# WebCanvas

## Benchmarking Web Agents in Online Environments

Yichen Pan*♠, Dehan Kong*♠, Sida Zhou*♠, Cheng Cui*♠, Yifei Leng♠, Bing Jiang♠, Hangyu Liu♠, Yanyi Shang♠, Shuyan Zhou♣, Tongshuang Wu♣, Zhengyang Wu♠

♠iMean AI,   ♣Carnegie Mellon University       mail: dehan@imean.ai

• Web agent evaluations are either static or operate within a limited environment.

• We eventually need to connect the agent with the real world with **Total Web Dynamics**.

• We craft and host your agentic data for you to connect your offline agent with the online environment!

### We build better interface for you to talk to your data

**50%** in annotation

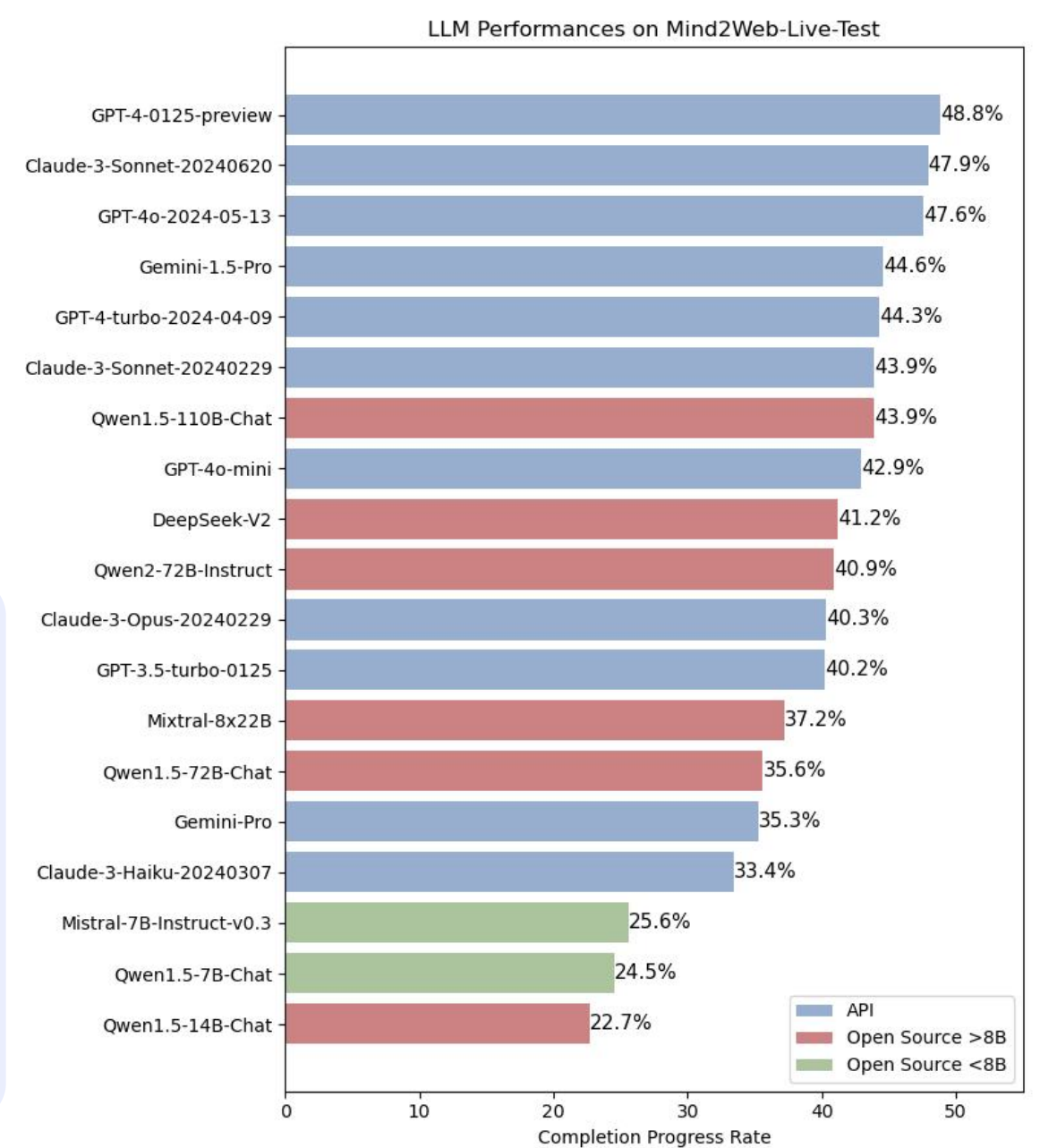**80%** in maintainance



Debug your trajectory

### Procedures to create your own benchmark online

→ Step1: Define your tasks.

→ Step2: Record your trajectories using iMean Builder (We support a broad range of action space, even with **data operations**).



→ Step3: Annotate the **Key Nodes** (static rule based evaluation for now).

→ Step4: Connect the data to **WebCanvas** by creating a challenge for the community (You can make it private during construction).



→ Step5: Evaluate your web agents online by easily integrating **WebCanvas** and accumulate insights of web agent performances as a community.



LLM Performances on Mind2Web-Live-Test

| Model | Completion Progress Rate |
|---|---|
| GPT-4-0125-preview | 48.8% |
| Claude-3-Sonnet-20240620 | 47.9% |
| GPT-4o-2024-05-13 | 47.6% |
| Gemini-1.5-Pro | 44.6% |
| GPT-4-turbo-2024-04-09 | 44.3% |
| Claude-3-Sonnet-20240229 | 43.9% |
| Qwen1.5-110B-Chat | 43.9% |
| GPT-4o-mini | 42.9% |
| DeepSeek-V2 | 41.2% |
| Qwen2-72B-Instruct | 40.9% |
| Claude-3-Opus-20240229 | 40.3% |
| GPT-3.5-turbo-0125 | 40.2% |
| Mixtral-8x22B | 37.2% |
| Qwen1.5-72B-Chat | 35.6% |
| Gemini-Pro | 35.3% |
| Claude-3-Haiku-20240307 | 33.4% |
| Mistral-7B-Instruct-v0.3 | 25.6% |
| Qwen1.5-7B-Chat | 24.5% |
| Qwen1.5-14B-Chat | 22.7% |

API / Open Source >8B / Open Source <8B

## Takeaways

• Better toolkits save you time in agentic data management.

• Key-node based metrics is essential for evaluating agents in the wild.

• Models finetuned on static datasets struggle to generalize in online environments a year later.

• Self reward doesn't help, but with human-labeled rewards as a reference, agents improve.

• Less capable models don't benefit from memory and ReAct in web tasks.

• Agent performance varies by domain, website and physical environment.

## Step Forward

| | Completion Rate | Task Success Rate | USD / Key Node Score |
|---|---|---|---|
| GPT-3.5-turbo | 42.5% | 17.3% | 0.092 |
| GPT-4o | 51.4% | 28.8% | 0.142 |
| GPT-4o mini | 42.9% | 21.5% | 0.004 |

Web agent efficiency should gain more attention!
GPT-4o mini is >30 times more cost efficient than GPT-4o

• Better observation: accurate and fast to compute

• Better conversion: human interface -> agent interface

• Dynamic evaluation functions

• Cloud environment: increase reliability

• Secured action of web agents

• Error handling and authorization

### OUR VISON IS TO CREATE THE ULTIMATE LIVE SANDBOX FOR WEB AGENT ADVANCEMENTS

⬇ ⬇ ⬇   WE WANT YOU TO JOIN WITH US AS A COMMUNITY!   ⬇ ⬇ ⬇

Paper📄          Github🐙          LiveDemo🔮