# Eliminating Contextual Prior Bias for Semantic Image Editing via Dual-Cycle Diffusion

Zuopeng Yang, Tianshu Chu, Xin Lin, Erdun Gao, Daqing Liu, Jie Yang, *Senior Member, IEEE*, and Chaoyue Wang

arXiv:2302.02394v3 [cs.CV] 5 Oct 2023

*Abstract*—The recent success of text-to-image generation diffusion models has also revolutionized semantic image editing, enabling the manipulation of images based on query/target texts. Despite these advancements, a significant challenge lies in the potential introduction of contextual prior bias in pretrained models during image editing, *e.g.*, making unexpected modifications to inappropriate regions. To address this issue, we present a novel approach called Dual-Cycle Diffusion, which generates an unbiased mask to guide image editing. The proposed model incorporates a Bias Elimination Cycle that consists of both a forward path and an inverted path, each featuring a Structural Consistency Cycle to ensure the preservation of image content during the editing process. The forward path utilizes the pre-trained model to produce the edited image, while the inverted path converts the result back to the source image. The unbiased mask is generated by comparing differences between the processed source image and the edited image to ensure that both conform to the same distribution. Our experiments demonstrate the effectiveness of the proposed method, as it significantly improves the D-CLIP score from $0.272$ to $0.283$. The code will be available at https://github.com/JohnDreamer/DualCycleDiffsion.

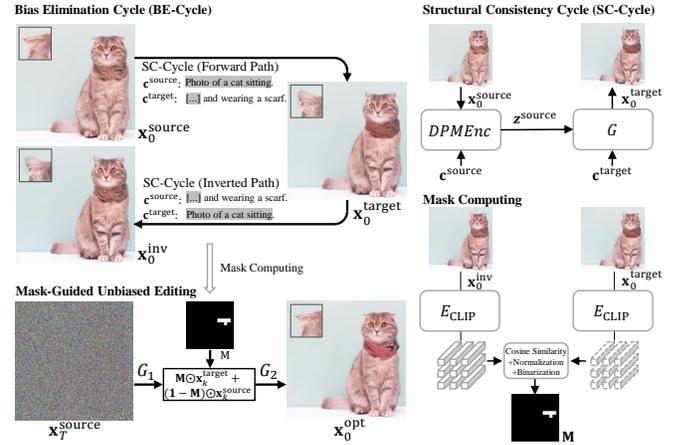*Index Terms*—Contextual prior bias, semantic image editing, Dual-Cycle Diffusion



Fig. 1. The overview of the proposed Dual-cycle Diffusion framework for the semantic image editing task. The pipeline and details of the Bias Elimination Cycle (BE-Cycle) and mask-guided unbiased editing are shown on the left. **The modifications on cat ears**, caused by **the contextual prior bias derived from the pre-trained model**, are illustrated in the left-top corner of the images. Given a source image, a source text, and a target text, we first leverage BE-Cycle to produce an unbiased mask, which is then used to guide image editing. On the right, the details of the Structural Consistency Cycle (SC-Cycle) [7] and the procedure of mask computing are shown. $\odot$ is the element-wise product.

## I. INTRODUCTION

SEMANTIC image editing is a critical and demanding problem within the field of image processing, with the objective of modifying an existing image based on a specified textual transformation query. Unlike conventional image editing techniques [42], such as those utilizing depth [38], skeleton [40], edge [39], heatmap [1], or segmentation [3], [37], the utilization of text [4]–[6], [43] to guide the editing process offers a more versatile and user-friendly approach to achieving desired modifications to the image. As illustrated in Fig 1, given an image of a cat and a query "*Photo of a cat sitting and wearing a scarf*", the aim is to add a scarf to the cat's neck while maintaining the integrity of the other

Zuopeng Yang, Tianshu Chu, and Jie Yang are with the Institute of Image Processing and Pattern Recognition, Department of Automation, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: {yzpeng, chutianshu, jieyang}@sjtu.edu.cn). Xin Lin is with the Institute of Artificial Intelligence and Blockchain, Guangzhou University, Guangzhou 510006, China (e-mail: linxj68@gmail.com). Erdun Gao is with the School of Mathematics and Statistics, The University of Melbourne, Melbourne, VIC 3010, Australia (e-mail: erdun.gao@student.unimelb.edu.au). Daqing Liu and Chaoyue Wang are with the JD Explore Academy, Beijing 100176, China (e-mail: liudq.ustc@gmail.com, chaoyue.wang@outlook.com).

regions of the image. Thus, this task can be regarded as an extension of text-conditioned image generation, with a precise and interpretable differentiation between regions that require modification and those that should be preserved.

In recent years, there has been substantial progress in text-conditioned generation techniques, with the advent of models such as DALL-E [19], Make-A-Scene [20], Imagen [21], Parti [25], and DALL-E2 [22]. These models are trained on enormous amounts of data collected from the Internet, resulting in significant improvements in the capacity for both textual semantic understanding and image synthesis. Among these models, diffusion-based models have garnered significant attention due to the iterative optimizations from random Gaussian noise to photo-realistic images.

As a pioneering method, SDEdit [13] achieves the editing through iterative denoising via a stochastic differential equation (SDE), which tends to modify the entire image. Prompt2prompt [24] localizes the edit by controlling the cross-attention maps to eliminate unexpected modifications only on synthesized images. DiffEdit [23] automatically generates a mask by contrasting the predictions of a diffusion model conditioned on different text prompts. However, it has been observed that it fails to produce masks focusing on the regions

requiring editing. Furthermore, to preserve the integrity of the image, CycleDiffusion [7] presents the DPMEncoder to encode the latent code.

Despite the benefits of utilizing pre-trained models for semantic image editing, a significant challenge arises in the form of prior bias introduced by distribution shifts between the training and target editing images. In this paper, we specifically focus on the contextual prior bias problem, which often leads to spurious correlations between different regions, thus resulting in undesired modifications in areas that were not intended to be altered, but still meet the contextual requirements of the target text. For instance, as exemplified in the top-left corner of Fig 1, when attempting to add a scarf to the cat's neck, the pre-trained model fails to maintain the shape of the ears.

To the best of our knowledge, our paper is one of the first studies to eliminate the contextual prior bias in text-guided semantic image editing. We introduce Dual-Cycle Diffusion, a method that effectively eliminates contextual prior biases by generating an unbiased mask to guide image editing. The method comprises a Bias Elimination Cycle (BE-Cycle) that consists of forward and inverted paths, each incorporating a Structural Consistency Cycle (SC-Cycle) to maintain image content integrity during the editing process. Given a target text, the forward path utilizes a pre-trained model to synthesize the edited image. The inverted path then converts the edited image back to the source image using a source text that describes the original image.

Additionally, the method allows for editing of both real images and synthetic images generated by other models, and the source text need not be the one used to generate the original image, so long as it adequately describes its main content. Specifically, the source text could be automatically produced by using an off-the-shelf captioning model [2] or obtained by making slight modifications, such as changing, adding, or removing several words from the target text. As shown in Fig 1, an example of this would be the removal of the phrase "and wearing a scarf. The resulting processed source image and edited image conform to the same distribution, enabling the automatic generation of an unbiased mask through the identification of differences in their features. This is achieved through an improved mask computing pipeline that utilizes the encoder of CLIP [10] to extract features from both images. The generated mask is then employed to guide unbiased semantic image editing. Comprehensive experiments demonstrate the proposed method's superiority, which significantly improves quantitative metrics and qualitative visual results.

## II. THE PROPOSED METHOD

### A. Preliminaries

**Diffusion probabilistic models (DPMs).** DPMs [8], [26]–[28] are a class of generative models that attempt to estimate the underlying data distribution $p(\mathbf{x})$ by gradually denoising a normally distributed variable. This can be seen as learning the inverse procedure of a fixed Markov Chain of length $T$. In the context that semantic image editing is an extension of the text-to-image generation task, our current focus is solely towards text-conditioned DPMs. During the inference stage, when given a text query $\mathbf{c}$, DPMs utilize a parametric mean estimator $\mu_\theta$ to progressively denoise a synthesized image sampled from white Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$. Using the reparameterization trick, we can sample $\mathbf{x}_{t-1} \sim \mathcal{N}(\mu_\theta(\mathbf{x}_t, \mathbf{c}, t), \sigma_t^2 \boldsymbol{I})$ as:

$$\mathbf{x}_{t-1} = \mu_\theta(\mathbf{x}_t, \mathbf{c}, t) + \sigma_t \boldsymbol{\epsilon}_t, \; t = 1, \cdots, T, \quad (1)$$

where $\sigma_t$ is the standard deviation in the $t$-th step and $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \boldsymbol{I})$. Note that the image generation direction proceeds from $\mathbf{x}_T$ to $\mathbf{x}_0$. To simplify notation, we use $\mathbf{z} := \mathbf{x}_T \oplus \boldsymbol{\epsilon}_{1:T}$ as the latent code throughout this paper, where $\oplus$ denotes concatenation. Consequently, the mapping from $\mathbf{x}_T$ to $\mathbf{x}_0$ can be represented as $\mathbf{x}_0 = G(\mathbf{z}, \mathbf{c})$.

**DPM-Encoder.** To obtain the latent code $\mathbf{z}$, Wu et al. [7] proposed DPM-Encoder, an invertible encoder for stochastic DPMs. In detail, stochastic DPMs define a posterior distribution $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ [8], [9] in the forward diffusion process. Then DPM-Encoder produces noisy images $\mathbf{x}_1, \cdots, \mathbf{x}_T$ from $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$ and computes the $\boldsymbol{\epsilon}_t$. Formally, the overall sampling process of DPM-Encoder $\mathbf{z} \sim q_{\text{DPMEnc}}(\mathbf{z}|\mathbf{x}_0, \mathbf{c}, \mu_\theta)$ can be defined as:

$$
\begin{aligned}
\mathbf{x}_1, \cdots, \mathbf{x}_T &\sim q(\mathbf{x}_{1:T}|\mathbf{x}_0), \\
\boldsymbol{\epsilon}_t &= (\mathbf{x}_{t-1} - \mu_\theta(\mathbf{x}_t, \mathbf{c}, t))/\sigma_t, \; t = 1, \cdots, T, \\
\mathbf{z} &:= \mathbf{x}_T \oplus \boldsymbol{\epsilon}_{1:T}.
\end{aligned}
\quad (2)
$$

### B. Structural Consistency Cycle

The goal of SC-Cycle [7] is to preserve the integrity of the image content during the editing process. The framework is shown in the top-right corner of Fig 1. Given a source image $\mathbf{x}_0^{\text{source}}$ and a source description $\mathbf{c}^{\text{source}}$, SC-Cycle first encodes the latent code $\mathbf{z}^{\text{source}}$ by a DPM-Encoder. Then the edited image $\mathbf{x}_0^{\text{target}}$ is generated by applying the mapping $G$ with the target text $\mathbf{c}^{\text{target}}$:

$$
\begin{aligned}
\mathbf{z}^{\text{source}} &\sim q_{\text{DPMEnc}}(\mathbf{z}|\mathbf{x}_0^{\text{source}}, \mathbf{c}^{\text{source}}, \mu_\theta), \\
\mathbf{x}_0^{\text{target}} &= G(\mathbf{z}^{\text{source}}, \mathbf{c}^{\text{target}}) = G(\mathbf{x}_T^{\text{source}} \oplus \boldsymbol{\epsilon}_{1:T}^{\text{source}}, \mathbf{c}^{\text{target}}).
\end{aligned}
\quad (3)
$$

The edited image $\mathbf{x}_0^{\text{target}}$ synthesized by SC-Cycle [7] preserves the details specific to the source image $\mathbf{x}_0^{\text{source}}$ by utilizing the latent code $\mathbf{z}^{\text{source}}$. This property leads to the similarity in structure between $\mathbf{x}_0^{\text{source}}$ and $\mathbf{x}_0^{\text{target}}$.

### C. Bias Elimination Cycle

BE-Cycle aims to automatically generate an unbiased mask by contrasting the generated images of the forward path and the inverted path. As illustrated in the top-left corner of Fig 1, each path comprises an SC-Cycle. In the forward path, $\mathbf{x}_0^{\text{target}}$ is synthesized following Eq. (3). Then in the inverted path, we attempt to restore the source image from $\mathbf{x}_0^{\text{target}}$, conditioned on the source text $\mathbf{c}^{\text{source}}$:

$$
\begin{aligned}
\mathbf{z}^{\text{target}} &\sim q_{\text{DPMEnc}}(\mathbf{z}|\mathbf{x}_0^{\text{target}}, \mathbf{c}^{\text{target}}, \mu_\theta), \\
\mathbf{x}_0^{\text{inv}} &= G(\mathbf{z}^{\text{target}}, \mathbf{c}^{\text{source}}) = G(\mathbf{x}_T^{\text{target}} \oplus \boldsymbol{\epsilon}_{1:T}^{\text{target}}, \mathbf{c}^{\text{source}}).
\end{aligned}
\quad (4)
$$

The forward and inverted paths of BE-cycle both rely on the same pre-trained model, resulting in a similar contextual prior bias in the generated images $\mathbf{x}_0^{\text{target}}$ and $\mathbf{x}_0^{\text{inv}}$. In other words, both images are conforming to the same distribution. Therefore, the contextual prior bias does not affect identifying

the regions that need to be edited by contrasting $\mathbf{x}_0^{\text{target}}$ and $\mathbf{x}_0^{\text{inv}}$.

To produce the mask, the first step is to extract the visual features of $\mathbf{x}_0^{\text{target}}$ and $\mathbf{x}_0^{\text{inv}}$ using a visual encoder $E_{\text{CLIP}}$ from a pre-trained CLIP model [10]:

$$\mathbf{F} = E_{\text{CLIP}}(\mathbf{x}_0), \ \mathbf{F} \in \mathbb{R}^{m \times n \times d}, \tag{5}$$

where $\mathbb{R}^{m \times n}$ is the spatial space, $d$ is the feature dimension, and $\mathbf{F}_{i,j,:}$ is the image grid feature with the grid's coordinate $(i, j)$, for $i \in \{1, 2, \cdots, m\}$ and $j \in \{1, 2, \cdots, n\}$. To measure the degree of change in each region, the cosine similarity between $\mathbf{F}_{i,j,:}^{\text{target}}$ and $\mathbf{F}_{i,j,:}^{\text{inv}}$ is calculated by:

$$\mathbf{S}_{i,j} = \cos \left\langle \mathbf{F}_{i,j,:}^{\text{target}}, \mathbf{F}_{i,j,:}^{\text{inv}} \right\rangle. \tag{6}$$

Finally, we can obtain the mask $\mathbf{M}$ by:

$$\mathbf{M}_{i,j} = \begin{cases} 1, & \text{if } \frac{\text{abs}(\mathbf{S}_{i,j}) - \min(\text{abs}(\mathbf{S}))}{\max(\text{abs}(\mathbf{S})) - \min(\text{abs}(\mathbf{S}))} > \delta, \\ 0, & \text{else}, \end{cases} \tag{7}$$

where $\delta$ is a threshold to control the binarization of the mask, $\max(\cdot)$, $\min(\cdot)$ obtain the maximal value, and the minimal value of a matrix, respectively, while $\text{abs}(\cdot)$ returns the absolute value of each element in a matrix or a scalar input. In practice, we set $\delta = 0.5$. To increase the resolution of the mask, the image is divided into $2 \times 2$ grids, each of which respectively generates a mask. All grid masks are then assembled to get the final mask according to their spatial positions. Finally, the mask is resized to match the spatial size of $\mathbf{x}_0$.

### D. Mask-Guided Unbiased Editing

In the final stage, the generated unbiased mask is used to guide the image editing process. The pipeline is illustrated in the left-bottom corner of Fig 1. As pointed in [11], the text-conditioned DPMs mainly rely on the text prompt to guide the sampling process at the early sampling stage, while gradually shifting towards visual features, as the generation continues. Therefore, we utilize the text prompt to guide the generation within the mask, while keeping the regions outside the mask as similar to the source image as possible. For convenience, we divide the mapping $G$ into two steps: $G(\cdot) = G_2(G_1(\cdot))$. In the first step, we obtain the image strongly influenced by the target text:

$$\mathbf{x}_k^{\text{target}} = G_1(\mathbf{x}_T^{\text{source}} \oplus \boldsymbol{\epsilon}_{(k-1):T}^{\text{source}}, \mathbf{c}^{\text{target}}), \tag{8}$$

where $k \in \{1, 2, \cdots, T\}$ indicates the $k$-th step of the diffusion process. Then we sample $\mathbf{x}_k^{\text{source}}$ from the posterior distribution $q(\mathbf{x}_k | \mathbf{x}_0^{\text{source}})$ and optimize the image by the mask:

$$\mathbf{x}_k^{\text{opt}} = \mathbf{M} \odot \mathbf{x}_k^{\text{target}} + (\mathbf{1} - \mathbf{M}) \odot \mathbf{x}_k^{\text{source}}, \tag{9}$$

where $\odot$ is the element-wise product. Finally, the edited image without the effect of contextual prior bias is synthesized by:

$$\mathbf{x}_0^{\text{opt}} = G_2(\mathbf{x}_k^{\text{opt}} \oplus \boldsymbol{\epsilon}_{1:k}^{\text{source}}, \mathbf{c}^{\text{target}}). \tag{10}$$

The unbiased mask only focuses on the regions that need editing, thus avoiding unexpected modification on the final edited image $\mathbf{x}_0^{\text{opt}}$, such as the cat ears.

TABLE I
QUANTITATIVE RESULTS ON THE ZERO-SHOT SEMANTIC IMAGE EDITING DATASET [7]. ∗ DENOTES THAT WE RAN ONLY 1 TRIAL FOR EACH HYPERPARAMETER COMBINATION, WHILE WE RAN 15 TRIALS IN THE REST EXPERIMENTS. BEST AND SECOND RESULTS ARE IN HIGHLIGHT.

| | Method | $\mathcal{S}_{\text{CLIP}} \uparrow$ | $\mathcal{S}_{\text{D-CLIP}} \uparrow$ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|
| LDM-400M | SDEdit [13] | 0.332 | 0.264 | 13.68 | 0.390 |
| | DiffEdit [23] | 0.323 | 0.208 | 18.50 | 0.588 |
| | CycleDiffusion [7] | 0.333 | 0.275 | 18.72 | 0.625 |
| | **Ours** | 0.333 | 0.281 | 19.12 | 0.635 |
| SD-v1-1 | SDEdit [13] | 0.339 | 0.248 | 15.50 | 0.498 |
| | DiffEdit [23] | 0.316 | 0.175 | 22.47 | 0.729 |
| | CycleDiffusion [7] | 0.331 | 0.262 | 21.98 | 0.731 |
| | **Ours** | 0.332 | 0.265 | 22.21 | 0.734 |
| SD-v1-4 | SDEdit [13] | 0.344 | 0.258 | 15.93 | 0.512 |
| | DiffEdit [23] | 0.318 | 0.172 | 22.70 | 0.731 |
| | CycleDiffusion [7] | 0.334 | 0.272 | 21.92 | 0.731 |
| | **Ours** | 0.337 | 0.283 | 22.08 | 0.730 |
| SD-v1-4 | SDEdit* [13] | 0.335 | 0.221 | 15.64 | 0.505 |
| | DiffEdit* [23] | 0.305 | 0.088 | 22.72 | 0.733 |
| | CycleDiffusion* [7] | 0.327 | 0.206 | 21.85 | 0.721 |
| | **Ours*** | 0.331 | 0.235 | 22.71 | 0.741 |

TABLE II
COMPARISONS BY USING DIFFERENT MASKS WITH/WITHOUT THE CLIP [10] ENCODER. THE BIASED/UNBIASED MASK IS GENERATED BY CONTRASTING THE INPUT AND OUTPUT OF THE FORWARD/INVERTED PATH IN THE BE-CYCLE. ALL THE EXPERIMENTS ARE CONDUCTED WITH THE SD-v1-4 PRE-TRAINED MODEL.

| | w/ CLIP Encoder | $\mathcal{S}_{\text{CLIP}} \uparrow$ | $\mathcal{S}_{\text{D-CLIP}} \uparrow$ | PSNR↑ | SSIM↑ |
|---|---|---|---|---|---|
| w/ Biased Mask | ✗ | 0.336 | 0.281 | 21.49 | 0.725 |
| | ✓ | 0.338 | 0.281 | 21.85 | 0.724 |
| w/ Unbiased Mask | ✗ | 0.337 | 0.282 | 21.90 | 0.725 |
| | ✓ | 0.337 | 0.283 | 22.08 | 0.730 |

## III. EXPERIMENTS

### A. Implementation Details

Experiments were conducted on the zero-shot semantic image editing dataset [7], which was specifically gathered for semantic image editing task. The dataset consists of a set of 150 tuples, each containing a source image, a source text, and a target text. For a fair comparison, DiffEdit [23] used the settings from its original paper and the remaining experiments followed the settings of CycleDiffusion [7] to use the DDIM sampler ($\eta = 0.1$) with 100 steps, set the classifier-free guidance scale of the encoding process as 1, and enumerate the classifier-free guidance scale of the decoding step as $\{1, 1.5, 2, 3, 4, 5\}$. To preserve image content, editing began with an image that was not fully noised. Therefore, we enumerated the step of adding noise as $\{85, 80, 75, 70, 60, 50\}$.[1] Then the optimization step $k$ was selected from $\{85, 80, 75, 70, 60, 50\}$. We ran 15 trials for each hyperparameter combination. To remove the effect of random noise in the mask computing process, we averaged all the masks generated with different hyperparameter combinations. The final results were automatically selected based on the directional CLIP score $\mathcal{S}_{\text{D-CLIP}}$.

**Metrics:** To evaluate editing performance over all comparison methods and our Dual-Cycle Diffusion, we adopted four metrics to evaluate edited images' faithfulness to source images and authenticity to target texts. They are **SSIM**, **PSNR**, CLIP score $\mathcal{S}_{\text{CLIP}}$ [10], and directional CLIP score $\mathcal{S}_{\text{D-CLIP}}$ [12]. **SSIM** is used to measure the similarity between

---

[1]It is the same with the early stop step of $\{15, 20, 25, 30, 40, 50\}$ in [7].

two images, and **PSNR** is used to quantify image quality. $\mathcal{S}_{\textbf{CLIP}}$ can assess the alignment between the generated image and the target text, while $\mathcal{S}_{\textbf{D-CLIP}}$ can evaluate the similarity between the images' and texts' changes. In addition, we conducted a user study to further evaluate the effectiveness of the model in eliminating contextual prior bias. The participants were instructed to vote for the images that had fewer unexpected modifications according to the target texts and source images, resulting in 1000 votes per method-to-method comparison. The results are shown in Table III.

### B. Comparisons with Existing Methods

To validate the effectiveness of the proposed Dual-Cycle Diffusion, we compared our model with other methods based on diffusion model, including SDEdit [13], DiffEdit [23], and CycleDiffusion [7]. To investigate the influence of data size, data quality, and training details on models' performance, several pre-trained text-to-image diffusion models are used: (1) LDM-400M, an LDM model [17] with 1.45B parameters, trained on LAION-400M [15], (2) SD-v1-1, a Stable Diffusion model [17] with 0.98B parameters, trained on LAION-5B [16], (3) SD-v1-4, finetuned from SD-v1-1 for improved aesthetics and classifier-free guidance sampling.

The quantitative results by the involved competitors are presented in Table I. The proposed method achieved the best or comparable scores on all metrics among the three pre-trained models, indicating its ability to enhance both edited images' faithfulness to source images and authenticity to target texts. Comparisons of the results of LDM-400M and SD-v1-4 suggest that large training data size can aid in comprehending image contents, leading to better preservation of image content when using SD-v1-4. Additionally, we can observe from a comparison of the SD-v1-1 and SD-v1-4 results that the improved classifier-free guidance sampling is useful for synthesizing edited images that are more authentic to target texts. Overall, using SD-v1-4 as the pre-trained model resulted in the best performance. Hence, all subsequent experiments and comparisons will be based on the SD-v1-4 model.

In Fig 2, we provide several visual comparisons of edited results generated by different methods. According to the visual comparisons, we can observe that the proposed method not only has accomplished modifications according to the text, but also has better image content preservation. For instance, as shown in the first row of Fig 2, the goal is to replace the sheep with a tiger. SDEdit made modifications to the entire image, while both DiffEdit and CycleDiffusion were able to maintain most parts of the source image. But they still unexpectedly added a door handle to the car, which was marked by a yellow box. This observation can also demonstrate the existence of contextual prior bias. Furthermore, the second row aims to add an apple. CycleDiffusion failed to preserve the shapes of the backpack and the apple already in the source image, while our method succeeded. DiffEdit even failed to produce another apple. The third row provides more results generated by the proposed method.

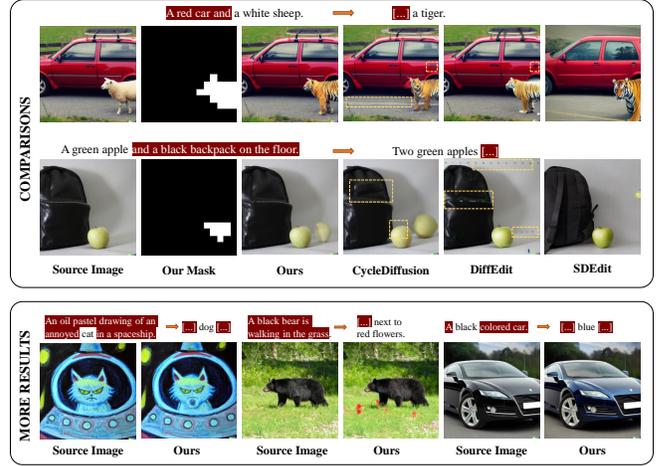| | CycleDiffusion | Ours w/ Biased Mask | Ours w/ Unbiased Mask |
|---|---|---|---|
| Votes (%) | 42.7 | 46.8 | *Reference* |



Fig. 2. The masks and edited samples generated by our proposed method. The first two rows show comparisons with two most representative baseline models: SDEdit [13], DiffEdit [23], and CycleDiffusion [7]. Among all samples, our method outperforms the other models in image content preservation. Yellow boxes mark unexpected modifications, and the masks reveal the specific edited areas. The last row shows more results generated by the proposed method.
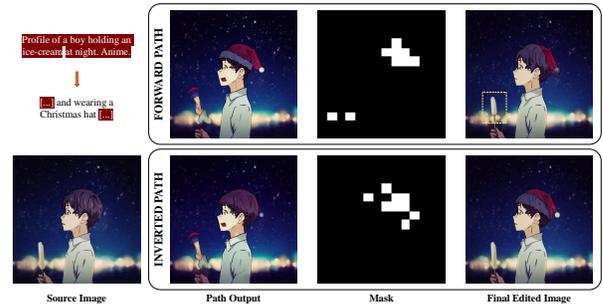


Fig. 3. The BE-Cycle's forward path and inverted path's masks, outputs, and final edited results. Specifically, the mask is generated by comparing the input and output of the forward or inverted path. To remove the effect of random noise in the mask computing process, we averaged all the masks generated with different hyperparameter combinations. From the comparisons, we can obtain an ice-cream similar to the source image with the unbiased mask.

### C. Ablation Study

In this section, we aim to first explore the effectiveness of the BE-Cycle and the improved mask computing pipeline, and then further demonstrate the existence of contextual prior bias. To this end, we conducted additional experiments based on the biased mask generated by contrasting $\mathbf{x}_0^{\text{source}}$ and $\mathbf{x}_0^{\text{target}}$ with the same process to calculate the unbiased mask. Additionally, we exhibited the results without the CLIP encoder. Here, we used SD-v1-4 as the pre-trained model. The quantitative results are reported in Table II. From the results, we can observe that all the models obtained similar $\mathcal{S}_{\textbf{CLIP}}$ and $\mathcal{S}_{\textbf{D-CLIP}}$ scores. However, the models, using the CLIP encoder, achieved better image content preservation, which demonstrates the effectiveness of the proposed mask computing pipeline. Then, we discuss the effects of the mask by comparing the models

using different masks. As seen from the results with the CLIP encoder, the model using the unbiased mask produced from the inverted path achieved a significant improvement in the **SSIM** and **PSNR** scores. This means that the inverted path increases the accuracy of identifying the regions that need to be edited. In other words, we can achieve better performance with fewer modifications. The reason why using the biased mask resulted in worse performance is the existence of contextual prior bias, which causes unexpected modifications. In short, the BE-Cycle is effective in eliminating the contextual prior bias derived from the pre-trained text-to-image models.

Additionally, we also depict the visual comparisons of the masks and outputs of the BE-Cycle's paths in Fig 3. The final edited images using different masks are also provided. Here, we averaged all the masks generated with different hyperparameter combinations to remove the noise's effect (*i.e.*, the boy's open mouth) in the mask computing process. From the results of the forward path, it can be observed that using a pre-trained model changed some image contents, such as the ice-cream's shape. This change is also reflected in the mask. Therefore, using the biased mask, the model shortened the ice-cream's length in the final edited image. In contrast, the unbiased mask produced from the inverted path can focus on the regions of the Christmas hat, which leads to synthesizing an ice-cream as similar to the source image as possible. In a word, the proposed method can achieve the goal of editing the image without the contextual prior bias's influence.

## IV. CONCLUSION

In this paper, we address the issue of contextual prior bias in semantic image editing and propose Dual-Cycle Diffusion to eliminate its effects by generating an unbiased mask to guide image editing. Our method employs SC-Cycle and BE-Cycle to eliminate distributional shifts between source and generated images from a pre-trained model. Leveraging CLIP's encoder to extract image visual contents, our model could produce an unbiased mask by identifying content differences in each region. This allows our model to focus on the regions that need editing without the effects of contextual prior bias. Through extensive experiments and ablation studies, we validate the superiority of Dual-Cycle Diffusion over other diffusion-based semantic image editing methods.

## REFERENCES

[1] H. Yan, H. Zhang, J. Shi, and J. Ma, "Texture brush for fashion inspiration transfer: A generative adversarial network with heatmap-guided semantic disentanglement," *IEEE TCSVT*, 2022.

[2] R. Mokady, A. Hertz, and A. H. Bermano, "Clipcap: Clip prefix for image captioning," *arXiv preprint arXiv:2111.09734*, 2021.

[3] S. Gu, J. Bao, H. Yang, D. Chen, F. Wen, and L. Yuan, "Mask-guided portrait editing with conditional gans," in *Proceedings of CVPR*, 2019, pp. 3436–3445.

[4] Z. Xu, T. Lin, H. Tang, F. Li, D. He, N. Sebe, R. Timofte, L. Van Gool, and E. Ding, "Predict, prevent, and evaluate: Disentangled text-driven image manipulation empowered by pre-trained vision-language model," in *Proceedings of CVPR*, 2022, pp. 18 229–18 238.

[5] B. Li, X. Qi, T. Lukasiewicz, and P. H. Torr, "Manigan: Text-guided image manipulation," in *Proceedings of CVPR*, 2020, pp. 7880–7889.

[6] T. Zhang, H.-Y. Tseng, L. Jiang, W. Yang, H. Lee, and I. Essa, "Text as neural operator: Image manipulation by text instruction," in *Proceedings of the 29th ACM MM*, 2021, pp. 1893–1902.

[7] C. H. Wu and F. De la Torre, "Unifying diffusion models' latent space, with applications to cyclediffusion and guidance," *arXiv preprint arXiv:2210.05559*, 2022.

[8] J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," *NeurIPS*, vol. 33, pp. 6840–6851, 2020.

[9] J. Song, C. Meng, and S. Ermon, "Denoising diffusion implicit models," in *ICLR*, 2021.

[10] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, "Learning transferable visual models from natural language supervision," in *ICML*. PMLR, 2021, pp. 8748–8763.

[11] Y. Balaji, S. Nah, X. Huang, A. Vahdat, J. Song, K. Kreis, M. Aittala, T. Aila, S. Laine, B. Catanzaro *et al.*, "ediffi: Text-to-image diffusion models with an ensemble of expert denoisers," *arXiv preprint arXiv:2211.01324*, 2022.

[12] O. Patashnik, Z. Wu, E. Shechtman, D. Cohen-Or, and D. Lischinski, "Styleclip: Text-driven manipulation of stylegan imagery," in *Proceedings of ICCV*, 2021, pp. 2085–2094.

[13] C. Meng, Y. He, Y. Song, J. Song, J. Wu, J.-Y. Zhu, and S. Ermon, "Sdedit: Guided image synthesis and editing with stochastic differential equations," in *ICLR*, 2021.

[14] X. Su, J. Song, C. Meng, and S. Ermon, "Dual diffusion implicit bridges for image-to-image translation," *arXiv preprint arXiv:2203.08382*, 2022.

[15] C. Schuhmann, R. Kaczmarczyk, A. Komatsuzaki, A. Katta, R. Vencu, R. Beaumont, J. Jitsev, T. Coombes, and C. Mullis, "Laion-400m: Open dataset of clip-filtered 400 million image-text pairs," in *NeurIPS Workshop Datacentric AI*, no. FZJ-2022-00923. Jülich Supercomputing Center, 2021.

[16] C. Schuhmann, R. Beaumont, R. Vencu, C. W. Gordon, R. Wightman, M. Cherti, T. Coombes, A. Katta, C. Mullis, M. Wortsman *et al.*, "Laion-5b: An open large-scale dataset for training next generation image-text models," in *Thirty-sixth Conference on NeurIPS Datasets and Benchmarks Track*, 2022.

[17] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of CVPR*, 2022, pp. 10 684–10 695.

[18] J. Zhang, P. Yang, W. Wang, Y. Hong, and L. Zhang, "Image editing via segmentation guided self-attention network," *SPL*, vol. 27, pp. 1605–1609, 2020.

[19] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, "Zero-shot text-to-image generation," in *ICML*. PMLR, 2021, pp. 8821–8831.

[20] O. Gafni, A. Polyak, O. Ashual, S. Sheynin, D. Parikh, and Y. Taigman, "Make-a-scene: Scene-based text-to-image generation with human priors," in *ECCV*, 2022.

[21] C. Saharia, W. Chan, S. Saxena, L. Li, J. Whang, E. Denton, S. K. S. Ghasemipour, B. K. Ayan, S. S. Mahdavi, R. G. Lopes *et al.*, "Photorealistic text-to-image diffusion models with deep language understanding," *arXiv preprint arXiv:2205.11487*, 2022.

[22] A. Ramesh, P. Dhariwal, A. Nichol, C. Chu, and M. Chen, "Hierarchical text-conditional image generation with clip latents," *arXiv preprint arXiv:2204.06125*, 2022.

[23] G. Couairon, J. Verbeek, H. Schwenk, and M. Cord, "Diffedit: Diffusion-based semantic image editing with mask guidance," in *ICLR*, 2023.

[24] A. Hertz, R. Mokady, J. Tenenbaum, K. Aberman, Y. Pritch, and D. Cohen-Or, "Prompt-to-prompt image editing with cross attention control," *arXiv preprint arXiv:2208.01626*, 2022.

[25] J. Yu, Y. Xu, J. Y. Koh, T. Luong, G. Baid, Z. Wang, V. Vasudevan, A. Ku, Y. Yang, B. K. Ayan *et al.*, "Scaling autoregressive models for content-rich text-to-image generation," *TMLR*, 2022.

[26] Y. Song and S. Ermon, "Generative modeling by estimating gradients of the data distribution," *NeurIPS*, vol. 32, 2019.

[27] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," in *ICLR*, 2020.

[28] D. Watson, W. Chan, J. Ho, and M. Norouzi, "Learning fast samplers for diffusion models by differentiating through sample quality," in *ICLR*, 2021.

[29] C. Lu, Y. Zhou, F. Bao, J. Chen, C. Li, and J. Zhu, "Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps," in *NeurIPS*, 2022.

[30] Q. Zhang and Y. Chen, "Fast sampling of diffusion models with exponential integrator," *arXiv preprint arXiv:2204.13902*, 2022.

[31] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," in *NeurIPS*, 2022.

[32] M. Chen and Z. Liu, "Edbgan: Image inpainting via an edge-aware dual branch generative adversarial network," *SPL*, vol. 28, pp. 842–846, 2021.

[33] S. S. Phutke and S. Murala, "Fasnet: Feature aggregation and sharing network for image inpainting," *SPL*, vol. 29, pp. 1664–1668, 2022.

[34] C. Han and J. Wang, "Face image inpainting with evolutionary generators," *SPL*, vol. 28, pp. 190–193, 2021.

[35] S. S. Phutke and S. Murala, "Diverse receptive field based adversarial concurrent encoder network for image inpainting," *SPL*, vol. 28, pp. 1873–1877, 2021.

[36] X. Xu, Y.-C. Chen, X. Tao, and J. Jia, "Text-guided human image manipulation via image-text shared space," *IEEE TPAMI*, 2021.

[37] H. Zheng, Z. Lin, J. Lu, S. Cohen, J. Zhang, N. Xu, and J. Luo, "Semantic layout manipulation with high-resolution sparse attention," *IEEE TPAMI*, 2022.

[38] G. Luo, Y. Zhu, Z. Weng, and Z. Li, "A disocclusion inpainting framework for depth-based view synthesis," *IEEE TPAMI*, vol. 42, no. 6, pp. 1289–1302, 2019.

[39] S. Xu, D. Liu, and Z. Xiong, "E2i: Generative inpainting from edge to image," *IEEE TCSVT*, vol. 31, no. 4, pp. 1308–1322, 2020.

[40] P. Zhang, L. Yang, X. Xie, and J. Lai, "Lightweight texture correlation network for pose guided person image generation," *IEEE TCSVT*, 2021.

[41] L. Zhang, H. Yang, T. Qiu, and L. Li, "Ap-gan: Improving attribute preservation in video face swapping," *IEEE TCSVT*, vol. 32, no. 4, pp. 2226–2237, 2021.

[42] Y. Liu, Q. Li, Q. Deng, and Z. Sun, "Towards spatially disentangled manipulation of face images with pre-trained stylegans," *IEEE TCSVT*, 2022.

[43] O. Avrahami, D. Lischinski, and O. Fried, "Blended diffusion for text-driven editing of natural images," in *Proceedings of CVPR*, 2022, pp. 18 208–18 218.