

## ZONE: Zero-Shot Instruction-Guided Local Editing

Shanglin Li<sup>1\*</sup>, Bohan Zeng<sup>1\*</sup>, Yutang Feng<sup>1\*</sup>, Sicheng Gao<sup>1</sup>, Xiuhui Liu<sup>1</sup>, Jiaming Liu<sup>2</sup>,  
 Lin Li<sup>2</sup>, Xu Tang<sup>2</sup>, Yao Hu<sup>2</sup>, Jianzhuang Liu<sup>4</sup>, Baochang Zhang<sup>1,3,5†</sup>

<sup>1</sup>Beihang University <sup>2</sup>Xiaohongshu Inc <sup>3</sup>Nanchang Institute of Technology, China  
<sup>4</sup>Shenzhen Institute of Advanced Technology, China <sup>5</sup>Zhongguancun Laboratory, China

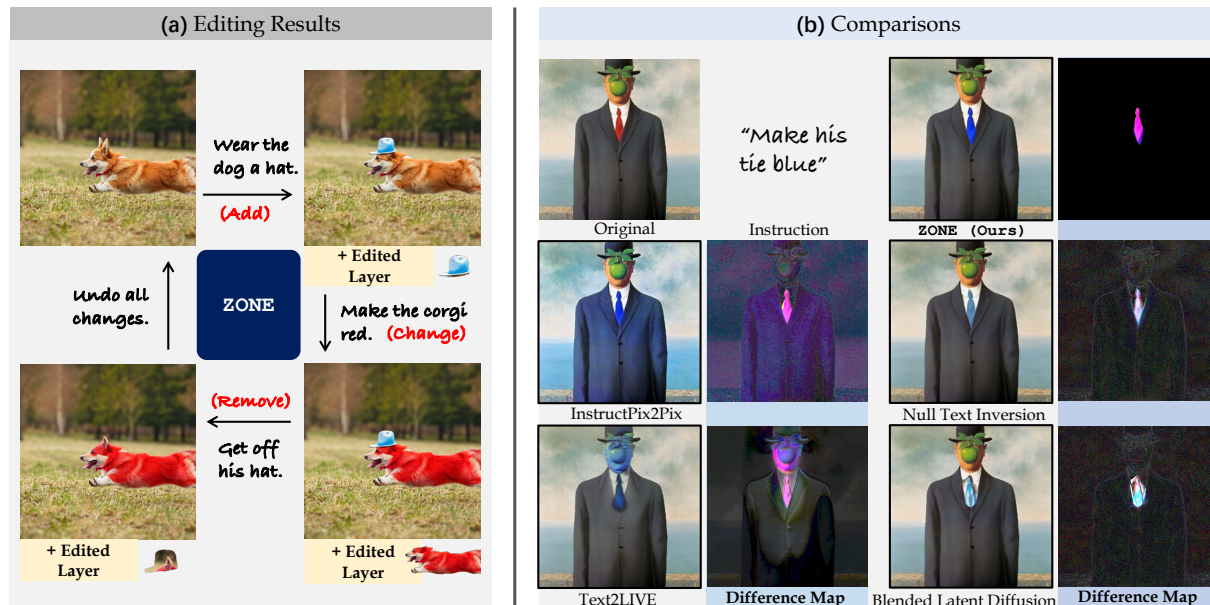


Figure 1. We propose ZONE, a zero-shot instruction-guided local editing approach. Our key idea is to edit and locate precise editing regions in an image with intuitive textual instructions. We demonstrate a multi-turn editing example in (a) and compare the difference maps between the edited image and the original image in (b) to highlight our method’s ability for local editing.

### Abstract

Recent advances in vision-language models like Stable Diffusion have shown remarkable power in creative image synthesis and editing. However, most existing text-to-image editing methods encounter two obstacles: First, the text prompt needs to be carefully crafted to achieve good results, which is not intuitive or user-friendly. Second, they are insensitive to local edits and can irreversibly affect non-edited regions, leaving obvious editing traces. To tackle these problems, we propose a Zero-shot instructiON-guided local image Editing approach, termed ZONE. We first convert the editing intent from the user-provided instruction (e.g., “make his tie blue”) into specific image editing regions through InstructPix2Pix. We then propose a Region-IoU

scheme for precise image layer extraction from an off-the-shelf segment model. We further develop an edge smoother based on FFT for seamless blending between the layer and the image. Our method allows for arbitrary manipulation of a specific region with a single instruction while preserving the rest. Extensive experiments demonstrate that our ZONE achieves remarkable local editing results and user-friendliness, outperforming state-of-the-art methods. Code is available at <https://github.com/lsl1001006/ZONE>.

### 1. Introduction

Large-scale vision-language models, such as Stable Diffusion [42], DALL-E 2 [41], and Imagen [45], have revolutionized text-guided image editing by bridging the gap between natural language and image content. Trained on vast

\*These authors contributed equally.

†Corresponding author: bczhang@buaa.edu.cn

visual and textual data, these methods harness generative power to manipulate appearance and style in natural images, offering a wide array of possibilities for enhancing and manipulating images in domains such as photography, advertising, and social media. These advancements have opened up new possibilities for text-guided image editing, making it increasingly important in various applications.

State-of-the-art (SOTA) image generative techniques [36, 41, 42, 53] predominantly concentrate on stylization, where the desired appearance is determined by a reference image or textual description, often leading to global image alterations [25, 29, 47]. However, these methods often lack straightforward local editing capabilities, and the precise localization of these edits typically needs additional input guidance, such as segmentation masks [1, 13, 32], making text-driven editing cumbersome and potentially limiting its scope. Recent description-guided works<sup>1</sup> like Prompt-to-Prompt [14], DiffEdit [8], and Text2LIVE [3] make noteworthy contributions to mask-free local edits, but they either require complex textual descriptions (*e.g.*, Prompt-to-Prompt requires word-to-word alignment between the source image caption and the edited image caption, and DiffEdit uses query and reference prompts) or need to specify the edited object (*e.g.*, Text2LIVE asks for multiple prompts), which are not user friendly. Instruction-guided editing methods<sup>2</sup> [5, 10, 55, 58] present more elegant characteristics in this regard. They eliminate the need for image-anchored descriptions, requiring only descriptions of the desired edits (*e.g.*, “make it snowy”), which facilitates concise and intuitive expression. However, these methods suffer from the over-edit problem, potentially distorting high-frequency details in non-edited regions (see Fig. 1 (b)).

To tackle these problems, we propose ZONE, a **Z**ero-shot **I**nstruction-guided **O**n-local image **E**ditin**G** approach. ZONE provides a more flexible and creative way to manipulate real images with layers.

Specifically, we leverage the pretrained instruction-guided model, InstructPix2Pix (IP2P) [5], for image editing. By exploring the attention mechanism of IP2P, we uncover the implicit associations between the editing locations and user-provided instructions in instruction-guided models. This allows us to identify the locations of the edited objects in instructions without the need for extra specification (*e.g.*, Stable Diffusion-based methods have to specify the tokens of the objects to edit). We further enhance this capability by proposing a Region-IoU scheme in conjunction with SAM [28], ensuring the mask refinement of the edited image layer. Our ZONE allows arbitrary image editing actions like “add”, “remove”, and “change”, all accomplished with intuitive instructions. Additionally, ZONE supports multi-turn local editing without affecting non-edited

regions, empowering high-fidelity local editing without any training or fine-tuning. Comprehensive experiments and user studies demonstrate that ZONE achieves remarkable results and user-friendliness in local image editing, outperforming existing SOTA methods.

To summarize, we make the following key contributions:

- We propose ZONE, a zero-shot image local editing method that enables users to edit localized regions of both real and synthetic images with simple instructions. ZONE preserves non-edited regions without loss and allows arbitrary manipulation of edited image layers.
- We reveal and exploit the different attention mechanisms between IP2P and Stable Diffusion when processing user instructions for image editing, with intuitive visual comparisons.
- We present a novel Region-IoU scheme and incorporate it with SAM for effective edited region refinement, and introduce a Fourier transform-based edge smoother to reduce the artifacts when compositing the image layers.
- Comprehensive experiments and user studies demonstrate that ZONE achieves high-fidelity local editing results without any auxiliary prompts, outperforming SOTA methods in photorealism and content preservation.

## 2. Related Work

### 2.1. Generative Models for Image Manipulation

Image manipulation is a fundamental process within the realm of computer vision, involving altering images with the aid of additional conditions like textual prompts, labels, masks, or reference images. Two mainstream editing methods include Generative Adversarial Networks (GANs) and Diffusion Models (DMs). Typical image manipulation tasks comprise image-to-image translation [7, 20, 26, 44, 47, 52, 59], super-resolution [12, 21, 30, 54], inpainting [19, 32, 39, 42], colorization [4, 31, 35, 51], and more. Although GAN-based methods excel when dealing with carefully curated data, they struggle with extensive and heterogeneous datasets [22, 23, 34]. To enhance generative expressiveness, [16, 17, 42, 48, 49] utilize DMs to achieve high-quality generation over diverse datasets. Recent research has yielded promising generation outcomes through the training or fine-tuning of large-scale text-to-image models [5, 18, 24, 33, 36, 41, 45, 53], as well as by harnessing CLIP [40] embeddings to guide image manipulation using textual prompts [9, 25, 29]. Some prior works [1, 2, 14, 38, 50] also demonstrate the zero-shot editing capability of pretrained DMs. Similarly, our method extensively exploits a pretrained DM’s generative capability to facilitate diverse and stylized image editing. However, we uniquely explore the implicit relationship between the DM’s editing regions during generation and the whole user instructions, enabling fine-grained layer-specific position-

<sup>1</sup>In this paper, we call them description-guided diffusion models.

<sup>2</sup>In this paper, we call them instruction-guided diffusion models.

ing.

## 2.2. Localized Image Editing

Several recent works have made attempts at localized image editing. Blend Diffusion [1] proposes a mask-guided method by blending edited regions with the other parts of the image at different noise levels along the diffusion process. Text2LIVE [3] introduces an RGBA layer generation approach with a CLIP-supervised generator for performing edits of objects in real images and videos. Prompt-to-Prompt [14] controls the spatial layouts of the image corresponding to the words in the prompt through cross-attention modification, enabling local edits by modifying textual prompts. Pix2Pix-Zero [38] preserves the structure of the original image with cross-attention guidance and applies an edit-direction embedding to make changes to localized objects. Instruction-based editing methods like IP2P [5] and MagicBrush [55] are trained or finetuned on triplet datasets to realize intuitive high-quality image editing based on user-provided instructions. PAIR-diffusion [13] allows editing the structure and appearance of each masked part in the original image independently. While these methods produce impressive results within their specific applications, they compromise on local image editing: instruction-guided methods [5, 55] and attention-based modifications [14, 38] introduce artifacts to non-edited regions, mask-based methods [1, 13] add complexity to user interactions, and CLIP-based methods [3, 38] sacrifice the flexibility of natural language editing. In contrast, our ZONE requires only a single instruction to achieve high-fidelity local image editing with an image layer.

## 2.3. Instruction-Guided Editing

Despite the significant progress of text-to-image models, most require detailed textual descriptions [36, 41–43, 45] to convey the desired image content, often falling short of user expectations for image editing. In contrast, direct instruction-guided modifications of target regions/attributes offer a more intuitive and convenient approach, such as “make the girl smile” and “give him a ball.” Recent advancements in instruction-guided editing and generation [5, 10, 37, 55, 57, 58] have made notable progress. For instance, IP2P [5] employs GPT-3 [6] and Prompt-to-Prompt [14] to synthesize an instruction-editing dataset, utilizes a pretrained Stable Diffusion model [42] for weight initialization, and trains a diffusion model specialized in instruction-guided editing. MagicBrush [55] fine-tunes IP2P using a real image dataset, thereby demonstrating a superior performance in instruction-guided editing. In this paper, we aim to leverage the instruction-editing capability of these pretrained instruction-guided diffusion models to eliminate the need for additional masks in previous local editing approaches [1, 2, 36], enabling flexible and high-fidelity local

editing based on a single user-provided instruction.

## 3. Preliminaries

**Diffusion Models.** Diffusion models [16, 46, 48] are probabilistic generative models founded on two complementary stochastic processes: *diffusion* and *denoising*. The *diffusion* process progressively adds different amounts of Gaussian noise to a clean image  $x_0$  towards Gaussian distribution  $x_T \sim \mathcal{N}(0, I)$  in  $T$  timesteps:  $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon$ , where  $\alpha_t$  defines the level of noise, and  $\epsilon \sim \mathcal{N}(0, I)$ .

In the *denoising* process, a neural network  $\epsilon_\theta$  is designed to predict the noise  $\epsilon$  for  $x_t$  to get a “cleaner” image gradually. This process is achieved by minimizing the denoising objective:  $\mathcal{L} = \mathbb{E}_{x_0, t, \epsilon} \|\epsilon - \epsilon_\theta(x_t, t)\|_2^2$ . Rombach et al. [42] introduce a latent diffusion model (LDM), which speeds up both processes by reducing images into a lower-dimensional latent space utilizing a variational auto-encoder [27]. This advancement has underpinned the achievements of Stable Diffusion, serving as the fundamental model for many diffusion-based works.

**InstructPix2Pix.** InstructPix2Pix [5] (IP2P) is a pioneering conditional diffusion model that edits images from user-provided instructions. Specifically, IP2P constructs an instruction dataset to fine-tune the pretrained Stable Diffusion. Given a target image  $x$ , an image condition  $c_I$ , and a textual instruction condition  $c_T$ , IP2P projects  $x$  to the latent  $z = \mathcal{E}(x)$  with a pretrained encoder  $\mathcal{E}$ , and then fine-tunes Stable Diffusion by minimizing the following objective:

$$\mathcal{L} = \mathbb{E}_{\mathcal{E}(x), \mathcal{E}(c_I), c_T, \epsilon \sim \mathcal{N}(0, 1), t} [\|\epsilon - \epsilon_\theta(z_t, t, \mathcal{E}(c_I), c_T)\|_2^2], \quad (1)$$

where the denoising network  $\epsilon_\theta$  accepts two input conditions and predicts the noise  $\epsilon$ . IP2P also finds it beneficial to perform classifier-free guidance [15] concerning both conditions, thus controlling the strength of edit by image guidance scale  $s_I$  and instruction guidance scale  $s_T$ :

$$\begin{aligned} \tilde{\epsilon}_\theta(z_t, c_I, c_T) = & \epsilon_\theta(z_t, \emptyset, \emptyset) \\ & + s_I \cdot (\epsilon_\theta(z_t, c_I, \emptyset) - \epsilon_\theta(z_t, \emptyset, \emptyset)) \\ & + s_T \cdot (\epsilon_\theta(z_t, c_I, c_T) - \epsilon_\theta(z_t, c_I, \emptyset)). \end{aligned} \quad (2)$$

At inference time, IP2P can modify an image with a user-provided instruction and trade-off the generated sample according to the strengths of the guidance image and the edit instruction through  $s_I$  and  $s_T$ .

## 4. Method

**Overview of ZONE.** We aim to make localized edits on an image with simple instructions. As depicted in Fig. 1 (a),

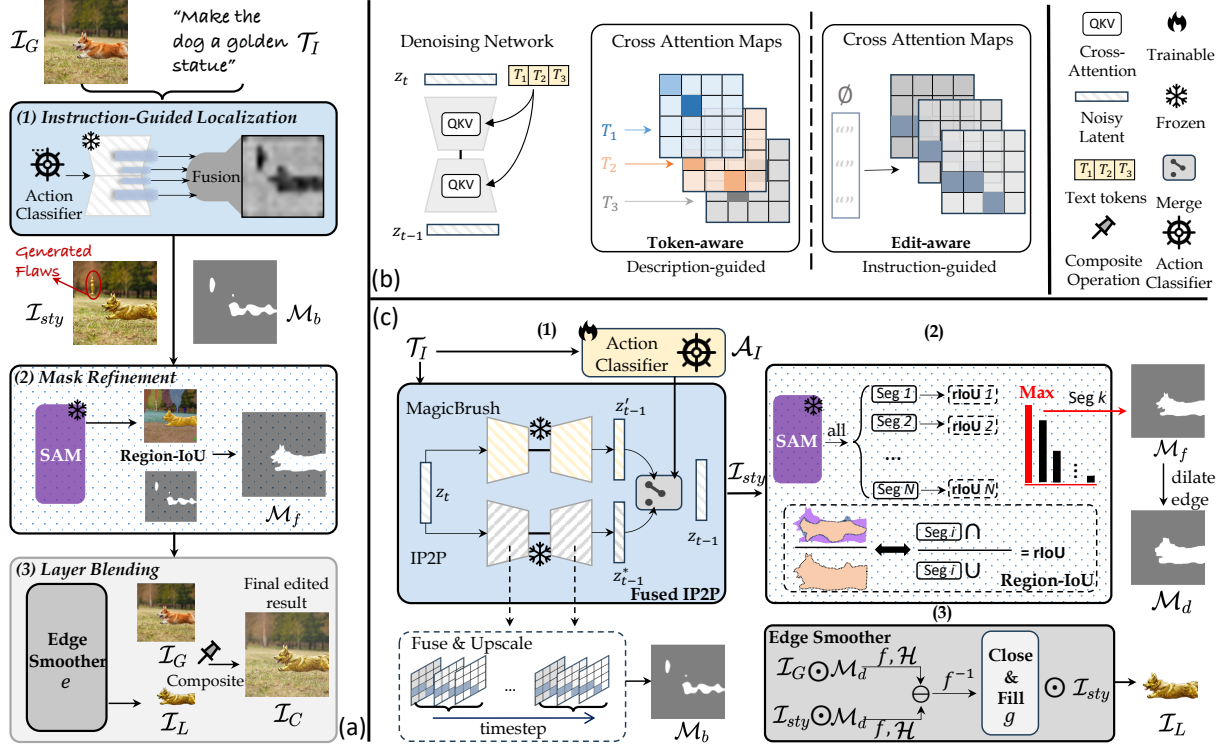


Figure 2. **Overview of ZONE.** (a) Three modules in ZONE. (b) The distinct difference between description-guided and instruction-guided diffusion models on cross-attention. The former usually follows a *token-aware* format, while the latter is *edit-aware*.  $\emptyset$  denotes the unconditional embeddings for null input. (c) Implementation details of the modules shown in (a).

such edits include performing three primary actions: (i) “add”: add an object to the image without specifying location with user-provided masks; (ii) “remove”: remove the object in the scene; (iii) “change”: change the style (*i.e.*, texture) of an existing object or replace the object with another object. Additionally, our method allows high-fidelity multi-turn edits with a series of instructions.

As outlined in Fig. 2 (a), our approach consists of the following steps: First, we train an action classifier for steering different editing requirements and concurrently generate and position the editing region using a fused IP2P, as detailed in Section 4.2 and Fig. 2 (c). Second, we devise a mask refinement module for an edited image layer in Section 4.3. Finally, in Section 4.4, we propose an FFT-based edge smoother for seamless blending of the edited image layer with the original image.

#### 4.1. Problem Statement

Given an RGB image  $\mathcal{I}_G \in \mathbb{R}^{3 \times H \times W}$  and a textual instruction  $\mathcal{T}_I$ , we aim to locate and edit image regions following  $\mathcal{T}_I$  and maintain the original non-edited regions. Inspired by Text2LIVE [3], we extract an edited layer  $\mathcal{I}_L$  with color and opacity that are composited over  $\mathcal{I}_G$ . As opposed to previous works [1, 3, 14, 38], we neither rely on any user-defined mask nor need non-intuitive prompt engineering, realizing

precise local editing and seamless layer blending.

#### 4.2. Instruction-Guided Localization

Many local editing methods require users to explicitly specify the object they want to edit with a prompt or a mask [1, 3, 8, 38]. This is not intuitive and often requires a certain learning cost. Our approach locates and edits the implicitly designated object from the user’s instruction. For example, a user-provided instruction like “make her old” can implicitly convey the user’s editing intent to modify the woman in the scene (*locate*) by making her appear older (*edit*).

As shown in Fig. 2 (b), our key finding is that the operational mechanisms of instruction-guided and description-guided diffusion models on cross-attention exhibit a distinct difference. Specifically, we empirically demonstrate that: (i) a description-guided model displays a *token-aware* characteristic on its cross-attention maps, associating each input text token with a corresponding spatial structure; (ii) an instruction-guided model’s cross-attention maps with unconditional embeddings share similar spatial features, demonstrating an *edit-aware* characteristic, being responsive to the overall editing intent.

Given a noisy latent  $z_t$  and a textual embedding  $c_T$ , the denoising UNet  $\epsilon_\theta$  predicts the noise  $\epsilon$  at each timestep  $t$ . The generation is conditioned on the textual prompt  $\mathcal{T}_I$

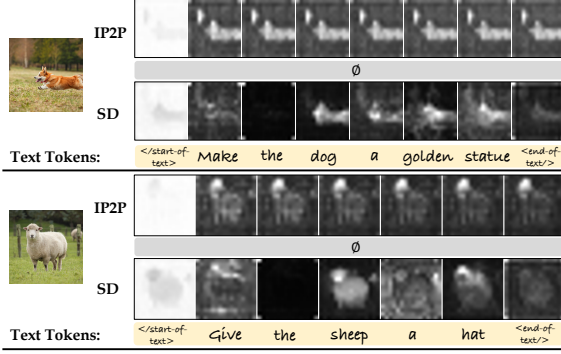


Figure 3. **Cross-attention map difference.** We average the cross-attention maps among all timesteps for each sample. IP2P shows consistency in the overall editing intent with unconditional embeddings  $\emptyset$ , while Stable Diffusion (SD) demonstrates a one-to-one correspondence with text tokens.

by computing cross-attention between the textual embedding  $c_T$  and the spatial features  $\phi(z_t)$ , and updates  $\phi(z_t)$  as  $\hat{\phi}(z_t)$ :

$$M = \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right), \hat{\phi}(z_t) = M \cdot V, \quad (3)$$

where the query  $Q = W_Q\phi(z_t)$ , key  $K = W_Kc_T$ , and value  $V = W_Vc_T$  are obtained with linear projections  $W_Q$ ,  $W_K$ , and  $W_V$ .  $M \in \mathbb{R}^{H' \times W' \times L}$  contains  $L$  cross-attention maps that are correlated to the similarity between  $Q$  and  $K$ . Typically  $H'$  and  $W'$  are  $1/32$  of the original image size  $H$  and  $W$  in Stable Diffusion. For the description-guided Stable Diffusion model, each token corresponds to its specific attention map  $M^l$  with text embeddings, where  $l \in \{1, 2, \dots, L\}$ . For the instruction-guided IP2P, we find its attention maps share a uniform characteristic with unconditional embeddings, concentrated directly at the edited location without token specification, as visualized in Fig. 3.

Based on this finding, we devise a simple yet effective localization module that semantically locates the edited region with instruction  $\mathcal{T}_I$ . Specifically, we first collect the attention maps of the denoising model of IP2P from all timesteps of the denoising process. Then, we average and resize the maps to obtain averaged attention maps  $\mathcal{M}_A \in \mathbb{R}^{H \times W \times L}$ . We observe that the first attention map  $\mathcal{M}_A^1$  primarily emphasizes the attention weights of the non-edited region. Subsequent attention maps  $\mathcal{M}_A^{2, \dots, L}$  shift attention towards the edited region. Therefore, we subtract the last cross-attention map from the first attention map and binarize each pixel  $(m, n)$  with a fixed threshold  $T$  to highlight the edited region and mitigate the background noise:

$$\mathcal{M}_b(m, n) = \begin{cases} 1, & \text{if } \mathcal{M}_A^1(m, n) - \mathcal{M}_A^L(m, n) < T, \\ 0, & \text{others,} \end{cases} \quad (4)$$

where  $T$  is empirically set to 128. This yields a rough, noise-filtered edited region mask  $\mathcal{M}_b$  most related to  $\mathcal{T}_I$  (see Fig. 2 (a)).

Moreover, we find that IP2P performs not as well as MagicBrush in the “remove” editing but preserves better object identity in terms of “add” and “change”. Therefore, we design a fused IP2P module with a trainable action classifier  $\mathcal{A}_I$ . As illustrated in Fig. 2 (c), we lock the weights of both IP2P and MagicBrush and use a pretrained action classifier  $\mathcal{A}_I$  to steer the denoising process based on  $\mathcal{T}_I$ :

$$z_{t-1} = (z_{t-1}^* + \beta \cdot z'_{t-1}) / (1 + \beta), \quad (5)$$

where  $z_{t-1}^*$  and  $z'_{t-1}$  are the denoised latents by IP2P and MagicBrush, respectively.  $\beta$  is a hyperparameter to control the guidance strength of MagicBrush on IP2P, empirically set to 0.2 if  $\mathcal{A}_I(\mathcal{T}_I)$  is classified to “remove” and 0.01 for other actions. This module generates a globally edited image  $\mathcal{I}_{sty}$  according to  $\mathcal{T}_I$ .  $\mathcal{I}_{sty}$  serves as the canvas, from which the edited region is cropped out to form a separate image layer in the following steps.

### 4.3. Mask Refinement

The location mask  $\mathcal{M}_b$  and  $\mathcal{I}_{sty}$  obtained in Section 4.2 are insufficient for precise local editing, since  $\mathcal{M}_b$  only indicates the general location of the edited region, as illustrated in Fig. 2 (a). An intuitive and effective mask refinement method is to use an off-the-shelf segmentation model. We leverage the Segment Anything Model (SAM) [28] to generate precise masks of the canvas  $\mathcal{I}_{sty}$  at various levels. However, we do not use SAM’s preset point or box prompts for segmentation selection, because these prompts could potentially lead to misselection or omission of SAM’s segmentation results due to IP2P’s over-edit problem (which is also reflected in  $\mathcal{M}_b$ , see  $\mathcal{I}_{sty}$  and  $\mathcal{M}_b$  in Fig. 2 (a)), resulting in a final mask that does not accurately reflect  $\mathcal{T}_I$ ’s editing intention. Therefore, we propose a Region-IoU (rIoU) scheme to obtain the accurate segmentation mask.

As depicted in Fig. 2 (c), by sending  $\mathcal{I}_{sty}$  to SAM, we extract all the possible instance segments  $\mathcal{S} = \{\mathcal{S}^j\}_{j=1}^N$ . Note that  $\mathcal{S}$  contains the segments from all levels of SAM’s segmentation. We define rIoU  $\mathcal{R}(j)$  as:

$$\mathcal{R}(j) = \frac{\text{area}(\mathcal{S}^j \cap \mathcal{M}_b)}{\text{area}(\mathcal{S}^j \cup \mathcal{M}_b)}, j = 1, 2, \dots, N. \quad (6)$$

If  $k = \arg \max_{j=1, 2, \dots, N} \{\mathcal{R}(j)\}$ , then we obtain the refined mask  $\mathcal{M}_f = \mathcal{S}^k$ . One example is shown in Fig. 2 (a) or (c).

### 4.4. Layer Blending

After the mask refinement, we obtain an edited image layer  $\mathcal{I}'_L = \mathcal{I}_{sty} \odot \mathcal{M}_f$ , which retains the color information of  $\mathcal{I}_{sty}$  within the region where  $\mathcal{M}_f = 1$ , with the rest being

Type	Methods	L1 ↓	L2 ↓	LPIPS ↓	CLIP-I ↑	CLIP-T ↑
Description-guided	DiffEdit [8]	<u>0.0426</u>	0.0099	<u>0.1695</u>	0.8947	0.2815
	Text2LIVE [3]	0.0511	<u>0.0075</u>	0.2176	0.9075	<b>0.3062</b>
	Pix2Pix-Zero [38]	0.1198	0.0342	0.4375	0.7679	0.2701
Instruction-guided	InstructPix2Pix [5]	0.0945	0.0274	0.2816	<u>0.9089</u>	0.2907
	MagicBrush [55]	0.0919	0.0378	0.2903	0.8959	0.2939
	ZONE (Ours)	<b>0.0146</b>	<b>0.0061</b>	<b>0.0441</b>	<b>0.9688</b>	<u>0.2969</u>

Table 1. **Quantitative evaluation.** We use L1 and L2 to gauge pixel-level structural similarity, LPIPS and CLIP-I to evaluate image quality, and CLIP-T to assess text-image semantic similarity. The best and the second best results are marked in **bold** and underline, respectively.



Figure 4. **Visualization and ablation.** The first 4 columns show the intermediate results related to the edge smoother. The last column compares the final edited results with and without the edge smoother.

transparent. A naïve way to get the final edited result  $\mathcal{I}_C$  is to stitch  $\mathcal{I}'_L$  and the original image  $\mathcal{I}_G$  at pixel-level. This fundamentally tackles the over-edit problem encountered in instruction-guided methods for local editing. Nevertheless, directly pasting  $\mathcal{I}'_L$  back to  $\mathcal{I}_G$  may result in noticeable artifacts, such as jagged edges and incomplete coverage of the edited region in the original image, as indicated by the yellow arrows in Fig. 4 (b).

We tackle this problem by designing a novel edge smoother with Fast Fourier Transform (FFT). Given the original image  $\mathcal{I}_G$ , the canvas  $\mathcal{I}_{sty}$ , and the refined location mask  $\mathcal{M}_f$ , we first dilate  $\mathcal{M}_f$  to  $\mathcal{M}_d$  to incorporate more edge information in  $\mathcal{I}_{sty}$  that may not be included in  $\mathcal{I}'_L$ . Then we get the dilated edited image layer  $\mathcal{I}_{L,d} = \mathcal{I}_{sty} \odot \mathcal{M}_d$  and the dilated original image layer  $\mathcal{I}_{G,d} = \mathcal{I}_G \odot \mathcal{M}_d$ , as shown in the second column of Fig. 4. The edge smoother  $e$  is defined by:

$$e(\mathcal{I}_{L,d}, \mathcal{I}_{G,d}) = g(f^{-1}(\mathcal{H}(f(\mathcal{I}_{L,d})) - \mathcal{H}(f(\mathcal{I}_{G,d})))), \quad (7)$$

where  $g$  is a composition of binarization and morphological closing and filling functions,  $f$  and  $f^{-1}$  represent FFT and inverse FFT, respectively, and  $\mathcal{H}$  is an ideal low-pass filter:

$$\mathcal{H}(f_s) = \begin{cases} f_s(c), & \text{if } \|c - c_0\|_2 \leq D_0, \\ 0, & \text{if } \|c - c_0\|_2 > D_0, \end{cases} \quad (8)$$

where  $f_s \in \mathbb{R}^{H \times W}$  is the frequency spectrum of the image

transformed by  $f$ ,  $c$  is the coordinate in  $f_s$ ,  $c_0$  is the center coordinate of  $f_s$ , and  $D_0$  is set empirically to 200 for a  $512 \times 512$  image. We use the edge smoother  $e$  to get the final mask  $\mathcal{M}_f^*$ .

As shown in the second column of Fig. 4, we observe that both  $\mathcal{I}_{G,d}$  and  $\mathcal{I}_{L,d}$  share similar low-frequency characteristics on non-edited regions (e.g., background), but they hold different low-frequency characteristics on the edited regions (e.g., hat and the shadow below it). Therefore, we can exclude the non-edited regions and retain the edited regions by subtracting the low-frequency components between  $\mathcal{I}_{L,d}$  and  $\mathcal{I}_{G,d}$  in the frequency domain:  $d_s = \mathcal{H}(f(\mathcal{I}_{L,d})) - \mathcal{H}(f(\mathcal{I}_{G,d}))$  and invert it back to the image domain to get the difference mask  $\mathcal{M}_{dm} = f^{-1}(d_s)$ . The final mask  $\mathcal{M}_f^*$  is then obtained by  $\mathcal{M}_f^* = g(\mathcal{M}_{dm}) = e(\mathcal{I}_{L,d}, \mathcal{I}_{G,d})$ . Finally, we get the final edited image layer  $\mathcal{I}_L$  by  $\mathcal{I}_L = \mathcal{I}_{sty} \odot \mathcal{M}_f^*$ , and the final edited result  $\mathcal{I}_C$  is acquired by compositing  $\mathcal{I}_G$  and  $\mathcal{I}_L$ . The intuitive visualization of these intermediate results are shown in Fig. 4.

The implementation details and more discussions can be found in the supplementary material.

## 5. Experiments

### 5.1. Experimental Setup

**Baselines.** We conduct comprehensive experiments for the local editing task by comparing ZONE with five state-of-the-art image editing methods that are capable of local editing: *Text2LIVE* [3], *DiffEdit* [8], *IP2P* [5], *Pix2Pix-Zero* [38], and *MagicBrush* [55]. The implementation of these methods can be found in the supplementary material.

**Datasets.** We randomly select and annotate 100 samples for evaluation, including 60 real images from the Internet and 40 synthetic images. To ensure the representativeness of the evaluation, we consider the diversity of scenes and objects in the sample selection. In particular, we divide the test set into three categories: 32 images for “add”, 54 for “change”, and 14 for “remove” actions. All these 100 images are listed in the supplementary material.

**Evaluation Metrics.** Following [5, 55], we perform qualitative and quantitative comparisons using a variety of eval-

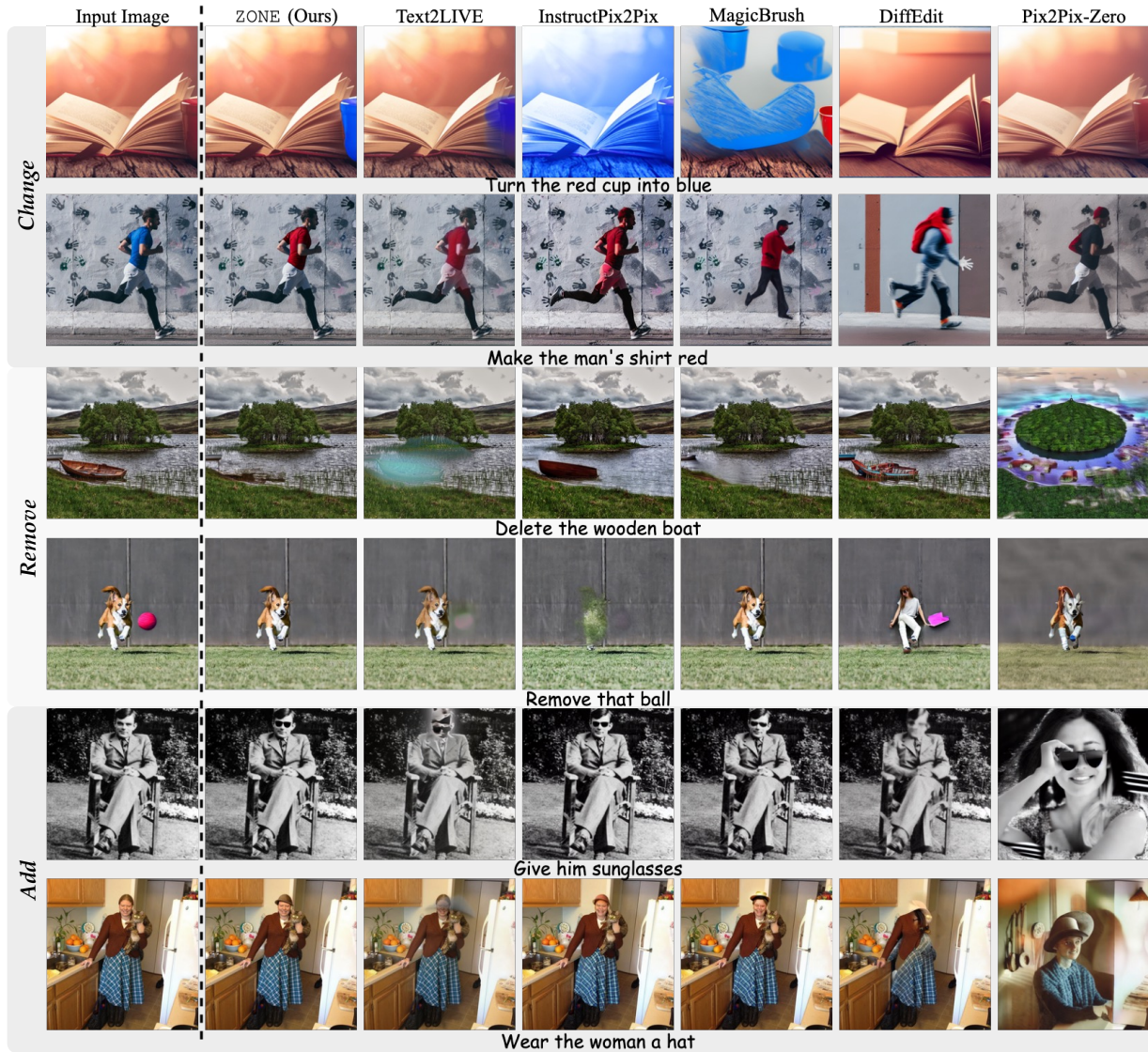


Figure 5. **Qualitative comparison.** We compare the editing efficacy of our ZONE with existing SOTA methods. The instructions (or instructions that are equivalent to the descriptions required by some baselines) used for editing are written below each row of the images.

uation metrics. Learned Perceptual Image Patch Similarity (LPIPS) [56] is used to quantify the perceptual similarity between the original and edited image. CLIP text-image similarity (CLIP-T) [11] is employed to assess the alignment between the edited image and its corresponding caption, and CLIP image similarity (CLIP-I) is used to evaluate the layout similarity and semantic correlation between the edited image and the original image, serving as a reliable indicator of the edited image’s quality. We also use L1 and L2 distances for pixel-level difference comparison.

## 5.2. Comparisons

**Quantitative Evaluation.** As shown in Table 1, we measure the models with the five metrics. The quantitative re-

sults indicate the following: (i) Our method significantly outperforms our counterparts on metrics related to image structure and quality, implying the efficacy of ZONE’s preservation of the non-edited regions. (ii) Text2LIVE performs best on CLIP-T, but the qualitative comparison in Fig. 5 does not support this result. We surmise that Text2LIVE performs better on this metric potentially due to its direct supervision by CLIP.

To quantify the stability of the edits, we divide the test set into three action groups: “change”, “add”, and “remove”. We then test the CLIP-I and CLIP-T metrics for each model and plot the CLIP curves in Fig. 6, where the performances of the same method on these actions are connected with lines of the same color. Our interpretation is as follows:

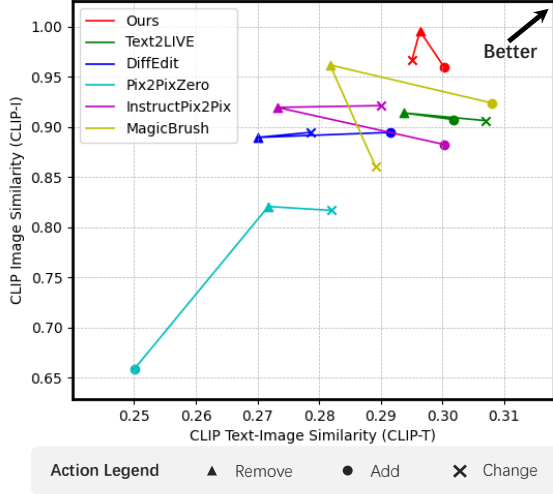


Figure 6. **Stability analysis.** We categorize the test set into three actions (“Remove”, “Add”, and “Change”) and calculate their respective CLIP-I and CLIP-T values. Our method achieves the best quality-stability trade-off for all actions.



Figure 7. **Detailed comparison.** We show a zoomed-in sample where ZONE effectively resolves the over-edit problem.

first, the shorter the projection of the line on the axis, the higher the semantic *stability* (i.e., maintaining similar performances under different editing instructions) of the image editing; second, if the curve is closer to the upper right corner, it indicates that the method’s editing *quality* is more superior. Our method achieves the best trade-off between *quality* and *stability*, demonstrating strong editing stability and representativeness.

**Qualitative Comparison.** In Fig. 5, we illustrate the editing results for the baselines and our method. We select six sets of images (including synthetic and real images) and group them based on actions. Our ZONE shows precise local editing capability while preserving the remaining pixels, this is especially important when there are perceptually important high-frequency details, such as faces, textures, or texts. A zoomed-in comparison is shown in Fig. 7. Both InstructPix2Pix and Text2LIVE introduce distortions to the non-edited areas during the editing process. For instance, InstructPix2Pix distorts the nearby clock and paints the orange outside of the basket red. In comparison, Text2LIVE maintains a better structure but generates a “barrel” of apples and introduces an obvious foggy effect to the image.

Methods	SR (%)	UPR (%)
DiffEdit [8]	27.1 ± 2.7	8.8
Text2LIVE [3]	33.0 ± 3.2	17.3
Pix2Pix-Zero [38]	19.2 ± 3.7	10.4
InstructPix2Pix [5]	59.8 ± 3.1	18.9
MagicBrush [55]	50.2 ± 2.9	18.0
ZONE (Ours)	<b>69.4 ± 3.5</b>	<b>26.6</b>

Table 2. **Human evaluation.** Our ZONE obtains the highest success rate (SR) and user preference rate (UPR).

Our method, however, can clearly distinguish between the edited region and the non-edited regions, demonstrating the best local editing efficacy.

### 5.3. Human Evaluation

Due to the lack of an effective metric to measure editing effects (mainly due to the absence of ground truth images after editing), the metrics mentioned in Section 5.1 alone are not sufficient to demonstrate the superiority of our method over existing ones. To further validate the editing effects of ZONE, in addition to the visual comparison in Fig. 5, we also conduct a human evaluation to calculate the success rate (SR) and user preference rate (UPR) of the edited images with the editing instructions. Table 2 shows a consistent preference for our method by users, as well as a dominant success rate over other methods.

Please refer to our supplementary material for more visualizations and details of this user study.

## 6. Conclusion

We present ZONE, a zero-shot instruction-guided local image editing approach, which leverages the localization capability within the pre-trained instruction-guided diffusion models. Our approach innovatively utilizes the editing intent regions inherent in the instructions, rather than focusing on individual tokens, eliminating the need for specific guidance. By integrating the Region-IoU scheme and FFT-based edge smoother with a pretrained segmentation model, ZONE effectively realizes precise local editing. Comprehensive experiments and user studies further demonstrate the superiority of ZONE over SOTA methods.

**Acknowledgements.** The work was supported by the National Key Research and Development Program of China (2023YFC3300029), Zhejiang Provincial Natural Science Foundation of China (LD24F020007), Beijing Natural Science Foundation (L223024), National Natural Science Foundation of China (62076016), “One Thousand Plan” projects in Jiangxi Province (Jxsg2023102268), Beijing Municipal Science & Technology Commission, Administrative Commission of Zhongguancun Science Park (Z231100005923035).



## References

- [1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *CVPR*, 2022. 2, 3, 4
- [2] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *TOG*, 2023. 2, 3
- [3] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *ECCV*, 2022. 2, 3, 4, 6, 8
- [4] Marc Górriz Blanch, Marta Mrak, Alan F Smeaton, and Noel E O'Connor. End-to-end conditional gan-based architectures for image colourisation. In *MMSPW*, 2019. 2
- [5] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *CVPR*, 2023. 2, 3, 6, 8
- [6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 3
- [7] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *CVPR*, 2018. 2
- [8] 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, 4, 6, 8
- [9] Katherine Crowson, Stella Biderman, Daniel Kornis, Dashiell Stander, Eric Hallahan, Louis Castricato, and Edward Raff. Vqgan-clip: Open domain image generation and editing with natural language guidance. In *ECCV*, 2022. 2
- [10] Alaaeldin El-Nouby, Shikhar Sharma, Hannes Schulz, Devon Hjelm, Layla El Asri, Samira Ebrahimi Kahou, Yoshua Bengio, and Graham W Taylor. Tell, draw, and repeat: Generating and modifying images based on continual linguistic instruction. In *CVPR*, 2019. 2, 3
- [11] Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators. *TOG*, 2022. 7
- [12] Sicheng Gao, Xuhui Liu, Bohan Zeng, Sheng Xu, Yanjing Li, Xiaoyan Luo, Jianzhuang Liu, Xiantong Zhen, and Baochang Zhang. Implicit diffusion models for continuous super-resolution. In *CVPR*, 2023. 2
- [13] Vidit Goel, Elia Peruzzo, Yifan Jiang, DeJia Xu, Nicu Sebe, Trevor Darrell, Zhangyang Wang, and Humphrey Shi. Pair-diffusion: Object-level image editing with structure-and-appearance paired diffusion models. *arXiv preprint arXiv:2303.17546*, 2023. 2, 3
- [14] 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, 3, 4
- [15] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3
- [16] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 2, 3
- [17] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *JMLR*, 2022. 2
- [18] Yi Huang, Jiancheng Huang, Yifan Liu, Mingfu Yan, Jiayi Lv, Jianzhuang Liu, Wei Xiong, He Zhang, Shifeng Chen, and Liangliang Cao. Diffusion model-based image editing: A survey. *arXiv preprint arXiv:2402.17525*, 2024. 2
- [19] Satoshi Iizuka, Edgar Simo-Serra, and Hiroshi Ishikawa. Globally and locally consistent image completion. *TOG*, 2017. 2
- [20] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, 2017. 2
- [21] Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung Park. Scaling up gans for text-to-image synthesis. In *CVPR*, 2023. 2
- [22] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *CVPR*, 2019. 2
- [23] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *CVPR*, 2020. 2
- [24] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *CVPR*, 2023. 2
- [25] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *CVPR*, 2022. 2
- [26] Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, and Jiwon Kim. Learning to discover cross-domain relations with generative adversarial networks. In *ICML*, 2017. 2
- [27] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 3
- [28] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 2, 5
- [29] Gihyun Kwon and Jong Chul Ye. Clipstyler: Image style transfer with a single text condition. In *CVPR*, 2022. 2
- [30] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 2
- [31] Jianxin Lin, Peng Xiao, Yijun Wang, Rongju Zhang, and Xi-angxiang Zeng. Diffcolor: Toward high fidelity text-guided image colorization with diffusion models. *arXiv preprint arXiv:2308.01655*, 2023. 2
- [32] 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
- [33] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jianjun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 2

- [34] Ron Mokady, Omer Tov, Michal Yarom, Oran Lang, Inbar Mosseri, Tali Dekel, Daniel Cohen-Or, and Michal Irani. Self-distilled stylegan: Towards generation from internet photos. In *SIGGRAPH*, 2022. 2
- [35] Kamyar Nazeri, Eric Ng, and Mehran Ebrahimi. Image colorization using generative adversarial networks. In *AMDO*, 2018. 2
- [36] Alexander Quinn 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. In *ICML*, 2022. 2, 3
- [37] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 3
- [38] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *SIGGRAPH*, 2023. 2, 3, 4, 6, 8
- [39] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context encoders: Feature learning by inpainting. In *CVPR*, 2016. 2
- [40] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 2
- [41] 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. 1, 2, 3
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1, 2, 3
- [43] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, 2023. 3
- [44] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *SIGGRAPH*, 2022. 2
- [45] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 1, 2, 3
- [46] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 3
- [47] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, et al. Styledrop: Text-to-image generation in any style. *arXiv preprint arXiv:2306.00983*, 2023. 2
- [48] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 3
- [49] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 2
- [50] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *CVPR*, 2023. 2
- [51] Hanzhang Wang, Deming Zhai, Xianming Liu, Junjun Jiang, and Wen Gao. Unsupervised deep exemplar colorization via pyramid dual non-local attention. *TIP*, 2023. 2
- [52] Tengfei Wang, Ting Zhang, Bo Zhang, Hao Ouyang, Dong Chen, Qifeng Chen, and Fang Wen. Pretraining is all you need for image-to-image translation. *arXiv preprint arXiv:2205.12952*, 2022. 2
- [53] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gungjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, et al. Scaling autoregressive models for content-rich text-to-image generation. *arXiv preprint arXiv:2206.10789*, 2022. 2
- [54] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. Resshift: Efficient diffusion model for image super-resolution by residual shifting. *arXiv preprint arXiv:2307.12348*, 2023. 2
- [55] Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing. *arXiv preprint arXiv:2306.10012*, 2023. 2, 3, 6, 8
- [56] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *CVPR*, 2018. 7
- [57] Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, et al. Hive: Harnessing human feedback for instructional visual editing. *arXiv preprint arXiv:2303.09618*, 2023. 3
- [58] Tianhao Zhang, Hung-Yu Tseng, Lu Jiang, Weilong Yang, Honglak Lee, and Irfan Essa. Text as neural operator: Image manipulation by text instruction. In *ACMMM*, 2021. 2, 3
- [59] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *ICCV*, 2017. 2