TRISHUL: Towards Region Identification and Screen Hierarchy Understanding for Large VLM based GUI Agents

Supplementary Material

1. Model Specifications and Endpoints

Since all our work leverages closed-source models like GPT-4V, GPT-4o, and Claude, we mention the model identifiers that we use for our API calls for clarity. For GPT-4V - "gpt-4-vision-preview", For GPT-4o - "gpt-4o-2024-08-06", and for Claude - "claude-3-5-sonnet-20241022". Unless otherwise noted, all experiments are conducted with a temperature setting of 0.0.

2. Hierarchical Screen Parsing Details

2.1. IoS Score

Similar to IoU score we define an IoS score as:

The IoS (Intersection over Size) score is a measure used to evaluate the overlap between two bounding boxes, typically in the context of object detection. It calculates the ratio of the intersection area between two boxes to the area of the first box. IoS(A, B) (also written as IOS_A), a score of 0.5 means 50% of A intersects with B.

2.2. Filtering Redundant Bounding boxes

The output of EasyOCR + SAM model combined is extremely cluttered (see 1) and contains numerous overlaps and false positive detections from both models. We deploy the following steps to parse the outputs of SAM and OCR (local elements as referred to in the main script) together.

• Generate GROI, Icon and Button Candidate Proposals Classify all SAM boxes based on A_{thresh-GROI}, A_{thresh-icon}, A_{thresh-button}. Let B represent the set of bounding boxes detected in the GUI.

Global Region of Interest (GROI) Candidates: The set of boxes with an area greater than the GROI threshold is given by:

$$GROI = \{b \in B \mid Area(b) > A_{thresh-GROI}\}\$$

Icon Candidates: The set of boxes with an area between the Button and Icon thresholds is defined as:

 $Icon = \{b \in B \mid A_{thresh-button} < Area(b) < A_{thresh-icon}\}$

Button Candidates: The set of boxes with an area less than the Button threshold is:

Button = { $b \in B$ | Area(b) < $A_{thresh-button}$ }

• **Remove False Positive Text Bounding Boxes:** Remove text boxes using a predefined dictionary, that are likely OCR mis-detections for icons. These texts usually contain only special characters or short, meaningless words. If a word contains one of these characters and has a length of less than ; 5 that text bbox is ignored.

Characters/Words to ignore:

- Remove Icons Inside or Overlapping Text Bounding Boxes: Remove icon bounding boxes that are either inside or intersect with text boxes, as they are likely to be text misidentified as icons by SAM.
- Filter Square-like Icon Bounding Boxes: Keep only icons that are roughly square-shaped, based on a specific aspect ratio range of [0.7, 1.3].
- Remove Redundant Icon and Button Bounding Boxes: Remove icon bounding boxes that are redundant, i.e., those that are inside or significantly overlap with IoS > 0.6with other icons or text boxes.

2.3. Non Max Suppression for GROIs

- **Reject boxes with low Information score** S If the current bounding box has an information score $S < S_{thresh}$ it is rejected. S_{thresh} is set to 10 for the ScreenPoint and Read task and 25 for action grounding task.
- **Reject Overlapping BBoxes:** If the current bounding box intersects with a previously selected bounding box with a higher Infrormation score and $IoS_{current} > IoS_{overlap-thresh}$ it is rejected. $IoS_{overlap-thresh}$ thresh is set to 0.5 for visual grounding task and 0 for ScreenPR task.
- Reject Contained BBoxes (Smaller GROIs Inside Larger): If the current bounding box is inside a previously selected bounding box with a higher Information

Score and if $IoS_{current} > IoS_{inside-thresh}$ it is rejected. $IoS_{overlap-thresh}$ thresh is set to 0.5 for visual grounding task and 0 for ScreenPR task.

• **Reject Engulfing Bboxes (Larger GROIs Inside Smaller):** If the current bounding box completely engulfs a bounding box with a higher Information Score then it is rejected.

2.4. GROI analysis for ScreenSpot and Visual-WebBench

We plot two additional statistics for the detected GROI's through our HSP block. In Figure 3 we plot the average number of GROIs per image, across the three different sub-categories of ScreenSpot and the full dataset of Visual-WebBench. We observe that GUI screenshots from mobiles have the lowest average number of GROIs per image. This is due to the fact that mobile regions are not semantically coherent, therefore lesser number of GROIs are generated.

In Figure 4 we plot the average of the total area covered by all GROIs in an image to the total area of the image. Mobile GUI screenshots have the least dense GROI coverage, due to the fact that we also detect fewer GROIs in mobile screenshots. These studies further validate the fact that GROIs are not as useful for mobile GUI's however they offer more benefit for PC and web based GUIs.



Figure 1. Candidate bounding boxes generated from SAM + OCR to the left and the corresponding HSP results (Icon + text + picture) to the right

You are a smart screen reader that outputs concise natural lanauage to answer auestions from users based on the area (box 1) pointed out by the user shown as a red dot on the screen. The red dot is inside the box 1 in the first image, at $(x,y) = (\{\cdot, \})$, where x and y are the normalized coordinates. Note: box 1 is the box with label 1 and box 2 is the box with label 2, box 1 is located inside box 2

Note: the first image shows what is inside box 2, and the second image shows the complete screen. Note: if the user asks about the location, focus on the layout, explain where box 1 is in box 2 and then explain where box 2 is in the overall complete screen.

Note: don't mention box 1, box 2 or the red dot in the output.

User question: (1) what is this? (2) where it is located in the screen? Your output should in format (1) ... (2) ..





Figure 3. Distribution of Number of GROIs per image for ScreenSpot and Visual WebBench



Figure 4. Distribution of Total GROI area / Image area for ScreenSpot and Visual WebBench

You are an expert User Interface (UI) navigation agent. You will receive a screenshot from a web, mobile, or PC interface, an instruction, and cropped Global Regions of Interest (GROIs) from the screenshot. Each GROI is annotated on the full image for context.

Your task:

1. Analyze each GROI individually using the full image for context.

 Briefly describe the content and functionality of each GROI.
 Identify the top 3 GROIs most likely to contain the UI element needed to complete the instruction successfully, listing them in order of relevance. 4. Incase of Overlap between 2 GROIs give preference to the smaller GROI.

Output format:

Description:

- ID 1: <Brief description of GROI ID 1> ID 2: <Brief description of GROI ID 2>
- ... (Repeat for each GROI)

Proposed GROIs

[<ID of the most relevant GROI>, <ID of the second most relevant GROI>, <ID of the third most relevant GROI>] Output the IDs as a list of Integers.

Figure 5. Prompt for instruction guided GROI Proposal generation

Task: Tor each instruction think carefully and follow these four steps in order. 1) Generate a brief description of the functionality and content of the UI screenshot provided. 2) For each 'text' type UI element. 2) For each 'text' type UI elements, and Its content provided in the input prompt, use the UI screenshot are reference to spell check the content, in case of errors, fix the errors in place. 2) For each "text' type UI elements append its functionality in context of the associated elements to the description. (e.g., "search' text is a search bar, "thtps://..." text link is a URL input field, "Jon" name text is an input field for name.) 2.) Classify the element as: 2.0.11 "text-or: (text contered if the text element is associated with a nearby "image" type element. (a.g., "search" text as search bar, "https://..." text link is a URL input field, "Join" nome texts is an input field for name.)
(a.g., "search" text as a search bar, "https://..." text link is a URL input field, "Join" nome texts is an input field for name.)
(a.g., "search" text as a search bar, "https://..." text link is a URL input field, "Join" nome text is an input field for name.)
(b.g., "search" text as a search bar, "https://..." text link is a URL input field, "Join" nome text is an input field for name.)
(c.g., "search" text as a search bar, "https://..." text link is a URL input field, "Join" nome text is an input field for name.)
(c.g., "search" text as a search bar, "https://..." text link is a URL input field, "Join" nome text is an input field for name.)
(c.g., "search" text as a search bar, "https://..."
(c.g., "text as the text as the text and the mapped to one "image" type element to max.
(c.g., "text-s", it text contable) if the "text element is NOT associated with any metry" image" type elements to any functionality.
(c.g., "fieth text-s", "type elements any one the description.
(c.g., "fieth text-s", "type elements any one the description.
(c.g., "fieth text-s", "type elements any one the description.
(c.g., "fieth text-s", "type elements any one the description in a start of "image" type elements the search text. How new text associated with then in the UI (e.g. Wifi icon with a wifi text at its bottom.), in a 1.1, "icon-s", "(icon-s", "icon-s", "icon-s", "icon-s", icon-search associated with any text element is search associated with any text element is associated with any text element is the UI (e.g. Wifi icon with a wifi text at its bottom.), in a 1.1, "icon-s", icon-search associated with any text element is associated with then any text element is associated with any text element is associated with any text element is associated with then any text element is associated If an icon does not have a bounding box due to imperfect image detection, do not try to generate its description, only focus on the annotated elements. 2.) Do not give redundant icon desc Note: The following things are very important for a successful execution of the task:\n 1.) Generate the output for all "text" elements first and then the "image" elements. 2.) Maintain spatial and semantic relationships between UI elements to correctly determine functionality. 3.) The ID tags associated with "image" type elements are unordered, so avoid sequencing blas when associating IDs in the prompt to IDs annotated in the screenshot. 4.) Match elements to ID tags using only the color of the ID tags, color of the bounding box, and their spatial proximity. 5.) The number of elements, bounding boxes, and IDs in the input must match the output exactly. Output Format:\n Description: Brief description of the content and functionality of the screenshol Annotations:\n nnatations:in [[D][X1, y1, X2, y2][text-p][D] of associated "image" element][Speil-corrected content for text type element]|n [[D][X1, y1, X2, y2][text-s][Speil-corrected content for text type element, context (if necessary]]]n [[D][X1, y1, X2, y2][text-s][Speil-corrected content for text type element][Functionality Description][n [[D][X1, y1, X2, y2][text-s][Speil-corrected content for text type element][Functionality Description]]n

(D) (Intrinsic local protocol and o Composition in the Composition in

You are an expert User Interface (UI) novigation agent. You will be given a screenshot of a UI of a mobile / desktop / web application. The screenshot will contain bounding box annotation for UI elements such icon/picture/text etc type elements. Additionally you will also be given an input prompt, containing information about the Hu Elements. In Your task is to, improve the prompt by describing the functionality of the various elements, correcting mistakes and establishing spatial and semantic relations between different types of elements tagged in the UI.)n

1.) There are 2 types on UI elements annotated onto the screenshot, "image" (generated using semantic segmentation) and "text" (generated using OCR). c) The data 2 yes of the elements during the schematory in mage (generated using elements during elements) and text (generated using durin, 2) foot if lement will have a bounding box and a numeric ID associated with it.
3) For "image" type U elements, each bounding box and its respective ID tag will share the same color, The ID will be at the top-left corner of the bounding box.
A) For "text" type U elements, the associated bounding box and its respective ID tag will share the same color, The ID will be at the top-left corner of the bounding box.
A) For "text" type U elements, the associated bounding box will be "ted" in color, the corresponding ID tag will NOT be annotated onto the image to avoid cluttering. In
5.) An input prompt with the description of all UI components in the following format will also be provided.

(Repeat for all UI elements)\n

Here are some examples

Input Description:\n

Task:\n

Figure 6. Prompt for SEED

5

You are a helpful AI assistant. You will be given a screenshot of a User Interface (UI) of a mobile / desktop / web application, along with an instruction. Your task is to find the correct element from the UI that will lead to a successful

completion of the instruction.\n Input Format:\n

- Input romation [ID][x1, y1, x2, y2][text-p][ID of associated "image" element][Content inside the text element]|n [ID][x1, y1, x2, y2][text-s][Content inside the "text" element][Functionality Description]\n [ID][x1, y1, x2, y2][text-a][Content inside the "text" element][Structural and Semantic description of icon][Functionality description]\n
- $[D][x_1, y_1, x_2, y_2][icon-y_1]$ based and Semantic description of icon][Functionality description][n [D][x_1, y_1, x_2, y_2][icon-y_1] [Structure] [Description of the image and context][n

[ID][x1, y1, x2, y2][irrelevant][]\n Input Description:\n

1.) Prompt:

1.1) There are eight types on UI elements annotated onto the screenshot, "text-a" (text actionable), "text-p" (text paired), "text-s" (text standalone), "icon-p" (icon paired), "icon-s" (icon standalone), "picture",

"irrelevant".\n Tretevant. (II) 1.2) For each element [ID] represents the associated ID tag, [x1, y1, x2, y2] are the endpoint coordinates of the rectangular bounding box associated with the element. 1.3) "text-p"; (text paired) are the "text" elements associated with a nearby "icon" type element. 1.4) "text-o"; (text actionable) are the "text" elements associated with a functionality. 1.5) "text-s"; (text standalone) are the "text" elements NOT associated with any nearby "icon" type elements or any functionality.

- 1.5) "text-SAM": (text SAM) are the "text" elements NOT detected by OCR but SAM Segmentation algorithm
- 1.6) "icon-p": (icon paired) are the "icon" elements associated with a nearby "text" type element.
 1.7) "icon-s": (icon standalone) are the "icon" elements NOT associated with a nearby "text" type element.

- 1.8) "picture": are the images embedded in the UI.1.9) "irrelevant": elements are false positive detections and should be ignored. 2.) Input Screenshot:

2.1) For "icon", "picture" "text-SAM" and "irrelevant" type UI elements, each bounding box and its respective ID tag will share the same color, The ID will be at the top-left or the bottom-left corner of the bounding box.\n

2.2) For "text-p", "text-a" and "text-s" type UI elements, the associated bounding box will be "red" in color, the corresponding ID tag will NOT be annotated onto the image to avoid cluttering. \n

Task:\n

Using the descriptions of the UI elements provided in the annotation, their classifications, and the screenshot as reference

correct element along with the reasoning behind their selection, clicking which will lead to the completion of the instruction in a single shot.

Note:\n

The following pointers are very important for a successful execution of the task:\n 1.) The selected element should lead to the completion of the instruction in a single shot.

- 2.) The correct element does not necessarily correspond to the element with the strongest string match, lay strong emphasis on the functionality requirement of the instruction. 3.) Correct semantic, structural, and functional understanding of elements should be established using the descriptors and image provided.
- Preference should always be given to icons and actionable text (e.g. texts associated with input fields, search bars, buttons, drop downs etc.) over static text.
 Pay attention on the action words in the instruction like "Search", "Input" "Select", "Enter" to determine the type of element required for the task.

Output Format:\n

{"id": <ID of the correct element>, "reasoning": <semantic and functional reasoning behind the elements selection>}

Do not output anything else

Figure 7. SoM grounding Prompt for ScreenSpot and VisualWebBench

You are an expert at completing instructions on Android phone screens.

You will be presented with two images. The first is the original screenshot. The second is the same screenshot with bounding box annotations.

You will additionally also be provided with a prompt with the descriptions of all the elements tagged on the screen

1. Input Screenshot Description:

2.1) For all "icon"/"picture" type UI elements, each bounding box and its respective ID tag will share the same color, The ID will be at the top-left or the bottom-left corner of the bounding box. 2.2) For all "text" type UI elements, the associated bounding box will be "red" in color, the corresponding ID tag will NOT be annotated onto the image to avoid cluttering. \n

- 2. Input Prompt Format:
- [ID][x1, y1, x2, y2][text-p][ID of paired "icon"/"picture" element][Content inside the text element]
- [ID][x1, y1, x2, y2][text-s][Content inside the "text" element]
- [ID][x1, y1, x2, y2][text-a][Content inside the "text" element][Functionality Description]
- [ID][x1, y1, x2, y2][icon-p][ID of paired "text" element][Structural and Semantic description of icon][Functionality description] [ID][x1, y1, x2, y2][icon-s][Structural and Semantic description of icon][Functionality description]
- [ID][x1, y1, x2, y2][picture][Description of the image and context]
- [ID][x1, y1, x2, y2][irrelevant][]
- 3. Annotation Description:
- -s stands for "static" -a stands for "actionable", -p stands for "paired".
- 4. Task:
 - If you decide to click somewhere, you should choose the numeric idx that is the closest to the location you want to click.
- The screenshot are most likely an intermediate step of this instruction, so in most cases there is no need to navigate back home.
- You should decide the action to continue this instruction.
- Here are the available actions:

rere dre the outside actions:
{"action_type": "click", "idx": <element_idx chosen from the second screen>, "bboxes": <bbox coordinates of the the chosen element as a singular list>}
{"action_type": "hove", "text": <the text to enter>}
{"action_type": "navigate_home"}
{"action_type": "navigate_back"}

- {"action_type": "scroll up"} {"action_type": "scroll down"} {"action_type": "scroll left"}

- {"action_type": "scroll right"}
- Your final answer must be in the above format.

Figure 8. Agentic task following prompt for AITW

You are an expert at completing instructions on websites. You will be given an instruction and two images. The first is the original screenshot.

The second is the same screenshot with some bounding box annotations. Additionally you will also be given a prompt giving details about the annotations on the screenshot. The format for the same is specified below. 1. Input Prompt Format:

[ID][x1, y1, x2, y2][text-p][ID of associated "image" element][Content inside the text element]

[10][x1, y1, x2, y2][text-p][Content inside the "text" element] [10][x1, y1, x2, y2][text-b][Content inside the "text" element] [10][x1, y1, x2, y2][text-b][Content inside the "text" element] [10][x1, y1, x2, y2][text-b][DI of associated "text" element][Structural and Semantic description] [10][x1, y1, x2, y2][text-b][DI of associated "text" element][Structural and Semantic description of icon][Functionality description] [10][x1, y1, x2, y2][text-b][Content inside the "text" element][Structural and Semantic description]

[ID][x1, y1, x2, y2][picture][Description of the image and context]

[ID][x1, y1, x2, y2][irrelevant][]
2. Input Screenshot Description:

2.1) For all icon/picture/irrelevant type UI elements, each bounding box and its respective ID tag will share the same color, The ID will be at the top-left or the bottom-left corner of the bounding box.\n 2.2) For all "text" type UI elements, the associated bounding box will be "red" in color, the corresponding ID tag will NOT be annotated onto the image to avoid cluttering. \n

3. History

You will also be provided a brief summary history of all actions taken so fa

Task:

Choose the best action to continue the instruction given the present screen. If you decide to click somewhere, Using the 2 UI screenshots and the annotation descriptors in the prompt as reference suggest the top 3 elements (text or icon) in decreasing order of relevance to the instruction, along with the reasoning behind their selection.

Here are the available actions:

"action_type": "CLICK", "lax": <element_idx of chosen element from the annotated screen>, "bbox": <bbox coordinate of chosen element from the annotated screen>}
("action_type": "TYPE", "lax": <element_idx oof chosen element from the annotated screen>, "bbox": <bbox coordinate of chosen element from the annotated screen>}
("action_type": "TYPE", "lax": <element_idx oof chosen element from the annotated screen>, "bbox": <bbox coordinate of chosen element from the annotated screen>}

Note: 1. **Important**: The sequence of screenshots are sampled from a pre-executed trajectory, and may not necessarily be a result of the previous action, choose the best action given the current screen. 2. Use {"action_type": "TYPE"} when a specific value needs to be typed into the input field

3. Output the answer as string and not json formatted

Output Format:

Action Reasoning: <Reasoning behind the selection of action>

UI Element Selection Reasoning [ID][x1, y1, x2, y2][element type]: <Semantic and Functional reasoning behind the selection of the most relevant UI component>

(Repeat for top 3 elements)

Action:

("action_type": "CLICK", "idx": <element_idx of chosen element from the annotated screen>, "bbox": <bbox coordinate of chosen element from the annotated screen>}
("action_type": "TYPE", "idx": <element_idx of chosen element from the annotated screen>, "bbox": <bbox coordinate of chosen element from the annotated screen>, "bbox":
("action_type": "TYPE", "idx": <element_idx of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox":
("bbox coordinate of chosen element from the annotated screen>, "bbox";
("bbox coordinate of chosen element from the annotated screen>, "bbox";
("bbox coordinate of chosen element from the annotated screen>,

Figure 9. Agentic task following prompt for Mind2Web