

Inversion-Free Image Editing with Natural Language

Sihan Xu^{1*} Yidong Huang^{1*} Jiayi Pan^{2†} Ziqiao Ma¹ Joyce Chai¹

¹University of Michigan ²University of California, Berkeley

<https://sled-group.github.io/InfEdit/>

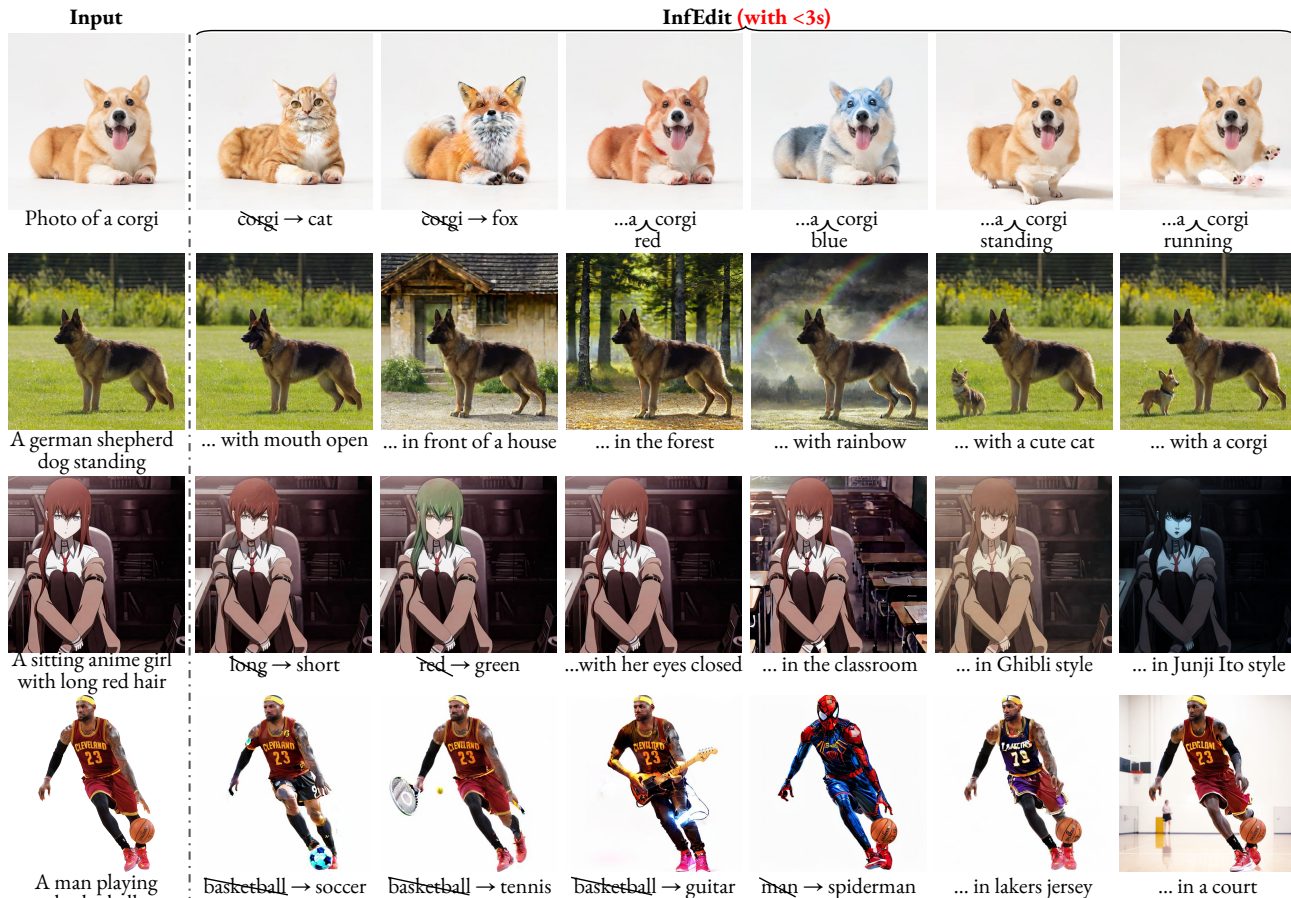


Figure 1. Our inversion-free editing (InfEdit) method demonstrates strong performance in various complex image editing tasks.

Abstract

Despite recent advances in inversion-based editing, text-guided image manipulation remains challenging for diffusion models. The primary bottlenecks include 1) the time-consuming nature of the inversion process; 2) the struggle to balance consistency with accuracy; 3) the lack of compatibility with efficient consistency sampling methods used in consistency models. To address the above issues, we start by asking ourselves if the inversion process can be eliminated for editing. We show that when the initial sample is known, a special variance schedule reduces the denoising step to the same form as the multi-step consistency sam-

pling. We name this Denoising Diffusion Consistent Model (DDCM), and note that it implies a virtual inversion strategy without explicit inversion in sampling. We further unify the attention control mechanisms in a tuning-free framework for text-guided editing. Combining them, we present inversion-free editing (InfEdit), which allows for consistent and faithful editing for both rigid and non-rigid semantic changes, catering to intricate modifications without compromising on the image’s integrity and explicit inversion. Through extensive experiments, InfEdit shows strong performance in various editing tasks and also maintains a seamless workflow (less than 3 seconds on one single A40), demonstrating the potential for real-time applications.

* Authors contributed equally to this work.

† Work done while the author was at the University of Michigan.

1. Introduction

Recent progress in image synthesis has been mostly driven by the development of Diffusion Models (DMs) [13, 29], which have outperformed traditional Generative Adversarial Networks (GANs) [8] and Variational Autoencoders (VAEs) [16] in various applications. A key factor in the wide success of DMs is their ability to incorporate diverse conditions, such as text [27], images [22, 35, 37], and even tactile input [36]. Building upon DMs, Consistency Models (CMs) [30] address the efficiency bottleneck by directly mapping noised samples along a trajectory to the same initial, promising *self-consistency*.

Enabling text-guided DMs for editing real images using natural language has presented significant challenges. Early methods typically require additional mask layers [2, 3, 6, 24] or training [5, 35, 39], which constrain their potential zero-shot application. Motivated by DDIM inversion [29], a prevailing paradigm of *inversion-based* editing has been established. The predominating methods along this line adopt *optimization-based inversion* [7, 17, 21] by aligning the forward source latents with the DDIM inversion trajectory. To address the issues of efficiency bottlenecks and far-from-ideal consistency, *dual-branch* methods [15, 34] have been introduced, which separate the source and target branches individually, and iteratively calibrate the trajectory of the target branch. However, inversion-based editing methods still face limitations in real-time and real-world language-guided image editing. Firstly, they typically rely on a lengthy inversion process to acquire the inversion branch as a series of anchors. Secondly, striking a balance between consistency and faithfulness remains challenging, even with extensive optimization or ways of calibrating the target branch. Lastly, these methods rely on variations of diffusion sampling, which are not compatible with the efficient consistency sampling using CMs.

To address the above challenges, we start by asking ourselves if the inversion process is really required for editing. We show that when the initial sample is known, there exists a special variance schedule such that the denoising step takes the same form as the multi-step consistency sampling. We name this Denoising Diffusion Consistent Model (DDCM), and note that it implies a sampling strategy that eliminates the inversion process. We further present Unified Attention Control (UAC), a tuning-free method that unifies attention control mechanisms for text-guided editing. Combining them, we present an inversion-free editing (InfEdit) framework that allows for consistent and faithful editing for both rigid and non-rigid semantic changes, catering to intricate modifications without compromising on the image’s integrity and explicit inversion. Through experiments, InfEdit shows strong performance in various editing tasks and also maintains a seamless workflow (less than 3s on one A40), demonstrating the potential for real-time editing.

2. Preliminaries

2.1. Diffusion Models

Diffusion models (DMs) [13] operate through a forward process that gradually adds Gaussian noises to data, described as follows:

$$z_t = \sqrt{\alpha_t}z_0 + \sqrt{1 - \alpha_t}\varepsilon \quad \varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (1)$$

where z_0 is a sample from the data distribution, $\alpha_{1:T}$ specify a variance schedule for $t \sim [1, T]$.

The training objective involves a parameterized noise prediction network, ε_θ , which aims to reverse the diffusion process. The training objective is to minimize the following loss based on a chosen metric function for measuring the distance between two samples $d(\cdot, \cdot)$:

$$\min_{\theta} \mathbb{E}_{z_0, \varepsilon, t} [d(\varepsilon, \varepsilon_\theta(z_t, t))] \quad (2)$$

Sampling from a diffusion model is an iterative process that progressively denoises the data. Following Eq (12) in Song et al. [29], the denoising step at t is formulated as:

$$\begin{aligned} z_{t-1} = & \sqrt{\alpha_{t-1}} \left(\frac{z_t - \sqrt{1 - \alpha_t} \varepsilon_\theta(z_t, t)}{\sqrt{\alpha_t}} \right) && \text{(predicted } z_0) \\ & + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \varepsilon_\theta(z_t, t) && \text{(direction to } z_t) \\ & + \sigma_t \varepsilon_t \quad \text{where } \varepsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) && \text{(random noise)} \end{aligned} \quad (3)$$

DDPM sampling [13] introduces a noise schedule σ_t so that Eq (3) becomes Markovian. By setting σ_t to vanish, DDIM sampling [29] results in an implicit probabilistic model with a deterministic forward process.

Following DDIM, we can use the function f_θ to predict and reconstruct \bar{z}_0 given z_t :

$$\bar{z}_0 = f_\theta(z_t, t) = (z_t - \sqrt{1 - \alpha_t} \cdot \varepsilon_\theta(z_t, t)) / \sqrt{\alpha_t} \quad (4)$$

Recently, Latent Diffusion Models (LDMs) [28] offer a new paradigm by operating in the latent space. The source latent z_0 is acquired by encoding a sample x_0 with an encoder \mathcal{E} , such that $z_0 = \mathcal{E}(x_0)$. So as to be reversed, the output can then be reconstructed by a decoder \mathcal{D} . This framework presents a computationally efficient way to generate high-fidelity images, as the diffusion process is conducted in a latent space with lower dimensions.

2.2. Consistency Models

Consistency models (CMs) [30] have recently been introduced, which greatly accelerate the generation process compared with previous DMs. One notable property of CMs is *self-consistency*, such that samples along a trajectory map to the sample initial. The key is a consistency function $f(z_t, t)$, which ensures a consistent distillation process by optimizing:

$$\min_{\theta, \theta^-; \phi} \mathbb{E}_{z_0, t} \left[d \left(f_\theta(z_{t_{n+1}}, t_{n+1}), f_{\theta^-}(\hat{z}_{t_n}^\phi, t_n) \right) \right] \quad (5)$$

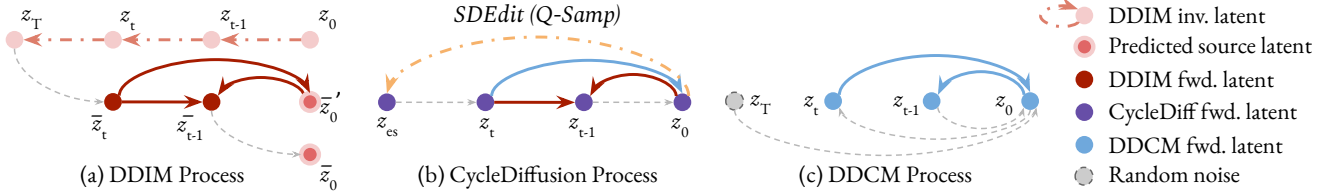


Figure 2. While DDIM is prone to reconstruction error and requires iterative inversion, DDCM accepts any random noise to start with. It introduces a non-Markovian forward process in which z_t directly points to the ground truth z_0 without neural prediction, and z_{t-1} does not depend on the previous step z_t like a consistency model.

in which f_θ denotes a trainable neural network that parameterizes these consistent transitions, while f_{θ^-} represents a slowly updated target model used for consistency distillation, with the update rule $\theta^- \leftarrow \mu\theta^- + (1 - \mu)\theta$ given a decay rate μ . The variable $\hat{z}_{t_n}^\phi$ denotes a one-step estimation of z_{t_n} from $z_{t_{n+1}}$.

Sampling in CMs is carried out through a sequence of timesteps $\tau_{1:n} \in [t_0, T]$. Starting from an initial noise \hat{z}_T and $z_0^{(T)} = f_\theta(\hat{z}_T, T)$, at each time-step τ_i , the process samples $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and iteratively updates the *Multistep Consistency Sampling* process:

$$\begin{aligned} \hat{z}_{\tau_i} &= z_0^{(\tau_{i+1})} + \sqrt{\tau_i^2 - t_0^2} \varepsilon \\ z_0^{(\tau_i)} &= f_\theta(\hat{z}_{\tau_i}, \tau_i) \end{aligned} \quad (6)$$

Latent Consistency Models (LCMs) [18] extend to accommodate a (text) condition c , which is crucial for text-guided image manipulation. Similarly, sampling in LCMs at τ_i starts with $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and updates:

$$\begin{aligned} \hat{z}_{\tau_i} &= \sqrt{\alpha_{\tau_i}} z_0^{(\tau_{i+1})} + \sigma_{\tau_i} \varepsilon, \\ z_0^{(\tau_i)} &= f_\theta(\hat{z}_{\tau_i}, \tau_i, c) \end{aligned} \quad (7)$$

2.3. Inversion-Based Image Editing with LDMs

DDIM inversion [29] is effective for unconditional diffusion applications, but lacks consistency with additional text or image conditions. As illustrated in Figure 2a, the predicted \hat{z}_0' deviates from the original source z_0 , cumulatively leading to undesirable semantic changes. This substantially restricts its use in image editing driven by natural language-guided diffusion.

To address this concern, various forms of inversion-based editing methods have been proposed. The predominant approaches utilize *optimization-based inversion* [7, 17, 21]. These methods aim to “correct” the forward latents guided by the source prompt (referred to as the source branch) by aligning them with the DDIM inversion trajectory. To tackle the efficiency bottlenecks and suboptimal consistency, very recent work has explored *dual-branch inversion* [15, 34]. These methods separate the source and target branches in the editing process: directly revert the source branch back to z_0 and iteratively calibrate the trajectory of the target branch. As shown in Figure 3a, they cal-

culate the distance between the source branch and the inversion branch (or directly sampled from q -sampling in [34]), and calibrate the target branch with this computed distance at each t .

3. Denoising Diffusion Consistent Models

We start with the following proposition.

Proposition 1 (Denoising Diffusion Consistent Models)

Consider a special case of Eq (3) when σ_t is chosen as $\sqrt{1 - \alpha_{t-1}}$ across all time t , the forward process naturally aligns with the *Multistep (Latent) Consistency Sampling*.

When $\sigma_t = \sqrt{1 - \alpha_{t-1}}$, the second term of Eq (3) vanishes:

$$\begin{aligned} z_{t-1} &= \sqrt{\alpha_{t-1}} \left(\frac{z_t - \sqrt{1 - \alpha_t} \varepsilon_\theta(z_t, t)}{\sqrt{\alpha_t}} \right) \quad (\text{predicted } z_0) \\ &+ \sqrt{1 - \alpha_{t-1}} \varepsilon_t \quad \varepsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (\text{random noise}) \end{aligned} \quad (8)$$

Consider $f(z_t, t; z_0) = (z_t - \sqrt{1 - \alpha_t} \varepsilon'(z_t, t; z_0)) / \sqrt{\alpha_t}$, where the initial z_0 is available (which is the case for image editing applications) and we replace the parameterized noise predictor ε_θ with ε' more generally. Eq (8) becomes

$$z_{t-1} = \sqrt{\alpha_{t-1}} f(z_t, t; z_0) + \sqrt{1 - \alpha_{t-1}} \varepsilon_t \quad (9)$$

which is in the same form as the *Multistep Latent Consistency Sampling* step in Eq (7).

In order to make $f(z_t, t)$ self-consistent so that it can be considered as a consistency function, i.e., $f(z_t, t; z_0) = z_0$, we can directly solve the equation and ε' can be computed without parameterization:

$$\varepsilon^{\text{cons}} = \varepsilon'(z_t, t; z_0) = \frac{z_t - \sqrt{\alpha_t} z_0}{\sqrt{1 - \alpha_t}} \quad (10)$$

As illustrated in Figure 2c, we arrive at a non-Markovian forward process, in which z_t directly points to the ground truth z_0 without neural prediction, and z_{t-1} does not depend on the previous step z_t like a consistency model. We name this *Denoising Diffusion Consistent Model* (DDCM).

3.1. DDCM for Virtual Inversion

We note that DDCM suggests an image reconstruction model without any explicit inversion operation, diverging from conventional DDIM inversion and its optimized or calibrated variations for image editing. It achieves the best

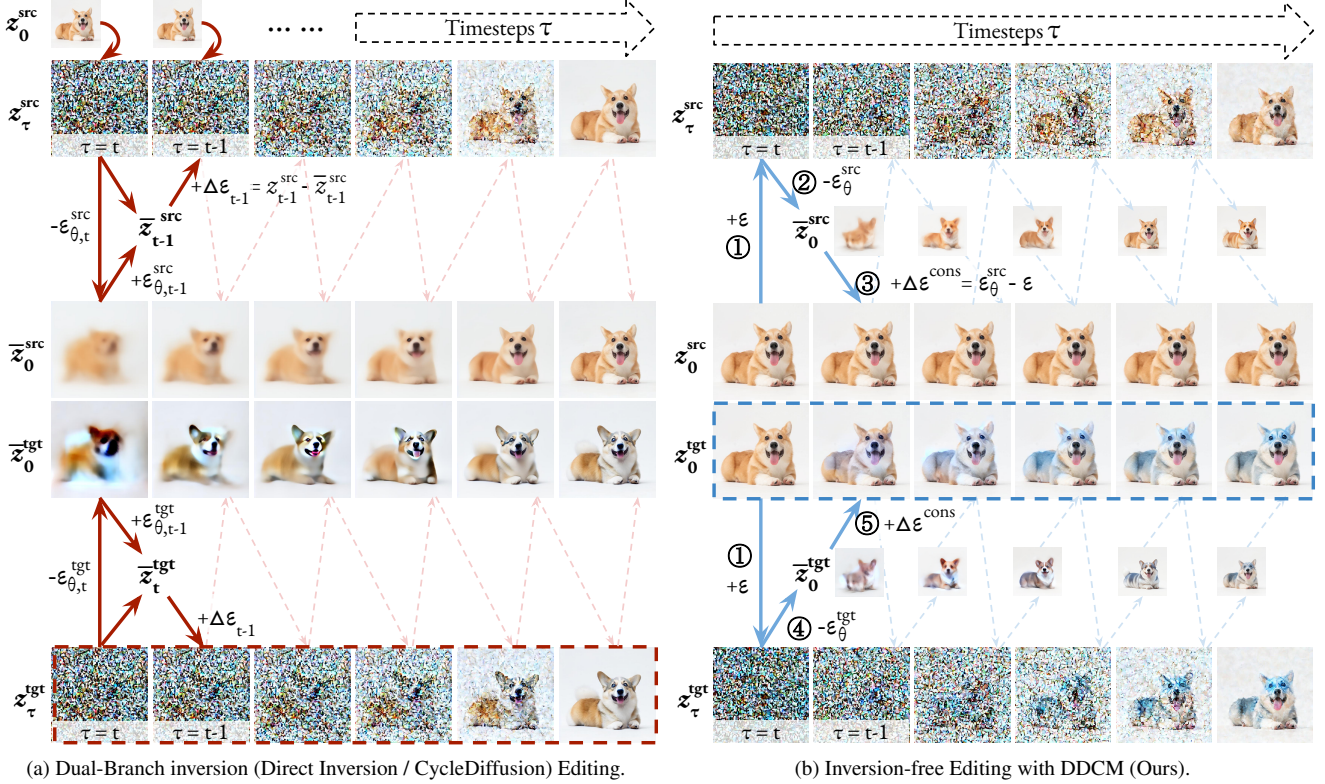


Figure 3. A comparative overview of the dual-branch inversion editing and inversion-free editing enabled by DDCM. While the former iteratively calibrates z_t^{tgt} in the target branch, inversion-free editing iteratively polishes the target branch initial z_0^{tgt} . In (b), we initialize z_0^{tgt} with z_0^{src} for visualization purposes, while in principle it can start from any random noise. The circled numbers correspond to Algorithm 2.

efficiency as it allows the forward process to start from any random noise and supports multi-step consistency sampling. On the other hand, it ensures exact consistency between original and reconstructed images as each step on the forward branch z_{t-1} only depends on the ground truth z_0 rather than the previous step z_t . Due to its inversion-free nature, we name this method *Virtual Inversion*. As outlined in Algorithm 1, $z = z_0$ is ensured throughout the process without parameterization.

Algorithm 1 DDCM Sampling for Virtual Inversion

- 1: **Input:**
 Sequence of timesteps $\tau_1 > \tau_2 > \dots > \tau_{N-1}$
 Reference initial input z_0
 - 2: Sample a random terminal noise $z_{\tau_1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 3: $\varepsilon_{\tau_1}^{\text{cons}} = (z_{\tau_1} - \sqrt{\alpha_{\tau_1}}z_0) / \sqrt{1 - \alpha_{\tau_1}}$
 - 4: $z = (z_{\tau_1} - \sqrt{1 - \alpha_{\tau_1}}\varepsilon_{\tau_1}^{\text{cons}}) / \sqrt{\alpha_{\tau_1}}$
 - 5: **for** $n = 2$ to $N - 1$ **do**
 - 6: Sample noise $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 7: $z_{\tau_n} = \sqrt{\alpha_{\tau_n}}z + \sqrt{1 - \alpha_{\tau_n}}\varepsilon$
 - 8: $\varepsilon_{\tau_n}^{\text{cons}} = (z_{\tau_n} - \sqrt{\alpha_{\tau_n}}z_0) / \sqrt{1 - \alpha_{\tau_n}}$
 - 9: $z = (z_{\tau_n} - \sqrt{1 - \alpha_{\tau_n}}\varepsilon_{\tau_n}^{\text{cons}}) / \sqrt{\alpha_{\tau_n}}$
 - 10: **end for**
 - 11: **Output:** z
-

3.2. DDCM for Inversion-Free Image Editing

Existing inversion-based editing methods are limited for real-time and real-world language-driven image editing applications. First, most of them still depend on a time-consuming inversion process to obtain the inversion branch as a set of anchors. Second, consistency remains a bottleneck given the efforts from optimization and calibration. Recall that dual-branch inversion methods perform editing on the target branch by iteratively calibrating the z_t^{tgt} with the actual distance between the source branch and the inversion branch at t , as is boxed in Figure 3a. While they ensure faithful reconstruction by leaving the source branch untouched from the target branch, the calibrated z_t^{tgt} does guarantee consistency from z_t^{src} in the source branch, as can be seen from the visible difference between z_0^{src} and z_0^{tgt} in Figure 3a. Third, all current inversion-based methods rely on variations of diffusion sampling, which are incompatible with efficient Consistency Sampling using LCMs.

DDCM offers an alternative to address these limitations, introducing an Inversion-Free Image Editing (InfEdit) framework. While also adopting a dual-branch paradigm, the key of our InfEdit method is to directly calibrate the initial z_0^{tgt} rather than the z_t^{tgt} along the branch, as is boxed

in Figure 3b. InfEdit starts from a random terminal noise $z_{\tau_1}^{\text{src}} = z_{\tau_1}^{\text{tgt}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. As shown in Figure 3b, the source branch follows the DDCM sampling process without explicit inversion, and we directly compute the distance $\Delta\varepsilon^{\text{cons}}$ between $\varepsilon^{\text{cons}}$ and $\varepsilon_{\theta}^{\text{src}}$ (the predicted noise to reconstruct a \bar{z}_0^{src}). For the target branch, we first compute the $\varepsilon_{\theta}^{\text{tgt}}$ to predict \bar{z}_0^{tgt} , and then calibrate the predicted target initial with the same $\Delta\varepsilon^{\text{cons}}$. Algorithm 2 outlines the mathematical details of this process, in which we slightly abuse the notation to define $f_{\theta}(z_t, t, \varepsilon) = (z_t - \sqrt{1 - \alpha_t}\varepsilon) / \sqrt{\alpha_t}$.

Algorithm 2 DDCM for inversion-free image editing

Input:

Conditional Diffusion/Consistency Model $\varepsilon_{\theta}(\cdot, \cdot, \cdot)$
Sequence of timesteps $\tau_1 > \tau_2 > \dots > \tau_{N-1}$
Reference initial input z_0^{src}
Source/target prompts as conditions $c^{\text{src}}, c^{\text{tgt}}$

- 1: Sample a random terminal noise $z_{\tau_1}^{\text{src}} = z_{\tau_1}^{\text{tgt}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 2: $\varepsilon_{\tau_1}^{\text{cons}} = (z_{\tau_1}^{\text{src}} - \sqrt{\alpha_{\tau_1}}z_0^{\text{src}}) / \sqrt{1 - \alpha_{\tau_1}}$
 - 3: $\varepsilon_{\tau_1}^{\text{src}}, \varepsilon_{\tau_1}^{\text{tgt}} = \varepsilon_{\theta}(z_{\tau_1}^{\text{src}}, \tau_1, c^{\text{src}}), \varepsilon_{\theta}(z_{\tau_1}^{\text{tgt}}, \tau_1, c^{\text{tgt}})$
 - 4: $z_0^{\text{tgt}} = f_{\theta}(z_{\tau_1}^{\text{tgt}}, \tau_1, \varepsilon_{\tau_1}^{\text{tgt}} - \varepsilon_{\tau_1}^{\text{src}} + \varepsilon_{\tau_1}^{\text{cons}})$
 - 5: **for** $n = 2$ to $N - 1$ **do**
 - 6: Sample noise $\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 7: ① $z_{\tau_n}^{\text{src}} = \sqrt{\alpha_{\tau_n}}z_0^{\text{src}} + \sqrt{1 - \alpha_{\tau_n}}\varepsilon$
 - 8: ① $z_{\tau_n}^{\text{tgt}} = \sqrt{\alpha_{\tau_n}}z_0^{\text{tgt}} + \sqrt{1 - \alpha_{\tau_n}}\varepsilon$
 - 9: ② $\varepsilon_{\tau_n}^{\text{src}} = \varepsilon_{\theta}(z_{\tau_n}^{\text{src}}, \tau_n, c^{\text{src}})$
 - 10: ③ $\varepsilon_{\tau_n}^{\text{cons}} = (z_{\tau_n}^{\text{src}} - \sqrt{\alpha_{\tau_n}}z_0^{\text{src}}) / \sqrt{1 - \alpha_{\tau_n}}$
 - 11: ④* $\varepsilon_{\tau_n}^{\text{tgt}} = \varepsilon_{\theta}(z_{\tau_n}^{\text{tgt}}, \tau_n, c^{\text{tgt}})$
 - 12: ⑤ $z_0^{\text{tgt}} = f_{\theta}(z_{\tau_n}^{\text{tgt}}, \tau_n, \varepsilon_{\tau_n}^{\text{tgt}} - \varepsilon_{\tau_n}^{\text{src}} + \varepsilon_{\tau_n}^{\text{cons}})$
 - 13: **end for**
 - 14: **Output:** z_0^{tgt}
 - 15: **Vanilla target noise prediction, no attention control.*
-

InfEdit addresses the current limitations of inversion-based editing methods. First, DDCM sampling allows us to abandon the inversion branch anchors required by previous methods, saving a significant amount of computation. Second, the current dual-branch methods calibrate z_t^{tgt} over time, while InfEdit directly refines the predicted initial z_0^{tgt} , without suffering from the cumulative errors over the course of sampling. Third, our framework is compatible with efficient Consistency Sampling using LCMs, enabling efficient sampling of the target image within very few steps.

4. Unifying Attention Control for Language-Guided Editing

InfEdit suggests a general inversion-free framework for image editing motivated by DDCM. In the realm of language-driven editing, achieving a nuanced understanding of the language condition and facilitating finer-grained interaction across modalities becomes a challenge. Hertz et al. [9] noticed that the interaction between the text and image modal-

ities occurs in the parameterized noise prediction network ε_{θ} , and opened up a series of attention control methods to compute a noise $\widehat{\varepsilon}_{\theta}^{\text{tgt}}$ that more accurately aligns with the language prompts. In the context of InfEdit specifically, attention control refines the original predicted target noise $\varepsilon_{\theta}^{\text{tgt}}$ (noted in ④ in Algorithm 2 and Figure 3b) with $\widehat{\varepsilon}_{\theta}^{\text{tgt}}$.

We follow [9] in terms of notation. Each basic block of the U-Net noise predictor contains a cross-attention module and a self-attention module. The spatial features are linearly projected into queries (Q). In cross-attention, the text features are linearly projected into keys (K) and values (V). In self-attention, the keys (K) and values (V) are also obtained from linearly projected spatial features. The attention mechanism [32] can be given as:

$$\text{Attention}(K, Q, V) = MV = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V \quad (11)$$

in which $M_{i,j}$ represents the attention map that determines the weight to aggregate the value of the j -th token on pixel i , and d denotes the dimension for K and Q .

Natural language specifies a wide spectrum of semantic changes. In the following, we describe how *rigid semantic changes*, e.g., those on the visual features and background, can be controlled via cross attention [9]; and how *non-rigid semantic changes*, e.g., those leading to adding/removing an object, novel action manners and physical state changes of objects, can be controlled via mutual self-attention [4]. We then introduce a Unified Attention Control (UAC) protocol for both types of semantic changes.

4.1. Cross-Attention Control

Prompt-to-Prompt (P2P) [9] observed that cross-attention layers can capture the interaction between the spatial structures of pixels and words in the prompts, even in early steps. This finding makes it possible to control the cross-attention for editing rigid semantic changes, simply by replacing the cross-attention map of generated images with that of the original images.

Global Attention Refinement At time step t , we compute the attention map M_t averaged over layers given the noised latent z_t and the prompt for both source and target branch. We drop the time step for simplicity and represent the source and target attention maps as M^{src} and M^{tgt} . To represent the common details, an alignment function $A(i) = j$ is introduced which signifies that the i^{th} word in the target prompt corresponds to the j^{th} word in the source prompt. Following Hertz et al. [9], we refine the target attention map by injecting the source attention map over the common tokens.

$$\text{Refine}(M^{\text{src}}, M^{\text{tgt}})_{i,j} = \begin{cases} (M^{\text{tgt}})_{i,j} & \text{if } A(j) = \text{None} \\ (M^{\text{src}})_{i,A(j)} & \text{otherwise} \end{cases} \quad (12)$$

This ensures that the common information from the source prompt is accurately transferred to the target, while the requested changes are made.

Local Attention Blends Besides global attention refinement, we adapt the blended diffusion mechanism from [1, 9]. Specifically, the algorithm takes optional inputs of target blend words w^{tgt} , which are words in the target prompt whose semantics need to be added; and source blend words w^{src} , which are words in the source prompt whose semantics need to be preserved. At time step t , we blend the noised target latent z_t^{tgt} following:

$$\begin{aligned}
 m^{\text{tgt}} &= \text{Threshold}[M_t^{\text{tgt}}(w^{\text{tgt}}), a^{\text{tgt}}] \\
 m^{\text{src}} &= \text{Threshold}[M_t^{\text{src}}(w^{\text{src}}), a^{\text{src}}] \\
 z_t^{\text{tgt}} &= (1 - m^{\text{tgt}} + m^{\text{src}}) \odot z_t^{\text{src}} + (m^{\text{tgt}} - m^{\text{src}}) \odot z_t^{\text{tgt}}
 \end{aligned}
 \tag{13}$$

in which m^{tgt} and m^{src} are binary masks obtained by calibrating the aggregated attention maps $M_t^{\text{tgt}}(w^{\text{tgt}})$, $M_t^{\text{src}}(w^{\text{src}})$ with threshold parameters a^{tgt} and a^{src} using threshold function:

$$\text{Threshold}(M, a)_{i,j} = \begin{cases} 1 & M_{i,j} \geq a \\ 0 & M_{i,j} < a \end{cases}
 \tag{14}$$

Scheduling Cross-Attention Control Applying cross-attention control throughout the entire sampling schedule will overly focus on spatial consistency, leading to an inability to capture the intended changes. Follow [9], we perform cross-attention control only in early steps before τ_c , interpreted as the cross-attention control strength:

$$\text{CrossEdit}(M^{\text{src}}, M^{\text{tgt}}, t) := \begin{cases} \text{Refine}(M^{\text{src}}, M^{\text{tgt}}) & t \geq \tau_c \\ M^{\text{tgt}} & t < \tau_c \end{cases}$$

4.2. Mutual Self-Attention Control

One key limitation of cross-attention control lies in its inability in non-rigid editing. Instead of applying controls over the cross-attention modules, MasaCtrl [4] observed that the layout of the objects can be roughly formed in the self-attention queries, covering the non-rigid semantic changes complying with the target prompt. The core idea is to synthesize the structural layout with the target prompt in the early steps with the original Q^{tgt} , K^{tgt} , V^{tgt} in the self-attention; and then to query semantically similar contents in K^{src} , V^{src} with the target query Q^{tgt} .

Controlling Non-Rigid Semantic Changes MasaCtrl suffers from the issue of undesirable non-rigid changes. As shown in Figure 8, MasaCtrl can lead to significant inconsistency from the source images, especially in terms of the composition of objects and when there are multiple objects and complex backgrounds. This is not surprising, as the target query Q^{tgt} is used throughout the self-attention control schedule. Instead of relying on the target prompts to

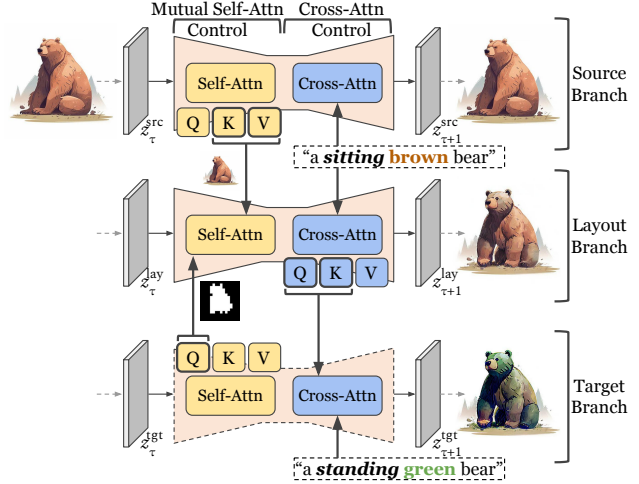


Figure 4. The proposed United Attention Control (UAC) framework to unify cross-attention control and mutual self-attention control. UAC introduces an additional layout branch as an intermediate to host the desired composition and structural information in the target image.

guide the premature steps, we form the structural layout with the source self-attention Q^{src} , K^{src} , V^{src} in the self-attention. We show in Section 5 that this design enables high-quality non-rigid changes while maintaining satisfying structural consistency.

Scheduling Mutual Self-Attention Control This mutual self-attention control is applied in the later steps after τ_s , interpreted as the mutual self-attention control strength:

$$\text{SelfEdit}(\{Q^{\text{src}}, K^{\text{src}}, V^{\text{src}}\}, \{Q^{\text{tgt}}, K^{\text{tgt}}, V^{\text{tgt}}\}, t) := \begin{cases} \{Q^{\text{src}}, K^{\text{src}}, V^{\text{src}}\} & t \geq \tau_s \\ \{Q^{\text{tgt}}, K^{\text{src}}, V^{\text{src}}\} & t < \tau_s \end{cases}
 \tag{15}$$

4.3. Unified Attention Control

To enable both rigid and non-rigid semantic changes within one unified framework is not trivial.

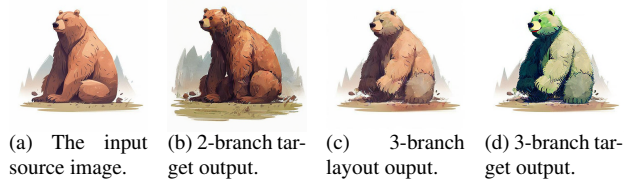


Figure 5. A comparison of the target branch outputs to edit “a sitting brown bear” to “a standing green bear”, involving both rigid and non-rigid semantic transformations. Random seed is fixed.

As is illustrated in Figure 5, the naïve combination of cross-attention control and mutual self-attention control sequentially would lead to a sub-optimal outcome in the original dual-branch setup, especially failing the global attention refinement. To address this issue, we introduce the Unified Attention Control (UAC) framework. UAC unifies cross-attention control and mutual self-attention control with an

additional latent *layout branch*, which serves as an intermediate to host the desired composition and structural information in the target image.

The UAC framework is detailed in Algorithm 3 and illustrated in Figure 4. During each forward step of the diffusion process, UAC starts with mutual self-attention control on z^{src} and z^{tgt} and assigns the output to the layout branch latent z^{lay} . Following this, cross-attention control is applied on M^{lay} and M^{tgt} to refine the semantic information for M^{tgt} . As is shown in Figure 5c, the layout branch output z_0^{lay} reflects the requested non-rigid changes (e.g., “standing”), while preserving the non-rigid content semantics (e.g., “brown”). The target branch output z_0^{tgt} (Figure 5d) builds upon the structural layout of the z_0^{lay} while reflecting the requested non-rigid changes (e.g., “green”).

Algorithm 3 Unified Attention Control Image Editing

1: **Input:**

Conditional Diffusion/Consistency Model $\varepsilon_\theta(\cdot, \cdot, \cdot)$
 Current timestep τ
 Reference initial input z_0^{src}
 Source/target prompts as conditions $c^{\text{src}}, c^{\text{tgt}}$
 Source/target blend words $w^{\text{src}}, w^{\text{tgt}}$
 Input latents $z_\tau^{\text{src}}, z_\tau^{\text{tgt}}, z_\tau^{\text{lay}}$

- 2: $\varepsilon^{\text{src}}, \{Q^{\text{src}}, K^{\text{src}}, V^{\text{src}}\}, M^{\text{src}} = \varepsilon_\theta(z_\tau^{\text{src}}, \tau, c^{\text{src}})$
 - 3: $\varepsilon^{\text{tgt}}, \{Q^{\text{tgt}}, K^{\text{tgt}}, V^{\text{tgt}}\}, M^{\text{tgt}} = \varepsilon_\theta(z_\tau^{\text{tgt}}, \tau, c^{\text{tgt}})$
 - 4: $\{\widehat{Q}^{\text{lay}}, \widehat{K}^{\text{lay}}, \widehat{V}^{\text{lay}}\} =$
 - 5: SelfEdit($\{Q^{\text{src}}, K^{\text{src}}, V^{\text{src}}\}, \{Q^{\text{tgt}}, K^{\text{tgt}}, V^{\text{tgt}}\}, \tau$)
 - 6: $\varepsilon^{\text{lay}}, M^{\text{lay}} = \varepsilon_\theta(z_\tau^{\text{lay}}, \tau, c^{\text{src}}; \{\widehat{Q}^{\text{lay}}, \widehat{K}^{\text{lay}}, \widehat{V}^{\text{lay}}\})$
 - 7: $\widehat{M}^{\text{tgt}} = \text{CrossEdit}(M^{\text{lay}}, M^{\text{tgt}}, \tau)$
 - 8: $\widehat{\varepsilon}^{\text{tgt}} = \varepsilon_\theta(z_\tau^{\text{tgt}}, \tau, c^{\text{tgt}}; \widehat{M}^{\text{tgt}})$
 - 9: $z_{\tau+1}^{\text{src}}, z_{\tau+1}^{\text{tgt}}, z_{\tau+1}^{\text{lay}} =$
 - 10: Sample($[z_\tau^{\text{src}}, z_\tau^{\text{tgt}}, z_\tau^{\text{lay}}], [\varepsilon^{\text{src}}, \widehat{\varepsilon}^{\text{tgt}}, \varepsilon^{\text{lay}}], \tau$)
 - 11: $m^{\text{tgt}} = \text{Threshold}[M_\tau^{\text{tgt}}(w^{\text{tgt}}), a^{\text{tgt}}]$
 - 12: $m^{\text{src}} = \text{Threshold}[M_\tau^{\text{src}}(w^{\text{src}}), a^{\text{src}}]$
 - 13: $z_{\tau+1}^{\text{tgt}} = (1 - m^{\text{tgt}} + m^{\text{src}}) \odot z_{\tau+1}^{\text{src}} + (m^{\text{tgt}} - m^{\text{src}}) \odot z_{\tau+1}^{\text{tgt}}$
 - 14: **Output:** $z_{\tau+1}^{\text{src}}, z_{\tau+1}^{\text{tgt}}, z_{\tau+1}^{\text{lay}}$
-

5. Experiments

5.1. Experiment Setups

Benchmarks We used established benchmarks to evaluate our proposed image editing method:

- **Language-Guided Image Editing.** We evaluate on the PIE-Bench introduced by Ju et al. [15], which assesses language-guided image editing in 9 different scenarios.
- **Image-to-Image (I2I) Translation.** We also evaluate on the I2I tasks at the scene-level (Summer \leftrightarrow Winter) and object-level (Horse \leftrightarrow Zebra) [40].

Evaluation Metrics We employ 4 distinct evaluation metrics to assess the generated image’s quality, the accuracy of the translation, the consistency against the source images, and the efficiency of the editing process.

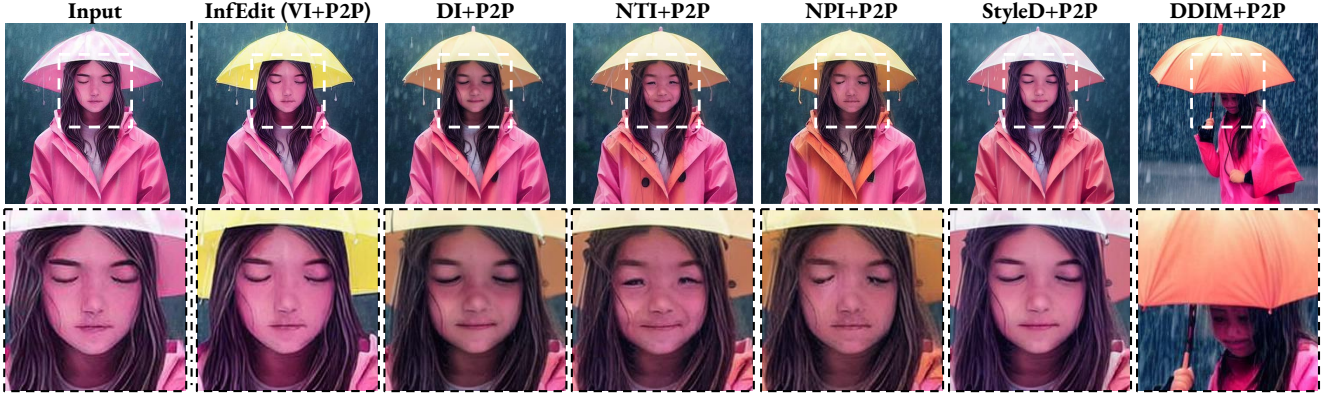
- **Image Quality.** We use the Fréchet Inception Distance (FID) [11] score, which compares the model outputs to real image distributions;
- **Translation Quality.** We use the CLIPScore [10] to quantify the semantic similarity of the generated image and target prompt with CLIP [26].
- **Translation Consistency.** We measure translation consistency using four different metrics: Mean Squared Error (MSE), Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM) [33], and Learned Perceptual Image Patch Similarity (LPIPS) [38].
- **Efficiency.** We directly compare the computation time on one A40 GPU for the inversion and forward process, as well as the number of sampling steps.

5.2. Inversion v.s. Inversion-Free Comparison

In this section, we present experiments to demonstrate that **inversion-free image editing (InfEdit) competes with the effectiveness of inversion-based methods, while also being significantly more efficient.** Recall that our InfEdit framework adopts Virtual Inversion (VI) derived from DDCM as the sampling framework, and takes any attention control for language-guided editing. We compare VI against other inversion-based methods on PIE-Bench, with 2 variants of InfEdit (VI+P2P and VI+UAC) for ablation. The inversion baselines we considered include DDIM [29], Null-Text (NT) [21], Negative Prompt (NP) [20], StyleDiffusion (StyleD) [17], CycleDiffusion (CycleD) [34], and Direct Inversion (DI) [15].

As depicted in Table 1, InfEdit competes with and often surpasses the effectiveness of inversion-based methods, especially in terms of background consistency. Additionally, InfEdit is significantly more efficient than inversion-based methods. On one hand, InfEdit does not require any inversion time. On the other hand, InfEdit generates high-quality target images with much fewer forward steps, and is even compatible with the LCM. We also present a qualitative example in Figure 6, demonstrating that with the same P2P attention control, VI allows faithful semantic changes as well as better consistency.

It is important to note a fundamental trade-off between reducing image editing distance and improving the faithfulness of image editing. StyleD [17], while demonstrating lower editing distances, exhibits limitations in effective editing, as evidenced by its scores in background preservation and CLIP similarity metrics. DI [15] and CycleD [34], surpassing InfEdit in structure distance, often fail to comply with editing instructions, leaving the source image untouched.



A girl with a pink yellow umbrella in the rain.

Figure 6. A qualitative example for ablation over inversion methods. With the same P2P attention control, InfEdit (VI+P2P) allows faithful semantic changes as well as better consistency.

Method		Structure	Background Preservation				CLIP Similarity		Efficiency (sec / #)		
Inverse	Edit	Distance $_{10^3}$ ↓	PSNR ↑	LPIPS $_{10^3}$ ↓	MSE $_{10^4}$ ↓	SSIM $_{10^2}$ ↑	Whole ↑	Edited ↑	Inverse Time ↓	Forward Time ↓	Steps ↓
DDIM	P2P	69.43	17.87	208.80	219.88	71.14	25.01	22.44	10.93 ± 0.01	12.79 ± 0.01	50
CycleD	P2P	6.06	28.25	43.96	25.85	85.61	23.68	20.87	N/A	4.55 ± 0.02	32
	NT	13.44	27.03	60.67	35.86	84.11	24.75	21.86	132.39 ± 7.69	12.90 ± 0.01	50
	NP	16.17	26.21	69.01	39.73	83.40	24.61	21.87	4.14 ± 0.00	12.78 ± 0.01	50
StyleD	P2P	11.65	26.05	66.10	38.63	83.42	24.78	21.72	810.17 ± 7.77	28.18 ± 1.30	50
	DI	11.65	27.22	54.55	32.86	84.76	25.02	22.10	16.83 ± 0.02	12.87 ± 0.01	50
	VI	14.22	27.52	47.98	34.17	85.05	24.89	22.03	N/A	4.50 ± 0.01	32
	VI*	15.61	26.64	55.85	41.15	84.66	24.57	21.69	N/A	2.60 ± 0.00	15
	VI*	13.78	28.51	47.58	32.09	85.66	25.03	22.22	N/A	2.22 ± 0.02	12

* Using the Latent Consistency Model (LCM) as the base model. Otherwise, Stable Diffusion (SD) v1.4 is adopted.

Table 1. Aggregated performances of different image inversion and editing methods on PIE-bench. We break InfEdit into the Virtual inversion (VI) and arbitrary choices of attention control mechanism and diffusion backbone. VI competes and even surpasses other inversion methods with the same P2P attention control. The integration of unified attention control (UAC) and the LCM backbone further enhances its performance. Notably, InfEdit runs about an order of magnitude faster than most of the baselines on one single A40.

5.3. Attention Control Comparison

In this section, we present experiments to demonstrate that **with unified attention control (UAC), InfEdit establishes state-of-the-art performance in terms of editing quality, consistency, and efficiency.** We compare UAC with other attention control baselines, especially Prompt-to-Prompt (P2P) [9], Plug-and-Play (PnP) [31], and Mutual Self-Attention Control (MasaCtrl) [4].

The comprehensive analysis across 9 distinct categories of editing tasks in PIE-Bench demonstrates the superior performance of InfEdit with UAC as the attention control mechanism, evidenced by enhanced editing quality in Figure 7a and improved image consistency in Figure 7b. Qualitative comparisons are provided in Figure 8, and more results are available in the Appendix.

5.4. Image-to-Image Translation Tasks

We further evaluate InfEdit (VI+UAC) in scene-level and object-level I2I translation tasks for more general compar-

isons. The baselines we considered include Text2LIVE [3], SDEdit [19], CycleD [34], NT [21], MasaCtrl [4], as well as the training-based state-of-the-art CycleNet [35]. As shown in Table 3, InfEdit strikes an effective balance between translation effects and consistency. Qualitative examples are shown in Figure 18 and 19 in Appendix.

5.5. Computational Efficiency Ablation

We finally study the editing efficiency of InfEdit. As shown in Table 1, InfEdit significantly outperforms the other baselines in terms of computational efficiency even without applying the LCM. We perform an ablation study to demonstrate that **InfEdit gains an advantage through its distinctive compatibility with latent consistency models, facilitating both efficient and high-quality image editing.** We compare the base diffusion backbones, Stable Diffusion (SD) 1.4 [28] and Latent Consistency Model (LCM) [18], by the CLIP Scores at different forward steps. Table 2 (also visualized in Figure 7c) shows that the LCM significantly

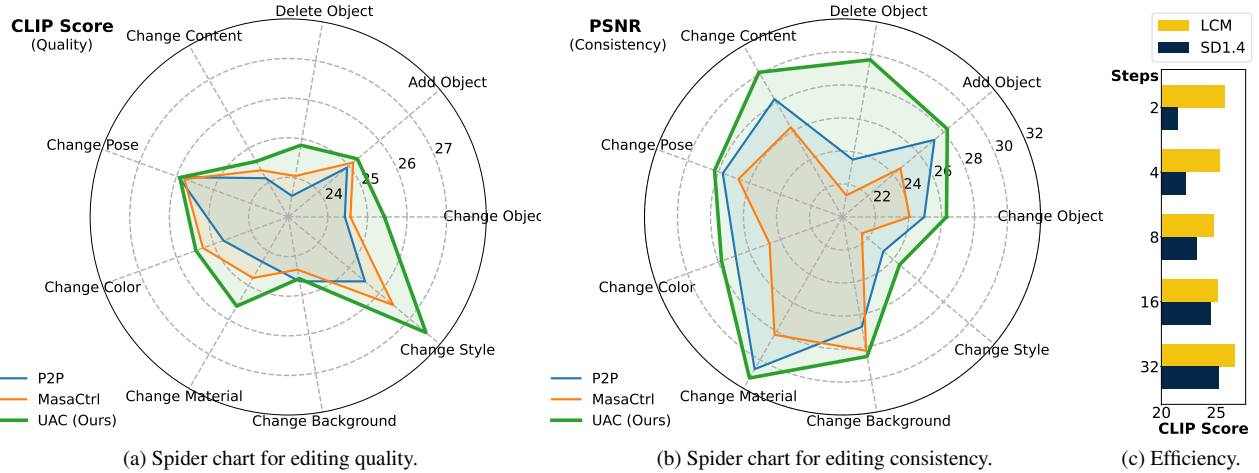


Figure 7. A comprehensive performance evaluation on the PIE-bench. We present spider charts of editing quality (CLIP Scores) and consistency (PSNR) across 9 editing tasks for Prompt-to-Prompt (P2P), MasaCtrl, and Unified Attention Control (UAC) methods. Accompanied by an analysis of editing efficiency for Stable Diffusion (SD) 1.4 and Latent Consistency Model (LCM) across different steps.



Figure 8. Qualitative comparisons of InfEdit (VI+UAC) against baselines. InfEdit attains editing goals with the best consistency.

outperforms SD 1.4 in editing quality, even with fewer sampling steps. While the CLIP Scores of SD incrementally improves from 21.47 to 25.18 as the number of forward steps increases from 2 to 32, LCM can achieve consistently higher CLIP Scores across varying step counts, showing superior speed in image editing.

5.6. Connecting InfEdit to LLMs

Based on the capability of InfEdit to edit images with natural language, we can further leverage large language

Method	CLIP Score				
	2 steps	4 steps	8 steps	16 steps	32 steps
InfEdit (VI+P2P)					
SD 1.4	21.47	22.17	23.16	24.45	25.18
LCM	25.76	25.33	24.76	25.08	26.68

Table 2. The use of LCM backbone in InfEdit allows a high CLIP Score even with fewer steps.

models (LLMs) to follow the image editing instructions. Through our experiments, we have validated the feasibility of prompting GPT-4 [25] to break down editing instructions into adequate source and target prompts for InfEdit. This

Task	Summer↔Winter (512 × 512)					Horse↔Zebra (512 × 512)				
Method	FID ↓	CLIP Sim ↑	LPIPS ↓	PSNR ↑	SSIM ↑	FID ↓	CLIP Sim ↑	LPIPS ↓	PSNR ↑	SSIM ↑
CycleNet	79.79	24.12	0.15	25.88	0.69	76.83	25.27	0.08	26.21	0.74
Text2LIVE	86.12	25.98	0.27	16.83	0.68	103.14	31.55	0.16	20.98	0.81
SDEdit	90.51	23.26	0.30	18.59	0.43	63.04	27.97	0.33	18.49	0.44
CycleD	84.52	24.40	0.24	21.66	0.68	41.17	29.09	0.29	19.41	0.61
NT+P2P	92.65	24.82	0.24	20.19	0.66	106.83	26.57	0.21	21.45	0.66
MasaCtrl	114.83	17.11	0.37	14.66	0.43	239.61	21.15	0.41	16.31	0.37
InfEdit	75.63	23.07	0.18	21.99	0.68	61.81	28.16	0.16	21.80	0.72

Table 3. Image2Image translation comparison. InfEdit methods achieve a favorable balance between consistency and translation quality.

allows users to give natural language instructions to control image editing, which improves the user experience. A Gradio demo is available on our project page.

6. Related Work

Image Manipulation with Diffusion Models Diffusion models (DMs) have achieved notable success in image generation [13, 29], with large-scale models pre-trained on text-to-image tasks [23, 28]. These models are increasingly adopted for image manipulation tasks, where DMs are augmented with additional conditions like text prompts [27] or images [35, 37] to generate the target image. The source image information is usually integrated into DMs through an inversion process [29] or via a side network [22, 37]. Additionally, mask-based methods have been proposed, utilizing either user-prompted or automatically generated masks [2, 6, 24], or augmentation layers [3], to facilitate more controlled and precise image manipulations. To enhance the consistency and quality of image edits, several techniques have been developed. Among these, attention control mechanisms [4, 9, 31] have emerged as a promising direction, especially when they are paired with inversion methods [15, 20, 21].

Inversion in Diffusion Models In DMs, real image editing methods usually rely on the inversion process, which produces a latent representation that can reconstruct the image through the generative process. Initially, SDEdits [19] was proposed, which adds random Gaussian noise to the source image as input but suffers from reconstruction quality. DDIM inversion [9] was then introduced for its deterministic mapping from latent space to image but is prone to errors accumulated in its multiple-step inversion process. Null-text inversion [21] used a null-text prompt for pivot tuning, improving real image editing but was time-consuming and not fully accurate. Negative prompt inversion [20] accelerated the inversion process by approximating the DDIM inversion, while sacrificing the reconstruction quality. CycleDiffusion [34], Editing-friendly Inversion [14] and Direct Inversion [15] use source latents from each inversion step as reference for editing the target

branch. However, these methods still struggle with the cumulative errors typical of the inversion process and tend to be slower overall due to the inherent need for inversion.

Attention-Control for Image Editing In the realm of zero-shot image editing, attention control becomes a pivotal technique to preserve consistency while manipulating visual content. Recent works like Prompt-to-Prompt (P2P) [9] and Plug-and-Play (PnP) [31] have contributed significantly to this field by replacing cross-attention and self-attention maps to maintain the original image layout and spatial information, thereby preserving the consistency during editing. In contrast, MasaCtrl [4] offers an alternative approach that enables the modification of layout and spatial attributes while safeguarding the semantic content inherent in the image, which addresses the limitation of conflating spatial edits with semantic preservation. P2Plus [17] extends the prompt-to-prompt paradigm by applying edits to both the text-conditional and unconditional branches during classifier-free guidance [12], thus offering a more comprehensive editing framework.

7. Conclusion

Recent advancements in inversion-based editing notwithstanding, text-guided image manipulation using diffusion models continues to be a challenge. The main challenges involve 1) the lengthy inversion process; 2) difficulties in maintaining both consistency and accuracy; 3) incompatibility with the efficient consistency sampling methods of consistency models. In response, we questioned whether it’s possible to bypass the inversion process in editing. Our findings reveal that with a known initial sample z_0 , a specific variance schedule σ can simplify the denoising step to a form akin to multi-step consistency sampling. This led to the development of the Denoising Diffusion Consistent Model (DDCM), which effectively introduces a virtual inversion strategy that eliminates the need for explicit inversion during sampling. Moreover, we present the Unified Attention Control (UAC) mechanisms as a tuning-free framework for text-guided editing. This integration forms the basis of our inversion-free editing approach, InfEdit, which

facilitates consistent and accurate editing across both rigid and non-rigid semantic transformations. InfEdit is adept at handling complex modifications without compromising the image’s integrity or requiring explicit inversion. Extensive experiments demonstrate its robust performance across a range of editing tasks and that it maintains a smooth workflow, completing tasks in under 3 seconds on a single A40. InfEdit unleashes the potential for real-time image editing applications.

Ethics Statement

While InfEdit offers promising advancements in image editing, it is crucial to consider its broader ethical, legal, and societal implications.

Copyright Infringement. As an advanced image editing tool, InfEdit could be used to modify and repurpose artists’ original works, raising concerns over copyright violations. It’s vital for practitioners to respect the rights of creators and maintain the integrity of the creative economy, ensuring adherence to licensing and copyright laws.

Deceptive Misuse. If exploited by nefarious entities, InfEdit’s capability to generate convincing image alterations could be used for misinformation, fraud, or identity theft. This necessitates responsible user guidelines and strong security protocols to prevent such misuse and safeguard against security threats.

Bias and Fairness. Furthermore, InfEdit builds upon pre-trained latent diffusion models and latent consistency models, which might carry inherent biases, leading to potential fairness issues. While the method is algorithmic and not pre-trained on large web-scale datasets, it’s important to recognize and mitigate any encoded biases in these pre-trained backbones to ensure fairness and ethical use.

By proactively addressing these concerns, we can leverage InfEdit’s capabilities responsibly, prioritizing ethical considerations, legal compliance, and the welfare of society. This approach is essential for advancing technology while safeguarding our community’s values and trust.

References

- [1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. [6](#)
- [2] Omri Avrahami, Thomas Hayes, Oran Gafni, Sonal Gupta, Yaniv Taigman, Devi Parikh, Dani Lischinski, Ohad Fried, and Xi Yin. Spatext: Spatio-textual representation for controllable image generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. [2](#), [10](#)
- [3] Omer Bar-Tal, Dolev Ofri-Amar, Rafail Fridman, Yoni Kasten, and Tali Dekel. Text2live: Text-driven layered image and video editing. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XV*, pages 707–723. Springer, 2022. [2](#), [8](#), [10](#)
- [4] Mingdeng Cao, Xintao Wang, Zhongang Qi, Ying Shan, Xiaohu Qie, and Yinqiang Zheng. Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 22560–22570, 2023. [5](#), [6](#), [8](#), [10](#)
- [5] Jooyoung Choi, Sungwon Kim, Yonghyun Jeong, Youngjune Gwon, and Sungroh Yoon. Ilvr: Conditioning method for denoising diffusion probabilistic models. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 14347–14356. IEEE Computer Society, 2021. [2](#)
- [6] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. In *The Eleventh International Conference on Learning Representations*, 2023. [2](#), [10](#)
- [7] Wenkai Dong, Song Xue, Xiaoyue Duan, and Shumin Han. Prompt tuning inversion for text-driven image editing using diffusion models. *arXiv preprint arXiv:2305.04441*, 2023. [2](#), [3](#)
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [2](#)
- [9] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or. Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Conference on Learning Representations*, 2023. [5](#), [6](#), [8](#), [10](#)
- [10] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7514–7528, 2021. [7](#)
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. [7](#)
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. [10](#)
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [2](#), [10](#)
- [14] Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddpm noise space: Inversion and manipulations, 2023. [10](#)
- [15] Xuan Ju, Ailing Zeng, Yuxuan Bian, Shaoteng Liu, and Qiang Xu. Direct inversion: Boosting diffusion-based editing with 3 lines of code. *arXiv preprint arXiv:2304.04269*, 2023. [2](#), [3](#), [7](#), [10](#)

- [16] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014. [2](#)
- [17] Senmao Li, Joost van de Weijer, Taihang Hu, Fahad Shahbaz Khan, Qibin Hou, Yaxing Wang, and Jian Yang. Stylediffusion: Prompt-embedding inversion for text-based editing. *arXiv preprint arXiv:2303.15649*, 2023. [2](#), [3](#), [7](#), [10](#)
- [18] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference, 2023. [3](#), [8](#)
- [19] 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. In *International Conference on Learning Representations*, 2021. [8](#), [10](#)
- [20] Daiki Miyake, Akihiro Iohara, Yu Saito, and Toshiyuki Tanaka. Negative-prompt inversion: Fast image inversion for editing with text-guided diffusion models. *arXiv preprint arXiv:2305.16807*, 2023. [7](#), [10](#)
- [21] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6038–6047, 2023. [2](#), [3](#), [7](#), [8](#), [10](#)
- [22] 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. *arXiv preprint arXiv:2302.08453*, 2023. [2](#), [10](#)
- [23] 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. In *International Conference on Machine Learning*, 2021. [10](#)
- [24] Alexander Quinn Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *International Conference on Machine Learning*, pages 16784–16804. PMLR, 2022. [2](#), [10](#)
- [25] OpenAI. Gpt-4 technical report, 2023. [9](#)
- [26] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. [7](#)
- [27] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *ArXiv*, abs/2204.06125, 2022. [2](#), [10](#)
- [28] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. [2](#), [8](#), [10](#)
- [29] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020. [2](#), [3](#), [7](#), [10](#)
- [30] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint arXiv:2303.01469*, 2023. [2](#)
- [31] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1921–1930, 2023. [8](#), [10](#)
- [32] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [5](#)
- [33] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing*, 13(4):600–612, 2004. [7](#)
- [34] Chen Henry Wu and Fernando De la Torre. A latent space of stochastic diffusion models for zero-shot image editing and guidance. In *ICCV*, 2023. [2](#), [3](#), [7](#), [8](#), [10](#)
- [35] Sihan Xu, Ziqiao Ma, Yidong Huang, Honglak Lee, and Joyce Chai. Cyclenet: Rethinking cycle consistent in text-guided diffusion for image manipulation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. [2](#), [8](#), [10](#)
- [36] Fengyu Yang, Jiacheng Zhang, and Andrew Owens. Generating visual scenes from touch. *International Conference on Computer Vision (ICCV)*, 2023. [2](#)
- [37] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. [2](#), [10](#)
- [38] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. [7](#)
- [39] Min Zhao, Fan Bao, Chongxuan Li, and Jun Zhu. Egsde: Unpaired image-to-image translation via energy-guided stochastic differential equations. *Advances in Neural Information Processing Systems*, 35:3609–3623, 2022. [2](#)
- [40] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017. [7](#)



Figure 9. Additional comparison on changing object tasks.

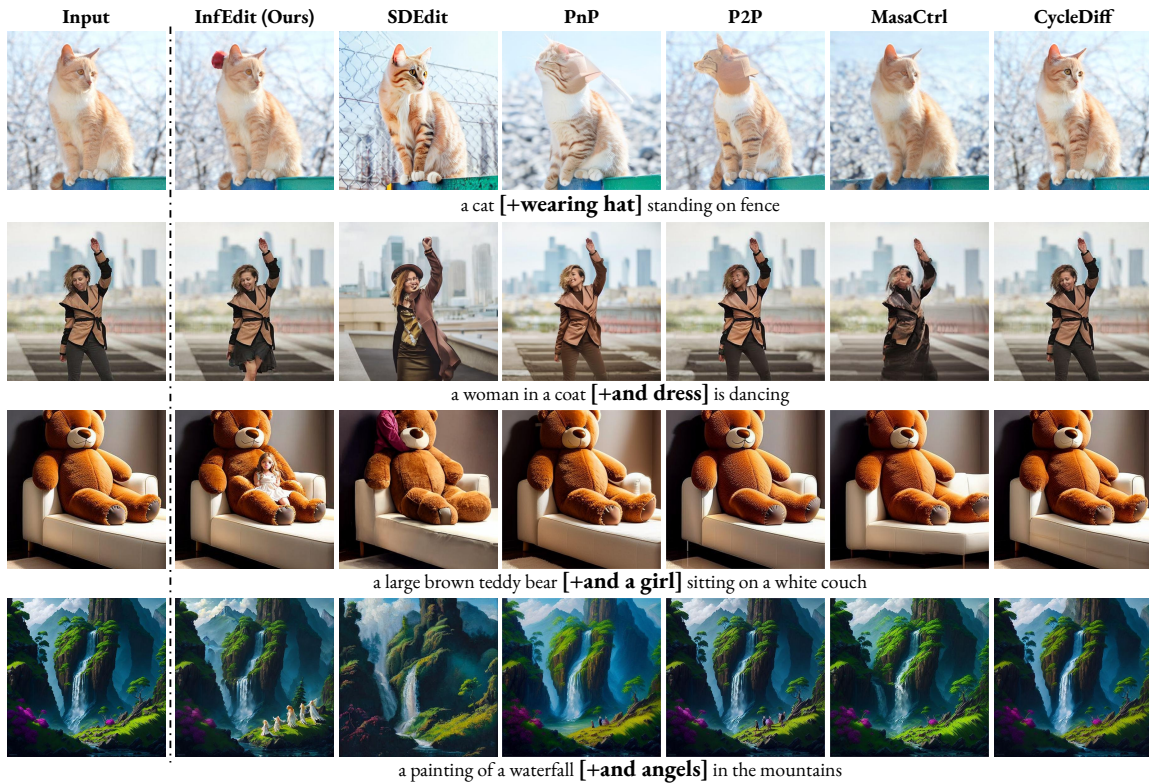


Figure 10. Additional comparison on adding object tasks.

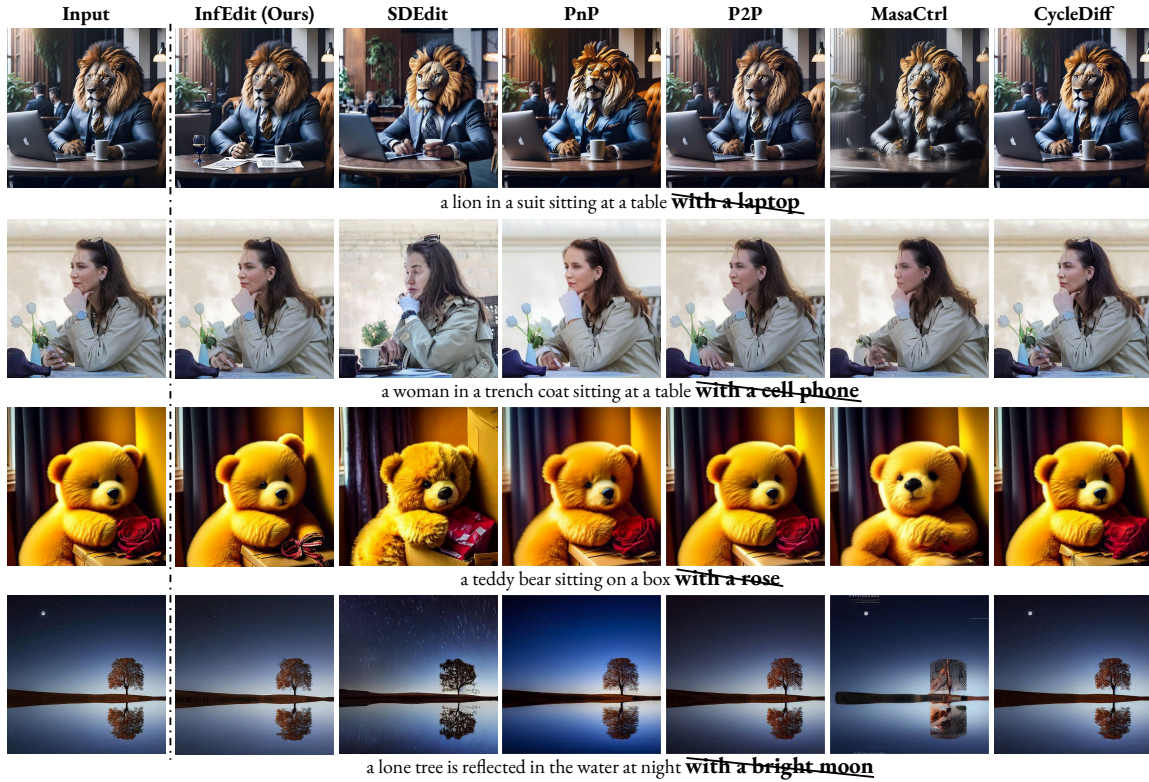


Figure 11. Additional comparison on deleting object tasks.

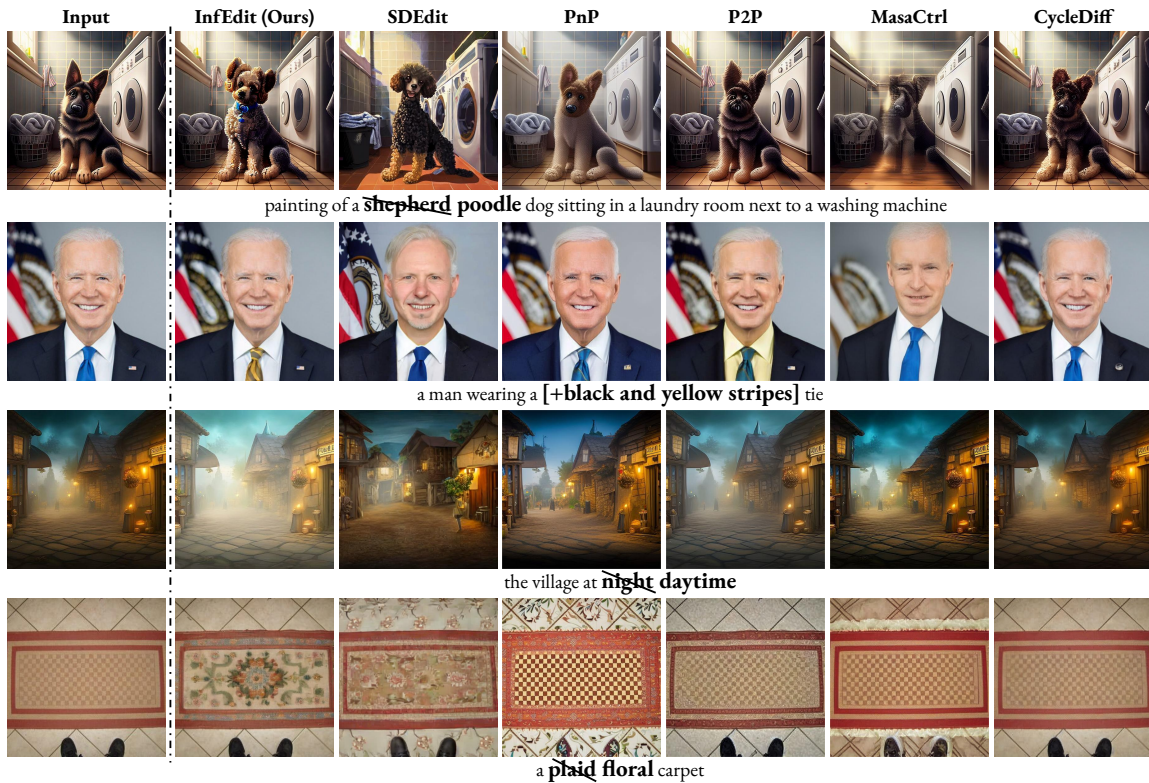


Figure 12. Additional comparison on changing content tasks.

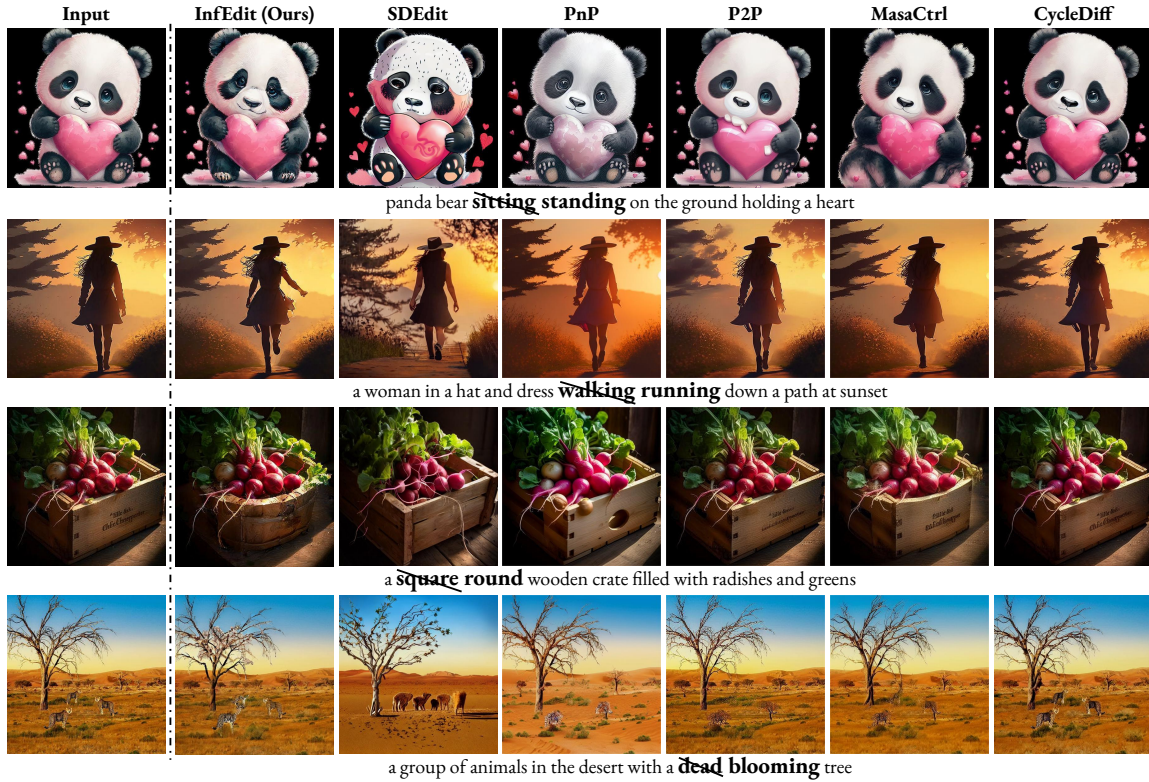


Figure 13. Additional comparison on changing pose tasks.

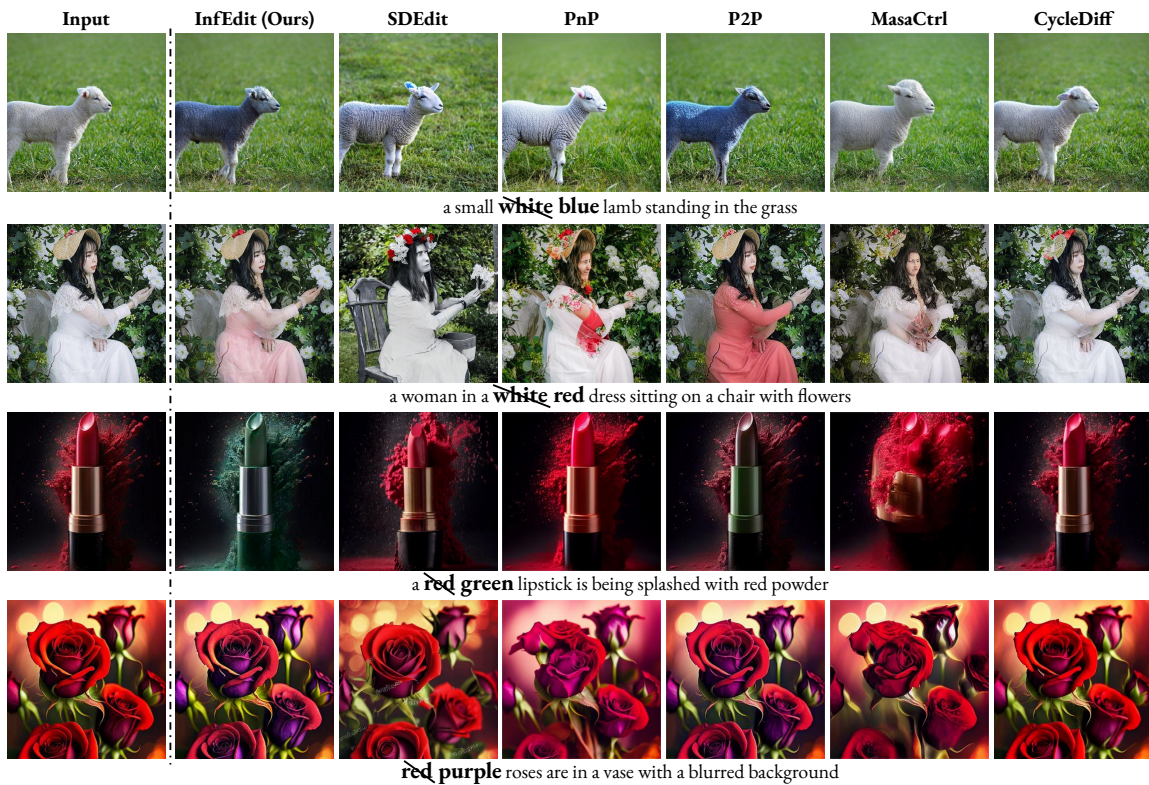


Figure 14. Additional comparison on changing color tasks.

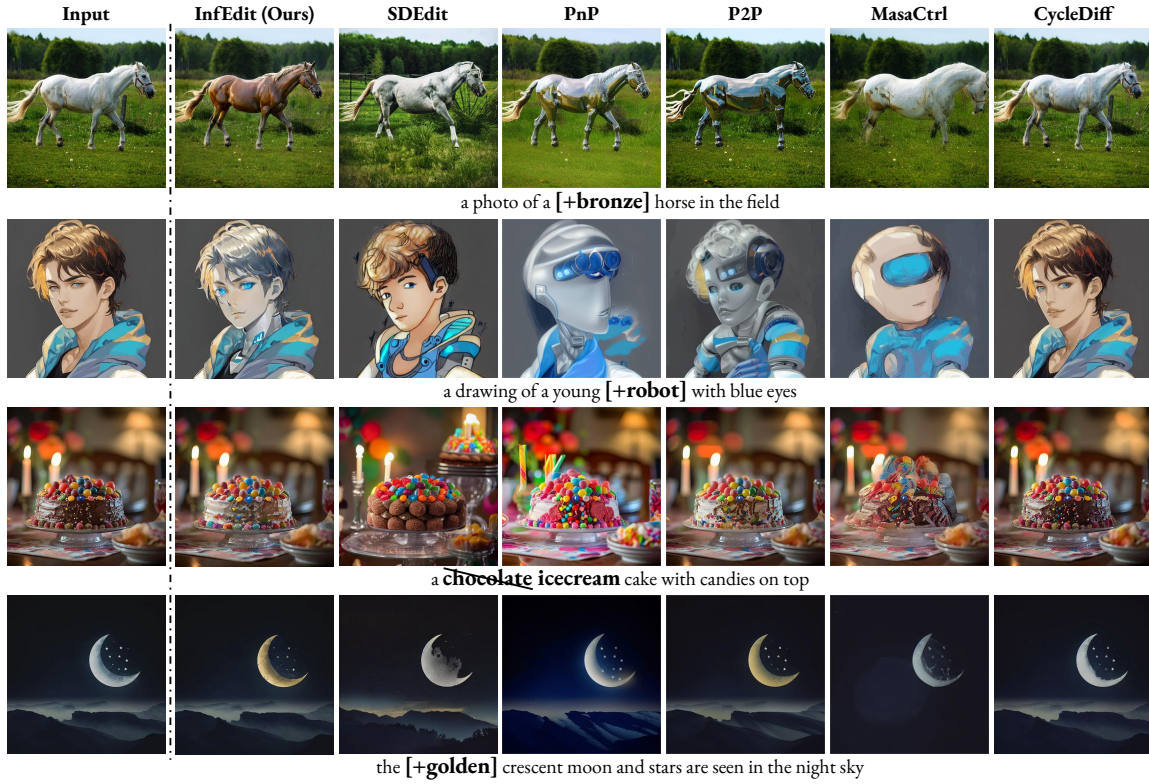


Figure 15. Additional comparison on changing material tasks.

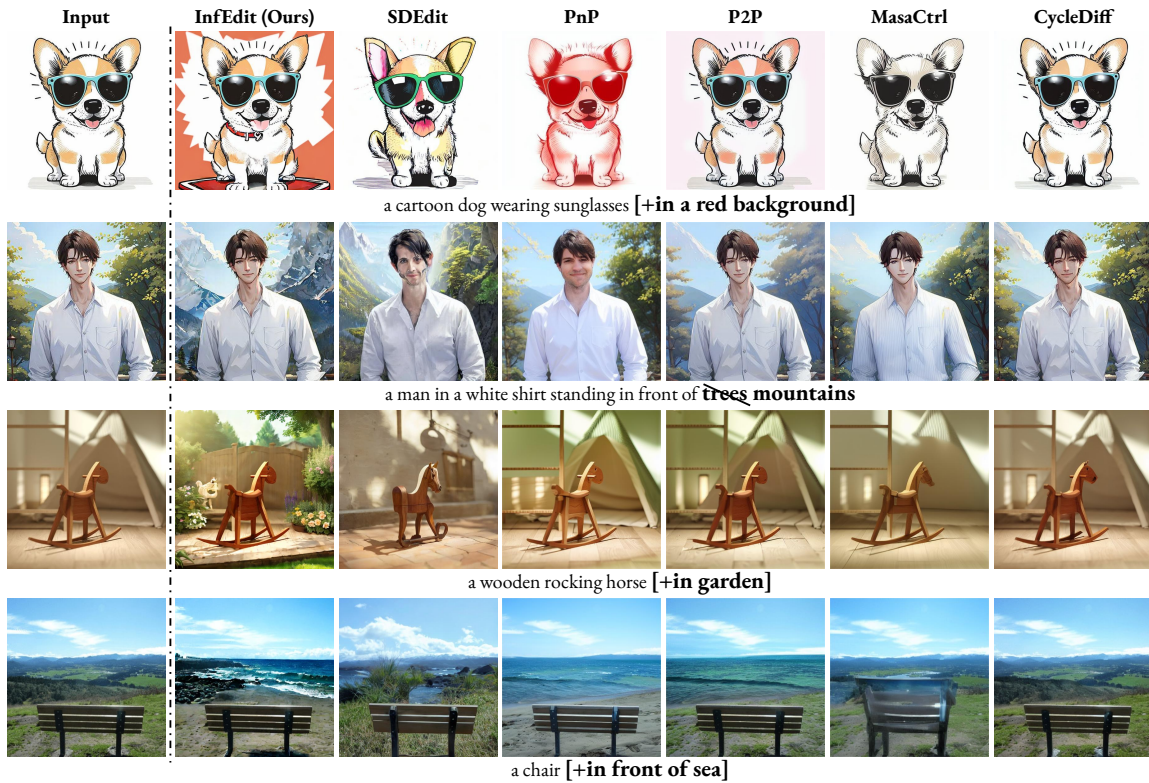


Figure 16. Additional comparison on changing background tasks.

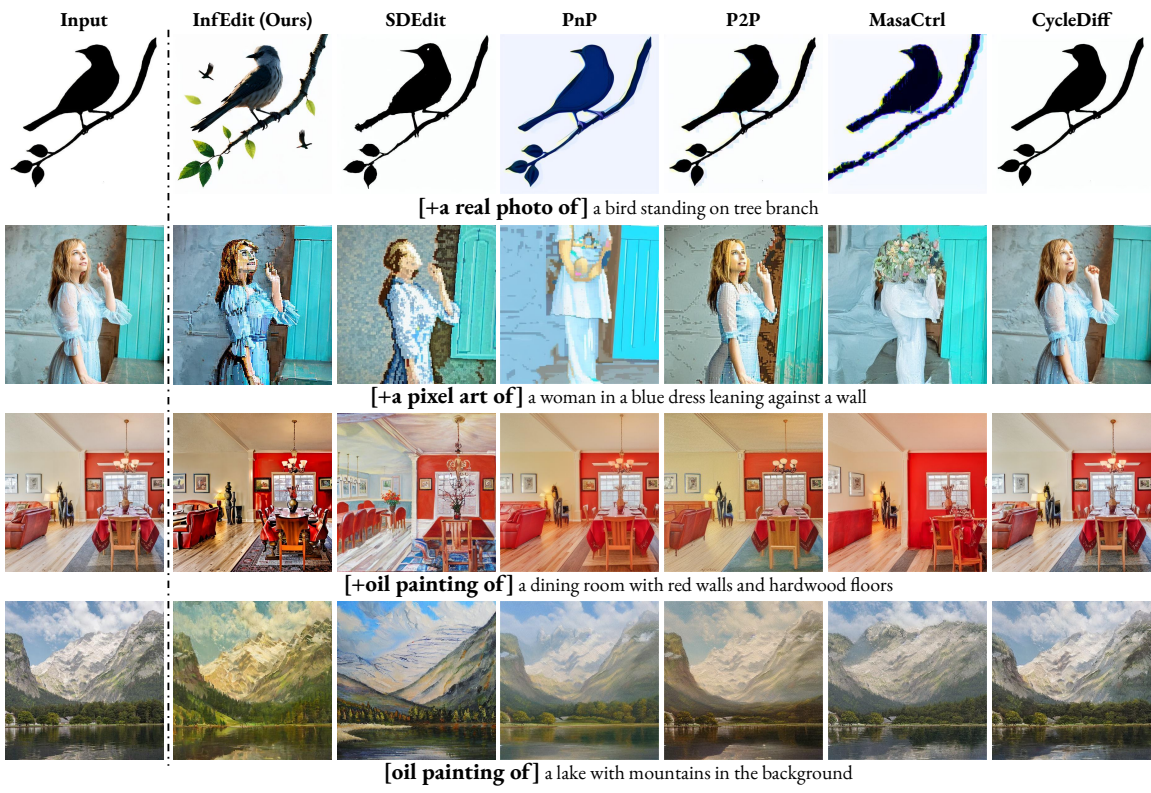


Figure 17. Additional comparison on changing style tasks.

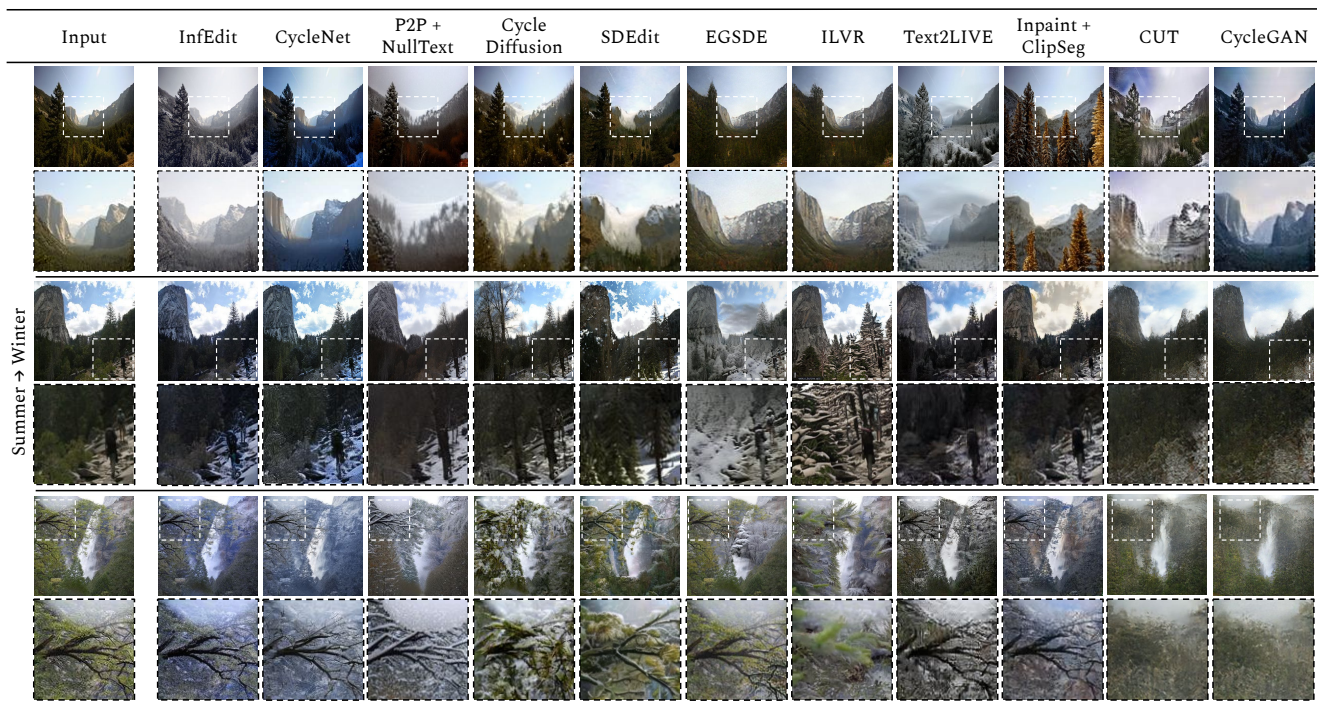


Figure 18. Qualitative comparison of InfEdit on Summer2Winter task with other baselines.

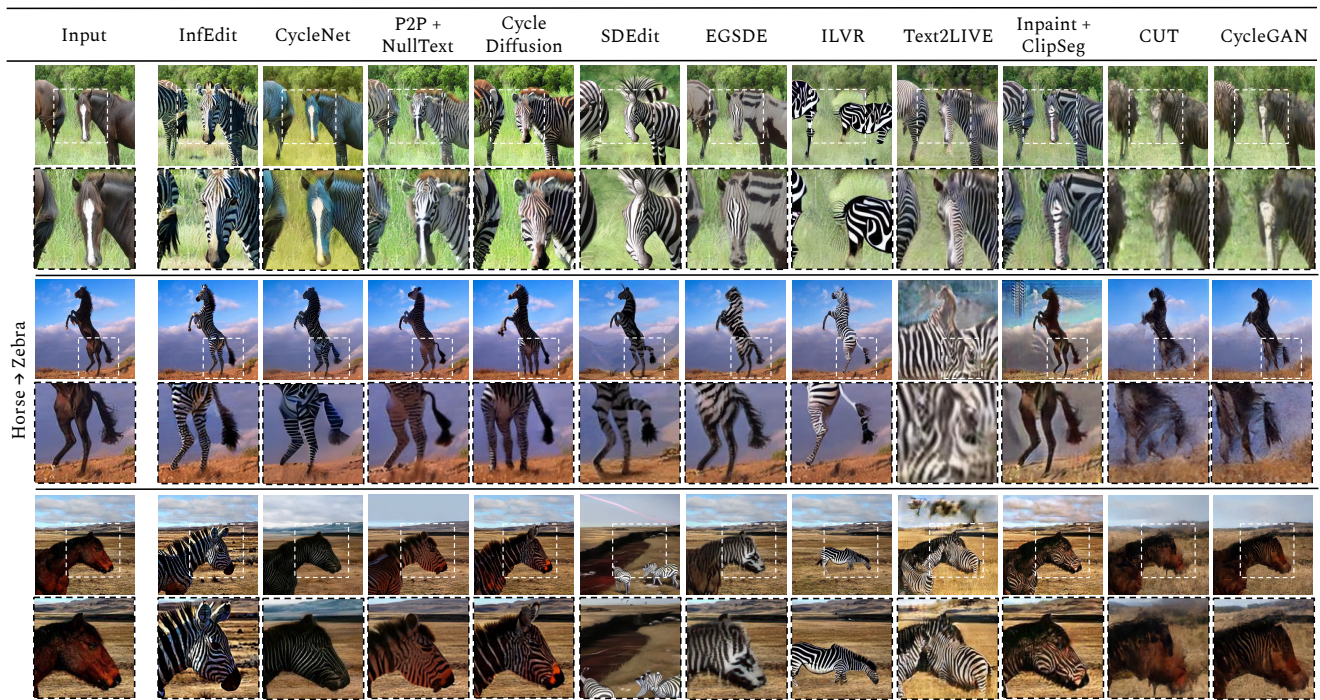


Figure 19. Qualitative comparison of InfEdit on Horse2Zebra task with other baselines.

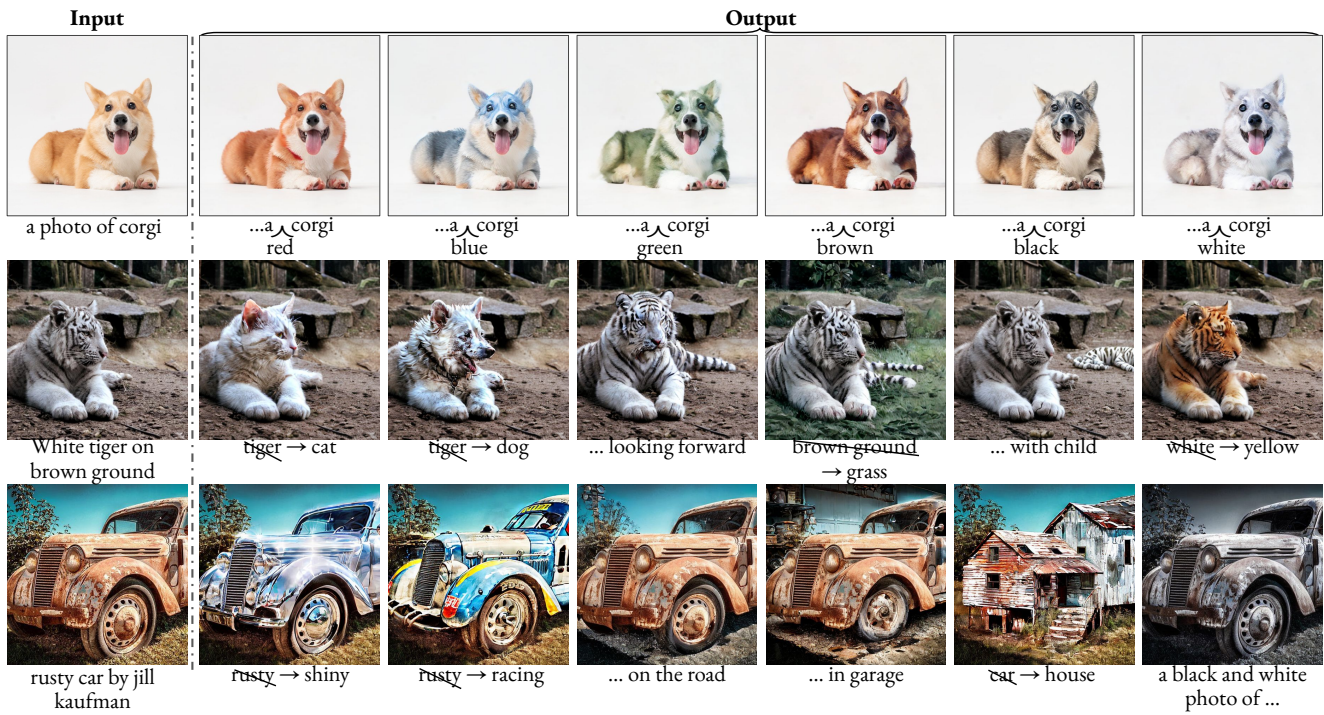


Figure 20. Additional results of InfEdit in various complex image editing tasks.

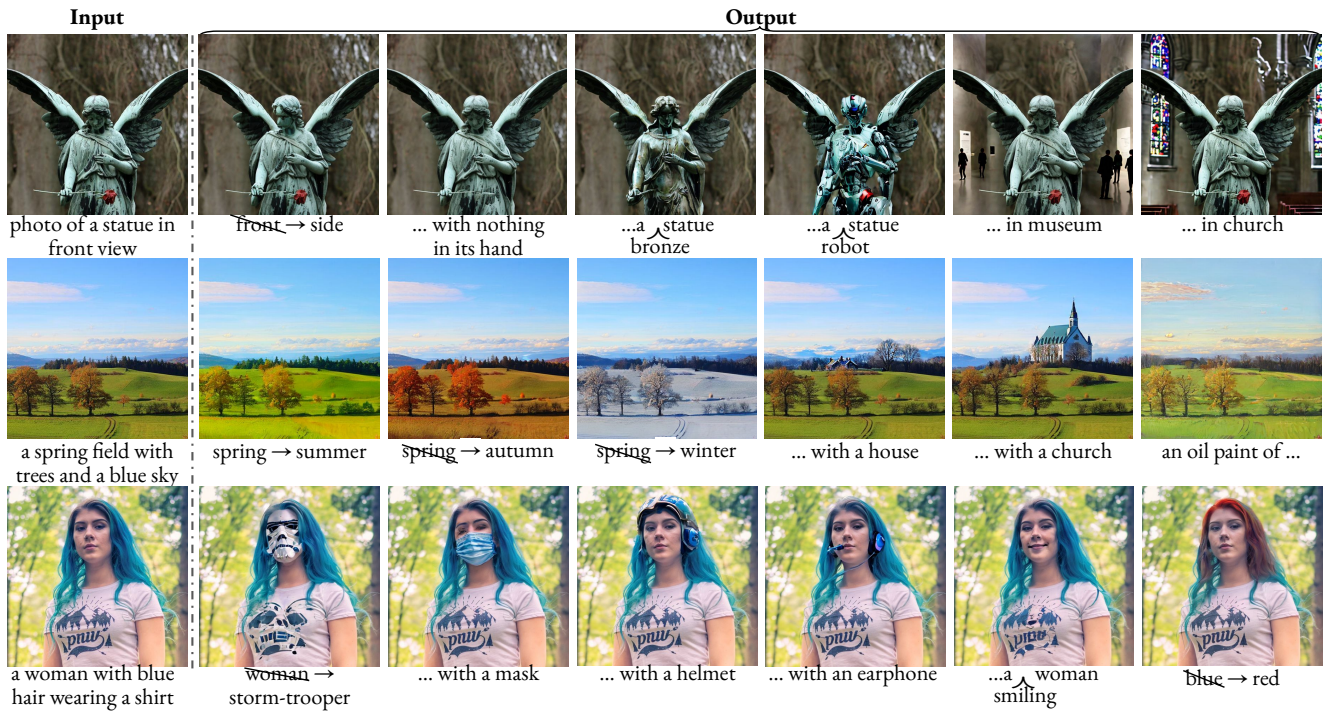


Figure 21. Additional results of InfEdit in various complex image editing tasks.

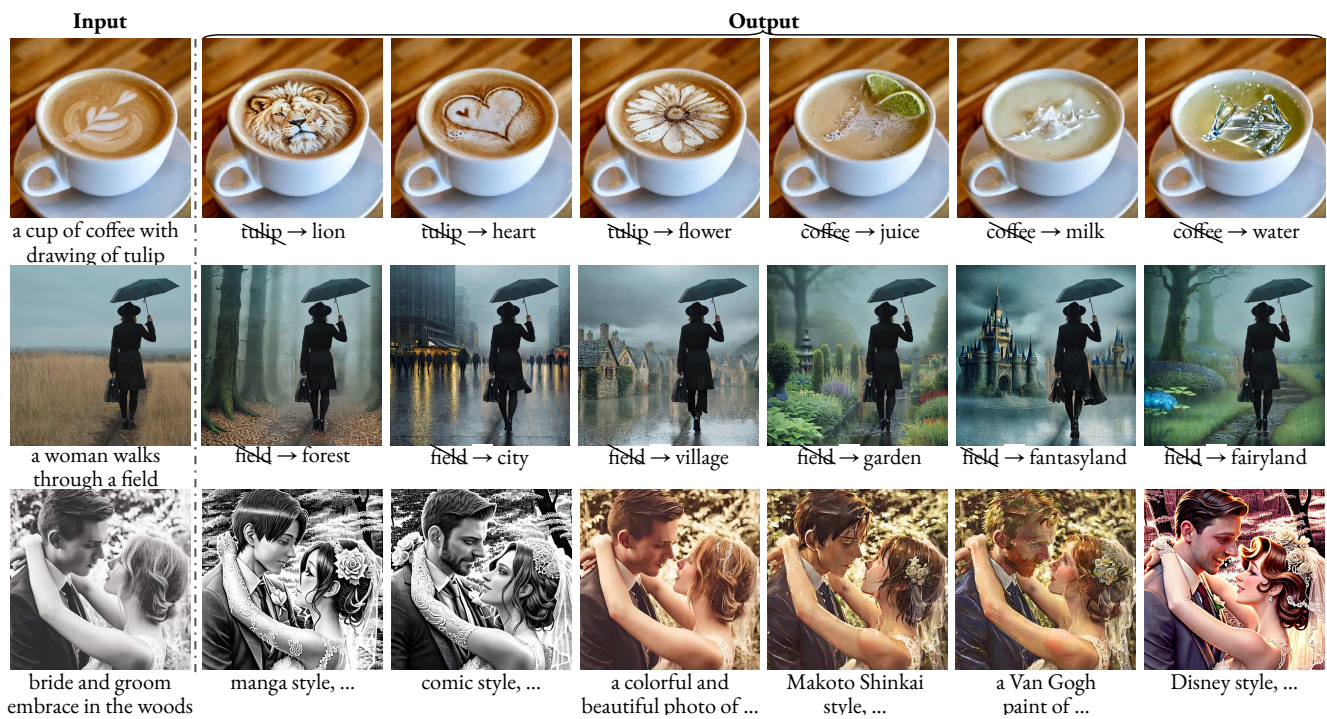


Figure 22. Additional results of InfEdit in various complex image editing tasks.

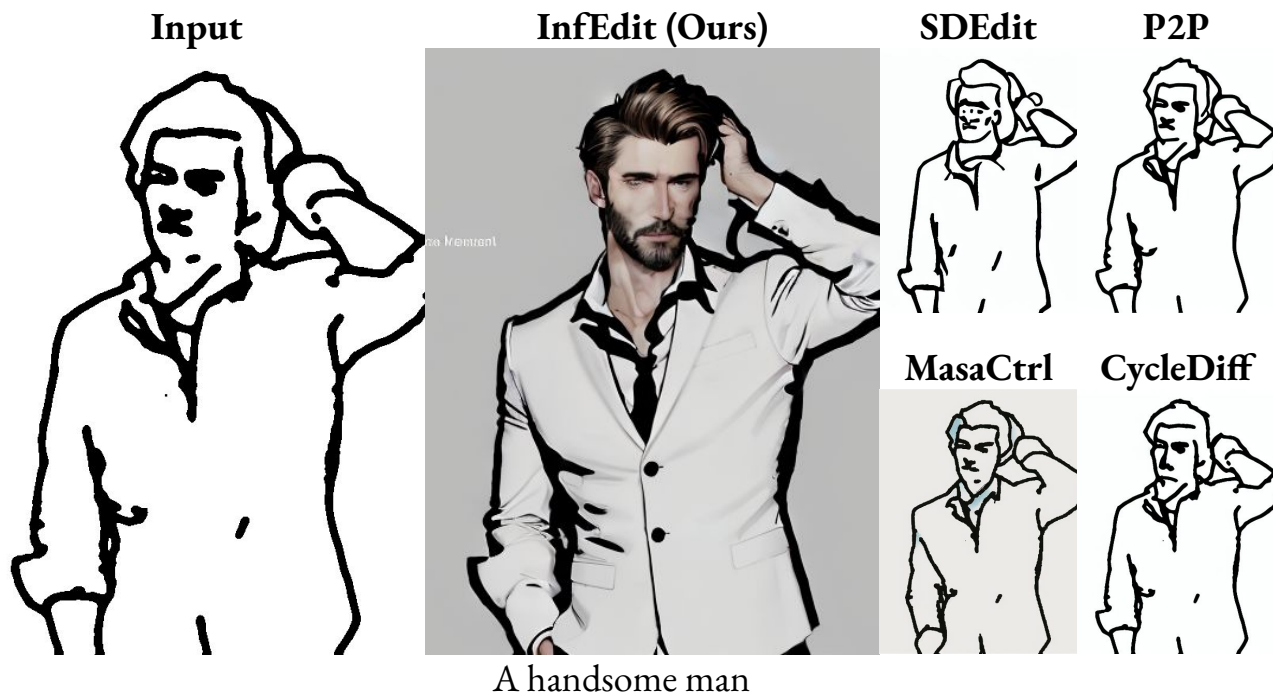


Figure 23. Additional results of InfEdit compared with other method.

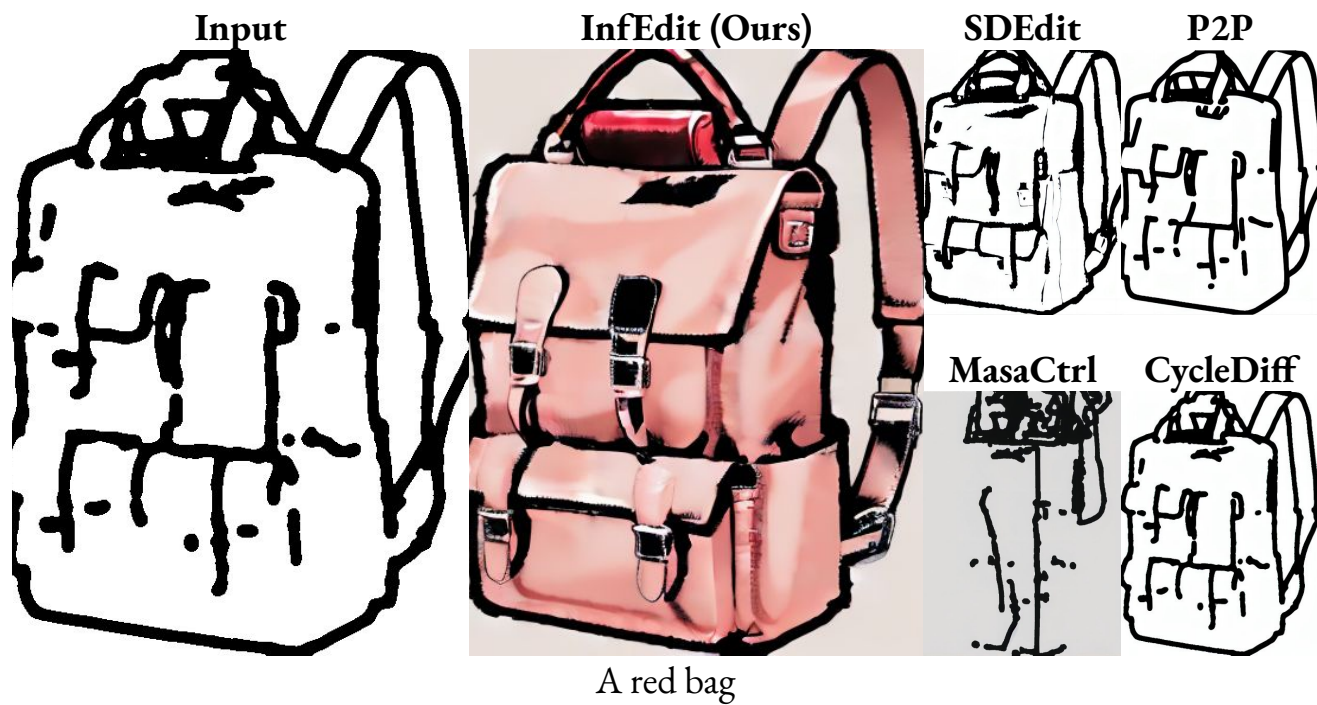


Figure 24. Additional results of InfEdit compared with other method.




Figure 25. Additional results of InfEdit in multi-modal editing.




Figure 26. Additional results of InfEdit in multi-modal editing.




 a anime girl with ~~green~~ orange hair.



 a anime girl with ~~shirt~~ lace skirt.



 a anime girl with ~~green~~ red eyes.



 a anime girl with ~~smile~~ angry face.




Figure 27. Multi-turn editing via InfEdit.