# Harnessing the Spatial-Temporal Attention of Diffusion Models for High-Fidelity Text-to-Image Synthesis

**Qiucheng Wu**[1*], **Yujian Liu**[1*], **Handong Zhao**[2], **Trung Bui**[2], **Zhe Lin**[2], **Yang Zhang**[3], **Shiyu Chang**[1]

[1]UC Santa Barbara, [2]Adobe Research, [3] MIT-IBM Watson AI Lab
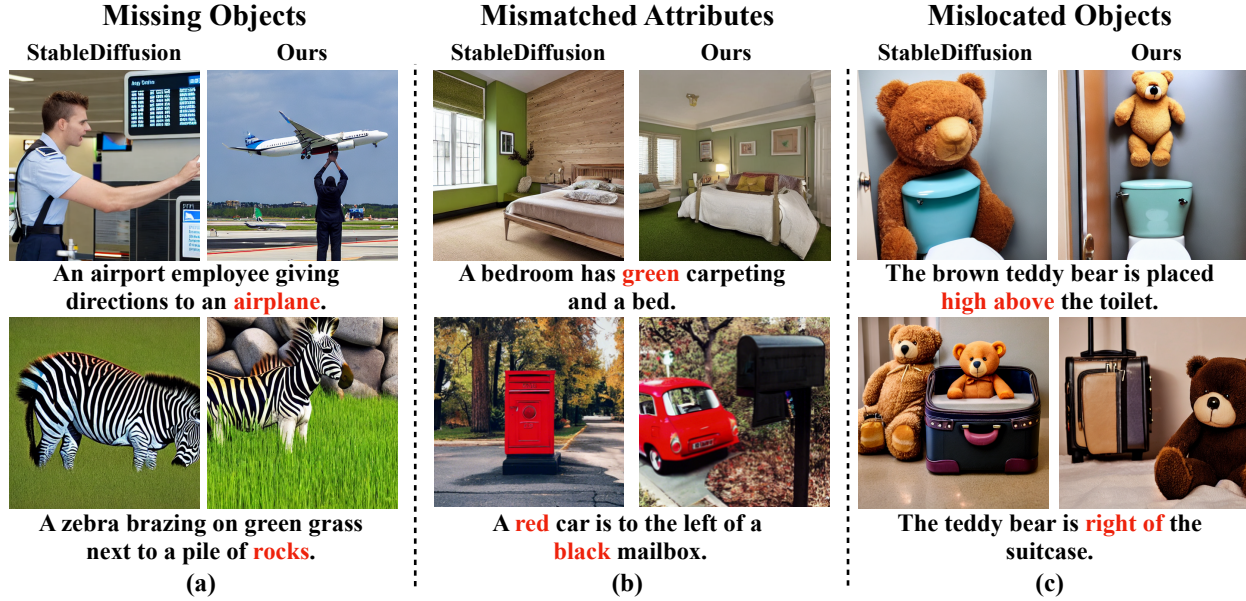
{qiucheng, yujianliu}@ucsb.edu

Figure 1. **Example generation by stable-diffusion-v1-4 and our method.** Stable diffusion model makes three types of errors that include *missing objects*, *mismatched attributes*, and *mislocated objects*. Errors are highlighted in red.

## Abstract

*Diffusion-based models have achieved state-of-the-art performance on text-to-image synthesis tasks. However, one critical limitation of these models is the low fidelity of generated images with respect to the text description, such as missing objects, mismatched attributes, and mislocated objects. One key reason for such inconsistencies is the inaccurate cross-attention to text in both the spatial dimension, which controls at what pixel region an object should appear, and the temporal dimension, which controls how different levels of details are added through the denoising steps. In this paper, we propose a new text-to-image algorithm that adds explicit control over spatial-temporal cross-attention in diffusion models. We first utilize a layout predictor to predict the pixel regions for objects mentioned in the text. We then impose spatial attention control by combining the attention over the entire text description and that over the local description of the particular object in the correspond-ing pixel region of that object. The temporal attention control is further added by allowing the combination weights to change at each denoising step, and the combination weights are optimized to ensure high fidelity between the image and the text. Experiments show that our method generates images with higher fidelity compared to diffusion-model-based baselines without fine-tuning the diffusion model. Our code is publicly available.[1]*

## 1. Introduction

Diffusion models [14, 19, 45, 46, 48] have recently revolutionized the field of image synthesis. Compared with previous generative models such as generative adversarial networks [1, 3, 11, 18] and variational autoencoders [21, 37, 38], diffusion models have demonstrated superior

---

[*]Equal contribution.
[1]https://github.com/UCSB-NLP-Chang/Diffusion-SpaceTime-Attn

performance in generating images with higher quality, more diversity, and better control over generated contents. Particularly, text-to-image diffusion models [2, 32, 35, 39, 41] allow generating images conditioned on a text description, which enables generation of creative images due to the expressiveness of natural language.

However, recent studies [10, 30] have revealed that one critical limitation of existing diffusion-model-based text-to-image algorithms is the *low fidelity* with respect to the text descriptions – the content of the generated image is sometimes at odds with the text description, especially when the description is complex. Specifically, typical errors made by stable diffusion models fall into three categories: *missing objects*, *mismatched attributes*, and *mislocated objects*. For example, in Fig. 1(a), stable diffusion model ignores the airplane even though it is mentioned in text; in Fig. 1(b), the model confuses "red car" and "black mailbox" and generates a red mailbox; in Fig. 1(c), the model locates teddy bear behind the toilet, despite the description "teddy bear is placed high above the toilet."

Such infidelity problems suggest that the cross-attention map on the text description may not be accurate. In particular, if we view the generation process of a diffusion model as a sequence of denoising steps, then the cross-attention on text descriptions can be considered as a function of both *spatial* (pixels) and *temporal* (denoising steps) information. Therefore, the inaccuracies of the cross-attention can result from the loose control over both the spatial and temporal dimensions. On one hand, spatial attention controls at what pixels the model should attend to each object and the corresponding attributes mentioned in the text. If the spatial attention is incorrect, the resulting images will have incorrect object locations or miss-associated attributes. On the other hand, temporal attention controls when the models should attend to different levels of details in the text. As previous works have revealed, diffusion models tend to focus on generating object outlines at earlier denoising steps and on details at later [50]. Thus loose control over the temporal aspect of attention can easily lead to overlooking certain levels of the object details. In short, to improve the fidelity of text-to-image synthesis, one would need to explicitly control both spatial and temporal attention to follow an accurate and optimal distribution.

In this paper, we propose a new text-to-image algorithm based on a pre-trained conditional diffusion model with explicit control over the spatial-temporal cross-attention map on text. The proposed algorithm introduces a layout predictor and a spatial-temporal attention optimizer. The layout predictor takes the text description as input and generates a spatial layout for each object mentioned in the text. Alternatively, the layout can also be provided by the user. Then the spatial-temporal attention optimizer imposes direct control over the spatial and temporal aspects of the attention

according to the spatial layout. In particular, for the spatial aspect, we parameterize the attention map such that the attention outputs in the designated pixel region for an object are a weighted combination of attention over the entire text description and that over the local description that specifically describes the corresponding object. In this way, we manage to emphasize the attention over the object descriptions. For the temporal aspect, we allow the combination weights to change across time and optimized according to a CLIP objective that measures the agreement between the generated images and the text description. In this way, we allow the attention to focus more on the entire description at the early stage and gradually shift to the detailed local descriptions as the denoising process proceeds. The entire pipeline resembles a typical painting process of a human, where each object's position is determined beforehand and the focus gradually shifts from global information to the local details of each object.

We conduct extensive experiments on datasets that contain real and template-based captions [6, 26] and our newly created synthetic dataset that contains complex text descriptions. Results show that our method generates images that better align with descriptions compared to other stable diffusion-based baselines. As shown in Fig. 1, our method effectively resolves the above-mentioned three errors. Particularly, controlling spatial attention locates objects at the desired position, and controlling temporal attention promotes the occurrence of objects with associated attributes. Our findings shed light on fine-grained control of diffusion models in text-to-image generation tasks.

## 2. Related Work

**Diffusion Models** Diffusion models are a class of generative models that have demonstrated state-of-the-art performance on image synthesis tasks [8, 14, 19, 45, 46, 48]. These models synthesize images by sampling a noisy image from the standard Gaussian distribution and iteratively denoising it back to a clean image. Their impressive performance has advanced research in multiple computer vision areas including inpainting [29, 40, 51], image editing [7, 12, 31, 50], super-resolution [15, 42], video synthesis [13, 16, 55], and applications beyond computer vision [22, 23, 52]. Among these, text-to-image diffusion models have gained significant attention [2, 35, 36, 41, 58]. Taking text descriptions as inputs, these models generate high quality images that are semantically aligned with the input text, which have led to many creative and artistic applications.

**Enhancing the Controllability of Text-to-Image Diffusion Models** While diffusion-based text-to-image models have shown promising results, recent works have highlighted cases where models fail to generate high-fidelity images with respect to the input text [10, 30]. To this end,
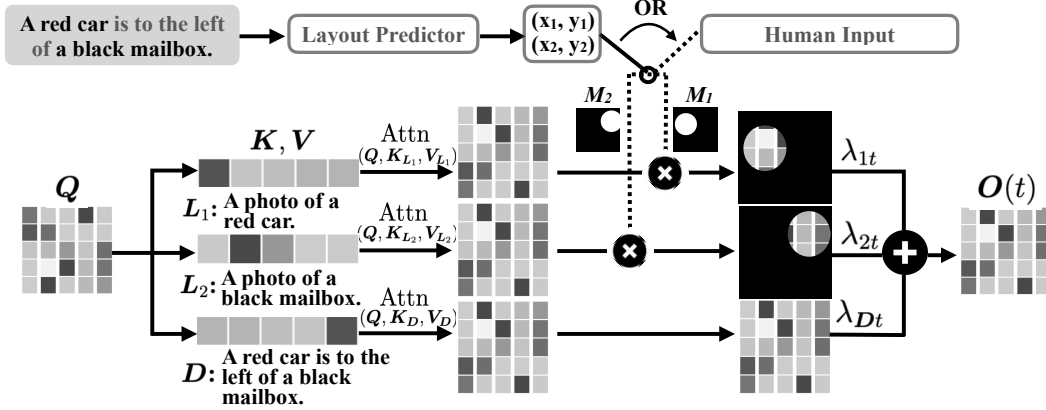
Figure 2. **Overview of our text-to-image generation pipeline at one denoising step.** Given input text $D$, we first parse it and extract all objects mentioned, constructing local descriptions $L_i$. Next, the layout predictor predicts the pixel region for each object in the text. The diffusion model attends to the global description $D$ and additionally attends to the local description $L_i$ in object $i$'s region. The final attention output is a weighted combination of attention to both global and local descriptions, where the combination weights sum up to 1 for each pixel and are optimized for each denoising step to achieve a high fidelity with $D$.

prior works have explored various ways to enhance the controllability of text-to-image diffusion models. One line of work enhances the controllability by improving diffusion models' ability to understand natural language, which includes using a more powerful text encoder that is separately trained on language modeling tasks [2, 41], incorporating linguistic structures in the text to guide the cross-attention between image and text [10], and decomposing a complex text description into multiple components that are easy to generate [27]. Another line of work conditions diffusion models' generation on auxiliary inputs such as object layout [2, 9, 24, 54] and silhouette [17, 33, 44, 56, 57]. By modifying diffusion models' attention operation according to these auxiliary information or directly fine-tuning diffusion models to take these auxiliary inputs, they are able to control the location and shape of the objects in the image. Finally, some work adds temporal aspect control on diffusion models by modifying the input text condition at each denoising step, which allows them to disentangle a desired attribute from other contents [50]. Different from prior works, our method imposes both spatial and temporal control in cross-attention layer. Moreover, our method does not require auxiliary inputs and does not fine-tune the diffusion model.

# 3. Methodology

## 3.1. Problem Formulation

We focus on the standard text-to-image problem. Given a text description, denoted as $D$, our goal is to generate an image that is consistent with $D$. For a concrete exposition, we will use an example description "*A red car is to the left of a black mailbox*" in the following. Our work aims to improve the fidelity of the generated image to text, which includes three requirements:

- **Object Fidelity:** The generated images should include all the objects mentioned in $D$. In our example, the generated image should contain a car and a mailbox.

- **Attribute Fidelity:** The attributes of each object in the image should match those in $D$. In our example, the car should be red and the mailbox should be black.

- **Spatial Fidelity:** The relative spatial positions of the object should match the description in $D$. In our example, the car should be on the left and the mailbox on the right.

We will tackle these problems using a pre-trained, fixed stable diffusion model.

## 3.2. Method Overview

To achieve high-fidelity text-to-image generation, we propose an algorithm consisting of the following four steps.
**Step 1: Object Identification**     We extract all the objects mentioned in $D$, denoted as $O_{1:N}$, by eliciting the noun phrases using spaCy.[2] $N$ is the total number of objects. In our sample, $O_1 = $"*a red car*" and $O_2 = $"*a black mailbox*".
**Step 2: Layout Prediction**     For each object $O_i$, we use a *layout predictor* to predict its pixel region, $\mathcal{R}_i$, which is a set of pixels roughly specifying where the object should be. The layout can also be provided by human users.
**Step 3: Local Description Generation**     For each object $O_i$, we generate a local text description, $L_i$, containing only that object information using a simple template. In our example, there are two local descriptions. $L_1 = $"*A photo of a red car*" and $L_2 = $"*A photo of a black mailbox*".
**Step 4: Attention Optimization**     During the generation process, we guide the diffusion model to combine attention to both the global description $D$ and the local descriptions

---

[2]https://spacy.io/.

$L_{1:N}$ according to the object layout. The attention combination weights are optimized for each input description.

The following subsections will provide further details about steps 2 and 4.

### 3.3. Layout Predictor

Our layout predictor is adapted from [53], which aims to predict the center coordinate $C_i = [X_i, Y_i]$ for each object $O_i$. Specifically, the layout predictor is a transformer that takes the text description $D$ as the input. At the output position where the input mentions the object $O_i$, we let the transformer output a set of Gaussian mixture model (GMM) parameters to fit the object's center coordinate $C_i$. Formally, denote $\boldsymbol{f}_i(\boldsymbol{D}; \boldsymbol{\theta})$ as the output of the layout predictor at the location of the mentioning object $O_i$. Then

$$\boldsymbol{f}_i(\boldsymbol{D}; \boldsymbol{\theta}) = \bigcup_{k=1}^{K} \{\boldsymbol{\mu}_{ik}, \boldsymbol{\Sigma}_{ik}, w_{ik}\}, \quad (1)$$

where $\boldsymbol{\mu}_{ik}$, $\boldsymbol{\Sigma}_{ik}$, and $w_{ik}$ denote the mean, covariance matrix and the prior probability of mixture $k$; $\boldsymbol{\theta}$ denotes the network parameters.

To improve the fidelity of the object positions in the predicted layout, we introduce a hybrid training objective including an *absolute position objective* and a *relative position objective*.

**Absolute Position Objective**   The absolute position objective provides direct supervision of the exact position of each object. Formally, we assume access to an image captioning dataset $\mathcal{D}_{\text{real}}$ with description $\boldsymbol{D}$ as well as the extracted labels of all the objects $\{O_i\}$ and their center coordinates $\{C_i\}$. Then the training objective is to minimize the negative log-likelihood of the ground-truth coordinates under the predicted GMM distribution, *i.e.*

$$\mathcal{L}_{\text{abs}}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{D}, \{C_i\} \sim \mathcal{D}_{\text{real}}} \Big[ \sum_{i=1}^{N} \ell_{\text{nll}}(\boldsymbol{C}_i; \boldsymbol{f}_i(\boldsymbol{D}; \boldsymbol{\theta})) \Big], \quad (2)$$

where $\ell_{\text{nll}}(\boldsymbol{C}_i; \boldsymbol{f}_i)$ denotes the negative log-likelihood of ground truth $\boldsymbol{C}_i$ under the GMM specified by $\boldsymbol{f}_i$.

**Relative Position Objective**   In many cases, the text description only mentions the relative positions of the objects, and thus it is more important to ensure the relative position of the predicted position is correct than the absolute position. To further enforce the fidelity of the relative positions, we introduce the following relative position objective.

To start with, we construct a synthetic dataset, $\mathcal{D}_{\text{syn}}$, which consists of text description $\boldsymbol{D}$ with explicit descriptions of relative spatial relations. We first randomly select $N$ objects with attribute modifiers. We then select $M$ pairs of objects to specify their relative spatial relation. For object pair $(O_i, O_j)$, their spatial relation, $R_{ij}$, is randomly drawn from "left of", "right of", "above" and "below". Finally,

we prompt GPT-3 [4] to generate the text description $\boldsymbol{D}$ that mentions all the objects and their relative positions. Our familiar "A red car is to the left of a black mailbox" is one such example. Further details are provided in Appendix B.

Next, we introduce a loss to penalize violating the relative position. If object $i$ is to the left of object $j$, we enforce that the rightmost mixture mean of object $i$ is to the left of the leftmost mixture mean of object $j$. Formally,

$$\ell_{\text{rel}}(R_{ij} = \text{"left"}; \boldsymbol{f}_i, \boldsymbol{f}_j) = \max\{\max_k \boldsymbol{\mu}_{ik}(0) - \min_k \boldsymbol{\mu}_{jk}(0), -\delta\}, \quad (3)$$

where $\boldsymbol{\mu}_{ik}(0)$ denotes the zeroth element (the $x$-coordinate) of $\boldsymbol{\mu}_{ik}$, and $\delta$ is a pre-specified margin. The $\ell_{\text{rel}}$ of the other three types of relations are defined similarly. The relative position loss is thus the aggregation of $\ell_{\text{rel}}$ across the synthetic dataset:

$$\mathcal{L}_{\text{rel}}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{D}, \{R_{ij}\} \sim \mathcal{D}_{\text{syn}}} \Big[ \sum_{i,j : R_{ij} \neq \emptyset} \ell_{\text{rel}}(R_{ij}; \boldsymbol{f}_i, \boldsymbol{f}_j) \Big], \quad (4)$$

where $\boldsymbol{f}_i$ is short for $\boldsymbol{f}_i(\boldsymbol{D}; \boldsymbol{\theta})$.

**Training and Inference**   To sum up, the final training objective is the combination of both:

$$\mathcal{L}_{\text{layout}}(\boldsymbol{\theta}) = \mathcal{L}_{\text{abs}}(\boldsymbol{\theta}) + \xi \mathcal{L}_{\text{rel}}(\boldsymbol{\theta}), \quad (5)$$

where $\xi$ is a hyperparameter. During inference time, we randomly draw the center coordinate $\boldsymbol{C}_i$ from the predicted GMM. The pixel region for object $O_i$, $\mathcal{R}_i$, is defined as a circular region centered at the drawn $\boldsymbol{C}_i$ with a *fixed radius* $r$. As would be shown in Appendix D, $\mathcal{R}_i$ only roughly regulates the position of the generated objects, and the actual size of the object can go beyond or below the size of $\mathcal{R}_i$. Thus a fixed $r$ is sufficient for this purpose.

### 3.4. Spatial-Temporal Attention Optimization

Recall that $\boldsymbol{D}$ is the global description and $\{L_i\}$ are local descriptions for each object. Our goal is to guide the diffusion model to attend to not only $\boldsymbol{D}$, but also $\boldsymbol{L}_i$ in object $i$'s region, $\mathcal{R}_i$, so that the model is more strongly prompted to generate the specific object as specified in the layout.

**Spatial-Temporal Attention**   Recall that the standard cross-attention [49] is defined as

$$\text{Attention}(\boldsymbol{Q}, \boldsymbol{K}, \boldsymbol{V}) = \text{softmax}\big(\boldsymbol{Q}\boldsymbol{K}^T / \sqrt{d}\big)\boldsymbol{V}, \quad (6)$$

where $\boldsymbol{Q} \in \mathbb{R}^{h \times w, d}$ is the queries vectors for each *pixel*, and $\boldsymbol{K}, \boldsymbol{V} \in \mathbb{R}^{l, d}$ are key and value vectors for each *text description token*; $h$ and $w$ represent the image height and width, $l$ represents the text description length; $d$ represents the dimension of each attention vector.

Since we have multiple text descriptions, we define $\boldsymbol{K}_D, \boldsymbol{V}_D$ as the key and value vectors for the global description $\boldsymbol{D}$, and $\boldsymbol{K}_{L_i}, \boldsymbol{V}_{L_i}$ for the local description $\boldsymbol{L}_i$. Further,

we introduce a set of binary mask matrices $\{M_i\}$ to indicate the region for each object, *i.e.*

$$M_i(x, y) = 1, \text{ if } (x, y) \in \mathcal{R}_i, \quad \text{and } 0, \text{ otherwise.} \quad (7)$$

Then the output of the attention layer of the denoising network at time $t$ is defined as

$$O(t) = \sum_{i=1}^{N} \lambda_{it} M_i \odot \text{Attention}(Q, K_{L_i}, V_{L_i})$$
$$+ \left(1 - \sum_{i=1}^{N} \lambda_{it} M_i\right) \odot \text{Attention}(Q, K_D, V_D), \quad (8)$$

where $\lambda_{it}$ are the attention combination weights, and $\odot$ denotes element-wise multiplication. Note that the combination weights are functions of time $t$, which is motivated by the observation that different denoising steps control the generation of different levels of details. By introducing this time dependency, we allow the diffusion model to focus more on the global description at earlier time steps and shift to local descriptions later.

**Optimization Objective**   All the attention combination weights, $\lambda = \{\lambda_{it}\}$, are determined by maximizing the consistency between the generated image and the text description as measured by the CLIP similarity [34]. We introduce two CLIP similarities, a global CLIP similarity that compares the entire image and the global description, and a set of local CLIP similarities that compare the images at each object region and the corresponding local descriptions. Formally, denote the generated image as $I(\lambda)$, and its local patch at each object's region as $I_{O_i}(\lambda)$.[3] Note that we adopt the *deterministic* PLMS denoising process [47] (*i.e.*, with $\sigma = 0$), so both $I$ and $I_{O_i}$ are deterministic functions of $\lambda$. Then the loss for optimizing $\lambda$ is given by

$$\mathcal{L}_{\text{attend}}(\lambda) = -\text{CLIP}(I(\lambda), D) - \gamma \sum_{i=1}^{N} \text{CLIP}(I_{O_i}(\lambda), L_i), \quad (9)$$

where $\text{CLIP}(\cdot)$ denotes the CLIP similarity, and $\gamma$ is a hyperparameter.

# 4. Experiments

We conduct experiments to evaluate our method's performance and generalizability. We also perform ablation study on important design choices of our method.

**Implementing Details:** We adopt RoBERTa-base [28] as the base model for layout predictor and use 5 mixtures in GMM. We use stable-diffusion-v1-4 [39] pre-trained on the laion dataset [43] and freeze it throughout all experiments. All generated images are in the size of $512 \times 512$. We use

---

[3]To standardize the image patch sizes, $I_{O_i}$ is obtained by cropping the original image with a minimum square that encompasses $\mathcal{R}_i$ and resizing it to 224×224.

PLMS sampler [47] to synthesize images with 50 denoising steps, and we use Adam [20] to optimize the layout predictor and attention combination weights. More details on hyperparameters and optimization are in Appendix A.

## 4.1. Evaluation on Fidelity of Generated Images

We first evaluate our method on object, attribute, and spatial fidelities as introduced in Sec. 3.1 using both objective and subjective metrics.

**Baselines:** We identify four baseline methods on generating images from complex text descriptions. (1) VANILLA-SD [39] is the pre-trained text-to-image stable diffusion model that directly generates images conditioned on the text description. (2) COMPOSABLE-DIFFUSION [27] is a diffusion-based compositional generation method. To generate images from a text description, the text is first decomposed into the conjunction of multiple components (*e.g.*, for the example in Sec. 3.1, the components are "A red car is to the left of a black mailbox." AND "A red car" AND "a black mailbox"). Each component is separately modeled by the diffusion model and then composed to generate an image by merging the outputs of each denoising step. (3) STRUCTURE-DIFFUSION [10] improves VANILLA-SD on generating images with better object and attribute fidelities. They do so by first extracting noun phrases at different levels of the parsing tree. The model then separately attends to each noun phrase and combines the attention outputs by taking their arithmetic mean. Different from our method, they perform cross-attention on the whole image instead of the region specific to an object. (4) PAINT-WITH-WORDS [2] assumes that users provide the pixel region of each object to be generated. Pixels in the region then increase their attention weights to the text that describes the corresponding object, and the amount of increase is determined by heuristic rules. For a fair comparison, we use the pixel region predicted by our layout predictor, and we report the performance with the ground truth region in Appendix E.

**Datasets and Metrics:** We conduct experiments on three datasets. (1) **MS-COCO** [25] contains photos taken by photographers and manually annotates the caption for each photo. We use the caption as the input text description. (2) **VSR** [26] is proposed for probing spatial understanding of vision-language models. Constructed from a subset of MS-COCO, it uses templates to generate captions that describe spatial relations in the image (*e.g.*, "The horse is to the left of the person."). (3) **GPT-synthetic** is the synthetic dataset introduced in Sec. 3.3. We manually check the test set to filter out sentences that do not conform to the specified spatial relation. Compared with VSR, this dataset has more diverse and complex text descriptions. The final dataset contains 500 descriptions, and we downsample MS-COCO and VSR to have the same size. Statistics including the number of objects and spatial relations are shown in Appendix B.

| | MS-COCO | | | VSR | | | GPT-synthetic | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Object (100) ↑ | Attribute (100) ↑ | Overall | Object (100) ↑ | Spatial (100) ↑ | Overall | Object (100) ↑ | Attribute (100) ↑ | Spatial (100) ↑ | Overall |
| Vanilla-SD | 66 | 62 | 42% | 64 | 66 | 48% | 60 | 58 | 60 | 38% |
| Composable-Diffusion | 53 | 48 | 30% | 63 | 60 | 18% | 63 | 60 | 47 | 28% |
| Structure-Diffusion | 60 | 63 | 48% | 54 | 64 | 32% | 58 | 54 | 58 | 32% |
| Paint-with-Words | **75** | **74** | 52% | 54 | 59 | 30% | 65 | 52 | 57 | 32% |
| Ours | **75** | **74** | – | **68** | **77** | – | **84** | **83** | **86** | – |

Table 1. **Subjective evaluation of our method and baselines. Best numbers** are in bold. Spatial relation is not available on MS-COCO because very few of its captions contain spatial relations. Attribute is not available on VSR as its captions do not consider attribute. Object, Attribute, and Spatial show the total score of 50 evaluations, where a model with the highest fidelity would achieve a score of 100. Overall denotes the percentage of generation of each method that is rated better than our method.

| | MS-COCO | | VSR | | GPT-synthetic | |
|---|---|---|---|---|---|---|
| | Object Recall | SPRel Precision | Object Recall | SPRel Precision | Object Recall | SPRel Precision |
| Vanilla-SD | 58.0% | – | 62.1% | 56.0% | 42.4% | 56.8% |
| Composable-Diffusion | 51.8% | – | 60.8% | 72.2% | 34.4% | 52.6% |
| Structure-Diffusion | 61.7% | – | 62.3% | 58.3% | 43.0% | 59.3% |
| Paint-with-Words | 57.3% | – | 63.8% | 47.1% | 45.9% | 53.5% |
| Ours | **69.6%** | – | **65.1%** | **75.0%** | **47.2%** | **66.7%** |

Table 2. **Automatic evaluation of our method and baselines.** SPRel Precision: Spatial Relation Precision. **Best numbers** are in bold. SPRel precision is not available on MS-COCO since most captions do not have an explicit spatial relation.

For automatic evaluations, we consider two metrics. (1) **Object Recall** measures the percentage of successfully synthesized objects over objects mentioned in the text. We use DETR [5] to detect objects in generated images. To calculate recall, we divide the number of detected objects in the text by the total number of objects in the text that also belong to one of the MS-COCO categories. It measures the **object fidelity** of generated images. (2) **Spatial Relation Precision** (SPRel Precision) measures the percentage of the correct spatial relations among all the relations whose corresponding objects are successfully synthesized. This metric measures the **spatial fidelity** of generated images. We consider the relation of left, right, above, and below because their correctness can be evaluated by comparing the bounding box centers. Qualitative examples of spatial relations beyond these four are shown in Appendix G. We also report the CLIP similarities between generated images and the input description in Appendix C. We observe that different methods have very close CLIP similarities, though our method still achieves competitive results.

**Results:** The results are shown in Table 2. As shown in the table, our method outperforms all baselines on both object recall and SPRel precision. Specifically, compared with baselines that do not have spatial attention control (all but Paint-with-Words), our method is significantly better on SPRel precision, which illustrates the effectiveness of spatial control and indicates that it is difficult for the stable diffusion model to generate correct spatial relations without guidance on the layout. Our method moves the burden of generating objects at the correct location from diffusion models to a separate layout predictor, allowing the whole pipeline to achieve better spatial fidelity. On the other hand, compared with Paint-with-Words that does not impose

the temporal attention control, our method achieves better object recall, showing that our optimized combination weights across the temporal dimension strike a better balance between local and global descriptions.

Fig. 3 demonstrates examples of text descriptions and images generated by our method and baselines. We observe that our method resolves three types of errors in baselines. First, our method alleviates the missing object issue (in the top panel). Baselines tend to focus on one object in the text and ignore other objects (*e.g.,* focusing on the elephant and ignoring the man in the second row), whereas our method generates all objects. Second, as shown in the middle panel, our method mitigates the mismatched attribute issue. Particularly, baselines struggle when multiple objects are mentioned in the text, where they mismatch the attribute and object (*e.g.,* red truck in the first row). Finally, the bottom panel shows examples where our method reduces the number of mislocated objects. Note that our method is effective both on the four relations and other spatial relations. More examples can be found in Appendix G.

**Subjective Evaluation:** To further evaluate the fidelity and quality of the generated images, we perform a subjective evaluation on Amazon Mechanical Turk. Specifically, we randomly sample 25 text descriptions for each dataset in Table 2. Each subject was presented with the text description and corresponding image generated by a single method or a pair of methods and asked the following four questions: (1) (*Object Fidelity*) Does the image contain all objects mentioned in the text? (2) (*Attribute Fidelity*) Are all synthesized objects consistent with their characteristics described in the text (*e.g.,* color and material)? (3) (*Spatial Fidelity*) Does the image locate all objects at the correct position such that the spatial relations in the text are satisfied (if an object in the relationship is missing, it is considered as an incorrect generation)? and (4) (*Overall*) Which image in the pair has higher fidelity with the text and has better quality? The first three questions are evaluated for each method individually with a score of 0, 1, or 2, where 2 denotes all objects/attributes/relations are correct and 0 denotes none of them is correct. The last question is evaluated on a pair of images, one generated by our method and the other by a baseline. Each generated image is evaluated by two subjects, so the total score of the first three questions is 100.
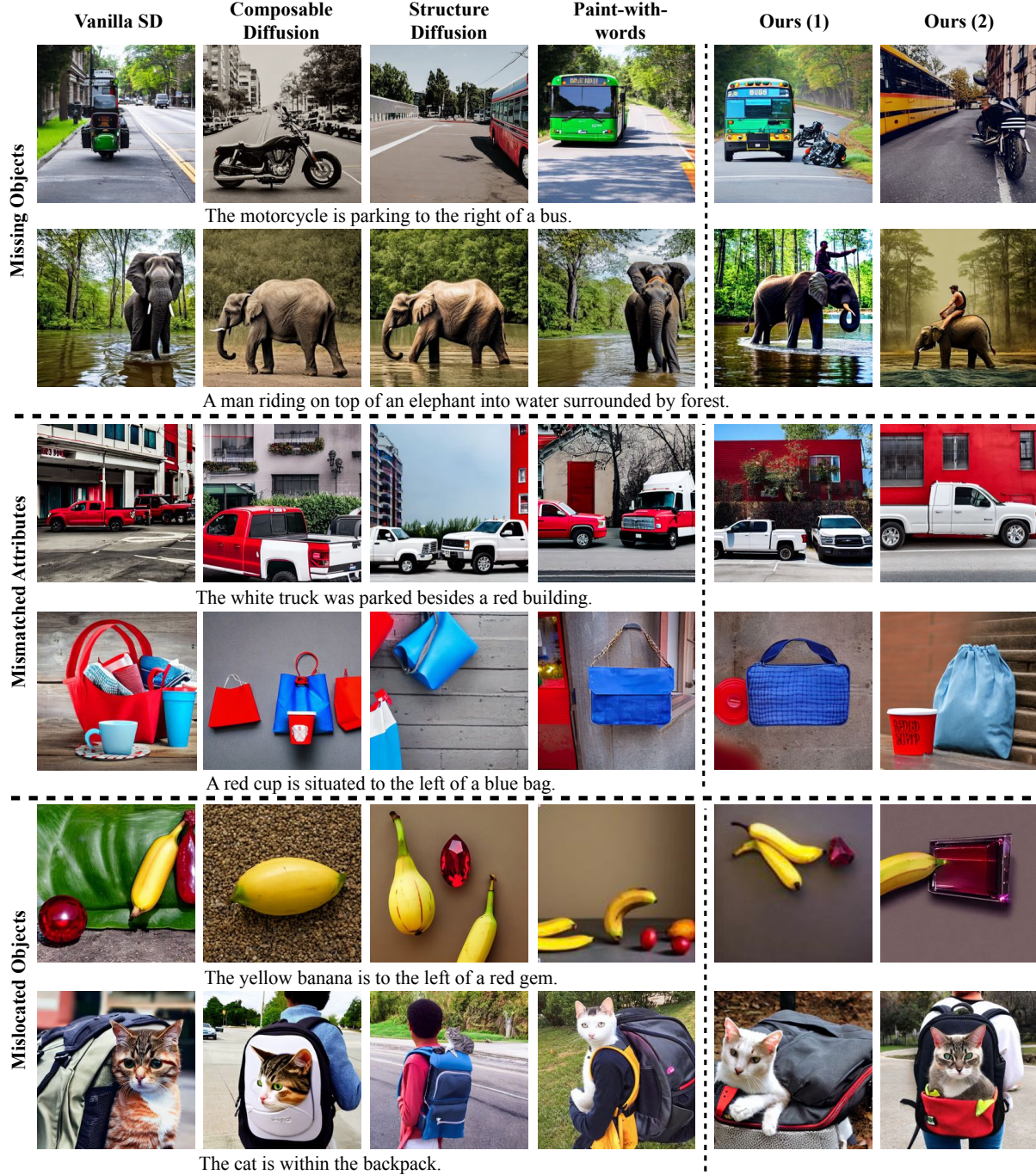
Figure 3. **Example images generated by our method and baselines.** Typical errors of baselines include missing objects, mismatched attributes, and mislocated objects. Ours (1)/(2) show the results with two different random seeds.

More details are in Appendix H.

Table 1 shows the results. Our method achieves significantly better performance in most cases, especially on GPT-synthetic dataset that contains multiple objects and relations in the same description. The results show that our method is effective at generating images with high fidelity without sacrificing perceptual quality.

## 4.2. Additional Analyses

**Performance on Each Text Complexity Level:** To explore the capability of our method, especially on complex text, we analyze the performance at each complexity level. Specifically, we use the GPT-synthetic dataset where multiple objects and relations can appear in the same description. We use the number of objects $N$ and number of spatial relations
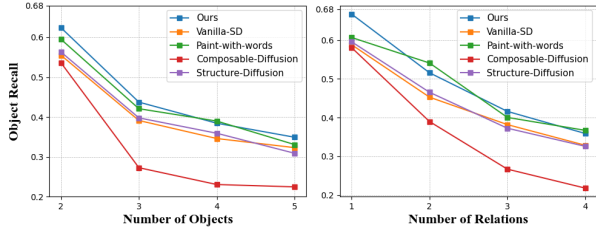
Figure 4. Performance when the number of objects and spatial relations in the text increase.
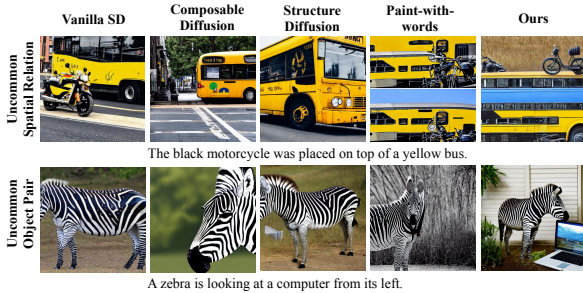


Figure 5. Example images generated by our method and baselines on novel object combinations.

$M$ introduced in Sec. 3.3 as proxies for text complexity and plot the performance of each method as the text becomes more complex. *i.e.,* we plot the performance for each value of $N$ regardless of $M$, and vise versa for $M$. As shown in Fig. 4, the performance of all methods decrease as the text becomes more complex, but our method still outperforms or is on par with others on complex descriptions. Note that PAINT-WITH-WORDS uses the same layout predictor as ours, which partially explains its strong performance.

**Generalizability to novel object combinations** is a critical requirement for text-to-image models. However, since MS-COCO dataset is collected from real photos, most of its captions contain object pairs and object-attribute pairs that are common in daily life, which are also more likely to overlap with stable diffusion's pre-training data. Thus we additionally evaluate our method on another synthetic dataset that is created similarly to Sec. 3.3, but contains *uncommon* object pairs, object-attribute pairs, and object spatial relations. Fig. 5 shows sample images generated by our method and baselines. It can be observed that our method successfully generates novel object pairs (*e.g.,* "zebra is looking at computer"). Performance on this dataset and more generated examples are shown in Appendix F. Overall, our method is able to generalize to novel object combinations.

### 4.3. Ablation Study

In this section, we will investigate the influence of two important steps in our method, namely, the layout predictor and the spatial-temporal attention. We first study how spatial-temporal attention affects our performance. To do

|  | Object Recall | SPRel Precision |
|---|---|---|
| Ours | **47.2%** | **66.7%** |
| No spatial control | 39.6% | 51.7% |
| No temporal control | 43.8% | 61.4% |
| No optimization | 41.0% | 55.6% |

Table 3. Ablation performance on GPT-synthetic test set.

that, we consider three variants of attention.

First, we explore the attention **without spatial control**. Concretely, the diffusion model still attends to both global description and local descriptions. However, instead of only attending to a local description in the pixel region of the object, the model now attends to all local descriptions in the *whole image*. The combination weights at each denoising step are optimized with only the global CLIP similarity in Sec. 3.4. Table 3 shows that its performance drops drastically, demonstrating the importance of spatial control.

Second, we explore the attention **without temporal control**, where combination weights remain the same for all denoising steps. Note that the weights can be different for each object and are optimized with $\mathcal{L}_{\text{attend}}(\boldsymbol{\lambda})$ in Sec. 3.4. The results in Table 3 show a significant degradation under both metrics, indicating the benefits of temporal dependency on attention.

Finally, we explore the attention **without optimization**. Specifically, we fix the combination weight $\lambda_{it} = \frac{1}{N}$ for all $i$ and $t$, where $N$ is the number of objects. This value is also the initialization point of $\{\lambda_{it}\}$ in our method. As shown in Table 3, the performance further drops compared to the no temporal control variant. The results indicate that optimizing the combination weights for a new text description is critical for generating images with high fidelity.

We further study the effects of our layout predictor, which includes the layout predictor trained with only one of the absolute and relative position objectives, the comparison with ground truth and user-provided pixel region, and a different strategy to construct the pixel region for an object. Results of these experiments are presented in Appendix E.

## 5. Conclusion

In this work, we study the text-to-image synthesis task based on diffusion models. We find existing methods lack explicit control on cross-attention in diffusion models, which leads to the generation of low-fidelity images. We propose an algorithm that imposes control on cross-attention in spatial and temporal aspects. Experiments show our method outperforms baselines in generating high-fidelity images. Further ablation study verifies the effectiveness of our spatial and temporal attention control. One limitation of our method is the reliance on a time-consuming optimization scheme, which takes around 10 minutes for each text-to-image generation. Future work may consider striking a better balance between performance and efficiency.

# References

[1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan, 2017. 1

[2] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. ediffi: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv*, 2022. 2, 3, 5

[3] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural image synthesis, 2018. 1

[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, 11

[5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers, 2020. 6

[6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C. Lawrence Zitnick. Microsoft coco captions: Data collection and evaluation server, 2015. 2

[7] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. 2

[8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 2

[9] Wan-Cyuan Fan, Yen-Chun Chen, DongDong Chen, Yu Cheng, Lu Yuan, and Yu-Chiang Frank Wang. Frido: Feature pyramid diffusion for complex scene image synthesis. *arXiv preprint arXiv:2208.13753*, 2022. 3

[10] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv*, 2022. 2, 3, 5

[11] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In Z. Ghahramani, M. Welling, C. Cortes, N. Lawrence, and K.Q. Weinberger, editors, *Neurips*, 2014. 1

[12] 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. 2

[13] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen

[14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Neurips*, 2020. 1, 2

[15] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *J. Mach. Learn. Res.*, pages 1–33, 2022. 2

[16] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022. 2

[17] Nisha Huang, Fan Tang, Weiming Dong, Tong-Yee Lee, and Changsheng Xu. Region-aware diffusion for zero-shot text-driven image editing. *arXiv preprint arXiv:2302.11797*, 2023. 3

[18] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation, 2017. 1

[19] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022. 1, 2

[20] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*, 2014. 5, 12

[21] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 1

[22] Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. *arXiv preprint arXiv:2009.09761*, 2020. 2

[23] Xiang Lisa Li, John Thickstun, Ishaan Gulrajani, Percy Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *arXiv preprint arXiv:2205.14217*, 2022. 2

[24] 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. *arXiv preprint arXiv:2301.07093*, 2023. 3

[25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 5, 11

[26] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial reasoning. *arXiv preprint arXiv:2205.00363*, 2022. 2, 5

[27] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B Tenenbaum. Compositional visual generation with composable diffusion models. In *ECCV*, 2022. 3, 5

[28] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv*, 2019. 5, 11, 12

[29] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *CVPR*, 2022. 2

[30] Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2, 2022. 2

[31] Chenlin Meng, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Image synthesis and

video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2

editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2

[32] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models, 2021. 2

[33] Dong Huk Park, Grace Luo, Clayton Toste, Samaneh Azadi, Xihui Liu, Maka Karalashvili, Anna Rohrbach, and Trevor Darrell. Shape-guided diffusion with inside-outside attention. *arXiv preprint arXiv:2212.00210*, 2022. 3

[34] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 5, 11, 12

[35] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2

[36] 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. 2

[37] Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. In *Neurips*, 2019. 1

[38] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. 1

[39] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 5, 11, 12

[40] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, 2022. 2

[41] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022. 2, 3

[42] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 2

[43] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv*, 2021. 5

[44] Jaskirat Singh, Stephen Gould, and Liang Zheng. High-fidelity guided image synthesis with latent diffusion models. *arXiv preprint arXiv:2211.17084*, 2022. 3

[45] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*. PMLR, 2015. 1, 2

[46] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 1, 2

[47] Jiaming Song, Arash Vahdat, Morteza Mardani, and Jan Kautz. Pseudoinverse-guided diffusion models for inverse problems. In *ICLR*, 2023. 5

[48] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021. 1, 2

[49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 4

[50] Qiucheng Wu, Yujian Liu, Handong Zhao, Ajinkya Kale, Trung Bui, Tong Yu, Zhe Lin, Yang Zhang, and Shiyu Chang. Uncovering the disentanglement capability in text-to-image diffusion models. *arXiv preprint arXiv:2212.08698*, 2022. 2, 3

[51] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. *arXiv preprint arXiv:2212.05034*, 2022. 2

[52] Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: a geometric diffusion model for molecular conformation generation, 2022. 2

[53] Cheng-Fu Yang, Wan-Cyuan Fan, Fu-En Yang, and Yu-Chiang Frank Wang. Layouttransformer: Scene layout generation with conceptual and spatial diversity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3732–3741, June 2021. 4

[54] Zuopeng Yang, Daqing Liu, Chaoyue Wang, Jie Yang, and Dacheng Tao. Modeling image composition for complex scene generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7764–7773, 2022. 3

[55] Sihyun Yu, Kihyuk Sohn, Subin Kim, and Jinwoo Shin. Video probabilistic diffusion models in projected latent space. *arXiv preprint arXiv:2302.07685*, 2023. 2

[56] Yu Zeng, Zhe Lin, Jianming Zhang, Qing Liu, John Collomosse, Jason Kuen, and Vishal M Patel. Scenecomposer: Any-level semantic image synthesis. *arXiv preprint arXiv:2211.11742*, 2022. 3

[57] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 3

[58] Yufan Zhou, Ruiyi Zhang, Changyou Chen, Chunyuan Li, Chris Tensmeyer, Tong Yu, Jiuxiang Gu, Jinhui Xu, and Tong Sun. Lafite: Towards language-free training for text-to-image generation. *arXiv preprint arXiv:2111.13792*, 2021. 2

## A. Implementation Details

To help reproduce our results, we include a comprehensive report of hyperparameters and the model architectures used in this work in Table 6.

The `Roberta-base` encoder [28] in layout predictor consists of 12 layers and 12 heads, with a hidden dimension size of 768. We fine-tune the model on GPT-synthetic dataset with relative position objective and MS-COCO dataset with absolute position objective for 100 epochs. The training batch size is 64, and the learning rate of encoder starts at 1e-6 and decays to 1e-8. The learning rate of the GMM output layer starts at 4e-5 and decays to 1e-8. We use the pre-trained `ViT-B/32` [34] to calculate CLIP similarities. For the diffusion model, we adopt the pre-trained `stable-diffusion-v1-4` [39] and stick with the default parameters. When optimizing the combination weights of cross-attention, the initial value is set to $1/N$, where $N$ is the number of objects. The weight is projected to $[-1, 2]$ after each gradient descent step to avoid extreme values.

## B. Details of GPT-synthetic Dataset and Dataset Statistics

In this section, we detail the process of creating the GPT-synthetic dataset and report the statistics of each dataset.

The GPT-synthetic dataset contains the 80 object categories in MS-COCO [25], and each description contains 2-5 objects and 1-4 relations. To create a text description with $N$ objects and $M$ relations, $N$ objects are first sampled without replacement from the same MS-COCO super-category (*e.g.,* $N$ objects from furniture), so that they are more likely to appear together in the same scene in real world. A color attribute is randomly assigned to each object with probability 0.5, and the assigned color is randomly sampled from a pre-defined list of colors. Among the $N$ objects, $M$ pairs are then sampled without replacement and randomly assigned a spatial relation from "left of," "right of," "above," and "below." We consider these four relations because they can be easily and reliably measured by comparing the center position, and we additionally check the relations to ensure no contradiction exists (*e.g.,* A is above B, B is above C, and C is above A). With specified objects and relations, GPT3 [4] is prompted to generate a sentence that mentions all objects and relations, given 5 demonstration examples. We specifically instruct GPT3 to generate diverse sentences. Table 7

|  | 2 objects | 3 objects | 4 objects | 5 objects |
|---|---|---|---|---|
| 1 relation | 200 | 50 | 0 | 0 |
| 2 relations | 0 | 50 | 50 | 0 |
| 3 relations | 0 | 0 | 50 | 50 |
| 4 relations | 0 | 0 | 0 | 50 |

Table 4. Statistics of GPT-synthetic dataset.

| | MS-COCO | | VSR | | GPT-synthetic | |
|---|---|---|---|---|---|---|
| | Global CLIP | Local CLIP | Global CLIP | Local CLIP | Global CLIP | Local CLIP |
| VANILLA-SD | 0.2890 | 0.2357 | **0.3071** | 0.2412 | **0.3006** | 0.2243 |
| COMPOSABLE-DIFFUSION | 0.2892 | **0.2397** | 0.2948 | 0.2394 | 0.2886 | 0.2327 |
| STRUCTURE-DIFFUSION | 0.2870 | 0.2339 | 0.2972 | 0.2395 | 0.2912 | 0.2396 |
| PAINT-WITH-WORDS | **0.2902** | 0.2391 | 0.2974 | 0.2394 | 0.2961 | **0.2418** |
| Ours | 0.2892 | 0.2375 | 0.3029 | **0.2415** | 0.2944 | 0.2403 |

Table 5. CLIP Similarity of our method and baselines.

shows the complete instruction, a sample demonstration, and a query that is used to generate a sentence. In practice, we manually write 20 demonstrations and randomly sample 5 for each generation.

Table 4 shows the number of text descriptions for the GPT-synthetic dataset. For the other two datasets, VSR contains 500 descriptions, and each description involves two objects and one spatial relation. MS-COCO contains 500 descriptions, including 200 descriptions with two objects, 150 descriptions with three objects, and 150 descriptions with four objects. There are no explicit spatial relations in MS-COCO descriptions. In general, GPT-synthetic contains the most complex text descriptions in terms of the number of objects and spatial relations.

## C. CLIP Similarity

We additionally use CLIP similarity [34] to measure how well the generated image aligns with the text description. Specifically, we consider CLIP similarity at two different granularities. **Global CLIP score** calculates the CLIP similarity between the whole text description and the whole image. On the other hand, **local CLIP score** calculates the similarity between the local object description and its corresponding bounding box in the image, if the object is detected. We use the local description defined in Sec. 3.2, *e.g.,* "A photo of a black mailbox."

The performance is presented in Table 5. We observe that different methods have very close global and local CLIP scores, which may be ascribed to the limited capability of vision-and-language models when text descriptions consist of multiple objects [59]. As such, the CLIP similarity is not sufficient to indicate which method performs better in our experiments, and we leverage the subjective evaluation in Sec. 4 for a more informative analysis.

## D. Layout Predictor Analyses

In this section, we analyze how the layout predictor affects our method from the following two aspects.

**Layout predictor helps synthesize correct spatial relations** We first demonstrate the position of each object predicted by our layout predictor and further show how the predicted position helps diffusion models generate objects with correct spatial relations. Fig. 6 illustrates two examples. Given the text description, the first column demon-

| Module | Attribute | Value |
|---|---|---|
| **Layout Predictor** | Model checkpoint | Roberta-base [28] |
| | Layers | 12 |
| | Heads | 12 |
| | Hidden dimension $d$ | 768 |
| | Training batch size | 64 |
| | Training epoch | 100 |
| | Learning rate | $1e-6 \rightarrow 1e-8$ for transformer layer<br>$4e-5 \rightarrow 1e-8$ for GMM |
| **Diffusion Model** | Model checkpoint | stable-diffusion-v1-4 [39] |
| | Sampling steps | 50 |
| | Sampling variance | 0.0 |
| | Resolution | $512 \times 512$ |
| | Latent channels | 4 |
| | Latent down-sampling factor | 8 |
| | Conditional guidance scale | 7.5 |
| **Attention Optimization** | Checkpoint for CLIP loss | ViT-B/32 [34] |
| | $\gamma$ | 5 |
| | Optimizer | Adam [20] |
| | Learning rate | 0.05 |
| | $\lambda_t$ initialization | $1/N$, where $N$ is object numbers |

Table 6. Hyperparameters and model architectures used in this paper.

| Instruction | Given several objects, write a sentence that describes the given objects. Additionally, if location relation between objects is specified, the sentence needs to contain sufficient information that reveals the relation. Try to generate sentences as diverse as possible and DO NOT simply state the object locations. |
|---|---|
| **Demonstrations** | Objects: silver car, green motorcycle, blue bus, yellow truck<br>Relation: silver car right of blue bus, yellow truck left of blue bus<br>Sentence: The blue bus was driving along the road, with a silver car positioned to its right and a yellow truck overtaken and left behind on the left side of the bus, while a green motorcycle zoomed past on the opposite lane. |
| **Query** | Objects: red sandwich, yellow carrot, brown hot dog, green cake<br>Relation: yellow carrot right of green cake<br>Sentence: |

Table 7. Complete instruction and example demonstration used to generate the GPT-synthetic Dataset.

strates images synthesized directly by vanilla stable diffusion model, where some objects are mislocated (*e.g.,* the spoon and bowl) and missed (*e.g.,* sandwich). The second column shows the predicted position of each object by our layout predictor, where it correctly locates objects according to the specified spatial relations (*e.g.,* the apple is placed beneath the sandwich). Finally, the last column shows images can be generated following the predicted layout, thus locating objects at correct position.

**Layout predictor does not restrict the object size** Once the position of each object is predicted, we define the pixel region of each object as a circle centered at the predicted coordinate with a radius $r$. We demonstrate two examples for images generated with different radii $r$ in Fig. 7. Generally, the object size increases as the radius increases. However, the circular region does not strictly bound the object, and objects can go beyond the region (*e.g.,* the cake in the first row is larger than predicted region). Moreover, the generation results are not highly sensitive to the choice of $r$, as shown by the similar outputs of $r = 0.15$ and $r = 0.2$. We thus fix $r = 0.2$ in our experiments.
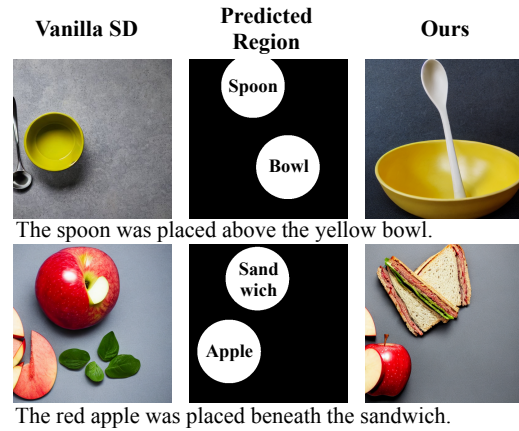


Figure 6. Layout predictor helps generate objects at correct locations. First column: Images generated by vanilla stable diffusion model. Second column: Pixel region generated by our layout predictor. Third column: Images generated by out method.

## E. Ablation Study for Layout Predictor

**Using ground truth location in place of layout predictor** In Sec. 4.1, we report the results based on our layout predic-
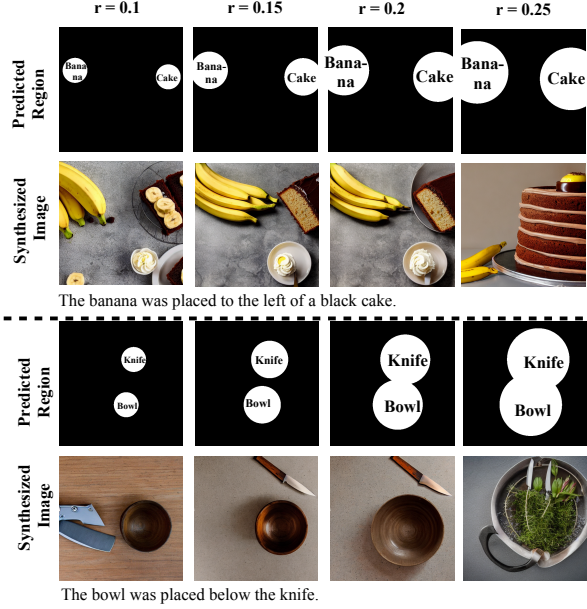
Figure 7. Example images generated by our method with different radii. For each panel, the top row is the predicted regions with different radii, and the bottom row is the generated images.

tor. We now investigate the performance when ground truth location is used in place of the layout predictor. Specifically, we use the center of the ground truth bounding box as the object position, and we use the same radius $r = 0.2$ to construct the pixel region. We evaluate the performance on VSR dataset since it contains the ground truth bounding box information. As can be observed in Table 8, providing ground truth location boosts the performance, especially for spatial relation precision, since the pixel region is guaranteed to preserve the correct spatial relations. We also evaluate the performance of PAINT-WITH-WORDS baseline when the ground truth position is given. It achieves 64.7% object recall and 58.3% SPRel precision, which is worse than the counterpart of our method.

**Training layout predictor with only absolute or relative position objective**    Our layout predictor is jointly trained with both absolution and relative position objectives (Sec. 3.3). We now explore our method's performance when the predictor is trained with only one of the objectives. The results are shown in Table 8. We observe that both settings lead to performance drop, indicating that both objectives are critical for an effective layout predictor. Moreover, removing relative position objective leads to significant performance degradation on SPRel precision, which demonstrates the importance of the objective.

**Hard versus Soft Threshold on Pixel Region**    We use a hard threshold to get the pixel region in Sec. 3.4. Here we explore a different strategy that produces a soft pixel region for an object. Specifically, we expand the pixel region of an

| | Object Recall | SPRel Precision |
|---|---|---|
| Ours | 65.1% | 75.0% |
| Ground truth position | 66.3% | 78.1% |
| No absolute position obj | 62.7% | 68.2% |
| No relative position obj | 63.5% | 58.6% |
| Soft pixel region | 64.3% | 69.8% |

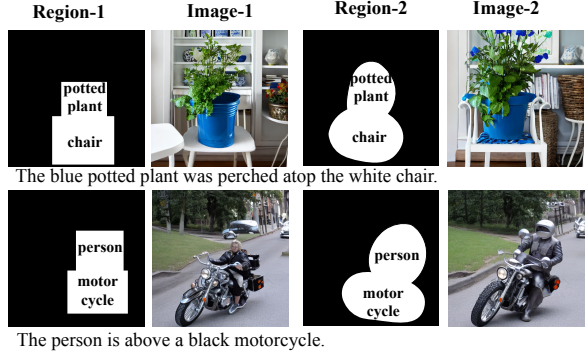Table 8. Ablation study for layout predictor on VSR dataset.



Figure 8. Example images generated from user provided region. The regions are shown in the first and third columns, and corresponding images are shown in the second and fourth columns.

object to the whole image but assign a smaller weight for pixels that are further away from the object. Formally, the output of the attention layer at time $t$ becomes

$$
\boldsymbol{O}(t) = \sum_{i=1}^{N} \lambda_{it} \boldsymbol{G}_i \odot \text{Attention}(\boldsymbol{Q}, \boldsymbol{K}_{\boldsymbol{L}_i}, \boldsymbol{V}_{\boldsymbol{L}_i}) \\
+ \left(1 - \sum_{i=1}^{N} \lambda_{it} \boldsymbol{G}_i\right) \odot \text{Attention}(\boldsymbol{Q}, \boldsymbol{K}_{\boldsymbol{D}}, \boldsymbol{V}_{\boldsymbol{D}}),
$$
(10)

where $\lambda_{it}, \boldsymbol{Q}, \boldsymbol{K}_{\boldsymbol{L}_i}, \boldsymbol{V}_{\boldsymbol{L}_i}, \boldsymbol{K}_{\boldsymbol{D}}, \boldsymbol{V}_{\boldsymbol{D}}$ are defined in Sec. 3.4, and $\boldsymbol{G}_i$ is a soft pixel region matrix for object $O_i$, with $\boldsymbol{G}_i(x, y) = g\left((x, y); \boldsymbol{C}_i, \sigma^2 \boldsymbol{I}\right) / g\left(\boldsymbol{C}_i; \boldsymbol{C}_i, \sigma^2 \boldsymbol{I}\right)$, where $g\left((x, y); \boldsymbol{C}_i, \sigma^2 \boldsymbol{I}\right)$ is the probability density of a 2D Gaussian distribution with mean $\boldsymbol{C}_i$ and covariance matrix $\sigma^2 \boldsymbol{I}$ at point $(x, y)$, $\boldsymbol{C}_i$ is the center coordinate of the object, and $\sigma$ is a hyperparameter. Intuitively, the combination weight of an object decreases as the pixel moves away from the object center, and the weights are normalized so that the object center has combination weight $\lambda_{it}$. The performance of this strategy is shown in Table 8, where it achieves a slightly worse performance compared to the hard threshold region.

Finally, we show examples in Fig. 8 where user provided region (possibly irregular) is given to our method. The generated images largely follow the provided layout, which demonstrates that our method can be adapted for image generation with better user interaction.

# F. Performance on Uncommon Combinations

To further test if our method can generate high-fidelity images for novel text descriptions, we demonstrate the per-

|                      | Object Recall | SPRel Precision |
|----------------------|---------------|-----------------|
| VANILLA-SD           | 39.8%         | 52.6%           |
| COMPOSABLE-DIFFUSION | 30.1%         | 50.8%           |
| STRUCTURE-DIFFUSION  | 40.1%         | 51.6%           |
| PAINT-WITH-WORDS     | 41.2%         | 54.7%           |
| Ours                 | **42.4%**     | **59.6%**       |

Table 9. Performance on text descriptions that contain uncommon object pairs, object-attribute pairs, and spatial relations.

formance of our method and baselines on a dataset that contains uncommon scenes. This uncommon synthetic dataset consists of 100 text descriptions, and it differs from the GPT-synthetic dataset in Sec. 3.3 from two aspects. (1) When sampling objects for a description, we remove the constraint that objects need to belong to the same super category. Sampling without this constraint can thus produce rare object pairs (*e.g.,* objects from food and vehicle can occur in the same description). (2) We manually check generated samples and only keep the ones that are unlikely to appear in real life. The result is demonstrated in Table 9. We observe that our method achieves the best result in terms of object recall and spatial relation precision, indicating that our method can better generalize to novel text descriptions. Some visual examples can be found in Fig. 5 and Fig. 9.
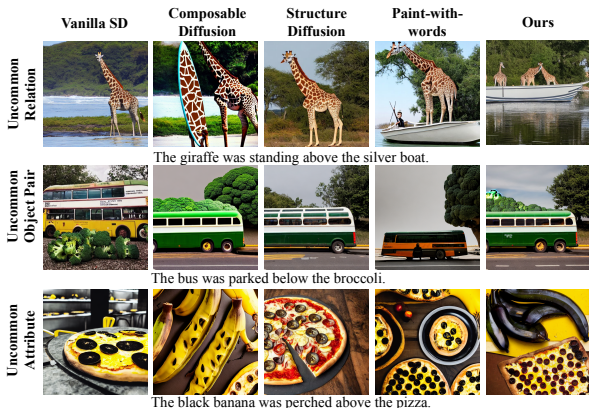


Figure 9. Example images generated by our method and baselines on uncommon relations, object pairs, and attributes.

## G. More Examples and Failure Cases

In this section, we provide more example images from our method and baselines. Then, we analyze the potential failure cases of our method.

**More Examples** We provide more example images from our method and baselines in Fig. 11. The results are consistent with Fig. 3, where our method generates images with high object, attribute, and spatial fidelities.

**Failure Cases** We present two failure cases in Fig. 10. For each example, we first show the predicted pixel region by our layout predictor and the corresponding generated im-

age (left two columns). We observe that these predicted positions tend to be at the edge of the image, which reduces the region of the object. We hypothesize that this will lead to insufficient attention to the corresponding object, and thus the object cannot be successfully synthesized (*e.g.,* the cell phone in the first row). We further demonstrate in the right two columns that moving pixel regions inside the image can resolve these failure cases, where the missing objects can be synthesized. Future work may consider adding the constraint that the predicted object center cannot locate at the edge of the image.
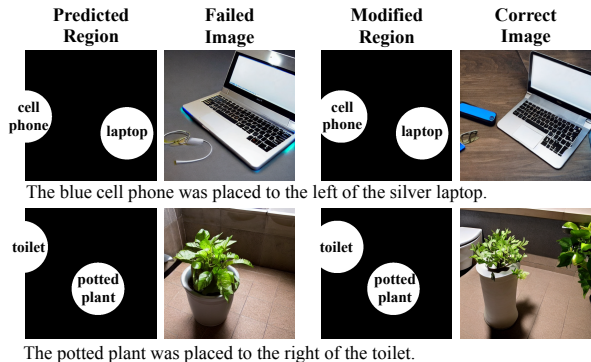


Figure 10. Failure Cases. The first two columns show the predicted region and the synthesized image, where some objects are missing. The last two columns demonstrate that modifying the pixel region can resolve the problem.

## H. Details of Subjective Evaluation

In this section, we provide further details of our subjective evaluation. We compare with all baselines on MS-COCO, VSR, and GPT-synthetic datasets. For each dataset, we randomly sample 25 text descriptions for evaluation. We evaluate on Amazon Mechanical Turk, and 85 workers participate in our study. During the subjective evaluation, the workers are asked four questions: (1) (*Object Fidelity*) Does the image contain all objects mentioned in the text? (2) (*Attribute Fidelity*) Are all synthesized objects consistent with their characteristics described in the text (*e.g.,* color and material)? (3) (*Spatial Fidelity*) Does the image locate all objects at the correct position such that the spatial relations in the text are satisfied (if an object in the relationship is missing, it is considered as an incorrect generation)? and (4) (*Overall*) Which image in the pair has higher fidelity with the text and has better quality? For the first three questions, we present the participant with a single image generated by one method, and ask the participant to rate the image using a score of 0, 1, or 2, where 2 denotes all objects/attributes/relations are correct and 0 denotes none of them is correct. For the last question, participant will see a pair of images, where one of them is generated by our method and the other one is generated by one baseline. The
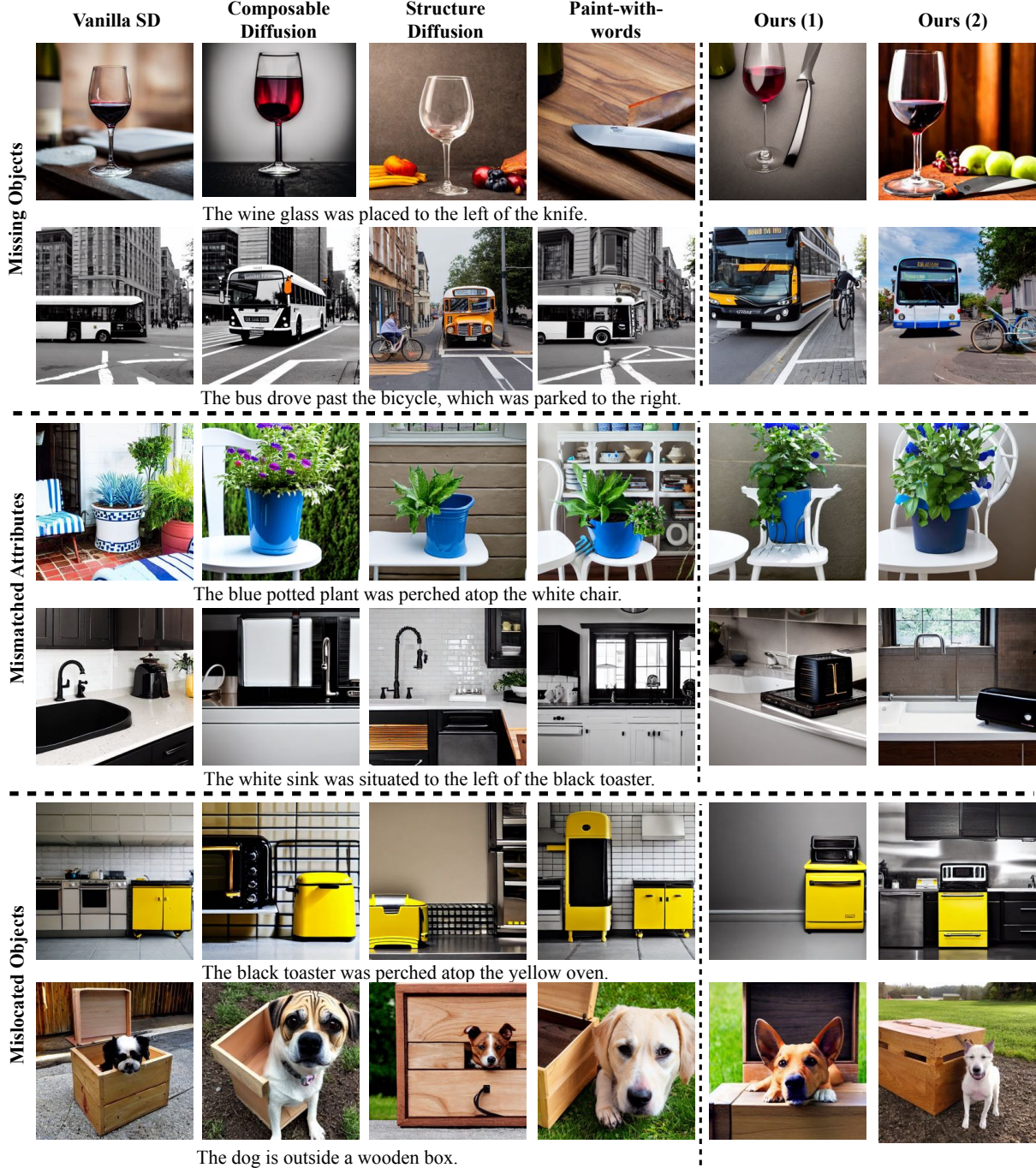
Figure 11. **Example images generated by our method and baselines.** Typical errors of baselines include missing objects, mismatched attributes, and mislocated objects. Ours (1)/(2) show the results with two different random seeds.

participant is then asked to select the better image in terms of overall fidelity and quality. The subjective evaluation interface is shown in Fig. 12 and Fig. 13. The subjective evaluation results are shown in Table 1. We also provide all generated images by our method and baselines in Figures 14, 15, and 16.

## References

[59] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visio-linguistic compositionality. In *CVPR*, 2022.

**Instructions:**

Please read the instructions carefully. Failure to follow the instructions may lead to rejection of your results. Your task will involve evaluating whether target objects have been successfully synthesized using AI models. First, you will see a text description that outlines the objects to be generated (e.g., "The bed is below the black cat."). Then you will see an image, which is generated based on the provided text by an AI algorithm. You will then be asked to evaluate if the generated image contains all the objects mentioned in text. You will use a scoring system ranging from 0 to 2, where 0 indicates all objects are incorrect or missing, 1 means some objects are incorrect or missing, and 2 means all objects are successfully generated. Notice that you should **only** rate if the objects are synthesized or not; you should disregard their inconsistencies with text description such as colors or relative positions (e.g., left/right).
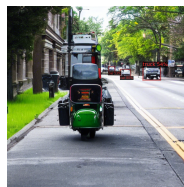
**Example:** We provide an example to help you understand how to evaluate the generated results. The text description is "**The bed is below the black cat.**"



We can observe that cat is successfully generated, but the bed is not. Therefore, this example is partially correct, and you should rate score 1. Again, notice that it synthesizes a white cat while the text says "black cat", but you should ignore the inconsistencies of color and position.

**Question:**

The text description is "**The motorcycle is parking to the right of a bus.**" Does this image contain the objects mentioned in the text description? Rate the generation results from 0 (all objects missed) to 2 (all objects are generated).



☐ 0
☐ 1
☐ 2

Figure 12. Instructions and an example question of the subjective evaluation on Amazon Mechanical Turk. The goal is to evaluate whether the generated images contain all specified objects (object fidelity). The interfaces for attribute fidelity and spatial fidelity are similar.
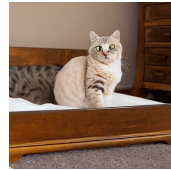
**Instructions:**

Please read the instructions carefully. Failure to follow the instructions will lead to the rejection of your results. In this task, you will be asked to judge and compare the quality of two AI-generated images. Specifically, you will first see a text description, which describes the desired content we want to generate (e.g., "The bed is below the white cat."). Then you will see two images, which are generated based on the provided text by different AI algorithms. You will then be asked to evaluate which image better follows the text description. When evaluating, you should consider the following aspects: (1) Does the synthesized image contain all objects mentioned in the text? (2) Does each object in the image follow the text description? (3) Does the image preserve the correct spatial relations mentioned in the text? (4) Does the image look real and natural? Then, you will choose the better image from the two candidate images.

**Example:** We provide an example to help you understand how to evaluate the generated results. The text description is "**The bed is below the white cat.**"
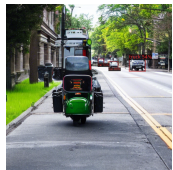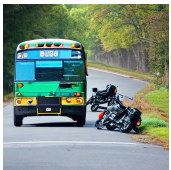


(A)                              (B)

You will evaluate based on the above criteria. First, it is important that the edited images faithfully synthesize the two objects "bed" and "cat" in the image. All methods generate a cat in the image. However, method A fails to generate high quality "bed", while method B generates a better bed. Second, both methods try to generate a white-colored cat. Third, both methods preserve the correct spatial relation that the cat is above a bed. Finally, the cat in method B looks more natural. Considering all the above analysis, method B is better.

**Question:**

The text description is "**The motorcycle is parking to the right of a bus.**" Which image in the pair has higher fidelity with the text and has better quality? Please give an overall evaluation based on the above criteria.



(A)                              (B)

☐ (A)
☐ (B)

Figure 13. Instructions and an example question of the subjective evaluation on Amazon Mechanical Turk. The goal is to compare two images generated by baselines and our method.
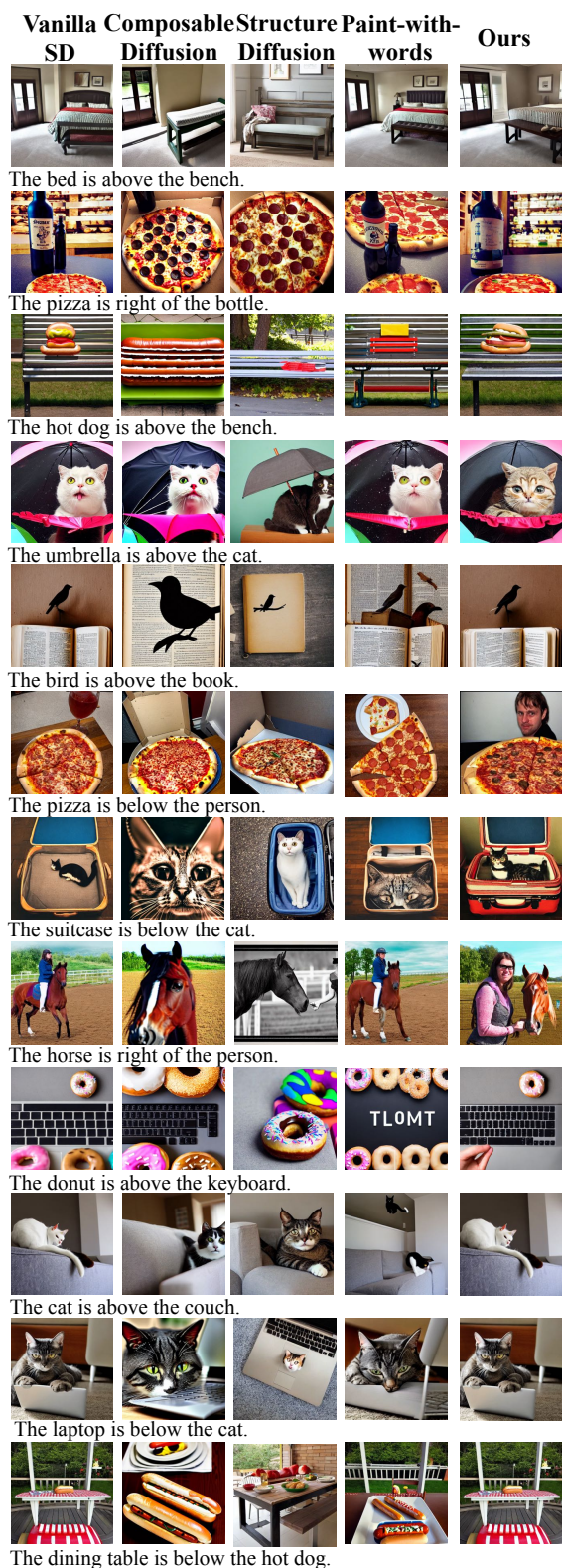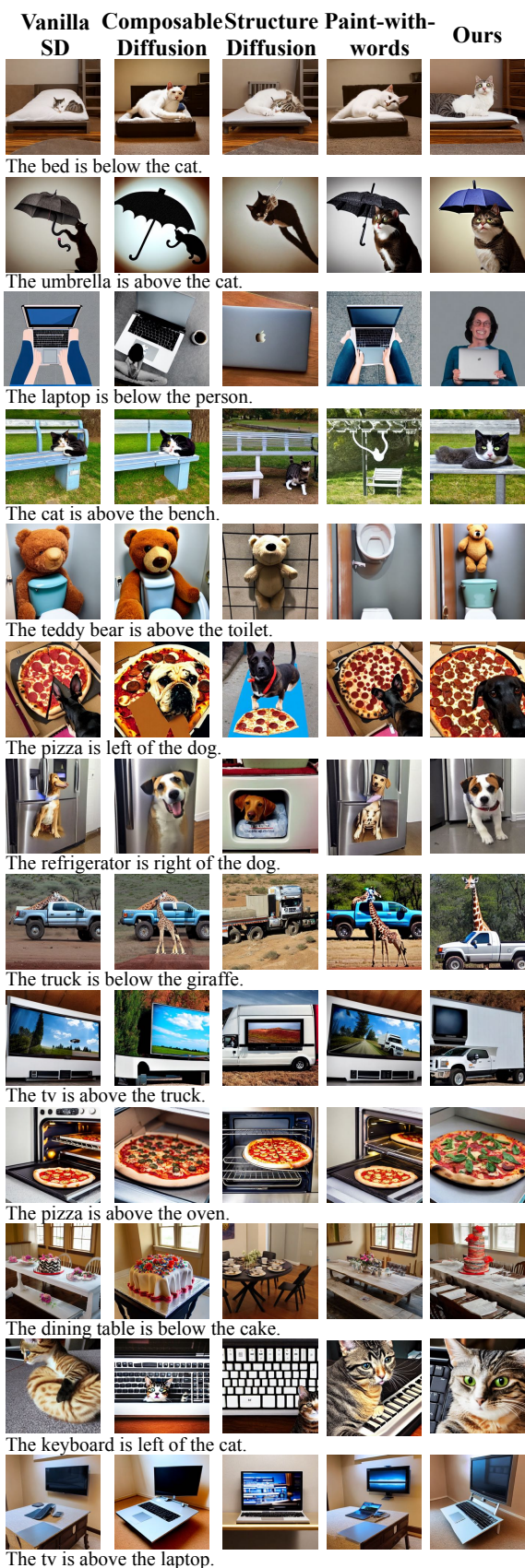
Figure 14. Generated images for subjective evaluation on VSR dataset.

| Vanilla SD | Composable Diffusion | Structure Diffusion | Paint-with-words | Ours | Vanilla SD | Composable Diffusion | Structure Diffusion | Paint-with-words | Ours |

The motorcycle is parking to the right of a bus.

The sink was installed above the white oven.

A red cup is situated to the left of a blue bag.

The red bottle was perched above the spoon on the counter.

The brown teddy bear is placed high above the toilet.

The person was standing to the left of the yellow snowboard.

The silver bed was situated to the right of the white couch.

The book was placed underneath the scissors.

The person is to the left of a tennis racket.

The black broccoli was placed to the left of the vibrant red orange.

The silver laptop was perched atop the green keyboard.

The sports ball is above the skateboard.

The red cup was placed to the left of the brown knife.

The wine glass was placed to the left of the knife.

The yellow chair was placed to the right of the white potted plant.

The donut was placed to the left of the apple.

The motorcycle is parking to the right of a bus.

The banana was placed atop the brown cake.

The person was standing beneath the black donut.

The donut was perched right of the black broccoli.

The brown pizza was placed on the plate, with a single stalk of broccoli to its right.

The blue boat was parked on the side of the road, with a motorcycle parked to its right.

The green suitcase was placed on the ground, and the umbrella was opened up above it to provide shade.

The green frisbee lay to the left of the blue tennis racket, ready to be picked up and thrown.

The person was holding a brown bowl in their right hand, while the bowl was positioned to their left.
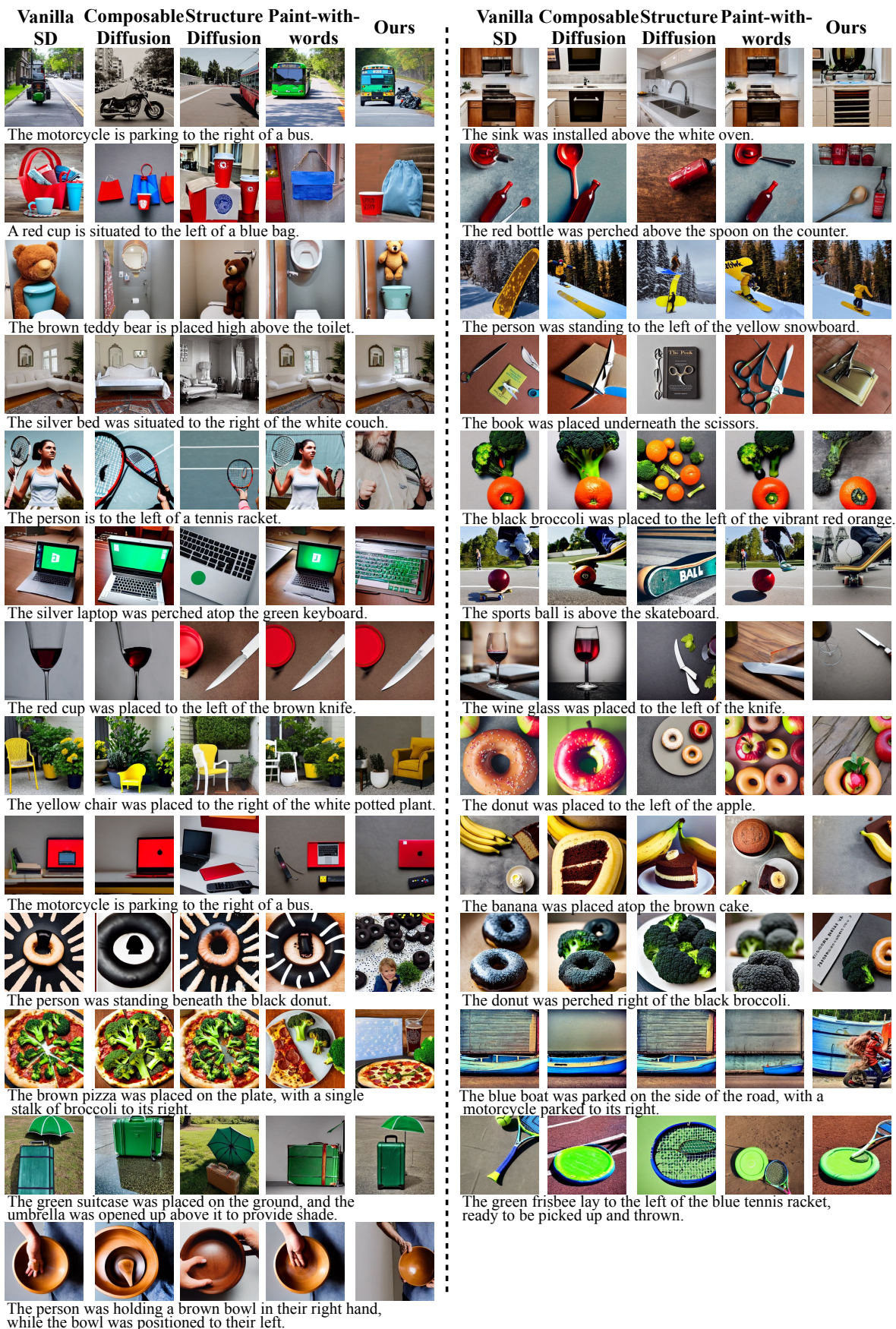
Figure 15. Generated images for subjective evaluation on GPT-synthetic dataset.
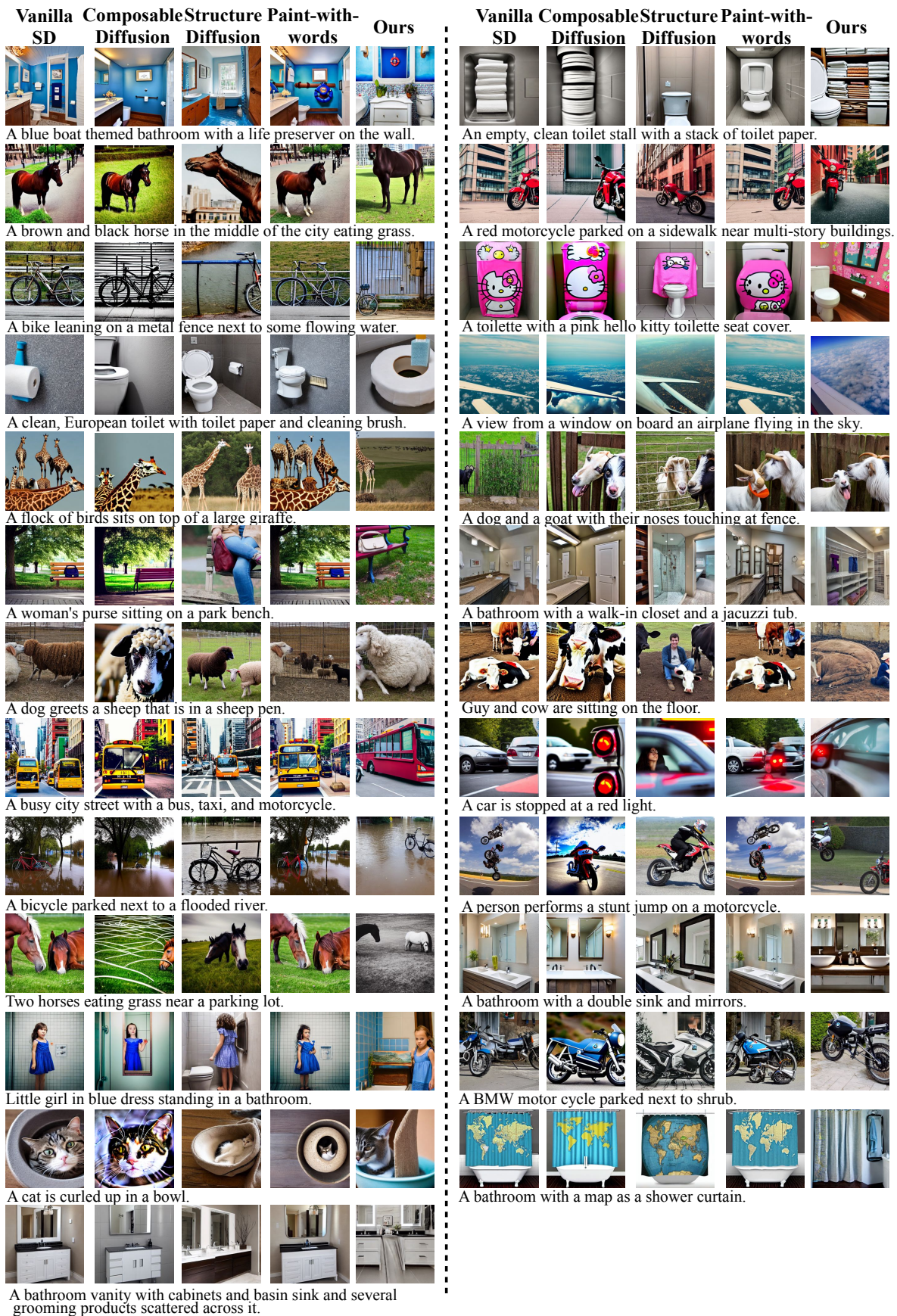
Figure 16. Generated images for subjective evaluation on MS-COCO dataset.