

Science-T2I: Addressing Scientific Illusions in Image Synthesis

Jialuo Li ¹, Wenhao Chai ², Xingyu Fu ³, Haiyang Xu ⁴, Saining Xie ¹

¹ New York University

² University of Washington

³ University of Pennsylvania

⁴ University of California, San Diego

Link: [Project Page](#) | [Leaderboard](#) | [Dataset & Model](#) | [Code](#)

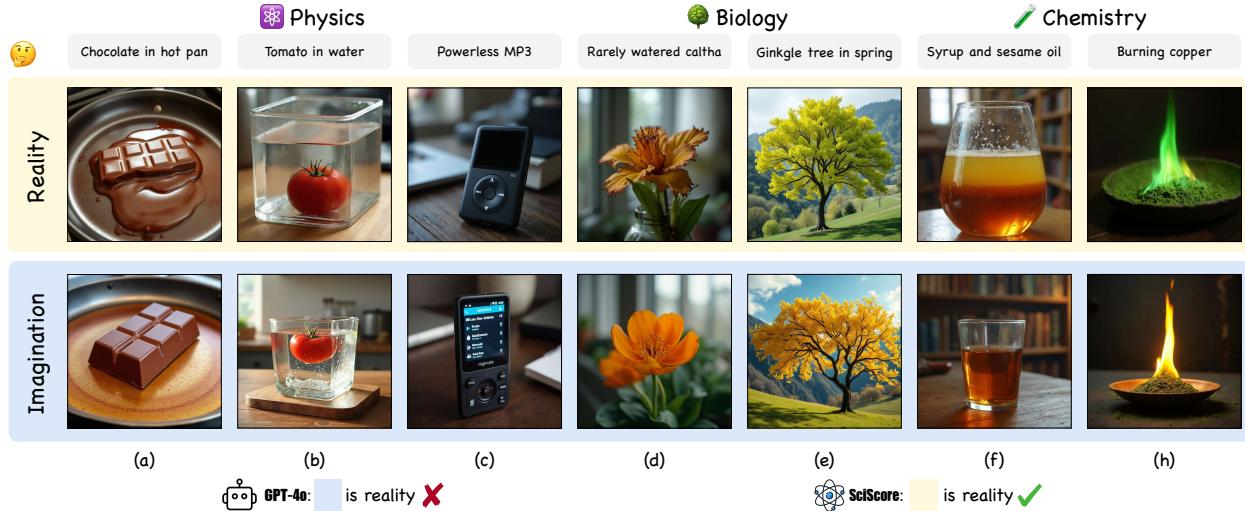


Fig. 1: Comparison between GPT-4o and SCISCORE. Given a prompt (in grey) requiring scientific knowledge, FLUX [1] model generates imaginary images (lower row) that are far from reality (upper row). Moreover, LMMs like GPT-4o [2] fail to identify the realistic image, whereas our end-to-end reward model SCISCORE succeeds. Notice that the prompts here are summarization of the real prompts that we used for illustration purposes.

We present a novel approach to integrating scientific knowledge into generative models, enhancing their realism and consistency in image synthesis. First, we introduce SCIENCE-T2I, an expert-annotated adversarial dataset comprising adversarial 20k image pairs with 9k prompts, covering wide distinct scientific knowledge categories. Leveraging SCIENCE-T2I, we present SCISCORE, an end-to-end reward model that refines the assessment of generated images based on scientific knowledge, which is achieved by augmenting both the scientific comprehension and visual capabilities of pre-trained CLIP model. Additionally, based on SCIENCE-T2I, we propose a two-stage training framework, comprising a supervised fine-tuning phase and a masked online fine-tuning phase, to incorporate scientific knowledge into existing generative models. Through comprehensive experiments, we demonstrate the effectiveness of our framework in establishing new standards for evaluating the scientific realism of generated content. Specifically, SCISCORE attains performance comparable to human-level, demonstrating a 5% improvement similar to evaluations conducted by experienced human evaluators. Furthermore, by applying our proposed fine-tuning method to FLUX, we achieve a performance enhancement exceeding 50% on SCISCORE.

1. Introduction

The quest to conceptualize the visual world and construct real world simulators has been a long-standing endeavor in the computer vision community [11, 21, 23, 33, 76, 77]. As articulated by [14], “The goal of image synthesis is to create, using the computer, a visual experience that is identical to what a viewer would experience when viewing a real environment.” In alignment with this vision, recent advances in generative modeling have notably improved the performance of image synthesis [8, 10, 54, 58, 62, 74]. While these advancements enable the generation of higher resolution, more aesthetically pleasing images with superior Frechet Inception Distance (FID) scores [1, 5, 54, 72], these models often produce superficial imitations rather than authentic representations of the real visual world [6, 22, 50, 51]. This limitation often arises from an inadequate understanding of the underlying scientific principles of realism, as demonstrated in the lower row of FLUX [1] generated images in Figure 1. Consequently, the images generated tend to mirror imaginative constructs, resulting in a noticeable gap between these creations and the tangible reality we inhabit.

This paper integrates scientific knowledge into image synthesis to bridge the gap between imagination and realism. We introduce SCIENCE-T2I, a comprehensive and expert-annotated dataset comprising over 20k adversarial image pairs and 9k prompts that span diverse fields such as physics, chemistry, and biology, and cover 16 unique scientific phenomena. Each data pair is collected in an adversarial setup, consisting of one image that accurately aligns with reality and another that does not, thereby facilitating preference modeling. To ensure quality and accuracy, all data were reviewed by human experts whose assessments were based on their professional expertise and consultation of an extensive knowledge base.

Leveraging SCIENCE-T2I, we further present SCISCORE, an end-to-end reward model infused with diverse expert-level scientific knowledge, designed to evaluate generated images as a science teacher would. Our results demonstrate that SCISCORE outperforms complex, prompt-engineering-reliant large multimodal models (LMMs) such as GPT-4o. Compared to GPT-4o, SCISCORE excels in capturing fine-grained visual details that LMMs often neglect as in Figure 1, and functions as a comprehensive end-to-end reward model – eliminating the dependence on language-guided inference processes, which can fail due to hallucinations.

Utilizing SCISCORE, we introduce a two-stage training methodology to develop an enhanced image synthesis model that conform to the realist with world knowledge. Specifically, we begin with supervised fine-tuning (SFT) on FLUX.1[dev][1] using SCIENCE-T2I. This initial phase is subsequently followed by an additional stage of online fine-tuning, where SCISCORE functions as the reward model and employs a masking strategy to improve the performance.

Our main contributions are summarized as follows:

- We introduce SCIENCE-T2I of over 9k prompts and 20k adversarial image pairs, annotated by experts to reflect reality, enabling the training of a language-guided reward model for text-to-image alignment with scientific knowledge.
- We propose an optimization strategy using the reward model SCISCORE to enhance diffusion-based generative models, showing improved alignment of generated images with reality on a quantitative benchmark.
- Extensive experiments show that our method outperforms the baseline by over 50%, marking a significant advancement in grounding the model in real-world scenarios.

2. Related Works

2.1. Physics Modeling in Generative Models

Integrating physical laws into generative models has become a vital area of research to enhance the realism and consistency of generated data across various domains, including image synthesis [45, 51], video generation [6, 9, 34, 50], and 3D modeling [25]. PhyBench [51] is a pioneering work that explores the incorporation of physical knowledge into current text-to-image (T2I) models by providing a comprehensive dataset designed to test physical commonsense across various domains. In the realm of text-to-video (T2V) models, benchmarks like VideoPhy [6] and PhyGenBench [50] evaluate whether current generative models can accurately simulate physical commonsense in real-world scenarios involving various material interactions. PhysComp [25] advances single-image 3D reconstruction by decomposing geometry into mechanical properties and enforcing static equilibrium. Our work differs by designing tasks as reasoning challenges, requiring models to understand and apply physical laws to generate accurate outputs. This approach pushes the boundaries of physical knowledge integration in generative models by emphasizing implicit reasoning over explicit description.

2.2. RL in Diffusion Models

Reinforcement learning (RL) has been effectively applied in diffusion models to enhance sample quality. For instance, VersaT2I [24] and DreamSync [64] simply use reject sampling. ReNO [18] focus on adapting a diffusion model during inference by purely optimizing the initial latent noise using a differentiable objective. Some other works [7, 67, 73] leverages DPO [57] or PPO [61] as optimization strategies. Our work differs by introducing a novel reward function that leverages physical commonsense to guide diffusion process, ensuring generated samples are physically plausible.

2.3. Benchmarking Image Synthesis Models

Standard metrics like FID [27], IS [59], LPIPS [78], and CLIPScore [26] are commonly used to assess image synthesis models. With model advancements, newer methods emphasize human evaluation and multimodal LLM-based assessment. HPSv2 [70], PickScore [37] and ImageReward [71] provide human preference annotations, while VQAScore [40], TIFA [29], VIEScore [39], LLMscore [49], and DSG [13] utilize VQA-style evaluations. For object attributes and relationships, benchmarks like T2I-CompBench [30] and CLIP-R-Precision [53] have been introduced. However, there are few benchmarks focusing on the physical commonsense. PhyBench [51] establishes a set of grading criteria and employs vision-language models to discretely score images. In contrast, we introduce SCISCORE, an end-to-end model designed to provide a more refined and continuous scoring mechanism for images.

3. Dataset: SCIENCE-T2I

We introduce SCIENCE-T2I, a novel dataset specifically designed to enhance text-to-image and multimodal models' understanding of underlying scientific principles. Unlike conventional datasets that focus on direct textual descriptions [15, 38, 44] and preference annotation [37, 70, 71], SCIENCE-T2I challenges models to perform implicit reasoning based on prompts that need scientific knowledge.

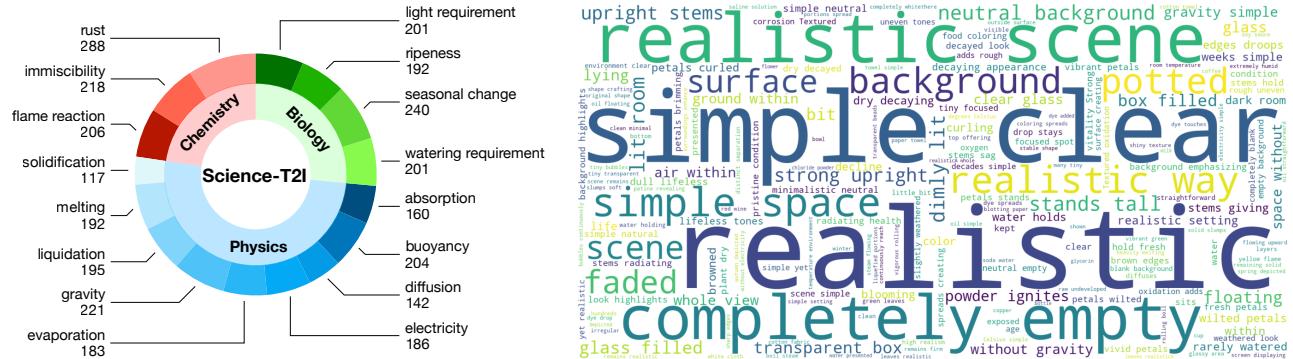


Fig. 2: Data statistics. (Left) SCIENCE-T2I is organized into three primary scientific fields: Chemistry, Biology, and Physics. Each field is divided into specific categories, with the numbers indicating the volume of implicit prompts collected for each category. (Right) Word cloud of structured prompt in SCIENCE-T2I.

Task Overview. As illustrated in Figure 2, SCIENCE-T2I consists of 16 tasks that require the model to infer or visualize concepts not explicitly stated in the prompts but rooted in underlying scientific principles. These tasks are inspired by existing research such as PhyBench [51] and Commonsense-T2I [22], as well as new concepts developed for this study. Each task is meticulously designed with the following objectives:

- **Rewriting Capability.** Tasks use prompts that allow flexible rephrasing, thereby enabling different expressions to effectively achieve the same visual meaning. For example, *an unripe apple* is able to be rephrased as *a green apple*, conveying the same visual concept. Further explanations and examples can be found in Appendix S1.
- **Scientific Knowledge Integration.** Tasks are based on established scientific principles in physics, chemistry, and biology, providing a clear and consistent framework. This approach reduces the ambiguity of commonsense knowledge, which can vary culturally or contextually.

For a comprehensive understanding of each task, detailed descriptions are provided in Appendix S2. Furthermore, beyond a classification based on scientific disciplines, the tasks can be categorized into two distinct groups:

- **Subject-oriented Tasks (ST)** require scientific reasoning to discern how inherent differences between subjects lead to varying visual features under identical conditions.
- **Condition-oriented Tasks (CT)** focus on how a single condition affects various subjects. Scientific reasoning in these tasks centers on the applied condition, not the subject’s individual properties.

This classification allows for a more nuanced understanding of model generalization across different tasks. Additional details regarding these definitions and the specific classification of each task can be found in Appendix S3.

Prompt Design. We classify prompts into two types: those requiring inference from scientific knowledge and their rewritten versions that utilize rewriting capabilities. Additionally, PhyBench [51] reveals that models often ignore these principles, focusing instead on descriptive text, indicating a third category based on description rather than inference. To clarify these concepts, we introduce specific terminologies:

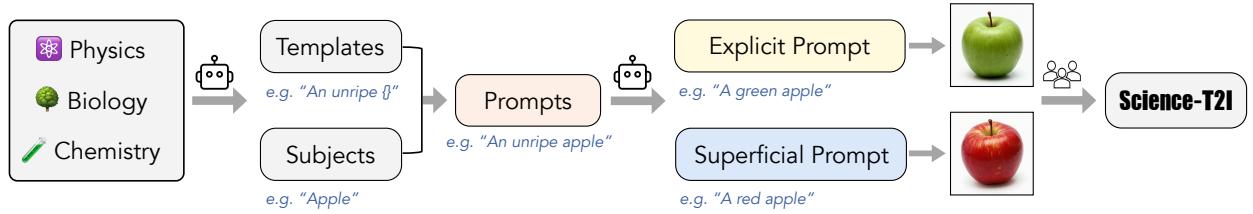


Fig. 3: Data curation pipeline. For each task, GPT-4o [2] first generates structured templates that capture the scientific principles while allowing for variability in objects or substances. These templates are used to create implicit prompts, which GPT-4o [2] then expands into explicit and superficial prompts, ultimately guiding the synthesis of corresponding explicit and superficial images.

- **Implicit Prompt (IP).** It contains specific terms or phrases that imply certain visual characteristics or phenomena requiring interpretative reasoning based on scientific knowledge. For example, the prompt "an unripe apple" suggests greenness without explicitly stating it.
- **Explicit Prompt (EP).** It reformulates the implicit prompt into a clear, descriptive statement that accurately reflects the intended image. For instance, the prompt "a green apple" directly conveys the immaturity.
- **Superficial Prompt (SP).** It provides an explicit interpretation but neglects scientific reasoning, focusing only on surface descriptions and simplistic interpretations. For example, interpreting "an unripe apple" as "a red apple" overlooks the implied maturity, leading to inaccuracies.

Data Curation. We leverage GPT-4o to generate templates and corresponding prompts during data curation. These outputs are then used to drive T2I models for image generation. The complete data curation pipeline is illustrated in Figure 3, with further details provided in Appendix S4.

4. Method: SCISCORE

While CLIP [55] effectively aligns textual and visual data, it struggles to accurately match implicit prompts with their corresponding images. To address this limitation, we introduce SCISCORE, a reward model fine-tuned on SCIENCE-T2I that extends CLIP’s architecture [55]. SCISCORE assesses the extent to which an image embodies the visual information derived from the scientific principles articulated within the prompt. In this section, we first define the reward mechanism for evaluating prompt-image compatibility (§4.1) and then detail the training methods used to optimize SCISCORE’s performance (§4.2).

4.1. Reward Modeling

SCISCORE extends the CLIP architecture [55] by independently encoding a text prompt x and an image y into a shared high-dimensional vector space using separate transformer encoders [65], E_{txt} and E_{img} . The reward is computed based on the alignment between textual and visual modalities, quantified by the inner product of their respective encoded representations and subsequently scaled by a learnable temperature parameter T :

$$r(y, x) = T \cdot \frac{E_{\text{txt}}(x) \cdot E_{\text{img}}(y)}{\|E_{\text{txt}}(x)\| \|E_{\text{img}}(y)\|}. \quad (1)$$

4.2. Training Techniques

For developing SCISCORE we employed a fine-tuning approach on the CLIP [55] using SCIENCE-T2I. Each training instance is structured as a tuple $(x_i, x_e, x_s, y_e, y_s)$, where x_i is the implicit prompt, x_e and x_s are the explicit and superficial prompts, respectively. Correspondingly, y_e and y_s denote the explicit and superficial images.

Predicted Preferences Calculation. Following preference modeling approaches in language from prior work [52, 63], the predicted preference $\hat{p}_{\text{img}}(x_a \succ x_b; y)$ for prompt x_a over prompt x_b for a given image y is calculated as:

$$\hat{p}_{\text{img}}(x_a \succ x_b; y) = \frac{\exp(r(y, x_a))}{\exp(r(y, x_b)) + \exp(r(y, x_a))} \quad (2)$$

Similarly, for a given prompt x , the predicted preference $\hat{p}_{\text{txt}}(y_a \succ y_b; x)$ for image y_a over image y_b is given by:

$$\hat{p}_{\text{txt}}(y_a \succ y_b; x) = \frac{\exp(r(y_a, x))}{\exp(r(y_a, x)) + \exp(r(y_b, x))} \quad (3)$$

Implicit Prompt Alignment (IPA). Preliminary experiments revealed that the pretrained CLIP model [55] tends to embed the implicit prompt in a manner similarly to the corresponding superficial prompt. To address this issue, we minimize the KL divergence between the target preference $p_{\text{txt}} = [1, 0]$ and the predicted preference $\hat{p}_{\text{txt}} = [\hat{p}_{\text{txt}}(y_e \succ y_s; x_i), \hat{p}_{\text{txt}}(y_s \succ y_e; x_i)]$. This effectively aligns the implicit prompt with the explicit image over the superficial image. The loss function is defined as:

$$\mathcal{L}_{\text{IPA}} = \sum_{j=1}^2 p_{\text{txt}_j} \left(\log p_{\text{txt}_j} - \log \hat{p}_{\text{txt}_j} \right) \quad (4)$$

Image Encoder Enhancement (IEE). To effectively handle reasoning tasks that involve fine-grained visual phenomena, it is imperative to enhance the capabilities of the image encoder. The objective of this enhancement is captured by the following loss function:

$$\mathcal{L}_{\text{IEE}} = \mathcal{L}_{\text{img}}^+ + \mathcal{L}_{\text{img}}^- \quad (5)$$

where $\mathcal{L}_{\text{img}}^+$ and $\mathcal{L}_{\text{img}}^-$ correspond to the losses associated with explicit and superficial image preferences, respectively. The explicit image loss $\mathcal{L}_{\text{img}}^+$ is defined as:

$$\mathcal{L}_{\text{img}}^+ = \sum_{j=1}^2 p_{\text{img}_j}^+ \left(\log p_{\text{img}_j}^+ - \log \hat{p}_{\text{img}_j}^+ \right), \quad (6)$$

where $p_{\text{img}}^+ = [1, 0]$ signifies a preference for the explicit image. The predicted probabilities are denoted by:

$$\hat{p}_{\text{img}}^+ = [\hat{p}_{\text{img}}(x_e \succ x_s; y_e), \hat{p}_{\text{img}}(x_s \succ x_e; y_e)], \quad (7)$$

Similarly, the superficial image loss $\mathcal{L}_{\text{img}}^-$ is defined as:

$$\mathcal{L}_{\text{img}}^- = \sum_{j=1}^2 p_{\text{img}_j}^- \left(\log p_{\text{img}_j}^- - \log \hat{p}_{\text{img}_j}^- \right), \quad (8)$$

where $p_{\text{img}}^- = [0, 1]$ indicates a preference for the superficial image. The predicted probabilities are given by:

$$\hat{p}_{\text{img}}^- = [\hat{p}_{\text{img}}(x_e \succ x_s; y_s), \hat{p}_{\text{img}}(x_s \succ x_e; y_s)] \quad (9)$$

The overall loss function integrates \mathcal{L}_{IPA} with \mathcal{L}_{IEE} as:

$$\mathcal{L} = \mathcal{L}_{\text{IPA}} + \lambda \mathcal{L}_{\text{IEE}}, \quad (10)$$

where λ is a hyper-parameter that controls the relative weight of the image encoder enhancement loss in relation to the implicit prompt alignment loss.

5. Two-Stage T2I Model Fine-Tuning

5.1. Supervised Fine-tuning (SFT)

Current post-training algorithms for diffusion models, such as those utilizing PPO [7, 19] and DPO [66, 72], have significantly advanced model fine-tuning. However, these methods are constrained by the requirement that the optimization objectives remain within the distribution of the pre-trained model. While this limitation is acceptable for tasks like aesthetic enhancement, which involve preferences among generated images, it poses challenges for applications requiring scientific reasoning. Preliminary experiments demonstrate that pre-trained models lack an understanding of scientific principles, as they are primarily trained on descriptive prompts paired with images. This shortcoming presents a significant obstacle for post-training techniques aimed at embedding scientific comprehension into diffusion models.

Our methodology begins with the supervised fine-tuning of a pre-trained model to enhance its scientific understanding, utilizing the SCIENCE-T2I. As illustrated by the experimental results in Table 3, FLUX [1] models consistently achieve superior performance in direct text-image alignment and exhibit a strong capacity for generating realistic styles, as evidenced by our preliminary experiments. Based on these observations, we adopt FLUX.1[dev][1] as our base model. Since FLUX [1] employs flow matching [46] framework, the SFT training objective is formulated as:

$$L_{\text{SFT}} = \mathbb{E}_{t, p_t(z|\epsilon), p(\epsilon)} \|v_\theta(z, t) - u_t(z|\epsilon)\|_2^2 \quad (11)$$

In this formulation, we adopt the same mathematical notation as presented in [17] to ensure consistency. For further clarity, additional details are provided in Appendix S15.

5.2. Online Fine-tuning (OFT)

After performing domain transfer using SFT, we apply an online fine-tuning approach for further model refinement with pipeline shown in Figure 4. Following the methodology proposed by DDPO [7], we conceptualize the denoising process within the diffusion model as a multi-step MDP:

$$s_t \triangleq (c, t, x_{1-t}), \pi_\theta(a_t | s_t) \triangleq p_\theta(x_{1-\Delta t-t} | c, t, x_{1-t}), \rho_0(s_0) \triangleq (p(c), \delta_0, \mathcal{N}(0, I)) \quad (12)$$

$$a_t \triangleq x_{1-\Delta t-t}, P(s_{t+\Delta t} | s_t, a_t) \triangleq (\delta_c, \delta_{t+\Delta t}, \delta_{x_{1-t-\Delta t}}), r(s_t, a_t) \triangleq \begin{cases} r(x_0, c) & \text{if } t = 1 \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

The mathematical notation and formulation above align with those used in DDPO [7], with slight adjustments to the timestamp notation for clarity. However, flow matching [46] is typically formulated

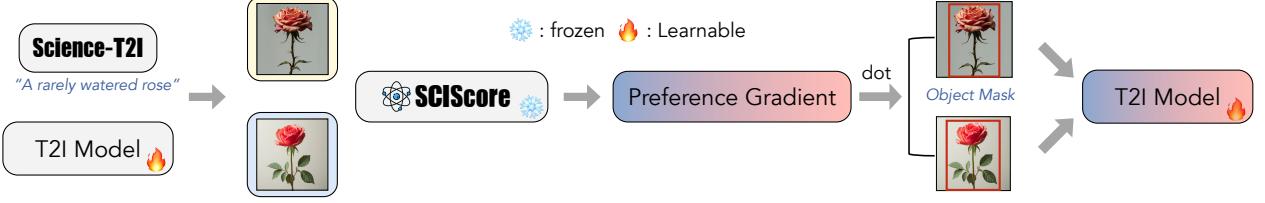


Fig. 4: Online fine-tuning pipeline. For each prompt, two images are generated to compute SCISCORE preference metric. Simultaneously, GroundingDINO [47] extracts segmentation masks from these images based on the prompts, which are then used to block gradient propagation in the corresponding spatial regions.

as an Ordinary Differential Equation (ODE), resulting in a deterministic process. This deterministic formulation complicates the computation of the policy $\pi_\theta(a_t | s_t)$:

$$\pi_\theta(a_t | s_t) = \delta(x_{1-\Delta t-t} - (x_{1-t} - v_\theta(s_t)\Delta t)) \quad (14)$$

In alignment with the discussion in [16], we can alternatively interpret flow matching [46] as a Stochastic Differential Equation (SDE), which is mathematically formulated as:

$$dx_t = \left(v_\theta(x_t, t) + \frac{\sigma_t^2}{2\beta_t\eta_t} \lambda_t \right) dt + \sigma_t dB_t \quad (15)$$

$$\eta_t = \left(\frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right), \quad \lambda_t = \left(v(x_t, t) - \frac{\dot{\alpha}_t}{\alpha_t} x_t \right) \quad (16)$$

where B_t denotes Brownian motion. By discretizing this equation while leveraging the rectified flow employed by FLUX [1], where $\alpha_t = t$ and $\beta_t = 1 - t$, we obtain:

$$\pi_\theta(a_t | s_t) = \mathcal{N}(a_t; \mu_\theta(s_t), \sigma_t^2 I) \quad (17)$$

$$\mu_\theta(s_t) = \frac{t\sigma_t^2 + 2(1-t)}{-2(1-t)} v_\theta(s_t) \Delta t + \frac{2(1-t) + \sigma_t^2 \Delta t}{2(1-t)} x_{1-t} \quad (18)$$

In this framework, the parameter σ_t is subject to manual configuration. Notably, setting $\sigma_t = 0$ simplifies the formulation to the deterministic case, as delineated in Equation 14. For the training objective, we adopt DPO as introduced by [56], with further technical details provided in Appendix S15. Specifically, given a condition (typically a prompt) c , we randomly sample two trajectories:

$$\sigma_w = \{s_0^w, a_0^w, s_{\Delta t}^w, a_{\Delta t}^w, \dots, s_1^w, a_1^w\}, \quad \sigma_l = \{s_0^l, a_0^l, s_{\Delta t}^l, a_{\Delta t}^l, \dots, s_1^l, a_1^l\} \quad (19)$$

Assuming that the reward satisfies $r(s_1^w, a_1^w) > r(s_1^l, a_1^l)$, the training objective is formulated as:

$$\mathcal{L} = \mathbb{E} \left[\log \rho \left(\beta \log \frac{\pi_\theta(a_k^l | s_k^l)}{\pi_{\text{ref}}(a_k^l | s_k^l)} - \beta \log \frac{\pi_\theta(a_k^w | s_k^w)}{\pi_{\text{ref}}(a_k^w | s_k^w)} \right) \right] \quad (20)$$

Subject-Based Masking Strategy. Considering the subject-oriented characteristics inherent to our scientific reasoning tasks, we employed a subject-based masking strategy during training. Specifically, we extract the subject from the input prompt and utilize GroundingDINO [47] to identify the bounding box around the subject. Subsequently, only the content within this bounding box is used for gradient backpropagation. Define mask corresponding to the box as \mathcal{M} , then the final training objective:

$$\mathcal{L} = -\mathbb{E} \left[\log \rho \left(\beta \log \frac{\mathcal{M}^w \odot \pi_\theta(a_k^w | s_k^w)}{\mathcal{M}^w \odot \pi_{\text{ref}}(a_k^w | s_k^w)} - \beta \log \frac{\mathcal{M}^l \odot \pi_\theta(a_k^l | s_k^l)}{\mathcal{M}^l \odot \pi_{\text{ref}}(a_k^l | s_k^l)} \right) \right].$$

6. Experiment: SCISCORE

6.1. Implementation Details

Training Setting. We fine-tune the CLIP-H model [32] using our framework on SCIENCE-T2I training set with both text and image encoder learnable. The experiment completes within one hour on 8 A6000 GPUs. For detailed hyperparameter configurations, please refer to Appendix S6.

Evaluation Setting. To thoroughly evaluate the model’s generalization ability across different environments, we have created two additional manually annotated test sets, both included in SCIENCE-T2I.

- **SCIENCE-T2I S** (671 tuples): This test set maintains the same stylistic attributes as the training set, emphasizing simplicity scene and focusing solely on the reasoning regions.
- **SCIENCE-T2I C** (227 tuples): This test set introduces additional and diverse scene settings both in the prompts and the corresponding images.

Additional details regarding their configuration and design are provided in Appendix S8.

Baseline Setup. We establish our baseline using three evaluation dimensions: VLMs, LMMs, and human assessments. For VLMs, we utilize CLIP-H [32], BLIPScore [42, 43] and SigLIP [75]. In the LMM category, we employ open-source models such as LLaVA-OV [41], Qwen2-VL [68] and InternVL [12], as well as the proprietary model GPT-4o-mini [2]. Additionally, we explore the use of CoT reasoning [69] during inference. Human evaluations involved 10 experts with science or engineering degrees. Further details regarding the baseline setup and evaluation process are provided in Appendix S9, while the experimental results are summarized in Table 1.

6.2. Results

Limited Performance of VLMs and LMMs. CLIP-H [55], BLIPScore [42] and SigLIP [75] demonstrate near-random accuracy across both test sets, underscoring their limitations in effectively distinguishing images when given only implicit prompts. While all LMMs, including GPT-4o-mini [2], despite being equipped with a vast knowledge base, fails to deliver satisfactory performance in these tasks. Notably, the application of CoT prompting [69] does not yield significant improvements in this context. A more thorough and detailed analysis of these limitations is provided in Appendix S11.

Table 1: Performance comparison of different models on SCIENCE-T2I S and SCIENCE-T2I C across different subjects, measured by accuracy in two-choice selection task. The best overall performance is highlighted with **bold values. To distinguish between model categories, green shading indicates the highest accuracy within VLMs, while purple shading shows the highest accuracy within LMMs.**

Model	SCIENCE-T2I S				SCIENCE-T2I C			
	Physics	Chemistry	Biology	Avg.	Physics	Chemistry	Biology	Avg.
Human Eval	87.67	75.85	95.29	87.01	84.71	85.40	89.14	86.02
Random Guess	50.00	50.00	50.00	50.00	50.00	50.00	50.00	50.00
CLIP-H [32]	55.08	52.38	55.88	54.69	56.56	44.44	76.67	59.47
BLIPScore [42]	50.35	43.08	59.86	55.00	49.78	60.00	58.33	51.54
SigLIP ViT-SO-14 [75]	59.63	53.17	55.94	57.23	61.48	51.11	70.00	61.67
Qwen2-VL-7B [68]	60.03	67.01	68.82	63.79	66.80	50.00	90.83	69.82
LLaVA-OV-7B [41]	68.22	57.82	64.71	65.05	74.59	50.00	76.67	70.26
InternVL2.5-8B [12]	67.80	62.24	84.41	70.79	73.77	65.56	85.83	75.33
GPT-4o mini [2]	61.97	73.81	86.76	70.83	69.29	70.00	90.00	74.78
GPT-4o mini+ CoT [69]	67.04	76.87	90.00	74.97	72.44	70.00	92.50	77.16
SCI SCORE (ours)	94.92	80.95	100.00	93.14	86.89	91.11	100.00	91.19

Generalization of SCI SCORE to Complex Scenes. Notably, SCI SCORE demonstrates strong performance not only on SCIENCE-T2I S, indicating its capability to handle scientific reasoning tasks in simple settings, but also on SCIENCE-T2I C, which features more complex scenes. This highlights the model’s ability to focus on specific regions of interest while effectively disregarding environmental distractions. Furthermore, it is noteworthy that SCI SCORE not only achieves but also surpasses human-level performance, with scores of 93.14 and 91.19 compared to human evaluation of 87.01 and 86.02.

Generalization of SCI SCORE across ST and CT. Given that all our predefined tasks can be classified as either ST or CT, we further investigate the generalization capabilities of SCI SCORE across these two categories. We evaluate the performance of SCI SCORE on SCIENCE-T2I S and SCIENCE-T2I C based on this categorization, as illustrated in Figure 5. The results reveal a significant performance gap between the two types of tasks, with nearly all failure cases occurring in ST. This outcome is anticipated, as CT rely on generalizable visual features (e.g., the absence of gravity implies floating objects), whereas ST are primarily driven by subject-specific characteristics, posing greater challenges. For novel or unseen subjects, the model struggles to identify appropriate visual content due to insufficient prior exposure to the subject’s distinctive attributes.

6.3. Ablation Study

Effect of IEE. To investigate the effect of IEE on the performance of SCI SCORE, comparative experiments are conducted. As shown in Table 2, there exists a trade-off between increasing the IEE loss rate and maintaining IPA loss. A lower IEE loss rate fails to enhance the image encoder’s ability to detect fine-grained details, whereas a higher IEE loss rate diminishes the focus on prompt alignment. We identified $\lambda = 0.25$ as the optimal for these objectives. Additionally, we provide further qualitative results and analysis in Appendix S12.

Table 2: Ablation study on different λ . The best performance is highlighted with **bold** values.

λ	SCIENCE-T2I S	SCIENCE-T2I C
0	93.14	88.99
0.1	92.85	90.75
0.5	92.85	91.19
0.75	93.14	88.99
0.25	93.14	91.19

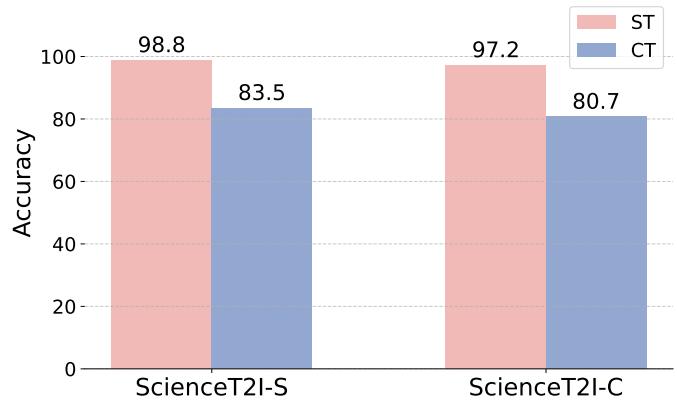


Fig. 5: Performance of SCISCORE in ST and CT.

6.4. Benchmarking T2I Generation

By leveraging the superior performance of SCISCORE, rather than relying on LMMs that require complex prompting techniques [51], we propose an end-to-end utilization of SCISCORE for benchmarking current T2I models on our predefined scientific reasoning tasks.

Three-Dimensional Evaluation. We assessed the scientific reasoning capabilities of current state-of-the-art T2I models through a three-dimensional evaluation. Specifically, we evaluated: the alignment between implicit prompts and (1) images generated from implicit prompts, (2) images generated from explicit prompts, and (3) images generated from superficial prompts. For each alignment evaluation, we selected one implicit, explicit, and superficial prompt forming a tuple from SCIENCE-T2I S and SCIENCE-T2I C, respectively. We generated two images per prompt using the T2I models and calculated the average SCISCORE. The experimental results are presented in Table 3.

Analysis: Explicit Prompt Alignment. The experiment results in Table 3 reveals that the FLUX series models [1] consistently outperform the Stable Diffusion series on explicit prompt alignment. In particular, SDv1.5 [58] exhibits a significant performance gap when compared to the other models. However, while SCISCORE effectively evaluates prompt-image alignment, its scores are inherently relative rather than absolute. Although it can discern which T2I model exhibits superior alignment with a given prompt, it cannot independently quantify alignment for a single T2I model. To comprehensively evaluate explicit prompt alignment, we include additional results and analysis in Appendix S14.

Analysis: Reasoning Capability. To evaluate the capacity of contemporary T2I models to reason from implicit prompts, we introduce a metric termed "Normalized Difference" (ND). The ND metric quantifies the degree to which images generated from implicit prompts resemble those generated from explicit or superficial prompts. Analysis of the data presented in Table 3 reveals a significant deficiency in the ability of current T2I models to effectively interpret implicit meanings. Specifically, these models predominantly generate outputs that correspond to the literal, surface-level aspects of the input prompts, rather than inferring or representing the underlying, implicit conceptualizations, such as those derived from scientific principles. This constraint is reflected in the observed ND scores, which for the majority of models fall below 50, yielding an average score of approximately 35.

Table 3: Performance of T2I Models on SCISCORE. ND is defined as $ND = (IP - SP) / (EP - SP)$. **Bold** values indicate the best performance, while underlined values represent the second-best performance.

T2I Model	SCIENCE-T2I S				SCIENCE-T2I C			
	SP	EP	IP	ND	SP	EP	IP	ND
Stable Diffusion v1.5 [58]	19.35	26.88	22.37	<u>40.11</u>	22.45	28.19	23.40	16.55
Stable Diffusion XL [54]	21.80	31.90	25.47	36.34	26.21	34.22	30.89	58.43
Stable Diffusion 3 [17]	18.99	32.53	22.31	24.52	24.01	34.65	27.88	36.37
FLUX.1[schnell] [1]	18.45	32.87	<u>24.43</u>	41.47	25.12	36.05	<u>29.66</u>	<u>41.54</u>
FLUX.1[dev] [1]	17.69	<u>32.85</u>	23.56	38.72	23.78	<u>34.70</u>	27.26	31.87

7. Experiment: Two-Stage T2I Model Fine-Tuning

7.1. Implementation Details

Training Setting. We first fine-tune FLUX.1[dev] [1] on SCIENCE-T2I using SFT in conjunction with LoRA [28] for 2,000 steps. This process generates LoRA weights intended for subsequent OFT. For the OFT phase, we randomly select 300 implicit prompts from SCIENCE-T2I to serve as the training set. During each epoch, 32 prompts are sampled, with each prompt paired with two images, and their corresponding SCISCORE is computed. Subject masks are extracted from the images using GroundingDINO [47]. Finally, the model undergoes fine-tuning for approximately 100 steps. Detailed training configurations can be found in Appendix S16.

Evaluation Setting. We construct two distinct prompt sets by extracting all implicit prompts from SCIENCE-T2I S and SCIENCE-T2I C. For evaluation, we generate five distinct images for each prompt and compute the average SCISCORE across these images to ensure robust results.

Relative Improvement Metric. To gain a more nuanced understanding of model performance beyond raw SCISCORE evaluation, we observe a consistent trend: SCISCORE achieved by the model when provided with explicit prompts invariably surpasses that obtained with implicit prompts. This principle allows us to establish an upper-bound performance estimate by calculating the average SCISCORE specifically for explicit prompts. To quantify the relative improvement achieved through finetuning, we introduce the **Relative Improvement** (RI) metric. Let $SciScore_B^{IP}$ and $SciScore_B^{EP}$ represent the SciScore of the base model under implicit and explicit prompts, respectively, and let $SciScore_F^{IP}$ denote the SciScore of the finetuned model under implicit prompts. The RI metric is then defined as:

$$RI = \frac{SciScore_F^{IP} - SciScore_B^{IP}}{SciScore_B^{EP} - SciScore_B^{IP}} \quad (21)$$

7.2. Results

Performance Gains from Two-stage Training. Table 4 substantiates the efficacy of both SFT and OFT in augmenting the performance of FLUX [1] on SCISCORE. Notably, the proposed methodology yields a substantial performance improvement, exceeding the baseline by over 50%, as evidenced

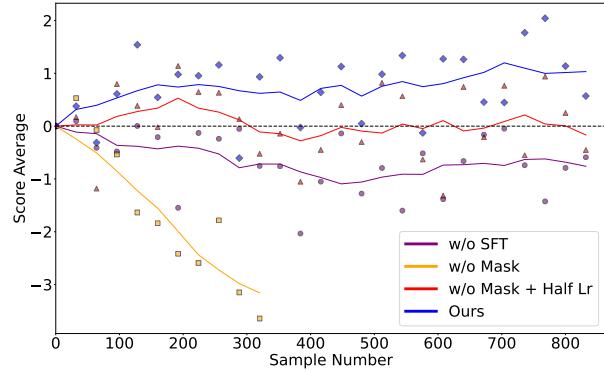


Fig. 6: Ablation study of two-stage training. The result illustrates the deviation from the initial baseline.

by the increase from 23.56 to 28.52 in SCIENCE-T2I S and from 27.26 to 30.11 in SCIENCE-T2I C. Furthermore, the results indicate that SFT contributes a greater performance gain relative to OFT. A comparative analysis of generated examples, presented in Figure 7, provides qualitative support for these quantitative findings.

Generalization to Complex Scenes. While SCIENCE-T2I mainly presents images and prompts in straightforward scenarios, the SCISCORE of finetuned FLUX [1] in SCIENCE-T2I C demonstrates a clear generalization to complex scene setups. This improvement indicates that the model has not merely memorized training examples. Rather, it has internalized underlying scientific principles through the finetuning process.

7.3. Ablation Study

In this section, we evaluate the individual contributions of each component within our proposed two-stage fine-tuning framework. Throughout the training process, we adhered to a consistent protocol for different settings: at each step, all implicit prompts from SCIENCE-T2I S were used as inputs for FLUX [1]. Then FLUX [1] generated two images per prompt, and the average SCISCORE was computed. Note that the results depicted in Figure 6 indicate the deviation from the baseline.

Necessity of SFT. In Figure 6, the blue line shows SFT performed before OFT, while the purple line illustrates the case without initial SFT. Both scenarios use identical configurations for OFT. The results demonstrate that initiating OFT with SFT leads to a more stable increase in SCISCORE. In contrast, OFT without preceding SFT does not improve SCISCORE. This discrepancy is likely due to the model’s limited ability to effectively learn from two suboptimal samples when SFT is not first applied. These observations highlight the critical role of starting with SFT to ensure the model trains within the distribution defined by the objective, facilitating effective OFT. Furthermore, these findings implicitly suggest the importance of a well-initialized base model for the subsequent online fine-tuning stage.

Masking Strategy As A Denoiser. Building upon an initial SFT checkpoint, we investigated the impact of masking strategy by testing two alternative configurations, represented by the red and yellow curves in Figure 6. Without masking (yellow curve), the performance became erratic, and the generated images began to exhibit signs of collapse. To address this, we attempted to stabilize

Table 4: SCISCORE on Various Methods. The best performance is highlighted with **bold** values.

Method	SCIENCE-T2I S		SCIENCE-T2I C	
	SCISCORE	RI	SCISCORE	RI
FLUX.1[dev]	23.56	/	27.26	/
+EP	32.85	/	34.70	/
+SFT	27.43	41.66	29.49	29.97
+SFT+OFT	28.52	53.39	30.11	38.31

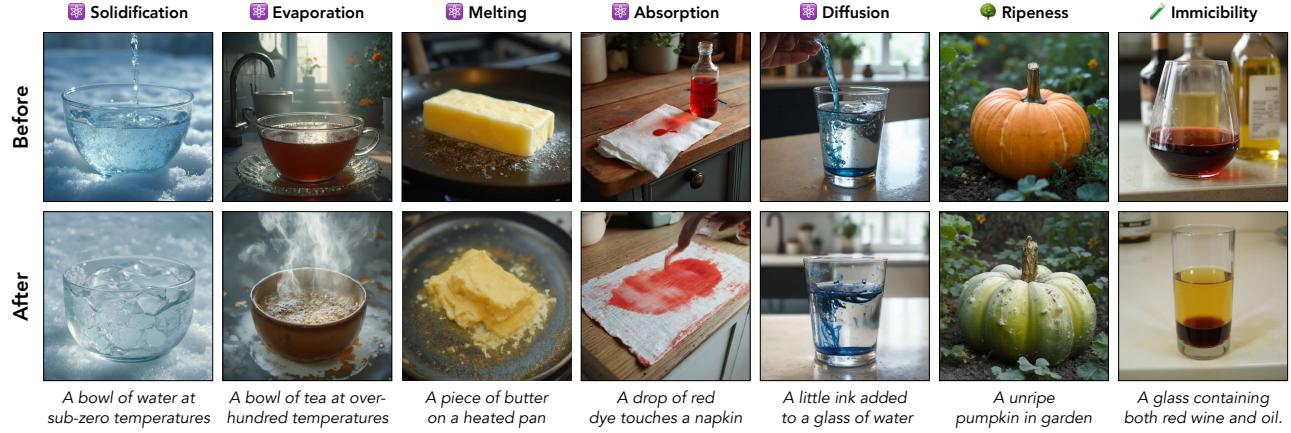


Fig. 7: Case study. The upper images are generated using the base FLUX.1[dev] [1], whereas the lower images are produced with our fine-tuning method. Each image pair utilizes an identical random seed to ensure consistency in comparison. Note that the displayed prompts are summaries of the original prompts.

training by halving the learning rate (red curve). While this adjustment prevented the collapse of the generated images, it did not lead to an increase of SCISCORE. This observation suggests that, without the masking strategy, the model tends to indiscriminately consider all features from the preferred images as equally important, effectively treating all features as ‘preferred’. However, only the visual features pertinent to the scientific principles contained in the prompt are truly relevant. This indiscriminate preference introduces substantial noise into the training process, hindering the model’s ability to learn effectively. In contrast, the model employing the masking strategy demonstrated a more stable increase on SCISCORE throughout training.

8. Conclusion

In summary, by leveraging our expert-annotated dataset, SCIENCE-T2I, which comprises over 20k adversarial image pairs and 9k prompts, we have developed a comprehensive framework for evaluating and enhancing image realism. Specifically, we introduce SCISCORE, a novel reward model designed to infuse scientific knowledge into image synthesis. Our results demonstrate that SCISCORE reaches human-level accuracy in aligning with scientific knowledge. Additionally, we propose a two-stage training framework for T2I models, utilizing SCISCORE as the reward model. This framework, which combines supervised fine-tuning with online fine-tuning, leads to significant performance improvements in generation tasks that require scientific reasoning.

Acknowledgments

We thank Alistair King for sharing insightful code, which was instrumental in our fine-tuning process. SX also acknowledges support from Open Path AI Foundation, Intel AI SRS, IITP grant funded by the Korean Government (MSIT) (No. RS-2024-00457882, National AI Research Lab Project), Amazon Research Award, and NSF Award IIS-2443404.

References

- [1] Flux. <https://blackforestlabs.ai/>.
- [2] Gpt-4o. <https://openai.com/index/hello-gpt-4o/>.
- [3] Ic-light. <https://openreview.net/pdf?id=u1cQYxRI1H>.
- [4] Laion-aesthetics. <https://laion.ai/blog/laion-aesthetics/>.
- [5] Dalle-3. <https://openai.com/index/dall-e-3/>.
- [6] Hritik Bansal, Zongyu Lin, Tianyi Xie, Zeshun Zong, Michal Yarom, Yonatan Bitton, Chen-fanfu Jiang, Yizhou Sun, Kai-Wei Chang, and Aditya Grover. Videophy: Evaluating physical commonsense for video generation. *arXiv preprint arXiv:2406.03520*, 2024.
- [7] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning, 2024. URL <https://arxiv.org/abs/2305.13301>.
- [8] Shidong Cao, Wenhao Chai, Shengyu Hao, and Gaoang Wang. Image reference-guided fashion design with structure-aware transfer by diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3525–3529, 2023.
- [9] Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stablevideo: Text-driven consistency-aware diffusion video editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 23040–23050, 2023.
- [10] Liang Chen, Shuai Bai, Wenhao Chai, Weichu Xie, Haozhe Zhao, Leon Vinci, Junyang Lin, and Baobao Chang. Multimodal representation alignment for image generation: Text-image interleaved control is easier than you think. *arXiv preprint arXiv:2502.20172*, 2025.
- [11] Wengling Chen and James Hays. Sketchygan: Towards diverse and realistic sketch to image synthesis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.
- [12] Zhe Chen, Jiannan Wu, Wenhui Wang, Weijie Su, Guo Chen, Sen Xing, Muyan Zhong, Qinglong Zhang, Xizhou Zhu, Lewei Lu, Bin Li, Ping Luo, Tong Lu, Yu Qiao, and Jifeng Dai. Internvl: Scaling up vision foundation models and aligning for generic visual-linguistic tasks, 2024. URL <https://arxiv.org/abs/2312.14238>.
- [13] Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-image generation. *arXiv preprint arXiv:2310.18235*, 2023.
- [14] Michael F Cohen and John R Wallace. *Radiosity and realistic image synthesis*. Morgan Kaufmann, 1993.
- [15] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.

- [16] Carles Domingo-Enrich, Michal Drozdal, Brian Karrer, and Ricky T. Q. Chen. Adjoint matching: Fine-tuning flow and diffusion generative models with memoryless stochastic optimal control, 2024. URL <https://arxiv.org/abs/2409.08861>.
- [17] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English, Kyle Lacey, Alex Goodwin, Yannik Marek, and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis, 2024. URL <https://arxiv.org/abs/2403.03206>.
- [18] Luca Eyring, Shyamgopal Karthik, Karsten Roth, Alexey Dosovitskiy, and Zeynep Akata. Reno: Enhancing one-step text-to-image models through reward-based noise optimization. *arXiv preprint arXiv:2406.04312*, 2024.
- [19] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models, 2023. URL <https://arxiv.org/abs/2305.16381>.
- [20] Hany Farid. Perspective (in)consistency of paint by text, 2022. URL <https://arxiv.org/abs/2206.14617>.
- [21] James A Ferwerda, Sumanta N Pattanaik, Peter Shirley, and Donald P Greenberg. A model of visual adaptation for realistic image synthesis. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 249–258, 1996.
- [22] Xingyu Fu, Muyu He, Yujie Lu, William Yang Wang, and Dan Roth. Commonsense-t2i challenge: Can text-to-image generation models understand commonsense?, 2024. URL <https://arxiv.org/abs/2406.07546>.
- [23] Donald P Greenberg, Kenneth E Torrance, Peter Shirley, James Arvo, Eric Lafortune, James A Ferwerda, Bruce Walter, Ben Trumbore, Sumanta Pattanaik, and Sing-Choong Foo. A framework for realistic image synthesis. In *Proceedings of the 24th annual conference on Computer graphics and interactive techniques*, pages 477–494, 1997.
- [24] Jianshu Guo, Wenhao Chai, Jie Deng, Hsiang-Wei Huang, Tian Ye, Yichen Xu, Jiawei Zhang, Jenq-Neng Hwang, and Gaoang Wang. Versat2i: Improving text-to-image models with versatile reward. *arXiv preprint arXiv:2403.18493*, 2024.
- [25] Minghao Guo, Bohan Wang, Pingchuan Ma, Tianyuan Zhang, Crystal Elaine Owens, Chuang Gan, Joshua B Tenenbaum, Kaiming He, and Wojciech Matusik. Physically compatible 3d object modeling from a single image. *arXiv preprint arXiv:2405.20510*, 2024.
- [26] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- [27] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.

- [28] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. URL <https://arxiv.org/abs/2106.09685>.
- [29] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering. *arXiv preprint arXiv:2303.11897*, 2023.
- [30] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. In *Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2023.
- [31] Ziwei Huang, Wanggui He, Quanyu Long, Yandi Wang, Haoyuan Li, Zhelun Yu, Fangxun Shu, Long Chan, Hao Jiang, Leilei Gan, and Fei Wu. T2i-factualbench: Benchmarking the factuality of text-to-image models with knowledge-intensive concepts, 2024. URL <https://arxiv.org/abs/2412.04300>.
- [32] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, July 2021. URL <https://doi.org/10.5281/zenodo.5143773>. If you use this software, please cite it as below.
- [33] Henrik Wann Jensen. *Realistic image synthesis using photon mapping*. AK Peters/crc Press, 2001.
- [34] Bingyi Kang, Yang Yue, Rui Lu, Zhijie Lin, Yang Zhao, Kaixin Wang, Gao Huang, and Jiashi Feng. How far is video generation from world model: A physical law perspective. 2024.
- [35] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022. URL <https://arxiv.org/abs/2206.00364>.
- [36] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022. URL <https://arxiv.org/abs/1312.6114>.
- [37] 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, 2023. URL <https://arxiv.org/abs/2305.01569>.
- [38] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- [39] Max Ku, Dongfu Jiang, Cong Wei, Xiang Yue, and Wenhui Chen. Viescore: Towards explainable metrics for conditional image synthesis evaluation. *arXiv preprint arXiv:2312.14867*, 2023.
- [40] Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Li, Yixin Fei, Kewen Wu, Xide Xia, Pengchuan Zhang, Graham Neubig, and Deva Ramanan. Evaluating and improving compositional text-to-visual generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5290–5301, 2024.
- [41] Bo Li, Yuanhan Zhang, Dong Guo, Renrui Zhang, Feng Li, Hao Zhang, Kaichen Zhang, Peiyuan Zhang, Yanwei Li, Ziwei Liu, and Chunyuan Li. Llava-onevision: Easy visual task transfer, 2024. URL <https://arxiv.org/abs/2408.03326>.

- [42] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation, 2022. URL <https://arxiv.org/abs/2201.12086>.
- [43] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models, 2023. URL <https://arxiv.org/abs/2301.12597>.
- [44] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. URL <https://arxiv.org/abs/1405.0312>.
- [45] Yunlong Lin, Tian Ye, Sixiang Chen, Zhenqi Fu, Yingying Wang, Wenhao Chai, Zhaohu Xing, Lei Zhu, and Xinghao Ding. Aglldiff: Guiding diffusion models towards unsupervised training-free real-world low-light image enhancement. *arXiv preprint arXiv:2407.14900*, 2024.
- [46] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling, 2023. URL <https://arxiv.org/abs/2210.02747>.
- [47] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, and Lei Zhang. Grounding dino: Marrying dino with grounded pre-training for open-set object detection, 2024. URL <https://arxiv.org/abs/2303.05499>.
- [48] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. URL <https://arxiv.org/abs/1711.05101>.
- [49] Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. Llmscore: Unveiling the power of large language models in text-to-image synthesis evaluation. *arXiv preprint arXiv:2305.11116*, 2023.
- [50] Fanqing Meng, Jiaqi Liao, Xinyu Tan, Wenqi Shao, Quanfeng Lu, Kaipeng Zhang, Yu Cheng, Dianqi Li, Yu Qiao, and Ping Luo. Towards world simulator: Crafting physical commonsense-based benchmark for video generation. *arXiv preprint arXiv:2410.05363*, 2024.
- [51] Fanqing Meng, Wenqi Shao, Lixin Luo, Yahong Wang, Yiran Chen, Quanfeng Lu, Yue Yang, Tianshuo Yang, Kaipeng Zhang, Yu Qiao, et al. Phybench: A physical commonsense benchmark for evaluating text-to-image models. *arXiv preprint arXiv:2406.11802*, 2024.
- [52] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. URL <https://arxiv.org/abs/2203.02155>.
- [53] Dong Huk Park, Samaneh Azadi, Xihui Liu, Trevor Darrell, and Anna Rohrbach. Benchmark for compositional text-to-image synthesis. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021.
- [54] 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, 2023. URL <https://arxiv.org/abs/2307.01952>.

- [55] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. URL <https://arxiv.org/abs/2103.00020>.
- [56] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. URL <https://arxiv.org/abs/2305.18290>.
- [57] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems*, 36, 2024.
- [58] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. URL <https://arxiv.org/abs/2112.10752>.
- [59] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016.
- [60] Ayush Sarkar, Hanlin Mai, Amitabh Mahapatra, Svetlana Lazebnik, D. A. Forsyth, and Anand Bhattad. Shadows don't lie and lines can't bend! generative models don't know projective geometry...for now, 2024. URL <https://arxiv.org/abs/2311.17138>.
- [61] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. URL <https://arxiv.org/abs/1707.06347>.
- [62] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. URL <https://arxiv.org/abs/2010.02502>.
- [63] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. URL <https://arxiv.org/abs/2009.01325>.
- [64] Jiao Sun, Deqing Fu, Yushi Hu, Su Wang, Royi Rassin, Da-Cheng Juan, Dana Alon, Charles Herrmann, Sjoerd van Steenkiste, Ranjay Krishna, et al. Dreamsync: Aligning text-to-image generation with image understanding feedback. In *Synthetic Data for Computer Vision Workshop@CVPR 2024*, 2023.
- [65] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL <https://arxiv.org/abs/1706.03762>.
- [66] 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, 2023. URL <https://arxiv.org/abs/2311.12908>.
- [67] 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.

- [68] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution, 2024. URL <https://arxiv.org/abs/2409.12191>.
- [69] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL <https://arxiv.org/abs/2201.11903>.
- [70] 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.
- [71] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation. *arXiv preprint arXiv:2304.05977*, 2023.
- [72] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Qimai Li, Weihan Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model, 2024. URL <https://arxiv.org/abs/2311.13231>.
- [73] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Weihan Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8941–8951, 2024.
- [74] Tian Ye, Sixiang Chen, Wenhao Chai, Zhaochu Xing, Jing Qin, Ge Lin, and Lei Zhu. Learning diffusion texture priors for image restoration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2524–2534, 2024.
- [75] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training, 2023. URL <https://arxiv.org/abs/2303.15343>.
- [76] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N. Metaxas. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.
- [77] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. Stackgan++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 41(8):1947–1962, 2018.
- [78] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018.

Supplementary Material

The supplementary material is structured as follows:

- Rational to prioritize rewriting capability in Section S1.
- Detailed task descriptions in Section S2.
- Detailed task classification in Section S3.
- Detailed data curation process in Section S4.
- Comparison with related benchmarks in Section S5.
- The training settings for SCISCORE in Section S6.
- Detailed benchmarking configuration in Section S7.
- Setup of SCIENCE-T2I S and SCIENCE-T2I C in Section S8.
- Detailed baseline setup and evaluation process for SCISCORE in Section S9.
- Additional results of SCISCORE in Section S10.
- Additional analysis of SCISCORE in Section S11.
- Qualitative analysis of IEE in Section S12.
- More results on benchmarking T2I model in Section S13.
- More results and analysis of explicit prompt alignment in Section S14.
- Details of two-stage training in Section S15.
- Detailed two-stage training settings in Section S16.
- Limitation and future direction in Section S17.
- Additional results of online fine-tuning in Section S18.

S1. Rationale to Prioritize Rewriting Capability

In the development of our study, we considered incorporating additional reasoning tasks, such as reflection-based tasks [20] that evaluate the consistency between objects and their reflections. However, these tasks present unique challenges that influenced our decision to exclude them.

Reflection-based tasks require the representation of precise geometric details to capture the relationships between objects and their reflections. Such intricate geometric information cannot be fully conveyed through textual descriptions alone. Consequently, current text-to-image generation models face difficulties in producing both correct and incorrect images for these tasks. This limitation hampers the creation of a consistent and valid dataset necessary for evaluating generative models on reflection-based generation.

Given these constraints, we prioritized tasks that can be effectively rephrased and consistently described using language-based prompts. This ensures generative models can interpret and generate the required images more reliably, thereby facilitating robust data collection and analysis. By focusing on linguistically describable tasks, we enhance the reproducibility and validity of our findings.

Tasks that lack this flexibility, particularly those requiring detailed geometric representation [60] and those subtle light-related features [3] beyond the capacity of textual prompts, are reserved for future exploration.

S2. Detailed Task Descriptions

In this section, we provide detailed descriptions of the tasks incorporated into our study. These tasks are designed to evaluate various biological, chemical, and physical phenomena presented in SCIENCE-T2I. Additionally, illustrative examples from SCIENCE-T2I are presented in Figure S12,S13,S14 to demonstrate the tasks.

- **Light Requirement (LR):** Plants change color and leaf size based on adequate or insufficient light exposure.
- **Watering Requirement (WR):** Plants exhibit differences in foliage health, wilting, and growth when receiving sufficient or inadequate water, leading to reduced growth.
- **Ripeness (RI):** Fruits alter their color and texture significantly when they are ripe compared to when they are unripe.
- **Seasonal Change (SC):** Plants display variations in leaf color, size, and blooming patterns across different seasons.
- **Flame Reaction (FR):** Chemical substances naturally produce their distinct flame colors vividly when burned.
- **Immiscibility (IM):** Two liquids either mix uniformly or separate into layers based on their chemical properties.
- **Rust (RU):** Metals appear shiny, smooth, and reflective before oxidation, and corroded, flaky, and brittle after rusting.
- **Absorption (AB):** A solid either soaks up a liquid or repels it, depending on their material properties.
- **Buoyancy (BU):** Substances either float on or sink in water based on their density relative to water.
- **Diffusion (DI):** When a small amount of liquid is added, it either disperses uniformly or remains separate.
- **Electricity (EL):** Electronic products change appearance, such as glowing or sparking, when electric current is applied.
- **Evaporation (EV):** Liquids boil and produce vapor when reaching boiling points; otherwise, they remain calm.
- **Gravity (GR):** Objects appear differently positioned when influenced by gravity versus in a gravity-free environment.
- **Liquidation (LI):** Air condenses into water droplets on surfaces cooled below room temperature.
- **Melting (ME):** Objects transition from solid to liquid, changing shape and structure upon reaching melting points.
- **Solidification (SO):** Liquids become solids, altering their form and texture when cooled below solidification points.

S3. Detailed Task Classification

During our investigation, particularly in the data curation phase, we observed that all the scientific phenomena involved can be uniformly represented using a **subject + condition** framework. Specifically, all tasks involve implicit prompts structured in this manner. For example, the prompt *an unripe apple* comprises the subject *apple* and the condition *unripe*; similarly, *a laptop without electricity* includes the subject *laptop* and the condition *without electricity*. Building on this observation, we identified that, for each task, the component requiring scientific reasoning can be closely associated either with the subject or with the condition. We classify these tasks as *subject-oriented* tasks and *condition-oriented* tasks, depending on the reasoning focus.

Subject-Oriented Tasks. In subject-oriented tasks, the necessity for scientific reasoning arises primarily from the subject's properties. In these tasks, different subjects under the same condition exhibit different visual features due to their inherent characteristics. For example, the *buoyancy* task is subject-oriented because different objects placed in water either float or sink depending on their densities relative to water, which is an intrinsic property of the subjects.

Condition-Oriented Tasks. In condition-oriented tasks, scientific reasoning is predominantly associated with the condition applied to the subject. In these tasks, varying subjects under a single condition consistently yield similar visual features. For instance, the *gravity* task is condition-oriented since many subjects behave similarly under identical gravity: floating in "without gravity" and resting on the ground in "normal gravity" scenarios.

Task Classification. As depicted in Figure S1, SCIENCE-T2I features a clear division: ten predefined tasks fall under the condition-oriented category, while the remaining six are classified as subject-oriented tasks, respectively.

Subject-oriented Task			Seasonal Change 240	Condition-oriented Task	Rust 288	Light Requirement 201		Water Requirement 201		Ripeness 192
	Immiscibility 218	Flame Reaction 206				Liquidation 195	Electricity 186	Melting 192	Gravity 221	
Buoyancy 204		Diffusion 142		Solidification 117						
	Absorption 160									

Fig. S1: Task classification of SCIENCE-T2I.

S4. Detailed Data Curation Process

In this section, we provide a detailed overview of our data curation process. We describe the methods used for generating subject-based prompts, synthesizing images, and establishing criteria for image selection.

Subject-Based Prompt. For each task, we first employ GPT-4o [2] to define a comprehensive set of templates for the implicit prompts. These templates act as structured frameworks that capture the essence of the reasoning required while allowing for variability in the objects or substances involved. Using the templates, GPT-4o [2] generated a variety of implicit prompts by inserting appropriate objects or substances into the placeholders. Then for each implicit prompt, we used GPT-4o [2] to generate the corresponding explicit prompt and superficial prompt. An illustration of this instruction process is provided in Figure S2.

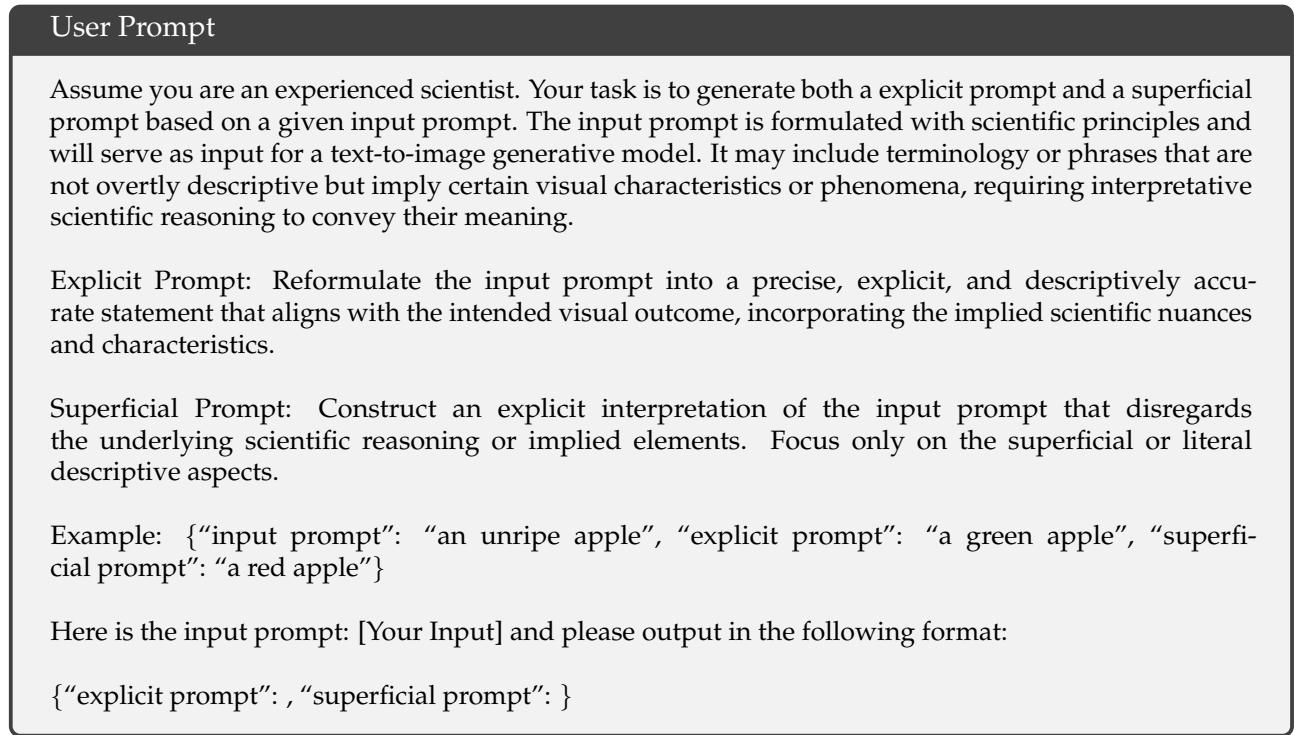


Fig. S2: Framework For Prompt Collection. This figure presents a detailed workflow for generating explicit and superficial prompts from implicit prompts using GPT-4o [2].

Synthetic Image Generation. The limited availability of images relevant to our specific scientific reasoning tasks within existing datasets and online resources necessitated the generation of synthetic images. However, we could not arbitrarily select a text-to-image model, as this choice directly affects both the quality of the generated data and the efficiency of data acquisition. Among the numerous advanced models available, our choice was informed by a comprehensive evaluation of several key factors. Below, we outline the primary considerations that guided our decision:

- **Descriptive Text-Image Alignment:** The core objective involves generating images that accurately reflect both explicit and superficial prompts. This necessitates a model with a robust capability to align textual descriptions with corresponding visual elements. Meanwhile, effective text-image alignment is also paramount for efficient data collection.
- **Realistic Style Consistency:** Our reasoning-based tasks are fundamentally grounded in scientific principles and real-world phenomena. Consequently, it is imperative that the generated images exhibit a style that reflects realism rather than abstract or cartoonish representations.

Based on these criteria, we conducted a qualitative evaluation of several state-of-the-art text-to-image models, including Stable Diffusion XL [54], Stable Diffusion 3 [17], DALLE 3 [5], and FLUX.1[dev] [1]. As illustrated in Figure S3, FLUX.1[dev] [1] consistently outperformed the other models in both text-image alignment and realistic style consistency. Therefore, FLUX.1[dev] [1] was selected as the model for synthetic image generation.



Fig. S3: Comparative Data Analysis. Models such as SDXL [54], SD3 [17], and DALLE3 [5] occasionally failed to align generated images accurately with the provided textual descriptions. Meanwhile, FLUX.1[dev] [1] demonstrated superior performance, producing the most realistic images among all evaluated models.

Criteria For Image Curation. As outlined in Section 3, the scientific principles inherent in the implicit prompt confer distinct visual features to the subject matter. During the image generation process for SCIENCE-T2I, particular emphasis was placed on the regions where these visual features are manifested. Our primary objective was to ensure that these regions accurately represent the concepts in alignment with the underlying scientific principles specified in the prompts. To achieve this, we established stringent criteria for the images, specifically: (1) minimizing noise and (2) preventing the introduction of irrelevant semantic information. As illustrated in Figure S12,S13,S14, we accomplish this by selecting data with the simplest possible backgrounds, such as solid colors. Additionally, we filter the data to ensure that the regions of interest are as large as possible, thereby maximizing the prominence of the visual features.

S5. Comparison with Related Benchmarks

In order to provide a comprehensive understanding of how SCIENCE-T2I distinguishes itself from other existing datasets, we present a detailed comparison in Table S1. The key distinguishing features and advantages of SCIENCE-T2I are as follows: (1) it enables the direct utilization of SCISCORE for benchmark T2I generation, significantly enhancing efficiency compared to approaches based on LMMs; (2) its test set is uniquely designed to also serve as a benchmark for evaluating LMMs, offering dual functionality; and (3) it includes a large-scale training set that not only supports the training of generative models but also facilitates advancements in multimodal research.

Table S1: Comparison with related benchmarks.

Benchmark	Type	Category	Training Set	Evaluation	
				Generation	LMM
Commonsense-T2I [22]	Commonsense	5	✗	✓	✗
T2I-FactualBench [31]	Commonsense	8	✗	✓	✗
PhyBench [51]	Science	31	✗	✓	✗
SCIENCE-T2I (Ours)	Science	16	✓	✓	✓

S6. Detailed Training Settings for SCISCORE

This section provides a overview of the hyper-parameter settings utilized during the training of SCISCORE. The key parameters, including batch size, learning rate, and optimizer configurations, are summarized in Table S2.

Table S2: Hyper-parameter settings used for training SCISCORE.

Hyper-parameters	SCISCORE
batch size	128
learning rate	2×10^{-6}
learning rate schedule	cosine
weight decay	0.3
training steps	600
warmup steps	150
optimizer	AdamW [48]
λ	0.25

Table S3: Configurations of T2I models mentioned in Section 6.4.

T2I Model	Guidance Scale	Inference Step
SDv1.5 [58]	7.5	50
SDXL [54]	5.0	50
SD3 [17]	7.0	28
FLUX.1[schnell] [1]	0.0	4
FLUX.1[dev] [1]	0.0	30

S7. Detailed Benchmarking Configuration

To facilitate equitable comparisons among the different T2I models, we standardized the output image resolution to 1024×1024 pixels for all models. Table S3 summarizes the configuration parameters used for each model, including the guidance scale and the number of inference steps.

S8. Setup of SCIENCE-T2I S and SCIENCE-T2I C

For the evaluation of SCISCORE compared with other VLMs and LMMs, two meticulously curated test sets are employed, each manually annotated and subjected to a stringent iterative review process by domain experts. This process involved cross-referencing the annotators' specialized knowledge with authoritative online sources to ensure accuracy and consistency. The validation procedure was repeated until unanimous consensus was achieved among all annotators, thereby enhancing the reliability of the test sets. These sets are strategically designed to evaluate the model's performance across varying levels of complexity and are characterized as follows:

- **SCIENCE-T2I S:** This test set closely replicates the stylistic and structural attributes of the training data. It emphasizes simplicity by focusing on specific regions and strictly adhering to the annotation criteria in Section S4. The goal of SCIENCE-T2I S is to assess the model's performance on data stylistically similar to its training set.
- **SCIENCE-T2I C:** This test set challenges the model in more complex scenarios, introducing contextual elements like explicit scene settings and diverse scenarios. Prompts in SCIENCE-T2I C may include phrases such as "in a bedroom" or "on the street," adding spatial and contextual variability. This complexity evaluates the model's ability to adapt to nuanced, less constrained environments.

S9. Detailed Baseline Setup and Evaluation Process for SCISCORE

This section provides detailed descriptions of the baseline setups employed to evaluate the performance of SCISCORE, as well as the evaluation process.

Evaluation Overview. The evaluation involves presenting one implicit prompt alongside two images: one aligned with the corresponding explicit prompt and the other with the superficial prompt. Models or humans are tasked with selecting the image that best corresponds to implicit prompt.

Vision-Language Models (VLMs). We employ three VLMs as baseline models, CLIP-H [32], BLIP-2 [43] and SigLIP [75]. The reward computation involves encoding the implicit input prompt and the input image using their respective text and image encoders. Subsequently, we apply the scoring mechanism described in Section 4.1 to evaluate the alignment between the text and image pairs.

Language Multimodal Models (LMMs). As a baseline for LMMs, we utilize open-source models like LLaVA-OV [41], Qwen2-VL [68], and InternVL [12], along with the proprietary model GPT-4o-mini [2]. For GPT-4o-mini [2], we conduct evaluations under two distinct settings: one without employing the CoT reasoning approach [69] and another incorporating CoT [69] to facilitate step-by-step reasoning. Specifically, we prompt LMMs to choose between two images by selecting either "the first" or "the second." Recognizing that the model may exhibit insensitivity to the order of image presentation, we mitigate this potential bias by conducting the evaluation twice, each time with the order of the input images reversed. We then compute the average accuracy across these two evaluations to obtain a more robust and reliable performance measure. The complete instruction set is detailed in Figure S4.

User Prompt

You will be presented with a textual prompt followed by two visual images. Your task is to critically analyze and compare both images, selecting the one that most accurately aligns with and represents the overall meaning of the given prompt. **First, you should imagine how an ideal image would look based on the prompt, and then describe both images in detail. Finally, combining your initial visualization with the descriptions of the two images, you should select the image that most effectively conveys the intended meaning of the prompt, providing a reasoned justification for your choice.**

Here is the input: {"prompt": [Your Input Prompt], "image-1": [Your Input Image], "image-2": [Your Input Image]}

Please output in the following format:

{'imagination': , 'description of image-1': , 'description of image-2': , 'justification for choice': , 'final choice': }

Fig. S4: Instruction For GPT Evaluation. Text segments in red are incorporated to facilitate CoT [69] reasoning.

Human Evaluation. To provide a human performance baseline, we collected data from 10 human evaluators, all of whom hold at least a college degree, primarily in science or engineering disciplines. This selection criterion ensures that the evaluators possess foundational scientific knowledge necessary to perform inference tasks. Moreover, it is critical to distinguish between human roles in evaluation and data curation. During curation, experts verify dataset accuracy using scientific literature and online knowledge bases. In contrast, evaluators in this study, distinct from curators, rely solely on their own knowledge to select the better option, which accounts for the absence of perfect performance.

S10. Additional Results of SCISCORE

In this section, we expand on the results presented in Table 1 by providing a more detailed breakdown of accuracy metrics for each category in Tables S4 and Table S5. This allows for a finer-grained evaluation of SCISCORE’s performance across different task categories. The abbreviations used in these tables are consistent with those outlined in Appendix S2. The extended results demonstrate that SCISCORE consistently outperforms baseline models across the majority of tasks. Furthermore, SCISCORE achieves perfect accuracy (100%) on several specific tasks, underscoring its effectiveness and robustness in diverse scenarios.

Table S4: Performance comparison on SCIENCE-T2I S across different categories. Best overall performance is in **bold**. Green highlights highest accuracy among VLMs, and blue highlights highest among LMMs.

Model	ME	DI	EL	SO	IM	EV	AB	LI	FR	SC	RI	RU	LR	WR	BU	GR
CLIP-H [32]	25.00	71.43	47.62	40.48	54.17	26.67	57.14	77.78	73.33	81.48	34.62	16.67	62.22	31.11	63.89	78.33
BLIPScore [42]	56.94	50.00	52.38	44.05	53.12	20.00	38.10	33.33	76.67	58.33	38.46	42.86	76.67	38.89	50.00	47.50
SigLIP ViT-SO-14 [75]	44.44	83.33	47.62	45.24	60.42	63.33	57.14	58.33	62.22	83.33	53.85	23.81	46.67	33.33	58.33	78.33
LLaVA-OV-7B [41]	36.11	75.00	77.38	55.95	55.21	100.00	38.10	95.83	48.89	78.70	46.15	59.52	51.11	72.22	45.83	92.50
Qwen2-VL-7B [68]	26.39	70.24	77.38	42.86	66.67	68.33	40.48	84.72	52.22	95.37	34.62	73.81	57.78	67.78	47.22	83.33
InternVL2.5-8B [12]	41.67	63.10	72.62	52.38	56.25	91.67	47.62	90.28	52.22	84.26	75.00	69.05	84.44	90.00	55.56	96.67
GPT-4o mini	36.11	77.38	82.14	35.71	65.63	100.00	33.33	76.39	58.89	97.22	53.85	95.24	96.67	83.33	56.94	71.31
GPT-4o mini+ CoT [69]	36.11	85.71	86.90	45.24	68.75	100.00	33.33	81.94	56.67	98.15	61.54	97.62	96.67	88.89	52.78	80.33
Human Eval	98.15	65.87	95.63	86.11	77.78	100.00	66.67	82.08	80.95	90.74	94.62	92.86	96.89	99.56	74.55	92.99
SCISCORE (ours)	100.00	97.62	100.00	90.48	68.75	100.00	71.43	100.00	97.78	100.00	100.00	100.00	100.00	100.00	66.67	98.33

Table S5: Performance comparison on SCIENCE-T2I C across different categories. Best overall performance is in **bold**. Green highlights highest accuracy among VLMs, and blue highlights highest among LMMs.

Model	ME	DI	EL	SO	IM	EV	AB	LI	FR	SC	RI	RU	LR	WR	BU	GR
CLIP-H [32]	66.67	78.57	21.43	57.14	50.00	0.00	64.29	66.67	46.67	88.89	75.00	35.71	80.00	60.00	58.33	75.00
BLIPScore [42]	58.33	50.00	28.57	42.86	62.50	50.00	29.17	60.00	75.00	54.17	57.14	53.33	46.67	62.50	40.00	
SigLIP ViT-SO-14 [75]	58.33	85.71	42.86	42.86	62.50	30.00	64.29	41.67	60.00	94.44	83.33	28.57	46.67	53.33	66.67	95.00
LLaVA-OV-7B [41]	45.83	71.43	64.29	71.43	59.38	100.00	57.14	79.17	36.67	77.78	79.17	53.57	63.33	86.67	75.00	100.00
Qwen2-VL-7B [68]	41.67	64.29	67.86	57.14	53.13	55.00	57.14	70.83	36.67	94.44	75.00	60.71	93.33	96.67	66.67	100.00
InternVL2.5-8B [12]	58.33	92.86	57.14	53.57	71.88	95.00	57.14	70.83	50.00	75.00	87.50	75.00	90.00	93.33	75.00	97.50
GPT-4o mini	67.65	67.86	64.29	50.00	68.75	90.00	50.00	75.00	53.33	88.89	87.50	89.29	100.00	83.33	54.17	97.50
GPT-4o mini+ CoT [69]	67.65	85.71	85.71	57.14	68.75	95.00	32.14	79.17	50.00	88.89	87.50	92.86	100.00	93.33	41.67	100.00
Human Eval	91.03	66.75	90.87	77.55	86.61	95.71	78.57	76.79	77.14	96.83	83.78	92.86	88.57	84.76	83.33	98.57
SCISCORE (ours)	100.00	85.71	85.71	92.86	81.25	100.00	71.43	100.00	100.00	100.00	100.00	92.86	100.00	100.00	41.67	100.00

S11. Additional Analysis

In this section, we present further in-depth analysis pertaining to the results and observations discussed in Section 6.2.

Performance of VLMs Approaches Random Guessing. The performance of CLIP-H [26], BLIP-Score [42] and SigLIP [75], is notably suboptimal, with accuracy levels hovering around 50% across both test sets. This is further corroborated by the ROC curves presented in Figure S6, which illustrate the performance of SCISCORE, CLIP-H [32], BLIPScore [42] and SigLIP [75]. The Area Under the Curve (AUC) scores (variable area in the figure) for these VLMs are relatively low, indicating that these models perform only marginally better than a random classifier. This limitation can be primarily attributed to the pretraining phase, where the majority of textual data are highly descriptive and explicitly reference their corresponding visual content. Consequently, during inference, the text encoder tends to rely heavily on these descriptive terms within the prompt. When a test prompt is associated with two images that both contain the main elements described in the prompt, the model struggles to differentiate between them effectively. This ambiguity leads to performance that is comparable to random guessing. In contrast, SCISCORE demonstrates superior efficacy, with a nearly optimal AUC score, indicating a high level of discriminative power and robustness in classification performance.



Fig. S5: Qualitative Results of IEE.
Images enclosed by green borders denote the correct selection in each pair.

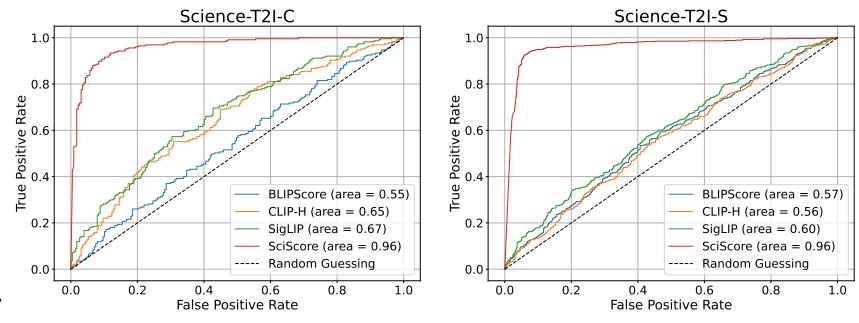


Fig. S6: ROC curves of CLIP, BLIPScore, SigLIP and SCISCORE.

Limitations of LMMs in Vision-Based Scientific Reasoning. Despite being equipped with an extensive knowledge base in their LLM backbone, current LMMs, including GPT-4o-mini [2], still struggle to achieve satisfactory performance in our vision-based scientific reasoning tasks, even when incorporating CoT [69] techniques. We posit that the primary reasons for this inadequate performance are twofold. First, the model exhibits a limited capacity to accurately capture and interpret the complex visual features inherent in scientific data, such as intricate diagrams, graphs, and microscopic images, which are crucial for tasks that rely heavily on visual information. This limitation hampers the model's ability to effectively integrate visual inputs with its existing knowledge base, leading to superficial or incorrect interpretations. Second, during the inference process, the model tends to generate reasoning chains that contain internal contradictions and inconsistencies, undermining the overall reliability and coherence of its scientific reasoning. These contradictory reasoning patterns within the CoT [69] framework suggest a fundamental challenge in maintaining logical consistency when processing and synthesizing information from visual sources, especially when dealing with complex or ambiguous data. To substantiate these claims, we present qualitative results in Figure S10, which illustrate specific

instances where GPT-4o-mini [2] fails to accurately interpret visual data and produces reasoning sequences that are internally conflicting and logically flawed.

SCISCORE Achieves Human-Level Performance. This enhanced efficacy can be primarily attributed to the inherent limitations in the specialized expertise of human evaluators. Although these evaluators typically possess undergraduate or advanced degrees and maintain a foundational understanding of relevant scientific domains, their knowledge bases are finite and often constrained by the boundaries of their specific areas of expertise. Such limitations can impede their ability to comprehensively assess all instances within diverse and extensive test sets, particularly when confronted with novel or interdisciplinary examples that lie outside their immediate knowledge scope. In contrast, SCISCORE leverages extensive contextual knowledge acquired from the training data, enabling it to generalize effectively and maintain consistent performance across diverse and challenging test scenarios.

S12. Qualitative Analysis of IEE

Qualitative results, as shown in Figure S5, demonstrate the effectiveness of incorporating IEE loss at an appropriate rate. The examples presented focus on the model’s ability to capture fine-grained and nuanced details. In the first two pairs, the task involves distinguishing between the frozen and liquid states of various liquids, which relies on subtle differences in transparency—frozen water exhibits lower transparency compared to liquid water. The third example pertains to a localized region within the image, where the model must determine whether the screen within this small region displays meaningful content. By incorporating IEE loss, the model enhances its visual discrimination and contextual analysis capabilities, enabling it to make more accurate and context-aware predictions.

S13. More Results on Benchmarking T2I Model

In this section, we further employ SCISCORE to benchmark additional T2I models. Due to budgetary constraints, our evaluation is limited to open-source models. The results are presented in Table S6.

Table S6: Performance of T2I Models on SCISCORE. **Bold** values indicate the best performance.

T2I Model	Size	SCIENCE-T2I S				SCIENCE-T2I C			
		SP	EP	IP	ND	SP	EP	IP	ND
Stable Diffusion 3.5	medium	20.40	34.29	24.11	26.71	24.67	36.14	29.32	40.54
	large	20.11	33.72	24.39	31.45	24.68	35.11	29.28	44.10
	turbo	19.37	31.57	22.71	27.38	23.86	33.49	27.38	36.55

S14. More Results on Explicit Prompt Alignment

While SCISCORE effectively evaluates the alignment between an implicit prompt and an image, it shares a common limitation inherent to all CLIP-based models [55]: the scores are only meaningful when comparing different pairs. In other words, SCISCORE can indicate that one prompt-image pair has better alignment than another but does not provide an absolute measure. To overcome this limitation in the context of the Explicit Prompt Alignment evaluation, we have developed systematic grading criteria for each tuple in SCIENCE-T2I S and SCIENCE-T2I C to assess alignment comprehensively. Inspired by PhyBench [51], our grading process is divided into two distinct aspects:

- **Main Subject Alignment (Scene Score, SS):** This aspect evaluates whether all descriptive visual content specified in the prompt is present in the corresponding image.
- **Implicit Visual Alignment (Reality Score, RS):** This aspect assesses whether the implicit visual elements, derived from underlying scientific principles present in the implicit prompt, are accurately represented in the image.

For illustrative purposes, we present several representative examples in Figure S8.

Evaluation Setup. After establishing the grading criteria, we selected all implicit prompts and their corresponding explicit prompts from SCIENCE-T2I S and SCIENCE-T2I C. Using T2I models, we generated two images for each explicit prompt. These images were then evaluated by GPT-4o-mini [2] following the instructions detailed in Figure S7. For each evaluation, GPT-4o-mini was provided with a single image and produced its corresponding SS and RS. However, within the two-tiered grading framework, it is redundant to assess RS if an image does not achieve a full score in SS, as reality grading assumes the presence of the main subject specified in the prompt. Consequently, we ask GPT-4o-mini to set RS as zero for any image that did not attain a full score in SS. The average SS and RS are summarized in Table S7. To further substantiate the effectiveness of the GPT-based evaluation, we examine the concordance between GPT assessments and human evaluations.

Table S7: Performance of T2I Models on Explicit Prompt Alignment. The Full Score (FS) is defined as $FS = SS + RS$. The Percentage of Expectation (PoE) is calculated by dividing the score by its expected value. **Bold** values indicate the best performance, while underlined values represent the second-best performance.

T2I Model	SCIENCE-T2I S				SCIENCE-T2I C			
	SS	PoE	FS	PoE	SS	PoE	FS	PoE
Stable Diffusion v1.5 [58]	1.298	64.90	2.470	49.40	1.261	63.05	2.446	48.92
Stable Diffusion XL [54]	<u>1.718</u>	<u>85.90</u>	3.510	70.20	1.679	83.95	3.360	67.20
Stable Diffusion 3 [17]	1.786	89.30	3.898	77.96	1.780	89.00	3.836	76.72
FLUX.1[schnell] [1]	1.730	86.50	<u>3.730</u>	<u>74.60</u>	<u>1.772</u>	<u>88.60</u>	<u>3.825</u>	76.50
FLUX.1[dev] [1]	1.720	86.00	3.641	72.82	1.702	85.10	3.676	73.52
Expectation	2.000	100.00	5.000	100.00	2.000	100.00	5.000	100.00

Relative Weakness in Scientific Scene Generation. The results in Table S7 reveal a relative weakness in the models' ability to generate outputs for complex scientific phenomena, as evidenced by the consistently lower average PoE of FS compared to SS. This suggests that the models exhibit weaker

User Prompt

Imagine you are an experienced scientist. Begin by evaluating the provided image using the specified scene composition criteria. If the image does not fully satisfy these criteria, assign a reality score of 0. However, if the scene meets all the criteria, proceed to assess its realism based on the given reality scoring guidelines, disregarding stylistic aspects and minor background details. Please first describe the image in detail and then adhere strictly to these criteria to ensure an accurate scoring of the image.

Here is the input: {"Prompt": [Your Input Prompt], "Scene Grading": [Your Input Scene Grading], "Reality Grading": [Your Input Reality Grading], "Image": [Your Input Image]}.

Please present your evaluation in the following format: {"description": , "scene score": , "reality score": }

Fig. S7: Image evaluation instruction. GPT generates SS and RS from the image, prompt, and grading criteria.

performance when generating images related to complex scientific phenomena compared to simpler subjects within prompts. A plausible explanation is that these phenomena often involve intricate features such as spatial relationships or uncommon object states (e.g., melting chocolate, a cup of frozen water), which are underrepresented in the models' pretraining data.

Concordance between GPT and Human. To assess the effectiveness of GPT-based evaluation methods, we designed an experiment aimed at demonstrating the alignment between GPT's judgments and those of human experts. Utilizing an established evaluation framework, we applied the same scoring methodology, which is detailed in Figure S7, to SCIENCE-T2I S and SCIENCE-T2I C. Human experts assigned scores based on these criteria, and after reaching consensus, we calculated the average scores. For explicit images, the human experts assigned an average SS of 2 and an average RS of 0; for superficial images, the average SS was 2 and the average RS was 3. Subsequently, we performed the same evaluation using GPT-4o-mini [2]. To quantify the correspondence between GPT-4o-mini's evaluations and those of the human experts, we calculated the human correspondence (HC) for both SS and RS. The human correspondence for SS is computed as:

$$HC_{SS} = \frac{SS}{2.0} \times 100 \quad (S22)$$

For RS, we used two separate formulas to compute the correspondence for explicit and superficial images. Specifically, for superficial images (SI) and explicit images (EI), HC for RS is calculated as:

$$HC_{RS}^{SI} = \left(1 - \frac{RS}{3.0}\right) \times 100, \quad HC_{RS}^{EI} = \frac{RS}{3.0} \times 100 \quad (S23)$$

The comparative results are presented in Table S8, which reveals that the high average agreement rate between GPT-4o-mini and human expert assessments demonstrates its reliability.

Table S8: Concordance Between GPT-4o-mini and Human Experts. The average agreement rate of over 80% demonstrates GPT-4o-mini's strong alignment with human expert assessments, highlighting its reliability.

Dataset	Input Type	SS	HC _{SS}	RS	HC _{RS}
SCIENCE-T2I S	EI	1.827	91.13	2.731	91.03
	SI	1.635	81.74	0.476	84.13
SCIENCE-T2I C	EI	1.855	92.73	2.490	83.00
	SI	1.630	81.50	0.636	78.79

Sampled Examples
<p>Example1: [</p> <p>"Prompt": "A transparent water-filled box holds a basketball, depicted realistically.",</p> <p>"Scene Grading": {</p> <ul style="list-style-type: none"> - 0 point: The picture does not feature a basketball inside a transparent box filled with water in any capacity. - 1 point: The picture shows a basketball, but it is not inside a transparent box. Alternatively, the basketball are in a transparent box, but there is no water present. - 2 points: The picture accurately depicts a basketball inside a transparent box filled with water. <p>},</p> <p>"Reality Grading": {</p> <ul style="list-style-type: none"> - 0 point: The basketball is completely sinking to the bottom of the water. - 1-2 point: The basketball is completely submerging in the water, but doesn't reach the bottom. Less mistakes will earn a higher score. - 3 points: The picture shows basketball floating on the surface of the water. <p>}]</p>
<p>Example2: [</p> <p>"Prompt": "A clear glass filled with water and oil, simple and realistic.",</p> <p>"Scene Grading": {</p> <ul style="list-style-type: none"> - 0 point: There is no glass or no liquid in the glass, or the scene is irrelevant (e.g., the focus is not on the glass or liquid at all). - 1 point: The glass contains liquid, but the focus on the liquid or the glass is unclear, or there are distracting elements in the scene. - 2 points: The glass is clearly depicted with some liquid in it, with no distractions, offering a simple, clear, and realistic depiction." <p>},</p> <p>"Reality Grading":{</p> <ul style="list-style-type: none"> - 0 points: Liquids are mixed or incorrectly positioned (e.g., water and oil blended or misplaced). - 1 point: Water and oil are present but with partial inaccuracies in separation or positioning (e.g., water floating on oil, blurred boundaries). - 2 points: Liquids are correctly positioned with visible separation (oil atop water), but minor deviations from realism exist (e.g., slight issues with clarity or texture). - 3 points: Fully realistic depiction with correct positioning (oil floating on water) and clear separation. <p>}]</p>

Fig. S8: Representative samples of grading criterion.

S15. Details of Two-Stage Training

In this section, we present a detailed overview of our two-stage training framework, which integrates SFT and masked online fine-tuning to enhance flow matching models.

Supervised Fine-tuning (SFT). Flow matching models [46] are continuous-time generative models that define a time-dependent velocity field $v(x_t, t)$ to transport samples from a noise distribution p_1 to data distribution p_0 over a time interval $t \in [0, 1]$. The transformation is governed by the ordinary differential equation (ODE):

$$\frac{dx_t}{dt} = v(x_t, t), \quad (\text{S24})$$

with the initial condition $x_1 \sim p_1$. The forward process is constructed as:

$$x_t = \alpha_t x_0 + \beta_t \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad (\text{S25})$$

where $\alpha_0 = 1$, $\beta_0 = 0$, $\alpha_1 = 0$, and $\beta_1 = 1$, ensuring the consistency of the marginal distributions with the initial and terminal conditions. The velocity field $v(x_t, t)$ is represented as the sum of two conditional expectations:

$$v(x, t) = \dot{\alpha}_t \mathbb{E}[x_* | x_t = x] + \dot{\beta}_t \mathbb{E}[\epsilon | x_t = x], \quad (\text{S26})$$

which can be approximated by the model $v_\theta(x, t)$ by minimizing the following training objective:

$$\mathcal{L}_{\text{SFT}}(\theta) := \mathbb{E}_{x_*, \epsilon, t} [\|v_\theta(x_t, t) - \dot{\alpha}_t x_* - \dot{\beta}_t \epsilon\|^2] \quad (\text{S27})$$

Direct Preference Optimization (DPO). RLHF aims to optimize a conditional distribution $p_\theta(x_0|c)$ such that the expected reward $r(c, x_0)$ is maximized, while simultaneously regularizing the KL-divergence from a reference distribution p_{ref} . This objective is formulated as:

$$\max_{p_\theta} \mathbb{E}_{c, x_0 \sim p_\theta(x_0|c)} [r(c, x_0)] - \beta \mathcal{D}_{\text{KL}} [p_\theta(x_0|c) \| p_{\text{ref}}(x_0|c)] \quad (\text{S28})$$

where the hyper-parameter β controls regularization. According to [56], the unique global optimal solution p_θ^* to this optimization problem is given by:

$$p_\theta^*(x_0|c) = p_{\text{ref}}(x_0|c) \exp\left(\frac{r(c, x_0)}{\beta}\right) / Z(c) \quad (\text{S29})$$

where $Z(c) = \sum_{x_0} p_{\text{ref}}(x_0|c) \exp\left(\frac{r(c, x_0)}{\beta}\right)$ is partition function. Then the reward function can be expressed as:

$$r(c, x_0) = \beta \log \frac{p_\theta^*(x_0|c)}{p_{\text{ref}}(x_0|c)} + \beta \log Z(c) \quad (\text{S30})$$

To model human preferences, the Bradley-Terry (BT) model is employed, which represents the probability of one outcome being preferred over another as:

$$p_{\text{BT}}(x_0^w \succ x_0^l | c) = \sigma(r(c, x_0^w) - r(c, x_0^l)) \quad (\text{S31})$$

where σ is the sigmoid function, x_0^w is the preferred outcome, and x_0^l is the less preferred one.. $r(c, x_0)$ can be parameterized by a neural network ϕ and estimated via maximum likelihood training for binary classification:

$$L_{\text{BT}}(\phi) = \mathbb{E}_{c, x_0^w, x_0^l} [\log \sigma(r_\phi(c, x_0^l) - r_\phi(c, x_0^w))] \quad (\text{S32})$$

By leveraging the relationship between the reward function and the optimal policy p_θ^* , the DPO objective is derived as:

$$\mathcal{L}_{\text{DPO}}(\theta) = -\mathbb{E}_{c, x_0^w, x_0^l} \left[\log \sigma \left(\beta \log \frac{p_\theta(x_0^w | c)}{p_{\text{ref}}(x_0^w | c)} - \beta \log \frac{p_\theta(x_0^l | c)}{p_{\text{ref}}(x_0^l | c)} \right) \right] \quad (\text{S33})$$

Choice of σ_t . We determine the value of σ_t by adhering to the methodology presented in [35]. Initially, we define the hyperparameters S_{churn} , S_{min} , S_{max} , and S_{noise} . Subsequently, we define γ_t as follows:

$$\gamma_t = \begin{cases} \min(S_{\text{churn}} \cdot \Delta t, \sqrt{2} - 1) & \text{if } t \in [S_{\text{min}}, S_{\text{max}}] \\ 0 & \text{otherwise,} \end{cases} \quad (\text{S34})$$

where Δt represents the timestep difference between consecutive sampling steps. Following this, we define σ_t by

$$\sigma_t = S_{\text{noise}} \cdot \sqrt{\gamma_t^2 + 2\gamma_t} \cdot (1 - t). \quad (\text{S35})$$

Pre-Training Subject Extraction. We integrate GroundingDINO [47] to facilitate the extraction of masks from images. To streamline the process, we initially employ LLM to identify and extract the relevant subjects from the training prompt set prior to the training phase. The extracted subjects are subsequently provided to GroundingDINO [47] during training to generate corresponding masks. These masks are then utilized to apply gradient masking.

Gradient Masking. The mask generated by GroundingDINO [47] is derived from the resolution of the RGB image. However, gradients are computed within the model's latent space, as detailed in LDM [58]. The connection between the RGB image and the latent space is facilitated by a pretrained Variational Autoencoder (VAE) [36], which inherently exhibits localist properties. Specifically, let the latent representation have dimensions (H_l, W_l, C_l) and the corresponding decoded image have dimensions (H, W, C) . If the mask extracted from the image is defined by the bounding box coordinates (x_1, y_1, x_2, y_2) , then the corresponding mask in the latent space is computed as:

$$\left(\frac{x_1}{H} \cdot H_l, \frac{y_1}{W} \cdot W_l, \frac{x_2}{H} \cdot H_l, \frac{y_2}{W} \cdot W_l \right) \quad (\text{S36})$$

This latent-space mask is subsequently applied to the gradients of the model to modulate the training process.

Padding Technique. Certain tasks require the careful consideration of positional relationships rather than solely the object's internal state. For example, in the *gravity* task, the object's position relative to the ground is critically important, making the use of the object mask alone insufficient for accurate analysis. To address this limitation, we extend the height and width dimensions of the mask by an additional 10%. This strategic padding ensures that the surrounding positional context is adequately captured, improving task performance and contextual understanding.

S16. Two-Stage Training Settings

In this section, we detail the hyper-parameter configurations employed in our two-stage training framework for the T2I model, which is presented in Table S9.

Table S9: Hyper-parameter settings for T2I model fine-tuning.

Hyper-parameters	SFT	OFT
batch size	32	8
learning rate	2×10^{-5}	3×10^{-4}
training steps	900	140
optimizer	AdamW [48]	AdamW [48]
gradient accumulation	8	2
LoRA rank	16	16
S_{churn}	/	0.1
S_{\min}, S_{\max}	/	$0, \infty$
S_{noise}	/	1.0
β	/	10

S17. Limitation and Future Direction

Limited Cross-Task Generalizability. Our method lacks generalizability across tasks due to the difficulty of transferring knowledge between distinct scientific domains. For instance, expertise in dynamics does not inherently facilitate comprehension of thermodynamics, even for humans. Consequently, performance improvements are confined to the 16 predefined tasks. Thus, integrating domain-specific data for extra training is essential to mitigate this limitation.

Limited Generalizability across ST. Since ST are heavily influenced by the subject, they present significant challenges for current T2I models, primarily due to the models' lack of "world knowledge". When encountering novel or unseen subjects, T2I models often fail to generate appropriate visual content, as they lack prior exposure to the subject's defining characteristics. In such cases, exploring how to incorporate the world knowledge of LLMs or leveraging additional data for training could be promising directions for future improvements.

S18. Additional Results of Online Fine-tuning

To further assess the effectiveness of the proposed algorithm, we conducted additional experiments utilizing different reward models. Specifically, we employed the LAION aesthetic predictor [4] and ImageReward [71] as the reward functions for our comprehensive evaluations. It is important to note that, in these experiments, we did not implement the masking strategy described in the main text.

Training Setting. All configurations align with those presented in Table S9, except for the specific settings detailed below. We fine-tuned the FLUX.1[schnell] [1] using 4 inference steps. During training with both reward models, each step involved sampling 128 images, utilizing a learning rate of 6×10^{-5} , and employing a gradient accumulation step size of 8. For training with the LAION aesthetic predictor [4], we conducted 164 training steps. In contrast, when training with ImageReward [71], we performed 550 training steps. Adhering to the configuration outlined in DDPO [7], the training prompt set comprised 45 distinct animal categories.

Evaluation Setting. The test prompt set consisted of an additional 10 animal categories not present in the training set. For each prompt, we generated 100 images and calculated the average reward assigned by the respective reward model, which served as our performance metric. The final experimental results, showcasing the average rewards achieved on the test set, are presented in Table S10. Additionally, qualitative examples are provided in Figure S9 for further analysis.

Table S10: Comparison of Average Rewards. The online fine-tuning approach consistently outperforms the baseline, demonstrating the effectiveness of the proposed algorithm.

Method	LAION [4]	ImageReward [71]
FLUX.1[schnell]	5.855	0.949
+OFT	6.074	1.023

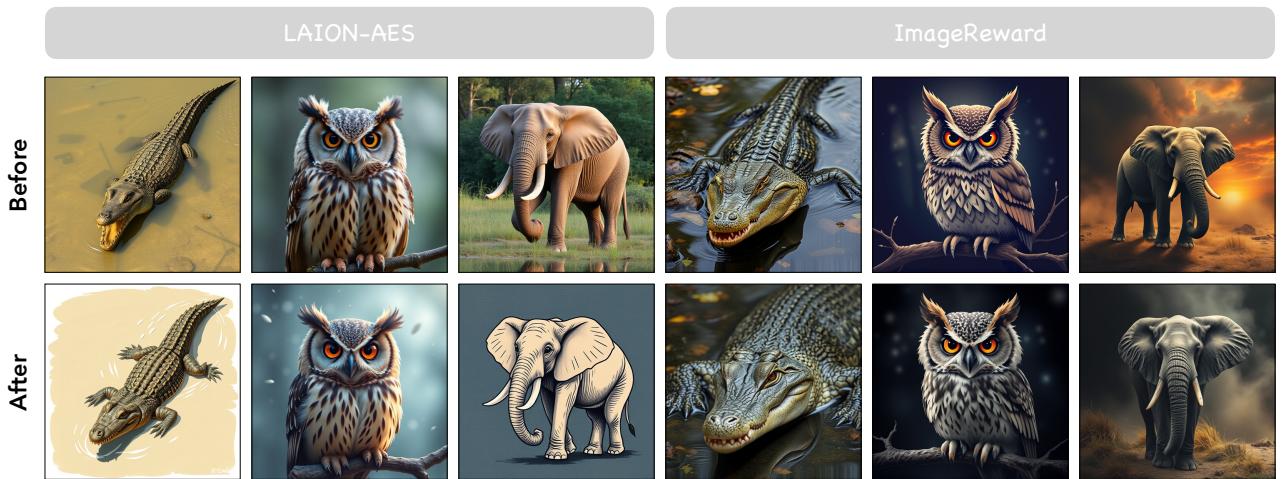
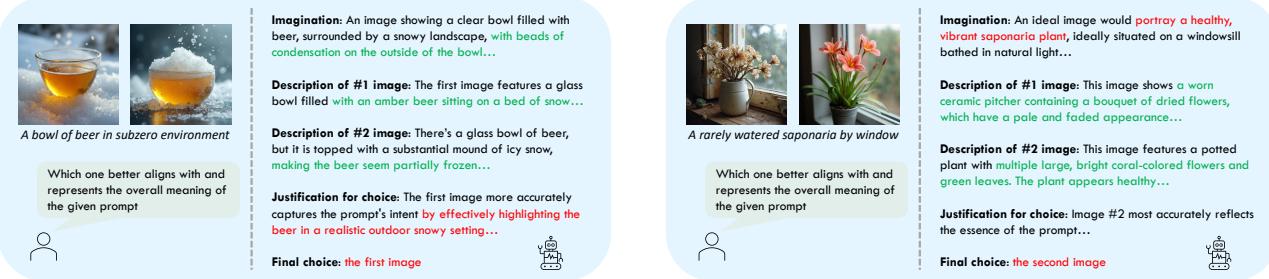
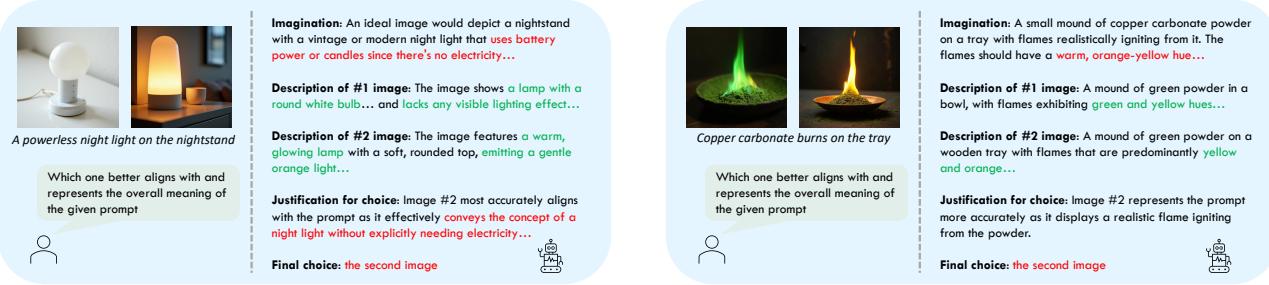
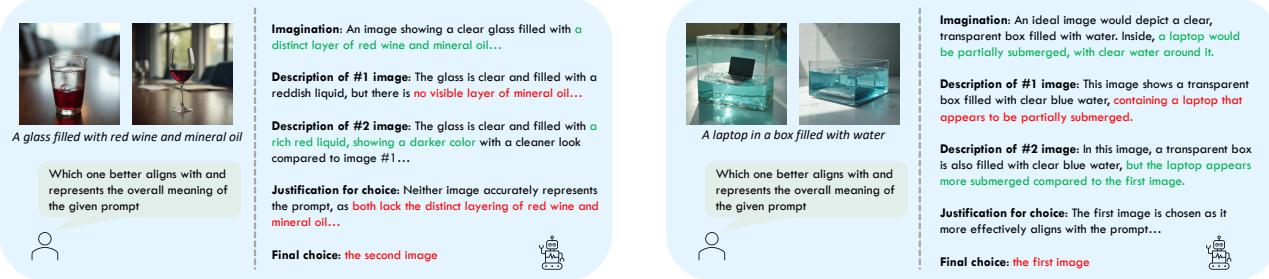


Fig. S9: Case study. The upper images are generated using the base FLUX.1[schnell] [1], whereas the lower images are produced with online fine-tuning method. Each image pair utilizes an identical random seed to ensure consistency in comparison. The prompts used here are "crocodile", "owl", and "elephant."



(a) Reasoning Failure. GPT-4o-mini [2] inaccurately infers the target image by misinterpreting the input prompt and neglecting the underlying scientific principles embedded within it. Instead of employing a systematic reasoning process, it relies predominantly on intuitive imagination.



(b) Visual Limitation. GPT-4o-mini [2] inaccurately describes the image, thereby impeding the reasoning process. Specifically, for tasks involving spatial relationships, it fails to make correct judgments, resulting in erroneous interpretations of positional dynamics within the visual content.

Fig. S10: Qualitative Failure Cases of GPT. In both cases, the CoT [69] reasoning approach from Figure S4 is applied, but errors in either interpretation or visual comprehension impact the final decision. Green text indicates correct inference, while red text marks errors.

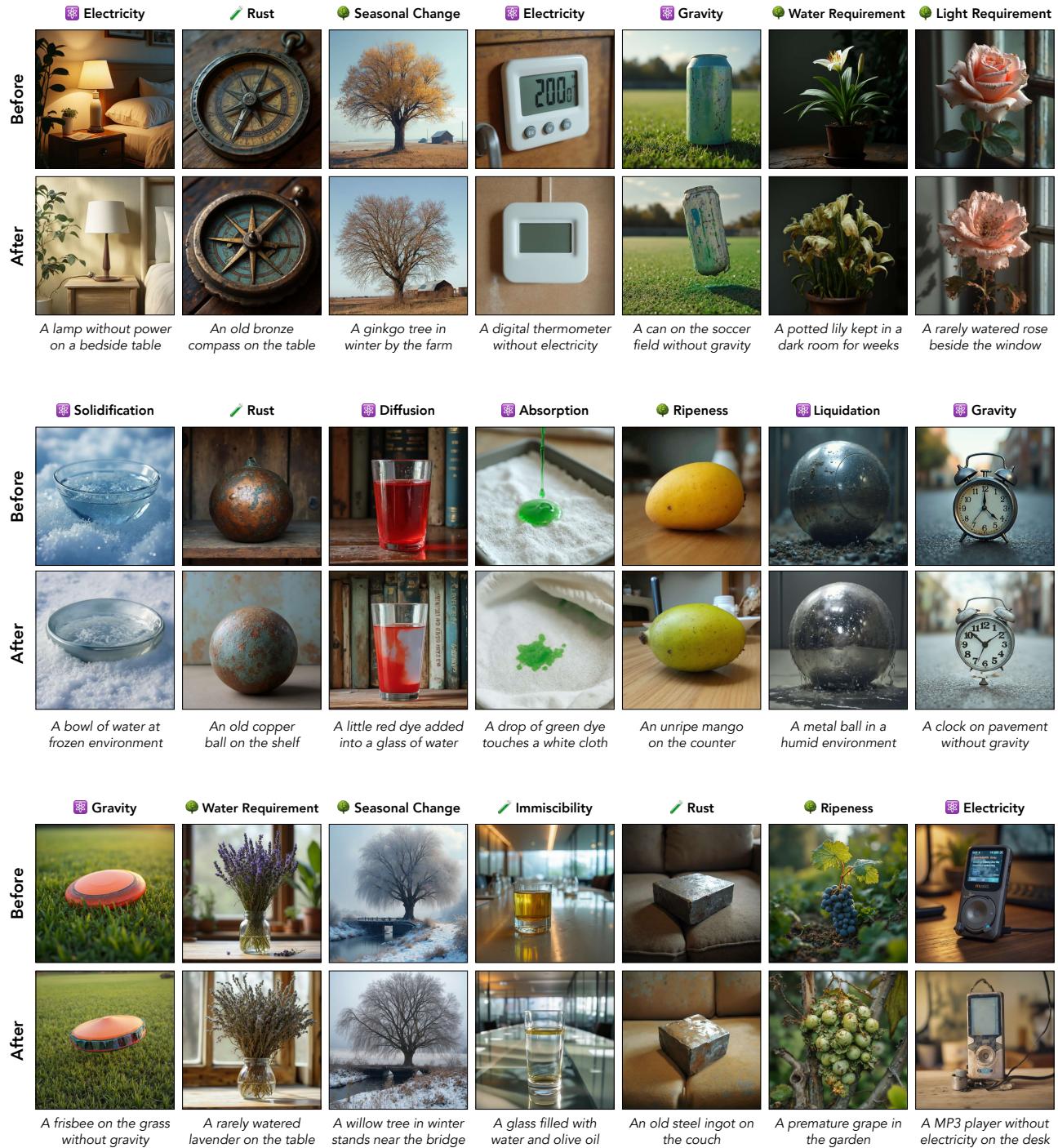


Fig. S11: Additional Generated Samples. Each pair of images is produced using the same random seed to ensure consistency.

<p> Buoyancy</p> <p>IP</p> <p>A transparent water-filled box holds an {apple} in a simple and empty background, depicted realistically.</p>		<p>A transparent box filled with water holds an {apple} floating on the surface, realistic. The background is completely empty.</p>	<p>EP</p>
<p> Diffusion</p> <p>IP</p> <p>A little bit {orange dye} added to a glass of water, simple and realistic.</p>		<p>A bit of {orange dye} spreads and diffuses into the glass of water, in a simple and realistic way.</p>	<p>EP</p>
<p> Gravity</p> <p>IP</p> <p>A {pillow} in a simple space without gravity, simple and realistic.</p>		<p>A realistic scene of a {pillow} floating in the air within a simple space. The background is completely empty.</p>	<p>EP</p>
<p> Melting</p> <p>IP</p> <p>A {butter stick} on a heated pan, simple and realistic.</p>		<p>A heavily melting {butter stick} on a heated pan, losing its original shape as liquefied portions spread into a glossy area.</p>	<p>EP</p>
<p> Flame Reaction</p> <p>IP</p> <p>A bit of {copper} powder ignites on a surface, in a simple and realistic way.</p>		<p>A bit of {copper} powder ignites into a {green} flame on a surface, the scene simple and realistic.</p>	<p>EP</p>
<p> Immiscibility</p> <p>IP</p> <p>A clear glass filled with {milk and oil}, simple and realistic.</p>		<p>A clean and distinct separation of layers is visible in the glass, with {milk} at the bottom and {oil} floating on top.</p> <p>A clear glass filled with {oil/milk}, straightforward and realistic.</p>	<p>EP</p>

Fig. S12: Several examples from SCIENCE-T2I. 'EP' denotes explicit prompts (yellow blocks), 'SP' denotes superficial prompts (blue blocks), and 'IP' denotes implicit prompts (grey blocks).

 Light Requirement

IP

A potted **{lily}** kept in a dark room for weeks, simple and realistic.



A potted **{lily}** sits in a dimly lit room, its petals wilted and curling with brown edges, while the stems sag.

EP

A potted **{lily}** stands tall in a dimly lit room, its vivid petals brimming with life and vitality. Strong, upright stems hold fresh petals.

SP

 Water Requirement

IP

A rarely watered **{rose}**, presented in a simple and realistic way.



A **{rose}** with wilted petals, curled and browned at the edges, droops from its stems, giving it a dry, decaying appearance.

EP

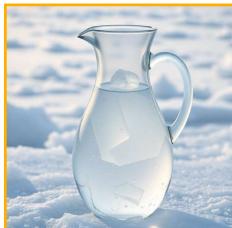
A blooming **{rose}** with vibrant petals stands tall on strong, upright stems, radiating health.

SP

 Solidification

IP

A **{carafe}** of **{water}** in a glacier, simple and realistic.



A **{carafe}** of frozen **{water}** in a glacier, simple and realistic.

EP

A **{carafe}** of fully liquid **{water}** in a glacier, simple and realistic.

SP

 Ripeness

IP

A unripe **{tomato}**, simple and realistic.



A green **{tomato}** with firm, smooth, and shiny skin is simple, clear, and realistic.

EP

A red **{tomato}**, making it simple and realistic.

SP

 Absorption

IP

A drop of **{blue dye}** touches a napkin, simple and realistic.



The **{blue dye}** spreads, creating a diffused blue stain on the napkin, simple and realistic.

EP

The **{blue dye}** drop stays as a tiny, focused spot on the napkin, creating a scene that's simple and realistic.

SP

 Chemistry --- Rust

IP

A **{iron hammer}** that has been exposed to oxygen for decades, simple and realistic.



The **{iron hammer}** has a look with a **{red rust}**, revealing its age and corrosion.

EP

A realistic **{iron hammer}** stands out against a completely blank background, simple and realistic.

SP

Fig. S13: Several examples from SCIENCE-T2I. 'EP' denotes explicit prompts (yellow blocks), 'SP' denotes superficial prompts (blue blocks), and 'IP' denotes implicit prompts (grey blocks).

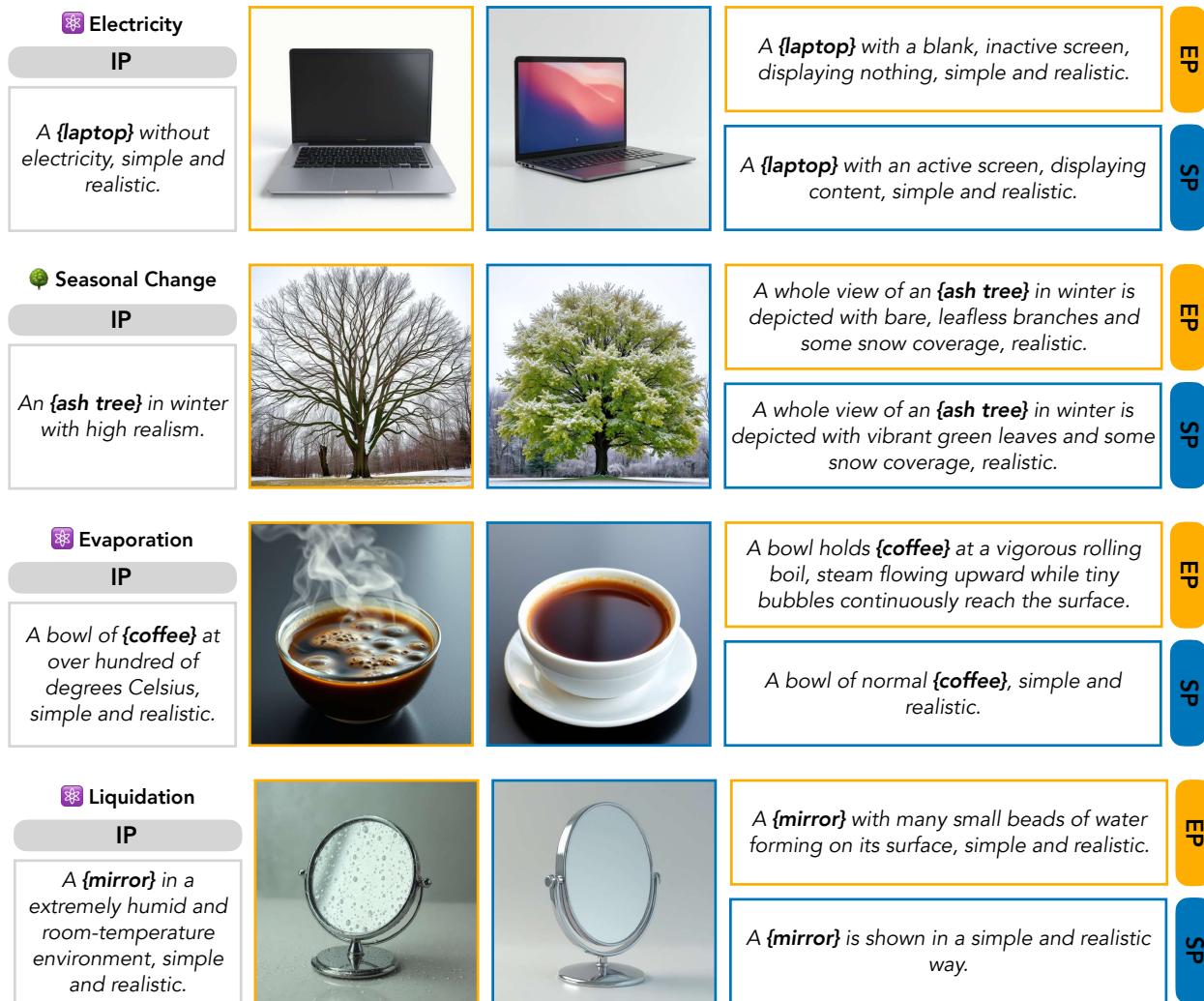


Fig. S14: Several examples from SCIENCE-T2I. 'EP' denotes explicit prompts (yellow blocks), 'SP' denotes superficial prompts (blue blocks), and 'IP' denotes implicit prompts (grey blocks).