# S WebMMU: A Benchmark for Multimodal Multilingual Website **Understanding and Code Generation**

**Anonymous ACL submission** 

#### Abstract

We present WebMMU, a multilingual benchmark that evaluates three core web tasks: (1) 004 website visual question answering, (2) code editing involving HTML/CSS/JavaScript, and (3) mockup-to-code generation. Unlike prior benchmarks that treat these tasks separately, WebMMU unifies them using expert-annotated, real-world web data to assess models' abilities in complex multi-step reasoning, precise element grounding, and functional UI comprehension and coding. Our evaluation shows 012 that while multimodal large language models 014 (MLLMs) perform well on basic information extraction, they struggle with reasoning and grounding, editing code to preserve functionality, and generating design-to-code that maintains hierarchy and supports multilingual con-019 tent. These findings reveal key limitations in current MLLMs and underscore the need for improved multimodal and cross-lingual reasoning to build future web agents capable of automating diverse web development tasks.

#### Introduction 1

001

011

024

037

041

The web is vital to daily life, enabling information access, shopping, and communication. Multimodal large language models (MLLMs) (Wang et al., 2024b; Hurst et al., 2024) that understand the Visual Web can help users extract information, support tasks like budget-conscious shopping, and handle multiple languages (Deng et al., 2024). They also show promise in automating web design and development, including front-end layout creation, user interface (UI) editing, and code generation (Anthropic, 2024). Unlike tasks focused only on text or images (Wang et al., 2024c; Yue et al., 2024), visual web understanding requires combining UI structure, layouts, text, interactivity, and visuals.

Existing benchmarks target specific aspects of web tasks but remain fragmented and lack comprehensive coverage. Website VQA datasets like WebQA (Chang et al., 2022) and WebSRC (Chen et al., 2021b) mainly focus on text retrieval, overlooking reasoning over UI structure, interactivity, and multilingual content. Recent web agent benchmarks evaluate online task completion (Koh et al., 2024; Deng et al., 2024; He et al., 2024), showing promise for agentic AI but are limited to artificial websites or lack fine-grained categorization (e.g. grounding, understanding, multi-step reasoning). In web development, design-to-code datasets such as Pix2Code (Beltramelli, 2018) and Web2Code benchmarks (Yun et al., 2024), as well as sketchbased datasets like Sketch2Code (Li et al., 2024), cover a limited variety of UIs and often fail to capture real web variability due to automated creation. Furthermore, current benchmarks lack multilingual and cross-domain generalization, limiting applicability beyond English and specific domains. These gaps motivate a unified benchmark integrating multiple web tasks with multimodal, reasoning, and cross-lingual capabilities for effective evaluation of AI in web development and advanced web understanding.

042

043

044

047

048

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

078

079

081

082

To address these challenges, we introduce Web-MMU (Figure 1), a multimodal, Multilingual, and MUlti-task benchmark for evaluating MLLMs on the Visual Web in four languages: English, Spanish, German, and French. WebMMU covers three core tasks: Website VQA (WebQA), which tests functional understanding, visual comprehension, and multi-step reasoning via visual questionanswering; Mockup2Code Generation, assessing design-to-code alignment for UI mockups and sketches, including both simple and complex nested layouts; and Web Code Editing, evaluating precise, context-aware HTML/CSS/JavaScript code editing for feature additions, UI tweaks, and bug fixes. The benchmark spans 20 domains such as shopping, booking, sports, and technology, ensuring wide real-world relevance.



Figure 1: **WebMMU Benchmark Overview.** WebMMU evaluates models on diverse web-based tasks: WebQA, Mockup2Code, and Web Code Editing. Covering 20 domains and four languages, it challenges models to answer visual questions requiring multi-step reasoning and action grounding. It also assesses design-to-code generation from website layouts of varying complexity and evaluates code editing for automated web development.

We benchmark state-of-the-art MLLMs across three core tasks, evaluating both open-source and closed-source models. Our results reveal significant challenges in action grounding and complex reasoning in the WebQA task, along with difficulties in structured layout understanding and accurate code generation for web development. While models (in particular, closed-source ones) exhibit strong general image understanding in WebQA, they struggle with complex reasoning, with most scoring below 50% and some as low as 2% (e.g., Fuyu-8B in English), alongside notable multilingual performance drops (Figure 2). In Web Code Editing, top-performing models like Gemini-2.0-Flash and Claude-3.5-Sonnet outperform open-source counterparts, yet still struggle with maintaining logical structure and syntactic correctness, highlighting the need for more structure-aware codeediting techniques, particularly for complex modifications. Similarly, in Mockup2Code, models such as OpenAI-o1 and Claude-3.5 achieve a high LLMas-Judge score (4/5) on simple layouts but fail with nested element structures, revealing limitations in UI hierarchy comprehension. These findings emphasize the need for improved multimodal align-

101

102 103

104

105

106

107

ment, UI-aware modeling, and cross-lingual robustness to bridge the gap between vision-language models and real-world web interaction. 108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

Our contributions are as follows:

- **Comprehensive Multi-Task Benchmark**: A unified evaluation suite encompassing website VQA, web design-to-code generation, and code editing tasks.
- Diverse, Expert-Annotated Multilingual Data: Fine-grained question-answer pairs, code edits, and UI design annotations across four languages, enabling comprehensive evaluation.
- **Findings**: MLLMs face challenges in multistep reasoning and grounding for WebQA, precise code editing, UI hierarchy understanding in web development, and multilingual generalization.

# 2 Related Work

Web Understanding and Agentic MLLMs. Multimodal learning has become central to web UI understanding, integrating visual, textual, and structural modalities to support both web com-

prehension and agentic navigation. Early work, 131 such as Screen2Words (Wu et al., 2021), parsed 132 web screenshots into UI elements, later influenc-133 ing MLLM pretraining(Lee et al., 2023). Recent 134 advances leverage patching strategies(Baechler 135 et al., 2024), grounding(Cheng et al., 2024), text-136 structural alignment(Xu et al., 2024; Bai et al., 137 2021), and context-aware UI representations(Kil 138 et al., 2024). These innovations have expanded 139 MLLM applications in web agents, enabling mod-140 els to navigate and manipulate websites based on 141 user instructions (Zheng et al.; Yoran et al., 2024; 142 Cheng et al., 2024). However, existing benchmarks 143 often rely on limited artificial websites(Deng et al., 144 2024; Zhou et al., 2023) or focus solely on En-145 glish data(He et al., 2024; Lù et al., 2024; Zhang 146 et al., 2024; Chen et al., 2024a), lacking diversity 147 and real-world complexity. WebMMU addresses 148 these gaps by incorporating real-world websites 149 and multilingual queries, requiring models to per-150 form complex reasoning and UI grounding, making 151 it a more comprehensive evaluation framework for MLLM-driven web understanding and navigation. 153

Visual Question Answering for Web. Progress 154 in web-based VQA has been driven by benchmarks 155 like WebSRC (Chen et al., 2021b), WebQA (Chang 156 157 et al., 2022), WebQuest (Wang et al., 2024a), VisualWebBench (Liu et al., 2024), and WebWalk-158 erQA (Wu et al., 2025) covering tasks such as 159 captioning, webpage QA, and element ground-160 ing. Compared to traditional VQA on natural im-161 ages (Yue et al., 2024), web-based VQA addition-162 ally requires understanding structured webpage lay-163 outs, the relationships between UI elements, and 164 their functional roles within web environments. 165 However, these benchmarks cover limited tasks, domains and languages. WebMMU addresses this gap by covering 20 domains in four languages and 168 adding fine-grained categories like action ground-169 ing, multi-step reasoning, and general understand-170 ing for more thorough evaluation. 171

Automatic Web Design and Development. 172 Code generation and editing have been widely stud-173 ied across programming languages, with bench-174 marks evaluating code generation (Chen et al., 175 2021a; Jimenez et al., 2024; Rodriguez et al., 177 2024b,a) and code editing based on natural language instructions (Guo et al., 2024; Tian et al., 178 2024). However, most previous studies focus on 179 general-purpose programming, neglecting web design and development. To bridge this gap, Gui et al. 181

(2024) and Yun et al. (2024) explore generating HTML/CSS from web screenshots. In contrast, WebMMU introduces Web Code Editing, which involves multilingual tasks for modifying a website's visual and functional features based on user instructions, better reflecting real-world web development use cases. Additionally, WebMMU includes Mockup2Code; unlike prior work (Jain et al., 2019; Barua et al., 2022) that relies on simplistic and artificial sketches drawn by researchers, our sketches are extracted from real-world websites, preserving complex element hierarchies through expert annotation. 182

183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

199

200

201

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

# 3 SWebMMU Benchmark

We introduce WebMMU, designed to evaluate MLLMs on real-world Visual Web tasks. In this section, we describe WebMMU's data collection, annotation process, and present an overview of benchmark tasks.

## 3.1 Data Collection and Annotation

Website Selection and Data Capture. To construct WebMMU, we curated a diverse set of webpage URLs from the FineWeb dataset (Penedo et al., 2024) and applied domain-specific heuristics to ensure coverage across 20 popular, contentrich, and feature-rich web domains (e.g., shopping, booking, technology). We selected webpages in four languages: English, German, French, and Spanish - considering linguistic diversity, annotator availability, and budget constraints. To capture full browsing sessions on a single webpage, we generated collages combining multiple snapshots taken at different scroll depths and interaction states within the page. A viewport-specific snapshot was retained alongside relevant HTML and assets (e.g., CSS, JavaScript). Selection strictly adhered to web crawling policies (e.g., robots.txt).

Annotation Process. Annotators were provided with webpage screenshots, corresponding HTML, and asset files and were tasked with three objectives: (1) generating open-ended and multiplechoice questions that capture real-world usage, including highlighting, clicking, and multi-step reasoning; (2) creating UI mockups of varying complexity and formats to support design-to-code workflows; and (3) formulating code edit requests that require programming expertise. A structured training phase ensured annotation consistency and quality. Further details on annotator guidelines are

281

231

233

234

241

243

244

245

246

248

249

250

252

254

255

263

264

265

267

269

271

272

273

276

277

278

given in the Appendix A.

**Quality Control and Annotator Demographics.** A 100% quality assurance framework was implemented in three stages: Trainer Review, where experienced annotators performed initial annotations; Primary QA (QA1), where independent specialists verified accuracy, completeness, and adherence to guidelines; and Secondary QA (QA2), ensuring consistency with expert-level annotation criteria. The dataset was annotated by 127 professionals across North America, South America, Europe, Africa, and Asia, representing diverse linguistic and domain expertise. English annotators primarily came from Asia, German and French from Europe, and Spanish from Latin America. Annotators held qualifications ranging from bachelor's to advanced degrees for specialized tasks and were compensated above fair market wages, ensuring ethical labor practices and high-quality results.

#### Tasks Overview 3.2

#### 3.2.1 Web Question Answering (WebQA)

The WebQA task in WebMMU evaluates models' ability to extract, integrate, and ground structured UI elements, numerical data, and graphical components from web screenshots while reasoning over hierarchical layouts, predicting actions, and ensuring spatial grounding. It consists of three categories: Agenctic Action, which focuses on web navigation and action execution without feedback from the environment, requiring models to understand UI elements like buttons, menus, and hyperlinks, identify elements (e.g., "Where can I find the coaching plans?"), and execute actions (e.g., "How can I save this drill?") while handling spatial grounding and distinguishing static vs. interactive elements across multilingual UIs; many of these tasks also require *coordinate-based reasoning* to localize UI components accurately. Multi-step Reasoning involves multi-step inference, numerical calculations, and comparisons across UI components (e.g., "If a customer were to buy all the camera models mentioned on the bottom of this page in Expert Camera Reviewstable, what would be the grand total?"), requiring models to integrate text, numerical values, and layout structures from structured web content, where hierarchical reasoning is essential despite being constrained to single-frame snapshots; and General Visual Comprehension, which assesses a model's ability to extract and synthesize structured and unstructured data from web

screenshots, including OCR-extracted text, images, graphical elements, and UI components (e.g., "How many brand logos are in the Featured Brands section?"), emphasizing semantic comprehension beyond standard OCR-based extraction.

# 3.2.2 Mockup2Code

The Mockup2Code task in WebMMU advances design-to-code by translating hand-drawn wireframes and high-fidelity digital mockups into structured code. Unlike text-based UI generation, it evaluates a model's ability to interpret spatial hierarchies and UI structures from visual inputs. The dataset includes low-fidelity sketches and digitally created mockups, challenging models to generalize across abstraction levels in web design while tackling component recognition, spatial alignment, and structured code synthesis. Unlike prior design-tocode datasets, WebMMU incorporates real-world web layouts, ensuring models generate syntactically correct and semantically meaningful code aligned with modern web development practices.

## 3.2.3 Web Code Editing

Web Code Editing is a novel task, which evaluates a model's ability to modify webpage code while preserving functional and structural integrity, given a screenshot, source code, and a user edit request. To perform well, models must complete three sub-tasks: (1) understand the provided inputs, including the webpage codebase, visual elements in the screenshot, and the requested modification; (2) identify the relevant code snippets that require modification; and (3) generate the appropriate HTML, CSS, or JavaScript edits to implement the requested change. These sub-tasks require an 314 advanced understanding of webpage development 315 and realistic code editing capabilities. The modi-316 fication requests span a broad range of visual and 317 functional changes. Visual edits include adjusting 318 font size and colors, repositioning elements, and 319 adding headers or footers. Functional modifica-320 tions involve adding interactive components such 321 as buttons or forms and enhancing user experience 322 with dynamic UI elements. The task is multilin-323 gual, aligning with the broader scope of WebMMU. 324 Given the length of webpage source code, models 325 are prompted to output only the necessary code dif-326 ferences rather than rewriting the entire codebase. 327 This improves both practicality and efficiency, ensuring that the generated edits remain concise and 329 targeted. More details on the prompt formulation 330

	En	Es	De	Fr	Total
Website Images	392	133	130	131	786
WebQA	1476	484	379	456	2795
Mockup2Code	180	93	85	78	436
Web Code Editing	165	75	67	68	375
Total	2213	785	661	733	4392

Table 1: **Dataset Statistics.** Language-wise dataset breakdown across tasks. We report the number of web images per language. English (En), Spanish (Es), German (De) and French (Fr).

Task	Metric	<b>Evaluation Details</b>
WebQA	LLM-as-	Measures accuracy; 0 (incorrect)
	Judge	/ 1 (correct).
Mockup2Code	LLM-as-	Assesses layout fidelity on a 1-5
	Judge	scale (layout, spacing, grid).
Cada Editina	BLEU, Tree-	Evaluates structural correctness
Code Editing	BLEU	by matching ground truth differ-
		ences.
	LLM-as-	Scores functional accuracy on
	Judge	a 1-5 scale (functional correct-
		ness).

Table 2:	Evaluation	Metrics	used in	WebMMU.

are provided in Appendix C.2.

331

332

333

334

338

341

342

343

351

354

357

### **3.3 Dataset Statistics**

WebMMU covers four languages: English, Spanish, German, and French (see Table 1). It contains 786 webpage images from domains like ecommerce, education, news, and finance. It includes 2795 WebQA samples, 436 Mockup2Code instances, and 375 Web Code Editing cases. Unlike previous datasets that focus on predefined UI layouts, WebMMU uses full-page snapshots, including dynamic content, nested structures, and multimodal dependencies. A small portion of images consist of multiple panels combined into a single image, capturing dense information and replicating browsing sessions.

#### 4 Evaluation

We evaluate state-of-the-art MLLMs across both closed-source and open-source categories. Model inference for WebQA, Mockup2Code, and Web Code Editing follows standardized prompts (Appendix D). Evaluation combines LLM-as-Judge (Zheng et al., 2023) scoring with established automatic metrics, as summarized in Table 2.

LLM-as-Judge is used to evaluate WebQA, where model responses receive binary correctness scores (0 or 1) based on predefined criteria for semantic accuracy and reasoning completeness (Appendix D.3). This structured approach ensures consistency and prevents arbitrary grading. Inspired by automated evaluation in image synthesis (Ku et al., 2023), Mockup2Code uses LLM-as-Judge, assessing the alignment between input sketches and rendered outputs across three key dimensions: layout structure, spacing, and grid consistency (Appendix D.2). Each aspect follows well-defined scoring guidelines, ensuring reproducible and fair assessments. For Web Code Editing, we evaluate both structural correctness and functional accuracy. The former is measured using BLEU (Papineni et al., 2002) and TreeBLEU (Gui et al., 2024), ensuring syntactic validity and adherence to coding conventions. The latter relies on LLM-as-Judge, where functional equivalence between reference and predicted edits is rated on a 1-5 scale. To avoid arbitrary scoring, rating criteria explicitly define correctness levels based on functional preservation and intended user modifications. Since web functionalities can be implemented in multiple ways, the evaluation accounts for semantically valid alternatives, preventing undue penalization of syntactically different but functionally correct edits. For all LLM-as-Judge evaluations, we use GPT40-1120, which has demonstrated strong alignment with human judgment and diverse scoring behavior (Feizi et al., 2025), ensuring robustness across tasks.

358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

383

384

385

387

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

## **5** Results

#### 5.1 WebQA Performance

Table 3 presents model accuracy for three question types. Closed-source models, such as Gemini 2.0 Flash and Claude 3.5 Sonnet, outperform opensource alternatives across all tasks but still struggle with agentic action, particularly in predicting spatial coordinates for interactive elements. Among open-source models, larger architectures (>30B parameters) like Qwen2VL-72B and Internvl2.5-38B perform better in general image understanding and UI recognition, while smaller models (<8B) exhibit poor generalization across tasks.

Performance varies by question type. General image understanding is easiest, relying mainly on visual recognition. Complex reasoning is harder, with most models scoring below 50% and some as low as 2% (e.g., Fuyu-8b in English), showing difficulties in retrieving and reasoning over structured webpage content. **Agentic action is the hardest, with top models rarely surpassing 10% accuracy**, as it requires precise spatial ground-

Model		English	ı		French			German			Spanish	
hour	Ť	ă,	Q	Ť	Ť	Q	Ť	ă,	Q	P	Ť	Q
Gemini2.0 Flash	44.3	1.2	59.2	41.6	9.0	52.8	18.2	12.8	29.1	46.1	12.0	36.1
Claude3.5 Sonnet	51.4	3.7	64.1	53.0	12.7	51.2	26.9	15.6	31.6	63.8	15.9	41.9
Phi3.5-VI-4b (Abdin et al., 2024)	8.90	1.80	31.60	2.20	6.90	39.00	8.40	13.00	23.90	3.00	10.20	32.00
UI-Tars-7b	19.30	8.10	<u>47.60</u>	7.70	8.90	<u>47.60</u>	7.80	14.30	<u>28.40</u>	20.90	14.00	38.80
Molmo-7b (Deitke et al., 2024)	12.30	<u>3.80</u>	32.90	7.00	7.50	47.60	8.30	13.70	31.90	15.10	10.30	32.00
Qwen2VL-7B (Wang et al., 2024b)	<u>18.00</u>	2.90	57.10	<u>10.10</u>	10.20	52.00	10.70	17.60	26.30	<u>19.30</u>	<u>14.00</u>	<u>36.50</u>
Fuyu-8b (Bavishi et al., 2023)	1.60	0.40	14.30	0.00	1.30	17.50	1.00	5.60	15.70	0.70	1.50	10.90
Internvl2.5-8b (Chen et al., 2024b)	16.30	1.90	46.30	11.00	13.30	40.00	7.40	16.00	25.90	13.80	11.90	31.10
Glm4V-9b (GLM et al., 2024)	15.30	8.10	41.80	11.40	13.90	48.10	14.70	13.80	25.00	21.60	<u>13.40</u>	35.60
Llama-3.2-11B-Vision (Dubey et al., 2024)	27.10	7.90	53.20	<u>11.60</u>	<u>11.30</u>	48.10	11.80	<u>14.30</u>	33.60	17.50	11.80	<u>37.90</u>
Pixtral-12b (Agrawal et al., 2024)	<u>27.10</u>	9.20	<u>44.90</u>	17.70	11.30	53.40	19.50	19.30	21.70	28.70	17.80	40.20
Internvl2.5-38b (Chen et al., 2024b)	<u>22.90</u>	<u>3.80</u>	59.30	20.90	15.30	65.70	18.00	20.10	39.70	36.20	14.90	41.40
Qwen2VL-72B (Wang et al., 2024b)	23.60	4.30	<u>53.70</u>	<u>16.90</u>	<u>13.90</u>	<u>54.50</u>	<u>15.30</u>	<u>17.50</u>	<u>36.20</u>	<u>29.10</u>	<u>12.70</u>	<u>41.00</u>

Table 3: Web VQA Performance. Model accuracy (%) by question type and language. 🛱 : Multi-step Reasoning, K: Agenctic Action, (2): General Visual Comprehension. Best and runner-up models per size category are **bold** and <u>underlined</u>. Model sizes: blue (<8B params), orange (8–12B), green (>12B), gray proprietary.



Figure 2: Failure Cases in WebQA for the top-performing open-source model (InternVL-38B). Most prominent errors occur in grounding actions (e.g., identifying clickable link coordinates) and multi-step reasoning tasks, such as detailed step-by-step calculations.

ing, such as recognizing interactive elements (e.g., 408 "About Me" in a menu) and predicting approxi-409 mate bounding boxes. While many models detect interactive parts, they struggle with localization, resulting in low scores. Error Analysis. Figure 2 reveals common failures: models often miscalculate numbers or fail in multi-step reasoning. In agentic action, inaccurate bounding boxes hurt performance. Multilingual generalization also remains weak despite resource-rich languages. These issues highlight the need for better spatial reasoning, numerical understanding, and cross-lingual adaptation to close the gap between vision-language models and real web interaction.

#### Mockup2Code Generation 5.2

Figure 3 evaluates the Mockup2Code task, reporting scores for each dimension and overall performance. Open MLLMs such as Phi3.5-VI, Fuyu-8B, and GLM4V-9B generally perform poorly across all metrics. Notably, Phi3.5-VI and Fuyu-8B score nearly 1 across all dimensions, indicating a complete failure on this task. Nevertheless, performance improves with model scale. For instance, Qwen2VL's score rises from 1.90 to 3.39 when scaling from 7B to 72B, while InternVL2.5 improves from 2.34 to 3.61 when scaling from 8B to 38B. Additionally, Pixtral-12B outperforms all 7B/8B models. Still, even the best open MLLMs struggle, especially with complex designs - InternVL2.5-38B, the highest performer, scores only 2.98 out of 422

423

424

425

426

427

428

429

430

431

432

433

434

435

436



Figure 3: Mockup2Code Performance. LLM-as-Judge evaluation scores for simple and complex UI mockups across three key dimensions: alignment, layout, and spacing, along with overall performance. Higher scores indicate better fidelity between the generated and reference web designs. Closed-source models outperform open-source alternatives, particularly in complex cases, yet challenges remain in high-fidelity code generation.

5. In contrast, proprietary models like Claude-3.5, Gemini-2.0-Flash, and OpenAI-o1 perform significantly better, particularly on simple UI designs, where they achieve LLM-as-Judge scores above 4. However, their performance declines in complex variants, with top scores reaching only 3.4 out of 5. Across all evaluation dimensions, both proprietary and large-scale open MLLMs struggle most with spacing, which requires accurately setting element dimensions and margins based on sketch input.

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

Case Analysis. Figure 5 shows both success and failure cases of the top model OpenAI-o1 on Mockup2Code. OpenAI-o1 handles simple, flat layouts well, even with moderate element variety and count. However, it struggles with nested structures, often misaligning and failing to preserve element hierarchy and spacing, especially when and <a> tags are nested within <div>. Similar issues occur in other samples and models, as seen in Figures 14 and 16.

#### 5.3 **Code Editing Performance**

Figure 4 shows Web Code Editing results evaluated by LLM-as-Judge (metrics in Table 7). Proprietary models achieve the highest functional accuracy, but only marginally outperform large opensource models, indicating both struggle to preserve functional correctness alongside syntactic consistency. Smaller models like Phi3.5-VI and

Fuyu-8b perform worst, often failing to generate valid code (score <1.5). Performance improves with size; Qwen2VL-72B and InternVL2.5-38B rival closed-source models. Yet, even the strongest exhibit clear limitations producing structurally correct edits that fully preserve functionality. Multilingual performance is stable for top models but more variable for smaller ones, reflecting challenges in adapting edits across languages. Crucially, all models – especially open-source – fail to **automatically** generate valid patch files for seamless source integration. Despite access to full source files, none produced patch content directly usable without manual fixes, making human oversight essential and highlighting a core challenge in automating web code edits.

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

485

486

487

490

491

492

494

495

#### 5.4 **Metric-Human Alignment**

We sampled 100 examples per task and enlisted 484 PhD students and researchers as annotators to evaluate alignment between human judgments and the automatic metric (LLM judge). For WebQA, humans agreed with the LLM judge in 89% of 488 cases. Most disagreements involved functional 489 questions, where the judge required exact bounding boxes, but humans were more lenient - accepting answers that correctly identified the clickable link location without a precise bounding box (e.g., the 493 "about us" link in the navbar). We consider the judge's stricter criteria correct since models were



Figure 4: **Performance on Code Edits.** LLM-as-Judge metric, on a scale of 1-5, used to evaluate functional correctness of code edits. All models, including closed-source models, struggle with the Web Code Editing task of WebMMU. Refer to Table 7 for full results, including BLEU and TreeBLEU scores, of all models.



Figure 5: Success (top) and failure (bottom) cases for Mockup2Code Generation from OpenAI-o1.

prompted to provide exact bounding boxes for such 496 questions (see Appendix C.1). For Mockup2Code 497 Spearman correlations were 0.39 (layout), 0.33 498 499 (spacing), and 0.46 (alignment), averaging 0.43 overall. Pearson correlations were slightly higher: 0.42 (layout), 0.41 (spacing), 0.48 (alignment), 501 with an overall average of 0.50. These correlations, while moderate, reflect the task's subjectivity and 503 support the reliability of the automatic evaluation. 504 For Web Code Editing, expert annotators validated 505 the LLM judge's assessments with 91% accuracy, demonstrating both the reliability of the evaluation and the validity of the associated judge's rationales assigned.

### 6 Conclusion

510

511 WebMMU evaluates MLLMs on a real-world, chal-512 lenging web question answering task and two code generation tasks: front-end design and code editing. Our tasks cover four languages and a wide variety of domains, sourced from human annotators. Our results show that Web VQA models struggle with interpreting complex UIs, reasoning, and multilingual generalization. Code editing models often generate syntactically correct but logically inconsistent code. UI generation models face a trade-off between precise element placement and preserving the original design's meaning. These challenges underscore the need for enhanced multimodal alignment, UI-aware architectures, and robust cross-lingual adaptation to develop future web agents capable of effectively performing a wide range of human tasks on the web.

513

514

515

516

517

518

519

520

521

522

523

524

525

526

## Limitations

528

553

555

557

563

564

567

568

569

570

571

572

573

576

While WebMMU provides a comprehensive evaluation of web-based AI reasoning and code genera-530 tion, it has several limitations. First, it is restricted to single-screenshot web reasoning, capturing static snapshots rather than supporting interactive envi-533 ronments or multi-turn navigation. Although multistep reasoning tasks are included, they rely solely 535 on single-image (including multiple panels of a 536 browsing session), limiting evaluation in dynamic web exploration. Second, linguistic coverage is 538 constrained to four languages: English, French, German, and Spanish; due to annotator availabil-540 ity, which may limit generalization to underrep-541 resented languages and regional web structures. 542 Third, while Mockup2Code and Web Code Editing cover core web technologies such as HTML, CSS, and JavaScript, modern frontend frameworks 545 like React, Angular, and Vue.js are not explicitly 546 evaluated. Finally, the automatic LLM judge met-547 ric, though reliable and fast, does not fully replicate human evaluation. Future work could explore improved automatic metrics or hybrid evaluation 550 approaches to better capture nuanced human judg-551 ments. 552

## Ethical Considerations

WebMMU is a benchmarking resource designed strictly for research purposes in multimodal and multilingual web understanding and generation. All tasks are created by human annotators using everyday web content and undergo thorough validation, so we do not anticipate misuse or harmful content. Compared to prior work, WebMMU expands evaluation across multiple languages, though coverage remains limited by annotator availability. To the best of our knowledge, the dataset contains no NSFW or harmful content. We commit to promptly removing any data upon valid requests once publicly released.

## References

Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, Alon Benhaim, Misha Bilenko, Johan Bjorck, Sébastien Bubeck, Martin Cai, Qin Cai, Vishrav Chaudhary, Dong Chen, Dongdong Chen, and 110 others. 2024. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*. Pravesh Agrawal, Szymon Antoniak, Emma Bou Hanna, Baptiste Bout, Devendra Chaplot, Jessica Chudnovsky, Diogo Costa, Baudouin De Monicault, Saurabh Garg, Theophile Gervet, Soham Ghosh, Amélie Héliou, Paul Jacob, Albert Q. Jiang, Kartik Khandelwal, Timothée Lacroix, Guillaume Lample, Diego Las Casas, Thibaut Lavril, and 23 others. 2024. Pixtral 12b. *Preprint*, arXiv:2410.07073. 577

578

579

580

581

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

- AI Anthropic. 2024. Claude 3.5 sonnet model card addendum. *Claude-3.5 Model Card*, 3(6).
- Gilles Baechler, Srinivas Sunkara, Maria Wang, Fedir Zubach, Hassan Mansoor, Vincent Etter, Victor Cărbune, Jason Lin, Jindong Chen, and Abhanshu Sharma. 2024. Screenai: A vision-language model for ui and infographics understanding. *arXiv preprint arXiv:2402.04615*.
- Chongyang Bai, Xiaoxue Zang, Ying Xu, Srinivas Sunkara, Abhinav Rastogi, Jindong Chen, and 1 others. 2021. Uibert: Learning generic multimodal representations for ui understanding. *arXiv preprint arXiv:2107.13731*.
- Somoy Subandhu Barua, Imam Mohammad Zulkarnain, Abhishek Roy, Md Golam Rabiul Alam, and Md Zia Uddin. 2022. Sketch2fullstack: Generating skeleton code of full stack website and application from sketch using deep learning and computer vision. *arXiv preprint arXiv:2211.14607.*
- Rohan Bavishi, Erich Elsen, Curtis Hawthorne, Maxwell Nye, Augustus Odena, Arushi Somani, and Sağnak Taşırlar. 2023. Introducing our multimodal models.
- Tony Beltramelli. 2018. pix2code: Generating code from a graphical user interface screenshot. In *Proceedings of the ACM SIGCHI symposium on engineering interactive computing systems*, pages 1–6.
- Yingshan Chang, Mridu Narang, Hisami Suzuki, Guihong Cao, Jianfeng Gao, and Yonatan Bisk. 2022. Webqa: Multihop and multimodal qa. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 16495–16504.
- Dongping Chen, Yue Huang, Siyuan Wu, Jingyu Tang, Liuyi Chen, Yilin Bai, Zhigang He, Chenlong Wang, Huichi Zhou, Yiqiang Li, and 1 others. 2024a. Guiworld: A dataset for gui-oriented multimodal llmbased agents. *arXiv preprint arXiv:2406.10819*.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, and 1 others. 2021a. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Xingyu Chen, Zihan Zhao, Lu Chen, Jiabao Ji, Danyang Zhang, Ao Luo, Yuxuan Xiong, and Kai Yu. 2021b. Websrc: A dataset for web-based structural reading comprehension. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 4173–4185.

- 634 635 636
- 63

64

- 64
- 64 64
- 6
- 6

6

- 654
- 6 6
- 6
- 6 6

6

66 66

66

667

67

673 674

675 676 677

- 678 679
- 6

6

6

6

68

588 589

- Zhe Chen, Jiannan Wu, Wenhai Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, and 1 others. 2024b. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24185–24198.
- Kanzhi Cheng, Qiushi Sun, Yougang Chu, Fangzhi Xu, Yantao Li, Jianbing Zhang, and Zhiyong Wu. 2024. Seeclick: Harnessing gui grounding for advanced visual gui agents. *arXiv preprint arXiv:2401.10935*.
- Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, YenSung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, and 32 others. 2024. Molmo and pixmo: Open weights and open data for state-of-the-art multimodal models. *arXiv preprint arXiv:2409.17146*.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2024.
  Mind2web: Towards a generalist agent for the web. Advances in Neural Information Processing Systems, 36.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, and et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Aarash Feizi, Sai Rajeswar, Adriana Romero-Soriano, Reihaneh Rabbany, Spandana Gella, Valentina Zantedeschi, and João Monteiro. 2025. Pairbench: A systematic framework for selecting reliable judge vlms. *arXiv preprint arXiv:2502.15210*.
- Team GLM, :, Aohan Zeng, Bin Xu, Bowen Wang, Chenhui Zhang, Da Yin, Dan Zhang, Diego Rojas, Guanyu Feng, Hanlin Zhao, Hanyu Lai, Hao Yu, Hongning Wang, Jiadai Sun, Jiajie Zhang, Jiale Cheng, Jiayi Gui, Jie Tang, and 40 others. 2024. Chatglm: A family of large language models from glm-130b to glm-4 all tools. *arXiv preprint arXiv:2406.12793*.
- Yi Gui, Zhen Li, Yao Wan, Yemin Shi, Hongyu Zhang, Yi Su, Shaoling Dong, Xing Zhou, and Wenbin Jiang. 2024. Vision2ui: A real-world dataset with layout for code generation from ui designs. *arXiv preprint arXiv:2404.06369*.
- Jiawei Guo, Ziming Li, Xueling Liu, Kaijing Ma, Tianyu Zheng, Zhouliang Yu, Ding Pan, Yizhi Li, Ruibo Liu, Yue Wang, and 1 others. 2024. Codeeditorbench: Evaluating code editing capability of large language models. *arXiv preprint arXiv:2404.03543*.
- Hongliang He, Wenlin Yao, Kaixin Ma, Wenhao Yu, Yong Dai, Hongming Zhang, Zhenzhong Lan, and Dong Yu. 2024. Webvoyager: Building an end-toend web agent with large multimodal models. *arXiv preprint arXiv:2401.13919*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*. 690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

717

718

719

720

721

722

723

724

725

726

727

728

729

730

732

735

736

737

738

739

740

741

- Vanita Jain, Piyush Agrawal, Subham Banga, Rishabh Kapoor, and Shashwat Gulyani. 2019. Sketch2code: transformation of sketches to ui in real-time using deep neural network. *arXiv preprint arXiv:1910.08930*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik R Narasimhan. 2024. Swe-bench: Can language models resolve real-world github issues? In *ICLR*.
- Jihyung Kil, Chan Hee Song, Boyuan Zheng, Xiang Deng, Yu Su, and Wei-Lun Chao. 2024. Dual-view visual contextualization for web navigation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14445–14454.
- Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram Duvvur, Ming Chong Lim, Po-Yu Huang, Graham Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and Daniel Fried. 2024. Visualwebarena: Evaluating multimodal agents on realistic visual web tasks. *arXiv preprint arXiv:2401.13649*.
- Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhu Chen. 2023. Viescore: Towards explainable metrics for conditional image synthesis evaluation. *arXiv preprint arXiv:2312.14867.*
- Kenton Lee, Mandar Joshi, Iulia Raluca Turc, Hexiang Hu, Fangyu Liu, Julian Martin Eisenschlos, Urvashi Khandelwal, Peter Shaw, Ming-Wei Chang, and Kristina Toutanova. 2023. Pix2struct: Screenshot parsing as pretraining for visual language understanding. In *International Conference on Machine Learning*, pages 18893–18912. PMLR.
- Ryan Li, Yanzhe Zhang, and Diyi Yang. 2024. Sketch2code: Evaluating vision-language models for interactive web design prototyping. *arXiv preprint arXiv:2410.16232*.
- Junpeng Liu, Yifan Song, Bill Yuchen Lin, Wai Lam, Graham Neubig, Yuanzhi Li, and Xiang Yue. 2024. Visualwebbench: How far have multimodal llms evolved in web page understanding and grounding? *arXiv preprint arXiv:2404.05955*.
- Xing Han Lù, Zdeněk Kasner, and Siva Reddy. 2024. Weblinx: Real-world website navigation with multiturn dialogue. *arXiv preprint arXiv:2402.05930*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th annual meeting of the Association for Computational Linguistics, pages 311–318.

801

Guilherme Penedo, Hynek Kydlíček, Loubna Ben allal, Anton Lozhkov, Margaret Mitchell, Colin Raffel, Leandro Von Werra, and Thomas Wolf. 2024. The fineweb datasets: Decanting the web for the finest text data at scale. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track.* 

743

744

745

746

747

748

750

751

752

753

754

756

762 763

764

765

766

767

771

772

773

774

775

776

781

790

791

793

794

795

796

797

- Juan Rodriguez, Xiangru Jian, Siba Smarak Panigrahi, Tianyu Zhang, Aarash Feizi, Abhay Puri, Akshay Kalkunte, François Savard, Ahmed Masry, Shravan Nayak, Rabiul Awal, Mahsa Massoud, Amirhossein Abaskohi, Zichao Li, Suyuchen Wang, Pierre-André Noël, Mats Leon Richter, Saverio Vadacchino, Shubbam Agarwal, and 24 others. 2024a. Bigdocs: An open and permissively-licensed dataset for training multimodal models on document and code tasks. *arXiv preprint arXiv:2412.04626*.
- Juan A. Rodriguez, Abhay Puri, Shubham Agarwal, Issam H. Laradji, Pau Rodriguez, Sai Rajeswar, David Vazquez, Christopher Pal, and Marco Pedersoli. 2024b. Starvector: Generating scalable vector graphics code from images and text. *Preprint*, arXiv:2312.11556.
- Runchu Tian, Yining Ye, Yujia Qin, Xin Cong, Yankai Lin, Yinxu Pan, Yesai Wu, Haotian Hui, Weichuan Liu, Zhiyuan Liu, and 1 others. 2024. Debugbench: Evaluating debugging capability of large language models. *arXiv preprint arXiv:2401.04621*.
- Maria Wang, Srinivas Sunkara, Gilles Baechler, Jason Lin, Yun Zhu, Fedir Zubach, Lei Shu, and Jindong Chen. 2024a. Webquest: A benchmark for multimodal qa on web page sequences. *arXiv preprint arXiv:2409.13711*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024b. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. arXiv preprint arXiv:2409.12191.
- Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyan Jiang, and 1 others.
  2024c. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. In The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track.
- Jason Wu, Xiaoyi Zhang, Jeff Nichols, and Jeffrey P Bigham. 2021. Screen parsing: Towards reverse engineering of ui models from screenshots. In *The 34th Annual ACM Symposium on User Interface Software and Technology*, pages 470–483.
- Jialong Wu, Wenbiao Yin, Yong Jiang, Zhenglin Wang, Zekun Xi, Runnan Fang, Deyu Zhou, Pengjun Xie, and Fei Huang. 2025. Webwalker: Benchmarking llms in web traversal. *arXiv preprint arXiv:2501.07572*.

- Linhui Xu, Yawen Zhang, Yawen Li, Yawen Zhang, and Yawen Li. 2024. Hierarchical multimodal finegrained modulation for visual grounding. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 1–9. ACM.
- Ori Yoran, Samuel Amouyal, Chaitanya Malaviya, Ben Bogin, Ofir Press, and Jonathan Berant. 2024. Assistantbench: Can web agents solve realistic and timeconsuming tasks? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 8938–8968.
- Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024.
  Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9556– 9567.
- Sukmin Yun, Haokun Lin, Rusiru Thushara, Mohammad Qazim Bhat, Yongxin Wang, Zutao Jiang, Mingkai Deng, Jinhong Wang, Tianhua Tao, Junbo Li, and 1 others. 2024. Web2code: A large-scale webpage-to-code dataset and evaluation framework for multimodal llms. *arXiv preprint arXiv:2406.20098*.
- Ziniu Zhang, Shulin Tian, Liangyu Chen, and Ziwei Liu. 2024. Mmina: Benchmarking multihop multimodal internet agents. *arXiv preprint arXiv:2404.09992*.
- Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. Gpt-4v (ision) is a generalist web agent, if grounded. In *Forty-first International Conference on Machine Learning*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue Ou, Yonatan Bisk, Daniel Fried, and 1 others. 2023. Webarena: A realistic web environment for building autonomous agents. *arXiv preprint arXiv:2307.13854*.

## A Human Annotator Instruction

## A.1 WebQA Annotations Guideline

You will be provided with screenshots of websites.848Your task is to create challenging questions that test849deep understanding and reasoning about the image850content. Each question should fall into one of the851three categories described below, and be designed852to encourage a detailed analysis of the screenshot.853Important Note: If a screenshot lacks sufficient854

858Agenctic ActionPurpose: Focus on the interac-859tive elements and navigation aspects of the website.860These questions should prompt the viewer to inter-861pret or locate specific functional elements, like but-862tons, menus, or links, and understand their purpose.863Example: "Where would a user click to access864their saved items?" Guidelines: Create questions865that require the viewer to understand how different866elements work or what actions they might trigger.867Avoid overly simple questions that don't involve868interaction or navigation. Do provide the bounding869box location or hint on how to navigate.

Multi-step Reasoning Purpose: These questions 870 should require multi-step thinking, involving the analysis of multiple parts of the image, compar-872 isons, or drawing inferences from the content. Example: "How does the timing of updates in dif-874 ferent news sources on this page provide insights into the event's coverage?" Guidelines: Formulate 876 questions that connect elements across the image or require interpretation of trends, relationships, or content hierarchy. These should not be answerable from a single part of the image. If answerable, then should be difficult e.g. solving a math question (see 881 example) or asking what will happen if the cart is doubled (see example).

General Visual Comprehension Purpose: Assess the viewer's ability to identify and comprehend basic information displayed in the image, such as titles, labels, or the overall structure. Example: "What is the main title or header of this page?" Guidelines: Keep these questions straightforward, focusing on textual or visual elements that convey the primary purpose or information displayed. Aim for questions that require attention to specific details rather than general impressions. Highlight the region of answer with bounding box if needed (upto your choice).

893

894

896

## A.2 Performing Code Editing on Websites

897 Understanding the Scope of Edits Before start898 ing, identify the specific task or issue with clarity
899 and precision. Ensure you fully understand the
900 requested visual or functional changes before pro901 ceeding. Examples of tasks by difficulty are out902 lined below

## **Basic Changes**

• Change the button color from blue to green.	904
• Fix a typo in the homepage headline.	905
• Remove the underlined style from all hyper- links.	906 907
• Add a border to images in the gallery section.	908
Intermediate Enhancements	909
• Replace the navigation bar font with 'Roboto' and ensure it matches the design mockup.	910 911
• Add a hover effect to all buttons, changing their background to light gray.	912 913
• Update the footer links to open in a new tab and add appropriate ARIA labels for accessi- bility.	914 915 916
• Create a consistent color scheme for all head- ings on the page.	917 918
Advanced Functional or Design Tasks	919
• Add a new section to the homepage to show- case recent blog posts, styled to match the website theme.	920 921 922
• Refactor the JavaScript for the carousel to improve performance and fix the sliding bug.	923 924
• Optimize the CSS for faster page load times by combining redundant rules and removing unused classes.	925 926 927
• Implement a lightbox feature for viewing images in the gallery.	928 929
• Create a visually engaging header with a full- width background image and overlay text for the homepage.	930 931 932
• Design a custom 404 error page with an ani- mated illustration and a link back to the home- page.	933 934 935
• Develop a visually interactive pricing table with hover effects to highlight selected options.	936 937 938
• Redesign the "About Us" section using a card layout for team member profiles, including images and bios.	939 940 941

998

999

1002

1003

1004

1005

1006

1007

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

988

989

1000 1001

• Update the contact form with a modern design, including floating labels and inline validation.

• Animate the scrolling experience for anchor links to smoothly transition between sections of the page.

**Key Principles:** a) Focus on Instructions. b) Only address the requested tasks and avoid unrelated changes unless explicitly instructed. c) Document Changes Clearly and d) For every modification, provide a clear record that includes:

• What was changed?

943

944

945

947

949

953

955

957

960

961

962

964

965

967

969

970

971

973

- Why was it changed?
- The location of the change (e.g., file name and line numbers, or element location in the inline HTML).

#### A.3 Performing Sketch Task

The distinction between simpler and more complex sketches typically depends on the number of components and the level of detail in the specifications. Simpler sketches usually have fewer elements (e.g., basic shapes, minimal labels), while complex sketches include multiple, interrelated components and detailed instructions (e.g., specifying dimensions, class names like 'div nav,' or explicit layout details). To differentiate, consider: **Simple**: Basic wireframes or mockups with minimal annotations (e.g., a rectangle representing a button). **Complex**: Detailed designs specifying attributes (e.g., 'button 200px wide, div with class="nav"") or involving hierarchical or nested components.

#### **B** Task Samples

Tables 4, 5, and 6 present representative exam-974 ples from the WebMMU dataset, covering WebQA, 975 Mockup2Code, and Web Code Editing tasks. The 976 WebQA task (Table 4) evaluates a model's ability to interact with webpage elements, recognize visual content, and perform complex reasoning based 979 on structured UI components. The Mockup2Code task (Table 5) illustrates how webpage screenshots are converted into structured HTML representa-983 tions, distinguishing between basic layout sketches and detailed UI component mappings. The Web 984 Code Editing task (Table 6) demonstrates automated HTML modifications, providing before-andafter visual transformations based on functional 987

and design-driven prompts. These task samples comprehensively showcase the challenges in webpage understanding, layout structuring, and automated UI refinement within the WebMMU benchmark.

# C Model Output Generation Prompts

# C.1 WebQA Task Completion Prompt

We present the prompt used for the WebQA task in Figure 6. The prompt instructs the model to analyze a website screenshot and provide a concrete answer to the given question. When the question requires identifying or interacting with specific elements on the screen, the model is asked to include the bounding box coordinates in its response.

# Web QA Inference

Analyze the website screenshot and provide a detailed answer to the question. If the question involves locating or interacting with specific elements on the screen, include the bounding box coordinates [x\_min, y\_min, x\_max, y\_max] in your response.

Figure 6: Prompt for Generating Output of WebQA task

# C.2 Web Code Editing Task Completion Prompt

This prompt guides a model in modifying the source code based on a modification instruction given by the user. The model outputs changes using the git diff format, highlighting additions and deletions with '+'s and '-'s respectively. This ensures clear and structured documentation of code edits. The prompt template can be seen in Figure 7.

## C.3 Mockup2Code Task Completion Prompt

The Mockup2Code task involves generating HTM-L/CSS code from an input sketch (see Figure 8). Given a visual layout, the model must produce accurate, well-structured HTML and CSS that replicate the design. The prompt guides the model to interpret elements, hierarchy, and styling for faithful image-to-code conversion.

# **D** Evaluation Prompts

This section provides details on the prompt formulations used throughout this work. These prompts102110221022

	Ť	<b>Q</b>	¢
V       V	How can I find more information about the player A.J. Dybantsa? Answer: Click on the "RIVALS150 ranking" at the lower left and se- lect "A.J. Dybantsa" at (x1:230.34, x2:297.32, y1:1049.92, y2:1083.07).	How many players are visible in the 4th image on the left side? A) 2 B) 1 C) 4 D) 3 Answer: D) 3	Which user pays less when subscribing annually vs. monthly? Answer: - Yearly: \$99.95 - Monthly: \$9.95 × 12 = \$119.40 - Savings: \$19.45
	¿En qué área puede el usuario hacer clic para ver publicaciones anteriores? A) Categorías B) Lo más visto C) Archivo del blog Respuesta: C) Archivo del blog.	¿A qué tipo de público está dirigido este blog de recetas? Respuesta: Personas interesadas en re- postería y cocina casera.	¿Cuál sería el efecto en la navegación si la sección "Lo más visto recientemente" estuviera al inicio? Respuesta: Aumentaría la accesibilidad de los artículos populares.
	Je cherche un soin pour le visage. Où dois-je aller ? Réponse: Dans la barre latérale sous "Soins visage par type de peau".	Quel est le prix total des articles si l'on exclut celui avec 5 étoiles ? A) 58,70 $\in$ B) 62,85 $\in$ C) 45,50 $\notin$ D) 51,90 $\in$ Réponse: B) 62,85 $\in$ .	Avec un budget de 15 €, quels produits puis-je acheter ? Réponse: - Masque Argile verte et menthe poivrée bio (5,90 €) - Masque Argile rose bio - Peaux sensibles (5,90 €)
	Comment filtrer les hôtels qui acceptent les animaux ? Réponse: Dans la section "Ausstat- tung", cliquez sur "Haustiere erlaubt".	Combien de chaises sont visibles dans l'image de "Sunflower Santa Maria Novella" ? Réponse: 4 chaises.	Quelles sont les meilleures options d'hébergement près de Via Faenza ? Réponse: "Sunflower Santa Maria Novella" avec une note de 8.9.

Table 4: **WebMMU VQA Task Samples.** This table presents diverse Visual Question Answering (VQA) task samples from the WebMMU dataset, categorized into three types: (1) Functional (interaction with webpage elements), (2) General Understanding (visual recognition within webpage images), and (3) Complex Reasoning (logical inference and numerical computation). Each row showcases an input webpage image alongside representative questions and answers.

guide the multimodal large language models in generating and evaluating responses across different tasks. The prompts are categorized based on their usage, including code modification, VQA evaluation, and UX scoring.

## D.1 WebQA Evaluation Prompt

1023

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1040

These prompts are used for evaluating model responses in VQA tasks. The model rates answers as 1 (Correct and Complete) or 0 (Incorrect or Irrelevant) based on factual accuracy and completeness. Example cases are provided to guide the evaluation. The prompt template can be seen in Figure 9.

## D.2 Mockup2Code Evaluation Prompt

The Mockup2Code evaluation task involves assessing the accuracy of an MLLM-generated website based on an input sketch (see Figure 10). The evaluation prompt directs the annotator to compare the AI-generated HTML/CSS output with the given visual layout, ensuring that the generated website 1041 accurately replicates the design in terms of struc-1042 ture, styling, and layout. The evaluation criteria 1043 focus on layout structure, spacing, proportions, and 1044 alignment, allowing for a detailed assessment of 1045 how closely the generated output matches the intended design. The goal is to evaluate the model's 1047 ability to interpret and transform the sketch into a 1048 functional, visually consistent website. 1049

## **D.3** Code Edit Evaluation Prompt

This prompt is used to evaluate model responses1051in code edition tasks. The model rates answers as10521-5 (5 refers to the most correct and complete, and10531 refers to incorrect or irrelevant) based on factual1054accuracy and completeness. Example cases guide1055the evaluation. The prompt template can be seen in1056Figure 11.1057

## Web Code Editing Generation Prompt

You are an expert web developer specializing in identifying and applying modifications to web code. You will receive a website's screenshot and a combination of it's HTML, CSS, and/or JavaScript code, formatted as follows:

- HTML Code: html\_code
- CSS Code: css\_code
- JavaScript Code: javascript\_code

You will also receive a modification prompt describing the required changes. Your task is to produce the necessary code modifications using 'git diff' format, even if some or all sections are missing. Follow these guidelines:

- 1. Input code: <input\_code>
- 2. Modification Prompt: <edit\_prompt>
- 3. Output Diff:
  - Use '+' for additions and '-' for deletions.
  - Modify only the relevant parts while preserving structure.
  - In case the code is missing, generate the necessary block of code from scratch.
  - Ensure readability and correctness in the modifications.

Only output the necessary diff; do not repeat the input code.

Figure 7: Web Code Editing generation prompt

## E Case Studies of Model Performance

1058

1059

1061

1062

1064 1065

1066

1067

## E.1 Case Studies for the Web Code Editing

We present case studies for the Web Code Editing task illustrating both success and failure examples. Figure 12 shows a success case where the Claude-3.5 model generates code that accurately follows the user's instructions. In contrast, Figure 13 highlights a failure case for the Gemini-2.0-Flash model, which overlooks key modifications requested by the user.

## E.2 Case Studies for the Mockup2Code

We provide several examples illustrating the perfor-1069 mance of different models on the Mockup2Code 1070 task, including both the input mockups and the generated outputs. Figure 14 shows failure cases 1072 where both the best closed-source model (OpenAI-1073 o1) and the best open-source model (Internvl2.5-1074 38b) struggled to accurately reproduce the designs. 1076 In contrast, Figure 15 highlights success cases for the OpenAI-o1 model, demonstrating its ability 1077 to handle both simple and complex mockups ef-1078 fectively. Additionally, Figures 16 and 17 present failure cases specifically for the open-source model 1080

Internvl2.5-8b and closed-source model OpenAI-1081o1, emphasizing areas where it underperforms on1082varying mockup complexities.1083

# Mockup2Code Generation Prompt

You are an expert website developer. Analyze the provided webpage sketch and generate a single, fully structured HTML file with embedded CSS that accurately reflects the design. The output must be a self-contained HTML document with internal <style> tags for CSS. Ensure all elements are structured exactly as seen in the sketch—no extra elements, no missing elements. HTML Requirements:

- **Components:** Include all necessary components such as headers, paragraphs, buttons, forms, and images, maintaining the correct hierarchy and placement.
- **Images:** Use images generated from https://placehold.co/ with exact dimensions matching the sketch, a neutral background color, and centered "Image" text. For example: <img src="https://placehold.co/300x200?text=Image&bg=cccccc" alt="Placeholder Image">
- Placeholder Text: Use Lorem Ipsum for placeholder text where needed.

## **CSS Requirements:**

- Implement CSS directly within the HTML file (inside a <style> block) to match the sketch, covering spacing, font sizes, colors, alignments, and element positioning.
- Use CSS Grid or Flexbox where appropriate to replicate the exact design layout.
- Apply styling for readability and interactive elements (e.g., fonts, button appearance).
- Ensure placeholder images maintain proper dimensions and design consistency.

# **Code Output:**

- Provide a single, complete HTML file with internal CSS (do not separate them into different files).
- Do not include explanations, comments, or any extra formatting outside the code itself.

Figure 8: Mockup2Code Generation Prompt: It takes input sketch and outputs HTML/CSS code of the given input

Input Image	<b>Basic Layout Sketch</b>	Detailed UI Representation
	Neder	
BIS Allows Banks to Hold 2% Top Rated		
of Reserves in Crypto		https://www.inter.com/particular         https://www.inter.com/particular         https://www.inter.com/particular           wd1x22:         jmg/stack_plastical_plagtical_glassical_glassical_glassical_plagtical_glassical_glassical
		wq176x381         wq176x381 <t< td=""></t<>
China page		(726-02) (726-02) (7276-02) (7
The art to be address of the art to be address		1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
acticate ingressments; Is activate ingressment activate activat		(D.5252) (D.668) (0.669)
	Coastin Agreenset	
Annue Alex Sun Mag Sers Grant Jap C C Makine (8 Sers		
Perdanar pro- NOV Schers Index		
িচনত Jonati from mo @ চা খ	Section 34	Section 36 Comment Texts o 640 X 54
Geri Art 200		Incyra 92 X 23 Incyra 92 X 23 Incyre Authoratia 600 X 23
	Comment Texts Comment Texts	Eprint         Comment Texts (p. 640 X 54           43.4.4         Epily as 2X x 20           repsy         struct rise a 640 X 54           struct rise of color 23         struct rise of color 23
Lat	Comment Texts	46 X 43     Theory a 32 X 23     Weapy Autor with a g62 X 23
Nan 1 2001 - Nan 9 2016 Mart	Comment Texts	Lotter         Control 1/2 = 0.01 × 0.01
So very such band and us very such initial attacky by band and bars taken and bars	Comment Texts	Logical Jack 24/3
ar may surry be reach from two lines and an an part of they are noted a part of values of a sure of a part of a markets and any of an almost . Page and stands you for channels plant and and all of theme years.	Comment Texts	Linge Comment lexits p 640 X 54
Contracted WeakWater Service W	Comment Texts	Express 12 X 22      Insign Express 12 X 22      Express 12 X 23      Express 12 X 23      Express 12 X 23      Comment Texts p 640 X 54
The other is a set of	Comment Texts	Image         Convertige         Sector           1mage         QUTEC 1015 x 600 x 53         Controllerit fexts p 640 x 54           6001         Controllerit fexts p 640 x 54
the first of a strained being companies. He gain will assign research in post track. Straining 	Comment Texts	Comment Texts p 640 x 54     A     A     A     A     A     A
Form bur at last to     The set of the	Comment Texts	9000 4070 118 a 600 X 23 601 48 X 48 Comment Texts p 640 X 54
Cite Saray and an and a second s	_	
- Transfer Memory + 1 - 1 - 1 - 2 Memory + 1 - 1 - 2 Memory + 1 - 2 Memory		
160		
PUBLICATIONS AND VIDEOS	Header	Logo (180x80) Logo
THE STOCKED BE WORLDW		NavBar(944x36)
7886 New Yorkenser of View Substantian Andreas Statements		
MBD dataset a trans Talawak bandangan per an mangkang antang sina ministra mangkang bandangan per ang sina per ang	Body	Main(1313x1070)
FABIO MEMBERS LOG-IN TO ACCESS ALL FABIO RESOURCES		
FABIG TECHNICAL NEWSLETTERS		
DOL 5 09         DESC 09         DESC 09           PARAMEN RAVE 701         ADMINIST PARAMENT RAVE 7011         ADMINIST PARAMENT RAVE 7011           CONTRS         CONTRST.         CONTRST.           - Optimum And and Contract Parks on the optimum rave rave rave rave rave rave rave rave	Footer	Quick Links Legals Contact (380/2260) (290/260)
term and a constraint of the second sec		
New - New - New -		

Table 5: **Mockup2Code Task Samples.** This table showcases examples from the Mockup2Code task, illustrating the transformation of webpage images into structured representations. Each row includes: (1) an Input Image (webpage screenshot), (2) a Simple Sketch (basic layout structure), and (3) a Complex Sketch (detailed UI components and text placements).



Table 6: Web Code Editing Task Samples. Code edition before and after screenshot of webpage.

```
Web QA Evaluation Prompt
```

```
examples = [
    {
       "INPUT": {
          "question": "What is the capital of France?",
          "model_answer": "Paris",
          "ground_truth": "Paris",
       },
       "OUTPUT": {
          "rating": 1,
          "rationale": "The model's answer matches the reference answer exactly."
       }
    },
    {
       "INPUT": {
          "question": "What is in the left of the image?",
          "model_answer": "A bus is in the left of the image.",
          "ground_truth": "A dog is in the left of the image.",
      "rating": 0,
          "rationale": "The model's answer is incorrect because the reference answer is 'A dog'."
       }
    },
    {
       "INPUT": {
          "question": "Where is the burger on the table? Tell me the coordinates.",
          "model_answer": "The burger is on the table.",
          "ground_truth": "The burger is on the table at (50, 10, 150, 60).",
       },
       "OUTPUT": {
          "rating": 0,
           "rationale": "The predicted answer is incomplete because it does not provide the
coordinates as requested in the question."
       }
    }
]
test_case = {
    "INPUT": {
       "question": question,
       "model_answer": model_answer,
"ground_truth": ground_truth
    }
}
You are evaluating a Visual Question Answering (VQA) system's response. Compare the model's answer with the ground
truth and rate its accuracy.
Rating Scale (1 or 0):
1 - Correct and Complete: - The predicted answer fully matches the ground truth. - No factual errors or missing details. -
Addresses the question with the correct level of specificity.
0 - Incorrect or Irrelevant: - Any factual errors or mismatches with the reference answer. - Does not address the question
properly. - Provides misleading or irrelevant information.
Examples for reference: json.dumps(examples, indent=4)
Question, Model Answer, and Ground Truth: json.dumps(test_case, indent=4)
You must provide your evaluation in the following JSON format (without any extra text): json.dumps("rating": 0 or 1,
"rationale": "[Brief explanation of why this rating was chosen]") ""
```

Figure 9: LLM-as-judge prompt for WebQA task using few-shot examples

#### Mockup2Code Evaluation Prompt

**Task Overview:** Your task is to evaluate the accuracy of an AI-generated website by comparing it against a provided input sketch. The AI-generated website is provided as an image rendering of the HTML/CSS output. Your goal is to assess how well this rendered image replicates the intended layout from the sketch. **Provided Inputs:** You will receive two images:

- 1. Input Sketch A wireframe illustrating the intended layout.
- 2. **Predicted AI-Rendered Website Image** A screenshot of the website generated from AI-created HTML/CSS based on the sketch.

Since the AI-generated website is provided as an image, your evaluation must be based entirely on visual accuracy, disregarding the underlying code implementation.

#### **Step 1: Detailed Description of Both Images**

For each image (**Input Sketch** and **AI-Rendered Website**), provide a highly-detailed breakdown based on the following categories. Ensure that descriptions follow the same format for both images to facilitate a precise comparison. **1. Identify All Structural Sections:** 

Describe in detail the overall structure of the webpage layout, covering the following:

- Header Does it contain a logo, navigation menu, search bar, or other elements?
- Navigation Bar Describe the menu items. How many items are there? Is the navigation horizontal or vertical?
- Main Content Area Identify distinct sections such as hero banners, text areas, images, or interactive components.
- Sidebars (if applicable) Is there a sidebar for additional navigation, filters, or widgets?
- Footer What content is present (e.g., links, social icons, contact information)?

For the AI-rendered website, note any differences compared to the sketch (e.g., missing sections, extra sections, missing items, misplaced content).

#### 2. List and Describe All Elements:

## List all key elements present in the Input Sketch and AI-Rendered Website:

- Text Elements Titles, paragraphs, labels, lists, captions.
- Images & Icons Identify all image placeholders and their intended placement.
- Buttons & Links Describe all interactive elements like CTAs, navigation links, or form buttons.
- Forms & Inputs Search bars, text fields, dropdowns, checkboxes, radio buttons, etc.
- Tables & Lists If present, describe their structure and formatting.

For the AI-rendered website, specify any elements that are missing, added, or incorrectly placed. **3. Layout & Positioning Details:** 

Describe and analyze the spatial arrangement of elements in both images:

- Column Structure Is the design single-column, multi-column, or grid-based?
- Alignment Are elements aligned left, center, right, or justified?
- Spacing & Proportions Are elements evenly spaced? Are margins, padding, and gaps consistent?
- Relative Proportions Are certain sections (e.g., hero banners, sidebars) larger than others?

For the AI-rendered website, describe any deviations from the sketch (e.g., elements' size differences, elements too large/small, uneven spacing, misalignments).

### Step 2: Evaluation of the AI-Rendered Website

After describing both images, evaluate the AI-generated website's accuracy using the following criteria. Assign a score from 1 to 5 for each.

#### 1. Layout Structure Accuracy (1-5):

Does the generated HTML structure strictly follow the wireframe in layout, hierarchy, and element grouping? This includes the correct placement, nesting, and semantic usage of standard structural elements: <header>, <nav>, <main>, <section>, <aside>, <article>, <footer>, <div>, and content containers like <img>, .

- $5 \rightarrow 100\%$  match. All elements are correctly placed, properly nested, fully grouped, and semantically accurate. No missing, misplaced, or extra elements.
- 4 → Mostly accurate, but minor structural inconsistencies exist (e.g., an unnecessary wrapper, slightly misplaced section, or minor redundancy). No missing elements.
- 3 → Some structural errors at least one missing or misused element, multiple misplaced sections, or noticeable grouping issues.
- $2 \rightarrow$  Major deviations multiple missing, misplaced, or incorrectly nested elements, affecting hierarchy and readability.
- $1 \rightarrow$  Severe structural failure multiple core sections are absent or completely misstructured, making the output unrecognizable compared to the wireframe.

### 2. Spacing & Proportions (1-5):

Do margins, paddings, and element dimensions (e.g., width, height, max-width, min-width, max-height, min-height, gap for flex/grid layouts) precisely match the wireframe?

- $5 \rightarrow 100\%$  correct. All elements have precise margins, paddings, widths, heights, and spacing. No deviations.
- $3 \rightarrow$  Noticeable discrepancies some elements are too large, too small, or unevenly spaced, affecting visual balance.
- $2 \rightarrow$  Significant spacing issues multiple elements have incorrect dimensions, margins, or paddings, leading to a visibly distorted layout.
- $1 \rightarrow$  Severe inaccuracies most elements have incorrect proportions or spacing, making the layout visually broken and inconsistent with the wireframe.

#### 3. Alignment & Grid Consistency (1-5):

Are elements precisely aligned according to the wireframe, following the expected grid/flex structure and ensuring uniform positioning?

- 5  $\rightarrow$  Perfect alignment. Every element follows the wireframe's grid, flex, or positioning structure exactly. No misalignments.
- $4 \rightarrow$  Mostly aligned, but minor deviations exist (e.g., slightly off-center text or small pixel variations in placement).
- $3 \rightarrow$  Some clear misalignments at least one noticeably off-grid or misplaced element that affects overall balance.
- $2 \rightarrow$  Major alignment issues, with multiple elements misaligned, overlapping, or not following the expected structure.
- 1 → Severe disorganization the output fails to follow the wireframe's grid or positioning, making the layout appear chaotic.

### **Final Score Calculation:**

```
"spacing_proportions": [1-5],
    "alignment_grid_consistency": [1-5]
},
    "final_score": [calculated average score],
    "reasoning": "[Concise evaluation highlighting key strengths and weaknesses]"
}
```

Figure 10: LLM-as-Judge input prompt: It evaluates the model output and the ground truth among some detailed criteria given in the prompt.

Model		Englis	h		Frenc	h		Germa	ın		Spanis	sh
	BLEU	TreeBLEU	LLM-as-Judge									
QwenVL-7B	9.02	28.91	2.48	4.11	22.17	1.81	5.41	24.02	2.18	7.22	14.19	2.31
Molmo-7B	1.98	11.91	1.93	3.12	5.77	1.62	1.01	12.30	1.70	1.82	4.21	1.64
Phi-3.5-VI	0.00	0.00	1.02	0.00	0.00	1.00	0.00	0.00	1.00	0.01	0.00	1.00
Fuyu-8B	0.02	0.09	1.07	0.00	0.00	1.06	0.00	0.00	1.06	0.00	1.11	1.03
InternVL-2.5-8B	10.46	25.96	2.50	6.61	14.03	2.07	9.68	23.40	2.34	5.73	15.03	2.13
Glm-4v-9B	6.09	21.74	2.50	4.75	15.23	2.31	4.60	17.00	2.21	5.12	7.38	2.37
Llava-OV-7B	8.08	27.98	2.41	3.32	16.76	2.06	5.78	17.63	2.09	4.42	11.34	2.07
Pixtral-12B	12.16	26.59	2.81	6.28	14.52	2.74	11.07	23.67	2.64	6.04	14.36	2.52
InternVL-2.5-38B	15.84	36.19	3.55	8.01	26.77	3.15	14.12	33.75	3.40	10.14	18.55	3.16
QwenVL-72B	16.00	38.38	3.27	9.40	25.34	3.13	14.16	30.41	3.01	10.36	19.97	3.04
Claude	22.80	38.92	4.04	16.57	24.66	3.69	20.61	32.61	4.00	13.65	22.79	3.64
Gemini-2-Flash	14.34	24.80	3.15	11.11	13.10	3.26	11.62	23.14	3.15	10.71	18.49	3.25
GPT-40	18.94	35.11	3.93	11.81	12.47	3.96	15.47	25.23	3.72	10.89	15.14	3.65

Table 7: Results of Web Code Editing on different languages.

#### Web Code Editing Evaluation Prompt

You are evaluating a system that generates HTML code based on a given task. Compare the predicted code with the ground truth code and rate its correctness based on functionality rather than exact syntax. If the code performs the intended task correctly, even if formatted differently or using a different approach, it should receive a high score. **Rating Scale:** 

- 5 PERFECT Fully achieves the required functionality as described in the reference output. May have differences in syntax or structure, but effectively performs the same task with no missing elements.
- 4 CORRECT BUT WITH MINOR ISSUES Achieves the intended functionality but has small flaws (e.g., slightly different behavior, minor inefficiencies).
- 3 PARTIALLY CORRECT Achieves part of the intended functionality but is missing key aspects or has notable issues.
- 2 MOSTLY INCORRECT Fails to accomplish most of the required functionality but shows some partial effort.
- 1 COMPLETELY INCORRECT The solution does not fulfill the required functionality or is entirely off-target.

```
Examples for reference:
examples = [
   {
      "INPUT": {
         "question": "Change the header's background color to blue."
         "model_answer":"+<style>header{background-color:blue;}</style> <header>Welcome</header>",
         "ground_truth": "<header style='background-color: blue;'>Welcome</header>"},
      "OUTPUT": {
         "rating": 5,
         "rationale": "The model answer correctly implements the change by ensuring the header
displays with a blue background. Despite using a style tag in the model answer versus inline
styling in the ground truth, both approaches deliver the exact intended functionality." }
   }]
Task for Evaluation:
   {
      "INPUT": {
         "question": "<question>",
         "model_answer": "<model_answer>"
         "ground_truth": "<ground_truth>"
      }
   }
Provide your evaluation in the following JSON format (using json delimiters, do not include any extra text):
   {
      "rating": "1 or 2 or 3 or 4 or 5",
      "rationale": "[Brief explanation of why this rating was chosen]"
   }
```

Figure 11: LLM-as-judge prompt for Web Code Editing task using few shot examples.



Figure 12: Code edit case study: **Success** (Claude-3.5). The generated code accurately addresses the instructions given as input.



Figure 13: Code edit case study: Failure (Gemini-2.0-Flash). The generated answer by the model skips main modifications requested by the user.



Figure 14: Examples of the **failure cases on the Mockup2Code task** for the best closed-source model (OpenAI-o1) and the best open-source model (InternVL2.5-38B).



Figure 15: Examples of the **success cases on the Mockup2Code task** for the best closed-source model (OpenAI-o1) for both simple and complex mockups.

Advanced AD and Large to a 20 million from the mark region of the mark region from the mark region for the		ana
<ul> <li>Second Second Sec</li></ul>	p         Image: Constraint of the second secon	202 202 202 202 202 202 202 202 202 202
Kee Meanin     Lee Meanin     Lee Meaning		ing and a set of the s
s ben and a set of the		ana an Mpa an Ma
Original Page	Mockup Image	Generation
T Score Evaluation on Simple	e Sketch	
the sketch. acing: 2. Uneven gaps and inc gnment: 2. Misaligned elemen arall Score: 2	consistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions
the sketch. scing: 2. Uneven gaps and inc gnment: 2. Misaligned elemen srall Score: 2 ternvl2.5-8b	consistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions n a disorganized layout.
the sketch. scing: 2. Uneven gaps and inc ignment: 2. Misaligned elemen erall Score: 2 termvl2.5-8b	consistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions n a disorganized layout.
the sketch. scing: 2. Uneven gaps and inc ignment: 2. Misaligned elemen erall Score: 2 termvl2.5-8b Versener en	consistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions n a disorganized layout.
the sketch. scing: 2. Uneven gaps and inc ignment: 2. Misaligned elemen erall Score: 2 termvl2.5-8b	consistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions n a disorganized layout.
the sketch. scing: 2. Uneven gaps and inc gmment: 2. Misaligned elemen erall Score: 2 termvl2.5-8b	ensistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions in a disorganized layout.
the sketch. scing: 2. Uneven gaps and inc sgment: 2. Misaligned elemen erall Score: 2 ternvl2.5-8b University of the second	ensistent spacing cause improper ts and an irregular grid result i	grouping and distorted proportions n a disorganized layout.

Figure 16: Examples of the **failure cases on the Mockup2Code task** for the open-source model (Internvl2.5-8b) for both simple and complex mockups.



Figure 17: Examples of the **failure cases on the Mockup2Code task** for the best closed-source model (OpenAI-o1) for both simple and complex mockups.