# `ImAgent`: A Unified Multimodal Agent Framework for Test-Time Scalable Image Generation

Kaishen Wang[*], Ruibo Chen[*], Tong Zheng, Heng Huang[†]

University of Maryland, College Park

## Abstract

*Recent text-to-image (T2I) models have made remarkable progress in generating visually realistic and semantically coherent images. However, they still suffer from randomness and inconsistency with the given prompts, particularly when textual descriptions are vague or underspecified. Existing approaches, such as prompt rewriting, best-of-N sampling, and self-refinement, can mitigate these issues but usually require additional modules and operate independently, hindering test-time scaling efficiency and increasing computational overhead. In this paper, we introduce `ImAgent`, a training-free unified multimodal agent that integrates reasoning, generation, and self-evaluation within a single framework for efficient test-time scaling. Guided by a policy controller, multiple generation actions dynamically interact and self-organize to enhance image fidelity and semantic alignment without relying on external models. Extensive experiments on image generation and editing tasks demonstrate that `ImAgent` consistently improves over the backbone and even surpasses other strong baselines where the backbone model fails, highlighting the potential of unified multimodal agents for adaptive and efficient image generation under test-time scaling.*

## 1. Introduction

Text-to-Image (T2I) models [33–36, 44, 47] have made remarkable progress in generating visually realistic and semantically coherent images from natural language descriptions. Despite these advancements, the quality of generated results often exhibits randomness and inconsistency with the given prompts [18, 62]. This limitation primarily stems from the strong dependency of existing T2I models on the clarity and specificity of textual inputs. When the prompt is vague or underspecified, the models tend to overlook key semantic elements, leading to images that deviate from the intended meaning and fail to fully capture the user's intent.

To mitigate this problem, prior research has proposed various strategies, such as prompt rewriting [1, 25, 29], best-of-N sampling [21, 42], classifier-free guidance [11, 37, 39], and self-revision or iterative refinement [14, 19, 41]. These approaches aim to either clarify the textual prompt or reduce the model's inherent randomness rather than finetuning the models, reflecting the underlying principle of test-time scaling—improving generation quality through additional inference-time computation.

While effective, these methods typically require additional components, such as a language model for prompt enhancement, a generative model for image synthesis, and a vision-language model for evaluation, thereby increasing memory consumption and computational cost. Moreover, they are often applied independently, requiring human intervention to determine the optimal approach for a given case, which substantially limits the efficiency of test-time scaling and leads to unnecessary computational overhead.

In this paper, we aim to construct a universal agent for image generation that can adaptively select the optimal action for a given case, allocate computational resources accordingly, and execute the selected action within the agent itself without relying on any external models. This design enables more efficient test-time scaling. Thanks to recent advances in unified multimodal models [6, 34, 54, 54, 60], which integrate text generation, image generation, and visual understanding within a single framework, this assumption becomes attainable. Building upon this foundation, we introduce `ImAgent`, a training-free unified multimodal agent designed to perform efficient test-time scaling for image generation.

Specifically, `ImAgent` is built around a policy controller that serves as its "brain", determining whether an image requires refinement based on the observation history and deciding which action should be invoked in the next step. Under its coordination, multiple predefined generation actions, such as *Prompt Enhancement with CoT*, *Image Detail Refinement*, and *Best-of-N Sampling*, operate collaboratively within a unified framework. Notably, all these actions are executed internally through the agent itself. Through this dynamic interaction, `ImAgent` trans-

---

[*]Equal Contribution
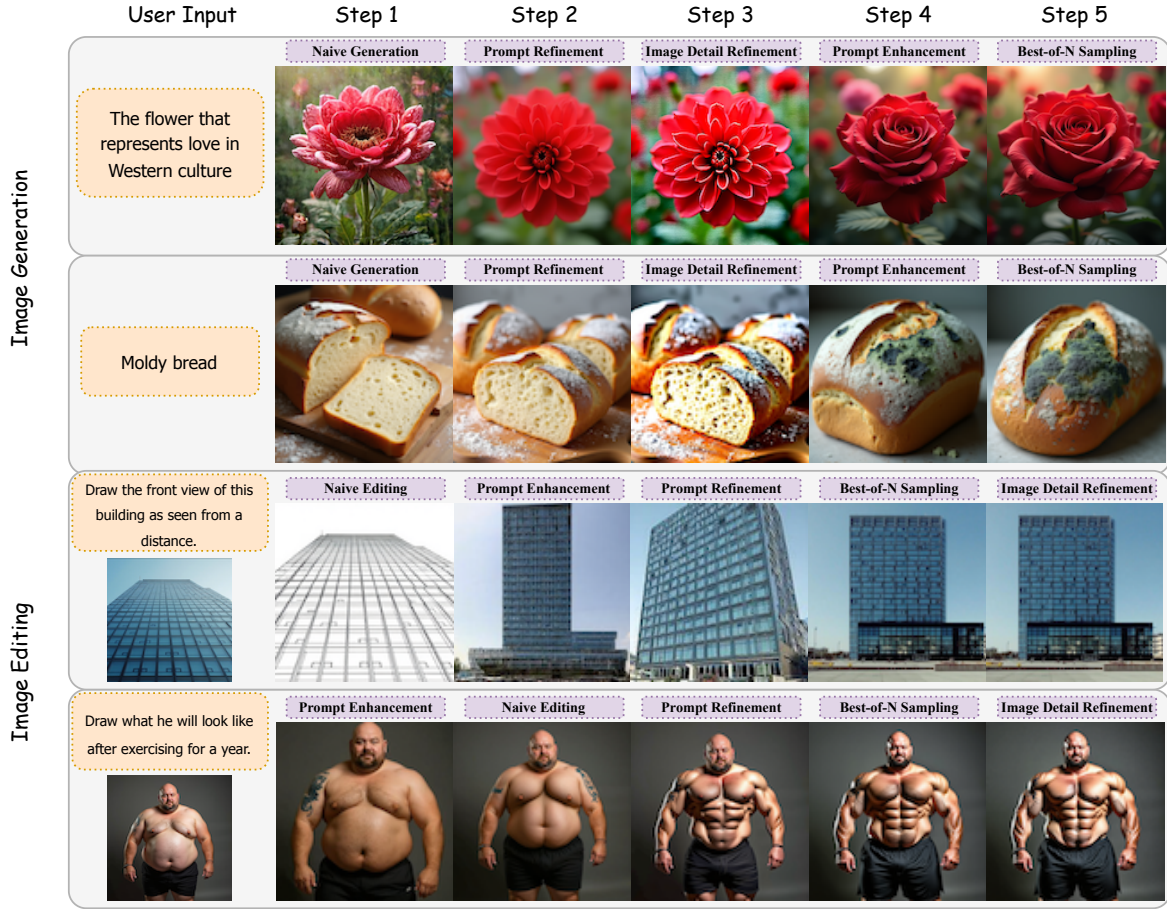[†]Corresponding Author
[‡]Preprint

Figure 1. Qualitative examples of `ImAgent` built upon Bagel for image generation and image editing.

forms what was previously a manually assembled pipeline into a self-organizing agent capable of reasoning, generation, and self-improvement, thereby enhancing the efficiency of test-time scaling.

We conduct experiments on two popular tasks, image generation and image editing, to evaluate the effectiveness of the proposed `ImAgent`. Extensive experimental results show that `ImAgent` achieves significant improvements over the backbone model, demonstrating its strong capability and the effectiveness of test-time scaling in enhancing image generation performance. The contributions of this paper are summarized as follows:

- We propose a unified multimodal agent framework for image generation, `ImAgent`, which integrates multiple generation actions and dynamically coordinates them through a policy controller.
- Extensive experiments on image generation and editing tasks demonstrate that `ImAgent` achieves substantial improvements over the backbone model, even outperform-

ing baselines that the backbone model fails to surpass, validating its effectiveness and in test-time scaling.

## 2. Related Work

### 2.1. Text-to-Image Generation

Text-to-Image (T2I) models have demonstrated remarkable ability in synthesizing visually realistic and semantically coherent images from natural language descriptions. The evolution of T2I models has followed several major paradigms, including diffusion-based models [7, 9, 33, 35, 36, 59], autoregressive models [44, 47, 52, 61, 69, 73], and more recently, unified multimodal models [6, 26, 49, 58, 60, 72, 74]. These approaches have significantly advanced image fidelity, diversity, and text-image alignment.

### 2.2. Unified Multimodal Models

Unified multimodal models refer to models that jointly possess multimodal understanding and generation capabilities

within a single framework. Unlike traditional pipelines that rely on separate models for text or image understanding, image generation, and image editing, these models integrate vision and language processing into a shared architecture, enabling seamless information exchange across modalities [49, 54, 54, 58, 74]. This convergence of understanding and generation not only improves efficiency and coherence but also paves the way for building more general-purpose, human-like agents capable of reasoning, creating, and self-correcting.

### 2.3. Generation Optimization Strategies

Despite the remarkable progress of T2I models, the quality of generated images still depends heavily on the clarity of prompts and the randomness inherent in the generation process. To alleviate these issues, various optimization strategies have been proposed. *Best-of-N sampling* [21, 48, 50] generates multiple candidates and selects the best one, reducing the influence of stochastic variations in model sampling. *Prompt rewriting* [48, 53, 64] reformulates or enriches the input text to provide more explicit guidance for the generator, improving semantic alignment and visual coherence. *classifier-free guidance* [11, 37, 39] controls the conditioning strength to balance fidelity and diversity, and *iterative self-refinement* [14, 19, 28, 41, 63] enables models to assess and revise their outputs through multiple rounds of generation. Although these techniques significantly enhance performance, they typically operate as independent modules with limited interaction between reasoning, generation, and evaluation—highlighting the need for unified multimodal frameworks that integrate these capabilities into a cohesive pipeline.

## 3. Method

We construct `ImAgent` upon unified multimodal models, which inherently possess both understanding and generation capabilities. Motivated by the observation that multimodal understanding is generally more reliable and easier to achieve than multimodal generation [67, 70], `ImAgent` leverages the model's strong understanding ability to adaptively select the optimal action for a given case that enhances the overall generation quality. This design enables `ImAgent` to fully exploit the potential of a single unified multimodal model without any additional training or external modules, thereby achieving efficient test-time scaling. In this section, we first introduce the overall agent framework of `ImAgent`, followed by a detailed description of its action space.

### 3.1. Agent Framework Overview

As shown in Figure 2, `ImAgent` dynamically selects the next action based on the current observation and executes the selected action until either the STOP action is chosen or

---

**Algorithm 1:** `ImAgent` for Image Generation.

**Input:** Initial user prompt $P_0$, action space $\mathcal{A} = \{a_1, \ldots, a_n\}$, maximum step number $T_{\max}$
**Output:** Final image $I^*$
1 Initialize prompt $p \leftarrow p_0$; image $I \leftarrow \varnothing$; observation history $\mathcal{O} \leftarrow []$
2 **for** $t = 1$ **to** $T_{\max}$ **do**
   `// 1. Reasoning and Action`
   `   Selection`
3  |  Use the policy controller $\pi_\theta$ to determine the next action:
4  |  $a_t \leftarrow \pi_\theta(a \mid s_t)$
5  |  **if** $a_t = STOP$ **then**
6  |  |  break `// The model decides the`
      |  |  `   generation is satisfactory`
   |  `// 2. Action Invocation`
7  |  Execute the selected action $a_t \in \mathcal{A}$:
8  |  $(P_{t+1}, I_{t+1}, o_t) \leftarrow f_{a_t}(P_t, I_t, \mathcal{O}_{t-1})$
9  |  where $o_t$ is the new observation (e.g., quality evaluation)
   |  `// 3. State Update`
10 |  Append $o_t$ to the observation history: $\mathcal{O}_t \leftarrow \mathcal{O}_{t-1} \cup \{o_t\}$
11 **return** Final image $I^* \leftarrow I$

---

the maximum step limit is reached. Formally, we define the state at the current $t$-th time step as:

$$s_t = \{P_0, I_0, P_t, I_t, \mathcal{O}_{t-1}\}, \qquad (1)$$

where $P_0$ and $I_0$ denote the initial user prompt and input image, $P_t$ and $I_t$ represent the current prompt and generated image at $t$-th step, and $\mathcal{O}_{t-1} = \{o_1, \ldots, o_{t-1}\}$ denotes the history of past observations, i.e., action history.

Then, `ImAgent` dynamically selects an action $a_t$ from the defined action space $\mathcal{A} = \{a_1, \ldots, a_n\}$ (detailed in Section 3.2) using a policy controller $\pi_\theta$ conditioned on the current state at $t$-th step:

$$a_t \sim \pi_\theta(a \mid s_t). \qquad (2)$$

Upon selecting $a_t$, the corresponding action function $f_{a_t}$ is invoked to produce an updated prompt $P_{t+1}$, an updated image $I_{t+1}$, and a new observation $o_t$:

$$(P_{t+1}, I_{t+1}, o_t) = f_{a_t}(P_t, I_t, \mathcal{O}_{t-1}). \qquad (3)$$

Then, the observation is updated as follows:

$$\mathcal{O}_t \leftarrow \mathcal{O}_{t-1} \cup o_t. \qquad (4)$$

Unlike conventional agent-based frameworks where observations directly influence state updates, in `ImAgent` the
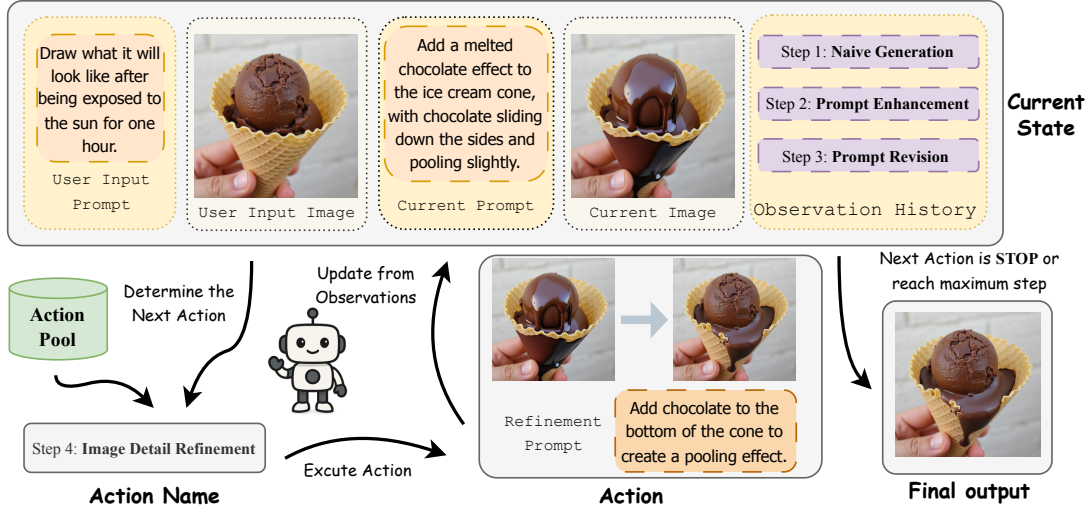
Figure 2. The overall architecture of `ImAgent`.

observation history is primarily used for action selection, while prompt and image updates are determined by the selected generation actions.

This iterative process continues until the agent selects the STOP action, which indicates that the generated result is satisfactory, or until the maximum number of steps $T_{\max}$ (set to 5 by default) is reached. Notably, both the policy controller and the action execution are implemented within a single unified multimodal model. The detailed algorithm for image generation is presented in Algorithm 1, while the algorithm for image editing follows a similar procedure and is provided in the Appendix.

### 3.2. Action Space

In this section, we introduce the action space $\mathcal{A}$ of `ImAgent`. The detailed definitions and usage of each action are provided in the Appendix.

**Naive Generation/Editing.** This action performs a one-shot image generation or editing operation directly based on the current prompt. It is typically used when the input description is simple, unambiguous, and requires no iterative refinement.

**Prompt Enhancement with CoT.** This action refines the input prompt by enriching vague or underspecified user queries with additional contextual and descriptive details. Since text-to-image (T2I) models are highly sensitive to prompt wording [1, 25, 29], more elaborate and specific prompts typically lead to higher-quality visual outputs. However, most T2I backbones are trained on simple captions or surface-level image descriptions, which limits their ability to reason over complex or compositional instructions. To address this limitation, we leverage the language

reasoning capability of the model's understanding module through Chain-of-Thought (CoT) prompting. By explicitly performing intermediate reasoning and elaboration, the agent transfers structured linguistic insights from the understanding domain to the generation process, thus producing more semantically aligned and visually coherent outputs.

**Prompt Revision Based on the Generated/Edited Image.** This action is triggered when the generated or edited image indicates that the current prompt is suboptimal. In this case, the unified multimodal model is prompted to analyze the discrepancy between the visual output and the intended semantics, and to self-revise the prompt accordingly. This self-correction mechanism allows the agent to iteratively refine textual descriptions based on visual feedback, thereby enhancing both semantic alignment and generation quality over successive iterations.

**Image Detail Refinement.** This action aims to correct minor imperfections in the generated or edited image when the input instruction is already satisfactory and the remaining issues stem from the generation module itself. It refines local visual details—such as textures, lighting, or small artifacts—without modifying the current prompt. By enhancing fine-grained fidelity while maintaining semantic consistency, this action improves the overall perceptual quality and realism of the generated results.

**Best-of-N Sampling.** This action mitigates the inherent stochasticity of text-to-image (T2I) generation. Unlike language models that typically produce stable, low-entropy outputs through supervised fine-tuning (SFT) and reinforcement learning (RL) optimization, T2I models exhibit substantial variance across different samples [21, 48, 50]. To

| Types | Model | Comm. | Comp. | Logical | Numerical | Causal | Overall |
|---|---|---|---|---|---|---|---|
| Gen | SD3-medium | 0.54 | 0.64 | 0.55 | 0.50 | 0.18 | 0.53 |
| | Sana-1.5 | 0.49 | 0.67 | 0.49 | 0.48 | 0.21 | 0.49 |
| | Lumina-T2I | 0.38 | 0.49 | 0.38 | 0.45 | 0.18 | 0.39 |
| | Omnigen | 0.43 | 0.60 | 0.51 | 0.47 | 0.34 | 0.48 |
| | LLM4GEN$_{SD1.5}$ | 0.55 | 0.48 | 0.55 | 0.39 | 0.45 | 0.51 |
| | ELLA$_{SD1.5}$ | 0.40 | 0.44 | 0.40 | 0.32 | 0.29 | 0.39 |
| | LlamaGen | 0.38 | 0.39 | 0.38 | 0.35 | 0.12 | 0.36 |
| | DALL-E-3 | 0.78 | 0.76 | 0.69 | 0.69 | 0.64 | 0.73 |
| | gpt-image-1 | 0.83 | 0.87 | 0.81 | 0.88 | 0.71 | 0.83 |
| Unified | EMU3 | 0.44 | 0.59 | 0.55 | 0.61 | 0.41 | 0.52 |
| | Show-o | 0.42 | 0.59 | 0.42 | 0.57 | 0.30 | 0.46 |
| | Lumina-Image 2.0 | 0.49 | 0.65 | 0.56 | 0.43 | 0.40 | 0.52 |
| Bagel | Vanilla | 0.47 | 0.65 | 0.53 | 0.68 | 0.40 | 0.54 |
| | ImAgent | **0.58** | **0.68** | **0.62** | **0.71** | **0.53** | **0.62** |
| Janus-Pro-7B | Vanilla | 0.45 | **0.60** | 0.46 | 0.46 | 0.36 | 0.47 |
| | ImAgent | **0.47** | 0.58 | **0.52** | **0.54** | **0.44** | **0.51** |

Table 1. Experimental results on **R2I-Bench** [4]. *Comm.* and *Comp.* denote the Commonsense and Compositional categories, respectively. Due to the inaccessibility of parts of the Mathematical and Concept-Mixing categories, these two categories are omitted. Vanilla represents the backbone model, while ImAgent denotes our agent built upon this model. The best performance between vanilla and ImAgent is highlighted in **bold**.

reduce this randomness, the agent generates $N$ candidate images and employs the unified multimodal model to evaluate their visual and semantic alignment. The image with the highest alignment score is selected as the final output.

**Termination (STOP).** This action signals the end of the iterative reasoning process when the agent determines that the current image and prompt have achieved satisfactory alignment. Instead of relying on a fixed number of steps, the agent autonomously decides to terminate based on its internal evaluation of visual quality.

## 4. Experiments

### 4.1. Experimental Setting

**Models.** We build our ImAgent on two of the most recent and powerful unified models, Bagel [6] and Janus-Pro-7B [5]. Both models possess versatile multimodal capabilities, including image understanding, image generation, which collectively form the foundation of our framework. For image generation, we leverage both Bagel and Janus-Pro-7B. However, since Janus-Pro-7B does not support image editing, we employ Bagel exclusively for the image editing tasks.

**Benchmarks.** We evaluate ImAgent on multiple benchmarks, including 3 benchmarks on image generation and 4 benchmarks for image editing tasks. For image generation, we use R2I-Bench [4], which is designed to rigorously assess reasoning-driven T2I generation; WISE [30], which covers six categories of image generation scenarios; and T2I-ReasonBench [43], which assesses the models' reasoning ability in generative tasks. For image editing, we
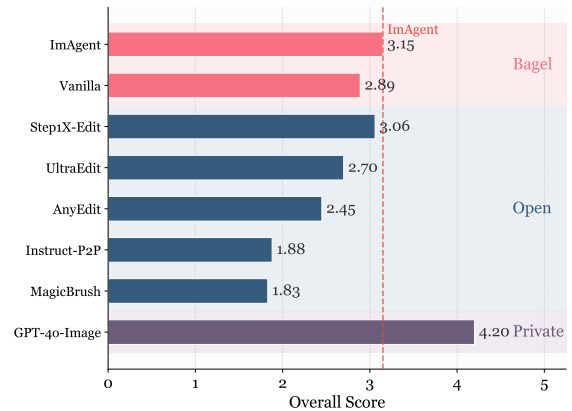


Figure 3. Experimental results on **ImgEdit-Bench** [65].

adopt GEdit-Bench [23], which contains both Chinese and English instruction-based editing tasks; RISEBench [71], which focuses on reasoning-informed visual editing across diverse reasoning types; KRIS-Bench [57], which evaluates reasoning capabilities over factual, conceptual, and procedural knowledge; and ImgEdit-Bench [65], which is used to evaluate image editing performance in terms of instruction adherence, editing quality, and detail preservation.

**Baselines.** For image generation, we employ SDv1.5 [36], SD3-Medium [36], SDXL [33], SD3.5-Medium [7], SD3.5-Large [7], PixArt-Alpha [3], FLUX.1-Dev [15], Sana-1.5 [59], Lumina-T2I [34], LLM4GEN$_{SD1.5}$ [22], ELLA$_{SD1.5}$ [12], LlamaGen [44], DALL-E-3 [27], gpt-image-1 [13], Omnigen [58], FLUX.1-schnell [16], andPlayground-v2.5 [20], as generation-based models. We also include Janus [54], VILA-U [56],

| Types | Model | Cultural | Time | Space | Biology | Physics | Chemistry | Overall |
|---|---|---|---|---|---|---|---|---|
| *Gen* | SDv1.5 | 0.34 | 0.35 | 0.32 | 0.28 | 0.29 | 0.21 | 0.32 |
| | SDXL | 0.43 | 0.48 | 0.47 | 0.44 | 0.45 | 0.27 | 0.43 |
| | SD3.5-large | 0.44 | 0.50 | 0.58 | 0.44 | 0.52 | 0.31 | 0.46 |
| | PixArt-Alpha | 0.45 | 0.50 | 0.48 | 0.49 | 0.56 | 0.34 | 0.47 |
| | playground-v2.5 | 0.49 | 0.58 | 0.55 | 0.43 | 0.48 | 0.33 | 0.49 |
| | FLUX.1-dev | 0.48 | 0.58 | 0.62 | 0.42 | 0.51 | 0.35 | 0.50 |
| *Unified* | Janus | 0.16 | 0.26 | 0.35 | 0.28 | 0.30 | 0.14 | 0.23 |
| | VILA-U | 0.26 | 0.33 | 0.37 | 0.35 | 0.39 | 0.23 | 0.31 |
| | Show-o | 0.28 | 0.40 | 0.48 | 0.30 | 0.46 | 0.30 | 0.35 |
| | Emu3 | 0.34 | 0.45 | 0.48 | 0.41 | 0.45 | 0.27 | 0.39 |
| | MetaQuery-XL | 0.56 | 0.55 | 0.62 | 0.49 | 0.63 | 0.41 | 0.55 |
| | GPT-4o | 0.81 | 0.71 | 0.89 | 0.83 | 0.79 | 0.74 | 0.80 |
| *Bagel* | Vanilla | 0.44 | 0.55 | 0.68 | 0.44 | 0.60 | 0.39 | 0.52 |
| | ImAgent | **0.63** | **0.63** | **0.72** | **0.59** | **0.69** | **0.53** | **0.63** |
| *Janus-Pro-7B* | Vanilla | 0.30 | 0.37 | 0.49 | 0.36 | 0.42 | 0.26 | 0.35 |
| | ImAgent | **0.44** | **0.46** | **0.55** | **0.45** | **0.50** | **0.27** | **0.44** |

Table 2. Experimental results on **WISE** [30]. Vanilla represents the backbone model, while `ImAgent` denotes our agent built upon this model. The best performance between vanilla and `ImAgent` is highlighted in **bold**.

| Types | Model | Idiom | | Textual | | Entity | | Scientific | | Overall | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Acc. | Qual. | Acc. | Qual. | Acc. | Qual. | Acc. | Qual. | Acc. | Qual. |
| *Gen* | FLUX.1-dev | 39.1 | 83.4 | 56.9 | 76.5 | 45.1 | 90.6 | 46.7 | 80.9 | 47.0 | 82.8 |
| | FLUX.1-schnell | 40.9 | 83.1 | 65.1 | 74.5 | 44.8 | 91.5 | 50.7 | 83.0 | 50.4 | 83.0 |
| | playground-v2.5 | 43.9 | 87.8 | 38.5 | 72.1 | 48.4 | 92.4 | 50.8 | 83.3 | 45.4 | 83.9 |
| | SD-3-Medium | 35.9 | 81.4 | 60.9 | 71.3 | 42.4 | 90.1 | 50.9 | 81.7 | 47.5 | 81.1 |
| | SD-3.5-Medium | 34.4 | 80.6 | 58.0 | 70.1 | 44.8 | 92.1 | 49.9 | 83.0 | 46.8 | 81.4 |
| | SD-3.5-Large | 35.6 | 85.3 | 62.2 | 75.4 | 46.6 | 92.6 | 52.9 | 84.5 | 49.3 | 84.4 |
| | gpt-image-1 | 75.7 | 94.5 | 86.9 | 97.6 | 77.5 | 96.6 | 74.7 | 94.3 | 78.7 | 95.8 |
| *Unified* | Emu3 | 33.1 | 82.9 | 33.7 | 68.7 | 33.8 | 85.2 | 40.1 | 77.0 | 35.2 | 78.5 |
| | Janus-Pro-7B | 25.5 | 78.0 | 37.2 | 70.9 | 38.5 | 87.6 | 44.9 | 77.8 | 36.5 | 78.6 |
| | show-o | 33.1 | 82.5 | 35.3 | 80.3 | 34.9 | 87.4 | 41.6 | 76.6 | 36.2 | 81.7 |
| | GoT | 29.7 | 76.4 | 30.6 | 70.7 | 31.0 | 86.2 | 36.8 | 76.3 | 32.0 | 77.4 |
| | Gemini-2.0 | 52.4 | 87.8 | 73.0 | 83.3 | 67.0 | 94.3 | 66.7 | 89.3 | 64.8 | 88.7 |
| *Bagel* | Vanilla | 30.2 | 85.7 | 36.6 | 68.4 | 45.0 | 94.7 | 54.4 | 87.5 | 41.6 | 84.1 |
| | ImAgent | **37.7** | **90.0** | **54.2** | **79.1** | **52.6** | **96.6** | **61.2** | **90.3** | **51.4** | **89.0** |
| *Janus-Pro-7B* | Vanilla | 25.5 | 78.0 | **37.2** | **70.9** | 38.5 | 87.6 | 44.9 | 77.8 | 36.5 | 78.6 |
| | ImAgent | **27.9** | **86.0** | 35.3 | 68.8 | **40.7** | **89.7** | **51.2** | **84.2** | **38.8** | **82.2** |

Table 3. Experimental results on **T2I-ReasonBench** [43]. Vanilla represents the backbone model, while `ImAgent` denotes our agent built upon this model. The best performance between vanilla and `ImAgent` is highlighted in **bold**.

Show-o [60], Janus-Pro-7B [5], Emu3 [52], Lumina-Image 2.0 [34], show-o [60], GoT [8] and MetaQuery-XL [32] as unified multimodal models.

For image editing, we evaluate both private and open baselines. The private baselines include Gemini-2.0 [10], Doubao [40], GPT-4o [31], and Gemini-2.0-Flash-pre [46]. The open baselines include Step1X-Edit [24], Instruct-Pix2Pix [2], MagicBrush [68], AnyEdit [66], Omni-Gen [58], EMU2 [45], Qwen-Image-Edit [55], FLUX.1-Kontext-Dev [17], Ovis-U1 [51], and Seedream-4.0 [38].

### 4.2. Quantitative Results

**Image Generation.** As shown in Table 2, `ImAgent` achieves outstanding performance on the WISE benchmark.

Compared to their respective backbone models, `ImAgent` built upon Bagel and Janus-Pro-7B consistently surpasses the vanilla counterparts, with performance improvements of 21.2% and 25.7%, respectively, demonstrating the effectiveness of our approach in enhancing image generation quality through coordinated reasoning and refinement. Notably, the vanilla Janus-Pro-7B underperforms compared to SDXL and EMU3 and performs on par with Show-o. In contrast, `ImAgent` constructed on Janus-Pro-7B surpasses all of them, highlighting its strong test-time scaling capability and generalization potential. Furthermore, on the reasoning benchmarks, `ImAgent` consistently delivers strong results. As shown in Tables 1 and 3, `ImAgent` outper-

| Types | Model | Temporal | Causal | Spatial | Logical | Overall |
|-------|-------|----------|--------|---------|---------|---------|
| *Private* | Gemini-2.0-Flash-pre | 10.6 | 13.3 | 11.0 | 2.3 | 9.4 |
| | Gemini-2.0 | 8.2 | 15.5 | 23.0 | 4.7 | 13.3 |
| | GPT-4o | 34.1 | 32.2 | 37.0 | 10.6 | 28.9 |
| *Open* | EMU2 | 1.2 | 1.1 | 0.0 | 0.0 | 0.5 |
| | OmniGen | 1.2 | 1.0 | 0.0 | 1.2 | 0.8 |
| | Step1X-Edit | 0.0 | 2.2 | 2.0 | 3.5 | 1.9 |
| | Qwen-Image-Edit | 4.7 | 10.0 | 17.0 | 2.4 | 8.9 |
| | FLUX.1-Kontext-Dev | 2.3 | 5.5 | 13.0 | 1.2 | 5.8 |
| | Ovis-U1 | 1.2 | 3.3 | 4.0 | 2.4 | 2.8 |
| | Seedream-4.0 | 12.9 | 12.2 | 11.0 | 7.1 | 10.8 |
| *Bagel* | Vanilla | 2.4 | 5.6 | 14.0 | 1.2 | 6.1 |
| | ImAgent | **17.6** | **15.6** | **16.0** | **2.4** | **13.1** |

Table 4. Experimental results on **RISEBench** [71]. Vanilla represents the backbone model, while `ImAgent` denotes our agent built upon this model. The best performance between vanilla and `ImAgent` is highlighted in **bold**.
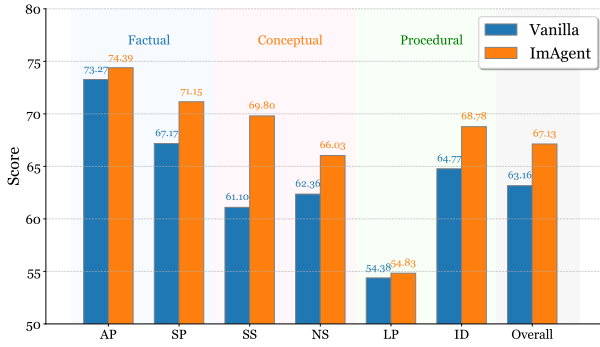


Figure 4. Experimental results on **KRISBench** [57].

forms the vanilla models across both backbone architectures. For instance, `ImAgent` achieves improvements of 14.8% and 7.5% across both backbones on R2I-Bench, respectively. Moreover, `ImAgent` built upon Bagel surpasses the commercial model Gemini-2.0 in overall quality on the T2I-ReasonBench, further demonstrating its effectiveness in reasoning-based image generation.

**Image Editing.** As shown in Table 4, vanilla Bagel achieves only 6.1 on RISEBench, underperforming compared to Qwen-Image-Edit and Seedream-4.0. In contrast, `ImAgent` achieves a 114.8% improvement over vanilla and surpasses both Qwen-Image-Edit and Seedream-4.0. Remarkably, `ImAgent` performs on par with Gemini-2.0, demonstrating its effectiveness and strong test-time scaling capability. This indicates that open-source models enhanced with `ImAgent` can achieve competitive performance with commercial models.

Furthermore, as shown in Figure 4, `ImAgent` also achieves excellent performance on KRISBench, attaining an overall score of 67.13, surpassing vanilla Bagel which scores 63.16. `ImAgent` performs well on ImgEdit-Bench (shown in Figure 3) as well, with vanilla achieving

2.89 and underperforming compared to Step1X-Edit, while `ImAgent` reaches 3.15, outperforming Step1X-Edit and achieving the best performance among open-source models. This further demonstrates the effectiveness of `ImAgent` in test-time scaling.

Moreover, `ImAgent` outperforms vanilla Bagel across both English and Chinese settings in GEdit-Bench, achieving improvements of 5.5% and 5.2%, respectively. Detailed results are provided in the Appendix.

### 4.3. Qualitative Analysis

**Image Refinement Flow.** As shown in Figure 1, we present several qualitative examples, including two for image generation and two for image editing, each executed for five steps. From an overall perspective, we can observe that the image quality improves progressively with each step, demonstrating that `ImAgent` refines the image gradually.

Specifically, consider a case where the user aims to generate an image based on the prompt "moldy bread." In the first step, `ImAgent` performs the *Naive Generation* action and produces an image containing fresh bread, which clearly deviates from the prompt. The policy controller then selects the *Prompt Refinement* action to improve the prompt using both the generated image and the user's original input. Although the updated image shows slight improvement, it remains unsatisfactory. To address this, the controller chooses the *Image Detail Refinement* action, leading to the appearance of mold on the bread. Subsequently, `ImAgent` applies the *Prompt Enhancement* action to further strengthen the textual guidance, resulting in an image that better aligns with the user's intent. Finally, recognizing that the prompt is now sufficiently clear but the visual result could still be improved, `ImAgent` executes the *Best-of-N Sampling* action to generate multiple candidates and select the most aligned one. The final output successfully captures the desired "moldy bread" concept.
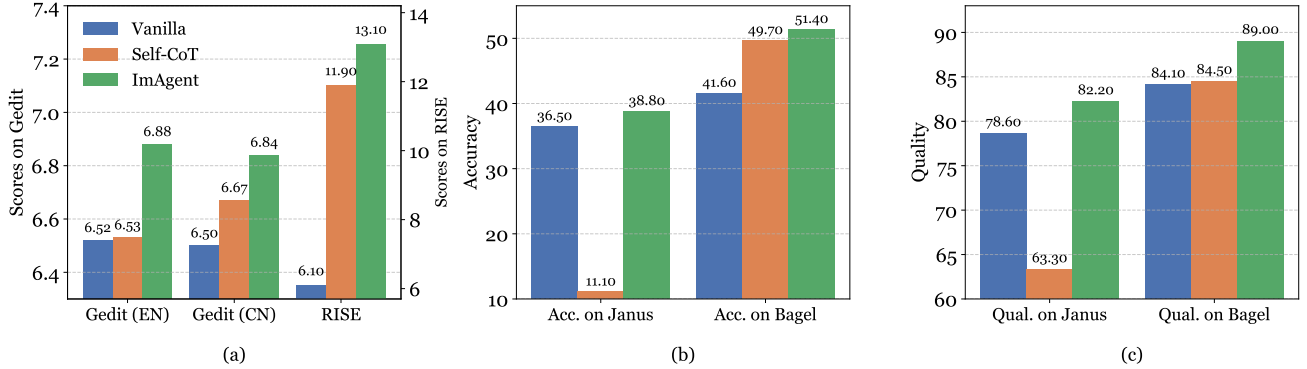
Figure 5. Comparison among the vanilla model, self-CoT, and `ImAgent`. (a) Results on Gedit-Bench (English and Chinese) and RISEBench using the Bagel backbone. (b) Accuracy and (c) Quality on Janus-Pro-7B and Bagel evaluated with T2I-ReasonBench.



Figure 6. An example of the policy controller making a decision.

**Detailed Policy Controller Decision.** We present a detailed case study of the policy controller's output to illustrate how it makes action decisions. As shown in Figure 6, which visualizes the decision-making process for the case discussed earlier, the policy controller selects the `prompt_refinement` action to enrich the textual prompt and improve semantic clarity for subsequent generation steps. We present more examples in the Appendix.

## 4.4. Ablation Studies

**Comparison with Self-CoT.** Since the vanilla Bagel model already incorporates self-CoT capabilities, we further compare `ImAgent` with self-CoT to evaluate the effectiveness of our approach. For the image editing task, as shown in Figure 5(a), `ImAgent` significantly outperforms both the vanilla and self-CoT variants on the Gedit-Bench across both English and Chinese settings. Furthermore, `ImAgent` also surpasses self-CoT in the image generation task. As illustrated in Figures 5(b) and (c), under the Janus-Pro-7B backbone, the performance of self-CoT notably drops compared to the vanilla model, whereas `ImAgent` consistently achieves superior results. Similarly, under the Bagel backbone, `ImAgent` maintains clear advantages over both the vanilla and self-CoT models. These findings suggest that the improvements brought by `ImAgent` stem not only from self-CoT–like reasoning refinement, but also from the different actions coordination enabled by the strong policy controller.

| Method | $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | Random | ImAgent |
|--------|-------|-------|-------|-------|-------|--------|---------|
| Score | 0.54 | 0.58 | 0.57 | 0.56 | 0.56 | 0.59 | **0.62** |

Table 5. Comparison of different action selection policies on R2I-Bench. $a_1$–$a_5$ represent the five actions (Section 3.2) in order.

**Effectiveness of the Whole `ImAgent`.** To evaluate the effectiveness of `ImAgent` in coordinately selecting and executing appropriate actions, we compare its performance with that of single-turn actions and a random action selection policy, where the $T_{max}$ for random selection is set the same as in `ImAgent`. As shown in Table 5, experiments are conducted on R2I-Bench [4]. All single-turn actions perform well compared to the naive generation action, which serves as the lower bound (0.54). Moreover, the random action selection policy achieves a score of 0.59, higher than all single-turn actions, indicating that naively constructing an agent is beneficial. Finally, `ImAgent` outperforms the random policy with a score of 0.62, demonstrating that it effectively and efficiently boosts test-time scaling in image generation.

## 5. Conclusion

In this paper, we presented `ImAgent`, a training-free unified multimodal agent designed to enhance image generation through efficient test-time scaling. By integrating multiple generation actions within a single framework and coordinating them via a policy controller, `ImAgent` dynamically selects and executes the most appropriate action for a given case. Extensive experiments on image generation and editing benchmarks demonstrate that `ImAgent` consistently improves over the backbone models and surpasses strong baseline methods, including scenarios where the backbone fails. These results highlight the potential of unified multimodal agents for adaptive, efficient, and high-quality image generation without relying on additional external models.

# References

[1] Stephen Brade, Bryan Wang, Mauricio Sousa, Sageev Oore, and Tovi Grossman. Promptify: Text-to-image generation through interactive prompt exploration with large language models. In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, pages 1–14, 2023. 1, 4

[2] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18392–18402, 2023. 6

[3] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart-$\sigma$: Weak-to-strong training of diffusion transformer for 4k text-to-image generation. In *European Conference on Computer Vision*, pages 74–91. Springer, 2024. 5

[4] Kaijie Chen, Zihao Lin, Zhiyang Xu, Ying Shen, Yuguang Yao, Joy Rimchala, Jiaxin Zhang, and Lifu Huang. R2i-bench: Benchmarking reasoning-driven text-to-image generation. *arXiv preprint arXiv:2505.23493*, 2025. 5, 8

[5] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025. 5, 6

[6] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 1, 2, 5

[7] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 2, 5

[8] Rongyao Fang, Chengqi Duan, Kun Wang, Linjiang Huang, Hao Li, Shilin Yan, Hao Tian, Xingyu Zeng, Rui Zhao, Jifeng Dai, et al. Got: Unleashing reasoning capability of multimodal large language model for visual generation and editing. *arXiv preprint arXiv:2503.10639*, 2025. 6

[9] Peng Gao, Le Zhuo, Dongyang Liu, Ruoyi Du, Xu Luo, Longtian Qiu, Yuhang Zhang, Chen Lin, Rongjie Huang, Shijie Geng, et al. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers. *arXiv preprint arXiv:2405.05945*, 2024. 2

[10] Google. Experiment with gemini 2.0 flash native image generation. https://www.google.com, 2025. Accessed: 2025-11-09. 6

[11] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 1, 3

[12] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 5

[13] Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, et al. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*, 2024. 5

[14] Jinwoo Jeon, JunHyeok Oh, Hayeong Lee, and Byung-Jun Lee. Iterative prompt refinement for safer text-to-image generation. *arXiv preprint arXiv:2509.13760*, 2025. 1, 3

[15] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. Accessed: 2025-11-02. 5

[16] Black Forest Labs. Flux. https://github.com/black-forest-labs/flux, 2024. 5

[17] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, Sumith Kulal, Kyle Lacey, Yam Levi, Cheng Li, Dominik Lorenz, Jonas Müller, Dustin Podell, Robin Rombach, Harry Saini, Axel Sauer, and Luke Smith. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space, 2025. 6

[18] Kyungmin Lee, Sangkyung Kwak, Kihyuk Sohn, and Jinwoo Shin. Direct consistency optimization for robust customization of text-to-image diffusion models. *Advances in neural information processing systems*, 37:103269–103304, 2024. 1

[19] Chunyi Li, Haoning Wu, Hongkun Hao, Zicheng Zhang, Tengchuan Kou, Chaofeng Chen, Lei Bai, Xiaohong Liu, Weisi Lin, and Guangtao Zhai. G-refine: A general quality refiner for text-to-image generation. In *Proceedings of the 32nd ACM International Conference on Multimedia*, pages 7375–7384, 2024. 1, 3

[20] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2. 5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv preprint arXiv:2402.17245*, 2024. 5

[21] Shufan Li, Konstantinos Kallidromitis, Akash Gokul, Arsh Koneru, Yusuke Kato, Kazuki Kozuka, and Aditya Grover. Reflect-dit: Inference-time scaling for text-to-image diffusion transformers via in-context reflection. *arXiv preprint arXiv:2503.12271*, 2025. 1, 3, 4

[22] Mushui Liu, Yuhang Ma, Zhen Yang, Jun Dan, Yunlong Yu, Zeng Zhao, Zhipeng Hu, Bai Liu, and Changjie Fan. Llm4gen: Leveraging semantic representation of llms for text-to-image generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 5523–5531, 2025. 5

[23] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, et al. Step1x-edit: A practical framework for general image editing, 2025. *URL https://arxiv. org/abs/2504.17761*. 5

[24] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, et al. Step1x-edit: A practical framework for general image editing. *arXiv preprint arXiv:2504.17761*, 2025. 6

[25] Tao Liu, Kai Wang, Senmao Li, Joost van de Weijer, Fahad Shahbaz Khan, Shiqi Yang, Yaxing Wang, Jian Yang, and Ming-Ming Cheng. One-prompt-one-story: Free-lunch

consistent text-to-image generation using a single prompt. *arXiv preprint arXiv:2501.13554*, 2025. 1, 4

[26] Haiying Luo, Tao He, and Zhang Yi. A stable mapping of nmode. *Artificial Intelligence Review*, 57(5):120, 2024. 2

[27] Shuailei Ma, Kecheng Zheng, Ying Wei, Wei Wu, Fan Lu, Yifei Zhang, Chen-Wei Xie, Biao Gong, Jiapeng Zhu, and Yujun Shen. Learning visual generative priors without text. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 8051–8061, 2025. 5

[28] Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, et al. Self-refine: Iterative refinement with self-feedback. *Advances in Neural Information Processing Systems*, 36:46534–46594, 2023. 3

[29] Wenyi Mo, Tianyu Zhang, Yalong Bai, Bing Su, Ji-Rong Wen, and Qing Yang. Dynamic prompt optimizing for text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26627–26636, 2024. 1, 4

[30] Yuwei Niu, Munan Ning, Mengren Zheng, Weiyang Jin, Bin Lin, Peng Jin, Jiaqi Liao, Chaoran Feng, Kunpeng Ning, Bin Zhu, et al. Wise: A world knowledge-informed semantic evaluation for text-to-image generation. *arXiv preprint arXiv:2503.07265*, 2025. 5, 6

[31] OpenAI. Introducing 4o image generation. https://openai.com/index/introducing-4o-image-generation, 2025. Accessed: 2025-11-09. 6

[32] Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang Wang, Zhiyang Xu, Jiuhai Chen, Kunpeng Li, Felix Juefei-Xu, et al. Transfer between modalities with metaqueries. *arXiv preprint arXiv:2504.06256*, 2025. 6

[33] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1, 2, 5

[34] Qi Qin, Le Zhuo, Yi Xin, Ruoyi Du, Zhen Li, Bin Fu, Yiting Lu, Jiakang Yuan, Xinyue Li, Dongyang Liu, et al. Lumina-image 2.0: A unified and efficient image generative framework. *arXiv preprint arXiv:2503.21758*, 2025. 1, 5, 6

[35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1 (2):3, 2022. 2

[36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 5

[37] Axel Sauer, Tero Karras, Samuli Laine, Andreas Geiger, and Timo Aila. Stylegan-t: Unlocking the power of gans for fast large-scale text-to-image synthesis. In *International conference on machine learning*, pages 30105–30118. PMLR, 2023. 1, 3

[38] Team Seedream, Yunpeng Chen, Yu Gao, Lixue Gong, Meng Guo, Qiushan Guo, Zhiyao Guo, Xiaoxia Hou, Weilin

Huang, Yixuan Huang, et al. Seedream 4.0: Toward next-generation multimodal image generation, 2025. *URL https://arxiv. org/abs/2509.20427*. 6

[39] Dazhong Shen, Guanglu Song, Zeyue Xue, Fu-Yun Wang, and Yu Liu. Rethinking the spatial inconsistency in classifier-free diffusion guidance. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9370–9379, 2024. 1, 3

[40] Yichun Shi, Peng Wang, and Weilin Huang. Seededit: Align image re-generation to image editing, 2024. *URL https://arxiv. org/abs/2411.06686*. 6

[41] Jaskirat Singh and Liang Zheng. Divide, evaluate, and refine: Evaluating and improving text-to-image alignment with iterative vqa feedback. *Advances in Neural Information Processing Systems*, 36:70799–70811, 2023. 1, 3

[42] Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025. 1

[43] Kaiyue Sun, Rongyao Fang, Chengqi Duan, Xian Liu, and Xihui Liu. T2i-reasonbench: Benchmarking reasoning-informed text-to-image generation. *arXiv preprint arXiv:2508.17472*, 2025. 5, 6

[44] Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan. Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint arXiv:2406.06525*, 2024. 1, 2, 5

[45] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiying Yu, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14398–14409, 2024. 6

[46] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023. 6

[47] Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling: Scalable image generation via next-scale prediction. *Advances in neural information processing systems*, 37:84839–84865, 2024. 1, 2

[48] Rui Tian, Mingfei Gao, Mingze Xu, Jiaming Hu, Jiasen Lu, Zuxuan Wu, Yinfei Yang, and Afshin Dehghan. Unigen: Enhanced training & test-time strategies for unified multimodal understanding and generation. *arXiv preprint arXiv:2505.14682*, 2025. 3, 4

[49] Shengbang Tong, David Fan, Jiachen Li, Yunyang Xiong, Xinlei Chen, Koustuv Sinha, Michael Rabbat, Yann LeCun, Saining Xie, and Zhuang Liu. Metamorph: Multimodal understanding and generation via instruction tuning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 17001–17012, 2025. 2, 3

[50] Claudio Mayrink Verdun, Alex Oesterling, Himabindu Lakkaraju, and Flavio P Calmon. Soft best-of-n sampling for model alignment. *arXiv preprint arXiv:2505.03156*, 2025. 3, 4

[51] Guo-Hua Wang, Shanshan Zhao, Xinjie Zhang, Liangfu Cao, Pengxin Zhan, Lunhao Duan, Shiyin Lu, Minghao Fu, Jianshan Zhao, Yang Li, and Qing-Guo Chen. Ovis-u1 technical report. *arXiv preprint arXiv:2506.23044*, 2025. 6

[52] Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan Zhang, Yueze Wang, Zhen Li, Qiying Yu, et al. Emu3: Next-token prediction is all you need. *arXiv preprint arXiv:2409.18869*, 2024. 2, 6

[53] Zhijie Wang, Yuheng Huang, Da Song, Lei Ma, and Tianyi Zhang. Promptcharm: Text-to-image generation through multi-modal prompting and refinement. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–21, 2024. 3

[54] Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified multimodal understanding and generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 12966–12977, 2025. 1, 3, 5

[55] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan Cai, and Zenan Liu. Qwen-image technical report, 2025. 6

[56] Yecheng Wu, Zhuoyang Zhang, Junyu Chen, Haotian Tang, Dacheng Li, Yunhao Fang, Ligeng Zhu, Enze Xie, Hongxu Yin, Li Yi, et al. Vila-u: a unified foundation model integrating visual understanding and generation. *arXiv preprint arXiv:2409.04429*, 2024. 5

[57] Yongliang Wu, Zonghui Li, Xinting Hu, Xinyu Ye, Xianfang Zeng, Gang Yu, Wenbo Zhu, Bernt Schiele, Ming-Hsuan Yang, and Xu Yang. Kris-bench: Benchmarking next-level intelligent image editing models. *arXiv preprint arXiv:2505.16707*, 2025. 5, 7

[58] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chaofan Li, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13294–13304, 2025. 2, 3, 5, 6

[59] Enze Xie, Junsong Chen, Yuyang Zhao, Jincheng Yu, Ligeng Zhu, Chengyue Wu, Yujun Lin, Zhekai Zhang, Muyang Li, Junyu Chen, et al. Sana 1.5: Efficient scaling of training-time and inference-time compute in linear diffusion transformer. *arXiv preprint arXiv:2501.18427*, 2025. 2, 5

[60] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 1, 2, 6

[61] Tianyi Xiong, Xiyao Wang, Dong Guo, Qinghao Ye, Haoqi Fan, Quanquan Gu, Heng Huang, and Chunyuan Li. Llava-critic: Learning to evaluate multimodal models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 13618–13628, 2025. 2

[62] Xingqian Xu, Jiayi Guo, Zhangyang Wang, Gao Huang, Irfan Essa, and Humphrey Shi. Prompt-free diffusion: Taking" text" out of text-to-image diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8682–8692, 2024. 1

[63] Zhengyuan Yang, Jianfeng Wang, Linjie Li, Kevin Lin, Chung-Ching Lin, Zicheng Liu, and Lijuan Wang. Idea2img: Iterative self-refinement with gpt-4v (ision) for automatic image design and generation. *arXiv preprint arXiv:2310.08541*, 2023. 3

[64] Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Retrieval-augmented multimodal language modeling. *arXiv preprint arXiv:2211.12561*, 2022. 3

[65] Yang Ye, Xianyi He, Zongjian Li, Bin Lin, Shenghai Yuan, Zhiyuan Yan, Bohan Hou, and Li Yuan. Imgedit: A unified image editing dataset and benchmark. *arXiv preprint arXiv:2505.20275*, 2025. 5

[66] Qifan Yu, Wei Chow, Zhongqi Yue, Kaihang Pan, Yang Wu, Xiaoyang Wan, Juncheng Li, Siliang Tang, Hanwang Zhang, and Yueting Zhuang. Anyedit: Mastering unified high-quality image editing for any idea. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 26125–26135, 2025. 6

[67] Jihai Zhang, Tianle Li, Linjie Li, Zhengyuan Yang, and Yu Cheng. Are unified vision-language models necessary: Generalization across understanding and generation. *arXiv preprint arXiv:2505.23043*, 2025. 3

[68] Kai Zhang, Lingbo Mo, Wenhu Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *Advances in Neural Information Processing Systems*, 36:31428–31449, 2023. 6

[69] Qian Zhang, Xiangzi Dai, Ninghua Yang, Xiang An, Ziyong Feng, and Xingyu Ren. Var-clip: Text-to-image generator with visual auto-regressive modeling. *arXiv preprint arXiv:2408.01181*, 2024. 2

[70] Xinjie Zhang, Jintao Guo, Shanshan Zhao, Minghao Fu, Lunhao Duan, Jiakui Hu, Yong Xien Chng, Guo-Hua Wang, Qing-Guo Chen, Zhao Xu, et al. Unified multimodal understanding and generation models: Advances, challenges, and opportunities. *arXiv preprint arXiv:2505.02567*, 2025. 3

[71] Xiangyu Zhao, Peiyuan Zhang, Kexian Tang, Xiaorong Zhu, Hao Li, Wenhao Chai, Zicheng Zhang, Renqiu Xia, Guangtao Zhai, Junchi Yan, et al. Envisioning beyond the pixels: Benchmarking reasoning-informed visual editing. *arXiv preprint arXiv:2504.02826*, 2025. 5, 7

[72] Tong Zheng, Lichang Chen, Simeng Han, R Thomas McCoy, and Heng Huang. Learning to reason via mixture-of-thought for logical reasoning. *arXiv preprint arXiv:2505.15817*, 2025. 2

[73] Tong Zheng, Hongming Zhang, Wenhao Yu, Xiaoyang Wang, Xinyu Yang, Runpeng Dai, Rui Liu, Huiwen Bao, Chengsong Huang, Heng Huang, et al. Parallel-r1: Towards

parallel thinking via reinforcement learning. *arXiv preprint arXiv:2509.07980*, 2025. 2

[74] Chunting Zhou, Lili Yu, Arun Babu, Kushal Tirumala, Michihiro Yasunaga, Leonid Shamis, Jacob Kahn, Xuezhe Ma, Luke Zettlemoyer, and Omer Levy. Transfusion: Predict the next token and diffuse images with one multi-modal model. *arXiv preprint arXiv:2408.11039*, 2024. 2, 3