# CogAgent: A Visual Language Model for GUI Agents

## Appendix

## 1. Details of Training Configurations

We report the detailed training settings of CogAgent in Table 7, and model configurations of CogAgent in Table 8.

| Configurations | Pre-train | Multi-task |
|---|---|---|
| Total steps | $60,000$ | $10,000$ |
| Warmup steps | 500 | 500 |
| Batch size | $4,608$ | $1,024$ |
| Learning rate | $2 \times 10^{-5}$ | |
| Learning rate decay | Cosine | |
| Weight decay | 0.05 | |
| Dropout ratio | 0.1 | |
| Adam $\epsilon$ | $1 \times 10^{-5}$ | |
| Adam $\beta$ | $(0.9, 0.95)$ | |

Table 7. Training settings of pre-training and multi-task fine-tuning.

| VLM decoder | |
|---|---|
| Architecture | Vicuna-1.5-7B + visual expert |
| Layers | 32 |
| Hidden size | $4,096$ |
| Attention heads | 32 |
| *Low-resolution visual encoder* | |
| Architecture | EVA2-CLIP-E |
| Input resolution | $224 \times 224$ |
| Patch size | $14 \times 14$ |
| *High-resolution visual encoder* | |
| Visual encoder | EVA2-CLIP-L |
| Input resolution | $1120 \times 1120$ |
| Patch size | $14 \times 14$ |
| *Cross Attention* | |
| Hidden size | $1,024$ |
| Attention heads | 32 |

Table 8. Model configurations of CogAgent.

## 2. Details of Evaluation Datasets

In this section, we will provide a detailed overview of the datasets used in our evaluations .

## 2.1. General VQA

- **VQAv2 [1].** VQAv2 is designed for visual question answering with natural images, covering a wide range of question types including yes/no, numerical counting, and more open-ended inquiries. The dataset comprised of a collection exceeding 200,000 images, paired with more than 1.1 million questions. Each question is paired with 10 answers annotated by different annotators.

- **OK-VQA [23].** OK-VQA (Outside Knowledge Visual Question Answering) dataset is constructed to evaluate visual question-answering skills that require external knowledge. Models need to combine image content and common sense to answer questions. The dataset includes 14,055 open-ended questions, each accompanied by 5 ground truth answers.

- **MM-Vet [41].** MM-Vet is designed to evaluate the overall capability of generalist visual language models in a zero-shot manner. It integrates 6 core VL capabilities to solve complex tasks (including recognition, OCR, knowledge, language generation, spatial awareness, and math), and explores 16 distinct integrations resulting from the combination of these capabilities. As for evaluation metrics, it utilizes a language model-based evaluator tailored for open-ended responses.

- **POPE [19].** POPE (Polling-based Object Probing Evaluation) is a dataset constructed to assess the object hallucination problem in large visual language models. It employs a polling-based object probing method, transforming hallucination assessment into a binary classification challenge. This is achieved by prompting large vision-language models (LVLMs) with straight-forward Yes-or-No queries regarding the objects in question (for example, "Is there a car in the image?"). Our evaluation is conducted under the dataset's most challenging setting: the adversarial setting.

## 2.2. Text-rich VQA

- **OCR-VQA [27].** OCR-VQA dataset comprises 207,572 images of book covers, paired with more than 1 million question-answer pairs. The questions inquire about book information including title, edition, year, author, and genre of the book, which requires text recognition and comprehension abilities.

- **TextVQA [34].** TextVQA is a benchmark of visual reasoning based on text in images. Models need to

| Task | Dataset | Description | Split | Metrics |
|------|---------|-------------|-------|---------|
| General VQA | VQAv2 | VQA on natural images. | test-dev | VQA Score($\uparrow$) |
| | OK-VQA | VQA on natural images requiring outside knowledge. | val | VQA Score ($\uparrow$) |
| | MM-Vet | Conversational style VQA on integrated capabilities. | test | GPT-4 score($\uparrow$) |
| | POPE | VQA for hallucination assessment. The adversarial setting is used. | test | F1 score($\uparrow$) |
| Text-rich VQA | OCR-VQA | VQA on images of book covers. | test | EM ($\uparrow$) |
| | TextVQA | VQA on natural images containing text. | test | VQA Score ($\uparrow$) |
| | ST-VQA | VQA on natural images requiring textual understanding. | test | ANLS ($\uparrow$) |
| | ChartQA | VQA about charts with visual and logical reasoning. | test | VQA Score ($\uparrow$) |
| | InfoVQA | VQA on infographics. | test | ANLS ($\uparrow$) |
| | DocVQA | VQA on document images. | test | ANLS ($\uparrow$) |
| GUI Agent | Mind2Web | Web behavior prediction given snapshots and historical actions. | test | step SR ($\uparrow$) |
| | AITW | Android behavior prediction given snapshots and historical actions. | test | Matching Score ($\uparrow$) |

Table 9. Summary of the evaluation benchmarks.

incorporate the textual information in the images and reason over it to answer TextVQA questions. It comprises a total of 28,408 images and 45,336 questions.

- **ST-VQA [4].** ST-VQA is designed to emphasize the significance of exploiting the semantic information present within images in textual form during the VQA process. It comprises tasks of diverse difficulties, for which recognizing the scene text and performing necessary reasoning is required to generate the answer. The dataset comprises 23,038 images sourced from multiple public datasets and 31,791 question-answer pairs.

- **ChartQA [24].** ChartQA is a benchmark of question-answering about logical and visual reasoning on charts. It consists of 20,882 charts curated from four different online sources, 9,608 manual written questions, as well as 23,111 questions automatically generated with T5 according to human-written chart summaries.

- **InfographicVQA(InfoVQA) [26].** The task of InfoVQA is to answer questions centering on a given infographic image. The answers to most questions can be extracted from the given documents, while the answers to a small percentage of questions are not extractive. There are 5K Images collected from the Internet and 30K manually annotated questions in the dataset.

- **DocVQA [25].** DocVQA focuses on question-answering given a document image. The answer for questions is often a span of text from the given documents. There are 12K images and 50K manually annotated questions in the datasets.

## 2.3. GUI Agent

CogAgent is evaluated on two GUI agent datasets, Mind2Web and Android in the Wild (AITW), corresponding to computer agent and smartphone agent respectively.

- **Mind2Web [10].** Mind2Web is designed to develop and evaluate web agents capable of executing intricate tasks on various websites based on language directions. While existing datasets for web agents commonly rely on simulated or overly simplified web data, Mind2Web utilizes real-world websites and is annotated by human annotators. It gathers data from 137 websites covering 31 domains, and collects over 2,000 open-ended tasks, each accompanied by a crowdsourced action sequence.

  In mind2web, evaluated agents are asked to accomplish a designated task on a chosen website by performing a sequence of actions. Each instance (i.e. a specific task) in Mind2Web contains a task description, action sequence, and webpage snapshots. Each action in the sequence is a (Target element, Operation) pair, and Operation includes Click, Type (with additional value), and Select (with additional value). Each action is paired with the concurrent webpage snapshots in a variety of formats including raw HTML code, DOM tree, screenshot, etc. As for CogAgent, we choose screenshot images as the input representation of websites; as for other language-based agents, HTML is chosen as the input representation.

  Following Deng et al. [10] and Zeng et al. [42], we formalize the problem as: first choose the target webpage element among top-k (k=10 or 50) candidates, then predict specific operations. The top-k candidates are provided by the candidate generation model in Deng et al. [10]. Step success rate (step SR) is reported on 3 out-of-domain test sets (cross-website, cross-domain,

cross-task) as metric. Only predictions with the same target element and operation as the ground truth are regarded as correct.

- **Android in the Wild (AITW) [31].** AITW is constructed to develop and evaluate Android device-control systems that are capable of understanding and acting upon human natural language instructions by directly manipulating the device's user interface. This dataset significantly surpasses others in its category in terms of size, encompassing 715k episodes across 30k distinct instructions, and covering four Android versions (v10–13). It also includes eight types of devices, ranging from Pixel 2 XL to Pixel 6, each with different screen resolutions. AITW consists of five subsets: GoogleApps, Install, WebShopping, General, and Single. The subsets have distinct tasks, while are in the same data format. Each episode (i.e. action sequence) comprises three components: a goal instruction provided in natural language, a user action sequence, and a corresponding screenshots sequence.

  As for screenshots, AITW only provides screenshot images and does not provide tree-based representations of UI. This is because a large portion of them in Smartphone applications are of low quality or even do not exist, and adopting the tree-based representation would strongly limit agents' applications. For visual agents, screenshots are provided to the agents in image format; for language-model-based agents evaluated by Zhan and Zhang [43], the textual representations of OCR and icons formatted in HTML syntax are provided. As for actions, AITW considers a variety of action types including tapping, swiping, typing, going home, going back, entering, etc. For each action, models are required to predict the exact action type; for tap, swipe and type, models are further required to predict the position, direction, and content to be typed, respectively. The detailed standard for computing matching scores is provided in Rawles et al. [31].

## 3. Derivation of Acceleration for High-Resolution Cross-Module

Suppose that $L_{I_{lo}}$, $L_{I_{hi}}$ and $L_T$ are the lengths of the low-resolution image, high-resolution image and text sequences. Let $H_{cross}, H_{dec}$ be the number of attention heads in cross-attention and self-attention, and $d_{cross}, d_{dec}$ be the dimension of each attention head.

If using our high-resolution cross-module, the computational complexity of attention is

$$T_{improved} = \mathbf{O}\big((L_{I_{lo}} + L_T)L_{I_{hi}}H_{cross}d_{cross} + (L_{I_{lo}} + L_T)^2 H_{dec}d_{dec}\big). \tag{5}$$

If not utilizing the high-resolution cross-module and directly substituting low-resolution images with high-resolution ones, the computational complexity would be

$$T_{original} = \mathbf{O}\big((L_{I_{hi}} + L_T)^2 H_{dec}d_{dec}\big). \tag{6}$$

The reduction factor of the computational complexity in attention, $T_{original}/T_{improved}$, equals to

$$\frac{(L_{I_{hi}} + L_T)^2 H_{dec}d_{dec}}{(L_{I_{lo}} + L_T)L_{I_{hi}}H_{cross}d_{cross} + (L_{I_{lo}} + L_T)^2 H_{dec}d_{dec}} \tag{7}$$

$$= \frac{L_{I_{hi}} + L_T}{L_{I_{lo}} + L_T} \frac{(L_{I_{hi}} + L_T)H_{dec}d_{dec}}{L_{I_{hi}}H_{cross}d_{cross} + (L_{I_{lo}} + L_T)H_{dec}d_{dec}} \tag{8}$$

$$= \frac{L_{I_{hi}} + L_T}{L_{I_{lo}} + L_T} \frac{(L_{I_{hi}} + L_T)\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}}}{L_{I_{hi}} + (L_{I_{lo}} + L_T)\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}}} \tag{9}$$

- **Case 1:** $L_{I_{lo}}, L_T \ll L_{I_{hi}}$.

  Given that $L_{I_{lo}}$ is much smaller than $L_{I_{hi}}$, when $L_T$ also satisfies being much smaller than $L_{I_{hi}}$, both $\frac{L_{I_{lo}}}{L_{I_{hi}}}$ and $\frac{L_T}{L_{I_{hi}}}$ become first-order small quantities. If conducting a 0-th order approximation for the complexity reduction factor, we obtain:

  $$\frac{T_{original}}{T_{improved}} = \frac{L_{I_{hi}}(1 + \frac{L_T}{L_{I_{hi}}})}{L_{I_{lo}} + L_T} \frac{L_{I_{hi}}(1 + \frac{L_T}{L_{I_{hi}}})\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}}}{L_{I_{hi}}(1 + \frac{L_{I_{lo}} + L_T}{L_{I_{hi}}}\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}})} \tag{10}$$

  $$= \frac{L_{I_{hi}}(1 + \frac{L_T}{L_{I_{hi}}})}{L_{I_{lo}} + L_T} \frac{(1 + \frac{L_T}{L_{I_{hi}}})\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}}}{1 + \frac{L_{I_{lo}} + L_T}{L_{I_{hi}}}\frac{H_{dec}d_{dec}}{H_{cross}d_{cross}}} \tag{11}$$

  $$\approx \frac{L_{I_{hi}}}{L_{I_{lo}} + L_T} \tag{12}$$

- **Case 2:** Our implementation.

  In our implementation, $d_{cross} = 32$, $H_{cross} = 32$, and we inherits $d_{dec} = 128$, $H_{dec} = 32$ from CogVLM-17B. Both high- and low-resolution encoders patchify images with $14 \times 14$-pixel patches, thus $L_{I_{hi}} = 6400$, $L_{I_{lo}} = 256$. Subsitute the values to Eq. (9), we obtain

  $$\frac{T_{original}}{T_{improved}} = \frac{6400 + L_T}{256 + L_T} \frac{4(6400 + L_T)}{6400 + 4(256 + L_T)} \tag{13}$$

  $$> \frac{6400 + L_T}{256 + L_T} \tag{14}$$

  That is to say, our method leads to at least $\frac{L_{I_{hi}} + L_T}{L_{I_{lo}} + L_T} = \frac{6400 + L_T}{256 + L_T} \times$ acceleration which is a stringent lower bound. Taking pre-training as an example, with $L_T$ typically below 512 in earlier stages, our high-resolution cross-module can get an over $25\times$ reduction in computational cost.

## 4. Performance Analysis on AITW

Android in the Wild (AITW) dataset is a large-scale dataset for Android device agents. Our analysis centers on instances where the CogAgent's responses diverge from the annotated ground truth, focusing on these discrepancies.

In our analysis, we conducted a sampling and manual verification of cases where there was a discrepancy between the model's responses and the labeled answers. We sampled instances and discovered that in 42% of these, the model's response actually represented an alternative correct method, rather than an error. This phenomenon is common due to the multiplicity of ways to accomplish a function on a smartphone. For example, many phones are equipped with both Google's app and a Google search bar, and searches can be conducted using either. During the use of smartphones, human users may exhibit simultaneous needs for both in-depth exploration and extensive browsing. For example, when browsing news, CogAgent might identify relevant results on the current page and select them immediately, while some users may prefer to scroll through the entire page first. Furthermore, the data collection for annotations might involve redundant or overly detailed steps, whereas CogAgent can directly find the most efficient path. Fig. 4 illustrates these typical examples of multiple correct answers.

Given the nature of smartphone usage, wherein accomplishing a task often involves multiple viable paths, it is advisable for future studies to explore a virtual environment with an evaluation of the accuracy of task completion following a sequence of actions, aligning more closely with real-world human usage scenarios.



Figure 4. Instances in the AITW dataset where CogAgent provides accurate responses that differ from the annotated ground truth.

## 5. Samples of Pre-train Data

Samples of pre-train data are listed in this section, including data for text recognition (Fig. 5), visual grounding (Fig. 6), and webpage-html pairs (Fig. 7).
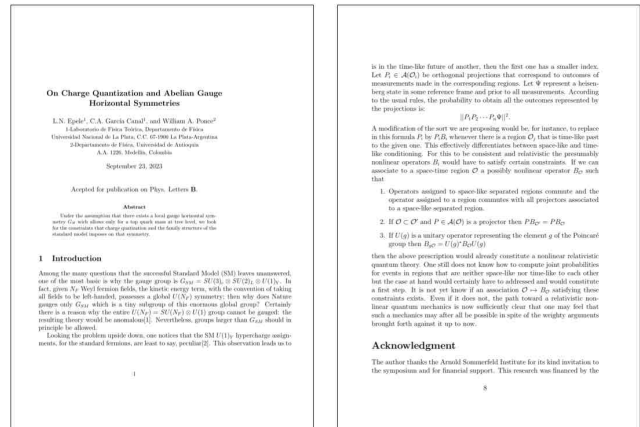


Figure 5. Samples of pre-train data for text recognition.

A bride [[592,179,866,956]] and groom [[717,160,862,882]] standing near the water [[425,426,990,898]] at sunset.
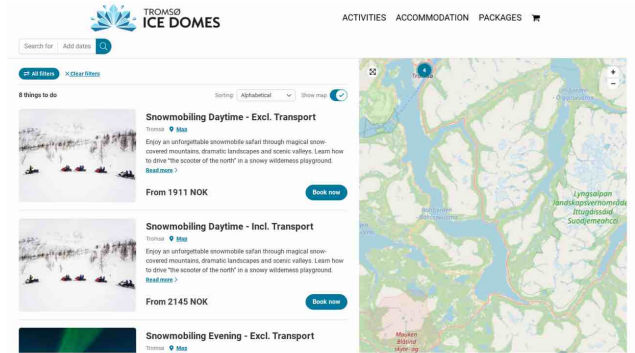


Chicken and rice dish [[301,160,997,792]] in a pan [[282,107,998,839]], with spoons [[091,476,630,953]] and plate [[002,004,258,374]] next to it.
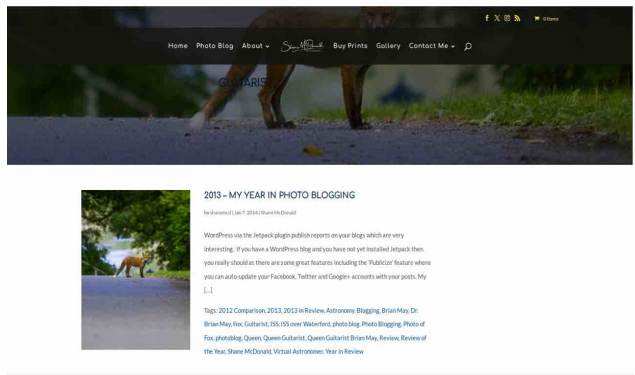


A bluebird [[302,413,640,752]] sitting on a branch coffee mugs [[279,588,677,804]].

Figure 6. Samples of pre-train data for visual grounding.



*[Box]* [[30, 281, 870, 445]]
*[HTML]* <div>\n <div>\n <div>\n <h2>\n <!-- -->\n <a>\n Snowmobiling Daytime - Excl. Transport\n </a>\n </h2>\n <a>\n <!-- -->\n </a>\n <div>\n <div>\n <!-- -->\n <span>\n Tromsø\n </span>\n <span>\n <a>\n <span>\n </span>\n <b>\n Map\n </b>\n <b>\n Map\n </b>\n </a>\n </span>\n </div>\n </div>\n </div>\n </div>\n <a>\n <div>\n <p>\n Enjoy an unforgettable snowmobile safari through magical snow-covered mountains, dramatic landscapes and scenic valleys. Learn how to drive "the scooter of the north" in a snowy wilderness playground.\n </p>\n <div>\n <span>\n Read more\n </span>\n <span>\n <span>\n </span>\n </span>\n </div>\n </div>\n </a>\n <!-- -->\n</div>



*[Box]* [[190, 510.203125, 1088.21875, 682.703125]]
*[HTML]* <p>\n WordPress via the Jetpack plugin publish reports on your blogs which are very interesting. If you have a WordPress blog and you have not yet installed Jetpack then you really should as there are some great features including the 'Publicize' feature where you can auto-update your Facebook, Twitter and Google+ accounts with your posts. My [...]\n</p>

Figure 7. Samples of webpage-html pairs.

## 6. Details of Fine-Tuning Data

### 6.1. Human annotation

To enable CogAgent to function as an agent on various apps and websites, we assembled a team of over ten annotators to collect and label screenshots. To ensure the quality of the annotations and prevent overly simplistic, homogenized content, the annotation process was divided into two phases.

In the first phase, annotators were required to independently select from a pool of apps and websites, and capture screenshots of pages strongly relevant to the main function of that app/website, with no less than 20 screenshots for each app/website. Subsequently, for each screenshot, the annotators would label them according to the following aspects:

1. **Five buttons.** List the names of five clickable buttons from the screenshots. The types of buttons should be diverse.

2. **Three clickable areas.** List three clickable areas from the screenshots, e.g. text boxes, search boxes, clickable images, hyperlinks, etc.

3. **Two questions extracting information from the image.** Ask questions involving the textual information in the images. For instance, "In what year did the user in the image register?" .

4. **One operation requirement.** What function of the webpage/app would you use in this interface? For example, adding the comment: "cool!". The requirement shouldn't be vague such as "adding comment".

In the second phase, annotators are required to provide grounding annotation for the questions and operational methods proposed in the first phase (the aforementioned parts 3 and 4). For example, for Figure 8, a possible annotation would be:

1. **Buttons**: Back; Search; Subscribed; Home; Library.

2. **Clickable areas**: Avatar in the middle top; Video preview in the middle bottom; Personal profile.

3. **Question 1**: Based on the page, how many followers does this author have?

   **Answer**: According to the personal profile at the top [[013,568,802,188]], this author has 4.97M followers.

   **Question 2**: Based on the page, how many videos has this author posted?

   **Answer**: According to the personal profile at the top [[013,568,802,188]], this author has posted 502 videos.
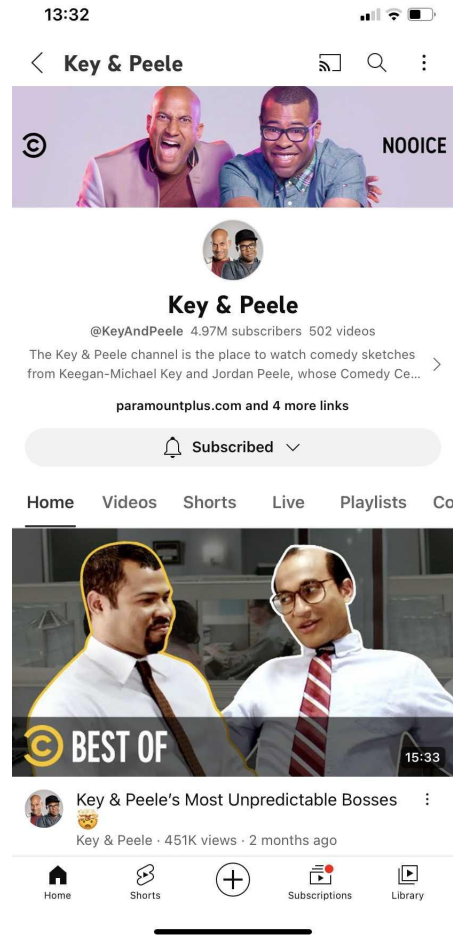


Figure 8. Samples of human-collected screenshot.

4. **Operation requirement**: Based on the page, write out the steps to complete the following function: Follow this author.

   **Answer**: Click on Subscribed [[049,826,728,078]] to follow this author.

### 6.2. Conversion of Agent Datasets

To convert Mind2Web to natural language with GPT4, we use the following prompt:

> Imagine that you are a robot operating a computer. Like how humans operate the computer, you can move the mouse, click with the mouse, or type some texts with the keyboard.
>
> **Your ultimate task is: "Find the lowest-priced round trip flight with hotel on May 2 from Kathmandu, Nepal KTM to Shanghai, China PVG and return on May 5. Book a double room and check out with the default flights.".**

You are given previous actions: (format: element → operation)

1. [link] Flight + Hotel → CLICK,
2. [textbox] Where from? → TYPE: KATHMANDU,
3. [div] Tribhuvan Intl Airport (KTM), Nepal → CLICK,
4. [textbox] Where to? → TYPE: SHANGHAI,
5. [div] Pudong Intl Airport (PVG), China → CLICK,
6. [span] Sat 8 Apr - Sat 15 Apr → CLICK,
7. [checkbox] 2 May 2023 → CLICK,
8. [checkbox] 5 May 2023 → CLICK,

Furthermore, you **secretly** know the following steps of action. **(But you have to pretend this is the result of your derivation and you don't know it before)**

1. [button] Search → CLICK,
2. [button] Recommended → CLICK,
3. [div] Lowest Price → CLICK,
4. [svg] → CLICK,
5. [button] Choose your room → CLICK,
6. [button] Book Business Double Room A → CLICK,
7. [svg] → CLICK,
8. [svg] → CLICK,
9. [button] Choose → CLICK

Please formulate an operational guide for future operations for solving the task. The guide includes:

1. Plan: A **multi-step future** plan that does not involve low-level operations **(start from current screen and action, DON'T include previous actions)**; steps indexed by numbers. **Be sure to pretend that you don't know the future interface and actions, and don't know the elements not existing on the current screen.**

2. Action: The specific **immediate** action that needs to be taken (Corresponding to the** very next action** [button] Search → CLICK at the box {"x_left": 0.876, "y_left": 0.308, "width": 0.063, "height": 0.034}), which should be detailed enough to directly control the operation, such as specifying the exact location (or box) to click, the exact content of the text to type, the item to select.

**Output Format:** A JSON dictionary strictly following the format:

{ "plan": "...<Your Plan Here>", "action": "...<Your Action Here>", "operation": "[button] Search → CLICK at the box "x_left": 0.876, "y_left": 0.308, "width": 0.063, "height": 0.034"} **Please do not output any content other than the JSON format. Don't mention elements that only appear in HTML such as svg, div.**

The returned answer from GPT-4 would be:

{ "plan": "1. After searching, you'll see a list of flight and hotel packages. 2. Review the recommended options, usually highlighted for convenience. 3. Focus on finding the lowest priced option. You might need to navigate through various sorting or filtering features to do this. 4. Once you've found the lowest priced option, select it to view more details. 5. In the detailed view, you'll have to choose your room type. Look for options that fit your requirements, like a business double room. 6. After selecting the room, proceed to finalize your booking. This might involve confirming your choice and entering any required information. 7. Review all the details carefully before finalizing the booking to ensure everything is correct.",

"action": "Click the 'Search' button to proceed with the search based on the entered criteria. This button is located towards the right side of the screen.",

"operation": "[button] Search → CLICK at the box {"x_left": 0.876, "y_left": 0.308, "width": 0.063, "height": 0.03}" }

# 7. Failure cases

The failure cases can be summarized into four categories (Fig. 9):

1. Incorrect action prediction;

2. Incorrect coordinate prediction;

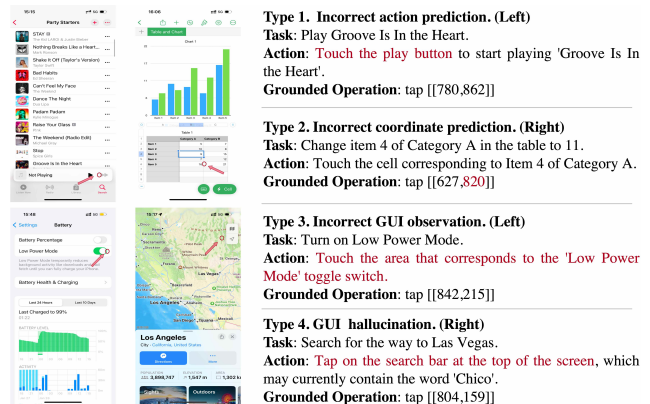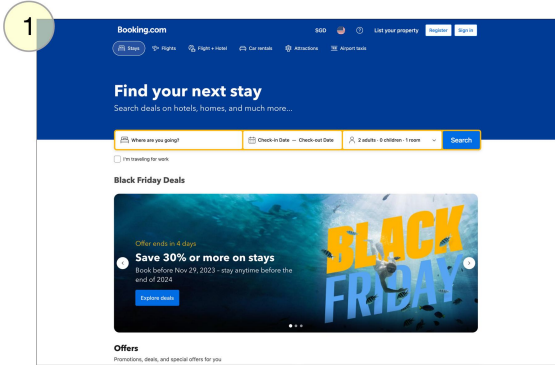3. Incorrect GUI observation;

4. Hallucination.



Type 1. **Incorrect action prediction. (Left)**
**Task**: Play Groove Is In the Heart.
**Action**: Touch the play button to start playing 'Groove Is In the Heart'.
**Grounded Operation**: tap [[780,862]]

Type 2. **Incorrect coordinate prediction. (Right)**
**Task**: Change item 4 of Category A in the table to 11.
**Action**: Touch the cell corresponding to Item 4 of Category A.
**Grounded Operation**: tap [[627,820]]

Type 3. **Incorrect GUI observation. (Left)**
**Task**: Turn on Low Power Mode.
**Action**: Touch the area that corresponds to the 'Low Power Mode' toggle switch.
**Grounded Operation**: tap [[842,215]]

Type 4. **GUI hallucination. (Right)**
**Task**: Search for the way to Las Vegas.
**Action**: Tap on the search bar at the top of the screen, which may currently contain the word 'Chico'.
**Grounded Operation**: tap [[804,159]]

Figure 9. Failure cases of GUI Agent.

# 8. More Generated Samples of CogAgent
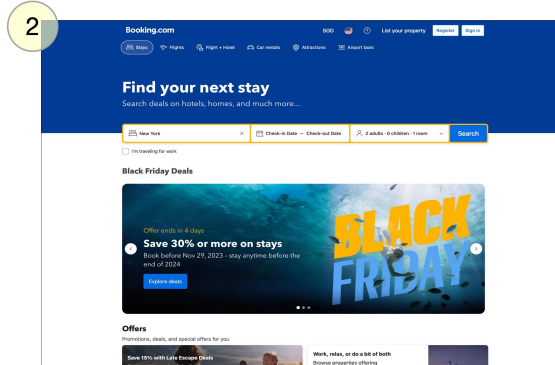
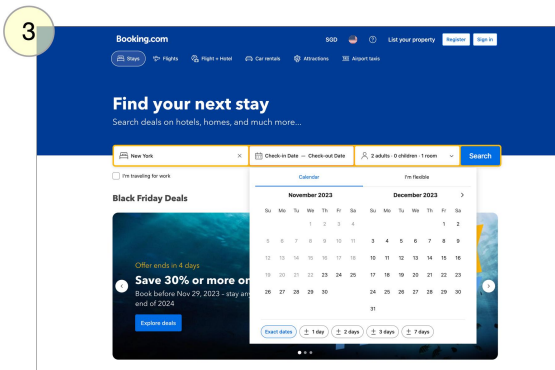To further demonstrate the capabilities of CogAgent, we provide additional generated samples in this section.

**Plan:** 1. Type 'New York' into the search box to specify the destination.
2. Select the correct New York location from the auto-suggested list.
3. Adjust the check-in and check-out dates to December 1st, respectively.
4. Execute the search by clicking the 'Search' button.
5. Once the results are displayed, sort them by price, selecting the lowest first.
6. Review the options and select the room offer that is the lowest priced.
7. Proceed to book the selected room following the website's booking procedure.
**Next Action:** Type 'New York' into the textbox located at the coordinates {'x_min': 0.175, 'y_min': 0.309, 'width': 0.28, 'height': 0.047} on the screen.
**Grounded Operation:** [textbox] Where are you going? -> TYPE: New York at the box [[175,309,455,356]]



**Plan:** To book the lowest priced room in New York on Dec 1st, the next steps would involve:
1. Adjusting the check-in and check-out dates to December 1st, respectively.
2. Executing a search for available rooms.
3. Sorting the search results by price, from low to high.
4. Selecting the room offer with the lowest price.
5. Proceeding to book the chosen room.
**Next Action:** Click on the 'Check-in Date' textbox at the specified coordinates to bring up the date selection interface.
**Grounded Operation:** [textbox] Check-in Date -> CLICK at the box [[438,309,549,357]]



**Plan:** 1. Confirm the checkout date by clicking on the appropriate date, which in this case is 1 December 2023.
2. Initiate the search for hotels by clicking the 'Search' button to find available rooms.
3. Once the results are displayed, sort the list of hotels by price, selecting the option to order them from the lowest to highest price.
4. Review the sorted list and select the hotel offering the lowest price.
5. Proceed to book the selected hotel room following the website's booking procedure.
**Next Action:** Click on the calendar box corresponding to the date 1 December 2023 to set it as the checkout date.
**Grounded Operation:** [gridcell] 1 December 2023 -> CLICK at the box [[739,564,776,611]]
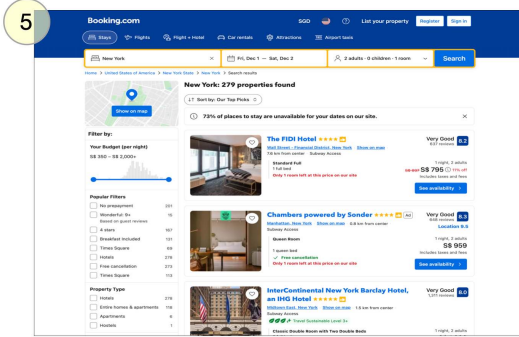


**Plan:** 1. Click the Search button to initiate the search with the current criteria.
2. Once the search results are displayed, sort the listings by price, selecting the option to order them from low to high.
3. Review the sorted list and select the room offer that is the lowest in price.
4. Click on the 'Book' button for the chosen offer to proceed with the reservation process.
**Next Action:** Click the Search button located at the bottom right of the search parameters box to execute the search for hotels in New York from December 1st to December 2nd for 2 adults and 0 children in 1 room.
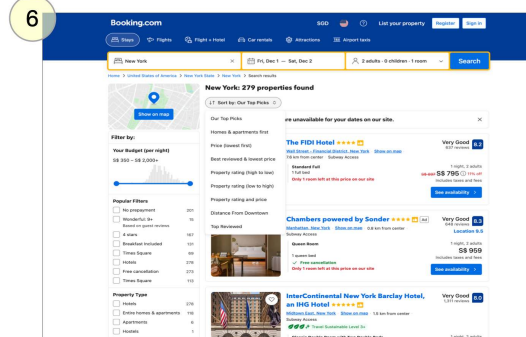**Grounded Operation:** [button] Search -> CLICK at the box [[798,367,875,424]]
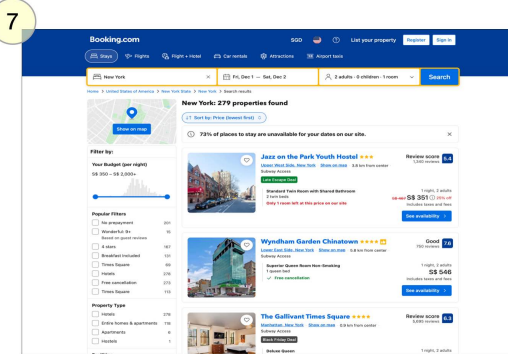
Figure 10. More generated samples of CogAgent.

**User:** What's the proper procedure for "**Book a room of the lowest price in New York on Dec 1st**"? (with grounding)
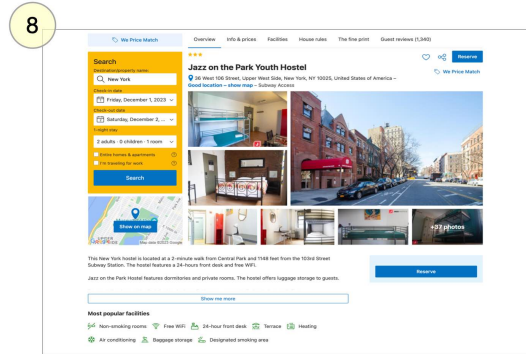


**Plan:** 1. Click on the 'Sort by: Our Top Picks' button to change the sorting criteria.
2. Select 'Price (lowest first)' from the sort options to arrange the listings by price in ascending order.
3. Identify the lowest-priced hotel listed after the sort operation.
4. Click on the 'See availability' button for the lowest-priced hotel to view available rooms.
5. Choose a room option that fits the requirements, such as a specific price range or room type.
6. Proceed with booking by clicking 'Reserve' and following the subsequent booking steps.
**Next Action:** Click at the box with coordinates starting approximately at 29.8% from the left of the screen and 27.2% from the top of the screen.
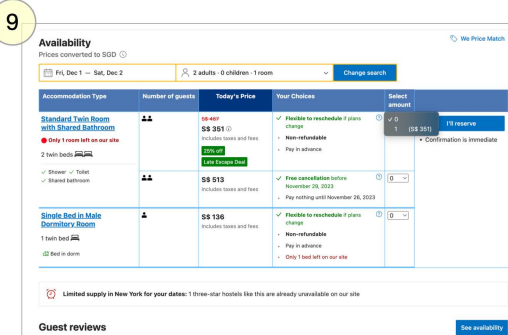**Grounded Operation:** [button] Sort by:Our Top Picks -> CLICK at the box [[298,272,420,307]]



**Plan:** 1. Click on the 'Price (lowest first)' option to sort the properties by price in ascending order.
2. Review the sorted list and identify the property with the lowest price.
3. Select the property by clicking on it.
4. Once on the property's page, choose a room by clicking on the 'See availability' button for that room.
5. Follow the booking process by selecting the date, room type, and any other necessary options, then proceed to reserve the room by clicking 'Reserve'.
**Next Action:** Click on the 'Price (lowest first)' sorting option located towards the top of the webpage, within the specified coordinates on the screen.
**Grounded Operation:** [li] Price (lowest first) -> CLICK at the box [[288,379,494,414]]



**Plan:** 1. Click on the 'Jazz on the Park Youth Hostel' to view more details about the property and to verify it meets all booking requirements including availability for the specified date, price, and amenities.
2. Review the detailed information provided after the click to ensure it matches the desired criteria.
3. If the property is suitable, proceed to reserve a room by following the booking instructions which may include selecting a room type, providing personal details, and making a payment.
4. Confirm the reservation and ensure that a confirmation message or email is received.
**Next Action:** Move the cursor over the box that represents the 'Jazz on the Park Youth Hostel' link with coordinates starting at approximately 38.9% from the left of the screen and 38.6% from the top of the screen.
**Grounded Operation:** [link] Jazz on the Park Youth Hostel Opener West Side, New Yor... -> CLICK at the box [[389,386,810,435]]



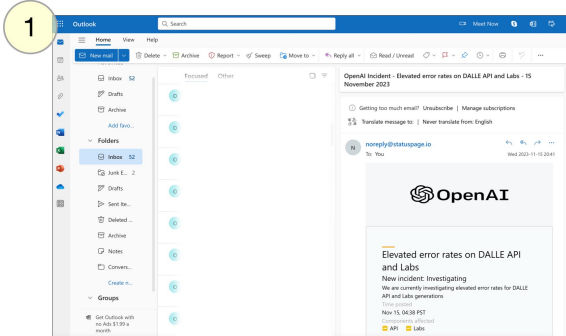**Plan:** After the current action, the next steps should include confirming the reservation details such as dates, number of guests, and room preferences. Then proceed to enter personal information required for booking. Following that, review all the details for accuracy, choose a payment method, and complete the booking process. Finally, wait for a confirmation message or email from the hotel to ensure the reservation has been successfully made.
**Next Action:** Click on the 'Reserve' button located towards the right side of the screen in a box with coordinates {'x_min': 0.781, 'y_min': 0.724, 'width': 0.106, 'height': 0.038}.
**Grounded Operation:** [button] Reserve -> CLICK at the box [[781,724,887,762]]



**Plan:** 1. Click on the 'Select amount' button for the room with the lowest price listed.
2. Choose the desired payment option, such as 'Pay in advance' if available and offered at a lower rate.
3. Confirm the selection and proceed to booking by clicking on 'Continue to book'.
4. Complete the booking process by entering personal and payment details as required.
5. Review the reservation summary and finalize the booking by confirming all details are correct and clicking 'Book now'.
**Next Action:** Click on the 'Select amount' button located towards the right side of the screen for the accommodation option with the lowest price displayed, which is the 'Standard Twin Room with Shared Bathroom' currently priced at $351 including taxes and fees.
**Grounded Operation:** [button] Select amount -> CLICK at the box [[739,272,906,318]]
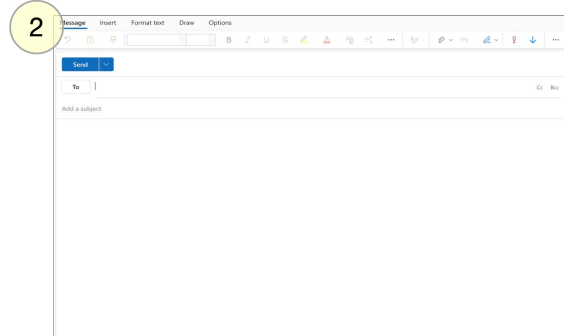
Figure 11. More generated samples of CogAgent.

**User:** Could you guide me through the steps to **"Write an email to cvpr2024@gmail.com, with subject of 'I'm CogAgent', and the main body of 'Here is CogAgent sending the email.'"** ?
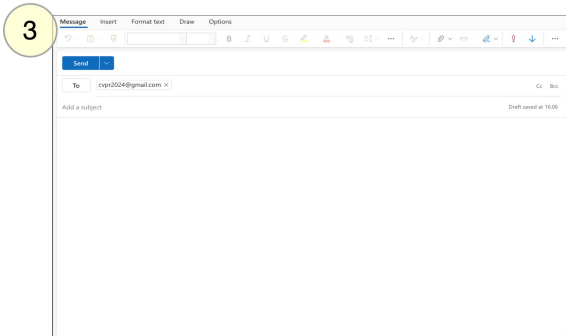


**CogAgent:**

**Plan:** 1. Click on the 'New email' button to start a new email.
2. Enter the recipient's email address in the 'To:' field.
3. Enter the subject in the 'Subject:' field.
4. Type the main body of the email in the 'Body' text area.
**Next Action:** Click on the 'New mail' button located at the top left corner of the screen, within the Outlook interface.
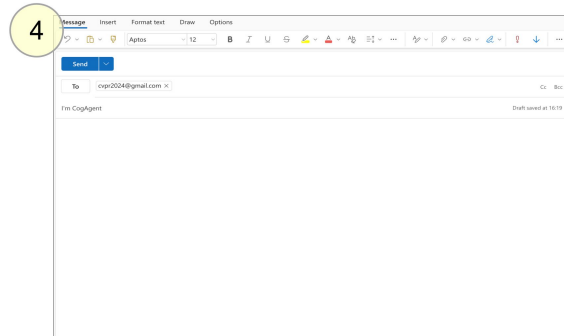


**CogAgent:**

**Plan:** 1. Type the email address into the 'To' field.
2. Enter the subject line in the 'Subject' field.
3. Use the keyboard to type the main body of the email.
**Next Action:** Move the cursor to the 'To' textbox, click on it to focus on it, and type 'cvpr2024@gmail.com'.



**CogAgent:**

**Plan:** 1. Type the subject line into the 'Subject' field.
2. Enter the main body of the email in the large text area below the subject line.
3. Use the 'Send' button to send the email.
**Next Action:** Click on the 'Subject' textbox, which is currently highlighted and type 'I'm CogAgent'. Then, use the keyboard to type the rest of the message in the large text area below it.



**CogAgent:**

**Plan:** 1. Click on the text area to bring up the keyboard if it's not already active.
2. Type the main body of the email, 'Here is CogAgent sending the email.'
3. Find and click the 'Send' button to send the email.
**Next Action:** Move the cursor to the text area where the message content needs to be typed and type 'Here is CogAgent sending the email.'
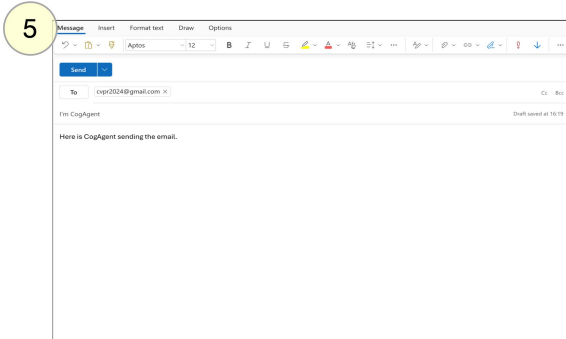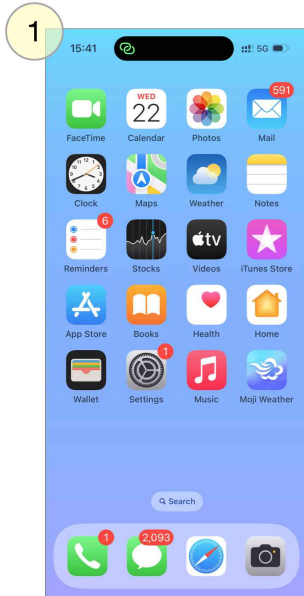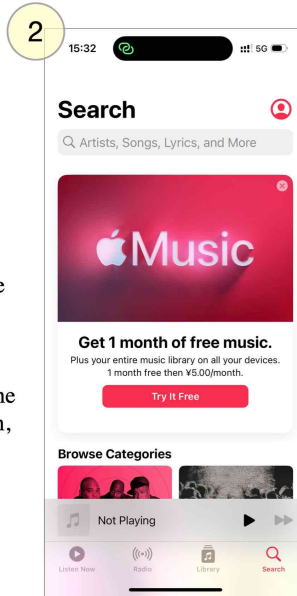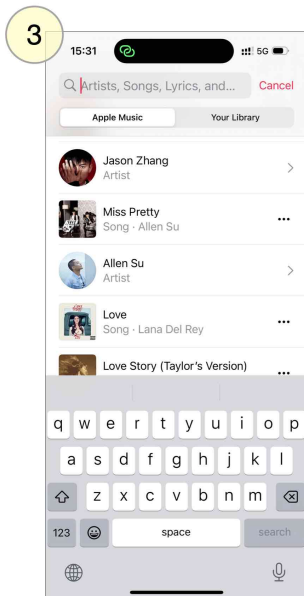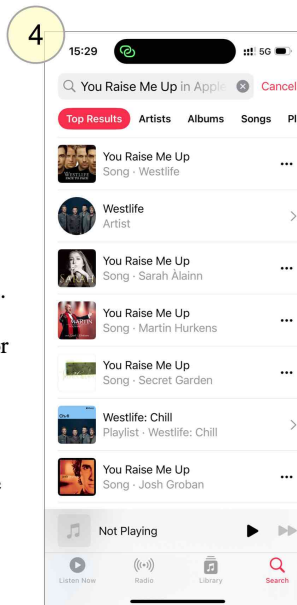


Figure 12. More generated samples of CogAgent.

**① CogAgent:**

**Plan:** 1. Locate the Music app icon.
2. Open the app to find the music library.
3. Search for 'You Raise Me Up' by using the search function within the app.
**Next Action:** Click on the Music app icon located in the second row from the bottom, third column from the left.

**② CogAgent:**

**Plan:** 1. Tap on the search bar at the top of the screen.
2. Enter the song title into the search field.
3. Execute the search by pressing the search button or enter key on the keyboard.
**Next Action:** Move the cursor to the search bar and click it to activate the text entry mode for typing in the song name 'You Raise Me Up'.

**③ CogAgent:**

**Plan:** 1. Tap on the search bar at the top of the screen.
2. Enter the text 'You Raise Me Up' into the search field.
3. Execute the search by pressing the search button or enter key on the keyboard.
**Next Action:** Click on the search bar where it says 'Artists, Songs, Lyrics, and...', then type 'You Raise Me Up'.

**④ CogAgent:**

**Plan:** 1. Identify the search bar at the top of the screen.
2. Review the list of songs that appear as a result of the search query.
**Next Action:** Look at the search results section just below the search bar to find the song 'You Raise Me Up'.

Figure 13. More generated samples of CogAgent.

**User:** What's shown in this picture? Describe in detail.

**CogAgent:** The screenshot shows a character standing in the game Genshin Impact, looking at a teleport waypoint. The HUD displays various elements: on the top left is a mini-map with icons indicating different locations; on the top right are character portraits and health bars for Raiden Shogun, Bennett, and Kaedehara Kazuha, along with an UID number; below these is a quest prompt 'Go to Merusea Village' with a distance indicator of 284m; in the center is the glowing Teleport Waypoint; at the bottom are the character's level (Lv. 90), health points (25123/3461), stamina (623), and other game UI elements like food, a key, and a compass.

**User:** I'm looking for guidance on how to "Go to the teleport waypoint"?

**CogAgent:** **Plan:** 1. Identify the direction to the teleport waypoint from the current location on the map.
2. Navigate through the environment or use in-game transportation methods to reach the destination. 3. Confirm arrival at the teleport waypoint and prepare for the next step of the journey within Genshin Impact.
**Next Action:** Move the cursor over the 'Teleport Waypoint' icon located towards the center right of the screen, just above the character's feet, and left-click to activate it.
**Grounded Operation:** [generic] Teleport Waypoint -> CLICK at the box [[506,497,686,547]]

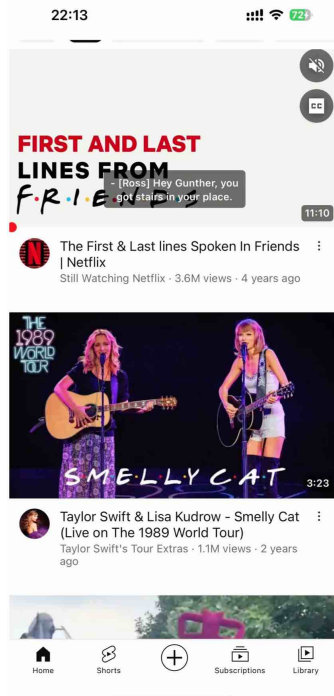**User:** Who are the teammates on the current mission?

**CogAgent:** The teammates on the current mission are Raiden Shogun, Bennett, and Kaedehara Kazuha.

**User:** I'm looking for guidance on how to "Find out how to view a detailed map"?

**CogAgent:** **Plan:** 1. Identify the map icon on the screen.
2. Click on the map icon to open the detailed map view.
**Next Action:** Move the cursor over the top left corner of the screen and click on the compass-like icon with 'N' inside it.



**User:** When did the Smelly Cat video post?

**CogAgent:** The Smelly Cat video post was made 2 years ago.

**User:** How long is the Smelly Cat video?

**CogAgent:** The Smelly Cat video is 3 minutes and 23 seconds long.

**User:** What should I do if I want to watch the Smelly Cat video (in the middle of the screen) and share it with my friend Rena?

**CogAgent:** **Plan:** 1. Play the Smelly Cat video by Taylor Swift & Lisa Kudrow to ensure it is fully loaded and ready for viewing.
2. Locate the share icon on the screen, which is typically represented by a symbol such as a paper airplane or an arrow pointing outwards from a box.
3. Tap the share icon to open sharing options.
4. Find and select the option to send via messaging app, like WhatsApp, to share the video with friends.
5. In the messaging app, locate the contact named Rena.
6. Send the video link to Rena through the messaging app.
**Next Action:** Touch the area in the middle of the screen where the 'Smelly Cat' thumbnail is displayed, corresponding to [Live on The 1989 World Tour] by Taylor Swift & Lisa Kudrow - Smelly Cat. This action is intended to start playing the video.
**Grounded Operation:** tap [[701,498]]

Figure 14. More generated samples of CogAgent.

**User**: Describe the image in detail.

**CogAgent**: The image is a screenshot of Twitter posts from the official Disney account. In the first post, there's an emoji-filled statement that reads 'So many iconic Disney and Pixar films were released in November! Take this quiz to find out how much you know about them.' This tweet was posted on 2023/11/15 at 05:02 and has received 46 reposts, 2 quotes, 319 likes, and 13 bookmarks. Below it, another tweet states 'What animal does Robin Hood disguise himself as when competing in the tournament of the golden arrow?' with multiple choice options like Duck, Pig, Stork, Wolf.

**User**: Which option has the highest number of votes?

**CogAgent**: The 'Stork' option has the highest number of votes with 61%.

**User**: How can I "Follow Disney's Twitter account"?

**CogAgent**: **Plan:** 1. Locate the 'Follow' button on the screen.
2. Click the 'Follow' button to follow Disney's Twitter account.
**Next Action:** Move the cursor over the 'Follow' button next to the Disney logo and click it.
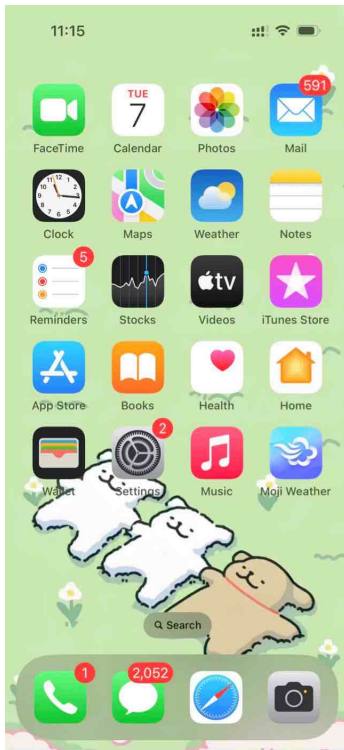
**User**: Describe the image in detail.

**CogAgent**: The image showcases a tweet from an account named 'Harvey' with the handle '@harveyyross'. The tweet contains an illustration of a young girl hugging a large dog. Above this, there is another text that reads: 'Some people will never understand how much someone can love a dog, but that's okay the dog knows...'. Below the illustration, there are statistics showing the tweet has been viewed 10.8K times, liked by 241 users, retweeted 6 times, and bookmarked by 6 users. Additionally, there is a reply to this tweet can 1 from an account named 'Farm Girl' with the handle '@far...', which says: 'Awww, that's so sweet. And yes,' indicating agreement o timesr appreciation for the original message.

**User**: Does the tweet have any reply? What is it?

**CogAgent**: Yes, there is a reply to the tweet from an account named 'Farm Girl' with the handle '@far...'. The reply reads: 'Awww, that's so sweet. And yes,' indicating agreement or appreciation for the original message.

Figure 15. More generated samples of CogAgent.

**User:** Describe the screen in detail.

**CogAgent:** The picture shows a mobile phone screen with various icons arranged in a grid. There are 24 app icons on the screen, including three cute cartoon bears at the bottom of the screen. These bears have different expressions and colors. In addition to the app icons, there is also a digital clock displaying "11:15" at the top left corner of the screen.

**User:** What apps do I have?

**CogAgent:** The apps on the phone are divided into two categories: system apps and third-party apps. System apps include "FaceTime", "Calendar", "Photos", "Mail", "Clock", "Maps", "Weather", "Notes", "Reminders", "Stocks", "Videos", "iTunes Store", "App Store", "Books", "Health", "Home", "Music", and "Moji Weather". Third-party apps include "Q Search" and other unidentified apps, which cannot be identified from their icons alone.

**User:** What should I do if I want to see photos taken 2 days ago?

**CogAgent:** You should tap on the "Photos" app, and it will display all your photos. If you want to see photos taken 2 days ago, you can select a date from the options at the top of the screen or use the search function in the Photos app to find specific photos.

**User:** If I want to take a photo, what should I do?

**CogAgent:** You should tap on the "Camera" app to take a photo.

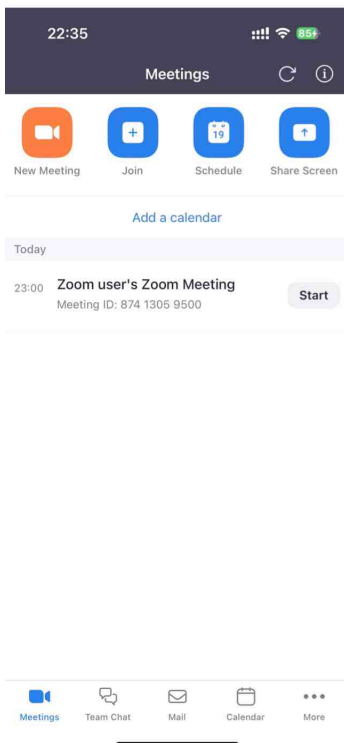**User:** What should I do if I want to call Susan?

**CogAgent:** You should tap on the "Phone" app to call Susan.

**User:** How can I know Susan's number?

**CogAgent:** You can call Susan by searching for her number in the phone's contact list or on your phone's address book. If you have saved her number, you can directly dial it from there.

**User:** What should I do if I want to download youtube?

**CogAgent:** You should tap on the "App Store" app, and then search for YouTube in it. Once you find the official YouTube app, you can download and install it to watch videos.

**User:** Is there a preferred method for "Join another zoom meeting with id 13572468"?

**CogAgent:** **Plan:** 1. Identify and select the 'Join' option for an existing or new meeting ID.
2. Enter the specific Meeting ID when prompted to join a particular Zoom meeting.
3. Confirm any additional required details such as passcodes if necessary, and then join the meeting successfully.
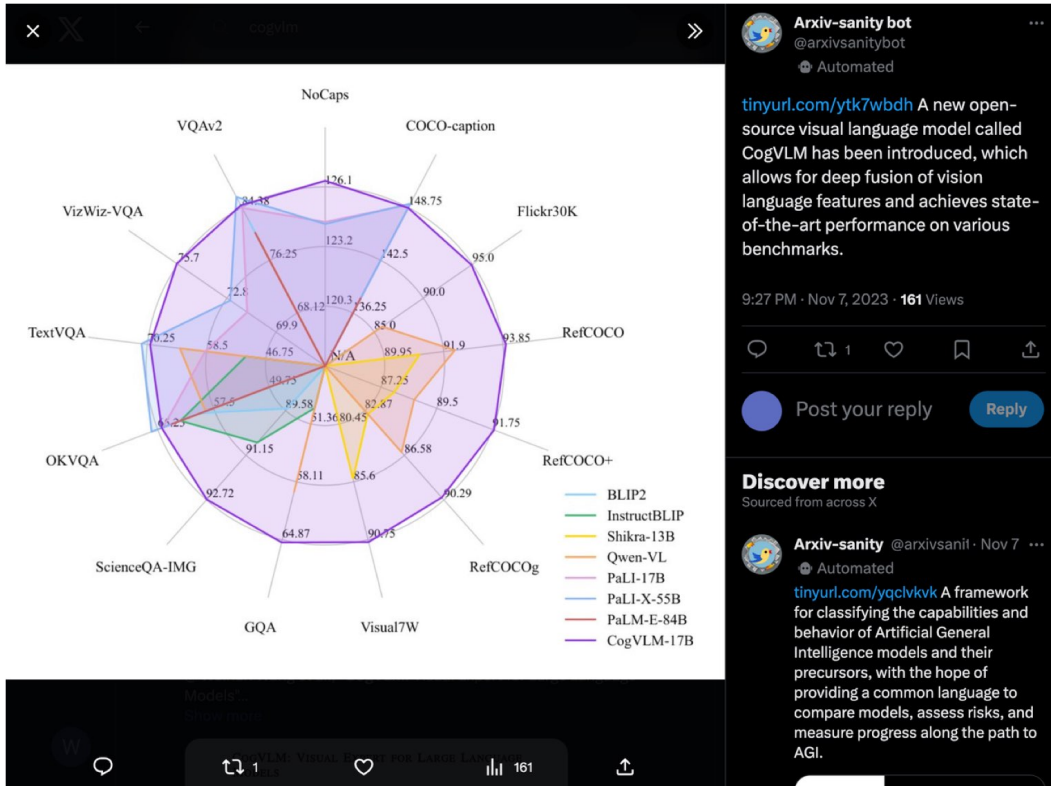**Next Action:** Tap on the 'Join' button located in the middle of the screen, which corresponds to the coordinates [x,y]=[0.456,0.18]. This action is intended to initiate the process of joining a Zoom meeting.
**Grounded Operation:** tap [[456,180]]

**User:** In how many minutes will the next meeting start? Please think step by step.

**CogAgent:** The next meeting will start in 23 minutes. This is determined by subtracting the current time (22:35) from the scheduled start time of the meeting (23:00).

Figure 16. More generated samples of CogAgent.

**User**: What's the performance on CogVLM-17B on Visual7W?

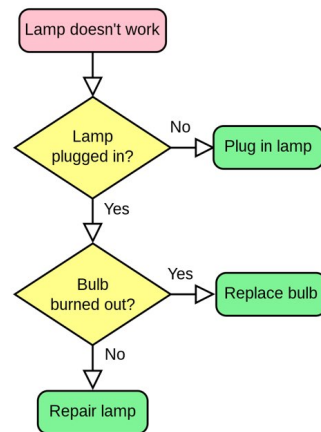**CogAgent**: The performance on CogVLM-17B on Visual7W is 90.75.

**User**: What's Arxiv-sanity's comment?

**CogAgent**: Arxiv-sanity's comment is: 'A new open-source visual language model called CogVLM has been introduced, which allows for deep fusion of vision language features and achieves state-of-the-art performance on various benchmarks.'



**User**: I am getting one funghi pizza and one prosciutto pizza. How much should I pay in total?
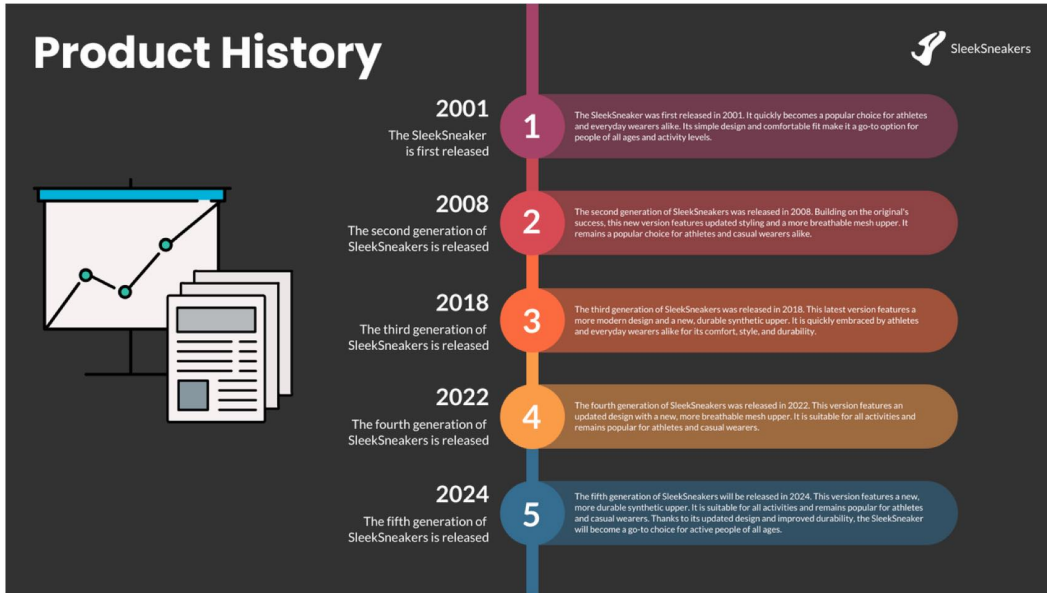
**CogAgent**: The funghi pizza costs 6,90\u20ac and the prosciutto pizza costs 7,90\u20ac. Therefore, the total cost for both pizzas would be 14,80\u20ac.



**User**: What is the solution if the lamp is plugged in and the bulb has burned out?

**CogAgent**: If the lamp is plugged in and the bulb has burned out, the solution is to 'Replace bulb'.

Figure 17. More generated samples of CogAgent.

**User**: Describe this image.

**CogAgent**: The image is a timeline showcasing the product history of 'SleekSneakers'. It starts from 2001, where 'The SleekSneaker' was first released, and ends in 2024, where the fifth generation of SleekSneakers was released. The timeline highlights key milestones, such as the introduction of the second generation of SleekSneakers in 2008, the fourth generation in 2022, and the fifth generation in 2024. Each milestone is accompanied by a brief description.

**User**: Read the content in 2022.

**CogAgent**: In 2022, the fourth generation of SleekSneakers was released. This version features an updated design with a new, more breathable mesh upper. It is suitable for all activities and remains popular for athletes and casual wearers.



**User**: In which country was this photo taken?

**CogAgent**: The photo was taken in Australia, as the Australian Open is one of the major tennis tournaments held there.



**User**: Who is the person on the left?

**CogAgent**: The person on the left is Keira Knightley.

Figure 18. More generated samples of CogAgent.

| method | top-1 err. | top-5 err. |
| --- | --- | --- |
| VGG [41] (ILSVRC'14) | - | 8.43[†] |
| GoogLeNet [44] (ILSVRC'14) | - | 7.89 |
| VGG [41] (v5) | 24.4 | 7.1 |
| PReLU-net [13] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |
| ResNet-152 | **19.38** | **4.49** |

Table 4. Error rates (%) of **single-model** results on the ImageNet validation set (except [†] reported on the test set).

**User**: Convert it to the markdown format.

**CogAgent**: | method | top-1 err. | top-5 err. |
| -------- | ------------ | ----------- |
| VGG [41] (ILSVRC'14) | - | 8.43 |
| GoogLeNet [44] (ILSVRC'14) | - | 7.89 |
| VGG [41] (v5) | 24.4 | 7.1 |
| PReLU-net [13] | 21.59 | 5.71 |
| BN-inception [16] | 21.99 | 5.81 |
| ResNet-34 B | 21.84 | 5.71 |
| ResNet-34 C | 21.53 | 5.60 |
| ResNet-50 | 20.74 | 5.25 |
| ResNet-101 | 19.87 | 4.60 |
| ResNet-152 | 19.38 | 4.49 |

### Dataset 1

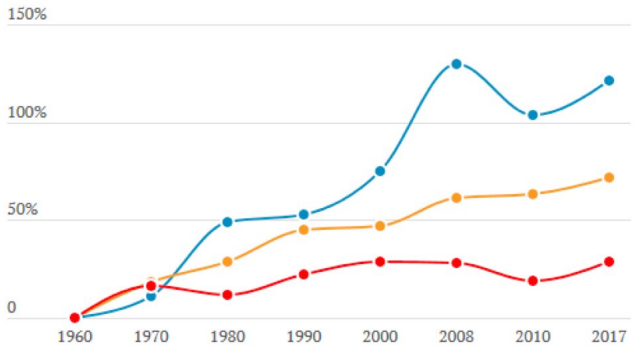| Department | Name | Salary |
| --- | --- | --- |
| Administration | Emily | $2,000 |
| Finance | Nick | $1,800 |
| Marketing | John | $2,600 |
| Marketing | Gloria | $1,800 |
| Administration | Jane | $2,500 |
| Finance | Max | $2,000 |
| Administration | Tyler | $3,000 |
| Finance | Wilkins | $2,500 |
| Marketing | Miles | $2,800 |
| Administration | Sabine | $1,900 |

**User**: Based on the table, how much is Wilkins earning?

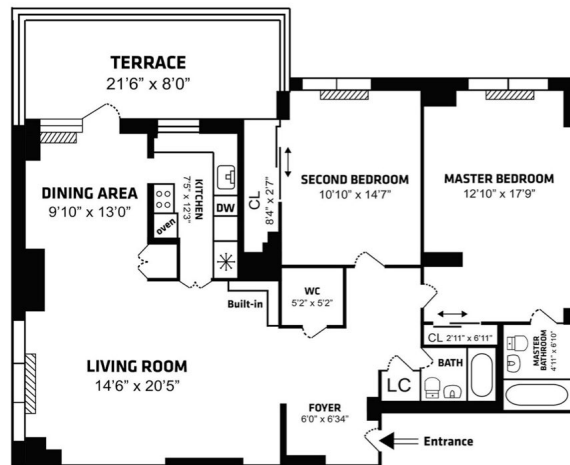**CogAgent**: Wilkins is earning $2,500.

**US Nationwide:**



**User**: Which year has the highest growth rate of median house price?

**CogAgent**: The year 2008 has the highest growth rate of median house price, which is 150%.



**User**: What are the appliances in the kitchen in this floorplan?

**CogAgent**: The kitchen in this floorplan has a built-in oven, a built-in dishwasher, and a built-in microwave.

Figure 19. More generated samples of CogAgent.