

## A. Datasets details

### A.1. Instruction-Tuning data

**Website.** We develop a parser using PyAutoGUI [5] and source websites from 22 representative scenarios such as Airbnb, Booking, AMD, and Apple, which covering shopping, technology, etc. For each scenario, we collect multiple screenshots to maximize annotation coverage. This process yields 22K screenshots with a total of 926K element annotations. After filtering out elements classified as `static text`, we retain 576K elements, averaging 26 elements per screenshot.

**Mobile.** We source mobile data from AMEX [1] annotations: (i) Element grounding and (ii) Functionalities. which covers 97K screenshots, with 885K elements and 178K functionalities.

**Desktop.** We collect 100 screenshots and 2,000 raw annotations from the OmniAct [4] Desktop training split, encompassing 15 software applications across iOS, Windows, and Linux desktops. Additionally, we augment these annotations using GPT-4o-assisted prompting, as detailed in the following section.

### A.2. Downstream tasks

**Mind2Web** [3] supports the development of generalist web agents capable of completing complex tasks on any website by following language instructions. The dataset aligns each HTML document with its corresponding webpage screenshot, featuring a training set of 7,775 actions and three test splits—test-task, test-website, and test-domain—verified for correct rendering and element visibility to ensure reliable evaluation across tasks, websites, and domains. It action space includes three actions: `CLICK`, `TYPE` and `SELECT`.

**AITW** [6] is an Android smartphone environment, which contains 30k instructions and 715K trajectories. We follow the setting by SeeClick [2], which divide the data by domains: General, Install, Google Apps, Single, Web Shopping. It action space includes 12 actions: `CLICK`, `TYPE`, `SELECT`, `SCROLL UP`, `SCROLL DOWN`, `SCROLL LEFT`, `SCROLL RIGHT`, `PRESS BACK`, `PRESS HOME`, `PRESS ENTER`, `STATUS TASK COMPLETE`, `STATUS TASK IMPOSSIBLE`.

**MiniWob** [7] comprises 2000 open-ended tasks from 137 real web environments, each with high-level instruction and action trajectory. It action space includes 2 actions: `CLICK`, `TYPE`.

## B. Settings

### B.1. Evaluation metrics

We follow the original SeeClick [2] settings for evaluation metrics: (i) ScreenSpot: Accuracy; (ii) Mind2Web: El-

ement accuracy, Operation F1-score, and Step success rate; (iii) AITW: Action matching scores; (iv) MiniWob: Success rate.

### B.2. Training details

We utilize 32 V100 GPUs for instruction-tuning, while downstream adaptation is conducted on 8 V100 GPUs. The batch size per GPU is set to 1, with gradient accumulation steps of 2. We use float16 precision for training. To enhance efficiency, we apply LoRA tuning with a rank of 64 and an alpha value of 128 for both the language model and visual encoder, resulting in 4% of the total learnable parameters. We leverage DeepSpeed Zero-2 and use the SDPA attention mechanism. The learning rate is configured to 1e-4. The maximum visual patch number is 1280. The instruction-tuning training duration is approximately two days.

**UI Connected Graph:** We apply the UI-Graph to both the visual encoder and language model with a masking ratio of 0.75, using cross-layer insertion at layer 14. During each iteration, a random ratio of visual tokens is masked. For inference usage, we uniformly sample tokens across each component, ensuring visibility of all components to the model.

**Interleaved Streaming:** In our streaming setting, we set up the history number as 2.

**Data Resampling:** To achieve data balance among existing datasets, probabilities are assigned to each dataset, with weights set as (Web:Mobile:Desktop: GUIAct-Web:GUIAct-Mobile) = (1:1:1: 1:1).

### B.3. More Discussions

**Challenges of Long-horizon task** Our UI-graph is crucial for long-horizon tasks significantly (*e.g.*, 20+ steps) to help reduce the total visual tokens (*e.g.*, 400+ tokens per frame). Instead of storing all screenshots, one potential solution is *streaming with truncated history* or *decomposing complex tasks into subtasks*.

**Replacing screenshots with visual captions.** We explored these ablation studies in AITZ [8], which provides visual captions. As shown in the table, captions do not improve the final action scores compared to using action-solely or with screenshot inputs. This is due to: (i) **Captions often overlook small but critical elements like buttons**, and (ii) **Controlling the captioning accuracy is hard due to model’s hallucinations**. Therefore, screenshots prove to be a more universal and effective form of history.

Models	History (length=2)	Match scores
Qwen2-VL-2B	Actions	59.4
Qwen2-VL-2B	Screen’s captions+Actions	58.4
Qwen2-VL-2B	Screenshots+Actions	<b>61.7</b>

**Relations between ShowUI and backbone model.** We selected Qwen2-VL-2B as the backbone for ShowUI based on

four key considerations: (i) high-resolution with enhanced OCR capabilities, (ii) a lightweight design, (iii) interleaved visual inputs, and (iv) long-context support.

We also considered Phi-3.5-V for ablation studies. As shown below, the Screenshot task is consistent with TextVQA, indicating that *a superior backbone model leads to improved performance, achieving efficient adaptation to GUI Agents.*

Models	Model size	TextVQA (General)	ScreenSpot (GUI)
Phi-3.5-V	4.2B	72.0	71.4
Qwen2-VL	2B	<b>79.7</b>	<b>75.1</b>

## B.4. Prompt templates

**GPT-4o Assisted Prompts.** Below, we display the prompt used for GPT-4o to augment the OmniAct original annotations, which mainly includes three types: (i) Appearance; (ii) Spatial-relationship; (iii) Situational.

You will receive a screenshot with a red bounding box surrounding the target element, along with the target element’s name. Analyze the screenshot and provide concise descriptive responses for each of the following dimensions.

- 1.Appearance:** Describe the target element’s color, shape, ocr and other visual characteristics.  
- Example: “A rectangular chat card with a blue background to ‘Ash.’”
- 2.Spatial:** Describe the target element’s position based on the contextual spatial relationship.  
- Example: “The element that positioned above the Clara.”
- 3.Situational:** Create an intent-oriented query related to the target element, considering how a user might interact with it.  
- Example: “Send a message to Ash.”

Please follow these guidelines:

- Do not confuse the red bounding box with the element itself.
- Provide responses as concise sentences (15 words or fewer).
- For each types, make the description specific enough to distinguish it from other elements.
- If a dimension does not apply, respond with “None.”
- Structure your response in JSON format as shown below:

```
{
  "appearance": "A rectangular chat card with a blue background with letter 'A'",
  "spatial": "The element that positioned above the Clara.",
  "situational": "Send a message to Ash."
}
```

**Grounding Prompts.** For grounding, we apply below prompt, which will be added in revision:

Based on the screenshot of the page, I give a text description and you give its corresponding location. The coordinate represents a clickable location [x, y] for an element, which is a relative coordinate on the screenshot, scaled from 0 to 1. {QUERY}

**Action README Template.** Below, we present the template for action navigation. Variables, marked in **Blue**, depend on specific scenarios, while actions used for loss calculation are highlighted in **Red**.

You are an assistant trained to navigate the {device} . Given a task instruction, a screen observation, and an action history sequence, output the next action and wait for the next observation.

Here is the action space:

# templated by action\_type with action description.

1. ‘CLICK’: Click on an element, value is the element to click and the position [x,y] is required.
2. ‘TYPE’: Type a string into an element, value is the string to type and the position [x,y] is not applicable.

...

Format the action as a dictionary with the following keys:

{‘action’:‘action\_type’, ‘value’:‘element’, ‘position’:[x,y]}

Position represents the relative coordinates on the screenshot and should be scaled to a range of 0-1.

Task: {task}

<past\_image\_1>{past\_action\_1}

...

<past\_image\_n>{past\_action\_n}

<image\_n+1>{action\_n+1}

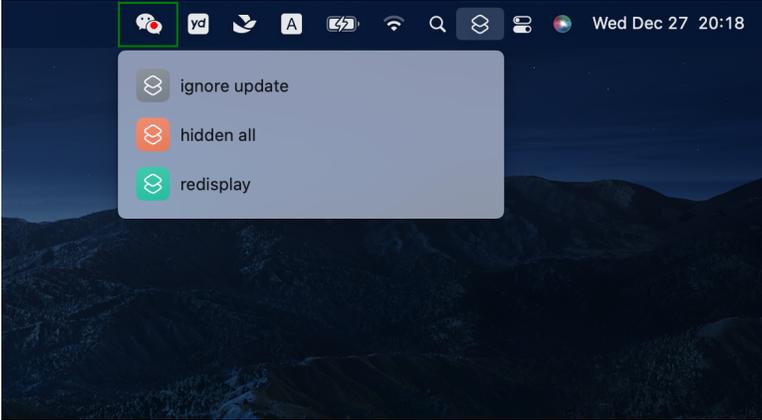
## B.5. Qualitative Examples

In Fig.1 and 2, we display several examples on Screenshot zero-shot grounding. We found that: with instruction tuning, ShowUI is able to perform some visual reasoning, such as it can distinguishes the correct operator among multiple abstract symbols or associate ‘view help’ with question mark, as shown in Fig.1 (b,e). Beside, we found in several failure cases, such as Fig.1 (d,f), there might have multiple possible clickable elements, leading to model confusion.

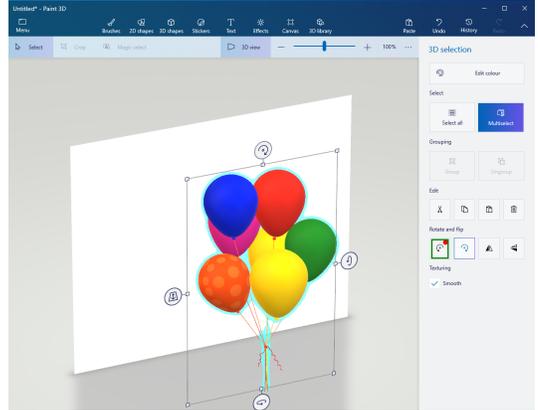
In Fig.4, we demonstrate more examples about how we leverage GPT4o to augment the original OmniAct-Desktop annotations with diverse queries based on Appearance, Spatial Relationships, and Intention.

## References

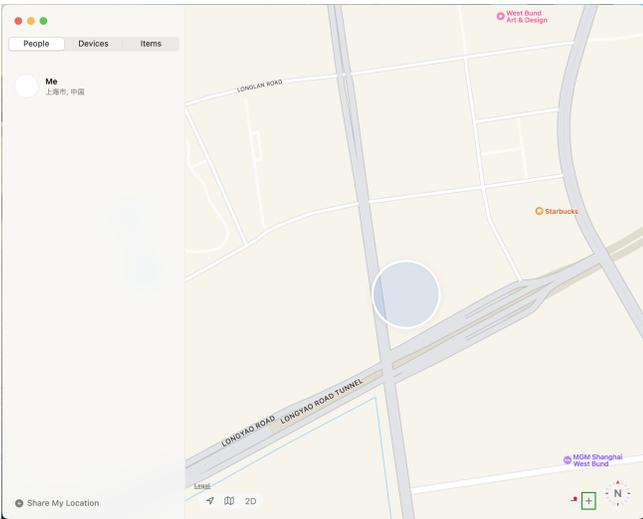
- [1] Yuxiang Chai, Siyuan Huang, Yazhe Niu, Han Xiao, Liang Liu, Dingyu Zhang, Peng Gao, Shuai Ren, and Hongsheng Li. Amex: Android multi-annotation expo dataset for mobile gui agents. *arXiv preprint arXiv:2407.17490*, 2024.
- [2] Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. Seeclck: Harnessing gui grounding for advanced visual gui agents. *ACL*, 2024.
- [3] Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems*, 36, 2024.
- [4] Raghav Kapoor, Yash Parag Butala, Melisa Russak, Jing Yu Koh, Kiran Kamble, Waseem Alshikh, and Ruslan Salakhutdinov. Omniact: A dataset and benchmark for enabling multimodal generalist autonomous agents for desktop and web. *arXiv preprint arXiv:2402.17553*, 2024.
- [5] PyAutoGUI. Pyautogui. 2024. <https://pyautogui.readthedocs.io/en/latest/>.



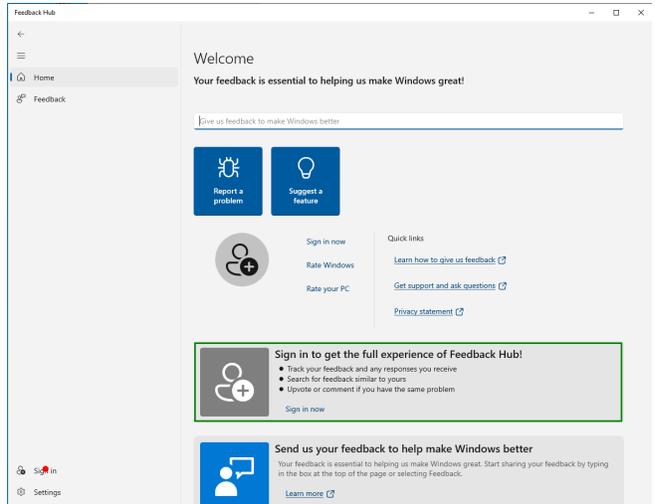
(a) ✓ Instruction: “Open wechat”. With instruction-tuning, ShowUI can recognize the appearance of the WeChat icon.



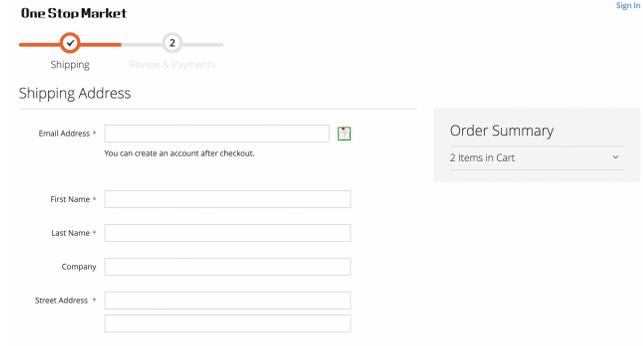
(b) ✓ Instruction: “Rotate left”. ShowUI distinguishes the correct operator among multiple abstract symbols.



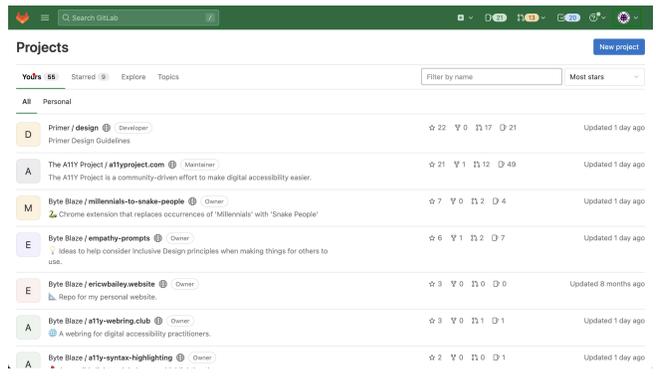
(c) ✗ Instruction: “Zoom in”. The model visually confuses the difference between zoom in and zoom out.



(d) ✗ Instruction: “Sign in”. There are two possible sign-in elements, but the query lacks sufficient information to determine the correct one.



(e) ✓ Instruction: “view help for email account”. ShowUI is able to associate “view help” with question mark clickable element.



(f) ✗ Instruction: “view my account”. ‘View my account’ could be interpreted as ‘Click Your Profile’ or ‘User Profile’ (top right), leading to confusion.

Figure 1. Case studies on Screenspot’s Desktop (a-d) and Web (e-f) grounding,

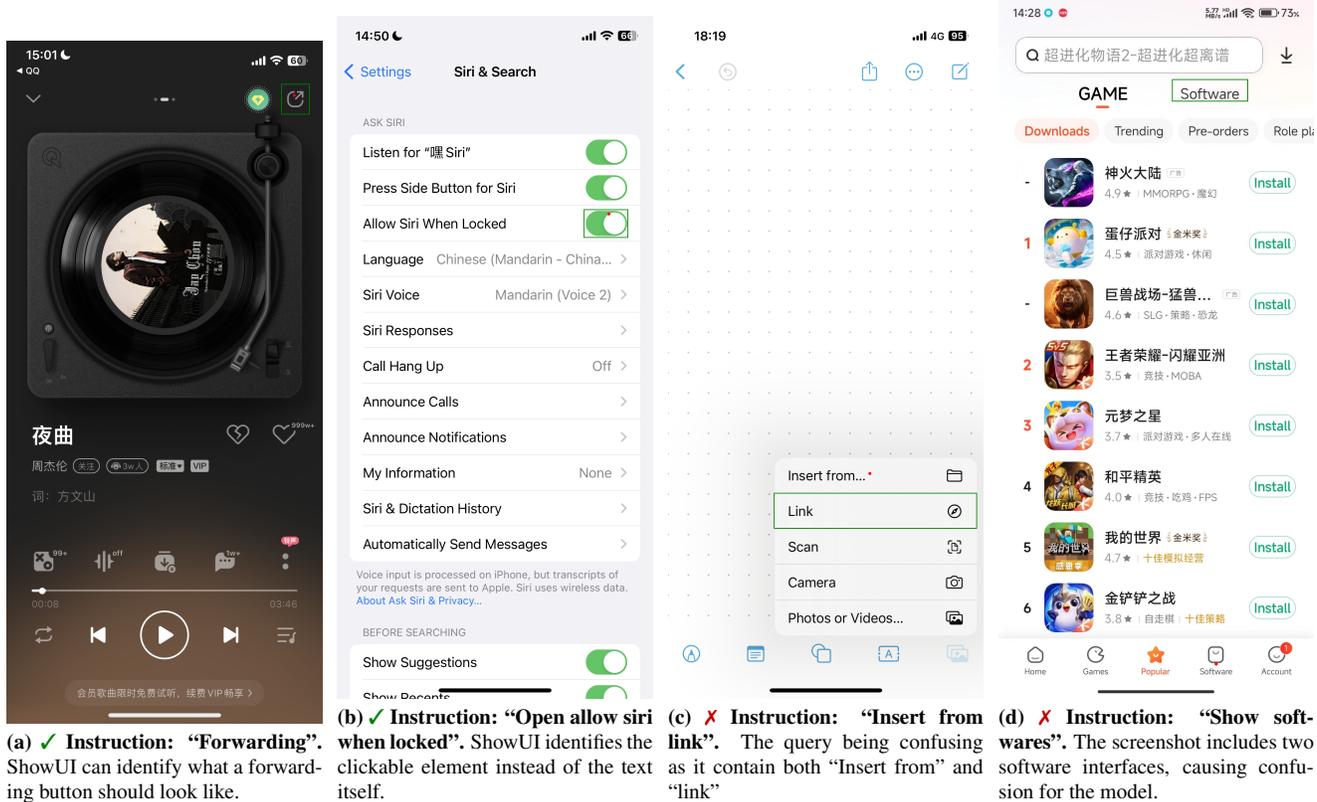
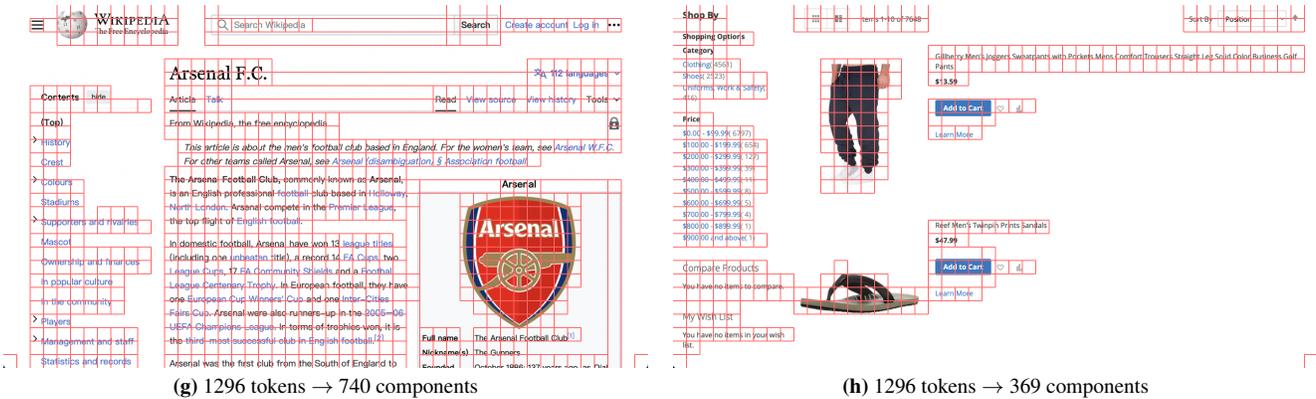
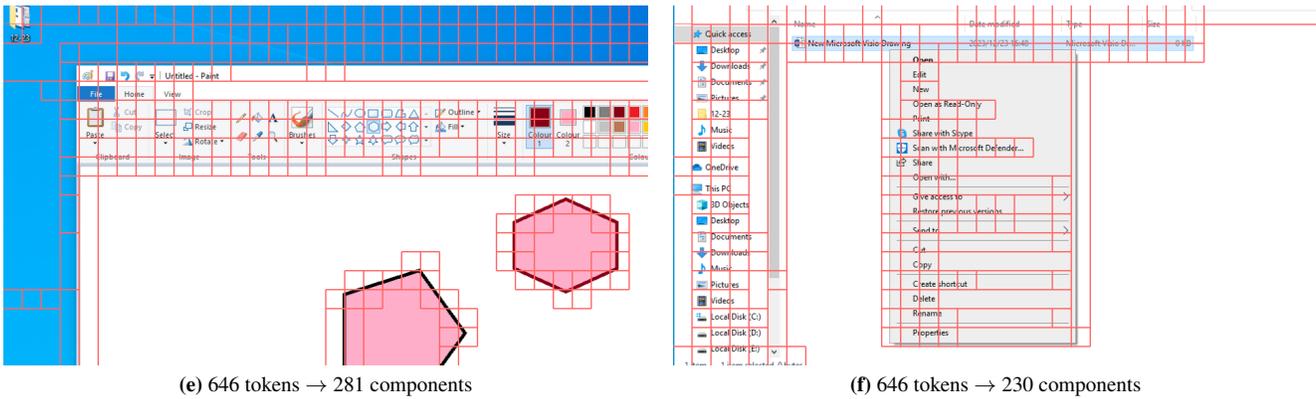
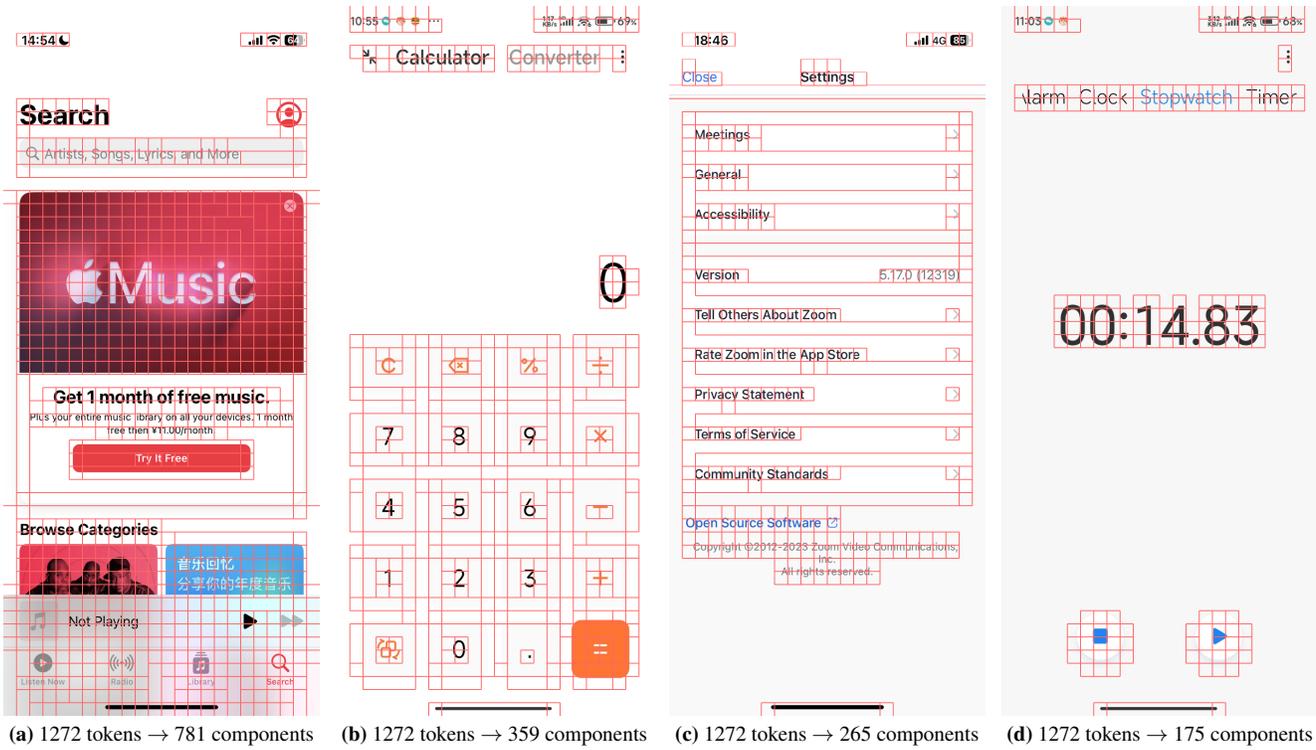


Figure 2. Case studies on Screenspot’s Mobile grounding.

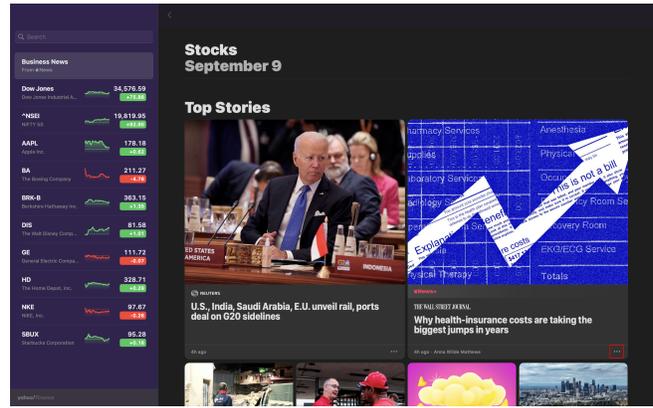
- [6] Christopher Rawles, Alice Li, Daniel Rodriguez, Oriana Riva, and Timothy Lillicrap. Android in the wild: A large-scale dataset for android device control. *arXiv preprint arXiv:2307.10088*, 2023.
- [7] Tianlin Shi, Andrej Karpathy, Linxi Fan, Jonathan Hernandez, and Percy Liang. World of bits: An open-domain platform for web-based agents. In *International Conference on Machine Learning*, pages 3135–3144. PMLR, 2017.
- [8] Jiwen Zhang, Jihao Wu, Yihua Teng, Minghui Liao, Nuo Xu, Xiao Xiao, Zhongyu Wei, and Duyu Tang. Android in the zoo: Chain-of-action-thought for gui agents. *arXiv preprint arXiv:2403.02713*, 2024.



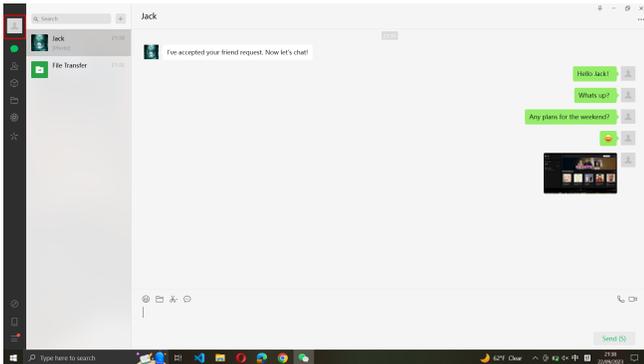
**Figure 3.** Illustration of our method constructs the **UI-connected graph** based on the informativeness of screenshots. (a–d) Mobile; (e–f) PC; (g–h) Web.



(a) Example of Weather. **Original:** ‘visibility’; **Appearance:** “A rectangular box with 28 km in white text.”; **Spatial:** “Positioned below ‘WIND’ and next to ‘PRESSURE’.”; **Intention:** “Check current fog or mist conditions.”



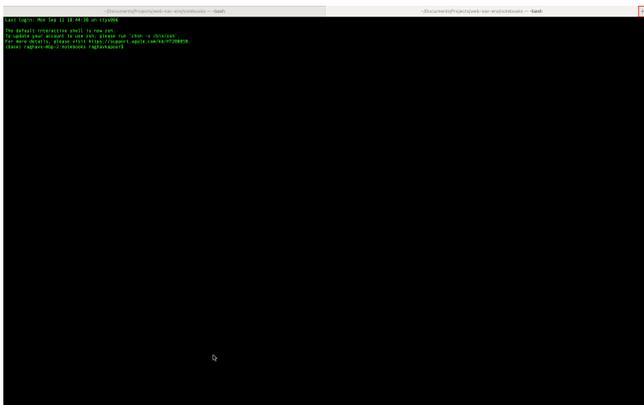
(b) Example of Stocks. **Original:** ‘Share option-health insurance’; **Appearance:** “Three vertical dots icon on a dark background.”; **Spatial:** “Located to the right of the health insurance headline.”; **Intention:** “Click to share the health insurance article.”



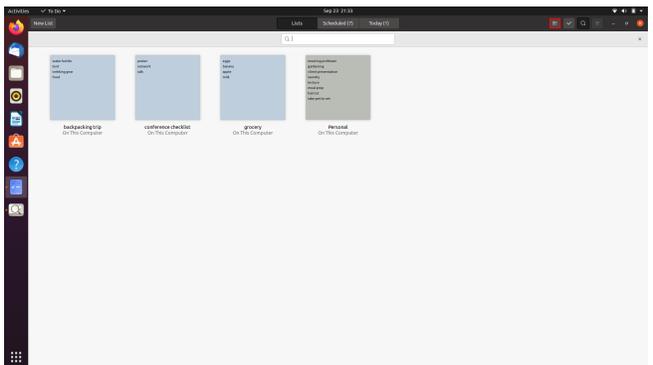
(c) Example of WeChat. **Original:** ‘expand\_profile’; **Appearance:** “A rounded gray button with a person icon.”; **Spatial:** “Located at the top-left corner of the chat pane.”; **Intention:** “Expand the contact’s profile view.”



(d) Example of VLC. **Original:** ‘Play’; **Appearance:** “White triangle icon within a black circular button.”; **Spatial:** “Located at the bottom left corner of the screen.”; **Intention:** “Click to play the video.”



(e) Example of Terminal. **Original:** ‘create\_new\_tab’; **Appearance:** “A small ‘+’ icon in a gray tab bar.”; **Spatial:** “Located at the far right of the tab bar.”; **Intention:** “Open a new terminal tab.”



(f) Example of Todo. **Original:** ‘view as list’; **Appearance:** “A gray, vertical button with a box and lines icon.”; **Spatial:** “Positioned at the top right beside the search bar.”; **Intention:** “Switch to list view.”

**Figure 4.** Illustration of how we augment the original OmniAct-Desktop annotations with diverse queries based on Appearance, Spatial Relationships, and Intention.