

DESIGNING, EVALUATING, AND IMPROVING LLM AGENTS AS SIMULATED USERS

Yuxuan Lu

Northeastern University | <https://yuxuan.lu> | lu.yuxuan@northeastern.edu

Jun 2026

INTRODUCTION

Yuxuan (Leo) Lu

<https://yuxuan.lu>

Ph.D. student, Northeastern University (Khoury) — advised by Prof. Dakuo Wang

Industry Experience: Applied Scientist Intern at Amazon (Sep 2024–present);
ML Research Intern at LinkedIn China (Jul 2022–May 2023).

Research Focus: I work at the intersection of NLP and HCI.

NLP / ML

Finetuned LLM Agent	ACL'26
WebServ	NeurIPS'26 [†]
Firefly	NeurIPS'26 [†]

HCI

UXAgent	CHI'25 LBW
+ full system	in submission
Agent A/B	CHI'26 Poster

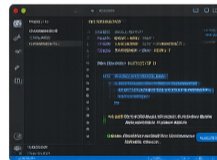
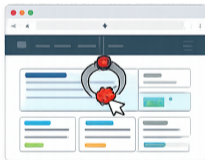
22 selected papers, 6 first-author (above; [†]under review).

16 co-authored — 11 accepted (ACL'26 ×3, ICLR'26, NAACL'24, EMNLP'23/'24, CHI'24/'25, CSCW'24/'26) and 5 under review (EMNLP'26 ×3, NeurIPS'26, +1 in submission).



opencode

 Claude Code



Greeting and quick check-in

- Shorten `under review (ACL'26)` to `ACL'26 under review`
- Even shorter: `ACL'26 submission`
- If you want to keep the full phrasing, split it across two lines inside the cell
- If needed, give the right column a bit more width with weighted X columns

I'd recommend this version first:

```
\begin{tblr}{
  colspec={X[1.4,1]X[0.9,r]},
  row{1}={font=\bfseries}
}
NLP & Venue/Status \\
SFT LLM Agent & ACL'26 under review \\
WebServ & ACL'26 under review \\
\end{tblr}
```

Why this helps:

- `ACL'26 under review` is shorter and reads more naturally on slides

Ask anything...

+

↑

Build GPT-5.4 Default

Review 3

Session changes

Unified Split Collapse all

```
files/talks/3_papers/talk.tex +8 -1
{0.48\linewidth} {0.48\linewidth}
26 unmodified lines
204 \centering 204 \centering
205
\includegraphics[wid \includegraphics[width
th=.8\linewidth] h=.8\linewidth]
{agents.png} {agents.png}
206 \end{figure} 206 \end{figure}
207 \end{frame} 207 \end{frame}
208
209 \begin{frame}[plain]
210 \begin{figure}[htp]
211 \centering
212
\includegraphics[width
h=.8\linewidth]
{opencode.png}
213 \end{figure}
214 \end{frame}
208 } 215 }
209
210 \begin{frame} 217 \begin{frame}{Today's
{Today's topic} topic}
211 one minute 218 one minute separate
separate slide slide
```

3 Changes

All files

- files
- talks
- 3_papers
- figures
- opencode.png A
- talk.pdf M
- talk.tex M



How can I train the crows in my neighborhood to bring me gifts?

This question does not make sense

This question should not be answered

Search results for: how to train crows to bring you gifts

Quotes



how to train crows to bring

Find in page



+ Add new quote

[How to Make Friends With Crows - PetHelpful](#)

If you did this a few times, your crows would learn your new place, but as I said, I'm not sure if they will follow or visit you there since it's probably not in their territory. The other option is simply to make new crow friends with the crows that live in your new neighborhood.

[Gifts From Crows | Outside My Window](#)

The partial piece of apple may have been left behind when the crow was startled rather than as a gift. If the crows bring bright objects you'll know for sure that it's a gift because it's not something they eat. Brandi Williams says: May 28, 2020 at 7:19 am.



[1] Gifts From Crows | Outside My Window (www.birdsoutsidemywindow.org)

Many animals give gifts to members of their own species but crows and other corvids are the only ones known to give gifts to humans.



Number of quote tokens left: 463

Number of actions left: 96



Done quoting! Write an answer

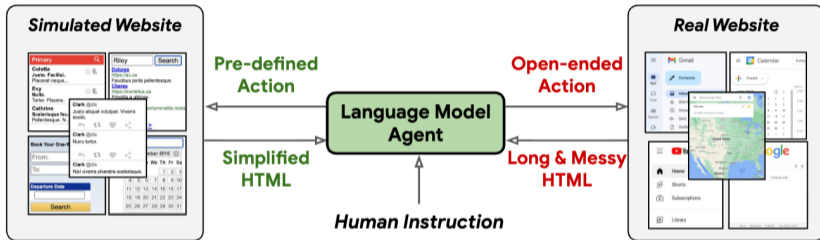
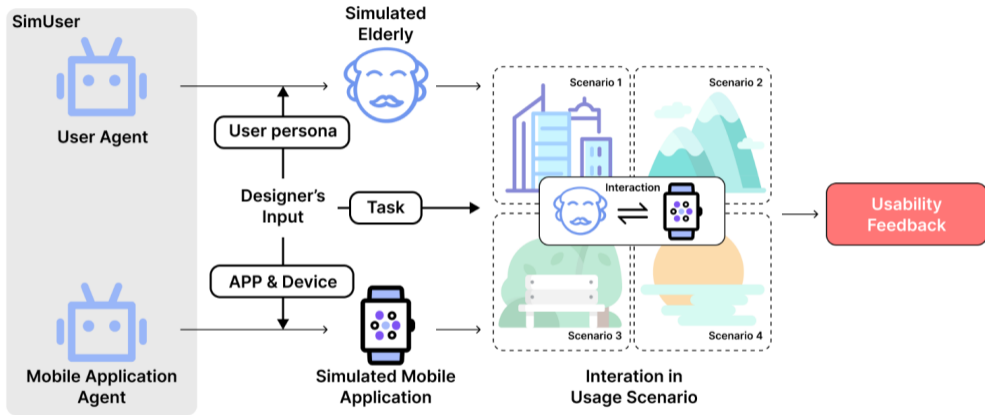
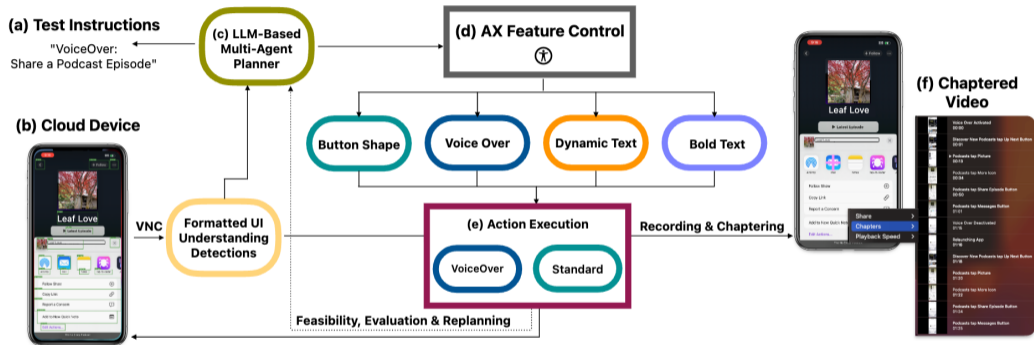


Figure 1: WebAgent, ICLR 2024







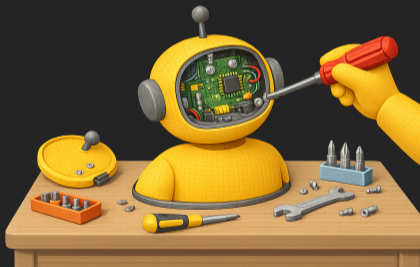
How can we design LLM agents as simulated users?



How can we evaluate LLM agents as simulated users?



How can we improve LLM agents as simulated users?



CASE STUDY 1: UXAGENT: A SYSTEM FOR SIMULATING USABILITY TESTING OF WEB DESIGN WITH LLM AGENTS

What's the worst nightmare for a researcher?

- Spending weeks on a study, only to find that the results are not significant
- Training a model for days, only to find a bug in your preprocessing pipeline
- The myth of “Reviewer 2”

- Spending weeks on a study, only to find that the results are not significant
- Training a model for days, only to find a bug in your preprocessing pipeline
- The myth of “Reviewer 2”
- Realizing that your experiment design is flawed one week before deadline

AS SOMEONE FROM BOTH NLP AND HCI, I THINK ...

	NLP	HCI
Experiment Subject	Models and Machines	Human Subjects
Experiment Design	Code and Data	Study Protocol
Experiment Cost	Money	Human Participants' Time
"Debugging" Method	Code Debugging	???

	NLP	HCI
Experiment Subject	Models and Machines	Human Subjects
Experiment Design	Code and Data	Study Protocol
Experiment Cost	Money	Human Participants' Time
"Debugging" Method	Code Debugging	???

Human Participants' Time is Valuable and Limited



Existing LLM Agent systems mostly works in **sandboxed environments**



**How can we better evaluate UX Research study design
before running the study?**

**How can we better evaluate usability testing study design
before running the study**

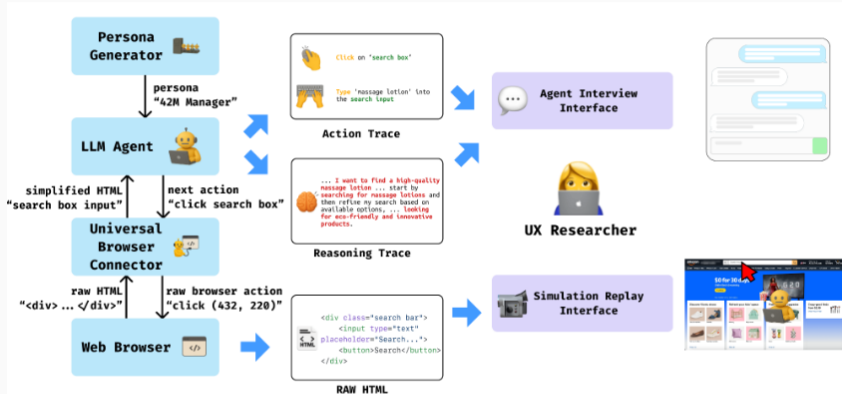


Figure 2: System Architecture of UXAgent

<https://www.youtube.com/watch?v=-2xpeJ04mRA>

- 1 Participant Recruitment
- 2 Survey
- 3 Review

Configure participant demographics Provide a questionnaire Confirm & Run

Recruitment Target Setting

URL of website being tested *

Number of Participants *

- +

Participant Task *

Example Persona *

Persona: Clara
Background:
Clara is a PhD student in Computer Science at a

(a) Participant Task Config

Demographics

Field Name	Value	Weight	Actions
Age	18-55	1	- + Remove Value Remove Field

+ Add Choice

Field Name	Value	Weight	Actions
Gender	Male	1	- + Remove Value
	Female	1	- + Remove Value
	Non-binary	1	- + Remove Value

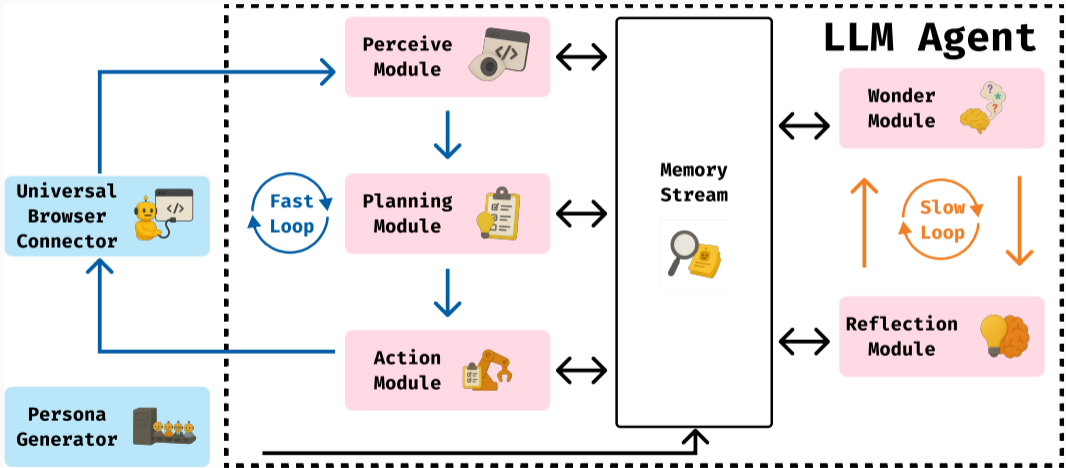
+ Add Choice Remove Field

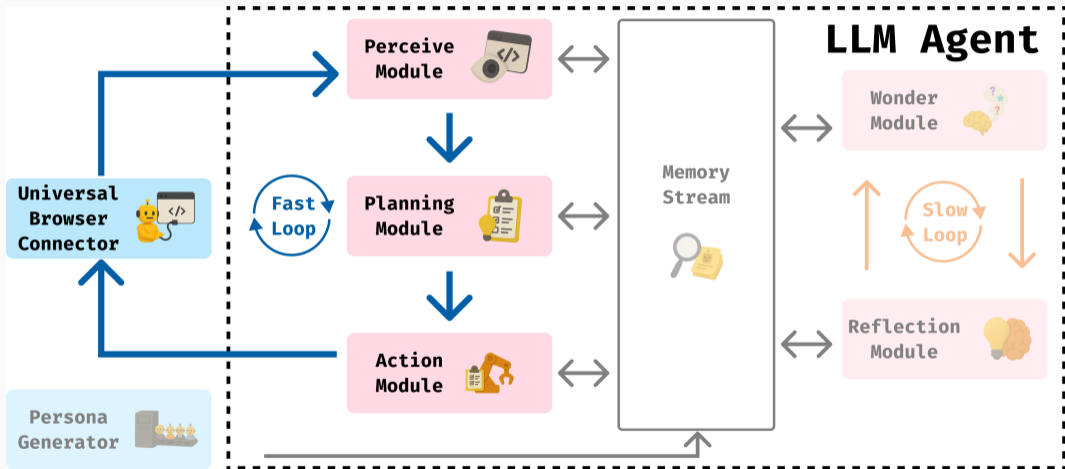
+ Add Field

Reset Form

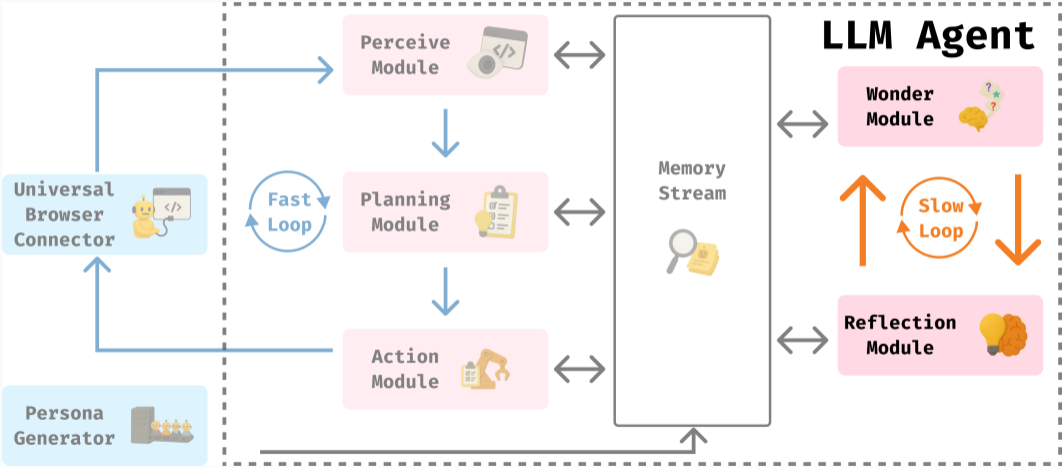
Next →

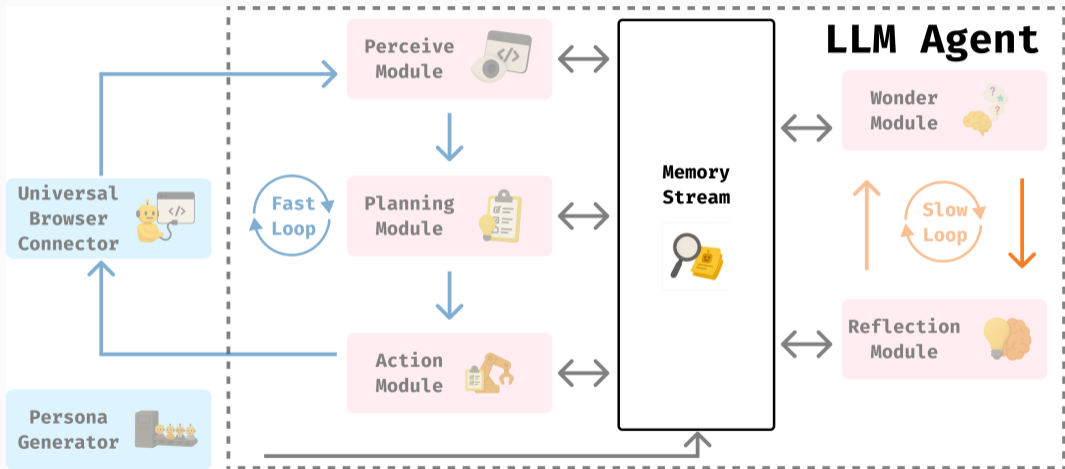
(b) Demographic Distribution Config






AGENT ARCHITECTURE - SLOW LOOP








Click on 'search box'

Type 'massage lotion' into the search input

Action Trace



... I want to find a high-quality massage lotion ... start by searching for massage lotions and then refine my search based on available options, ... looking for eco-friendly and innovative products.

Reasoning Trace



```
<div class="search bar">
  <input type="text"
  placeholder="Search...">
  <button>Search</button>
</div>
```

RAW HTML

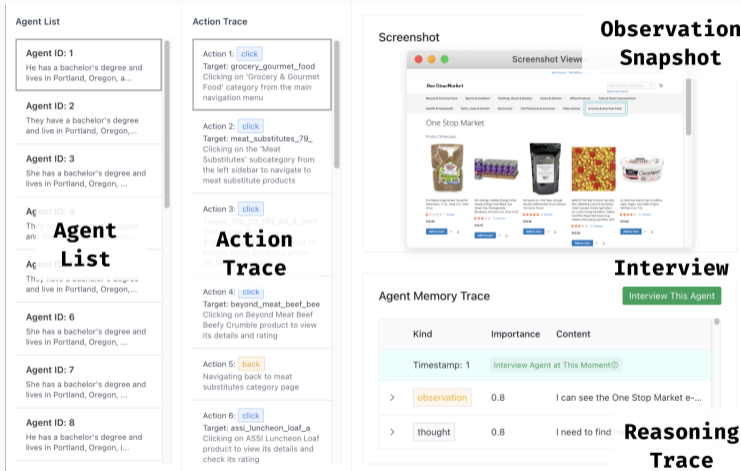


Figure 3: Result Viewer Interface

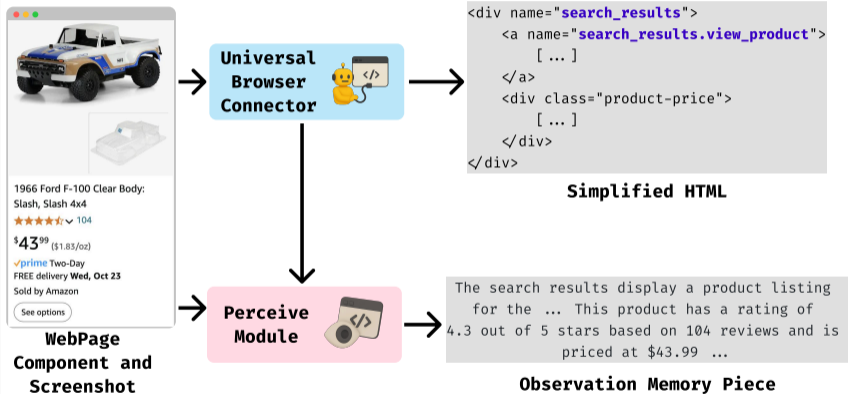


Figure 4: Universal Browser Connector

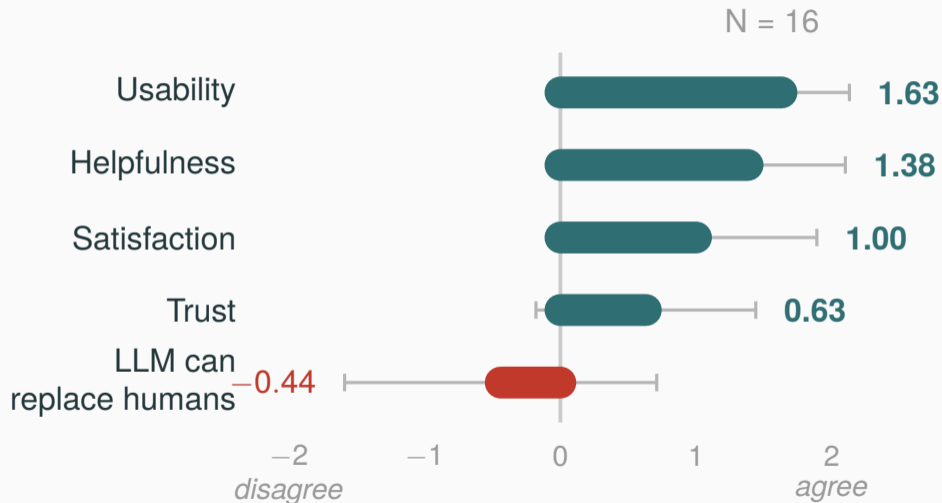
UXAgent Evaluation: can simulated agents help UX researchers?

16 UX researchers · 20 simulated sessions · ~30 min / session



Task: buy a meat substitute (top-rated, \$100–200), stop at checkout

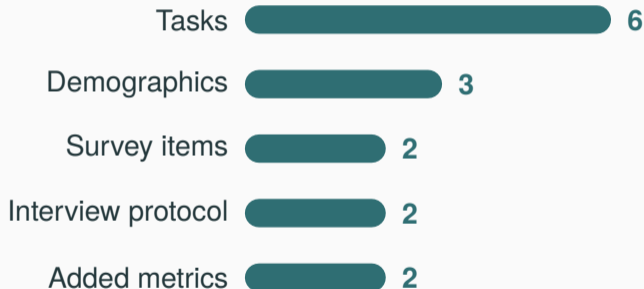
POST-STUDY RATINGS (LIKERT -2 TO 2)



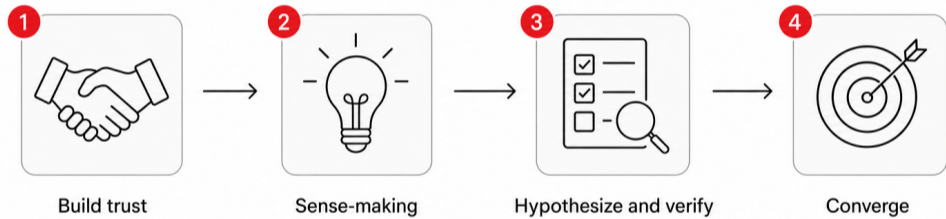
Researchers find UXAgent usable and helpful — but reject it as a *replacement* for humans.

14/16

proposed further study-design changes after using UXAgent



TYPICAL ITERATION CADENCE



“

*I trust it pretty much because it's aligned
with the reasoning.*

— P3

Researchers cross-check **reasoning** ↔ **action** ↔ **screenshot** before trusting the data.
Trust rated $M = 0.62$ (SD 0.81).

convincing

“The reasoning of the actions is very reasonable.”

— P7

not convincing

“I just do not think a real human participant will have actions like that.”

— P6

...*but* “users always surprise me” (P8) — real users are unpredictable too, so some surprising agent behavior may feel *authentic* rather than less.

“

I'm really afraid that researchers only rely on this and overtrust such data.

— P1

Most disagreed that agents can **replace** humans ($M = -0.44$).
Position: **augment** human studies — best for early, low-cost iteration.

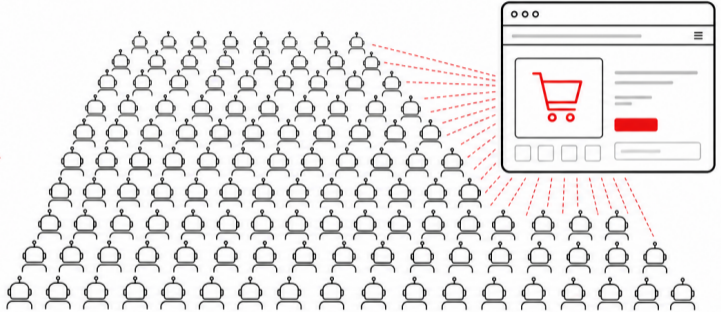
CASE STUDY 2: SIMULATING LARGE-SCALE A/B TESTING STUDY WITH LLM AGENTS

UXAgent: LLM agents simulate users at *small scale* to help UX researchers.

UXAgent: LLM agents simulate users at *small scale* to help UX researchers.

But real design decisions ship through *large-scale A/B testing*.

Can simulated users scale up?



THREE BOTTLENECKS IN A/B TESTING



No lightweight piloting



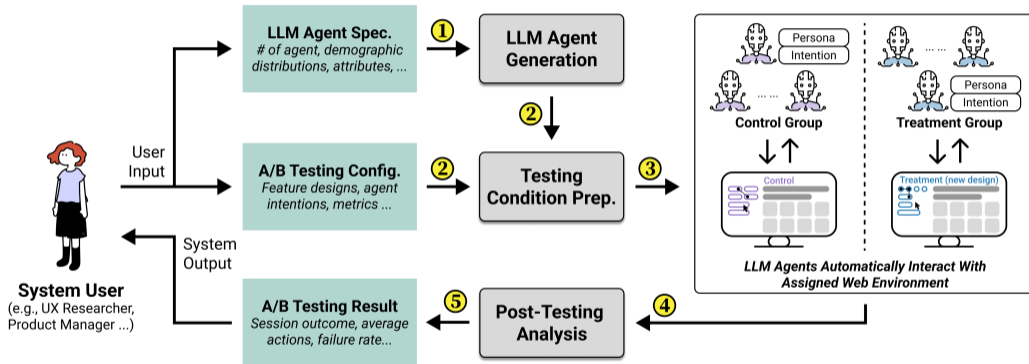
Scarce user traffic



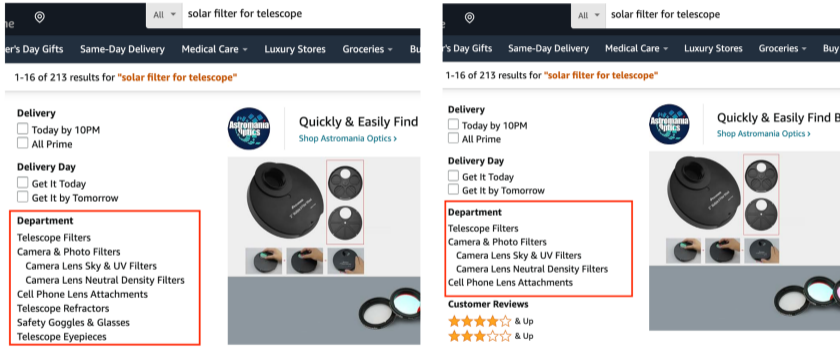
Slow feedback, weeks

From a formative study with 6 A/B-testing practitioners — many promising ideas never get piloted before launch.

**Can we deploy thousands of persona-driven LLM agents
on the *live* website
— to pilot an A/B test before spending real user traffic?**

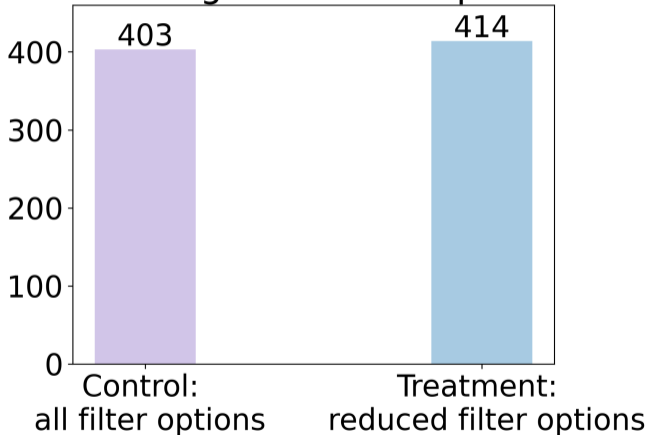


Plug-and-play with existing agent stacks (Claude computer-use, ReAct).

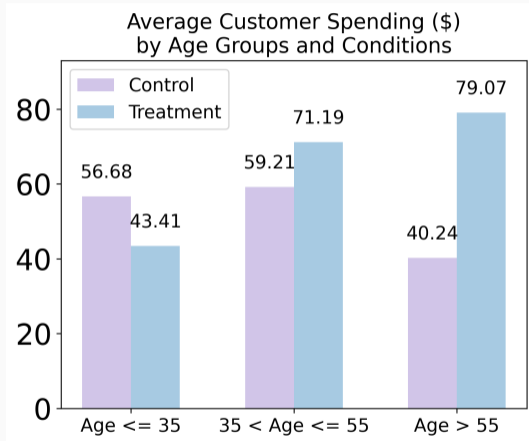
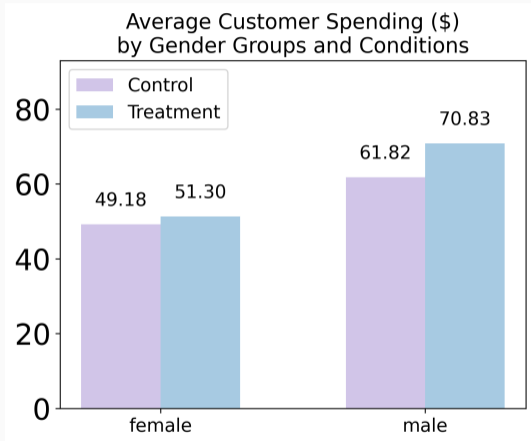


Control: full filter panel vs. **Treatment:** reduced panel (similarity ranking)
1,000 agents (500/condition) — $\approx 10\%$ of a real A/B test.

of LLM Agents made a purchase



414 vs. 403 purchases — a modest but statistically reliable increase.



Effects are stronger for **male** and **older** personas.

Agent outcomes align *directionally* with a parallel large-scale human A/B experiment.

**Simulated users can scale —
and point in the right direction.**

**Simulated users can scale —
and point in the right direction.**

But “directionally aligned” \neq “accurate”.

How faithful are these agents, really?

CASE STUDY 3: USE REAL ONLINE CUSTOMER BEHAVIOR DATA TO EVALUATE AND IMPROVE

CASE STUDY 3: USE REAL ONLINE CUSTOMER BEHAVIOR DATA TO EVALUATE AND IMPROVE

TASK & METHOD

- We focus on the **human behavior simulation task**
 - Generate the next user action based on the context and past actions.
 - Specifically, in the online shopping scenario.

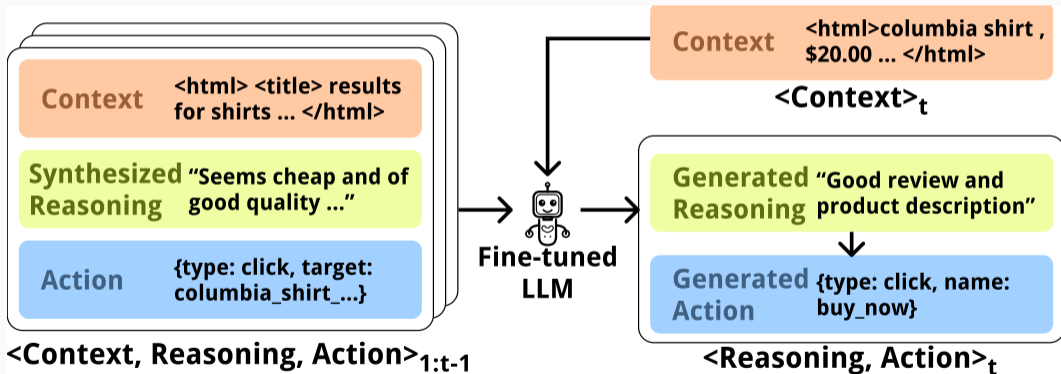
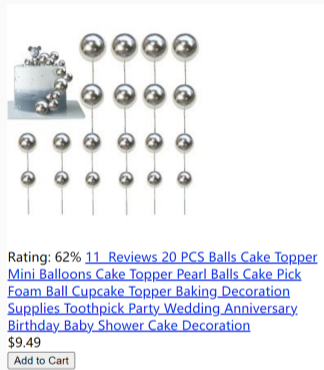


Figure 5: Overview of the next action prediction task.

- Collected from a real-world online shopping platform.
- 31,865 sessions from 3,526 users
- 230,965 user actions
- 4,432 purchases, 27,433 terminations



Rendered page

Context: simplified HTML

JS/CSS stripped, page structure kept; each interactable element gets a unique name (e.g. `product_form.add_to_cart`).

- **Action** is defined as the next raw browser action conducted by the user.
 - Generalizable to other domains beyond online shopping.
- `click` (click on an element)
- `type_and_submit` (type text and submit a form by hitting enter)
- `terminate` (user ends the session by closing the browser window)

- **Reasoning** is defined as a natural language sentence that describes the reasoning behind an action.
 - *“I want to find a comfortable piece of clothing, so I’m looking for options with high ratings.”*
- Enhances the explainability of the model.
- Not present in existing datasets.
- Synthesize with teacher LLM

- To enhance LLMs' accuracy in simulating human behavior, we finetune them on the task.
 - **Input:** $\langle Context, Reasoning, Action \rangle_{1:t-1} + \langle Context \rangle_t$
 - **Output:** $\langle Reasoning, Action \rangle_t$
- Training:
 - Entire session is inputted as a whole.
 - Minimize the loss of the predicted action and reasoning tokens
- Inference:
 - Input the context, past actions and corresponding reasoning.
 - Output the next action and reasoning.

CASE STUDY 3: USE REAL ONLINE CUSTOMER BEHAVIOR DATA TO EVALUATE AND IMPROVE

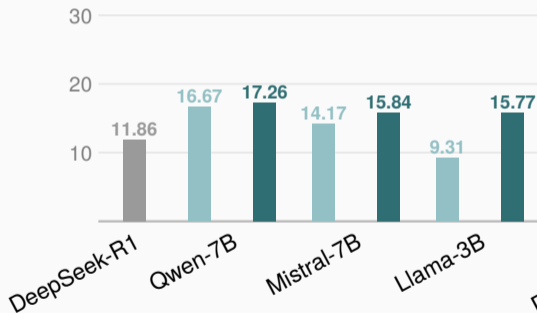
EVALUATION AND EXPERIMENTS

- Next Action Generation
 - Exact Match
 - Predicted action is only considered correct if both the action type (click, terminate, etc.) and the action attribute (the click target / the input text) match the ground truth.

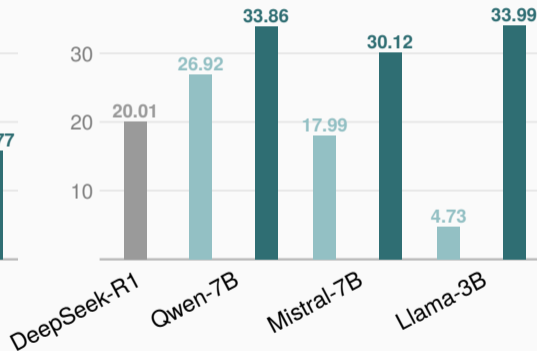
- Next Action Generation
 - Exact Match
 - Predicted action is only considered correct if both the action type (click, terminate, etc.) and the action attribute (the click target / the input text) match the ground truth.
- Shopping Outcome Prediction
 - Essentially predicting the last action based on the session history.
 - One of `click` on a `buy_now` button or `terminate` the session.
 - F1 score

■ prompted baseline ■ SFT ■ SFT + reasoning

Next Action — Accuracy (%)



Session Outcome — F1 (%)



Prompting is not all you need: small fine-tuned models beat the best prompted model, and synthesized *reasoning* lifts outcome F1 further.

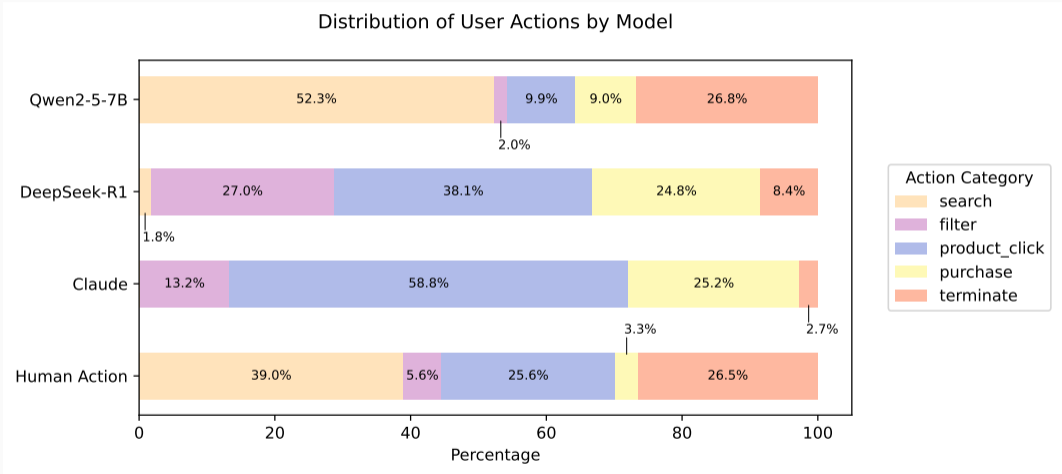


Figure 6: Distribution of the action types in the dataset.

Thank you!
Questions?

