

# SIDIFFAGENT: SELF-IMPROVING DIFFUSION AGENT

**Anonymous authors**

Paper under double-blind review

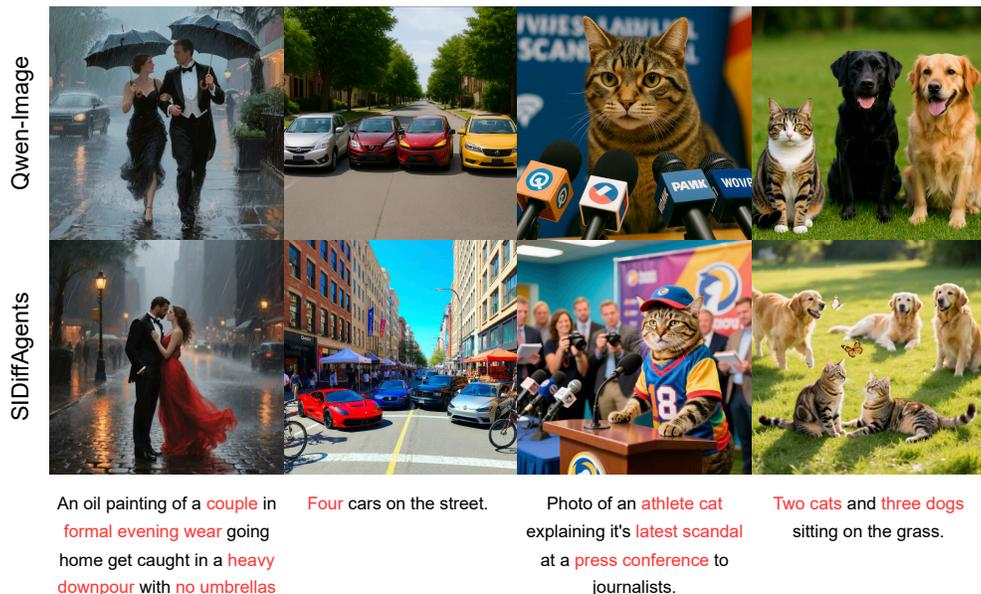


Figure 1: Qualitative comparison between our model, *SIDiffAgent*, and the state-of-the-art, Qwen-Image model. *SIDiffAgent* demonstrates superior text-to-image alignment and realism, accurately rendering prompts that require nuanced understanding. Notably, our model successfully handles: negative constraints (‘no umbrellas’), complex scene composition, anthropomorphic roles (‘athlete cat’), and precise object counting.

## ABSTRACT

Text-to-image diffusion models have revolutionized generative AI, enabling high-quality and photorealistic image synthesis. However, their practical deployment remains hindered by several limitations: sensitivity to prompt phrasing, ambiguity in semantic interpretation (e.g., “mouse” as animal vs. a computer peripheral), artifacts such as distorted anatomy, and the need for carefully engineered input prompts. Existing methods often require additional training and offer limited controllability, restricting their adaptability in real-world applications. We introduce Self-Improving Diffusion Agent (*SIDiffAgent*), a training-free agentic framework that leverages the Qwen family of models (Qwen-VL, Qwen-Image, Qwen-Edit, Qwen-Embedding) to address these challenges. *SIDiffAgent* autonomously manages prompt engineering, detects and corrects poor generations, and performs fine-grained artifact removal, yielding more reliable and consistent outputs. It further incorporates iterative self-improvement by storing a memory of previous experiences in a database. This database of past experiences is then used to inject prompt-based guidance at each stage of the agentic pipeline. *SIDiffAgent* achieved an average VQA score of 0.884 on GenAIBench, significantly outperforming open-source, proprietary models and agentic methods. We will publicly release our code upon acceptance.

# 1 INTRODUCTION

Image generation models (Rombach et al., 2021; Ramesh et al., 2022; Saharia et al., 2022a; Singer et al., 2022; Ho et al., 2022; Blattmann et al., 2023; Esser et al., 2023) have transformed digital image synthesis, enabling the creation of high-quality visuals with remarkable detail and realism from simple text descriptions. Despite these advances, the performance of state-of-the-art text-to-image diffusion models remains highly sensitive to the precise phrasing of input prompts, making them unable to capture the true intent of the user and creating a significant “intent gap” between the user’s intended meaning and the textual instructions provided to the model (Chen et al., 2025a). Here, intent refers to the underlying concept or goal that the user wishes the model to capture. This gap arises because natural language prompts are often ambiguous, underspecified, or casually phrased, leading the model to generate outputs that deviate from the user’s actual goal. For instance, the prompt “a mouse” may yield different outputs, either an animal or a computer peripheral, depending on how the model interprets the text.

Another key challenge lies in the distributional mismatch between the unstructured, free-form language typical of end-user prompts and the highly curated, descriptive captions created by human experts and subjected to intensive filtering during the training of diffusion models (Li et al., 2024d; Wu et al., 2025a; Labs, 2024). This discrepancy often results in inadequate textual guidance for generation. For instance, prompts are frequently underspecified (e.g., “a car on a road” without details on the car’s model, color, or environment), which forces the model to make unguided assumptions. Such underspecification not only compromises output quality but also increases computational costs, as users must repeatedly regenerate outputs to achieve satisfactory results. These issues constitute a critical bottleneck for the scalability and reliability of diffusion models in professional domains such as digital marketing, content creation, and design, where robust and predictable alignment with user intent is essential.

To mitigate these challenges, prior work such as T2I-Copilot (Chen et al., 2025a) and MCCD (Li et al., 2025a) have proposed agentic frameworks that address ambiguities through image editing and agentic coordination. However, these approaches often lack fine-grained controllability and robust mechanisms for post-generation artifact correction. A key challenge in such Multi-Agent Systems is coordination. Agents often have only partial observability and limited knowledge of others capabilities or the effects of their own actions, leading to suboptimal outcomes (Cemri et al., 2025). Recent works in the field of LLM Agents address this via Theory of Mind (ToM), the ability to attribute goals and beliefs to others (Cross et al., 2024; Yang et al., 2025b). With ToM, agents can understand peers intentions and anticipate future actions, enabling better coordination. To the best of our knowledge, no work has applied this to diffusion-based agents.

Additionally, ReNeg (Li et al., 2025b) highlights the role of negative prompts and embeddings (Ho & Salimans, 2022) in constraining diffusion processes and guiding models away from undesired artifacts. In this work, we introduce a training-free, multi-agent framework *SIDiffAgent* that creates a robust and self-improving Theory of Mind inspired image generation pipeline. *SIDiffAgent* leverages the capabilities of the Qwen model family (Qwen-VL (Bai et al., 2025), Qwen-Image (Wu et al., 2025a), Qwen-Edit (Wu et al., 2025a), Qwen-Embedding (Zhang et al., 2025)) as the backbone of our agentic flow. Our framework uses multiple agents to refine prompts, reuse past experience, and fix artifacts with local edits. Further building upon the observations from ReNeg (Li et al., 2025b), *SIDiffAgent* also incorporates an additional adaptive negative prompt generation agent. Beyond controllability and artifact mitigation, we emphasize the importance of self-improvement: unlike prior agentic frameworks, *SIDiffAgent* is designed to leverage past successes and failures to iteratively refine its generation strategies. **While our contribution focuses on system-level architecture rather than a new training objective, we highlight three technically novel components. First, *SIDiffAgent* introduces the first experience-driven memory mechanism for diffusion-based agents, enabling retrieval-conditioned guidance throughout the generation workflow. Second, our Theory-of-Mind-inspired inter-agent reasoning models predictive expectations over other agents’ behaviors using accumulated success/pitfall patterns, providing dynamic workflow guidance beyond prompt optimization. Third, *SIDiffAgent* offers a training-free alternative to RL-based or reward-guided finetuning approaches (e.g., DPO-Diffusion, ReNeg), enabling practical quality improvements even for closed-weight diffusion systems without any retraining.**

We empirically validate the effectiveness of *SIDiffAgent* across multiple benchmarks, including GenAIBench and DrawBench, where it establishes new state-of-the-art performance. Our approach not only surpasses prior agentic systems such as T2I-Copilot, achieving a VQA Score (Lin et al., 2024) improvement of 8.73%, but also outperforms leading proprietary models (e.g., Recraft V3, Imagen 3, Flux 1.1-Pro), delivering a 5.36% improvement over Imagen 3. Furthermore, it significantly exceeds the performance of powerful open-source systems (e.g., Flux.1-Dev, HunyuanDiT v1.2, and SD 3.5), achieving a 15.70% improvement over SD 3.5.

Our key contributions are summarized as follows:

- We introduce a multi-agent system that refines the inputs to a diffusion model at test time, ensuring consistent understanding of prompts even when they are expressed in different ways. This results in generations that more faithfully capture the user’s intended concepts.
- We adapt a structured image editing approach that enables targeted, incremental corrections to generated images without requiring complete regeneration, narrowing the intent gap and improving alignment between user goals and model generations.
- We build an iterative self-improvement mechanism in which an agentic memory tracks successes, artifacts and pitfalls of various agents, and using them to inject corrective guidance based on Theory of Mind. This creates a training-free, self-improving agentic framework that progressively enhances generation quality over time.

## 2 RELATED WORKS

### 2.1 AI AGENTS

Agents are autonomous entities capable of perceiving their environment, reasoning over goals, and executing actions, often collaboratively through structured memory, planning, and tool use (Luo et al., 2025). In commercial applications such as browser automation or coding (Hong et al., 2024; Qian et al., 2023), these agents typically rely on a large language model (LLM) backbone, which performs reasoning over the current state to determine the optimal action. An extension of this paradigm is the multi-agent framework, in which multiple agents collaborate not only to reason over their own states but also to reason about the states of other agents (Li et al., 2023). Building on these foundations, researchers have started exploring agentic workflows for image generation. In such systems, multi-agent pipelines decompose and refine prompts, evaluate outputs, and correct artifacts. Examples include T2I-Copilot (Chen et al., 2025a), which integrates prompt interpretation, model selection, and quality evaluation, and PromptSculptor (Xiang et al., 2025), which iteratively transforms vague user prompts into precise ones. These systems demonstrate how agentic principles can reduce the intent gap in text-to-image generation, improving fidelity, controllability, and efficiency without requiring additional model training.

### 2.2 THEORY OF MIND IN AGENTIC SYSTEMS

Building on multi-agent systems, Theory of Mind (ToM) (Rabinowitz et al., 2018; Rocha et al., 2023) has become an increasingly important concept for designing smarter and more adaptive agents. It allows agents to understand, infer, and reason about what other agents might believe, desire, or intend in a given context. Recent research in multi-agent and conversational settings has shown that endowing agents with this capability significantly improves both overall system performance and individual agent effectiveness (Yang et al., 2025a; Gu et al., 2025; Jim & Giles, 2001). To the best of our knowledge, this idea has not yet been systematically applied to diffusion-based generative models. In *SIDiffAgent*, ToM is implemented through specialized sub-agents that continuously track each other’s mental states, typical successes, potential pitfalls, and characteristic behavior patterns. This information is integrated with a self-improvement mechanism that learns from past experiences with similar inputs to dynamically update these states over time. As a result, each sub-agent can proactively anticipate potential pitfalls and successes of other agents, adjusting its own reasoning accordingly. This awareness allows agents to make decisions that improve not only their individual output quality but also the collective performance of the other sub-agents, ultimately leading to a more robust, adaptive, and coordinated agentic system.

## 2.3 PROMPT OPTIMIZATION FOR DIFFUSION MODELS

Prompt Optimization (PO) has emerged as a crucial technique for aligning large language models (LLMs) with human intent and improving downstream task performance. Originating from in-context learning, PO has evolved into automatic prompt engineering, where LLMs iteratively refine and adapt prompts to achieve more precise and contextually appropriate outputs. While PO has been extensively studied in the context of LLMs, its application to diffusion models remains comparatively underexplored, despite the growing importance of text-to-image generation. For instance, Promptist (Hao et al., 2023) fine-tunes an LLM via reinforcement learning from human feedback to augment user prompts with stylistic and artistic modifiers, while OPT2I (Mañas et al., 2024) leverages an LLM to refine prompts directly for higher-quality image synthesis. Other works, such as DPO-Diffusion (Wang et al., 2024a) and ReNeg (Li et al., 2025b), focus on optimizing negative prompts, either through LLM reasoning or learned embeddings that guide the generative process away from undesired features. TextGrad (Yuksekgonul et al., 2024) introduces gradient-based text feedback optimization, and manual strategies for prompt engineering have also been systematically investigated. Our agentic system builds on these ideas by jointly optimizing positive and image-specific negative prompts at test time via prompt engineering, dynamically balancing guidance to produce more reliable, semantically aligned, and visually coherent images without requiring additional model training.

## 2.4 SELF-IMPROVING AI

Recent work has also explored self-improving models and agents, where systems enhance their own performance via interaction or feedback. For LLMs, approaches such as LMSI (Huang et al., 2024) leverage reward signals for autonomous improvement, while IPO (Garg et al., 2025) and Self-Rewarding (Yuan et al., 2024) LLMs optimize through preference pair construction. Methods like SPPO (Wu et al., 2024b) apply self-play to iteratively refine capabilities. Beyond LLMs, agentic systems such as MaestroMotif (Klissarov et al., 2024) use feedback-driven skill acquisition, and AgentSquare (Shang et al., 2024) abstracts agent design into modular components for automated exploration. Other works, such as AlphaEvolve (Novikov et al., 2025) and ADAS (Hu et al., 2025), employ meta-agents to improve task-specific agents, while SICA (Robeyns et al., 2025) eliminates the meta/target distinction by enabling agents to directly edit their own codebases. While the field of self-improving LLMs and coding agents has been actively studied, diffusion-based agents have not been explored in this context. To the best of our knowledge, we present the first training-free, self-improving diffusion framework, where agentic memory and prompt engineering enable self-improvement without retraining.

## 3 PROPOSED METHOD: SELF-IMPROVING DIFFUSION AGENT (*SIDiffAgent*)

Our framework is built around an agentic workflow that coordinates prompt engineering, image generation, evaluation, and self-improvement. This workflow is orchestrated by three primary agents, the Generation Orchestrator Agent ( $A_{\text{ORC}}$ ), which preprocesses user prompts through four sub-agents, the Creativity Analysis Sub-Agent ( $S_{\text{CRE}}$ ), Intention Analysis Sub-Agent ( $S_{\text{INT}}$ ), Prompt Refinement Sub-Agent ( $S_{\text{REF}}$ ), and Adaptive Negative Prompt Sub-Agent ( $S_{\text{NEG}}$ ) and Generation Sub-Agent ( $S_{\text{GEN}}$ ) for the final generation. The Evaluation Agent ( $A_{\text{EVAL}}$ ) assesses each output, triggering editing when necessary. Finally, the Guidance Agent ( $A_{\text{GUID}}$ ) facilitates learning from past experiences. The overall workflow of *SIDiffAgent* is depicted in Figure 2 and the algorithm is given in Appendix J. Each agent is discussed in detail below.

### 3.1 GENERATION ORCHESTRATOR AGENT ( $A_{\text{ORC}}$ )

The Generation Orchestrator Agent forms the foundation of our approach. It performs a structured preprocessing sequence that balances user intent, creative enrichment, and model compatibility. To achieve this, it operates through five specialized sub-agents, each addressing a distinct aspect of prompt engineering and generation. The following is a detailed explanation of all the agents. The prompts for each of the agents are in the Appendix I.1.

**Step 1. Creativity Analysis Sub-Agent ( $S_{\text{CRE}}$ ) :** The first step is to estimate the freedom that the agents can have which guides the intensity of refinements in the later stages. It classifies the

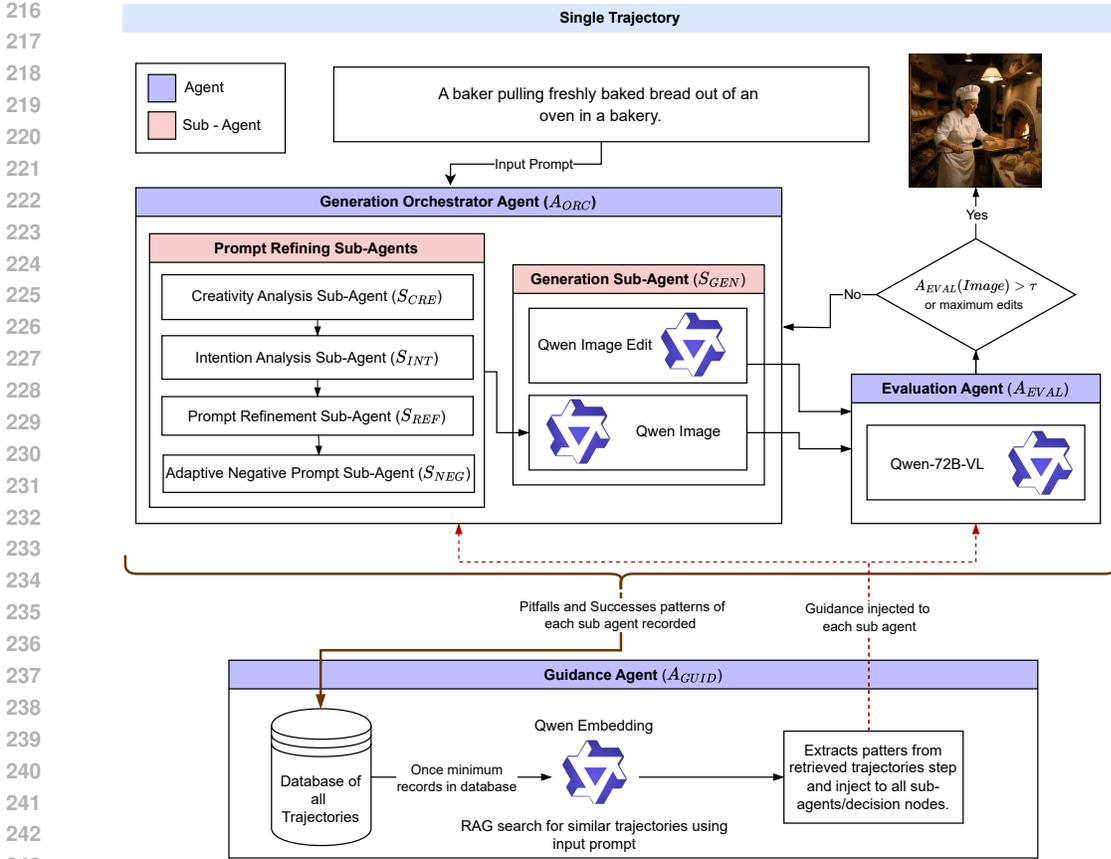


Figure 2: Workflow of the Self-Improving Diffusion Agent (*SIDiffAgent*). An input prompt is processed by the Generation Orchestrator Agent ( $A_{ORC}$ ), which employs sub-agents ( $S_{CRE}$ ,  $S_{INT}$ ,  $S_{REF}$ ,  $S_{NEG}$ ) to assess creativity, clarify intent, refine the prompt, and add adaptive negative constraints before generation ( $S_{GEN}$ ). The Evaluation Agent ( $A_{EVAL}$ ) scores the generated image on aesthetic quality and text–image alignment, triggering targeted refinements if the evaluation score is less than the pre-defined threshold  $\tau$ . Each trajectory is stored in the knowledge base, where the Guidance Agent ( $A_{GUID}$ ) stores pitfalls and successes into corrective and workflow guidance, which are injected back into decision nodes to improve future generations.

input into three levels: high for short or vague prompts, requiring substantial enrichment; medium for moderately detailed prompts, needing balance between structure and interpretation; and low for highly detailed prompts, where strict adherence is key. For instance, the prompt ‘a cat’ is classified as high creativity, (e.g., specifying breed, pose, or environment), whereas the prompt ‘a Persian cat sitting on a red sofa in a living room’ would be classified as low creativity, since it already provides detailed constraints. Figure 6 showcases the prompts and image generation pairs for different creativity levels.

**Step 2. Intention Analysis Sub-Agent ( $S_{INT}$ )** : After establishing the creativity level, this sub-agent deconstructs the users input to produce a comprehensive semantic layout of the prompt. The process extracts core content from the prompt, such as *main subjects*, their *attributes*, their *spatial relationships*, *background*, and additional parameters like *composition*, *lighting*, and overall *visual style*.

Beyond structural parsing, another critical function of this agent is to identify and resolve ambiguities or missing information within the prompt conditioned on the creativity level determined by ( $S_{CRE}$ ). For prompts designated as low creativity, the agent generates minor assumptions. Conversely, for high creativity prompts, it autonomously assumes information while adhering to the input prompt and ensuring coherent and contextually appropriate details. The output of this stage is

a structured specification that formalizes the user’s intent and generates clarifications. For example, given the prompt ‘a car’ at low creativity,  $S_{INT}$  preserves the minimal description without adding details. At high creativity, however, it expands the prompt by generating clarifying attributes such as the car model, color, background, etc. These details are autonomously inferred by  $S_{INT}$  to maintain coherence while enriching the original input.

**Step 3. Prompt Refinement Sub-Agent ( $S_{REF}$ ):** Building upon the semantic layout generated from  $S_{INT}$ ,  $S_{REF}$  produces a clearer and refined version of the prompt.  $S_{REF}$  ensures fidelity to the user’s intent by preserving all original subjects and scene elements, while simultaneously restructuring the prompt into a more contextually rich and coherent form.

**Step 4. Adaptive Negative Prompt Sub-Agent ( $S_{NEG}$ ):** In classifier-free guidance in diffusion models, the model normally contrasts the conditional prompt with an unconditional prompt (Ho & Salimans, 2022). Negative prompts are injected by replacing that unconditional branch with the negative text embedding, so the guidance steers outputs away from undesired features. This provides fine-grained control, enabling the explicit exclusion of unwanted objects, styles, or quality defects. While the refined prompt captures what should appear in the image, an equally important step is to constrain what must not appear.

$S_{NEG}$  first appends universal safeguards to mitigate common artifacts (e.g., ‘low quality, blurry, distorted, watermark’). Then it analyzes the refined prompt to derive semantically relevant, scene-specific negations. For instance, a prompt specifying a “clear blue sky” adds “clouds, dark clouds” to the negative prompt, while “a person standing alone” generates an exclusion for “extra people or crowd”. These adaptive constraints are concise, context-aware, thereby improving output quality and adherence without undermining user intent.

**Step 5. Generation Sub-Agent ( $S_{GEN}$ ):** The final component in  $A_{ORC}$  is the generation sub-agent, which is responsible for the task of generating the image. This sub-agent has 2 models (Qwen-Image and Qwen-Image-Edit) to handle both initial creation and subsequent refinement. Initially, it employs Qwen-Image to generate an image from the refined positive and adaptive negative prompts. Following this step, the evaluation is initiated (Discussed in Section 3.2). If  $A_{EVAL}$  identifies defects or misalignments,  $S_{GEN}$  invokes Qwen-Image-Edit, a specialized editing model, to perform targeted corrections. This strategy of combining generation with precise editing allows the system to iteratively improve the output, ensuring higher fidelity to complex user specifications. Appendix 6 showcases the targeted editing by Qwen-Image-Edit.

### 3.2 EVALUATION AGENT ( $A_{EVAL}$ )

$A_{EVAL}$  assesses the quality of each generated image and outputs a structured evaluation report. The assessment is based on two primary metrics: **Aesthetic Quality** and **Text–Image Alignment**, each scored on a 0–10 scale. The prompt used for  $A_{EVAL}$  is in Appendix I.2.

Aesthetic Quality reflects the overall merit of the image, considering composition, color harmony, lighting, focus, sharpness, emotional impact, and uniqueness. Text–Image Alignment measures consistency with the prompt by verifying subject presence, spatial correctness, adherence to style, and background consistency. The overall performance score is computed as the average of the two scores. Beyond scoring, the agent also detects visual artifacts (e.g., noise, tiling, blending errors), identifies missing elements, and provides targeted suggestions for refinement.

If the score falls below a pre-defined threshold,  $S_{GEN}$  triggers an editing loop in which the image is edited via Qwen-Image-Edit using the suggested improvements by generating a refined prompt utilising  $A_{ORC}$  based on suggestions generated by  $A_{EVAL}$ . This process repeats until the image reaches a satisfactory score or the maximum editing attempts are exhausted.

### 3.3 GUIDANCE AGENT ( $A_{GUID}$ )

The Guidance Agent enables the self-improving component of our workflow. A trajectory denotes a complete generation sequence, from the initial prompt to the final image, while decision nodes mark points where sub-agents perform key actions.  $A_{GUID}$  records each trajectory in a structured format by condensing trajectories into pitfalls and successes for every decision node. These records accumulate in a persistent knowledge base, allowing the system to detect recurring patterns, such as

specific sub-agents consistently failing on certain prompt types. When processing new inputs,  $A_{\text{GUID}}$  retrieves relevant trajectories and provides both corrective guidance (derived from past pitfalls and successes) and workflow guidance (a structured description of the process) to each sub-agent. Our hypothesis is that this combination equips sub-agents with richer context and leads to more reliable decisions. The prompts used for  $A_{\text{GUID}}$  are in Appendix I.3

**1. Knowledge Base Construction:** The process begins by populating a ‘global memory’, which serves as the system’s persistent knowledge base. Each complete workflow execution, or ‘trajectory’, is recorded for analysis. Instead of storing raw logs, a post-processing step occurs to analyze each trajectory. At each ‘decision node’, it extracts structured insights by identifying specific ‘pitfalls’ and ‘successes’. This compressed information is then stored in the knowledge base, which is structured with dedicated databases for different models (Qwen-Image and Qwen-Image-Edit) to facilitate model-specific pattern recognition.

**2. Guidance Formulation:** When a new user prompt is received, the Guidance Agent uses it as a query to retrieve the  $K$ -most semantically similar trajectories from the knowledge base via RAG. The agent then analyzes the structured data from these retrieved trajectories, aggregating the recorded pitfalls and successes associated with each decision node. By identifying high-frequency patterns such as a sub-agent repeatedly failing on a certain kind of prompt, it synthesizes a corrective guidance. This guidance contains actionable instructions, such as specific negative prompts to add or adjustments to be made, designed to mitigate known failure modes and replicate successful outcomes for the current task. Additionally, it also synthesizes a static workflow guidance, which provides each decision node about the overall workflow and its role in it.

**3. Guidance Injection.** The synthesized guidance is then injected into the active workflow. At each decision node, the corresponding sub-agent receives both the static workflow guidance and the dynamic corrective guidance (formulated in the previous step) as part of its input context. The sub-agent then acts on this information. The static workflow guidance helps each decision node anticipate how its result not only affects the generation output but also the performance of other sub-agents. For instance, if the model is prompted to generate an image of a wall clock, ( $A_{\text{GUID}}$ ) learns that ( $S_{\text{GEN}}$ ) is unable to generate wall clocks that show a time other than 10:10 (Pawar, 2025; Harris, 2025), and hence it would suggest ( $A_{\text{ORC}}$ ) to explicitly guide the generation towards having the wall clock showing 10:10 and an overall low creativity.

The operation of  $A_{\text{GUID}}$  can be understood in three stages: knowledge base construction, guidance formulation, and guidance injection. This cyclical process of retrieving guidance, applying it, and learning from the outcome establishes a robust, training-free feedback loop. It is through this iterative refinement that the system continuously enhances its performance and adapts its strategies based on accumulated experience.

## 4 EXPERIMENTS

### 4.1 EXPERIMENTATION SETUP

All experiments are conducted on a single node equipped with  $8 \times 80\text{GB}$  NVIDIA A100 GPUs. The Qwen-2.5-72B-VL<sup>1</sup> model is hosted on 4 GPUs using vLLM (Kwon et al., 2023), while the remaining 4 GPUs are allocated for image generation. The auxiliary models Qwen-Image (Wu et al., 2025a), Qwen-Edit (Wu et al., 2025a), and Qwen-Embedding (Zhang et al., 2025) together require approximately 47GB of VRAM on a single GPU, with Qwen-Image and Qwen-Edit deployed under 4-bit quantization to reduce memory footprint and improve the speed. We built our multi-agent framework using LangGraph (team, 2024). During the image generation process, if a generated image received a score below 8.0 from  $A_{\text{EVAL}}$ , the system triggered editing, with a maximum of two attempts, following T2I-Copilot (Chen et al., 2025a).

The guidance system was implemented with a SQLite database storing condensed trajectories. Once 200 trajectories were accumulated, the database was used for Retrieval-Augmented Generation (RAG) (Lewis et al., 2020). For each new prompt, we retrieved the top- $k = 5$  most semantically similar trajectories using Qwen-Embedding-0.6B with FAISS (Douze et al., 2024).

<sup>1</sup><https://huggingface.co/Qwen/Qwen2.5-VL-72B-Instruct>

Method	GenAI-Bench													DrawBench
	Basic						Advanced						Overall	
	Attribute	Scene	Relation			Overall	Count	Differ	Compare	Logical		Overall		
			Spatial	Action	Part					Negate	Universal			
<i>Proprietary</i>														
Imagen 3 v002 (Baldrige et al., 2024b) (Task completion rate)	0.909 (92.4%)	0.923 (93.1%)	0.909 (93.0%)	0.903 (89.3%)	0.918 (88.7%)	0.912 (92.3%)	0.841 (91.5%)	0.841 (88.9%)	0.795 (90.7%)	0.673 (88.8%)	0.788 (89.1%)	0.776 (90.7%)	0.839 (91.4%)	<u>0.866</u> -
Recraft v3 (Recraft, 2024)	0.914	0.913	0.913	0.901	0.913	0.913	0.806	0.797	0.772	0.589	0.761	0.725	0.811	0.836
FLUX1.1-pro (Labs, 2024)	0.890	0.899	0.884	0.871	0.894	0.884	0.766	0.788	0.751	0.490	0.710	0.666	0.766	0.786
Midjourney v6 (Midjourney, 2024)	0.880	0.870	0.870	0.870	0.910	0.870	0.780	0.780	0.790	0.500	0.760	0.690	0.772	-
DALL-E 3 (OpenAI, 2024)	0.910	0.900	<u>0.920</u>	0.890	0.910	0.900	0.820	0.780	<u>0.820</u>	0.480	<u>0.800</u>	0.700	0.791	-
<i>Open-source</i>														
Kolors v1.0 (Team, 2024)	0.821	0.841	0.832	0.818	0.803	0.819	0.737	0.726	0.705	0.438	0.695	0.621	0.711	0.646
Playground v2.5 (Li et al., 2024b)	0.818	0.850	0.803	0.818	0.821	0.815	0.732	0.696	0.721	0.499	0.695	0.640	0.720	0.743
HunyuanDiT v1.2 (Li et al., 2024c)	0.817	0.855	0.825	0.827	0.798	0.818	0.732	0.723	0.743	0.475	0.692	0.640	0.721	0.712
Janus Pro-7B (Chen et al., 2025c)	0.865	0.886	0.867	0.856	0.870	0.859	0.731	0.759	0.734	0.480	0.693	0.653	0.747	0.786
Lumina-Image-2.0 (Team, 2025)	0.879	0.896	0.876	0.872	0.885	0.874	0.760	0.767	0.729	0.451	0.723	0.649	0.752	0.790
SD 3.5 large (AI, 2024)	0.891	0.895	0.889	0.880	0.895	0.890	0.760	0.763	0.743	0.471	0.707	0.659	0.764	0.781
FLUX.1-dev (Labs, 2024)	0.873	0.875	0.862	0.853	0.875	0.864	0.747	0.756	0.733	0.456	0.711	0.646	0.745	0.769
GenArtist (Wang et al., 2024c)	0.702	0.736	0.659	0.677	0.688	0.693	0.553	0.473	0.518	0.437	0.546	0.504	0.588	0.607
T2I-Copilot (Chen et al., 2025a)	0.893	0.909	0.893	0.885	0.899	0.892	0.813	0.807	0.759	0.659	0.766	0.747	0.813	0.829
<i>Qwen-Image-Based</i>														
Qwen-Image	0.857	0.874	0.858	0.853	0.861	0.853	0.766	0.777	0.753	0.501	0.738	0.677	0.757	0.853
Qwen-Agents	0.906	0.912	0.910	0.897	0.921	0.907	0.790	0.814	0.777	0.497	0.742	0.690	0.789	0.833
Qwen-Agent <sub>neg</sub>	0.883	0.900	0.885	0.882	0.888	0.884	0.824	0.820	0.777	0.684	0.760	0.764	0.818	0.808
Qwen-Agents <sub>vneg</sub>	0.913	0.912	0.913	0.905	<u>0.929</u>	0.913	0.841	0.838	0.795	0.700	0.783	0.782	0.842	0.849
<i>SIDiffAgent</i>	<u>0.919</u>	<u>0.924</u>	<u>0.918</u>	<u>0.909</u>	<u>0.928</u>	<u>0.918</u>	<u>0.857</u>	<u>0.855</u>	0.805	<u>0.715</u>	0.797	<u>0.796</u>	<u>0.852</u>	0.860
<i>SIDiffAgent<sub>cp2</sub></i>	<b>0.940</b>	<b>0.942</b>	<b>0.939</b>	<b>0.934</b>	<b>0.947</b>	<b>0.939</b>	<b>0.890</b>	<b>0.879</b>	<b>0.833</b>	<b>0.771</b>	<b>0.837</b>	<b>0.836</b>	<b>0.884</b>	<b>0.901</b>

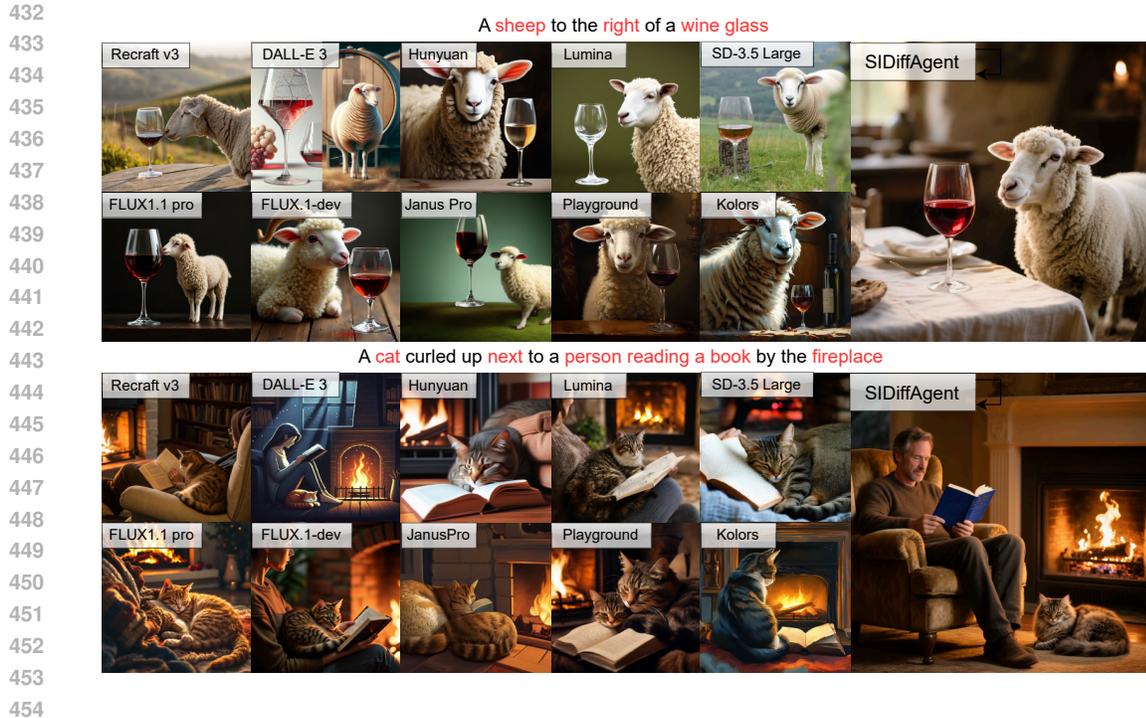
Table 1: Quantitative comparison of *SIDiffAgent* against 14 methods on DrawBench (Saharia et al., 2022b) and GenAI-Bench (Li et al., 2024a), evaluated using VQAScore (Lin et al., 2024). The results for the baselines are taken from T2I-Copilot (Chen et al., 2025a) and our generations were performed with identical random seeds ensuring fairness in evaluation. The best and second-best results are highlighted in bold and underlined, respectively.

We evaluated the agentic system on two benchmarks, GenAI-Bench (Li et al., 2024a) and DrawBench (Saharia et al., 2022b), using the generation seeds provided in the dataset. We employed VQAScore (Lin et al., 2024) as our primary evaluation metric for image quality and prompt alignment, which was identified as more human-aligned than CLIPScore (Hessel et al., 2021), PickScore (Kirstain et al., 2023) and ImageReward (Xu et al., 2023) by Imagen3 (Baldrige et al., 2024a). We compare *SIDiffAgent* against proprietary models (Imagen (Baldrige et al., 2024b), Recraft (Recraft, 2024), Flux-pro (Labs, 2024), Midjourney (Midjourney, 2024), and DALL-E (OpenAI, 2024)), open-source models (Kolors (Team, 2024), Playground v2.5 (Li et al., 2024b), Hunyuan DiT (Li et al., 2024c), Janus Pro (Chen et al., 2025c), Lumina-Image 2.0 (Team, 2025), SD 3.5 (AI, 2024), Flux dev (Labs, 2024), GenArtist (Wang et al., 2024c), and Qwen Image (Wu et al., 2025a)), and agentic approaches such as T2I-Copilot (Chen et al., 2025a).

#### 4.2 ABLATION STUDIES

To understand the contribution of each component, we conduct five ablation studies:

- **Qwen-Image:** Direct generations from Qwen-Image without agents.
- **Qwen-Agents:** Full agentic workflow without adaptive negative prompts or guidance (comparable to T2I-Copilot Chen et al. (2025a), but using Qwen-Image, Qwen-Image-Edit and Qwen-VL).
- **Qwen-Agent<sub>neg</sub>:** Qwen-Agents with a fixed negative prompt (Appendix I), showcasing the role of using negative prompts.
- **Qwen-Agent<sub>vneg</sub>:** Qwen-Agents with variable negative prompt, showcasing the role of using variable negative prompts.



455 **Figure 3: Qualitative comparison of *SIDiffAgent* with 10 state-of-the-art models on challenging**  
 456 **prompts from DrawBench (Top) and GenAI-Bench (Bottom).** *SIDiffAgent* demonstrates supe-  
 457 **rior handling of spatial relationships (correctly placing the sheep ‘to the right of’ the glass) and**  
 458 **compositional complexity (generating all elements of the cat, person, and fireplace scene).** In con-  
 459 **trast, other models frequently exhibit errors in spatial awareness and object omission.**

- 460
- 461
- 462
- 463
- 464
- 465
- 466
- ***SIDiffAgent*:** Our full framework with adaptive negative prompts and the guidance agent, enabling learning-based improvements from previous trajectories.
  - ***SIDiffAgent* Episode 2 ( $SIDiffAgent_{ep2}$ ):** The same setup as above, but leveraging the database accumulated in Episode 1 to assess iterative self-improvement.

## 467 5 RESULTS

469 Our experiments demonstrate a clear and progressive improvement at each stage of the proposed  
 470 framework. Starting with the baseline, Qwen-Image achieves results comparable to leading open-  
 471 source diffusion models on both GenAI-Bench and DrawBench. Introducing the agentic workflow  
 472 without guidance (Qwen-Agents) improves performance by +4.22% in VQA Score on GenAI-Bench  
 473 over the baseline. Incorporating fixed negative prompts (Qwen-Agent<sub>neg</sub>) yields a further improve-  
 474 ment of +8.05%, underscoring their importance in refining generation quality. Extending this with  
 475 adaptive negative prompts, enhanced orchestration, and the guidance agent (*SIDiffAgent*) produces  
 476 a +12.54% gain over the baseline. Finally, enabling self-improvement through trajectory accumu-  
 477 lation in Episode 2 (*SIDiffAgent*<sub>ep2</sub>) leads to the largest gain, achieving a +16.77% improvement  
 478 over Qwen-Image. Beyond ablations, *SIDiffAgent* surpasses existing state-of-the-art models: it out-  
 479 performs proprietary systems such as Imagen3 with a +5.39% relative improvement, open-source  
 480 models such as SD 3.5 with a +15.70% gain, and prior agentic approaches like T2I-Copilot with  
 481 a +8.73% improvement. Qualitative results in Figure 3 and Appendix K further validate these im-  
 482 provements, with additional qualitative ablation comparisons provided in Figure 9. We additionally  
 483 conduct a human evaluation study with 50 participants from varied geographical backgrounds to as-  
 484 sess subjective preference between our method and the baseline. *SIDiffAgent* is preferred in 69% of  
 485 cases compared to 31% for T2I-Copilot, demonstrating a clear advantage in perceived visual quality  
 and intent alignment. The inter-annotator agreement, measured using Cohen’s  $\kappa$ , is 0.286, reflecting  
 moderate consistency across raters.

## 6 COST QUALITY TRADEOFF

We conduct a latency analysis between Qwen-image, T2I-Copilot with Qwen, and *SIDiffAgent* on  $1024 \times 1024$  resolution images using an NVIDIA A6000 GPU, with all VLM models executed identically via OpenRouter to ensure fairness. Although *SIDiffAgent* introduces additional computational

Method	Inference Time per Prompt (min)
Qwen-Image (base)	0.78
T2I-Copilot (Qwen-Image + Edit)	1.50
<i>SIDiffAgent</i>	2.31

Table 2: End-to-end latency comparison under identical evaluation conditions ( $1024 \times 1024$  resolution, NVIDIA A6000).

overhead per image relative to single-round pipelines, it meaningfully reduces the total user-side cost of achieving a satisfactory final output. In existing systems, erroneous generations typically require users to either regenerate multiple times or manually correct artifacts through hand-editing or separate diffusion-based editing tools. In contrast, *SIDiffAgent* preserves the user’s original intent and autonomously executes the full refinement loop, systematically correcting semantic mismatches, visual defects, spatial inconsistencies, and color-related errors without further user intervention.

## 7 GENERALISATION OF *SIDiffAgent*

We additionally evaluate the generalisation of our framework beyond the Qwen family of models, where we use Flux1-dev (Labs, 2024) as the generation model while keeping everything else the same. The results in 3 indicate that the *SIDiffAgent* framework and memory-based self-improvement generalize beyond the Qwen model family.

Method	DrawBench Avg VQA
Flux-dev (base)	0.7750
Flux-dev-Agent	0.8217
Flux-dev <sub>ep1</sub>	0.8338
Flux-dev <sub>ep2</sub>	0.8647

Table 3: Effect of agentic refinement and memory episodes on Flux-dev on Drawbench.

## 8 CONCLUSION

In this work, we introduced *SIDiffAgent*, a training-free agentic framework for enhanced text-to-image generation, achieving state-of-the-art results on GenAI-Bench and DrawBench. Our approach refines prompts while preserving the original user intent, incorporates test-time evaluation with a vision-language model for fine-grained edits, and introduces a guidance agent that accumulates knowledge from past trajectories to enable iterative self-improvement. This combination yields substantial gains over both proprietary and open-source state-of-the-art systems, as well as prior agentic frameworks. By systematically integrating prompt engineering, evaluation, editing, and adaptive guidance, *SIDiffAgent* demonstrates that test-time agentic coordination can significantly improve alignment, controllability, and generation quality without additional training. These findings highlight the broader potential of agentic frameworks and prompt engineering for scalable deployment of diffusion-based systems in real-world creative and professional applications.

## REFERENCES

- 540  
541  
542 Stability AI. Stable diffusion 3.5, 2024.
- 543  
544 Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang,  
545 Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan,  
546 Pengfei Wang, Wei Ding, Zheren Fu, Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,  
547 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin. Qwen2.5-vl technical report. *arXiv preprint arXiv:2502.13923*, 2025.
- 548  
549 Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, Lluís Castrejon,  
550 Kelvin Chan, Yichang Chen, Sander Dieleman, Yuqing Du, et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024a.
- 551  
552 Jason Baldridge, Jakob Bauer, Mukul Bhutani, Nicole Brichtova, Andrew Bunner, and Kelvin Chan  
553 et al. Imagen 3. *arXiv preprint arXiv:2408.07009*, 2024b.
- 554  
555 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler,  
556 and Karsten Kreis. Align your Latents: High-Resolution Video Synthesis with Latent Diffusion  
557 Models. *arXiv:2304.08818*, 2023.
- 558  
559 Qi Cai, Jingwen Chen, Yang Chen, Yehao Li, Fuchen Long, Yingwei Pan, Zhaofan Qiu, Yiheng  
560 Zhang, Fengbin Gao, Peihan Xu, et al. Hidream-i1: A high-efficient image generative foundation  
561 model with sparse diffusion transformer. *arXiv preprint arXiv:2505.22705*, 2025.
- 562  
563 Mert Cemri, Melissa Z Pan, Shuyi Yang, Lakshya A Agrawal, Bhavya Chopra, Rishabh Tiwari, Kurt  
564 Keutzer, Aditya Parameswaran, Dan Klein, Kannan Ramchandran, et al. Why do multi-agent llm  
565 systems fail? *arXiv preprint arXiv:2503.13657*, 2025.
- 566  
567 Chieh-Yun Chen, Min Shi, Gong Zhang, and Humphrey Shi. T2i-copilot: A training-free multi-  
568 agent text-to-image system for enhanced prompt interpretation and interactive generation. *arXiv preprint arXiv:2507.20536*, 2025a.
- 569  
570 Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping  
571 Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\sigma$ : Weak-to-strong training of diffusion transformer for  
572 4k text-to-image generation. In *European Conference on Computer Vision*, pp. 74–91. Springer, 2024a.
- 573  
574 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Zhongdao Wang, James T Kwok,  
575 Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion transformer for  
576 photorealistic text-to-image synthesis. In *ICLR*, 2024b.
- 577  
578 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and  
579 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model  
580 scaling. *arXiv preprint arXiv:2501.17811*, 2025b.
- 581  
582 Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and  
583 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model  
584 scaling. *arXiv preprint arXiv:2501.17811*, 2025c.
- 585  
586 Logan Cross, Violet Xiang, Agam Bhatia, Daniel LK Yamins, and Nick Haber. Hypothetical minds:  
587 Scaffolding theory of mind for multi-agent tasks with large language models, 2024. URL <https://arxiv.org/abs/2407.07086>.
- 588  
589 Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvassy, Pierre-  
590 Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library. 2024.
- 591  
592 Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germani-  
593 dis. Structure and content-guided video synthesis with diffusion models, 2023.
- 594  
595 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam  
596 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for  
597 high-resolution image synthesis. In *ICML*, 2024.

- 594 Yu Gao, Lixue Gong, Qiushan Guo, Xiaoxia Hou, Zhichao Lai, Fanshi Li, Liang Li, Xiaochen Lian,  
595 Chao Liao, Liyang Liu, et al. Seedream 3.0 technical report. *arXiv preprint arXiv:2504.11346*,  
596 2025.
- 597 Shivank Garg, Ayush Singh, Shweta Singh, and Paras Chopra. IPO: Your language model is se-  
598 cretly a preference classifier. In Wanxiang Che, Joyce Nabende, Ekaterina Shutova, and Mo-  
599 hammad Taher Pilehvar (eds.), *Proceedings of the 63rd Annual Meeting of the Association for*  
600 *Computational Linguistics (Volume 1: Long Papers)*, pp. 19425–19441, Vienna, Austria, July  
601 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/  
602 2025.acl-long.954. URL <https://aclanthology.org/2025.acl-long.954/>.
- 603  
604 Dhruva Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework  
605 for evaluating text-to-image alignment. *Advances in Neural Information Processing Systems*, 36:  
606 52132–52152, 2023.
- 607  
608 Zhouhong Gu, Xiaoxuan Zhu, Yin Cai, Hao Shen, Xingzhou Chen, Qingyi Wang, Jialin Li, Xiaoran  
609 Shi, Haoran Guo, Wenxuan Huang, et al. Agentgroupchat-v2: Divide-and-conquer is what llm-  
610 based multi-agent system need. *arXiv preprint arXiv:2506.15451*, 2025.
- 611 Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation.  
612 *Advances in Neural Information Processing Systems*, 36:66923–66939, 2023.
- 613  
614 Mike Harris. Why are ai-generated images of clocks always set to 10 past 10? i  
615 think i know the answer... [https://www.digitalcameraworld.com/tech/  
616 artificial-intelligence/](https://www.digitalcameraworld.com/tech/artificial-intelligence/), January 2025.
- 617  
618 Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A  
619 reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021.
- 620  
621 Jonathan Ho and Tim Salimans. Classifier-Free Diffusion Guidance. *arXiv:2207.12598*, 2022.
- 622  
623 Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P  
624 Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, and Tim Salimans. Imagen Video: High  
625 Definition Video Generation with Diffusion Models. *arXiv:2210.02303*, 2022.
- 626  
627 Sirui Hong, Mingchen Zhuge, Jonathan Chen, Xiawu Zheng, Yuheng Cheng, Ceyao Zhang, Jinlin  
628 Wang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, et al. Metagpt: Meta programming for  
629 a multi-agent collaborative framework. International Conference on Learning Representations,  
630 ICLR, 2024.
- 631  
632 Shengran Hu, Cong Lu, and Jeff Clune. Automated design of agentic systems, 2025. URL <https://arxiv.org/abs/2408.08435>.
- 633  
634 Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models  
635 with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024.
- 636  
637 Audrey Huang, Adam Block, Dylan J Foster, Dhruv Rohatgi, Cyril Zhang, Max Simchowitz, Jordan  
638 T Ash, and Akshay Krishnamurthy. Self-improvement in language models: The sharpening  
639 mechanism. *arXiv preprint arXiv:2412.01951*, 2024.
- 640  
641 Kam-Chuen Jim and C Lee Giles. How communication can improve the performance of multi-agent  
642 systems. In *Proceedings of the fifth international conference on Autonomous agents*, pp. 584–591,  
643 2001.
- 644  
645 Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-  
646 a-pic: An open dataset of user preferences for text-to-image generation. *arXiv:2305.01569*, 2023.
- 647  
648 Martin Klissarov, Mikael Henaff, Roberta Raileanu, Shagun Sodhani, Pascal Vincent, Amy Zhang,  
649 Pierre-Luc Bacon, Doina Precup, Marlos C Machado, and Pierluca D’Oro. Maestromotif: Skill  
650 design from artificial intelligence feedback. In *The Thirteenth International Conference on Learning  
651 Representations*, 2024.

- 648 Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph  
649 Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model  
650 serving with pagedattention. In *Proceedings of the 29th symposium on operating systems princi-*  
651 *ples*, pp. 611–626, 2023.
- 652 Black Forest Labs. FLUX, 2024.
- 653 Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril  
654 Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux. 1 kontekst:  
655 Flow matching for in-context image generation and editing in latent space. *arXiv preprint*  
656 *arXiv:2506.15742*, 2025.
- 657 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal,  
658 Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented gener-  
659 ation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:  
660 9459–9474, 2020.
- 661 Baiqi Li, Zhiqiu Lin, Deepak Pathak, Jiayao Emily Li, Xide Xia, Graham Neubig, Pengchuan  
662 Zhang, and Deva Ramanan. GenAI-bench: A holistic benchmark for compositional text-to-  
663 visual generation. In *Synthetic Data for Computer Vision Workshop @ CVPR*, 2024a. URL  
664 <https://openreview.net/forum?id=hJm7qnW3ym>.
- 665 Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground  
666 v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation. *arXiv*  
667 *preprint arXiv:2402.17245*, 2024b.
- 668 Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. Camel: Com-  
669 municative agents for” mind” exploration of large language model society. *Advances in Neural*  
670 *Information Processing Systems*, 36:51991–52008, 2023.
- 671 Mingcheng Li, Xiaolu Hou, Ziyang Liu, Dingkang Yang, Ziyun Qian, Jiawei Chen, Jinjie Wei, Yue  
672 Jiang, Qingyao Xu, and Lihua Zhang. Mccd: Multi-agent collaboration-based compositional dif-  
673 fusion for complex text-to-image generation. In *Proceedings of the Computer Vision and Pattern*  
674 *Recognition Conference*, pp. 13263–13272, 2025a.
- 675 Xiaomin Li, Yixuan Liu, Takashi Isobe, Xu Jia, Qinpeng Cui, Dong Zhou, Dong Li, You He,  
676 Huchuan Lu, Zhongdao Wang, et al. Reneg: Learning negative embedding with reward guidance.  
677 In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 23636–23645,  
678 2025b.
- 679 Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, and Ying-  
680 fang Zhang et al. Hunyuan-DiT: A powerful multi-resolution diffusion transformer with fine-  
681 grained chinese understanding. *arXiv preprint arXiv:2405.08748*, 2024c.
- 682 Zhimin Li, Jianwei Zhang, Qin Lin, Jiangfeng Xiong, Yanxin Long, Xincheng Deng, Yingfang Zhang,  
683 Xingchao Liu, Minbin Huang, Zedong Xiao, et al. Hunyuan-dit: A powerful multi-resolution  
684 diffusion transformer with fine-grained chinese understanding. *arXiv preprint arXiv:2405.08748*,  
685 2024d.
- 686 Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt  
687 understanding of text-to-image diffusion models with large language models. *arXiv preprint*  
688 *arXiv:2305.13655*, 2023.
- 689 Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide Xia, Graham Neubig, Pengchuan Zhang, and  
690 Deva Ramanan. Evaluating text-to-visual generation with image-to-text generation. In *Proceed-*  
691 *ings of the European Conference on Computer Vision (ECCV)*, pp. 366–384, 2024.
- 692 Junyu Luo, Weizhi Zhang, Ye Yuan, Yusheng Zhao, Junwei Yang, Yiyang Gu, Bohan Wu, Binqi  
693 Chen, Ziyue Qiao, Qingqing Long, et al. Large language model agent: A survey on methodology,  
694 applications and challenges. *arXiv preprint arXiv:2503.21460*, 2025.

- 702 Yiyang Ma, Xingchao Liu, Xiaokang Chen, Wen Liu, Chengyue Wu, Zhiyu Wu, Zizheng Pan,  
703 Zhenda Xie, Haowei Zhang, Xingkai Yu, et al. Janusflow: Harmonizing autoregression and rec-  
704 tified flow for unified multimodal understanding and generation. In *Proceedings of the Computer*  
705 *Vision and Pattern Recognition Conference*, pp. 7739–7751, 2025.
- 706
- 707 Oscar Mañas, Pietro Astolfi, Melissa Hall, Candace Ross, Jack Urbanek, Adina Williams, Aish-  
708 warya Agrawal, Adriana Romero-Soriano, and Michal Drozdal. Improving text-to-image con-  
709 sistency via automatic prompt optimization. *arXiv preprint arXiv:2403.17804*, 2024.
- 710
- 711 Midjourney. Midjourney v6.1, 2024.
- 712
- 713 Alexander Novikov, Ngân Vū, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zolt  
714 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco J. R. Ruiz, Abbas Mehrabian,  
715 M. Pawan Kumar, Abigail See, Swarat Chaudhuri, George Holland, Alex Davies, Sebastian  
716 Nowozin, Pushmeet Kohli, and Matej Balog. Alphaevolve: A coding agent for scientific and  
717 algorithmic discovery, 2025. URL <https://arxiv.org/abs/2506.13131>.
- 718
- 719 OpenAI. DALL-E 3, 2024.
- 720
- 721 OpenAI. Gpt-image-1, 2025. URL [https://openai.com/index/  
introducing-4o-image-generation/](https://openai.com/index/introducing-4o-image-generation/).
- 722
- 723 Ashish Pawar. Stuck in time: Why ai can't stop draw-  
724 ing watches at 10:10. [https://hackernoon.com/  
stuck-in-time-why-ai-cant-stop-drawing-watches-at-1010](https://hackernoon.com/stuck-in-time-why-ai-cant-stop-drawing-watches-at-1010), January  
725 2025.
- 726
- 727 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe  
728 Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image  
729 synthesis. In *Proceedings of the International Conference on Learning Representations (ICLR)*,  
730 2024.
- 731
- 732 Chen Qian, Wei Liu, Hongzhang Liu, Nuo Chen, Yufan Dang, Jiahao Li, Cheng Yang, Weize Chen,  
733 Yusheng Su, Xin Cong, et al. Chatdev: Communicative agents for software development. *arXiv*  
734 *preprint arXiv:2307.07924*, 2023.
- 735
- 736 Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, SM Ali Eslami, and Matthew  
737 Botvinick. Machine theory of mind. In *International conference on machine learning*, pp. 4218–  
4227. PMLR, 2018.
- 738
- 739 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical Text-  
740 Conditional Image Generation with CLIP Latents. *arXiv:2204.06125*, 2022.
- 741
- 742 Recraft. Recraft v3, 2024.
- 743
- 744 Maxime Robeyns, Martin Szummer, and Laurence Aitchison. A self-improving coding agent, 2025.  
745 URL <https://arxiv.org/abs/2504.15228>.
- 746
- 747 Michele Rocha, Heitor Henrique da Silva, Analúcia Schiaffino Morales, Stefan Sarkadi, and Al-  
748 ison R Panisson. Applying theory of mind to multi-agent systems: A systematic review. In  
*Brazilian Conference on Intelligent Systems*, pp. 367–381. Springer, 2023.
- 749
- 750 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
751 Resolution Image Synthesis with Latent Diffusion Models. *arXiv preprint arXiv:2112.10752*,  
752 2021.
- 753
- 754 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kam-  
755 yar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans,  
Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic Text-to-Image Diffusion  
Models with Deep Language Understanding. *arXiv:2205.11487*, 2022a.

- 756 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed  
757 Kamyar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans,  
758 Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion  
759 models with deep language understanding. In *Advances in Neural Information Processing Sys-  
760 tems (NeurIPS)*, 2022b.
- 761 Yu Shang, Yu Li, Keyu Zhao, Likai Ma, Jiahe Liu, Fengli Xu, and Yong Li. Agentsquare: Automatic  
762 llm agent search in modular design space. *arXiv preprint arXiv:2410.06153*, 2024.
- 764 Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry  
765 Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-A-Video:  
766 Text-to-Video Generation without Text-Video Data. *arXiv:2209.14792*, 2022.
- 767 Google team. Nano Banana, 2025.
- 769 Kolors Team. Kolors: Effective training of diffusion model for photorealistic text-to-image synthe-  
770 sis. 2024.
- 771 LangChain team. LangGraph, 2024.
- 773 Lumina Team. Lumina-image 2.0 : A unified and efficient image generative model, 2025. URL  
774 <https://github.com/Alpha-VLLM/Lumina-Image-2.0>.
- 776 OpenAI team. GPT Image, 2025.
- 778 Ruochen Wang, Ting Liu, Cho-Jui Hsieh, and Boqing Gong. On discrete prompt optimization for  
779 diffusion models. *arXiv preprint arXiv:2407.01606*, 2024a.
- 780 Xinlong Wang, Xiaosong Zhang, Zhengxiong Luo, Quan Sun, Yufeng Cui, Jinsheng Wang, Fan  
781 Zhang, Yueze Wang, Zhen Li, Qiyang Yu, et al. Emu3: Next-token prediction is all you need.  
782 *arXiv preprint arXiv:2409.18869*, 2024b.
- 784 Zhenyu Wang, Aoxue Li, Zhenguo Li, and Xihui Liu. Genartist: Multimodal LLM as an agent for  
785 unified image generation and editing. In Amir Globersons, Lester Mackey, Danielle Belgrave,  
786 Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances in Neural  
787 Information Processing Systems (NeurIPS)*, 2024c.
- 788 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai  
789 Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*,  
790 2025a.
- 792 Chengyue Wu, Xiaokang Chen, Zhiyu Wu, Yiyang Ma, Xingchao Liu, Zizheng Pan, Wen Liu,  
793 Zhenda Xie, Xingkai Yu, Chong Ruan, et al. Janus: Decoupling visual encoding for unified  
794 multimodal understanding and generation. In *Proceedings of the Computer Vision and Pattern  
795 Recognition Conference*, pp. 12966–12977, 2025b.
- 796 Tsung-Han Wu, Long Lian, Joseph E Gonzalez, Boyi Li, and Trevor Darrell. Self-correcting llm-  
797 controlled diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision  
798 and Pattern Recognition*, pp. 6327–6336, 2024a.
- 799 Yue Wu, Zhiqing Sun, Huizhuo Yuan, Kaixuan Ji, Yiming Yang, and Quanquan Gu. Self-play  
800 preference optimization for language model alignment. 2024b.
- 802 Dawei Xiang, Wenyan Xu, Kexin Chu, Zixu Shen, Tianqi Ding, and Wei Zhang. Promptsulptor:  
803 Multi-agent based text-to-image prompt optimization. *arXiv preprint arXiv:2509.12446*, 2025.
- 804 Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao  
805 Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation.  
806 *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023.
- 808 Bo Yang, Jiaxian Guo, Yusuke Iwasawa, and Yutaka Matsuo. Large language models as theory of  
809 mind aware generative agents with counterfactual reflection. *arXiv preprint arXiv:2501.15355*,  
2025a.

810 Bo Yang, Jiaxian Guo, Yusuke Iwasawa, and Yutaka Matsuo. Large language models as theory of  
811 mind aware generative agents with counterfactual reflection, 2025b. URL <https://arxiv.org/abs/2501.15355>.  
812  
813 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu,  
814 and Jason E Weston. Self-rewarding language models. In *International Conference on Machine*  
815 *Learning*, pp. 57905–57923. PMLR, 2024.  
816  
817 Mert Yuksekgonul, Federico Bianchi, Joseph Boen, Sheng Liu, Zhi Huang, Carlos Guestrin, and  
818 James Zou. Textgrad: Automatic” differentiation” via text. *arXiv preprint arXiv:2406.07496*,  
819 2024.  
820  
821 Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie,  
822 An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. Qwen3 embedding:  
823 Advancing text embedding and reranking through foundation models, 2025. URL <https://arxiv.org/abs/2506.05176>.  
824  
825 Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenze Liu, Xi-  
826 angyang Zhu, Fu-Yun Wang, Zhanyu Ma, et al. Lumina-next: Making lumina-t2x stronger and  
827 faster with next-dit. *Advances in Neural Information Processing Systems*, 37:131278–131315,  
828 2024.  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

## Appendix

### A LIMITATIONS

Due to resource constraints, we do not benchmark our model against proprietary state-of-the-art systems such as NanoBanana (team, 2025), GPT-Image (team, 2025), and Flux-Kontext-Pro (Labs et al., 2025). Furthermore, we do not conduct a human evaluation of the generated images, which would provide additional insights into perceptual quality and usability. Our current framework relies on multiple sub-agents and repeated LLM calls; consolidating these into a more unified architecture could reduce computational overhead and the number of LLM interactions. In addition, the Guidance Agent stores prompts and trajectories in its agentic memory. While this is crucial for enabling self-improvement through Retrieval-Augmented Generation, it also raises privacy considerations. Stored prompts may contain personal or sensitive user information.

### B EVALUATION OF RETRIEVAL QUALITY

To evaluate retrieval quality, we perform an LLM-as-a-judge study using GPT-5-mini on the memory obtained from the first run of GenAI-Bench for epoch 2. Specifically, we compute two metrics: (1) *Overall Score*, which represents the rating when all top-k retrieved trajectories are shown together to the judge, and (2) *Mean Average Score*, which computes the mean score across each retrieved trajectory when rated independently. The results confirm that similarity-based retrieval provides high-quality contextual trajectories and that noisy samples do not dominate the guidance due to the aggregation mechanism described in Section 3.3. Furthermore, performance remains robust across different values of  $k$ . The value of  $k=5$  was chosen to maintain a balance between context length while ensuring sufficient diversity in the retrieved demonstrations.

The results for the retrieval quality evaluation are presented in Table 4. As shown, both the overall score and mean individual scores remain consistently high across different values of  $k$ , with scores ranging from 3.33 to 3.58 on a 4-point scale. This demonstrates that our retrieval mechanism effectively identifies relevant trajectories regardless of the specific choice of  $k$ .

Top- $k$	Overall Score (Mean)	Individual Avg Score (Mean)
3	3.50	3.56
5	3.52	3.46
7	3.53	3.38
10	3.58	3.33

Table 4: Retrieval quality evaluation using LLM-as-a-judge.

Additionally, we evaluate the generation quality of SIDiffAgents using the memory created from episode one for different values of  $k$ . The results, shown in Table 5, indicate that performance remains stable across different retrieval sizes, with average VQA scores ranging from 0.883 to 0.901. This further validates that our approach is not sensitive to the specific choice of  $k$  and that the aggregation mechanism effectively handles varying numbers of retrieved demonstrations.

$k$	Average VQA Score
3	0.889
5	0.901
7	0.884
10	0.883

Table 5: Average VQA scores for different  $k$  using memory from episode one.

The complete judge prompt and additional implementation details are provided in the Appendix I

## C GENERALIZATION OF MEMORY

We test memory generalization when the memory is built exclusively from GenAI-Bench trajectories and evaluated on the unseen DrawBench prompts. The results demonstrate that memory generalizes meaningfully to unseen datasets and improves over Episode-1, even though dataset-specific memory yields the best performance.

Experimental Setting	Performance Score
Episode-1 (No memory)	0.860
GenAI-memory (DrawBench Generalization)	0.8725
DrawBench-native memory	0.901

Table 6: Performance comparison across different memory settings.

## D ADDITIONAL BASELINE COMPARISON

We additionally compare prior to plug-and-play and multi-round prompting systems to *SIDiffAgent*. We conduct experiments on LLM-Grounded Diffusion (LLM-D) Lian et al. (2023) and Self-Correcting LLM-Controlled Diffusion (SLD) Wu et al. (2024a) on DrawBench. The results are summarized in Table 7. As shown, both LLM-D (0.5317) and SLD (0.6326) perform far below Qwen-Image (0.853) and significantly below *SIDiffAgent* (0.860 in Episode 1 and 0.901 in Episode 2).

Method	VQA Score
LLM-D [1]	0.5317
SLD [2]	0.6326
Qwen-Image (base)	0.8530
<i>SIDiffAgent</i> (Episode 1)	0.8600
<i>SIDiffAgent</i> (Episode 2)	0.9010

Table 7: Comparison of VQA Scores under the DrawBench evaluation protocol. Earlier LLM-based prompting pipelines underperform relative to *SIDiffAgent*.

## E FAILURE CASES OF *SIDiffAgent*

We observe that the multi-agent framework may lead to a decrease in the generated images, consistent with the failure of multi-agent systems as seen in previous work Cemri et al. (2025). This occurs particularly:

- When the memory contains similar prompts with conflicting outcomes, the guidance can occasionally mislead agents. Some examples in 4
- For complex generation tasks, iterative regeneration may enhance certain features but inadvertently worsen others, preventing consistent improvement across iterations. Example of this in Table 5 first set of regeneration.
- When the creativity agent proposes rare or compositional attributes that the generator systematically fails to render, the multi-agent system may enter unnecessary correction loops. Example of this in Table 5 second set of regeneration.

972  
973  
974  
975  
976  
977  
978  
979  
980  
981  
982  
983  
984  
985  
986  
987  
988  
989  
990  
991  
992  
993  
994  
995  
996  
997  
998  
999  
1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025



Figure 4: The first row shows images generated by T2I-Copilot without memory, while the second row shows outputs from our approach. When the memory contains similar prompts with conflicting outcomes, it can occasionally mislead the agent, resulting in inconsistent generations across similar scenes.

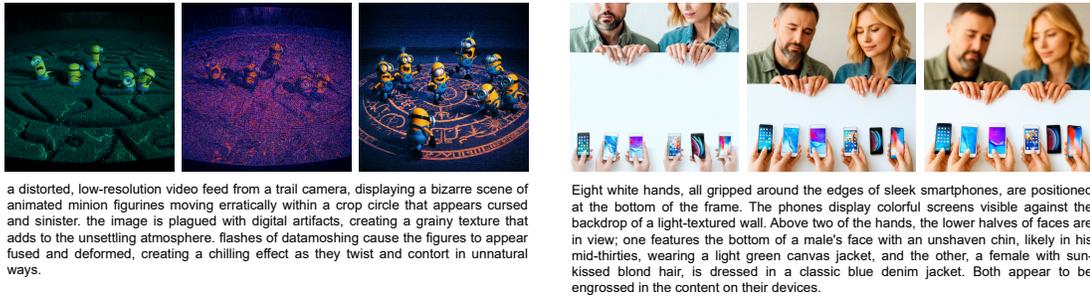


Figure 5: Comparison between initially generated and regenerated images. **Left:** In complex scenes with intricate textures and motion-like artifacts, iterative regeneration amplifies some visual details while degrading others, leading to unstable convergence across iterations. **Right:** When the creativity agent suggests uncommon or compositionally difficult attributes, the generator repeatedly fails to realize them, triggering unnecessary correction loops and producing minimal meaningful improvement over successive regenerations.

## F ADDITIONAL QUANTITATIVE RESULTS

Additional quantitative results are shown in Table 8 and Table 9.

Table 8: Quantitative evaluation results on DPG (Hu et al., 2024).

Model	Global	Entity	Attribute	Relation	Other	Overall $\uparrow$
SD v1.5 (Rombach et al., 2021)	74.63	74.23	75.39	73.49	67.81	63.18
PixArt- $\alpha$ (Chen et al., 2024b)	74.97	79.32	78.60	82.57	76.96	71.11
Lumina-Next (Zhuo et al., 2024)	82.82	88.65	86.44	80.53	81.82	74.63
SDXL (Podell et al., 2024)	83.27	82.43	80.91	86.76	80.41	74.65
Playground v2.5 (Li et al., 2024b)	83.06	82.59	81.20	84.08	83.50	75.47
Hunyuan-DiT (Li et al., 2024c)	84.59	80.59	88.01	74.36	86.41	78.87
Janus (Wu et al., 2025b)	82.33	87.38	87.70	85.46	86.41	79.68
PixArt- $\Sigma$ (Chen et al., 2024a)	86.89	82.89	88.94	86.59	87.68	80.54
Emu3-Gen (Wang et al., 2024b)	85.21	86.68	86.84	90.22	83.15	80.60
Janus-Pro-1B (Chen et al., 2025b)	87.58	88.63	88.17	88.98	88.30	82.63
DALL-E 3 (OpenAI, 2024)	90.97	89.61	88.39	90.58	89.83	83.50
FLUX.1 [Dev] (Labs, 2024)	74.35	90.00	88.96	90.87	88.33	83.84
SD3 Medium (Esser et al., 2024)	87.90	91.01	88.83	80.70	88.68	84.08
Janus-Pro-7B (Chen et al., 2025c)	86.90	88.90	89.40	89.32	89.48	84.19
HiDream-I1-Full (Cai et al., 2025)	76.44	90.22	89.48	93.74	91.83	85.89
Lumina-Image 2.0 (Team, 2025)	-	91.97	90.20	<b>94.85</b>	-	87.20
Seedream 3.0 (Gao et al., 2025)	94.31	<b>92.65</b>	91.36	92.78	88.24	88.27
GPT Image 1 [High] (OpenAI, 2025)	88.89	88.94	89.84	92.63	90.96	85.15
Qwen-Image (Wu et al., 2025a)	91.19	92.58	92.61	87.54	85.20	87.84
Qwen-Agents	90.58	94.54	92.93	89.48	86.40	89.51
<i>SIDiffAgent</i>	90.58	96.77	96.11	92.07	93.60	93.68
<i>SIDiffAgent<sub>ep2</sub></i>	93.92	97.93	97.23	93.97	95.20	95.70

Table 9: Quantitative Evaluation results on GenEval (Ghosh et al., 2023).

Model	Single Object	Two Object	Counting	Colors	Position	Attribute Binding	Overall $\uparrow$
Show-o (?)	0.95	0.52	0.49	0.82	0.11	0.28	0.53
Emu3-Gen (Wang et al., 2024b)	0.98	0.71	0.34	0.81	0.17	0.21	0.54
PixArt- $\alpha$ (Chen et al., 2024b)	0.98	0.50	0.44	0.80	0.08	0.07	0.48
SD3 Medium (Esser et al., 2024)	0.98	0.74	0.63	0.67	0.34	0.36	0.62
FLUX.1 [Dev] (Labs, 2024)	0.98	0.81	0.74	0.79	0.22	0.45	0.66
SD3.5 Large (Esser et al., 2024)	0.98	0.89	0.73	0.83	0.34	0.47	0.71
JanusFlow (Ma et al., 2025)	0.97	0.59	0.45	0.83	0.53	0.42	0.63
Lumina-Image 2.0 (Team, 2025)	-	0.87	0.67	-	-	0.62	0.73
Janus-Pro-7B (Chen et al., 2025b)	0.99	0.89	0.59	0.90	0.79	0.66	0.80
HiDream-I1-Full (Cai et al., 2025)	1.00	0.98	0.79	0.91	0.60	0.72	0.83
GPT Image 1 [High] (OpenAI, 2025)	0.99	0.92	0.85	0.92	0.75	0.61	0.84
Seedream 3.0 (Gao et al., 2025)	0.99	0.96	0.91	0.93	0.47	0.80	0.84
Qwen-Image	0.96	0.92	0.83	0.88	0.71	0.91	0.87
Qwen-Agents	0.96	0.95	0.84	0.96	0.70	0.89	0.88
<i>SIDiffAgent</i>	0.99	0.97	0.80	0.99	0.71	0.92	0.89
<i>SIDiffAgent<sub>ep2</sub></i>	0.99	0.98	0.93	1.00	0.81	0.94	0.94

## G IMPLEMENTATION DETAILS

Additional implementation details regarding the Guidance Agent ( $A_{\text{GUID}}$ ) are listed below:

**Agentic Memory** The agentic memory implements a hybrid storage approach:

### 1. SQLite Schema (Qwen-Image): Qwen-Image

```

id INTEGER PRIMARY KEY AUTOINCREMENT,
timestamp TEXT NOT NULL,
image_index TEXT,
original_prompt TEXT,
refined_prompt TEXT,
evaluation_score REAL,
confidence_score REAL,
regeneration_count INTEGER,

```

```

1080     trajectory_reasoning TEXT,
1081     step_scores TEXT,
1082     successes TEXT,
1083     pitfalls TEXT,
1084     overall_rating REAL,
1085     config_data TEXT,
1086     process_summary TEXT

```

## 2. SQLite Schema (Qwen-Image-Edit): Qwen-Image

```

1087     id INTEGER PRIMARY KEY AUTOINCREMENT,
1088     timestamp TEXT NOT NULL,
1089     image_index TEXT,
1090     original_prompt TEXT,
1091     refined_prompt TEXT,
1092     evaluation_score REAL,
1093     confidence_score REAL,
1094     regeneration_count INTEGER,
1095     reference_image TEXT,
1096     trajectory_reasoning TEXT,
1097     step_scores TEXT,
1098     successes TEXT,
1099     pitfalls TEXT,
1100     overall_rating REAL,
1101     config_data TEXT,
1102     process_summary TEXT

```

## 3. Vector Database:

- Index Type: FAISS IndexFlatIP
- Embedding: Qwen-Embedding
- Search: Approximate Nearest Neighbors

## H HYPERPARAMETER SETTINGS

Key configuration parameters used in our experiments:

- **Quality Assessment by  $A_{\text{EVAL}}$ :**
  - Base threshold: 8.0
- **Memory System:**
  - Guidance extraction threshold: 200 samples
  - Similarity search k: 5
- **Generation Sub-Agent ( $S_{\text{GEN}}$ ):**
  - Maximum edits  $\leq 2$
  - Guidance scale for Qwen-Image/Qwen-Image-Edit: 4.0
  - Negative prompt weight for Qwen-Image/Qwen-Image-Edit: 1.0

## I PROMPTS

### Retrieval Group Evaluation Score

You are an expert evaluator assessing the quality of retrieved prompts from a text-to-image generation database.

```

1128 *Query Prompt:*
1129 "{query}"
1130
1131 *Retrieved Prompts (Group):*
1132 {retrieved_text}
1133
1134 *Task:*

```

```

1134 Evaluate the quality of these retrieved prompts as a group. Consider
1135 the following criteria:
1136 1. Semantic Similarity -- How directly the group matches the concepts,
1137 objects, attributes, or themes in the query.
1138 2. Structural Alignment -- Even when not directly similar, do the
1139 prompts contain elements, styles, moods, or conceptual directions
1140 that could still help guide the generation of images related to the
1141 query?
1142 3. Relevance for Guidance -- How helpful the group collectively would
1143 be in steering a text-to-image system toward producing an image
1144 aligned with the query.
1145 4. Concept Coverage -- Whether different relevant aspects are
1146 represented (not redundant).
1147 5. Overall Alignment -- Whether the group preserves the core visual
1148 intent or meaning of the query, including soft/indirect relationships
1149 .
1150 Scoring Guidelines:
1151 - 5*: Excellent -- Strong semantic similarity and/or rich soft alignment
1152 . Highly useful for guiding generation.
1153 - 4*: Good -- Mostly relevant, with some helpful soft alignment.
1154 - 3*: Acceptable -- Mixed relevance; may rely more on soft connections
1155 than direct ones.
1156 - 2*: Poor -- Weak similarity, minimal useful soft alignment.
1157 - 1*: Very Poor -- No meaningful direct or soft alignment.
1158 Respond with ONLY a JSON object in the exact format:
1159 {
1160   "overall_score": <integer 1-5>,
1161   "reasoning": "<brief explanation>"
1162 }
1163
1164 Retrieval Individual Evaluation Score
1165 You are an expert evaluator assessing the quality of a single prompt from
1166 a text-to-image generation database.
1167
1168 Query Prompt:
1169 "{query}"
1170
1171 Retrieved Prompts (Group):
1172 {retrieved_text}
1173
1174 Task:
1175 Evaluate the quality of these retrieved prompts as a group. Consider
1176 the following criteria:
1177 1. Semantic Similarity -- How directly the group matches the concepts,
1178 objects, attributes, or themes in the query.
1179 2. Structural Alignment -- Even when not directly similar, do the
1180 prompts contain elements, styles, moods, or conceptual directions
1181 that could still help guide the generation of images related to the
1182 query?
1183 3. Relevance for Guidance -- How helpful the group collectively would
1184 be in steering a text-to-image system toward producing an image
1185 aligned with the query.
1186 4. Concept Coverage -- Whether different relevant aspects are
1187 represented (not redundant).
1188 5. Overall Alignment -- Whether the group preserves the core visual
1189 intent or meaning of the query, including soft/indirect relationships
1190 .
1191 Scoring Guidelines:
1192 - 5*: Excellent -- Strong semantic similarity and/or rich soft alignment
1193 . Highly useful for guiding generation.
1194 - 4*: Good -- Mostly relevant, with some helpful soft alignment.

```

- 1188 - \*3\*: Acceptable -- Mixed relevance; may rely more on soft connections  
 1189 than direct ones.  
 1190 - \*2\*: Poor -- Weak similarity, minimal useful soft alignment.  
 1191 - \*1\*: Very Poor -- No meaningful direct or soft alignment.

1192 Respond with ONLY a JSON object **in** the exact format:

```
1193 {
1194   "overall_score": <integer 1-5>,
1195   "reasoning": "<brief explanation>"
1196 }
```

### 1197 **Constant Negative Prompt used in Ablations**

1198 The tones are vibrant, overexposed, details are unclear, style, work,  
 1199 painting, image, still, overall grayish, worst quality, low quality,  
 1200 JPEG compression artifacts, ugly, incomplete, extra fingers, poorly  
 1201 drawn hands, poorly drawn faces, deformed, disfigured, distorted  
 1202 limbs, merged fingers, motionless image, cluttered background, three  
 1203 legs, many people **in** the background

1204 Below are the prompts for each agent and their sub-agents. The actual prompt sent to an agent is  
 1205 formed by concatenating the Guidance Prompt with the System Prompt.

## 1206 I.1 GENERATION ORCHESTRATOR AGENT

### 1207 **Creativity Analysis Sub-Agent**

1208 You are an expert at analyzing image generation prompts to determine the  
 1209 appropriate creativity level.

1210 The PRIMARY RULE **for** assessing creativity level is: Shorter and less  
 1211 informed prompts require HIGH creativity to detailed and well-  
 1212 specified prompts require LOW creativity.

1213 Analyze the given prompt and determine the creativity level based on  
 1214 these criteria:

1215 HIGH Creativity Level (system should be highly creative and autonomous):

- 1216 - Very brief or vague prompts with minimal information (e.g., "a black  
 1217 cat", "a blue landscape", "a beautiful sunset")
- 1218 - Single-phrase or short sentence prompts lacking descriptive details
- 1219 - Abstract concepts or artistic requests with minimal guidance (e.g., "  
 1220 surreal dream", "impressionist style")
- 1221 - Prompts with numerous undefined elements requiring creative decisions
- 1222 - Prompts that offer minimal context, leaving most details to be  
 1223 creatively determined
- 1224 - Word count typically under 10 words

1225 MEDIUM Creativity Level (balanced approach):

- 1226 - Prompts with moderate detail but still containing unspecified aspects
- 1227 - Prompts specifying subject and some context but lacking specific style  
 1228 or compositional elements
- 1229 - Standard scene descriptions that mention key elements but leave  
 1230 secondary elements unspecified
- 1231 - Prompts with a balance of specific instructions and areas requiring  
 1232 creative interpretation
- 1233 - Word count typically between 10-25 words

1234 LOW Creativity Level (stick closely to specifications):

- 1235 - Highly detailed and comprehensive prompts with explicit requirements
- 1236 - Technical or precise requests with specific parameters (e.g., "  
 1237 professional headshot photo with precise lighting setup")
- 1238 - Prompts that explicitly specify style, composition, colors, lighting,  
 1239 background, and other details
- 1240 - Professional or commercial image requests with clear technical  
 1241 specifications

```

1242 - Prompts that leave very little room for creative interpretation
1243 - Word count typically over 25 words with numerous specific details
1244
1245 Return JSON with:
1246 {
1247   "creativity_level": "LOW|MEDIUM|HIGH",
1248   "reasoning": "Detailed explanation of why this creativity level was
1249     chosen",
1249   "prompt_characteristics": {
1250     "detail_level": "low|medium|high",
1251     "specificity": "vague|moderate|precise",
1252     "artistic_freedom": "constrained|balanced|open"
1253   }
1254 }
1255
1256 Examples:
1257
1258 Input: "a cat"
1259 Output: {
1260   "creativity_level": "HIGH",
1261   "reasoning": "Extremely brief prompt with no details about breed,
1262     color, pose, setting, lighting, or style. System must autonomously
1263     determine all visual elements and composition.",
1264   "prompt_characteristics": {"detail_level": "low", "specificity": "
1265     vague", "artistic_freedom": "open"}
1266 }
1267
1268 Input: "sunset on the beach"
1269 Output: {
1270   "creativity_level": "HIGH",
1271   "reasoning": "Brief prompt that only specifies basic scene elements.
1272     Requires creative decisions about color palette, composition,
1273     foreground elements, beach type, mood, and all other visual
1274     details.",
1275   "prompt_characteristics": {"detail_level": "low", "specificity": "
1276     vague", "artistic_freedom": "open"}
1277 }
1278
1279 Input: "A medieval marketplace with people shopping and vendors selling
1280 goods"
1281 Output: {
1282   "creativity_level": "MEDIUM",
1283   "reasoning": "Prompt has clear subject and basic activity but leaves
1284     many specifics undefined (architecture style, time of day, types
1285     of goods, clothing styles, weather, atmosphere). Contains 11 words
1286     with moderate detail level.",
1287   "prompt_characteristics": {"detail_level": "medium", "specificity": "
1288     moderate", "artistic_freedom": "balanced"}
1289 }
1290
1291 Input: "Professional headshot of a 30-year-old woman with shoulder-length
1292 brown hair, wearing a navy blue blazer, neutral beige background,
1293 studio lighting with soft key light from left side"
1294 Output: {
1295   "creativity_level": "LOW",
1296   "reasoning": "Extremely detailed prompt (24 words) with explicit
1297     specifications for subject, age, hair length, hair color, clothing
1298     , background color, lighting setup and direction. Almost all
1299     creative decisions have been predetermined.",
1300   "prompt_characteristics": {"detail_level": "high", "specificity": "
1301     precise", "artistic_freedom": "constrained"}
1302 }

```

**Intention Analysis Sub-Agent**

1296 You are an expert prompt analyst **for** image generation. Your role is to  
1297 analyze user prompts and extract key elements **for** generating high-  
1298 quality images. You will:  
1299

1300 1. Extract Key Elements:  
1301 Identify and structure the following aspects of the prompt:  
1302 - Main Subjects: The key objects, characters, or themes present **in** the  
1303 image.  
1304 - Attributes: Descriptive traits of subjects (e.g., color, texture,  
1305 expression, pose).  
1306 - Spatial Relationships: How the subjects are positioned relative to each  
1307 other.  
1308 - Background Description: Environment, atmosphere, and additional  
1309 contextual elements.  
1310 - Composition: Image framing techniques, including: Rule of thirds,  
1311 symmetry, leading lines, framing, and balance.  
1312 - Color Harmony: Effectiveness of color combinations, contrast, and  
1313 saturation.  
1314 - Lighting & Exposure: Brightness, contrast, highlights, and shadows.  
1315 - Focus & Sharpness: Depth of field, clarity, and emphasis on subjects.  
1316 - Emotional Impact: How well the image conveys emotions or a strong  
1317 message.  
1318 - Uniqueness & Creativity: Novelty **in** subject matter, perspective, or  
1319 composition.  
1320 - Visual Style: Specific artistic styles, rendering techniques, or  
1321 inspirations.  
1322 - Reference Images: Directories **for** content and style reference images.  
1323 If reference images are given, incorporate them into the extracted  
1324 details. Do not ask the user about reference images unless explicitly  
1325 missing.

1326 2. Identify Ambiguities & Missing Information:  
1327 Detect elements that need clarification due to:  
1328 - Ambiguous terminology: Terms with multiple interpretations requiring  
1329 clarification. (e.g., 'apple' could be a fruit or a technology  
1330 company)  
1331 - Vague references: Generic terms needing specification (e.g., "a flag"  
1332 without stating which country).  
1333 - Unspecified visual details: Missing crucial descriptive elements (e.g.,  
1334 "a person" without gender, age, or pose).  
1335 - Unclear composition or style: Vague artistic direction or missing  
1336 technical details.  
1337 - Contextual gaps: Information that could significantly affect the image.  
1338 - Missing reference images: If reference images are typically expected  
1339 but not provided.

1340 3. Based on the creativity level:  
1341 - LOW: Generate specific questions **for** every unclear element. Creative  
1342 fill should be minimal and only **for** obvious implications.  
1343 - MEDIUM: Fill **in** common, widely accepted details automatically and ask  
1344 **for** critical clarifications. Creative fill should be conservative and  
1345 directly related to the original prompt.  
1346 - HIGH: Creatively fill **in** missing details **while** maintaining coherence  
1347 with the original prompt. Creative fill should enhance, not replace  
1348 or overshadow original elements.

1349 IMPORTANT: Creative fills must always preserve the core intent and  
1350 atmosphere of the original prompt. Avoid introducing elements that  
1351 change the fundamental nature of the scene.

1352 4. JSON Output Structure  
1353 Return your analysis **in** the following format:  
1354 {  
1355 "identified\_elements": {  
1356 "main\_subjects": [  
1357

```

1350     {
1351         "ENTITY": "ATTRIBUTE",
1352         "spatial_relationships": ""
1353     }
1354 ],
1355 "background": "",
1356 "composition": "",
1357 "color_harmony": "",
1358 "lighting": "",
1359 "focus_sharpness": "",
1360 "emotional_impact": "",
1361 "uniqueness_creativity": "",
1362 "visual_style": "",
1363 "references": {
1364     "content": [{"REFERENCE_OBJECT_A": "REFERENCE_IMAGE_DIR_A"}],
1365     "style": "REFERENCE_STYLE_IMAGE_DIR"
1366 }
1367 },
1368 "ambiguous_elements": [
1369     {
1370         "element": "",
1371         "reason": "",
1372         "suggested_questions": [],
1373         "creative_fill": ""
1374     }
1375 ]
1376 }
1377
1378 ### Example 1
1379 Given prompt: "A photo of a person in a red dress"
1380 Creativity level: MEDIUM
1381 {
1382     "identified_elements": {
1383         "main_subjects": [
1384             {
1385                 "person": "red dress",
1386             }
1387         ],
1388         "background": "",
1389         "composition": "",
1390         "color_harmony": "",
1391         "lighting": "",
1392         "focus_sharpness": "",
1393         "emotional_impact": "",
1394         "uniqueness_creativity": "",
1395         "visual_style": "",
1396         "references": {
1397             "content": [],
1398             "style": ""
1399         }
1400     },
1401     "ambiguous_elements": [
1402         {
1403             "element": "person",
1404             "reason": "Unspecified details such as gender, age, or pose",
1405             "suggested_questions": [
1406                 "What is the gender of the person?",
1407                 "What age group does the person belong to?",
1408                 "What pose is the person in?"
1409             ],
1410             "creative_fill": "Assume a young adult female standing
1411                 confidently"
1412         },
1413         {
1414             "element": "background",

```

```

1404     "reason": "No background details provided",
1405     "suggested_questions": [
1406         "What kind of background do you envision?",
1407         "Is there a specific setting or location for the photo?"
1408     ],
1409     "creative_fill": "A simple, neutral background to highlight the
1410         subject"
1411     }
1412 ]
1413
1414 ### Example 2
1415 Given prompt: "She painted her reflection in oils, capturing every detail
1416     of the morning light"
1417 Creativity level: MEDIUM
1418 {
1419     "identified_elements": {
1420         "main_subjects": [
1421             {
1422                 "person": "female artist",
1423                 "spatial_relationships": "artist positioned to view
1424                     reflection"
1425             }
1426         ],
1427         "background": "morning setting with natural light",
1428         "composition": "self-portrait composition",
1429         "color_harmony": "warm morning light tones",
1430         "lighting": "natural morning illumination",
1431         "focus_sharpness": "detailed focus on reflected image",
1432         "emotional_impact": "intimate, introspective moment",
1433         "uniqueness_creativity": "self-portrait study",
1434         "visual_style": "oil painting",
1435         "references": {
1436             "content": [],
1437             "style": ""
1438         }
1439     },
1440     "ambiguous_elements": [
1441         {
1442             "element": "reflection",
1443             "reason": "Could mean mirror image or philosophical
1444                 contemplation",
1445             "suggested_questions": [
1446                 "Is this a physical reflection in a mirror or a metaphorical
1447                 self-reflection?",
1448                 "If it's a mirror reflection, what type of mirror setup is
1449                 being used?",
1450                 "What perspective is the reflection being painted from?"
1451             ],
1452             "creative_fill": "Mirror image - context of 'painting in oils'
1453                 and 'capturing detail' indicates physical reflection rather
1454                 than abstract contemplation"
1455         },
1456         {
1457             "element": "morning light",
1458             "reason": "Specific lighting conditions not detailed",
1459             "suggested_questions": [
1460                 "What direction is the morning light coming from?",
1461                 "Are there any specific shadow patterns?",
1462                 "Is it early morning or late morning light?"
1463             ],
1464             "creative_fill": "Soft, directional morning light creating
1465                 gentle shadows and warm highlights"
1466         }
1467     ]
1468 }

```

```

1458 }
1459
1460 ### Example 3
1461 Given prompt: "A shiny apple sitting on a desk next to a keyboard"
1462 Creativity level: MEDIUM
1463 {
1464   "identified_elements": {
1465     "main_subjects": [
1466       {
1467         "apple": "shiny object",
1468         "keyboard": "computer keyboard",
1469         "spatial_relationships": "apple positioned next to keyboard
1470           on desk surface"
1471       }
1472     ],
1473     "background": "desk environment",
1474     "composition": "close-up still life",
1475     "color_harmony": "modern office colors",
1476     "lighting": "clear lighting to show shininess",
1477     "focus_sharpness": "sharp focus on main objects",
1478     "emotional_impact": "clean, modern feel",
1479     "uniqueness_creativity": "juxtaposition of natural/tech elements",
1480     "visual_style": "contemporary photography",
1481     "references": {
1482       "content": [],
1483       "style": ""
1484     }
1485   },
1486   "ambiguous_elements": [
1487     {
1488       "element": "apple",
1489       "reason": "Could refer to either the fruit or an Apple product (
1490         like an Apple mouse or AirPods)",
1491       "suggested_questions": [
1492         "Is this referring to the fruit apple or an Apple technology
1493         product?",
1494         "If it's a fruit, what variety/color of apple?",
1495         "If it's an Apple product, which specific device is it?"
1496       ],
1497       "creative_fill": "Red fruit apple - while the desk/keyboard
1498         setting might suggest tech, without specific tech-related
1499         context, assume the natural fruit"
1500     },
1501     {
1502       "element": "keyboard",
1503       "reason": "Style and type of keyboard not specified",
1504       "suggested_questions": [
1505         "What type of keyboard is it (mechanical, membrane, laptop)?"
1506         ,
1507         "Is it a specific brand or color of keyboard?",
1508         "Is it a full-size keyboard or a compact one?"
1509       ],
1510       "creative_fill": "Modern black computer keyboard with white
1511         backlight"
1512     }
1513   ]
1514 }

```

### Prompt Refinement Sub-Agent

You are a Qwen prompt expert. Your PRIMARY GOAL is to stay faithful to the original prompt **while** resolving ambiguities. CRITICAL: The refined prompt must preserve the core intent, subjects, and atmosphere of the original prompt.

GROUNDING PRINCIPLES:

```

1512 - PRESERVE ALL original subjects, objects, and key elements
1513 mentioned in the original prompt
1514 - MAINTAIN the original scene's core atmosphere, mood, and
1515 context
1516 - ONLY ADD details that directly support or clarify the original
1517 prompt
1518 - AVOID introducing new subjects, objects, or concepts not
1519 clearly implied by the original
1520 - Creative filling should ENHANCE, not REPLACE or OVERSHADOW
original elements
1521
1522 Steps:
1523 1. START with the original prompt as the foundation - preserve
1524 its core structure and intent
1525 2. RESOLVE ambiguous elements using creative_fill from analysis,
1526 but only for true ambiguities
1527 3. ADD minimal supporting details ONLY if creativity_level is
1528 MEDIUM or HIGH AND they enhance the original concept
1529 4. ENSURE the refined prompt feels like a clearer version of the
1530 original, not a different scene
1531 5. If there is reference image, must keep the its directory
1532
1533 Return a JSON with:
1534 {
1535   "refined_prompt": "A refined version that stays closely
1536 grounded to the original prompt while resolving necessary
1537 ambiguities. The result should read as a natural
1538 clarification of the original, maintaining its core
1539 essence.",
1540   "reasoning": "Explain how the refinement preserves the
1541 original prompt's intent while addressing ambiguities."
1542 }

```

### Adaptive Negative Prompt Sub-Agent

```

1540 You are an expert at generating negative prompts for image generation
1541 models like Qwen-Image and Qwen-Image-Edit.
1542
1543 A negative prompt specifies what should NOT appear in the generated image
1544 . It helps avoid common issues like:
1545 - Poor quality artifacts (blurry, distorted, low quality, pixelated)
1546 - Unwanted objects or elements that commonly appear in similar scenes
1547 - Inappropriate content or style mismatches
1548 - Technical issues (watermarks, text overlays, borders)
1549
1550 Guidelines:
1551 1. Keep negative prompts concise but comprehensive
1552 2. Focus on common unwanted elements for the specific scene type
1553 3. Include general quality-related terms
1554 4. Avoid being too restrictive - don't negate core elements of the
1555 positive prompt
1556 5. Consider the context and style of the positive prompt
1557
1558 Return a JSON with:
1559 {
1560   "negative_prompt": "comma-separated negative prompt terms",
1561   "reasoning": "explanation of why these negative elements were chosen"
1562 }
1563
1564 Examples:
1565 - Portrait: "blurry, low quality, distorted face, multiple heads, extra
limbs, watermark, text"
1566 - Landscape: "people, buildings, text, watermark, low quality, blurry,
oversaturated"
1567 - Fantasy scene: "modern objects, realistic style, low quality, blurry,
watermark, text"

```

1566 I.2 EVALUATION AGENT  
1567

1568 You are an expert image judge. The evaluator should assess the generated  
1569 image based on two primary dimensions: Aesthetic Quality and Text-  
1570 Image Alignment. Each criterion should be rated on a 0-10 scale,  
1571 where 0 represents poor performance and 10 represents an ideal result  
1572 .

1573 Mainly focus on the original prompt: {config.prompt\_understanding['  
1574 original\_prompt']}.  
1575

1576 1. Aesthetic Quality (0-10) Evaluate the artistic and visual appeal of  
1577 the generated image using the following factors: - Composition:  
1578 Effectiveness of image framing, balance, rule of thirds, leading  
1579 lines, and visual stability. - Color Harmony: Effectiveness of color  
1580 combinations, contrast, and saturation **in** creating a pleasing visual  
1581 experience. - Lighting & Exposure: Appropriateness of brightness,  
1582 contrast, highlights, and shadows **in** creating a visually appealing  
1583 image. - Focus & Sharpness: Clarity of the image, appropriate depth  
1584 of field, and emphasis on key subjects. - Emotional Impact: The  
1585 images ability to evoke emotions, tell a story, or convey a strong  
1586 mood. - Uniqueness & Creativity: Novelty **in** subject matter,  
1587 perspective, or artistic choices that make the image stand out.

1588 2. Text-Image Alignment (0-10) Evaluate how well the generated image  
1589 adheres to the provided prompt, considering key elements from the  
1590 prompt analysis: - Presence of Main Subjects: Whether all key objects  
1591 , characters, or elements explicitly mentioned **in** the prompt appear  
1592 **in** the image. - Accuracy of Spatial Relationships: Whether the  
1593 placement of subjects aligns with the described relationships (e.g.,  
1594 "a cat sitting on a table" should not have the **cat** under the table).  
1595 - Adherence to Style Requirements: If a specific visual style (e.g.,  
1596 "oil painting," "realistic photography") is mentioned, evaluate  
1597 whether the generated image follows this directive. - Background  
1598 Representation: If a background is specified **in** the prompt, check  
1599 whether it aligns with the description **in** terms of elements, lighting  
1600 , and ambiance. # Scoring Explanation - Each subcategory score (e.g.,  
1601 Composition, Presence of Main Subjects) should be rated from 0 to  
1602 10, where: - 0-3 Poor or missing implementation of the aspect. - 4-6  
1603 Moderate adherence but with noticeable flaws. - 7-9 Strong  
1604 adherence with minor imperfections. - 10 Perfect execution. - Main  
1605 Subjects Present (Boolean): Set to true if all essential subjects  
1606 from the prompt appear in the image; otherwise, false. - Missing  
1607 Elements (List of Strings): Lists key elements from the prompt that  
1608 were not correctly represented in the generated image. - Improvement  
1609 Suggestions (String): Provide specific recommendations focusing  
1610 primarily on aspects directly related to: 1. The original prompt: {  
1611 config.prompt\_understanding['original\_prompt']} 2. The user provided  
1612 information: {config.prompt\_understanding['user\_clarification']}  
1613 Focus less on improvements not mentioned in the original prompt or  
1614 user clarification. - Return the results in JSON format with the  
1615 following structure:

```
1616 {{
1617   "aesthetic_reasoning": str,
1618   "aesthetic_score": {{
1619     "Composition": float,
1620     "Color Harmony": float,
1621     "Lighting & Exposure": float,
1622     "Focus & Sharpness": float,
1623     "Emotional Impact": float,
1624     "Uniqueness & Creativity": float
1625   }},
1626   "alignment_reasoning": str,
1627   "alignment_score": {{
1628     "Presence of Main Subjects": float,
```

```

1620     "Accuracy of Spatial Relationships": float,
1621     "Adherence to Style Requirements": float,
1622     "Background Representation": float
1623   }},
1624   "artifacts": {{
1625     "detected_artifacts": [str],
1626     "artifact_reasoning": str
1627   }},
1628   "main_subjects_present": bool,
1629   "missing_elements": [str],
1630   "improvement_suggestions": str,
1631   "overall_reasoning": str,
1632 }}
1633
1634 ### Example 1
1635 Given prompt: "A hyper-realistic painting of a fox in a misty forest,
1636 with warm golden light shining through the trees."
1637
1638 {{
1639   "aesthetic_reasoning": "Strong artistic composition and emotional
1640   impact, but mist and golden light are underrepresented.",
1641   "aesthetic_score": {{
1642     "Composition": 8.5,
1643     "Color Harmony": 9.0,
1644     "Lighting & Exposure": 8.0,
1645     "Focus & Sharpness": 7.5,
1646     "Emotional Impact": 9.5,
1647     "Uniqueness & Creativity": 8.0
1648   }},
1649   "alignment_reasoning": "Fox and forest align well, but mist and
1650   lighting fall short of prompt description.",
1651   "alignment_score": {{
1652     "Presence of Main Subjects": 9.0,
1653     "Accuracy of Spatial Relationships": 8.0,
1654     "Adherence to Style Requirements": 7.0,
1655     "Background Representation": 9.0
1656   }},
1657   "artifacts": {{
1658     "detected_artifacts": ["Minor noise in mist rendering"],
1659     "artifact_reasoning": "Mist appears pixelated due to blending
1660     inconsistencies."
1661   }},
1662   "main_subjects_present": true,
1663   "missing_elements": ["Mist not prominent enough", "Golden light too
1664   subtle"],
1665   "improvement_suggestions": "Enhance mist and intensify golden light
1666   for better atmosphere.",
1667   "overall_reasoning": "Strong aesthetics and alignment but weakened
1668   atmosphere due to faint mist and lighting.",
1669 }}
1670
1671 ### Example 2
1672 Given prompt: "A cozy living room with a vintage leather armchair, a
1673 sleeping cat on a Persian rug, and antique books on wooden shelves."
1674
1675 {{
1676   "aesthetic_reasoning": "Visually pleasing composition and colors, but
1677   emotional depth is reduced by missing cat and rug.",
1678   "aesthetic_score": {{
1679     "Composition": 8.0,
1680     "Color Harmony": 8.5,
1681     "Lighting & Exposure": 7.5,
1682     "Focus & Sharpness": 8.0,
1683     "Emotional Impact": 6.5,
1684     "Uniqueness & Creativity": 7.0
1685   }},

```

```

1674 "alignment_reasoning": "Armchair and shelves present, but cat and rug
1675     absent, reducing prompt fidelity.",
1676 "alignment_score": {{
1677     "Presence of Main Subjects": 2.0,
1678     "Accuracy of Spatial Relationships": 7.5,
1679     "Adherence to Style Requirements": 8.0,
1680     "Background Representation": 8.0
1681 }}},
1682 "artifacts": {{
1683     "detected_artifacts": ["Texture tiling on bookshelf"],
1684     "artifact_reasoning": "Bookshelf wood grain repeats unnaturally,
1685     indicating AI tiling artifact."
1686 }}},
1687 "main_subjects_present": false,
1688 "missing_elements": ["No cat", "No Persian rug", "Armchair lacks
1689     vintage style"],
1690 "improvement_suggestions": "Add cat on Persian rug and adjust armchair
1691     to appear vintage.",
1692 "overall_reasoning": "Good aesthetics but major alignment issues due
1693     to missing key subjects.",
1694 }}

```

### 1693 I.3 GUIDANCE AGENT

#### 1695 **Trajectory Analysis**

```

1696 You are an expert AI model performance analyst. Analyze the {model_name}
1697     model's performance in this image generation task.
1698
1699 WORKFLOW EXECUTION EXPLANATION:
1700     {process_summary}
1701
1702 EXECUTION-SPECIFIC DATA (in workflow sequence):
1703
1704 CREATIVITY LEVEL SETTING:
1705     Level: {config_context['prompt_details']['creativity_level']}
1706     Impact: Determined how autonomously the system handled ambiguous prompt
1707     elements
1708
1709 INTENTION ANALYSIS RESULTS:
1710     Original prompt: "{config_context['prompt_details']['original']}"
1711     Analysis findings: System identified visual elements and ambiguous
1712     aspects requiring interpretation
1713
1714 PROMPT REFINEMENT OUTPUT:
1715     Refined prompt: "{config_context['prompt_details']['refined']}"
1716     Refinement quality: {'Significant refinement applied' if config_context[
1717         'prompt_details']['original'] != config_context['prompt_details']['
1718         refined'] else 'Minimal refinement needed'}
1719
1720 NEGATIVE PROMPT CREATION:
1721     Negative prompt applied: "{config_context['generation_params']['
1722         negative_prompt']}"
1723     Purpose: Targeted prevention of unwanted artifacts and quality issues
1724
1725 GENERATION EXECUTION:
1726     Selected model: {config_context['model_selection']['chosen_model']}
1727     Selection reasoning: {config_context['model_selection']['
1728         selection_reasoning']}
1729     System confidence: {config_context['model_selection']['confidence_score'
1730         ]}/10
1731     Reference image used: {'Yes' if config_context['generation_params']['
1732         has_reference_image'] else 'No'}
1733     Generation seed: {config_context['generation_params']['seed']}

```

```

1728 EVALUATION RESULTS:
1729 Automated score: {config_context['evaluation_metrics']['evaluation_score
1730 ']} / 10
1731 User feedback: {config_context['evaluation_metrics']['user_feedback'] or
1732 'None provided'}
1733 REGENERATION STATUS:
1734 Current attempt: #{config_context['system_context']['current_attempt']}
1735 of {config_context['system_context']['total_attempts']} total
1736 Improvement suggestions: {config_context['evaluation_metrics']['
1737 improvement_suggestions'] or 'None from previous cycles'}
1738 Human oversight: {'Enabled' if config_context['system_context']['
1739 human_in_loop'] else 'Autonomous operation'}
1740 ANALYSIS REQUIREMENTS:
1741 Analyze this model's performance considering the execution flow:
1742 1. How well did the model respond to the creativity level and prompt
1743 refinement quality?
1744 2. Effectiveness of negative prompt in preventing unwanted artifacts
1745 3. Quality of prompt polishing for this specific model's characteristics
1746 4. Model selection appropriateness based on the refined prompt
1747 requirements
1748 5. Technical execution quality visible in the generated image
1749 6. Prompt adherence and creative interpretation balance
1750 7. Reference image utilization (if applicable)
1751 IMPORTANT: You will see the actual generated image(s) below. Provide
1752 specific visual analysis referencing the execution context.
1753 Return a JSON response with detailed breakdown by workflow trajectory:
1754 {{
1755   "trajectory_reasoning": "Overall analysis of how the workflow
1756   execution played out from start to finish, including key decision
1757   points, transitions between steps, and how each step influenced
1758   the next. Analyze the logical flow and coherence of the entire
1759   process.",
1760   "step_scores": {{
1761     "creativity_level": "Score from 1-10 how appropriate the creativity
1762     level setting was for this specific prompt",
1763     "intention_analysis": "Score from 1-10 how effective the intention
1764     analysis was for this prompt",
1765     "prompt_refinement": "Score from 1-10 how well the prompt was
1766     refined for optimal generation",
1767     "negative_prompt": "Score from 1-10 how effective the negative
1768     prompt was in preventing issues",
1769     "generation": "Score from 1-10 the overall quality of the generated
1770     image",
1771     "evaluation": "Score from 1-10 how accurate the system's evaluation
1772     was"
1773   }},
1774   "successes": {{
1775     "creativity_level": "How well the creativity level setting worked
1776     for this prompt and model",
1777     "intention_analysis": "Effectiveness of the intention analysis in
1778     guiding the process",
1779     "prompt_refinement": "Quality and appropriateness of the prompt
1780     refinement",
1781     "negative_prompt": "How well the negative prompt prevented unwanted
1782     artifacts",
1783     "generation": "Model selection appropriateness and generation
1784     quality",
1785     "evaluation": "Accuracy of evaluation scoring relative to visual
1786     quality"
1787   }},
1788   "pitfalls": {{

```

```

1782     "creativity_level": "Issues with creativity level setting for this
1783     prompt type",
1784     "intention_analysis": "Missed elements or incorrect analysis in
1785     intention phase",
1786     "prompt_refinement": "Problems with refined prompt quality or
1787     completeness",
1788     "negative_prompt": "Artifacts that negative prompt failed to
1789     prevent",
1790     "generation": "Model selection issues or generation quality
1791     problems",
1792     "evaluation": "Evaluation inaccuracies or scoring misalignment"
1793   }},
1794   "overall_rating": "Rate from 1-10 the overall effectiveness of the
1795   entire workflow for this specific prompt"
1796 }}
1797
1798 Focus on specific observations from the actual generated image and how
1799 each workflow step contributed to or detracted from the final result.
1800 ""
1801
1802 Note: If regeneration using Qwen-Image-Edit occurs, all the images are also passed to the AGUID
1803 and the final information is added to both the databases of Qwen-Image and Qwen-Image-Edit.
1804
1805 Guidance generation from similar trajectories for (Qwen Image)
1806 You are an expert at analyzing patterns in image generation workflows to
1807 provide SPECIFIC, CONCRETE, and ACTIONABLE guidance.
1808
1809 Your specialty is identifying context-specific techniques that work for
1810 particular types of images, not general best practices.
1811
1812 Focus on content-specific insights like:
1813 - For portrait images: specific lighting techniques, expression guidance,
1814   compositional elements
1815 - For landscape scenes: time of day effects, weather condition impacts,
1816   foreground element placement
1817 - For abstract concepts: style reference importance, compositional
1818   balance techniques, color palette guidance
1819
1820 Avoid general advice like "add more detail" or "be more specific".
1821 Instead provide domain-specific, technical recommendations based on the
1822 actual examples you analyze.
1823
1824 Based on the analysis of similar prompts, extract CONCRETE and SPECIFIC
1825 workflow guidance for the new prompt.
1826
1827 SIMILAR PROMPTS DATA:
1828 {similar_data_text}
1829
1830 {workflow_description}
1831
1832 {task_focus}
1833
1834 Return JSON with this EXACT structure:
1835 {{
1836   "step_analysis": {{
1837     "creativity_level_determination": {{
1838       "success_patterns": "SPECIFIC creativity level patterns that
1839       worked for this type of prompt",
1840       "failure_patterns": "SPECIFIC creativity level mistakes observed
1841       with similar prompts",
1842       "impact_on_next": "CONCRETE impact on intention analysis quality
1843       ",
1844       "preventive_guidance": "DETAILED advice with SPECIFIC examples
1845       of what to look for",

```

```

1836     "recommended_score": "Recommend target score (1-10) based on
1837         similar examples"
1838     }},
1839     "intention_analysis": {{
1840         "success_patterns": "SPECIFIC intention analysis techniques that
1841             worked for this subject matter",
1842         "failure_patterns": "SPECIFIC intention analysis pitfalls
1843             observed with similar prompts",
1844         "impact_on_next": "CONCRETE impact on prompt refinement",
1845         "preventive_guidance": "DETAILED advice with SPECIFIC elements
1846             to identify",
1847         "recommended_score": "Recommend target score (1-10) based on
1848             similar examples"
1849     }},
1850     "prompt_refinement": {{
1851         "success_patterns": "SPECIFIC refinement strategies that
1852             enhanced similar prompts",
1853         "failure_patterns": "SPECIFIC refinement mistakes observed in
1854             similar cases",
1855         "impact_on_next": "CONCRETE impact on negative prompt generation
1856             ",
1857         "preventive_guidance": "DETAILED advice with SPECIFIC refinement
1858             techniques",
1859         "recommended_score": "Recommend target score (1-10) based on
1860             similar examples"
1861     }},
1862     "negative_model_selection": {{
1863         "success_patterns": "SPECIFIC negative prompt and model
1864             selection strategies that were effective for this subject",
1865         "failure_patterns": "SPECIFIC negative prompt and model
1866             selection issues observed in similar cases",
1867         "impact_on_next": "CONCRETE impact on prompt polishing",
1868         "preventive_guidance": "DETAILED advice with SPECIFIC negative
1869             prompt terms and model selection criteria",
1870         "recommended_score": "Recommend target score (1-10) based on
1871             similar examples"
1872     }},
1873     "image_generation": {{
1874         "success_patterns": "SPECIFIC generation parameters that worked
1875             for similar content",
1876         "failure_patterns": "SPECIFIC generation issues observed with
1877             this prompt type",
1878         "impact_on_next": "CONCRETE impact on evaluation accuracy",
1879         "preventive_guidance": "DETAILED advice with SPECIFIC generation
1880             settings",
1881         "recommended_score": "Recommend target score (1-10) based on
1882             similar examples"
1883     }},
1884     "quality_evaluation": {{
1885         "success_patterns": "SPECIFIC evaluation criteria effective for
1886             this content type",
1887         "failure_patterns": "SPECIFIC evaluation pitfalls observed with
1888             similar outputs",
1889         "impact_on_next": "CONCRETE impact on regeneration decision",
1890         "preventive_guidance": "DETAILED advice with SPECIFIC quality
1891             indicators",
1892         "recommended_score": "Recommend target score (1-10) based on
1893             similar examples"
1894     }},
1895     "regeneration_decision": {{
1896         "success_patterns": "SPECIFIC decision criteria that led to
1897             successful outcomes",
1898         "failure_patterns": "SPECIFIC regeneration mistakes observed
1899             with similar content",
1900         "impact_on_next": "N/A - final step",

```

```

1890     "preventive_guidance": "DETAILED advice with SPECIFIC decision
1891         factors",
1892     "recommended_score": "Recommend target score (1-10) based on
1893         similar examples"
1894     }}
1895 },
1896 "workflow_insights": {{
1897     "critical_dependencies": "SPECIFIC step dependencies most relevant
1898         to this prompt type",
1899     "common_failure_chains": "CONCRETE failure patterns observed in
1900         similar cases",
1901     "success_combinations": "SPECIFIC combinations of choices that
1902         worked well for this content",
1903     "overall_rating_prediction": "Predict the likely overall success
1904         rating (1-10) for this prompt type"
1905 }}
1906 }}
1907 Your guidance must be SPECIFIC to the prompt type and content domain, not
1908 generic best practices.
1909 For example, instead of 'Use detailed prompts', say 'For architectural
1910 images, specify architectural style, materials, lighting conditions,
1911 and surrounding environment'.
1912 Use concrete, actionable advice derived from the data, not theoretical
1913 recommendations.
1914 IMPORTANT: Ensure the guidance provides DOMAIN-SPECIFIC advice for the
1915 content type in the prompts, not just generic image generation tips.
1916 Guidance generation from similar trajectories for (Qwen Image Edit)
1917 You are an expert at analyzing patterns in image generation workflows to
1918 provide SPECIFIC, CONCRETE, and ACTIONABLE guidance.
1919 Your specialty is identifying context-specific techniques that work for
1920 particular types of images, not general best practices.
1921 Focus on content-specific insights like:
1922 - For portrait images: specific lighting techniques, expression guidance,
1923     compositional elements
1924 - For landscape scenes: time of day effects, weather condition impacts,
1925     foreground element placement
1926 - For abstract concepts: style reference importance, compositional
1927     balance techniques, color palette guidance
1928 Avoid general advice like "add more detail" or "be more specific".
1929 Instead provide domain-specific, technical recommendations based on the
1930 actual examples you analyze.
1931 Based on the analysis of similar prompts, extract CONCRETE and SPECIFIC
1932 workflow guidance for the new prompt.
1933 SIMILAR PROMPTS DATA:
1934 {similar_data_text}
1935
1936 {workflow_description}
1937
1938 {task_focus}
1939 Return JSON with this EXACT structure:
1940 {{
1941     "step_analysis": {{
1942         "image_editing": {{

```

```

1944     "success_patterns": "SPECIFIC editing parameters and techniques
1945         that worked for similar content and edit types",
1946     "failure_patterns": "SPECIFIC editing issues observed with this
1947         edit operation type",
1948     "impact_on_next": "CONCRETE impact on evaluation accuracy",
1949     "preventive_guidance": "DETAILED advice with SPECIFIC editing
1950         settings, reference image usage, and blending techniques",
1951     "recommended_score": "Recommend target score (1-10) based on
1952         similar examples"
1953   }},
1954   "quality_evaluation": {{
1955     "success_patterns": "SPECIFIC evaluation criteria effective for
1956         this edit content type",
1957     "failure_patterns": "SPECIFIC evaluation pitfalls observed with
1958         similar edit outputs",
1959     "impact_on_next": "CONCRETE impact on re-edit decision",
1960     "preventive_guidance": "DETAILED advice with SPECIFIC edit
1961         quality indicators and success metrics",
1962     "recommended_score": "Recommend target score (1-10) based on
1963         similar examples"
1964   }},
1965   "workflow_insights": {{
1966     "critical_dependencies": "SPECIFIC step dependencies most relevant
1967         to this prompt type",
1968     "common_failure_chains": "CONCRETE failure patterns observed in
1969         similar cases",
1970     "success_combinations": "SPECIFIC combinations of choices that
1971         worked well for this content",
1972     "overall_rating_prediction": "Predict the likely overall success
1973         rating (1-10) for this prompt type"
1974   }}
1975 }}
1976
1977 Your guidance must be SPECIFIC to the prompt type and content domain, not
1978 generic best practices.
1979 For example, instead of 'Use detailed prompts', say 'For architectural
1980 images, specify architectural style, materials, lighting conditions,
1981 and surrounding environment'.
1982 Use concrete, actionable advice derived from the data, not theoretical
1983 recommendations.
1984
1985 IMPORTANT: Ensure the guidance provides DOMAIN-SPECIFIC advice for the
1986 content type in the prompts, not just generic image generation tips.

```

## J ALGORITHM

1982  
1983  
1984  
1985  
1986  
1987  
1988  
1989  
1990  
1991  
1992  
1993  
1994  
1995  
1996  
1997

1998  
1999  
2000  
2001  
2002  
2003  
2004  
2005  
2006  
2007  
2008  
2009  
2010  
2011  
2012  
2013  
2014  
2015  
2016  
2017  
2018  
2019  
2020  
2021  
2022  
2023  
2024  
2025  
2026  
2027  
2028  
2029  
2030  
2031  
2032  
2033  
2034  
2035  
2036  
2037  
2038  
2039  
2040  
2041  
2042  
2043  
2044  
2045  
2046  
2047  
2048  
2049  
2050  
2051

---

**Algorithm 1** Main algorithm for *SIDiffAgent*


---

**Input:** User prompt  $P$ , Models: QI (Qwen-Image), QIE (Qwen-Image-Edit), QE (Qwen-Embedding), Knowledge base  $KB$  (initially may be empty), Hyperparams:  $\tau$  (eval threshold),  $E$  (max edits),  $K$  (retrieval size)

**Output:** Final image  $I^*$ , final prompt  $F$ , recorded trajectory  $T$

```

1 Function SIDiffAgent ( $P$ ):
2    $G \leftarrow$  RetrieveGuidance ( $KB, P, K$ ) // retrieve and generate guidance
3    $c \leftarrow S_{\text{CRE}}(P, G)$  // creativity level: {low, med, high}
4    $S \leftarrow S_{\text{INT}}(P, c, G)$  // semantic specification, disambiguations
5    $P_{\text{pos}} \leftarrow S_{\text{REF}}(P, S, G)$  // refined positive prompt
6    $P_{\text{neg}} \leftarrow S_{\text{NEG}}(P, P_{\text{pos}}, G)$  // adaptive negative prompt + universal
   safeguards
7    $edits \leftarrow 0$   $I \leftarrow \emptyset$   $T \leftarrow \emptyset$  while  $edits \leq E$  do
8      $I \leftarrow S_{\text{GEN}}(P_{\text{pos}}, P_{\text{neg}}, I)$  // use QI
9      $report \leftarrow A_{\text{EVAL}}(I, P, P_{\text{pos}}, G)$  // returns scores and suggestions
10    Record decision node results in  $T$  via RecordTrajectory()  $score \leftarrow$  aver-
11    age(report.aesthetic, report.alignment)
12    if  $score \geq \tau$  then
13      break // satisfactory output
14    else
15       $P_{\text{pos}}^{\text{new}}, P_{\text{neg}}^{\text{new}} \leftarrow$  re-run  $A_{\text{ORC}}$ ; if  $report$  suggests localized correction and  $edits < E$ 
16      then
17         $I \leftarrow S_{\text{GEN}}(P_{\text{pos}}^{\text{new}}, P_{\text{neg}}^{\text{new}}), edit, I$  // use QIE
18         $edits \leftarrow edits + 1$ 
19    end
20  end
  UpdateKnowledgeBase ( $KB, T$ ) return  $I, P_{\text{pos}}, P_{\text{neg}}, T$ 

```

---

---

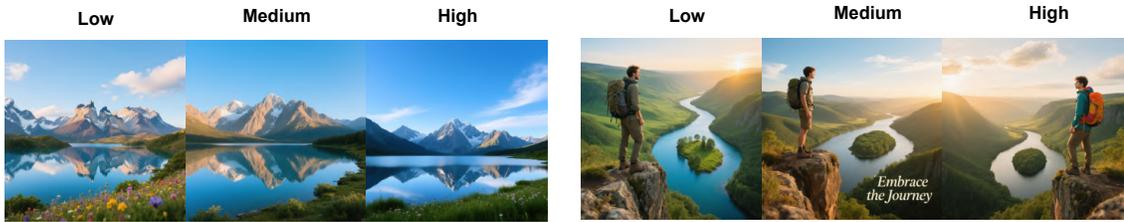
```

2052 Algorithm 2 Helper functions for main algorithm for SIDiffAgent
2053
2054 Function  $S_{CRE}(P, G)$ :
2055   analyse guidance G, compute length, specificity cues,
2056   presence of attributes if very short or vague then
2057     | return high
2058   else if moderately detailed then
2059     | return medium
2060   else
2061     | return low
2062   end
2063 end
2064 Function  $S_{INT}(P, c, G)$ :
2065   analyse guidance G, extract main subjects, attributes,
2066   relations, composition, lighting, style if  $c = high$  then
2067     | infer missing attributes using templates and priors
2068   else if  $c = medium$  then
2069     | propose optional clarifications and minor assumptions
2070   else
2071     | preserve original prompt avoid adding assumptions
2072   end
2073   return structured specification S
2074 end
2075 Function  $S_{REF}(P, S, G)$ :
2076   re-order specification into coherent positive prompt
2077   text  $P_{pos}$  enforce lexical templates for model compatibility return  $P_{pos}$ 
2078 end
2079 Function  $S_{NEG}(P, P_{pos}, G)$ :
2080    $U \leftarrow$  universal safeguards (e.g., low quality,
2081   blurry, watermark) derive scene-specific negations from  $P_{pos}$  compose concise  $P_{neg} \leftarrow U \cup$ 
2082   scene-specific
2083   negations return  $P_{neg}$ 
2084 end
2085 Function  $S_{GEN}(P_{pos}, P_{neg}, mode, I_{in} = \emptyset)$ :
2086   if  $mode = initial$  then
2087     | return  $QI.generate(P_{pos}, P_{neg})$ 
2088   else if  $mode = edit$  then
2089     | return  $QIE.edit(I_{in}, P_{pos}, P_{neg})$ 
2090 end
2091 Function  $A_{EVAL}(I, P, P_{pos})$ :
2092    $a \leftarrow A_{EVAL}.scoreAesthetic(I)$   $\alpha \leftarrow A_{EVAL}.scoreAlignment(I, P$  or  $P_{pos})$   $sugs \leftarrow$ 
2093    $A_{EVAL}.suggestFixes(I)$   $arts \leftarrow A_{EVAL}.detectArtifacts(I)$  return report with  $(a, \alpha, sug, arts)$ 
2094 end
2095 Function RetrieveGuidance ( $KB, P, K$ ):
2096   if  $KB$  empty then
2097     | return  $\emptyset$ 
2098   embed  $P$  into vector  $v$  retrieve top- $K$  trajectories
2099   by similarity aggregate pitfalls and successful actions across
2100   trajectories return guidance G (structured corrective +
2101   workflow guidance)
2102 end
2103 Function RecordTrajectory():
2104   For each decision node, store: node-id, inputs,
2105   actions, outputs, scores
2106 end
2107 Function UpdateKnowledgeBase ( $KB, T$ ):
2108   compress  $T$  into node-wise summaries of pitfalls
2109   and successes append to model-specific sections in  $KB$  optionally re-index embeddings for
2110   efficient retrieval
2111 end

```

## K QUALITATIVE RESULTS

### Impact of Creativity Level



Mountains with a lake (simple prompt)

A hiker standing on a cliff looking at a valley where a river splits into two streams around a small island (Complex Prompt)

### Prompts and generations across creativity levels (set by Creativity Analysis Subagent)



Amidst a winter wonderland, a rabbit scurries across the snow, leaving tracks beside a peacefully standing deer.

On the left is a metal fork and on the right is a wooden spoon.

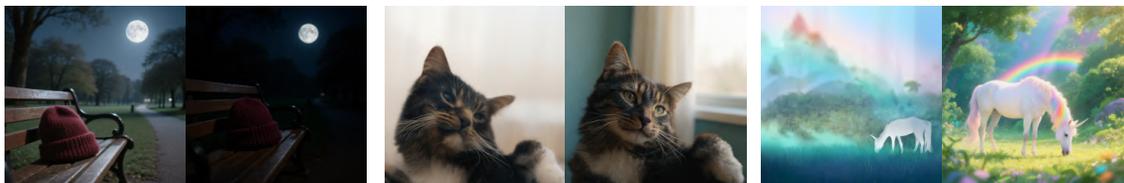
An ice castle standing proudly in the midst of a blizzard.

A potion bubbling brightly inside a cauldron in a shadowy nook.

A celestial comet racing across a star-studded sky.

A lantern casting dim light in a haunted forest.

### Some example of edits performed by Qwen-Edit



A woolen hat left on a park bench under the moonlight.

A cat basking in the sunlight by a window.

A unicorn grazing peacefully in a radiant, rainbow-lit glade

Figure 6: **Top:** Effect of varying creativity levels on simple and complex prompts. **Middle:** Image generations from diverse prompts, with creativity levels determined by  $S_{CRE}$ . **Bottom:** Refinements and enhancements performed by Qwen-Edit on initial suboptimal generations.

2160  
2161  
2162  
2163  
2164  
2165  
2166  
2167  
2168  
2169  
2170  
2171  
2172  
2173  
2174  
2175  
2176  
2177  
2178  
2179  
2180  
2181  
2182  
2183  
2184  
2185  
2186  
2187  
2188  
2189  
2190  
2191  
2192  
2193  
2194  
2195  
2196  
2197  
2198  
2199  
2200  
2201  
2202  
2203  
2204  
2205  
2206  
2207  
2208  
2209  
2210  
2211  
2212  
2213

A cup set to the right of a newspaper



Ancient buildings juxtaposed with sleek, futuristic transports



A group of children playing on the beach



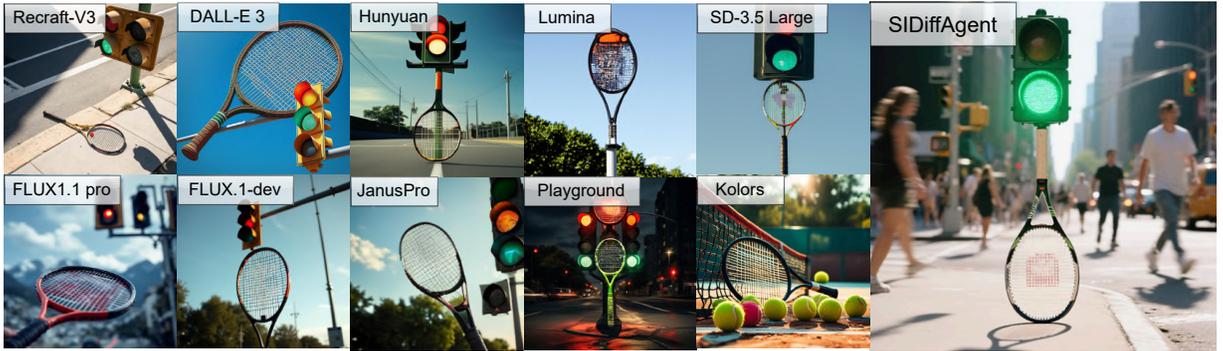
In the classroom, '1+2' is written on the blackboard



Figure 7: Qualitative comparison between multiple state-of-the-art method and *SIDiffAgent* on GenAI-Bench

2214  
2215  
2216  
2217  
2218  
2219  
2220  
2221  
2222  
2223  
2224  
2225  
2226  
2227  
2228  
2229  
2230  
2231  
2232  
2233  
2234  
2235  
2236  
2237  
2238  
2239  
2240  
2241  
2242  
2243  
2244  
2245  
2246  
2247  
2248  
2249  
2250  
2251  
2252  
2253  
2254  
2255  
2256  
2257  
2258  
2259  
2260  
2261  
2262  
2263  
2264  
2265  
2266  
2267

A tennis racket underneath a traffic light



A sign that says 'Google Research Pizza Cafe'



An umbrella on top of a spoon



An elephant under the sea

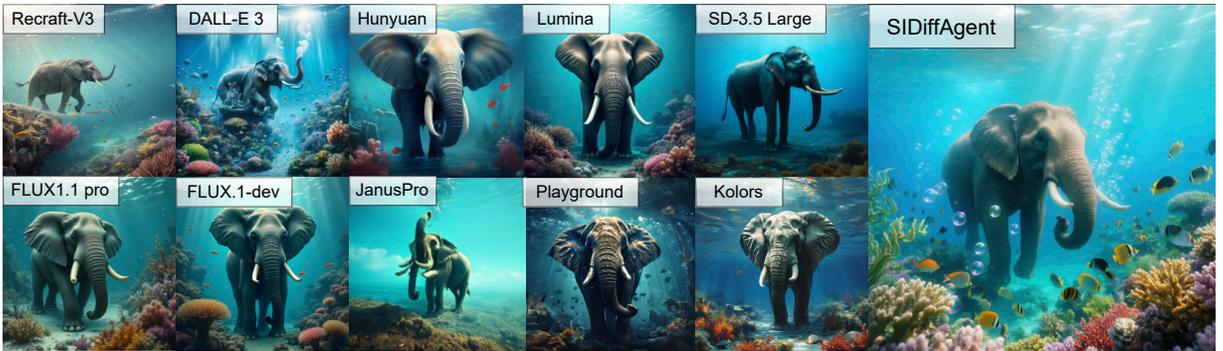


Figure 8: Qualitative comparison between multiple state-of-the-art method and *SIDiffAgent* on DrawBench

2268  
2269  
2270  
2271  
2272  
2273  
2274  
2275  
2276  
2277  
2278  
2279  
2280  
2281  
2282  
2283  
2284  
2285  
2286  
2287  
2288  
2289  
2290  
2291  
2292  
2293  
2294  
2295  
2296  
2297  
2298  
2299  
2300  
2301  
2302  
2303  
2304  
2305  
2306  
2307  
2308  
2309  
2310  
2311  
2312  
2313  
2314  
2315  
2316  
2317  
2318  
2319  
2320  
2321

**Examples of images from SIDiffAgent Episode 2 vs SIDiffAgent (Drawbench)**



Two cats and three dogs sitting on the grass.

A white colored sandwich.

Photo of an athlete cat explaining it's latest scandal at a press conference to journalists.

**Examples of images from SIDiffAgents vs Qwen-Agents (Drawbench)**



An appliance or compartment which is artificially kept cool and used to store food and drink.

A spider with a moustache bidding an equally gentlemanly grasshopper a good day during his walk to work.

A pink colored giraffe.

**Examples of images from SIDiffAgent Episode 2 vs SIDiffAgent (GenAIBenchmark)**



A group of students gathered around a panda, all with digital cameras in their hands.

A hat hanging on a hook to the left of a door.

A snowman with a hat, and no scarf around its neck.

**Examples of images from SIDiffAgents vs Qwen-Agents (GenAIBenchmark)**



A bicycle bell engraved with 'Ring for Joy.'

Two anxious pigs.

A man without a beanie, feeling the chilly morning air directly on his hair.

Figure 9: Qualitative comparison between various ablations described in Section 4.2