

Detail++: Training-Free Detail Enhancer for Text-to-Image Diffusion Models

Lifeng Chen¹ Jiner Wang¹ Zihao Pan¹ Beier Zhu^{1, 2} Xiaofeng Yang^{1, 2} Chi Zhang^{1†}

¹AGI Lab, Westlake University, ²Nanyang Technological University.

<https://detail-plus-plus.github.io/>



Figure 1. **A comparison between our method and current state-of-the-art generative models.** The mainstream models often suffer from issues such as semantic overflow, complex attribute mismatching, and style blending. Even Flux, the leading generative model under the DiT framework, struggles to overcome these challenges. In contrast, our method, Detail++, based on SDXL, achieves highly accurate semantic binding in a *training-free* way.

Abstract

Recent advances in text-to-image (T2I) generation have led to impressive visual results. However, these models still face significant challenges when handling complex prompts—particularly those involving multiple subjects with distinct attributes. Inspired by the human drawing process, which first outlines the composition and then incrementally adds details, we propose Detail++, a training-free framework that introduces a novel Progressive Detail Injection (PDI) strategy to address this limitation. Specifically, we decompose a complex prompt into a sequence of simplified sub-prompts, guiding the generation process in stages. This staged generation leverages the inherent layout-controlling capacity of self-attention to first ensure global composition, followed by precise refinement. To achieve accurate binding between attributes and corresponding subjects, we exploit cross-attention mechanisms and further introduce a Centroid Alignment Loss at test time to reduce binding noise and enhance attribute consistency.

Extensive experiments on T2I-CompBench and a newly constructed style composition benchmark demonstrate that Detail++ significantly outperforms existing methods, particularly in scenarios involving multiple objects and complex stylistic conditions.

1. Introduction

In recent years, text-to-image (T2I) generation [13, 27, 42, 48, 50] techniques have advanced remarkably, enabling the creation of high-quality, detail-rich images from textual descriptions and demonstrating considerable potential in art design, advertising, entertainment, and education [11, 28, 30]. Despite these improvements, current methods still struggle when handling complex descriptions that involve multiple subjects. A well-known challenge, termed “detail binding,” emerges when descriptions assign distinct sets of attributes to multiple subjects. This often causes models to misassociate attributes across subjects, resulting in semantic overflow, incorrect attribute matching, and undesired style blending, as shown in Fig. 1.

†Corresponding author.

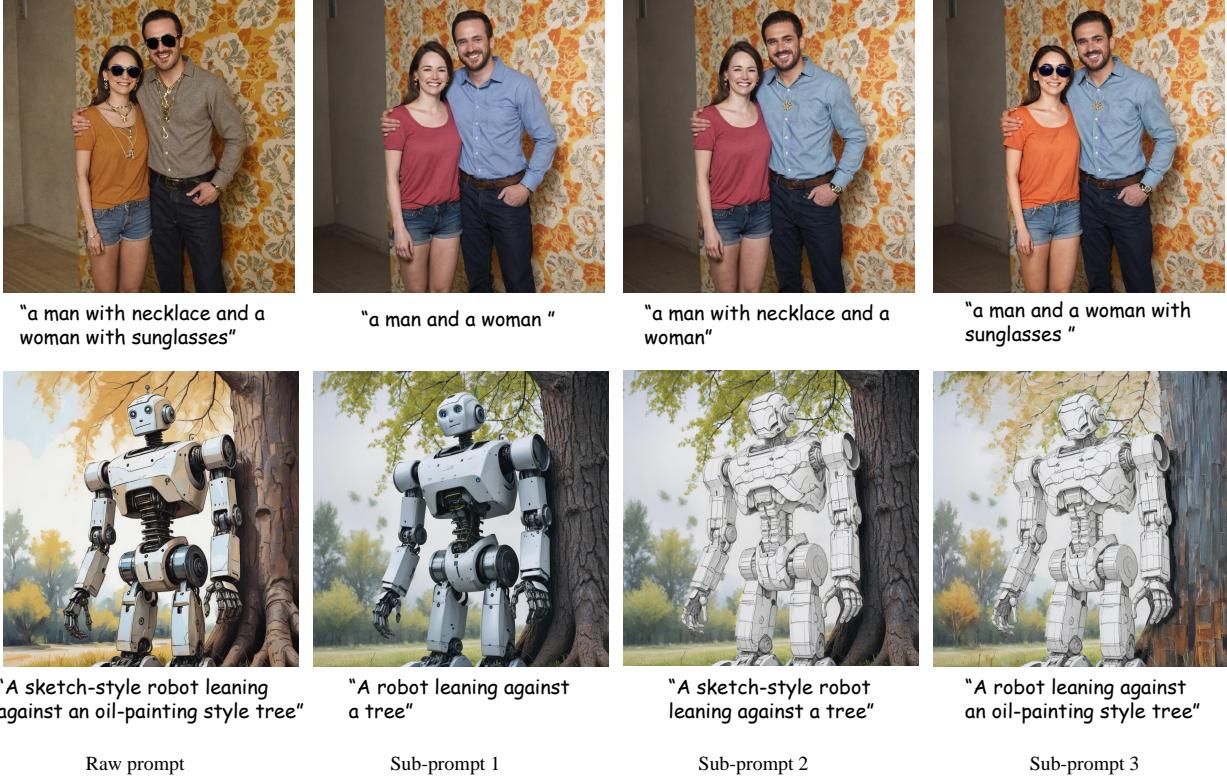


Figure 2. **The basic process of Detail++.** As shown in the first column, generating complex prompts in a single branch often results in inaccurate or blended attribute assignments. For example, attributes such as “sunglasses” and “necklace” may be mistakenly applied to the wrong subject. Our method addresses this challenge through a progressive approach: we first ignore all complex modifiers to produce a rough generation base, then systematically inject details to ensure each attribute is precisely added to its corresponding subject region. The prompt displayed below each image indicates the input used for that specific branch, and the second row demonstrates the method’s effectiveness in style combination scenarios. Note that all four branches here are generated in parallel.

Regarding the above challenges, our observation is that current state-of-the-art models attempt to render all elements and their attributes from a prompt simultaneously, leading to imperfect detail binding. Drawing inspiration from human artistic practice—where artists typically begin by sketching the basic spatial layout and outlines before gradually refining details—we progressively inject attributes to accurately integrate specific details into their corresponding subject regions, thereby achieving more semantically coherent image generation. To achieve our goal, we propose **Detail++**, a multi-branch image generation framework where different simplified versions of the prompt work together to produce accurate image content. The process begins with an initial denoising branch using the original prompt, during which we extract and share self-attention maps across other following branches to maintain consistent layout. The second denoising branch employs the simplest sub-prompt—containing only the subject and its basic layout—to create a foundation without fine details. In subsequent parallel branches, specific attributes from each sub-prompt are progressively injected into the image, with

all branches maintaining the same overall layout. To ensure proper binding of details to their correct subjects, we introduce an Accumulative Latent Modification strategy that uses cross-attention maps to generate binary masks, precisely locating different subject regions for targeted detail injection. Additionally, we incorporate Centroid Alignment Loss during test time to optimize subject-specific cross-attention maps, addressing map inaccuracies and further improving detail injection precision. The whole process is illustrated in Fig. 2.

Our proposed Detail++ was evaluated on the widely used T2I-CompBench benchmark [23] as well as on our custom style composition benchmark. The results demonstrate that Detail++ consistently outperforms existing methods, particularly in scenarios involving multiple style bindings. In qualitative evaluations, **Detail++** produces high-quality images without additional fine-tuning [15, 22–24] or predefined layout information [15, 32, 33, 63, 66, 67], underscoring the practicality of our approach. The main contributions of this paper are as follows:

- We innovatively propose **Detail++**, a training-free multi-

- branch framework. It improves the accuracy of existing T2I models when facing complex prompts involving multiple subjects through progressive attribute injection.
- We developed a self-attention map sharing mechanism that utilizes the properties of self-attention map in U-Net to control the image layout, thus creating a consistent base for accurate attribute binding.
 - We further introduced a test-time optimization strategy based on Centroid Alignment Loss, which further improved the accuracy of attribute binding.
 - Extensive experiments demonstrate the advantages of Detail++ compared to current SOTA methods, and we will subsequently open-source the code.

2. Related work

2.1. Text-to-Image Diffusion Models

Diffusion Models have become the dominant approach in text-to-image (T2I) generation [2, 12, 13, 27, 42, 45, 48]. Models like Stable Diffusion [48], DALL-E 2 [46], andImagen [50] generate high-quality images by progressively denoising noise. They leverage neural networks to map text to visuals, often using pre-trained text encoders like CLIP [44] to enhance text comprehension and guide generation. However, even the most advanced models, such as Flux [27], still struggle to generate images that closely align with the text prompts. This has led to a line of studies focused on improving T2I alignment.

2.2. T2I Alignment

As mentioned earlier, T2I baseline models struggle to precisely follow text conditions, often facing issues such as attribute overflow (where attributes spill over onto unmentioned subjects), mismatching (where attributes are incorrectly matched), and blending (where attributes from different subjects mix onto a single subject). To address these challenges, several solutions have been proposed [7, 14, 17, 26, 40, 54, 55]. Among them, some methods based on test-time optimization [1, 6, 10, 31, 40, 47, 64, 65] aim to align the cross-attention maps of attributes with the subject tokens, but they are often slow and difficult to optimize for noise efficiently. Other methods based on layout [15, 33, 43, 59, 62, 63] rely on users or LLMs to predefined the layout, often involving complex intermediate steps. There are also approaches that need fine-tuning [9, 15, 22, 24, 60], such as ELLA [22], which enhances the model’s understanding of text by integrating LLMs into T2I generation models, but requires significant modifications to the baseline model, leading to high training costs. While these methods have made notable progress in addressing attribute mismatching, they still struggle with overflow and blending. Recently, a method called ToME [21] has provided a new approach to addressing overflow by merging token embeddings gener-

ated by the CLIP text encoder. However, the issue of attribute blending remains largely unresolved. Our method, Detail++, effectively avoids these three problems by progressively injecting attributes, allowing for accurate binding of different types of attributes, including object, color, texture, and style.

2.3. Image Editing

Image editing manipulates specific regions while preserving overall fidelity. Traditional text-driven image editing methods [29, 39, 56] have shown promise by combining GANs [16, 41] with CLIP. Recently, a large number of new methods [2, 3, 25, 37, 57] based on diffusion models have emerged to address various aspects of image editing tasks in a training-free manner, most of which [5, 18, 34, 38, 51] leverage the attention mechanism [52]. For example, Prompt-to-Prompt (P2P) [18] adjusts cross-attention maps to enable a variety of editing operations. MasaCtrl [5] realizes rigid editing while maintaining the overall textures and identity in a mutual attention way. Building on these works, FPE [34] distinguishes the different functionalities between self-attention and cross-attention in stable diffusion models. Our method draws inspiration from these works and further extends them to a progressive editing framework.

3. Preliminary

This section introduces the basics of the attention mechanism, which serves as a key component of our proposed method.

Attentions in U-Net. The U-Net architecture [49] has become a cornerstone in many state-of-the-art diffusion. U-Net in diffusion models integrates both self-attention and cross-attention mechanisms, each serving unique purposes in image generation. Cross-attention focuses on extracting the semantic elements of the prompt, ensuring the generated images are consistent with text prompt [18]. The attention matrix in cross-attention can be defined as:

$$M_{\text{cross}} = \text{Softmax} \left(\frac{Q_{\text{img}} K_{\text{text}}^{\top}}{\sqrt{d_k}} \right), \quad (1)$$

where Q_{img} is the query embedding derived from the U-Net’s internal features, K_{text} is the key embedding originating from the textual embeddings of the prompt, and d_k is the dimensionality of the key vectors. On the other hand, self-attention primarily relates to the spatial layout and structural details of the image by computing relationships within the image feature space [34]. This attention map can be represented as:

$$M_{\text{self}} = \text{Softmax} \left(\frac{Q_{\text{img}} K_{\text{img}}^{\top}}{\sqrt{d_k}} \right), \quad (2)$$

where both Q_{img} and K_{img} are the query and key matrices derived from the internal features within the U-Net. The

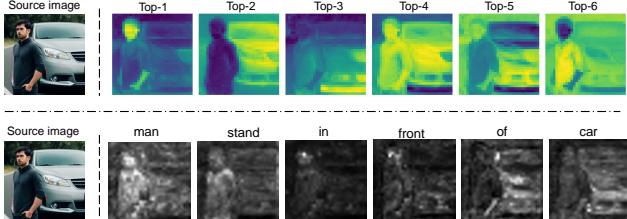


Figure 3. Attention visualization. The visualization of the self-attention map and cross cross-attention map of the prompt “*man stand in front of car*”. The self-attention map is visualized by displaying the top-6 components obtained after SVD [53]. The cross-attention map visualizations correspond to each token in the prompt.

self-attention map provides an internal relational structure that is crucial for preserving spatial consistency. To better understand the difference between self-attention and cross-attention, we visualize the attention maps of them separately in Fig. 3. Our proposed method builds upon these foundational attention mechanisms, specifically modifying and extending them to improve the handling of complex prompts.

4. Method

In this section, we provide an overview of our method, **Detail++**. Our core idea is to first ignore all complex, hard-to-generate details to obtain a rough generation base, and then progressively add more details in a painting-like manner. Given a complex original prompt, we automatically generate a series of simplified prompts, as described in Sec. 4.1. To ensure that images generated from these prompts with different attributes maintain the same layout, we share the self-attention map obtained during the generation of the original prompt with all sub-prompt generation processes (Sec. 4.1). This shared self-attention map provides a stable structural foundation. Controlled by our Accumulative Latent Modification strategy (Sec. 4.1), correct details are progressively injected into the corresponding subject regions. At each timestep, cross-attention maps for each subject keyword are used to create precise masks that guide the selective attribute injection. Moreover, in Sec. 4.2, we introduce a Centroid Alignment Loss to achieve more accurate and focused cross-attention map during test time, ensuring that the added details are strictly confined to their intended subject regions. The overall workflow is illustrated in Fig. 4.

4.1. Progressive Detail Injection

Prompt decomposition. To handle the complexity of detailed prompts, the first step in our method is to simplify the input prompt. Specifically, we use a language model, such as spaCy [20] or ChatGPT [4] to decompose the original complex prompt p_0 into a sequence of progressively simpli-

fied sub-prompts, denoted as $\mathcal{P} = \{p_0, p_1, \dots, p_{n-1}, p_n\}$, where p_1 represents the most simplified version by removing all modifiers. This simplification ensures that all attributes are eliminated in the p_1 branch, providing a clear base for progressive detail injection. Sub-prompts p_2, p_3, \dots, p_n each add an additional modifier to p_1 . The choice of decomposition is explained in more detail in the Appendix 9.

Furthermore, for each pair of prompt $p_{i+1}(i > 0)$ and p_1 , we also use the language model to identify the subject q_i to which the newly added attribute in p_{i+1} corresponds. These subjects are represented by the set $\mathcal{Q} = \{q_1, \dots, q_{n-2}, q_{n-1}\}$. For example, given the prompt $p_0 = “a red teddy bear wearing a green tracksuit”$, the set \mathcal{P} and subject set \mathcal{Q} would be as follows:

$$\mathcal{P} = \left\{ \begin{array}{l} p_0 : “a red teddy bear wearing a green tracksuit”, \\ p_1 : “a teddy bear wearing a tracksuit”, \\ p_2 : “a red teddy bear wearing a tracksuit”, \\ p_3 : “a teddy bear wearing a green tracksuit” \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : “teddy bear”, \\ q_2 : “tracksuit” \end{array} \right\}.$$

Shared self-attention map. Once we have the set of sub-prompts, the next step is to use these prompts to generate images cooperatively. Since previous work [5, 34, 38, 51] has shown that the self-attention map in U-Net stores the layout information of the image, we let the network generate images based on each sub-prompt, while achieving a consistent layout by sharing the self-attention map. Specifically, we cache the U-Net self-attention maps in the first branch and reuse them in subsequent branches. In this way, we make the images that different sub-prompts generated consistent with the layout of the original prompt. It is observed that the early stages of denoising are most critical for constructing the overall layout of the image, so we only use this strategy during the initial S steps of the denoising process. For the remaining $T - S$ steps, we allow each sub-prompt’s U-Net models to predict noise independently to ensure the image fidelity, as is done in the standard diffusion process. In this way, we can keep the layout of the images generated by the different sub-prompts consistent with the original prompt.

Accumulative Latent Modification. In order to precisely bind the newly added attributes in each sub-prompt to their corresponding subjects, we further develop latent-level subject masking for each attribute injection. Specifically, we extract a binary mask based on the cross-attention map corresponding to each subject q_i (Sec. 4.1) from a sub-prompt, which highlights the relevant regions of each subject in the image. This process involves normalizing the cross-attention maps and applying a thresholding operation

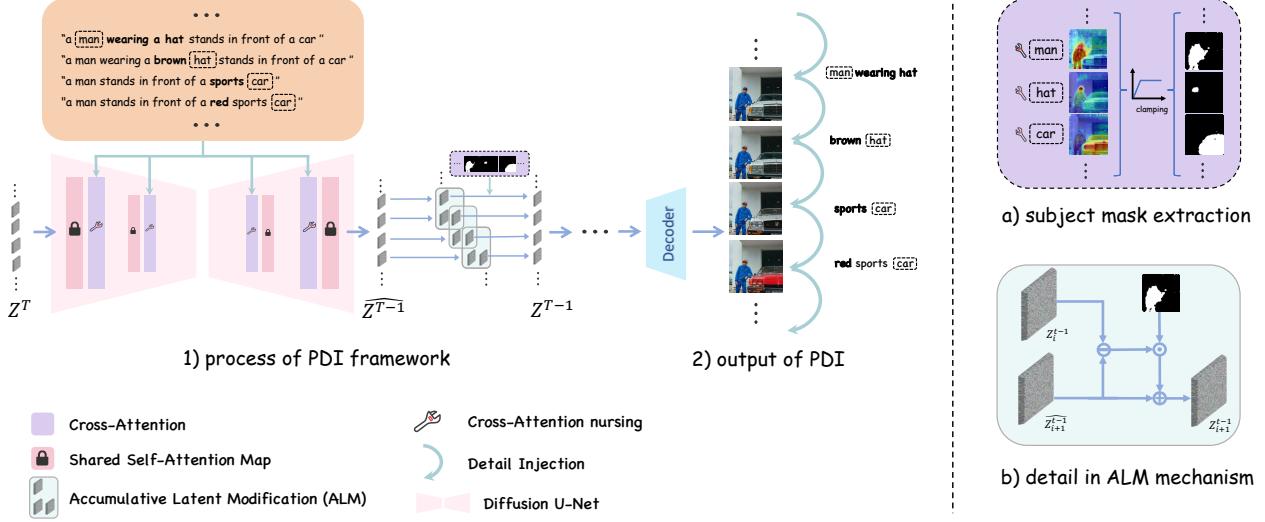


Figure 4. Overview of Detail++. Our method consists of a large framework, Progressive Detail Injection (PDI), and test-time attention nurturing based on Centroid Alignment Loss. The sub-prompts, decomposed by spaCy, are together passed through U-Net for parallel inference. In each denoising step, the resulting batch of latents undergoes Accumulative Latent Modification (ALM), modifying the only region current adding attribute corresponding to. Note that, in our framework, the self-attention maps for all branches are unified, which ensures a consistent layout and avoids conflicts among editing effects.

to create a distinct binary mask for each subject q_i . The binary mask extraction for subject q_i is given by the following formula:

$$B_i = \mathbf{1} \left[\frac{\overline{M}_i - \min(\overline{M}_i)}{\max(\overline{M}_i) - \min(\overline{M}_i)} > \tau \right], \quad (3)$$

where B_i is the binary mask of subject q_i , and \overline{M}_i represents its averaged attention map across layers. 15. τ is a threshold parameter. The indicator function $\mathbf{1}[\cdot]$ converts the normalized attention values to a binary mask, where values above τ are set to 1 and others to 0. $\min(\overline{M}_i)$ and $\max(\overline{M}_i)$ represent the minimum and maximum values of the averaged attention map respectively.

Binary masks can restrict latent modifications to a specific region. Specifically, for the i -th branch, the modified latent representation is computed as:

$$z_{i+1}^{t-1} = z_i^{t-1} + B_i \odot (\widehat{z}_{i+1}^{t-1} - z_i^{t-1}), \quad (4)$$

where \widehat{z}_{i+1}^{t-1} represents the denoised latent in $i+1$ -th branch, z_i^{t-1} is the latent from the previous branch, B_i is the binary semantic mask for the current subject, and z_{i+1}^{t-1} is the modified latent incorporating correct attribute. The operator \odot denotes element-wise multiplication.

For example, if a region's mask value b is 1, this signifies that the region corresponds to the target subject. In this case, the new latent z_{i+1} generated with the additional prompt detail is applied to this region, ensuring our framework integrates the added details precisely within the designated area.

Conversely, if b is 0 for a given region, indicating that the area is unrelated to the current attribute, the model reuses the latent z_i from the previous branch, thereby excluding the new attribute's influence from unrelated regions. This selective application prevents cross-subject interference and preserves the distinction of each subject's features across the image. Thus, the entire process of progressive attention substitution is outlined in Algorithm 1.

4.2. Centroid Alignment Loss

When generating binary masks from cross-attention maps, we observe a critical issue in existing cross-attention mechanisms: attention activations tend to be dispersed and unfocused, causing a subject's attention to spread into other subject regions, resulting in attribute injection errors and degraded generation quality. To address this problem, we propose a test-time optimization method by introducing a Centroid Alignment Loss to refine cross-attention maps. Our key insight is that in ideal conditions, a subject's cross-attention map should form concentrated activations at the center of the subject region, rather than being scattered across the entire image space. Based on this, we design a mechanism that encourages each subject's salient attention point (the point with largest attention value) to align with the centroid of its attention distribution.

Specifically, we first calculate the centroid $p_{\text{centroid}}(q_i)$ for each subject q_i 's attention map, where q_i represents the subject word to which the attribute in the i -th branch is being added. This centroid serves as the weighted center of the cross-attention map corresponding to q_i 's token in the

Algorithm 1 Progressive Detail Injection

Require: A source prompt p_0 , and sub-prompts derived from it p_1, \dots, p_n , DM represents the diffusion models to predict the noise of next step. T is the total denoising timesteps.

Ensure: Input $\mathcal{P} = \{p_0, p_1, \dots, p_n\}$ to model in parallel.

- 1: $z_0^T, z_1^T, \dots, z_n^T \leftarrow z^T \sim \mathcal{N}(0, \mathcal{I})$; \triangleright Initialized with same latent.
- 2: $\mathcal{Z}^T \leftarrow \{z_0^T, z_1^T, \dots, z_n^T\}$
- 3: **for** $t = T, \dots, T - S + 1$ **do**
- 4: $\mathcal{Z}^{t-1} \leftarrow \text{DM}(\mathcal{Z}^t, \mathcal{P}, t, \overline{M}_{self}^t)$; $\triangleright \overline{M}_{self}^t$ is from first branch.
- 5: **for** $i = 1, 2, \dots, n - 1$ **do**
- 6: $\widehat{z}_{i+1}^{t-1} \leftarrow z_{i+1}^{t-1}$
- 7: $z_{i+1}^{t-1} \leftarrow z_i^{t-1} + B_i \odot (\widehat{z}_{i+1}^{t-1} - z_i^{t-1})$; $\triangleright B_i$ is the binary mask consistent with Sec 4.2.3.
- 8: **end for**
- 9: **end for**
- 10: **for** $t = T - S, \dots, 1$ **do**
- 11: $\mathcal{Z}^{t-1} = \text{DM}(\mathcal{Z}^t, \mathcal{P}, t)$;
- 12: **end for**
- 13: **return** \mathcal{Z}^0

prompt:

$$p_{\text{centroid}}(q_i) = \frac{1}{\sum_{h,w} \overline{M}_i(h,w)} \left[\frac{\sum_{h,w} w \cdot \overline{M}_i(h,w)}{\sum_{h,w} h \cdot \overline{M}_i(h,w)} \right], \quad (6)$$

where h and w represent the height and width spatial coordinates respectively, and $\overline{M}_i(h,w)$ is the normalized attention value at that position.

Next, we define the Centroid Alignment Loss by minimizing the Euclidean distance between the attention centroid and the brightest point $p_{\max}(q_i)$, encouraging a more focused attention distribution:

$$L_{\text{align}} = \sum_i \|p_{\text{centroid}}(q_i) - p_{\max}(q_i)\|^2. \quad (7)$$

We combine this loss term with the entropy loss used in ToME[21] to form the total loss function L_{total} . The entropy loss is defined as:

$$L_{\text{ent}} = - \sum_{m \in \overline{M}_i} m \log(m), \quad (8)$$

where m represents the attention values in the normalized cross-attention map \overline{M}_i for subject q_i . Our total loss function becomes $L_{\text{total}} = L_{\text{align}} + \lambda L_{\text{ent}}$, where λ is a trade-off hyperparameter. During each diffusion step t , we update the latent variable through gradient descent:

$$\mathbf{z}'_t \leftarrow \mathbf{z}_t - \alpha_t \cdot \nabla_{\mathbf{z}_t} L_{\text{total}}, \quad (9)$$

where α_t is the step size.

Experimental results demonstrate that introducing the Centroid Alignment Loss not only produces more concentrated cross-attention maps but also significantly improves the precision of attribute injection. The optimized attention maps generate more accurate binary masks, ensuring that detailed attributes are correctly injected into their corresponding subject regions, effectively resolving attribute confusion issues in complex prompts.

5. Experiments

5.1. Experimental Setups

Evaluation Benchmarks and Metrics. We conduct extensive experiments on the T2I-CompBench[23], a widely-used benchmark designed specifically for text-to-image generation tasks involving complex compositional prompts. Specifically, we adopt subsets in T2I-CompBench designed to evaluate attributes such as color, shape, texture, because they directly reflect our method's capacity for semantic alignment. To further validate the subjective quality of generated images and assess human preference scores, we incorporate the ImageReward[58] model, a learned metric widely accepted for reflecting human perceptual judgments in image generation tasks.

Existing benchmarks typically lack an effective evaluation mechanism for attributes related to the combination of multiple artistic styles, which facing the problems of style blend. To bridge this gap and provide a more comprehensive evaluation, we propose a novel benchmark called the Style Composition Benchmark (SCB). SCB consists of 300 carefully designed prompts that cover a wide range of distinguished artistic styles. More details can be found in Appendix 11.

Our evaluation method for SCB independently measures the alignment between each element in the image and its corresponding style descriptor. Specifically, it involves parsing the prompt of each generated image to identify style descriptors associated with different semantic components. Each component in the generated image is first cropped out using Grounding DINO [36], based on the component words parsed from the prompt. We then calculate the component-style correspondence for each segmented element using the CLIP-Score [19, 44] model. The final style score for each generated image is computed as the average of the style-matching scores across all segmented semantic elements. Therefore, based on the deviation from the reference images, our new metric can efficiently assess whether style blending has occurred.

Implementation Details. Our method is built upon the SDXL [42] baseline and utilizes the SpaCy [20] NLP toolkit for modifier recognition and automatic prompt management. The different branches of our model are processed in parallel, significantly reducing runtime. The self-attention

Table 1. Quantitative comparison of detail binding performance across multiple metrics. Our proposed Detail++ outperforms all baselines on color, texture, shape and style binding. The highest scores are highlighted in blue, and second-highest in green.

Method	Train	BLIP-VQA \uparrow			Human-preference \uparrow			$SCB \uparrow$
		Color	Texture	Shape	Color	Texture	Shape	
SD1.5[48]	✓	0.4719	0.4334	0.3898	-0.360	-0.689	-0.483	0.2196
Composable Diffusion[35]	✗	0.4063	0.3645	0.3299	-0.209	-1.043	-0.813	0.2032
Structured Diffusion[14]	✗	0.4990	0.4900	0.4218	-0.025	-0.325	0.169	0.2207
GORS[23]	✓	0.6603	0.6287	0.4785	0.017	-0.417	1.115	-
ELLA v1.5[22]	✓	0.6911	0.6308	0.4938	-	-	-	-
SDXL[42]	✓	0.6369	0.5637	0.5408	0.733	0.124	1.184	0.2488
PixArt- α [8]	✓	0.6886	0.7044	0.5582	0.242	-0.788	0.165	0.2398
Attention Regulation[64]	✗	0.5860	0.5173	0.4672	0.268	-0.603	1.118	0.2248
Ranni xl[15]	✓	0.6893	0.6325	0.4934	-	-	-	-
ToMe[21]	✗	0.6583	0.6371	0.5517	-0.028	-0.851	0.454	0.2554
Detail++(Ours)	✗	0.7389	0.7241	0.5582	1.773	0.300	1.424	0.2687

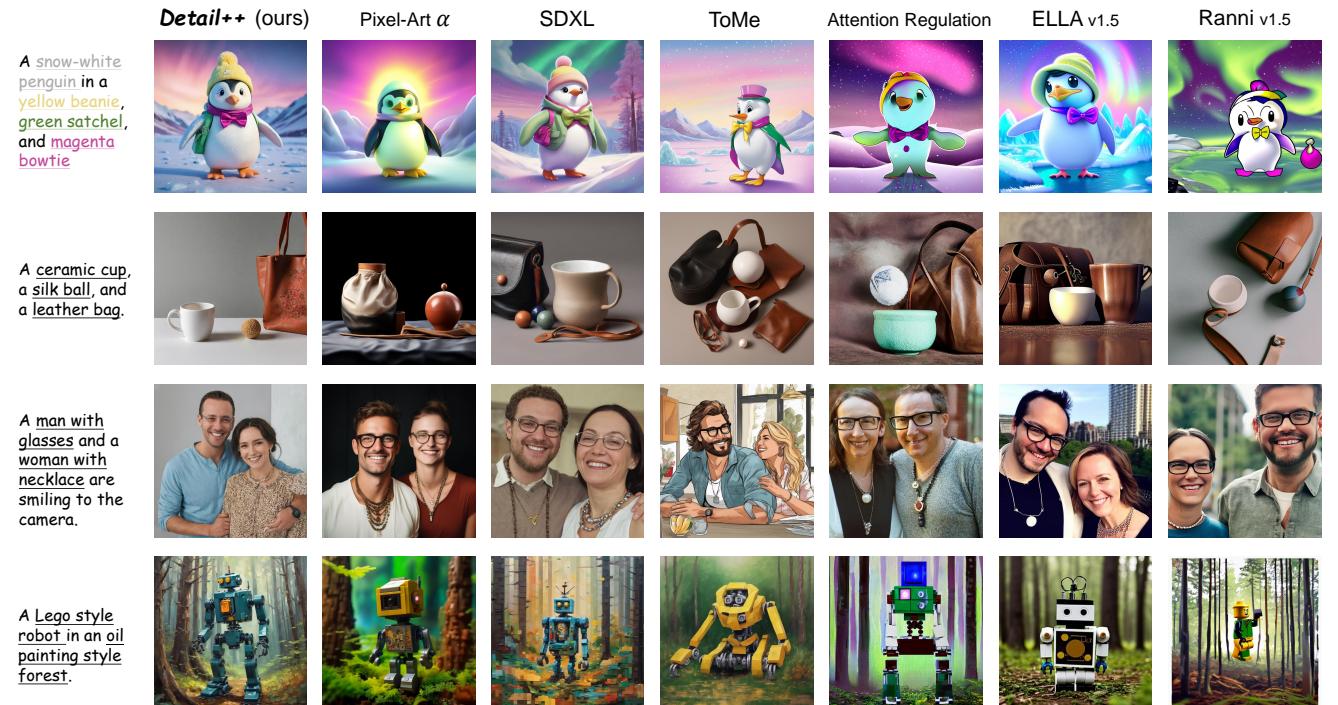


Figure 5. Qualitative comparison of various methods on complex prompts with multiple attribute types (object, color, texture, and style). Our method effectively prevents attribute overflow, complex attribute mismatching, and style blending while maintaining high visual fidelity.

map sharing mechanism is activated at 80% of the total denoising steps ($S = 0.8T$), employing attention sharing with map sizes of 32×32 from all blocks. In test-time optimization, we set the $\lambda = 1$ to achieve a good balance between these two losses. More experiment details can be found in

appendix.

Comparison Methods. To comprehensively evaluate the performance of our proposed approach, we conduct comparative analyses against several methods. In addition to comparing against the established baseline method,

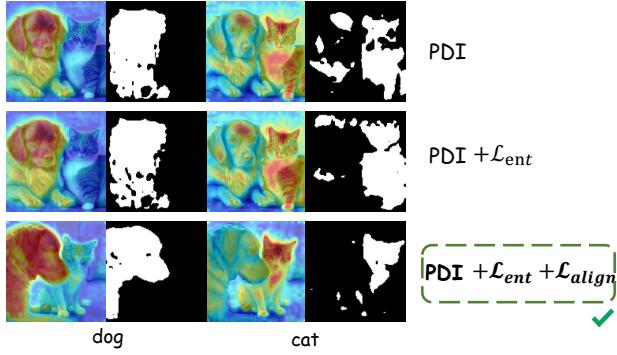


Figure 6. Ablation study of different optimization terms in binary mask extraction. The combined use of both L_{ent} and L_{align} produces more accurate and focused subject masks compared to using only L_{ent} or no optimization.



Figure 7. Visual results of different optimization strategies for the prompt “a dog wearing sunglasses and a cat wearing a necklace”. All results in 3 rows are generated based on the same random seed. The combined use of both L_{ent} and L_{align} achieves more precise detail injection, preventing attributes from incorrectly spilling into adjacent subject regions.

SDXL [42], we include multiple advanced methods explicitly designed to enhance detail binding and compositionality: Ranni [15], an approach introduces a semantic panel as an intermediary between textual prompts and image generation; ELLA [22], known for leveraging LLMs to enhance alignment between textual prompts and generated imagery; Attention Regulation [64], a sophisticated approach focusing on dynamically regulating attention weights to boost semantic accuracy and reduce visual ambiguities; ToME [21], a recent method based on efficient token merging, demonstrating notable improvements in compositional clarity and visual fidelity.

5.2. Experimental Results

Quantitative Comparison. As shown in Table 1, our method demonstrates superior or comparable performance

Table 2. Ablation study of different optimization terms on T2I-CompBench. The combined use of both L_{ent} and L_{align} performs best.

PDI	\mathcal{L}_{ent}	\mathcal{L}_{align}	BLIP-VQA		
			Color	Texture	Shape
	✗	✗	0.6369	0.5637	0.5408
	✓	✗	0.6942	0.6836	0.5507
	✓	✓	0.7034	0.6995	0.5521
	✓	✓	0.7389	0.7241	0.5582

in terms of color, texture, shape binding capabilities on T2I-CompBench [23] and style binding ability in our proposed SCB benchmark compared to other methods. Additionally, using the ImageReward [58] score as a proxy for human preference, our method (Detail++) exhibits stronger human preference alignment.

Qualitative Comparison. Fig. 5 presents a qualitative comparison that illustrates how different methods handle four distinct prompt types: complex color binding, complex texture binding, object binding, and style composition. From the first and second rows, it is evident that for color and texture attributes, when the number of subjects is large, attribute mismatching becomes more prominent. In the third row, for object-type attributes, semantic leakage often occurs, causing unmodified subjects to be influenced by the modifiers of other subjects, as seen when the necklace modifies the man, while the woman is also affected by the sunglasses. In the fourth row, we observe that existing mainstream methods struggle to handle multiple style modifiers applied to different components, often resulting in style blending. Our method, however, is the only one that successfully generates the Lego-style robot while maintaining the pure oil painting background, without any blending between the two. The visual results clearly demonstrate the superiority of our method (Detail++) in preventing attribute overflow, complex attribute binding, and style blending, which is consistent with our quantitative results in Table 1.

Ablation Study. As shown in Table 2, we conducted a quantitative analysis of different components and loss terms. It can be observed that using only L_{ent} [21] results in a slight improvement in the experimental outcomes. However, when the proposed L_{align} is added, the results show a significant enhancement. Furthermore, we visualize the generated binary masks of subjects when generating the prompt “a dog wearing sunglasses and a cat wearing a necklace” in Fig. 6. In the middle set of images, we observe that adding the entropy loss helps the cross-attention map converge to some extent. However, the highlighted region for the ‘cat’ token still includes the ‘dog’s forehead and neck area, which can lead to imprecise attribute addition. This is evident in the second row of Fig. 7, where

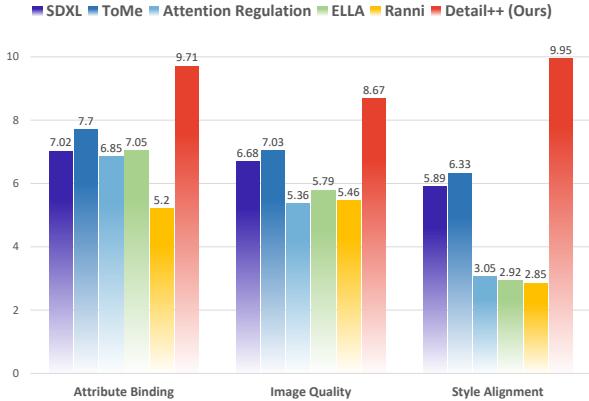


Figure 8. User study results evaluating different methods across three key metrics. Detail++ consistently outperforms all baseline methods in attribute binding, image quality, and style alignment, receiving significantly higher scores particularly in style alignment.

adding the necklace to the cat also inadvertently affects the dog’s neck area. Our proposed centroid alignment loss, however, further refines the initial branch generation, making the subjects’ cross-attention maps more focused and accurate, greatly reducing the overlap in binary masks across different subjects, which also can be found the generated image in Fig. 7.

5.3. User Study

To comprehensively evaluate the effectiveness of our proposed method, we conducted a user study with 145 participants in Fig. 8. The study consisted of 12 questions per questionnaire, covering three key aspects: Attribute Binding, Image Quality, and Style Binding. To ensure reliable and unbiased evaluation, we constructed randomized image pools for each comparative model (SDXL [42], ToMe [21], Attention Regulation [64], ELLA [22], Ranni [15] and our method, where each pool contained 120 images generated from 30 diverse prompts using 4 different random seeds. For each question in the survey, images were randomly selected from these pools based on the corresponding prompts. The results demonstrate that our method consistently outperforms all baseline approaches across all three metrics.

6. Conclusion

In this work, we present Detail++, a novel training-free approach that addresses the challenge of accurate detail binding in text-to-image generation with complex prompts. Our Progressive Detail Injection framework combines self-attention map sharing and Accumulative Latent Modification to ensure proper attribute assignment while maintaining consistent layouts. The introduced Centroid Alignment Loss further improves binding precision by optimizing cross-attention maps. Extensive

experiments on public benchmarks demonstrate that Detail++ significantly outperforms existing methods in preventing attribute overflow, mismatching, and style blending, while preserving high image fidelity. As a plug-and-play module compatible with current diffusion models, Detail++ represents an important step toward more controlled and semantically accurate text-to-image generation without requiring additional training.

References

- [1] Aishwarya Agarwal, Srikrishna Karanam, KJ Joseph, Apoorv Saxena, Koustava Goswami, and Balaji Vasan Srinivasan. A-star: Test-time attention segregation and retention for text-to-image synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2283–2293, 2023. 3
- [2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Qinsheng Zhang, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 3
- [3] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 18392–18402, 2023. 3
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020. 4
- [5] Mingdeng Cao, Xiantao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yingqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 22560–22570, 2023. 3, 4
- [6] 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. 3
- [7] Chieh-Yun Chen, Chiang Tseng, Li-Wu Tsao, and Hong-Han Shuai. A cat is a cat (not a dog!): Unraveling information mix-ups in text-to-image encoders through causal analysis and embedding optimization. *arXiv preprint arXiv:2410.00321*, 2024. 3
- [8] 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. 7, 1

- [9] Kai Chen, Enze Xie, Zhe Chen, Yibo Wang, Lanqing Hong, Zhenguo Li, and Dit-Yan Yeung. Geodiffusion: Text-prompted geometric control for object detection data generation. *arXiv preprint arXiv:2306.04607*, 2023. 3
- [10] 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*, 36:6048–6069, 2023. 3
- [11] Nassim Dehouche and Kullathida Dehouche. What’s in a text-to-image prompt? the potential of stable diffusion in visual arts education. *Heliyon*, 9(6), 2023. 1
- [12] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 3
- [13] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024. 1, 3
- [14] 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. 3, 7, 1
- [15] 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. 2, 3, 7, 8, 9, 1
- [16] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on Graphics (TOG)*, 41(4):1–13, 2022. 3
- [17] Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image generation with rich text. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7545–7556, 2023. 3
- [18] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 3
- [19] Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation via shared attention. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4775–4785, 2024. 6
- [20] Matthew Honnibal. spacy 2: Natural language understanding with bloom embeddings, convolutional neural networks and incremental parsing. (*No Title*), 2017. 4, 6
- [21] Taihang Hu, Linxuan Li, Joost van de Weijer, Hongcheng Gao, Fahad Shahbaz Khan, Jian Yang, Ming-Ming Cheng, Kai Wang, and Yaxing Wang. Token merging for training-free semantic binding in text-to-image synthesis. *Advances in Neural Information Processing Systems*, 37:137646–137672, 2025. 3, 6, 7, 8, 9, 1, 2
- [22] 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, 3, 7, 8, 9, 1
- [23] 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. 2, 6, 7, 8
- [24] Dongzhi Jiang, Guanglu Song, Xiaoshi Wu, Renrui Zhang, Dazhong Shen, Zhuofan Zong, Yu Liu, and Hongsheng Li. Comat: Aligning text-to-image diffusion model with image-to-text concept matching. *Advances in Neural Information Processing Systems*, 37:76177–76209, 2024. 2, 3
- [25] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2426–2435, 2022. 3
- [26] Jeeyung Kim, Erfan Esmaeili, and Qiang Qiu. Text embedding is not all you need: Attention control for text-to-image semantic alignment with text self-attention maps. *arXiv preprint arXiv:2411.15236*, 2024. 3
- [27] Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024. 1, 3
- [28] Mingkun Lei, Xue Song, Beier Zhu, Hao Wang, and Chi Zhang. Stylestudio: Text-driven style transfer with selective control of style elements. *arXiv preprint arXiv:2412.08503*, 2024. 1
- [29] Bowen Li, Xiaojuan Qi, Thomas Lukasiewicz, and Philip HS Torr. Manigan: Text-guided image manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7880–7889, 2020. 3
- [30] Jialu Li, Yuanzhen Li, Neal Wadhwa, Yael Pritch, David E Jacobs, Michael Rubinstein, Mohit Bansal, and Nataniel Ruiz. Unbounded: A generative infinite game of character life simulation. *arXiv preprint arXiv:2410.18975*, 2024. 1
- [31] Yumeng Li, Margret Keuper, Dan Zhang, and Anna Khoreva. Divide & bind your attention for improved generative semantic nursing. *arXiv preprint arXiv:2307.10864*, 2023. 3
- [32] 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. 2
- [33] 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. 2, 3
- [34] Bingyan Liu, Chengyu Wang, Tingfeng Cao, Kui Jia, and Jun Huang. Towards understanding cross and self-attention in stable diffusion for text-guided image editing. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 7817–7826, 2024. 3, 4
- [35] 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. 7, 1
- [36] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Qing Jiang, Chunyuan Li, Jianwei Yang,

- Hang Su, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. In *European Conference on Computer Vision*, pages 38–55. Springer, 2024. 6
- [37] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 3
- [38] Jisu Nam, Heesu Kim, DongJae Lee, Siyoon Jin, Seungryong Kim, and Seunggyu Chang. Dreammatcher: appearance matching self-attention for semantically-consistent text-to-image personalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8100–8110, 2024. 3, 4
- [39] Seonghyeon Nam, Yunji Kim, and Seon Joo Kim. Text-adaptive generative adversarial networks: manipulating images with natural language. *Advances in neural information processing systems*, 31, 2018. 3
- [40] Geon Yeong Park, Jeongsol Kim, Beomsu Kim, Sang Wan Lee, and Jong Chul Ye. Energy-based cross attention for bayesian context update in text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 36: 76382–76408, 2023. 3
- [41] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2085–2094, 2021. 3
- [42] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1, 3, 6, 7, 8, 9
- [43] 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. 3
- [44] 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. 3, 6
- [45] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 3
- [46] 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. 3
- [47] Royi Rassin, Eran Hirsch, Daniel Glickman, Shauli Ravfogel, Yoav Goldberg, and Gal Chechik. Linguistic binding in diffusion models: Enhancing attribute correspondence through attention map alignment. *Advances in Neural Information Processing Systems*, 36:3536–3559, 2023. 3
- [48] 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, 3, 7
- [49] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5–9, 2015, proceedings, part III* 18, pages 234–241. Springer, 2015. 3
- [50] 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, 3
- [51] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 3, 4
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 3
- [53] Michael E Wall, Andreas Rechtsteiner, and Luis M Rocha. Singular value decomposition and principal component analysis. In *A practical approach to microarray data analysis*, pages 91–109. Springer, 2003. 4
- [54] Tianyi Wei, Dongdong Chen, Yifan Zhou, and Xingang Pan. Enhancing mmdit-based text-to-image models for similar subject generation. *arXiv preprint arXiv:2411.18301*, 2024. 3, 2
- [55] Feize Wu, Yun Pang, Junyi Zhang, Liyanu Pang, Jian Yin, Baoquan Zhao, Qing Li, and Xudong Mao. Core: Context-regularized text embedding learning for text-to-image personalization. *arXiv preprint arXiv:2408.15914*, 2024. 3
- [56] Weihao Xia, Yujiu Yang, Jing-Hao Xue, and Baoyuan Wu. Tedigan: Text-guided diverse face image generation and manipulation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2256–2265, 2021. 3
- [57] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22428–22437, 2023. 3
- [58] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagewerard: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023. 6, 8
- [59] 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. 3

- [60] 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. 3
- [61] 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. 3, 4
- [62] Tianjun Zhang, Yi Zhang, Vibhav Vineet, Neel Joshi, and Xin Wang. Controllable text-to-image generation with gpt-4. *arXiv preprint arXiv:2305.18583*, 2023. 3
- [63] Xinchen Zhang, Ling Yang, Yaqi Cai, Zhaochen Yu, Kai-Ni Wang, Ye Tian, Minkai Xu, Yong Tang, Yujiu Yang, Bin Cui, et al. Realcompo: Balancing realism and compositionality improves text-to-image diffusion models. *Advances in Neural Information Processing Systems*, 37:96963–96992, 2024. 2, 3
- [64] Yang Zhang, Teoh Tze Tzun, Lim Wei Hern, and Kenji Kawaguchi. Enhancing semantic fidelity in text-to-image synthesis: Attention regulation in diffusion models. In *European Conference on Computer Vision*, pages 70–86. Springer, 2024. 3, 7, 8, 9, 1
- [65] Yasi Zhang, Peiyu Yu, and Ying Nian Wu. Object-conditioned energy-based attention map alignment in text-to-image diffusion models. In *European Conference on Computer Vision*, pages 55–71. Springer, 2024. 3
- [66] Peiang Zhao, Han Li, Ruiyang Jin, and S Kevin Zhou. Loco: Locally constrained training-free layout-to-image synthesis. *arXiv preprint arXiv:2311.12342*, 2023. 2
- [67] Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion: Controllable diffusion model for layout-to-image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22490–22499, 2023. 2

7. Limitation and Discussion

The proposed Progressive Detail Injection (PDI) framework addresses detail binding issues in text-to-image generation through a training-free approach, yet it exhibits certain limitations. The method heavily relies on the quality and accuracy of initial Self-Attention maps—if the layout generation in the early phase is suboptimal, subsequent attribute injection processes may not fully rectify these issues, potentially constraining the final generation quality.

8. Implementation Details

Method details. Our method is implemented on SDXL. We extract the cross attention maps from all layers with resolution of 32×32 of all blocks in U-Net, as we found this is comparatively fine-grained and accurate. Also, our all experiments are conducted in a single Nvidia H-20 GPU.

Baseline methods implementation. For quantitative comparison in Table 1, we use official implementation of Stable diffusion 1.5 [48], Stable diffusion XL [42], Structured Diffusion Guidance [14], Composable diffusion [35], PixelArt- α [8], ELLA [22], ToME [21], Attention Regulation [64]. Since the SDXL checkpoint of Ranni [15] is not open-sourced, thus we directly refer to the score in their paper.

9. Sub-prompts Management

In this section, we compare the effects of different prompt management granularities on generation quality. For example, given the original prompt $p_{\text{ori}} = "a red dog with sunglasses and a blue cat with a necklace."$, we can manage the sub-prompts in two ways: 1) most complex first, or 2) simplest first. Regarding granularity, we can manage them in two ways: a) one subject per branch, or b) one attribute per branch. This results in a total of four possible configuration combinations. These are:

Config. A (1 + a):

$$\mathcal{P} = \left\{ \begin{array}{l} p_0 : "a red dog with sunglasses and a blue cat with a necklace.", \\ p_1 : "a dog and a cat.", \\ p_2 : "a red dog with sunglasses and a cat.", \\ p_3 : "a dog and a blue cat with a necklace." \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : "dog", \\ q_2 : "cat" \end{array} \right\}.$$

Config. B (1+b):

$$\mathcal{P} = \left\{ \begin{array}{l} p_0 : "a red dog with sunglasses and a blue cat with a necklace.", \\ p_1 : "a dog and a cat.", \\ p_2 : "a red dog and a cat.", \\ p_3 : "a dog with sunglasses and a cat.", \\ p_4 : "a dog and a blue cat.", \\ p_5 : "a dog and a cat with a necklace." \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : "dog", \\ q_2 : "dog", \\ q_3 : "cat", \\ q_4 : "cat" \end{array} \right\}.$$

Config. C (1 + a):

$$\mathcal{P} = \left\{ \begin{array}{l} p_1 : "a dog and a cat.", \\ p_2 : "a red dog with sunglasses and a cat.", \\ p_3 : "a dog and a blue cat with a necklace." \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : "dog", \\ q_2 : "cat" \end{array} \right\}.$$

Config. D (2+b):

$$\mathcal{P} = \left\{ \begin{array}{l} p_1 : "a dog and a cat.", \\ p_2 : "a red dog and a cat.", \\ p_3 : "a dog with sunglasses and a cat.", \\ p_4 : "a dog and a blue cat.", \\ p_5 : "a dog and a cat with a necklace." \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : "dog", \\ q_2 : "dog", \\ q_3 : "cat", \\ q_4 : "cat" \end{array} \right\}.$$

Therefore, we conduct plenty of experiments to test which one can best reflect our method’s efficiency, as shown in Table A2. Another tokenization approach is accumulative prompts, where the previously added attribute continues to appear in the prompt of the next branch. This is referred to as the **accUmu. prompt**:

$$\mathcal{P} = \left\{ \begin{array}{l} p_1 : "a dog and a cat.", \\ p_2 : "a red dog and a cat.", \\ p_3 : "a red dog with sunglasses and a cat.", \\ p_4 : "a red dog with sunglasses and a blue cat.", \\ p_5 : "a red dog with sunglasses and a blue cat with a necklace." \end{array} \right\},$$

$$\mathcal{Q} = \left\{ \begin{array}{l} q_1 : "dog", \\ q_2 : "dog", \\ q_3 : "cat", \\ q_4 : "cat" \end{array} \right\}.$$

However, as shown in the second-to-last row of Table A2, the quantification results are not ideal. We hypothesize that this is due to modifier overflow in the longer prompts [21], which degrades the efficiency of our method.

Table A1. Time Complexity of various methods. The highest scores are highlighted in blue, and second-highest in green.

Method	Inference Steps	Time Cost	Color	Texture	Shape
SDXL	50	30s	0.6275	0.5637	0.5408
ToMe	50	87s	0.6583	0.6371	0.5517
Attention Regulation	50	83s	0.5860	0.5173	0.4672
Detail++(Ours)	50	31s	0.7389	0.7241	0.5582

Table A2. Different configurations' effect.

Config	1/2	a/b	BLIP-VQA			Time Consumption	Memory Consumption
			Color	Texture	Shape		
A	1	a	0.7286	0.7205	0.5573	30s	16.21G
B	1	b	0.7389	0.7241	0.5582	31s	20.10G
C	2	a	0.7238	0.7156	0.5567	30s	14.26G
D	2	b	0.7325	0.7192	0.5580	30s	18.14G
accumu. prompt			0.7314	0.7163	0.5541	30s	18.14G

10. Time complexity

Time complexity. As shown in Table. A2, we compare the time cost with the other two most recent training-free methods [21, 54]. Thanks to parallel inference, our method achieves a similar runtime to the baseline while delivering better results. In contrast, other methods typically require LLMs [15] for local inference or rely on test-time optimization during the denoising process, often resulting in suboptimal time consumption.

11. SCB details

Style Composition Benchmark(SCB) construction. Aiming to comprehensively evaluate text-to-image generation under multi-style composition scenarios, We adopt a unified prompt format—“A [type of style] style [subject] in a [type of style] style [background].” As shown in Table A3, the “[type of style]” column lists the common style categories used for prompt construction. while [subject] and [background] are randomly selected from the corresponding word pool in the table. Notably, to better evaluate the ability of different models to maintain style consistency, we ensure that the [subject] and [background] within the same prompt do not share the same style.

SCB evaluation metrics. We propose a novel style alignment evaluation methodology that employs object detection to segment composite images into subject and background components. Each segment is independently analyzed against its intended style description using CLIP score metrics. The final averaged score reflects the overall fidelity of the composite image to its intended style descriptors across all elements. As illustrated in Fig. A1 A2,

our method effectively evaluates style alignment across composite images. In Fig. A1, the robot incorporates oil painting stylistic elements while the forest background exhibits traces of sketch characteristics, resulting in a lower CLIP score (0.2344) that indicates reduced alignment due to style contamination. In contrast, Fig. A2 demonstrates clearer stylistic delineation—a sketch-style robot (0.2818) against an oil painting forest background (0.2493)—yielding a higher average alignment score (0.2656). This granular approach enables precise measurement of style fidelity across different image elements. Noticed that, why we use a cropping way instead of a segmentation way, is the Clip cannot accurately recognize the images with partial vacancy.

12. More Visualization Results

In this section, we provide additional generation outputs of our method, demonstrating that it can produce images with high aesthetic visual quality, which can be seen in Fig. A3.

Table A3. Style Composition Benchmark

Type of style	Lego; Oil-painting; Cyberpunk; Sketch; Pixel-Art; Watercolor; Graffiti
Subject	Dog; Cat; Robot; Car; Unicorn
Background	forest; Space; Desert; City; Ruins

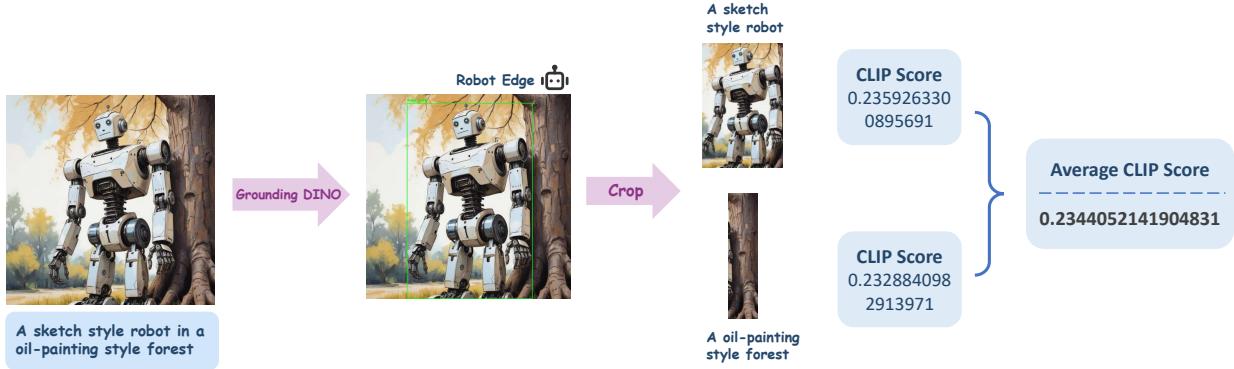


Figure A1. Pipeline of measure the style accuracy in circumstances of blending.

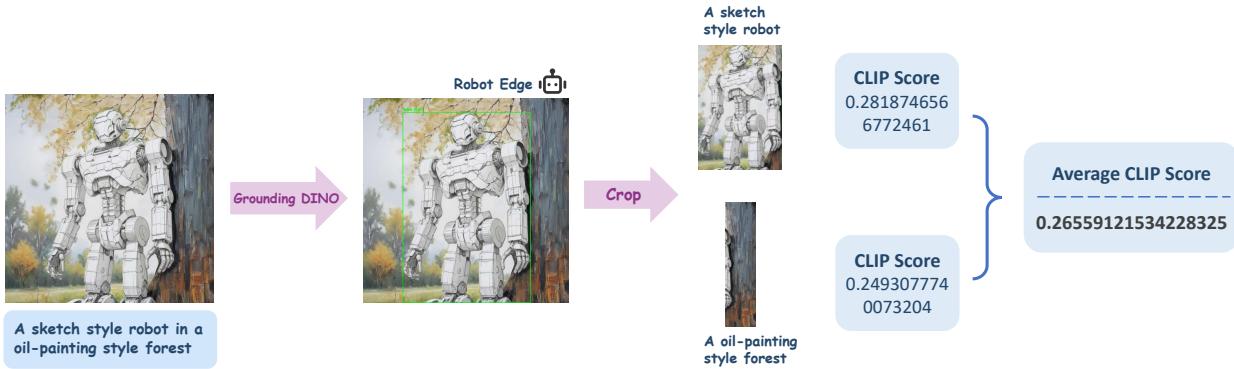


Figure A2. Pipeline of measure the style accuracy in circumstances for our methods.

13. Comparison between SA and CA mechanism

In this section, we present key concepts and techniques that serve as the foundation for our proposed method. These preliminary components provide a clearer understanding of the motivations behind our approach and the technical mechanisms leveraged in its design.

14. Why first 80% timesteps self-extension?

Why SA map? While many recent studies [5, 34, 38, 51] indicate that self-attention maps carry structural information and thus have strong layout-preserving capabilities, the cross-attention map-based editing approach [18] still has a significant impact. So, why choose to replace the self-attention map rather than the cross-self attention map for progressive editing? Here, we visualize the difference between replacing self-attention maps (SA) and cross-attention maps (CA) at various steps. As shown in Fig. A6, editing based on the self-attention map tends to provide stronger editing capabilities, consistent with the findings in FPE [34].

We hypothesize that the relatively weaker effect of edit-

ing based on cross-attention arises because it only replaces or influences the attention map of the edited object. In contrast, editing based on the self-attention map is more like conditional generation under the prior assumption of a diffusion model’s layout, which resembles the effect described in [61]. In this way, our framework, **Detail++**, ensures that each editing step is guided by an accurate detail binding relation prompt, enabling precise detail binding and preventing issues such as overflow, mismatching, and blending.

Why 80%? A higher percentage of self-attention map extension across all branches results in better control over the cohesive layout. However, increasing the number of steps diminishes the original generation capability of the current branch, which may lead to images with lower fidelity. This phenomenon is illustrated in Fig. A7. As observed, when S is very small, the layout of each branch’s output is inconsistent, leading to structural inconsistency. On the other hand, when S is very large, such as when $S = T$, the image quality degrades, especially at the edges of the subjects, where artifacts appear, as indicated in the marked areas of the image.



A pale peach-colored fox mascot with sleek fur, a **metallic silver** jacket, a **neon purple** scarf, and a **mint-green** hat sits atop a dreamy turquoise log in the Forest of Dreams.



A magical deer, with iridescent fur of mint green, lavender and antlers adorned with **glittering golden** vines, stands in an enchanted forest of **coral pink** trees under a **turquoise** sky.



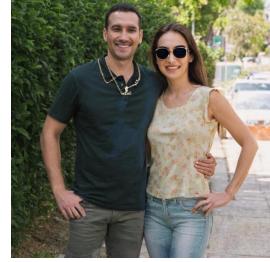
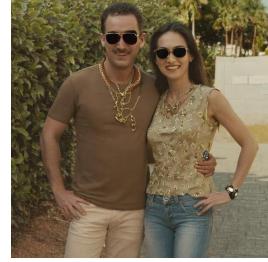
red teddy bear wearing a **green tracksuit** and a **yellow hat** is eating a **blue cake**.



A cream-colored bunny mascot with floppy ears, wearing a **pastel teal scarf**, **rose gold sunglasses**, and a **lavender backpack**, standing in a glowing coral-pink meadow under a twilight sky.



A **dog wearing sunglasses** and a **cat wearing necklace**.



Full body shot: a **man wearing necklace** and a **woman wearing sunglasses**

Figure A3. Our methods can generate images with accurate detail binding also high fidelity.

15. Cross-attention Layers Selection

In this section, we compare different layer selections for binary mask generation. We visualize cross-attention maps from various layers (see Fig. A4) and provide additional detailed visualizations in Fig. A5. It is evident that extracting the 32×32 resolution cross-attention map from the down, up, and mid blocks produces the cleanest and most accurate subject mask.

16. Combination with Other Methods

Our method can also be successfully integrated with other commonly used techniques in AIGC, such as Control-net [61]. In this section, we demonstrate the efficiency of our approach when combined with Control-net, as shown in Fig. A8.

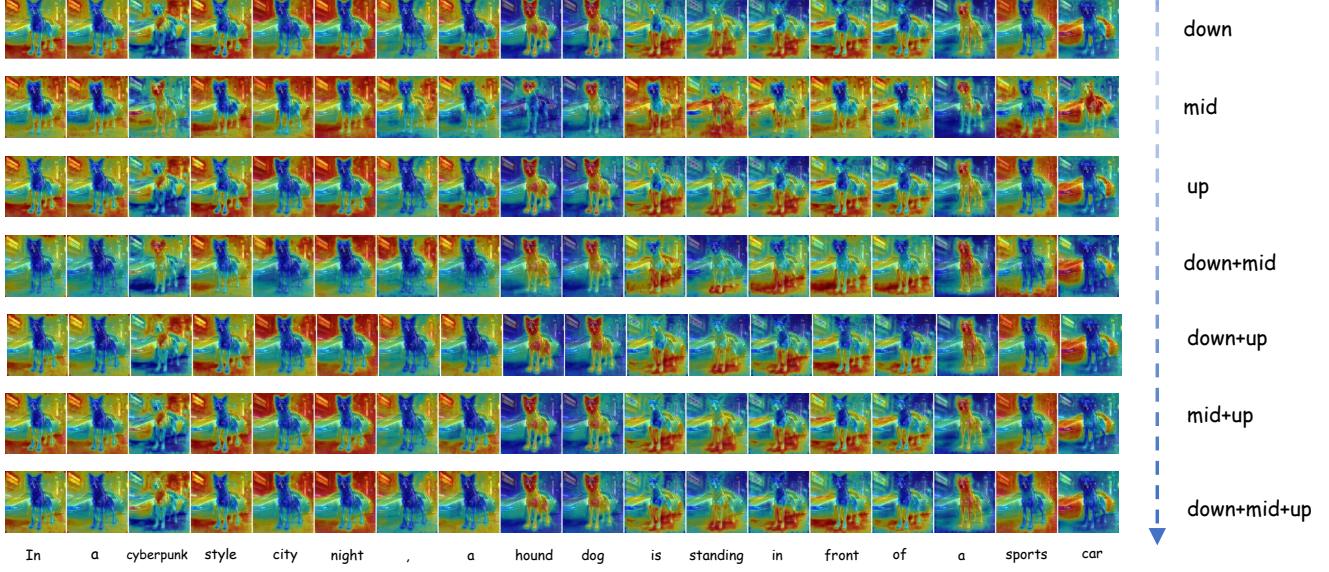


Figure A4. Visualization of cross-attention maps at different layers. “Down” denotes the cross-attention layers in the first downsampling block, while “Up” indicates those in the second upsampling block. Other layers with a 64×64 cross-attention map require significantly more memory and are thus omitted here to avoid excessive GPU usage.



Figure A5. Comparison of the effects of different cross-attention layer selections in binary mask extraction. It can be observed that extracting the 32×32 resolution cross-attention map from the down, up, and mid blocks yields the cleanest and most accurate subject mask.

Figure A6. Comparative analysis of attention map types in image generation. As indicated by the marked areas, replacing the cross-attention map results in weaker editing effects, making it less effective at accurately assigning attributes to the subject. In contrast, replacing the self-attention map allows for more precise attribute assignment in the same scenario.

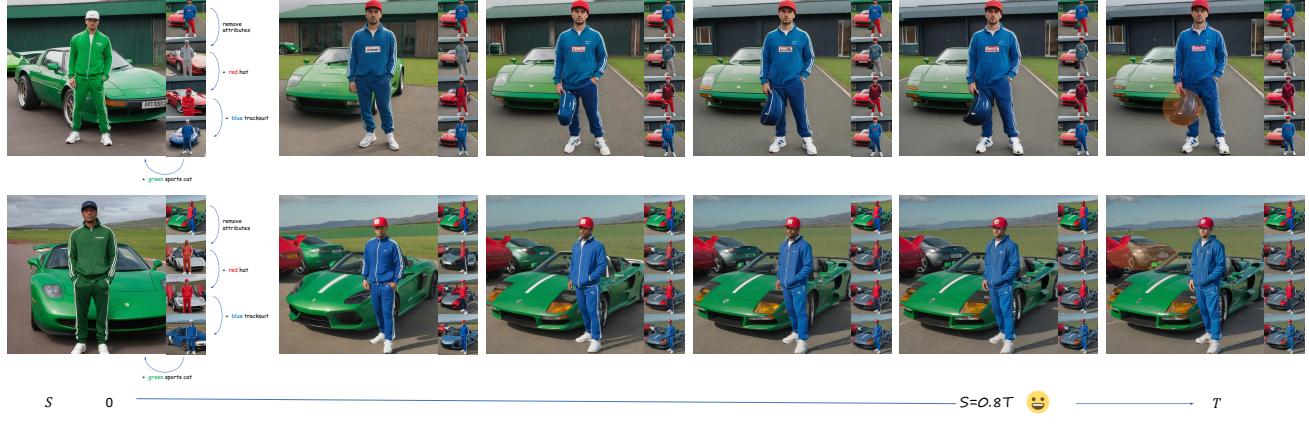


Figure A7. Comparative analysis of different attention extension timesteps. Both rows of images are generated from the prompt “A man wearing a red hat and blue tracksuit is standing in front of a green sports car.” The subfigures to the right of each main image represent the progressive generation process, where layout changes can be observed as the process advances.



Figure A8. Our methods can generate images with accurate detail binding.