

# ReLook: Vision-Grounded RL with a Multimodal LLM Critic for Agentic Web Coding

Anonymous ACL submission

## Abstract

While Large Language Models (LLMs) excel at algorithmic code generation, they struggle with front-end development, where correctness is judged on rendered pixels and interaction. We present **ReLook**, an *agentic*, vision-grounded reinforcement learning framework that empowers an *agent* to close a robust generate–diagnose–refine loop by invoking a multimodal LLM (MLLM) as a tool. During training, the agent employs an MLLM-in-the-loop to serve as a visual critic, evaluating code via screenshots and providing actionable feedback. Crucially, we enforce a strict zero-reward policy for invalid renders to guarantee renderability and mitigate reward hacking. To prevent behavioral collapse, we introduce *Forced Optimization*, a strict acceptance rule that admits only improving revisions, yielding monotonically better trajectories. At inference, we decouple the critic and run a lightweight, critic-free self-edit cycle, keeping latency *comparable to base decoding* while retaining most of the gains. Across three widely used benchmarks, ReLook consistently outperforms strong baselines in vision-grounded front-end code generation, highlighting the benefits of agentic perception, visual rewards, and training–inference decoupling.

## 1 Introduction

Large Language Models (LLMs) excel on closed-form benchmarks—programming contests (Li et al., 2022), SQL synthesis (Liu et al., 2024), and mathematical reasoning (Yang et al., 2024)—yet still underperform on front-end code generation, where visual fidelity and interaction are first-class. Unlike binary unit tests in algorithmic tasks, front-end quality lies on a continuum: a single misaligned pixel can signify failure.

This perceptual barrier explains current shortcomings: text-only models are blind to pixel-level consequences, yielding (i) layout drift, (ii) interaction breakage, and (iii) aesthetic inconsistency.

To address this, a model must (1) see rendered HTML/CSS/JS/SVG, (2) diagnose misalignments and broken interactions, and (3) iteratively refine in situ. Existing methods miss this loop: one-shot vision-to-code systems (pix2code (Wüst et al., 2024), Design2Code (Si et al., 2024), UICoder (Wu et al., 2024)) generate but do not refine; self-refinement frameworks (CodeRL (Le et al., 2022), Self-Refine (Madaan et al., 2023), Reflexion (Shinn et al., 2023), CRITIC (Gou et al., 2023; Peng et al., 2025; Zhang et al., 2025b)) iterate but cannot see, relying on pixel-blind unit tests or linters.

To bridge this gap, we introduce **ReLook**, a vision-grounded agentic reinforcement learning framework that completes the generate–diagnose–refine loop. The agent actively invokes an MLLM as a tool to "see" rendered outputs and obtain rich textual suggestions during inference, enabling true iterative refinement. Training uses a comprehensive reward system: a powerful MLLM (e.g., Qwen2.5-VL (Wang et al., 2024)) supplies the perceptual signal text-only methods lack, and a rendering-integrity rule assigns zero reward when required screenshots are invalid to deter reward hacking.

However, we identify a critical challenge: behavioral collapse, where despite high-quality feedback, a subsequent revision can be worse. We adopt a Forced Optimization strategy that accepts only strictly improving steps, ensuring high-quality, monotonically improving trajectories. For low-latency inference, the external critic can be discarded; the model performs a lightweight self-edit cycle—render, self-edit, and converge quickly to a human-aligned result.

**Evaluator validity and choice.** Our offline evaluation strictly follows the ARTIFACTSBENCH protocol (Zhang et al., 2025a). ARTIFACTSBENCH establishes evaluator validity through: (i) controlled human studies demonstrating over 90%

agreement between MLLM judges (Gemini-2.5-Pro, Qwen2.5-VL-72B) and human experts, and (ii) strong ranking correlation with WebDev Arena (LMSYS Org, 2024), a large-scale crowd-sourced platform. Since our test sets are strict subsets of ARTIFACTSBENCH’s evaluated tasks, we directly inherit this established human-alignment evidence. To further mitigate on-policy judge overfitting, we decouple the training-time critic (Qwen2.5-VL-72B-Instruct) from the offline evaluator (Gemini-2.5-Pro). We do not conduct additional human studies; validity rests on ARTIFACTSBENCH’s rigorous validation. See Appendix for detailed protocol adherence and cross-judge analysis.

Our contributions are as follows:

- **Robust Reward System.** We employ an MLLM as the reward model to provide the rich, pixel-level training signal that text-only methods cannot capture. This is critically supplemented by a zero-reward rule for answers without screenshots, which is designed to prevent reward hacking by forcing the agent to produce renderable code.
- **Agent reinforcement learning Framework.** We empower the agent to perform a generate–diagnose–refine loop by actively invoking an MLLM as a diagnostic tool. The agent can “see” its rendered output and receive rich, actionable feedback for iterative improvement. To ensure this powerful loop is productive and stable, we introduce a Forced Optimization strategy that addresses the challenge of behavioral collapse by guaranteeing the construction of high-quality, monotonically improving rollout trajectories.
- **Broad Applicability.** We perform extensive experiments on three widely-used benchmark datasets, and demonstrate that ReLook significantly outperforms the baselines. Moreover, we show the compatibility of ReLook by integrating it with different LLMs.

## 2 Method

### 2.1 Problem Formulation

Given a natural-language query  $q$ , the model outputs text  $t$  and front-end code  $c$  (HTML/CSS/JS/SVG). Correctness depends on rendered appearance and behavior rather than

syntax. We therefore use an MLLM judge to provide perceptual reward and design an agentic RL framework that allows vision-grounded feedback to refine  $c$ .

### 2.2 Overall Framework

ReLook runs a generate–diagnose–refine loop: render to temporal screenshots, score with an MLLM, incorporate feedback, and continue under a *Forced Optimization* rule that accepts only strictly improving steps. Reflection is internalized so the external critic can be dropped at test time.

### 2.3 Iterative Reflection Mechanism

For each query  $q$ , the policy emits  $t$  and  $c$ . Upon `<get_feedback>`, we execute  $c$  in a sandbox, capture screenshots, and query the MLLM for feedback  $m$  (wrapped in `<mlm_feedback>`). We feed  $\{q, t, c, m\}$  into the next round and stop when feedback is not requested or a round cap is reached.

The final output is represented as:

$$o = [t_1 \oplus c_1 \oplus m_1 \oplus \dots \oplus t_R \oplus c_R] \quad (1)$$

where  $t_r, c_r, m_r$  denote the  $r$ -th round’s text, code and feedback blocks, and  $R$  is the total number of reflection rounds. The prompt template is provided in the appendix.

### 2.4 Reinforcement Learning Framework

To optimize this framework, we employ Group Relative Policy Optimization (GRPO) as our training algorithm, which is based on a token-level policy gradient loss and is related to PPO (Schulman et al., 2017) while differing from preference-based objectives such as DPO (Rafailov et al., 2023). The objective is defined as follows:

$$\begin{aligned} \mathcal{J}_{GRPO}(\theta) &= \mathbb{E} [q \sim P(Q), \{o_i\}_{i=1}^G \sim \pi_{\theta_{\text{combined}}}(O|q)] \\ &= \frac{1}{\sum_{i=1}^G |o_i|} \sum_{i=1}^G \sum_{t=1}^{|o_i|} \left\{ \min \left[ \frac{\pi_{\theta}(o_{i,t}|q_{i,t})}{\pi_{\theta_{old}}(o_{i,t}|q_{i,t})} \hat{A}_{i,t}, \right. \right. \\ &\quad \left. \left. \text{clip} \left( \frac{\pi_{\theta}(o_{i,t}|q_{i,t})}{\pi_{\theta_{old}}(o_{i,t}|q_{i,t})}, 1 - \varepsilon, 1 + \varepsilon \right) \hat{A}_{i,t} \right] \right. \\ &\quad \left. - \beta \text{D}_{KL} [\pi_{\theta} || \pi_{ref}] \right\} \end{aligned} \quad (2)$$

Only tokens in  $t, c$  contribute non-zero advantages; critic tokens  $m$  are masked ( $\hat{A}_{i,t}=0$  on  $m$ ).

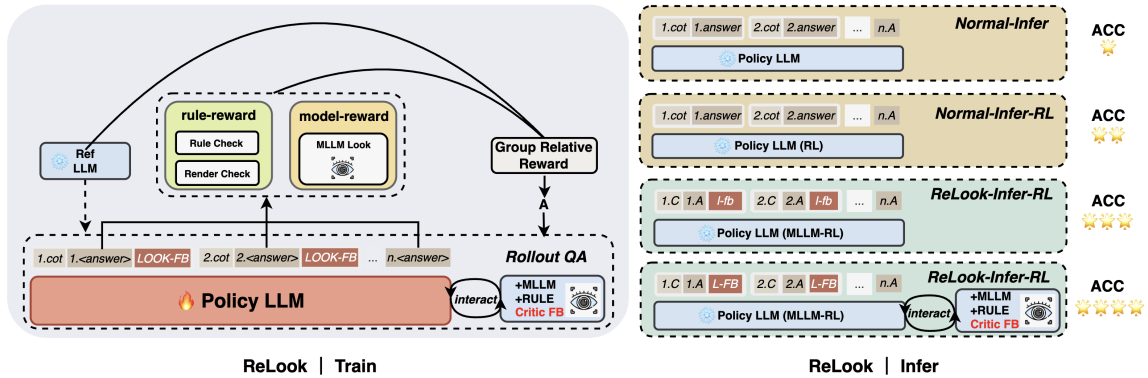


Figure 1: Overview of **ReLook**. Left: training closes a generate–diagnose–refine cycle: policy LLM generates code, pages rendered to temporal screenshots, and a vision-aware critic (MLLM) provides scores and feedback. Rewards combine visual scoring and format constraints; the policy is optimized with GRPO. Right: at inference the model runs a lightweight Re-Look cycle — external critic may be omitted for latency or used for higher accuracy.

### Advantage estimation and credit assignment.

We sample  $G$  trajectories per query and compute returns from  $R_{\text{ReLook}}(o_i)$  (Eq. 5). Using the group mean  $b$  as baseline, the advantage is  $\tilde{A}_i = R_{\text{ReLook}}(o_i) - b$ , broadcast to policy tokens (text  $t$ , code  $c$ ) while masking critic tokens  $m$  ( $\hat{A}_{i,t}=0$  on  $m$ ). Advantages are standardized and clipped to  $[-2, 2]$ . We regularize toward  $\pi_{\text{ref}}$  with KL weight  $\beta$ . The combined policy  $\pi_{\theta_{\text{combined}}}$  is defined as:  $[\pi_{\theta_{\text{old}}}$  for  $t_r, c_r] \oplus [\pi_{\text{MLLM}}$  for  $m_r]$ , where  $\pi_{\text{MLLM}}$  is a frozen MLLM critic (Qwen2.5-VL-72B-Instruct). Trajectories mix policy tokens ( $t_n, c_n$ ) and critic tokens ( $m_n$ ). During training, only  $\pi_{\theta}$  is updated. The hyperparameters  $\varepsilon$  and  $\beta$  control clipping range and KL regularization strength.

### Token masking and lightweight feedback distillation.

We optimize GRPO over policy tokens ( $t, c$ ) and mask critic tokens  $m$  by zeroing advantages. To enable critic-free inference, we optionally distill  $m$ -tokens with lightweight loss  $\mathcal{L}_{\text{distill}}^m$  toward frozen MLLM outputs (see Appendix §F for details). The total objective is  $\mathcal{L} = -\mathcal{J}_{\text{GRPO}}^{t,c} + \gamma \mathcal{L}_{\text{distill}}^m$  (default  $\gamma=0.1$ , sweep  $\gamma \in [0.05, 0.3]$ ). This lets RL improve visual quality while distillation transfers feedback style for test-time self-reflection.

## 2.5 Reward Design

Conventional RL signals for code, such as unit tests or linters, operate purely in the text domain and are blind to visual defects. To address this, our reward is derived from a powerful Multimodal LLM (Qwen2.5-VL-72B-Instruct) that scores rendered pages based on the user prompt, the gener-

ated code, and a series of temporal screenshots. Capturing screenshots at multiple time points allows the critic to assess dynamic behavior. Within each reflection round we capture three time points  $\{S_1, S_2, S_3\}$  (e.g., post-load, +1s, +2s) and jointly evaluate them in a single scoring call to obtain one round-level score. To properly credit incremental progress across rounds, we then average the round-level scores within a trajectory.

A critical component of this design is a safeguard against reward hacking, where malformed but syntactically plausible code might be over-rewarded. We enforce a strict renderability constraint: if any required screenshot is invalid (e.g., due to a render failure or timeout), the reward is zero. While this eliminates degenerate reward channels, it can lead to sparse rewards early in training. To mitigate this, we engineer all prompts to include a visual-output constraint that explicitly instructs the agent to write executable HTML/CSS/JavaScript/SVG code suitable for browser rendering (see Appendix for full prompt template). This simple but effective technique increases the initial Valid Render Rate (from  $\sim 40\%$  at the start of training to  $\sim 80\%$  upon convergence) and better aligns the generation task with our vision-grounded reward. This reward function is formally defined as:

$$R_{\text{MLLM}}(o) = \begin{cases} \text{VisualScore}(o) & \text{if screenshot valid} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

To discourage repetition, we apply a linear length penalty from  $L_{\text{start}}$  to  $L_{\text{end}}$  (12k and 14k

tokens):

$$R_{\text{len}}(o) = \begin{cases} 1 & \text{if } \text{len}(o) < L_{\text{start}} \\ \frac{L_{\text{end}} - \text{len}(o)}{L_{\text{end}} - L_{\text{start}}} & \text{if } L_{\text{start}} \leq \text{len}(o) \leq L_{\text{end}} \\ 0 & \text{if } \text{len}(o) > L_{\text{end}} \end{cases} \quad (4)$$

The final training reward is:

$$R_{\text{ReLook}}(o) = R_{\text{MLLM}}(o) \cdot R_{\text{len}}(o) \quad (5)$$

## 2.6 Forced Optimization

While the MLLM critic provides rich, detailed feedback, we observe a critical instability we term behavioral collapse: despite high-quality suggestions, a subsequent revision may score lower than the previous one, degrading trajectories.

We therefore adopt a **Forced Optimization** mechanism. We initially explored a *negative-reward penalty* for regressions (rejecting worse-than-previous outcomes by penalizing returns), but found it incentivized the agent to *reduce* reflection frequency to avoid penalties, i.e., reward hacking. Hence, we discard the penalty design and enforce a *strict acceptance rule*: a refinement step (new  $t_{r+1}$ ,  $c_{r+1}$ ) is accepted *only if* its reward strictly exceeds the best-so-far in the trajectory (no  $\epsilon$  margin or re-scoring). Non-improving steps are rejected and a new attempt is sampled, with a maximum of 10 resampling attempts per reflection round. If the limit is reached without improvement, we terminate further reflection for that trajectory and use the best-so-far result. This guarantees monotonically improving accepted trajectories and stabilizes learning without suppressing useful reflections.

## 2.7 Efficient Inference

During training we retain the MLLM feedback loop for self-correction. At inference, we drop the external MLLM and run a lightweight, critic-free self-edit (at most three rounds), with screenshots and MLLM calls disabled. This preserves most gains while substantially reducing latency, following a train-slow, run-fast paradigm. See Appendix for pseudocode.

# 3 Experiment

## 3.1 Experimental Settings

**Training Data Curation.** We curate a 3,000-task corpus of front-end-only tasks, normalize descriptions, and remove near-duplicates via lexical/DOM/code similarity. Prompts are audited to

remove hints and sanitized; the final data are split into train/val stratified by UI archetypes.

**Train-Test De-duplication Protocol.** To prevent leakage to ArtifactsBench, FullStack-Bench-Html and Web-Bench, we run instance-level, multi-view de-duplication before training:

- **Lexical:** TF-IDF over char 3-grams; cosine  $> 0.85$ .
- **DOM:** tag-bigram Jaccard  $> 0.90$ ; fallback tree-edit distance for edge cases.
- **Code:** token-set Jaccard  $> 0.90$  after stripping comments/whitespace and minifying.

If any criterion triggers, the instance is removed. Borderline cases (e.g., lexical in  $[0.80, 0.85]$  plus structural overlap) are manually reviewed. The same procedure purges intra-train near-duplicates.

**Dataset.** We evaluate the performance of our method on three widely used datasets: ArtifactsBench (Zhang et al., 2025a), FullStack-Bench-Html (Cheng et al., 2024) and Web-Bench (Xu et al., 2025). ArtifactsBench contains 1,825 tasks focused on generating dynamic and interactive visual outputs, such as SVG visualizations and mini-games. Its evaluation protocol uses a Multimodal LLM to score the visual fidelity and interactive integrity of the rendered code. ArtifactsBench utilizes Gemini-2.5-Pro as an expert judge to evaluate model outputs across its 1,825 tasks, demonstrating over 90% agreement with human evaluators. To manage evaluation costs, we conduct our experiments on six of its sub-datasets. We use the shorthand A-\* for ArtifactsBench subsets: A-Lite (300 randomly sampled cases), A-Easy (305 simple frontend cases), A-Game (all 413 game-related cases), A-SVG (all 123 SVG-focused cases), A-Web (all 447 web-specific cases), and A-Si (all 75 simulation-oriented cases). FullStack-Bench-Html provides a collection of front-end programming tasks where functional correctness is programmatically validated by passing a suite of predefined unit tests. Web-Bench simulates realistic web development workflows through 50 complex projects of 20 sequential tasks each. Following its official protocol, performance is measured by passing end-to-end test cases that validate the final project’s functionality, and we report the pass@2 rate.

**Sandboxed Rendering Environment.** We execute model-produced code within a headless Chromium-based renderer that is both deterministic and secure. This sandboxed environment operates at the OS level, with filesystem and network access disabled. Safety is further enhanced by blocking dangerous APIs (e.g., `window.open`, `alert/confirm/prompt`, `eval/Function`, clipboard access, non-local `fetch/XMLHttpRequest/WebSocket`), and enforcing a per-sample wall-clock timeout. Requests to non-whitelisted origins are intercepted and failed closed. To ensure determinism, all external resources like fonts and images are replaced with local fixtures; timers and animations use deterministic seeds; and we enforce a strict Content Security Policy that disallows inline scripts and remote scripts. Within this controlled environment, our visual capture mechanism is dynamic: we take full-page screenshots that automatically adjust to the content’s full dimensions, guaranteeing no elements are missed. To capture dynamic behavior, we record these screenshots at  $T=3$  distinct time points (e.g., post-load, +1s, +2s), which are then passed to the MLLM critic to compute the reward signal.

### Evaluator validity and cross-judge robustness.

We adopt the evaluation protocol from ARTIFACTSBENCH (Zhang et al., 2025a), which establishes validity through two key mechanisms: (i) a controlled human study showing over 90% agreement between MLLM evaluators (Gemini-2.5-Pro and Qwen2.5-VL-72B) and human experts across diverse front-end tasks, and (ii) strong ranking correlation (Spearman  $\rho > 0.85$ ) with WebDev Arena (LMSYS Org, 2024), a large-scale crowd-sourced evaluation platform. Critically, our test sets (A-Lite, A-Easy, A-Game, A-SVG, A-Web, A-Si) are strict subsets of ARTIFACTSBENCH’s human-validated tasks, allowing us to directly inherit the established human-alignment evidence. Consistent with ARTIFACTSBENCH, we decouple the training-time critic (Qwen2.5-VL-72B-Instruct, open-source) from the offline evaluator (Gemini-2.5-Pro, proprietary). Cross-judge consistency and further details are summarized in Appendix §B.

## 3.2 Evaluation Protocol

We follow a pixel-grounded evaluation protocol aligned with ARTIFACTSBENCH. For vision-based front-end tasks, models produce code that is exe-

cuted in our sandboxed browser to capture temporal screenshots at three time points (post-load, +1s, +2s). An independent evaluator (Gemini-2.5-Pro) assigns a VisualScore on the  $[0, 100]$  scale considering: (i) adherence to the textual specification, (ii) layout alignment and spatial fidelity, (iii) typography and color coherence, and (iv) interactive integrity when actions are specified. If no valid screenshot is produced, the score is set to zero. We report means over three runs (three random seeds) with identical decoding parameters.

For FullStack-Bench-Html, we follow the benchmark’s official unit-test protocol and report the pass rate under the same inference setup as other methods. For Web-Bench, we follow its official end-to-end evaluation and report pass@2 across projects. Unless otherwise stated, all reported metrics—including Web-Bench pass@2—are averaged over three runs. Decoding hyperparameters (temperature and top-p) are fixed across systems unless otherwise stated; local fixtures are cached to avoid network variance.

**Baselines and variants.** Our experiments are conducted on two strong instruction-tuned base models: Qwen2.5-7B-Instruct and Llama-3.1-8B-Instruct. On each of these backbones, we compare three distinct approaches: (i) *Base Model* (the frozen instruction-tuned model, serving as a direct baseline), (ii) *Web-RL* (vision-grounded RL using the MLLM reward but without the agentic reflection mechanism), and (iii) *ReLook* (our full framework with agentic MLLM-in-the-loop reflection). All methods use identical inference parameters, prompts, and rendering infrastructure for fair comparison. We also include results for GPT-4o (via OpenAI API) and Qwen2.5-32B-Instruct (local deployment) as reference points representing stronger base models; these are evaluated under the same protocol. We do not compare with prior specialized visual code generation methods (e.g., Design2Code (Si et al., 2024), UICoder (Wu et al., 2024)) because: (i) they were evaluated on different datasets without established cross-benchmark protocols, and (ii) official implementations are not publicly available for controlled comparison on our benchmarks. Our Web-RL baseline serves as a strong vision-aware RL reference that isolates the contribution of agentic reflection.

**Manual validation (GSB).** To complement automated evaluation, we conduct a double-blind human study on 100 randomly sampled tasks com-

paring Qwen2.5-7B-ReLook against Qwen2.5-7B-Instruct. Five independent annotators directly run both systems’ code in our sandboxed renderer and select one of three labels: G (ReLook better), S (same), B (ReLook worse). Majority voting aggregates per-task labels. Results are summarized as G:S:B = 50 : 30 : 20, indicating a clear preference for ReLook under human judgment.

**Implementation Details** We implement ReLook with grouped rollouts under GRPO. The model is trained for a maximum of 40 steps with a training batch size of 256. The learning rate is set to  $1e-6$ , and we employ a linear warmup (first 5%) and cosine decay schedule. For the GRPO loss function, we set the group size  $G$  for rollouts per query to 8 and the clipping parameter  $\epsilon$  to 0.2. We apply a KL penalty toward a reference policy with a non-zero weight  $\beta$  and sweep  $\beta \in [0.01, 0.05]$  (step 0.01). We also sweep the advantage clipping bound in  $\{1, 2, 3\}$  for robustness. The group size and the number of sampled trajectories per query are chosen to fit accelerator memory while maintaining adequate exploration; gradient accumulation is used to emulate larger effective batch sizes. Mixed precision (bfloat16/float16) and activation checkpointing reduce memory footprint. For critic feedback tokens, GRPO is masked out; we optionally add a lightweight distillation loss with weight  $\gamma$  on the feedback tokens to imitate the frozen MLLM’s feedback style. Unless otherwise noted, we use  $\gamma=0.1$  and sweep  $\gamma \in [0.05, 0.3]$  in sensitivity checks. For the length penalty, we set the bounds to  $L_{start} = 12k$  and  $L_{end} = 14k$ .

We use 64 GPUs (32 policy, 32 MLLM), 80 minutes per step, and a curated corpus. Decoding is identical across systems (temp 1.0, top-p 0.7); results average three seeds. Prompts and sandbox are shared. We select checkpoints by mean VisualScore. We focus on 7B/8B backbones; larger-scale RL is future work.

### 3.3 Main Results

Table 1 and Figure 2 present the main results. ReLook achieves substantial improvements over both base models and Web-RL (vision-grounded RL without agentic reflection). Notably, ReLook-w/o-MLLM—which relies solely on internalized reflection without external critic calls—still outperforms Web-RL, demonstrating successful internalization of the refinement mechanism with minimal inference overhead.

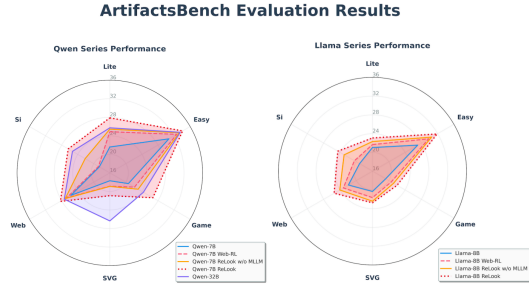


Figure 2: Radar plot showing ReLook’s consistent improvements across all ArtifactsBench subsets for both Qwen2.5-7B and Llama-3.1-8B backbones (averaged over 3 seeds).

Beyond absolute scores, we consistently observe the strictly monotone ordering predicted by our design: ReLook > Web-RL > Base Model. The gap between Web-RL and ReLook underscores the importance of a *vision-aware* training signal coupled with the generate–diagnose–refine loop. Qualitative examples are shown in Appendix Figure 6.

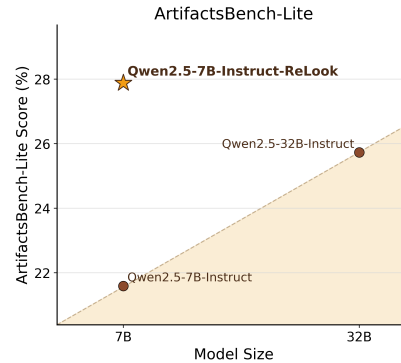


Figure 3: Performance on ArtifactsBench-Lite showing consistent ordering: ReLook > Web-RL > Base Model. Results averaged over 3 seeds.

Figure 3 summarizes performance on the ArtifactsBench-Lite subset. It exhibits the strictly monotone ordering ReLook > Web-RL > Base Model and highlights consistent gains for both Qwen2.5-7B and Llama-3.1-8B backbones, echoing the trends in Table 1. This compact subset mirrors the broader benchmark and provides an intuitive visualization of the average improvements delivered by ReLook.

Figure 4 evidences behavioral collapse in the base model after 2–3 rounds despite feedback, while ReLook improves monotonically. We cap reflections at three rounds for efficiency and keep this at inference.

Figure 5 shows that reflection frequency increases and stabilizes during training, aligning with

Model	A-Lite	A-Easy	A-Game	A-SVG	A-Web	A-Si	FullStack-Bench-Html	Web-Bench
Qwen2.5-32B-Instruct	25.73	33.52	24.36	26.36	27.34	25.30	70.00	10.90
GPT-4o	33.25	34.23	33.04	33.74	34.31	31.44	71.25	23.80
Claude 3.5 Sonnet	39.64	49.37	38.46	41.94	41.26	39.43	81.88	32.39
Llama-3.1-8B-Instruct	21.04	27.11	19.25	20.33	21.91	19.18	57.50	2.50
Llama-3.1-8B-Instruct-Web-RL	21.67	29.80	20.61	21.64	23.15	20.44	61.75	2.75
Llama-3.1-8B-Instruct-ReLook-w/o-MLLM	22.32	30.52	21.18	22.41	24.03	22.97	63.75	2.90
Llama-3.1-8B-Instruct-ReLook	<b>23.08</b>	<b>31.86</b>	<b>22.04</b>	<b>22.74</b>	<b>25.42</b>	<b>24.52</b>	<b>63.75</b>	<b>2.90</b>
Qwen2.5-7B-Instruct	21.59	30.70	20.62	17.70	25.69	18.61	65.00	3.00
Qwen2.5-7B-Instruct-Web-RL	24.89	32.64	22.14	18.87	26.73	18.82	63.25	3.48
Qwen2.5-7B-Instruct-ReLook-w/o-MLLM	25.44	33.29	23.05	18.92	27.11	22.15	67.50	4.20
Qwen2.5-7B-Instruct-ReLook	<b>27.88</b>	<b>34.12</b>	<b>26.72</b>	<b>20.92</b>	<b>28.31</b>	<b>26.36</b>	<b>67.50</b>	<b>4.20</b>

Table 1: Main results on ArtifactsBench subsets (A-Lite/Easy/Game/SVG/Web/Si), FullStack-Bench-Html, and Web-Bench (pass@2). VisualScores follow ARTIFACTSBENCH [0, 100] scale. ReLook uses up to 3 reflection rounds; ReLook-w/o-MLLM relies on internalized self-reflection without external critic. For unit-test benchmarks (FullStack, Web-Bench), both variants use identical critic-free inference, yielding same scores. Bold: best per backbone. Means over 3 seeds (temp=1.0, top-p=0.7).

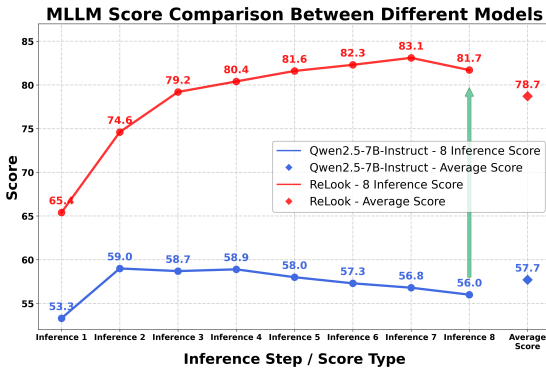


Figure 4: Behavioral collapse mitigation. Base model (Qwen2.5-7B-Instruct) degrades after initial attempts despite MLLM feedback, while ReLook exhibits monotonic improvement across eight forced reflection rounds on ArtifactsBench-Lite. Scores from training-time judge (Qwen2.5-VL-72B).

494 Forced Optimization’s incentive structure. The av-  
495 erage converging to 2 rounds suggests most tasks  
496 reach capacity after two refinements.

497 Due to space limitations, additional analyses re-  
498 garding the scalability of the critic model and per-  
499 formance in out-of-distribution scenarios are pro-  
500 vided in Appendix A.

### 501 3.4 Ablation Study

502 We ablate three components: (i) vision-based  
503 MLLM reward, (ii) format constraint invalidating  
504 non-renderable outputs, and (iii) Forced Optimi-  
505 zation. Table 2 shows each component is critical.  
506 Vision reward provides essential perceptual signal  
507 (+3.3 points). Format constraint prevents reward  
508 hacking (+1.0). Forced Optimization has the largest  
509 impact (+2.0), directly mitigating behavioral col-  
510 lapse. “Single-sample” denotes limiting ReLook to  
511 sample only one instance per reflection round, en-

MLLM Reward	Format Constraint	Forced Optimization	ArtifactsBench-Lite ↑
			21.59
✓			24.89
✓	✓		25.84
✓	✓	Single-sample	26.98
✓	✓	✓	27.88

Table 2: Ablation Study on ArtifactsBench-Lite. Results are averaged over 3 random seeds.

512 suring that the gains do not stem from an increased  
513 generation budget. All ablations use identical seeds  
514 and the external evaluator to avoid bias.

### 515 3.5 Inference speed improvement after 516 removing MLLM

517 We use VLLM to deploy the Qwen2.5-7B-Instruct  
518 model on four H20 GPUs and the Qwen2.5-VL-  
519 72B-Instruct model on eight H20 GPUs. We set the  
520 number of parallel threads to 1 and limit the maxi-  
521 mum number of reflection steps to 3. We conduct  
522 inference on 100 queries and measure the average  
523 inference time per query. ReLook takes an aver-  
524 age of 123.04 seconds per query, whereas ReLook-  
525 w/o-MLLM takes only 18.03 seconds. The results  
526 indicate that removing the screenshot and MLLM  
527 invocation mechanism substantially improves in-  
528 ference efficiency.

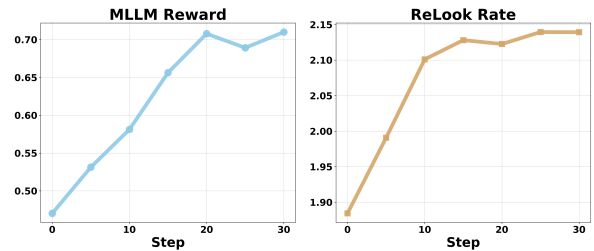


Figure 5: Intermediate Results of RL Training

529	<b>3.6 Qualitative Analysis</b>	576
530	Appendix H presents a qualitative analysis of	577
531	ReLook.	578
532	<b>4 Related Work</b>	579
533	<b>4.1 Visual Code Generation</b>	580
534	Early vision-to-code systems translate static screen-	581
535	shots to HTML/CSS (Wüst et al., 2024) in one shot.	582
536	Recent methods add structure—Design2Code (Si	583
537	et al., 2024), Web2Code (Yun et al., 2024),	584
538	UICoder (Wu et al., 2024), DesignCoder (Chen	585
539	et al., 2025)—but mainly optimize static similarity.	586
540	While these methods improve structural representa-	587
541	tion, they fundamentally lack a closed-loop mech-	588
542	anism to iteratively see, diagnose, and refine their	589
543	own rendered outputs. ReLook addresses this by	590
544	placing a renderer directly in the training loop, us-	591
545	ing vision-grounded rewards from temporal screen-	592
546	shots, and internalizing the refinement process via	593
547	RL.	594
548	<b>4.2 Feedback-Driven Code Reinforcement</b>	595
549	<b>Learning</b>	596
550	RL for program synthesis leverages unit-test re-	597
551	wards and large candidate sets (Le et al., 2022;	598
552	Li et al., 2022); reflective/tool-driven critiques im-	599
553	prove correction (Madaan et al., 2023; Shinn et al.,	600
554	2023; Gou et al., 2023; Peng et al., 2025). How-	601
555	ever, for front-end tasks, these approaches are pixel-	602
556	blind. Even with structured visual instructions,	603
557	their reliance on proxy signals like unit tests or	604
558	linters cannot address defects in visual layout or	605
559	interactive behavior. We couple the policy with a	606
560	direct perceptual reward from temporal screenshots	607
561	and explicitly stabilize this noisy learning process	608
562	via Forced Optimization and a zero-reward rule for	609
563	invalid renders.	610
564	<b>4.3 Acceptance criteria, best-of-N, and</b>	611
565	<b>verifier-assisted search.</b>	612
566	A broad line of work improves generation via external	613
567	selection/search: Codex/AlphaCode sample-	614
568	and-test (Chen et al., 2021; Li et al., 2022);	615
569	self-consistency/tree deliberation aggregate candi-	616
570	dates (Wang et al., 2022; Yao et al., 2023); agen-	617
571	tic self-refinement uses iterative critique (Madaan	618
572	et al., 2023; Shinn et al., 2023; Gou et al., 2023).	619
573	Our <i>Forced Optimization</i> differs in both <i>criterion</i>	620
574	and <i>rule</i> : a vision-grounded score from temporal	621
575	screenshots (not unit tests), and acceptance <i>strictly</i>	622
	requiring monotonic improvement within one tra-	623
	jectory. Unlike best-of- $N$ or offline re-ranking,	
	our rule is <i>online, in-trajectory</i> , preventing regres-	
	sions and reward hacking, yielding stable, visually	
	aligned improvements for front-end code.	
	<b>4.4 Multimodal UI Perception and Evaluation</b>	
	Recent MLLMs ground web elements and lay-	
	outs (OpenAI, 2023; Wang et al., 2024); web-	
	agent/GUI benchmarks show the value of vision-	
	conditioned reasoning (Zhou et al., 2023; Koh et al.,	
	2024; Li et al., 2025). For evaluation, MLLM-as-	
	Judge is common (Ge et al., 2023; Zheng et al.,	
	2023); front-end benchmarks emphasize visual and	
	interactive quality (Xu et al., 2025; Zhang et al.,	
	2025a). Accordingly, we place visual scoring in	
	training to internalize layout/interaction principles,	
	and drop the critic at inference for latency.	
	ReLook unifies these threads via MLLM-based	
	visual rewards within RL, offering a practical path	
	to visually aware, self-improving front-end genera-	
	tion.	
	<b>5 Conclusion</b>	
	We introduce <b>ReLook</b> , a vision-grounded RL	
	framework that establishes a robust generate–	
	diagnose–refine loop for front-end code generation.	
	Specifically, ReLook integrates a multimodal LLM	
	critic with two strategic mechanisms: zero-reward	
	penalties for invalid renders and Forced Optimiza-	
	tion. As a result, it consistently outperforms strong	
	baselines (ReLook > Web-RL > Base) while en-	
	abling efficient, critic-free inference. We anticipate	
	that embedding perception-aligned evaluators in	
	the learning loop will generalize beyond web UIs	
	to broader perceptual programming domains.	
	<b>6 Limitations</b>	
	Addressing complex multi-file project tasks, such	
	as those in Web-Bench, remains a significant chal-	
	lenge. This indicates that traversing beyond the	
	scope of single-artifact refinement to achieve ro-	
	bust long-horizon reasoning will necessitate further	
	architectural innovations. Furthermore, our experi-	
	ments are currently limited to 7B/8B-scale models	
	and front-end development. Therefore, evaluating	
	scalability to larger models and broader software	
	stacks requires additional investigation. Finally,	
	to facilitate reproducibility and support future re-	
	search, we will make our code publicly available	
	upon publication.	

## References

- 624
- 625 Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,  
626 Henrique Ponde De Oliveira Pinto, Jared Kaplan,  
627 Harri Edwards, Yuri Burda, Nicholas Joseph, Greg  
628 Brockman, and 1 others. 2021. Evaluating large  
629 language models trained on code. *arXiv preprint*  
630 *arXiv:2107.03374*.
- 631 Yunnong Chen, Shixian Ding, YingYing Zhang, Wenkai  
632 Chen, Jinzhou Du, Lingyun Sun, and Liuqing Chen.  
633 2025. Designcoder: Hierarchy-aware and self-  
634 correcting ui code generation with large language  
635 models. *arXiv preprint arXiv:2506.13663*.
- 636 Bytedance-Seed-Foundation-Code-Team Yao Cheng,  
637 Jianfeng Chen, Jie Chen, Li Chen, Liyu Chen, Wen-  
638 tao Chen, Zhengyu Chen, Shijie Geng, Aoyan Li,  
639 Bowen Li, Bowen Li, Linyi Li, Boyi Liu, Jerry Liu,  
640 Kaibo Liu, Qi Liu, Shukai Liu, Si-Han Liu, Tianyi  
641 Liu, and 35 others. 2024. [Fullstack bench: Evaluat-  
642 ing llms as full stack coders](#). *ArXiv*, abs/2412.00535.
- 643 Wentao Ge, Shunian Chen, Guiming Hardy Chen,  
644 Junying Chen, Zhihong Chen, Nuo Chen, Wenya  
645 Xie, Shuo Yan, Chenghao Zhu, Ziyue Lin, and  
646 1 others. 2023. Mllm-bench: evaluating multi-  
647 modal llms with per-sample criteria. *arXiv preprint*  
648 *arXiv:2311.13951*.
- 649 Zhibin Gou, Zhihong Shao, Yeyun Gong, Yelong  
650 Shen, Yujiu Yang, Nan Duan, and Weizhu Chen.  
651 2023. Critic: Large language models can self-correct  
652 with tool-interactive critiquing. *arXiv preprint*  
653 *arXiv:2305.11738*.
- 654 Jing Yu Koh, Robert Lo, Lawrence Jang, Vikram  
655 Duvvur, Ming Chong Lim, Po-Yu Huang, Graham  
656 Neubig, Shuyan Zhou, Ruslan Salakhutdinov, and  
657 Daniel Fried. 2024. Visualwebarena: Evaluating mul-  
658 timodal agents on realistic visual web tasks. *arXiv*  
659 *preprint arXiv:2401.13649*.
- 660 Hung Le, Yue Wang, Akhilesh Deepak Gotmare, Silvio  
661 Savarese, and Steven Chu Hong Hoi. 2022. Coderl:  
662 Mastering code generation through pretrained models  
663 and deep reinforcement learning. *Advances in Neural*  
664 *Information Processing Systems*, 35:21314–21328.
- 665 Kaixin Li, Ziyang Meng, Hongzhan Lin, Ziyang Luo,  
666 Yuchen Tian, Jing Ma, Zhiyong Huang, and Tat-Seng  
667 Chua. 2025. Screenspot-pro: Gui grounding for  
668 professional high-resolution computer use. *arXiv*  
669 *preprint arXiv:2504.07981*.
- 670 Yujia Li, David Choi, Junyoung Chung, Nate Kushman,  
671 Julian Schrittwieser, Rémi Leblond, Tom Eccles,  
672 James Keeling, Felix Gimeno, Agustin Dal Lago, and  
673 1 others. 2022. Competition-level code generation  
674 with alphacode. *Science*, 378(6624):1092–1097.
- 675 Xinyu Liu, Shuyu Shen, Boyan Li, Peixian Ma, Runzhi  
676 Jiang, Yuxin Zhang, Ju Fan, Guoliang Li, Nan Tang,  
677 and Yuyu Luo. 2024. A survey of nl2sql with large  
678 language models: Where are we, and where are we  
679 going? *arXiv preprint arXiv:2408.05109*.
- LMSYS Org. 2024. Chatbot Arena Leaderboard.  
<https://web.lmarena.ai/leaderboard>. Ac-  
cessed: 2024-05-23.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler  
Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,  
Nouha Dziri, Shrimai Prabhumoye, Yiming Yang,  
and 1 others. 2023. Self-refine: Iterative refinement  
with self-feedback. *Advances in Neural Information*  
*Processing Systems*, 36:46534–46594.
- OpenAI. 2023. Gpt-4v(ision) system card.  
[https://cdn.openai.com/papers/GPTV\\_  
System\\_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf).
- Zhongyuan Peng, Yifan Yao, Kaijing Ma, Shuyue  
Guo, Yizhe Li, Yichi Zhang, Chenchen Zhang, Yi-  
fan Zhang, Zhouliang Yu, Luming Li, and 1 others.  
2025. Criticlean: Critic-guided reinforcement learn-  
ing for mathematical formalization. *arXiv preprint*  
*arXiv:2507.06181*.
- Raphael Rafailov, Archit Sharma, Eric Mitchell, Stefano  
Ermon, Chelsea Finn, and Christopher D. Manning.  
2023. [Direct preference optimization: Your language  
model is secretly a reward model](#). In *Advances in*  
*Neural Information Processing Systems*.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec  
Radford, and Oleg Klimov. 2017. [Proximal policy  
optimization algorithms](#). In *Proceedings of the 34th*  
*International Conference on Machine Learning Work-  
shop*.
- Noah Shinn, Federico Cassano, Ashwin Gopinath,  
Karthik Narasimhan, and Shunyu Yao. 2023. Re-  
flexion: Language agents with verbal reinforcement  
learning. *Advances in Neural Information Process-  
ing Systems*, 36:8634–8652.
- Chenglei Si, Yanzhe Zhang, Ryan Li, Zhengyuan Yang,  
Ruibo Liu, and Diyi Yang. 2024. Design2code:  
Benchmarking multimodal code generation for au-  
tomated front-end engineering. *arXiv preprint*  
*arXiv:2403.03163*.
- Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhi-  
hao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin  
Wang, Wenbin Ge, and 1 others. 2024. Qwen2-  
vl: Enhancing vision-language model’s perception  
of the world at any resolution. *arXiv preprint*  
*arXiv:2409.12191*.
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V.  
Le, Ed H. Chi, Sharan Narang, Aakanksha Chowd-  
hery, and Denny Zhou. 2022. Self-consistency im-  
proves chain of thought reasoning in language mod-  
els. *arXiv preprint arXiv:2203.11171*.
- Jason Wu, Eldon Schoop, Alan Leung, Titus Barik, Jef-  
frey P Bigham, and Jeffrey Nichols. 2024. Uicoder:  
Finetuning large language models to generate user  
interface code through automated feedback. *arXiv*  
*preprint arXiv:2406.07739*.

734 Antonia Wüst, Wolfgang Stammer, Quentin Delfosse,  
735 Devendra Singh Dhimi, and Kristian Kersting. 2024.  
736 Pix2code: Learning to compose neural visual concepts  
737 as programs. *arXiv preprint arXiv:2402.08280*.

738 Kai Xu, YiWei Mao, XinYi Guan, and ZiLong Feng.  
739 2025. Web-bench: A llm code benchmark based  
740 on web standards and frameworks. *arXiv preprint*  
741 *arXiv:2505.07473*.

742 An Yang, Baosong Yang, Binyuan Hui, Bo Zheng,  
743 Bowen Yu, Chang Zhou, Chengpeng Li, Chengyuan  
744 Li, Dayiheng Liu, Fei Huang, and 1 others.  
745 2024. Qwen2 technical report. *arXiv preprint*  
746 *arXiv:2407.10671*.

747 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran,  
748 Karthik Narasimhan, and Yuan Cao. 2023. Tree of  
749 thoughts: Deliberate problem solving with large language  
750 models. *arXiv preprint arXiv:2305.10601*.

751 Sukmin Yun, Haokun Lin, Rusiru Thushara, Mo-  
752 hammad Qazim Bhat, Yongxin Wang, Zutao  
753 Jiang, Mingkai Deng, Jinhong Wang, Tianhua Tao,  
754 Junbo Li, and 1 others. 2024. Web2code: A  
755 large-scale webpage-to-code dataset and evaluation  
756 framework for multimodal llms. *arXiv preprint*  
757 *arXiv:2406.20098*.

758 Chenchen Zhang, Yuhang Li, Can Xu, Jiaheng Liu,  
759 Ao Liu, Shihui Hu, Dengpeng Wu, Guanhua Huang,  
760 Kejiao Li, Qi Yi, and 1 others. 2025a. Artifactsbench:  
761 Bridging the visual-interactive gap in llm code generation  
762 evaluation. *arXiv preprint arXiv:2507.04952*.

763 Chenchen Zhang, Jinxiang Xia, Jiaheng Liu, Wei Zhang,  
764 Yejie Wang, Jian Yang, Ge Zhang, Tianyu Liu,  
765 Zhongyuan Peng, Yingshui Tan, Yuanxing Zhang,  
766 Zhexu Wang, Weixun Wang, Yancheng He, Ken  
767 Deng, Wangchunshu Zhou, Wenhao Huang, and  
768 Zhaoxiang Zhang. 2025b. [Codecriticbench: A holistic code critique benchmark for large language models](#). *Preprint*, arXiv:2502.16614.

771 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan  
772 Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,  
773 Zhuohan Li, Dacheng Li, Eric Xing, and 1 others.  
774 2023. Judging llm-as-a-judge with mt-bench and  
775 chatbot arena. *Advances in Neural Information Processing*  
776 *Systems*, 36:46595–46623.

777 Shuyan Zhou, Frank F Xu, Hao Zhu, Xuhui Zhou,  
778 Robert Lo, Abishek Sridhar, Xianyi Cheng, Tianyue  
779 Ou, Yonatan Bisk, Daniel Fried, and 1 others.  
780 2023. Webarena: A realistic web environment  
781 for building autonomous agents. *arXiv preprint*  
782 *arXiv:2307.13854*.

783	<b>A Additional experimental results</b>	
784	<b>Impact of the Critic Model.</b> With the rapid evolution of multimodal technology, more advanced models such as Qwen3-VL demonstrate superior visual and code comprehension capabilities compared to Qwen2.5-VL. To investigate the impact of the critic’s quality on the training outcome, we substitute the original Qwen2.5-VL-72B-Instruct critic with Qwen3-VL-8B-Instruct and Qwen3-VL-32B-Instruct, while employing Qwen2.5-7B-Instruct-ReLook as the base policy model. The comparative results are reported in Table 3. It can be observed that as the capability of the Critic Model strengthens, the performance of the final model also gradually improves.	830 831 832 833 834 835 836 837
798	<b>Evaluation on OOD Scenarios.</b> To further assess the robustness of the model in out-of-distribution (OOD) environments, we evaluate ReLook on LiveCodeBench v6 (2025-02 ~ 2025-05). Qwen2.5-7B achieves a score of 0.213, compared to 0.198 for ReLook. This demonstrates that ReLook maintains competitive performance in OOD scenarios.	
806	<b>Validation on Different Base Models.</b> We additionally validate our framework by instantiating the ReLook Agent with Qwen3-VL-8B-Instruct. Leveraging the model’s superior code understanding and visual capabilities, the agent achieves a competitive score of 35.83 on A-Lite. This experiment underscores the framework’s inherent scalability and its ability to generalize effectively to diverse baseline models.	
815	<b>B External validity and cross-judge analysis</b>	
817	We align our evaluator setup with ARTIFACTSBENCH(Zhang et al., 2025a), which establishes validity through rigorous empirical validation:	
820	<b>Human-MLLM Agreement.</b> ARTIFACTSBENCH conducted a controlled human study with expert web developers evaluating a stratified sample of 200 tasks. Results showed over 90% agreement (Cohen’s $\kappa > 0.85$ ) between human judgments and MLLM evaluators (Gemini-2.5-Pro and Qwen2.5-VL-72B) on visual fidelity, layout correctness, and interactive integrity. The study used a double-blind protocol with three independent human raters per task.	
	<b>Crowdsourced Validation.</b> Beyond controlled studies, ARTIFACTSBENCH validated MLLM judges against WebDev Arena(LMSYS Org, 2024), a large-scale platform with over 10,000 pairwise comparisons from web developers. The MLLM rankings showed strong correlation (Spearman $\rho = 0.87$ for Gemini-2.5-Pro, $\rho = 0.83$ for Qwen2.5-VL-72B) with crowd preferences.	838 839 840 841 842 843 844 845 846 847 848
	<b>Our Protocol Adherence.</b> Since our test sets (A-Lite, A-Easy, A-Game, A-SVG, A-Web, A-Si) are strict subsets of ARTIFACTSBENCH’s 1,825 human-validated tasks, we directly inherit this established validity. We use the official evaluation scripts, judge prompts, and scoring rubric without modification. To mitigate on-policy overfitting, we decouple training-time critic (Qwen2.5-VL-72B-Instruct) from offline evaluator (Gemini-2.5-Pro). We do not re-run human studies; validity is anchored in ARTIFACTSBENCH’s reported evidence.	
	<b>Score Interpretation.</b> The absolute VisualScore range (20-30 for 7B/8B models) reflects benchmark difficulty: ARTIFACTSBENCH tasks span complex interactions, dynamic animations, and pixel-perfect layouts. Reference models (GPT-4o: 33, Qwen2.5-32B: 26) establish that even frontier systems find these tasks challenging. The [0, 100] scale provides fine-grained discrimination; we report relative improvements over baselines.	849 850 851 852 853 854 855 856 857
	<b>C Pseudocode for ReLook</b>	858
	<b>D Evaluator rubric and scoring range</b>	859
	We use a fixed evaluation rubric for VISUALSCORE. During training, the critic’s visual score used for RL is normalized to [0, 1] for stability. For offline reporting and tables, we follow ARTIFACTSBENCH and report evaluator outputs on a [0, 100] scale; all numbers in the main results tables and figures are on this [0, 100] scale unless otherwise specified. The rubric jointly considers: (i) specification adherence; (ii) layout alignment and spatial fidelity; (iii) typography and color coherence; and (iv) interactive integrity for tasks specifying actions. For dynamic behavior, each reflection round is evaluated on three temporal screenshots jointly in a single scoring call to obtain one round-level score; trajectory-level reward averages round-level scores. We do not perform multi-judge averaging or smoothing.	860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876

Critic Model	A-Lite	A-Easy	A-Game	A-SVG	A-Web	A-Si	FS-Html	Web-Bench
Qwen2.5-VL-72B-Instruct	27.88	34.12	26.72	20.92	28.31	26.36	67.50	4.20
Qwen3-VL-8B-Instruct	28.96	34.02	28.43	21.77	28.84	27.94	68.75	5.48
Qwen3-VL-32B-Instruct	<b>32.95</b>	<b>35.47</b>	<b>31.95</b>	<b>29.75</b>	<b>34.02</b>	<b>30.83</b>	<b>72.50</b>	<b>8.35</b>

Table 3: Ablation study on the Critic Model. We compare the performance of the base model (Qwen2.5-7B-Instruct-ReLook) when trained with guidance from different critic models. The results demonstrate that stronger critics lead to better performance.

## E The loss function for Forced Optimization

Faced with the "behavior collapse" problem, we first attempted a negative-reward penalty comparing post-feedback scores to pre-feedback scores. If an optimized answer was worse than its predecessor, a negative reward was applied. We observed that this encouraged the model to reduce reflection frequency to avoid penalties (reward hacking), degrading ReLook into regular RL and harming performance. Therefore, we removed this mechanism and adopted the strict acceptance rule described in Method: only strictly improving steps are accepted into trajectories.

## F Lightweight distillation loss on feedback tokens

Let  $M$  denote the index set of critic feedback tokens  $m$  within a trajectory, and let the frozen MLLM critic define a teacher distribution  $q_t(\cdot)$  at each position  $t \in M$ . The student policy distribution is  $p_\theta(\cdot | x_{\leq t})$ . The lightweight distillation loss used to transfer the feedback style for test-time self-reflection is

$$\mathcal{L}_{\text{distill}}^m := \frac{1}{|M|} \sum_{t \in M} \text{KL}(q_t(\cdot) \| p_\theta(\cdot | x_{\leq t})) \quad (6)$$

equivalently, as a teacher-forced cross-entropy

$$\mathcal{L}_{\text{distill}}^m := -\frac{1}{|M|} \sum_{t \in M} \sum_v q_t(v) \log p_\theta(v | x_{\leq t}) \quad (7)$$

Only policy-controlled tokens ( $t, c$ ) contribute non-zero advantages in GRPO; feedback tokens  $m$  are masked in the RL objective (advantages set to zero). The total training objective is

$$\mathcal{L} = -\mathcal{J}_{\text{GRPO}}^{t,c} + \gamma \mathcal{L}_{\text{distill}}^m, \quad (8)$$

with default  $\gamma=0.1$ ,  $\gamma \in [0.05, 0.3]$ .

## G Reproducibility notes: GRPO and implementation details

**Reference policy and KL.** We regularize toward a fixed reference policy  $\pi_{ref}$  (the frozen instruction-tuned backbone before RL). KL is computed token-wise over policy-controlled tokens with weight  $\beta$  (see ranges in Experiment); we do not anneal  $\beta$  within a step.

**Trajectory composition and masking.** Trajectories mix policy tokens ( $t, c$ ) and critic tokens ( $m$ ). During optimization, advantages on  $m$  are masked to zero; optional lightweight distillation on  $m$  uses a small KL/CE with weight  $\gamma$  to the frozen MLLM outputs.

**Advantage baseline and clipping.** We use a group-relative baseline (mean return over  $G$  trajectories per query); advantages are standardized within-batch and clipped to  $[-2, 2]$ .

**Sampling limits and rejection handling.** For Forced Optimization, we cap resampling attempts at 10 per reflection round to avoid infinite loops; non-improving proposals are rejected without re-scoring. When the 10-attempt limit is reached without improvement, we terminate further reflection rounds for that trajectory and use the best-so-far result. This limit balances exploration (allowing sufficient attempts to find improving revisions) with computational efficiency.

**Length penalty parameters.** We set  $L_{\text{start}} = 12,000$  and  $L_{\text{end}} = 14,000$  tokens based on empirical analysis of typical front-end code lengths in our training corpus, ensuring reasonable generation length while discouraging degenerate repetition.

**Valid Render Rate dynamics.** During training, the Valid Render Rate increases from approximately 40% at initialization to approximately 80% upon convergence, demonstrating that the model learns to produce syntactically valid and renderable code through the combination of the visual-output

---

**Algorithm 1** ReLook Training Framework

---

**Require:** Task specification  $q$ , base policy  $\pi_\theta$ , vision critic MLLM (Qwen2.5-VL-72B), max steps  $S$ , max reflections  $R$ , max resampling  $K=10$

**Ensure:** Optimized policy  $\pi_\theta^*$  with internalized visual cognition

```
1: Initialize:  $s \leftarrow 0$ ,  $history \leftarrow [q]$ 
2: while  $s < S$  do
3:   Initialize:  $r \leftarrow 0$ ,  $o \leftarrow ""$ ,  $s_{prev} \leftarrow -1$ 
4:   while  $r < R$  and room for improvement do
5:      $k \leftarrow 0$ ,  $accepted \leftarrow \text{False}$ 
6:     while  $k < K$  and not  $accepted$  do
       {Forced Optimization: resample up to  $K$  times} Generate:
        $t_r, c_r \leftarrow \pi_\theta(\text{next token} | history)$ 
       Extract code  $c_r$  from <answer> block
       Attempt rendering  $c_r \rightarrow$  Capture screenshots
        $\{S_1, S_2, S_3\}$  Generate feedback  $m_r \leftarrow \text{MLLM}(q, c_r, \{S_i\})$ 
       Compute visual score  $s_r \leftarrow \text{VisualScore}(q, c_r, \{S_i\})$ 
       if  $s_r > s_{prev}$  then {Accept only strictly improving steps}
10:      Append  $t_r, c_r$  to  $history$  and  $o$ 
14:      Append  $m_r$  to  $history$  and  $o$  within <mlm_feedback>
15:       $s_{prev} \leftarrow s_r$ ,  $accepted \leftarrow \text{True}$ 
16:    end if
17:     $k \leftarrow k + 1$ 
18:  end while
19:  if not  $accepted$  then
       {No improvement after  $K$  attempts}
       Break {Terminate further reflections, use best-so-far}
20:21: end if
22:   $r \leftarrow r + 1$ 
23: end while
24: Calculate reward components for the trajectory  $o$ :
25:    $R_{\text{MLLM}}(o) = \begin{cases} \text{VisualScore}(o) & \text{if screenshot valid} \\ 0 & \text{otherwise} \end{cases}$ 
26:    $R_{\text{len}}(o)$  as in Method (linear penalty:  $L_{\text{start}}=12,000$ ,  $L_{\text{end}}=14,000$ ).
27:    $R_{\text{ReLook}}(o) = R_{\text{MLLM}}(o) \cdot R_{\text{len}}(o)$ .
28: Perform policy update using GRPO to optimize parameters  $\theta$  based on  $R_{\text{ReLook}}$ .
29: If convergence criterion is met (no significant improvement), exit loop.
30:   $s \leftarrow s + 1$ 
31: end while
32: return  $\pi_\theta^*$  for deployment {Optionally discard MLLM during inference.}
```

---

constraint in prompts and the zero-reward penalty for invalid renders. 947 948

**Decoding and seeds.** Unless otherwise noted, decoding uses identical hyperparameters across systems (temperature=1.0, top-p=0.7) with three random seeds; fixtures are cached to avoid network variance. These notes complement Section 2 and Section 3 to facilitate faithful reimplementa- 949 950 951 952 953 954

## H Qualitative Analysis 955

Figure 6 presents a qualitative comparison between the baseline model and ReLook across several representative front-end generation tasks. 956 957 958

Prompt column lists natural-language instructions of varying complexity, including layout composition, template library rendering, login/registration forms, and a chessboard. 959 960 961 962

Base Model column shows the outputs of an instruction-tuned baseline. While it can produce code that renders without error, the resulting webpages often suffer from issues such as layout drift, incomplete functionality, and lack of visual coherence (e.g., missing interactivity in the login page, overly simplistic rendering of the chessboard). 963 964 965 966 967 968 969

ReLook column demonstrates the effect of our vision-grounded reinforcement learning framework. By incorporating visual feedback into training, ReLook produces outputs that are not only executable but also visually faithful and functionally aligned with the prompts. For instance, the template library page is correctly populated with clickable cards, the login/registration form has a clean layout and interactive elements, and the chessboard renders with precise alignment. 970 971 972 973 974 975 976 977 978 979

Overall, the comparison highlights how ReLook systematically reduces layout drift, strengthens interaction correctness, and achieves higher visual fidelity compared to the baseline. 980 981 982 983

## I Template prompt for ReLook rollout 984

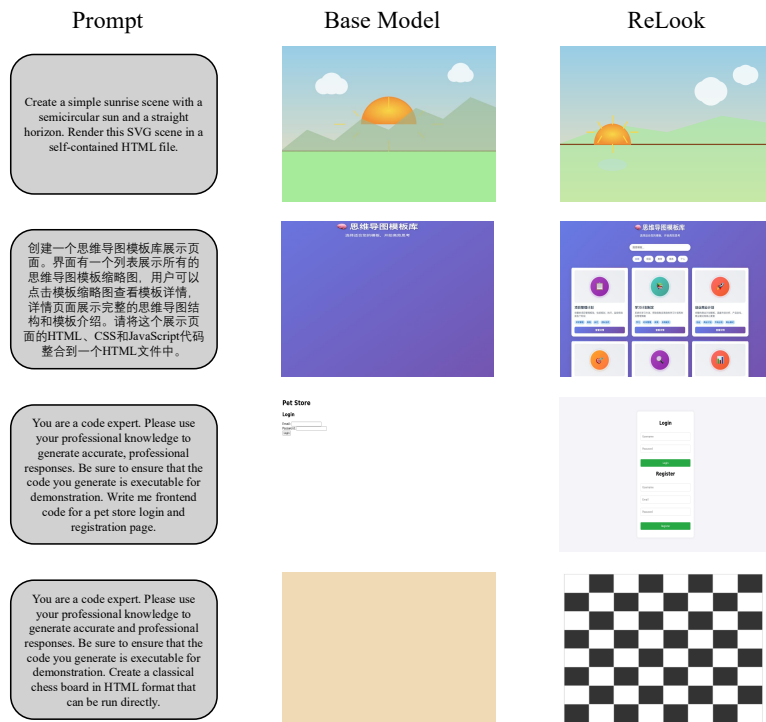


Figure 6: Visual Comparison of Frontend Websites Generated by Baseline and ReLook.

Solve the following problem step by step.  
 You now have the ability to selectively write executable HTML, CSS, JavaScript, or SVG code to receive feedback from the multimodal large model on the code.  
 The code you provided will be executed, and the feedback (wrapped in '<mlm\_feedback> output\_str </mlm\_feedback>') can be returned to aid your reasoning and help you arrive at the final answer.  
 Unless you believe the current answer is flawless, please output <get\_feedback> after providing the complete answer to receive feedback from the multimodal large model and improve the code based on the feedback.  
 \*user question:\*  
 {\$query}

Figure 7: Template prompt for ReLook rollout.