# Benchmarks and Tree Search for Multimodal LLM Web Agents
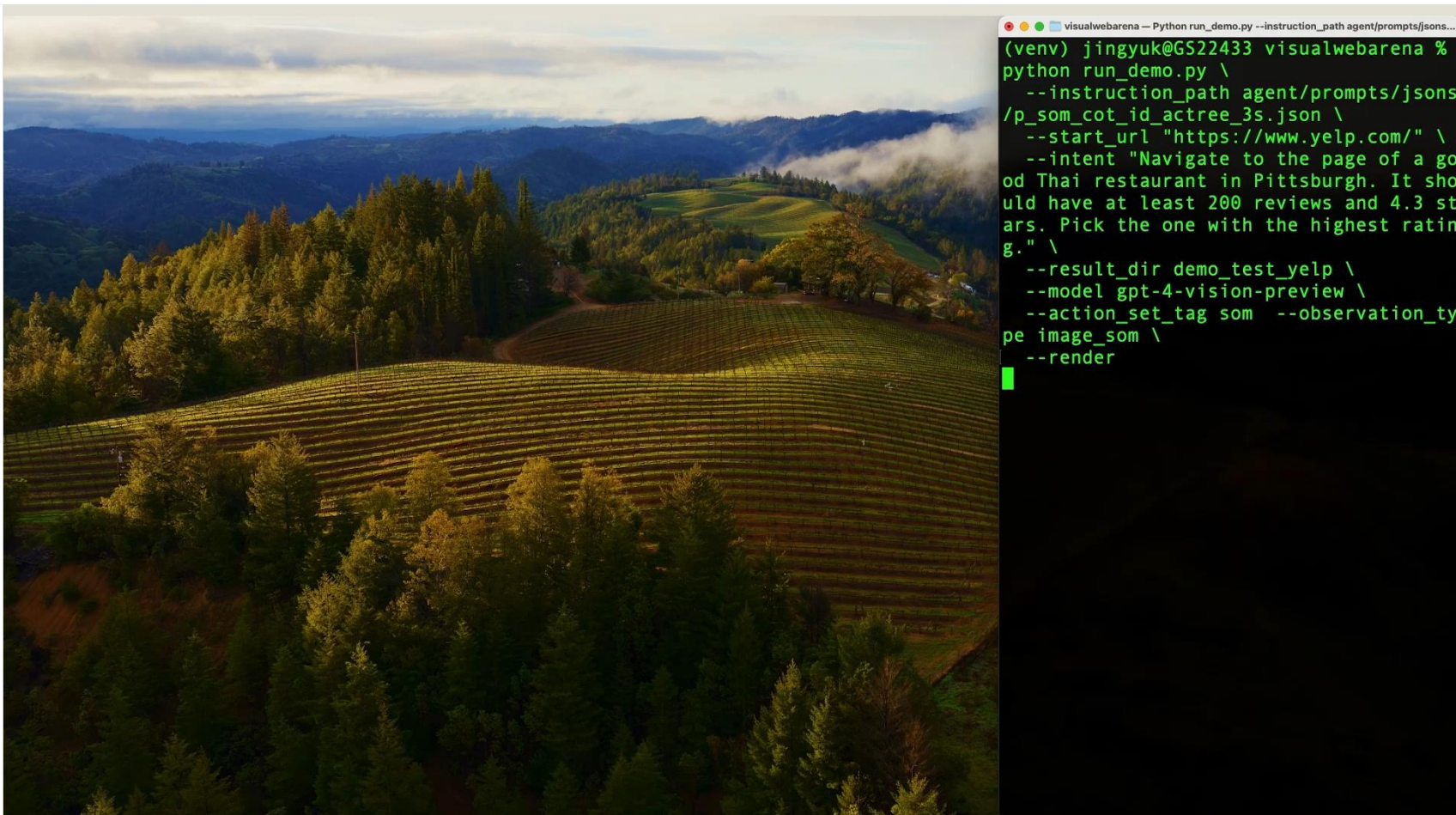
Daniel Fried

Language Technologies Institute

Carnegie Mellon University

**Task**: Navigate to the page of a good Thai restaurant in Pittsburgh. It should have at least 200 reviews and 4.3 stars. Pick the one with the highest rating.

# Why Web Agents?

# Why Web Agents?

# Web Agent Benchmarks

# Simulators (with simplified sites and tasks)

## World of Bits: An Open-Domain Platform for Web-Based Agents

Tianlin (Tim) Shi [1,2]   Andrej Karpathy [2]   Linxi (Jim) Fan [1]   Jonathan Hernandez [2]   Percy Liang [1]



‣ Simplified tasks, but comes with a simulator (can act, explore do RL).

‣ Introduced in 2017, remained challenging for some time afterward!

# Real Sites and Tasks (but without simulators)

**MIND2WEB: Towards a Generalist Agent for the Web**

Xiang Deng*    Yu Gu    Boyuan Zheng    Shijie Chen
Samuel Stevens    Boshi Wang    Huan Sun*    Yu Su*
The Ohio State University
https://osu-nlp-group.github.io/Mind2Web

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

▸ 2000 crowdsourced tasks and trajectories from ~100 real and diverse websites.

▸ Can perform reference-based evaluation, but lacks a simulator to allow agents to act freely.

# Simulators with Real-World Sites



Shuyan Zhou    Frank Xu    Jing Yu Koh



**VisualWebArena Sites**

**WebArena** (Zhou*, Xu* et al., ICLR 2024)
Standalone, self-hostable web environments



"Help me make a post selling this item and navigate to it. Price it at $10 cheaper than the most similar item on the site."

"Navigate to the comments section of the latest image post in the /f/Art subreddit that contains animals."

"Buy the cheapest color photo printer and send it to Emily's place (as shown in the image)."

Webpage        Task Specification

LLM / VLM Agent

click [1602]

**VisualWebArena** (Koh et al., ACL 2024)
Benchmark for *multimodal* web agents

# Reproducible Environments

POMDP environment: $\mathcal{E} = \langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{T} \rangle$

Observations $\mathcal{O}$

Actions $\mathcal{A}$

| Action Type $a$ | Description |
|---|---|
| click [elem] | Click on element elem. |
| hover [elem] | Hover on element elem. |
| type [elem] [text] | Type text on element elem. |
| press [key_comb] | Press a key combination. |
| new_tab | Open a new tab. |
| tab_focus [index] | Focus on the i-th tab. |
| tab_close | Close current tab. |
| goto [url] | Open url. |
| go_back | Click the back button. |
| go_forward | Click the forward button. |
| scroll [up\|down] | Scroll up or down the page. |
| stop [answer] | End the task with an optional output. |

Execution-based evaluation (reward) function: $r(\mathbf{a}, \mathbf{s})$

# Execution-Based Evaluation

| Webpage / Input Image(s) | Example Intent | Reward Function $r(s, a)$ Implementation |
|---|---|---|
|  | What is the ISIN of the company that occupies the largest portion in Warren Buffet's portfolio? Answer using the information from the Wikipedia site in the second tab. | `exact_match(â, "US0378331005")` |
|  | Add something like what the man is wearing to my wish list. | `url="/wishlist"`<br>`locator(".wishlist .product-image-photo")`<br>`eval_vqa(s, "Is this a polo shirt? (yes/no)", "yes")`<br>`eval_vqa(s, "Is this shirt green? (yes/no)", "yes")` |

# **VisualWebArena:** Task Distribution



Distribution of Tasks Across Sites



Distribution of Tasks by Difficulty

**Task:** "Please add to my shopping cart all the items from this page that can connect these devices from the two images."

Jing Yu Koh

# Building Multimodal LLM Agents

# (Multimodal) LLMs as Agents



ReAct prompting [Yao et al. 2022] with Set-of-Marks visual representation [Yang. 2014]

**Task:** Make a reservation at Pusadee's Garden for 2 people on the earliest date for dinner. Use my name JY Koh and phone number 650-555-5555.

```
(venv) jingyuk@GS22433 visualwebarena %
python run_demo.py \
  --instruction_path agent/prompts/jsons
/p_som_cot_id_actree_3s.json \
  --start_url "https://www.google.com/"
\
  --intent "Make a reservation at Pusade
e's Garden for 2 people on the earliest
date at any time. Use my name JY Koh and
 phone number 650-555-5555." \
  --result_dir demo_test_yelp \
  --model gpt-4-vision-preview \
  --action_set_tag som  --observation_ty
pe image_som \
  --render
```

# (Multimodal) LLMs as Agents



Success Rates of GPT-4 on VWA

# Common Failure Modes

- **Failures in visual processing**
  - Clicking the wrong item
  - Identifying specific items in complex webpages
  - Spatial reasoning ("what are the prices of products in the first row?")

- **Long horizon reasoning and planning**
  - Getting stuck in loops
  - Correctly performing tasks but undoing them

# Exponential Error Compounding in Agents

| Accuracy @ k steps: | | | | |
|---|---|---|---|---|
| **1 (single step)** | **5** | **10** | **30** | **50** |
| 90% | 59.05% | 34.87% | 4.24% | 0.52% |
| 95% | 77.38% | 59.87% | 21.46% | 7.69% |
| 99% | 95.10% | 90.44% | 73.97% | 60.50% |
| 99.9% | 99.50% | 99.00% | 97.04% | 95.12% |
| 99.99% | 99.95% | 99.90% | 99.70% | 99.50% |

# Local Decisions; Global Consequences

*"Add this and coconut milk to my cart"*



click [**31**]

0.5

0.2

...

...

0.2

click [**24**]

0.1

✓

# Local Decisions; Global Consequences

*"Add this and coconut milk to my cart"*



**0.5**

**0.2**

**...**

**...**

**0.2**

**0.1**

✓

# Search By Repeated Sampling

1. Sample actions from the language model until [STOP].



See also Pan et al. 2024,
Autonomous Evaluation and Refinement of Digital Agents

# Search By Repeated Sampling

1. Sample actions from the language model until [STOP].
2. Evaluate the resulting trajectory
3. Repeat?



✗

See also Pan et al. 2024,
Autonomous Evaluation and Refinement of Digital Agents

# Search By Repeated Sampling

1. Sample actions from the language model until [STOP].
2. Evaluate the resulting trajectory
3. Repeat?

See also Pan et al. 2024,
Autonomous Evaluation and Refinement of Digital Agents

# Search By Repeated Sampling

1. Sample actions from the language model until [STOP].
2. Evaluate the resulting trajectory
3. Repeat?

See also Pan et al. 2024,
Autonomous Evaluation and Refinement of Digital Agents

# Search By Repeated Sampling



Oracle vs Predicted Scores

Repeated
sampling helps!

- But the space is exponentially large. Can we guide exploration?
- Key idea of our approach: apply value function to intermediate nodes.

# Our Method: Tree Search

- Best-first search algorithm
- Ingredients:
  - Baseline agent to propose actions.
  - Way to backtrack in the environment.
  - A **value function** to score and rerank candidate states.
    - In this work, we prompt GPT-4o to act as an evaluator.



v = 0.35

v = 0.55

Koh et al. 2024,
Tree Search for Language Model Agents

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

## GPT-4o Agent

## GPT-4o Agent + Search

Starting State

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

## GPT-4o Agent

①

## GPT-4o Agent + Search

①

v = 0.5

Starting State

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

① Step sequence

**Backtracking**

v = 1.0
State values

## GPT-4o Agent

① 

## GPT-4o Agent + Search

①

v = 0.5

Starting State

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

① Step sequence

→ Backtracking

v = 1.0 State values

**GPT-4o Agent**

**GPT-4o Agent + Search**

Starting State

v = 0.5

v = 0.45

v = 0.45

v = 0.45

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 → 2

**GPT-4o Agent + Search**

Starting State

1 v = 0.5

2 v = 0.45

v = 0.45

v = 0.45

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 2

**GPT-4o Agent + Search**

Starting State

1 v = 0.5

2 v = 0.45

v = 0.4

...

v = 0.45

v = 0.45

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

## GPT-4o Agent

## GPT-4o Agent + Search

Starting State

v = 0.5
v = 0.45
v = 0.4
v = 0.45
v = 0.45

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

① Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

**GPT-4o Agent + Search**

Starting State

v = 0.5

v = 0.45

v = 0.4

v = 0.45

v = 0.5

v = 0.45

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

① Step sequence    $v = 1.0$   State values

→ Backtracking

## GPT-4o Agent

## GPT-4o Agent + Search

Starting State

$v = 0.5$   $v = 0.45$   $v = 0.4$

$v = 0.45$   $v = 0.5$

$v = 0.45$

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

① Step sequence

v = 1.0 State values

→ Backtracking

## GPT-4o Agent

Failure

## GPT-4o Agent + Search

Starting State

v = 0.5

v = 0.45

v = 0.4

v = 0.45

v = 0.5

v = 0.55

v = 0.45

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"

**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✖ **Failure**

**GPT-4o Agent + Search**

Starting State

1 v = 0.5
2 v = 0.45 ... v = 0.4
3 v = 0.45
4 v = 0.5
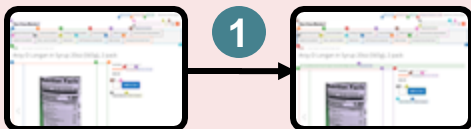5 v = 0.55 v = 0.3

v = 0.45

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
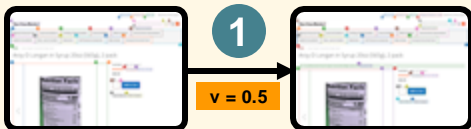
**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 2 3 4 ✗ **Failure**

**GPT-4o Agent + Search**

Starting State

v = 0.5 1

v = 0.45 2 v = 0.4

v = 0.45 3 v = 0.5 4 v = 0.55 5 v = 0.3

v = 0.45 6

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
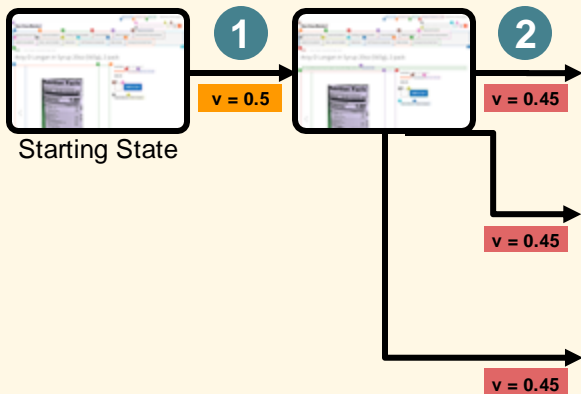
**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✖ **Failure**

**GPT-4o Agent + Search**

Starting State

v = 0.5 → 1 → v = 0.45 → 2 → v = 0.4

v = 0.45 → 3 → v = 0.45 → 4 → v = 0.5 → 5 → v = 0.55 → v = 0.3

v = 0.45 → 6 → v = 0.55

v = 0.5

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
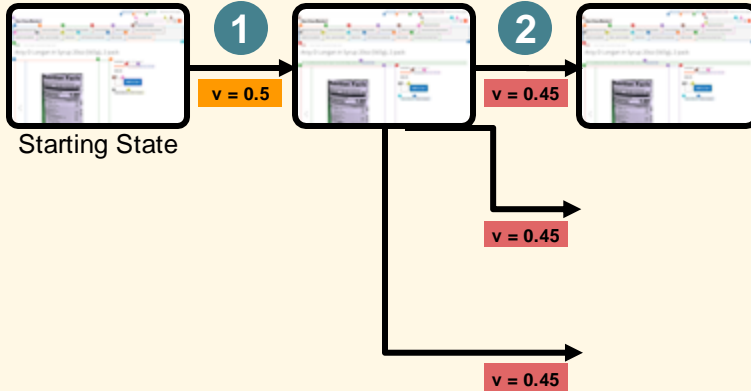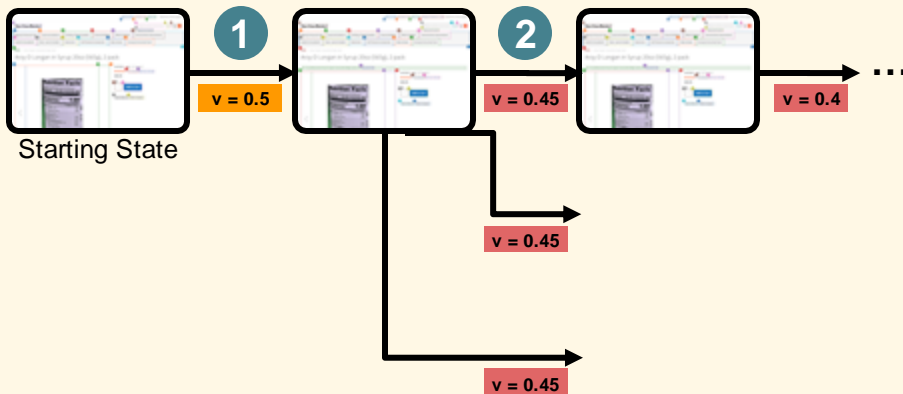
**Legend**

1 — Step sequence

Backtracking

v = 1.0 — State values

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✖ **Failure**

**GPT-4o Agent + Search**

1 — v = 0.5

Starting State

2 — v = 0.45 ... v = 0.4

3 — v = 0.45 4 — v = 0.5 5 — v = 0.55 ... v = 0.3
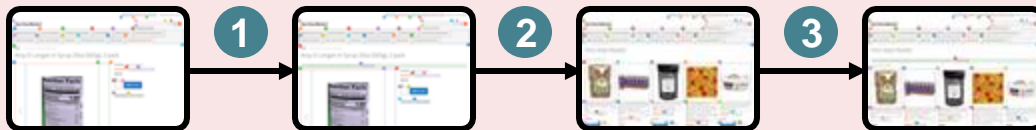
6 — v = 0.45 7 — v = 0.55

v = 0.5

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
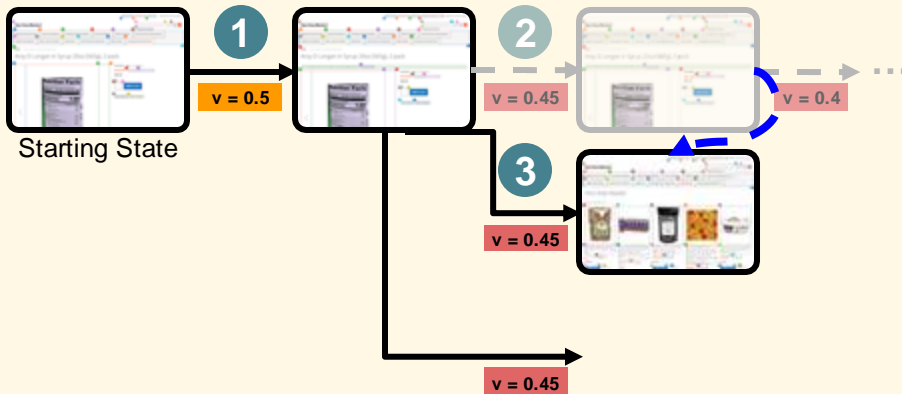
**Legend**

- ① Step sequence
- → Backtracking
- $v = 1.0$ State values

**GPT-4o Agent**

Failure

**GPT-4o Agent + Search**

Starting State

$v = 0.5$

$v = 0.45$ $v = 0.4$

$v = 0.45$ $v = 0.5$ $v = 0.55$ $v = 0.3$

$v = 0.45$ $v = 0.55$ $v = 0.68$

$v = 0.5$

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
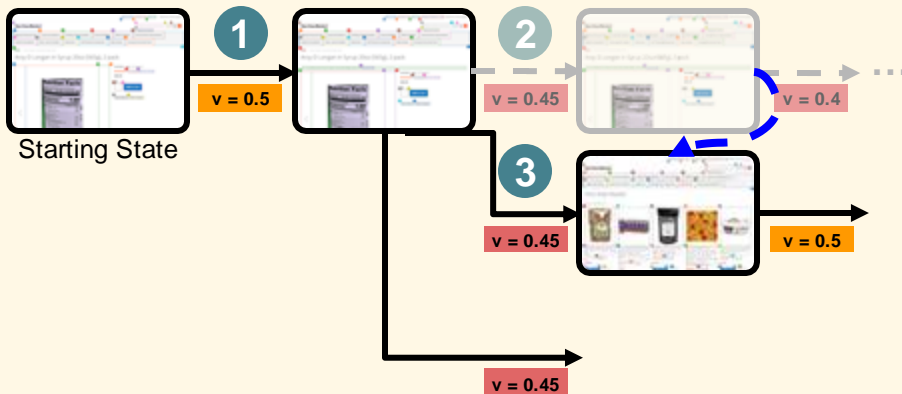
**Legend**

1 Step sequence

v = 1.0 State values

→ Backtracking

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✖ **Failure**

**GPT-4o Agent + Search**

Starting State

1 v = 0.5

2 v = 0.45

v = 0.4

3 v = 0.45

4 v = 0.5

5 v = 0.55

v = 0.3

6 v = 0.45

7 v = 0.55

8 v = 0.68

v = 0.5

**Legend**

1 Step sequence

$v = 1.0$ State values

→ Backtracking

**Task Instruction ($I$):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
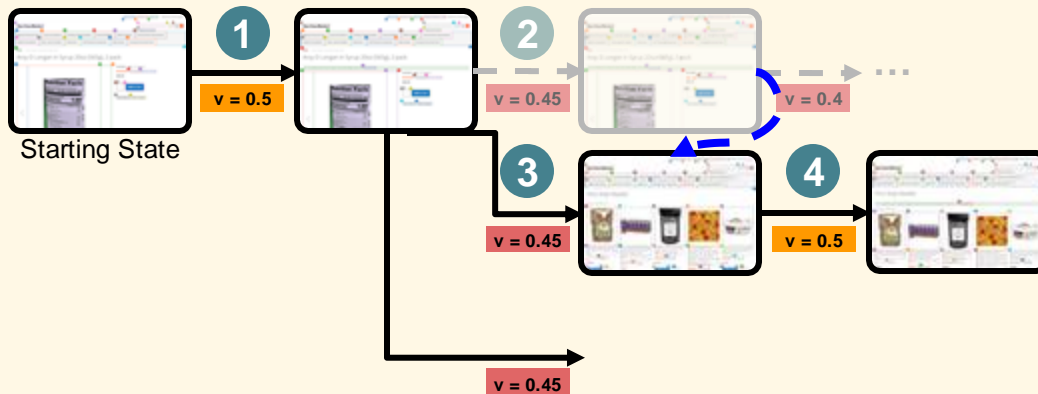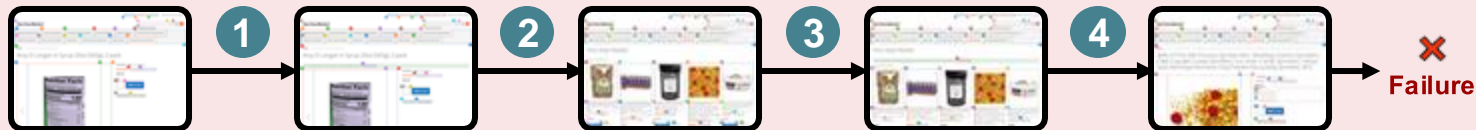
**GPT-4o Agent**

Starting State → 1 → 2 → 3 → 4 → **Failure**

**GPT-4o Agent + Search**

Starting State

1 → $v = 0.5$

2 → $v = 0.45$ ... $v = 0.4$

3 → $v = 0.45$

4 → $v = 0.5$

5 → $v = 0.55$ ... $v = 0.3$

6 → $v = 0.45$

7 → $v = 0.55$

8 → $v = 0.68$ → ... $v = 0.2$

$v = 0.5$

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
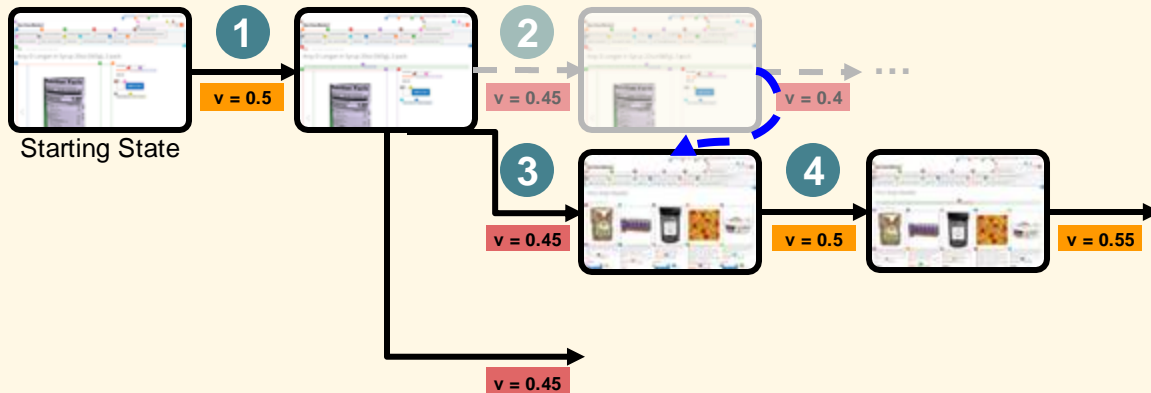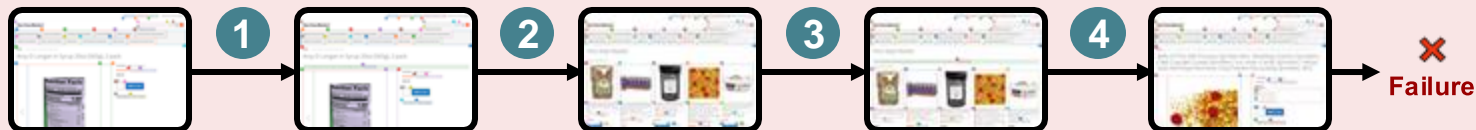
**Legend**

1 Step sequence

$v = 1.0$ State values

→ Backtracking

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✖ **Failure**

**GPT-4o Agent + Search**

Starting State

1 $v = 0.5$

2 $v = 0.45$ $v = 0.4$

3 $v = 0.45$ 4 $v = 0.5$ 5 $v = 0.55$ $v = 0.3$

6 $v = 0.45$ 7 $v = 0.55$ 8 $v = 0.68$ $v = 0.2$

9 $v = 0.5$

**Task Instruction ( $I$ ):** "Can you add this and the other canned fruit (of the same brand) that looks like this, but red instead of brown to the comparison page?"
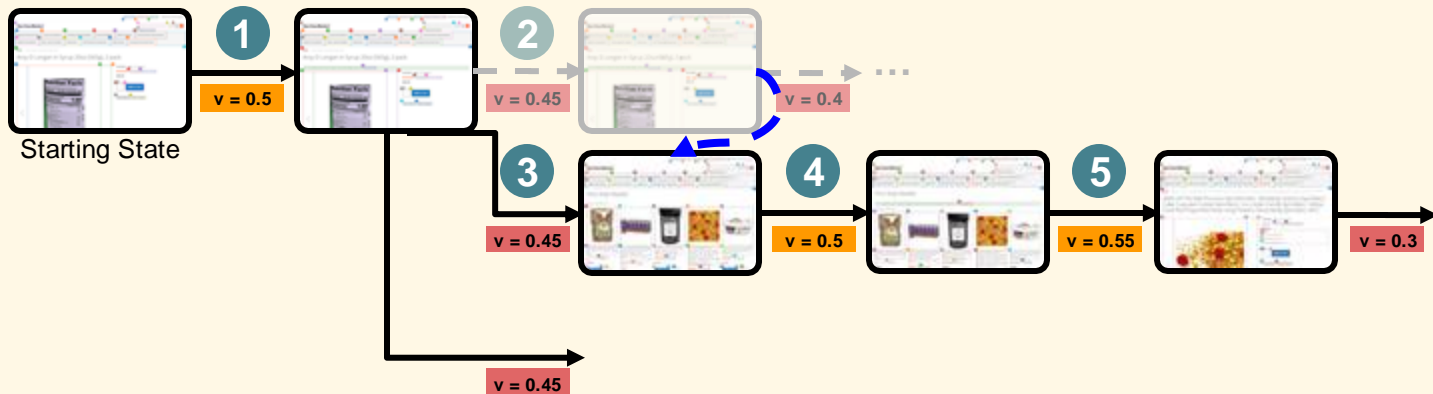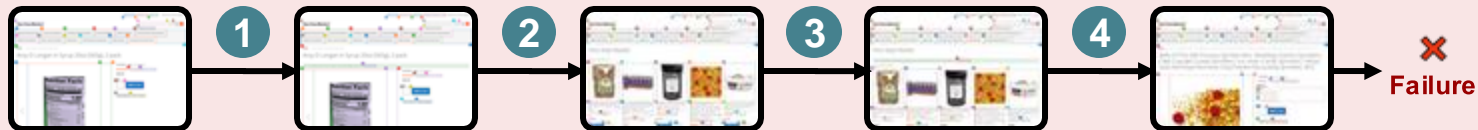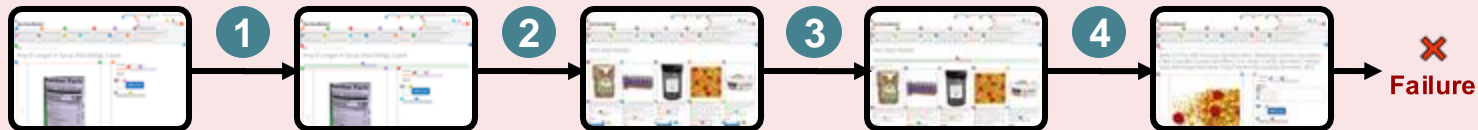
**Legend**

1 Step sequence

$v = 1.0$ State values

→ Backtracking

**GPT-4o Agent**

1 → 2 → 3 → 4 → ✗ **Failure**

**GPT-4o Agent + Search**

Starting State

1 → $v = 0.5$

2 $v = 0.45$ $v = 0.4$

3 $v = 0.45$ 4 $v = 0.5$ 5 $v = 0.55$ $v = 0.3$

6 $v = 0.45$ 7 $v = 0.55$ 8 $v = 0.68$ $v = 0.2$

9 $v = 0.5$ 10 $v = 1.0$ ✅ **Success**

# Results



VWA Success Rate
- Baseline
- w/ Search

| Llama-3-Instruct-70B | GPT-4o |
| --- | --- |
| 7.6% / 16.7% | 18.9% / 26.4% |

WA Success Rate
- Baseline
- w/ Search

| Llama-3-Instruct-70B | GPT-4o |
| --- | --- |
| 7.6% / 10.1% | 15.0% / 19.2% |

# Ablations

| Depth $d$ | Branch $b$ | SR (↑) | Δ |
|:---:|:---:|:---:|:---:|
| 0 | 1 | 24.5% | 0% |
| 1 | 3 | 26.0% | +6% |
| | 5 | 32.0% | +31% |
| 2 | 3 | 31.5% | +29% |
| | 5 | 35.0% | +43% |
| 3 | 5 | 35.5% | +45% |
| 5 | 5 | **37.0%** | +51% |

Success rate (SR) and relative change over the baseline (Δ) on a subset of 200 VWA tasks with varying search depth ($d$) and branching factor ($b$). $d = 0$ indicates no search is performed. All methods use a max search budget $c = 20$.

# Ablations

- Having a good value function is essential!
- There is still a lot of headroom for improving both the base agent policy, and the value function

| Value Function | SR ($\uparrow$) |
|---|---|
| None (no search) | 24.5% |
| LLaVA-v1.6-34B | 30.0% |
| GPT-4o (no SC) | 28.5% |
| GPT-4o | 37.0% |
| Groundtruth | 43.5% |

Table 3: Success rate of the GPT-4o agent with different value functions.

# Qualitative Results



**Task Instruction ($I$):** "I recall seeing this exact item on the site, help me find the most recent post of it. I recall seeing it in either the Collectibles or Antiques section."

## GPT-4o Agent + Search

Starting State

1 — v = 0.5
2 — v = 0.53 — v = 0.45 ...
3 — v = 0.53 — 4 — v = 0.63 — v = 0.2 ...
5 — v = 0.63 — 6 — v = 1.0 — ✅ Success

**Legend:** ① Search sequence  - ➤ Backtracking  v = 1.0 State values

# Limitations

- Search is slow
  - We implemented backtracking in a relatively naive way (store actions in a queue, take them again to get to the original state)
  - See Chen et al. 2024, *When is Tree Search Useful?*

- Dealing with destructive actions
  - Some things on the web are very difficult to undo, e.g., ordering an item

# Future Work

- Search as a policy improvement function.

- What's the value of value functions?

- What if we don't have a perfect simulator?

- Search to improve safety.

# Adversarial Attacks on Multimodal Agents

**Chen Henry Wu, Jing Yu Koh, Ruslan Salakhutdinov, Daniel Fried, Aditi Raghunathan**
Carnegie Mellon University
{chenwu2,jingyuk,rsalakhu,dfried,aditirag}@cs.cmu.edu

# Collaborators


Vikram Duvvur

Po Yu Huang

Lawrence Jang

Jing Yu Koh

Ming Chong Lim

Robert Lo

Stephen McAleer

Graham Neubig

Russ Salakhutdinov

Frank Xu

Shuyan Zhou

# Thanks!

{dfried,jingyuk,rsalakhu}@cs.cmu.edu

jykoh.com/vwa // jykoh.com/search-agents

**Reddit**                    **Wikipedia**

**Task:** "What is the 2022 total nominal GDP of the area that produces most sugarcane in the year of 2021? (in billion)?"

# Results

| Benchmark | Agent Model | Max Actions | No Search | + Search | Relative Change |
|---|---|---|---|---|---|
| VisualWebArena | GPT-4o + SoM [1] | 30 | 19.8% | - | - |
| | Llama-3-70B-Instruct [1] | | 9.8% | - | - |
| | Llama-3-70B-Instruct (ours) | 5 | 7.6% | 16.7% | +119.7% |
| | GPT-4o + SoM (ours) | | 18.9% | **26.4%** | +39.7% |
| WebArena | GPT-4o [2] | 30 | 13.1% | - | - |
| | GPT-4 + Reflexion [3] | | 15.6% | - | - |
| | AutoWebGLM [4] | | 18.2% | - | - |
| | AutoEval [3] | | 20.2% | - | - |
| | BrowserGym (GPT-4) [5] | | 23.5% | - | - |
| | SteP [6] | | **35.8%** | - | - |
| | GPT-4o (ours) | 5 | 15.0% | 19.2% | +28.0% |

# Baseline Agents

| Model Type | LLM Backbone | Visual Backbone | Inputs | Success Rate (↑) | | | |
|---|---|---|---|---|---|---|---|
| | | | | **Classifieds** | **Reddit** | **Shopping** | **Overall** |
| Text-only | LLaMA-2-70B | | Acc. Tree | 0.43% | 1.43% | 1.29% | 1.10% |
| | Mixtral-8x7B | | | 1.71% | 2.86% | 1.29% | 1.76% |
| | Gemini-Pro | - | | 0.85% | 0.95% | 3.43% | 2.20% |
| | GPT-3.5 | | | 0.43% | 0.95% | 3.65% | 2.20% |
| | GPT-4 | | | 5.56% | 4.76% | 9.23% | 7.25% |
| Caption-augmented | LLaMA-2-70B | BLIP-2-T5XL | Acc. Tree + Caps | 0.00% | 0.95% | 0.86% | 0.66% |
| | Mixtral-8x7B | BLIP-2-T5XL | | 1.28% | 0.48% | 2.79% | 1.87% |
| | GPT-3.5 | LLaVA-7B | | 1.28% | 1.43% | 4.08% | 2.75% |
| | GPT-3.5 | BLIP-2-T5XL | | 0.85% | 1.43% | 4.72% | 2.97% |
| | Gemini-Pro | BLIP-2-T5XL | | 1.71% | 1.43% | 6.01% | 3.85% |
| | GPT-4 | BLIP-2-T5XL | | 8.55% | 8.57% | 16.74% | 12.75% |
| Multimodal | IDEFICS-80B-Instruct | | Image + Caps + Acc. Tree | 0.43% | 0.95% | 0.86% | 0.77% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 4.29% | 8.15% | 6.04% |
| | GPT-4V | | | 8.12% | 12.38% | **19.74%** | 15.05% |
| Multimodal (SoM) | IDEFICS-80B-Instruct | | Image + Caps + SoM | 0.85% | 0.95% | 1.07% | 0.99% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 3.81% | 7.73% | 5.71% |
| | GPT-4V | | | **9.83%** | **17.14%** | 19.31% | **16.37%** |
| Human Performance | - | - | Webpage | 91.07% | 87.10% | 88.39% | 88.70% |

# Baseline Agents

| Model Type | LLM Backbone | Visual Backbone | Inputs | Success Rate (↑) | | | |
|---|---|---|---|---|---|---|---|
| | | | | **Classifieds** | **Reddit** | **Shopping** | **Overall** |
| Text-only | LLaMA-2-70B | - | Acc. Tree | 0.43% | 1.43% | 1.29% | 1.10% |
| | Mixtral-8x7B | | | 1.71% | 2.86% | 1.29% | 1.76% |
| | Gemini-Pro | | | 0.85% | 0.95% | 3.43% | 2.20% |
| | GPT-3.5 | | | 0.43% | 0.95% | 3.65% | 2.20% |
| | GPT-4 | | | 5.56% | 4.76% | 9.23% | 7.25% |
| Caption-augmented | LLaMA-2-70B | BLIP-2-T5XL | Acc. Tree + Caps | 0.00% | 0.95% | 0.86% | 0.66% |
| | Mixtral-8x7B | BLIP-2-T5XL | | 1.28% | 0.48% | 2.79% | 1.87% |
| | GPT-3.5 | LLaVA-7B | | 1.28% | 1.43% | 4.08% | 2.75% |
| | GPT-3.5 | BLIP-2-T5XL | | 0.85% | 1.43% | 4.72% | 2.97% |
| | Gemini-Pro | BLIP-2-T5XL | | 1.71% | 1.43% | 6.01% | 3.85% |
| | GPT-4 | BLIP-2-T5XL | | 8.55% | 8.57% | 16.74% | 12.75% |
| Multimodal | IDEFICS-80B-Instruct | | Image + Caps + Acc. Tree | 0.43% | 0.95% | 0.86% | 0.77% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 4.29% | 8.15% | 6.04% |
| | GPT-4V | | | 8.12% | 12.38% | **19.74%** | 15.05% |
| Multimodal (SoM) | IDEFICS-80B-Instruct | | Image + Caps + SoM | 0.85% | 0.95% | 1.07% | 0.99% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 3.81% | 7.73% | 5.71% |
| | GPT-4V | | | **9.83%** | **17.14%** | 19.31% | **16.37%** |
| Human Performance | - | - | Webpage | 91.07% | 87.10% | 88.39% | 88.70% |

# Baseline Agents

| Model Type | LLM Backbone | Visual Backbone | Inputs | Success Rate (↑) | | | |
|---|---|---|---|---|---|---|---|
| | | | | **Classifieds** | **Reddit** | **Shopping** | **Overall** |
| Text-only | LLaMA-2-70B | - | Acc. Tree | 0.43% | 1.43% | 1.29% | 1.10% |
| | Mixtral-8x7B | | | 1.71% | 2.86% | 1.29% | 1.76% |
| | Gemini-Pro | | | 0.85% | 0.95% | 3.43% | 2.20% |
| | GPT-3.5 | | | 0.43% | 0.95% | 3.65% | 2.20% |
| | GPT-4 | | | 5.56% | 4.76% | 9.23% | 7.25% |
| Caption-augmented | LLaMA-2-70B | BLIP-2-T5XL | Acc. Tree + Caps | 0.00% | 0.95% | 0.86% | 0.66% |
| | Mixtral-8x7B | BLIP-2-T5XL | | 1.28% | 0.48% | 2.79% | 1.87% |
| | GPT-3.5 | LLaVA-7B | | 1.28% | 1.43% | 4.08% | 2.75% |
| | GPT-3.5 | BLIP-2-T5XL | | 0.85% | 1.43% | 4.72% | 2.97% |
| | Gemini-Pro | BLIP-2-T5XL | | 1.71% | 1.43% | 6.01% | 3.85% |
| | GPT-4 | BLIP-2-T5XL | | 8.55% | 8.57% | 16.74% | 12.75% |
| Multimodal | IDEFICS-80B-Instruct | | Image + Caps + Acc. Tree | 0.43% | 0.95% | 0.86% | 0.77% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 4.29% | 8.15% | 6.04% |
| | GPT-4V | | | 8.12% | 12.38% | **19.74%** | 15.05% |
| Multimodal (SoM) | IDEFICS-80B-Instruct | | Image + Caps + SoM | 0.85% | 0.95% | 1.07% | 0.99% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 3.81% | 7.73% | 5.71% |
| | GPT-4V | | | **9.83%** | **17.14%** | 19.31% | **16.37%** |
| Human Performance | - | - | Webpage | 91.07% | 87.10% | 88.39% | 88.70% |

# Baseline Agents

| Model Type | LLM Backbone | Visual Backbone | Inputs | Success Rate (↑) | | | |
|---|---|---|---|---|---|---|---|
| | | | | **Classifieds** | **Reddit** | **Shopping** | **Overall** |
| Text-only | LLaMA-2-70B | | Acc. Tree | 0.43% | 1.43% | 1.29% | 1.10% |
| | Mixtral-8x7B | | | 1.71% | 2.86% | 1.29% | 1.76% |
| | Gemini-Pro | - | | 0.85% | 0.95% | 3.43% | 2.20% |
| | GPT-3.5 | | | 0.43% | 0.95% | 3.65% | 2.20% |
| | GPT-4 | | | 5.56% | 4.76% | 9.23% | 7.25% |
| Caption-augmented | LLaMA-2-70B | BLIP-2-T5XL | Acc. Tree + Caps | 0.00% | 0.95% | 0.86% | 0.66% |
| | Mixtral-8x7B | BLIP-2-T5XL | | 1.28% | 0.48% | 2.79% | 1.87% |
| | GPT-3.5 | LLaVA-7B | | 1.28% | 1.43% | 4.08% | 2.75% |
| | GPT-3.5 | BLIP-2-T5XL | | 0.85% | 1.43% | 4.72% | 2.97% |
| | Gemini-Pro | BLIP-2-T5XL | | 1.71% | 1.43% | 6.01% | 3.85% |
| | GPT-4 | BLIP-2-T5XL | | 8.55% | 8.57% | 16.74% | 12.75% |
| Multimodal | IDEFICS-80B-Instruct | | Image + Caps + Acc. Tree | 0.43% | 0.95% | 0.86% | 0.77% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 4.29% | 8.15% | 6.04% |
| | GPT-4V | | | 8.12% | 12.38% | **19.74%** | 15.05% |
| Multimodal (SoM) | IDEFICS-80B-Instruct | | Image + Caps + SoM | 0.85% | 0.95% | 1.07% | 0.99% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 3.81% | 7.73% | 5.71% |
| | GPT-4V | | | **9.83%** | **17.14%** | 19.31% | **16.37%** |
| Human Performance | - | - | Webpage | 91.07% | 87.10% | 88.39% | 88.70% |

# Baseline Agents

| Model Type | LLM Backbone | Visual Backbone | Inputs | Success Rate (↑) | | | |
|---|---|---|---|---|---|---|---|
| | | | | **Classifieds** | **Reddit** | **Shopping** | **Overall** |
| Text-only | LLaMA-2-70B | | Acc. Tree | 0.43% | 1.43% | 1.29% | 1.10% |
| | Mixtral-8x7B | | | 1.71% | 2.86% | 1.29% | 1.76% |
| | Gemini-Pro | - | | 0.85% | 0.95% | 3.43% | 2.20% |
| | GPT-3.5 | | | 0.43% | 0.95% | 3.65% | 2.20% |
| | GPT-4 | | | 5.56% | 4.76% | 9.23% | 7.25% |
| Caption-augmented | LLaMA-2-70B | BLIP-2-T5XL | Acc. Tree + Caps | 0.00% | 0.95% | 0.86% | 0.66% |
| | Mixtral-8x7B | BLIP-2-T5XL | | 1.28% | 0.48% | 2.79% | 1.87% |
| | GPT-3.5 | LLaVA-7B | | 1.28% | 1.43% | 4.08% | 2.75% |
| | GPT-3.5 | BLIP-2-T5XL | | 0.85% | 1.43% | 4.72% | 2.97% |
| | Gemini-Pro | BLIP-2-T5XL | | 1.71% | 1.43% | 6.01% | 3.85% |
| | GPT-4 | BLIP-2-T5XL | | 8.55% | 8.57% | 16.74% | 12.75% |
| Multimodal | IDEFICS-80B-Instruct | | Image + Caps + Acc. Tree | 0.43% | 0.95% | 0.86% | 0.77% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 4.29% | 8.15% | 6.04% |
| | GPT-4V | | | 8.12% | 12.38% | **19.74%** | 15.05% |
| Multimodal (SoM) | IDEFICS-80B-Instruct | | Image + Caps + SoM | 0.85% | 0.95% | 1.07% | 0.99% |
| | CogVLM | | | 0.00% | 0.48% | 0.43% | 0.33% |
| | Gemini-Pro | | | 3.42% | 3.81% | 7.73% | 5.71% |
| | GPT-4V | | | **9.83%** | **17.14%** | 19.31% | **16.37%** |
| Human Performance | - | - | Webpage | 91.07% | 87.10% | 88.39% | 88.70% |

# Ablations

- Search helps more for medium difficulty (4-9 actions to solve) tasks
- May be related to our own compute limitations: we fixed the max search depth to be 5 in our experiments
- Increasing the depth is likely to help hard tasks

| Difficulty | No Search | Search | Δ |
|---|---|---|---|
| easy | 34.2% | 42.3% | +24% |
| medium | 12.7% | 22.2% | +75% |
| hard | 10.2% | 14.9% | +47% |

Table 3: Success rates and relative change (Δ) of the GPT-4o agent on VWA tasks of different action difficulty levels.

# Analysis

- Consistent gains across all site types
- Value function is already fairly general for web tasks

| Website | No Search | Search | Δ |
|---|---|---|---|
| Classifieds | 18.4% | 26.5% | +44% |
| Reddit | 17.1% | 20.5% | +20% |
| Shopping | 20.0% | 29.0% | +45% |
| Overall | 18.9% | 26.4% | +40% |

Table 4: Success rates and relative change (Δ) of the GPT-4o agent on VWA websites.

| Website | No Search | Search | Δ |
|---|---|---|---|
| CMS | 11.0% | 16.5% | +50% |
| Map | 21.1% | 25.8% | +22% |
| Shopping | 24.0% | 28.1% | +17% |
| Reddit | 7.9% | 10.5% | +33% |
| Gitlab | 10.2% | 13.3% | +30% |
| Overall | 15.0% | 19.2% | +28% |

Table 5: Success rates and relative change (Δ) of the GPT-4o agent on WA websites.

# Value Model via Prompting

▶ Self-consistency chain-of-thought prompting (adapted from Pan et al. 2024), with 20 samples and values ranging from 0 to 1

system_message:
You are an expert in evaluating the performance of a web navigation agent. The agent is designed to help a human user navigate a website to complete a task. Given the user's intent, the agent's action history, the final state of the webpage, and the agent's response to the user, your goal is to decide whether the agent's execution is successful or not. If the current state is a failure but it looks like the agent is on the right track towards success, you should also output as such.

There are three types of tasks: