

# Image-Editing Specialists: A Multi-Objective Approach for Diffusion Models

Elior Benarous  
Harvard University  
ETH Zürich  
[ebenarous@ethz.ch](mailto:ebenarous@ethz.ch)

Yilun Du  
Harvard University  
Google DeepMind

Heng Yang  
Harvard University



## Abstract

[Pending abstract] We present a method for training image-editing diffusion models based on textual instructions and visual prompts to capture subtle user preferences without the need for long descriptive prompts. It achieves so while requiring neither a large and curated training data or human feedback. Our method largely improves the alignment with instructions in two ways. First, it better localizes the region to edit, largely reducing the modification of regions in the image that are unrelated to the instruction. Second, it captures fine nuances in the desired edit by leveraging a visual prompt, which considerably simplifies the user’s efforts to attain a highly specific edit. We show that our model can perform intricate edits in complex image like busy road scenes, and can also perform sim2real transfer that greatly helps applications where real data is scarce.

## 1. Introduction

Text-to-image (T2I) generative models have achieved remarkable success in creating visually appealing images based on text prompts [25, 44, 47], largely due to advancements in aligning captions with images [1, 42].

Beyond the remarkable generative capabilities of T2I models, instructional image editing has become one of the most practical approach to semantic modification [11, 54]. Unlike traditional image editing [12, 22, 28, 33, 38, 63], which necessitates descriptive captions for both the input and modified images, instruction-based image editing re-

quires only natural-language instructions. This method is more straightforward, as it requires specifying only the elements that need alteration in the original image, without involving the remaining unrelated attributes.

Although instructional editing has gained popularity in consumer applications, it has yet to show effective reusability of its samples as data augmentation in downstream model training. Simple-to-generate synthetic samples have shown to serve as pretraining data to teach early-training features [7, 8, 27]. Recently, some studies focused on leveraging T2I models as sources of high quality data augmentation. Some works generate diverse samples by applying naive perturbation to the noisy latents [62, 78], the conditioning embeddings [80] or throughout the denoising process [19]. Others do so by designing, optimizing, or reformulating prompts [6, 16, 51, 56, 57, 74]. Further, Azizi et al. [3] finetune a T2I model to reduce the gap between real data from a specific domain and class-conditional generated samples. Still, none leverage instructional editing for data augmentations that are closely aligned with the user’s intentions. In our work, we focus on a method to train instruction-based diffusion models to generate samples that are of high precision and follow instructions with high fidelity. We show that such specialized models can serve as highly useful training data generators, particularly in data-scarce domains, where the collection might be costly or timely.

To become truly useful for downstream applications, we identify two key criteria that an image-editing model should satisfy.

**Structural Alignment:** It should ensure high-fidelity

preservation of content in the initial image that is unrelated to the editing request. Current state-of-the-art image-editing models without segmentation or semantic masks in the forward pass struggle to perform precise edits on the subject specified by the text prompt, while leaving unrelated areas of the image intact. Notably, the training data for InstructPix2Pix was created with the Prompt-to-Prompt method [22]. Hence, it suffers from the same limitations as Prompt-to-Prompt, including imperfect preservation of background pixels.

**Semantic Alignment:** It should enable fine-grained control over the visual aspects of desired modifications. One effective approach is to utilize both textual and image prompts to express nuanced stylistic preferences. This strategy captures subtleties that are difficult to articulate with text alone and may elude language models. While a text prompt can convey the broad theme of an edit, specifying the precise nuances of its interpretation proves more challenging and impractical. Consequently, it is preferable to enhance user alignment with their envisioned edits without resorting to lengthy prompts, which contradicts the essence of instruction-based systems. Most works on visual prompting for image generation have been targeted at style transfer, where an image’s style is edited across the whole frame [23, 58, 65]. Other recent studies have focused on finetuning pre-trained T2I models for subject-driven editing from a set of reference images [21, 49]. However their methods require unique identifiers to capture the edit nuance. In our work, we bring our attention to localized edits conditioned on a visual prompt, and limit the complexity of our text prompts to simple instructions. We highlight the relevance of visual prompts, which further alleviate the user’s burden to convey the desired edit when combined with instruction-styled text prompts.

The standard diffusion model training objective presents inherent challenges for image editing because it enforces uniformly reconstructing the entire image, without distinguishing between regions that require preservation and those that need modification. Achieving this balance with a standard reconstruction loss necessitates an oracle to generate precise input-output pairs. However, such oracles are often biased, lack scalability, or produce low-quality samples that require further pruning using various filtering techniques [11, 13, 54, 75]. To address these limitations, we frame the problem as an alignment issue, leveraging the capabilities of the InstructPix2Pix model [11]. Although this model can generate edits, it requires refined training to better align with human semantic and structural preferences, especially in complex scene edits. Our reinforcement learning (RL) framework enhances this alignment and circumvents the need for an oracle, as the method is based on self-play. By achieving closer alignment with human preferences, the model is more likely to generate natural-

looking edits with better visual appeal, making it a suitable method for high-quality data augmentation.

Accordingly, we introduce a novel RL approach designed to specialize image editing models into producing edits with styles highly aligned with visual prompts while accurately preserving the original structure in non-pertinent areas of the image. Unlike traditional reinforcement learning with human feedback (RLHF) methods, our approach bypasses the need for human annotators by employing AI models to provide feedback that accurately simulates human preferences in terms of structure and semantics. We demonstrate our method’s effectiveness in editing complex scenes, such as outdoor city street environments for autonomous driving applications, and its utility in data augmentation for downstream robot manipulation policy training. Our experiments show that it surpasses current state-of-the-art methods in structural preservation and instruction alignment. Our contributions are summarized as follow:

1. We design a framework to train image-editing models with RLAIF
2. We craft a multi-objective approach tailored to learning visual nuances beyond simple textual instructions
3. We showcase the quality of the samples by using them as effective training data for data-scarce applications

## 2. Related Works

**Text-guided Image Editing.** In prior work, text-guided image editing methods fall in three categories: architectural tricks that require no training, few-steps optimization for each sample, and large scale finetuning. In the first category, Prompt-to-Prompt (P2P) [22] manipulates attention maps in the diffusion model to control the layout and content of the editing. Plug-and-Play (PNP) [63] injects self-attention maps and spatial features to improve the structure. Still, these remain limited in their ability to edit complex scenes. Others leverage a segmentation mask [2, 12, 34, 40, 55, 66, 69] or semantic mask [37, 39, 73, 76] during the forward pass, which imposes a greater inconvenience on users by requiring them to supply these additional masks.

In the second category, Null-Text Inversion [38] optimizes the null-text embedding during the inversion of an input image. Imagic [28] finetunes model weights and embeddings to align with the input image and the edit text prompt. RB-Modulation [48] leverages stochastic optimal controller to align content and style with visual prompts. Still, these are slow at inference because optimization is required for each generation.

Lastly, other methods adopt the standard reconstruction training like [11, 13, 53, 54, 75]. InstructPix2Pix [11] generates a large synthetic dataset of image pairs with P2P, coupled with instruction-based text prompts. Similarly, Emu Edit [54] expands the dataset with P2P, screens training samples with semantic and structural filters techniques,

and leverages multi-task training. SuTI [13] leverages fine-tuned experts version of Imagen [50] to create high quality samples for the editing model to learn from. MagicBrush [75] assembles an image editing dataset synthesized with DALL-E 2 [44] and manually prune samples. Alchemist [53] uses a rendering tool targeted to modifying material attributes. However, these methods require an oracle to generate images, which requires some curation, and may still infuse its limitations.

**Reinforcement Learning for Diffusion.** Aligning model outputs with human preferences has seen a wide success in the field of language modeling. For objectives that are complex to define explicitly, a popular strategy is reinforcement learning with human feedback (RLHF) [4, 15, 41, 59], where we first teach a reward function to capture output preferences, and leverage reinforcement learning algorithms like proximal policy optimization [52] to finetune models with such rewards.

In the field of diffusion models, several works study the use of human feedback for T2I generation. [31] collects human annotations, and perform maximum likelihood training where the reward is naively used as a weight. [68] designs a reward model that better captures finegrained human preferences. [9, 17] show that diffusion models can be trained with RL using a reward model judging images’ aesthetics [70]. Specific to instructional image-editing, HIVE [77] extends large-dataset supervised training by collecting human feedback on edits and performing off-policy RLHF training. However, these methods need a reward model trained on large-scale human annotation. Not only is this process cumbersome, but it also provides supervision of limited quality. First, it quickly becomes hard for humans to evaluate at a finegrained level the preservation faithfulness of pixels unrelated to the edit instruction. Second, semantic alignment remains vague as it is only compared to the short instruction prompts, where different individuals could easily disagree on the specific interpretation.

Inspired by advances with RLAIF [5, 30], we alleviate the need for humans-in-the-loop and opt for a method where AI models provide the preference supervision tailored to solving the two issues above. Further, we train on-policy by leveraging the framework established by D3PO [71], posit-ing that using online samples would lead to better results.

### 3. Method

In this section, we describe the custom objective designed to obtain parallel supervision for the semantic and structural alignment. In Sec. 3.1, we describe how to alleviate the need for a reward model. Then, we explain in Sec. 3.2 how we design our two separate objectives. Finally, in Sec. 3.3, we present the modified architecture to intake the additional visual prompt conditioning and its modified score

estimate formulation for classifier-free guidance with three conditionings.

#### 3.1. Reinforcement Learning Training of Diffusion Models

Our model should learn from a reward that captures the structural and semantic alignment. For such, the reward function must intake the input image to be edited  $I_{in}$ , the instruction prompt  $c_T$ , the target style image  $I_{sty}$ , and compare those to the generated edit  $I_{gen}$ .

Most RLHF methods train a reward model to then train a downstream model. However, Direct Preference Optimization (DPO) [43] showed that preference ranking can be used to train language models and circumvent reward models. [64] showed that this could be extended to diffusion models. In our work, we leverage the framework introduced by D3PO [71], which expands that of DPO into a multi-step Markov Decision Process (MDP).

Given a pair of outputs  $(y_1, y_2) \sim \pi_{ref}(y|x)$  generated from a reference pre-trained model  $\pi_{ref}$ , we denote the preference as  $y_w \succ y_l|x$  and store the raking tuple  $(x, y_w, y_l)$  in dataset  $\mathcal{D}$ , where  $y_w$  and  $y_l$  are the prefered and dispreferred samples respectively. Following the Bradley-Terry model [10], the human preference distribution  $p^*$  can be expressed by using a reward function  $r^*$  as:

$$p^*(y_1 \succ y_2 | x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))} \quad (1)$$

A parametrized reward model  $r_\phi$  can then be trained via maximum likelihood estimation to approximate  $r^*$  with:

$$\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \rho(r_\phi(x, y_w) - r_\phi(x, y_l))] \quad (2)$$

where  $\rho$  is the logistic function. Prior works in RL have for objective to optimize a distribution such that its associated reward is maximized, allthewhile regularizing this distribution with the KL divergence to remain similar to its initial reference distribution:

$$\max_{\pi_\theta} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_\theta(y|x)} [r_\phi(x, y)] - \beta \mathbb{D}_{KL} [\pi_\theta(y|x) \| \pi_{ref}(y|x)] \quad (3)$$

where  $\beta$  controls the deviation between  $\pi_\theta$  and  $\pi_{ref}$ . This distribution takes the following for optimal solution:

$$\pi_r(y | x) = \frac{1}{Z(x)} \pi_{ref}(y | x) \exp\left(\frac{1}{\beta} r(x, y)\right) \quad (4)$$

where  $Z(x) = \sum_y \pi_{ref}(y|x) \exp\left(\frac{1}{\beta} r(x, y)\right)$  is the partition function. Reorganizing Eq. 4, we obtain the expression for the reward as a function of its associated optimal policy.

$$r(x, y) = \beta \log \frac{\pi_r(y | x)}{\pi_{ref}(y | x)} + \beta \log Z(x) \quad (5)$$

Substituting the parametrized reward function and policy for their optimal counterparts, we reintegrate that expression into Eq. 2. With the change of variables, the loss function is now expressed over policies rather than over reward functions. This closed form avoids having to train a reward model, but rather allows us to directly optimize the model.

$$\mathcal{L}_{\text{DPO}}(\pi_\theta; \pi_{\text{ref}}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \log \rho \left( \beta \log \frac{\pi_\theta(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_\theta(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (6)$$

Extending this to diffusion models, we note a key difference to the described framework. The output is not generated from a single forward pass, but rather a sequential process. To address this, we pose the  $T$ -horizon MDP formulation, adapted from [60], for the  $T$ -timesteps long denoising process.

$$\begin{aligned} s_t &= (\mathbf{x}_{T-t}, \mathbf{c}, t) & P_0(s_0) &= (\mathcal{N}(\mathbf{0}, \mathbf{I}), p(\mathbf{c}), \delta_0) \\ \mathbf{a}_t &= \mathbf{x}_{T-t-1} & P(\mathbf{s}_{t+1} | \mathbf{s}_t, \mathbf{a}_t) &= (\delta_{\mathbf{x}_{T-t-1}}, \delta_c, \delta_{t+1}) \\ r(\mathbf{s}_t, \mathbf{a}_t) &= r((\mathbf{x}_{T-t}, \mathbf{c}, t), \mathbf{x}_{T-t-1}) \\ \pi(\mathbf{a}_t | s_t) &= p_\theta(\mathbf{x}_{T-t-1} | \mathbf{x}_{T-t}, \mathbf{c}, t) \end{aligned}$$

where  $p_\theta(\mathbf{x}_{0:T}|\cdot)$  is a T2I diffusion model,  $\delta$  is the Dirac delta distribution, and  $\mathbf{c}$  is the conditioning distributed according to  $p(\mathbf{c})$ . Note that we disregard  $r$  as our method circumvents it. With such, we treat the denoising process as a sequence of observations and actions:  $\sigma = \{s_0, a_0, s_1, a_1, \dots, s_{T-1}, a_{T-1}\}$ . Since we can only judge the denoised output, we would need to update  $\pi_\theta(\sigma) = \prod_t^T \pi_\theta(s_t, a_t)$ , which is intractable. Following [71], we assume that if the final output of a sequence is preferred over that of another sequence, then any state-action pair of the winning sequence is preferred over that of the losing sequence. Hence, we determine the preferred sequence by sampling an initial state  $s_0 = s_0^w = s_0^l$ , generating two independent sequences, and ranking their output. Accordingly, we express the objective at a certain timestep as:

$$\mathcal{L}_t(\pi_\theta) = -\mathbb{E}_{(s_t, \sigma_w, \sigma_l)} \left[ \log \rho \left( \beta \log \frac{\pi_\theta(a_t^w | s_t^w)}{\pi_{\text{ref}}(a_t^w | s_t^w)} - \beta \log \frac{\pi_\theta(a_t^l | s_t^l)}{\pi_{\text{ref}}(a_t^l | s_t^l)} \right) \right] \quad (7)$$

### 3.2. Multi-objective Joint Training

Now that we know how a preference ranking is converted into a loss that can train our diffusion model, we explain how we rank two generated sample based off their structural alignment with the input image, and their semantic alignment with the text and visual prompts. Our approach splits each objective into two separate scores.

**Structural score.** To obtain a score for how well the structure of the input image has been preserved, we leverage a monocular depth estimation model [72]. Given a pair

of input and edited images, we compute the depth map of each image independently, and define our structural score as the  $L_1$  distance between the two resulting depth maps.

$$\mathcal{L}_{\text{struct}} = \frac{1}{h \cdot w} \sum_{i,j}^{h,w} |f_\phi(I_{\text{in}})_{i,j} - f_\phi(I_{\text{gen}})_{i,j}|$$

where  $f_\phi$  is the depth model. With this metric, we capture any missing, additional or deformed elements in the edit compared to the input.

**Semantic score.** Contrary to the structural score, the semantic alignment should only be measured inside the region of the generated image where the edit is expected. Accordingly, we isolate the relevant regions with a text-conditioned segmentation model to locate the element to be edited. Here, we use grounded-SAM2 [35, 45, 46]. Same goes for the visual prompt, whose style may not be covering the entire frame. Subsequently, we determine the semantic alignment score by computing the distance between the embeddings of relevant patches in the generated and target style images. Additionally, we find that adding the standard pixel-space reconstruction objective on the instruction-irrelevant pixels, which should remain identical, enforces better localization and sharper bounds of the region of the image that must be edited. This serves as a regularizer to prevent the style from spreading over the whole generated frame to unrelated regions. Accordingly, our semantic score is:

$$\begin{aligned} \mathcal{L}_{\text{sem}} &= D(m_{\text{in}} \odot I_{\text{in}}, m_{\text{sty}} \odot I_{\text{sty}}, f_\theta) \\ &\quad + \lambda \cdot (1 - m_{\text{in}}) \odot \|I_{\text{in}} - I_{\text{gen}}\|_2 \end{aligned}$$

where  $D(\cdot)$  is a distance metric,  $m_{\text{og}}$  and  $m_{\text{sty}}$  are the binary segmentation map obtained from the input and style images respectively,  $\odot$  defines the element-wise multiplication, here the cosine distance,  $f_\theta$  is an encoder, and  $\lambda$  is a hyperparameter to weigh the influence of the pixel reconstruction objective relative to the semantic one. We find empirically that  $\lambda = 0.05$  achieves the ideal balance, and that DreamSim [18], amongst other encoders, captures features best aligned with this task (see Sec. 4.2 for ablations).

Similar to [79], we obtain *advantages* [61] by normalizing the scores on a per-batch basis using the mean and variance of each training batch. We then combine the distinct advantages into a unique score, with their relative contribution weighed by a hyperparameter  $\alpha$ .

$$\mathcal{L}_{\text{total}} = \hat{A}_{\text{struct}} + \alpha \cdot \hat{A}_{\text{sem}} , \text{ where } \hat{A} = \frac{\mathcal{L} - \mu_{\mathcal{L}}}{\sqrt{\sigma_{\mathcal{L}}^2 + \epsilon}}$$

Finally, we rank two sequences according to the score of their respective output.

### 3.3. Architecture for Multiple Conditionings

Our architecture is based on that of InstructPix2Pix [11], itself is adapted from Stable Diffusion [47] with one main

architecture modification. That is, they add input channels to the first convolutional layer to intake the encoded input image. We extend this and add more input channels to the first convolutional layer to intake both the input and style images. The weights that operate on the newly added input channels are initialized to zero. In practice, we find that applying a cross-attention layer before feeding the visual prompt improves performance, as it helps identify regions of both images relevant to the instruction. Here, the query is the linear projection of the concatenation of  $\mathcal{E}(I_{in})$  and  $\mathcal{E}(I_{sty})$ , and the key and query are obtained by projecting the CLIP encoding of the instruction prompt.

Accordingly, we adapt the score estimate formulation as follows. Classifier-free guidance (CFG) [24] shifts probability mass where an implicit classifier assigns high likelihood to the conditioning, improving visual quality and faithfulness of samples. The standard CFG estimate is:

$$\tilde{e}_\theta(z_t, c) = e_\theta(z_t, \emptyset) + s \cdot (e_\theta(z_t, c) - e_\theta(z_t, \emptyset))$$

InstructPix2Pix adds a conditioning on the input image to be edited on top of the text instruction, further disentangling the score estimate yields:

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

In our case, we find it beneficial to extend CFG with respect to the additional visual prompt, and use the following (see Appendix 6.1 for the complete derivation):

$$\begin{aligned} \tilde{e}_\theta(z_t, c_{I_{in}}, c_{I_{sty}}, c_T) &= e_\theta(z_t, \emptyset, \emptyset, \emptyset) \\ &\quad + s_{I_{in}} \cdot (e_\theta(z_t, c_{I_{in}}, \emptyset, \emptyset) - e_\theta(z_t, \emptyset, \emptyset, \emptyset)) \\ &\quad + s_{I_{sty}} \cdot (e_\theta(z_t, c_{I_{in}}, c_{I_{sty}}, \emptyset) - e_\theta(z_t, c_{I_{in}}, \emptyset, \emptyset)) \\ &\quad + s_T \cdot (e_\theta(z_t, c_{I_{in}}, c_{I_{sty}}, c_T) - e_\theta(z_t, c_{I_{in}}, c_{I_{sty}}, \emptyset)) \end{aligned}$$

During training, we alternate between 5 style images that relate to a same text prompt to be used as conditioning. This helps prevent overfit on spurious cues.

## 4. Experiments

**incomplete** This section presents the experimental results and ablation studies of our technical choices, demonstrating the effectiveness of our method. We apply the default guidance scale parameters ( $s_{I_{in}}$  and  $s_T$ ) from InstructPix2Pix for a fair comparison.

We evaluate our method focusing on the ability to perform located edits in complex scenes. Accordingly, we choose crowded street images from the Oxford RobotCar dataset [36] as our input images to be edited. We train our method to learn 7 types of edits separately, all related to modifying the road, as it is guaranteed to remain visible on every frame without requiring complex curation of the input



Figure 1. Qualitative comparison. Our method outperforms its counterparts by significantly editing the image while sharply preserving the structure of regions unrelated to the text prompt. We omit visualizing the conditioning style image as the other methods cannot leverage it. Visualize more samples in Sec. 6

images (see full list in Sec. 6.2). The edit requests range from realistic weather modification like adding snow on the road, to stylistic material modification like turning the road into wood. We conduct evaluations for each edit type on a set of 3500 images. We also showcase the ability to capture subtle instruction nuances beyond the text prompt by training two models with the same text instruction but two different styles of visual conditionings. Additionally, we display the ability of our model to produce edits that serve as training data in a data-scarce domain.

### 4.1. Baseline Comparisons

We compare our model’s performance against that of InstructPix2Pix (IP2P) [11], HIVE [77], MagicBrush (MBrush) [75], all based on the same v1.5 of stable diffusion.

We first compare the results qualitatively in Figure 3. We notice that InstructPix2Pix and MagicBrush struggle to precisely locate the region to be edited while leaving the rest intact, and to generate naturalistic edits with unrealistic prompts. Conversely, HIVE often prefers maintaining high fidelity with the input image at the cost of minimalistic edits. Our model achieves a better balance between excessive and insufficient editing, resulting in edits that are both highly aligned with the prompts and remain faithful to the input image. Further, our method reduces the impact of biases associated with concepts mentioned in the text prompt. For instance, we suppress InstructPix2Pix’s hallucinations of a forest when generating an image with the word ‘wood’ mentioned in the instruction.

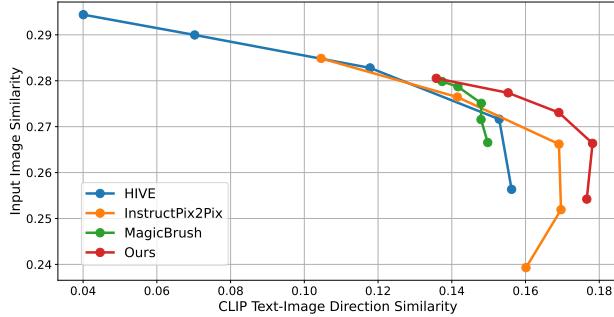


Figure 2. Comparison between instructional editing models. We plot the trade-off between consistency with the input image (Y-axis) and consistency with the edit (X-axis). For both metrics, higher is better. For all methods, we fix the same parameters as in [11] and vary the  $s_{I_{in}} \in [1.0, 2.0]$

Method	$\text{Depth}_{\text{out}} \downarrow$	$\text{Depth}_{\text{all}} \downarrow$	$L_2 \downarrow$	$\text{STY}_{\text{in}} \uparrow$	$\text{CLIP}_{\text{txt}} \uparrow$
IP2P	35.93	50.52	0.064	0.264	0.238
MBrush	46.06	60.69	0.054	0.250	0.234
HIVE	21.56	29.47	<b>0.036</b>	0.248	0.224
Ours	<b>14.16</b>	<b>18.39</b>	0.062	<b>0.282</b>	<b>0.251</b>

Table 1. Comparison of structure preservation through predicted monocular depth mask alignment and reconstruction metrics (left), and semantic alignment with both the text and visual prompts (right). For the text alignment, we use a descriptive prompt to capture all information on the image.

We also quantitatively analyze the tradeoff between the fidelity with the input image, and the alignment with the text instruction. The former is measured through the cosine similarity of image patch embeddings that are outside of the mask locating the element to be edited. We rely on grounded SAM2 to produce high-quality masks, which we find surpass human-made masks provided in public benchmarks [75], regarding contour precision. With such masking, our reported metrics are better targeted to the regions of interest. We capture both high and low-level visual features by taking the average score across three encoders, namely DINO, CLIP, and DreamSim. The text alignment is measured through the directional CLIP similarity [20], which denotes how much the change in descriptive text captions agrees with the change in the images. This evaluation operates on the entire frame, i.e. no masking operation is involved. Both metrics are antagonistic, increasing the desired edit strength will reduce the output’s faithfulness to the input image. Still, we find that when comparing our method with its counterparts in Fig. 2, our results have notably higher directional similarity values for the same image consistency. We also confirm our qualitative insights, noticing that HIVE tends to prioritize preserving the original image over executing the requested edits, whereas our method

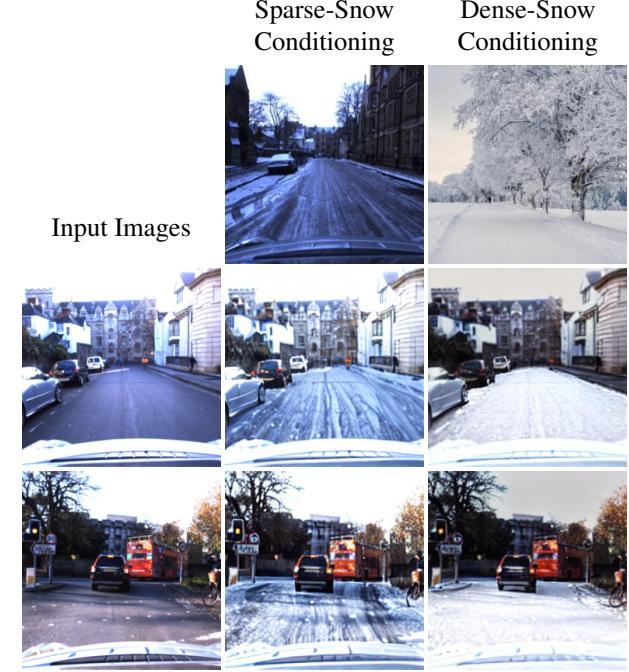


Figure 3. Visualization of the impact of the visual prompt on generated samples when provided the text instruction “add snow on the road”. Our method effectively captures semantic nuances beyond that described in the text prompt.

respects closely the editing directive. In Table 1, we also find that HIVE better reconstructs the regions that should be kept intact. Still, our method outperforms the others on the depth mask alignment metrics, both measured on these regions and on the whole images. This suggests that, while our model has some imprecision on a finegrained level, it successfully preserves the fundamental structure, which is crucial for visual coherence and perceived realism.

We then demonstrate our model’s capability to interpret subtle details beyond what is specified in a text prompt. To illustrate this, we train two distinct models to add a dense and a sparse layer of snow on the road, even when these specifics were not included in the text prompt. Fig. 3 reveals that our model excels at reproducing the visual style of the concept hinted at in the text prompt, while respecting the spatial composition of the input image during the edit. For instance, the sparse snow generated by our model mirrors the stripe-like pattern of the conditioning image and aligns with the road’s direction in the input images. This underscores the benefit of leveraging a visual prompt to infuse fine-grained nuances without requiring an extensive text prompt.

We quantitatively validate these observations in Table 1. Our evaluation focuses on two key aspects: visual semantic alignment between regions-of-interest, and text-image alignment scores. For the visual alignment, we compute the

Method	HPSv2 $\uparrow$	ImageReward $\uparrow$	PickScore $\uparrow$
IP2P	21.85	-0.314	18.92
MBrush	21.92	-0.454	18.84
HIVE	21.61	-0.567	18.89
Ours	<b>22.12</b>	<b>-0.134</b>	<b>18.94</b>

Table 2. Comparison between instructional editing models on human-preference-aligned visual scoring metrics.

cosine similarity between masked regions in the generated and conditioning images, averaging across DINO, CLIP, and DreamSim embeddings. The text-image alignment is measured through CLIP similarity ( $CLIP_{txt}$ ) between the edited image and output caption. We find that our model outperforms its counterparts, both in the visual and text-image alignments, further confirming our method’s efficacy to increase instruction fidelity.

Finally, in Table 2, we compare the scores obtained from T2I synthesis preference prediction models, namely HPSv2 [67], ImageReward [70], and PickScore [29], which are all trained to emulate human preferences. We find that our model surpasses all other baselines across the different edits. This confirms our model’s superior ability to preserve the essential structural features of the input image, which are crucial for perceived realism, while also aligning effectively with the specified editing instructions.

## 4.2. Ablation Study

**Classifier-free-guidance score** How does the parameter  $s_{sty}$  influence the generated samples

**Information Encoders** Different encoders capture different information in their representations, guiding the learning process with distinct models leads to different outputs. Other works commonly use DINO and CLIP. However, recently, Dreamsim [18] showed to outperform those encoders in alignment with human preferences. Similarly, we find that this encoder captures the best balance between high and low level features compared to DINO and CLIP. We here only analyze qualitative results, since all semantic visual alignment metrics are based on these encoders

## 4.3. Sim2Real Edits

Finally, we showcase the ability for our model to produce complex edits that are hard to explain with a text prompt, and that effectively serve as downstream model training data. Robotics policies are known to be very brittle and perform poorly out-of-distribution. Further, due to the complexity of collecting real-world data in those robot learning tasks, many works use a simulation environment to collect large-scale data to pre-train their robot policy [32]. Still, the gap between the simulation domain and real domain is

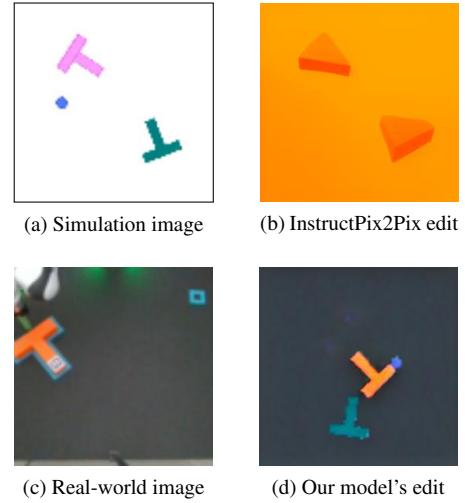


Figure 4. Comparison between the base InstructPix2Pix model and our specialized model on Sim-2-Real edits. Our method captures both the texture of the table, the color of the foreground object, and adds some slight shading around, like there is in the real world.

generally large. Hence, models often require a layer of fine-tuning on real-world data to perform correctly, which still requires costly and lengthy data collection in real-world. Here, we aim to train our diffusion model to perform image editing to modify said simulation images into samples closely resemble the real-world domain. This requires only a few images from the real world to serve as target style conditioning for the editing process, rather than tens of thousands. We select a robot manipulation task referred to as Push-T, first introduced in [14], where a T-shaped block is placed randomly on a surface and a autonomous arm must push it to align its pose with that of a target, drawn somewhere else on that same surface.

As shown in Fig. 4, we see that our model can perform highly realistic edits, far beyond what the base model can produce. These are generated in two steps, first by modifying the background to match that of the table, and second by changing the color of the T-block being pushed. In this toy example, we find that the input images are too unrealistic for the depth model to produce coherent depth maps, hindering the training process. Accordingly, we enforce structural alignment through the alignment of Canny edge maps rather than depth maps.

pending

## 5. Conclusion

In our paper, we introduce a novel method for instructional editing with neither a large dataset of curated images or human feedback. Rather, we leverage other AI models to

provide a balanced supervision to align generated images with what humans would like. Specifically, we showcase its ability to produce sharp edits and preserve instruction-irrelevant regions in an image. Further, we show that our model can interpret fine-grained visual nuances beyond a text prompt, largely alleviating the need for complex and highly descriptive textual prompts. We find that combining these two yields generated samples with higher perceived realism compared to methods trained for pixel-perfect reconstruction of fine textures. Lastly, we present an example application in a robot manipulation setting, where the T2I generative model serves as data augmentation, and significantly reduces the required amount of real-world data.

While our method demonstrates impressive performance, we have also identified limitations inherent to its training method. Types of edits are limited to local ones, where a clear segmentation can be drawn around the shape of the region to be edited. As shown with the nighttime edit examples, a trick is to create a mask that captures the entire image when teaching a model to do whole-image modifications. Still, some more automated adaptability is desirable. Further, the edits are limited to ones that do not modify the global shape of the element to be modified, but rather modify its texture or surface. Some interesting future work could include adding some more flexible bounds in the masking and depth alignment operations (by ignoring the top quantile for depth, or voluntarily enlarging the segmentation mask for semantics) to allow for edits that involve adding, removing, or modifying the shape of objects a bit. Another limitation is that the semantic alignment is dependent on what the encoder can capture. Also, the method relies on the zero-shot performance of the base model. Hence, the biases of such encoder and diffusion model might be infused in the final outputs. Nevertheless, these limitations can be addressed by substituting the encoder or generative prior with suitable alternatives in a plug-and-play fashion.

## References

- [1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning, 2022. 1
- [2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, page 18187–18197. IEEE, 2022. 2
- [3] Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mohammad Norouzi, and David J. Fleet. Synthetic data from diffusion models improves imagenet classification, 2023. 1
- [4] Yuntao Bai, Andy Jones, Kamal Ndousse, Amanda Askell, Anna Chen, Nova DasSarma, Dawn Drain, Stanislav Fort, Deep Ganguli, Tom Henighan, Nicholas Joseph, Saurav Kadavath, Jackson Kernion, Tom Conerly, Sheer El-Showk, Nelson Elhage, Zac Hatfield-Dodds, Danny Hernandez, Tristan Hume, Scott Johnston, Shauna Kravec, Liane Lovitt, Neel Nanda, Catherine Olsson, Dario Amodei, Tom Brown, Jack Clark, Sam McCandlish, Chris Olah, Ben Mann, and Jared Kaplan. Training a helpful and harmless assistant with reinforcement learning from human feedback, 2022. 3
- [5] Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, Carol Chen, Catherine Olsson, Christopher Olah, Danny Hernandez, Dawn Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson, Ethan Perez, Jamie Kerr, Jared Mueller, Jeffrey Ladish, Joshua Landau, Kamal Ndousse, Kamile Lukosuite, Liane Lovitt, Michael Sellitto, Nelson Elhage, Nicholas Schiefer, Noemi Mercado, Nova DasSarma, Robert Lasenby, Robin Larson, Sam Ringer, Scott Johnston, Shauna Kravec, Sheer El Showk, Stanislav Fort, Tamara Lanham, Timothy Telleen-Lawton, Tom Conerly, Tom Henighan, Tristan Hume, Samuel R. Bowman, Zac Hatfield-Dodds, Ben Mann, Dario Amodei, Nicholas Joseph, Sam McCandlish, Tom Brown, and Jared Kaplan. Constitutional ai: Harmlessness from ai feedback, 2022. 3
- [6] Hritik Bansal and Aditya Grover. Leaving reality to imagination: Robust classification via generated datasets, 2023. 1
- [7] Manel Baradad, Jonas Wulff, Tongzhou Wang, Phillip Isola, and Antonio Torralba. Learning to see by looking at noise, 2022. 1
- [8] Elior Benarous, Sotiris Anagnostidis, Luca Biggio, and Thomas Hofmann. Harnessing synthetic datasets: The role of shape bias in deep neural network generalization, 2023. 1
- [9] Kevin Black, Michael Janner, Yilun Du, Ilya Kostrikov, and Sergey Levine. Training diffusion models with reinforcement learning, 2024. 3
- [10] Ralph Allan Bradley and Milton E. Terry. Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika*, 39(3/4):324–345, 1952. 3
- [11] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. 1, 2, 4, 5, 6
- [12] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactr: Tuning-free mutual self-attention control for consistent image synthesis and editing, 2023. 1, 2
- [13] Wenhui Chen, Hexiang Hu, Yandong Li, Nataniel Ruiz, Xuhui Jia, Ming-Wei Chang, and William W. Cohen. Subject-driven text-to-image generation via apprenticeship learning, 2023. 2, 3
- [14] Cheng Chi, Zhenjia Xu, Siyuan Feng, Eric Cousineau, Yilun Du, Benjamin Burchfiel, Russ Tedrake, and Shuran Song. Diffusion policy: Visuomotor policy learning via action diffusion. *The International Journal of Robotics Research*, 2024. 7
- [15] Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences, 2023. 3
- [16] Lisa Dunlap, Alyssa Umino, Han Zhang, Jiezhi Yang, Joseph E. Gonzalez, and Trevor Darrell. Diversify your vision datasets with automatic diffusion-based augmentation, 2023. 1
- [17] Ying Fan, Olivia Watkins, Yuqing Du, Hao Liu, Moonkyung Ryu, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, Kangwook Lee, and Kimin Lee. Dpok: Reinforcement learning for fine-tuning text-to-image diffusion models, 2023. 3
- [18] Stephanie Fu, Netanel Tamir, Shobhita Sundaram, Lucy Chai, Richard Zhang, Tali Dekel, and Phillip Isola. Dreamsim: Learning new dimensions of human visual similarity using synthetic data, 2023. 4, 7
- [19] Yunxiang Fu, Chaoqi Chen, Yu Qiao, and Yizhou Yu. Dreamda: Generative data augmentation with diffusion models, 2024. 1
- [20] Rinon Gal, Or Patashnik, Haggai Maron, Gal Chechik, and Daniel Cohen-Or. Stylegan-nada: Clip-guided domain adaptation of image generators, 2021. 6
- [21] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 2
- [22] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control, 2022. 1, 2
- [23] Amir Hertz, Andrey Voynov, Shlomi Fruchter, and Daniel Cohen-Or. Style aligned image generation via shared attention, 2024. 2
- [24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance, 2022. 5
- [25] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. 1
- [26] Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(24):695–709, 2005. 12

- [27] Hirokatsu Kataoka, Kazushige Okayasu, Asato Matsumoto, Eisuke Yamagata, Ryosuke Yamada, Nakamasa Inoue, Akio Nakamura, and Yutaka Satoh. Pre-training without natural images, 2021. 1
- [28] 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, 2023. 1, 2
- [29] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Maitana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation, 2023. 7
- [30] Harrison Lee, Samrat Phatale, Hassan Mansoor, Thomas Mesnard, Johan Ferret, Kellie Lu, Colton Bishop, Ethan Hall, Victor Carbune, Abhinav Rastogi, and Sushant Prakash. Rlaif vs. rlhf: Scaling reinforcement learning from human feedback with ai feedback, 2024. 3
- [31] Kimin Lee, Hao Liu, Moonkyung Ryu, Olivia Watkins, Yuqing Du, Craig Boutilier, Pieter Abbeel, Mohammad Ghavamzadeh, and Shixiang Shane Gu. Aligning text-to-image models using human feedback, 2023. 3
- [32] Xuanlin Li, Kyle Hsu, Jiayuan Gu, Karl Pertsch, Oier Mees, Homer Rich Walke, Chuyuan Fu, Ishikaa Lunawat, Isabel Sieh, Sean Kirmani, Sergey Levine, Jiajun Wu, Chelsea Finn, Hao Su, Quan Vuong, and Ted Xiao. Evaluating real-world robot manipulation policies in simulation, 2024. 7
- [33] Jun Hao Liew, Hanshu Yan, Daquan Zhou, and Jiashi Feng. Magicmix: Semantic mixing with diffusion models, 2022. 1
- [34] Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, and Sanja Fidler. Editgan: High-precision semantic image editing, 2021. 2
- [35] Shilong Liu, Zhaoyang Zeng, Tianhe Ren, Feng Li, Hao Zhang, Jie Yang, Chunyuan Li, Jianwei Yang, Hang Su, Jun Zhu, et al. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*, 2023. 4
- [36] Will Maddern, Geoff Pascoe, Chris Linegar, and Paul Newman. 1 Year, 1000km: The Oxford RobotCar Dataset. *The International Journal of Robotics Research (IJRR)*, 36(1):3–15, 2017. 5
- [37] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations, 2022. 2
- [38] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models, 2022. 1, 2
- [39] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models, 2023. 2
- [40] 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, 2022. 2
- [41] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. 3
- [42] 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. 1
- [43] Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model, 2024. 3
- [44] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. 1, 3
- [45] Nikhila Ravi, Valentin Gabeur, Yuan-Ting Hu, Ronghang Hu, Chaitanya Ryali, Tengyu Ma, Haitham Khedr, Roman Rädle, Chloe Rolland, Laura Gustafson, Eric Mintun, Junting Pan, Kalyan Vasudev Alwala, Nicolas Carion, Chao-Yuan Wu, Ross Girshick, Piotr Dollár, and Christoph Feichtenhofer. Sam 2: Segment anything in images and videos, 2024. 4
- [46] Tianhe Ren, Shilong Liu, Ailing Zeng, Jing Lin, Kun-chang Li, He Cao, Jiayu Chen, Xinyu Huang, Yukang Chen, Feng Yan, Zhaoyang Zeng, Hao Zhang, Feng Li, Jie Yang, Hongyang Li, Qing Jiang, and Lei Zhang. Grounded sam: Assembling open-world models for diverse visual tasks, 2024. 4
- [47] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 1, 4
- [48] Litu Rout, Yujia Chen, Nataniel Ruiz, Abhishek Kumar, Constantine Caramanis, Sanjay Shakkottai, and Wen-Sheng Chu. Rb-modulation: Training-free personalization of diffusion models using stochastic optimal control, 2024. 2
- [49] 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, 2023. 2
- [50] 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, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. 3
- [51] Mert Bulent Sarıyıldız, Kartek Alahari, Diane Larlus, and Yannis Kalantidis. Fake it till you make it: Learning transferable representations from synthetic imagenet clones, 2023. 1
- [52] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms, 2017. 3
- [53] Prafull Sharma, Varun Jampani, Yuanzhen Li, Xuhui Jia, Dmitry Lagun, Fredo Durand, William T. Freeman, and Mark Matthews. Alchemist: Parametric control of material properties with diffusion models, 2023. 2, 3

- [54] Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh, and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks, 2023. 1, 2
- [55] Yichun Shi, Xiao Yang, Yangyue Wan, and Xiaohui Shen. Semanticstylegan: Learning compositional generative priors for controllable image synthesis and editing, 2022. 2
- [56] Joonghyuk Shin, Minguk Kang, and Jaesik Park. Fill-up: Balancing long-tailed data with generative models, 2023. 1
- [57] Jordan Shipard, Arnold Wiliem, Kien Nguyen Thanh, Wei Xiang, and Clinton Fookes. Diversity is definitely needed: Improving model-agnostic zero-shot classification via stable diffusion, 2023. 1
- [58] Kihyuk Sohn, Nataniel Ruiz, Kimin Lee, Daniel Castro Chin, Irina Blok, Huiwen Chang, Jarred Barber, Lu Jiang, Glenn Entis, Yuanzhen Li, Yuan Hao, Irfan Essa, Michael Rubinstein, and Dilip Krishnan. Styledrop: Text-to-image generation in any style, 2023. 2
- [59] Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford, Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. 3
- [60] R.S. Sutton and A.G. Barto. Reinforcement learning: An introduction. *IEEE Transactions on Neural Networks*, 9(5): 1054–1054, 1998. 4
- [61] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. Policy gradient methods for reinforcement learning with function approximation. In *Advances in Neural Information Processing Systems*. MIT Press, 1999. 4
- [62] Brandon Trabucco, Kyle Doherty, Max Gurinas, and Ruslan Salakhutdinov. Effective data augmentation with diffusion models, 2023. 1
- [63] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation, 2022. 1, 2
- [64] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization, 2023. 3
- [65] Haofan Wang, Matteo Spinelli, Qixun Wang, Xu Bai, Zekui Qin, and Anthony Chen. Instantstyle: Free lunch towards style-preserving in text-to-image generation, 2024. 2
- [66] Su Wang, Chitwan Saharia, Ceslee Montgomery, Jordi Pont-Tuset, Shai Noy, Stefano Pellegrini, Yasumasa Onoe, Sarah Laszlo, David J. Fleet, Radu Soricut, Jason Baldridge, Mohammad Norouzi, Peter Anderson, and William Chan. Imagen editor and editbench: Advancing and evaluating text-guided image inpainting, 2023. 2
- [67] Xiaoshi Wu, Yiming Hao, Keqiang Sun, Yixiong Chen, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score v2: A solid benchmark for evaluating human preferences of text-to-image synthesis, 2023. 7
- [68] Xiaoshi Wu, Keqiang Sun, Feng Zhu, Rui Zhao, and Hongsheng Li. Human preference score: Better aligning text-to-image models with human preference, 2023. 3
- [69] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model, 2022. 2
- [70] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imagereward: Learning and evaluating human preferences for text-to-image generation, 2023. 3, 7
- [71] Kai Yang, Jian Tao, Jiafei Lyu, Chunjiang Ge, Jiaxin Chen, Qimai Li, Weihan Shen, Xiaolong Zhu, and Xiu Li. Using human feedback to fine-tune diffusion models without any reward model, 2024. 3, 4
- [72] Lihe Yang, Bingyi Kang, Zilong Huang, Zhen Zhao, Xiaogang Xu, Jiashi Feng, and Hengshuang Zhao. Depth anything v2, 2024. 4
- [73] Hu Ye, Jun Zhang, Sibo Liu, Xiao Han, and Wei Yang. Ipadapter: Text compatible image prompt adapter for text-to-image diffusion models, 2023. 2
- [74] Jianhao Yuan, Francesco Pinto, Adam Davies, and Philip Torr. Not just pretty pictures: Toward interventional data augmentation using text-to-image generators, 2024. 1
- [75] Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing, 2024. 2, 3, 5, 6
- [76] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2
- [77] Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, Caiming Xiong, and Ran Xu. Hive: Harnessing human feedback for instructional visual editing, 2024. 3, 5
- [78] Yifan Zhang, Daquan Zhou, Bryan Hooi, Kai Wang, and Jiashi Feng. Expanding small-scale datasets with guided imagination, 2023. 1
- [79] Yinan Zhang, Eric Tzeng, Yilun Du, and Dmitry Kislyuk. Large-scale reinforcement learning for diffusion models, 2024. 4
- [80] Yongchao Zhou, Hshmat Sahak, and Jimmy Ba. Training on thin air: Improve image classification with generated data, 2023. 1

## 6. Appendix

### 6.1. Derivation for Classifier-free Guidance with Three Conditionings

We introduce separate guidance scales like ip2p to enable separately trading off the strength of each conditioning. The modified score estimate for our model is derived as follows. Our generative model learns  $P(z|c_T, c_{I_{sty}}, c_{I_{in}})$ , which corresponds to the probability distribution of the image latents  $z = \mathcal{E}(x)$  conditioned on an input original image  $c_{I_{in}}$ , a target style image  $c_{I_{sty}}$ , and a text instruction  $c_T$ . We arrive at our particular classifier-free guidance formulation by expressing the conditional probability as follows:

$$\begin{aligned} P(z|c_T, c_{I_{sty}}, c_{I_{in}}) &= \frac{P(z, c_T, c_{I_{sty}}, c_{I_{in}})}{P(c_T, c_{I_{sty}}, c_{I_{in}})} \\ &= \frac{P(c_T|c_{I_{sty}}, c_{I_{in}}, z)P(c_{I_{sty}}|c_{I_{in}}, z)P(c_{I_{in}}|z)P(z)}{P(c_T, c_{I_{sty}}, c_{I_{in}})} \end{aligned} \quad (8)$$

Diffusion models estimate the score [26] of the data distribution, i.e. the derivative of the log probability. Taking the logarithm of the expression above yields the following:

$$\begin{aligned} \log(P(z|c_T, c_{I_{sty}}, c_{I_{in}})) &= \log(P(c_T|c_{I_{sty}}, c_{I_{in}}, z)) \\ &\quad + \log(P(c_{I_{sty}}|c_{I_{in}}, z)) \\ &\quad + \log(P(c_{I_{in}}|z)) \\ &\quad + \log(P(z)) \\ &\quad - \log(P(c_T, c_{I_{sty}}, c_{I_{in}})) \end{aligned} \quad (9)$$

Taking the derivative and rearranging, we obtain:

$$\begin{aligned} \nabla_z \log(P(z|c_T, c_{I_{sty}}, c_{I_{in}})) &= \nabla_z \log(P(z)) \\ &\quad + \nabla_z \log(P(c_{I_{in}}|z)) \\ &\quad + \nabla_z \log(P(c_{I_{sty}}|c_{I_{in}}, z)) \\ &\quad + \nabla_z \log(P(c_T|c_{I_{sty}}, c_{I_{in}}, z)) \end{aligned} \quad (10)$$

### 6.2. Training Details

Depending on the specialization, we train the base Instruct-Pix2Pix model over 15-40 steps, with a batch size of 512 samples. We posit that the biggest varying factor is how far out-of-distribution the requested edit is, as the KL-Divergence regularization limits the weights of the model from shifting too quickly too far from its initialization. For example, adding thick snow on the road needed 14 steps to converge, whereas turning the road into wood needed 43. For a fair comparison of our method, we do not conduct any hyperparameter optimization, and keep the same RL training hyperparameters as in D3PO.

We train our model to perform 7 types of edits of the road separately, these are adding snow (dense and sparse), rain, and sand on top; and changing it into gold and wood. We include a last one operating on the entire image, which

is changing the time to nighttime. To train the latter, we handcraft the mask in the supervision to select the entire frame.

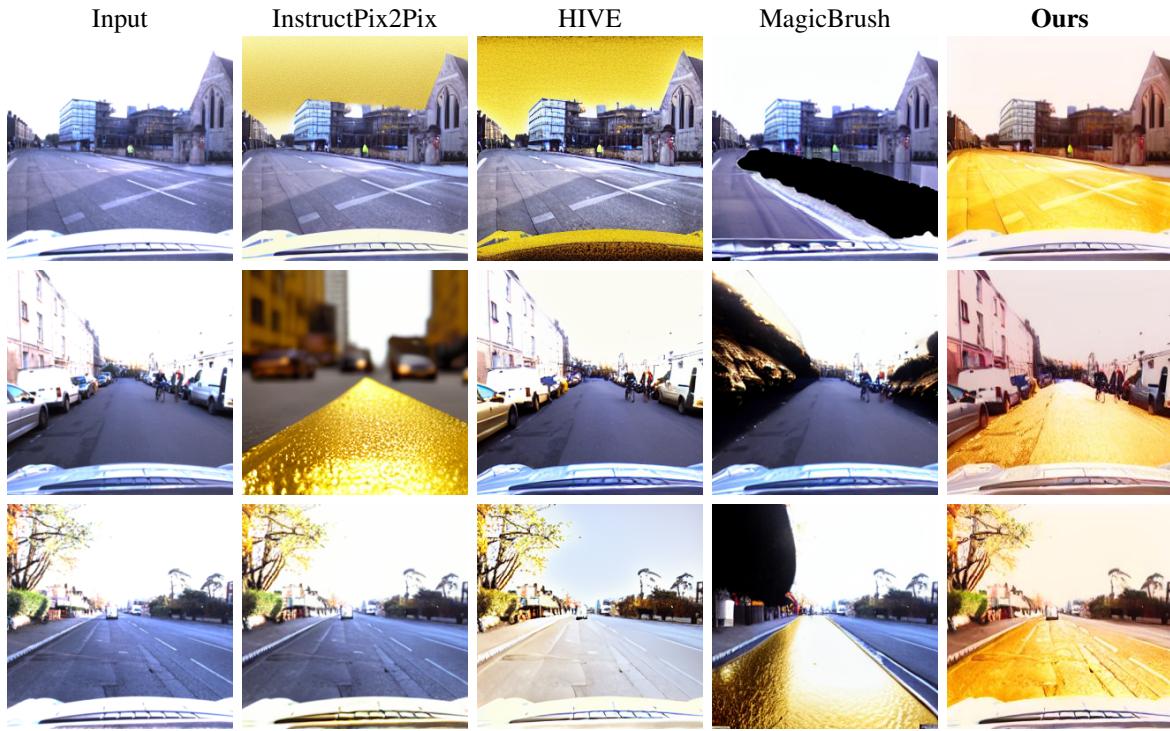


Figure 5. More comparisons between our method and other baselines when given the textual instruction prompt: “*change the road into gold*”

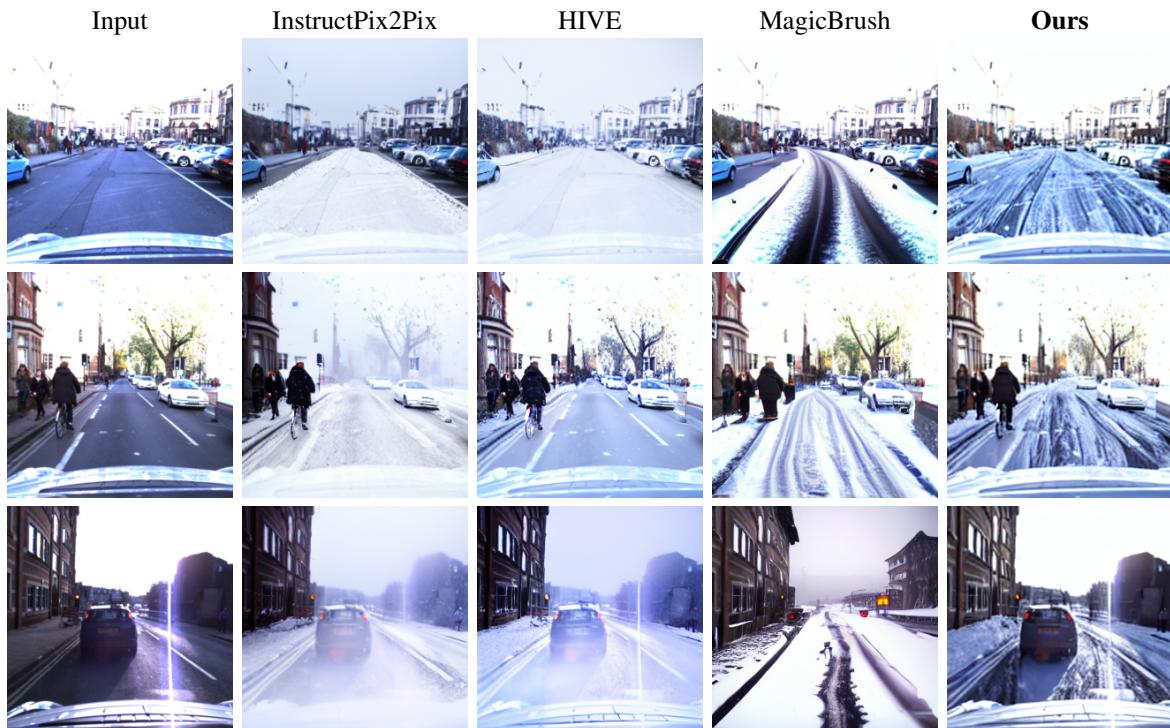


Figure 6. More comparisons between our method and other baselines when given the textual instruction prompt: “*add snow on the road*”



Figure 7. More comparisons between our method and other baselines when given the textual instruction prompt: “*change the time to nighttime*”



Figure 8. More comparisons between our method and other baselines when given the textual instruction prompt: “*add rain on the road*”

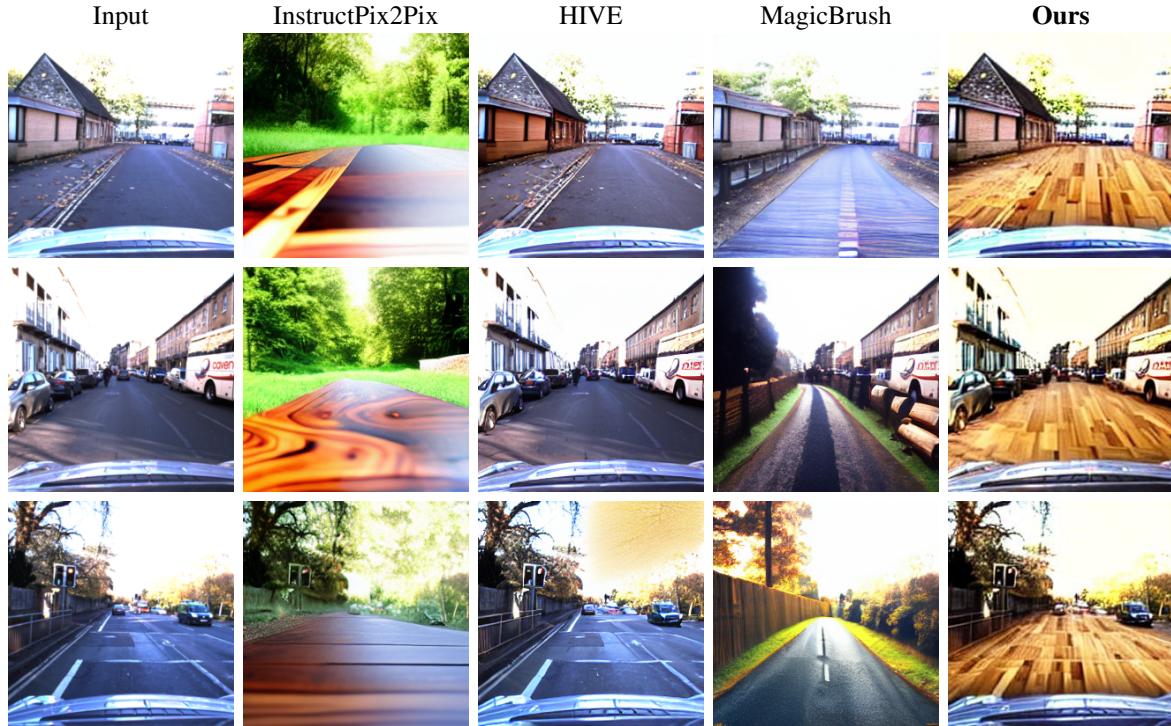


Figure 9. More comparisons between our method and other baselines when given the textual instruction prompt: “*change the road into wood*”

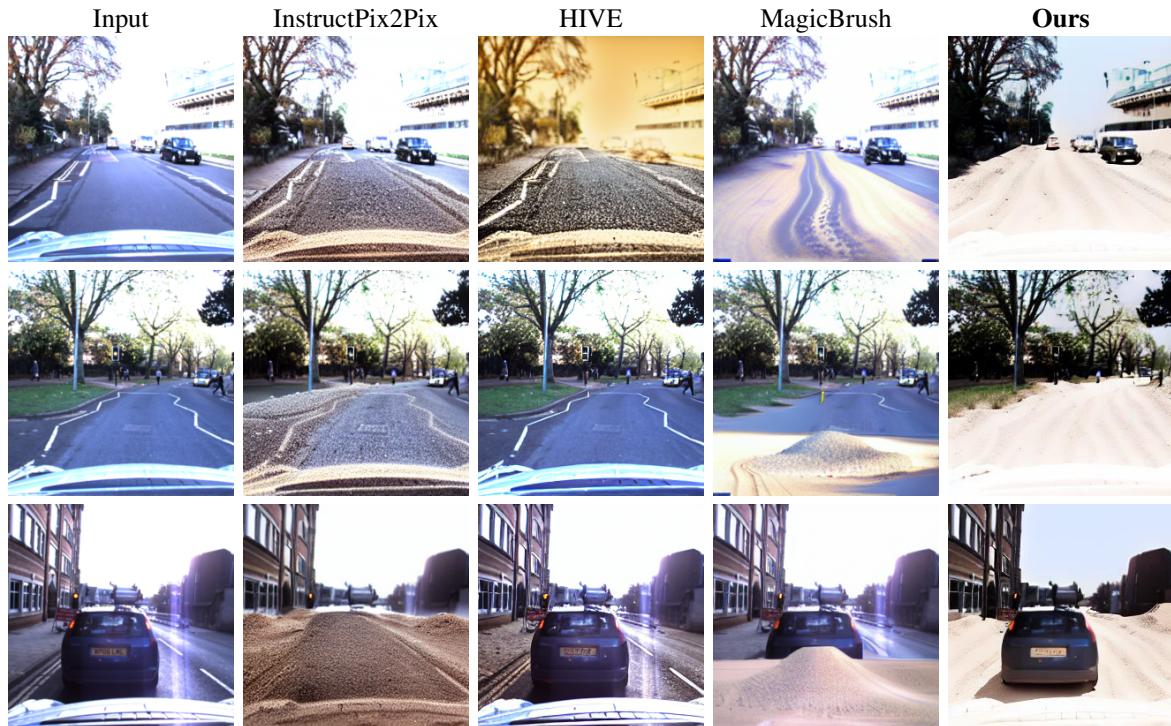


Figure 10. More comparisons between our method and other baselines when given the textual instruction prompt: “*add sand on the road*”



Figure 11. More comparisons between our method and other baselines when given the textual instruction prompt: “add snow on the road”