

# INFINITEWEB: Scalable Web Environment Synthesis for GUI Agent Training

Ziyun Zhang<sup>1\*</sup> Zezhou Wang<sup>2\*</sup> Xiaoyi Zhang<sup>3†</sup>  
Zongyu Guo<sup>3</sup> Jiahao Li<sup>3</sup> Bin Li<sup>3</sup> Yan Lu<sup>3</sup>

<sup>1</sup>Peking University <sup>2</sup>Nanjing University  
<sup>3</sup>Microsoft Research Asia

## Abstract

GUI agents that interact with graphical interfaces on behalf of users represent a promising direction for practical AI assistants. However, training such agents is hindered by the scarcity of suitable environments. We present INFINITEWEB, a system that automatically generates functional web environments at scale for GUI agent training. While LLMs perform well on generating a single webpage, building a realistic and functional website with many interconnected pages faces challenges. We address these challenges through unified specification, task-centric test-driven development, and a combination of website seed with reference design image to ensure diversity. Our system also generates verifiable task evaluators enabling dense reward signals for reinforcement learning. Experiments show that INFINITEWEB surpasses commercial coding agents at realistic website construction, and GUI agents trained on our generated environments achieve significant performance improvements on OSWorld, Online-Mind2Web, and MobileWorld, demonstrating the effectiveness of the proposed system. Our code and generated websites are publicly available at <https://github.com/microsoft/FIVE-UI-Evol>.

## 1 Introduction

GUI agents, autonomous systems that interact with graphical user interfaces to complete tasks on behalf of users, have emerged as a promising direction for building practical AI assistants (Xie et al., 2024; Zhou et al., 2024). Recent advances (Hong et al., 2024; Qin et al., 2025) have demonstrated that vision-language models can be end-to-end trained with reinforcement learning algorithms as GUI agents to understand screenshots, reason about UI elements, and execute human-like actions

\*: Equal contribution and work done during the internship at Microsoft Research Asia. †: Project lead.

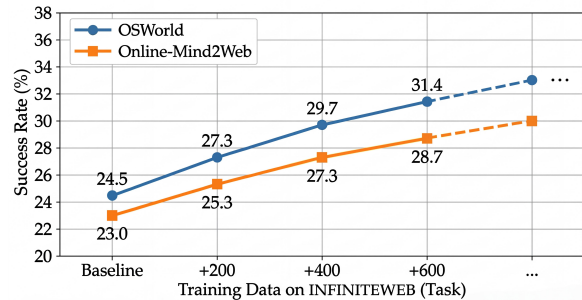


Figure 1: GUI agent performance improves with more training data generated by INFINITEWEB. Dashed lines indicate potential for further scaling.

to automate tasks in digital world. However, training such agents remains challenging due to the scarcity of suitable environments.

Existing GUI agent benchmarks, such as Mini-WoB++ (Liu et al., 2018), WebArena (Zhou et al., 2024), and OSWorld (Xie et al., 2024), provide valuable testbeds but suffer from fundamental limitations in **scale** and **diversity** as training environments. These benchmarks are manually constructed, requiring significant human effort to design websites or download applications, define tasks, and create evaluation criteria. As a result, they contain only tens to hundreds of applications, insufficient for training agents that can generalize across the vast diversity of real-world websites. Although recent work (Sun et al., 2025; Xu et al., 2024; Xie et al., 2025a) proposes synthesizing tasks or trajectories, these approaches still operate within the same benchmark environments, limiting model training on a small set of specific applications.

A natural question arises: *Can we automatically generate environments for GUI agent training?* While large language models (LLMs) have shown remarkable code generation capabilities (Chen et al., 2022; Si et al., 2025; Jimenez et al., 2023), especially for web frontend (Leviathan et al., 2026), directly applying them to generate complete, func-

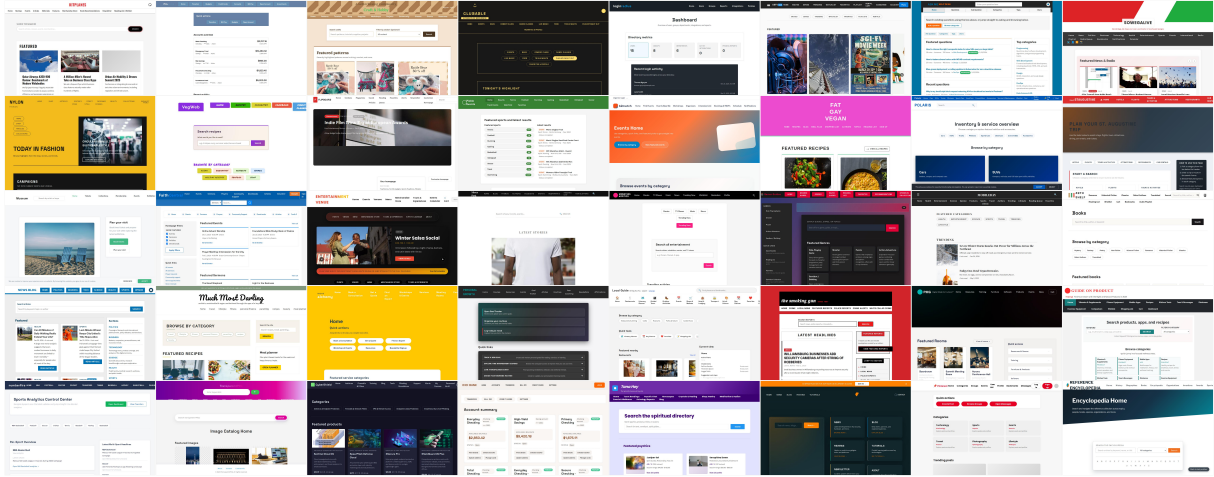


Figure 2: A showcase of diverse websites generated by INFINITEWEB, spanning various categories including e-commerce, booking platforms, social media, and more.

tional websites faces three critical challenges.

Generating such environments presents three intertwined challenges. First, **consistency**: While LLMs perform well on generating a single webpage, a realistic website comprises multiple interconnected pages sharing data, visual styles, and backend interfaces. LLMs generating pages independently often produce incompatible implementations, different backend interface signatures, conflicting data formats, or inconsistent state management, breaking the cross-page interactions essential for realistic websites. Second, **correctness**: website functionalities require multiple coordinated steps, but LLM-generated code frequently contains functional bugs that compound over long-horizon tasks, causing incorrect reward signals that can destabilize reinforcement learning. Third, **diversity**: LLMs tend to produce repetitive task patterns and homogeneous visual styles, risking agent overfitting to specific interaction patterns rather than learning generalizable skills.

In this paper, we present INFINITEWEB, an agentic system that automatically generates functional web environments at scale for GUI agent training, addressing aforementioned challenges.

For **consistency**, we propose **Unified Specification**: rather than generating pages independently, we first derive a complete set of data models and interfaces from user tasks, then generate all pages according to this shared specification, ensuring the realistic cross-page interactions. To ensure **correctness**, inspired by the classic software engineering practice (Williams et al., 2003), we introduce the **task-centric test-driven development (TCTDD)**

approach, where test cases are firstly derived from task specifications and then code is iteratively refined until all task-relevant tests pass. For **diversity**, our system addresses it from both functional and visual dimensions: functionally, by taking a *website seed* (a brief description) and generating tasks specifically designed to match that seed. Visually, by providing reference design images, we use vision-language models to extract characteristics and generate websites that match the target style. It enables the leverage the millions of visually distinct websites available in resources like Common Crawl (Common Crawl Foundation, 2024) as an abundant source of diverse designs.

Furthermore, to support RL-based training, our system is designed to generate **verifiable task evaluators** along with the website and tasks, which tracks key task-related variables during agent running, enabling dense reward signals for reinforcement learning. We conduct systematical analysis on our system from two aspects: generated website quality and the effect to training GUI agent as simulated environment. The results demonstrate the superiority of our system as an environment synthesis system.

We summarize our contributions as follows and we will release the artifacts of this work to further contribute to the research community:

- We propose INFINITEWEB, the first system specifically designed for generating functional web environments with verifiable evaluators for GUI agent training at scale (Figure 2).
- Experiments demonstrate that our system surpasses advanced coding agents in building realis-

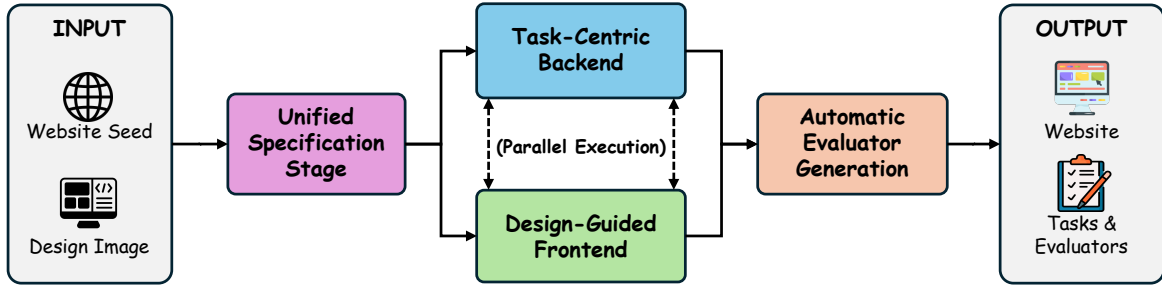


Figure 3: Overview of INFINITEWEB. Given a website seed and design image, our system produces a functional website with tasks and evaluators through four stages: the **Unified Specification Stage** generates tasks and derives data models and interfaces; the **Task-Centric Backend** and **Design-Guided Frontend** execute in parallel; and the **Evaluator Generation** creates task-specific evaluators for dense reward signals.

tic web environments on WebGen-Bench, achieving superior performance in both visual and functional quality.

- Training on our generated environments significantly improves GUI agent performance: +6.9% on OSWorld, +5.7% on Online-Mind2Web, and +3.9% on MobileWorld, demonstrating the realism and quality of simulated environments produced by our system and their cross-platform transferability.

## 2 Related Work

**GUI Agent Benchmarks.** While there are benchmarks evaluating separate ability of GUI Agents like UI element grounding (Li et al., 2025a; Liu et al., 2025) or UI understanding (Wang et al., 2025), end-to-end evaluating GUI agents requires interactive environments. Early work such as MiniWoB++ (Liu et al., 2018) introduced simplified web interaction tasks, demonstrating the potential of reinforcement learning for web automation. Subsequent benchmarks have increased realism and complexity: WebArena (Zhou et al., 2024) provides self-hosted websites for autonomous agent evaluation, OSWorld (Xie et al., 2024) extends to full desktop environments across multiple operating systems, and Mind2Web (Deng et al., 2023) offers large-scale web task annotations. However, these benchmarks share a fundamental limitation: they are manually constructed, requiring significant human effort to design environments, define tasks, and create evaluators. This limits their scale and diversity, potentially leading to agent overfitting. Our work addresses this bottleneck by automatically generating functional web environments at scale.

### LLM-based Code and Website Generation.

Large language models have shown remarkable code generation capabilities, from solving competitive programming problems (Li et al., 2022) to generating complete applications. Recent work has explored UI-to-code generation: Design2Code (Si et al., 2025) benchmarks the conversion of visual designs to front-end code, while WebGen-Bench (Lu et al., 2025) evaluates end-to-end website generation from natural language descriptions. However, a key challenge remains: LLM-generated code frequently contains bugs. CodeT (Chen et al., 2022) addresses this by generating tests alongside code to filter incorrect solutions. Our approach builds on this insight but differs in a crucial way: rather than attempting to verify all generated code, we focus on *task-centric correctness*, ensuring only the functionality required for specific user tasks is bug-free, making the verification problem tractable.

### Synthetic Environment and Data Generation.

Procedural generation has proven valuable for training robust agents. Cobbe et al. (2020) demonstrated that procedurally generated game levels significantly improve reinforcement learning generalization. In the GUI agent domain, recent work has explored synthetic data generation: WebSailorV2 (Li et al., 2025b) uses synthetic trajectories and scalable RL to train web agents, UI-Evol (Zhang et al., 2025) automatically evolves knowledge for computer use agents, while AgentSynth (Xie et al., 2025a) synthesizes long-horizon desktop tasks from atomic subtasks. These approaches focus on generating *training data* (action trajectories) within existing environments. In contrast, our work generates complete, functional *environments* themselves, including websites, tasks, and automatic evaluators, addressing the environment

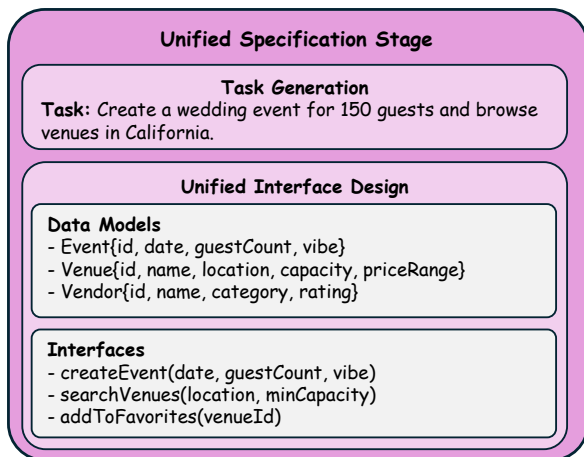


Figure 4: Unified Specification Stage. Given a website seed and design image, this stage generates realistic tasks, then derives shared interface design consisting of data models and programming interfaces across pages.

scalability problem at its source.

### 3 Method

#### 3.1 Overview

Figure 3 illustrates our system pipeline. Our system takes a website seed (e.g., “online bookstore website”) and a design image as input, and outputs a fully functional website along with tasks that can be done in the website and corresponding automatic evaluators. Both website seeds and design images are extracted from Common Crawl to provide diverse visual and functional references (details in Appendix B).

Our pipeline consists of four main stages, with the backend and frontend executing in parallel. First, the **Unified Specification Stage** generates tasks and derives unified data models and interfaces, ensuring **consistency** and **functional diversity**. Second, the **Task-Centric Backend** uses TCTDD to validate business logic, ensuring **correctness** of task-relevant functionality. Third, the **Design-Guided Frontend** extracts visual features from design images to guide page generation, ensuring **visual diversity**. Fourth, **Evaluator Generation** produces task-specific evaluators with **dense reward signals** for reinforcement learning.

#### 3.2 Unified Specification Stage

This stage addresses the **consistency** challenge while enabling **functional diversity**. A functional website typically consists of multiple pages that share data and interfaces. When generating pages independently, LLMs often produce inconsistent

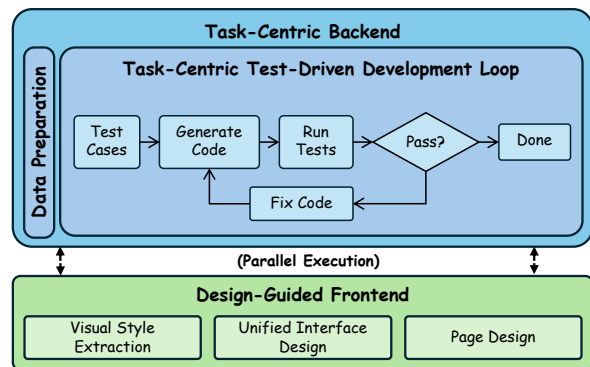


Figure 5: Task-Centric Backend and Design-Guided Frontend in parallel. The backend uses TCTDD to iteratively generate and validate business logic. The frontend extracts visual styles and generates pages.

implementations. Our key insight is that *everything should be derived from tasks*: by first generating tasks specific to the website seed, then deriving unified data models and interfaces from them, we ensure all pages share identical specifications while tasks naturally vary across different website seeds.

**Task Generation.** Given a website seed (e.g., “online bookstore”), we prompt an LLM to generate realistic user tasks specific to that website seed. This ensures **functional diversity**: a booking website generates reservation-related tasks (e.g., “book a hotel room for next weekend”), while an e-commerce site generates shopping tasks (e.g., “find and purchase a laptop under \$500”). Each task represents a complete user goal that varies in complexity and covers different aspects of the website’s functionality.

**Unified Interface Design.** From the generated tasks, we derive three unified specifications that all pages share. First, we extract *data models*: if tasks involve searching products, viewing details, and making purchases, we derive entities such as Product, Cart, and Order with their attributes and relationships. Second, we perform *preliminary architecture planning* to identify all pages required (e.g., homepage, search results, product details, cart, checkout) and their primary functions. Third, we derive a unified set of *programming interfaces*: each task step implies one or more interface calls, and crucially, these interface specifications are shared across all pages, ensuring identical parameters and data formats.

The interfaces are designed to be *user-facing*: the system automatically classifies parameters into system-managed (e.g., userId,

sessionId, managed internally) and user-provided (e.g., productId, quantity). For example, the original interface `addToCart(userId, sessionId, productId, quantity)` is wrapped as `addToCart(productId, quantity)`, with system parameters automatically retrieved from `localStorage`. This unified interface design ensures that all pages use identical API signatures and data formats, enabling seamless cross-page interactions.

With the unified specification stage complete (tasks, data models, interfaces), we now turn to generating business logic and frontend pages in two parallel pipelines.

### 3.3 Task-Centric Backend

This stage addresses the **correctness** challenge. LLM-generated code frequently contains bugs, making naively synthesized environments unsuitable for agent learning. Our key insight is to adopt *task-centric correctness* as the correctness criterion. Since agents interact only with a narrow, task-induced subspace determined by task specifications and their policies, correctness outside this subspace does not contribute to the learning signal or policy optimization. Rather than enforcing full functional correctness over the entire website, we focus on ensuring that only the functionalities required for the target tasks are correct. This alignment allows correctness verification and refinement to be focused on task-relevant execution paths, which we operationalize through TCTDD.

**Data Preparation.** We generate concrete data instances that populate the website, ensuring consistency with both data models and tasks. For example, if a task requires finding products under \$50, we ensure the generated product catalog contains such items. Placeholder resources (e.g., image URLs) are replaced with real, context-relevant content via external APIs.

**Task-Centric Test-Driven Development.** We adopt the TCTDD approach to ensure correctness of task-relevant functionality. TCTDD works as follows: based on task specifications and generated data, test cases and implementation code are generated in parallel; then tests are run and iteratively fixed until all pass.

Test cases and implementation code use the same pre-generated data, ensuring consistency. For example, if the generated data contains products priced at \$29.99 and \$45.00, the test will verify

that exactly these products are returned for an “under \$50” query. When tests fail, we provide the LLM with the failing test case, expected vs. actual output, and relevant code segment. The LLM generates a fix and re-tests, continuing until all tests pass or a maximum iteration limit is reached.

### 3.4 Design-Guided Frontend

This stage is designed to address the **visual diversity** challenge. LLMs tend to generate websites with similar visual styles. Our key insight is to draw *design images* as references from abundant and visually diverse website screenshot in the real world. Given a reference design image, we extract visual characteristics and generate pages that match the guidance.

**Visual Style Extraction.** We employ a vision-language model to decode visual attributes from the design image, establishing a global style constraint. Specifically, we extract the **color system** (e.g., primary and neutral palettes), **typography hierarchy** including font families and weights, **spacing rules**, and **component styling** such as button patterns. These extracted structural style specifications serve as a consistent visual specification for all subsequent generation steps.

**Page Design.** Building on the unified specification, we conduct detailed architectural design for each page in parallel. This step defines specific **functional requirements** including content blocks and interaction flows, determines **routing logic** via URL parameters, and establishes **responsive layouts** defined by grid systems and breakpoints.

**Page Realization.** We first generate a unified page framework containing the shared header, footer, and CSS variables for the website, based on the extracted visual features, ensuring consistent styling across all pages. Then, for each page, we generate the HTML structure, CSS styles, and a JavaScript UI layer that connects elements to the backend SDK using a data-attribute-driven pattern (e.g., `data-populate`, `data-action`). Finally, we inject an initialization script into the homepage that writes the generated data to `localStorage`, which is the browser’s built-in persistent key-value storage. It enables data persistence across pages without a backend server.

Method	Instruction Categories			Test Case Categories			Overall
	Content Pres.	User Inter.	Data Mgmt.	Functional	Data Display	Design Valid.	
Bolt.diy	83.2	59.9	62.8	58.4	84.1	41.7	67.0
Claude-Code	87.9	70.1	67.6	67.3	87.5	61.1	74.3
Codex	89.8	79.2	75.6	72.8	<b>96.2</b>	76.4	81.2
<b>Ours</b>	<b>91.5</b>	<b>83.8</b>	<b>82.7</b>	<b>80.9</b>	94.1	<b>82.8</b>	<b>85.6</b>

Table 1: Category-wise evaluation results on WebGen-Bench (%). Instruction Categories classify the website functionality type. Test Case Categories classify the evaluation type. Results are averaged over three runs.

Method	Chrome	GIMP	Calc	Impress	Writer	Multi	OS	Thunder.	VLC	VSCode	Overall
Computer-use-preview	36.9	34.6	10.6	25.4	30.4	10.8	45.8	46.7	29.4	47.8	26.0
Claude-3.7-Sonnet	41.2	34.6	8.5	29.7	39.1	10.8	50.0	33.3	35.3	43.5	27.1
Doubao-1.5-thinking-0428	44.4	46.2	13.0	31.8	39.1	4.8	30.4	66.7	23.5	56.5	28.1
Claude-4-Sonnet	36.9	46.2	17.0	36.2	43.5	9.7	37.5	66.7	38.5	60.9	31.2
OpenCUA-32B	40.5	55.1	13.5	30.7	39.1	9.7	52.2	44.4	25.0	52.8	29.7±1.1
UI-TARS-1.5-7B	22.9	51.9	11.7	29.7	39.1	3.8	34.8	26.7	34.2	63.0	24.5±1.2
+ 200 tasks	34.8	61.5	10.0	27.7	34.8	8.0	41.7	40.0	25.5	55.1	27.3±0.9
+ 400 tasks	35.5	69.2	10.6	29.8	37.7	9.0	48.6	35.5	33.3	62.3	29.7±1.1
+ 600 tasks	36.9	69.2	12.8	29.8	47.8	9.7	45.8	40.0	35.3	65.2	<b>31.4±1.0</b>

Table 2: Results on OSWorld under 15 maximum steps by domain (%). The lower section shows UI-TARS-1.5-7B trained with tasks from INFINITEWEB-generated websites. Calc/Impress/Writer refer to LibreOffice applications. Multi = Multi-Apps, Thunder. = Thunderbird. Standard deviation computed over three runs.

### 3.5 Automatic Evaluator Generation

A critical requirement for GUI agent training is automatically evaluating whether a task has been successfully completed. Our system automatically generates task-specific evaluators by leveraging existing state variables and code instrumentation.

The evaluators leverage two types of variables: *existing variables* representing state naturally stored by the application (e.g., cart contents, user preferences), and *instrumentation variables* that are explicitly added checkpoints tracking task-specific progress. For instrumentation variables, we identify the key steps required for each task’s completion and record progress in `localStorage` when the corresponding functions execute. For example, a “search and purchase” task might track: search query submitted, product viewed, item added to cart, and checkout completed.

Based on these variables, we generate a JavaScript evaluator function that checks variables to determine task completion, capable of assessing partial completion rather than only binary success/failure. This enables *dense reward signals* for reinforcement learning: agents receive partial credit based on completed steps, facilitating more effective learning for complex multi-step tasks.

## 4 Experiments

We evaluate INFINITEWEB on three dimensions: (1) functional correctness of generated websites, (2) visual quality, and (3) effectiveness for GUI agent training.

### 4.1 Experimental Setup

We evaluate on four benchmarks and assess visual quality through pairwise comparisons. For fair comparison, all website generation methods use GPT-5 as the backbone LLM with reasoning effort set to “high”. Implementation details including generation hyperparameters and agent training configuration are provided in Appendix B. We also provide manual evaluation detailed in Appendix E.

**WebGen-Bench.** WebGen-Bench (Lu et al., 2025) evaluates functional correctness of LLM-generated websites through agent-based task execution on 101 websites. Each website is generated with a set of predefined tasks that it must support. For evaluation, an LLM agent is presented with a user task and attempts to complete it by interacting with the generated website. Each task outcome is classified as: *Passed* if the task is fully completed with correct results, *Partial* if the agent makes progress but does not complete the task entirely, or *Failed* if the task cannot be accomplished

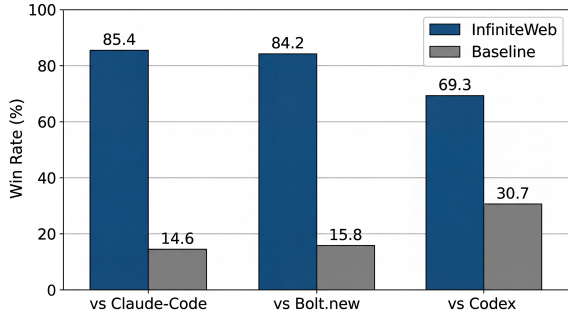


Figure 6: LLM-as-Judge visual quality evaluation. Each pair shows win rates for ours (left) vs baseline (right).

due to missing functionality or errors. We report three metrics: Passed rate, Partial rate, and the Overall score. Since the original WebGen-Bench does not include design images, we match each test website with a design image extracted from Common Crawl based on website category. This design image is provided to all methods as input to enable fair comparison.

**LLM-as-Judge Visual Quality.** We assess visual quality through LLM-as-Judge pairwise comparisons (Zheng et al., 2023) on 200 generated websites. For each website, we capture a full-page screenshot and present it alongside the reference design image to GPT-5. The model is prompted to evaluate which implementation better matches the target design across five dimensions: (1) visual layout similarity, (2) color scheme matching, (3) typography and spacing, (4) component arrangement and structure, and (5) overall aesthetic consistency. The model outputs one of three judgments: our method wins, the baseline wins, or tie. We report win rates for each pairwise comparison, where higher percentages indicate stronger visual fidelity to the design reference.

**Online-Mind2Web.** Online-Mind2Web (Xue et al., 2025) extends the original Mind2Web benchmark (Deng et al., 2023) to evaluate web agents on live websites, testing their ability to complete realistic tasks on real-world web pages. Unlike static benchmarks with cached HTML snapshots, Online-Mind2Web requires agents to interact with actual deployed websites, introducing challenges such as dynamic content loading, varying page layouts, and real network latency. We use this benchmark to measure **in-domain generalization**: whether training on our synthetic websites improves performance on real-world web interactions that the agent has never seen during

training.

**OSWorld.** OSWorld (Xie et al., 2024) is a benchmark for evaluating GUI agents on real desktop applications across diverse domains including web browsers, office suites (Calc, Impress, Writer), media players (VLC), code editors (VS Code), and email clients (Thunderbird). We use this benchmark to measure **out-of-domain transfer**: whether training on synthetic web environments transfers to real desktop application tasks. Specifically, we adopt **OSWorld-Verified** (Xie et al., 2025b), a refined version with improved task quality and evaluation robustness.

**MobileWorld.** MobileWorld (Kong et al., 2025) benchmarks autonomous mobile agents in interactive environments. We use this benchmark to measure **mobile transfer**: whether training on our generated web environments, which use responsive layout design, transfers to mobile GUI tasks without additional fine-tuning.

## 4.2 Website Functional Correctness

We compare against three representative approaches for AI-powered website generation: **Codex** (v0.46.0) (Chen et al., 2021), OpenAI’s coding assistant agent; **Claude-Code** (v2.0.0) (Anthropic, 2025), Anthropic’s coding assistant agent; and **Bolt.diy** (v0.0.7) (StackBlitz Labs, 2024), an open-source AI website builder from StackBlitz. All methods are given the same website seed and a homepage design image. The prompt template for baselines is provided in Appendix F.

Table 1 presents the functional correctness results on WebGen-Bench. Our method achieves the highest overall score of 85.6%, significantly outperforming all baselines. We report performance across two classification schemes: *Instruction Categories* (Content Presentation, User Interaction, Data Management) that classify the type of website functionality being tested, and *Test Case Categories* (Functional Testing, Data Display Testing, Design Validation Testing) that classify the type of evaluation being performed. Our method achieves the best performance in Functional Testing (80.9%) and Design Validation (82.8%), demonstrating particularly strong advantages on the most challenging task categories. Detailed results with statistical significance tests are provided in Appendix C.1.

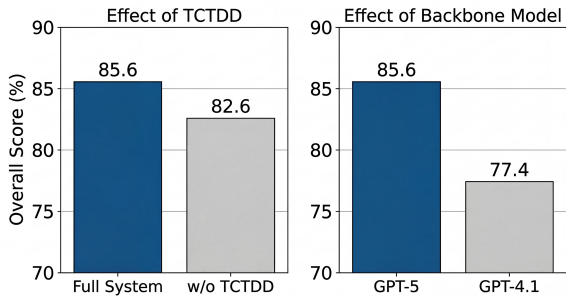


Figure 7: Ablation study results on WebGen-Bench. Left: Effect of TCTDD validation loop. Right: Effect of backbone model.

### 4.3 LLM-as-Judge Visual Quality

We compare the same websites generated in Section 4.2 (ours vs. Codex, Claude-Code, and Bolt.diy). Figure 6 shows the pairwise comparison results. Our method consistently outperforms all baselines (69–85% win rate). Human evaluation confirms 91% agreement with automated assessments (Appendix D).

### 4.4 Effectiveness for Agent Training

The ultimate goal of INFINITEWEB is to provide training environments for GUI agents. We generate 600 tasks spanning diverse website categories (e-commerce, social media, booking platforms, etc.) and use them to train UI-TARS-1.5-7B (Qin et al., 2025). Training uses GRPO (Group Relative Policy Optimization) (Shao et al., 2024) with our dense reward signals from code instrumentation, enabling the agent to receive partial credit for intermediate progress rather than binary success/failure. We then evaluate on Online-Mind2Web (in-domain), OSWorld (out-of-domain), and MobileWorld (mobile transfer) benchmarks.

As shown in Figure 1, training on our generated environments leads to substantial improvements: +6.9% on OSWorld (24.5%  $\rightarrow$  31.4%), +5.7% on Online-Mind2Web (23.0%  $\rightarrow$  28.7%), and +3.9% on MobileWorld (6.4%  $\rightarrow$  10.3%). Table 2 shows the per-domain breakdown on OSWorld, where improvements are observed across most application categories. This suggests that skills acquired from training on web environments can transfer beyond web tasks to both desktop and mobile applications. Appendix A provides case studies analyzing this transfer. We also compare against training on Claude-Code-generated websites using the same setup; results in Appendix C.4 show that data qual-

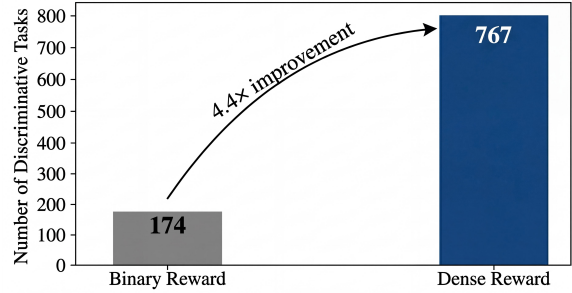


Figure 8: Number of discriminative tasks for GRPO training. Dense reward enables learning from 4.4 $\times$  more tasks by providing partial credit for intermediate steps.

ity, not quantity, drives the improvement.

The improvement scales with the amount of training data, suggesting that generating more diverse environments could yield further gains.

### 4.5 Generated Environment Quality

To evaluate the quality of our generated environments, we compare the success rate and average successful steps of two agents on InfiniteWeb and OSWorld: UI-TARS-1.5-7B and Agent S2 (Agashe et al., 2025), a multi-agent system using GPT-4.1 as planner and UI-TARS-72B for grounding. Table 3 shows the results.

Env.	Agent	Score	Steps
InfiniteWeb	UI-TARS-1.5-7B	7.4	10.3
	Agent S2	14.1	10.9
OSWorld	UI-TARS-1.5-7B	24.5	9.0
	Agent S2	27.3	9.3

Table 3: Agent performance on InfiniteWeb and OSWorld. Score is the average task completion rate (%). Steps is the average steps for successful tasks.

**Higher Difficulty.** Compared to OSWorld, InfiniteWeb is markedly more challenging: agents achieve 2–3 $\times$  lower scores, and successful tasks require longer trajectories, suggesting increased task complexity.

**Better Discriminability.** Performance on InfiniteWeb is more sensitive to agent capability, resulting in a 6.7 percentage point gap between Agent S2 and UI-TARS, compared to 2.8 on OSWorld.

### 4.6 Ablation Studies

Having established the effectiveness of our full system, we now examine the contribution of individual components through ablation studies. Figure 7 shows the results.

**Effect of TCTDD.** Removing the TCTDD validation loop reduces the overall score by 5.0 points. This confirms that iterative test-driven refinement is crucial for achieving high functional correctness, even when using a strong backbone model. Notably, even without TCTDD, our method still achieves 80.6%, comparable to Codex, showing that our base architecture is itself competitive.

**Effect of Backbone Model.** Replacing GPT-5 with GPT-4.1 reduces the score by 8.2 points (85.6  $\rightarrow$  77.4). Even with GPT-4.1, our method still outperforms Claude-Code using GPT-5 (75.8%), showing that our approach remains competitive even with a weaker backbone model.

**Effect of Dense Reward.** Our instrumentation system enables dense reward signals by tracking intermediate task steps. To evaluate its impact on reinforcement learning, we run UI-TARS-1.5-7B on 4,000 generated tasks with 4 trajectories per task and compare the number of discriminative tasks where GRPO can effectively learn, i.e., tasks where at least one trajectory in a group receives different scores. As shown in Figure 8, dense reward enables learning from 767 tasks compared to 174 with binary reward, a 4.4 $\times$  increase. This demonstrates that dense reward substantially expands the effective training signal by providing partial credit for intermediate progress, thereby improving training data efficiency.

#### 4.7 Generation Efficiency

We analyze the computational cost of website generation. On average, generating a single website consumes approximately 0.36M input tokens and 0.34M output tokens. Using GPT-5 batch processing pricing (\$0.625/M input, \$5.00/M output), this translates to approximately \$1.93 per website. The median generation time is approximately 20 minutes per website with our API configuration, though this is highly dependent on API response speed and rate limits. Since each website is generated independently, multiple websites can be generated in parallel to increase throughput.

## 5 Conclusion

We presented INFINITEWEB, a system that aims to generate functional web environments for GUI agent training, addressing consistency through unified interface design, correctness through task-centric test-driven development, and diversity

through website seed variation and design image guidance. Our system surpasses commercial coding agents in this scenario and experiment results demonstrate its advantages for training GUI agents. By releasing our system and generated datasets, we hope to support future research in building more capable and generalizable GUI agents.

## Limitations

Our work has several limitations that suggest directions for future research.

**Single-Website Scope.** Our current tasks operate within individual websites. Cross-website tasks, such as comparing prices across multiple shopping sites or aggregating information from different sources, represent an important real-world scenario not yet covered. While we observe improvement on multi-application tasks in OSWorld (3.8%  $\rightarrow$  9.7%), suggesting that single-website training already enables some cross-application transfer, generating coordinated multi-website environments with cross-site tasks remains a valuable direction for future work.

**Mobile Evaluation.** While our generated websites use responsive layout design and we demonstrate positive transfer to MobileWorld (+3.9%), our current system only generates web environments. Generating native mobile application environments for mobile GUI agent training is a direction for future work.

**Generation Cost.** Generating a complete website environment requires multi-stage LLM calls, including task generation, architecture design, code generation, and test validation. The median generation time is approximately 20 minutes per website, at a cost of \$1.93 per site using GPT-5 batch pricing. While this cost is orders of magnitude lower than manual benchmark construction (e.g., OSWorld required approximately 1,800 person-hours for 369 tasks) and generation accounts for only 6.9% of total training time due to embarrassingly parallel execution (600 websites in approximately 125 minutes), further reducing per-website cost and latency remains an engineering direction.

## References

Saaket Agashe, Kyle Wong, Vincent Tu, Jiachen Yang, Ang Li, and Xin Eric Wang. 2025. [Agent s2: A compositional generalist-specialist framework for computer use agents](#). *Preprint*, arXiv:2504.00906.

- Anthropic. 2025. Claude code. <https://claude.ai/code>.
- Bei Chen, Fengji Zhang, Anh Nguyen, Daoguang Zan, Zeqi Lin, Jian-Guang Lou, and Weizhu Chen. 2022. *Codet: Code generation with generated tests*. Preprint, arXiv:2207.10397.
- 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. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. 2020. Leveraging procedural generation to benchmark reinforcement learning. In *International Conference on Machine Learning (ICML)*.
- Common Crawl Foundation. 2024. Common crawl. <https://commoncrawl.org>.
- Xiang Deng, Yu Gu, Boyuan Zheng, Shijie Chen, Sam Stevens, Boshi Wang, Huan Sun, and Yu Su. 2023. Mind2web: Towards a generalist agent for the web. *Advances in Neural Information Processing Systems (NeurIPS)*.
- Wenyi Hong, Weihang Wang, Qingsong Lv, and 1 others. 2024. Cogagent: A visual language model for gui agents. In *CVPR*.
- Carlos E Jimenez, John Yang, Alexander Wettig, Shunyu Yao, Kexin Pei, Ofir Press, and Karthik Narasimhan. 2023. Swe-bench: Can language models resolve real-world github issues? *arXiv preprint arXiv:2310.06770*.
- Qiyu Kong, Xu Zhang, Zhenyu Yang, Nolan Gao, Chen Liu, Panrong Tong, Chenglin Cai, Hanzhang Zhou, Jianan Zhang, Liangyu Chen, Zhidan Liu, Steven Hoi, and Yue Wang. 2025. *Mobileworld: Benchmarking autonomous mobile agents in agent-user interactive, and mcp-augmented environments*. Preprint, arXiv:2512.19432.
- Yaniv Leviathan, Dani Valevski, Matan Kalman, Danny Lumen, Eyal Segalis, Eyal Molad, Shlomi Pasternak, Vishnu Natchu, Valerie Nygaard, Srinivasan Venkatachary, James Manyika, and Yossi Matias. 2026. *Generative ui: Llms are effective ui generators*. Preprint, arXiv:2604.09577.
- Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo, Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng Chua. 2025a. Screenspot-pro: Gui grounding for professional high-resolution computer use. In *Proceedings of the 33rd ACM International Conference on Multimedia*, pages 8778–8786.
- Kuan Li, Zhongwang Zhang, Huifeng Yin, Rui Ye, Yida Zhao, Liwen Zhang, Litu Ou, Dingchu Zhang, Xixi Wu, Jialong Wu, and 1 others. 2025b. Websailor-v2: Bridging the chasm to proprietary agents via synthetic data and scalable reinforcement learning. *arXiv preprint arXiv:2509.13305*.
- Yujia Li, David Choi, Junyoung Chung, and 1 others. 2022. Competition-level code generation with alpha-code. *Science*, 378:1092–1097.
- Evan Zheran Liu, Kelvin Guu, Panupong Pasupat, Tianlin Shi, and Percy Liang. 2018. Reinforcement learning on web interfaces using workflow-guided exploration. In *International Conference on Learning Representations (ICLR)*.
- Xinyi Liu, Xiaoyi Zhang, Ziyun Zhang, and Yan Lu. 2025. Ui-e2i-synth: Advancing gui grounding with large-scale instruction synthesis. *arXiv preprint arXiv:2504.11257*.
- Zimu Lu, Yifan Yang, Haoran Ren, Haocheng Hou, Hangyu Xiao, Ke Wang, Wenji Shi, Aojun Zhou, Minghao Zhan, and Hongsheng Li. 2025. Webgenbench: Evaluating llms on generating interactive and functional websites from scratch. *arXiv preprint arXiv:2505.03733*.
- Yujia Qin, Yining Ye, Junjie Fang, Haoming Wang, Shihao Liang, Shizuo Tian, Junda Zhang, Jiahao Li, Yunxin Li, Shijue Huang, and 1 others. 2025. Uitars: Pioneering automated gui interaction with native agents. *arXiv preprint arXiv:2501.12326*.
- 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. *Deepseekmath: Pushing the limits of mathematical reasoning in open language models*. Preprint, arXiv:2402.03300.
- Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang, Ruibo Liu, and Diyi Yang. 2025. *Design2code: Benchmarking multimodal code generation for automated front-end engineering*. Preprint, arXiv:2403.03163.
- StackBlitz Labs. 2024. *Bolt.diy*. Accessed: 2025-04-22.
- Qiushi Sun, Kanzhi Cheng, Zichen Ding, Chuanyang Jin, Yian Wang, Fangzhi Xu, Zhenyu Wu, Chengyou Jia, Liheng Chen, Zhoumianze Liu, and 1 others. 2025. Os-genesis: Automating gui agent trajectory construction via reverse task synthesis. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5555–5579.
- Xuehui Wang, Zhenyu Wu, JingJing Xie, Zichen Ding, Bowen Yang, Zehao Li, Zhaoyang Liu, Qingyun Li, Xuan Dong, Zhe Chen, and 1 others. 2025. Mmbench-gui: Hierarchical multi-platform evaluation framework for gui agents. *arXiv preprint arXiv:2507.19478*.
- Laurie Williams, E. Michael Maximilien, and Mladen Vouk. 2003. Test-driven development as a defect-reduction practice. In *14th International Symposium on Software Reliability Engineering (ISSRE)*.

- Jingxu Xie, Dylan Xu, Xuandong Zhao, and Dawn Song. 2025a. Agentsynth: Scalable task generation for generalist computer-use agents. *arXiv preprint arXiv:2506.14205*.
- Tianbao Xie, Mengqi Yuan, Danyang Zhang, Xinzhuang Xiong, Zhennan Shen, Zilong Zhou, Xinyuan Wang, Yanxu Chen, Jiaqi Deng, Junda Chen, Bowen Wang, Haoyuan Wu, Jixuan Chen, Junli Wang, Dunjie Lu, Hao Hu, and Tao Yu. 2025b. [Introducing osworld-verified](#). *xlang.ai*.
- Tianbao Xie, Danyang Zhang, Jixuan Chen, Xiaochuan Li, Siheng Zhao, Ruisheng Cao, Toh Jing Hua, Zhoujun Cheng, Dongchan Shin, Fangyu Lei, and 1 others. 2024. Osworld: Benchmarking multimodal agents for open-ended tasks in real computer environments. *arXiv preprint arXiv:2404.07972*.
- Yiheng Xu, Dunjie Lu, Zhennan Shen, Junli Wang, Zekun Wang, Yuchen Mao, Caiming Xiong, and Tao Yu. 2024. Agenttrek: Agent trajectory synthesis via guiding replay with web tutorials. *arXiv preprint arXiv:2412.09605*.
- Tianci Xue, Weijian Qi, Tianneng Shi, Chan Hee Song, Boyu Gou, Dawn Song, Huan Sun, and Yu Su. 2025. [An illusion of progress? assessing the current state of web agents](#). *Preprint*, arXiv:2504.01382.
- Ziyun Zhang, Xinyi Liu, Xiaoyi Zhang, Jun Wang, Gang Chen, and Yan Lu. 2025. [Ui-evol: Automatic knowledge evolving for computer use agents](#). *Preprint*, arXiv:2505.21964.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, and 1 others. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. In *NeurIPS*.
- Shuyan Zhou, Frank F. Xu, Hao Zhu, Xuhui Zhou, Robert Lo, Abishek Sridhar, Xianyi Cheng, Yonatan Bisk, Daniel Fried, Uri Alon, and 1 others. 2024. Webarena: A realistic web environment for building and evaluating autonomous agents. In *International Conference on Learning Representations (ICLR)*.

## A Case Studies and Analysis

### A.1 Cross-Domain Transfer Analysis

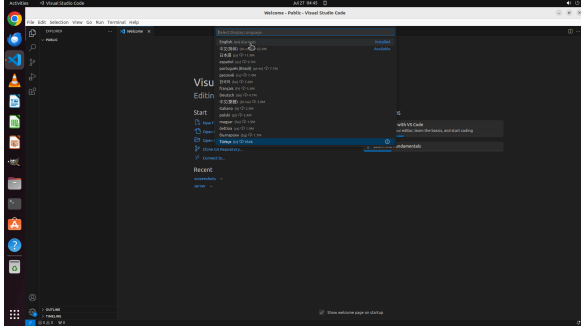
To understand why website training improves performance across all OSWorld domains, we analyzed execution traces from the multi-run experiments. To eliminate cases attributable to random variation, we applied strict filtering and focused on “strong positive transfer” cases, where the baseline failed consistently across all repeated runs while the trained model succeeded consistently. Analyzing the baseline failure patterns revealed three universal GUI interaction capabilities that website training develops:

**Exploration Persistence.** The trained model persists in exploring alternatives when initial attempts fail, rather than prematurely giving up. In one VS Code task requiring language change to Arabic, the baseline browsed the language list with PageDown, concluded “Arabic is not in the visible range,” and terminated after only 5 steps. The trained model continued for 15 steps, trying multiple approaches (typing “Arabic”, scrolling, clearing and retrying) until successfully locating and selecting the option. Figure 9 shows the comparison.

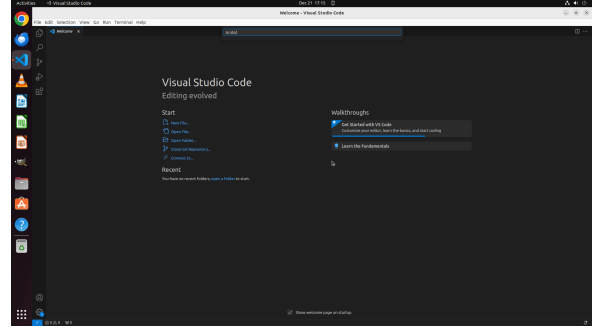
**Flow Completeness.** The trained model executes complete task workflows instead of stopping part-way. For a Spotify installation task, the baseline opened Ubuntu Software Center, searched for “Spotify,” and called done after 4 steps, without clicking Install. The trained model completed the full 13-step flow: search, click Install, enter password for authentication, wait for installation progress, and verify completion. Figure 10 illustrates this difference.

**Loop Avoidance.** The trained model avoids getting stuck in repetitive action cycles. In an email attachment task in Thunderbird, the baseline successfully attached a file in 4 steps but then became confused about task completion, unsure what to do next after adding the attachment, it entered a futile loop of repeatedly opening the file picker and canceling for the remaining 11 steps. The trained model completed the same task cleanly in 5 steps. Figure 11 demonstrates this pattern.

These capabilities are domain-agnostic: avoiding loops, completing workflows, and persisting through obstacles apply equally to image editors, office suites, and system utilities. Website environments, with their diverse interaction patterns



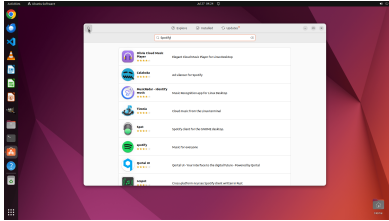
(a) Baseline (Step 4): Language list visible, no Arabic found



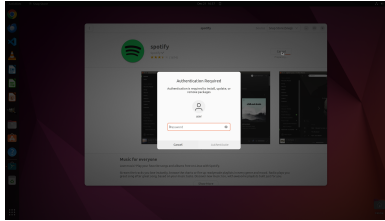
(b) Trained (Step 5): Types "Arabic" to search

Step	Baseline	Trained
1-3	Open command palette, search "Configure Display Language"	Same as baseline
4	PageDown to browse list	Type "Arabic" in search
5	<b>Give up:</b> "Arabic not visible"	Scroll, clear, retry
6-15	-	Continue exploring alternatives

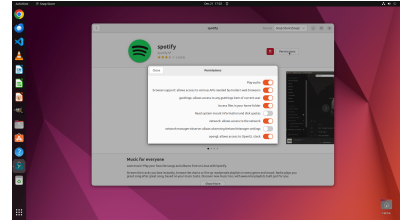
Figure 9: VS Code language change task: Exploration Persistence.



(a) Baseline Step 4: Search results, task ends here



(b) Trained Step 5: Auth dialog after Install



(c) Trained Step 12: Configuring permissions

Step	Baseline	Trained
1-3	Open Software Center, search "Spotify"	Same as baseline
4	<b>Done:</b> Search results displayed	Click Spotify app
5	-	Click Install button
6-7	-	Enter password, authenticate
8-13	-	Wait for installation, configure permissions

Figure 10: Spotify installation task: Flow Completeness. Baseline stops at search results; trained completes full installation.

and multi-step transactions, effectively train these transferable behaviors.

## A.2 Automatic Evaluator Generation

InfiniteWeb automatically generates dense reward evaluators that provide proportional rewards for partial task completion. Figure 12 shows an evaluator for the task "Subscribe to newsletter with weekly specials" from a restaurant website.

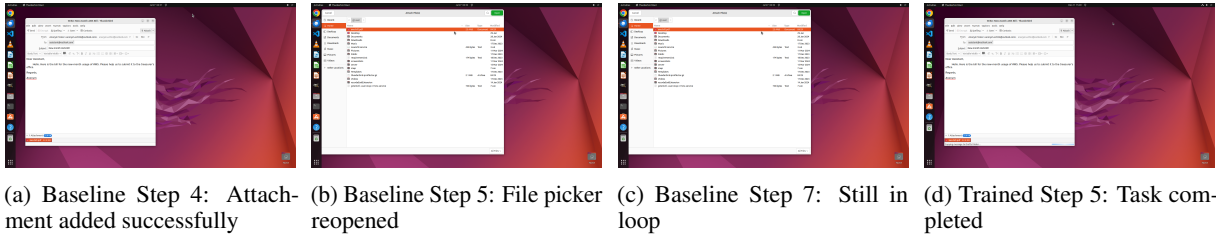
The evaluator uses weighted checkpoints that enable dense reward signals for GRPO training. Each checkpoint validates a different aspect of task completion: (1) *instrumentation flags* that track whether the agent performed required actions, (2) *data consistency* that verifies records were properly

created, and (3) *confirmation state* that ensures the full workflow completed. The weighted sum allows partial credit, an agent that initiates but fails to complete a task still receives proportional reward. This design prevents shortcuts (directly manipulating localStorage fails instrumentation checks) while providing richer training signals than sparse 0/1 rewards.

## A.3 TCTDD Validation and Auto-Fix

The TCTDD validation loop automatically detects and fixes implementation errors. Table 4 shows an example from a B2B industrial equipment website where one test initially failed.

This iterative process ensures that the generated



(a) Baseline Step 4: Attachment added successfully (b) Baseline Step 5: File picker reopened (c) Baseline Step 7: Still in loop (d) Trained Step 5: Task completed

Step	Baseline	Trained
1–4	Attach file successfully (a–c)	Same as baseline
5	<b>Loop starts:</b> Reopen file picker (b)	<b>Done:</b> Save to drafts (d)
6–15	Repeat open/cancel 11 times (c shows step 7)	–

Figure 11: Thunderbird email attachment task: Loop Avoidance. After successfully attaching the file (a), baseline becomes confused and enters a futile loop (b, c), while trained model completes cleanly (d).

Iteration 1: Test Execution	
✓	Task 1: Request quote for three forklifts
✓	Task 2: Add two electric forklifts to cart
×	Task 3: Schedule demo for product <i>Error: Demo request should be submitted</i>
✓	Task 4–10: (passed)
<b>Result: 9/10 passed, 1 failed</b>	
Auto-Fix: LLM Analysis and Repair	
The LLM identifies that <code>submitDemoRequest()</code> returns undefined instead of a success object, and generates a corrected implementation.	
Iteration 2: Re-validation	
✓	Task 1–10: All tests passed
<b>Result: 10/10 passed</b>	

Table 4: TCTDD validation loop example. The system detects a failing test, uses an LLM to analyze and fix the implementation, then re-validates until all tests pass.

business logic correctly implements all required functionality. In our experiments, most websites require 1–3 iterations to pass all tests, with a maximum of 8 iterations allowed.

## B Data Collection and Implementation

### Website Seed and Design Image Extraction.

We sample web pages from Common Crawl. For each sampled page, we render it in a headless browser and capture a full-page screenshot as the design image. We then use an LLM to analyze the visual content of the screenshot, generating a concise natural language description as the website seed, while filtering out pages that violate robots.txt or contain illegal content.

**Generation Hyperparameters.** We use the following configuration: temperature 0.7, maximum output tokens 32,000, task count range 8–10 per

website, maximum 12 pages per website, and maximum 8 iterations for TCTDD validation loop.

**Agent Training.** We post-train UI-TARS-1.5-7B using GRPO. The training configuration includes: learning rate  $1e-6$ , AdamW optimizer with bf16 precision, gradient clipping at 1.0, global batch size 16, PPO epochs 1, clip ratio 0.2–0.3, and discount factor  $\gamma = 0.95$ . For rollout, we use 128 parallel environments, sample 8 trajectories per task, set maximum 15 steps per episode, and use temperature 1.0 for sampling.

## C Experimental Details and Results

**Baseline Implementation.** For the direct prompting baselines, we use GPT-5 with high reasoning effort as the backbone model. The prompt specifies website seed, required functionality, technical requirements (up to 12 pages, localStorage, reference design image), and code standards. The full prompt template is provided in Appendix F.

### C.1 WebGen-Bench Results

Table 5 presents the detailed results on WebGen-Bench across three independent runs. Welch’s t-tests confirm that InfiniteWeb significantly outperforms all baselines: vs Bolt.diy ( $t=14.81$ ,  $p<0.001$ ), vs Claude-Code ( $t=6.33$ ,  $p<0.01$ ), and vs Codex ( $t=6.57$ ,  $p<0.05$ ).

Table 6 shows the ablation study results. Both ablations show significant degradation: using GPT-4.1 instead of GPT-5 ( $t=7.70$ ,  $p<0.01$ ) and removing TCTDD validation ( $t=2.82$ ,  $p<0.05$ ).

### C.2 Online-Mind2Web Results

Table 7 presents the results on Online-Mind2Web across three independent runs, broken down by task

### Generated Dense Reward Evaluator Example

```

const checkpoints = [];

// Load instrumentation and domain data
const attempt = JSON.parse(localStorage.getItem('task10_subscriptionAttempt'));
const confirmed = JSON.parse(localStorage.getItem('task10_subscriptionConfirmedOnSite'));
const subs = JSON.parse(localStorage.getItem('newslettersubscriptions') || '[]');

// CP1 (0.35): Subscription attempt with valid email and weekly specials enabled
const cp1 = attempt && isTestEmail(attempt.email) && attempt.wantsWeeklySpecials;
checkpoints.push({ passed: cp1, weight: 0.35 });

// CP2 (0.30): Subscription record matches attempt in storage
const subRecord = subs.find(s => s.id === attempt?.newsletterSubscriptionId);
const cp2 = subRecord && subRecord.email === attempt.email;
checkpoints.push({ passed: cp2, weight: 0.30 });

// CP3 (0.35): On-site confirmation references the subscription
const cp3 = confirmed?.newsletterSubscriptionId === attempt?.newsletterSubscriptionId;
checkpoints.push({ passed: cp3, weight: 0.35 });

// Dense reward: sum of weighted checkpoints (returns 0.0 to 1.0)
return checkpoints.reduce((sum, cp) => sum + (cp.passed ? cp.weight : 0), 0);

```

Figure 12: A generated dense reward evaluator with weighted checkpoints. Each checkpoint validates a different aspect: (1) user action tracking via instrumentation (weight 0.35), (2) data record consistency verification (weight 0.30), and (3) confirmation state validation (weight 0.35). Partial task completion yields proportional rewards, e.g., completing only the subscription attempt earns 0.35 points. This enables more effective GRPO training compared to sparse 0/1 rewards.

Method	Task Completion (%)
Bolt.diy	67.1 ± 1.8
Claude-Code	75.8 ± 2.4
Codex	80.8 ± 0.4
InfiniteWeb (Ours)	85.6 ± 1.2

Table 5: Detailed results on WebGen-Bench with standard deviation over three runs.

Configuration	Task Completion (%)
GPT-4.1 (vs GPT-5)	77.4 ± 1.4
w/o TCTDD Validation	82.6 ± 1.4
InfiniteWeb (Full)	85.6 ± 1.2

Table 6: Ablation study results on WebGen-Bench.

difficulty. Welch’s t-tests comparing against baseline show: 200 tasks ( $t=2.96, p<0.05$ ), 400 tasks ( $t=4.47, p<0.05$ ), and 600 tasks ( $t=6.58, p<0.01$ ).

### C.3 Appearance Win Rate

Table 8 presents the appearance comparison results across three independent runs.

Training	Easy	Medium	Hard	Overall
Orig.	46.9	18.9	5.3	<b>23.0±0.9</b>
200	54.3	18.2	7.9	<b>25.3±1.0</b>
400	48.2	27.3	5.3	<b>27.3±1.4</b>
600	56.8	23.1	9.2	<b>28.7±1.2</b>

Table 7: Results on Online-Mind2Web by difficulty (%). 600/400/200 denote InfiniteWeb with different training task counts. Orig. is the original baseline. Standard deviation computed over three runs.

### C.4 Comparison with Claude-Code-Generated Websites

To disentangle the effect of data quality from data quantity, we compare training on InfiniteWeb-generated websites against training on Claude-Code-generated websites using the same 600-task GRPO setup. Both use GPT-5 with reasoning effort set to “high” for generation. All results are averaged over three runs.

Training on Claude-Code-generated websites yields marginal or even negative improvements, while InfiniteWeb achieves significant gains across all three benchmarks. This demonstrates that data quality, not quantity, drives the improvement. The quality gap stems from context management: In-

Comparison	Win Rate (%)
Ours vs Claude-Code	85.4 $\pm$ 0.5
Ours vs Bolt.diy	84.2 $\pm$ 0.8
Ours vs Codex	69.3 $\pm$ 0.7

Table 8: Appearance win rate comparison. Win rate indicates how often InfiniteWeb-generated websites are judged visually closer to the reference design image. Standard deviation computed over three runs.

Training Data	Mind2Web	OSWorld	MobileWorld
Baseline	23.0	24.5	6.4
Claude-Code	23.7 (+0.7)	25.2 (+0.7)	5.1 (-1.3)
InfiniteWeb	<b>28.7</b> (+5.7)	<b>31.4</b> (+6.9)	<b>10.3</b> (+3.9)

Table 9: Comparison of training on InfiniteWeb vs. Claude-Code-generated websites (%). Both use 600 tasks with identical GRPO training. Mind2Web refers to Online-Mind2Web.

finiteWeb provides each generation step with only directly relevant specifications (interface definitions, data models), while Claude-Code accumulates the full codebase in context, leading to cross-page inconsistencies, reduced scalability (4.9 vs. 12.6 pages per website), and less complete business logic as the context grows.

## D Human Evaluation for Visual Quality

To validate the reliability of our automated visual evaluation (Section 4), we conducted a human verification study. We randomly sampled 100 comparison cases across all three baseline comparisons (InfiniteWeb-Codex, InfiniteWeb-Claude, and InfiniteWeb-Bolt), with approximately equal representation from each. Human evaluators were presented with the reference design image and two website screenshots (A and B), and asked to determine which implementation more closely matches the reference design.

The human judgments achieved a 91% agreement rate with the automated GPT-5 evaluations, indicating that the automated visual quality assessment is highly reliable and well-aligned with human perception. The disagreement cases primarily involved subtle differences where both implementations were reasonably close to the reference design, making the distinction less clear-cut.

## E Human Verification of Task and Evaluator Quality

To validate the quality of generated tasks and automatic evaluators, we conducted a manual veri-

fication study. We randomly sampled 100 tasks from the generated websites and had human evaluators assess: (1) whether the task description is clear and executable on the generated website, and (2) whether the automatic evaluator correctly determines task completion.

Of the 100 sampled tasks, 95 passed human verification, confirming that our system generates high-quality tasks with reliable automatic evaluators.

Additionally, we analyzed the TCTDD validation loop statistics. Among all generated websites, only 1.5% remained unfixed after the maximum 8 TCTDD iterations, demonstrating the effectiveness of our iterative test-driven approach in ensuring functional correctness.

## F Prompts

## Prompt Template for Baseline Website Generation

Generate a full static {website\_type} website using HTML/CSS/JavaScript in the current directory.

### **BASIC REQUIREMENTS:**

- Create up to 12 pages
- Homepage must be index.html
- Use browser's localStorage to store data
- Use design\_image.png as the visual design reference
- The website must implement at least these functions: {function\_requirements}

### **FILE STRUCTURE:**

- Each HTML page should have its own CSS file (e.g., index.html → index.css)

### **DESIGN MATCHING:**

- Carefully analyze design\_image.png and extract: color scheme, typography, layout patterns, spacing system
- Accurately reproduce the visual style to ensure design consistency across all pages

### **HTML STANDARDS:**

- Use semantic HTML5 elements: <header>, <main>, <footer>, <nav>, <section>, <article>

### **CSS STANDARDS:**

- Use modern CSS features: Flexbox, Grid, CSS Custom Properties
- Implement interactive states: hover effects and smooth transitions

### **JAVASCRIPT STANDARDS:**

- Use modern ES6+ syntax: template literals, arrow functions, const/let, destructuring

### **FUNCTIONALITY REQUIREMENTS:**

- Ensure all specified user tasks can be completed end-to-end
- Fully implement each page's core functionality, not just static displays
- Beyond the explicitly required functions, add other common features appropriate for this website seed
- When page parameters are missing, provide reasonable default content

### **DATA QUALITY:**

- Ensure temporal data alignment: dates should be logically consistent
- Generate diverse data with sufficient variety to support different scenarios
- Create realistic, professional content appropriate for the website seed

Figure 13: Prompt template for baseline website generation. Variables {website\_type} and {function\_requirements} are filled based on the input specification.

## Prompt: Task Generation

You are a UX researcher. Generate {task\_count\_range} realistic user tasks for a {website\_type}.

### IMPORTANT REQUIREMENTS:

1. This is a mock website, so tasks should NOT depend on any external services like email authentication.
2. Each task MUST contain between {min\_steps}-{max\_steps} detailed steps for proper complexity.
3. Tasks should be suitable for RL model training, requiring multiple decisions and interactions.

### Each task should:

- Represent a SPECIFIC user goal with MEASURABLE success criteria
- Contain {min\_steps}-{max\_steps} DETAILED action steps
- Include CLEAR decision criteria (e.g., “select the cheapest option”, “choose items with 4+ stars”)
- Specify EXACT targets (e.g., “add 3 items under \$50”, “find products with free shipping”)
- Use CONCRETE values and thresholds (prices, quantities, ratings, dates)
- Cover different aspects of the website functionality

### Task specificity requirements:

- BAD: “Compare products and select the best one”
- GOOD: “Compare two laptops and select the one with more RAM under \$1000”
- BAD: “Search for headphones and add to cart”
- GOOD: “Search for wireless headphones under \$200 with 4+ star rating and add the first result to cart”

### Step detail requirements (FOCUS ON ACTIONS, NOT VERIFICATION):

- Specific navigation actions (e.g., “Navigate to the homepage”)
- Clear element interactions (e.g., “Click the search button in the header”)
- Precise data entry (e.g., “Type ‘wireless headphones’ in the search field”)
- Selection actions (e.g., “Select ‘Blue’ from the color dropdown”)
- Page transitions (e.g., “Click on the product image to open details page”)

### AVOID these types of steps:

- Verification steps (e.g., “Verify the page loaded”)
- Validation steps (e.g., “Validate the price is correct”)
- Confirmation steps (e.g., “Ensure the button is visible”)

### Return JSON format:

```
{ "tasks": [ { "id": "task_1", "name": "...",  
  "description": "...", "steps": ["..."]} ] }
```

Figure 14: Prompt for automatic task generation from website seed.

## Prompt: Primary Architecture Design

Design a complete website architecture for a {website\_type}.

### User Tasks that the website must support:

{tasks\_text}

Based on these tasks, design a COMPLETE architecture with ALL pages needed:

1. All pages needed for the website (maximum {max\_pages} pages)
2. Primary functions each page should provide
3. Keep it simple and focused on user needs
4. DO NOT include authentication/login pages
5. DO NOT consider multi-user scenarios
6. This is for single-user use only

### Return JSON format:

```
{ "all_pages": [ { "name": "Page Name", "filename": "page.html" } ],  
  "pages": [ { "name": "Page Name", "filename": "page.html",  
    "description": "Brief description",  
    "primary_functions": ["Function 1", "Function 2"]} ] }
```

### Requirements:

- Include all pages needed to complete the user tasks
- Each page should have clear, focused responsibilities
- Use descriptive filenames (e.g., index.html, products.html, cart.html)
- Primary functions should be high-level user actions
- Ensure all task steps can be completed with the designed pages
- index.html must be contained and as the homepage

Figure 15: Prompt for designing primary website architecture based on user tasks.

### Prompt: Data Model Extraction

You are a data architect. Analyze the user tasks and extract ALL data entities and fields needed.

**Website Seed:** {website\_seed}

**User Tasks:** {tasks\_json}

**Website Architecture Pages:** {pages\_json}

For each task, identify:

1. Core entities directly mentioned (e.g., Product, Cart)
2. Supporting entities needed for functionality
3. All necessary fields for each entity
4. Relationships between entities

#### IMPORTANT REQUIREMENTS:

- This is for SINGLE-USER agent training only - NO multi-user support needed
- DO NOT include User entity or userId/sessionId fields
- DO NOT include authentication-related entities
- Extract ALL entities needed, not just the minimal set
- Include all fields necessary for the tasks
- Specify data types for each field
- Identify primary keys (but NO foreign keys to User)
- Specify data\_pre\_generation\_num for each entity: "many", "few", or "none"
  - "many": Generate 10-20 items (for catalog entities like Product, Category)
  - "few": Generate 3-5 items (for limited entities like Brand, Department)
  - "none": No pre-generation needed (for runtime entities like Cart, Order)
- Provide storage\_key for localStorage (lowercase plural form)

#### Return JSON format:

```
{
  "entities": [
    {
      "name": "Product",
      "storage_key": "products",
      "fields": [
        {
          "name": "id",
          "type": "string",
          "primary_key": true
        },
        {
          "name": "price",
          "type": "number",
          "required": true
        }
      ],
      "data_pre_generation_num": "many"
    }
  ],
  "relationships": [
    {
      "from": "CartItem",
      "to": "Product",
      "type": "belongs_to",
      "field": "productId"
    }
  ]
}
```

Figure 16: Prompt for extracting data models from user tasks.

## Prompt: Interface Design

You are a software architect. Design comprehensive interfaces for both user tasks AND page functionality.

**Website Seed:** {website\_seed}

**User Tasks:** {tasks\_json}

**Data Models:** {data\_models\_json}

**Website Pages and Functions:** {pages\_info}

### IMPORTANT REQUIREMENTS:

1. Design USER-FACING interfaces that will be directly called from UI
2. This is for SINGLE-USER agent training - NO userId, sessionId parameters
3. System state (cart, session) should be managed internally, not passed as parameters

### CRITICAL: Design interfaces for TWO purposes:

#### A. TASK EXECUTION INTERFACES - For user tasks:

- What information must be shown BEFORE the user can act (display interfaces)
- What action the user performs (action interfaces)
- What feedback/results need to be shown AFTER the action (result interfaces)

#### B. PAGE FUNCTIONALITY INTERFACES - For each page's primary\_functions:

- Review EVERY primary\_function in the Website Pages list
- Ensure there's an interface to support EACH function
- Examples: "Navigate to featured product categories" → needs getCategory()

#### Additional requirements:

- Interfaces should handle complete operations (e.g., addToCart handles cart creation if needed)
- Do NOT create unnecessary CRUD, but DO create interfaces needed for page display
- For interfaces that get data for display, return user-friendly fields

#### Return JSON format:

```
{
  "interfaces": [
    {
      "name": "addToCart",
      "description": "Add a product to cart",
      "parameters": [
        {
          "name": "productId",
          "type": "string"
        }
      ],
      "returns": {
        "type": "object",
        "properties": {
          "success": {
            "type": "boolean"
          }
        }
      },
      "relatedTasks": ["task_1"],
      "helperFunctions": [
        {
          "name": "_getOrCreateCart",
          "description": "Internal helper",
          "visibility": "private"
        }
      ]
    }
  ]
}
```

Figure 17: Prompt for designing user-facing interfaces based on tasks and data models.

## Prompt: Interface Wrapping

You are a software architect analyzing interface parameters for a {website\_type}.

Your task: Identify parameters that should be hidden from user-facing interfaces and generate wrapped versions.

**ORIGINAL INTERFACES:** {original\_interfaces\_json}

**EXISTING DATA MODELS:** {data\_models\_json}

**PARAMETER CLASSIFICATION RULES:**

**1. SYSTEM-MANAGED PARAMETERS (should be hidden):**

- User identity: userId, guestId, sessionId, currentUser
- System context: cartId, deviceId, timestamp, requestId
- Authentication: authToken, userRole, permissions, isAuthenticated
- Environment: locale, timezone, region, currency

**2. USER-PROVIDED PARAMETERS (should remain exposed):**

- Business data: productId, quantity, rating, comment
- User selections: selectedSize, color, filters
- User input: searchQuery, address, paymentDetails

**ANALYSIS CRITERIA:**

- Ask: "Would a user type this value into a form or select it from a UI?"
- If YES → Keep as parameter (user-provided)
- If NO → Hide and manage through state (system-managed)

**EXAMPLE TRANSFORMATION:**

Original: addToCart(userId, guestId, productId, quantity, selectedSize)

Wrapped: addToCart(productId, quantity, selectedSize)

State Needed: UserSession with currentUserId/currentGuestId

**Return JSON format:**

```
{
  "wrapped_interfaces": [
    {
      "name": "addToCart",
      "parameters": [
        {
          "name": "productId",
          "type": "string"
        }
      ],
      "state_data_models": [
        {
          "name": "UserSession",
          "fields": [
            {
              "name": "currentUserId",
              "type": "string"
            }
          ]
        }
      ],
      "implementation_mapping": [
        {
          "wrapped_function": "addToCart",
          "parameter_mapping": {
            "userId": "_getSession().currentUserId"
          }
        }
      ]
    }
  ]
}
```

Figure 18: Prompt for wrapping interfaces to hide system-managed parameters.

## Prompt: Architecture Design

You are a web architect. Design complete website architecture based on user tasks and interfaces.

**Website Seed:** {website\_seed}

**User Tasks:** {task\_summary\_json}

**Primary Architecture (initial design):** {primary\_arch\_json}

**Available Interfaces:** {interface\_summary\_json}

**Data Models:** {data\_summary\_json}

### IMPORTANT:

- This is for SINGLE-USER agent training - NO authentication/login pages needed
- The interfaces provided are USER-FACING interfaces (no userId/sessionId parameters)
- System state is managed automatically through localStorage

### Design Requirements:

1. Use EXACTLY the pages from primary architecture - do not add or remove pages
2. Assign appropriate interfaces to each page based on functionality
3. Use URL parameters for navigation (NOT localStorage for page data)
4. Define incoming parameters (what parameters the page accepts)
5. Define outgoing connections (what pages this page navigates to)
6. Specify access methods for each page
7. Design header and footer navigation links

### Access Method Guidelines:

- "navigation": Accessible through header/footer navigation
- "url\_param": Accessible through URL parameters from other pages
- "direct\_link": Accessible through direct links in content
- "form\_submission": Accessible after form submission

### Return JSON format:

```
{ "all_pages": [{"name": "Home", "filename": "index.html"}],  
  "pages": [{"name": "Home", "filename": "index.html",  
             "assigned_interfaces": ["searchProducts"],  
             "incoming_params": [],  
             "outgoing_connections": [{"target": "product.html",  
                                       "params": {"id": "productId"}}],  
             "access_methods": [{"type": "navigation"}]},  
  "header_links": [{"text": "Home", "url": "index.html"}]}
```

Figure 19: Prompt for designing complete website architecture with page navigation.

### Prompt: Page Functionality Design

You are a senior web functional designer. Design the functional aspects and workflows of a webpage.

**Website Seed:** {website\_seed}

**Page Architecture:** {page\_spec\_json}

**Available Data Models:** {data\_dict\_json}

**Assigned Interfaces for This Page:** {interface\_details\_json}

**Navigation Information:** {navigation\_info}

**DESIGN REQUIREMENTS:**

1. Create an engaging, specific page title
2. Write a rich, detailed description of the page
3. Design core features based on the assigned interfaces
4. Define user workflows that utilize the interfaces
5. Specify user interactions (clicks, forms, navigation)
6. Describe state logic using URL parameters (NOT localStorage)
7. Create functional components that use the interfaces

**IMPORTANT GUIDELINES:**

- Use ONLY the assigned interfaces for this page
- Navigation uses URL parameters (e.g., product.html?id=123)
- Focus on functionality, not visual appearance
- Components should be functional, not presentational
- Each component should have clear data binding and event handlers
- Output should not involve any static data or hardcoded values

**Return JSON format:**

```
{ "title": "Page title", "description": "Page description",  
  "page_functionality": { "core_features": ["Feature 1"],  
    "user_workflows": ["Workflow step"],  
    "interactions": ["Click action"],  
    "state_logic": "URL parameter handling"},  
  "components": [{"id": "search-form", "type": "search-form",  
    "functionality": "Handles product search",  
    "data_binding": ["Product"],  
    "event_handlers": ["onSubmit"]}]}  
}
```

Figure 20: Prompt for designing page functionality and components.

## Prompt: Design Image Analysis

You are a senior UI/UX design analyst. Analyze the provided design image to extract all visual characteristics.

**Website Seed:** {website\_seed}

### ANALYSIS TASKS:

#### 1. Visual Features Analysis:

- Identify overall visual style (modern, minimalist, vibrant, corporate, etc.)
- Describe visual hierarchy and focal points
- Note use of whitespace and visual breathing room

#### 2. Color Scheme Extraction:

- Primary colors (main brand colors)
- Secondary colors (supporting colors)
- Accent colors (for CTAs, highlights)
- Neutral colors (backgrounds, text, borders)
- Provide exact hex color values

#### 3. Layout Characteristics:

- Grid system (12-column, custom, etc.)
- Layout patterns (sidebar, centered, full-width)
- Section organization and alignment principles

#### 4. UI Patterns: Button styles, card designs, form elements, navigation patterns

#### 5. Typography: Font families, size hierarchy, font weights, line heights

#### 6. Spacing System: Base unit, padding/margin patterns, component spacing

#### 7. Interaction Hints: Hover states, transitions, animation suggestions

#### Return JSON format:

```
{
  "visual_features": {
    "overall_style": "modern minimalist",
    "color_scheme": {
      "primary": ["#hex"],
      "accent": ["#hex"]
    },
    "layout_characteristics": {
      "grid_system": "12-column"
    },
    "ui_patterns": [
      {
        "pattern_type": "button",
        "characteristics": {
          "shape": "rounded"
        }
      }
    ],
    "typography": {
      "font_families": {
        "heading": "Inter"
      }
    },
    "spacing_system": {
      "base_unit": "8px"
    }
  }
}
```

Figure 21: Prompt for analyzing design image to extract visual characteristics.

## Prompt: Layout Design

You are a senior UI/UX designer. Create a thoughtful, detailed layout for existing components.

### DESIGN DNA (extracted from design image):

- Visual Style: {visual\_style}
  - Grid System: {grid\_system}
  - Layout Pattern: {layout\_pattern}
  - Spacing System: {spacing\_system\_json}
- PAGE CONTEXT:** Website Seed: {website\_seed}, Page: {page\_name}

### Components to Layout: {components\_list}

#### STEP 1: Choose Layout Strategy Combination

For each dimension, provide reasoning and make a choice:

1. **Content Arrangement:** linear-flow, grid-based, asymmetric, centered-focus, masonry, split-screen, sidebar-content, magazine-layout
2. **Component Grouping:** functional-clusters, visual-zones, priority-based, workflow-aligned, data-centric
3. **Space Allocation:** equal-distribution, primary-focus, golden-ratio, thirds-rule, flexible-grid
4. **Content Density:** spacious, balanced, compact, variable
5. **Visual Flow:** top-down, z-pattern, f-pattern, circular, focal-center

**STEP 2:** Describe each component's layout using natural language (position, size, relationships)

**STEP 3:** Describe overall layout picture

#### Return JSON format:

```
{
  "chosen_strategies": {
    "content_arrangement": {
      "reasoning": "...",
      "choice": "grid-based"
    },
    "overall_layout_description": "Description of full layout",
    "component_layouts": [
      {
        "id": "search-form",
        "layout_narrative": "Position and size description",
        "visual_prominence": "primary"
      }
    ]
  }
}
```

Figure 22: Prompt for designing component layouts based on design analysis.

### Prompt: Page Framework Generation

You are a senior web developer. Analyze the provided design image and generate a complete HTML framework with header and footer that matches the visual style.

**Website Seed:** {website\_seed}

**Header Navigation Links:** {header\_links\_json}

**Footer Links:** {footer\_links\_json}

**Design Analysis Context:** {design\_context}

**Requirements:**

1. ANALYZE THE DESIGN IMAGE to extract: visual style, color palette, typography, layout patterns, spacing
2. Create a complete HTML framework matching the design (reusable for all pages)
3. Only include header, footer, and main content area (id="content")
4. Header matching the design's header style with provided navigation links
5. Footer matching the design's footer style with provided footer links
6. Modern, semantic HTML5 structure

**CSS Requirements:**

- Extract exact colors from the design image
- Match typography from the design
- Replicate spacing and sizing
- Create CSS variables for the design system

**CRITICAL:**

- Use English only
- Do NOT include interactive elements without corresponding links
- SVG files are not allowed in the framework

**Return JSON format:**

```
{"framework_html": "HTML with header/footer",  
 "framework_css": "CSS replicating the design"}
```

Figure 23: Prompt for generating page framework (header/footer) from design image.

### Prompt: HTML Page Generation

You are a senior web developer. Generate the main content HTML for a {website\_type} website page with UI JavaScript.

**Page Information:** {page\_design\_json}

**Navigation Information:** {page\_architecture\_json}

**Framework HTML:** {framework\_html}

**Data Dictionary:** {data\_dict\_json}

**Page-Specific SDK Interfaces:** {page\_interfaces\_json}

**Requirements:**

1. Generate ONLY content for <main id="content"> section
2. Call interfaces as WebsiteSDK.functionName() - they are SYNCHRONOUS
3. Handle incoming\_params: Extract URL parameters this page expects
4. Implement outgoing\_connections: Navigate to other pages with correct parameters
5. Add data attributes: data-populate, data-action, data-component

**UI JavaScript Requirements:**

1. Initialize page when DOM is ready
2. Extract URL parameters for incoming\_params
3. Call SDK methods based on data-populate attributes
4. Set up event listeners based on data-action attributes
5. Implement navigation with correct parameters
6. Always call WebsiteSDK.methodName() directly (no method extraction)

**CRITICAL:** Call SDK interfaces with positional arguments only. Use only relative .html URLs for internal navigation.

**Return:** {"html\_content": "Complete HTML page with UI JavaScript"}

Figure 24: Prompt for generating HTML pages with integrated UI JavaScript.

### Prompt: CSS Page Generation

You are a senior web developer. Generate CSS styles for the page based on its HTML structure.

**Page Design:** {page\_design\_json}

**Page Layout:** {page\_layout\_json}

**Design Analysis:** {design\_analysis\_json}

**Framework CSS (build upon this):** {framework\_css}

**Generated HTML (style this content):** {html\_content}

**Requirements:**

1. Include complete framework CSS - no abbreviations
2. Style the content area and page-specific components
3. Follow the design analysis color scheme and typography
4. Implement the layout specifications (grid, spacing, etc.)
5. Ensure responsive design with proper breakpoints
6. Use CSS variables defined in framework CSS
7. Add hover states and transitions for interactive elements
8. Use modern CSS features (flexbox, grid, custom properties)

**CRITICAL:** Put this at the VERY TOP of css\_content:

```
[hidden] { display: none !important; }
```

**Return:** {"css\_content": "Complete CSS including framework and page-specific styles"}

Figure 25: Prompt for generating CSS styles based on HTML structure and design analysis.

### Prompt: Data Generation

You are a data generator specializing in realistic website data. Generate comprehensive, realistic data based on the EXACT data dictionary specifications.

**Website Seed:** {website\_seed}

**User Tasks Context:** {tasks\_json}

**Data Dictionary Structure:** {data\_types\_info\_json}

**CRITICAL CONSTRAINTS:**

1. **Use data\_type\_name as JSON key:** Use the exact value from "data\_type\_name" field
2. **Use EXACT field names:** Only fields defined in fields dictionary
3. **Follow field types:** string, number, boolean, array, datetime as specified
4. **Intelligent Volume Decision:** Based on generation\_type:
  - "many": Generate substantial amount approaching max\_items
  - "few": Generate small representative set (20-30% of max\_items)
5. **No extra fields:** Do NOT add fields not in the dictionary

**IMAGE URL REQUIREMENTS:** Use ONLY real, working image services:

- Unsplash: [https://images.unsplash.com/photo-\[ID\]?w=800&h=600](https://images.unsplash.com/photo-[ID]?w=800&h=600)
- Picsum: [https://picsum.photos/800/600?random=\[1-1000\]](https://picsum.photos/800/600?random=[1-1000])

**DATA QUALITY:** Generate realistic, diverse content appropriate for the website seed. Ensure data relationships are logical and consistent.

**Return JSON format:**

```
{"static_data": {"products": [{"field1": "value"}],  
  "categories": [{"id": "cat_1", "name": "Category"}]}}
```

Figure 26: Prompt for generating realistic website data based on data models.

### Prompt: Backend Implementation Generation

You are an expert JavaScript developer. Generate a complete business logic implementation.

**Website Seed:** {website\_seed}

**Tasks:** {tasks\_json}

**Data Models:** {data\_models\_json}

**Interfaces:** {interfaces\_json}

**REQUIREMENTS:**

1. Implement ALL core interfaces specified
2. Add helper functions as needed (prefix with \_ for private)
3. Use localStorage for ALL data persistence (browser-compatible)
4. NO DOM operations, NO window/document references (except localStorage)
5. Must work in both browser and Node.js environments
6. All data must be JSON serializable for localStorage
7. Implement interfaces with positional arguments only

**STRUCTURE:**

```
const localStorage = (function() { ... })(); // polyfill
class BusinessLogic {
  constructor() { this._initStorage(); }
  _initStorage() { /* init localStorage tables */ }
  _getFromStorage(key) { /* retrieve data */ }
  _saveToStorage(key, data) { /* persist data */ }
  addToCart(productId, quantity) { /* implementation */ }
}
module.exports = BusinessLogic;
```

**Return:** {"code": "javascript code here"}

Figure 27: Prompt for generating business logic implementation.

### Prompt: Backend Test Generation

You are an expert test engineer. Generate flow-based integration tests for the business logic.

**Website Seed:** {website\_seed}

**Tasks:** {tasks\_json}

**Interfaces:** {interfaces\_json}

**Generated Data:** {generated\_data\_json}

**CRITICAL REQUIREMENTS:**

1. Use Generated Data ONLY in setupTestData() for initial localStorage population
2. NEVER hardcode expected return values - always extract from actual API responses
3. Chain API calls properly: Call API, capture response, extract needed values for next calls
4. Test complete user flows, not individual functions
5. Focus on happy path (successful scenarios)
6. Must run in Node.js environment
7. Test ALL tasks provided

**CORRECT Flow Testing Example:**

```
const addResult = this.logic.addToCart(userId, productId, 2);
const actualCartId = addResult.cartId; // Extract from response
const cartData = this.logic.getCart(actualCartId); // Use actual ID
this.assert(cartData.total > 0, 'Total should be positive');
```

**Return:** {"code": "javascript test code"}

Figure 28: Prompt for generating flow-based integration tests.

### Prompt: Evaluator Generation

You are a QA engineer. Create evaluators to check if users complete tasks successfully.

**Website Seed:** {website\_seed}

**Tasks to evaluate:** {tasks\_json}

**Cross-Page States Structure:** {cross\_page\_states\_json}

**Generated Data Structure:** {data\_structure\_json}

For each task, create an evaluator that:

- Uses cross-page states stored in localStorage to determine completion
- Uses data structure knowledge to create precise validation logic
- References exact field names and data types from the data structure
- Provides clear evaluation criteria and logic
- Uses JavaScript logic to check task completion status

**Guidelines:**

- Use localStorage.getItem() to access both cross-page states and static data
- Parse JSON data when retrieving complex objects from localStorage
- Check for null/undefined values before accessing object properties
- Use realistic validation logic based on the actual data structure

**Return JSON format:**

```
{
  "evaluators": [
    {
      "task_id": "task_1",
      "name": "Evaluator Name",
      "description": "What this evaluator checks",
      "localStorage_variables": ["selectedProductId", "products"],
      "evaluation_logic": "const products = JSON.parse(...); ..."
    }
  ]
}
```

Figure 29: Prompt for generating task completion evaluators.

### Prompt: Instrumentation Analysis

You are analyzing JavaScript business logic to determine what instrumentation variables are needed to evaluate task completion.

**TASKS TO EVALUATE:** {tasks\_json}

**CURRENT BUSINESS LOGIC:** {code\_snippet}

**EXISTING LOCALSTORAGE VARIABLES:** {existing\_storage\_vars\_json}

**DATA STORAGE KEYS:** {storage\_keys\_json}

**ANALYSIS REQUIREMENTS:** For each task, determine:

1. What operations must occur for the task to be considered complete?
2. Can we use existing localStorage variables to determine completion?
3. If NOT, what new instrumentation variables are needed?

**INSTRUMENTATION GUIDELINES:**

- Only add variables if existing localStorage is insufficient
- Use naming convention: taskN\_actionDescription (e.g., task1\_searchCompleted)
- Specify which function should set the variable and under what condition
- Be conservative - only add instrumentation if truly necessary

**Return JSON:**

```
{
  "requirements": [
    {
      "task_id": "task_1",
      "needs_instrumentation": true,
      "required_variables": [
        {
          "variable_name": "task1_searchCompleted",
          "set_in_function": "searchNeighborhoods",
          "set_condition": "After successful search"
        }
      ]
    }
  ]
}
```

Figure 30: Prompt for analyzing instrumentation requirements for task tracking.

### Prompt: Instrumentation Code Generation

You are adding instrumentation variables to JavaScript business logic for task completion tracking.

**ORIGINAL CODE:** {original\_code}

**INSTRUMENTATION SPECIFICATIONS:** {instrumentation\_specs\_json}

**INSTRUCTIONS:** For each instrumentation variable:

1. Find the specified function in the code
2. Add localStorage.setItem() call at the appropriate location based on set\_condition
3. Wrap instrumentation code in try-catch to ensure non-invasive behavior
4. Use the exact variable\_name and value\_to\_set from specifications

**CRITICAL REQUIREMENTS:**

- DO NOT change any original functionality
- DO NOT modify function signatures or return values
- Instrumentation code must be wrapped in try-catch
- Only add localStorage.setItem() calls as specified
- Preserve all existing code structure and comments
- Place instrumentation BEFORE the return statement

**Return:** Complete instrumented business\_logic.js code

Figure 31: Prompt for generating instrumented code with tracking variables.

### Prompt: Instrumentation Evaluator Generation

You are generating evaluators to check if users completed tasks successfully.

**TASKS:** {tasks\_json}

**INSTRUMENTATION VARIABLES AVAILABLE:** {var\_mapping\_json}

**BUSINESS LOGIC IMPLEMENTATION:** {business\_logic\_code}

**WEBSITE DATA:** {website\_data\_json}

**INSTRUCTIONS:** For each task, create an evaluator based on the instrumentation plan:

**Case 1: Tasks with needs\_instrumentation=true**

- Use the instrumentation\_variables specific to that task
- Validate the variable values match expected values

**Case 2: Tasks with needs\_instrumentation=false**

- Use the existing\_variables to infer task completion
- Check the ACTUAL data structure from the business logic implementation

All evaluators must:

- Check if the variables exist in localStorage
- Use the EXACT data structure from the business logic implementation
- Return true if the task is completed, false otherwise

**Return JSON:**

```
{
  "evaluators": [
    {
      "task_id": "task_1",
      "name": "...",
      "localStorage_variables": ["var1", "var2"],
      "evaluation_logic": "// JavaScript returning boolean"
    }
  ]
}
```

Figure 32: Prompt for generating evaluators with instrumentation support.