

See, Think, Act: Teaching Multimodal Agents to Effectively Interact with GUI by Identifying Toggles



Zongru Wu¹, Rui Mao¹, Zhiyuan Tian¹, Pengzhou Cheng¹, Tianjie Ju¹, Zheng Wu¹,
Lingzhong Dong¹, Haiyue Sheng², *Zhuosheng Zhang^{1*}*, *Gongshen Liu^{1*}*

¹School of Computer Science, Shanghai Jiao Tong University

²School of Foreign Languages, Beijing Institute of Technology

wuzongru@sjtu.edu.cn



上海交通大学

SHANGHAI JIAO TONG UNIVERSITY

CVPR
JUNE 3-7, 2026



DENVER
COLORADO



Motivation: The GUI Interaction Bottleneck



Omitted Toggling (False Negative)

Activate 9:00 AM alarm

Action: finished()

Omitting the required toggle action results in task failure.

Excessive Toggling (False Positive)

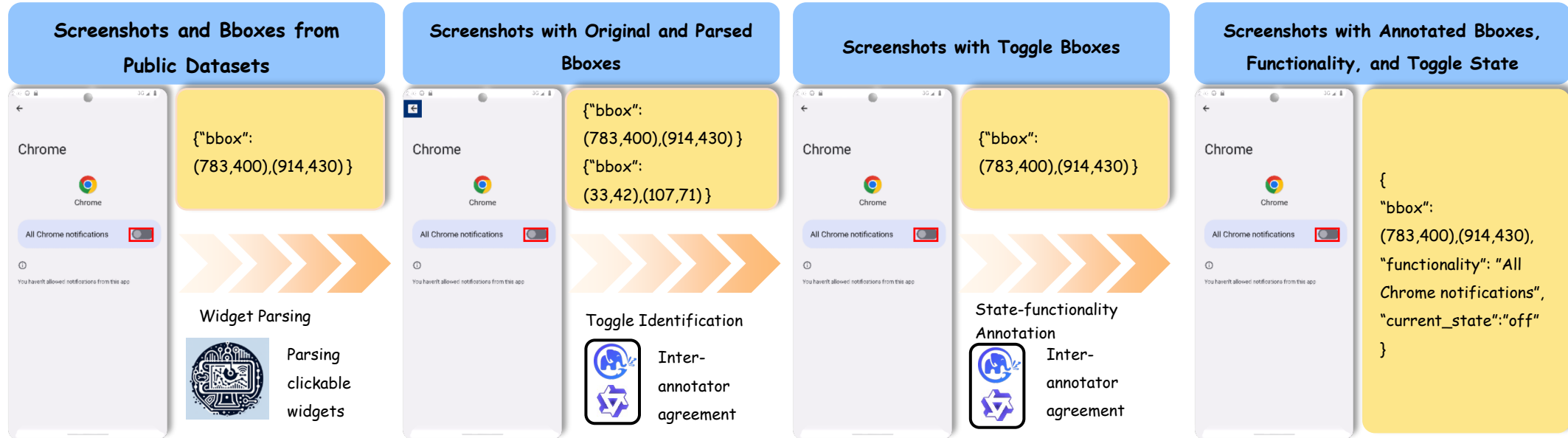
Turn off 9:00 AM alarm

Action: click(852,358)

Excessive toggling that should not occur results in task failure.

Reliable Toggle Execution is Missing

- GUI interaction, especially **toggle control**, is ubiquitous but error-prone for current multimodal agents.
- Typical errors include: **Omitted** Toggling (False Negative) and **Excessive** Toggling (False Positive).

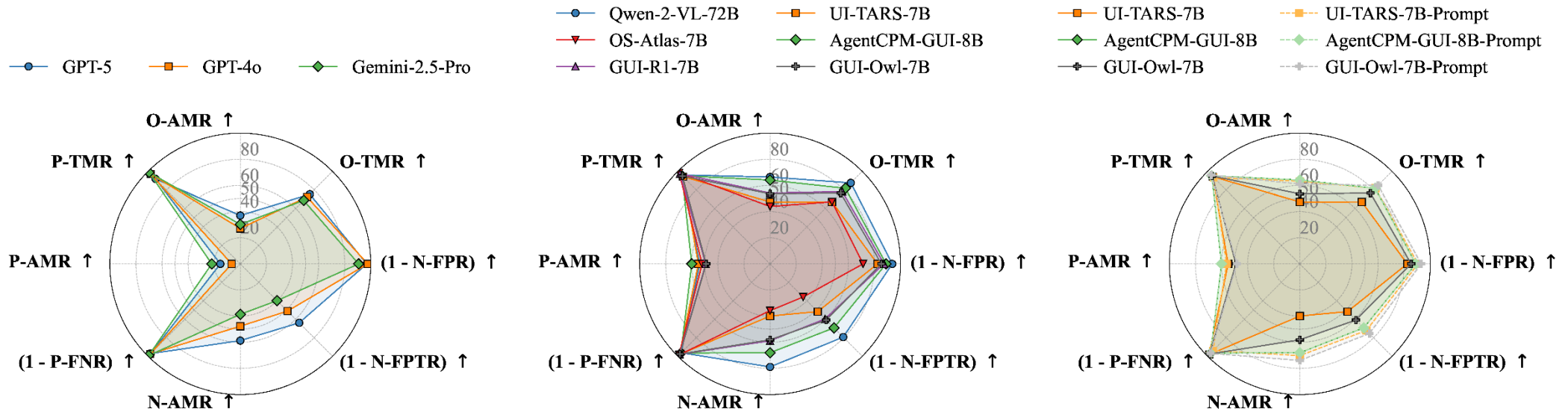


- **Widget Parsing:** Parse all clickable elements from screenshots.
- **Toggle Identification:** Identify GUI toggles from clickable widgets with inter-annotator agreement.
- **State-functionality Annotation:** Label toggle state and functionality with inter-annotator agreement.

Split	AITW [12]	RICOSCA [7]	OS-Atlas [16]	AMEX [2]	AndroidWorld [13]	GUIAct-Mobile [3]	Total
Train	68,380	4,144	496	444	130	58	73,652
Test	7,558	496	60	56	12	2	8,184



Evaluation of Existing Agents



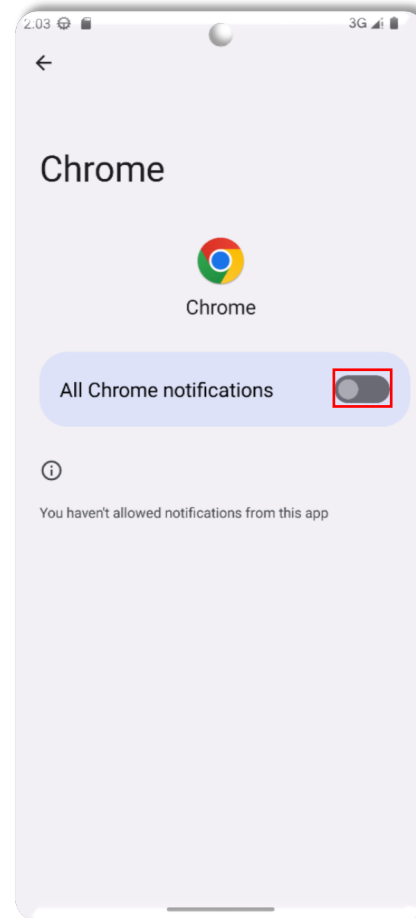
(a) Proprietary MLLMs

(b) Further-trained open-source MLLMs

(c) Comparison with prompt engineering

- Proprietary MLLM agents perform poorly, even GPT-5 yields less than 50% O-AMR (accuracy).
- Open-source MLLM agents remain unsatisfactory, few exceeding 60% O-AMR.
- All agents exhibit a **togglng bias** towards **false positive** togglng.
- Prompting to focus on toggle state during reasoning offers no fundamental improvement.

- Train agents to follow the structured reasoning for toggling:
 - **Perceiving**: Identify the *current state* from screenshot.
 - **Analyzing**: Infer the *desired state* from user query.
 - **Deciding**: Determine the final action based on the *comparison* of the current and desired states.
- General GUI reasoning is preserved with original format.



Positive Reasoning

Turn on `All Chrome notifications`

Perceiving: According to the screenshot, the switch `All Chrome notifications` is currently `Off`.

Analyzing: While according to the goal, it should be `On`.

Deciding: Switch needs to be toggled. Executing: `click on location setting`.

Action: click(850, 415)

Negative Reasoning

Turn off `all Chrome notifications`

Perceiving: According to the screenshot, the switch `All Chrome notifications` is currently `Off`.

Analyzing: While according to the goal, it should be `Off`.

Deciding: No toggle needed. Executing: `set the task as completed`.

Action: finished()



Performance on State Control Benchmark



Model	O-TMR↑	O-AMR↑	P-TMR↑	P-AMR↑	P-FNR↓	N-AMR↑	N-FPTR↓	N-FPR↓
<i>Zero-shot</i>								
Qwen-2-VL-72B	87.59	66.42	96.21	<u>53.89</u>	3.69	78.96	20.67	6.33
GUI-R1-7B	78.27	54.14	<u>97.58</u>	49.32	2.03	58.97	40.37	12.63
OS-Atlas-7B	67.16	43.95	98.51	52.10	1.27	35.80	64.10	28.67
UI-TARS-7B	67.14	47.45	94.33	54.94	<u>1.71</u>	39.96	48.29	17.62
AgentCPM-GUI-8B	<u>81.74</u>	<u>64.08</u>	95.38	60.04	3.32	<u>68.11</u>	<u>30.69</u>	<u>11.07</u>
GUI-Owl-7B	76.58	53.57	94.99	48.97	2.32	58.16	39.15	14.66
<i>w/ StaR-style Prompting</i>								
OS-Atlas-7B	73.52 ^{↑6.36}	50.07 ^{↑6.12}	96.77 ^{↓1.74}	49.88 ^{↓2.22}	2.96 ^{↑1.69}	50.27 ^{↑14.47}	49.62 ^{↓14.48}	22.21 ^{↓6.46}
UI-TARS-7B	81.18 ^{↑14.04}	<u>62.89</u> ^{↑15.44}	90.98 ^{↑3.35}	<u>54.40</u> ^{↑0.54}	8.55 ^{↑6.84}	<u>71.38</u> ^{↑31.42}	<u>27.54</u> ^{↓20.75}	<u>9.38</u> ^{↑8.24}
AgentCPM-GUI-8B	<u>82.14</u> ^{↑0.40}	64.43 ^{↑0.35}	<u>95.36</u> ^{↓0.02}	59.95 ^{↓0.09}	<u>3.59</u> ^{↑0.27}	68.91 ^{↑0.80}	30.28 ^{↓0.41}	10.58 ^{↓0.49}
GUI-Owl-7B	84.27 ^{↑7.69}	60.92 ^{↑7.35}	94.06 ^{↓0.93}	47.36 ^{↓1.61}	4.55 ^{↑2.23}	74.49 ^{↑16.33}	23.68 ^{↓15.47}	7.48 ^{↓7.18}
<i>w/ StaR Training</i>								
OS-Atlas-7B	96.13 ^{↑28.97}	79.72 ^{↑35.77}	95.77 ^{↓2.74}	62.95 ^{↑10.85}	<u>4.23</u> ^{↑2.96}	96.48 ^{↑60.68}	3.52 ^{↓60.58}	1.52 ^{↓27.15}
UI-TARS-7B	95.82 ^{↑28.68}	77.86 ^{↑30.41}	95.11 ^{↑0.78}	59.19 ^{↑4.25}	4.89 ^{↑3.18}	<u>96.53</u> ^{↑56.57}	<u>3.45</u> ^{↓44.84}	1.34 ^{↓16.28}
AgentCPM-GUI-8B	<u>95.98</u> ^{↑14.24}	<u>79.00</u> ^{↑14.92}	94.50 ^{↓0.88}	<u>60.53</u> ^{↑0.49}	5.50 ^{↑2.18}	97.46 ^{↑29.35}	2.54 ^{↓28.15}	0.95 ^{↓10.12}
GUI-Owl-7B	<u>95.99</u> ^{↑19.41}	77.60 ^{↑24.03}	<u>95.65</u> ^{↑0.66}	58.87 ^{↑9.90}	4.35 ^{↑2.03}	96.33 ^{↑38.17}	3.67 ^{↓35.48}	<u>1.56</u> ^{↓13.10}

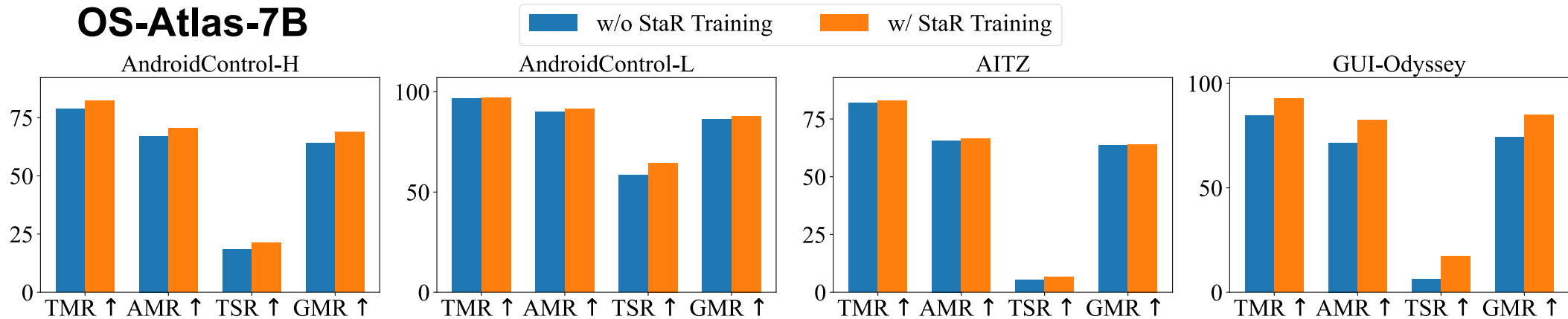
- StaR achieves substantial (over 30%) overall improvements.
- StaR improves negative-instruction accuracy and reduces false positives.



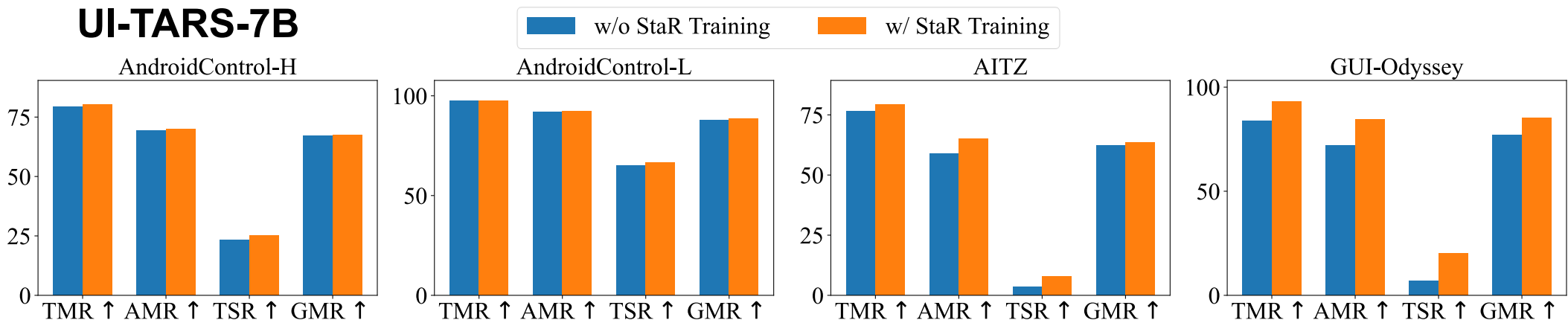
Performance on GUI Agentic Benchmark



- OS-Atlas-7B**



- UI-TARS-7B**



- StaR consistently improves or preserves general agentic performance.



Summary



💡 For reliable GUI interaction, agents must **See** the current state and **Think** before they **Act**.

■ Contributions

- **Benchmark:** We construct a state control benchmark, revealing the unreliability of existing multimodal agents in GUI toggle control.
- **Method:** To fix the unreliable toggling bias, we propose **StaR**, a reasoning method that enables agents to **perceive** the current state, **infer** the desired state, and **act** accordingly.
- **Result:** Boosts toggle accuracy by over 30% while preserving general GUI agentic performance.

■ Source

- **Paper:** <https://arxiv.org/abs/2509.13615>
- **Code:** <https://github.com/ZrW00/StaR>
- **Benchmark:** https://huggingface.co/datasets/ZrW00/StaR_state_control_benchmark



Paper



Code



Benchmark

Thanks For Listening

