

MCCD: Multi-Agent Collaboration-based Compositional Diffusion for Complex Text-to-Image Generation

Mingcheng Li^{1,2*} Xiaolu Hou^{1,2*} Ziyang Liu³ Dingkang Yang^{1,2§} Ziyun Qian^{1,2}
 Jiawei Chen^{1,2} Jinjie Wei^{1,2} Yue Jiang^{1,2} Qingyao Xu^{1,2} Lihua Zhang^{1,2,4,5§}

¹Academy for Engineering and Technology, Fudan University ²Cognition and Intelligent Technology Laboratory (CIT Lab)

³School of Future Science and Engineering, Soochow University, Suzhou, China

⁴Jilin Provincial Key Laboratory of Intelligence Science and Engineering, Changchun, China

⁵Engineering Research Center of AI and Robotics, Ministry of Education, Shanghai, China

{mingchengli21, xlhou23}@m.fudan.edu.cn, dkyang20@fudan.edu.cn

Abstract

Diffusion models have shown excellent performance in text-to-image generation. Nevertheless, existing methods often suffer from performance bottlenecks when handling complex prompts that involve multiple objects, characteristics, and relations. Therefore, we propose a Multi-agent Collaboration-based Compositional Diffusion (MCCD) for text-to-image generation for complex scenes. Specifically, we design a multi-agent collaboration-based scene parsing module that generates an agent system comprising multiple agents with distinct tasks, utilizing MLLMs to extract various scene elements effectively. In addition, Hierarchical Compositional diffusion utilizes a Gaussian mask and filtering to refine bounding box regions and enhance objects through region enhancement, resulting in the accurate and high-fidelity generation of complex scenes. Comprehensive experiments demonstrate that our MCCD significantly improves the performance of the baseline models in a training-free manner, providing a substantial advantage in complex scene generation.

1. Introduction

Recently, diffusion models [7, 30, 32, 38] have shown significant advancements in Text-to-Image (T2I) generation, such as Stable Diffusion [30], Imagen [31] and DALL-E 2/3 [2, 29]. However, despite their noteworthy performance in generating realistic images consistent with text prompts, these models have large limitations in processing complex textual prompts and generating complex scenes, leading to unsatisfactory image generation [10, 20, 23]. Therefore, T2I models require powerful spatial perceptions that can

precisely align multiple objects with different attributes and complex relationships involved in compositional prompts.

Some studies attempt to introduce additional conditions to solve the problems above, which can be divided into two parts: (i) spatial information-based methods [19, 20, 27, 41, 43], and (ii) feedback-based methods. Spatial information-based methods utilize additional spatial information (e.g., layouts and boxes) as conditions to enhance the compositionality of T2I generation. For example, GLIGEN [19] introduces trainable gated self-attention layers to integrate spatial inputs based on the pre-trained stable diffusion models. ReCo [41] utilizes additional sets of positional tokens for T2I generation to achieve effective region control and fine-tune the pre-trained T2I models. Feedback-based methods [9, 15–17, 33, 36] use the generated images as feedback to optimize the T2I generation. For instance, DreamSync [33] utilizes a visual question answer model and an aesthetic quality evaluation model to recognize fine-grained discrepancies between the generated image and the textual input, thus enhancing the semantic alignment capabilities of the T2I model. GORS [15] fine-tunes pre-trained T2I models leveraging generated images aligned to the text prompt and text-image alignment reward-weighted loss. Parrot [17] jointly optimizes the T2I model with a multi-reward optimization strategy to improve the quality of image generation. However, the above methods suffer from the following limitations: (i) lacking fine-grained and precise spatial information guidance, resulting in unrealistic spatial locations and relations in the generated images. (ii) Difficulty in obtaining high-quality image feedback to effectively optimize the image generation. (iii) Fine-tuning of the T2I model (e.g., stable diffusion) results in a large amount of computational and time overheads.

To address the above-mentioned problems, we propose a Multi-agent Collaboration-based Compositional Diffu-

Corresponding authors. ^{}Equal contributions.

sion (MCCD) for high-quality text-to-image generation and complex scene generation. Our novelty stems from three core contributions: **(i)** We propose a multi-agent collaboration-based scene parsing module that constructs multiple agents with different tasks to implement collaboration in forward thought chain reasoning and backward feedback processes to precisely parse key scene elements. **(ii)** Furthermore, Hierarchical Compositional diffusion is proposed to achieve sufficient interaction among parsed multiple scene elements and accurately generate complex scenes that match the text prompt. **(iii)** Comprehensive qualitative and quantitative experiments demonstrate that our method significantly improves the performance of the baseline models in a training-free manner, which has large advantages.

2. Related Work

2.1. Text-to-Image Generation

Text-to-Image (T2I) generation, *i.e.*, text-conditional image synthesis, has been a key research hotspot in the field of multimodal learning [2, 29]. Numerous efforts have been devoted to generating visually natural and realistic images. Generative Adversarial Networks (GANs) are typical T2I models that utilize adversarial training between the generator and the discriminator to produce images that are as close as possible to the real images. In recent years, inspired by the application of Auto-Regressive Models (ARMs) in the field of text generation, many works have achieved favorable results in the field of T2I generation utilizing ARMs, such as CogView[8] and DALL-E 2/3 [2, 29]. Despite the progress achieved by the above studies, they still have many limitations, such as unstable training, difficult convergence, and unidirectional bias, which lead to poorer quality of image generation and lower generalizability. Due to the natural fit of the inductive bias to the image data, diffusion models [7, 13, 30, 32, 38] are now widely used for T2I generation and significantly improve image generation quality and fidelity. GLIDE [25] utilizes the pre-trained CLIP model [28] to achieve semantic alignment between the text prompts and the generated images during the image sampling process. Recent advances in T2I diffusion models have significantly improved the quality and realism of image generation in recent years, such as SDXL [26], DALL-E 2/3 [2, 29] and ContextDiff [40]. In recent years, Large Language Models (LLMs) have been widely used in many tasks due to their powerful comprehension and reasoning capabilities [5, 6, 18, 34, 37, 45]. Many studies use Multimodal Large Language Models (MLLMs) in text-to-image (T2I) generation tasks and achieve performance gains [11, 12, 14, 22, 27, 39, 44]. For example, RPG [39] utilizes the Chain-of-Thought (CoT) of MLLMs to extract layouts from text prompt to enhance T2I generation. LMD [20] utilizes MLLMs to enhance the compositional gener-

ation of diffusion models by generating images grounded on bounding box layouts from the MLLMs. However, the above methods have the following limitations:(1) Simply using MLLMs to process text prompts without sufficiently exploiting and utilizing their powerful comprehension and inference capabilities. (2) Only utilizing MLLMs as a layout generator to control image synthesis, neglecting the extraction of other important scene elements. To address the above problems, we propose a Multi-agent Collaboration-based scene Parsing (MCP) module, which constructs a multi-agent system consisting of multiple agents with various divisions of labor, and utilizes multi-agent collaboration and interaction to achieve adequate extraction and parsing of key scene elements in text prompt, thereby facilitating the subsequent scene generation process.

2.2. Compositional Diffusion Generation

In recent years, many methods have been introduced to improve compositional T2I generation [19, 20, 24, 41, 42] to enhance the capabilities of diffusion models in terms of attribute binding, object relations, and numeracy. For example, ReCo [41] and GLIGEN [19] introduce a location-aware adapter in diffusion models to improve the spatial plausibility of the generated images. LMD [20] utilizes LLM to generate scene layouts and designs a controller to bootstrap the pre-trained diffusion model. RPG [39] denoises each subregion in parallel and applies a post-processing step of resizing and concatenation for high-quality compositional generation. The T2I-Adapter [24] facilitates compositional T2I generation by controlling the semantic structure through some high-level features of the image. However, these methods can only implement coarse compositional control, leading to unsatisfactory image generation results, especially when dealing with complex prompts. Therefore, we propose hierarchical compositional diffusion, which utilizes more precise control to progressively refine image synthesis.

3. Methodology

3.1. Overall Framework

As shown in Figure 1, given a complex text prompt containing multiple objects and relations, the goal of MCCD is to produce realistic and high-quality images. Our MCCD is a training-free framework with the following workflow: (i) Multi-agent collaboration-based scene parsing module utilizes a multi-agent collaborative approach to parse individual elements in a text prompt. (2) Hierarchical compositional diffusion is used to interact with scene elements and achieve high-quality image generation using dynamic integration, regional enhancement, and latent space smoothing.

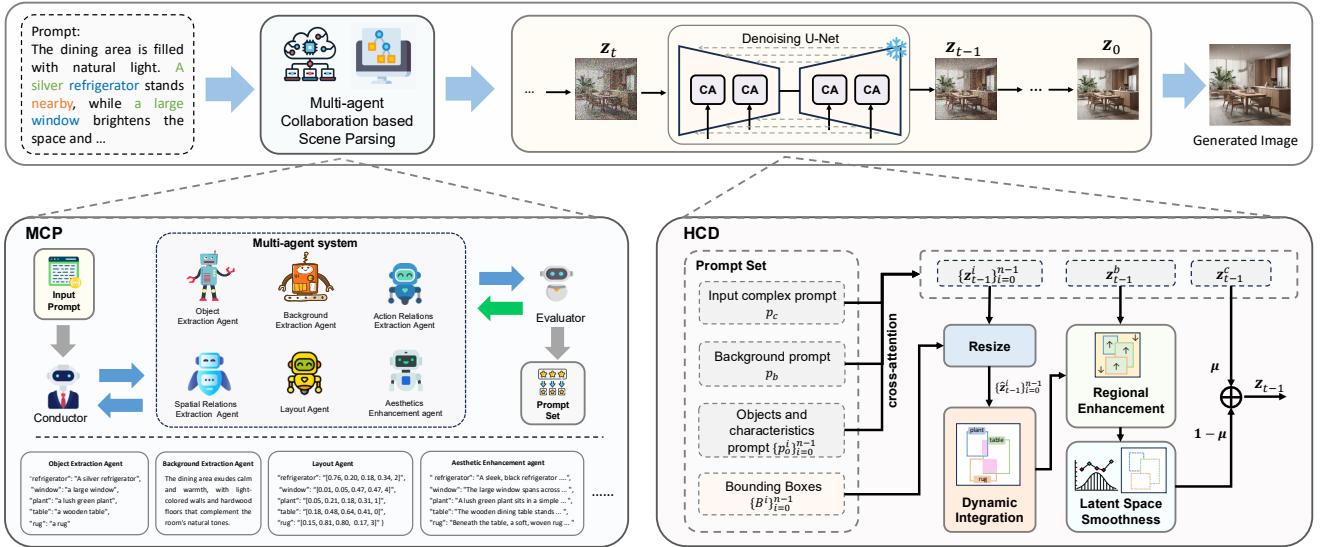


Figure 1. **The overall framework of the proposed MCCD.** MCCD consists of two core components: Multi-agent Collaboration-based scene Parsing (MCP) module and Hierarchical Compositional Diffusion (HCD) module. In MCP, the blue and green arrows indicate forward CoT reasoning and backward feedback processes, respectively

3.2. Multi-agent Collaboration based Scene Parsing

Based on the sufficient consideration of the composition of complex scenes, we categorize the scene elements into several parts: objects and their characteristics, backgrounds, action relations, spatial relations, and layouts. To generate complex scenes, we need to explore all the key elements in the text prompt thoroughly. The previous method [39] utilizes the Chain-of-Thought (CoT) capability of MLLMs to extract objects and layouts from the input text prompt to a certain extent. Despite some success, when dealing with more complex scenarios, the CoT-based approach often fails to sufficiently parse the intricate relations in the scene. Furthermore, unidirectional CoT reasoning lacks error correction mechanisms, leading to inaccurate elemental parsing and image generation. Therefore, we propose a Multi-agent Collaboration-based Scene Parsing (MCP) module that splits the scene parsing process into multiple sub-stages and constructs corresponding specialized agents for each stage, and generates a multi-agent system simultaneously. In the system, each agent dynamically collaborates and interacts with other agents while strictly performing its specialized tasks, thus maximizing the advantages of teamwork. This paradigm ensures that MLLMs accurately parse multiple elements contained in text prompts, providing significant advantages over traditional CoT reasoning.

The constructed multi-agent system consists of six agents whose definitions and tasks are: (1) object extraction agent, whose task is to extract objects and their characteristics from text prompts. (2) Background extraction agent, whose task is to extract the object-independent background

descriptions. (3) Action relations extraction agent, whose task is to capture the action relations between the objects. (4) Spatial relations extraction agent, whose task is to capture the spatial relations between the objects. (5) Layout agent, whose task is to conceptualize the layout of objects for a reasonable composition. (6) Aesthetics enhancement agent, whose task is to beautify the objects' characteristics to enhance the image aesthetics. Our goal is to integrate the outputs of multiple agents to ultimately generate the objects and characteristics prompt, background prompt, and bounding boxes of all objects. Each agent is defined by an MLLM and is associated with a prompt template while allowing access to an external knowledge base.

Moreover, we construct a conductor \mathcal{C} and an evaluator \mathcal{E} to coordinate multiple agents for effective collaboration. Specifically, the conductor dynamically directs different agents to construct a forward CoT in iterations and passes candidate answers to the external program execution environment. The evaluator then uses feedback signals to trigger a backward feedback process. These two processes are elaborated below.

Forward CoT Reasoning. In forward CoT reasoning, the conductor \mathcal{C} adaptively determines the entire set of participating agents, casting the procedure as a sequential decision-making problem in which each possible action corresponds to a particular agent selection. We define the input prompt as \mathcal{P} and a set of predefined agents as $\xi = \{\mathcal{A}_{\phi_1}, \mathcal{A}_{\phi_2}, \dots, \mathcal{A}_{\phi_n}\}$, where n is the total number of agents and ϕ_i is the configuration of the agent i -th. The set of output for the t -th inference step is denoted

as \mathcal{O}_t , and the state is denoted as $\mathcal{S}_t = (\mathcal{P}, \mathcal{O}_t, t)$. With the powerful prompt learning capability of MLLM, agents can achieve the same functionality in a training-free manner as decision-making agents that require large amounts of data for training in traditional reinforcement learning, with significant advantages in terms of efficiency and overhead. Thus, we take the conductor as a policy function for selecting an agent, denoted as:

$$\mathcal{C}(a | s) = P_r \{\mathcal{A}_{\phi_t} = a | S_t = s\}. \quad (1)$$

The agent selection strategy can be translated into the design of the prompt template, which requires prompt engineering to realize the optimal strategy. Each step of forward CoT reasoning is denoted as:

$$\mathcal{A}_{\phi_{i_t}} = \mathcal{C}(S_t), \quad (2)$$

$$o = \mathcal{A}_{\phi_{i_t}}(\mathcal{P}, \mathcal{O}_t), \quad (3)$$

$$\mathcal{O}_{t+1} = \mathcal{O}_t \cup \{o\}, \quad (4)$$

where $\mathcal{A}_{\phi_{i_t}}$ represents the selected i_t -th agent at step t and o denotes the output of the selected agent. After reaching the maximum step T , the forward process is aborted and all outputs are integrated into the result \mathcal{R} .

Backward feedback. The backward feedback strategy enables a multi-agent system to adjust collaborative behavior by using the evaluator's evaluation of the forward response. The execution order of the agents is defined as $\eta = \{\mathcal{A}_{\phi_{i_1}}, \mathcal{A}_{\phi_{i_2}}, \dots, \mathcal{A}_{\phi_{i_n}}\}$, where i_t denotes the index of the agent at step t . The backward feedback process begins with an external feedback f_{raw} , usually provided by the program execution environment, denoted as

$$f_{raw} = \text{execution}(\mathcal{R}). \quad (5)$$

The raw feedback is then evaluated by evaluator \mathcal{E} to obtain an initial signal: $(f_0, sf_0) = \mathcal{E}(f_{raw})$, where f_0 indicates whether the backward process needs to continue or not, and sf_0 indicates the localization of the error in the backward feedback. If the result \mathcal{R} in the forward reasoning is correct, then f_0 is set to false and the whole reasoning process ends. Otherwise, the conductor \mathcal{C} initiates a backward feedback process that updates the response by back-propagating from the last agent. At step t backward, the state update is represented as follows:

$$(f_t, sf_t) \leftarrow \text{feedback} \left(\mathcal{A}_{\phi_{i_{T-t+1}}}, \mathcal{P}, \mathcal{O}_t, f_{t-1} \right), \quad (6)$$

$$\mathcal{O}_{t+1} = \mathcal{O}_t \cup \{sf_t\}. \quad (7)$$

The backward feedback process continues to be executed iteratively until the feedback signal indicates that the agent has made a mistake or until all the agents have undergone reflection. Once this occurs, the forward process will then be executed again. This process is repeated continuously until a reasonable answer is produced.

3.3. Hierarchical Compositional Diffusion

Previous region-based diffusion methods, such as [39], divide images into multiple complementary regions. It takes into account spatial relations to a certain extent, but it cannot cope with more complex scenes due to the non-overlapping and independent properties of regions. Additionally, [35] imposes constraints on cross-attention maps so that the positions and sizes of objects are as consistent as possible with the bounding box, but ignores the consideration of the overlaps between the bounding boxes and smoothness near the bounding box, leading to unrealistic scene synthesis. To this end, we propose a Hierarchical Compositional Diffusion (HCD) module that performs progressive interaction of multiple elements obtained from scene parsing. Specifically, we dynamically balance the overlapping regions between multiple objects using the Gaussian mask and a regional enhancement strategy is used to enlarge the discrepancy between objects and background in the latent space. Moreover, Gaussian smoothing is used to enhance the smoothness around the bounding box.

As shown in Figure 1, the MCP parses a complex prompt with n objects into multiple prompts, *i.e.*, the input complex prompt p_c , the object prompts $\{p_o^i\}_{i=0}^{n-1}$ and background prompt p_b . At each timestep, we feed each prompt into the denoising network in parallel, using cross-attention layers to generate the corresponding latent representations.

$$z_{t-1} = \text{Softmax} \left(\frac{(W_Q \cdot \phi(z_t))(W_K \cdot \psi(p))}{\sqrt{d}} \right) (W_V \cdot \psi(p)), \quad (8)$$

where image latent z_t is the query and each prompt is the key and the value. W_Q, W_K, W_V are linear projections and d is the latent projection dimension of the keys and queries. The latent representations of complex prompt, all object prompts, and background prompt are denoted as z_{t-1}^c , $\{z_{t-1}^i\}_{i=0}^{n-1}$, and z_{t-1}^b . We use bilinear interpolation to resize the latent representation of the object based on the dimensions of the bounding box, denoted as follows:

$$\hat{z}_{t-1}^i = R(z_{t-1}^i, B^i). \quad (9)$$

Complex scenes often contain multiple objects with intricate positions and action relationships across multiple objects, resulting in a large number of overlapping regions. To this end, we design a dynamic integration mechanism based on a depth-aware Gaussian mask to achieve adaptive and smooth feature fusion in overlapping regions. In general, the center region of the bounding box contains the core features of the object, while the features in the edge region gradually decrease in importance. To keep as many core features as possible in the overlapping regions of the bounding box, for each bounding box (x_0, y_0, w, h) , we construct a Gaussian mask matrix that contains weights that gradually

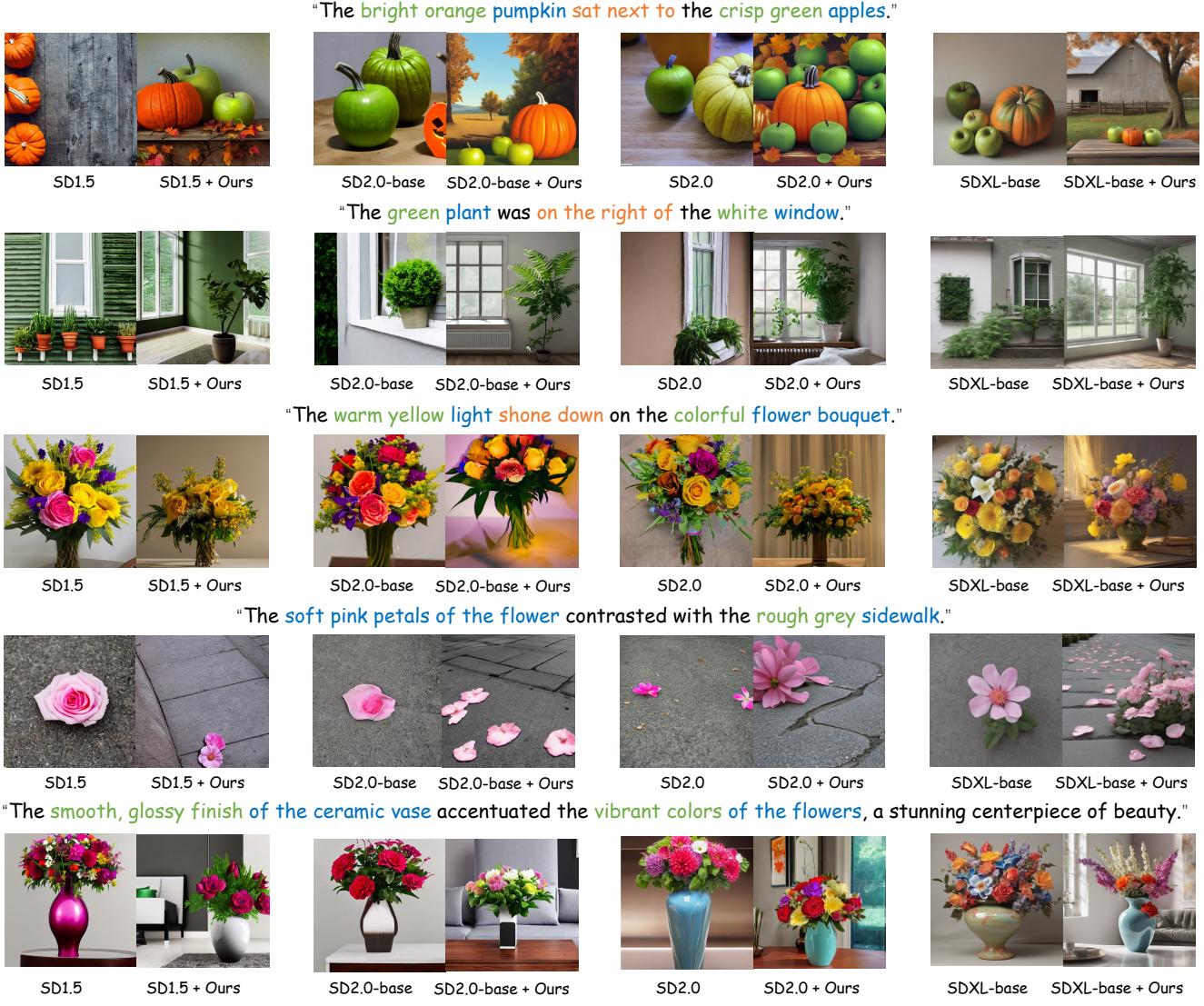


Figure 2. **Qualitative results of MCCD improving diffusion models.** MCCD enhances the attribute binding and spatial relationships of the base diffusion models. The generated results have reasonable backgrounds and detailed textures with great aesthetics and realism.

decay from the center to the edges, represented as follows:

$$M(x, y) = \exp\left(-\frac{(x - \mu_x)^2 + (y - \mu_y)^2}{2\sigma^2}\right), \quad (10)$$

where $(\mu_x, \mu_y) = (x_0 + w/2, y_0 + h/2)$ is the center of the mask matrix, and the $\sigma = \max(w, h)/2$ is the standard deviation controlling the width of the Gaussian distribution. Furthermore, an object’s layer depth d indicates its front and back positions in the image. Objects with smaller layer depth appear closer in the image, so greater weights should be assigned during feature fusion. To achieve a smooth transition between foreground and background, we calculate

late continuous and eased layer depths:

$$w_i = \frac{1}{1 + \exp(\alpha \cdot (d_i - \frac{n-1}{2}))}, \quad (11)$$

where $d \in \{0, 1, \dots, n-1\}$ is the layer depth of the object, and α is the smoothness control parameter. Given any coordinates (x, y) of the overlapping regions of the bounding boxes, their dynamically fused features are represented as,

$$\hat{z}_{t-1}^i(x, y) = \frac{\sum_{i=1}^m w_i \cdot M_i(x, y) \cdot z_{t-1}^i(x, y)}{\sum_{i=1}^m w_i \cdot M_i(x, y)}, \quad (12)$$

where m is the number of overlapping regions. Then, we concatenate the latent representations of all object prompts

Table 1. **Evaluation results on T2I-CompBench.** MCCD shows the best performance in terms of Attribute Binding, Object Relationship, and Complex. Basic data is derived from [15].

Models	Attribute Binding			Object Relationship		Complex↑
	Color↑	Shape↑	Texture↑	Spatial↑	Non-Spatial↑	
Composable Diffusion [23]	0.4063	0.3299	0.3645	0.0800	0.2980	0.2898
Structured Diffusion [10]	0.4990	0.4218	0.4900	0.1386	0.3111	0.3355
Attn-Exct v2 [3]	0.6400	0.4517	0.5963	0.1455	0.3109	0.3401
GORS [15]	0.6603	0.4785	0.6287	0.1815	0.3193	0.3328
DALL-E 2 [29]	0.5750	0.5464	0.6374	0.1283	0.3043	0.3696
PixArt- α [4]	0.6886	0.5582	0.7044	0.2082	0.3179	0.4117
SD1.5 [30]	0.3134	0.2954	0.3772	0.1087	0.3015	0.2994
SD1.5 + MCCD	0.3508	0.3193	0.4026	0.1462	0.3078	0.3054
SD2.0-base [30]	0.4498	0.3584	0.4319	0.1140	0.3045	0.3021
SD2.0-base + MCCD	0.4823	0.3702	0.4621	0.1613	0.3106	0.3146
SD2.0 [30]	0.4852	0.4066	0.4591	0.1490	0.2905	0.3173
SD2.0 + MCCD	0.5090	0.4170	0.4921	0.1785	0.3123	0.3250
SDXL-base [2]	0.5744	0.4705	0.4907	0.1971	0.3009	0.3130
SDXL-base + MCCD	0.6278	0.4832	0.5647	0.2350	0.3132	0.3348

based on the positions of the bounding boxes to achieve control over positional relations. The area not covered by the bounding box is filled by the latent representation of the background. This process is represented as:

$$\mathbf{z}'_{t-1} = \text{Concat}(\{\mathbf{M}_{B^i} \cdot \hat{\mathbf{z}}_{t-1}^i\}_{i=0}^{n-1}, \{\mathbf{M}_{-(\bigcup_{i=0}^{n-1} B^i)} \cdot \mathbf{z}_{t-1}^b\}), \quad (13)$$

where \mathbf{M}_{B^i} and $\mathbf{M}_{-(\bigcup_{i=0}^{n-1} B^i)}$ are masks denoting the region within the bounding box B^i and the background region outside all bounding boxes, respectively.

To ensure that objects are generated within the designated bounding box regions, we implement a regional enhancement mechanism that emphasizes the latent representation of the bounding box regions and simultaneously suppresses the influence of surrounding background areas.

$$\mathbf{z}'_{t-1}+ = \lambda_{\text{pos}} \cdot \mathbf{M}_{B^i} \cdot (\text{Max}(\hat{\mathbf{z}}_{t-1}^i) - \mathbf{z}'_{t-1}), \quad (14)$$

$$\mathbf{z}'_{t-1}- = \lambda_{\text{neg}} \cdot \left(\mathbf{M}_{-(\bigcup_{i=0}^{n-1} B^i)} \right) \cdot (\mathbf{z}'_{t-1} - \text{Min}(\hat{\mathbf{z}}_{t-1}^i)), \quad (15)$$

where λ_{pos} and λ_{neg} are both set to 0.2.

The features inside and outside the bounding box have large discrepancies, resulting in an unsmooth excess of the generated image near the bounding box. Therefore, we employ Gaussian filtering to smooth the features near the

bounding box in the latent space. First, we define a two-dimensional Gaussian kernel, denoted as follows:

$$G_\sigma(i, j) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{i^2 + j^2}{2\sigma^2}\right), \quad (16)$$

where $\sigma = 1.0$ is the standard deviation that controls the spread of the Gaussian filtering. Then, assuming an arbitrary feature with coordinates (x, y) near the bounding box, we perform Gaussian filtering on it, denoted as follows:

$$z_{t-1}^{\text{smooth}}(x, y) = \sum_{i=-k}^k \sum_{j=-k}^k z'_{t-1}(x+i, y+j) \cdot G_\sigma(i, j). \quad (17)$$

Furthermore, to preserve the overall semantics and distribution of the image as much as possible, the weighted sum of z_{t-1}^c and z_{t-1}^{smooth} is used to produce the final denoised output that maintains consistency between image and text.

$$z_{t-1} = \mu \cdot z_{t-1}^c + (1 - \mu) \cdot z_{t-1}^{\text{smooth}}, \quad (18)$$

where μ is the trade-off weight set to 0.8.



Figure 3. **Ablation results of MCCD.** The poor results after removing the critical components prove that each component is crucial.



"The soft, warm glow of the campfire illuminated the faces of the hikers, as they roasted marshmallows and swapped stories. A starry sky stretches over a distant mountain range. Nearby, a tent suggests that the group is on a camping adventure."

"The dining area is filled with natural light. A silver refrigerator stands nearby, while a large window brightens the space and frames a lush green plant sitting just in front of it. The dining area features a wooden table, and a rug adds warmth to the room, creating a cozy atmosphere."

Figure 4. **Additional qualitative results.** MCCD can handle complex text prompts with multiple objects and attribute binding relationships effectively, generating reasonable bounding box layouts and producing aesthetically pleasing images with high realism.

4. Experiments

4.1. Datasets and Evaluation Metrics

Since our framework is training-free, no dataset is required for training. To comprehensively evaluate the performance of different methods, we selected six metrics from three categories in the T2I-CompBench benchmark [15], namely

Attribute Binding (Color, Shape, Texture), Object Relationship (Spatial, Non-Spatial), and Complex. The three metrics of Attribute Binding are evaluated by disentangled BLIP-VQA. The spatial metric is determined by detecting the object through UniDet and comparing the center of the object's bounding box in the image to determine the spatial relationship. CLIP-Score is applied for the evaluation

of the Non-Spatial metric. The Complex metric is the average score of disentangled BLIP-VQA, CLIP-Score, and UniDet. Each metric corresponds to 300 prompts, and each prompt yields 10 images for evaluation with diverse seeds.

4.2. Implementation Details

Our MCCD is both generic and extensible, allowing us to seamlessly integrate a wide range of MLLM architectures and diffusion models into the framework. In our experiments, we use GPT-4o-mini [1] to construct multiple agents in MCP. The pre-trained diffusion models in HCD consist of: SDv1.5 [30], SDv2 [30], SDv2-base [30], and SDXL-base [26]. The inference process is set to 20 steps, with the classifier-free guidance scale set to 7.0. We have carefully designed task templates and selected high-quality context examples to ensure optimal performance. All experiments are conducted on a single A800 GPU.

4.3. Main Results

Qualitative Results. Figure 2 shows the results of MCCD on multiple diffusion models. We have the following observations (i) Attribute Binding: Attributes (*e.g.*, color, quantity) are incorrectly bound to objects in the generated results of the base diffusion models, while each attribute is correctly bound to the corresponding object after applying MCCD. (ii) Spatial location relationships are often incorrect in the results generated from the base diffusion models. For example, objects may be incorrectly placed in unnatural locations, or spatial relationships between objects may be confused. The spatial relationship confusion problem is solved after applying MCCD. (iii) Realism and aesthetics: Many of the generated results of the base diffusion models lack suitable backgrounds and detailed texture structures, while MCP provides detailed, realistic, and aesthetically pleasing descriptions of objects and backgrounds for the complex prompts, which makes the results highly aesthetic.

Figure 4 shows the results of complex prompts with multiple objects and attribute binding by MCCD. Through in-depth understanding and parsing of the input complex prompts, MCCD can effectively recognize and extract the objects and their interrelationships in the text prompt and then generate a reasonable and accurate bounding-box layout. The position and size of each bounding box closely correspond to the elements (*e.g.*, objects, colors, shapes, etc.) in the complex prompt, ensuring the accurate positioning and realistic rendering of each element in the image. In addition, MCCD fully captures the connection between the object and its surroundings, resulting in an image that exhibits a high degree of realism and aesthetics. In this way, MCCD can transform complex and layered text prompts into clearly structured and highly realistic images, perfectly combining text and visual expression.

Quantitative Results. Table 1 shows the results of mul-

tiple open source reproducible models and MCCD in T2I-CompBench. We can draw the following conclusions: (i) Overall, the best results are achieved by applying MCCD to SDXL, with an overall metric improvement of 9.04%, which is significantly better than the baseline models. (ii) Compared to other diffusion models, the application of MCCD resulted in the highest overall improvement in metrics of 9.04%, demonstrating the general scalability of MCCD. (iii) Benefiting from MCP, the key components of MCCD, the Spatial metric of the diffusion model is increased by up to 41.49%, and the Complex metric is increased by up to 8.73%, which effectively improves the spatial position relationship, semantic expressiveness, and image fidelity of the generated images.

4.4. Ablation Study

To verify the necessity of each component, we conduct full-scale ablation studies. The results are shown in Figure 3. (i) Firstly, MCP is removed from MCCD. The poor attribute binding and spatial location relationship indicate that it is crucial to dynamically invoke multiple agents to generate reasonable bounding boxes, objects, and background descriptions before image generation. (ii) In addition, we remove the complex prompt in HCD, and the image generation results that are inconsistent with the prompt descriptions suggest that the complex prompt is indispensable for maintaining consistency between image and text. (iii) Then, Regional Enhancement is removed from HCD. The phenomenon of poor semantic representation of the objects and the background in the image suggests that regional enhancement is crucial. (iv) Additionally, we remove Dynamic Integration, and the image appears to be partially missing in the occluded object’s part, presenting an anomalous structure. (v) Finally, we remove Latent Space Smoothness, and the image exhibits unnatural transitions where the bounding box connects to the background, indicating that smoothness is essential for natural transitions between the bounding boxes and between the bounding box and the background in the image.

5. Conclusion

In this paper, we propose Multi-Agent Collaboration-based Compositional Diffusion (MCCD) for generating high-quality complex scenes. Specifically, we design a multi-agent collaboration-based scene parsing module to fully extract the scene elements contained in the text by constructing a multi-intelligentsia system based on MLLMs. In addition, we propose a hierarchical compositional diffusion that utilizes dynamic integration of overlapping regions, regional enhancement, and latent space smoothness to generate realistic and aesthetically pleasing images. Comprehensive experiments prove the advantages of our method.

MCCD: Multi-Agent Collaboration-based Compositional Diffusion for Complex Text-to-Image Generation

Supplementary Material

6. The Prompt design of MCP

We provide detailed descriptions and prompt template implementations for the conductor, evaluator, and all agents. The text enclosed within the curly braces denotes placeholders that will be dynamically populated during runtime based on the input text prompt and the agent's output.

6.1. Conductor

The role of a conductor is highly specialized and significant, which necessitates a more intricate prompt design compared to other agents. The task of the conductor is to coordinate all the agents you manage so that they can work together to solve the problem. The prompt template for a conductor is shown in Table 2.

6.2. Evaluator

The evaluator's task is to exercise critical thinking to assess the truthfulness and reasonableness of the agent's output. If it is not reasonable, then make recommendations for modification. The prompt template for the evaluator is given in Table 3.

6.3. Agent System

In this section, we provide an in-depth overview of the individual agents involved in our MCCD. Each agent is assigned a specific role and domain knowledge related to problem-solving.

Object extraction agent. The object extraction agent's task is to extract key entities and their corresponding characteristics from the text input prompt. The prompt template for the object extraction agent is shown in Table 4.

Background extraction agent. The task of the background extraction agent is to extract the background from this complex prompt. The extracted background is required not to contain any object and its characteristics, but only a description of the scene as a whole. The prompt template for the background extraction agent is shown in Table 5.

Action relations extraction agent. The task of the action relation extraction agent is to extract the action relations between objects in the scene, such as holding, sitting, and so on, to bind multiple objects. Its prompt template is illustrated in Table 6.

Spatial relations extraction agent. The task of the spatial relations extraction agent is to extract the spatial relations between objects in the scene, such as left, beside, and so on, to fully extract the spatial information of the scene. Its prompt template is illustrated in Table 7.

Table 2. The prompt template for the conductor.

Conductor

[INST] <SYS>

You are the leader of an agent system for text parsing in complex scenes. Now, you need to coordinate all the agents you manage so that they can work together to solve the problem. Next, you are given a specific text prompt, and your goal is to select the agents you think are best suited to solicit insights and suggestions. Generally speaking, the parsing of complex scenes includes several processes: object extraction, background extraction, relation extraction, layout extraction, and aesthetic optimization. Different text prompts may correspond to different processes, so you need to select the corresponding agent to solve the problem dynamically. </SYS>

<USER>

The text prompt is: {text prompt}.

Remember, based on the capabilities of different agents and the current status of the problem-solving process, you need to decide which agent to consult next. The agents' capabilities are described as follows: {agent info}.

Agents that have already outputted their answers include: {outputted agents} .

Please select an agent to consult from the remaining agents {remaining agents}.

Remember, the agent must choose from the existing list above.

Note that you must complete the workflow within the remaining {remaining steps} steps.

You should output the name of the agent directly. The next agent is: </USER> [INST]

Layout agent. The task of the layout agent is to layout the scene by generating a bounding box for each object. Its prompt template is illustrated in Table 8.

Aesthetics enhancement agent. The task of the aesthetic enhancement agent is to perform the role of an aesthetic guide to optimize the description of the object characteristics, thus enhancing the artistic and aesthetic qualities of the image. Its prompt template is displayed in Table 9.

7. Case Study of MCP

Figures 5 and 6 show two cases of MCP workflows. As shown in the figures, the conductor organizes multiple

Table 3. The prompt template for the evaluator.

<i>Evaluator</i>
[INST] <SYS>
Your task is to exercise critical thinking to assess the truthfulness and reasonableness of the agent's output. If it is not reasonable, then make recommendations for modification.
Output format: {"Result": "Evaluation results, with a value of right or wrong", "Problem": "If there is a problem, describe it in detail, otherwise, the value is null", "Modification Suggestion": "If the result is incorrect, describe the proposed change, otherwise, the value is null"}. </SYS>
<USER>
The input prompt is described as {input prompt}. The output of all agents is {outputs}.
Please evaluate the reasonableness of the agents' outputs. If they are not reasonable, please state your suggestions for modification.
</USER> [INST]

Table 4. The prompt template for the object extraction agent.

<i>Object extraction agent</i>
[INST] <SYS>
As an object extraction agent, you extract key entities and their corresponding characteristics from the text input prompt.
Extract multiple object and characteristic pairs if multiple characteristics of an entity describe different parts of a person, such as the head, clothes/body, and underwear. To ensure numeric accuracy, objects with the same class name (e.g., five apples) will be separately assigned to different regions.
The output format is {object ₁ :characteristic ₁ , object ₂ :characteristic ₂ , ..., object _n :characteristic _n } </SYS>
<USER>
The input prompt is described as: {input prompt}. You are supposed to refer to the output of other agents: {outputs}.
Please output the extracted objects and their characteristics in the text prompt.
</USER> [INST]

agents in an orderly manner according to the input prompt to achieve a well-collaborated result. Figure 5 is a case where each agent correctly outputs the result. When a problem occurs during the system's execution of a task, the Evaluator can identify it through a backward feedback process and suggest modifications to resolve it quickly. For example, in the case of Figure 6, when the behavior of a certain agent does not meet the expectation, the Evaluator analyzes that the problem occurs in the layout agent and suggests a modification. The backward feedback process identifies the agent's error and regenerates the prompt set. This mechanism of feedback and correction provides a high degree of flexibility and self-adaptation for the system, thus enhancing the robustness and intelligence of the whole system.

Table 5. The prompt template for the background extraction agent.

<i>Background extraction agent</i>
[INST] <SYS>
As an object extraction agent, your task is to extract the background from this complex prompt. The extracted background is required not to contain any object and its characteristics, but only a description of the scene as a whole.
</SYS>
<USER>
The input prompt is described as: {input prompt}. The outputs given by other agents are as follows: {outputs}, please refer to them carefully.
Please extract and output the background in the text prompt.
</USER> [INST]

Table 6. The prompt template for the action relations extraction agent.

<i>Action relation extraction agent</i>
[INST] <SYS>
As an action relation extraction agent, your task is to extract the action relations between objects in the scene, such as holding, sitting, and so on, to fully extract the spatial information of the scene.
The output format is {(object ₁ , action relation ₁ , object ₂), ..., (object _n , action relation _n , object _m)} </SYS>
<USER>
The input prompt is described as: {input prompt}. The outputs given by other agents are as follows: {outputs}.
Please extract and output the action relations between objects in the text prompt.
</USER> [INST]

uator can identify it through a backward feedback process and suggest modifications to resolve it quickly. For example, in the case of Figure 6, when the behavior of a certain agent does not meet the expectation, the Evaluator analyzes that the problem occurs in the layout agent and suggests a modification. The backward feedback process identifies the agent's error and regenerates the prompt set. This mechanism of feedback and correction provides a high degree of flexibility and self-adaptation for the system, thus enhancing the robustness and intelligence of the whole system.

Table 7. The prompt template for the spatial relations extraction agent.

<i>Spatial relation extraction agent</i>
[INST] <SYS>
As a spatial relation extraction agent, your task is to extract the spatial relations between objects in the scene, such as “left” and “beside”, to fully capture the spatial information of the scene.
The output format is $\{(object_1, \text{ spatial relation}_1, object_2), \dots, (object_n, \text{ spatial relation}_n, object_n)\}$
</SYS>
<USER>
The input prompt is described as: {input prompt}. You should refer to the output of other agents: {outputs}.
Please extract and output the spatial relations between objects in the text prompt.
</USER> [INST]

Table 8. The prompt template for the layout agent.

<i>Layout agent</i>
[INST] <SYS>
You are a layout agent, and your task is to generate the bounding boxes for the objects. The following rules must be strictly followed during generation. A layout denotes a set of “object: bounding box” items. “object” means any object name, which starts the object name with “a” or “an” if possible. “bounding box” is formulated as $[x, y, w, h, d]$, where “ x, y ” denotes the top left coordinate of the bounding box, “ w ” denotes the width, “ h ” denotes the height, and “ d ” indicates the order of the object’s front and back position in the image, starting from 0, the smaller d indicates the object is more forward. The top-left corner has coordinates $[0, 0]$. The bottom-right corner has coordinates $[1, 1]$.
The output format: “ $\{object_1: \text{ bounding box}_1, object_2: \text{ bounding box}_2, \dots, object_n: \text{ bounding box}_n\}$ ”
</SYS>
<USER>
The input prompt is described as: {input prompt}. You should refer to the output of other agents: {outputs}.
Please generate and layout according to the task description and rules.
</USER> [INST]

Table 9. The prompt template for the layout agent.

<i>Aesthetics enhancement agent</i>
[INST] <SYS>
As an aesthetic enhancement agent, you serve as an aesthetic guide, refining the descriptions of an object’s characteristics to amplify its artistic and aesthetic appeal. This involves thoughtful consideration of composition, color balance, texture, and other key elements contributing to the image’s overall impact.
</SYS>
<USER>
The input prompt is: {input prompt}. You should refer to the output of other agents: {outputs}.
Please generate a glorified characterization as required.
</USER> [INST]

8. Additional Ablation Results of MCP

To demonstrate the scalability and generalizability of MCCD, we perform comprehensive ablation experiments on MLLMs used in MCP based on the SDXL-Base model. We select GPT-4o-mini [1], GPT-4o [1], and LLaVa-1.5-7b [21] as the MLLMs. The qualitative results are shown in Figure 7. According to the experimental results, compared to the errors in object attribute binding and quantity generation of the SDXL-Base model, MCCD can generate correct image content with high fidelity and aesthetic quality after utilizing different MLLMs. Specifically, these models can accurately recognize and bind the attributes of objects while exhibiting high stability and consistency in the number and layout of objects, which significantly improves the quality and usability of the generated results. This suggests that MCCD can fully utilize the potential of these advanced MLLMs to generate images with high visual realism and artistic value in complex scenes.

9. Additional Qualitative Analysis

Figure 8 shows the results of complex text prompts with multiple objects and attribute bindings by MCCD. Through an in-depth understanding and parsing of the input detailed text description, MCCD can effectively recognize and extract the objects and their interrelationships in the text, generating a reasonable and accurate bounding box layout. The position and size of each bounding box closely correspond to the elements (e.g., objects, colors, shapes, etc.) in the textual descriptions, ensuring the accurate positioning and realistic rendering of each element in the image. In addition, MCCD not only focuses on the physical positional relationship of the object but also fully captures the interaction be-

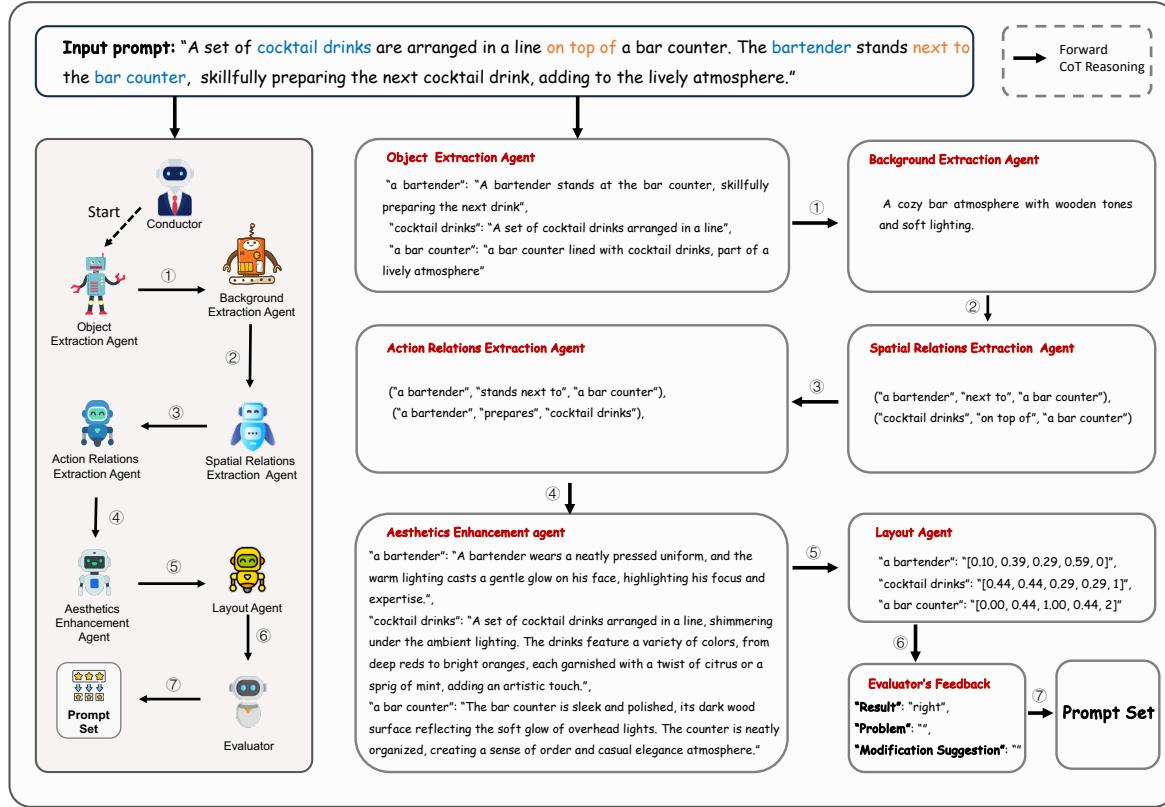


Figure 5. A case to illustrate the workflow of MCP.

tween the object and its surroundings, resulting in an image that exhibits a high degree of realism and aesthetics. In this way, MCCD can transform complex and layered textual information into clearly structured and visually rich images, perfectly combining text and visual expression. MCCD can generate reasonable bounding box layouts, resulting in attractive images with high realism.

10. Discussion of Broader Impact

MCCD, as a training-free T2I approach with excellent performance, is capable of transforming complex and layered text information into clearly structured and visually rich images, perfectly combining text and visual representation. However, it also has potential negative impacts. For example, it may be used to generate scenarios involving immoral or illegal activities that can harm society. Additionally, automatically generated images may raise concerns about intellectual property and copyright.

11. Discussion of Limitations and Future Work

The proposed MCCD serves as a training-free plug-in that can be adapted to any Diffusion-based T2I methods, using a hierarchical compositional generative paradigm to enhance

the quality of complex scene generation for the model. A potential problem is that the MCCD inference overhead is somewhat affected by the number of objects in a complex text prompt. As the number of bounding boxes increases, the inference time also increases. In the future, we will focus on designing efficient performance optimization strategies to improve the inference speed of the method.

Acknowledgements

This work is supported in part by the Shanghai Municipal Science and Technology Committee of Shanghai Outstanding Academic Leaders Plan (No. 21XD1430300), and in part by the National Key R&D Program of China (No. 2021ZD0113503).

References

- [1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023. 8, 11
- [2] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce

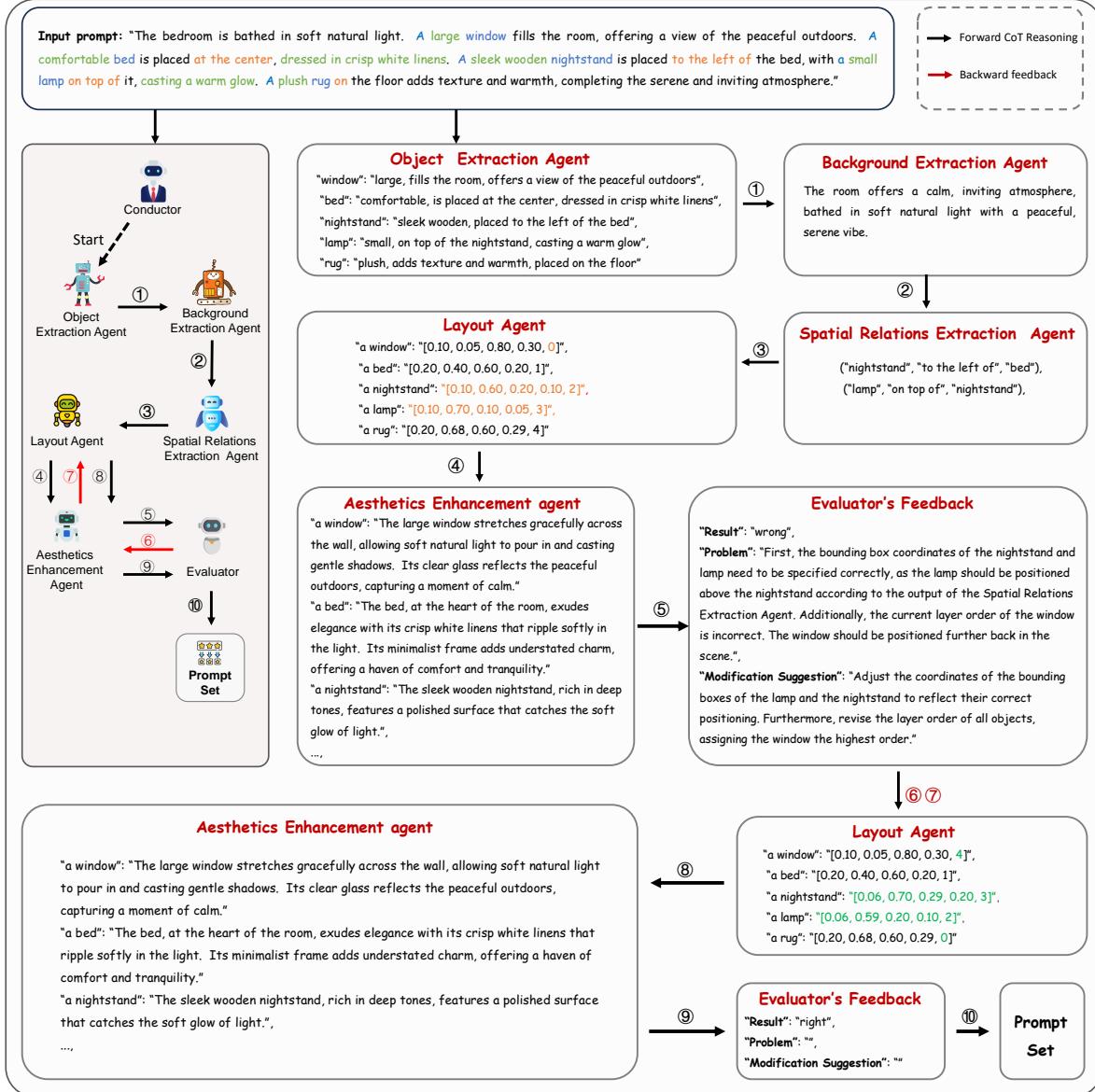


Figure 6. A case to illustrate the workflow of MCP. In the agents' outputs, orange text indicates errors, and green text indicates corrections.

- Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023. 1, 2, 6
- [3] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Transactions on Graphics (TOG)*, 42(4):1–10, 2023. 6
- [4] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart-alpha: Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 6
- [5] Wei-Ge Chen, Irina Spiridonova, Jianwei Yang, Jianfeng Gao, and Chunyuan Li. Llava-interactive: An all-in-one demo for image chat, segmentation, generation and editing. *arXiv preprint arXiv:2311.00571*, 2023. 2
- [6] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023. 2
- [7] Prafulla Dharwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1, 2
- [8] Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng,

"A woman in a pink shirt and jeans holds a white umbrella in the rain."



"A glass vase and a metallic watering can are placed beside each other, both filled with colorful flowers."

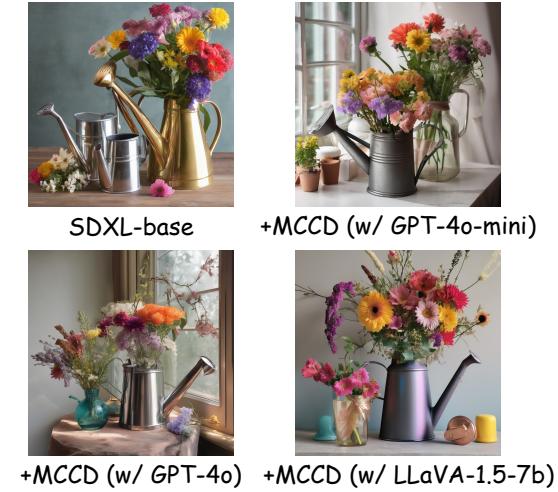
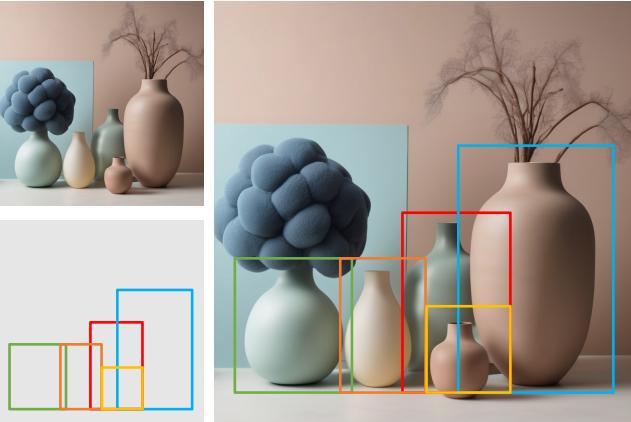
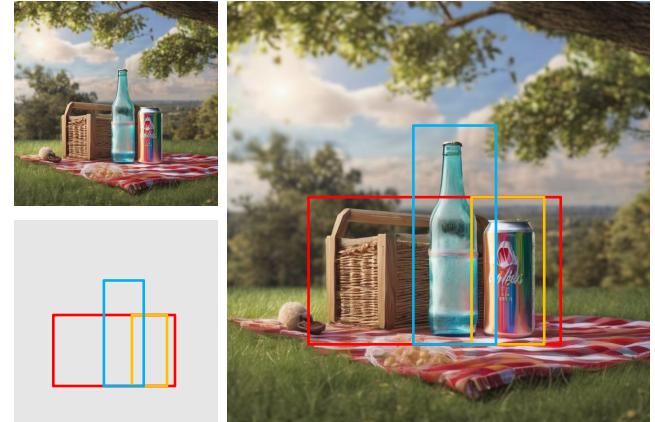


Figure 7. Ablation results of MLLMs in MCP.



"Five uniquely colored and shaped vases are arranged on a white surface, with two in soft pink, two in subtle green, and one in warm white."



"A glass water bottle is positioned at the center, with an aluminum can to its right. Behind them, a wooden basket adds a rustic touch, enhancing the natural setting. They are set on a lush, green grassy field, creating a serene and peaceful atmosphere."

Figure 8. Additional qualitative analysis.

Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, et al. Cogview: Mastering text-to-image generation via transformers. *Advances in neural information processing systems*, 34:19822–19835, 2021. 2

[9] Guian Fang, Zutao Jiang, Jianhua Han, Guangsong Lu, Hang Xu, and Xiaodan Liang. Boosting text-to-image diffusion models with fine-grained semantic rewards. *arXiv preprint arXiv:2305.19599*, 5, 2023. 1

[10] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022. 1, 6

- and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*, 2022. 1, 6
- [11] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large language models. *Advances in Neural Information Processing Systems*, 36, 2024. 2
- [12] Hanan Gani, Shariq Farooq Bhat, Muzammal Naseer, Salman Khan, and Peter Wonka. Llm blueprint: Enabling

- text-to-image generation with complex and detailed prompts. *arXiv preprint arXiv:2310.10640*, 2023. 2
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [14] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment. *arXiv preprint arXiv:2403.05135*, 2024. 2
- [15] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-combench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36:78723–78747, 2023. 1, 6, 7
- [16] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback. *arXiv preprint arXiv:2302.12192*, 2023.
- [17] Seung Hyun Lee, Yinxiao Li, Junjie Ke, Innfarm Yoo, Han Zhang, Jiahui Yu, Qifei Wang, Fei Deng, Glenn Entis, Junfeng He, et al. Parrot: Pareto-optimal multi-reward reinforcement learning framework for text-to-image generation. In *European Conference on Computer Vision*, pages 462–478. Springer, 2025. 1
- [18] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. In *International conference on machine learning*, pages 19730–19742. PMLR, 2023. 2
- [19] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22511–22521, 2023. 1, 2
- [20] Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models. *arXiv preprint arXiv:2305.13655*, 2023. 1, 2
- [21] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 26296–26306, 2024. 11
- [22] Mushui Liu, Yuhang Ma, Yang Zhen, Jun Dan, Yunlong Yu, Zeng Zhao, Zhipeng Hu, Bai Liu, and Changjie Fan. Llm4gen: Leveraging semantic representation of llms for text-to-image generation. *arXiv preprint arXiv:2407.00737*, 2024. 2
- [23] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*, pages 423–439. Springer, 2022. 1, 6
- [24] Chong Mou, Xiantao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 4296–4304, 2024. 2
- [25] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 2
- [26] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 2, 8
- [27] Leigang Qu, Shengqiong Wu, Hao Fei, Liqiang Nie, and Tat-Seng Chua. Layoutllm-t2i: Eliciting layout guidance from llm for text-to-image generation. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 643–654, 2023. 1, 2
- [28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2
- [29] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022. 1, 2, 6
- [30] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1, 2, 6, 8
- [31] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in neural information processing systems*, 35:36479–36494, 2022. 1
- [32] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 1, 2
- [33] Jiao Sun, Deqing Fu, Yushi Hu, Su Wang, Royi Rassin, Da-Cheng Juan, Dana Alon, Charles Herrmann, Sjoerd van Steenkiste, Ranjay Krishna, et al. Dreamsync: Aligning text-to-image generation with image understanding feedback. In *Synthetic Data for Computer Vision Workshop@ CVPR 2024*, 2023. 1
- [34] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerzel, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022. 2
- [35] Jinheng Xie, Yuexiang Li, Yawen Huang, Haozhe Liu, Wentian Zhang, Yefeng Zheng, and Mike Zheng Shou. Boxdiff: Text-to-image synthesis with training-free box-constrained diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7452–7461, 2023. 4

- [36] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imageward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36, 2024. [1](#)
- [37] Aiyuan Yang, Bin Xiao, Bingning Wang, Borong Zhang, Ce Bian, Chao Yin, Chenxu Lv, Da Pan, Dian Wang, Dong Yan, et al. Baichuan 2: Open large-scale language models. *arXiv preprint arXiv:2309.10305*, 2023. [2](#)
- [38] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *ACM Computing Surveys*, 56(4): 1–39, 2023. [1](#), [2](#)
- [39] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and CUI Bin. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multi-modal llms. In *Forty-first International Conference on Machine Learning*, 2024. [2](#), [3](#), [4](#)
- [40] Ling Yang, Zhilong Zhang, Zhaochen Yu, Jingwei Liu, Minkai Xu, Stefano Ermon, and CUI Bin. Cross-modal contextualized diffusion models for text-guided visual generation and editing. In *The Twelfth International Conference on Learning Representations*, 2024. [2](#)
- [41] Zhengyuan Yang, Jianfeng Wang, Zhe Gan, Linjie Li, Kevin Lin, Chenfei Wu, Nan Duan, Zicheng Liu, Ce Liu, Michael Zeng, et al. Reco: Region-controlled text-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14246–14255, 2023. [1](#), [2](#)
- [42] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 3836–3847, 2023. [2](#)
- [43] Xinchen Zhang, Ling Yang, Yaqi Cai, Zhaochen Yu, Jiake Xie, Ye Tian, Minkai Xu, Yong Tang, Yujiu Yang, and Bin Cui. Realcompo: Dynamic equilibrium between realism and compositionality improves text-to-image diffusion models. *arXiv preprint arXiv:2402.12908*, 2024. [1](#)
- [44] Shanshan Zhong, Zhongzhan Huang, Weushao Wen, Jinghui Qin, and Liang Lin. Sur-adapter: Enhancing text-to-image pre-trained diffusion models with large language models. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 567–578, 2023. [2](#)
- [45] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. [2](#)