

MS-DIFFUSION: MULTI-SUBJECT ZERO-SHOT IMAGE PERSONALIZATION WITH LAYOUT GUIDANCE

Anonymous authors

Paper under double-blind review



Figure 1: Representative outputs showcase the capabilities of MS-Diffusion in typical applications. The MS-Diffusion framework facilitates personalization across both single-subject scenarios (the upper panel) and multi-subject contexts (the lower panel). Notably, while preserving the intricacies of subject detail, MS-Diffusion achieves a marked enhancement in textual fidelity.

ABSTRACT

Recent advancements in text-to-image generation models have dramatically enhanced the generation of photorealistic images from textual prompts, leading to an increased interest in personalized text-to-image applications, particularly in multi-subject scenarios. However, these advances are hindered by two main challenges: firstly, the need to accurately maintain the details of each referenced subject in accordance with the textual descriptions; and secondly, the difficulty in achieving a cohesive representation of multiple subjects in a single image without introducing inconsistencies. To address these concerns, our research introduces the MS-Diffusion framework for layout-guided zero-shot image personalization with multi-subjects. This innovative approach integrates grounding tokens with the feature resampler to maintain detail fidelity among subjects. With the layout guidance, MS-Diffusion further improves the cross-attention to adapt to the multi-subject inputs, ensuring that each subject condition acts on specific areas. The proposed multi-subject cross-attention orchestrates harmonious inter-subject compositions while preserving the control of texts. Comprehensive quantitative and qualitative experiments affirm that this method surpasses existing models in both image and text fidelity, promoting the development of personalized text-to-image generation.

1 INTRODUCTION

Recent advancements in text-to-image (T2I) diffusion methodologies (Rombach et al., 2022; Saharia et al., 2022; Betker et al., 2023) have propelled the field to new heights, realizing unprecedented levels of photorealism while demonstrating a refined ability to conform to textual prompts. These achievements have spurred the development of a broad spectrum of applications, most prominently in the domain of personalized T2I (P-T2I) models, which are tasked with the complex undertaking of assimilating and regenerating novel visual concepts or subjects across diverse contexts with a heightened demand for conceptual and compositional fidelity. Despite fine-tuning-based techniques such as DreamBooth (Ruiz et al., 2023) and Textual Inversion (Gal et al., 2023) yield results with considerable accuracy, they necessitate extensive resources for tuning individual instances and for the storage of multiple models, which renders them less feasible for widespread application. To circumvent these resource-intensive requirements, fine-tuning-free alternatives have come to the fore.

Single-subject driven personalization methods, IP-Adapter (Ye et al., 2023) and ELITE (Wei et al., 2023) for instance, introduce a specialized cross-attention mechanism that distinctly processes text and image features, thereby affording the possibility to employ reference images directly as visual prompts within the model. Furthermore, recent works have employed multi-subject driven customization methodologies to concatenate visual concepts with textual prompts, offering a glimpse of the potential in techniques like SSR-Encoder (Zhang et al., 2024), λ -ECLIPSE (Patel et al., 2024), Emu2 (Sun et al., 2023), and KOSMOS-G (Pan et al., 2023). These models harness identity data and amalgamate it with text via cross-attention, exhibiting proficiency in adjusting textures. Nevertheless, above zero-shot personalization methods encounter limitations, notably in adapting a pressing question pertains to the congruence of granular details between the subject depicted in the synthesized imagery and its corresponding subjects, along with the degree of semblance between the content of the generated image and associated textual descriptions. The challenge is further amplified in scenarios requiring the personalization of multiple subjects. Especially, the challenge of ensuring harmonious representation when multiple subjects are incorporated—specifically, the elucidation of whether the resultant image manifests any discordant elements or deleterious interactions in accordance with textual directives and multi-subject referential controls. As illustrated in Figure 2, multi-subject personalization methods frequently incur notable detail inaccuracies in a fine-tuning-free framework and often lead to subject neglect, subject overcontrol, and subject-subject conflict issue within the generated images.

To confront these identified challenges, we are the **first** to introduce the layout-guided zero-shot image personalization with multiple subjects (MS-Diffusion) framework, which consolidates the accommodation of multiple subjects, the incorporation of zero-shot learning capabilities, the provision of layout guidance, and the preservation of the foundational model’s parameters. Firstly, we design the grounding resampler to extract the subject detailed features and integrate them with grounding information containing entities and boxes. As an image projection module, the proposed grounding resampler can enhance the subject fidelity while appending semantic and positional priors. Secondly, we propose a novel cross-attention mechanism for multiple subjects, which confines subjects to represent themselves in specific areas. This confluence not only facilitates the efficacious infusion of multi-subject data into the model but also mitigates conflicts between text and image subject control conditions. Such an approach culminates in the refined granularity of control over the image’s multi-subject composition. The experimental results demonstrate our method consistently outperforms the state-of-the-art approaches on all the benchmarks. **We conclude the previous P-T2I works and provide an overall comparison in Table 1.** The contributions can be summarized as follows:

- We introduce a layout-guided, zero-shot multi-subject image personalization framework within the diffusion model paradigm, designated as ‘MS-Diffusion’. This innovation streamlines the complex process of preserving detailed subject references. Moreover, it seamlessly integrates multiple subjects into a coherent and harmonious personalized image.
- The ‘Grounding Resampler’ is advanced as a novel feature refinement mechanism. This construct enriches the detail extraction from images by ascertaining the correlative content and fusing it with box embeddings that demarcate the anticipated spatial zones for each subject. Additionally, we introduce a specialized multi-subject cross-attention mechanism, confronting and rectifying prevalent complications in multi-subject personalization, including inadvertent subject neglect, disproportionate subject dominance, and internecine subject conflicts.

Table 1: **An overview of previous studies of P-T2I tasks.** MS-Diffusion is the **first** approach to support multi-reference zero-shot P-T2I generation with layout guidance and base model freezing.

Method	Zero Shot	Multi Subject	Base Model Freezing	MLLM Free	Layout Guidance
Textual Inversion (Gal et al., 2023)	✗	✗	✓	✓	✗
DreamBooth (Ruiz et al., 2023)	✗	✗	✗	✓	✗
ELITE (Wei et al., 2023)	✓	✗	✗	✓	✗
BLIP-Diffusion (Li et al., 2023a)	✓	✗	✗	✗	✗
IP-Adapter (Ye et al., 2023)	✓	✗	✓	✓	✗
Emu2 (Sun et al., 2023)	✓	✓	✗	✗	✗
Kosmos-G (Pan et al., 2023)	✓	✓	✓	✗	✗
λ-ECLIPSE (Patel et al., 2024)	✓	✓	✓	✗	✗
SSR-Encoder (Zhang et al., 2024)	✓	✓	✓	✓	✗
MS-Diffusion	✓	✓	✓	✓	✓

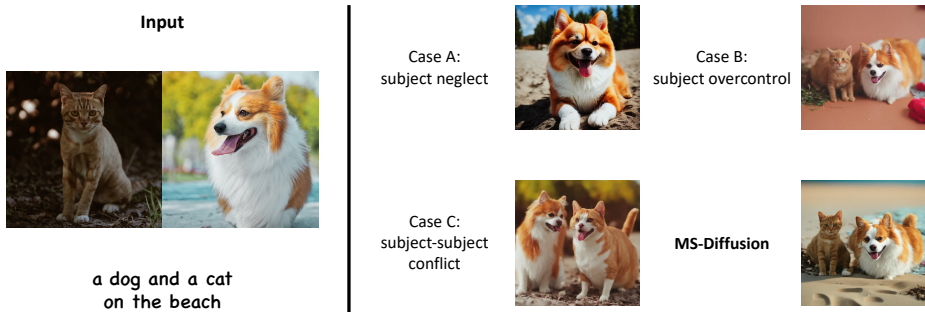


Figure 2: **Challenges inherent in multi-subject personalization approaches.** Through the explicit layout guidance, MS-Diffusion addresses these challenges by directing subject conditioning to specific areas, simultaneously maintaining high image fidelity.

- The capabilities of 'MS-Diffusion' are empirically substantiated through its ability to synthesize a broader spectrum of images with notable fidelity. The paper further delineates comprehensive ablation studies, underpinning the rationale behind our design decisions and affirming the efficacy of our proposed approach.

2 RELATED WORK

2.1 TEXT-TO-IMAGE GENERATION

Text-to-image generative models (Saharia et al., 2022; Bao et al., 2023; Esser et al., 2024; Podell et al., 2023) are capable of producing high-quality images using user-provided text prompts. In recent times, diffusion-based models have shown strong performance in text-to-image tasks. Stable Diffusion (Rombach et al., 2022) proposes conducting the diffusion process in latent space rather than pixel space, which reduces the sampling steps without compromising image quality. Kandinsky (Razhigaev et al., 2023) takes both text embedding and image embedding as conditions to generate images more controllably. DALL-E-3 (Betker et al., 2023) recaptions the training data pairs and utilizes T5-XXL (Chung et al., 2022) as the text encoder to strengthen the prompt-following ability. StableCascade (Pernias et al., 2023) presents a cascaded architecture to leverage outputs of front stages as priors, further reducing the latent space. PixArt-α (Chen et al., 2023a) also employs a large T5 text encoder and replaces the original U-Net backbone with a transformer (Peebles & Xie, 2023). These models focus on the basic text-to-image ability and cannot handle the situation when users provide specific subjects.

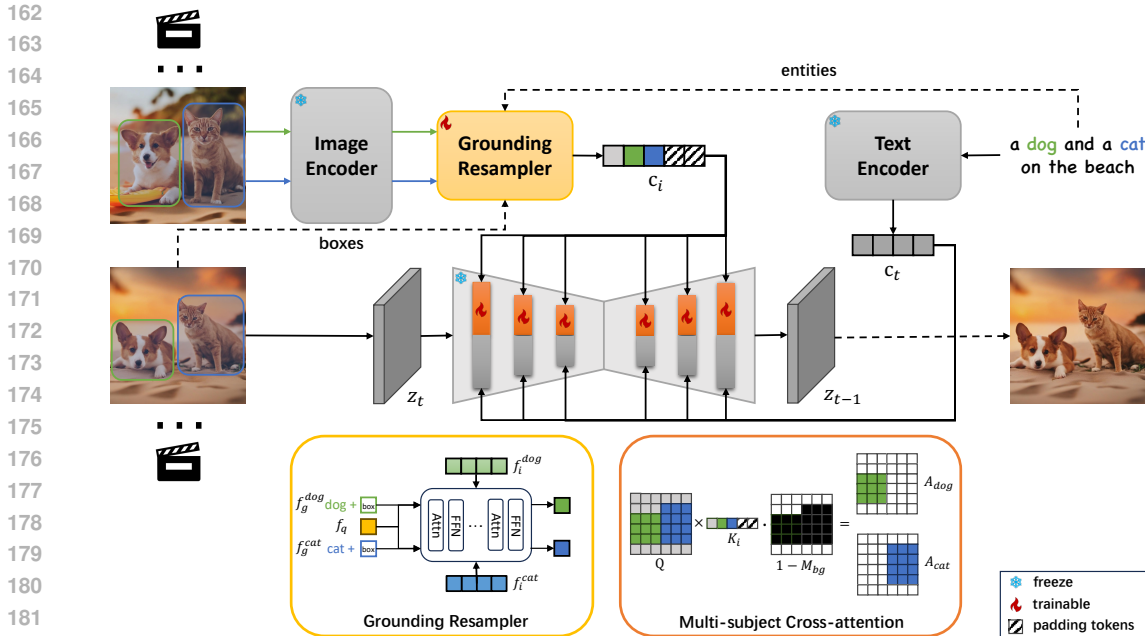


Figure 3: **The overall pipeline of MS-Diffusion.** It introduces two pivotal enhancements to the model: the grounding resampler and multi-subject cross-attention mechanisms. Firstly, the grounding resampler adeptly assimilates visual information, correlating it with specific entities and spatial constraints. Subsequently, the cross-attention mechanism facilitates precise interactions between the image condition and the diffusion latent within the multi-subject attention layers. Throughout the training phase, all components of the pre-existing diffusion model remain frozen.

To finely control text-to-image generation, some diffusion models support users in providing layout guidance. Layout Diffusion (Zheng et al., 2023) and GLIGEN (Li et al., 2023c) input positions and labels of bounding boxes into the diffusion model and train it to learn the layout information. DenseDiffusion (Kim et al., 2023) develops a training-free method and modulates the attention maps in the inference phase. Instance Diffusion (Wang et al., 2024b) and MIGC (Zhou et al., 2024) extend the layout-conditioned diffusion to the instance level, enabling the model to generate multiple objects with precise quantities. While layout-guided diffusion models have robust controllability, they cannot reference specific concepts, which is emphasized in personalized text-to-image generation.

2.2 TEXT-TO-IMAGE PERSONALIZATION

Text-to-image personalization (Han et al., 2023; Qiu et al., 2023; Chen et al., 2024; Shi et al., 2024; Hu et al., 2024) has attracted much attention in the community for its powerful referencing capability of both text and image prompts. Textual inversion (Gal et al., 2023) and DreamBooth (Ruiz et al., 2023) utilize an identifier in the text to bind the visual concept through fine-tuning. IP-Adapter (Ye et al., 2023) proposes a zero-shot personalized model by projecting the image embedding to the cross-attention layers. InstanceID (Wang et al., 2024a) develops an approach for identity personalization, replacing the image encoder with a face encoder and employing ControlNet (Zhang et al., 2023) to integrate the face landmarks. To narrow the gap between image and text prompts, Kosmos-G (Pan et al., 2023) and λ -ECLIPSE (Patel et al., 2024) conduct multi-modal training to unite the inputs by text-image interleaving. SSR-Encoder (Zhang et al., 2024) design a query network to extract a single subject from images with multiple subjects for personalization.

Though past research in this field has significantly enhanced the ability to reference single subjects, few studies have explored zero-shot multi-subject personalized models. Moreover, existing related works struggle to address conflicts in personalized generation with multiple subjects and generate bad results, which is precisely the focus of our work, to discuss and resolve these issues.

216 3 METHOD

217
218 The pipeline of MS-Diffusion is shown in Figure 3. Through an improved data construction strategy,
219 we can get multiple subject images together with the corresponding entities and layouts as input. We
220 propose a grounding resampler to separately extract the image features and integrate them with phrase
221 embedding and box embedding for condition enhancement. Inside the cross-attention layers, we
222 further introduce the masked cross-attention to guide the generation with layout priors and alleviate
223 conflicts of multiple subjects. The training needs no pre-trained weights to be optimized and remains
224 plug-and-play in various base models.

225 3.1 PRELIMINARIES

226
227 **Stable Diffusion with Image Prompt.** As a widely used diffusion model, Stable Diffusion
228 (SD) (Rombach et al., 2022) conducts the diffusion process in the latent space. Given an image
229 and a text prompt, SD encodes them into latent code \mathbf{z} and condition embedding \mathbf{c}_t utilizing
230 VAE (van den Oord et al., 2017) and CLIP (Radford et al., 2021) text encoder, respectively. In
231 zero-shot image personalization architectures like IP-Adapter (Ye et al., 2023), images can also be
232 considered a condition of the diffusion model. Specifically, a subject image is encoded to image
233 embeddings by an image encoder and then projected into the original condition space of the diffusion
234 model denoted as \mathbf{c}_i . For a timestep t which is uniformly sampled from a fixed range, the model θ
235 predicts the noise ϵ_θ and is optimized through the objective:

$$236 \mathcal{L}_{IP} = \mathbb{E}_{\mathbf{z}, \mathbf{c}, \epsilon, t} \left[\|\epsilon - \epsilon_\theta(\mathbf{z} | \mathbf{c}_t, \mathbf{c}_i, t)\|_2^2 \right] \quad (1)$$

237 where $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. In this work, we employ \mathcal{L}_{IP} with SDXL (Podell et al., 2023) as the pre-trained
238 model, which contains two CLIP text encoders and additional condition inputs besides \mathbf{c} and t ,
239 omitted for brevity.

240
241 **Cross-attention.** In IP-Adapter Ye et al. (2023), both \mathbf{c}_i and \mathbf{c}_t are integrated into the U-Net
242 backbone through cross-attention layers:

$$243 \text{Attn}(\mathbf{Q}, \mathbf{K}_i, \mathbf{K}_t, \mathbf{V}_i, \mathbf{V}_t) = \gamma \cdot \mathbf{z}_{img} + \mathbf{z}_{txt} = \gamma \cdot \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}_i^T}{\sqrt{d}} \right) \mathbf{V}_i + \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}_t^T}{\sqrt{d}} \right) \mathbf{V}_t \quad (2)$$

244 where $\mathbf{Q} = \mathbf{z}W_q$, $\mathbf{K}_i = \mathbf{c}_iW_k^i$, $\mathbf{K}_t = \mathbf{c}_tW_k^t$, $\mathbf{V}_i = \mathbf{c}_iW_v^i$, $\mathbf{V}_t = \mathbf{c}_tW_v^t$, and W_q, W_k, W_v are
245 corresponding projection weight matrices. And d represents the dimensionality of the key vectors
246 Note that the key and value matrix of \mathbf{c}_i and \mathbf{c}_t are independent of each other to decouple conditions
247 of different modalities. Previous studies (Hertz et al., 2023; Tang et al., 2023) have found that
248 attention maps $\mathbf{A} = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}} \right)$ can reflect the attribution relation between generated
249 images and conditions, which means that they determine the effect of condition controls.

250 3.2 DISCUSSION ON MULTI-SUBJECT IMAGE PERSONALIZATION

251
252 A widely used method for achieving multi-subject image personalization involves training individual
253 models for each subject, followed by their integration. Tuning-based methods (Gu et al., 2023;
254 Avrahami et al., 2023; Liu et al., 2023b; Kumari et al., 2023) with improvements on visual conflicts
255 can produce impressive multi-subject personalized images. However, zero-shot methods eliminate
256 the need for individual subject tuning or the merging of different combinations, significantly reducing
257 costs and enhancing the practicality of multi-subject personalization. This makes research into
258 zero-shot multi-subject image personalization both necessary and promising.

259
260 Most of the relevant studies (Ma et al., 2024; Gu et al., 2023; Liu et al., 2023b; Xiao et al., 2024)
261 focus on mitigating visual conflicts in text cross-attentions of the base model. While modifying
262 text cross-attention can be effective, it presents certain limitations. First, adjustments to text cross-
263 attention can directly impact the control over text conditions. Second, text cross-attention does not
264 directly dictate the areas of influence for image conditions; rather, it exerts an indirect influence on
265 image conditions by shaping the image layout generated by the diffusion model. This indirect control
266 may result in low performance and increased uncertainty.

3.3 DATA CONSTRUCTION

In the field of multi-subject personalization, creating a robust dataset architecture is a significant challenge, especially when no pre-existing dataset includes a variety of reference subjects with their validated truths. Our method starts with applying a Named Entity Recognition (NER) protocol to textual data to extract relevant entities. These entities are then used within a detection framework to define the corresponding bounding boxes. This step generates training samples that encompass a range of [*subjects, entities, spatial layouts*].

Previous studies have mostly created training examples from stand-alone images which is essentially a reconstruction task, leading the resulting models to favor replication, often resulting in 'copy-and-paste' artifacts (Chen et al., 2023b). To address this issue, our enhanced approach involves extracting subjects from a single video frame and using another frame from the same sequence as a reference for the truth. This technique effectively separates the personalized references from the target images. Due to possible variations in subjects between frames, we use a specialized subject-matching algorithm to ensure accurate matching. We provide a detailed description of this data processing pipeline in Section A.

3.4 GROUNDING RESAMPLER

Different from text embedding, image embedding generally contains more information and is sparser, making projection into the condition space challenging. Leveraging embeddings from all image patches primarily control the condition, but the pooled output from the image encoder tends to omit many details. Different from Flamingo (Alayrac et al., 2022) and IP-Adapter (Ye et al., 2023), we propose the integration of a grounding resampler that functions as an alternative form of image projector. Utilizing a set of learnable tokens, a *resampler* queries and distills pertinent information from the image features. Specifically, with an image embedding f_i and a learnable query f_q , the resampler comprises several attention layers:

$$\text{RSAttn} = \text{Softmax} \left(\frac{\mathbf{Q}(f_q) \mathbf{K}^T([f_i, f_q])}{\sqrt{d}} \right) \mathbf{V}([f_i, f_q]) \quad (3)$$

where $[f_i, f_q]$ denotes the concatenation of the image embedding f_i and the learnable query f_q . The architecture incorporates fully connected feedforward networks (FFNs), analogous to those utilized in standard vision transformers (Dosovitskiy et al., 2021).

As detailed in 3.3, we can obtain entities of multiple referenced subjects and their target area boxes in the generated image. We present to initialize the query f_q with grounding tokens f_g derived from text embedding of entities and Fourier embedding of boxes. Entities are related to the semantic information of images and boxes indicate the areas where the subjects are supposed to be. It would be helpful for the resampler to extract the image features appropriately and the cross-attention layers to condition the generation finely. To prevent the model from becoming dependent on the grounding tokens during inference, we randomly replace them with the original learnable queries in the training. For the input of n subjects, the projection processes of different subject images do not affect each other. The resulting n queries will be concatenated and input into the subsequent model as \mathbf{c}_i with $N = n * n_t$ tokens, where n_t is the token quantity per subject.

3.5 MULTI-SUBJECT CROSS-ATTENTION

In scenarios involving the generation of multiple subjects, challenges frequently arise that are not exclusive to personalization tasks. These include discordances between subjects and their backgrounds and amongst the subjects themselves. A viable solution to mitigate such conflicts leverages attention masks, contingent upon the availability of layout priors. The incorporation of attention masks within cross-attention mechanisms facilitates the exclusion of padding tokens from the condition, thus minimizing their impact.

To confine the context of each subject to a designated spatial domain, we propose an enhancement to the conventional attention mask, denoted as \mathbf{M} . This adjustment involves the bilateral neglect of

tokens within both the query and key matrices, applied specifically for the j th subject as follows:

$$\mathbf{M}_j(x, y) = \begin{cases} 0 & \text{if } [x, y] \in B_j \\ -\infty & \text{if } [x, y] \notin B_j \end{cases} \quad (4)$$

Here, B_j signifies the coordinate set of bounding boxes pertaining to the j th subject. By this means, the conditional image latent $\hat{\mathbf{z}}_{img}$ is derived through:

$$\hat{\mathbf{z}}_{img} = \text{Softmax} \left(\frac{\mathbf{Q}\mathbf{K}_i^T}{\sqrt{d}} + \mathbf{M} \right) \mathbf{V}_i \quad (5)$$

Herein, \mathbf{M} represents the amalgamation of all subject-specific masks, $\text{Concat}(\mathbf{M}_0, \dots, \mathbf{M}_n)$. In this way, the model ensures each subject to be represented in a certain area, thus resolving the issues of subject neglect and conflict in Figure 2.

However, an inherent limitation arises when a query patch token is ubiquitously masked across all referenced subjects or remains unmasked (in instances of overlapping bounding boxes), thereby diminishing the intended efficacy of multi-subject cross-attention. To counteract this, we introduce dummy tokens initialized at random preceding the image tokens to symbolize the background. This strategy is instrumental in ensuring that text conditions predominantly govern areas devoid of any guided layout, thereby solving the subject overcontrol issue in Figure 2. Following the acquisition of \mathbf{A} , we apply \mathbf{M}_{bg} to seamlessly mask these tokens within $\hat{\mathbf{z}}_{img}$, as illustrated:

$$\mathbf{z}_{img} = (1 - \mathbf{M}_{bg}) \cdot \hat{\mathbf{z}}_{img} \quad (6)$$

where \mathbf{M}_{bg} is articulated as a binary mask, with elements within the subject bounding boxes designated as zero. In contrast to the methods discussed in Section 3.2, our proposed multi-subject cross-attention directly manages the image conditions by employing masked image cross-attention in the targeted areas. While addressing multi-object conflicts, our method ensures that text conditions remain unaffected, as evidenced by the significantly higher text adherence capability shown in Table 2. Notably, certain studies have sought to resolve these conflicts through the application of objectives on attention maps within the cross-attention mechanism. A series of rigorous experiments have been conducted to substantiate our design, with the details elucidated in Section 4.4.

4 EXPERIMENTS

4.1 EXPERIMENT SETUP

Datasets. For training, we utilize an in-house video dataset that contains 3.6M video clips. For evaluation, we measure the single-subject and multi-subject performance on DreamBench (Ruiz et al., 2023) and MS-Bench, respectively. DreamBench includes 30 subjects and 25 prompts and we preset all input boxes to [0.25, 0.25, 0.75, 0.75]. To thoroughly assess the performance of multi-subject personalization, we have established a new evaluation standard, MS-Bench, which includes 40 subjects and 1148 combinations, each combination having up to 6 prompts, totaling 4488 distinct test samples. Details of datasets are provided in Section A and Section B.

Evaluation metrics. Following previous works, we measure the performance through image and text fidelity. To assess image fidelity, we employ cosine similarity measures between generated images and subject images within CLIP (Radford et al., 2021) and DINO (Caron et al., 2021) spaces, referred to as CLIP-I and DINO, respectively. For text fidelity, we utilize cosine similarity between generated images and text prompts in CLIP space, denoted as CLIP-T. In multi-subject personalization, using the average fidelity to reflect image fidelity is insufficient, as it fails to reveal cases of subject neglect. We further employ the product of multi-subject DINO, denoted as M-DINO, to indicate whether each subject has been recreated in the results.

Baselines. For single-subject personalization, we compare our model with methods mentioned in Table 1. Emu2 (Sun et al., 2023), Kosmos-G (Pan et al., 2023), and λ -ECLIPSE (Patel et al., 2024) are all MLLM-based methods, while λ -ECLIPSE is reported to have better performance. Therefore, we select SSR-Encoder (Zhang et al., 2024) and λ -ECLIPSE (Patel et al., 2024) as baselines for multi-subject personalization. Considering the fairness, the qualitative results of MS-Diffusion are generated without any fine-tuning on benchmarks. The implementation details of MS-Diffusion and these methods are contained in Section C.

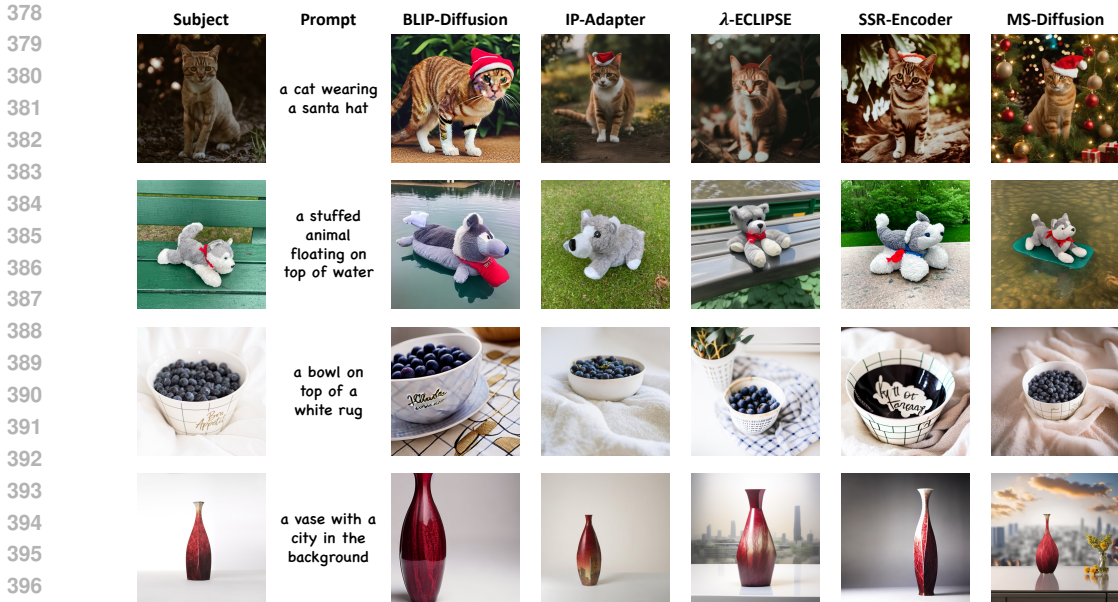


Figure 4: Qualitative results of MS-Diffusion and baselines in single-subject personalization.

Table 2: **Quantitative comparison on MS-Diffusion and baselines.** **Bold** and underline represent the highest and second-highest metrics, respectively. For single-subject, the results for IP-Adapter, Emu2, and Kosmos-G were obtained from λ-ECLIPSE (Patel et al., 2024), while the rest are reported in the corresponding papers. * denotes the model is fine-tuned on DreamBench.

Method	Single-subject			Multi-subject			
	CLIP-I	DINO	CLIP-T	CLIP-I	DINO	M-DINO	CLIP-T
Textual Inversion	0.780	0.569	0.255	-	-	-	-
DreamBooth	0.803	0.668	<u>0.305</u>	-	-	-	-
BLIP-Diffusion*	0.805	0.670	0.302	-	-	-	-
λ-ECLIPSE*	0.796	<u>0.682</u>	0.304	-	-	-	-
MS-Diffusion*	0.805	0.702	0.313	-	-	-	-
ELITE	<u>0.771</u>	<u>0.621</u>	<u>0.293</u>	-	-	-	-
BLIP-Diffusion	0.779	0.594	0.300	-	-	-	-
IP-Adapter	0.810	0.613	0.292	-	-	-	-
Emu2	0.765	0.563	0.273	-	-	-	-
Kosmos-G	0.822	0.618	0.250	-	-	-	-
λ-ECLIPSE	0.783	0.613	0.307	<u>0.724</u>	0.419	0.094	<u>0.316</u>
SSR-Encoder	<u>0.821</u>	0.612	<u>0.308</u>	0.725	0.425	<u>0.107</u>	0.303
MS-Diffusion	0.792	0.671	0.321	0.698	0.425	0.108	0.341

4.2 SINGLE-SUBJECT COMPARISON

In the single-subject comparison, a detailed examination is carried out utilizing both qualitative and quantitative comparisons to gauge the performance of different methodologies. On the qualitative front, as shown in Figure 4, MS-Diffusion shows an exceptional ability to generate single-subject images with high fidelity and detail retention. In quantitative results provided in Table 2, MS-Diffusion also achieves competitive scores, with obviously the highest DINO and CLIP-T scores at 0.671 and 0.321 respectively on zero-shot scenarios, and leading CLIP-I score of 0.792. For tuning methods, MS-Diffusion outperforms baselines in all metrics. As discussed in DreamBooth (Ruiz et al., 2023), DINO more accurately captures the similarity in details between the results and the



Figure 5: Qualitative results of MS-Diffusion and baselines in multi-subject personalization.

Table 3: **Ablation study of MS-Diffusion.** RS, GRS, MCA, TAL, IAL, and LG represent resampler, grounding resampler, multi-subject cross-attention, text attention loss, image attention loss, and layout guidance, respectively.

Method	Single-subject			Multi-subject			
	CLIP-I	DINO	CLIP-T	CLIP-I	DINO	M-DINO	CLIP-T
MS-Diffusion	0.792	0.671	0.321	0.698	0.425	0.108	0.341
w/o RS	0.775	0.583	0.320	0.680	0.372	0.082	0.336
w/o GRS	0.777	0.646	0.320	0.681	0.389	0.090	0.331
w/o MCA	0.798	0.662	0.312	0.693	0.422	0.100	0.309
w/o LG w/ IAL	0.761	0.577	0.284	0.675	0.377	0.080	0.305
w/o LG w/ IAL&TAL	0.809	0.660	0.293	0.687	0.413	0.093	0.316

labels, whereas CLIP-I may exhibit high scores in situations of background overfitting, resulting in a clear advantage for DINO, but a slight disadvantage for CLIP-I of MS-Diffusion.

4.3 MULTI-SUBJECT COMPARISON

From a qualitative perspective in Figure 5, MS-Diffusion manages to maintain natural interactions among subjects in generated images while ensuring each subject retains its distinctiveness and recognizability. Quantitatively, results in Table 2 demonstrate the strength of MS-Diffusion in DINO, M-DINO, and CLIP-T. Unlike in single-subject personalization, there is a larger gap in text fidelity between MS-Diffusion and the baselines in multi-subject personalization, demonstrating that MS-Diffusion not only effectively generates the multiple subjects outlined in the text but also excellently preserves the text control capabilities inherent to SD. Additionally, the image fidelity of MS-Diffusion is comparable, highlighting its superior ability to retain details, particularly significant as low text fidelity is commonly associated with overfitting.

4.4 ABLATION STUDY

Module ablation. We conduct an ablation experiment on the proposed two modules, grounding resampler (GRS) and multi-subject cross-attention (MCA), to validate their effects. For GRS, we replace it with a linear projection layer and a normal resampler (Alayrac et al., 2022; Ye et al., 2023).



Figure 6: **Visualization results of module ablation.** Models without RS or GRS have an obvious decrease in the detail-preserving capability. For the multi-subject generation, the model without MCA cannot handle the subject conflicts.

Results in Table 3 indicate that the resampler-like image projector significantly enhances the details, as evidenced by DINO being obviously higher than the linear projector. Moreover, The substantial improvement in multi-object image fidelity by GRS reflects the critical role of the information carried by grounding tokens in multi-object generation. As a key module for resolving conflicts, removing MCA results in a noticeable degradation of text fidelity, especially in multi-subject generation. The combined use of both modules ensures that MS-Diffusion maintains high image and text fidelity simultaneously. We provide visualization results regarding the module ablation in Figure 6. As clearly reflected by the qualitative examples, GRS enhances the details and MCA handles the conflicts.

Layout guidance. As mentioned in Section 3.5, we have explored the indispensable role of explicit layout guidance (LG), including grounding tokens and MCA. A straightforward approach to implicitly utilizing layout involves incorporating an attention loss during training. Besides the image attention loss (IAL), we also introduce text attention loss (TAL) to training by setting the original cross-attention layers trainable. The detailed loss definition is provided in Section G. As illustrated in Table 3, an objective to guide the image cross-attention helps the personalization hardly at all. TAL has somewhat resolved the conflict issues, but its performance is inferior to MS-Diffusion while introducing additional training parameters. We consider the inclusion of LG necessary and rational, not merely for the performance enhancements it offers, but also because it effectively resolves the various multi-object generation issues highlighted in Figure 2.

5 CONCLUSION

This study makes a significant contribution to the field of P-T2I diffusion models with the development of MS-Diffusion. This zero-shot framework excels at capturing intricate subject details and smoothly blending multiple subjects into a single coherent image. Equipped with the innovative Grounding Resampler and Multi-subject Cross-attention mechanisms, our model effectively overcomes common multi-subject personalization issues, such as subject neglect and conflict. Extensive ablation studies underscore MS-Diffusion’s enhanced performance in image synthesis fidelity compared to existing models. It stands as a groundbreaking approach for P-T2I applications that are free from the need for fine-tuning and require layout guidance.

REFERENCES

- 540
541
542 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
543 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford,
544 Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob L. Menick,
545 Sebastian Borgeaud, Andy Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski,
546 Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karén Simonyan. Flamingo: a visual
547 language model for few-shot learning. In *NeurIPS*, 2022.
- 548 Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene:
549 Extracting multiple concepts from a single image. In *SIGGRAPH Asia*, pp. 96:1–96:12. ACM,
550 2023.
- 551 Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth
552 words: A vit backbone for diffusion models. In *CVPR*, pp. 22669–22679. IEEE, 2023.
- 553 James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang
554 Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer*
555 *Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2:3, 2023.
- 556 Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image
557 editing instructions. In *CVPR*, pp. 18392–18402. IEEE, 2023.
- 558 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and
559 Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, 2021.
- 560 Hong Chen, Yipeng Zhang, Simin Wu, Xin Wang, Xuguang Duan, Yuwei Zhou, and Wenwu Zhu.
561 Disenbooth: Identity-preserving disentangled tuning for subject-driven text-to-image generation.
562 In *ICLR*. OpenReview.net, 2024.
- 563 Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James T.
564 Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- α : Fast training of diffusion transformer
565 for photorealistic text-to-image synthesis. *CoRR*, abs/2310.00426, 2023a.
- 566 Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor:
567 Zero-shot object-level image customization. *CoRR*, abs/2307.09481, 2023b.
- 568 Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li,
569 Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language
570 models. *arXiv preprint arXiv:2210.11416*, 2022.
- 571 Jiankang Deng, Jia Guo, Niannan Xue, and Stefanos Zafeiriou. Arcface: Additive angular margin
572 loss for deep face recognition. In *CVPR*, pp. 4690–4699. Computer Vision Foundation / IEEE,
573 2019.
- 574 Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas
575 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit,
576 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale.
577 In *ICLR*. OpenReview.net, 2021.
- 578 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
579 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, Dustin Podell, Tim Dockhorn, Zion English,
580 and Robin Rombach. Scaling rectified flow transformers for high-resolution image synthesis. In
581 *ICML*. OpenReview.net, 2024.
- 582 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and
583 Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using
584 textual inversion. In *ICLR*. OpenReview.net, 2023.
- 585 Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao,
586 Rui Zhao, Shuning Chang, Weijia Wu, Yixiao Ge, Ying Shan, and Mike Zheng Shou. Mix-of-
587 show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. In
588 *NeurIPS*, 2023.
- 589
590
591
592
593

- 594 Ligong Han, Yinxiao Li, Han Zhang, Peyman Milanfar, Dimitris N. Metaxas, and Feng Yang. Svdiff:
595 Compact parameter space for diffusion fine-tuning. In *ICCV*, pp. 7289–7300. IEEE, 2023.
596
- 597 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-
598 to-prompt image editing with cross-attention control. In *ICLR*. OpenReview.net, 2023.
599
- 600 Hexiang Hu, Kelvin CK Chan, Yu-Chuan Su, Wenhui Chen, Yandong Li, Kihyuk Sohn, Yang Zhao,
601 Xue Ben, Boqing Gong, William Cohen, et al. Instruct-imagen: Image generation with multi-modal
602 instruction. In *CVPR*, pp. 4754–4763, 2024.
- 603 Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri,
604 and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *CVPR*, pp.
605 6007–6017. IEEE, 2023.
- 606 Yunji Kim, Jiyoung Lee, Jin-Hwa Kim, Jung-Woo Ha, and Jun-Yan Zhu. Dense text-to-image
607 generation with attention modulation. In *ICCV*, 2023.
608
- 609 Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloé Rolland, Laura Gustafson, Tete
610 Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross B. Girshick.
611 Segment anything. In *ICCV*, pp. 3992–4003. IEEE, 2023.
- 612 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept
613 customization of text-to-image diffusion. In *CVPR*, pp. 1931–1941. IEEE, 2023.
614
- 615 Dongxu Li, Junnan Li, and Steven C. H. Hoi. Blip-diffusion: Pre-trained subject representation for
616 controllable text-to-image generation and editing. In *NeurIPS*, 2023a.
617
- 618 Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image
619 pre-training with frozen image encoders and large language models. In *ICML*, volume 202 of
620 *Proceedings of Machine Learning Research*, pp. 19730–19742. PMLR, 2023b.
- 621 Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li,
622 and Yong Jae Lee. Gligen: Open-set grounded text-to-image generation. *CVPR*, 2023c.
623
- 624 Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker:
625 Customizing realistic human photos via stacked ID embedding. In *CVPR*, pp. 8640–8650. IEEE,
626 2024.
- 627 Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei
628 Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for
629 open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023a.
630
- 631 Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli Zhao,
632 Jingren Zhou, and Yang Cao. Customizable image synthesis with multiple subjects. In *NeurIPS*,
633 2023b.
- 634 Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personalized
635 text-to-image generation without test-time fine-tuning. In *SIGGRAPH*, pp. 25. ACM, 2024.
636
- 637 Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhui Chen, and Furu Wei. Kosmos-g:
638 Generating images in context with multimodal large language models. *CoRR*, abs/2310.02992,
639 2023.
- 640 Maitreya Patel, Sangmin Jung, Chitta Baral, and Yezhou Yang. λ -eclipse: Multi-concept personalized
641 text-to-image diffusion models by leveraging CLIP latent space. *CoRR*, abs/2402.05195, 2024.
642
- 643 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *ICCV*, pp.
644 4172–4182. IEEE, 2023.
- 645 Pablo Pernias, Dominic Rampas, Mats Leon Richter, Christopher Pal, and Marc Aubreville.
646 Würstchen: An efficient architecture for large-scale text-to-image diffusion models. In *The
647 Twelfth International Conference on Learning Representations*, 2023.

- 648 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
649 Penna, and Robin Rombach. SDXL: improving latent diffusion models for high-resolution image
650 synthesis. *CoRR*, abs/2307.01952, 2023.
- 651 Zeju Qiu, Weiyang Liu, Haiwen Feng, Yuxuan Xue, Yao Feng, Zhen Liu, Dan Zhang, Adrian Weller,
652 and Bernhard Schölkopf. Controlling text-to-image diffusion by orthogonal finetuning. In *NeurIPS*,
653 2023.
- 654 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
655 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever.
656 Learning transferable visual models from natural language supervision. In *ICML*, volume 139 of
657 *Proceedings of Machine Learning Research*, pp. 8748–8763. PMLR, 2021.
- 658 Anton Razzhigaev, Arseniy Shakhmatov, Anastasia Maltseva, Vladimir Arkhipkin, Igor Pavlov,
659 Ilya Ryabov, Angelina Kuts, Alexander Panchenko, Andrey Kuznetsov, and Denis Dimitrov.
660 Kandinsky: An improved text-to-image synthesis with image prior and latent diffusion. In *EMNLP*,
661 pp. 286–295. Association for Computational Linguistics, 2023.
- 662 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
663 resolution image synthesis with latent diffusion models. In *CVPR*, pp. 10674–10685. IEEE,
664 2022.
- 665 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
666 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*,
667 pp. 22500–22510. IEEE, 2023.
- 668 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kam-
669 yar Seyed Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, Jonathan
670 Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with
671 deep language understanding. In *NeurIPS*, 2022.
- 672 Jing Shi, Wei Xiong, Zhe Lin, and Hyun Joon Jung. Instantbooth: Personalized text-to-image
673 generation without test-time finetuning. In *CVPR*, pp. 8543–8552, 2024.
- 674 Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Zhengxiong Luo, Yueze Wang,
675 Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models
676 are in-context learners. *CoRR*, abs/2312.13286, 2023.
- 677 Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus
678 Stenertorp, Jimmy Lin, and Ferhan Ture. What the DAAM: interpreting stable diffusion using cross
679 attention. In *ACL*, pp. 5644–5659. Association for Computational Linguistics, 2023.
- 680 Aäron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning.
681 In *NeurIPS*, pp. 6306–6315, 2017.
- 682 Qixun Wang, Xu Bai, Haofan Wang, Zekui Qin, and Anthony Chen. Instantid: Zero-shot identity-
683 preserving generation in seconds. *arXiv preprint arXiv:2401.07519*, 2024a.
- 684 Xudong Wang, Trevor Darrell, Sai Saketh Rambhatla, Rohit Girdhar, and Ishan Misra. Instancediffu-
685 sion: Instance-level control for image generation. In *CVPR*, pp. 6232–6242, 2024b.
- 686 Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Grounding
687 diffusion with token-level supervision. *CoRR*, abs/2312.03626, 2023.
- 688 Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. ELITE: encoding
689 visual concepts into textual embeddings for customized text-to-image generation. In *ICCV*, pp.
690 15897–15907. IEEE, 2023.
- 691 Zhichao Wei, Qingkun Su, Long Qin, and Weizhi Wang. Mm-diff: High-fidelity image personalization
692 via multi-modal condition integration. *CoRR*, abs/2403.15059, 2024.
- 693 Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer:
694 Tuning-free multi-subject image generation with localized attention. *International Journal of
695 Computer Vision*, pp. 1–20, 2024.

702 Hu Ye, Jun Zhang, Sib0 Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt
703 adapter for text-to-image diffusion models. *CoRR*, abs/2308.06721, 2023.
704

705 Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
706 diffusion models. In *IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris,*
707 *France, October 1-6, 2023*, pp. 3813–3824. IEEE, 2023.

708 Yuxuan Zhang, Yiren Song, Jiaming Liu, Rui Wang, Jinpeng Yu, Hao Tang, Huaxia Li, Xu Tang,
709 Yao Hu, Han Pan, et al. Ssr-encoder: Encoding selective subject representation for subject-driven
710 generation. In *CVPR*, pp. 8069–8078, 2024.
711

712 Guangcong Zheng, Xianpan Zhou, Xuewei Li, Zhongang Qi, Ying Shan, and Xi Li. Layoutdiffusion:
713 Controllable diffusion model for layout-to-image generation. In *CVPR*, pp. 22490–22499. IEEE,
714 2023.

715 Dewei Zhou, You Li, Fan Ma, Xiaoting Zhang, and Yi Yang. Migc: Multi-instance generation
716 controller for text-to-image synthesis. In *CVPR*, pp. 6818–6828, 2024.
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A TRAINING DATASET CONSTRUCTION PIPELINE

Figure 7 illustrates the construction pipeline of the training dataset. Firstly, we randomly select two frames from a video clip, one as the reference and the other as the ground truth. Both frames are captioned by BLIP-2 (Li et al., 2023b). Secondly, we utilize a NER model¹ to extract entities from the caption. Entities and images are then input into Grounding DINO (Liu et al., 2023a) to obtain the boxes, which are parts of the final input of the model. Taking the boxes as prompts of SAM (Kirillov et al., 2023), we can further obtain segmentation masks to extract subjects from the reference image. Since the entities in different frames can be different, we design a subject matcher, which finds the correspondence between the frames by conducting Hungarian Algorithm on the entity image embeddings. Frames in a video typically contain the same entities but exhibit clear differences in details such as angles and poses. This makes them highly suitable as training data for personalized image generation models, which can help mitigate the model’s tendency to copy-and-paste.

Our dataset comprises 2.8M general scenario videos and 0.8M product demonstration videos, where the former covers more scenarios and the latter has more clear subjects. 2-5 frames for each are adopted in the training. In practice, there may only be 1-2 subjects successfully matched. To ensure that the training data contains a sufficient number of reference subjects, for the targets where matching fails, we directly use the corresponding parts of the ground truth as references. Subjects that are too small, too large, or have imbalanced proportions are filtered out. Each training sample can have up to 4 subjects, and we pad in the ones with fewer than 4.

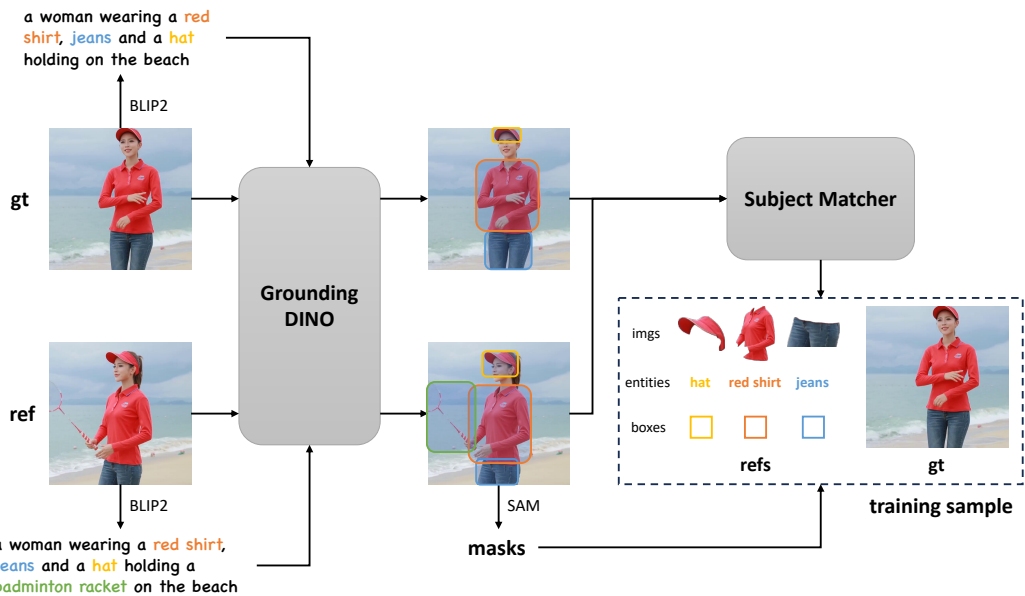


Figure 7: **Data construction pipeline of our work.** For the input of two frames, we can get subject images, entities, and boxes. Note that the entities and boxes are from the ground truth frame since they indicate the information in the generated result.

B DETAILS OF MS-BENCH

To construct MS-Bench, we collect subjects from previous studies (Ruiz et al., 2023; Gal et al., 2023; Kumari et al., 2023), the Internet², and an internal dataset that does not overlap with the training set. MS-Bench contains four data types and 13 combination types with two or three subjects. We provide the details in Table 4. Each combination type other than those related to the scene has 6 prompt variations. There are 1148 combinations and 4488 evaluation samples, where entities and

¹<https://spacy.io/>

²<https://unsplash.com/>

Table 4: **Explanation of MS-Bench.** Each combination type has preset prompts and boxes. [S] represents prompt variations about the scene, including "in a room", "in the jungle", "in the snow", "on the beach", "on the grass", and "on a cobblestone street".

Type	Prompt	Boxes
living+living living+object object+object	a {0} and a {1} [S]	[0.00, 0.25, 0.50, 0.75] [0.50, 0.25, 1.00, 0.75]
living+upwearing	a {0} wearing a {1} [S]	[0.25, 0.25, 0.75, 0.75] [0.25, 0.00, 0.75, 0.25]
living+midwearing living+wholewearing	a {0} wearing a {1} [S]	[0.25, 0.25, 0.75, 0.75] [0.25, 0.25, 0.75, 0.75]
midwearing+downwearing	a woman wearing a {0} and a {1} [S]	[0.25, 0.25, 0.75, 0.60] [0.25, 0.60, 0.75, 1.00]
living+scene object+scene	a {0} with a {1} in the background	[0.25, 0.25, 0.75, 0.75] [0.00, 0.00, 1.00, 1.00]
living+living+living object+object+object	a {0}, a {1}, and a {2} [S]	[0.00, 0.25, 0.35, 0.75] [0.35, 0.25, 0.65, 0.75] [0.65, 0.25, 1.00, 0.75]
living+object+scene	a {0} and a {1} with a {2} in the background	[0.00, 0.25, 0.50, 0.75] [0.50, 0.25, 1.00, 0.75] [0.00, 0.00, 1.00, 1.00]
upwearing+midwearing+ downwearing	a woman wearing a {0}, a {1}, and a {2} [S]	[0.25, 0.00, 0.75, 0.25] [0.25, 0.25, 0.75, 0.60] [0.25, 0.60, 0.75, 1.00]

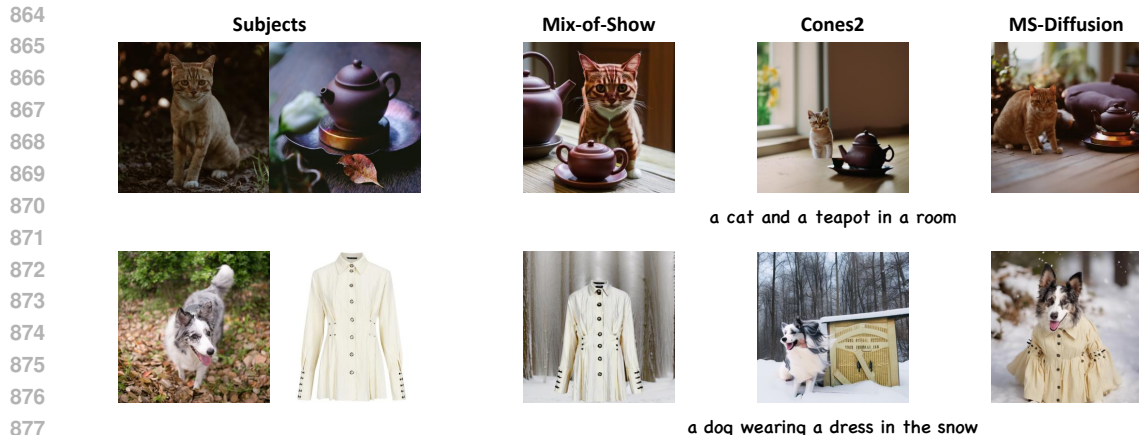
boxes are subject categories and preset layouts. Compared to other multi-subject benchmarks, our MS-Bench ensures that the model performance can be reflected comprehensively in abundant cases.

C EXPERIMENT SETTINGS

Training and Inference. The pre-trained model employed in MS-Diffusion is Stable Diffusion XL (SDXL) (Podell et al., 2023). Implemented by Pytorch 2.0.1 and Diffusers 0.23.1, our model is trained on 16 A100 GPUs for 120k steps with a batch size of 8 and a learning rate of 1e-4. Following the training of IP-adapter (Ye et al., 2023), we set $\gamma = 1.0$ in cross-attention layers and dropped the text and image condition using the same probability. To ensure the model is not dependent on the grounding tokens (Section 3.4), we also randomly drop them with a probability of 0.1. We generate five images for each sample during the inference, with unconditional guidance scale and γ set to 7.5 and 0.6, respectively, to get better results.

Comparative methods. Here we provide the details of the baselines compared in qualitative and quantitative experiments:

- **BLIP-Diffusion** (Li et al., 2023a) utilizes BLIP-2 (Li et al., 2023b) to unite the text and image embeddings. We implement it in the qualitative comparison using Diffusers.
- **IP-Adapter** (Ye et al., 2023) also uses image prompt as the condition. We run qualitative samples on their official code with the scale set to 0.5 recommended by the paper. Considering fairness, we use the result of SDXL in Table 2.
- **SSR-Encoder** (Zhang et al., 2024) design a query network to extract the specified subject of a single image, which enables it to finish multi-subject generation. We leverage it as one of the baselines in multi-subject personalization. For performance comparison, we



879 **Figure 8: Qualitative comparison of Mix-of-Show, Cones2, and MS-Diffusion.** MS-Diffusion
 880 outperforms Mix-of-Show and Cones2 in details preserving in the first row. Mix-of-Show and Cones2
 881 struggle in handling the interactions between subjects.

882

883

884 employed the official code provided, alongside the default hyperparameters specified in the
 885 code repository.

- 886 • λ -ECLIPSE (Patel et al., 2024) trains an independent multi-modal encoder and employs
 887 Kandinsky as the generative backbone. Since it outperforms other MLLM-based approaches,
 888 we choose it to be the representative. To facilitate a comparative analysis of performance, the
 889 study utilized the officially provided code, in conjunction with the default hyperparameters
 890 delineated within the corresponding code repository.

891 D MORE RESULTS OF SINGLE-SUBJECT PERSONALIZATION

892

893

894 Additional qualitative results on DreamBench are provided in Figure 12. MS-Diffusion shows
 895 excellent text fidelity in all subjects while keeping subject details, especially the living ones (dogs). It
 896 can be noticed that some elements in the background (the third line and the fourth line) also occur
 897 in the results (the grass and the teapot holder) since the entire images are referenced during the
 898 generation. Their scope of action depends on the input bounding box. In practical applications, using
 899 masked images as a condition is recommended.

900 One of the limitations in detail preservation is the insufficient capability of the image encoder. We pro-
 901 vide some uncommon examples in Figure 13. MS-Diffusion utilizes the CLIP image encoder, which
 902 results in the loss of some details in uncommon and complex cases. However, it still significantly
 903 outperforms state-of-the-art methods benefiting from the proposed grounding resampler.

904 E MORE RESULTS OF MULTI-SUBJECT PERSONALIZATION

905

906

907 We provide additional multi-subject personalized images based on MS-Bench in Figure 14. The
 908 results encompass various combination types, fully demonstrating the generalizability and robustness
 909 of MS-Diffusion. When the scene changes freely according to the text, the details of the subject are
 910 preserved without being affected. In addition to common parallel combinations, MS-Diffusion also
 911 performs well in personalized generation for combinations with certain overlapping areas, such as
 912 "living+midwearing" and "object+scene".

913 F COMPARISON WITH TUNING-BASED METHODS

914

915

916 While zero-shot methods like MS-Diffusion can decrease the tuning cost, they may also suffer from
 917 performance degradation compared to tuning-based approaches due to the limitations of pre-training
 scale. However, MS-Diffusion indicates comparable performance in single-subject quantitative results



Figure 9: **Qualitative comparison of Break-A-Scene and MS-Diffusion.** MS-Diffusion gets comparable results by extracting subjects from a single image. Break-A-Scene tends to overfit the input image (the bird pose in the first row and the white creature sitting on the black bowl in the second row).

in Table 2. Since the proposed MS-Bench can be too large for tuning-based methods, we provide qualitative comparison results with Mix-of-Show (Gu et al., 2023), Cones2 (Liu et al., 2023b), and Break-A-Scene (Avrahami et al., 2023) in Figure 8 and Figure 9. The results show that MS-Diffusion achieves comparable results to tuning-based methods. Mix-of-Show and Cones2 face certain issues when handling multiple subjects with complex interactions (e.g., the example of a dog wearing a dress). Break-A-Scene tends to overfit the original image’s interactions (e.g., the pose of birds and the white creature sitting on the black bowl). While avoiding these issues, MS-Diffusion requires no test-time tuning and only one subject image during inference, unlike the tuning-based methods that need additional time and multiple images for tuning.

G LAYOUT GUIDANCE

Cross-attention maps can intuitively reflect the condition-image attribution relation (Tang et al., 2023; Hertz et al., 2023). Recent works (Wang et al., 2023; Wei et al., 2024) have studied utilizing an objective on the cross-attention maps in multi-subject generation. The objective exists only in training, considered implicit and insufficient to handle multi-subject conflicts. We have provided the performance comparison between our explicit layout guidance and attention loss in Section 4.4. For a single cross-attention layer, the attention loss of the j th subject is calculated by:

$$\mathcal{L}_{am}^j = \left(1 - \frac{\sum_{[x,y] \in B_j} \mathbf{A}_{[x,y],j}}{\sum_{[x,y]} \mathbf{A}_{[x,y],j}} \right)^2 \tag{7}$$

where $[x, y]$ corresponds to a latent token in \mathbf{Q} and B_j is the bounding box of the j th subject. This objective aims to promote the activation of attention maps within specific boxes. We average \mathcal{L}_{am}^j across layers and subjects and set its weight to 0.01 in the final loss. To validate the text attention loss, we also optimize the text cross-attention layers in training, increasing approximately 70% in learnable parameters.

Although our model provides explicit layout guidance, it still significantly differs from layout-based diffusion. Firstly, the information of boxes in the grounding resampler is prior, and its conditional effect is relatively weak. We have also reduced the model’s reliance on this input by randomly dropping grounding tokens. Secondly, the multi-subject cross-attention only exists in image cross-attention, inherently controlling the action of image conditions in specified areas, but cannot determine the whole generation of the diffusion model. Our goal is not to develop a method that fully supports layout control but to utilize layout information to guide the model in resolving conflicts in multi-subject generation.



Figure 10: **Text-image attribution analysis of MS-Diffusion.** We average the attention maps corresponding to the subjects and translate them to normalized heat maps.

To further explore the layout control capability of MS-Diffusion, we provide qualitative results in Figure 15. It can be demonstrated that MS-Diffusion can generate images that adhere to layout conditions, even in the case of two instances of the same category. However, the generated positions are not entirely accurate, especially in *"a cat and a cat on the grass"*, illustrating that the layout condition is relatively weak compared to text and image prompts in the personalization task.

1000 H INTEGRATION WITH CONTROLNET

1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

In the realm of text-to-image diffusion models, a notable application is the enhancement of structural control within image generation. Our MS-Diffusion maintains the original network architecture unchanged, thus ensuring full compatibility with existing controllable tools. Consequently, this allows for the generation of images that are not only prompted by images but also governed by additional conditions. By integrating our MS-Diffusion with established controllable mechanisms such as ControlNet, we demonstrate the capacity to produce images under varied structural directives. Figure 16 displays an array of images synthesized with image prompts coupled with distinct structural controls (depth, canny edge, openpose). This seamless cooperation between our method and these tools facilitates the creation of highly controllable images without necessitating fine-tuning.

I MULTI-SUBJECT INTERACTION

Unlike image editing (Kawar et al., 2023; Brooks et al., 2023), personalized image generation features a high degree of freedom, enabling effective handling of interactions among different elements. Benefiting from its architecture design that does not impact the base model, MS-Diffusion successfully inherits the multi-subject interaction capabilities of the base model. As illustrated in Figure 17, MS-Diffusion can flexibly manage interactions between reference subjects, even when there is overlap among these objects, thereby demonstrating its potential in practical applications.

J HUMAN AND ANIME PERSONALIZATION

Human and anime subjects are popular in the use of personalized model. We provide results of MS-Diffusion on human and anime subjects in Figure 18. Some research (Wang et al., 2024a; Li et al., 2024) has explored human personalization in text-to-image diffusion models. By utilizing a



Figure 11: **Qualitative examples when applying pseudo layout guidance during the inference.** In this figure, *pse-lg* and *pse* respectively indicate whether layout prior is used.

face encoder (Deng et al., 2019) and training on the face dataset, MS-Diffusion can also be extended to a personalized model for these subjects.

K TEXT-IMAGE ATTRIBUTION ANALYSIS

In MS-Diffusion, our focus is primarily on resolving conflicts between subjects without altering the control mechanism of the text. While some approaches (Wang et al., 2023; 2024b; Zhou et al., 2024; Kim et al., 2023) in non-personalized text-to-image tasks address multi-object generation conflicts through text cross-attention, this inevitably requires tuning the diffusion model’s parameters, thereby affecting the plug-and-play nature, which is not preferred by us. As demonstrated in the text-image attribution analysis presented in Figure 10, the control of multiple objects by text in our model is also quite evident. This may be related to the explicit layout guidance for subjects, since the images and text condition jointly in the generation process. We also attempted to control text cross-attention using the same mechanism as in Section 3.5, but no differences were observed in the results.

L SUBJECT INTERPOLATION

The blending of two distinct subject representations to yield composite subjects with hybrid characteristics is feasible through the navigation of the embedding space linking the subjects. As depicted in Figure 19, linear interpolations are conducted among dog and hat representations, subsequently rendering the interpolated subject in an unaccustomed context. The visualization reveals a natural gradation of subject appearance along the interpolated trajectory that harmonizes with the surrounding environment. This technique proves beneficial when applied in subject fusion and style transfer.

M LIMITATIONS

There are certain limitations in MS-Diffusion. The box-based indication of positions lacks precision, making it challenging to work effectively when the interaction between subjects is stronger. Moreover, the model requires explicit layout input during inference, and generating complex scenes becomes difficult. Though MS-Diffusion beats SOTA personalized diffusion methods in both single-subject and multi-subject generation, it still suffers from the influence of background in subject images.

1080 We explore a solution for explicit layout needs during inference. MS-Diffusion supports using the text
1081 cross-attention maps as the pseudo layout guidance. Specifically, as indicated in Figure 10, since the
1082 text cross-attention maps can reflect the area of each text token, we can replace the layout prior with
1083 them during the inference. In practice, we set a threshold to extract masks from text cross-attention
1084 maps and apply them after T denoising steps. Before T , we experiment with completely disabling
1085 layout guidance or using a rough box as a layout prior. The results are presented in Figure 11.
1086 Although disabling layout guidance experiences a decline in subject consistency, it still demonstrates
1087 that explicit layout guidance during inference can be optimized. One direction for exploration is to
1088 enable the model to learn the layout during training.

1089 1090 N SOCIETAL IMPACT

1091
1092 As an image personalization method, MS-Diffusion aims to customize images based on user-provided
1093 subjects without fine-tuning. Additionally, the multi-subject reference capability of MS-Diffusion
1094 allows users to freely combine and re-create different concepts. However, MS-Diffusion can also be
1095 used to generate deceptive images, especially those involving subject combinations that would not
1096 exist in reality, an issue that remains to be addressed in the future.

1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

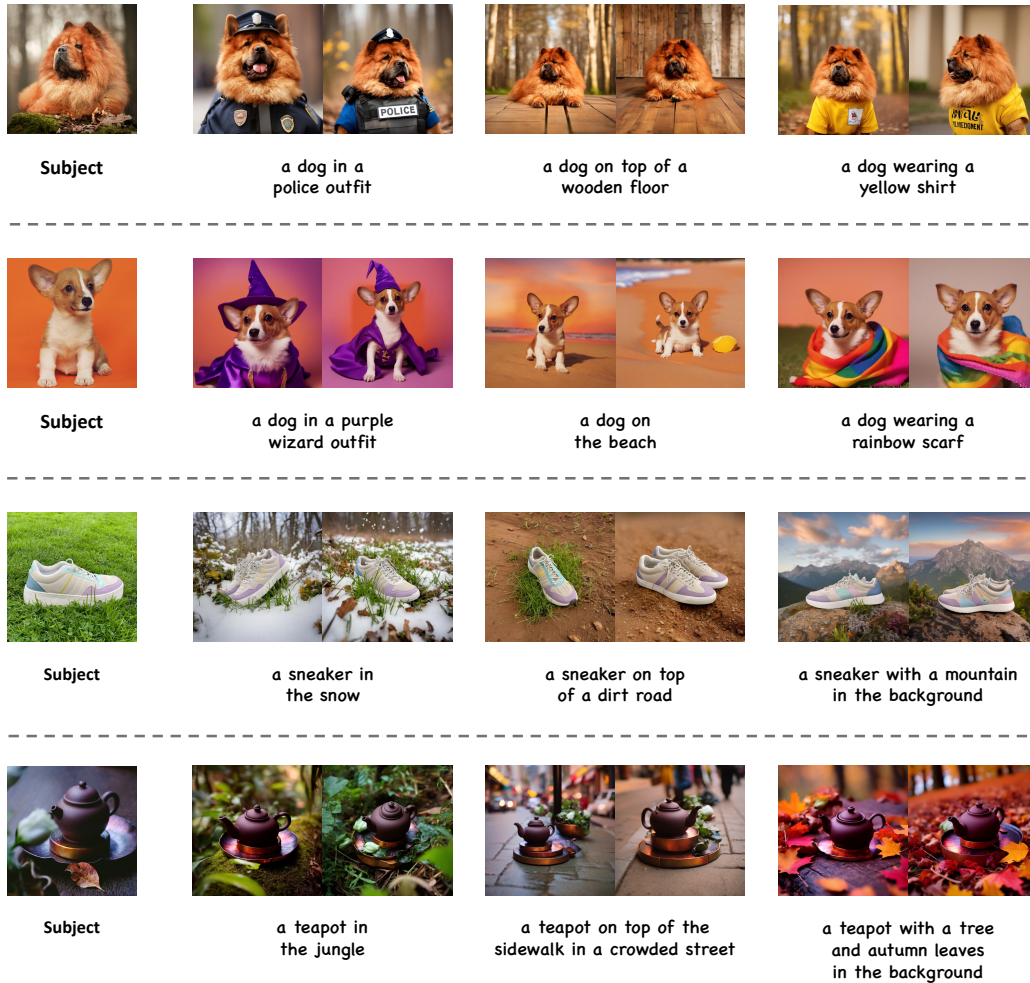


Figure 12: Additional qualitative results of MS-Diffusion in single-subject personalization.

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241



Figure 13: **Qualitative comparison of MS-Diffusion and zero-shot personalized SOTAs on un-common subjects.** Though losing some details, MS-Diffusion outperforms other SOTAs obviously.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

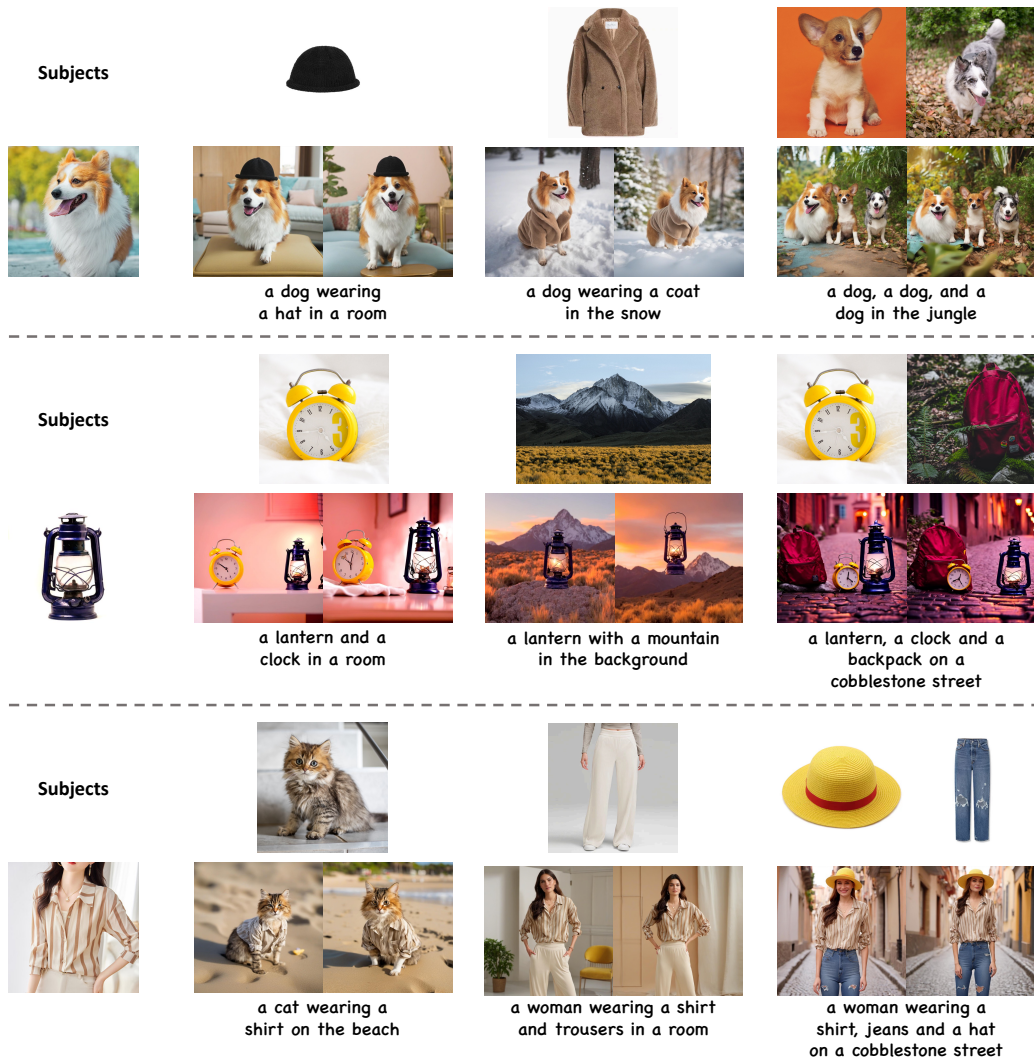


Figure 14: Additional qualitative results of MS-Diffusion in multi-subject personalization.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

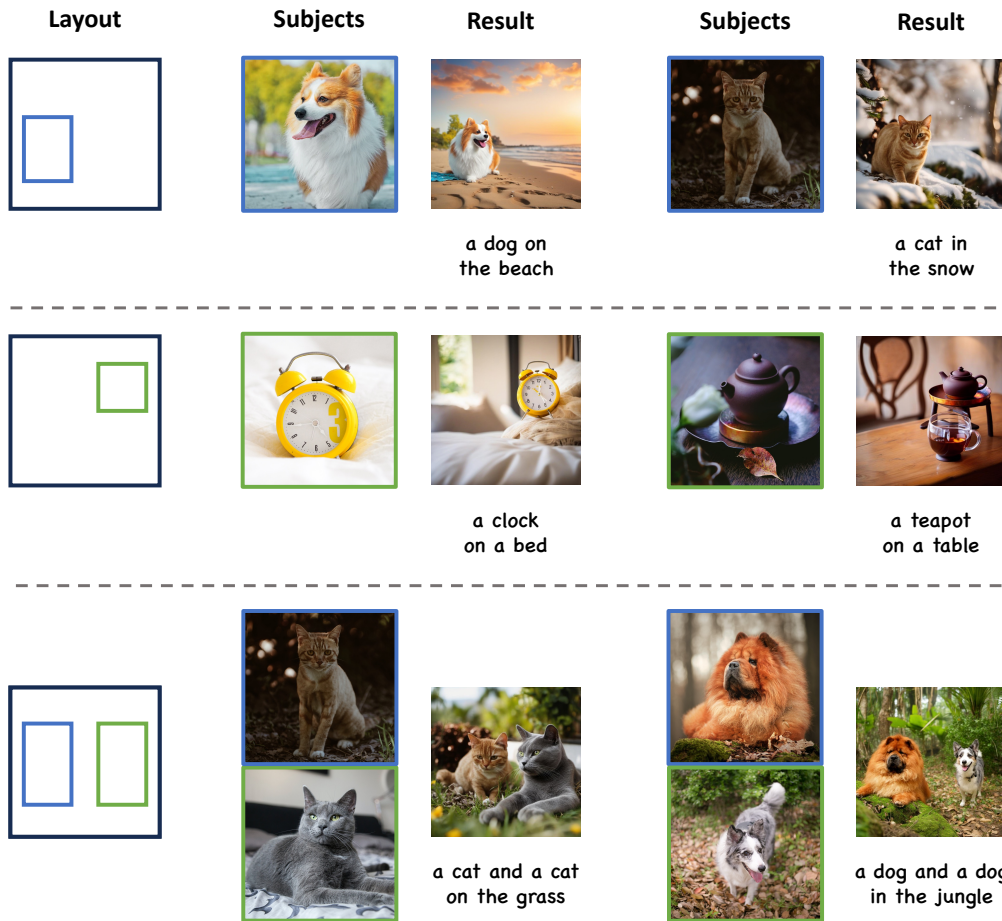


Figure 15: **Qualitative examples of MS-Diffusion about the layout control ability.** Bounding boxes of different colors correspond to subjects with different color borders.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

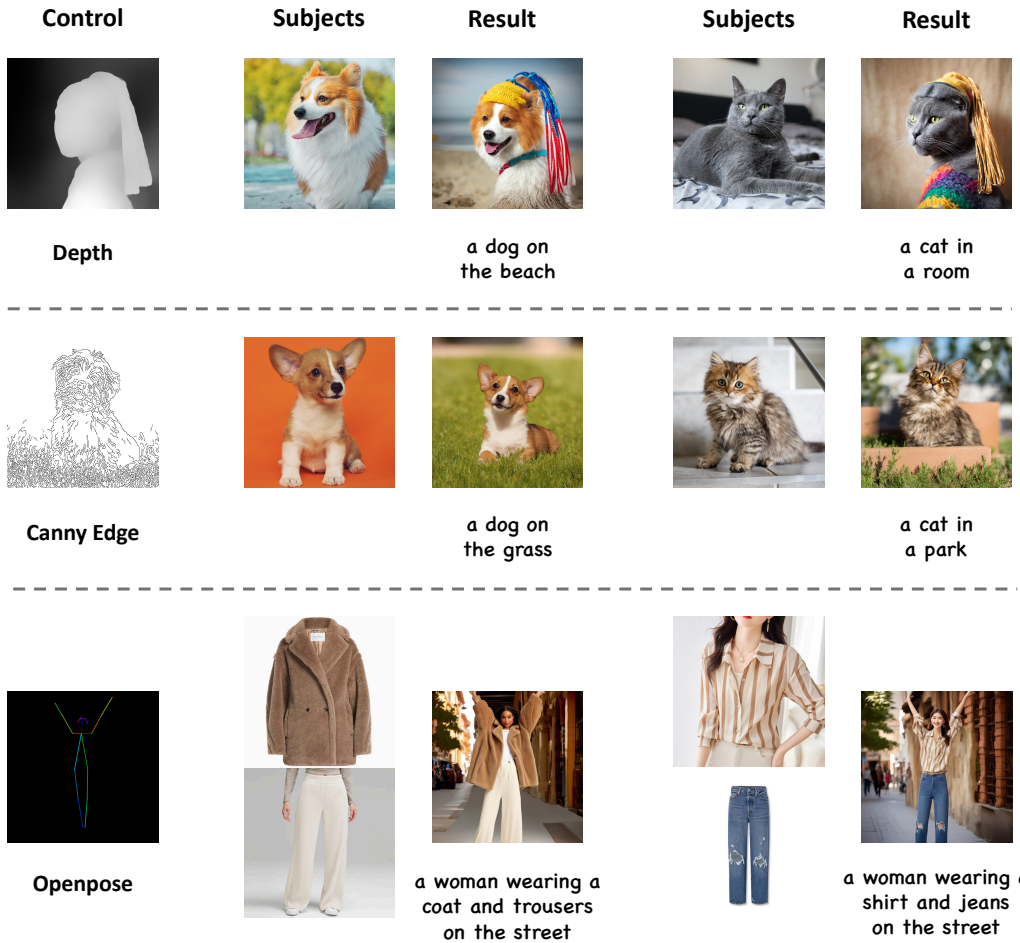


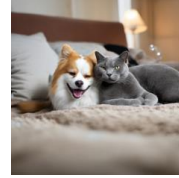
Figure 16: **Generative results when integrating different control conditions.** The integrated ControlNets are composed of depth, canny edge, and openpose.

1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457

Subjects



Interaction Subjects Generation



a dog and a cat playing on the beach

a dog and a cat lying together



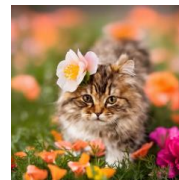
a dog wearing a t-shirt

a dog painted on a t-shirt



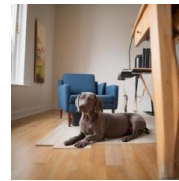
a woman holding a cat in the snow

a woman lying on the grass with a cat on her



a cat holding a flower

a cat with a flower on its head



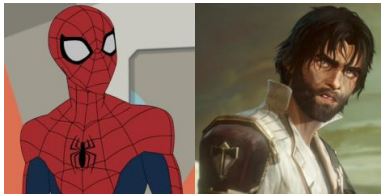
a dog on a chair in a room

a dog under a chair in a room

Figure 17: **Examples of prompts with complex interaction of multiple subjects.** MS-Diffusion can generate high-quality images following both the subjects and prompts.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

Subjects



Results



a man walking on the road



an anime man on a crowded street



an anime man and an anime man fighting together



a man talking with an anime man wearing red hat



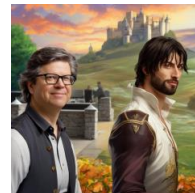
a man standing in the church



an anime man flying between buildings



an anime man sitting besides an anime man



a man and an anime man with a castle in the background

Figure 18: **Personalized results of MS-Diffusion on human and anime subjects.** MS-Diffusion can generate high-quality images in both single-subject and multi-subject personalization for humans and anime.

1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565

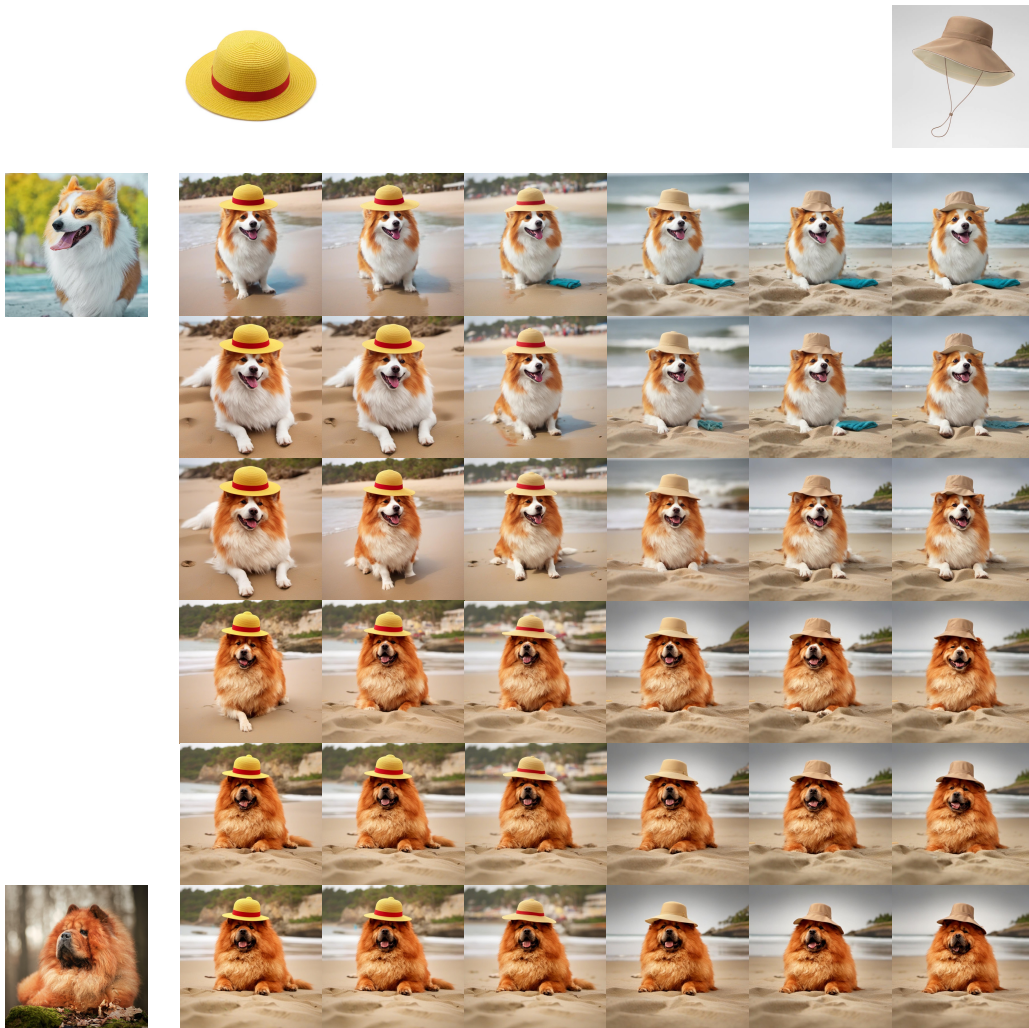


Figure 19: **Subjects interpolation in multi-subject generation.** We select two dogs and two hats to conduct linear interpolation with the text set to "a dog wearing a hat on the beach".