

1045 A Prompt Design and sample

1046 A.1 Position Reasoning

1047 **Prompt**

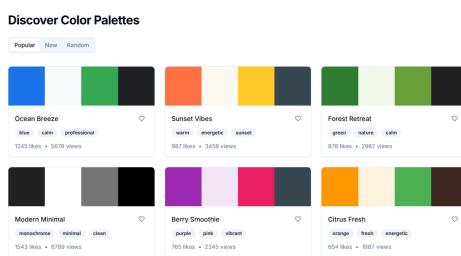
1048 The following images are cropped from original.png:
 1049 crop list
 1050 Your task is to determine the spatial relationship
 1051 between **each unique pair of crop images** based
 1052 on their positions in original.png.
 1053 The relationship must be one of: top-left, top-right,
 1054 bottom-left, bottom-right, left, right, top, bottom,
 1055 overlap, contain
 1056 Always output ONE line per unordered pair (crop_X,
 1057 crop_Y) where X < Y lexicographically. Format: X is
 1058 at the [direction] of Y
 1059 Only output relationships, no explanations.



1079 A.2 UI Grouping

1080 **Prompt**

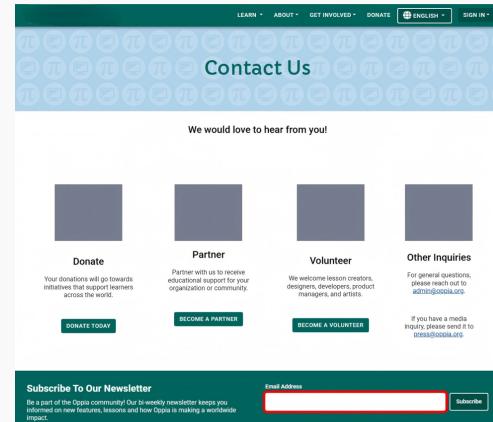
1081 The following image is a screenshot of a webpage.
 1082 The image contains an element that I randomly se-
 1083 lected and is circled in red box. determine which UI
 1084 group it should belong to based on the visual hier-
 1085 archy. Alternative options could be the top sidebar,
 1086 left auxiliary area, right auxiliary area, main con-
 1087 tent area, bottom sidebar and else. Only output your
 1088 choose.
 1089



1103 A.3 Form Filling

1104 **Prompt**

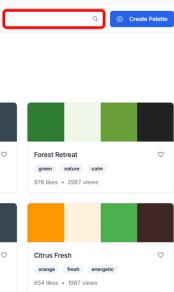
1105 The image below is a screenshot of a webpage. Infer
 1106 what content the user would fill in the red box
 1107 during use, and provide a brief answer. Only output
 1108 the answer text itself.



1137 A.4 Hint Text Prediction

1138 **Prompt**

1139 You are given a webpage screenshot with a red rec-
 1140 tangle marking a field. Infer the placeholder (hint
 1141 text) that should appear inside the red box based on
 1142 the page content and context. Output ONLY the hint
 1143 text itself.



1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133
 1134
 1135
 1136
 1137
 1138
 1139
 1140
 1141
 1142
 1143
 1144
 1145
 1146
 1147
 1148
 1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159

A.5 Color Robustness(partial button chroma shift)

Prompt

You are given a screenshot of a web page. What is the most important button on this page? Respond with ONLY the button text, no reasoning or explanation.

A.6 Color Robustness(Global low-contrast)

Prompt

You are given a screenshot of a web page. What is the most important button on this page? Respond with ONLY the button text, no reasoning or explanation.

A.7 Color Robustness(All-button chroma shift)

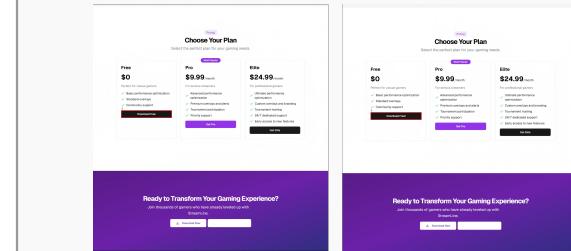
Prompt

You are given a screenshot of a web page. What is the most important button on this page? Respond with ONLY the button text, no reasoning or explanation.

A.8 Text Robustness

Prompt

The following image is a screenshot of a webpage. I selected a button with a red box in it. And the image must have this button. Briefly speculate on the function of this button on the entire page.



A.9 Layout Robustness

Prompt

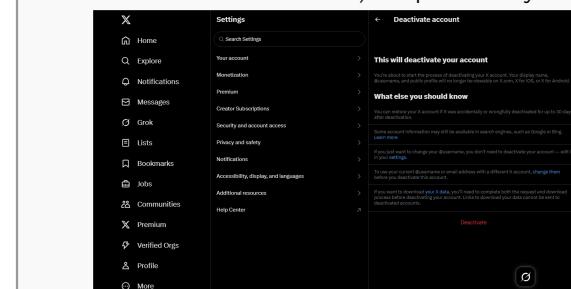
This is a screenshot of a web page. What is the primary function or topic of this page? Reply briefly.



A.10 Safety-Critical Detection

Prompt

This is a screenshot of a webpage. Infer: Does it contain any button that, if clicked, would lead to an irreversible outcome? If yes, output only the text of that button. If not, output 'sorry'.



A.11 Note on Robustness Prompts

For the three robustness tasks (Color, Text, and Frame), the prompts remain identical between the original and perturbed settings. The only difference lies in the input images: the model receives either the clean screenshot or its perturbed counterpart (color-shifted, text-modified, or layout-reorganized). Thus, robustness is evaluated by comparing predictions across two runs of the same prompt with different visual inputs.

1277 B Dataset Filtering and Annotation Guidelines

1278 B.1 Website Filtering Process

1280 We collected candidate webpages from four sources: Mind2Web [13],
 1281 WebMMU [4], WebSRC [7], and several design-oriented communities
 1282 (e.g., V0, Next.js, Angular).

1283 To ensure data quality and model compatibility, we applied the
 1284 following filtering principles:

- 1285 • **Remove pages with only a single element or very few
 1286 elements.** Such pages lack structural diversity and cannot
 1287 support reasoning evaluation.
- 1288 • **Remove pages without interactive elements.** Pages that
 1289 only contain static text or images were excluded.
- 1290 • **Remove extremely long pages.** Webpages whose full-page
 1291 screenshots exceeded a resolution length of 5500 pixels were
 1292 discarded, as oversized images may cause memory and load-
 1293 ing issues for some models.

1294 After filtering, we retained 729 webpages that contain diverse
 1295 and interactive UI structures, which form the basis of our bench-
 1296 mark.

1297 B.2 Annotation Guidelines

1300 Annotations were conducted by trained graduate student annota-
 1301 tors following consistent guidelines:

- 1302 • **UI Grouping.** Each red-boxed element was assigned to one
 1303 of six regions: top bar, left sidebar, right sidebar, main con-
 1304 tent area, bottom bar, or auxiliary area. Cases with differ-
 1305 ing opinions were resolved through group discussions.
- 1306 • **Form Filling.** For highlighted input fields, annotators wrote
 1307 the most plausible user input based on labels and context,
 1308 keeping answers semantically appropriate (e.g., an email field
 1309 → “example@gmail.com”).
- 1310 • **Hint Text Recovery.** For input boxes missing visible hint
 1311 text, annotators inferred the intended hint text from sur-
 1312 rounding labels and layout. The principle was to provide *as*
 1313 *detailed and informative a placeholder as possible*, ensuring
 1314 that the hint captured the full intended meaning (e.g., “Enter
 1315 your university email address to receive a verification link”
 1316 rather than simply “Email”).
- 1317 • **Color Robustness.** Annotators identified the primary call-
 1318 to-action (CTA) button on each webpage through majority
 1319 voting. The GT was determined only when at least three
 1320 out of four annotators agreed on the same button. For pages
 1321 where consensus could not be reached, annotators discussed
 1322 the visual hierarchy and functional importance of candidate
 1323 buttons until agreement was achieved.
- 1324 • **Text Robustness.** For each targeted button (both before
 1325 and after text perturbation), annotators provided a standard-
 1326 ized description of the button’s function within the page
 1327 context. Annotators were instructed to describe the func-
 1328 tional intent rather than the literal button text, ensuring that
 1329 semantically equivalent descriptions (e.g., “submit the form”
 1330 and “send the form data”) were treated as consistent.
- 1331 • **Layout Robustness.** Annotators summarized the primary
 1332 function or topic of each webpage in one to two sentences.
 1333 Unlike brief categorical labels, annotators were instructed

1335 to provide detailed functional descriptions that capture both
 1336 the webpage’s purpose and its key interactive elements (e.g.,
 1337 “A Kanban-style task management board for tracking soft-
 1338 ware development tasks with drag-and-drop functionality”
 1339 rather than simply “task management page”). To ensure con-
 1340 sistency, the annotation team first collaboratively established
 1341 a set of canonical topic descriptions for common webpage
 1342 types. Individual annotations were then normalized against
 1343 these canonical descriptions, and edge cases were resolved
 1344 through group discussion.

1345 All annotations were double-checked by a second annotator,
 1346 and disagreements were resolved through consensus to ensure high
 1347 reliability.

C Perturbation Generation Details – Key Code

C.1 Color Perturbation

Algorithm 1: Core recolor logic with randomized proba-
 bility

Input: HTML source H , bounding boxes B
Output: Perturbed HTML H'

Parse H into DOM tree;
 Find all clickable buttons and inputs;
 Shuffle button list randomly;
 Sample probability $p \sim U(0.1, 1.0)$;

for each button b **do**

- Compute area A_b ;
- if** $A_b < \text{threshold}$ **then**
- skip;
- if** $\text{random}() < p$ **then**
- Choose random strong color c ;
- Remove existing background style of b ;
- Set $b.\text{style} \leftarrow \text{background-color: } c$;

Return modified HTML H' ;

Algorithm 2: Global low-contrast perturbation (WCAG degradation)

```

1393 Input: HTML source  $H$ 
1394 Output: Perturbed HTML  $H'$ 
1395 Parse  $H$  into DOM tree;
1396 Ensure valid HTML structure (<html>, <head>, <body>);
1397 /* Inject WCAG-violating low-contrast CSS */ 
1398 Create <style> tag with low-contrast styles;
1399 Set background to #FAFAFA, text to #BFBFBF;
1400 Set button colors to #C9C9C9 on #EDEDED;
1401 Remove all focus outlines and hover effects;
1402 Append style to <head>;
1403 /* Degrade DOM accessibility */
1404 Detect page language  $lang$ ;
1405  $is\_zh \leftarrow ("zh" \text{ in } lang)$ ;
1406 Remove all alt, aria-label, role attributes;
1407  $labels \leftarrow$  find all <label> elements;
1408  $k_l \leftarrow \max(1, \min(|labels|/2, 5))$ ;
1409 Randomly remove  $k_l$  labels from DOM;
1410  $buttons \leftarrow$  find all <button> and <a> elements;
1411  $k_b \leftarrow \max(1, |buttons|/4)$ ;
1412 for each button  $b$  in random sample of size  $k_b$  do
1413   Add icon-only class to hide text visually;
1414 for each input  $inp$  in <input>, <textarea>, <select> do
1415   if placeholder is None then
1416      $inp.placeholder \leftarrow$  weak hint based on  $is\_zh$ ;
1417   Remove required attribute;
1418 return modified HTML  $H'$ ;

```

C.2 Text Perturbation

Algorithm 3: Button text perturbation and DOM write-back

```

1431 Input: Button text  $t$ 
1432 Output: Perturbed text  $t'$ 
1433 Define perturbation strategies:
1434   (1) Replace characters with look-alikes (e→3, o→0, l→1, ...);
1435   (2) Shuffle characters if length > 3;
1436   (3) Reverse string;
1437   (4) Insert spaces between characters;
1438   (5) Replace with keyword (e.g., "Submit") if trigger word
1439     detected;
1440   (6) Append punctuation (! or ·);
1441 Randomly select strategy  $s$ ;
1442 Apply  $s$  to  $t$ , obtain  $t'$ ;
1443 if  $t' = t$  then
1444   append marker to  $t'$ 
1445 Update button innerText =  $t'$ ;

```

C.3 Layout Perturbation

Algorithm 4: DOM-level frame perturbations (wrapper / role / noise)

```

1451 Input: DOM tree  $T$ 
1452 Output: Perturbed DOM tree  $T'$ 
1453 Wrapper Injection: Wrap selected node with new <div>
1454   container;
1455 Role Replacement: Replace <button> or <input
1456   type=submit> with <div role=button> keeping label
1457   text;
1458 Redundant Nodes: Insert hidden <div> nodes at the end
1459   of DOM;
1460 Difficulty levels:
1461   • Easy: few redundant nodes
1462   • Medium: wrapper injection + 1 role replacement +
1463     moderate redundant nodes
1464   • Hard: deeper wrappers + multiple replacements + many
1465     redundant nodes
1466 Apply operators according to difficulty;
1467 Return  $T'$ ;

```

C.4 Evaluation Metrics

We adopt evaluation metrics tailored to the characteristics of each subtask.

Categorical reasoning tasks. For **position relationship reasoning** (relative spatial relationship) and **UI group classification**, we employ accuracy as the primary metric:

$$\text{Accuracy} = \frac{1}{N} \sum_{i=1}^N \mathbf{1}\{\hat{y}_i = y_i\}, \quad (1)$$

where N is the number of samples, \hat{y}_i is the model prediction, and y_i is the ground truth.

Form understanding tasks. For **form filling** and **hint-text recovery**, outputs are free-form text with potential semantic variation. We thus compute embedding-based similarity:

$$\text{Sim}_{\text{emb}}(x, y) = \frac{\langle f(x), f(y) \rangle}{\|f(x)\| \cdot \|f(y)\|}, \quad (2)$$

where $f(\cdot)$ is a sentence embedding model. This graded metric captures semantic alignment beyond surface lexical match.

Robustness tasks. For **text robustness**, we employ embedding similarity, comparing predictions under original versus perturbed webpages. For **layout robustness**, which involves global page topic consistency, we use TF-IDF cosine similarity:

$$\text{Sim}_{\text{tfidf}}(x, y) = \frac{\sum_j \text{tfidf}_{x,j} \cdot \text{tfidf}_{y,j}}{\|\text{tfidf}_x\| \cdot \|\text{tfidf}_y\|}. \quad (3)$$

For **color robustness**, accuracy is used:

$$\text{accuracy} = \frac{TP}{TP + FN}. \quad (4)$$

1509 *Safety critical detection.* For **safety critical detection**, recall
 1510 is employed as webpages are curated to always contain safety-
 1511 relevant elements, and false negatives represent critical failures.

1512 Importantly, our robustness evaluation can also be regarded as
 1513 an *ablation-style analysis*: by presenting models with identical tasks
 1514 under clean and perturbed webpages, we isolate the influence of
 1515 specific perturbations (color, text, layout).

1518 D Fine-tuning Details

1519 We select Qwen2.5-VL-7B-Instruct as the target model for fine-
 1520 tuning. The fine-tuning focused on three benchmark dimensions
 1521 where performance gaps were most pronounced and supervision
 1522 signals relatively tractable: pairwise position relationship reasoning,
 1523 UI group classification, and color robustness.

1525 D.1 Training Data Construction

1526 **Position Relationship Reasoning.** We decomposed the original
 1527 four-element spatial reasoning questions in WebRSSBench into
 1528 pairwise samples. Each training sample consists of a full-page screenshot,
 1529 two cropped element images, and a ground-truth spatial relation label.
 1530 The ground truths were automatically computed based on element
 1531 bounding box coordinates from the HTML source. This
 1532 decomposition ensures that each sample corresponds to a single,
 1533 unambiguous spatial relation, simplifying the learning objective
 1534 while preserving the core reasoning challenge.

1535 **UI Group Classification.** Each training sample consists of a webpage
 1536 screenshot with a randomly selected element highlighted by a red
 1537 bounding box, paired with a ground-truth UI group label (one of:
 1538 top bar, left sidebar, right sidebar, main content area, bottom
 1539 bar, or auxiliary area) annotated by human experts following
 1540 the guidelines described in Appendix ??.

1541 **Color Robustness.** Each training sample consists of a webpage
 1542 screenshot with 10–30% of buttons recolored using the partial
 1543 chroma shift perturbation (Algorithm 1), paired with the ground-
 1544 truth CTA button label determined by human majority voting. We
 1545 selected this perturbation level as it creates a scenario where color-
 1546 based shortcuts become unreliable for only a subset of elements,
 1547 encouraging the model to rely on structural and semantic cues
 1548 rather than visual saliency alone.

1550 D.2 Data Splits and Leakage Control

1551 **Data Splits.** For all three tasks, datasets were split into training
 1552 and evaluation sets:

- 1554 • **Position Relationship Reasoning:** Training and evalua-
 1555 tion samples were derived from the pairwise decomposition
 1556 of WebRSSBench position reasoning questions.
- 1557 • **UI Group Classification:** Samples were split into training
 1558 (`uigroup_train.jsonl`) and evaluation (`uigroup_eval.jsonl`)
 1559 sets.
- 1560 • **Color Robustness:** Samples were split into training (`train.jsonl`)
 1561 and evaluation (`test.jsonl`) sets.

1562 **Leakage Control.** To prevent data leakage, we performed *site-
 1563 level splitting*: all samples derived from the same source website
 1564 were assigned exclusively to either the training or evaluation set,
 1565 with no overlap at the website or template level. This ensures that

1566 the model cannot exploit memorized site-specific layouts or visual
 1567 patterns during evaluation.

1568 D.3 Training Configuration

1569 **LoRA Configuration.** We adopted LoRA (Low-Rank Adaptation)
 1570 with the following hyperparameters: rank $r = 16$, scaling factor
 1571 $\alpha = 32$, and dropout rate 0.05. LoRA was applied to the following
 1572 projection layers: q_{proj} , k_{proj} , v_{proj} , o_{proj} , $\text{gate}_{\text{proj}}$, up_{proj} ,
 1573 $\text{down}_{\text{proj}}$, w_1 , w_2 , w_3 , and W_{pack} .

1574 **Hyperparameters.** Training hyperparameters for each task are
 1575 summarized in Table 3.

1576 **Table 3: Fine-tuning hyperparameters for each task.**

Hyperparameter	Position	UI Group	Color
Base Model	Qwen2.5-VL-7B	Qwen2.5-VL-7B	Qwen2.5-VL-7B
Batch size	1	4	1
Gradient accumulation	8	4	8
Effective batch size	8	16	8
Epochs	2	3	2
Learning rate	2e-4	2e-5	2e-4
Warmup ratio	0.03	0.03	0.03
Precision	BF16	BF16	BF16
Gradient checkpointing	Yes	No	Yes

1577 **Image Preprocessing.** To ensure training stability and prevent
 1578 out-of-memory errors, we applied the following image preprocess-
 1579 ing:

- 1580 • Images with maximum dimension exceeding 1024 pixels were
 1581 downsampled proportionally.
- 1582 • Extremely small images (< 8 pixels) or images with extreme
 1583 aspect ratios (> 200 : 1) were padded using letterbox filling
 1584 to a minimum size of 32 × 32 pixels.

1585 **Reproducibility.** All experiments were conducted with a fixed
 1586 random seed of 42 (`set_seed(42)`) for data shuffling and model
 1587 initialization. We also enabled `torch.backends.cuda.matmul.allow_tf32`
 1588 and `torch.backends.cudnn.benchmark` for training acceleration.

1589 **Training Schedule.** We trained for a fixed number of epochs (2
 1590 epochs for position reasoning and color robustness, 3 epochs for UI
 1591 group) with epoch-level checkpoint saving. The final checkpoint
 1592 was selected for evaluation. No early stopping was applied.

1593 **Computational Resources.** Fine-tuning was performed on NVIDIA
 1594 GPUs using the `device_map="auto"` setting from Hugging Face
 1595 Transformers, which automatically distributes model parameters
 1596 across available devices. Each fine-tuning run was completed within
 1597 several hours depending on dataset size.

1598 D.4 Transfer Across Perturbations

1599 For color robustness, we fine-tuned exclusively on the 10–30% par-
 1600 tial button chroma shift setting. This design choice was intentional:
 1601 by training on partial perturbations where only a subset of buttons
 1602 are recolored, the model is forced to develop robust CTA identifi-
 1603 cation strategies that do not rely solely on color saliency. We hypo-
 1604 thesize that this approach encourages better generalization com-
 1605 pared to training on 100% perturbations, which may lead the model
 1606 to overfit to uniformly altered color distributions.

1625 D.5 Fine-tuning Scripts

1626 We provide the core training configurations below for reproducibility.

1629 Listing 1: Fine-tuning configuration for Position Relationship Reasoning

```
1632 # Model and LoRA setup
1633 model = Qwen2_5_VLForConditionalGeneration.from_pretrained(
1634     MODEL_ID, torch_dtype=torch.bfloat16, device_map="auto"
1635 )
1636 lora_config = LoraConfig(
1637     r=16, lora_alpha=32, lora_dropout=0.05, bias="none",
1638     task_type="CAUSAL_LM",
1639     target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
1640     "gate_proj", "up_proj", "down_proj", "w1", "w2", "w3
1641     ", "W_pack"]
1642 )
1643 model = get_peft_model(model, lora_config)
1644
1645 # Training arguments
1646 args = TrainingArguments(
1647     per_device_train_batch_size=1,
1648     gradient_accumulation_steps=8,
1649     num_train_epochs=2,
1650     learning_rate=2e-4,
1651     warmup_ratio=0.03,
1652     save_strategy="epoch",
1653     eval_strategy="epoch",
1654     bf16=True,
1655     gradient_checkpointing=True,
```

1656 Listing 2: Fine-tuning configuration for UI Group Classification

```
1657 # LoRA configuration
1658 lora_config = LoraConfig(
1659     r=16, lora_alpha=32, lora_dropout=0.05, bias="none",
1660     task_type="CAUSAL_LM",
1661     target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
1662     "gate_proj", "up_proj", "down_proj"]
1663 )
1664
1665 # Training arguments
1666 training_args = TrainingArguments(
1667     per_device_train_batch_size=4,
1668     gradient_accumulation_steps=4,
1669     num_train_epochs=3,
1670     learning_rate=2e-5,
1671     warmup_ratio=0.03,
1672     evaluation_strategy="steps",
1673     save_strategy="epoch",
1674     bf16=True,
```

1675 Listing 3: Fine-tuning configuration for Color Robustness

```
1676 # LoRA configuration (same as Position task)
1677 lora_config = LoraConfig(
1678     r=16, lora_alpha=32, lora_dropout=0.05, bias="none",
1679     task_type="CAUSAL_LM",
1680     target_modules=["q_proj", "k_proj", "v_proj", "o_proj",
1681     "gate_proj", "up_proj", "down_proj", "w1", "w2", "w3
1682     ", "W_pack"]
1683 )
1684
1685 # Training arguments
1686 args = TrainingArguments(
1687     per_device_train_batch_size=1,
1688     gradient_accumulation_steps=8,
```

```
1689 learning_rate=2e-4,
1690 warmup_ratio=0.03,
1691 save_strategy="epoch",
1692 eval_strategy="epoch",
1693 bf16=True,
1694 gradient_checkpointing=True,
```

1695 D.6 Position Relationship Reasoning

1696 Listing 4: Fine-tuning script for position reasoning

```
1697 processor = AutoProcessor.from_pretrained(MODEL_ID,
1698     trust_remote_code=True)
1699 model = Qwen2_5_VLForConditionalGeneration.
1700     from_pretrained(
1701     MODEL_ID, torch_dtype=torch.bfloat16, device_map
1702     ="auto", trust_remote_code=True
1703 )
1704 model.config.use_cache = False
1705 if hasattr(model, "enable_input_require_grads"):
1706     model.enable_input_require_grads()
1707
1708 model = get_peft_model(model, LoraConfig(
1709     r=16, lora_alpha=32, lora_dropout=0.05, bias="
1710     none",
1711     task_type="CAUSAL_LM",
1712     target_modules=["q_proj", "k_proj", "v_proj", "o_
1713     proj",
1714     "gate_proj", "up_proj", "down_proj
1715     ", "w1", "w2", "w3", "W_pack"]
1716 ))
1717
1718 args = TrainingArguments(
1719     per_device_train_batch_size=1,
1720     gradient_accumulation_steps=8,
1721     num_train_epochs=2,
1722     learning_rate=2e-4,
1723     warmup_ratio=0.03,
1724     logging_steps=20,
1725     bf16=True,
1726     gradient_checkpointing=True,
1727     dataloader_pin_memory=True,
```

1728 D.7 Color Robustness

1729 Listing 5: Fine-tuning script for color robustness

```
1730 model.config.use_cache = False
1731 if hasattr(model, "enable_input_require_grads"):
1732     model.enable_input_require_grads()
1733
1734 lora_config = LoraConfig(
1735     r=16, lora_alpha=32, lora_dropout=0.05, bias="
1736     none",
1737     task_type="CAUSAL_LM",
1738     target_modules=["q_proj", "k_proj", "v_proj", "o_
1739     proj",
1740     "gate_proj", "up_proj", "down_proj
1741     ", "w1", "w2", "w3", "W_pack"]
```

```

1741 )
1742     model = get_peft_model(model, lora_config)
1743
1744     args = TrainingArguments(
1745         per_device_train_batch_size=1,
1746         gradient_accumulation_steps=8,
1747         num_train_epochs=2,
1748         learning_rate=2e-4,
1749         warmup_ratio=0.03,
1750         logging_steps=20,
1751         save_strategy="epoch",
1752         eval_strategy="epoch",
1753         bf16=True,
1754         gradient_checkpointing=True,
1755     )
1756
1757

```

D.8 UI Group

Listing 6: Fine-tuning script for UI group classification

```

train_data = load_dataset("json", data_files=os.path
    .join(DATA_DIR,"uigroup_train.jsonl"))["train"]
eval_data = load_dataset("json", data_files=os.path
    .join(DATA_DIR,"uigroup_eval.jsonl"))["train"]

lora_config = LoraConfig(
    r=16, lora_alpha=32, lora_dropout=0.05, bias="none",
    task_type="CAUSAL_LM",
    target_modules=["q_proj", "k_proj", "v_proj",
        "o_proj",
        "gate_proj", "up_proj", "down_proj
        ", "w1", "w2", "w3", "W_pack"]
)
model = get_peft_model(model, lora_config)

training_args = TrainingArguments(
    per_device_train_batch_size=4,
    per_device_eval_batch_size=4,
    gradient_accumulation_steps=4,
    num_train_epochs=3,
    learning_rate=2e-5,
    warmup_ratio=0.03,
    logging_steps=10,
    evaluation_strategy="steps",
    save_strategy="epoch",
    bf16=True,
    report_to="none",
)

```

E Robustness and fail cases

The following table 4 shows the ratio of unchanged output of the model before and after our interference.

The following cases 5, 6, 7 are also derived from our interference and the pre-and post output of the model.

Table 4: Model performance on robustness tasks (Self compare).

Model	Color Robustness	Text Robustness	Layout Robustness
GPT-5	75.88	65.13	57.65
claude-4-Sonnet	70.54	54.14	62.47
Gemini 2.5-Pro	74.06	70.33	47.87
Pixtral-Large	52.86	48.27	62.69
InternVL3-78B	74.21	66.59	64.52
Qwen 2.5-VL-72B	74.39	69.83	75.18
Qwen2.5-VL-32B	64.69	71.24	58.63
Intern-S1	78.35	74.86	62.49
Pixtral-12B	61.24	72.86	71.81
Llama4-Scout17B	76.69	64.23	61.32
Qwen2.5-VL-7B	71.72	62.77	63.72
Qwen-vl-plus	75.12	69.47	70.78

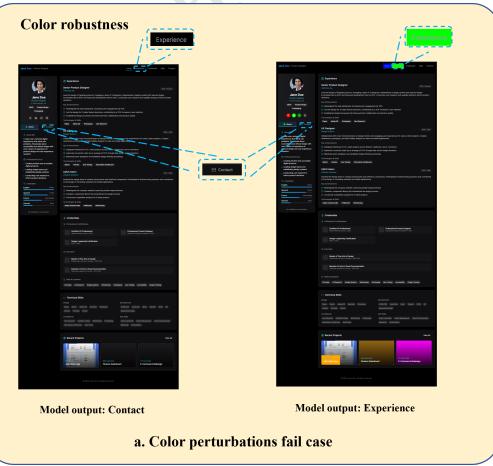


Figure 5: Output examples of MLLMs on WebRSSBench when facing different perturbations.

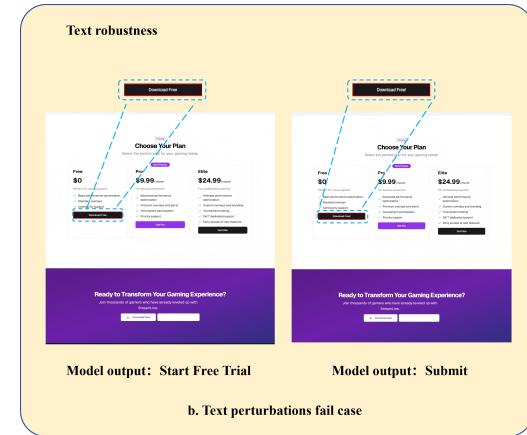


Figure 6: Output examples of MLLMs on WebRSSBench when facing different perturbations.

1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855

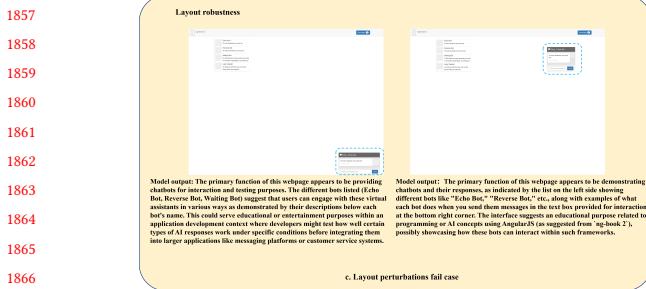


Figure 7: Output examples of MLLMs on WebRSSBench when facing different perturbations.

F LLM Usage

Large Language Models (LLMs) were used to integrate paragraph clarity and sentence grammar correction, while ensuring that sentences have high readability.

Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009