

INSTRUCTBRUSH: LEARNING ATTENTION-BASED VISUAL INSTRUCTION FOR IMAGE EDITING

Anonymous authors

Paper under double-blind review



Figure 1: In terms of abstract and fine-grained edits, text-guided methods have difficulty accurately analogy the edits exhibited in reference image pairs, even with the help of multimodal large language models. In contrast, our method can better capture these edit concepts and apply them to edits of new images.

ABSTRACT

Diffusion-based image editing methods have garnered significant attention in image editing. However, despite encompassing a wide range of editing priors, these methods are helpless when handling editing tasks that are challenging for users to accurately describe. We propose *InstructBrush*, an inversion method for instruction-based image editing methods to bridge this gap. It extracts editing effects from example image pairs as editing instructions to guide the editing of new images. Two key techniques are introduced into *InstructBrush*, *Attention-based Instruction Optimization* and *Transformation-oriented Instruction Initialization*, to address the limitations of the previous method in terms of inversion effects and instruction generalization. To explore the ability of visual prompt editing methods to guide image editing in open scenarios, we establish a **Transformation-Oriented Paired Benchmark (TOP-Bench)**. Quantitatively and qualitatively, our approach achieves superior performance in editing and is more semantically consistent with the target editing effects. The code and benchmark will be released upon acceptance.

1 INTRODUCTION

Recently developed diffusion-based image editing methods Hertz et al. (2022); Tumanyan et al. (2023); Brooks et al. (2023); Xu et al. (2024) enable users to effortlessly achieve their editing goals using natural language prompts. While they have garnered significant attention owing to their flexibility and versatility in image editing, they still face challenges when dealing with editing tasks that are difficult for users to describe. Specifically, while guiding image editing with language is natural and straightforward, it becomes particularly challenging when users wish to apply analogous manipulation on others’ finished edits or image transformations implemented by other tools, as shown in Figure 1. In such cases, using text or a single image as a condition to guide diffusion models for editing is quite difficult. It makes sense to provide a pair of example images to demonstrate this editing effect.

This motivates the demand for the problem of *visual prompt editing*. Similar to image analogies Jacobs et al. (2001), this problem learns an edit concept from image pairs, and subsequently applies it to edit new images. These image pairs that provide information about image transformations are also called *visual prompts*, which serves as a valuable replacement when language is imprecise in describing specific editing concepts.

One way to implement visual prompt editing is visual in-context learning Yang et al. (2024); Gu et al. (2024). It constructs the visual prompt as well as the input image and prediction noise as a grid-like input, and then uses the inpainting diffusion model to model the task as an inpainting task to predict the output. Although this paradigm can learn general image transformations by analogy with visual prompts, its performance is slightly inferior for the specific task of image editing. In addition, due to the limitation of grid input, it cannot be applied to the editing of high-resolution images. To address these issues, visual instruction inversion Nguyen et al. (2023) replaces the inpaint diffusion model with the instruction-based editing model Brooks et al. (2023); Geng et al. (2023) to improve the performance on image editing tasks while supporting high-resolution image editing. It uses the text inversion method Gal et al. (2022a) to invert the editing concepts revealed by visual prompts into the feature space of text instructions to guide the editing of new input images, but it struggles with the editing effects for two reasons: 1) Inverting instructions in textual space limits their representational ability. Since the text encoder is aligned on text-image pairs with rough descriptions, it is challenging to provide specific representations of the image editing details Chen et al. (2023d). 2) Its semantic-level instruction initialization introduces editing-irrelevant content from visual prompts, hence limiting the generalization of the instruction in generalized scenarios.

To bridge these gaps, we introduce *InstructBrush*, an instruction inversion-based method for visual prompt editing by leveraging the instruction-based image editing model. In contrast to the previous methods, we propose the *Attention-based Instruction Optimization*. It improve the representation ability of instruction guidance by localizing and learning the editing concepts represented by visual prompts in the cross-attention layer of the diffusion model. To introduce semantic-level guidance related to editing, we introduce the *Transformation-oriented Instruction Initialization*. It ingeniously separates editing-related information from the content of visual prompts and incorporates it into the learned instructions. This effectively mitigates the risk of previous method Nguyen et al. (2023) compromising instruction generalization by introducing irrelevant content information, and promotes semantic alignment of the instruction with the objectives.

To investigate the ability of the visual prompt editing methods in guiding image editing in diverse scenarios, we establish **Transformation-Oriented Paired Benchmark (TOP-Bench)**. This benchmark comprises a total of 750 images, encompassing 25 distinct editing effects, with each effect having 10 pairs of training data and 5 pairs of testing data. The creation of this benchmark not only helps to evaluate the potential of existing methods in guiding image editing, but also paves the way for further research in visual prompt editing. Qualitatively and quantitatively, our method surpasses the existing methods in terms of performance and demonstrates greater semantically consistency with the target editing effects.

In summary, our contributions are threefold:

- We introduce *InstructBrush*, a novel solution to visual prompt editing, which extracts the editing concepts from exemplar image pairs for the subsequent image editing task.

- We propose the *Attention-based Instruction Optimization*, which is optimized within the feature space of the cross-attention, improving the representation ability of instruction guidance, and the *Transformation-oriented Instruction Initialization* to ingeniously introduce semantic-level guidance related to editing.
- We establish **Transformation-Oriented Paired Benchmark (TOP-Bench)** for visual prompt editing to assess its adaptability across diverse scenarios. Both qualitatively and quantitatively, our approach achieves more robust editing and is more semantically consistent with the target editing effects.

2 RELATED WORK

Instruction-based Image Editing. Text-guided diffusion models Nichol et al. (2021); Ramesh et al. (2022); Saharia et al. (2022); Rombach et al. (2022); Podell et al. (2023); Betker et al. (2023); Dai et al. (2023) have taken the world by storm. By leveraging the robust generative priors of these models, InstructPix2Pix (IP2P) Brooks et al. (2023) makes the initial attempt to use a triplet dataset for training a model that edits images based on instructions, achieving intuitive and user-friendly instruction-based image editing. HIVE Zhang et al. (2023b) incorporates reward learning from human feedback to fine-tune IP2P for instruction editing that is more aligned with user preferences. MagicBrush Zhang et al. (2023a) constructs a large-scale manually annotated dataset to fine-tune IP2P, greatly improving the effect in real image editing. Several existing methods, such as InstructDiffusion Geng et al. (2023) and Emu Edit Sheynin et al. (2023) extend instruction-based editing methods to new visual tasks, demonstrating its potential as a universal framework for visual tasks. Recently, some efforts Fu et al. (2023); Huang et al. (2023a) leverage Multimodal Large Language Models (MLLMs) to enhance the performance of instructions, facilitating more accurate editing. Other efforts Simsar et al. (2023); Guo & Lin (2023); Li et al. (2023a) concentration flexible and high-fidelity local editing, addressing the limitations of instruction-based editing in processing local details of images. Additionally, instruction-based image editing has been extended to 3D Chen et al. (2023b) and video Xing et al. (2023) editing tasks, showcasing its tremendous application value.

Visual In-context Learning. In-context learning Brown et al. (2020), which originated from the field of natural language processing (NLP), has been promoted as a learning paradigm. This paradigm enables the execution of a given task on a sample query after learning the task from a set of examples. VisualPrompting Bar et al. (2022) first introduced the concept of visual contextual learning. It uses an inpainting-based approach with grid-like inputs and has shown remarkable results in many tasks. Subsequent works Wang et al. (2023a;b); Fang et al. (2024) broaden the application areas of the framework, such as keypoint detection Wang et al. (2023a), image denoising Wang et al. (2023a), image segmentation Wang et al. (2023b) and 3D point cloud Fang et al. (2024). Recent works Wang et al. (2024); Chen et al. (2023c) introduce in-context learning on diffusion models to accomplish various visual tasks, but they require guidance from textual instructions. Yang et al. (2024); Gu et al. (2024) models visual transformations as a diffusion-based inpainting problem. However, it still requires grid images as input, which poses a significant burden when processing high-resolution images. Unlike these methods, Visii Nguyen et al. (2023) focuses on editing tasks. It inverses exemplar image pairs into a text instruction within an instruction-based image editing model, replacing textual instructions to guide the editing of new images. Our approach similarly focuses on image editing based on instruction inversion and achieves more robust editing and generalization ability to new scenarios.

3 PRELIMINARIES

Latent Diffusion Models. Stable Diffusion (SD), a variant of the latent diffusion model (LDM) Rombach et al. (2022), serves as a text-guided diffusion model. To generate high-resolution images while enhancing computational efficiency in the training process, it employs a pre-trained variational autoencoder (VAE) encoder $\mathcal{E}(\cdot)$ to map images into latent space and perform an iterative denoising process. Subsequently, the predicted images is mapping back into pixel space through the pre-trained VAE decoder $\mathcal{D}(\cdot)$. For each denoising step, the simplified optimization objective is defined as follows:

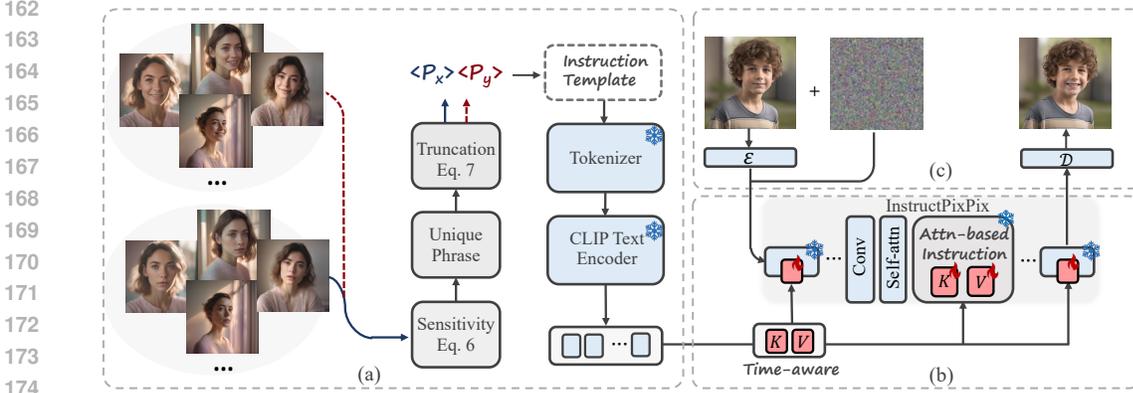


Figure 2: **The Framework of InstructBrush.** InstructBrush inverts instructions from exemplar image pairs by proposing novel (a) *Transformation-oriented Instruction Initialization* and (b) *Attention-based Instruction Optimization* modules. After optimization, the learned instructions are used to guide the editing of new images (c).

$$L_{LDM}(\theta) := \mathbb{E}_{\mathcal{E}(x), \epsilon, t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \tau_{\theta}(c))\|_2^2 \right]. \quad (1)$$

In this process, the text description c is first tokenized into textual embeddings by a Tokenizer. The textual embeddings are then passed through the CLIP text encoder $\tau_{\theta}(\cdot)$ to obtain text conditions. The resulting text conditions are used to guide the diffusion denoising process.

InstructPix2Pix. InstructPix2Pix (IP2P) Brooks et al. (2023) is an instruction-guided image editing method. After encoding the input image c_I using the VAE encoder, IP2P concatenates the noisy latent z_t with the encoded latent $\mathcal{E}(c_I)$ in the first convolutional layer of SD. Subsequently, it uses a generated triplet dataset to perform instruction tuning Wei et al. (2021) on the improved network. This method maximizes the utilization of SD’s powerful generative prior, thereby enabling stunning image editing based on human instructions c_T . The simplified denoising optimization objective is defined by:

$$L_{IP2P}(\theta) := \mathbb{E}_{\mathcal{E}(x), \mathcal{E}(c_I), c_T, \epsilon, t} \left[\|\epsilon - \epsilon_{\theta}(z_t, t, \mathcal{E}(c_I), c_T)\|_2^2 \right]. \quad (2)$$

The dual conditional framework of IP2P employs both input image I and text instruction t for guidance, achieved through an enhanced classifier-free guidance (CFG) strategy Ho & Salimans (2022). The improved CFG incorporates two distinct guidance scales, s_T and s_I , adjustable to balance guidance strength between text and image conditions. It learns the score estimate predicted by the network corresponding to a single denoising step as follows:

$$\begin{aligned} \tilde{e}_{\theta}(z_t, c_I, c_T) &= e_{\theta}(z_t, \emptyset, \emptyset) \\ &\quad + s_I \cdot (e_{\theta}(z_t, c_I, \emptyset) - e_{\theta}(z_t, \emptyset, \emptyset)) \\ &\quad + s_T \cdot (e_{\theta}(z_t, c_I, c_T) - e_{\theta}(z_t, c_I, \emptyset)). \end{aligned} \quad (3)$$

4 METHOD

The pipeline of *InstructBrush* is demonstrated in Figure 2. Based on the instruction-based image editing methods Brooks et al. (2023), *InstructBrush* inverts exemplar image pairs as editing instructions and applies them to editing new images. It proposes novel *Attention-based Instruction Optimization* and *Transformation-oriented Instruction Initialization* modules. The former introduces the editing instruction into the cross-attention layers of the instruction-based image editing model and directly optimizes the Keys and Values corresponding to the instruction within these layers, facilitating more effective instruction inversion (Section 4.1). The latter introduces semantic-level guidance related to editing, ingeniously separates editing-related information from the content of visual prompts and incorporates it into the learned instructions. This effectively promotes semantic alignment of the instruction with the objectives. (Section 4.2).

4.1 ATTENTION-BASED INSTRUCTION OPTIMIZATION

Inspired by Textual Inversion Gal et al. (2022a), The current instruction inversion method Nguyen et al. (2023) optimizes the embeddings of the text encoder using image pairs, aiming to represent the transformation effects between image pairs in textual space. However, the text encoder is trained on text-image pairs with rough descriptions, and its feature space is prone to losing the detailed representation of the image Chen et al. (2023d). Therefore, it is difficult to achieve the requirement of only optimizing the instruction that represents the target transformation in this space. Instead, we focus on optimizing the features in cross-attention layer of the diffusion model. These features are projected from textual embeddings to representations consistent with image features, enabling a more precise representation of image transformation details Hertz et al. (2022); Simsar et al. (2023). As a result, we introduce an attention-based instruction optimization that optimizes editing instructions in the image feature space of the cross-attention layers in the diffusion model, fostering more effective instruction inversion.

Attention-based Instruction. Considering a single-head cross-attention, let Q be the query, K, V be the keys and values from the instruction, respectively, the cross-attention is given by:

$$\text{Attention}(Q, K, V) = \text{Softmax}\left(\frac{QK^T}{\sqrt{d'}}\right)V. \quad (4)$$

Here, $K, V \in \mathbb{R}^{l \times d}$, where l represents the token length of the instruction, and d represents the feature dimension, the value of which depends on the position of the cross-attention layer in the U-Net framework. We optimize the features $\gamma_K, \gamma_V \in \mathbb{R}^{m \times d}$ with a length of $m \in l$ in the key and value corresponding to the first m tokens of the text instruction. Because after linear projection, instruction embeddings transform from text embedding to image features, exhibiting stronger image representation capabilities. To optimize the feature embeddings of the editing instruction, our optimization objective is derived from the simplified least squares error in Eq. 2:

$$\gamma = \arg \min \mathbb{E}_{\mathcal{E}(x), \mathcal{E}(c_I), c_T, \epsilon, t} \left[\|\epsilon - \epsilon_\theta(z_t, t, \mathcal{E}(c_I), c_T)\|_2^2 \right]. \quad (5)$$

Here, $\gamma = \{\gamma_K, \gamma_V\}_{1 \dots n}$ represents the features of keys and values from the first m tokens of the text instruction in all n cross-attention layers. The value of m corresponds to the number of text tokens used for instruction initialization, as described in Section 4.2.

Time-aware Instruction (Optional). In the text-guided diffusion models, the denoising process focuses on image generation from low-frequency structure to high-frequency details Daras & Dimakis (2022); Zhang et al. (2023c). We believe that a similar property also exists in instruction-based editing models, where different denoising processes primarily focus on distinct transformations. We confirm this view in Figure 15. Therefore, we divide the instruction optimization equally into j parts based on denoising time steps, emphasizing instruction learning within the editing-related denoising time steps. Now, we have $\gamma = \{\gamma_K, \gamma_V\}_{1 \dots n}^j$, where j is 5 by default. In this way, the learned instructions can capture more details of transformations, which can guide the editing of new images more robustly.

4.2 TRANSFORMATION-ORIENTED INSTRUCTION INITIALIZATION

Concept inversion Gal et al. (2022a); Voynov et al. (2023); Zhang et al. (2023c) uses the semantic class word (e.g., dog, cat) for initialization, providing prior information for the target concept learning. However, instruction inversion requires learning a sentence as an instruction that describes the *image transformation*. Manually initializing a sentence based on the transformation of reference image pairs is not only laborious but also subjective. The existing work Nguyen et al. (2023) utilizes the caption method Wen et al. (2023) to obtain the caption of after-editing images in the training set as the start point of the optimization. Despite the introduction of transformation-related prior knowledge, it simultaneously introduce editing-irrelevant content information about the training scenario, hindering the generalization of instruction to new scenarios. In addition, although existing multimodal large language models (MLLM) can directly compare two images to obtain a description of the differences, the daunting model size and lack of prior knowledge in professional vocabulary have caused certain obstacles in its practical application. In contrast, our approach extracts transformation-related information in a simpler and more effective way. Specifically, we first extract *unique phrases* that differentiate the images before and after editing as editing-related priors. Subsequently, we incorporate them into the *instruction template* for instruction initialization.

Unique Phrase Extraction. Given a set of image pairs $\{\{x\}, \{y\}\}$, where $\{x\}$ and $\{y\}$ represent the image sets before and after editing, for a single set $\{x\}$, we search for the fixed-length phrase set $P_x = \{\langle p_1 \rangle, \dots, \langle p_r \rangle\}$ with the highest cosine similarity between image and text features. Here, $\langle p_i \rangle$ represents a text phrase from a vocabulary set, which can be customized according to the task domain or use a public vocabulary set pha (2022). And r represents the adjustable number of phrases to form the caption, which is set to 5 by default. Subsequently, we compare the feature similarity between P_x and the image sets $\{x\}$, $\{y\}$ respectively, and then measure the difference in feature similarity of the same phrase with the two sets as the *sensitivity*. This process can be represented as follows:

$$sens_i(\langle p_i \rangle) = sim(\langle p_i \rangle, \{x\}) - sim(\langle p_i \rangle, \{y\}) \quad (6)$$

Here, $sens_i$ denotes the sensitivity of the i th phrase in P_x and sim denotes the CLIP feature similarity. We identify the phrase with maximum sensitivity as the *unique phrase* $\langle p_x \rangle$ of the set $\{x\}$. However, there exist certain edits whose editing-related information cannot be recognized. To avoid the unique phrase containing editing-irrelevant information, we define the truncation conditions:

$$\langle p_x \rangle = \begin{cases} \langle p_x \rangle & \text{if } sens(\langle p_x \rangle) \geq \eta \\ \emptyset & \text{otherwise,} \end{cases} \quad (7)$$

where η represents a constant that controls the truncation of unique phrases, set to 0.15 by default.

Instruction Template. With the above method, we can get the unique phrases $\langle p_x \rangle$ and $\langle p_y \rangle$ for sets $\{x\}$ and $\{y\}$. Then we incorporate them into the instruction template. The form of the instruction template is strictly aligned with the base model’s editing instructions to maximize the use of the textual prior. For example, we use “turn $\langle p_x \rangle$ into $\langle p_y \rangle$ ” as a starting point for instruction optimization. Note that when $\langle p_y \rangle = \emptyset$, we use *None* instruction for initialization and optimize fixed-length features for Keys and Values in cross-attention. Although the initialized instruction is not sufficient to express the target editing effect, it can introduce prior knowledge of transformation, aiding the semantics of learned instruction to be close to the target.

5 TRANSFORMATION-ORIENTED PAIRED BENCHMARK

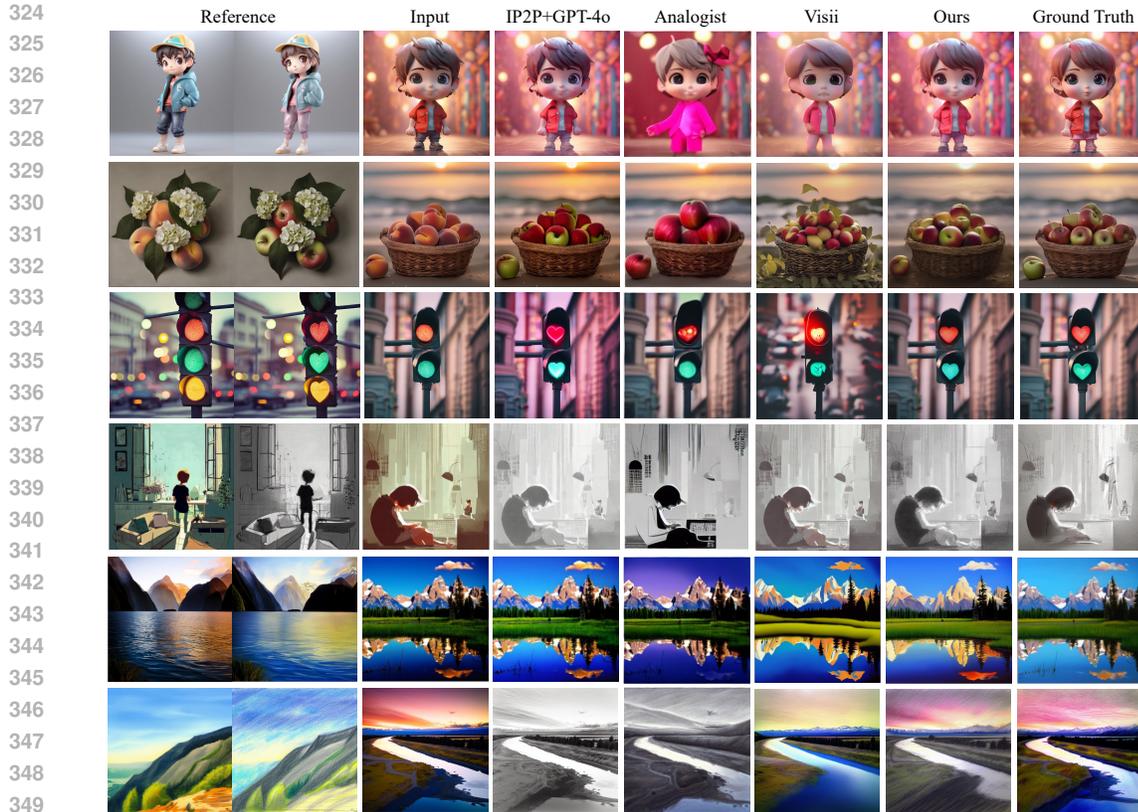
To investigate the editing capabilities of various instruction inversion methods in open scenarios and facilitate a fair comparison of these methods, We establish a benchmark named *TOP-Bench* (Transformation-Oriented Paired Benchmark), which can be utilized for both qualitative and quantitative evaluations. Our benchmark contains few-shot rather than one-shot datasets because of the effect of image transformation that is difficult to fully visualize with a single image pair. It spans 25 datasets corresponding to different editing effects. It covers a wide range of editing categories and scenarios, allowing for division from multiple dimensions. Each dataset consists of 10 pairs of training images and 5 pairs of testing images, totaling 750 images. Additionally, we provide text instructions aligned with the transformation effects for each dataset. Please refer to the Supplementary for data acquisition and detailed introduction.

To further analyze the advantages of our method, we categorize the benchmark into two different categories: TOP-Global and TOP-Local, corresponding to datasets of 14 global editing effects and 11 local editing effects, respectively. We compare the quantitative results of different methods in these two categories to validate the effectiveness of our method.

6 EXPERIMENTS

In this section, we present qualitative and quantitative results. The implementation details of our method are detailed in Appendix C. Since the effects of image transformations are difficult to visualize completely from individual image pairs, we focus on the analysis of experiments in the few-shot setting in this chapter. Additional one-shot experiments are shown in Appendix F.1.

Metrics. We use several objective evaluation metrics on the benchmark. Specifically, we employ full-reference quality metrics PSNR, SSIM Wang et al. (2004), LPIPS Zhang et al. (2018), CLIP image similarity score and DINO score to assess the consistency between the generated images and the ground truth, quantifying the image editing capabilities of each method. In addition, we measure



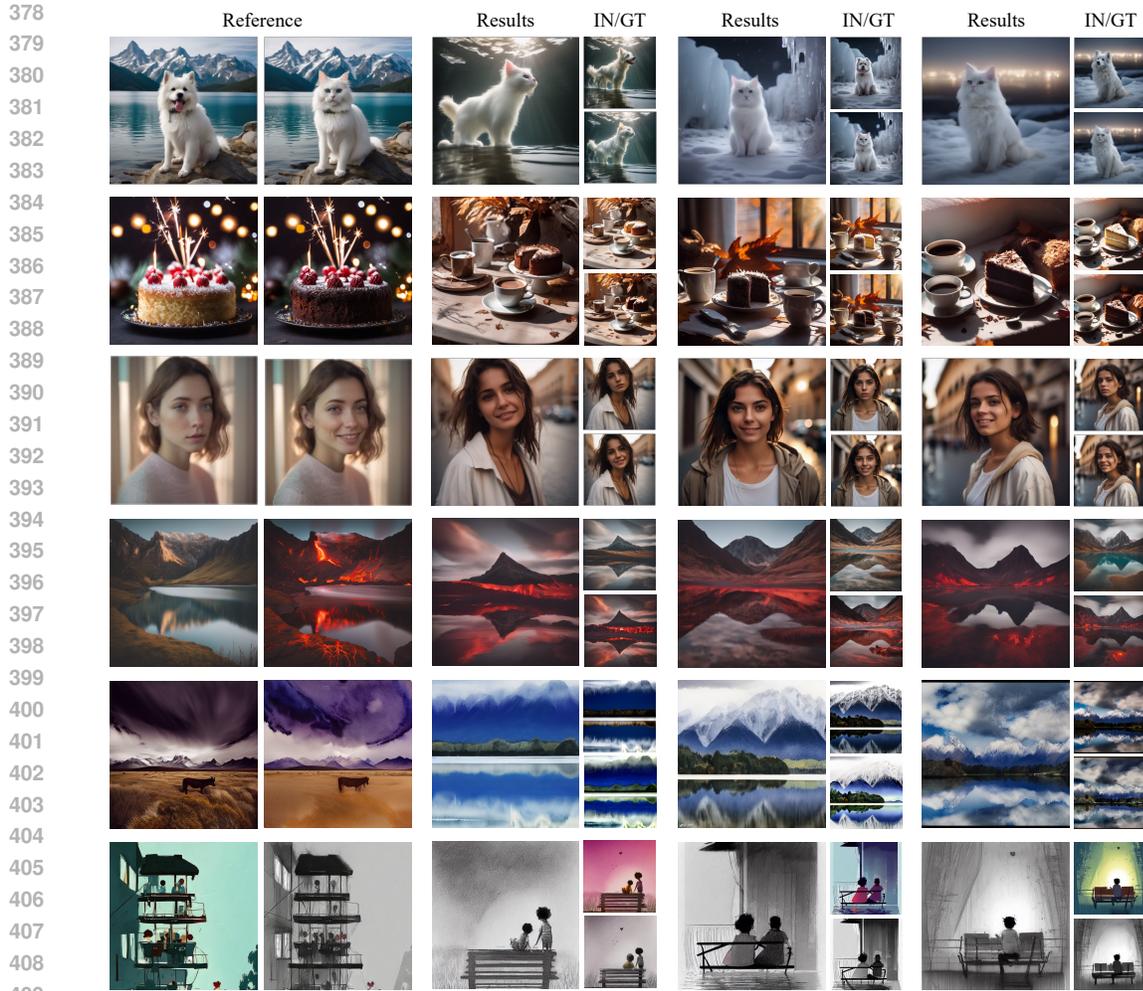
351 **Figure 3: Qualitative Comparisons with Existing Methods.** Our method achieves superior per-
 352 formance in both local and global image editing. It effectively avoids introducing editing-irrelevant
 353 information from the training images, showing better instruction generalization.

354 the CLIP directional similarity Gal et al. (2022b) between image pairs to evaluate the semantic
 355 alignment between the editing direction of each method and the target. Specifically, we measure the
 356 consistency between the average editing direction from the input images to the generated images and
 357 the average direction of the training image pairs, see Appendix C for more details. Additionally, we
 358 compared the runtime of different methods, see the appendix for more details.

359 **Compared Methods.** We compare our *InstructBrush* with the state-of-the-art competitor Visii
 360 Nguyen et al. (2023) and Analogist Gu et al. (2024). Considering that the multimodal large language
 361 model (MLLM) can compare the differences between reference image pairs to obtain editing instruc-
 362 tions, which can be used as the input of the text-guided editing model, we introduce GPT-4o-based
 363 IP2P Brooks et al. (2023) for comparison. We use an image resolution of 512×512 for comparison
 364 with other methods. For Visii and Analogist, we utilize its official implementation, while for IP2P,
 365 we employ its Diffusers von Platen et al. (2022) version. All experiments are conducted following
 366 the official recommended configurations.

367 6.1 COMPARISONS

368 **Qualitative Comparisons.** We use our *TOP-Bench* to evaluate the results of different methods.
 369 For instruction inversion methods, we employ 10 reference before-and-after editing image pairs
 370 to optimize the instruction for each editing effect. For IP2P and Analogist, we leverage GPT-
 371 4o to compare differences between reference image pairs to obtain textual captions for editing.
 372 Subsequently, we present a comparison of the results of editing the test images in Figure 3. Since
 373 the editing effects of IP2P and Analogist are mainly affected by the text conditional prior of the
 374 diffusion model rather than directly from the reference image pairs, they show certain deviations
 375 when analogizing the editing effects of the reference images. IP2P employs text instructions to
 376 guide image editing. Although such text-based editing approach cannot accurately extract the editing
 377 concepts between image pairs, the generalization ability of text and the model priors of IP2P ensure



410 **Figure 4: More Visualization Results of Our Method.** Our method demonstrates robust performance
411 on both local and global editing. And it does not introduce scene information of the training image
412 when editing new images, which reflects the instruction generalization of our method.
413

414
415 the quality of the generated images. In contrast, Analogist leverages the priors of the inpainting
416 diffusion model, and compared to IP2P, it has a lesser understanding of the editing instructions.
417 Additionally, the extra structural constraints imposed on the attention further exacerbate its lower
418 adherence to the instructions. For example, suboptimal results are observed in the local edits from
419 row 1 to row 3. Although Visii optimizes instructions to learn the target editing concept and solves
420 the problem of IP2P not being able to specifically represent image changes using text instructions
421 alone, its content-oriented initialization reduces the instructions generalization. It can easily introduce
422 content information in the training image during the instruction editing process, as shown in rows 2
423 and 3. In addition, the limitations of optimization space also make it difficult to accurately learn target
424 editing concepts. By contrast, our *InstructBrush* demonstrates superior editing performance. Fig. 4
425 illustrates more qualitative results obtained by our method, which demonstrates robust performance
426 on both local and global editing. And it does not introduce scene information of the training image
427 when editing new images, which reflects the instruction generalization of our method. we present
428 additional results in Appendix F, encompassing one-shot and real-world images evaluations ¹, which
429 further substantiate the robust performance of our method.

429 **Quantitative Comparisons.** We conduct a detailed quantitative evaluation of these methods on
430 TOP-Local, TOP-Global, and overall TOP-Bench. As shown in Table 1, the editing performance
431

¹We conduct one-shot evaluations, real-world images evaluations and visualization of applications.

Table 1: **Quantitative Results.** We measure the performance of our method against several other methods based on average PSNR, SSIM, LPIPS, CLIP direction score (CLIP-D), CLIP image similarity score (CLIP-I), and DINO score. Our approach offers a basic version without the use of Time-aware Instructions as well as a complete version that utilizes them.

Datasets	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP-D \downarrow	CLIP-I \uparrow	DINO \uparrow
TOP-Global	IP2P+GPT-4o Brooks et al. (2023)	13.53	0.4878	0.3884	0.6897	0.8676	0.8959
	Analogist Gu et al. (2024)	12.28	0.3860	0.3993	0.8025	0.8047	0.8800
	Visii Nguyen et al. (2023)	15.87	0.4947	0.3866	0.3938	0.8471	0.8767
	Ours (basic)	17.51	0.5509	0.3005	0.2814	0.8541	0.8812
	Ours	18.66	0.5842	0.2526	0.2798	0.9127	0.9354
TOP-Local	IP2P+GPT-4o Brooks et al. (2023)	17.33	0.7016	0.2738	0.6188	0.8844	0.9072
	Analogist Gu et al. (2024)	14.03	0.5120	0.3334	0.8559	0.8535	0.9021
	Visii Nguyen et al. (2023)	18.76	0.7157	0.2585	0.2560	0.8695	0.9024
	Ours (basic)	22.36	0.8115	0.1337	0.3032	0.8887	0.9238
	Ours	23.26	0.8297	0.1143	0.3576	0.9100	0.9486
TOP-Bench	IP2P+GPT-4o Brooks et al. (2023)	15.20	0.5819	0.3380	0.6585	0.8750	0.9009
	Analogist Gu et al. (2024)	13.05	0.4414	0.3703	0.8260	0.8262	0.8897
	Visii Nguyen et al. (2023)	17.14	0.5919	0.3303	0.3332	0.8569	0.8880
	Ours (basic)	19.64	0.6656	0.2271	0.2910	0.8693	0.8999
	Ours	20.68	0.6922	0.1918	0.3140	0.9115	0.9412

of our method surpasses that of other methods at both the editing effects and semantic alignment. In addition, compared to the results of Visii, our method shows a more significant improvement on TOP-Local than on TOP-Global. This is because in local editing tasks, training images contain more editing-irrelevant scene information. The content-oriented initialization of Visii introduces them to the initialized instructions, posing a greater obstacle to optimization. On the contrary, our transformation-oriented instruction initialization method can accurately capture the transformations between image pairs and use them for initialization, thus improving instruction generalization.

6.2 ABLATION STUDIES

Attention-based Instruction Ablation. The use of attention-based instruction aims to avoid the limitation of CLIP space on the representation ability of target transformations and achieve a more accurate representation of image transformation details. The metrics PSNR, SSIM, and LPIPS are calculated between the output and the ground truth to evaluate the editing performance. We report results in Table 2 and observe that adopting attention-based instructions replaced with CLIP space-based instructions effectively improves the editing performance of the instructions. Additionally, we also observe in Figure 5 that compared to inversion in CLIP space, optimizing instruction in attention space has shown significant improvements in editing.

Transformation-oriented Instruction Initialization Ablation. Content-oriented initialization methods introduce irrelevant content information from the training images, thereby interfering with the optimization process. As depicted in Figure 5, the use of the content-oriented initialization method results in the leakage of content information from the training image into the edited image. By enabling instruction initialization to prioritize image changes over image content, it not only enhances the editing capabilities of learned instructions, but also aligns the edited image with the target transformation in terms of semantic information, which is confirmed in Table 2.

Time-aware Instruction Ablation (Optional). The use of time-aware instructions facilitates instruction optimization by allowing instructions to focus on learning different transformations at different denoising time steps. Table 2 explicitly shows that the use of time-aware instruction helps to improve the editing effect. The same result is confirmed in Figure 5. *Note that in the second row of Figure 5, the reason the time-aware instruction does not show significant improvement is because this type of editing is relatively simple, and using other modules is sufficient to achieve such editing effects. The fine-grained facial glasses editing in the first row demonstrates the importance of this setting for fine-grained editing. Additionally, we present more qualitative ablation results in Figure 12, visualizing the importance of the module’s design.* However, the use of this module significantly increases the optimization time. Therefore, this module is discarded in the basic version of our method in exchange for shorter optimization time.



Figure 5: **Visualization Results of Ablation Study.** We visualize the independent effects of our proposed attention-based instruction, time-aware instruction, and transition-oriented instruction initialization on the results, intuitively highlighting the importance of these configurations.

Table 2: **Ablation Study.** We validate the independent impact of our proposed attention-based instruction, time-aware instruction, and transformation-oriented instruction initialization on results, emphasizing the importance of these configurations.

Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP-D \downarrow	CLIP-I \uparrow	DINO \uparrow
w/o Attn	19.56	0.6709	0.2179	0.3124	0.8455	0.8757
w/o Init	20.22	0.6841	0.2018	0.3747	0.8993	0.9231
w/o Time (Ours-fast)	19.64	0.6656	0.2271	0.2910	0.8693	0.8999
Ours	20.68	0.6922	0.1918	0.3140	0.9115	0.9412

7 LIMITATIONS AND CONCLUSION

Although optimization-based methods represented by our framework are easier to learn editing concepts from multiple pairs of reference images, they increase the time cost in training. In addition, our initialization method is limited by the vocabulary used to search for unique phrases. If the phrase is not present in the vocabulary, our initialization method will initialize using *None* instruction, which will not introduce any editing prior.

Our method extracts editing effects from image pairs for editing tasks that are difficult for users to describe. It introduces a new instruction optimization and initialization method, achieving better instruction optimization and generalization. Numerous experiments have demonstrated the advantages of our method. In the future, we will apply our method to more powerful instruction-based image editing models for more robust editing performance. We hope that this work will stimulate more research and serve as a prior extraction method to aid in the training of downstream tasks.

REFERENCES

- Clip-interrogator. <https://github.com/pharmapsychotic/clip-interrogator>, 2022.
- Yuval Alaluf, Elad Richardson, Gal Metzer, and Daniel Cohen-Or. A neural space-time representation for text-to-image personalization. *ACM Transactions on Graphics (TOG)*, 42(6):1–10, 2023.
- Moab Arar, Rinon Gal, Yuval Atzmon, Gal Chechik, Daniel Cohen-Or, Ariel Shamir, and Amit H. Bermano. Domain-agnostic tuning-encoder for fast personalization of text-to-image models. In *SIGGRAPH Asia 2023 Conference Papers*, pp. 1–10, 2023.
- Amir Bar, Yossi Gandelsman, Trevor Darrell, Amir Globerson, and Alexei Efros. Visual prompting via image inpainting. *Advances in Neural Information Processing Systems*, 35:25005–25017, 2022.
- James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*. <https://cdn.openai.com/papers/dall-e-3.pdf>, 2(3):8, 2023.

- 540 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image
541 editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
542 *Recognition*, pp. 18392–18402, 2023.
- 543 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,
544 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are
545 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 546 Li Chen, Mengyi Zhao, Yiheng Liu, Mingxu Ding, Yangyang Song, Shizun Wang, Xu Wang, Hao
547 Yang, Jing Liu, Kang Du, et al. Photoverse: Tuning-free image customization with text-to-image
548 diffusion models. *arXiv preprint arXiv:2309.05793*, 2023a.
- 549 Minghao Chen, Junyu Xie, Iro Laina, and Andrea Vedaldi. Shap-editor: Instruction-guided latent 3d
550 editing in seconds. *arXiv preprint arXiv:2312.09246*, 2023b.
- 551 Tianqi Chen, Yongfei Liu, Zhendong Wang, Jianbo Yuan, Quanzeng You, Hongxia Yang, and
552 Mingyuan Zhou. Improving in-context learning in diffusion models with visual context-modulated
553 prompts. *arXiv preprint arXiv:2312.01408*, 2023c.
- 554 Xi Chen, Lianghua Huang, Yu Liu, Yujun Shen, Deli Zhao, and Hengshuang Zhao. Anydoor:
555 Zero-shot object-level image customization. *arXiv preprint arXiv:2307.09481*, 2023d.
- 556 Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon
557 Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. Emu: Enhancing image generation
558 models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807*, 2023.
- 559 Giannis Daras and Alex Dimakis. Multiresolution textual inversion. In *NeurIPS 2022 Workshop on*
560 *Score-Based Methods*, 2022.
- 561 Zhongbin Fang, Xiangtai Li, Xia Li, Joachim M Buhmann, Chen Change Loy, and Mengyuan Liu.
562 Explore in-context learning for 3d point cloud understanding. *Advances in Neural Information*
563 *Processing Systems*, 36, 2024.
- 564 Tsu-Jui Fu, Wenze Hu, Xianzhi Du, William Yang Wang, Yinfei Yang, and Zhe Gan. Guid-
565 ing instruction-based image editing via multimodal large language models. *arXiv preprint*
566 *arXiv:2309.17102*, 2023.
- 567 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel
568 Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual
569 inversion. *arXiv preprint arXiv:2208.01618*, 2022a.
- 570 Rinon Gal, Or Patashnik, Haggai Maron, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or.
571 Stylegan-nada: Clip-guided domain adaptation of image generators. *ACM Transactions on*
572 *Graphics (TOG)*, 41(4):1–13, 2022b.
- 573 Rinon Gal, Moab Arar, Yuval Atzmon, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or.
574 Encoder-based domain tuning for fast personalization of text-to-image models. *ACM Transactions*
575 *on Graphics (TOG)*, 42(4):1–13, 2023.
- 576 Zigang Geng, Binxin Yang, Tiankai Hang, Chen Li, Shuyang Gu, Ting Zhang, Jianmin Bao, Zheng
577 Zhang, Han Hu, Dong Chen, et al. Instructdiffusion: A generalist modeling interface for vision
578 tasks. *arXiv preprint arXiv:2309.03895*, 2023.
- 579 Zheng Gu, Shiyuan Yang, Jing Liao, Jing Huo, and Yang Gao. Analogist: Out-of-the-box visual
580 in-context learning with image diffusion model. *ACM Transactions on Graphics (TOG)*, 43(4):
581 1–15, 2024.
- 582 Qin Guo and Tianwei Lin. Focus on your instruction: Fine-grained and multi-instruction image
583 editing by attention modulation. *arXiv preprint arXiv:2312.10113*, 2023.
- 584 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-
585 to-prompt image editing with cross-attention control. In *The Eleventh International Conference on*
586 *Learning Representations*, 2022.

- 594 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
595 2022.
- 596
- 597 Yuzhou Huang, Liangbin Xie, Xintao Wang, Ziyang Yuan, Xiaodong Cun, Yixiao Ge, Jiantao Zhou,
598 Chao Dong, Rui Huang, Ruimao Zhang, et al. Smartedit: Exploring complex instruction-based
599 image editing with multimodal large language models. *arXiv preprint arXiv:2312.06739*, 2023a.
- 600
- 601 Ziqi Huang, Tianxing Wu, Yuming Jiang, Kelvin CK Chan, and Ziwei Liu. Reversion: Diffusion-
602 based relation inversion from images. *arXiv preprint arXiv:2303.13495*, 2023b.
- 603
- 604 Chuck Jacobs, D Salesin, N Oliver, A Hertzmann, and AB Curless. Image analogies. In *Proceedings*
605 *of Siggraph*, pp. 327–340, 2001.
- 606
- 607 Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. Direct inversion: Boosting
608 diffusion-based editing with 3 lines of code. In *The Twelfth International Conference on Learning*
Representations, 2023.
- 609
- 610 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
611 based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577,
612 2022.
- 613
- 614 Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolu-
615 tional neural networks. *Advances in neural information processing systems*, 25, 2012.
- 616
- 617 Shanglin Li, Bohan Zeng, Yutang Feng, Sicheng Gao, Xuhui Liu, Jiaming Liu, Li Lin, Xu Tang,
618 Yao Hu, Jianzhuang Liu, et al. Zone: Zero-shot instruction-guided local editing. *arXiv preprint*
arXiv:2312.16794, 2023a.
- 619
- 620 Zhen Li, Mingdeng Cao, Xintao Wang, Zhongang Qi, Ming-Ming Cheng, and Ying Shan. Photomaker:
621 Customizing realistic human photos via stacked id embedding. *arXiv preprint arXiv:2312.04461*,
622 2023b.
- 623
- 624 Jie Liang, Hui Zeng, Miaomiao Cui, Xuansong Xie, and Lei Zhang. Ppr10k: A large-scale portrait
625 photo retouching dataset with human-region mask and group-level consistency. In *Proceedings of*
the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 653–661, 2021.
- 626
- 627 Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for
628 editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference*
on Computer Vision and Pattern Recognition (CVPR), pp. 6038–6047, 2023.
- 629
- 630 Saman Motamed, Danda Pani Paudel, and Luc Van Gool. Lego: Learning to disentangle and
631 invert concepts beyond object appearance in text-to-image diffusion models. *arXiv preprint*
arXiv:2311.13833, 2023.
- 632
- 633 Thao Nguyen, Yuheng Li, Utkarsh Ojha, and Yong Jae Lee. Visual instruction inversion: Image
634 editing via image prompting. In *Thirty-seventh Conference on Neural Information Processing*
Systems, 2023.
- 635
- 636 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,
637 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with
638 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 639
- 640 Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu.
641 Zero-shot image-to-image translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, pp.
642 1–11, 2023.
- 643
- 644 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
645 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
646 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 647
- Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.

- 648 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
649 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference*
650 *on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 651
652 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
653 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
654 text-to-image diffusion models with deep language understanding. *Advances in Neural Information*
655 *Processing Systems*, 35:36479–36494, 2022.
- 656 Shelly Sheynin, Adam Polyak, Uriel Singer, Yuval Kirstain, Amit Zohar, Oron Ashual, Devi Parikh,
657 and Yaniv Taigman. Emu edit: Precise image editing via recognition and generation tasks. *arXiv*
658 *preprint arXiv:2311.10089*, 2023.
- 659 Enis Simsar, Alessio Tonioni, Yongqin Xian, Thomas Hofmann, and Federico Tombari. Lime:
660 Localized image editing via attention regularization in diffusion models, 2023.
- 661
662 Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for
663 text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer*
664 *Vision and Pattern Recognition (CVPR)*, pp. 1921–1930, 2023.
- 665 Yael Vinker, Andrey Voynov, Daniel Cohen-Or, and Ariel Shamir. Concept decomposition for visual
666 exploration and inspiration. *ACM Transactions on Graphics (TOG)*, 42(6):1–13, 2023.
- 667
668 Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul,
669 Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022.
- 670
671 Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. $p+$: Extended textual condi-
672 tioning in text-to-image generation. *arXiv preprint arXiv:2303.09522*, 2023.
- 673
674 Xinlong Wang, Wen Wang, Yue Cao, Chunhua Shen, and Tiejun Huang. Images speak in images: A
675 generalist painter for in-context visual learning. In *Proceedings of the IEEE/CVF Conference on*
676 *Computer Vision and Pattern Recognition*, pp. 6830–6839, 2023a.
- 677 Xinlong Wang, Xiaosong Zhang, Yue Cao, Wen Wang, Chunhua Shen, and Tiejun Huang. Seggpt:
678 Segmenting everything in context. *arXiv preprint arXiv:2304.03284*, 2023b.
- 679
680 Zhendong Wang, Yifan Jiang, Yadong Lu, Pengcheng He, Weizhu Chen, Zhangyang Wang, Mingyuan
681 Zhou, et al. In-context learning unlocked for diffusion models. *Advances in Neural Information*
682 *Processing Systems*, 36, 2024.
- 683 Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from
684 error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612,
685 2004.
- 686 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
687 Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International*
688 *Conference on Learning Representations*, 2021.
- 689
690 Yuxiang Wei, Yabo Zhang, Zhilong Ji, Jinfeng Bai, Lei Zhang, and Wangmeng Zuo. Elite: Encoding
691 visual concepts into textual embeddings for customized text-to-image generation. *arXiv preprint*
692 *arXiv:2302.13848*, 2023.
- 693
694 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein.
695 Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery.
696 *arXiv preprint arXiv:2302.03668*, 2023.
- 697
698 Zhen Xing, Qi Dai, Zihao Zhang, Hui Zhang, Han Hu, Zuxuan Wu, and Yu-Gang Jiang. Vidiff: Trans-
699 lating videos via multi-modal instructions with diffusion models. *arXiv preprint arXiv:2311.18837*,
700 2023.
- 701
702 Sihan Xu, Yidong Huang, Jiayi Pan, Ziqiao Ma, and Joyce Chai. Inversion-free image editing with
703 language-guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer*
704 *Vision and Pattern Recognition*, pp. 9452–9461, 2024.

- 702 Yifan Yang, Houwen Peng, Yifei Shen, Yuqing Yang, Han Hu, Lili Qiu, Hideki Koike, et al.
 703 Imagebrush: Learning visual in-context instructions for exemplar-based image manipulation.
 704 *Advances in Neural Information Processing Systems*, 36, 2024.
- 705
 706 Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt
 707 adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023.
- 708
 709 Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated
 710 dataset for instruction-guided image editing. *arXiv preprint arXiv:2306.10012*, 2023a.
- 711
 712 Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable
 713 effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on
 computer vision and pattern recognition*, pp. 586–595, 2018.
- 714
 715 Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang,
 716 Silvio Savarese, Stefano Ermon, et al. Hive: Harnessing human feedback for instructional visual
 717 editing. *arXiv preprint arXiv:2303.09618*, 2023b.
- 718
 719 Yuxin Zhang, Weiming Dong, Fan Tang, Nisha Huang, Haibin Huang, Chongyang Ma, Tong-
 720 Yee Lee, Oliver Deussen, and Changsheng Xu. Prospect: Prompt spectrum for attribute-aware
 personalization of diffusion models. *ACM Transactions on Graphics (TOG)*, 42(6):1–14, 2023c.
- 721
 722 Ruoyu Zhao, Mingrui Zhu, Shiyin Dong, Nannan Wang, and Xinbo Gao. Catversion: Concatenating
 723 embeddings for diffusion-based text-to-image personalization. *arXiv preprint arXiv:2311.14631*,
 2023.

724 725 726 A APPENDIX

727 728 B ADDITIONAL RELATED WORK

729
 730 **Diffusion-based Prompt Inversion.** The diffusion-based prompt inversion methods aim to learn the
 731 text prompt from a handful of images describing concepts, thereby guiding the generation of diffusion
 732 models. Textual Inversion Gal et al. (2022a) learns text embeddings corresponding to pseudo-words
 733 to represent the target concepts. The pseudo-words can be combined with free text to guide the
 734 generation of images containing target concepts. Based on their research, some works Daras &
 735 Dimakis (2022); Voynov et al. (2023); Zhang et al. (2023c); Alaluf et al. (2023); Zhao et al. (2023)
 736 explore the effects of different inversion spaces on prompt inversion. Other works Gal et al. (2023);
 737 Wei et al. (2023); Arar et al. (2023); Chen et al. (2023a); Ye et al. (2023); Li et al. (2023b) train
 738 an image encoder based on text inversion to achieve generation guided by a given reference image.
 739 Additionally, ReVersion Huang et al. (2023b) focuses on learning the relation between objects through
 740 contrastive learning. PEZ Wen et al. (2023) inverts hard prompts by projecting learned embeddings
 741 onto adjacent interpretable word embeddings, providing a new solution for image captioning. Vinker
 742 et al. (2023) decomposes a visual concept, allowing users to explore hidden sub-concepts of the
 743 object of interest. Lego Motamed et al. (2023) uses carefully designed prompt learning methods to
 744 learn abstract concepts that are entangled with the subject from few samples. These methods focus
 745 on learning concepts to guide image generation, while our study aims to learn the transformations
 between image pairs to guide image editing.

746 747 C EXPERIMENTAL SETTINGS

748
 749 **Implementation Details.** The implementation is based on one NVIDIA Tesla V100 GPU. We use
 750 public vocabulary set pha (2022) to search unique phrases for instruction initialization. Afterward,
 751 based on pre-trained IP2P Brooks et al. (2023), we optimize the features of the keys and values
 752 corresponding to approximately 10 initialization instruction tokens. Note that our method is not
 753 limited to IP2P and can also be applied to other instruction-based editing models Geng et al. (2023);
 754 Zhang et al. (2023a); Sheynin et al. (2023). We divide the learned instructions into 5 parts according
 755 to the denoising time step, and optimize each part with 1000 steps using a learning rate of 0.001 and
 a batch size of 1, respectively, for a total of 5000 steps. The whole training process takes about 20

minutes. During both training and inference, we adopt a text guidance scale $s_T = 7.5$ and an image guidance scale $s_I = 1.5$. And we use the Euler ancestral sampler with denoising variance schedule Karras et al. (2022) with a sampling step of $T = 20$ during the inference process.

Evaluaiion Metrics. We use six objective evaluation metrics on the benchmark. Specifically, we employ full-reference quality metrics PSNR, SSIM Wang et al. (2004), and LPIPS Zhang et al. (2018) to assess the consistency between the generated images and the ground truth, quantifying the image editing capabilities of each method. Among them, higher PSNR indicates more similarity between the results and the ground truth; higher SSIM indicates that the results are structurally more similar to the ground truth; we implement the evaluation of LPIPS based on AlexNet Krizhevsky et al. (2012), and smaller LPIPS indicates that the results has a better features similarity between the results and the ground truth. We also use CLIP image similarity score and DINO score to assess the consistency between the generated images and the ground truth. In addition, we measure the CLIP directional similarity Gal et al. (2022b) between image pairs to evaluate the semantic alignment between the editing direction of each method and the target. Specifically, we measure the consistency between the average editing direction from the input images to the generated images and the average direction of the training image pairs. **The CLIP image directional similarity Parmar et al. (2023); Nguyen et al. (2023); Gal et al. (2022b) is defined as follows:**

$$1 - \cos(\Delta_{x \rightarrow y}, \Delta_{x' \rightarrow y'}), \quad (8)$$

where $\Delta_{x \rightarrow y}$ is the CLIP direction from the input image to the result image, and $\Delta_{x' \rightarrow y'}$ is the CLIP direction between the reference images.

Num	Name	Instruction	Editing Type	
			Local	Global
1	boy2girl	"make boy and dog into a girl and cat"	✓	
2	midnight	"make it nighttime"		✓
3	sea painting	"turn it into a painting"		✓
4	sketch style	"make the image a pencil sketch"		✓
5	summer	"make it summer"		✓
6	wallpaper	"make it snow"		✓
7	charcoal	"turn it into a charcoal drawing"		✓
8	glasses	"add a pair of glasses"	✓	
9	painting	"Make it a painting"		✓
10	painting snow	"make it snow"		✓
11	pencil sketch	"as a pencil sketch"		✓
12	purple	"make the sky a deep purple"		✓
13	snow	"have it snow"		✓
14	watercolor	"as a watercolor painting"		✓
15	4dboy	"Turn the boy into a girl"	✓	
16	apple	"Turn peaches into apples"	✓	
17	cake	"Make it a chocolate cake"	✓	
18	cloud kitty	"Make the cat into a bear"	✓	
19	dog2cat	"Make the dog into a cat"	✓	
20	juice	"Make it a lemonade"	✓	
21	lava	"Turn it into lava"		✓
22	rain	"Turn the rain into snow"		✓
23	read books	"Make newspapers into books"	✓	
24	smile	"Add a smile"	✓	
25	traffic lights	"make it a heart-shaped light"	✓	

Table 3: **Benchmark Presentation.** The benchmark has a total of 25 editing effects, evenly covering both local and global editing.

D BENCHMARK CONSTRUCTION

In recent years, there has been rapid development in text-guided image editing methods. The evaluation of image editing effectiveness has also evolved. Initially, the editing effect is solely

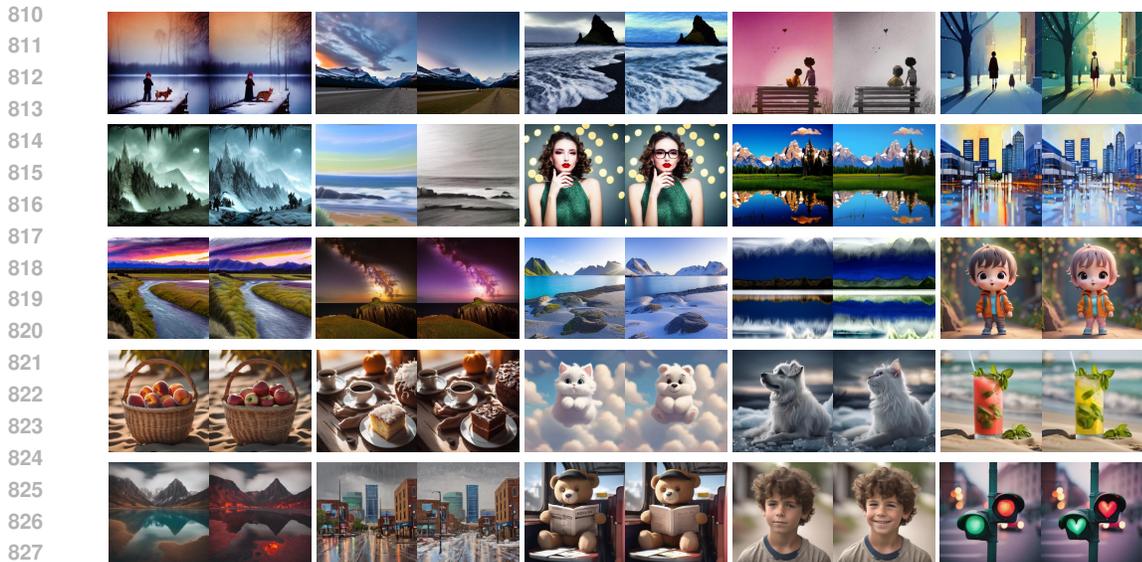


Figure 6: **Visualization of Our Benchmark.** Our benchmark spans 25 datasets corresponding to different editing effects. It covers a wide range of editing categories and scenarios, allowing for division from multiple dimensions. Each dataset consists of 10 pairs of training images and 5 pairs of testing images, totaling 750 images. We show a pair of before-and-after transformation examples for each editing effect.

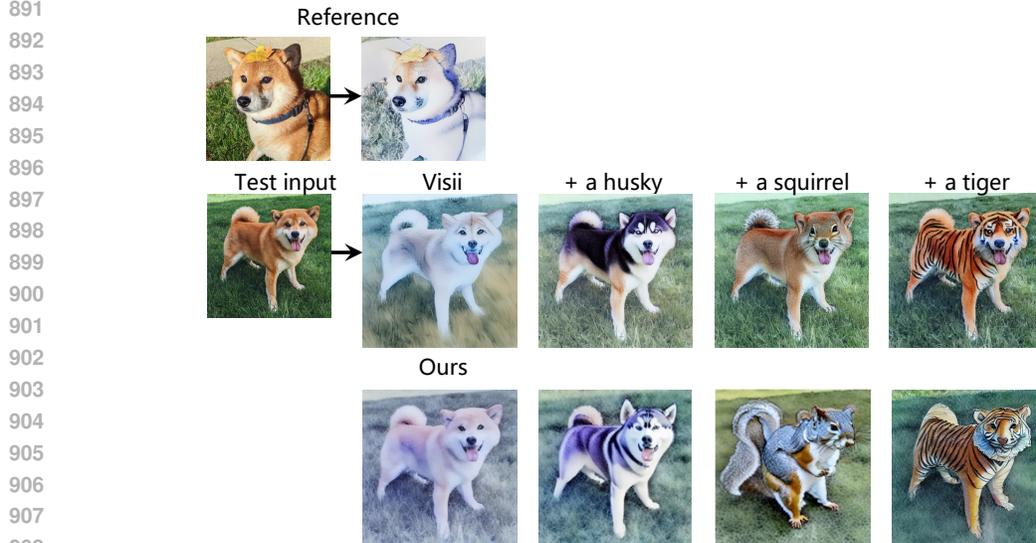
evaluated through qualitative presentations and user study Hertz et al. (2022); Mokady et al. (2023), which led to significant subjectivity. Subsequently, PNP Tumanyan et al. (2023) establishes a benchmark for text-guided image editing, which assesses the performance of text-based image editing methods using text-image and image-image feature similarity scores. Later, Direct Inversion Ju et al. (2023) introduces a more robust benchmark for text-guided image editing methods, comprising 700 images and 10 editing types, and utilizes 8 evaluation metrics for an objective and comprehensive assessment. Although these benchmarks are widely used by existing text-guided image editing methods, however, the lack of paired training data prevents them from being applicable to the instruction inversion methods. Visii Nguyen et al. (2023) utilizes the filtered dataset of IP2P Brooks et al. (2023) for evaluation. However, despite being filtered by CLIP similarity, The overall quality of the IP2P training data is still poor, which is reflected in the quality and fidelity of the images before and after their editing. Furthermore, the dataset of IP2P contains fewer pairs of data for the same editing type, which hinders an accurate assessment of the performance of the instruction inversion method under the few-shot setting.

To investigate the editing capabilities of various instruction inversion methods in open scenarios and facilitate a fair comparison of these methods, We establish a benchmark named *TOP-Bench* (**T**ransformation-**O**riented **P**aired **B**enchmark), which can be utilized for both qualitative and quantitative evaluations. Our benchmark spans 25 datasets corresponding to different editing effects. It covers a wide range of editing categories and scenarios, allowing for division from multiple dimensions. Each dataset consists of 10 pairs of training images and 5 pairs of testing images, totaling 750 images. Additionally, we provide text instructions aligned with the transformation effects for each dataset.

In order to obtain paired data representing image editing, we refer to the IP2P method of generating data and utilize the existing image editing method P2P Hertz et al. (2022) to directly generate paired data before and after editing. For different editing effects, some of them completely replicate the training set of IP2P, i.e., using the image caption as well as the editing instructions are from the training set of IP2P, and the same settings of IP2P are used to generate and filter the high-quality data of the present method so as to represent the editing of the scene in the domain; while for some editing effects, we generate them through the SDXL-based P2P, while the image caption as well as the editing instructions are obtained based on GPT-4 to represent the editing of the out-of-domain scene. *TOP-Bench* provides paired before and after editing data. It is suitable for the evaluation of instruction inversion methods. At the same time, *TOP-Bench* can be segmented in multiple



885 **Figure 7: Image Tone Modification.** Our InstructBrush can extract various image tones from a
886 handful of data pairs and apply them to new images. The images, from left to right, show three
887 reference image pairs for optimization, the input images, the corresponding editing results, and the
888 ground truth.



909 **Figure 8: Hybrid Instruction.** Our method can combine the visual instruction that represent
910 a particular style with different textual instructions in order to jointly guide the image editing.
911 In contrast, the visual instruction of Visii is forgotten in the process of combining with textual
912 instructions.

913
914
915
916 dimensions to comprehensively evaluate the performance of instruction reversal methods. A detailed
917 presentation of the datasets representing the different editing effects within TOP-Bench and their
categorization is shown in Figure 6 and Table 3.

Table 4: **Extra Ablation Study.** We ablate the impact of optimizing the first m tokens initialized in K , V (Ours) and the impact of optimizing all tokens.

Setting	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP-D \downarrow	CLIP-I \uparrow	DINO \uparrow
All tokens	19.66	0.6660	0.2254	0.3214	0.9379	0.9114
Ours	20.68	0.6922	0.1918	0.3140	0.9115	0.9412

E APPLICATIONS

Image Retouching. Image retouching is the process of changing or improving the quality of an image. This involves enhancing colors, removing imperfections, adjusting lighting, or making other edits to improve the overall appearance of an image. Implementing image retouching using the instruction-based image editing models is challenging because the vast majority of image retouching transformations are difficult to describe using textual instructions. Our method helps in this task. Given paired data before and after image retouching, our InstructBrush can extract editing instructions representing this image transformation from the prior of generative model. The aligned instructions obtained through this process facilitate training for downstream tasks. As shown in Figure 7, our InstructBrush can extract image tones based on three image pairs that represent tonal transformations from PPR10K Liang et al. (2021) and apply these transformations to new images.

Hybrid Instruction. Hybrid instruction allows the use of textual and visual instructions together to guide the editing of an image. Adopting the approach mentioned in Nguyen et al. (2023), we concatenate the embeddings representing visual and textual instructions for guided editing. For visual instruction optimization, we disabled time-aware instruction optimization to achieve better hybrid instruction results. As shown in Figure, our method can combine the visual instruction that represent a particular style with different textual instructions in order to jointly guide the image editing. In contrast, the visual instruction of Visii Nguyen et al. (2023) which represents style is forgotten in the process of combining with textual instructions.

F ADDITIONAL EXPERIMENTS

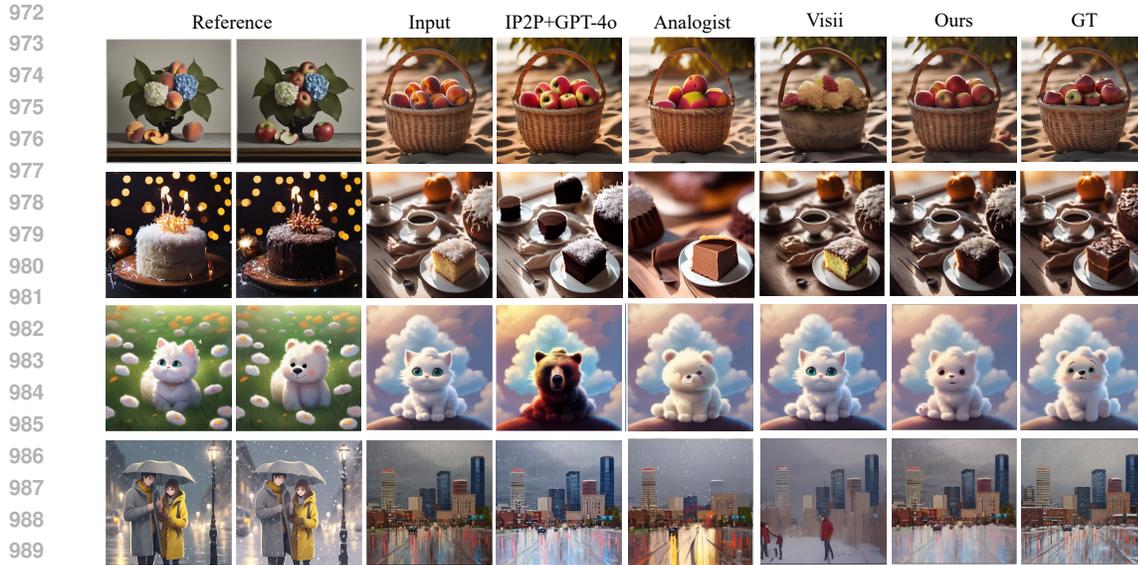
F.1 ONE-SHOT EDITING

To further demonstrate the advantages of our method, we test the quantitative results of different methods under 1-shot setting. We use the first pair of images from each training sets within the benchmark as our 1-shot training data pair. All settings were kept the same as in previous experiments. As shown in Table 5. Our method outperforms the other methods under 1-shot and has the same trend as the few-shot quantitative experiments. Compared to the results of Visii, our method still shows a more significant improvement on TOP-Local than on TOP-Global for the 1-shot setting. This shows in local editing tasks, training images contain more editing-irrelevant scene information. The content-oriented initialization of Visii introduces them to the initialized instructions, posing a greater obstacle to optimization. Our method can accurately capture the transformations between image pairs and use them for initialization, thus improving instruction generalization.

The 1-shot comparative experiment is shown in Figure 9. It proves that the editing effect of our method under 1-shot is better than that of Visii, while the latter easily leaks the content of the reference image pair in the result. This further verifies the advantage of the transformation-oriented initialization design of our method. Additionally, our method demonstrates more consistent editing results with reference image pairs compared to IP2P, which further demonstrates the advantage of providing image pairs for image editing. Additional 1-shot results are shown in Figure 10, which further confirms the editing effect and generalization capabilities of our method.

F.2 EXTRA ABLATION STUDY

In Table 4, we ablate the impact of optimizing the first m tokens initialized in K , V and the impact of optimizing all tokens. The results show that optimizing the first m tokens, which is our current



991 Figure 9: **One-shot Qualitative Comparisons with Existing Methods.** Our method achieves superior performance in both local and global image editing in one-shot setting. It effectively avoids introducing editing-irrelevant information from the training images, showing better instruction generalization.

992 Table 5: **Quantitative Results for One-shot.** We measure the average PSNR, SSIM, LPIPS, and CLIP direction scores of several methods in different editing tasks. In 1-shot settings, our method demonstrates significant superiority over other methods. We highlight in red the percentage of our method that exceeds Visii.

993
994
995
996
997
998
999

Datasets	Method	PSNR \uparrow	SSIM \uparrow	LPIPS \downarrow	CLIP-D \downarrow	CLIP-I \uparrow	DINO \uparrow
TOP-Global	IP2P+GPT-4o Brooks et al. (2023)	13.53	0.4878	0.3884	0.6897	0.8676	0.8959
	Analogist Gu et al. (2024)	12.28	0.3860	0.3993	0.8025	0.8047	0.8800
	Visii Nguyen et al. (2023)	16.01	0.5071	0.3692	0.2909	0.8560	0.8834
	Ours	17.79	0.5761	0.2748	0.3008	0.9009	0.9276
TOP-Local	IP2P+GPT-4o Brooks et al. (2023)	17.33	0.7016	0.2738	0.6188	0.8844	0.9072
	Analogist Gu et al. (2024)	14.03	0.5120	0.3334	0.8559	0.8535	0.9021
	Visii Nguyen et al. (2023)	19.73	0.7293	0.2309	0.5736	0.8794	0.9136
	Ours	23.08	0.8270	0.1172	0.4790	0.9422	0.9764
TOP-Bench	IP2P+GPT-4o Brooks et al. (2023)	15.20	0.5819	0.3380	0.6585	0.8750	0.9009
	Analogist Gu et al. (2024)	13.05	0.4414	0.3703	0.8260	0.8262	0.8897
	Visii Nguyen et al. (2023)	17.65	0.6049	0.3083	0.4153	0.8663	0.8967
	Ours	20.11	0.6865	0.2055	0.3792	0.9191	0.9491

1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014

method’s setting, yields better effects. We attribute this to the fact that optimizing tokens with semantic information is sufficient and brings more generalization to the optimized instructions.

1015 F.3 MORE VISUALIZATION RESULTS

1016
1017
1018

Visualization Results Testing on Our Benchmark. We show more visualization results of our method applied to local and global editing in Figure 16 and Figure 17.

1019
1020
1021
1022
1023
1024
1025

Visualization Results Testing on Real-world Images. We test the performance of our method on real-world images. These data are obtained from the website as well as the PIE-Bench Ju et al. (2023). As shown in Figure 11, Our method can still achieve various editing tasks well even on real-world-images. This further validates the generalization ability of our method.

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

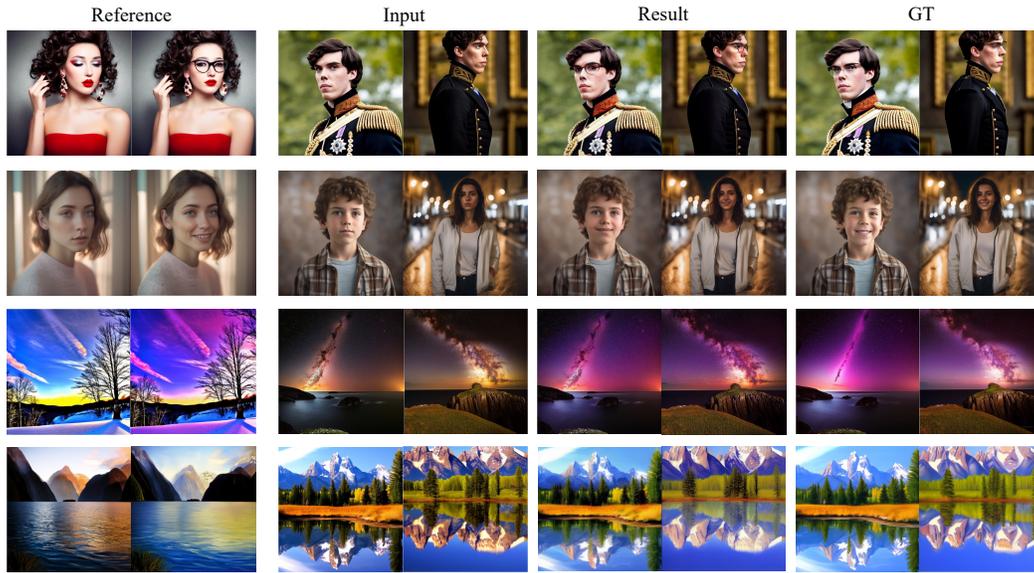
1041

1042

1043

1044

1045



1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

Figure 10: **One-shot Visualization Results of Our Method.** Our method demonstrates robust performance on both local and global editing in one-shot setting. Even if the training and test scenes are quite different, our method can well extract the target editing effect from the training pairs and apply it to new images. This further verifies the generalization ability of our method.

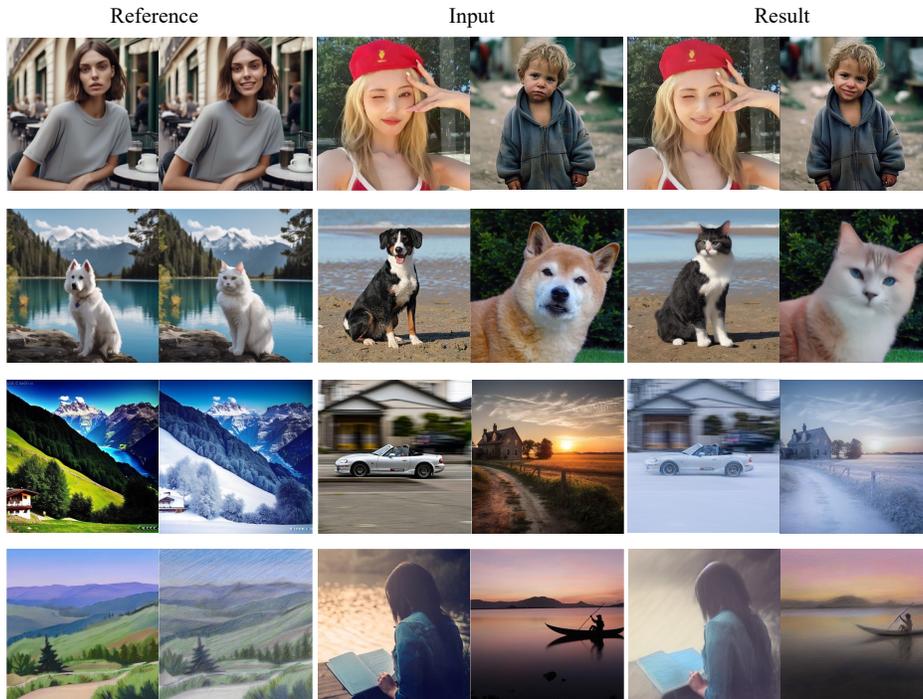


Figure 11: **Testing of our method on real-world images.** Our method demonstrates robust performance on both local and global editing. Our method also works well for editing real-world images. This further verifies the generalization ability of our method.

1080
 1081
 1082
 1083
 1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095
 1096
 1097
 1098
 1099
 1100
 1101
 1102
 1103
 1104
 1105
 1106
 1107
 1108
 1109
 1110
 1111
 1112
 1113
 1114
 1115
 1116
 1117
 1118
 1119
 1120
 1121
 1122
 1123
 1124
 1125
 1126
 1127
 1128
 1129
 1130
 1131
 1132
 1133

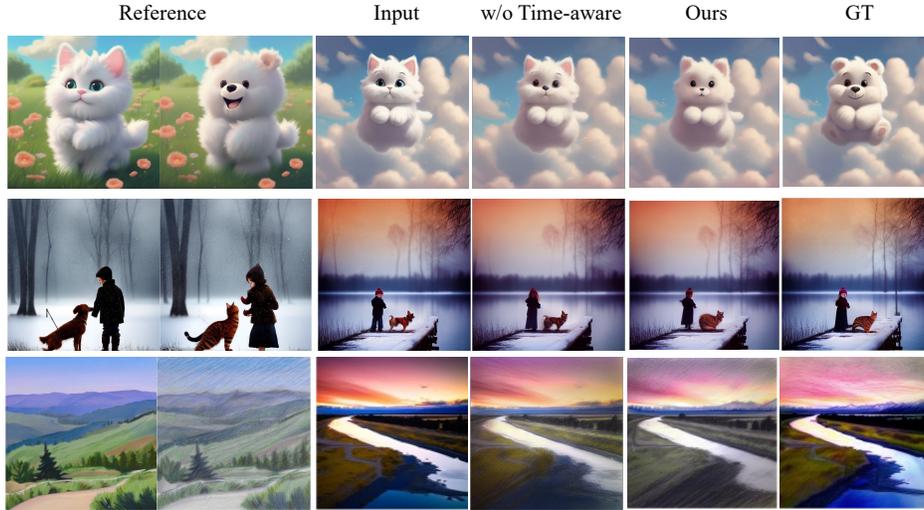


Figure 12: **Visualization effect of the Time-aware optimization method.** We provide qualitative results on the ablation of the time-aware optimization method to visualize the effectiveness of this design.

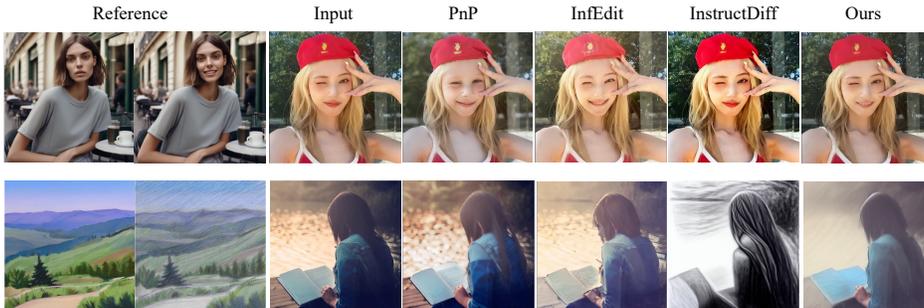


Figure 13: **Comparison with text-guided editing methods.** Our method is compared with text-guided image editing methods Tumanyan et al. (2023); Xu et al. (2024); Geng et al. (2023), and our method can more accurately compare the editing effect of the reference image pair, as shown in the figure. All text-guided editing methods use GPT-4o to obtain editing prompts.

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

1154

1155

1156

1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

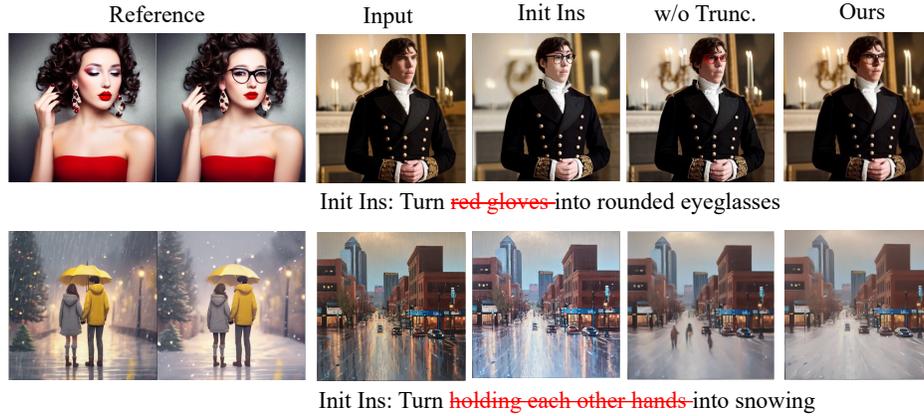


Figure 14: **Visualization of the initialization instruction.** We separately visualized the editing results of the initialization instructions obtained through our Transformation-oriented initialization method and the instructions learned through our Attention-based optimization method. Additionally, we provided the instructions before truncation to verify the effect of using truncation. Below each set of pictures, we present our initialization instructions, with the unique phrases that have been truncated indicated by red strikethroughs.



Figure 15: **Visualization of Applying Time-aware Instructions to Various Denoising Steps.** Example: $T = 800$ represents the application of our time-aware instruction before the denoising time step of 800 (steps 1000 to 800), while the *None* instruction is applied to the denoising process after 800 steps (steps 800 to 0). Therefore, $T = 1000$ indicates the input image, and $T = 0$ indicates our full implementation. The visualization results show that in the early denoising stages, the editing focuses on coarse information such as colors (rows 2 and 3); in the later stages, the editing focuses on detailed information such as textures and facial expressions (rows 1 and 3).

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

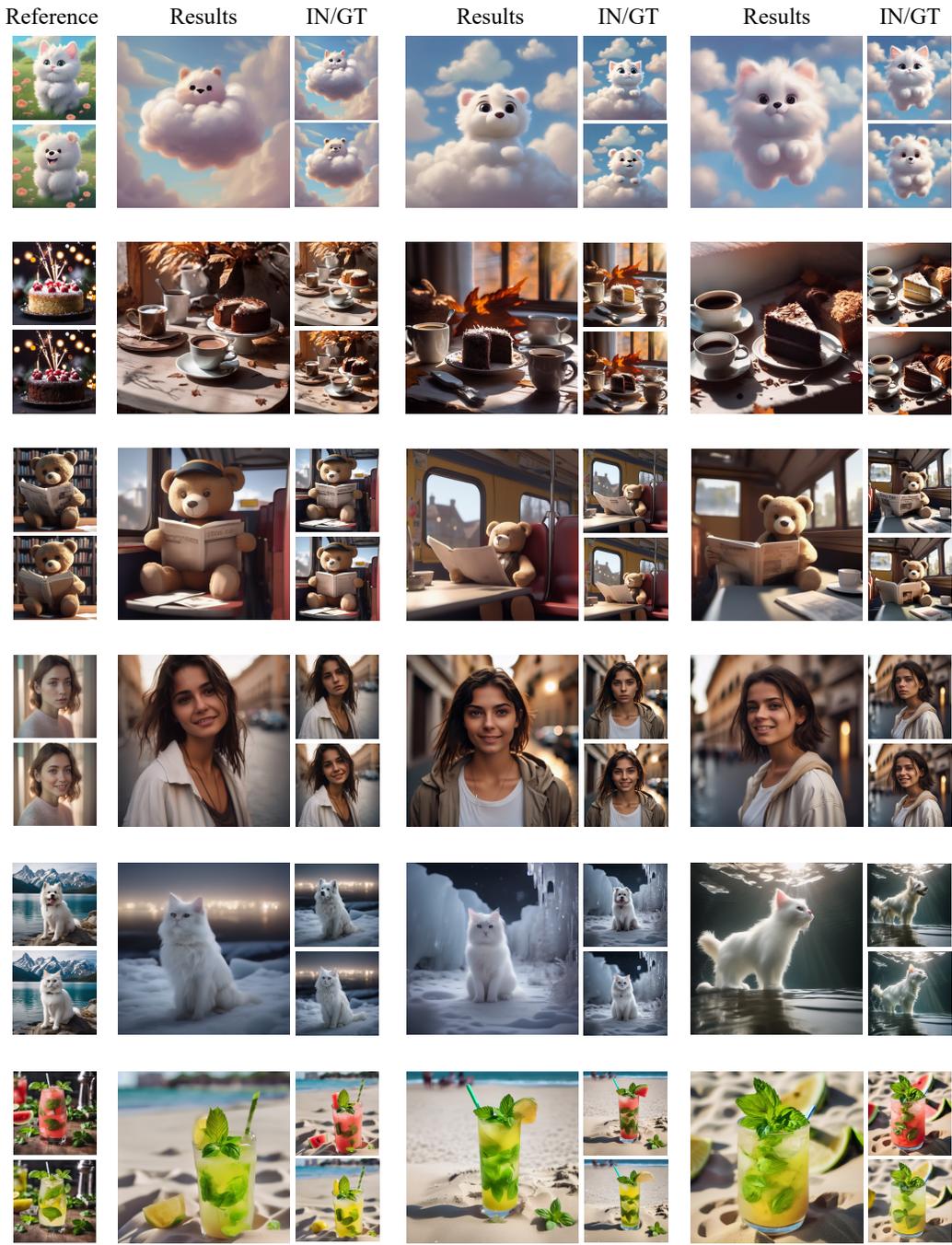


Figure 16: **More Visualization Results of Our Method for Local Editing.** Our method shows robust performance in local editing. Moreover, it does not introduce the scene information of the training image when editing a new image, which reflects the instructive generality of our method.

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

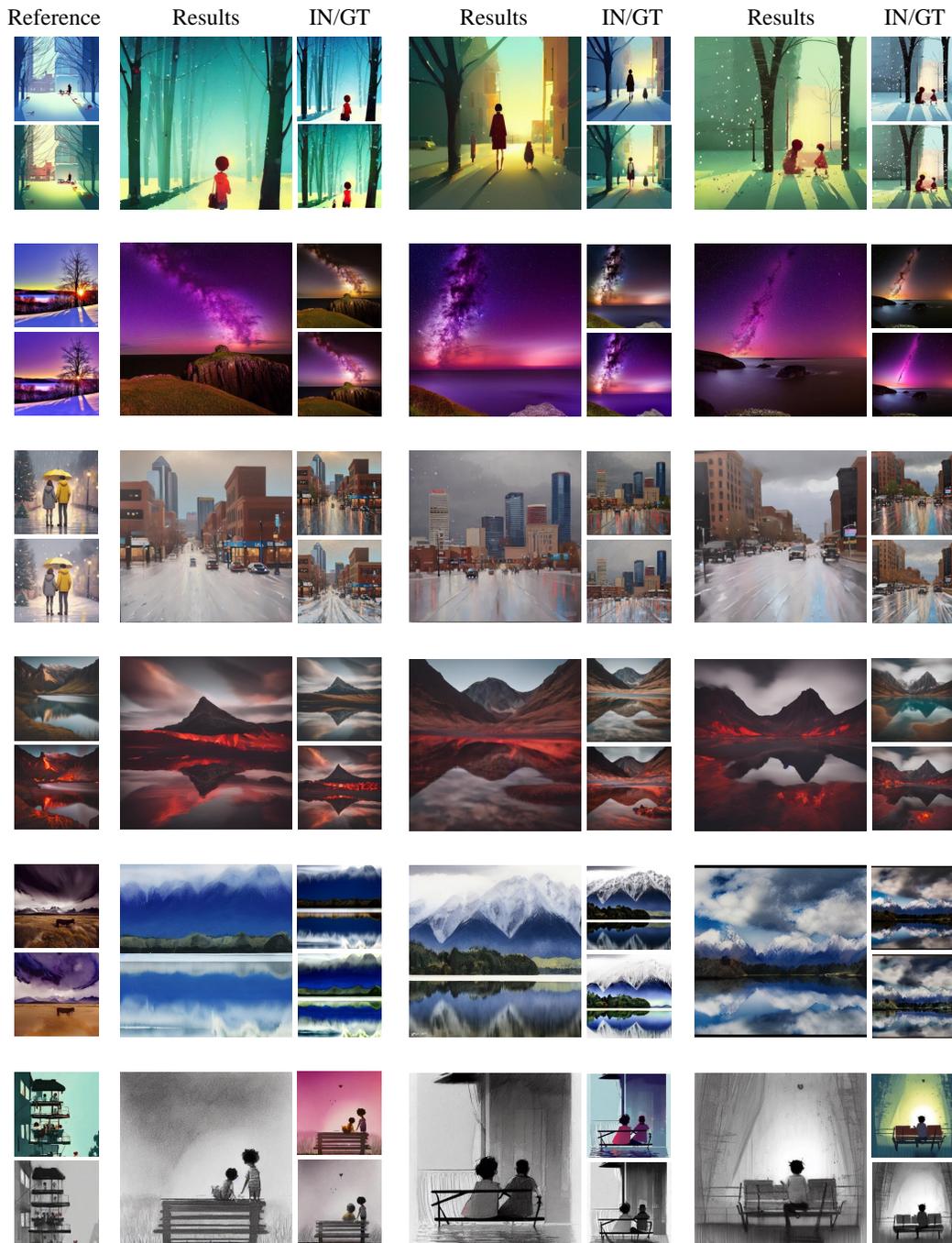


Figure 17: **More Visualization Results of Our Method for Global Editing.** Our method shows robust performance in global editing. Moreover, it does not introduce the scene information of the training image when editing a new image, which reflects the instructive generality of our method.