

RealGen: Photorealistic Text-to-Image Generation via Detector-Guided Rewards

Junyan Ye^{1,2*}, Leiqi Zhu^{1,3*}, Yuncheng Guo¹, Dongzhi Jiang⁴,
Zilong Huang², Yifan Zhang⁵, Zhiyuan Yan⁶, Haohuan Fu⁵, Conghui He¹, Weijia Li^{2†}
¹Shanghai AI Lab, ²Sun Yat-Sen University, ³Nanjing University,
⁴CUHK MMLab, ⁵Tsinghua University, ⁶Peking University



Figure 1. **Images generated by our proposed RealGen.** RealGen achieves superior photorealism and enhanced details, outperforming both powerful T2I models, such as Qwen-Image, and specialized photorealistic models, like FLUX-Krea.

Abstract

With the continuous advancement of image generation technology, advanced models such as GPT-Image-1 and Qwen-Image have achieved remarkable text-to-image consistency and world knowledge. However, these models still fall short in photorealistic image generation. Even on simple T2I tasks, they tend to produce "fake" images with distinct AI artifacts, often characterized by "overly smooth skin" and "oily facial sheens". To recapture the original goal of "indistinguishable-from-reality" generation, we propose RealGen, a photorealistic text-to-image framework. RealGen integrates an LLM component for prompt optimization and a diffusion model for realistic image generation. Inspired by adversarial generation, RealGen introduces a "Detector Reward" mechanism, which quantifies artifacts and assesses realism using both semantic-level and

feature-level synthetic image detectors. We leverage this reward signal with the GRPO algorithm to optimize the entire generation pipeline, significantly enhancing image realism and detail. Furthermore, we propose RealBench, an automated evaluation benchmark employing Detector-Scoring and Arena-Scoring. It enables human-free photorealism assessment, yielding results that are more accurate and aligned with real user experience. Experiments demonstrate that RealGen significantly outperforms general models like GPT-Image-1 and Qwen-Image, as well as specialized photorealistic models like FLUX-Krea, in terms of realism, detail, and aesthetics. The code is available at <https://github.com/yejy53/RealGen>.

1. Introduction

Image generation has undergone a significant evolution from GANs [11] to Diffusion models [9, 29], leading to a genera-

*Equal contribution. †Corresponding author.

tion of powerful models such as GPT-Image-1 [23], Nano-Banana [6], and Qwen-Image [38]. These models demonstrate exceptional capabilities in areas like precise attribute control and complex text rendering. However, recent research has focused on enhancing prompt fidelity and leveraging world knowledge for complex content generation. Despite these advancements, even existing powerful generative models like FLUX.1 Pro [16] and Qwen-Image [38] still tend to produce images that lack photorealism, such as "overly smooth skin" and "oily facial sheens", as shown in Fig. 2. This arguably deviates from the original goal of image generation: *"to produce visuals that are indistinguishable from reality,"* which is precisely the focus of our work.

To pursue higher photorealism, existing research has attempted to combine generative models with reinforcement learning (RL) and human preference scores [30, 32]. For instance, methods like DanceGRPO [43] and FlowGRPO [21] utilize reward models such as PickScore [15] and HPSv2 [40] to align generation with human preferences. However, trained on human preference data may introduce prior biases, such as an excessive preference for colors like red or purple. More critically, human preference scores do not directly equate to photorealism. A candid snapshot might be highly realistic but receive a low score due to a lack of aesthetic appeal, inadvertently encouraging the model to favor artistic or anime styles. Meanwhile, approaches like FLUX-Krea [18] employ large-scale, human-curated, high-quality image samples combined with TPO to enhance realism. However, this method is not only cost-prohibitive but also highly dependent on the subjective preferences of annotators. Therefore, the critical bottleneck in leveraging RL for photorealism lies in *how to establish an objective, scalable, and human-free metric for image realism.*

Inspired by adversarial generation, a promising approach is to quantify photorealism using synthetic image detectors. The core rationale is that *a more realistic image should possess fewer AI artifacts and thus be more difficult for a detector to classify as "Fake."* In recent years, as generative models have rapidly advanced, corresponding image detectors have evolved in parallel, achieving robust detection performance. Furthermore, some MLLM-based methods have extended this capability to semantic-level, interpretable detection, such as analyzing "AI-feel" skin textures. Therefore, we posit that two distinct classes of detectors can be employed to measure image realism: one for analyzing visible, semantic-level artifacts, and another for deep feature-level artifacts. As shown in Figure 2, during the RL process, we define these detectors as the reward model to guide the generator to "escape" detection, which effectively reduces artifacts and enhances image realism.

In addition, the quality of image synthesis is highly correlated with the complexity of the text prompt [2]. Expert users, through sophisticated "prompt engineering," can pro-

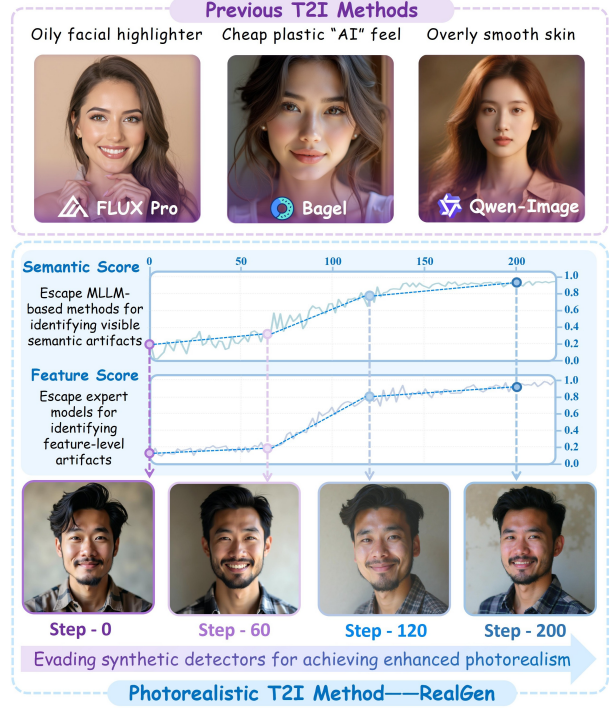


Figure 2. **From Synthetic Artifacts to Photorealism.** Contrasting the common "fake-feel" AI artifacts in previous T2I methods, our proposed RealGen achieves enhanced photorealism by progressively evading semantic and feature-level detectors.

duce photorealistic, cinematic-quality images. Conversely, simple prompts from casual users often yield low-quality results that lack realism. This discrepancy arises because T2I models often rely on complex prompt structures; simple prompts provide low information entropy, causing the model to default to its high-frequency priors, which results in images lacking detail and specificity [55]. Therefore, automatically rewriting or enriching the user's input directive presents another critical pathway to enhancing the quality and realism of the subsequently generated images [33].

In summary, we propose **RealGen**, a framework for photorealistic text-to-image generation. The framework includes a Large Language Model (LLM) component for optimizing user prompts and a diffusion model dedicated to generating realistic images. We utilize detector models as rewards, optimizing both the LLM and the generative model via Generalized Reinforcement Policy Optimization (GRPO) algorithm [21, 43] to effectively reduce artifacts in the generated images. Specifically: (1) With the generative model held fixed, we train the LLM to optimize the input prompts, using the detector model to score the final synthesized image; this encourages the LLM to generate richer and more effective prompts. (2) While maintaining text input consistency, we optimize the diffusion model itself, enabling it to "escape" the detectors and generate more realistic and detailed images.

Furthermore, we also propose **RealBench**, a new benchmark for evaluating the photorealism of generated images.

RealBench contains a diverse set of text prompts, supplemented by high-quality, real-world image references. Addressing the lack of effective evaluation methods beyond manual human assessment, RealBench employs two automated evaluation protocols: Detector-Scoring and Arena-Scoring. First, Detector-Scoring utilizes multiple synthetic image detectors (including those held out from the reward) to score the images; images with fewer AI artifacts receive higher scores. Second, for Arena-Scoring, we adapt the LMArena methodology, using several different MLLMs to simulate user preferences in pairwise comparisons, selecting the result that appears more realistic. Outputs from each model are put into pairwise "battles" against other models or real-world images for at least 3000 random pairings to determine a final win rate. The inclusion of real-world data not only enhances scoring stability but also validates the effectiveness of this adversarial scoring approach. Our main contributions are as follows:

- We propose **RealGen**, a text-to-image generator capable of producing highly convincing photorealistic images. It leverages a Detector Reward-guided GRPO post-training to escape detector identification, thereby reducing artifacts and enhancing image realism and detail.
- We introduce **RealBench**, a new benchmark for evaluating photorealism that achieves human-free automated scoring through Detector-Scoring and Arena-Scoring.
- RealGen significantly outperforms both general image models (like GPT-Image-1, Qwen-Image) and specialized realistic models (like FLUX-Krea) in realism, details, and aesthetics on the T2I task.

2. Related Work

2.1. Image Generation Models

Image generation models have made significant progress in recent years [1, 3, 8, 46], demonstrating exceptional performance across diverse downstream tasks [19, 48, 50]. Representative models such as Stable Diffusion [9, 26, 28], FLUX [16, 17], Emu [5, 31, 36] and DALL-E [27] demonstrating powerful text-to-image generation capabilities. As research gradually shifts toward multimodal generative models, these models achieve understanding and generation through unified architectures [4, 25, 42, 45]. Although powerful models like GPT-Image-1 [23], Bagel [7] and OmniGen2 [39] can generate precise objects and complex text, challenges remain in generating realistic images, especially with the strong oily appearance of human faces. RealGen focuses on enhancing realistic image generation, optimizing image realism and detail.

2.2. Reinforcement Learning and Human Preferences in Image Generation

Recently, with the development of reinforcement learning, more methods have begun to explore its application in the Diffusion image generation domain [43, 45]. For example, DiffusionDPO [32] and FlowGRPO [21] use human preference data as a reward function to generate images that align with human aesthetics. However, reward models based on human preferences can introduce prior biases, such as HPSv2 [40] favoring red-toned images and PickScore [15] preferring purple ones. While human-captured images are generally realistic, they do not always score highly in terms of preference. SRPO [30] has made some progress in enhancing image realism but still lacks in aesthetic quality. FLUX-Krea [18], through large-scale manual collection of high-quality image samples and TPO reinforcement learning, optimizes the realism of generated images. However, this method is limited by the costs of data labeling and the personal preferences of annotators. In contrast, RealGen does not rely on human preference data. Instead, it uses advanced AIGC detection models to guide the generative model away from AI artifacts at both semantic and feature levels, further improving image realism.

2.3. Synthesis detection for Image generation

As the realism of image generation continuously increases, corresponding detection techniques have rapidly evolved [20, 49]. From early detectors like CNNSpot [35] advanced methods leveraging powerful pre-trained vision models (e.g., OmniAID [12], Effort [44]), these detectors have demonstrated success rates exceeding 80% [49, 56]. Recently, detector based on MLLMs have emerged, such as FakeVLM [37] and LEGION [14], which not only achieve good detection performance but also enhance the explainability of artifact detection. Since synthetic image detection and generation are a "cat-and-mouse" game, some studies have explored using detector feedback to optimize generation quality. However, these methods are currently more limited to post-processing techniques, such as Inpainting, to optimize already generated images [53, 54]. In contrast, our work combines reinforcement learning and synthetic detectors to directly optimize the generative model itself.

3. RealGen

3.1. Model architecture

As illustrated in Figure 3, the architecture of RealGen comprises two core components: a LLM for understanding and refining user intent, and a diffusion model for realistic image synthesis. First, the LLM receives the initial user instruction and performs "thought and planning": it expands the short prompt into a longer, more diverse text description by adding rich details. Subsequently, the diffusion model utilizes this

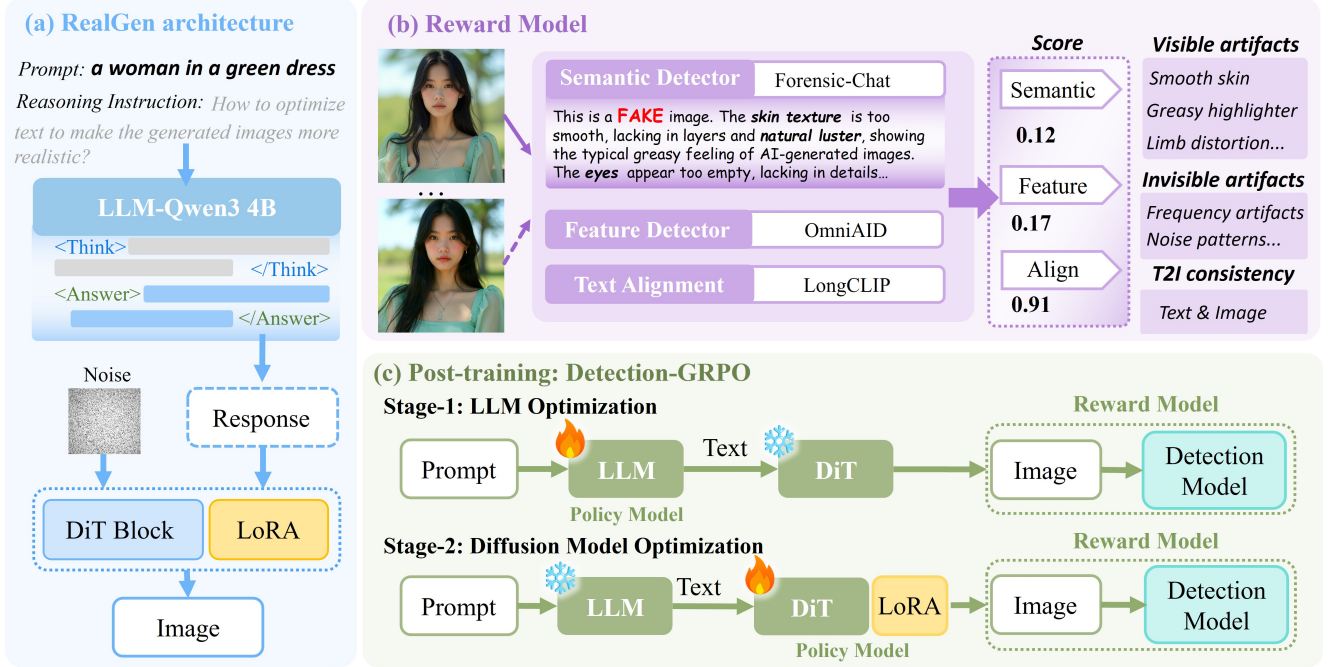


Figure 3. **Overview of the RealGen Method.** (a) The architecture of RealGen, consisting of an LLM component and a Diffusion component. (b) Our detector-based reward model, which evaluates images based on visible artifacts, feature-level artifacts, and text-image alignment. (c) The two-stage post-training process guided by this reward model, which respectively optimizes the LLM and Diffusion components.

refined text as its input condition, executing the denoising and decoding process to generate the final image.

For implementation, we employ Qwen-3 4B [47] as the base LLM. For image generation, we utilize the advanced pre-trained diffusion model, FLUX.1-dev dev [16], integrated with fine-tuned LoRA layers. Both the LLM and the diffusion model first undergo specialized Supervised Fine-Tuning (SFT) as a cold-start phase, followed by Reinforcement Learning optimization via GRPO. The subsequent sections will detail our designed detection reward function and the RL training process.

3.2. Detection Reward

To steer the model optimization towards high-fidelity realism, the design of the reward function is critical, as it must accurately quantify authenticity. We adopt a "detection-as-reward" paradigm inspired by adversarial generation, designing a multi-objective reward function. As shown in Fig. 3(b), this function combines detectors at two distinct levels—semantic and feature—to penalize both perceptible artifacts and imperceptible synthesis traces, respectively.

Semantic Detector: We employ Forensic-Chat [20], a generalizable and interpretable detector optimized from Qwen2.5-VL-7B [34]. It assesses authenticity by analyzing image content (e.g., smooth greasy skin, artifacts in faces/hands, unnatural background blur). We define the semantic reward R_{semantic} as the normalized probability of its

output "real" token probability:

$$R_{\text{semantic}} = \text{softmax}([L(\text{"fake"}, \text{"Fake"}), L(\text{"real"}, \text{"Real"})])_1 \quad (1)$$

Feature Detector: We utilize the advanced expert detector OmniAID [12], which achieves stable and accurate detection by being pre-trained on large-scale real and synthetic datasets. Feature-level artifacts are primarily associated with frequency artifacts and abnormal noise patterns. We define the feature-level reward (R_{feature}) as one minus its output "fake" probability, i.e., $1 - P_{\text{OmniAID}}(\text{fake})$.

Text Alignment: Finally, we include the Long-CLIP [52] score as an auxiliary reward (R_{align}) to maintain text-image alignment. This prevents the model from sacrificing fidelity to the input prompt in its pursuit of realism.

For the GRPO process, we combine the three reward functions defined above: $\{R_{\text{semantic}}, R_{\text{feature}}, R_{\text{align}}\}$. For a batch of N samples $\{I_i\}_{i=1}^N$, we first evaluate each sample I_i to obtain its raw scores $\{r_i^{\text{sem}}, r_i^{\text{feat}}, r_i^{\text{align}}\}$. Since each reward function operates on a different scale and distribution, we first normalize the scores for each reward dimension within the batch. Then, we sum these normalized scores to fuse them into a single advantage function $A(I_i)$, as defined in Equation (2):

$$A(I_i) = \sum_{k \in \{\text{sem}, \text{feat}, \text{align}\}} \frac{r_i^k - \text{Mean}(\{r_j^k\}_{j=1}^N)}{\text{Std}(\{r_j^k\}_{j=1}^N)} \quad (2)$$

3.3. Post-training: Detection-GRPO

LLM Optimization. In the first stage, we optimize the LLM for user intent refinement, while the parameters of the image generation model remain frozen. In this setup, the LLM π_θ acts as the policy network. We first conduct SFT as a cold-start to teach it the "think-plan-then-generate" pattern required by the system prompt. During the RL phase, for a given text input x , the policy LLM samples and rewrites it, generating N optimized prompts $\{y_i\}_{i=1}^N$. Subsequently, the frozen diffusion model generates an image I conditioned on y . This image I is then evaluated by our multi-objective reward function to obtain the reward $R(I, y)$. In this context, each trajectory $\{o_i\}_{i=1}^N$ corresponds to the text sequence $y_i = (y_{1,i}, \dots, y_{T,i})$ sampled by the LLM. We update the LLM's parameters by maximizing the GRPO objective function defined in Equation (3).

$$\mathcal{J}(\theta) = \mathbf{E}_{c \sim \mathcal{C}, \{o_i\}_{i=1}^N \sim \pi_{\theta_{\text{old}}}(\cdot|c)} \left[\frac{1}{N} \sum_{i=1}^N \frac{1}{T} \sum_{t=1}^T \left(\mathcal{J}_{GRPO} - \beta \mathbf{D}_{\text{KL}}(\pi_\theta \| \pi_{\theta_{\text{ref}}}) \right) \right], \quad (3)$$

$$\mathcal{J}_{GRPO} = \min(\rho_{t,i} A_i, \text{clip}(\rho_{t,i}, 1 - \varepsilon, 1 + \varepsilon) A_i)$$

where A_i is the fused advantage function obtained from Equation (2) and $\rho_{t,i}$ denotes the importance sampling ratio between policies, which is defined in Equation (4).

$$\rho_{t,i} = \frac{\pi_\theta(y_{t,i}|x, y_{<t,i})}{\pi_{\theta_{\text{old}}}(y_{t,i}|x, y_{<t,i})} \quad (4)$$

The optimized LLM learns to explore and generate superior prompts, for instance: (1) adding rich scenic details, (2) naturally incorporating "flaws and imperfections," and (3) appending specific auxiliary words (e.g., "shot on iPhone") that aid the diffusion model in producing realistic images.

Diffusion Model Optimization. In the second stage, we optimize the diffusion model for photorealistic image generation, while keeping the parameters of the LLM component frozen. We employed a diffusion model RL framework Flow-GRPO [21], however, due to the stochastic nature of RL, images generated through limited denoising steps or full-trajectory exploration tend to be noisy and blurry [30]. This unstable process subsequently misleads the judgment of detection models, which are typically trained on datasets lacking exposure to such noisy or ambiguous artifacts. To address this training-inference inconsistency, for a given input x , the policy model p_θ executes a complete denoising inference $(x_T, x_{T-1}, \dots, x_0)$ for evaluation. Simultaneously, several consecutive steps Δt from this process is designated for random exploration and obtain N trajectories $(x_{T',i}, x_{T'-1,i}, \dots, x_{T'-\Delta t,i})$, which effectively balance RL explorativeness with training efficiency. We optimize the policy model through GRPO process the same as stage-1,

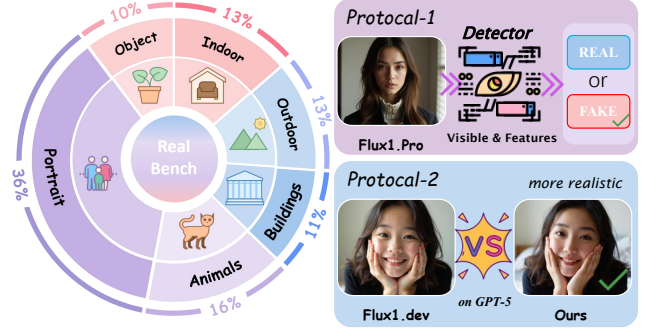


Figure 4. **Overview of the RealBench.** The left shows the categorical data composition. The right details its evaluation protocol.

as defined in Equation (5), and the probability ratio $\rho_{t,i}$ is given by:

$$\rho_{t,i} = \frac{p_\theta(x_{t-1,i}|x_{t,i}, c)}{p_{\theta_{\text{old}}}(x_{t-1,i}|x_{t,i}, c)} \quad (5)$$

SDE sampling provides stochasticity to the reverse process. During GRPO training, a portion of short texts are processed by the text-rewrite LLM, while the rest bypass this component. This strategy enhances the model's generalization robustness across prompts of varying lengths. Finally, detector reward optimization enables diffusion models to significantly reduce semantic and feature-level artifacts, consequently generating more realistic images.

4. RealBench

4.1. Dataset overview

We introduce RealBench, a benchmark platform designed to comprehensively evaluate the photorealism of T2I synthesized images. As illustrated in Fig. 4, RealBench features a meticulously curated dataset of 1000 high-quality, real-world images and their corresponding captions, sourced from the internet and free photography websites¹. This dataset spans seven distinct categories. Recognizing that "Portrait" is one of the most common and challenging categories in user T2I prompts, we have significantly increased its proportion while ensuring diversity across other categories. Unlike existing benchmarks focusing on instruction following (e.g., GenEval [10], DPG-Bench [13]) or human preference (e.g., HPD V2 [41]), RealBench focuses exclusively on assessing the photorealism of generated results. RealBench comprises two key evaluation protocols: Protocol-1, detector-based realism quantification (Detector-Scoring); and Protocol-2, arena-style preference evaluation (Arena-Scoring).

4.2. Detector-Based Realism Quantification

The core motivation for this protocol is that more photorealistic images should better evade advanced synthetic de-

¹<https://pixabay.com/>; <https://www.pexels.com/>

Table 1. Evaluation of image generation ability on RealBench. **Bold** indicates the best result, and underlined denotes the second best. * indicates that we used an LLM component for prompt optimization.

Model	Detector-Scoring				Arena-Scoring		Other metrics			
	Forensic-chat	OmniAID	Effort	GPT 5-Prompt	VS Real	VS Others	PickScore	CLIP	HPSv2.1	HPSv3
<i>Closed-Sourced T2I Model</i>										
FLUX-Pro [16]	57.45	21.55	20.94	50.14	18.20	-	23.68	86.85	30.79	<u>12.78</u>
Nano-Banana [6]	46.43	<u>31.02</u>	<u>11.74</u>	<u>73.19</u>	42.17	-	23.67	84.02	30.90	12.95
SeedDream 3.0 [16]	<u>63.47</u>	23.73	8.91	79.92	<u>36.40</u>	-	23.92	88.61	31.48	12.02
GPT-Image-1 [23]	75.63	33.59	5.70	70.98	33.71	-	<u>23.89</u>	<u>88.57</u>	<u>31.25</u>	12.48
<i>Open-Sourced T2I Model</i>										
SDXL [26]	43.37	24.44	8.44	23.82	9.22	35.20	23.02	84.65	28.44	9.87
SD-3.5-Large [9]	48.63	20.02	17.45	76.46	23.82	55.58	23.46	<u>88.23</u>	30.68	12.06
FLUX.1-dev [16]	40.91	21.32	14.85	43.03	12.61	43.60	<u>23.59</u>	86.33	<u>31.21</u>	<u>13.58</u>
FLUX.1-Kontext [17]	37.20	20.68	10.68	10.76	3.65	20.40	22.80	84.20	30.08	11.61
Echo-4o [51]	38.86	16.93	15.71	17.86	6.34	34.35	23.53	89.25	31.69	12.27
Bagel [7]	39.47	17.77	14.47	21.32	5.47	22.95	23.39	87.92	31.19	12.43
Qwen-Image [38]	57.47	36.82	17.10	65.03	18.25	47.35	21.97	82.83	25.50	8.15
SRPO [30]	64.14	<u>40.09</u>	24.73	81.29	40.70	64.30	23.55	86.16	29.66	12.43
FLUX-Krea [18]	57.10	32.44	18.42	79.44	37.60	66.40	23.77	87.60	30.75	12.50
Ours	<u>70.59</u>	37.85	<u>31.71</u>	<u>92.79</u>	<u>43.41</u>	<u>74.80</u>	23.58	86.80	31.87	13.61
Ours*	80.84	47.20	38.35	96.73	50.15	84.85	21.75	87.69	28.24	11.11

tectors, thus lowering their probability of being classified as "fake." Therefore, we utilize the probability of an image being deemed "real" as its photorealism score. In our assessment, we deploy detectors at both the visible semantic layer and the invisible feature layer. First, we include Forensic-Chat [20] and OmniAID [12], which are consistent with our reward function. At the semantic level, we also employ a leading closed-source MLLM (GPT-5 [24]) as a discriminator, guided by strict prompts (see supplementary material) to focus on common artifact regions. At the feature layer, we additionally use Effort [44], an unrelated expert synthetic detector known for its strong generalization in assessing image authenticity. The inclusion of these two reward-independent detectors helps ensure the comprehensiveness and robustness of the evaluation.

4.3. Arena-Style Preference Evaluation

We observe that it is inherently difficult for preference models, expert models, or even MLLMs to provide reliable, absolute quantitative scores for a single image. Therefore, inspired by LM-Arena², we adopt a pairwise comparison protocol, termed Arena-Scoring. In this protocol, we employ GPT-5 as the "judge model" to simulate user preferences. During evaluation, the judge model is presented with two semantically similar images corresponding to the identical text prompt (e.g., Model A's output vs. Model B's output, or Model A's output vs. a real image) and is compelled to make a forced-choice decision, selecting the one it perceives as more realistic. Finally, the images generated by

each model undergo at least 3000 random pairwise "battles" against outputs from other models and real-world data. We then calculate the final win rate. The inclusion of real images in the comparison pool not only enhances scoring stability but also allows us to validate the plausibility and reliability of the MLLM judge's preferences.

5. Experiments

5.1. Experimental Setup

Implement Details. Our RealGen model employs Qwen3-4B-Instruct [47] as the LLM component and FLUX.1-dev [16] as the base image generator. The training data for SFT and RL is primarily sourced from the real image subset of HPD v3 [22]. All experiments are conducted on 8 H200 GPUs. For the first stage of GRPO, we set the batch size to 32; for the second stage, the batch size is 12. The entire RL training process iterates for approximately 230 steps. To ensure a fair evaluation, we assessed model effectiveness on two datasets independent of the HPD v3 training data: our newly proposed RealBench and the "Photo" subset of HPD v2 [41].

Comparison Method. We conduct extensive comparisons between RealGen and current T2I models. These baselines include general generative models, such as the closed-source Flux.1 Pro [16], GPT-Image-1 [23], and Nano Banana [6], as well as open-source models like Flux.1 dev [16], Bagel [7], and Qwen-Image [38]. For a fair comparison, we also benchmark against methods specifically employing RL optimization for human preference or realism, such as Flux-Krea [18]

²<https://lmarena.ai/leaderboard/text-to-image>



Figure 5. Qualitative comparison of different methods on RealBench.

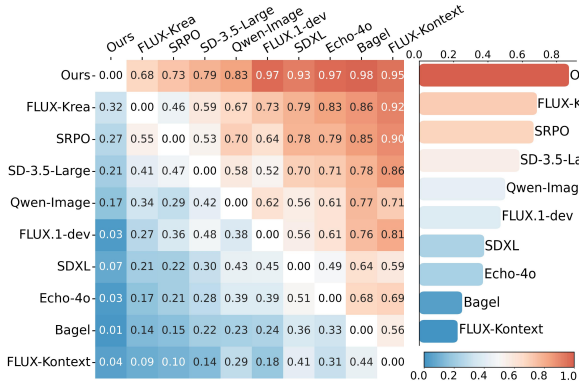


Figure 6. Pairwise Realism Comparison Matrix of Open-Sourced Text-to-Image Models.

and SRPO [30]. To ensure a fair comparison, all baselines utilize their official default settings.

Evaluation Protocol. We employ a comprehensive evaluation protocol composed of three main categories of metrics: (1) Detector-Scoring, (2) Arena-Scoring, and (3) Other Metrics. The detailed configurations for the first two categories are described in Section 4 (RealBench). The third category, Other Metrics, includes preference alignment scores such as Pick-Score [15], HPSv2.1 [41], and HPSv3 [22], as well as a LongCLIP [52] image-text alignment metric. It is worth noting that HPSv3, unlike its predecessors, incorporates real-world images into its training data, allowing it to more comprehensively reflect human preferences for both photorealism and overall quality.

5.2. Main Results

Table 1 presents the evaluation results of different methods on RealBench. Our proposed RealGen outperforms existing leading T2I models across multiple key photorealism metrics,

regardless of whether the LLM prompt rewrite component is used. For instance, images generated by general-purpose open-source T2I models, such as FLUX.1-kontext and Bagel, are easily identified by synthetic detectors, exposing obvious AI artifacts. Methods specialized for photorealism, like Flux-Krea and SRPO, achieve sub-optimal performance, but still a significant gap remains compared to our method. Crucially, strong performance on held-out evaluators (GPT-5 and Effort, Table 1), which were excluded from Detector-Reward training, validates that RealGen learned generalizable realism rather than merely overfitting to our reward models.

In the arena-style pairwise comparisons, we first analyze the "battle" results against Real images. RealGen demonstrates a significant advantage in this comparison, achieving a win rate approaching 50%, which suggests its outputs are capable of being confused with reality. In contrast, 8 of the 13 competing models achieved win rates below 30% against real images, clearly exposing their lack of photorealism. This stark disparity also validates the effectiveness of our Arena-Scoring as a reliable arbiter for realism. Furthermore, Figure 6 displays the win-rate matrix from the model-vs-model battles. The matrix confirms that RealGen achieves the highest overall win rate, indicating it is consistently selected as the more realistic option when compared directly against its peers.

Fig. 5 illustrates the qualitative visual comparisons. Base models like FLUX-dev and Bagel tend to produce images with excessive oiliness and unnatural highlights. Qwen-Image often generates overly smooth skin, exhibiting a typical AI artifact. Beyond a distinct "AI plastic" feel, GPT-Image-1 also shows an unnatural warm color cast skewed towards yellow-green. In contrast, the results from RealGen are superior in both texture and detail, appearing visually closer to actual photographs.

Table 2. Evaluation on the "Photo" subset of HPDv2 dataset [41].

Model	Detector-Scoring		Aesthetic Scoring	
	Forensic-Chat	OmniAID	HPSv2.1	HPSv3
FLUX-Pro [16]	66.57	40.13	28.57	11.64
Nano-Banana [6]	37.44	33.09	29.37	12.23
GPT-Image-1 [23]	63.75	27.56	29.73	12.19
SDXL [26]	44.53	32.22	25.92	7.27
FLUX.1 Kontext [17]	41.25	30.66	29.06	11.43
Bagel [7]	37.19	30.01	29.70	11.99
Qwen-Image [38]	48.85	32.17	27.60	9.44
SRPO [30]	62.65	48.23	27.88	11.11
FLUX-Krea [18]	58.00	44.96	28.76	11.39
Ours	71.34	56.93	30.18	13.10

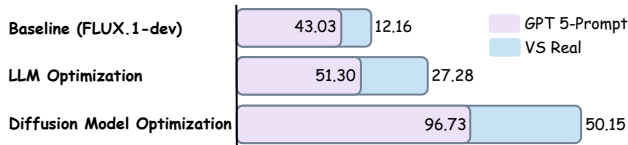


Figure 7. Ablation experiments on different key components.



Figure 8. Illustration of the key advantages of synthetic images.

To further validate the generalization ability of RealGen, we conducted an additional evaluation on the "photo" subset of HPDv2, with results presented in Table 2. The results clearly demonstrate that images synthesized by RealGen not only achieve exceptionally high detector-based realism scores but also rank among the top in human aesthetic preference scores. This provides evidence that our method can significantly enhance photorealism while simultaneously maintaining high aesthetic quality. The supplementary material contains further quantitative and qualitative results, such as detailed artifact analysis from Forensic-Chat and data from human evaluations.

Table 3. Ablation experiments on different reward functions.

Method	Effort	GPT 5	VS Real	CLIP
Baseline (FLUX.1-dev)	14.85	43.03	12.61	86.33
+ PickScore [15]	12.75	45.15	22.96	86.17
+ HPSv2.1 [41]	11.46	39.00	19.12	86.24
+ Detector-Score	31.71	92.79	43.41	86.80



Figure 9. Visualization of the effects of different reward functions.

5.3. Ablation experiments

As shown in Fig. 7, we conducted a progressive ablation study to explore the impact of the LLM optimization component and the Diffusion optimization component. For quantitative analysis, we employed two reward-independent metrics: the GPT-prompt Score and vs. Real images Arena-Scoring. The quantitative results indicate that applying only the Phase 1 LLM prompt optimization already leads to more realistic and detailed images by generating richer and more diverse prompts. Building on this, the subsequent inclusion of the Diffusion model optimization further improves image quality. The visualization in Fig. 8 shows the effect of component optimization: the LLM adds realistic descriptions to enrich detail, while optimizing the Diffusion model further enhances the image realism of the person, and glass.

Furthermore, we explored the distinct impacts of different reward functions on model optimization, comparing human preference rewards against our Detector-score. As shown in Table 3, when evaluated on multiple metrics held-out from the optimization objective, our Detector-Score demonstrates a clear advantage, proving it guides the model towards superior photorealism. The qualitative results in Figure 9 intuitively explain this gap: reward functions like PickScore and HPSv2.1 tend to bias the model towards cartoonish or artistic styles and still produce "oily" textures on portraits. In contrast, our proposed Detector-Score consistently leads to more realistic and photorealistic results.

6. Conclusion

To address the photorealism gap in current T2I models, we introduced RealGen. Our core contribution is the "Detector Reward" mechanism, inspired by adversarial generation, which utilizes both semantic and feature-level detectors to quantify image realism. By leveraging this reward signal with the GRPO algorithm, we successfully optimized the entire generation pipeline, enhancing image realism and detail; **this effectively realizes a "Detection for Generation" paradigm in the RL era.** Furthermore, we proposed RealBench, an automated benchmark employing Detector-Scoring and Arena-Scoring, which enables human-free photorealism evaluation. We hope our work will inspire further advancements in photorealistic image synthesis.

References

- [1] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. 3
- [2] Tingfeng Cao, Chengyu Wang, Bingyan Liu, Ziheng Wu, Jinhui Zhu, and Jun Huang. Beautifulprompt: Towards automatic prompt engineering for text-to-image synthesis. *arXiv preprint arXiv:2311.06752*, 2023. 2
- [3] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-alpha: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 3
- [4] Jiu hai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025. 3
- [5] Yufeng Cui, Honghao Chen, Haoge Deng, Xu Huang, Xinghang Li, Jirong Liu, Yang Liu, Zhuoyan Luo, Jinsheng Wang, Wenxuan Wang, et al. Emu3. 5: Native multimodal models are world learners. *arXiv preprint arXiv:2510.26583*, 2025. 3
- [6] Google DeepMind. Gemini 2.5 pro. Accessed: 2025-11-11, 2025. 2, 6, 8
- [7] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pre-training. *arXiv preprint arXiv:2505.14683*, 2025. 3, 6, 8
- [8] Haoge Deng, Ting Pan, Haiwen Diao, Zhengxiong Luo, Yufeng Cui, Huchuan Lu, Shiguang Shan, Yonggang Qi, and Xinlong Wang. Autoregressive video generation without vector quantization. *arXiv preprint arXiv:2412.14169*, 2024. 3
- [9] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 1, 3, 6
- [10] Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:52132–52152, 2023. 5
- [11] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1
- [12] Yuncheng Guo, Junyan Ye, Chenjue Zhang, Hengrui Kang, Haohuan Fu, Conghui He, and Weijia Li. Omniaid: Decoupling semantic and artifacts for universal ai-generated image detection in the wild, 2025. 3, 4, 6
- [13] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 5
- [14] Hengrui Kang, Siwei Wen, Zichen Wen, Junyan Ye, Weijia Li, Peilin Feng, Baichuan Zhou, Bin Wang, Dahua Lin, Linfeng Zhang, et al. Legion: Learning to ground and explain for synthetic image detection. *arXiv preprint arXiv:2503.15264*, 2025. 3
- [15] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. *Advances in neural information processing systems*, 36:36652–36663, 2023. 2, 3, 7, 8
- [16] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 2, 3, 4, 6, 8
- [17] Black Forest Labs. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space. 2025. 3, 6, 8
- [18] Sangwu Lee, Erwann Millon, Titus Ebbecke, Will Beddow, Le Zhuo, Iker García-Ferrero, Liam Esparraguera, Mihai Petrescu, Gian Saß, Gabriel Menezes, and Victor Perez. Flux.1 krea [dev]. <https://github.com/krea-ai/flux-krea>, 2025. 2, 3, 6, 8
- [19] Weijia Li, Jun He, Junyan Ye, Huaping Zhong, Zhi-meng Zheng, Zilong Huang, Dahua Lin, and Conghui He. Crossviewdiff: A cross-view diffusion model for satellite-to-street view synthesis. *arXiv preprint arXiv:2408.14765*, 2024. 3
- [20] Kaiqing Lin, Zhiyuan Yan, Ruoxin Chen, Junyan Ye, Ke-Yue Zhang, Yue Zhou, Peng Jin, Bin Li, Taiping Yao, and Shouhong Ding. Seeing before reasoning: A unified framework for generalizable and explainable fake image detection. *arXiv preprint arXiv:2509.25502*, 2025. 3, 4, 6
- [21] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv preprint arXiv:2505.05470*, 2025. 2, 3, 5
- [22] Yuhang Ma, Xiaoshi Wu, Keqiang Sun, and Hongsheng Li. Hpsv3: Towards wide-spectrum human preference score. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15086–15095, 2025. 6, 7
- [23] OpenAI. Gpt-4o. <https://openai.com/index/introducing-4o-image-generation>, 2025. 2, 3, 6, 8

- [24] OpenAI. Gpt-5 system card. <https://cdn.openai.com/gpt-5-system-card.pdf>, 2025. 6
- [25] Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang Wang, Zhiyang Xu, Jiahai Chen, Kunpeng Li, Felix Juefei-Xu, et al. Transfer between modalities with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025. 3
- [26] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 3, 6, 8
- [27] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2): 3, 2022. 3
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3
- [29] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1
- [30] Xiangwei Shen, Zhimin Li, Zhantao Yang, Shiyi Zhang, Yingfang Zhang, Donghao Li, Chunyu Wang, Qinglin Lu, and Yansong Tang. Directly aligning the full diffusion trajectory with fine-grained human preference. *arXiv preprint arXiv:2509.06942*, 2025. 2, 3, 5, 6, 7, 8
- [31] Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8871–8879, 2024. 3
- [32] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238, 2024. 2, 3
- [33] Linqing Wang, Ximing Xing, Yiji Cheng, Zhiyuan Zhao, Donghao Li, Tiankai Hang, Jiale Tao, Qixun Wang, Ruihuang Li, Comi Chen, et al. Promptenhancer: A simple approach to enhance text-to-image models via chain-of-thought prompt rewriting. *arXiv preprint arXiv:2509.04545*, 2025. 2
- [34] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. 4
- [35] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8695–8704, 2020. 3
- [36] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yuezeng Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 3
- [37] Siwei Wen, Junyan Ye, Peilin Feng, Hengrui Kang, Zichen Wen, Yize Chen, Jiang Wu, Wenjun Wu, Conghui He, and Weijia Li. Spot the fake: Large multimodal model-based synthetic image detection with artifact explanation. *arXiv preprint arXiv:2503.14905*, 2025. 3
- [38] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. 2, 6, 8
- [39] Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yuezeng Wang, Wanli Li, Xiyan Jiang, Yexin Liu, Junjie Zhou, et al. Omnigen2: Exploration to advanced multimodal generation. *arXiv preprint arXiv:2506.18871*, 2025. 3
- [40] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis. *arXiv preprint arXiv:2306.09341*, 2023. 2, 3
- [41] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis, 2023. 5, 6, 7, 8
- [42] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 3
- [43] Zeyue Xue, Jie Wu, Yu Gao, Fangyuan Kong, Lingting Zhu, Mengzhao Chen, Zhiheng Liu, Wei Liu, Qiushan Guo, Weilin Huang, et al. Dancegrp: Unleashing grp on visual generation. *arXiv preprint arXiv:2505.07818*, 2025. 2, 3
- [44] Zhiyuan Yan, Jiangming Wang, Zhendong Wang, Peng Jin, Ke-Yue Zhang, Shen Chen, Taiping Yao, Shouhong Ding, Baoyuan Wu, and Li Yuan. Effort: Efficient orthogonal modeling for generalizable ai-generated image detection. *arXiv preprint arXiv:2411.15633*, 2, 2024. 3, 6
- [45] Zhiyuan Yan, Kaiqing Lin, Zongjian Li, Junyan Ye, Hui Han, Zhendong Wang, Hao Liu, Bin Lin, Hao Li, Xue Xu, et al. Can understanding and generation truly benefit together—or just coexist? *arXiv e-prints*, pages arXiv–2509, 2025. 3
- [46] Zhiyuan Yan, Junyan Ye, Weijia Li, Zilong Huang, Shenghai Yuan, Xiangyang He, Kaiqing Lin, Jun He, Conghui He, and Li Yuan. Gpt-imgeval: A comprehensive benchmark for diagnosing gpt4o in image generation. *arXiv preprint arXiv:2504.02782*, 2025. 3
- [47] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint arXiv:2505.09388*, 2025. 4, 6
- [48] Junyan Ye, Jun He, Weijia Li, Zhutao Lv, Jinhua Yu, Haote Yang, and Conghui He. Skydiffusion: Street-to-satellite image synthesis with diffusion models and bev paradigm. *arXiv e-prints*, pages arXiv–2408, 2024. 3

- [49] Junyan Ye, Baichuan Zhou, Zilong Huang, Junan Zhang, Tianyi Bai, Hengrui Kang, Jun He, Honglin Lin, Zihao Wang, Tong Wu, et al. Loki: A comprehensive synthetic data detection benchmark using large multimodal models. *arXiv preprint arXiv:2410.09732*, 2024. 3
- [50] Junyan Ye, Jun He, Xiang Zhang, Yi Lin, Honglin Lin, Conghui He, and Weijia Li. Satellite image synthesis from street view with fine-grained spatial textual guidance: A novel framework. *IEEE Geoscience and Remote Sensing Magazine*, 2025. 3
- [51] Junyan Ye, Dongzhi Jiang, Zihao Wang, Leqi Zhu, Zhenghao Hu, Zilong Huang, Jun He, Zhiyuan Yan, Jinghua Yu, Hongsheng Li, et al. Echo-4o: Harnessing the power of gpt-4o synthetic images for improved image generation. *arXiv preprint arXiv:2508.09987*, 2025. 6
- [52] Beichen Zhang, Pan Zhang, Xiaoyi Dong, Yuhang Zang, and Jiaqi Wang. Long-clip: Unlocking the long-text capability of clip. In *European conference on computer vision*, pages 310–325. Springer, 2024. 4, 7
- [53] Lingzhi Zhang, Yuqian Zhou, Connelly Barnes, Sohrab Amirghodsi, Zhe Lin, Eli Shechtman, and Jianbo Shi. Perceptual artifacts localization for inpainting. In *European Conference on Computer Vision*, pages 146–164. Springer, 2022. 3
- [54] Lingzhi Zhang, Zhengjie Xu, Connelly Barnes, Yuqian Zhou, Qing Liu, He Zhang, Sohrab Amirghodsi, Zhe Lin, Eli Shechtman, and Jianbo Shi. Perceptual artifacts localization for image synthesis tasks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7579–7590, 2023. 3
- [55] Xiaofeng Zhang, Aaron Courville, Michal Drozdal, and Adriana Romero-Soriano. The intricate dance of prompt complexity, quality, diversity, and consistency in t2i models. *arXiv preprint arXiv:2510.19557*, 2025. 2
- [56] Mingjian Zhu, Hanting Chen, Qiangyu Yan, Xudong Huang, Guanyu Lin, Wei Li, Zhijun Tu, Hailin Hu, Jie Hu, and Yunhe Wang. Genimage: A million-scale benchmark for detecting ai-generated image. *Advances in Neural Information Processing Systems*, 36:77771–77782, 2023. 3