

# ScreenCoder: Advancing Visual-to-Code Generation for Front-End Automation via Modular Multimodal Agents

Anonymous ACL submission

## Abstract

Automating the transformation of user interface (UI) designs into front-end code holds significant promise for accelerating software development and democratizing design workflows. While multimodal large language models (MLLMs) can translate images to code, they often fail on complex UIs, struggling to unify visual perception, layout planning, and code synthesis within a single monolithic model, which leads to frequent perception and planning errors. To address this, we propose ScreenCoder, a modular multi-agent framework that decomposes the task into three interpretable stages: grounding, planning, and generation. By assigning these distinct responsibilities to specialized agents, our framework achieves significantly higher robustness and fidelity than end-to-end approaches. Furthermore, ScreenCoder serves as a scalable data engine, enabling us to generate high-quality image-code pairs. We use this data to fine-tune open-source MLLM via a dual-stage pipeline of supervised fine-tuning and reinforcement learning, demonstrating substantial gains in its UI generation capabilities. Extensive experiments demonstrate that our approach achieves state-of-the-art performance in layout accuracy, structural coherence, and code correctness.

## 1 Introduction

Automating front-end engineering is a critical step toward efficient software development, and recent large language models (LLMs) have advanced the generation of code directly from text instructions (Qwen, 2025; Bolt, 2025). However, this text-based approach faces significant limitations. Generating detailed UIs requires verbose prompts to capture structure and styling, struggles to specify fine-grained visual design like spacing or alignment, and fundamentally deviates from practical design workflows that begin with visual sketches, not paragraphs of text. Relying solely on textual

input is therefore sub-optimal for real-world deployment and often fails to capture the full visual intent.

To bridge this gap, multimodal large language models (MLLMs) offer the promise of directly interpreting UI design images and translating them into code (Yang et al., 2023). While conceptually appealing, our analysis reveals that current MLLMs struggle with this task as it requires a unified set of capabilities, visual understanding, structural layout planning, and domain-specific code synthesis, that they are not holistically designed for. Empirically, this leads to two recurring failure modes: (1) perception errors, where components are missed or misclassified, and (2) planning errors, where components are placed incorrectly or violate layout constraints.

To address these limitations, we propose **ScreenCoder**, a modular multi-agent framework that decomposes the UI-to-code task into three interpretable stages: grounding, planning, and generation. The grounding agent leverages a multimodal large language model to localize and semantically label key UI regions. The planning agent then constructs a hierarchical layout tree using domain knowledge of web layout systems. Finally, the generation agent produces HTML and CSS code via adaptive prompt-based synthesis, incorporating both layout context and optional user instructions to support interactive design. This decomposition introduces architectural modularity, enabling more robust component recognition, layout planning, and structured code generation than end-to-end black-box methods. Experiments show that our framework achieves state-of-the-art performance in layout fidelity, structural coherence, and generation quality.

Beyond inference, our framework acts as a **scalable data engine**. This is crucial because training on raw web data is often infeasible, as its length and noise from dependencies and scripts destabi-

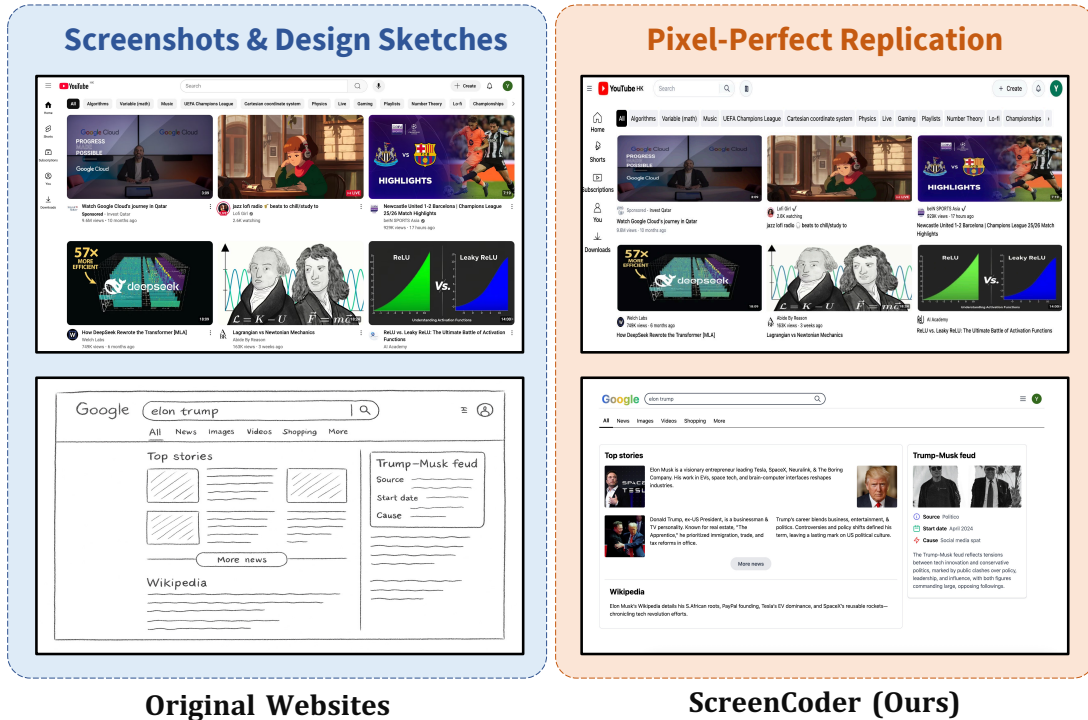


Figure 1: **ScreenCoder** accurately transforms website screenshots and design sketches into pixel-perfect front-end code. The figure showcases a variety of inputs on the left, including high-fidelity screenshots and a low-fidelity design sketch. The right column displays the corresponding webpages rendered from our model’s generated code, demonstrating its high-fidelity replication capabilities.

084 lize training and prevent models from learning the  
 085 core visual-to-code mapping (Si et al., 2025). To  
 086 address the challenge, we leverage ScreenCoder  
 087 to create **Screen-10K**, a new large-scale training  
 088 dataset of 10,000 high-quality image-code pairs,  
 089 curated by filtering an initial crawl of 50,000 web-  
 090 pages. We use Screen-10K to significantly enhance  
 091 open-source MLLM via a dual-stage supervised  
 092 fine-tuning and reinforcement learning pipeline.  
 093 Furthermore, to facilitate a more rigorous eval-  
 094 uation of modern models, we introduce **Screen-**  
 095 **Bench**, a challenging new benchmark composed of  
 096 1,000 high-quality, up-to-date websites that reflect  
 097 contemporary web design. Our framework thus  
 098 provides a practical path for both scalable dataset  
 099 creation and robust model alignment. To sum up,  
 100 our contributions are as follows:

- We conduct a systematic investigation into the limitations of existing MLLMs on UI-to-code tasks and propose a novel modular multi-agent framework that decomposes the complex UI-to-code generation task into three interpretable stages: grounding, planning, and generation, significantly outperforming existing end-to-end multimodal models in layout

084 fidelity and structural coherence. 109

- Leveraging our framework as a scalable data engine, we introduce Screen-10K, a new large-scale dataset containing 10,000 high-quality image-code pairs, which addresses a critical bottleneck in the field by providing a substantial resource for training more capable UI-to-code generation models. 110-116
- To facilitate more rigorous and relevant evaluation, we present ScreenBench, a new challenging benchmark of 1,000 diverse and contemporary web designs. ScreenBench provides a more accurate measure of model performance on real-world tasks compared to existing benchmarks. 117-123

## 2 Related Work 124

### 2.1 Multimodal Large Language Models 125

Multimodal Large Language Models (MLLMs) integrate vision and text to enable joint reasoning. Early models like VisualGPT (Chen et al., 2022) and Frozen (Tsimpoukelli et al., 2021) pioneered this by using pre-trained LLMs to decode 126-130

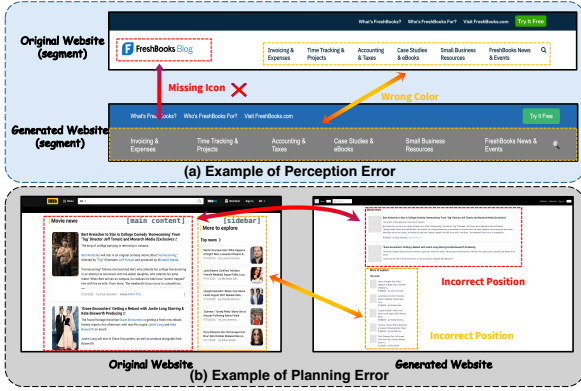


Figure 2: **Analysis of Common MLLM Failure Modes in UI-to-Code Generation.** We identify two primary error categories: (a) Perception Errors, where the model fails to accurately interpret visual details, leading to missing icons or incorrect colors, and (b) Planning Errors, where the model fails to correctly reason about the spatial layout, resulting in elements being placed in the wrong positions.

visual features. Architectural innovations, such as Flamingo’s (Alayrac et al., 2022) gated cross-attention and BLIP-2’s (Li et al., 2023) Q-Former, further improved vision-language alignment. Modern systems like Gemini 2.5 (Google, 2024) and GPT-4o (OpenAI, 2024) have dramatically scaled these capabilities, excelling at complex multimodal tasks and enabling applications like website generation from images (Zhu et al., 2023). However, despite their impressive general-purpose abilities, these models struggle with domain-specific structured generation, such as UI-to-code synthesis. This limitation stems from a lack of inductive biases for spatial layout and hierarchical planning, as well as a monolithic architecture that hinders the injection of task-specific knowledge.

## 2.2 Visual-to-Code Generation

Early visual-to-code methods used CNNs and LSTMs to translate UI screenshots into domain-specific languages (DSLs) (Beltramelli, 2018), which offered limited real-world applicability (Xu et al., 2021). Subsequent research shifted towards generating general-purpose HTML/CSS (Chen et al., 2018) and improving component recognition, layout understanding (Cizotto et al., 2023), interaction-awareness (Xiao et al., 2024), and multi-page generation (Wan et al., 2024). Alternative approaches have included OCR-based techniques (Nguyen and Csallner, 2015) and object detection for screen parsing (Wu et al., 2021). More recently, divide-and-conquer methods (Wan et al.,

2025; Wu et al., 2025; Gui et al., 2025b,a) have emerged, but often depend on heuristic-based segmentation. Despite these advances, prior works are often brittle, rely on synthetic data, and lack the interpretability needed to jointly model complex visual semantics, layouts, and coding patterns. In contrast, our approach introduces a modular multi-agent framework that decomposes the task into interpretable sub-tasks: grounding, planning, and generation. This allows for explicit reasoning and the integration of domain-specific priors. Our system also functions as a scalable data engine to train future MLLMs, addressing the scarcity of high-quality, realistic image-code datasets.

## 3 Motivation: Why MLLMs Fail at UI-to-Code Generation?

To motivate our approach, we analyze why state-of-the-art multimodal large language models (MLLMs) struggle with the direct, end-to-end generation of code from UI screenshots. While models like GPT-4o exhibit powerful general visual reasoning, our analysis of their outputs on real-world webpages reveals two primary failure modes, as illustrated in Figure 2: perception failures and planning failures.

First, **perception failures** stem from the model’s inability to accurately interpret visual elements. This materializes in two ways: (1) **element omission**, where smaller or less salient components like icons and secondary text are completely missed, and (2) **element distortion**, where an element is recognized but its attributes (e.g., text content, color) are incorrect, or its type is misidentified (e.g., an input field rendered as static text). Second, **planning failures** involve the inability to organize perceived elements into a coherent spatial and hierarchical structure. Even if all components are identified, they are often assembled incorrectly, leading to (1) **element misarrangement**, where components are placed in the wrong positions, and (2) **hierarchical incoherence**, where the model produces a “flat” layout that fails to infer the nested, DOM-like structure essential for modern, responsive web design. This highlights a lack of inductive bias for fundamental front-end layout conventions.

Our analysis indicates these failures arise not from an intrinsic inability to perform any single sub-task, but from the burden of a monolithic, end-to-end approach that overloads a single model with granular perception, complex spatial planning, and

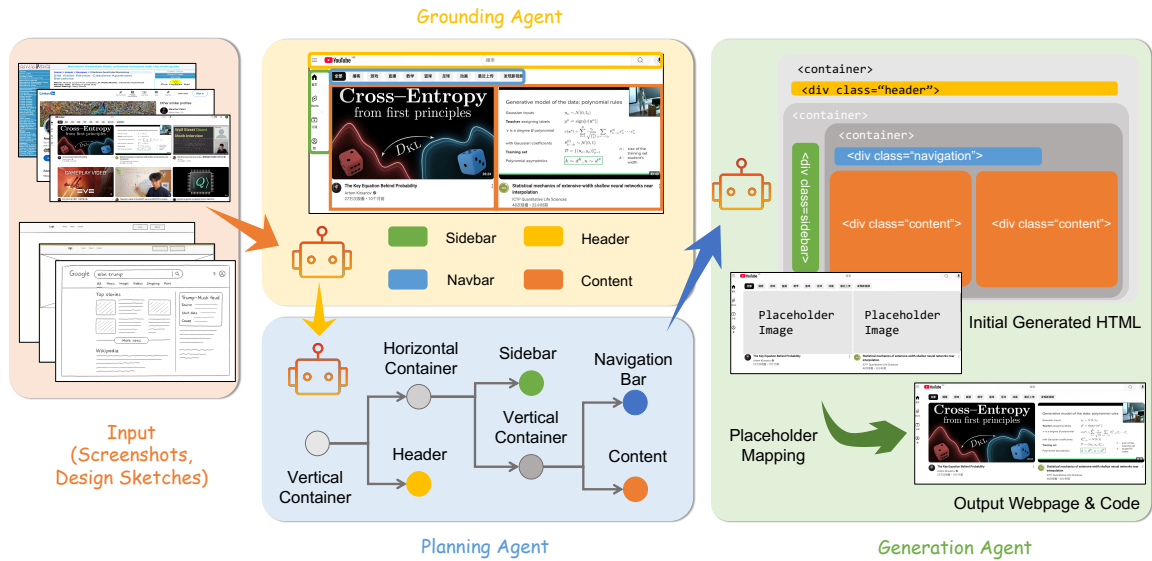


Figure 3: **Overview of ScreenCoder.** Given UI screenshots or design sketches as input, the Grounding Agent first detects and labels key components (e.g., header, navbar, sidebar, content). The Planning Agent organizes these components into a hierarchical layout using front-end engineering priors. The Generation Agent synthesizes initial HTML code with placeholders, followed by content mapping to produce the final webpage and code.

structured code synthesis simultaneously. This task also demands deep domain knowledge of front-end development, such as layout systems and component hierarchies, which general-purpose MLLMs lack. Inspired by agentic workflows that decompose complex problems, we hypothesize that separating these distinct challenges will yield more robust results. This insight directly motivates our modular, multi-agent framework. By explicitly decoupling the task into **grounding** (to address perception failures), **planning** (to address structural failures), and **generation** (to focus on code synthesis), we enable each agent to specialize. Crucially, this modularity allows us to inject specific domain knowledge at each stage, such as established layout conventions in the planning agent, thereby mitigating the failure modes inherent in a single, end-to-end process.

## 4 Method

We propose a modular, multi-agent framework for UI-to-code generation that decomposes the complex task into three sequential agents: grounding, planning, and generation. This design is explicitly motivated by the failure modes identified in Section 2; each agent is specialized to address a distinct sub-problem, allowing the system to leverage both visual understanding and structured reasoning in a coordinated manner. The grounding agent targets *perception errors* by accurately iden-

tifying UI components. The planning agent tackles *planning errors* by constructing a coherent layout hierarchy. Finally, the generation agent translates this structured plan into high-fidelity code. Our overall framework is shown in Figure 3.

### 4.1 Grounding Agent: Overcoming Perception Errors

The grounding agent serves as the perceptual front-end of our framework, tasked with detecting and semantically labeling major structural components to overcome the *perception errors* (element omission and distortion) common in end-to-end models. This design choice, assigning explicit labels like sidebar, header, and navigation, is crucial for enabling interactive, language-driven design, as it allows users and downstream agents to reference and manipulate specific components via natural language (e.g., “resize the sidebar”).

To accomplish this, the agent employs a Multimodal Large Language Model (MLLM) queried with prompts such as “Where is the sidebar?” or “Locate the header area.” The MLLM returns a set of grounded regions as bounding boxes and their associated labels:

$$\mathcal{B} = \{(b_i, l_i) \mid l_i \in \mathcal{L}\}_{i=1}^N, \quad (1)$$

where  $\mathcal{L} = \{\text{sidebar, header, navigation}\}$ . Here, each  $b_i = (x_i, y_i, w_i, h_i)$  is a bounding box in pixel coordinates. Unlike traditional object detection,

---

**Algorithm 1** Visual-to-Structural Tree Mapping

---

```
1: Input: Layout dict  $L$ , Dimensions  $W, H$ 
2: Output: Layout Tree  $\mathcal{T}$ 
3:  $\mathcal{T} \leftarrow \text{createNode}(\text{'root'}, \text{style}=\{$ 
     $\text{'position': 'relative', 'width': '100\%'}, \dots\})$ 
4: for label  $l$ , box  $b_l$  in  $L$  do
5:    $N_l \leftarrow \text{createNode}(l)$ 
6:    $(x\%, y\%, w\%, h\%) \leftarrow (b_l.x/W, b_l.y/H,$ 
     $b_l.width/W, b_l.height/H) \times 100$ 
7:    $N_l.\text{style} \leftarrow \{\text{'pos': 'abs', 'left': } x\%, \dots \}$ 
8:   if  $l$  contains subdivisions then
9:      $N_l.\text{is\_grid\_container} \leftarrow \text{true}$ 
10:   end if
11:    $\text{AddChild}(\mathcal{T}, N_l)$ 
12: end for
13: return  $\mathcal{T}$ 
```

---

269 this MLLM-based approach allows grounding to  
270 be flexibly guided by text, making the system ex-  
271 tensible to new UI concepts.

272 To ensure robustness, the agent performs sev-  
273 eral post-processing operations: (1)Deduplication  
274 and Conflict Resolution, using class-specific non-  
275 maximum suppression (NMS) to filter multiple  
276 detections for the same label and retain the most  
277 confident one; (2)Fallback Recovery, invoking a  
278 heuristic based on spatial priors if a key compo-  
279 nent is missed (e.g., a wide, short box at the top is  
280 likely a header); and (3)Main Content Inference,  
281 which robustly defines the primary content area  
282 by inferring it as the largest rectangular area not  
283 overlapping any detected component. The final  
284 output is a layout dictionary which provides the  
285 semantic and spatial foundation for the next stage.  
286 Unlike traditional object detectors which require  
287 costly retraining, our MLLM-based approach is in-  
288 herently extensible. The system can be adapted to  
289 recognize new UI concepts simply by expanding  
290 the textual label set  $\mathcal{L}$ , offering a flexible path for  
291 future domain adaptations.

## 292 4.2 Planning Agent: Correcting Planning 293 Errors

294 The Planning Agent mitigates common MLLM fail-  
295 ures in spatial reasoning, such as component mis-  
296 arrangement and hierarchical incoherence. Instead  
297 of a generative approach, it uses a novel, determin-  
298 istic Visual-to-Structural Tree Mapping algorithm1  
299 to programmatically translate unstructured visual  
300 information into a well-formed layout. The algo-  
301 rithm embeds key domain knowledge from modern  
302 web development by converting the flat 2D canvas  
303 of bounding boxes into a DOM-like tree, the stan-  
304 dard hierarchical data structure for web pages. This  
305 deliberately trades the unconstrained flexibility of

generative models for structural integrity and pre- 306  
dictability, ensuring the generated code faithfully 307  
mirrors the source screenshot’s layout. 308

309 First, our algorithm establishes the global lay-  
out by creating tree nodes for primary components 310  
(e.g., header, sidebar), converting absolute pixel 311  
coordinates into responsive percentage-based val- 312  
ues, and using absolute positioning to preserve the 313  
macro-structure. It then recursively handles inter- 314  
nal component layouts by injecting domain knowl- 315  
edge, designating parent regions with children as 316  
CSS Grid containers to arrange nested elements 317  
with high fidelity. This process yields an inter- 318  
pretable layout tree that serves as a blueprint, effec- 319  
tively decoupling the abstract planning of the UI 320  
structure from the concrete task of code generation, 321  
embodying the principle of separation of concerns. 322

## 323 4.3 Generation Agent: High-Fidelity Code 324 Synthesis

325 The generation agent translates the hierarchical lay-  
out tree  $\mathcal{T}$  into executable HTML and CSS. It tra- 326  
verses the tree and, for each node, uses a large lan- 327  
guage model to generate code based on an *adaptive* 328  
*prompt*. This prompt combines the component’s 329  
semantic label, its structural context from the tree, 330  
and optional user instructions. This approach pro- 331  
vides the LLM with rich context, guiding it to pro- 332  
duce code that is not only visually correct but also 333  
structurally sound and responsive to user intent. 334  
The generated code snippets for each component 335  
are then assembled according to the tree structure, 336  
preserving the hierarchy and layout defined by the 337  
planning agent. This closes the loop from visual 338  
perception to structured, interactive code synthesis. 339

## 340 4.4 Placeholder Mapping

341 To restore visual fidelity, we introduce a final place-  
holder mapping stage that replaces generic image 342  
placeholders with their original visual assets. The 343  
process begins by using a UI Element Detection 344  
(UIED)(Xie et al., 2020) model to extract all vi- 345  
sual elements (e.g., icons, images) from the source 346  
screenshot. These elements are then partitioned 347  
according to the semantic regions defined by the 348  
Planning Agent (e.g., header, sidebar). 349

350 Within each region, we solve an optimal assign-  
ment problem to match the detected elements to 351  
the placeholders. We construct a cost matrix based 352  
on the negative Complete IoU (CIoU) between the 353  
placeholder boxes and the original element boxes, 354  
after applying a localized affine transformation to 355

correct for minor rendering discrepancies. This bipartite matching problem is solved using the Hungarian Algorithm to find the optimal one-to-one mapping. Finally, the matched image patches are cropped from the original screenshot and inserted into the generated code, restoring the full visual content of the UI.

## 5 Enhancing MLLMs with Scalable Data Generation and Dual-Stage Post-Training

Beyond its inference capabilities, our framework serves a crucial role as a scalable data engine, addressing a fundamental challenge in training visual-to-code models. Directly training on raw, crawled web data is often infeasible (Si et al., 2025). Real-world code is typically long and noisy, replete with complex dependencies, external links, and irrelevant scripts that make training unstable and hinder the model’s ability to learn the core mapping from visual structure to clean, self-contained code.

To overcome this, we leverage our engine to generate Screen-10K, a large-scale dataset of 10,000 high-quality image-code pairs. This dataset was curated by initially crawling 50,000 webpages and applying a rigorous, automated filtering process to retain only the most valuable, well-structured examples. This clean dataset provides the stable foundation necessary for our dual-stage post-training pipeline. First, we perform supervised fine-tuning (SFT) on a 9,000-pair subset to align the model’s visual understanding with correct code syntax, establishing a strong baseline. The remaining 1,000 pairs are then used in a subsequent reinforcement learning (RL) stage to further optimize for high visual fidelity.

The reinforcement learning stage is based on Group Relative Policy Optimization (GRPO) (Shao et al., 2024). We optimize the policy  $\pi_\theta$  to maximize the expected reward over our RL dataset:

$$\max_{\theta} \mathbb{E}_{(x,y) \sim \pi_{\theta}} [\mathcal{R}(x, y)], \quad (2)$$

where  $x$  is the input image and  $y$  is the generated code.

To directly optimize for visual fidelity, our reward function  $\mathcal{R}(x, y)$  is based on the pixel-level similarity between the original screenshot and the webpage rendered from the generated code. For each output  $y$ , we first render it to produce an image,  $\text{Render}(y)$ . The reward is then defined as the negative Mean Squared Error (MSE) between the

original image  $x$  and the rendered output. Since RL seeks to maximize reward, using a negative error term incentivizes the policy to minimize the pixel-wise difference:

$$\mathcal{R}(x, y) = -\text{MSE}(x, \text{Render}(y)) \quad (3)$$

where the MSE between two images  $I_1$  and  $I_2$  of height  $H$  and width  $W$  is calculated as:

$$\text{MSE}(I_1, I_2) = \frac{1}{H \times W} \sum_{i=1}^H \sum_{j=1}^W \|I_1(i, j) - I_2(i, j)\|_2^2 \quad (4)$$

This reward function provides a strong, holistic signal that guides the model toward generating code that achieves a pixel-perfect visual replication of the original design.

## 6 Experiments

We evaluate our framework from two complementary perspectives: (1) the visual fidelity and semantic consistency of the generated webpages, and (2) its effectiveness as a scalable data engine for fine-tuning multimodal large language models (MLLMs). We assess the quality of the generated code by comparing the rendered output against ground-truth screenshots using the metrics and benchmarks detailed below.

### 6.1 Experimental Setup

**Datasets.** To rigorously evaluate modern UI-to-code models, we introduce **ScreenBench**, a new benchmark of 1,000 high-quality image-code pairs. This benchmark addresses the limitations of existing datasets like Design2Code (Si et al., 2025), which, with its 484 samples from older websites, primarily tests for textual content reproduction over structural complexity. In contrast, ScreenBench is substantially larger and sourced from contemporary web applications, featuring the complex, nested layouts (e.g., CSS Grid/Flexbox) that define modern web design. We also adopt Design2Code (Si et al., 2025) for evaluation.

**Evaluation Metrics.** We follow the methodology of Design2Code (Si et al., 2025) and evaluate visual similarity using both high-level and low-level metrics. For high-level assessment, we compute the CLIP similarity (Radford et al., 2021) between the rendered output and the reference screenshot. For low-level evaluation, we extract OCR-based visual blocks from both images, align them using text

Model	ScreenBench					Design2Code				
	Block	Text	Position	Color	CLIP	Block	Text	Position	Color	CLIP
GPT-4o	0.745	0.835	0.725	0.702	0.775	0.845	0.962	0.903	0.881	0.917
GPT-4V	0.721	0.815	0.701	0.682	0.758	0.831	0.955	0.895	0.872	0.905
Gemini-2.5-Pro	0.741	0.825	<u>0.752</u>	0.695	0.788	0.841	0.969	0.901	0.879	0.908
LLaVA 1.6-7B	0.635	0.830	0.544	0.592	0.727	0.736	0.910	0.729	0.816	0.802
DeepSeek-VL-7B	0.680	0.773	0.570	0.614	0.732	0.718	0.824	0.702	0.720	0.843
Qwen2.5-VL	0.723	0.828	0.613	0.632	0.762	0.822	0.951	0.815	0.831	0.893
Seed1.5-VL	0.727	<u>0.852</u>	0.742	<u>0.729</u>	0.783	0.829	<u>0.968</u>	<u>0.915</u>	<u>0.897</u>	0.911
DCGen	0.731	0.831	0.713	0.699	0.767	0.836	<u>0.958</u>	<u>0.885</u>	<u>0.865</u>	0.901
LaTCoder	0.735	0.834	0.718	0.703	0.772	0.839	0.960	0.890	0.869	0.904
LayoutCoder	0.733	0.836	0.717	0.706	0.762	0.841	0.965	0.891	0.876	0.909
Websight-8B	0.678	0.768	0.554	0.606	0.748	0.755	0.903	0.767	0.785	0.859
ScreenCoder (Agentic)	<b>0.768</b>	<b>0.857</b>	<b>0.755</b>	<b>0.734</b>	<b>0.812</b>	<b>0.865</b>	<b>0.975</b>	<b>0.925</b>	<b>0.908</b>	<b>0.922</b>
ScreenCoder (Finetuned)	<u>0.758</u>	0.841	0.742	0.718	<u>0.791</u>	<u>0.849</u>	0.968	0.913	0.886	<u>0.915</u>

Table 1: **Evaluation results on the ScreenBench and Design2Code benchmarks.** For each metric, the best result is in **bold** and the second best is underlined.

449 similarity, and then measure four key aspects based  
450 on these matched elements: block reproduction ac-  
451 curacy, textual consistency, spatial alignment, and  
452 color similarity.

453 **Baselines.** We benchmark our approach against  
454 a comprehensive suite of state-of-the-art models.  
455 This includes leading proprietary MLLMs (GPT-  
456 4o (OpenAI, 2024), GPT-4V (OpenAI, 2023), and  
457 Gemini-2.5-Pro (Google, 2024)); a range of open-  
458 source MLLMs (LLaVA 1.6-7B (Liu et al., 2023),  
459 DeepSeek-VL-7B (Lu et al., 2024), Qwen2.5-  
460 VL (Bai et al., 2025), and Seed1.5-VL (Guo et al.,  
461 2025)); and specialized UI-to-code methods (DC-  
462 Gen (Wan et al., 2025), Websight-8B (Laurençon  
463 et al., 2024)), LaTCoder (Gui et al., 2025a) and  
464 LayoutCoder (Wu et al., 2025). Our method is im-  
465 plemented on the open-source Qwen2.5-VL-32B  
466 model. We show other implementation details in  
467 Appendix C.1 and C.2.

## 468 6.2 Main Results

469 As shown in Table 1, our ScreenCoder framework,  
470 in both its agentic and fine-tuned variants, con-  
471 sistentlly surpasses all open-source baselines. The  
472 primary agentic model achieves state-of-the-art per-  
473 formance, outperforming even top proprietary sys-  
474 tems on the challenging ScreenBench with a Block  
475 score of 0.755. Furthermore, our fine-tuned model  
476 secures the second-best results on many metrics,  
477 and achieve comparable performance with close-  
478 source models. These results validate our frame-  
479 work’s dual utility as both a high-performing infer-  
480 ence system and an effective data engine for signif-

481 icantly enhancing open-source MLLMs. Besides  
482 automatic evaluation, we also conduct human ex-  
483 pert evaluation, whose setting and result are shown  
484 in Appedix A.

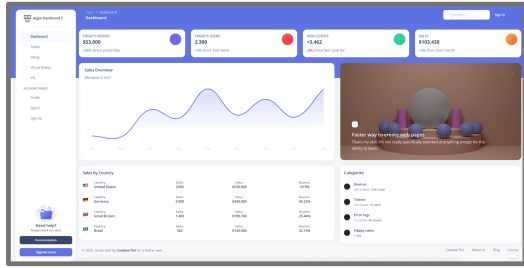
## 485 6.3 Qualitative Analysis

486 Figure 4 provides a qualitative comparison that  
487 highlights the practical advantages of our method.  
488 Leading end-to-end MLLMs, such as Qwen-2.5-  
489 VL and GPT-4o, demonstrate significant perception  
490 and planning failures. They struggle to replicate the  
491 target design, resulting in distorted layouts, incor-  
492 rect component arrangements, and a general loss of  
493 styling information. In stark contrast, ScreenCoder  
494 produces a high-fidelity result that closely mirrors  
495 the original website’s appearance and organization.  
496 More qualitative results are shown in Figure 1 and  
497 Appendix.

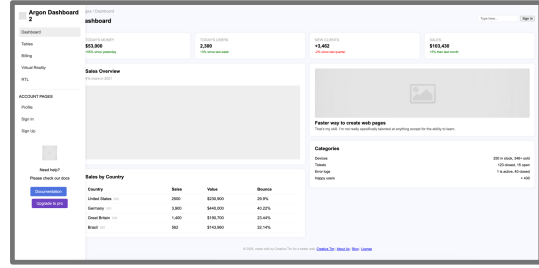
## 498 6.4 Ablation Study

499 **Effect of Dual-Stage Training.** To isolate the con-  
500 tributions of our training pipeline, we conducted a  
501 stage-wise ablation study (Table 2). Starting with  
502 the base Qwen2.5-VL model, the application of  
503 Supervised Fine-Tuning (SFT) yielded significant  
504 improvements, particularly in spatial layout aware-  
505 ness (‘Position’). The subsequent Reinforcement  
506 Learning (RL) stage further refined the model’s ca-  
507 pabilities, providing incremental but crucial boosts  
508 to achieve our final results.

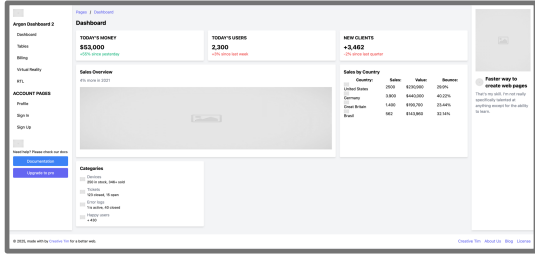
509 **Data Quality vs. Authenticity.** We further  
510 validated the necessity of our synthetic data en-  
511 gine by training the base model on 10,000 raw



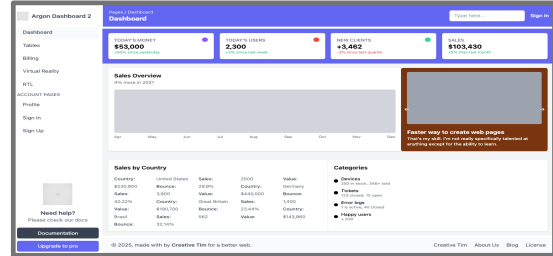
Original Website



Qwen-2.5-VL



GPT-4o



ScreenCoder (ours)

Figure 4: **Qualitative comparison of UI-to-code generation.** While leading MLLMs fail to accurately replicate the target website’s layout, styling, and component structure, our method, ScreenCoder, produces a high-fidelity result that closely matches the original design in both appearance and organization.

Training Stage	ScreenBench					Design2Code				
	Block	Text	Pos.	Color	CLIP	Block	Text	Pos.	Color	CLIP
Base Model (Qwen2.5-VL)	0.723	0.828	0.613	0.632	0.762	0.822	0.951	0.815	0.831	0.893
+ SFT (Screen-10K)	0.741	0.835	0.705	0.681	0.784	0.842	0.963	0.890	0.870	0.908
+ RL (Final)	<b>0.758</b>	<b>0.841</b>	<b>0.742</b>	<b>0.718</b>	<b>0.791</b>	<b>0.849</b>	<b>0.968</b>	<b>0.913</b>	<b>0.886</b>	<b>0.915</b>

Table 2: **Impact of our dual-stage training pipeline (SFT followed by RL).** SFT provides the foundational structural understanding, while RL refines pixel-level alignment.

Data Source	ScreenBench					Design2Code				
	Block	Text	Pos.	Color	CLIP	Block	Text	Pos.	Color	CLIP
Base Model	0.723	0.828	0.613	0.632	0.762	0.822	0.951	0.815	0.831	0.893
Real-10K (Raw Data)	0.719	0.832	0.605	0.628	0.755	0.825	0.947	0.811	0.825	0.891
Screen-10K (Ours)	<b>0.741</b>	<b>0.835</b>	<b>0.705</b>	<b>0.681</b>	<b>0.784</b>	<b>0.842</b>	<b>0.963</b>	<b>0.890</b>	<b>0.870</b>	<b>0.908</b>

Table 3: **Validation of the Data Engine.** Training on raw web data (Real-10K) degrades structural performance due to noise, whereas our synthetic data (Screen-10K) significantly improves it.

HTML/screenshot pairs collected from the same websites (Real-10K). As shown in Table 3, training on raw web data resulted in negligible gains or even performance degradation (e.g., ScreenBench Block score dropped from 0.723 to 0.719). This confirms that the noise in raw web code such as mini-fied classes and redundant wrappers hinders visual-structural alignment. In contrast, our cleaned, synthetic Screen-10K data provides a far stronger training signal, demonstrating that code quality is more critical than code authenticity.

## 6.5 Scalability and the Cost-Quality Trade-Off.

Our modular framework offers a flexible trade-off between computational cost and output quality via test-time scaling. Users can select smaller, faster models for quick, low-fidelity drafts or larger, more powerful models for high-fidelity, production-ready code. While a high-quality generation takes tens

of seconds, this is an acceptable trade-off for its intended use as a workflow accelerator where initial code quality is the priority. This versatility allows the system to be adapted for different stages of the development process. Future work can further enhance this flexibility. We envision a system where developers can interactively refine the output, dedicating more compute only to the specific components that require changes.

## 7 Conclusion

We present ScreenCoder, a modular multi-agent framework for UI-to-code generation that addresses key limitations of end-to-end models, and also functions as a scalable data engine to improve MLLMs via dual-stage post-training. Experiments demonstrate state-of-the-art performance in visual fidelity and code correctness, offering a practical solution and foundation for front-end automation.

## 549 Limitations

550 We discuss two main limitations to our work.

551 **Scope of Dynamic Interactivity.** ScreenCoder  
552 is currently optimized for the high-fidelity recon-  
553 struction of visual layouts and responsive styling,  
554 effectively bridging the gap between design and  
555 front-end implementation. As such, the system fo-  
556 cuses on generating static HTML and CSS. While  
557 it successfully handles interactive states (e.g., hover  
558 effects, responsive resizing), the automated genera-  
559 tion of complex client-side logic, data fetching, or  
560 backend integrations remains outside the current  
561 scope. Extending the agentic framework to infer  
562 and implement intricate behavioral logic represents  
563 a promising direction for future work.

564 **Approximation of Non-Standard Geometries.**  
565 To prioritize the generation of clean, maintain-  
566 able, and standard-compliant code, our Planning  
567 Agent operates within the constraints of the stan-  
568 dard web box model. Consequently, when encoun-  
569 tering highly artistic or experimental designs fea-  
570 turing non-rectangular geometry (e.g., significant  
571 skewing or irregular arbitrary shapes), the system  
572 is designed to approximate these elements into a ro-  
573 bust grid-based structure. While this design choice  
574 ensures that the output remains functional and pro-  
575 fessionally structured, it may result in a simplifica-  
576 tion of extreme stylistic deviations in favor of code  
577 validity.

## 578 References

579 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,  
580 Antoine Miech, Iain Barr, Yana Hasson, Karel  
581 Lenc, Arthur Mensch, Katherine Millican, Malcolm  
582 Reynolds, and 1 others. 2022. Flamingo: a visual  
583 language model for few-shot learning. *Advances in*  
584 *neural information processing systems*, 35:23716–  
585 23736.

586 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wen-  
587 bin Ge, Sibao Song, Kai Dang, Peng Wang, Shi-  
588 jie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu,  
589 Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei  
590 Wang, Wei Ding, Zheren Fu, Yiheng Xu, and 8 oth-  
591 ers. 2025. [Qwen2.5-vl technical report](#). *Preprint*,  
592 arXiv:2502.13923.

593 T. Beltramelli. 2018. pix2code: Generating code from  
594 a graphical user interface screenshot. In *Proceed-*  
595 *ings of the ACM SIGCHI Symposium on Engineering*  
596 *Interactive Computing Systems*, pages 1–6.

597 Bolt. 2025. [Introduction to bolt](#). Accessed: 2025-07-  
598 23.

C. Chen, T. Su, G. Meng, Z. Xing, and Y. Liu. 2018. 599  
From ui design image to gui skeleton: a neural ma- 600  
chine translator to bootstrap mobile gui implementa- 601  
tion. In *Proceedings of the 40th International Con-* 602  
*ference on Software Engineering*, pages 665–676. 603

Jun Chen, Han Guo, Kai Yi, Boyang Li, and Mohamed 604  
Elhoseiny. 2022. Visualgpt: Data-efficient adapta- 605  
tion of pretrained language models for image caption- 606  
ing. In *Proceedings of the IEEE/CVF Conference* 607  
*on Computer Vision and Pattern Recognition*, pages 608  
18030–18040. 609

A. A. J. Cizotto, R. C. T. de Souza, V. C. Mariani, 610  
and L. dos Santos Coelho. 2023. Web pages from 611  
mockup design based on convolutional neural net- 612  
work and class activation mapping. *Multimedia Tools* 613  
*and Applications*, pages 1–27. 614

Google. 2024. [Gemini api](#). Accessed: 2024-06-06. 615

Yi Gui, Zhen Li, Zhongyi Zhang, Guohao Wang, Tian- 616  
peng Lv, Gaoyang Jiang, Yi Liu, Dongping Chen, 617  
Yao Wan, Hongyu Zhang, Wenbin Jiang, Xuanhua 618  
Shi, and Hai Jin. 2025a. [Latcoder: Converting web-](#) 619  
[page design to code with layout-as-thought](#). In *Pro-* 620  
*ceedings of the 31st ACM SIGKDD Conference on* 621  
*Knowledge Discovery and Data Mining V.2*, KDD 622  
'25, page 721–732, New York, NY, USA. Association 623  
for Computing Machinery. 624

Yi Gui, Yao Wan, Zhen Li, Zhongyi Zhang, Dongping 625  
Chen, Hongyu Zhang, Yi Su, Bohua Chen, Xing 626  
Zhou, Wenbin Jiang, and Xiangliang Zhang. 2025b. 627  
[Uicopilot: Automating ui synthesis via hierarchical](#) 628  
[code generation from webpage designs](#). In *Proceed-* 629  
*ings of the ACM on Web Conference 2025*, WWW 630  
'25, page 1846–1855, New York, NY, USA. Associa- 631  
tion for Computing Machinery. 632

Dong Guo, Faming Wu, Feida Zhu, Fuxing Leng, 633  
Guang Shi, Haobin Chen, Haoqi Fan, Jian Wang, 634  
Jianyu Jiang, Jiawei Wang, Jingji Chen, Jingjia 635  
Huang, Kang Lei, Liping Yuan, Lishu Luo, Pengfei 636  
Liu, Qinghao Ye, Rui Qian, Shen Yan, and 178 oth- 637  
ers. 2025. [Seed1.5-vl technical report](#). *Preprint*, 638  
arXiv:2505.07062. 639

Hugo Laurençon, Léo Tronchon, and Victor Sanh. 2024. 640  
[Unlocking the conversion of web screenshots into](#) 641  
[html code with the websight dataset](#). *Preprint*, 642  
arXiv:2403.09029. 643

Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 644  
2023. Blip-2: Bootstrapping language-image pre- 645  
training with frozen image encoders and large lan- 646  
guage models. In *International conference on ma-* 647  
*chine learning*, pages 19730–19742. PMLR. 648

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae 649  
Lee. 2023. Visual instruction tuning. *arXiv preprint* 650  
arXiv:2304.08485. 651

Haoyu Lu, Wen Liu, Bo Zhang, Bingxuan Wang, Kai 652  
Dong, Bo Liu, Jingxiang Sun, Tongzheng Ren, Zhu- 653  
oshu Li, Hao Yang, Yaofeng Sun, Chengqi Deng, 654

655	Hanwei Xu, Zhenda Xie, and Chong Ruan. 2024.	J. Wu, X. Zhang, J. Nichols, and J. P. Bigham. 2021.	708
656	<a href="#">Deepseek-vl: Towards real-world vision-language understanding</a> . <i>Preprint</i> , arXiv:2403.05525.	Screen parsing: Towards reverse engineering of ui models from screenshots. In <i>The 34th Annual ACM Symposium on User Interface Software and Technology</i> , pages 470–483.	709
657			710
658	T. A. Nguyen and C. Csallner. 2015. Reverse engineering mobile application user interfaces with remaui (t).		711
659	In <i>2015 30th IEEE/ACM International Conference on Automated Software Engineering (ASE)</i> , pages 248–259.	Jingyu Xiao, Yuxuan Wan, Yintong Huo, Zhiyao Xu, and Michael R. Lyu. 2024. <a href="#">Interaction2code: How far are we from automatic interactive webpage generation?</a> <i>ArXiv</i> , abs/2411.03292.	713
660			714
661			715
662			716
663	OpenAI. 2023. <a href="#">Gpt-4v(ision) system card</a> . Accessed: 2025-07-24.	Mulong Xie, Sidong Feng, Zhenchang Xing, Jieshan Chen, and Chunyang Chen. 2020. <a href="#">Uied: a hybrid tool for gui element detection</a> . In <i>Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, ESEC/FSE 2020</i> , page 1655–1659, New York, NY, USA. Association for Computing Machinery.	717
664			718
665	OpenAI. 2024. <a href="#">Hello gpt-4o</a> . Accessed: 2024-06-06.		719
666	Qwen. 2025. <a href="#">Qwen3 blog post</a> . Accessed: 2025-07-23.		720
667	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. <a href="#">Learning transferable visual models from natural language supervision</a> . <i>arXiv preprint arXiv:2103.00020</i> .		721
668			722
669			723
670			724
671			725
672			726
673	Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, Y.K. Li, Y. Wu, and Daya Guo. 2024. <a href="#">Deepseekmath: Pushing the limits of mathematical reasoning in open language models</a> . <i>arXiv preprint arXiv:2402.03300</i> .	Y. Xu, L. Bo, X. Sun, B. Li, J. Jiang, and W. Zhou. 2021. image2emmet: Automatic code generation from web user interface image. <i>Journal of Software: Evolution and Process</i> , 33(8):e2369.	727
674			728
675			729
676			730
677			731
678			732
679	Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2025. <a href="#">Design2Code: Benchmarking multimodal code generation for automated front-end engineering</a> . In <i>Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)</i> , pages 3956–3974, Albuquerque, New Mexico. Association for Computational Linguistics.	Zhengyuan Yang, Linjie Li, Kevin Lin, Jianfeng Wang, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. 2023. <a href="#">The dawn of lmms: Preliminary explorations with gpt-4v(ision)</a> . <i>ArXiv</i> , abs/2309.17421.	733
680			734
681			735
682			736
683			737
684			738
685			739
686			740
687			741
688	Maria Tsimpoukelli, Jacob L Menick, Serkan Cabi, SM Eslami, Oriol Vinyals, and Felix Hill. 2021. Multimodal few-shot learning with frozen language models. <i>Advances in Neural Information Processing Systems</i> , 34:200–212.	Yaowei Zheng, Richong Zhang, Junhao Zhang, Yanhan Ye, Zheyang Luo, Zhangchi Feng, and Yongqiang Ma. 2024. <a href="#">Llamafactory: Unified efficient fine-tuning of 100+ language models</a> . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , Bangkok, Thailand. Association for Computational Linguistics.	742
689			743
690			744
691			
692			
693	Yuxuan Wan, Yi Dong, Jingyu Xiao, Yintong Huo, Wenxuan Wang, and Michael R. Lyu. 2024. <a href="#">Mrweb: An exploration of generating multi-page resource-aware web code from ui designs</a> . <i>ArXiv</i> , abs/2412.15310.	Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. 2023. <a href="#">Minigpt-4: Enhancing vision-language understanding with advanced large language models</a> . <i>ArXiv</i> , abs/2304.10592.	741
694			742
695			743
696			744
697			
698	Yuxuan Wan, Chaozheng Wang, Yi Dong, Wenxuan Wang, Shuqing Li, Yintong Huo, and Michael Lyu. 2025. <a href="#">Divide-and-conquer: Generating ui code from screenshots</a> . <i>Proceedings of the ACM on Software Engineering</i> , 2(FSE):2099–2122.		
699			
700			
701			
702			
703	Fan Wu, Cuiyun Gao, Shuqing Li, Xin-Cheng Wen, and Qing Liao. 2025. <a href="#">Mllm-based ui2code automation guided by ui layout information</a> . <i>Proceedings of the ACM on Software Engineering</i> , 2(ISSTA):1123–1145.		
704			
705			
706			
707			

## Appendix

### A Human Evaluation

To validate our automatic metrics, we conducted a human evaluation study assessing two critical dimensions: the **perceived visual fidelity** of the generated webpages and the **practical utility** of ScreenCoder in a real-world development workflow.

#### A.1 Pairwise Comparison of Visual Fidelity

**Setup.** We recruited six Ph.D. students with web development experience to serve as expert annotators. Following the standard methodology for evaluating generative models, we performed a pairwise comparison study. In each trial, annotators were shown an original UI screenshot alongside two rendered webpages—one generated by **ScreenCoder (Agentic)** and one by a strong baseline (**GPT-4o**). They were asked to vote for the page that was more visually similar to the original ("A is better," "B is better," or "Tie"). To ensure reliable judgments, each comparison was evaluated by three annotators, and a winner was declared based on a majority vote ( $\geq 2$ ).

**Results.** The results, summarized in Table 4, reveal a strong human preference for our method. Annotators judged ScreenCoder’s output as superior to GPT-4o’s in **65%** of cases, while finding it inferior in only **11%**. This outcome confirms that the improvements captured by our automatic metrics translate into a noticeably better perceptual experience, suggesting that our modular approach produces layouts that are more structurally coherent and visually faithful.

Table 4: Pairwise comparison of visual fidelity between ScreenCoder and GPT-4o. Results show the percentage of times each outcome was chosen by human annotators.

Outcome	Win Rate (%)
ScreenCoder Wins	65%
GPT-4o Wins (Baseline)	11%
Tie	24%

#### A.2 Workflow Usefulness Study

**Setup.** To measure ScreenCoder’s impact on developer productivity, we designed a task-based study. Six participants with front-end experience were divided into an **Experimental Group (n=3)**,

Table 5: Results of the workflow usefulness study. We compare the performance of developers using ScreenCoder against a control group using GPT-4o.

Metric	ScreenCoder (Experimental)	GPT-4o (Control)
Avg. Completion Time (min)	8.5	18.7
Avg. UI Similarity (out of 5.0)	4.6	3.4
Tasks Timed Out (out of 6)	0	2

who used a tool powered by ScreenCoder, and a **Control Group (n=3)**, who used their preferred method (all chose GPT-4o). Each participant was tasked with converting two UI screenshots (one simple, one complex) into functional webpages, with a 20-minute time limit for each task. We measured both the completion time and the quality of the final output, which was rated for UI similarity on a 5-point Likert scale by two independent expert judges.

**Results.** As detailed in Table 5, ScreenCoder provides a significant boost to development efficiency. On average, the ScreenCoder group completed tasks **2.2 times faster** than the control group (8.5 minutes vs. 18.7 minutes). Notably, all participants using ScreenCoder finished comfortably within the time limit, whereas two tasks in the control group timed out. Furthermore, the quality of the output was substantially higher for the ScreenCoder group, which achieved an average similarity score of **4.6/5.0**, compared to **3.4/5.0** for the control group (with a strong inter-judge Pearson correlation of 0.78). These findings demonstrate that ScreenCoder acts as a powerful accelerator in a human-in-the-loop process, enabling developers to build UIs faster and with greater accuracy.

### B Additional Qualitative Comparisons

To further demonstrate the effectiveness of our framework, this appendix provides an extended set of qualitative results. The following figures showcase side-by-side comparisons between the webpages generated by our method, ScreenCoder, and those produced by leading baseline models.

These examples were selected to cover a diverse range of website styles and layout complexities. They serve to highlight the common failure modes of existing end-to-end models—such as incorrect spatial arrangements, missing UI elements, and distorted styling—and illustrate how ScreenCoder’s modular, agentic approach successfully overcomes these challenges. As shown in the comparisons, our method consistently produces webpages with

825 significantly higher visual fidelity and structural co-  
826 herence, more closely matching the original design  
827 intent of the source screenshots.

## 828 C Implementation Detail

### 829 C.1 Prompt Template

830 The prompt templates are shown in Figure 13 and  
831 14.

### 832 C.2 Training Details

833 Our training and experiments were conducted using  
834 the Llama Factory codebase on a high-performance  
835 cluster of 16 NVIDIA A100 GPUs. We adopted a  
836 two-stage training pipeline. For the SFT stage, we  
837 fine-tuned the base model on 9,000 samples from  
838 our Screen-10K dataset. We adopt the LLaMa Fac-  
839 tory (Zheng et al., 2024) as the code base. The  
840 model was optimized using the AdamW optimizer  
841 with a cosine learning rate schedule, a peak learn-  
842 ing rate of  $2 \times 10^{-5}$ , a weight decay of 0.01, and  
843 a warmup ratio of 10%. In the subsequent RL  
844 stage, we used the remaining 1,000 samples from  
845 Screen-10K. We employed the Group Relative Pol-  
846 icy Optimization (GRPO) algorithm, continuing  
847 from the SFT checkpoint. The reward function was  
848 based on the negative Mean Squared Error (MSE)  
849 between the rendered output and the ground-truth  
850 image, directly optimizing for visual fidelity.

## 851 D LLM Usage Statement

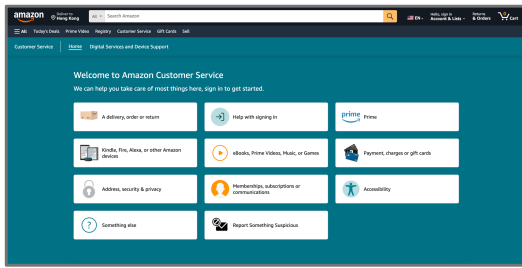
852 We used LLMs as an assistive tool for improv-  
853 ing grammar and clarity in the text, and for code  
854 completion. The core research, including the ex-  
855 perimental design and final implementation, was  
856 conceived and executed entirely by the authors.

## 857 E Failure Case Analysis and Robustness 858 Check

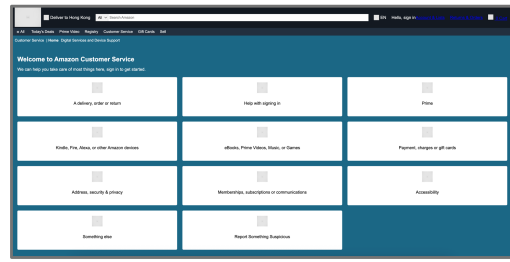
859 In this section, we provide a qualitative analysis  
860 of failure modes, specifically focusing on “Hard  
861 Cases” involving non-standard, artistic layouts that  
862 deviate from the standard rectangular box model.  
863 We compare the performance of ScreenCoder  
864 against the strong end-to-end baseline, Qwen2.5-  
865 VL.

866 As illustrated in Figures 15 to 17, we test the  
867 models on a design featuring significant CSS trans-  
868 formations (e.g., skewY) and overlapping elements.  
869 This analysis highlights a critical trade-off between  
870 pixel imitation and structural integrity:

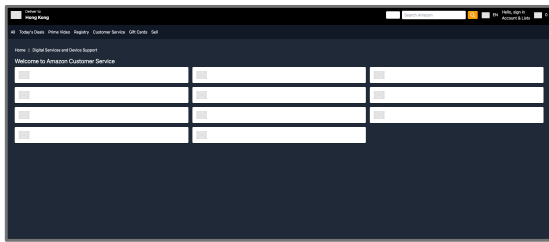
- 871 • **Catastrophic Failure in Baselines** 871  
872 **(Qwen2.5-VL):** Monolithic MLLMs tend to 872  
873 prioritize pixel-level visual matching. When 873  
874 encountering skewed or rotated elements, the 874  
875 baseline attempts to replicate the geometry 875  
876 using absolute positioning and incorrect 876  
877 transform approximations. This results 877  
878 in “spatial hallucination,” where elements 878  
879 suffer from Z-index errors, text overlaps, 879  
880 and fragmented coordinate placement. The 880  
881 resulting code, while attempting to look like 881  
882 the input, is functionally unusable. 882
  
- 883 • **Graceful Degradation in ScreenCoder** 883  
884 **(Ours):** Our modular approach prioritizes 884  
885 valid DOM structure. The Planning Agent 885  
886 maps the non-rectangular visual input to 886  
887 the nearest robust rectangular grid (CSS 887  
888 Flexbox/Grid). While this results in a loss 888  
889 of the specific artistic nuance (the skew an- 889  
890 gle is removed), the system fails *gracefully*. 890  
891 The generated code remains clean, responsive, 891  
892 and structurally valid. This demonstrates that 892  
893 ScreenCoder prefers simplification over bro- 893  
894 ken complexity, ensuring the output is always 894  
895 a viable starting point for developers. 895



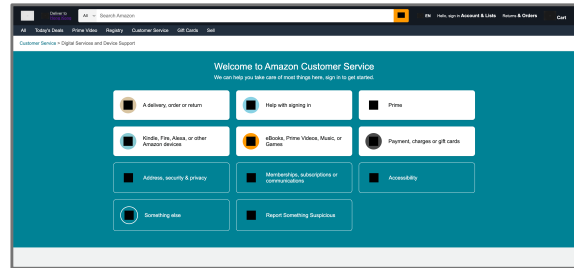
Original Website



Generated by Doubao-1.5-thinking-vision-pro

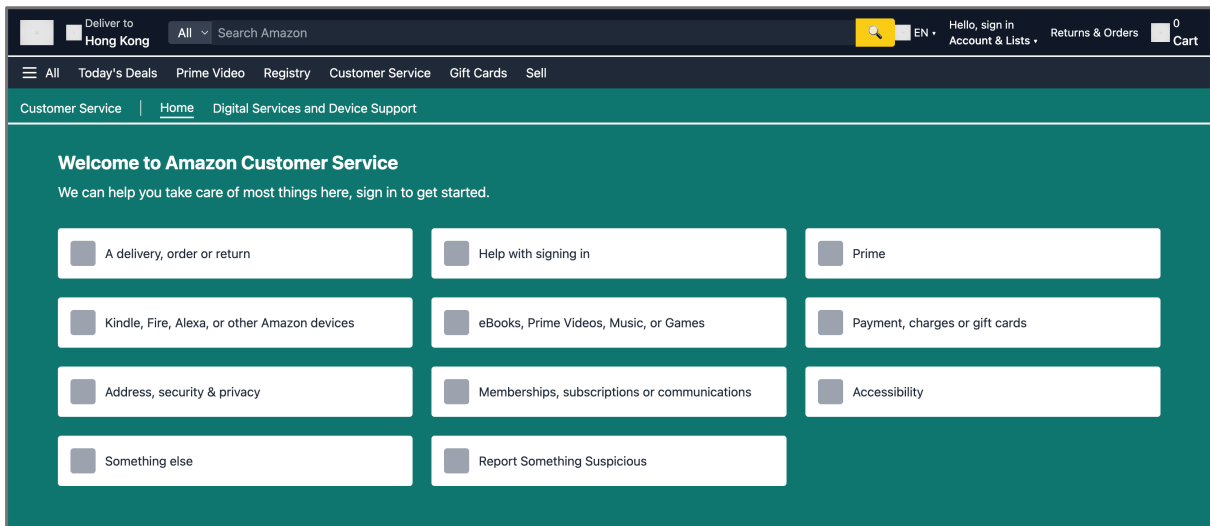


Generated by Qwen-2.5-VL



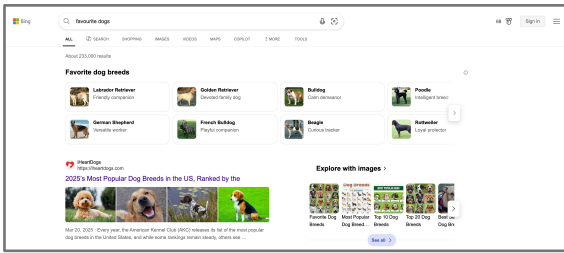
Generated by Gemini-2.5-pro

Figure 5: Qualitative comparison between our proposed method and various baselines.

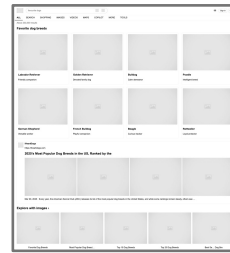


Generated by ScreenCoder (ours)

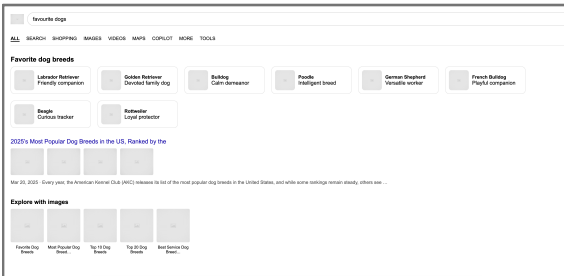
Figure 6: Qualitative comparison between our proposed method and various baselines.



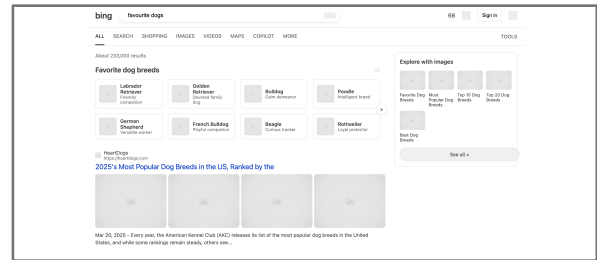
Original Website



Generated by Doubao-1.5-thinking-vision-pro

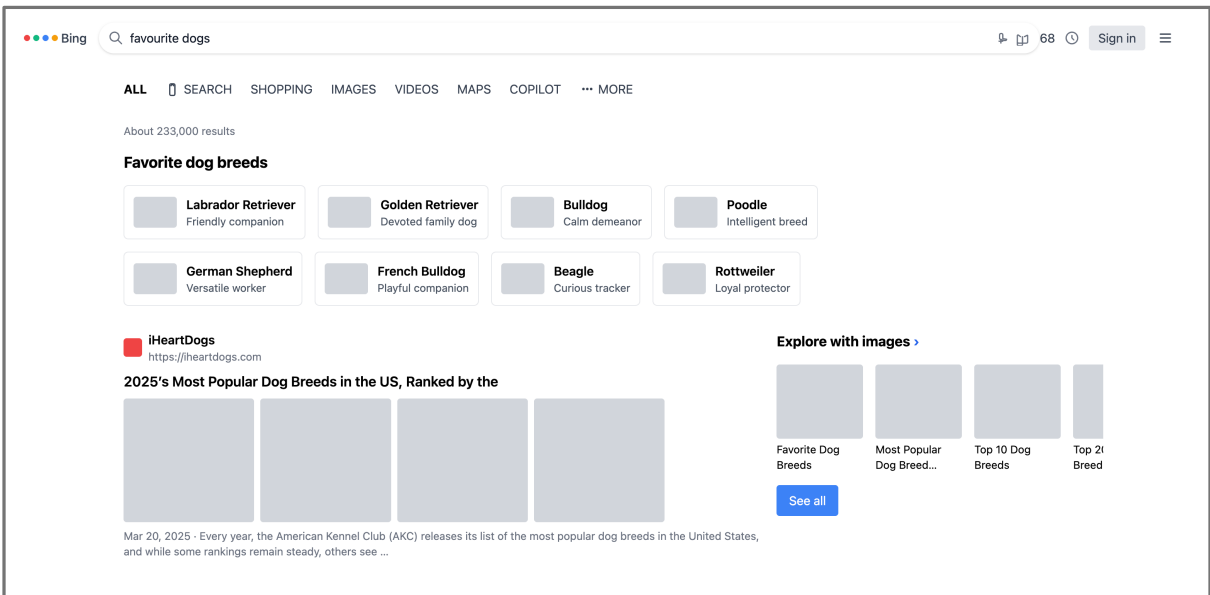


Generated by GPT-4o



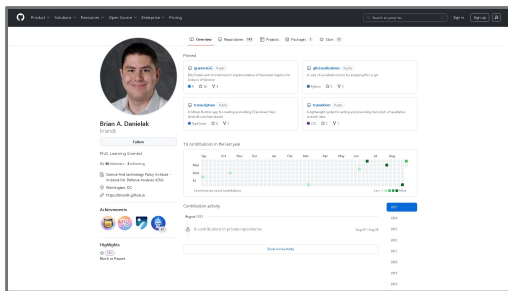
Generated by Gemini-2.5-pro

Figure 7: Qualitative comparison between our proposed method and various baselines.

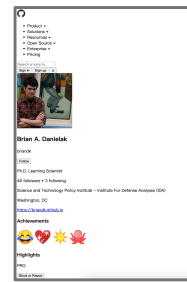


Generated by ScreenCoder (ours)

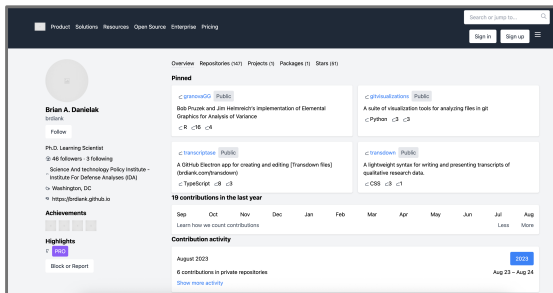
Figure 8: Qualitative comparison between our proposed method and various baselines.



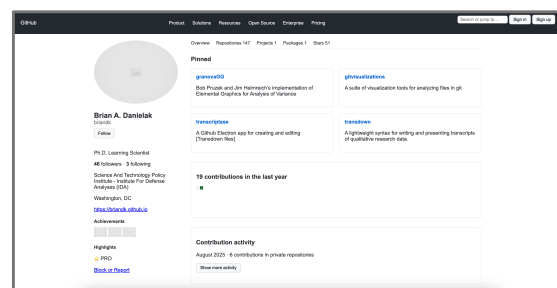
Original Website



Generated by Doubao-1.5-thinking-vision-pro

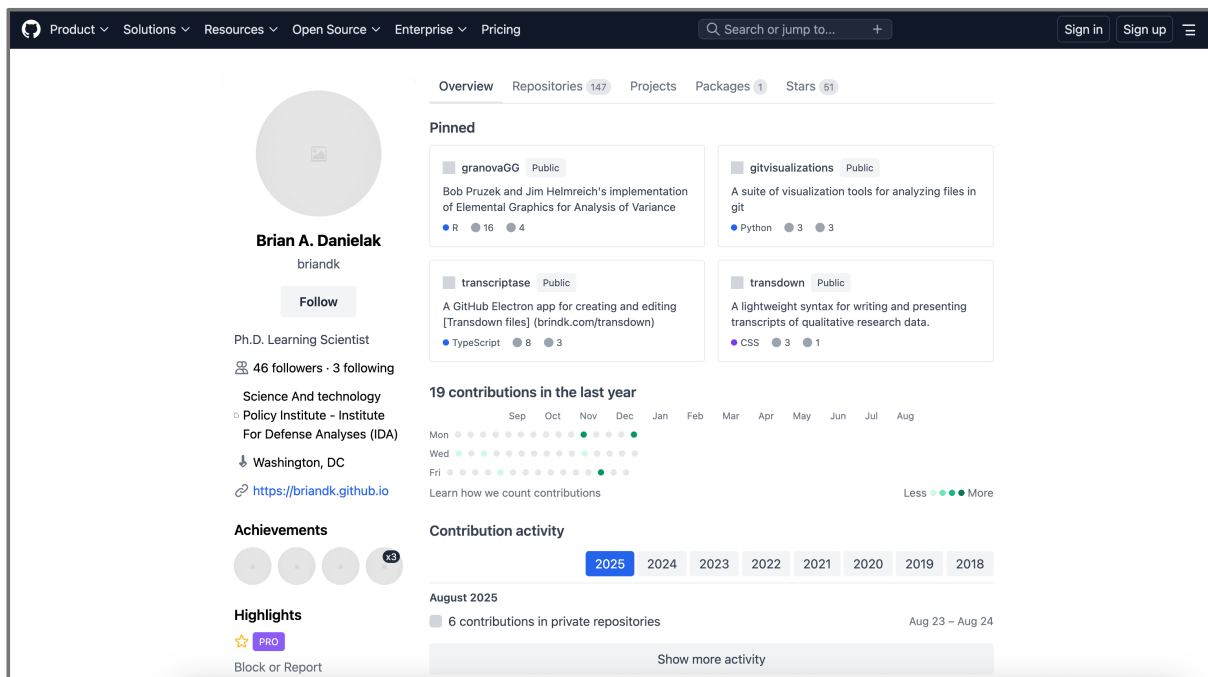


Generated by Qwen-2.5-VL



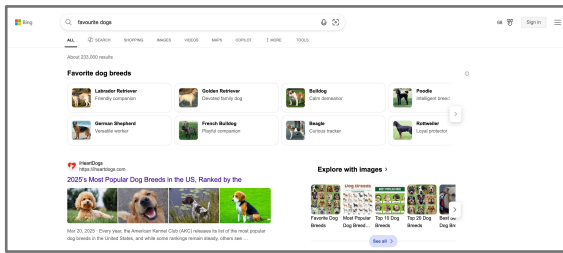
Generated by GPT-4o

Figure 9: Qualitative comparison between our proposed method and various baselines.

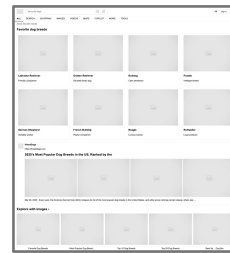


Generated by ScreenCoder (ours)

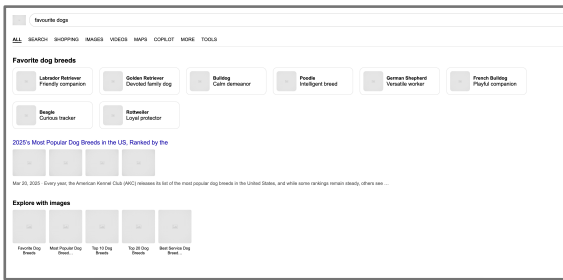
Figure 10: Qualitative comparison between our proposed method and various baselines.



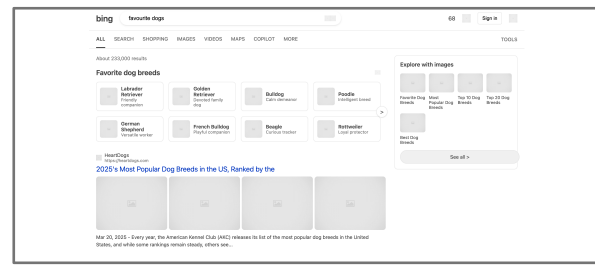
Original Website



Generated by Doubao-1.5-thinking-vision-pro

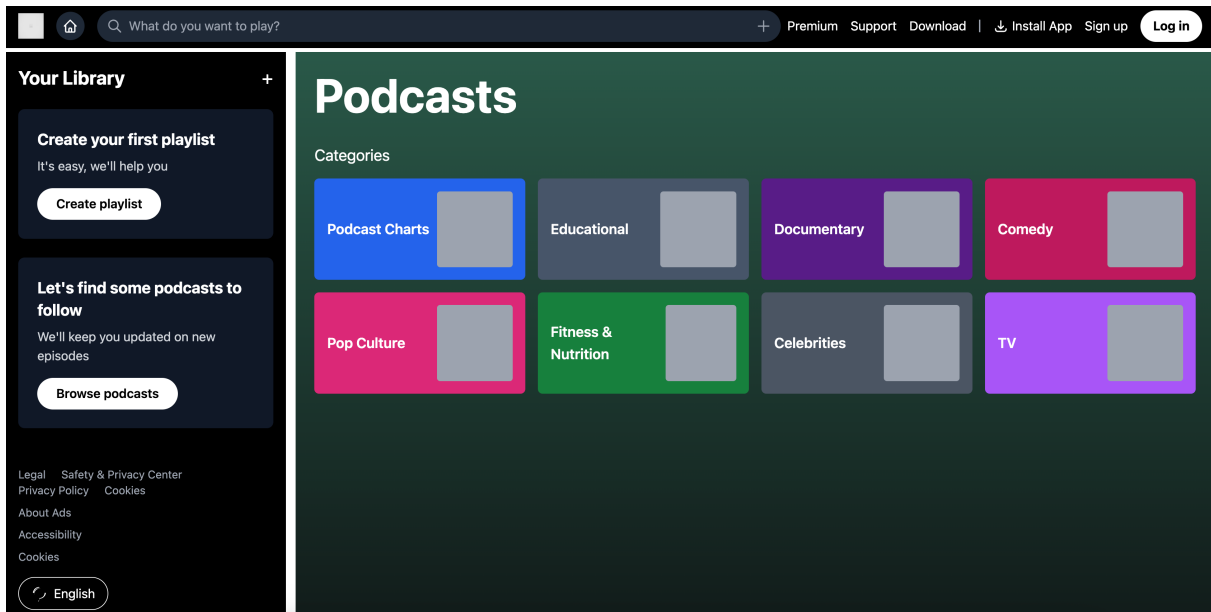


Generated by GPT-4o



Generated by Gemini-2.5-pro

Figure 11: Qualitative comparison between our proposed method and various baselines.



Generated by ScreenCoder (ours)

Figure 12: Qualitative comparison between our proposed method and various baselines.

## Prompt for Grounding

### Seed-1.5-thinking-pro-vision

""Only return the bounding boxes of the 'header', 'footer', 'main content', 'sidebar', and 'navigation' in this webpage screenshot. Please only return the corresponding bounding boxes in the format of 'name: <bbox>x1 y1 x2 y2</bbox>' for each line without extra information and comments!!! Note: 1. All text information, images and other content should be framed, don't miss any information; 2. The areas should not overlap; 3. Output a label and the corresponding bounding box for each line. You can decide whether to include some regions or not.""

### Qwen-2.5-VL

""Only return the bounding boxes of the 'header', 'footer', 'main content', 'sidebar', and 'navigation' in this webpage screenshot. Please only return the corresponding bounding boxes in the format of 'header: <bbox>x1 y1 x2 y2</bbox>\nfooter: <bbox>x1 y1 x2 y2</bbox>\nmain content: <bbox>x1 y1 x2 y2</bbox>\nsidebar: <bbox>x1 y1 x2 y2</bbox>\nnavigation: <bbox>x1 y1 x2 y2</bbox>\n' for each line without extra information and comments!!! Note: 1. All text information, images and other content should be framed, don't miss any information; 2. The areas should not overlap; 3. Output a label and the corresponding bounding box for each line. You can decide whether to include some regions or not. Also, output the image size in the format of "width: <width> height: <height>\n".""

Figure 13: Prompt Templates.

## Prompt for Generation

`{"sidebar": f"\"This is a screenshot of a container. Here is the user's additional instruction: {user_instruction[\"sidebar\"]}\". Please fill in a complete HTML and Tailwind CSS code to accurately reproduce the given container. Please ensure that all block layouts, icon styles, sizes, and text information are consistent with the original screenshot, based on the user's additional conditions. Below is the code template to fill in:`

```
<div>
```

```
your code here
```

```
</div>
```

`Only return the code within the <div> and </div> tags.\"\"\",`

`"header": f"\"This is a screenshot of a container. Here is the user's additional instruction: {user_instruction[\"header\"]}\". Please fill in a complete HTML and Tailwind CSS code to accurately reproduce the given container. Please ensure that all blocks' relative positions, layout, text information, and colors within the bounding box are consistent with the original screenshot, based on the user's additional conditions. Below is the code template to fill in:`

```
<div>
```

```
your code here
```

```
</div>
```

`Only return the code within the <div> and </div> tags.\"\"\",`

`"navigation": f"\"This is a screenshot of a container. Here is the user's additional instruction: {user_instruction[\"navigation\"]}\". Please fill in a complete HTML and Tailwind CSS code to accurately reproduce the given container. Please ensure that all blocks' relative positions, text layout, and colors within the bounding box are consistent with the original screenshot, based on the user's additional conditions. Please use the same icons as in the original screenshot. Below is the code template to fill in:`

```
<div>
```

```
your code here
```

```
</div>
```

`Only return the code within the <div> and </div> tags.\"\"\",`

`"main content": f"\"This is a screenshot of a container. Here is the user's additional instruction: {user_instruction[\"main content\"]}\". Please fill in a complete HTML and Tailwind CSS code to accurately reproduce the given container. Please replace the images in the original screenshot with solid gray blocks of the same size; text inside the images does not need to be recognized. Please ensure that all blocks' relative positions, layout, text information, and colors within the bounding box are consistent with the original screenshot, based on the user's additional conditions. Below is the code template to fill in:`

```
<div>
```

```
your code here
```

```
</div>
```

`Only return the code within the <div> and </div> tags.\"\"\",`

`"footer ": f"\"This is a screenshot of a container. Here is the user's additional instruction: {user_instruction[\"footer\"]}\". Please fill in a complete HTML and Tailwind CSS code to accurately reproduce the given container. Please ensure that all blocks' relative positions, text layout, and colors within the bounding box are consistent with the original screenshot, based on the user's additional conditions. Please use the same icons as in the original screenshot. Below is the code template to fill in:`

```
<div>
```

```
your code here
```

```
</div>
```

`Only return the code within the <div> and </div> tags.\"\"\"}`

`user_instruction: dictionary for storing user-defined instructions for each region ("header", "footer", "sidebar", "navigation", "main content").`

Figure 14: Prompt Templates.



Figure 15: **Failure Case Part 1: Input Design.** The input features a non-standard layout with a skewed container and rotated text.



Figure 16: **Failure Case Part 2: Qwen2.5-VL (Baseline).** The baseline suffers from *spatial hallucination*. It prioritizes pixel matching via fragile absolute positioning, leading to Z-index errors and overlapping text.

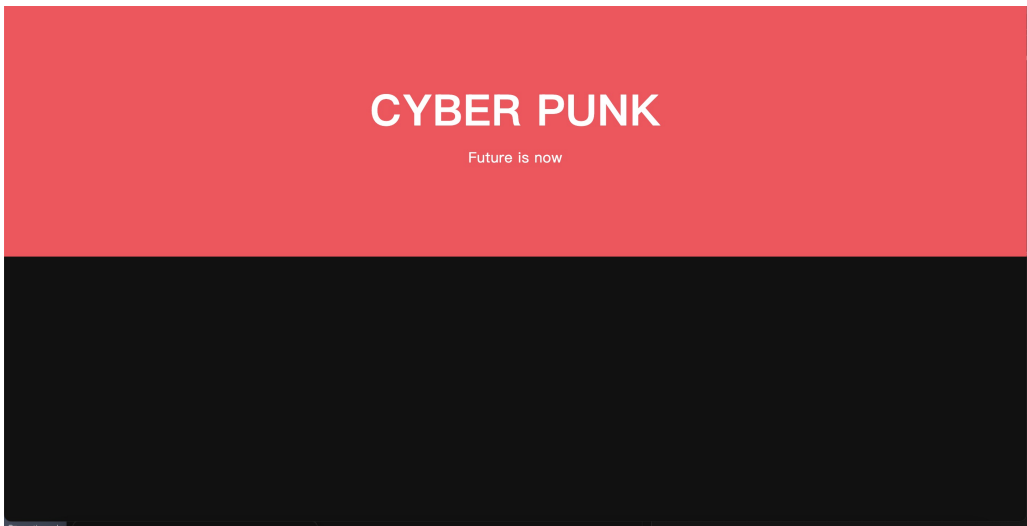


Figure 17: **Failure Case Part 3: ScreenCoder (Ours).** ScreenCoder simplifies the geometry to a standard rectangular grid (Graceful Degradation). While it misses the artistic skew, it maintains structural integrity, generating clean and valid code.