

# Draw Like an Artist: Complex Scene Generation with Diffusion Model via Composition, Painting, and Retouching

Minghao Liu<sup>1</sup>, Le Zhang<sup>2</sup>, Yingjie Tian<sup>1\*</sup>, Xiaochao Qu<sup>2</sup>, Luoqi Liu<sup>2</sup>, Ting Liu<sup>2\*</sup>

<sup>1</sup>University of Chinese Academy of Sciences, Beijing, China

<sup>2</sup>MT Lab, Meitu Inc., Beijing, China

## Abstract

Recent advances in text-to-image diffusion models have demonstrated impressive capabilities in image quality. However, complex scene generation remains relatively unexplored, and even the definition of ‘complex scene’ itself remains unclear. In this paper, we address this gap by providing a precise definition of complex scenes and introducing a set of **Complex Decomposition Criteria (CDC)** based on this definition. Inspired by the artist’s painting process, we propose a training-free diffusion framework called **Complex Diffusion (CxD)**, which divides the process into three stages: composition, painting and retouching. Our method leverages the powerful chain-of-thought capabilities of large language models (LLMs) to decompose complex prompts based on CDC and to manage composition and layout. We then develop an attention modulation method that guides simple prompts to specific regions to complete the complex scene painting. Finally, we inject the detailed output of the LLM into a retouching model to enhance the image details, thus implementing the retouching stage. Extensive experiments demonstrate that our method outperforms previous SOTA approaches, significantly improving the generation of high-quality, semantically consistent, and visually diverse images for complex scenes, even with intricate prompts.

## 1 Introduction

Recently, diffusion models [33; 34; 43; 18] have represented a significant advancement in text-to-image generation, showcasing impressive capabilities in producing high-quality images from textual descriptions. However, despite their remarkable performance, these models face substantial challenges when tasked with generating complex scenes. Specifically, as illustrated in Figure 1, when prompts involve multiple entities, intricate spatial positions, and conflicting relationships, the models often encounter issues such as entity

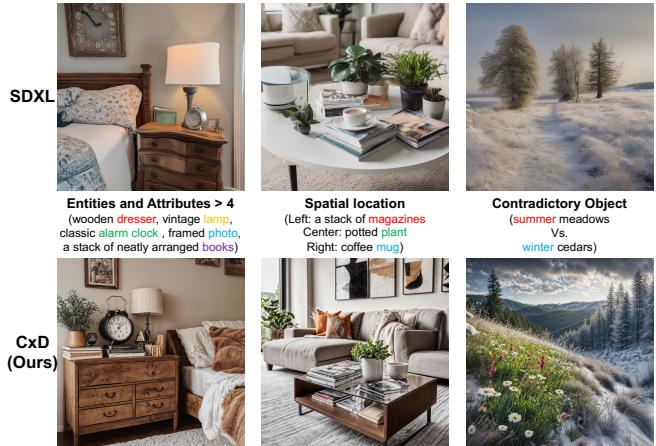


Figure 1: Limitations of pre-trained diffusion models in complex scene generation.

omissions, spatial inconsistencies, and overall disharmony in the generated images.

Several works [23; 39; 29; 5; 40] have sought to address these challenges by incorporating additional layouts or boxes to decompose complex scene relationships. For instance, LAW-diffusion [40] integrates layout configurations into the synthesis process to enhance the coherence of object relationships in complex scenes. Similarly, some approaches [19; 36; 42] utilize prompt-aware attention guidance to improve compositional text-to-image synthesis. Dense Diffusion [19], for example, adjusts intermediate attention maps based on layout conditions to support detailed captions, while RPG [42] employs MLLM as a global planner to decompose complex images into simpler tasks within subregions, introducing complementary regional diffusion for region-specific compositional generation. Despite these advancements, gaps remain when dealing with highly complex scene prompts, and the definition of a ‘complex scene’ continues to be somewhat ambiguous.

In parallel, traditional art processes for creating complex scenes involve a meticulous three-stage approach: composition, painting, and retouching [1; 8; 12; 32]. Artists begin by sketching the overall layout and positioning of elements (composition), followed by detailed painting where the

\*Corresponding authors.

main features are developed (painting), and finally, they refine the artwork by adding details and correcting imperfections (retouching). This methodical approach ensures that every aspect of the scene is carefully considered and harmonized. Adapting this artistic process to model-based generation might provide a solution for complex scene generation.

To address these issues, we propose a novel training-free diffusion framework called **Complex Diffusion (CxD)**, which draws inspiration from the artistic creation process and divides the scene generation into three stages: composition, painting, and retouching. We start by leveraging the powerful reasoning abilities of chain-of-thought in a Large Language Model (LLM) to decompose and compose complex prompts according to **Complex Decomposition Criteria (CDC)**. Next, CxD adapts the cross-attention layer to handle the simplified prompts and compositions derived from the LLM decomposition. Finally, we use a retouching model (ControlNet[46]) to enhance the details of complex scenes based on the attributes obtained during the LLM decomposition process. Comprehensive experiments across multiple tasks demonstrate that our proposed method consistently outperforms previous SOTA approaches, achieving significant improvements in generating high-quality, semantically consistent, and visually diverse images for complex scenes, even when conditioned on intricate textual prompts. In summary, our contributions can be outlined as follows:

- **Definition and Criteria.** We provide a clear experimental definition of complex scenes and introduce **Complexity Decomposition Criteria (CDC)** to effectively manage complex prompts.
- **CxD Framework:** Drawing inspiration from the artistic creation process, we propose a training-free **Complex Diffusion (CxD)** framework that divides the generation of complex scene images into three stages: composition, painting, and retouching.
- **Validation and Performance.** Extensive experiments demonstrate that CxD generates high-quality, consistent, and diverse images of complex scenes, even when dealing with intricate prompts

## 2 Related Work

### 2.1 Complex Sense Generation

Text-to-image(T2I) generation has made significant strides, but generating complex scenes remains challenging. Transformer-based models [45; 18] have improved spatial layout modeling and data efficiency, enhancing texture, structure, and relational accuracy. GAN-based models [13] leverage a generator-discriminator framework to produce images aligned with textual descriptions, but complex scenes with multiple object categories demand a highly proficient discriminator [22; 15; 21]. Recently, diffusion models [7; 14; 31] have gained popularity over GANs for their high-quality, diverse outputs. Notably, LAW-Diffusion [40] enhances object relationship coherence in complex scenes by integrating layout configurations into the synthesis process.

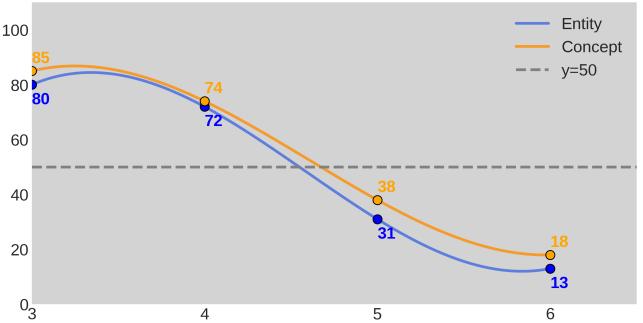


Figure 2: Performance trends of the SD XL model with varying numbers of entities and concepts

### 2.2 Compositional Diffusion Generation

Composition is vital in painting, and recent studies have explored compositional diffusion models to aid this process[23; 44; 17]. Some methods, like ControlNet [46] and T2I-Adapter [26], enhance controllability by training additional modules, but this adds extra costs. Others, like Composable Diffusion [25] MultiDiffusion [2], and Dense Diffusion [19], manipulate latent or cross-attention maps to control the model without extra training, often using bounding boxes to guide composition [5; 39]. However, these approaches struggle with complex prompts, leading to incomplete or distorted images. To address this, we propose an efficient, training-free method that maintains image controllability for complex prompts without additional costs.

### 2.3 LLM for Image generation

Large language models (LLMs) have revolutionized AI research, exemplified by models like ChatGPT [37] with strong language comprehension and reasoning. Leveraging LLMs, diffusion models improve text-image alignment and image quality [38; 44; 27]. LLMs can control layout by assigning locations based on prompts [6; 24], and models like Layout-GPT [10] enhance this by providing retrieval samples. Ranni [11] adds a semantic panel with multiple attributes, while RPG [42] uses LLMs for image composition planning. However, most methods rely on LLMs' inherent abilities or simple prompts. In contrast, we introduce complex decomposition criteria (CDC) to guide LLMs in helping diffusion models understand complex text prompts.

## 3 Complex Scene in Pre-Trained Diffusion

As AI-generated content (AIGC) progresses, image generation, particularly complex scene generation, has become a significant focus and research challenge [20; 16]. Many studies[40; 42] describe complex scenes using vague terms like “multiple”, “different”, and “diverse”, without clearly defining what constitutes “complexity.” This ambiguity can lead to biased evaluations and ineffective solutions. Thus, a precise definition of “complexity” in scene generation is crucial.

## Inspiration and Theme Stage



The tree has a thick trunk with rugged bark and a wide canopy of vibrant green leaves. Nearby, sunflowers stand tall with large, bright yellow petals and dark brown centers. Above this picturesque landscape, two colorful hot air balloons soar high, their fabric adorned with vibrant patterns in shades of red, blue, yellow, and green. A flock of geese flies in a V-formation across the serene sky, while a family enjoys a picnic on a cozy blanket spread out on lush green grass under the shade of a large, leafy tree. Completing the scene, a breathtaking canyon nestled between majestic mountains features a crystal-clear river winding through the rugged terrain.

## Composition and Layout Stage

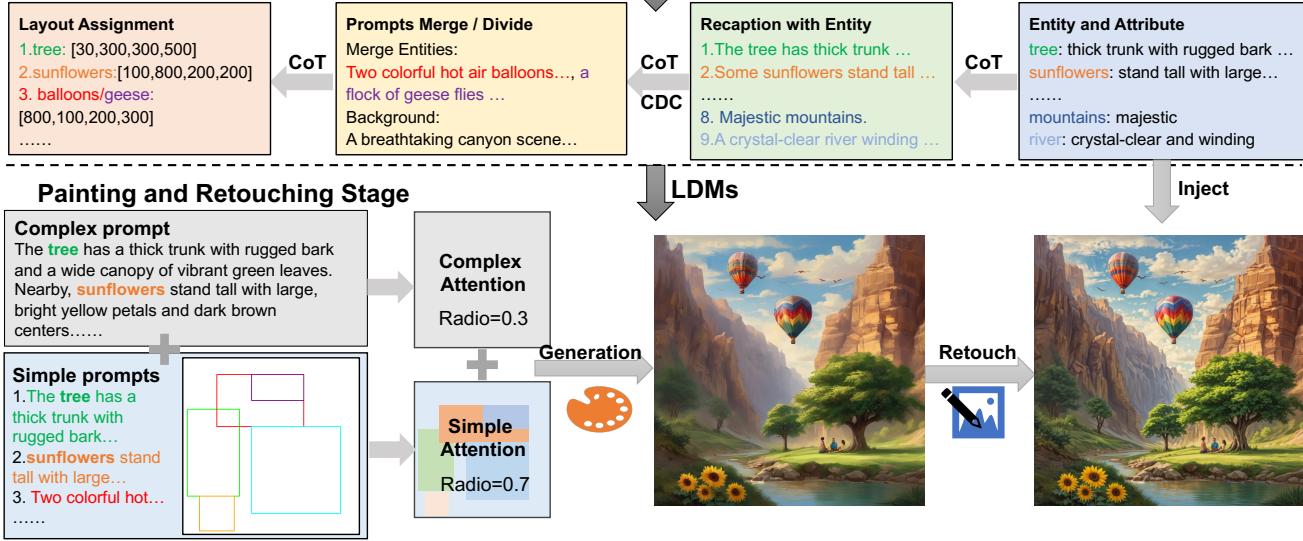


Figure 3: Overview of CxD framework for complex scene image generation.

### 3.1 Definition of ‘Complex’ in Pre-Trained Diffusion

Although some studies attempt to define “complexity” in terms of the number of entities, their attributes, and relationships [20; 16; 40; 42], our experiments suggest these factors must be considered collectively. Generating multiple entities alone is not as challenging for diffusion models as combining these entities with specific attributes and relationships, which often results in omissions, inconsistencies, and disharmony.

Building on previous research[16], we redefine complexity in scene generation by considering four key factors: the number of entities, attributes, spatial positioning, and relationships among entities. Notably, relationships here refer to associations or conflicts rather than spatial arrangements. Our findings indicate that prompts containing conflicting entities are particularly challenging for diffusion models, often leading to inconsistent and visually unsatisfactory results, as shown in Figure 1.

First, we tested diffusion models with prompts containing three to six positively correlated entities. As shown in Figure 2, when the number of entities reached five, nearly 70% of the images showed omissions or disharmony, worsening with six entities. Next, we examined the effect of adding attributes. Next, we examined the impact of adding attributes. When the total number of entities and attributes was four, the results were similar to those with four entities alone, achieving a 74% success rate. However, when the combined number increased to five, the success rate dropped to 38%, indicating that the model struggles with prompts involving more

than four concepts. (entities plus attributes). Lastly, we explored spatial positioning and relationships. Even with two entities, accurately capturing spatial relationships proved difficult, with success rates falling below 50%. For conflicting relationships (e.g., desert vs. rainforest and summer meadows vs. winter cedars), the success rate plummeted to 10%.

In summary, scenes with fewer than five concepts and no spatial or conflicting relationships are simple, while those exceeding these criteria are considered complex.

### 3.2 Complex scene Decomposition Criteria (CDC)

Based on our definition, we propose Complex Scene Decomposition Criteria (**CDC**) to help both humans and LLMs simplify complex prompts:

- **Identify Conflicting Entities:** If conflicting entities are present, classify the scene as complex and separate these entities into different prompts.
- **Check for Spatial Relationships:** If spatial relationships are involved, classify the scene as complex and split entities with positional relationships into different prompts, maintaining their spatial context.
- **Evaluate Number of Concepts:** If a prompt contains more than four concepts, decompose it by entities, ensuring each prompt retains as many attributes as possible without exceeding four per entity.

This approach ensures that complex prompts are effectively simplified for more accurate and consistent image generation.

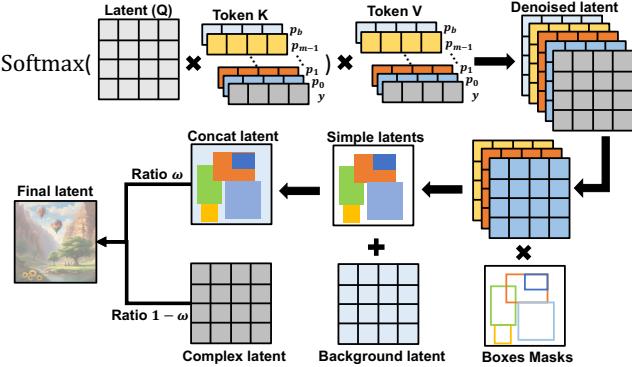


Figure 4: The demonstration of each sampling step in CxD.

## 4 Method:CxD

In this section, we present our training-free framework, **CxD**, which mirrors the artist’s drawing process by dividing complex scene generation into three stages: Composition, Painting, and Retouching, as shown in Figure 3. Starting with a complex scene prompt, we use the Chain-of-Thought (CoT) approach in Large Language Models (LLMs) for composition. The LLM extracts entities and attributes, rephrases with entities, merges them, divides the background based on Complex Decomposition Criteria (CDC), and assigns layouts. CxD then computes and combines complex and simple cross-attention maps at each sampling step. Finally, attributes extracted by the LLM are injected into ControlNet tile [46] for detailed retouching. Below, we detail the methods used in these three stages.

### 4.1 Composition and layout generation with LLMs

#### Entities extraction

Upon receiving a complex scene prompt  $y$  from the user, we leverage the advanced language understanding and reasoning capabilities of the LLM to extract the entities  $E$  and corresponding attributes  $A$  from the prompt. This process can be described as follow:

$$\{E_i\}_{i=0}^n = E_0, E_1, \dots, E_n = LLM_{extract\_enti}(y) \quad (1)$$

$$\begin{aligned} \{A_i\}_{i=0}^n &= \{A_0, A_1, \dots, A_n\} = \{(a_0^0, \dots, a_0^j), \dots, (a_0^n, \dots, a_n^k)\} \\ &= LLM_{extract\_attr}(y) \end{aligned} \quad (2)$$

where  $n$  denotes the number of entities in the complex scene prompt, with  $E_i$  representing the  $i$ -th entity.  $A_i$  denotes the set of attributes corresponding  $E_i$ , and  $a_i^j$  refers to the  $j$ -th attribute of the  $i$ -th entity. It is important to note that the number of attributes for different entities is not necessarily equal, so  $j$  may differ from  $k$ . Furthermore, the number of entities  $n$  and the attributes for each entity  $j$  are not predefined hyperparameters but are determined dynamically through LLM heuristics.

#### Prompts reception

Inspired by RPG [42], which utilizes LLMs to recaption prompts and plan for region divisions with chain-of-thought

(CoT). We also employ LLM to recaption prompt into sub-prompts based on the extracted entities  $E$  and corresponding attributes  $A$ . These sub-prompts are designed to be as consistent as possible with the relevant description in the original complex prompt. This process can be denoted as:

$$\{\hat{y}_i\}_{i=0}^n = \{\hat{y}_0, \dots, \hat{y}_n\} = LLM_{recaption}(E_i, A_i, y) \quad (3)$$

where  $\hat{y}_i$  represents the sub-prompt describing the  $i$ -th entity, with each entity corresponding to a distinct sub-prompt.

#### Prompts merge or divide

After recaptioning, the sub-prompts have been simplified a lot compared to the original complex prompt. However, we cannot guarantee that all sub-prompts will be sufficiently simple for the generative model, as some may still be relatively complex. In addition, some sub-prompts may be very simple on their own, even when combined, the overall prompt might still be straightforward for the generative model. To ensure the quality and efficiency of image generation, we use LLM to merge or split sub-prompts based on the Complex Decomposition Criteria (CDC). The results of merging or splitting are documented as simple prompts:

$$\{p_i\}_{i=0}^m = \{p_0, p_1, \dots, p_m\} = LLM_{md}(\hat{y}) \quad (4)$$

where  $m$  is the number of simple prompts,  $m \leq n$ .  $LLM_{md}$  denotes the operation of merging and dividing prompt.

For all simple prompts, we instruct the LLM to filter out background prompt  $p_b$ , which is image background and do not participate in layout assignment.

#### Layout assignment

Except for the background prompt, all simple prompts are assigned layouts by the LLM to complete the final composition. These layouts are specified with coordinates in the  $(x, y, width, height)$  format. Specifically, the LLM prioritizes the positional and size relationships between different entities mentioned in the prompt to assign bounding boxes. For entities without specified relationships, the LLM uses its common sense knowledge to determine appropriate bounding boxes. The bounding boxes are denoted as  $\{B_i\}_{i=0}^{m-1} = \{B_0, B_1, \dots, B_{m-1}\}$ . The area not covered by these layouts is reserved for the background prompt  $\{BG\}$ .

Finally, we sort the layouts assigned by the LLM in descending order of their area size and adjust the order of the corresponding simple prompts accordingly. This approach aligns with the artist’s practice of prioritizing the main subject first and also helps prevent smaller objects from being obscured by larger ones when entities overlap during image generation.

### 4.2 Cross-Attention Modulation

As analyzed in the previous section, diffusion models tend to be less effective with complex scenarios involving more than four concepts. To address this challenge, we modulate the cross-attention to adapt to composition generated by LLM, facilitating efficient handling of complex scene prompts, as shown in Figure 4.



A sprawling cyberpunk metropolis at night, where towering skyscrapers are adorned with intricate, animated billboards and holographic displays. The streets are a labyrinth of elevated walkways and neon-lit alleyways, bustling with diverse crowds in futuristic attire. Above, a lot of flying vehicles weaves through the sky, casting dynamic shadows on the sleek, reflective surfaces below.



An opulent underwater restaurant inside a giant coral reef. The restaurant has panoramic glass walls for a 360-degree view of vibrant marine life and colorful corals. Elegant tables with fine dining setups and soft, romantic lighting create an inviting atmosphere. A chandelier made of bioluminescent sea anemones illuminates the room. The reef outside is bustling with marine species like schools of fish.



A vibrant tropical island with a white sand beach and turquoise water. Palm trees sway in the breeze, and colorful birds are seen in the lush jungle. A waterfall cascades into a serene lagoon, surrounded by rocks and plants. A thatched-roof hut on stilts overlooks the scene, with a bright sky and setting sun in the background.

Figure 5: Qualitative comparison between CxD and SOTA text-to-image models

### Prompt batch process

Since the complex scene prompt is decomposed into various simple prompts by the LLM, we sample the same latent  $z_t$  as the query at each timestep to ensure image consistency. Additionally, we construct a prompt batch consisting of the complex prompt  $y$ , simple prompts  $\{p_i\}_{i=0}^{m-1}$  and background prompt  $p_b$ . At each timestep, this prompt batch is fed into the denoising network, which manipulates the cross-attention layers to generate different latents in parallel: the complex latent  $z^c$ , simple latents  $\{z^i\}_{i=0}^{m-1}$  and background latent  $z^b$ . The process can be formulated as follows:

$$z_{t-1}^i = \text{Softmax}\left(\frac{(W_Q \cdot \phi(z_t))(W_K \cdot \psi(p_i))}{\sqrt{d}}\right)(W_V \cdot \psi(p_i)) \quad (5)$$

where  $W_Q$ ,  $W_K$ ,  $W_V$  are linear projections, and  $d$  is the latent projection dimension of the keys and queries.  $\phi(z_t)$  and  $\psi(p_i)$  denote the embeddings of the latent  $z_t$  and the simple prompt  $p_i$  at the  $t$ -th timestep, respectively.

### Cross-attention enhancement modulation

To avoid missing concepts involving numerous entities and attributes and to enhance details, we propose attention enhancement modulation. Specifically, we resize the bounding boxes generated by the LLM to match the shape of the latent

$z_{t-1}^i$  as follows:

$$\hat{B} = \text{Resize}(B_i, l_{scale}) \quad (6)$$

where  $l_{scale}$  represents the scaling factor for the latent.

For each latent, we emphasize the region corresponding to bounding box  $\hat{B}_i$  while minimizing the influence of the surrounding area. This process can be described as:

$$z_{t-1}^i + = \lambda_{pos} \cdot M_{\hat{B}}^i \cdot (\text{Max}(z_{t-1}^i) - z_{t-1}^i) \quad (7)$$

$$z_{t-1}^i - = \lambda_{neg} \cdot (\sim M_{\hat{B}}^i) \cdot (z_{t-1}^i - \text{Min}(z_{t-1}^i)) \quad (8)$$

where  $\lambda_{pos}$  and  $\lambda_{neg}$  are hyperparameters that control the extent of modulation.  $M_{\hat{B}}$  and  $(\sim M_{\hat{B}}^i)$  are masks indicating the regions within and outside the bounding box  $\hat{B}^i$ , respectively.

After modulating the results, we concatenate the results of all simple prompts' denoising latents according to the area of the bounding boxes to achieve control over the positional relationship. The areas not covered by the bounding boxes are filled with the background denoising latent results. We define this process as:

$$z_{t-1}^{cat} = \text{Concat}(\{z_{t-1}^i(\hat{B}_i)\}_{i=0}^{m-1}, \{z_{t-1}^b(\sim \hat{B})\}) \quad (9)$$

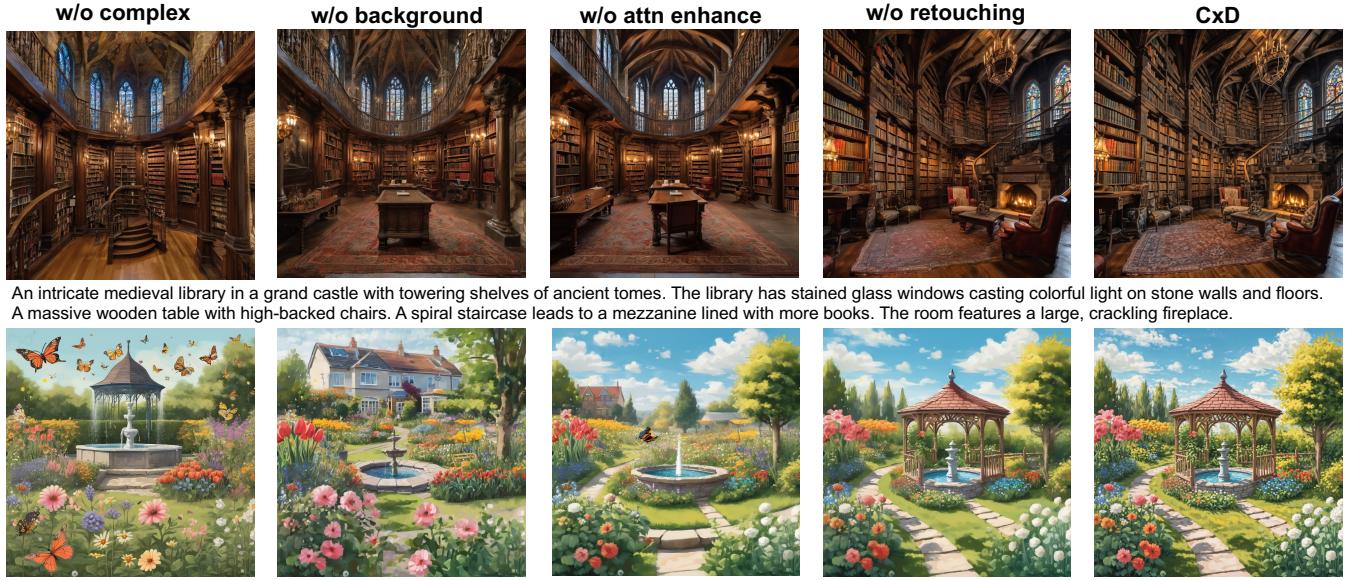


Figure 6: Visualization of CxD ablation study

where  $z_{t-1}^b$  denotes the background denoising latent and  $(\sim \hat{B})$  represents region uncovered by the bounding boxes.

Finally, to ensure a seamless transition between regions and harmonious blending of the background with the entities, we blend the complex prompt latent  $z_{t-1}^c$  with the concatenated latent  $z_{t-1}^{cat}$  using a weighted sum. This approach integrates both elements effectively to produce the final denoised output.

$$z_{t-1} = \omega \cdot z_{t-1}^{cat} + (1 - \omega) \cdot z_{t-1}^c \quad (10)$$

where  $\omega$  is the weight used to balance the contributions of the complex prompt and the simple prompts.

To tackle the challenges of complex scenes, we decompose complex prompts into simpler ones to manage concept overload. Bounding boxes from the LLM help create precise latent representations for each simple prompt, ensuring accurate positional control. Generating each latent independently minimizes conflicts between entities. In summary, CxD effectively addresses the issues related to complex scenes.

### 4.3 Retouching with ControlNet tile Model

Our method effectively generates images that align with the descriptions of complex prompts. However, when the number of entities and attributes exceeds the capacity of the pre-trained diffusion model, some local details unrelated to the complex prompt may be lost or blurred. To address this, we employ retouching models to refine the results, much like an artist adding finishing touches to a painting. We supply the entities and attributes extracted by the LLM as details to a ControlNet [46] extension—ControlNet-tile model—which enhances the image by correcting defects and incorporating new details. After applying ControlNet-tile, the image retains its original semantics but gains enhanced clarity in details and

textures. So far, we have completed the creation of a complex scene image through three stages—composition, painting, and retouching—much like the process an artist would follow.

## 5 Experiment

### 5.1 Experiment Setting

For our CxD framework, we utilize the open-source LLaMA-2 [35] 13B version as our large language model (LLM) and the Stable Diffusion XL [28] version as our pre-trained diffusion model. However, CxD is designed to be a general and extensible framework, capable of integrating various LLM architectures. All experiments in this study were conducted on an NVIDIA RTX 3090 GPU. Generating a complex scene image with CxD takes approximately 2 minutes, including the time required for processing complex prompts with the LLM. We have carefully crafted task-aware templates and high-quality in-context examples to leverage the chain-of-thought (CoT) capabilities of the LLM effectively.

### 5.2 Qualitative Assessment

We evaluated CxD’s performance against various complexity indicators, including the number of concepts, spatial locations, and conflicting relationships. Figure 1 compares results from the SD XL [28] model and CxD. The top row shows SD XL struggling with high complexity, including distortion and inaccuracies in spatial positioning when handling prompts with five entities and attributes. It also tends to ignore one side of entity conflicts. In contrast, CxD effectively manages high complexity, precise spatial arrangements, and conflicting entities, producing consistently harmonious and visually appealing images.

Table 1: Evaluation results on T2I-CompBench. We denote the best score in blue, and the second-best score in green.

Model	Attribute Binding			Object Relationship		Complex
	Color	shape	Texture	Spatial	Non-Spatial	
Stable Diffusion V1.4[31]	0.3765	0.3576	0.4156	0.1246	0.3079	0.3080
Stable Diffusion v2 [31]	0.5065	0.4221	0.4922	0.1342	0.3096	0.3386
Composable Diffusion[25]	0.4063	0.3299	0.3645	0.0800	0.2980	0.2898
Structured Diffusion[9]	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[16]	0.6603	0.4785	0.6287	0.1815	0.3193	0.3328
DALL-E 2[30]	0.5750	0.5464	0.6374	0.1283	0.3043	0.3696
SDXL[28]	0.6369	0.5408	0.5637	0.2032	0.3110	0.4091
PixArt- $\alpha$ [4]	0.6886	0.5582	0.7044	0.2082	0.3179	0.4117
ConPreDiff[41]	0.7019	0.5637	0.7021	0.2362	0.3195	0.4184
RPG[42]	0.8335	0.6801	0.8129	0.4547	0.3462	0.5408
CxD(Ours)	0.8562	0.6533	0.8563	0.6241	0.5426	0.6713

We compared CxD with previous state-of-the-art text-to-image models, including SDXL [28], LDM+ [24], DALLE-3 [22], and RPG [42]. LDM+ and RPG use LLMs for composition assistance. As shown in Figure 5, SDXL and LDM+ struggle with complex prompts, resulting in images that do not fully meet prompt expectations. While DALLE-3 and RPG effectively capture overall content, they sometimes miss local details in complex prompts (e.g., the red part in Figure 5). In contrast, CxD decomposes complex prompt into simple prompts, ensuring no entities or attributes are omitted. Consequently, CxD excels in managing both overall semantics and local details, demonstrating its effectiveness in handling complex scenes.

### 5.3 Quantitative Experiments

We compared our CxD with previous SOTA text-to-image models using the T2I-Compbench benchmark [16]. As shown in Table 1, Our CxD consistently outperforms all others in both general text-to-image generation and complex generation, with RPG coming in second. This highlights the superiority of our approach in handling complex scene generation tasks. We evaluated our CxD model against previous SOTA text-to-image models using the T2I-Compbench benchmark [16]. As shown in Table 1, our model sets a new SOTA benchmark in most tasks, particularly excelling in object relationships and complex scenarios, and significantly outperforms the second-best method. This exceptional performance can be attributed to the strong alignment of these tasks with our proposed Complex Decomposition Criteria (CDC), demonstrating the superiority of our approach in addressing complex scene generation.

### 5.4 Ablation Study

We evaluate each component of our CxD framework: (a) Complex prompt latent, (b) Background prompt latent, (c) Attention enhancement modulation, and (d) Image retouching, as shown in Figure 6. The first column shows images

without complex prompt latent, resulting in disjointed and inconsistent outputs. The second column, lacking background prompt latent, displays backgrounds that do not meet prompt requirements. The third column, without attention enhancement modulation, results in obscured entities. The fourth column, missing modification, produces images with blurred details due to too many entities. The final column shows our CxD framework’s output, preserving semantics and enhancing details, highlighting the importance of each CxD component in generating complex scenes.

## 6 Conclusion

In this paper, we propose CxD, a training-free diffusion framework designed to tackle the challenges of complex scene generation. We define ‘complex scenes’ with precision and provide a set of Complex Decomposition Criteria (CDC) for both humans and large language models (LLMs) to effectively handle complex scene prompts. The CxD framework divides the generation process into three stages—composition, painting, and retouching—mirroring the traditional artist’s approach to drawing. Our experimental results demonstrate that CxD performs well in generating complex scenes. Future work will focus on integrating additional modal data as input conditions to further enhance controllability.

## References

- [1] Rudolf Arnheim. *Art and visual perception: A psychology of the creative eye*. Univ of California Press, 1954.
- [2] Omer Bar-Tal, Lior Yariv, Yaron Lipman, and Tali Dekel. Multidiffusion: Fusing diffusion paths for controlled image generation. 2023.
- [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.
- [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.
- [5] Minghao Chen, Iro Laina, and Andrea Vedaldi. Training-free layout control with cross-attention guidance. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 5343–5353, 2024.
- [6] Jaemin Cho, Abhay Zala, and Mohit Bansal. Visual programming for step-by-step text-to-image generation and evaluation. *Advances in Neural Information Processing Systems*, 36, 2024.
- [7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021.
- [8] Arthur Wesley Dow. *Composition: A series of exercises in art structure for the use of students and teachers*. Univ of California Press, 2023.
- [9] 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.
- [10] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Arjun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutpt: Compositional visual planning and generation with large language models. *Advances in Neural Information Processing Systems*, 36, 2024.
- [11] Yutong Feng, Biao Gong, Di Chen, Yujun Shen, Yu Liu, and Jingren Zhou. Ranni: Taming text-to-image diffusion for accurate instruction following. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4744–4753, 2024.
- [12] Ernst Hans Gombrich and EH Gombrich. *The story of art*, volume 12. Phaidon London, 1995.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- [15] Tianyu Hua, Hongdong Zheng, Yalong Bai, Wei Zhang, Xiao-Ping Zhang, and Tao Mei. Exploiting relationship for complex-scene image generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 1584–1592, 2021.
- [16] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation. *Advances in Neural Information Processing Systems*, 36:78723–78747, 2023.
- [17] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.
- [18] Manuel Jahn, Robin Rombach, and Björn Ommer. High-resolution complex scene synthesis with transformers. *arXiv preprint arXiv:2105.06458*, 2021.
- [19] Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image generation with attention modulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7701–7711, 2023.
- [20] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A Shamma, et al. Visual genome: Connecting language and vision using crowd-sourced dense image annotations. *International journal of computer vision*, 123:32–73, 2017.
- [21] Hanbit Lee, Youna Kim, and Sang-Goo Lee. Multi-scale contrastive learning for complex scene generation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 764–774, 2023.
- [22] Hanbit Lee, Sang-Goo Lee, Jaehui Park, and Junho Shim. Improving complex scene generation by enhancing multi-scale representations of gan discriminators. *IEEE Access*, 11:43067–43079, 2023.
- [23] 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.
- [24] 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.
- [25] 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.
- [26] Chong Mou, Xintao 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*, volume 38, pages 4296–4304, 2024.
- [27] Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhui Chen, and Furu Wei. Kosmos-g: Generating im-

- ages in context with multimodal large language models. *arXiv preprint arXiv:2310.02992*, 2023.
- [28] 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.
  - [29] 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.
  - [30] 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.
  - [31] 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.
  - [32] Jean-Marie Schaeffer. Art of the modern age: Philosophy of art from kant to heidegger. 2023.
  - [33] 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.
  - [34] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020.
  - [35] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
  - [36] Ruichen Wang, Zekang Chen, Chen Chen, Jian Ma, Haonan Lu, and Xiaodong Lin. Compositional text-to-image synthesis with attention map control of diffusion models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 5544–5552, 2024.
  - [37] Tianyu Wu, Shizhu He, Jingping Liu, Siqi Sun, Kang Liu, Qing-Long Han, and Yang Tang. A brief overview of chatgpt: The history, status quo and potential future development. *IEEE/CAA Journal of Automatica Sinica*, 10(5):1122–1136, 2023.
  - [38] Tsung-Han Wu, Long Lian, Joseph E Gonzalez, Boyi Li, and Trevor Darrell. Self-correcting llm-controlled diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6327–6336, 2024.
  - [39] Jinpeng 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.
  - [40] Binbin Yang, Yi Luo, Ziliang Chen, Guangrun Wang, Xiaodan Liang, and Liang Lin. Law-diffusion: Complex scene generation by diffusion with layouts. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 22669–22679, 2023.
  - [41] Ling Yang, Jingwei Liu, Shenda Hong, Zhilong Zhang, Zhilin Huang, Zheming Cai, Wentao Zhang, and Bin Cui. Improving diffusion-based image synthesis with context prediction. *Advances in Neural Information Processing Systems*, 36, 2024.
  - [42] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and CUI Bin. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *Forty-first International Conference on Machine Learning*, 2024.
  - [43] 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.
  - [44] 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.
  - [45] Zuopeng Yang, Daqing Liu, Chaoyue Wang, Jie Yang, and Dacheng Tao. Modeling image composition for complex scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7764–7773, 2022.
  - [46] 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.