

# LEDITS++: Limitless Image Editing using Text-to-Image Models

Manuel Brack<sup>1,2\*†</sup> Felix Friedrich<sup>1,3\*</sup> Katharina Kornmeier<sup>1\*</sup> Linoy Tsaban<sup>4</sup>  
 Patrick Schramowski<sup>1,2,3</sup> Kristian Kersting<sup>1,2,3</sup> Apolinário Passos<sup>4</sup>  
<sup>1</sup>TU Darmstadt, <sup>2</sup>DFKI, <sup>3</sup>hessian.AI, <sup>4</sup>Huggingface  
 {brack, friedrich}@cs.tu-darmstadt.de

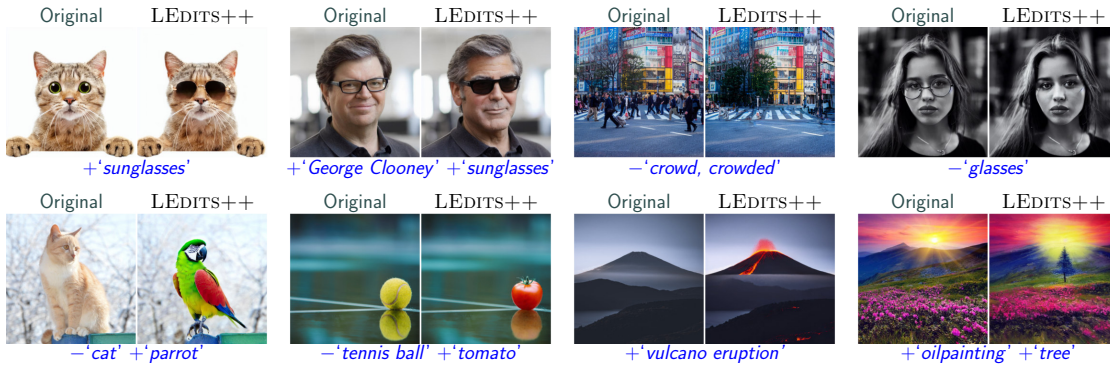


Figure 1. LEDITS++ facilitates versatile image-to-image editing. Several complex cases are available now.

## Abstract

Text-to-image diffusion models have recently received increasing interest for their astonishing ability to produce high-fidelity images from solely text inputs. Subsequent research efforts aim to exploit and apply their capabilities to real image editing. However, existing image-to-image methods are often inefficient, imprecise, and of limited versatility. They either require time-consuming fine-tuning, deviate unnecessarily strongly from the input image, and/or lack support for multiple, simultaneous edits. To address these issues, we introduce LEDITS++, an efficient yet versatile and precise textual image manipulation technique. LEDITS++’s novel inversion approach requires no tuning nor optimization and produces high-fidelity results with a few diffusion steps. Second, our methodology supports multiple simultaneous edits and is architecture-agnostic. Third, we use a novel implicit masking technique that limits changes to relevant image regions. We propose the novel TEdBench++ benchmark as part of our exhaustive evaluation. Our results demonstrate the capabilities of LEDITS++ and its improvements over previous methods.

## 1. Introduction

Text-to-image diffusion models (DM) have garnered recognition for their ability to generate high-quality images from textual descriptions. A growing body of research has recently been dedicated to utilizing these models for manipulating real images.

However, several barriers prevent many real-world applications of diffusion-based image editing. Current methods often entail computationally expensive model tuning or other optimization, presenting practical challenges [6, 18, 28, 30, 44]. Additionally, existing techniques have the proclivity to induce profound changes to the original image [17, 26], often resulting in completely different images. Lastly, all these approaches are inherently constrained when editing multiple (arbitrary) concepts simultaneously. We tackle these problems by introducing LEDITS++<sup>1</sup>, a diffusion-based image editing technique addressing these limitations.

LEDITS++<sup>2</sup> offers a streamlined approach for textual image editing, eliminating the need for extensive parameter tuning. To this end, we derive image inversion for a more efficient diffusion sampling algorithm to a) drastically reduce computational resources and b) guarantee perfect image re-

\*Equal contribution †Partially as research intern at Adobe

Proceedings of the 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

<sup>1</sup>LEDITS++ stands for *Limitless Edits* with sde-dpm-solver++.

<sup>2</sup><https://huggingface.co/spaces/leditsplusplus/project>

construction. Thus, we overcome computational obstacles and avoid changes in the edited image in the first place. Furthermore, we use a novel implicit masking approach to semantically ground each edit instruction to its relevant image region. This further optimizes changes to the image by retaining the overall image composition and object identity. Additionally, LEDITS++ is the only method to date to facilitate easy and versatile image editing by supporting multiple simultaneous instructions without causing undue interference. Finally, its lightweight architecture-agnostic nature ensures compatibility with both latent and pixel-based diffusion models, providing high accessibility.

In this work, we establish the methodical benefits of LEDITS++ and demonstrate that this intuitive, lightweight approach offers sophisticated semantic control for image editing. Specifically, we contribute by (i) devising a formal definition of LEDITS++ while (ii) deriving perfect inversion for a more efficient diffusion sampling method, (iii) qualitatively and empirically demonstrating its efficiency, versatility, and precision, (iv) providing an exhaustive empirical comparison to concurrent works with automatic and human user metrics, and thereby (v) introducing **Textual Editing Benchmark++** (TEdBench++), a more holistic and coherent testbed for evaluating textual image manipulation.

## 2. Background

Recently, large-scale, text-guided DMs have enabled versatile applications in image generation [3, 33, 37]. Especially latent diffusion models [31, 34] have gained attention for their computational efficiency. Below, we discuss related work for efficient, versatile image manipulation with DMs.

**Diffusion Sampling.** Generating outputs with DMs requires multiple iterative denoising steps that constitute the main bottleneck at inference. Commonly used sampling methods such as DDPM [15] or DDIM [42] require tens or hundreds of steps to produce high-quality samples. Consequently, numerous works have been dedicated to speeding up the sampling process without loss in quality. Distillation efforts progressively reduce the number of required steps through further training [25, 27]. Other works focus on improving the sampling itself, e.g. using high-order ODE-solvers [22, 23, 46]. Such solvers can be readily combined with pre-trained DMs at inference to lower the number of denoising steps. With LEDITS++, we derive perfect image inversion with the DPM-Solver++, allowing image editing in as few as 20 total steps.

**Semantic Control during Diffusion.** While text-to-image DMs generate new, astonishing images, fine-grained control over the generative process remains challenging. Minor changes to the text prompt lead to entirely differ-

ent outputs. Wu et al. [45] studied concept disentanglement using linear combinations of text embeddings to gain semantic control. Methods like Prompt-to-Prompt [14] and other works [8, 30] utilize the DM’s attention layers to attribute pixels to tokens from the text prompt. Dedicated operations on the attention maps enable more control over the generated images. Other works have focused on the noise estimates of DMs [5, 20] providing semantic control over the generation process. With LEDITS++, we now enable fine-grained semantic control for manipulating real images, going beyond purely generative applications.

**Real Image Editing.** Since DMs’ rise in popularity for text-to-image generation, they have also been explored for (real) image-to-image editing. As a first, simple approach, SDEdit added noise to the image for an intermediate step in the diffusion process [26]. While lightweight, the resulting image diverges substantially from the input as it is (partially) regenerated. Inpainting allows to keep the change small by having a user provide additional masks to restrict changes to certain image regions [2, 29]. Yet, user masks are costly or often simply unavailable. Other works have thus explored semantically grounded approaches using cross-attention instead to better control image manipulation [7, 28, 30]. In contrast, LEDITS++ leverages both attention- and noise-based masking to obtain fine-grained masks, enabling strong semantic control over real images.

Another important aspect of image manipulation methods is the required tuning and overall runtime. Instruct-Pix2Pix continues training a DM at scale to enable image editing capabilities [6]. Finetuning instead on each individual input to constrain the generation on the real image has shown helpful [18, 44] but not computationally efficient. Consequently, recent works have largely relied on inverting the deterministic DDIM sampling process [42] to save computational resources. DDIM inversion identifies an initial noise vector that results in the input image when denoised again. However, faithful reconstructions are only obtained in the limit of small steps, thus requiring large numbers of inversion steps. Moreover, small errors will still incur at each timestep, often accumulating into meaningful deviations from the input, requiring costly error correction through optimization [28, 30]. Recently, Huberman-Spiegelglas *et al.* proposed an inversion technique [17] for the DDPM sampler [15] to address the limitations of DDIM inversion. LEDITS++ provides the same guarantees of perfect inversion with even further reduced runtime alongside an edit-friendly latent space, enabling more versatility.

## 3. Image Editing with Text-to-Image Models

Before devising the methodology of LEDITS++, let us first motivate the desired features and use cases. Specifically, we aim for efficiency, versatility, and precision. The goal

is to provide a method that enables a fast exploratory workflow for image editing in which a user can iteratively interact with the model and explore various edits. Consequently, LEDITS++ produces outputs quickly with no tuning or optimization to not disrupt the creative process. Further, arbitrary editing instructions and combinations thereof are supported to facilitate a wide range of image manipulations (e.g., complex multi-editing). Lastly, we provide precise and sophisticated semantic control over the image editing. Each of the (potentially multiple) edit instructions can be steered individually, and changes are automatically restricted to relevant image regions. Importantly, with LEDITS++ we prioritize compositional robustness.

### 3.1. Guided Diffusion Models

Let us first define some general background for diffusion models. DMs iteratively denoise a Gaussian distributed variable to produce samples of a learned data distribution. Let’s consider a diffusion process that gradually turns an image  $x_0$  into Gaussian noise.

$$x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}n_t, \quad t = 1, \dots, T \quad (1)$$

where  $n_t$  are iid normal distributed vectors and  $\beta_t$  a variance schedule. The diffusion process is equivalently expressed as

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t \quad (2)$$

where  $\alpha_t = 1 - \beta_t$ ,  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$  and  $\epsilon_t \sim \mathcal{N}(0, \mathbf{I})$ . Importantly, all  $\epsilon_t$  are *not* statistically independent. Instead, consecutive pairs  $\epsilon_t, \epsilon_{t-1}$  are strongly dependent, which will be relevant later. To generate an (new) image  $\hat{x}_0$  the reverse diffusion process starts from random noise  $x_T \sim \mathcal{N}(0, \mathbf{I})$  which can be iteratively denoised as

$$x_{t-1} = \hat{\mu}_t(x_t) + \sigma_t z_t, \quad t = T, \dots, 1 \quad (3)$$

Here  $z_t$  are iid standard normal vectors, and common variance schedulers  $\sigma_t$  can be expressed in the general form

$$\sigma_t = \eta \beta_t \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}$$

where  $\eta \in [0, 1]$ . In this formulation,  $\eta = 0$  corresponds to the deterministic DDIM [42] and  $\eta = 1$  to the DDPM scheme [15]. Lastly, in these cases, we have  $\hat{\mu}_t(x_t) =$

$$\frac{\sqrt{\bar{\alpha}_{t-1}}x_t - \sqrt{1 - \bar{\alpha}_t}\hat{\epsilon}_\theta(x_t)}{\sqrt{\bar{\alpha}_t}} + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}\hat{\epsilon}_\theta(x_t)$$

Here  $\hat{\epsilon}_\theta(x_t)$  is an estimate of  $\epsilon_t$  produced by our neural network DM with learned parameters  $\theta$ , commonly implemented as a U-Net [35]. For text-to-image generation, the model is conditioned on a text prompt  $p$  to produce images faithful to that prompt. The DM is trained to produce the

noise estimate  $\hat{\epsilon}_\theta(x_t)$  needed for iteratively sampling  $\hat{x}_0$  (Eq. 3). For text-conditioned DMs,  $\hat{\epsilon}_\theta$  is calculated using specific guidance techniques.

Most DMs rely on classifier-free guidance [16], a conditioning method using a purely generative diffusion model, eliminating the need for an additional classifier. During training, the text conditioning  $c_p$  is randomly dropped with a fixed probability, resulting in a joint model for unconditional and conditional objectives. During inference, the score estimates for the  $\epsilon$ -prediction are adjusted so that:

$$\hat{\epsilon}_\theta(x_t, c_p) := \hat{\epsilon}_\theta(x_t) + s_g(\hat{\epsilon}_\theta(x_t, c_p) - \hat{\epsilon}_\theta(x_t)) \quad (4)$$

with guidance scale  $s_g$  and  $\hat{\epsilon}_\theta$  defining the noise estimate with parameters  $\theta$ . Intuitively, the unconditioned  $\epsilon$ -prediction is pushed in the direction of the conditioned one, with  $s_g$  determining the extent of the adjustment.

### 3.2. LEDITS++

With the fundamentals established, the methodology of LEDITS++ can now be broken down into three components: (1) efficient image inversion, (2) versatile textual editing, and (3) semantic grounding of image changes.

**Component 1: Perfect Inversion.** Utilizing text-to-image models for editing real images requires conditioning the generation on the input image. Recent works have largely relied on inverting the sampling process to identify  $x_T$  that will be denoised to the input image  $x_0$  [28, 30]. Inverting the DDPM scheduler is generally preferred over DDIM inversion since the former can be achieved in fewer timesteps and with no reconstruction error [17].

However, there exist more efficient schemes than DDPM for sampling DMs that greatly reduce the required number of steps and consequently DM evaluations. We here propose a more efficient inversion method by deriving the desired inversion properties for such a scheme. As demonstrated by Song *et al.*[43], DDPM can be viewed as a first-order stochastic differential equation (SDE) solver when formulating the reverse diffusion process as an SDE. This SDE can be solved more efficiently—in fewer steps—using a higher-order differential equation solver, here *dpm-solver++* [23]. The reverse diffusion process from Eq. 3 for the second-order sde-dpm-solver++ can be written as

$$x_{t-1} = \hat{\mu}_t(x_t, x_{t+1}) + \sigma_t z_t, \quad t = T, \dots, 1 \quad (5)$$

where now

$$\sigma_t = \sqrt{1 - \bar{\alpha}_{t-1}}\sqrt{1 - e^{-2h_{t-1}}}$$

and higher-order  $\hat{\mu}_t$  depends now on  $x$  from two timesteps,

$x_t$  and  $x_{t+1}$ . Such that  $\hat{\mu}_t(x_t, x_{t+1}) =$

$$\frac{\sqrt{1 - \bar{\alpha}_{t-1}}}{\sqrt{1 - \bar{\alpha}_t}} e^{-h_{t-1}x_t} + \sqrt{\bar{\alpha}_{t-1}}(1 - e^{-2h_{t-1}})\hat{\epsilon}_\theta(x_t) \\ + 0.5\sqrt{\bar{\alpha}_{t-1}}(1 - e^{-2h_{t-1}})\frac{-h_t}{h_{t-1}}(\hat{\epsilon}_\theta(x_{t+1}) - \hat{\epsilon}_\theta(x_t))$$

with

$$h_t = \frac{\ln(\sqrt{\bar{\alpha}_t})}{\ln(\sqrt{1 - \bar{\alpha}_t})} - \frac{\ln(\sqrt{\bar{\alpha}_{t+1}})}{\ln(\sqrt{1 - \bar{\alpha}_{t+1}})}$$

For the detailed derivation of the solver and proof of faster convergence, we refer the reader to the relevant literature [22, 23]. Based on the above, we now devise our inversion process. Given an input image  $x_0$  we construct an auxiliary reconstruction sequence of noise images  $x_1, \dots, x_T$  as

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\tilde{\epsilon}_t \quad (6)$$

where  $\tilde{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I})$ . Contrary to Eq. 2, the  $\tilde{\epsilon}_t$  are now statistically *independent*, which is a desirable property for image editing [17]. Lastly, the respective  $z_t$  for the inversion can be derived from Eq. 5 as

$$z_t = \frac{x_{t-1} - \hat{\mu}_t(x_t, x_{t+1})}{\sigma_t}, \quad t = T, \dots, 1 \quad (7)$$

with  $\hat{\mu}$  and  $\sigma_t$  as defined above. We base our implementation on the multistep variant of sde-dpm-solver++, which only requires one evaluation of the DM at each diffusion timestep by reusing the estimates from the previous step. The number of timesteps can be reduced further by stopping the inversion at an intermediate step  $t < T$  and starting the generation at that step. Empirically, we observed that  $t \in [0.9T, 0.8T]$  usually produces edits of the same fidelity as  $t = T$ , supporting observations in previous work [17, 26] that earlier timesteps are less relevant to the edit.

**Component 2: Textual Editing.** After creating our reconstruction sequence  $x_1, \dots, x_T$  and calculating the respective  $z_t$ , we now edit the image by manipulating the noise estimate  $\hat{\epsilon}_\theta$  based on a set of edit instructions  $\{e_i\}_{i \in I}$ . We devise a dedicated guidance term for each concept  $e_i$  based on conditioned and unconditioned estimates. Let us formally define LEDITS++’s guidance by starting with a single editing prompt  $e$ . We compute

$$\hat{\epsilon}_\theta(x_t, c_e) := \hat{\epsilon}_\theta(x_t) + \gamma(x_t, c_e) \quad (8)$$

with guidance term  $\gamma$ . Consequently, setting  $\gamma = 0$  will reconstruct the input image  $x_0$ . We construct  $\gamma$  to push the unconditioned score estimate  $\hat{\epsilon}_\theta(x_t)$ —i.e. the input image reconstruction—away from/towards the edit concept estimate  $\hat{\epsilon}_\theta(x_t, c_e)$ , depending on the guidance direction:

$$\gamma(x_t, c_e) = \phi(\psi; s_e, \lambda)\psi(x_t, c_e) \quad (9)$$

where  $\phi$  applies an edit guidance scale  $s_e$  element-wise, and  $\psi$  depends on the edit direction:  $\psi(x_t, c_e) =$

$$\begin{cases} \hat{\epsilon}_\theta(x_t, c_e) - \hat{\epsilon}_\theta(x_t) & \text{if pos. guidance} \\ -(\hat{\epsilon}_\theta(x_t, c_e) - \hat{\epsilon}_\theta(x_t)) & \text{if neg. guidance} \end{cases} \quad (10)$$

Thus, changing the guidance direction is reflected by the direction between  $\hat{\epsilon}_\theta(x_t, c_e)$  and  $\hat{\epsilon}_\theta(x_t)$ . The term  $\phi$  identifies those dimensions of the image and respective  $\hat{\epsilon}_\theta$  that are relevant to a prompt  $e$ . Consequently,  $\phi$  returns 0 for all irrelevant dimensions and a scaling factor  $s_e$  for the others. We describe the construction of  $\phi$  in detail below. Larger  $s_e$  will increase the effect of the edit, and  $\lambda \in (0, 1)$  reflects the percentage of the pixels selected as relevant by  $\phi$ . Notably, for a single concept  $e$  and uniform  $\phi = s_e$ , Eq. 8 generalizes to the classifier-free guidance term in Eq. 4.

For multiple  $e_i$ , we calculate  $\gamma_t^i$  as described above, with each defining their own hyperparameter values  $\lambda^i, s_e^i$ . The sum of all  $\gamma_t^i$  results in

$$\hat{\gamma}_t(x_t, c_{e_i}) = \sum_{i \in I} \gamma_t^i(x_t, c_{e_i}) \quad (11)$$

**Component 3: Semantic Grounding.** The masking term  $\phi$  (Eq. 9) is the intersection (pointwise product) of binary masks  $M^1$  and  $M^2$  combined with scaling factor  $s_e$ :

$$\phi(\psi; s_{e_i}, \lambda) = s_{e_i} M_i^1 M_i^2 \quad (12)$$

where  $M_i^1$  is a binary mask generated from the U-Net’s cross-attention layers and  $M_i^2$  is a binary mask derived from the noise estimate. Intuitively, each mask is an importance map, where  $M_i^1$  is more strongly grounded than  $M_i^2$ , but of significantly coarser granularity. Therefore, the intersection of the two yields a mask both focused on relevant image regions and of fine granularity. With LEDITS++, we empirically demonstrate that these maps can also capture regions of an image relevant to an editing concept that is not already present. Specifically for multiple edits, calculating a dedicated mask for each edit prompt ensures that the corresponding guidance terms remain largely isolated, limiting interference between them.

Formally, at each time step  $t$ , a U-Net forward pass with editing prompt  $e_i$  is performed to generate cross-attention maps for each token of the editing prompt. All cross-attention maps of the smallest resolution (e.g.,  $16 \times 16$  for SD) are averaged over all heads and layers, and the resulting maps are summed over all editing tokens, resulting in a single map  $A_t^{e_i} \in \mathbb{R}^{16 \times 16}$ . Importantly, we utilize the same U-Net evaluation  $\hat{\epsilon}_\theta(x_t, c_e)$  already performed in Eq. 10 to produce  $M^1$  with minimal overhead. Each map  $A_t^{e_i}$  is up-sampled to match the size of  $x_t$ . Cross-attention mask  $M^1$  is derived by calculating the  $\lambda$ -th percentile of



Figure 2. Comparison of image editing methods. (top) LEDITS++ is the only method to restrict edits to the tree leaves and position of the car. (bottom) Ours is the only approach faithfully executing all three edits and keeping changes minimal. (Best viewed in color)

Method	Reconstruction Error (RMSE) ↓	Execution Time (s) ↓	Variation/ Sampling	Semantic Grounding	Multi-Editing
SDEdit [26]	0.81 ±0.07	2.10 ±0.02	✓	✗	✗
Imagic [18]	0.58 ±0.12	349.98 ±0.45	✓	✗	✗
Vanilla DDIM Inversion	0.22 ±0.10	37.23 ±0.04	✗	✗	✗
Pix2Pix-Zero [30]	0.20 ±0.09	56.78 ±0.14	(✓)	✓	✗
DiffEdit [9]	0.13 ±0.03	27.65 ±0.03	✓	✓	✗
Edit-friendly DDPM [17]	<b>0.00</b>	10.36 ±0.05	✓	✗	✗
LEDITS++ (Ours)	<b>0.00</b>	<b>1.78</b> ±0.03	✓	✓	✓

Table 1. Comparing key properties for diffusion-based image editing techniques, with LEDITS++ offering clear methodological benefits. Due to LEDITS++’s efficient perfect inversion, it is the fastest and error-free method. At the same time, its methodology is the only enabling versatility in terms of *variation*, *semantic grounding*, and *multi-editing*. Subscript numbers indicate standard deviation.

up-sampled  $A_t^{e_i}$  and

$$M_t^1 = \begin{cases} 1 & \text{if } |A_t^{e_i}| \geq \eta_\lambda(|A_t^{e_i}|) \\ 0 & \text{else} \end{cases} \quad (13)$$

where  $\eta_\lambda(|\cdot|)$  is the  $\lambda$ -th percentile. By definition,  $M^1$  only selects image regions that correlate strongly with the editing prompt, and  $\lambda$  determines the size of this selected region.

The fine-grained mask  $M^2$  is calculated based on the guidance vector  $\psi$  of noise estimates derived in Eq. 10. The difference between unconditioned and conditioned  $\hat{\epsilon}_\theta$ , generally captures outlines and object edges of  $x_t$ . Consequently, the largest absolute values of  $\psi$  provide meaningful segmentation information of fine granularity for  $M^2$

$$M^2 = \begin{cases} 1 & \text{if } |\psi| \geq \eta_\lambda(|\psi|) \\ 0 & \text{else} \end{cases} \quad (14)$$

In general, threshold  $\lambda$  should correspond to the performed edit. Changes affecting the entire image, such as style transfer, should choose smaller  $\lambda$  ( $\rightarrow 0$ ), whereas edits targeting specific objects or regions should use  $\lambda$  proportional to the region’s prominence in the image.

## 4. Properties of LEDITS++

With the fundamentals of LEDITS++ established, we next showcase its unique properties and capabilities.

**Efficiency.** First off, LEDITS++ offers substantial performance improvements over other image editing methods. In Tab. 1, we provide a qualitative runtime comparison, with all methods being implemented for Stable Diffusion (SD) 1.5 [34]. As a parameter-free approach, LEDITS++ does not require any computationally expensive fine-tuning or optimization. Consequently, LEDITS++ is orders of magnitude faster than methods like Imagic [18] or Pix2Pix-Zero [30]. Further, we only need to invert the same number of diffusion steps used at inference, which results in significant runtime improvements over the standard DDIM inversion (21x). In addition to efficient inversion, we use a recent, fast scheduler that generally requires fewer total steps, further boosting performance. This way, LEDITS++ is six times faster than recent DDPM inversion [17] and on par with fast but poor-quality SDEdit [26].

**Versatility.** In addition to its efficiency, LEDITS++ remains versatile, enabling sheer limitless editing possibilities. In Fig. 1, we showcase a broad range of edit types. LEDITS++ facilitates fine-grained edits (adding/removing



Figure 3. Exemplary edit performed with LEDITS++ in only 25 diffusion steps with SD1.5. We apply a complex, compounded edit and ground each to a semantically reasonable image region.

glasses) and holistic changes such as style transfer (painting/sketch). Furthermore, object removal and replacement facilitate even more image editing tasks. Importantly, the overall image composition is preserved in all cases. To our knowledge, LEDITS++ is the only diffusion-based image editing method inherently supporting multiple edits in isolation, which allows for more complex image manipulation. Fig. 2 highlights LEDITS++ benefits over previous methods. Our method produces the highest edit fidelity and is the only approach capable of faithfully executing multiple, simultaneous instructions. Moreover, LEDITS++ also makes the least changes to unrelated objects and the overall background and composition of the image.

Lastly, the editing versatility benefits from the stochastic nature of the perfect but non-deterministic inversion. LEDITS++ provides meaningful image variations by resampling  $\tilde{\epsilon}_t$ . Additionally, the visual expression of each concept in the edited image scales monotonically with  $s_e$ , and the direction and magnitude of each concept can be varied freely. We present examples of both features in App. B.

**Precision.** Furthermore, LEDITS++’s methodology keeps edits concise and avoids unnecessary deviations from the input image (Fig. 2). First, the perfect inversion will reconstruct the exact input image if no edit is applied (cf. Sec. 3.2). Consequently, we already improve on faithfulness to the input image even before applying any edits. This benefit over other methods is highlighted by the reconstruction error in Tab. 1. Second, implicit masking will semantically ground each edit to relevant image regions. This is specifically important for editing multiple concepts at the same time. While other methods only utilize one prompt for all instructions, LEDITS++ isolates edits from each other (Eq. 11). Thus, we get dedicated masks for each concept as shown in Fig. 3. This design ensures that each instruction (e.g., red mask for ‘cherry blossom’) will be only applied where necessary. Subsequently, we provide further evidence for the efficacy of LEDITS++’s masking approach.

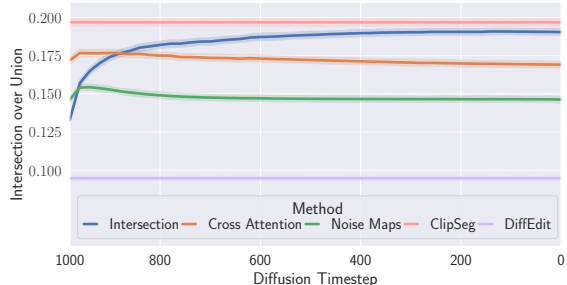


Figure 4. Semantic segmentation quality of LEDITS++. We show the intersection over union (higher is better) for COCO panoptic segmentation. The intersection masks outperform each by a clear margin, close to the CLIPSeg reference. (Best viewed in color)

## 5. Semantically Grounded Image Editing

Cross-attention maps of DMs have been used extensively to ground regions of interest during image generation semantically [7, 8, 14, 30]. Nonetheless, these have not been combined with noise-based masks so far and thus lack fine granularity. Hence, we empirically evaluate the quality of implicit masks, i.e., attention maps  $M^1$  and noise maps  $M^2$  (Eq. 13 and 14) in the LEDITS++ setup. We use a broad segmentation task for common objects as a proxy to measure the performance of implicit masks in identifying relevant image areas from edit instructions. Specifically, we utilize segmentation masks from the COCO panoptic segmentation challenge [19]. For each unique object in an image, we retrieve the masks  $M^1$ ,  $M^2$ , and their intersection per diffusion step. We use the (semantic) class label (e.g. ‘person’ or ‘TV’) as editing concept  $e$ . We consider masks at each of 50 total diffusion steps without actually editing the input image. Furthermore, we approximate mask threshold  $\lambda$  based on the relative size of an object’s bounding box.

**Concise masks with LEDITS++.** Fig. 4 shows implicit masking as a reliable means to identify relevant image regions. Importantly, the intersection of both cross-attention masks  $M^1$  and noise maps  $M^2$  clearly outperforms each separate mask. The overall performance is even similar to a dedicated CLIPSeg model [24], despite LEDITS++ masks being implicitly calculated at inference with only minimal overhead. At the same time, LEDITS++’s masking is superior to DiffEdit’s [3]. Consequently, our method’s intersection of cross-attention masks and noise maps provides strong semantic grounding while being efficient during image manipulation to ensure precise editing.

## 6. Image Editing Evaluation

Let us now compare LEDITS++ to current SOTA methods for image manipulation on two benchmarks.

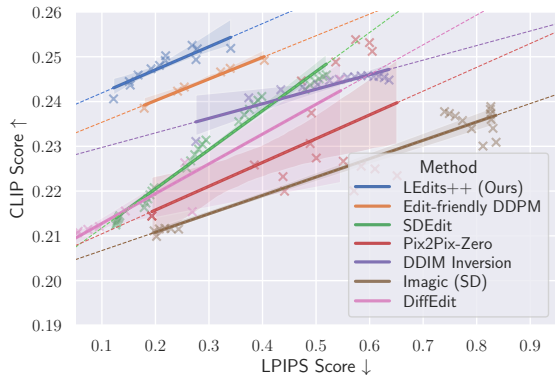


Figure 5. Comparison of instruction-alignment vs. image similarity trade-off for different editing methods. Results were reported for simultaneous manipulation of three facial attributes on CelebA. We plot CLIP scores (higher is better) of the target attributes against LPIPS similarity (lower is better). LEdits++ clearly outperforms all competing methods. (Best viewed in color)

### 6.1. Editing Multiple Concepts

First, we investigate the complex task of performing multiple edits simultaneously. We rely on a well-established setup for semantic image manipulation [5] to evaluate multi-conditioned attribute manipulation in facial images. In our experiment, we consider 100 images from the CelebA dataset [21]. For each image, we simultaneously edit three attributes out of a set of five, leading to ten total combinations of edit concepts. Further, we perform each edit across ten different seeds, resulting in 10,000 evaluated images for each method and hyperparameter setting, over 1M images in total. As measures for comparison, we employ CLIP and LPIPS scores. CLIP measures the text-to-image similarity of the edit instruction to the edited image, and LPIPS measures the image-to-image similarity of the real to the edited image. This way, we assess the trade-off between the versatility of edits (CLIP) and the precision of those manipulations (LPIPS). We implement all methods based on SD1.5 and provide more details in App. C.

**LEdits++ outperforms competing methods.** Fig. 5 shows the resulting CLIP vs. LPIPS plots for all methods. The top left corner represents the ideal editing method with maximum edit alignment without deviating from the initial image. Generally, one can observe a natural trade-off between versatility and precision for all methods, i.e., higher image-to-text alignment comes at the expense of lower similarity to the original image. LEdits++ is closest to the ideal region and thus clearly outperforms the other methods. In particular, the outputs remain close to the original image (low LPIPS scores), thanks to the precise implicit masking. At the same time, it faithfully performs the edits (high CLIP scores) due to the dedicated, isolated editing for each concept. The depicted scores reflect our qualitative inspections

	TEdBench		TEdBench++	
	SR $\uparrow$	LPIPS $\downarrow$	SR $\uparrow$	LPIPS $\downarrow$
Imagic w/ SD1.5	0.55	0.56	0.58	0.57
LEdits++ w/ SD1.5	<b>0.75</b>	<b>0.28</b>	<b>0.79</b>	<b>0.30</b>
Imagic w/ Imagen [18]	0.83	0.59	—	—
LEdits++ w/ SD-XL	<b>0.84</b>	<b>0.33</b>	<b>0.87</b>	<b>0.34</b>

Table 2. Success rate (SR) and LPIPS scores on the original TeDBench [18] and our revised version (TEdBench++). We compare Imagic to LEdits++ based on different DMs and find the latter to outperform on both metrics and benchmarks.

for Pix2Pix-Zero and Imagic on such complex manipulations (cf. Fig. 2). We observed that these methods often break—either failing to perform all three edits and/or drastically altering the input image. Only edit-friendly DDPM [17] and LEdits++ reliably achieve the maximum average CLIP score of over 0.25. This value seems to represent an upper bound according to our manual investigations, as each attribute is edited correctly for all input images, and higher scores are not observed. Despite being computationally very efficient, LEdits++ faithfully executes each edit instruction while keeping the changes to the input low, highlighting the method’s versatility and precision.

### 6.2. TEdBench(++)

Next, we investigate the versatility of LEdits++’s editing capabilities by running the Textual Editing Benchmark (TEdBench [18]), a collection of 100 input images paired with textual edit instructions. However, we observed a variety of inconsistencies in TEdBench and a lack of relevant editing tasks. Therefore, we propose TEdBench++ (Fig. 6a and App. D), a more challenging revised benchmark now containing 120 entries in total.<sup>3</sup> We addressed misspellings and rephrased ambiguous and inconclusive instructions. In addition to resolving these issues, we added instructions targeting challenging types of image manipulations previously not included in TEdBench: multi-conditioning, object/concept removal, style transfer, and complex replacements (Fig. 6a). We provide more details in App. D.

We compare LEdits++ on TEdBench(++) to one of the strongest editing methods, Imagic [18] with Imagen [37]. Since both are not publicly available, we can only compare to this specific combination of DM and editing method using Kawar *et al.*’s [18] curated outputs for TEdBench. Additionally, we, therefore, cannot combine LEdits++ with Imagen[37] and instead use a similarly advanced diffusion model, SD-XL [32]. However, to not only compete for the best fidelity outputs but focus the evaluation on methodological differences—not the pre-trained DM—we also compare both methods implemented with SD1.5. We provide further details in App. C.

<sup>3</sup>[https://huggingface.co/datasets/AIML-TUDA/TEdBench\\_plusplus](https://huggingface.co/datasets/AIML-TUDA/TEdBench_plusplus)



(a) Novel challenging examples of TEDBench++ and LEDITS++ applied, showcasing the versatility of supported edits.

(b) Qualitative comparison of LEDITS++ and Imagic on TEDBench, clearly highlighting the performance improvement.

Figure 6. Benchmark examples for LEDITS++ and Imagic on TEDBench(+). (Best viewed in color)

**LEDITS++ edits images reliably.** We first asked users to assess the overall success of edits, i.e., if an edit instruction was faithfully realized for a given input image. The results in Tab. 2 show that LEDITS++ outperforms Imagic on TEDBench despite a greatly reduced runtime (Tab. 1). The difference is even stronger when comparing both methods on the same pre-trained DM, i.e., SD1.5. The high success rate on TEDBench++ (87%) and the examples shown in Fig. 1 and 6a once again highlight LEDITS++’s versatility. Overall, our proposed method can reliably perform a diverse set of editing instructions for real images.

**High-quality edits with LEDITS++.** While investigating both methods’ performance we observed a substantial difference in edit quality. The examples in Fig. 6b particularly highlight the discrepancy in compositional robustness and object coherence. Hence, we also assessed both methods’ editing quality on TEDBench(++). We focus on samples where both methods performed a successful edit, i.e., were labeled as successful by users. We show the perceptual similarity (LPIPS) to the input image in Tab. 2. One can observe that the LPIPS scores for LEDITS++ are much lower than for Imagic, empirically supporting the qualitative examples in Fig. 6b. When manually inspecting the generated images, we often found Imagic to generate a completely *new* image based on the edit instruction, entirely disregarding the input image (cf. App. Fig. 11a).

## 7. Discussion

Let us now discuss open research questions and limitations.

**Model Dependency.** While LEDITS++ achieves impressive results on a large variety of image manipulation tasks, there are external factors to consider. Since the method is architecture-agnostic, it can be easily used with any DM. At

the same time, the general editing quality strongly depends on the overall capabilities of the underlying pre-trained DM. Naturally, more capable models will also enable better edits. But, at times, specific editing instructions may fail because the used DM does not have a decent representation of the targeted concept to begin with. One example is the model failing to edit a giraffe to be *sitting* since the underlying DM generally fails to generate this pose (cf. App. F). This effect can also clearly be seen in Tab. 2, with the editing success rate of a method varying strongly between DMs. Although the same image editing method is employed (LEDITS++), the more capable SD-XL variant outperforms the weaker SD1.5 model. Nonetheless, this means that the architecture-agnostic LEDITS++ will benefit from increasingly powerful DMs.

**Coherence Trade-offs.** Next to the benefits of LEDITS++’s semantic grounding, there are also downsides to this approach. Overall, implicit masking limits changes to relevant portions of the image and achieves strong coherence with the original image composition. Yet, the object and its identity within the masked area may change based on various factors. Generic prompts, like “a standing cat” (cf. App. F), do not contain detailed information about this specific object (“cat”). Thus, an edit with this prompt does not guarantee to preserve object identity, particularly for strong hyperparameters. We observed that fine-tuning approaches like Imagic make the opposite trade-off, better preserving the object identity while changing the background and image composition substantially (cf. App. F). A potential remedy for a loss in object coherence with LEDITS++, is more descriptive edit prompting, e.g. using textual inversion [12].

Lastly, the automatically-inferred implicit masks allow for easy use of LEDITS++ without users tediously providing masks. Nonetheless, user intentions are diverse and



cannot always be automatically inferred. Sometimes, individual user masks provide better control over the editing process. Such user masks can be easily integrated into LEDITS++ (cf. App. F), wherefore we encourage future research in this promising direction.

**Societal Impact.** LEDITS++ is an easy-to-use image editing technique that lowers the barrier for users and puts them in control for fruitful human-machine collaboration. Yet, the underlying text-to-image models offer both promise and peril, as highlighted by prior research [4, 11]. The (societal) biases within these models will also impact image editing applications [11]. Moreover, image manipulation can also be used adversarially to generate inappropriate [39] or fake content. Hence, we advocate for a cautious deployment of generative models together with image editing methods.

## 8. Conclusion

We introduced LEDITS++, an efficient yet versatile and precise method for textual image manipulation with diffusion models. It facilitates the editing of complex concepts in real images. Our approach requires no finetuning nor optimization, can be computed extremely efficiently, and is architecture agnostic. At the same time, it perfectly reconstructs an input image and uses implicit masking to limit changes to relevant image regions, thus editing precisely. Our large experimental evaluation confirms the efficiency, versatility, and precision of LEDITS++ and its components, as well as its benefits over several related methods.

**Acknowledgements** We gratefully acknowledge support from the BMBF (Grant No 01IS22091). This work benefited from the ICT-48 Network of AI Research Excellence Center “TAILOR” (EU Horizon 2020, GA No 952215), the Hessian research priority program LOEWE within the project WhiteBox, the HMWK cluster projects “Adaptive Mind” and “Third Wave of AI”, and from the NHR4CES.

## References

- [1] Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene: Extracting multiple concepts from a single image. *SIGGRAPH Asia*, 2023. 16
- [2] Omri Avrahami, Ohad Fried, and Dani Lischinski. Blended latent diffusion. *arXiv:2206.02779*, 2022. 2
- [3] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, Tero Karras, and Ming-Yu Liu. eDiff-I: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv:2211.01324*, 2022. 2, 6
- [4] Federico Bianchi, Pratyusha Kalluri, Esin Durmus, Faisal Ladhak, Myra Cheng, Debora Nozza, Tatsunori Hashimoto, Dan Jurafsky, James Zou, and Aylin Caliskan. Easily Accessible Text-to-Image Generation Amplifies Demographic Stereotypes at Large Scale. In *Proceedings of ACM Conference on Fairness, Accountability, and Transparency (FAcT)*, 2023. 9, 11
- [5] Manuel Brack, Felix Friedrich, Dominik Hintersdorf, Lukas Struppek, Patrick Schramowski, and Kristian Kersting. Segat: Instructing text-to-image models using semantic guidance. In *Proceedings of the Advances in Neural Information Processing Systems: Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2023. 2, 7
- [6] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2
- [7] 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. *arXiv:2304.08465*, 2023. 2, 6
- [8] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models. *ACM Trans. Graph.*, 42, 2023. 2, 6
- [9] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. 5
- [10] Zoe De Simone, Angie Boggust, Arvind Satyanarayan, and Ashia Wilson. What is a Fair Diffusion Model? Designing Generative Text-To-Image Models to Incorporate Various Worldviews. *arXiv preprint arXiv:2309.09944*, 2023. 11
- [11] Felix Friedrich, Manuel Brack, Lukas Struppek, Dominik Hintersdorf, Patrick Schramowski, Sasha Luccioni, and Kristian Kersting. Fair Diffusion: Instructing Text-to-Image Generation Models on Fairness. *arXiv preprint arXiv:2302.10893*, 2023. 9, 11
- [12] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*, 2022. 8, 16
- [13] Rohit Gandikota, Joanna Materzynska, Jaden Piotto-Kaufman, and David Bau. Erasing concepts from diffusion

- models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2023. 11
- [14] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2023. 2, 6
- [15] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Proceedings of the Advances in Neural Information Processing Systems: Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2020. 2, 3
- [16] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv:2207.12598*, 2022. 3
- [17] Inbar Huberman-Spiegelglas, Vladimir Kulikov, and Tomer Michaeli. An edit friendly ddp noise space: Inversion and manipulations. *arXiv preprint arXiv:2304.06140*, 2023. 1, 2, 3, 4, 5, 7, 12, 13, 14
- [18] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 1, 2, 5, 7, 12, 13, 14, 15
- [19] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2014. 6
- [20] Nan Liu, Shuang Li, Yilun Du, Antonio Torralba, and Joshua B. Tenenbaum. Compositional visual generation with composable diffusion models. In *Proceedings of European Conference on Computer Vision (ECCV)*, 2022. 2
- [21] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, December 2015. 7
- [22] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ODE solver for diffusion probabilistic model sampling in around 10 steps. In *NeurIPS*, 2022. 2, 4
- [23] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*, 2022. 2, 3, 4
- [24] Timo Lüddecke and Alexander Ecker. Image segmentation using text and image prompts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 6
- [25] Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthesizing high-resolution images with few-step inference. *arXiv:2310.04378*, 2023. 2
- [26] 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 *Proceedings of the International Conference on Learning Representations (ICLR)*, 2022. 1, 2, 4, 5, 13, 14
- [27] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik P. Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. 2
- [28] 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 (CVPR)*, 2023. 1, 2, 3, 13
- [29] 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 *Proceedings of the International Conference on Machine Learning (ICML)*. PMLR, 2022. 2, 11
- [30] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *SIGGRAPH*, 2023. 1, 2, 3, 5, 6, 12, 13, 14
- [31] Pablo Pernias, Dominic Rampas, Mats L. Richter, Christopher J. Pal, and Marc Aubreville. Wuerstchen: An efficient architecture for large-scale text-to-image diffusion models. *arXiv:2306.00637*, 2023. 2
- [32] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv:2307.01952*, 2023. 7
- [33] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with CLIP latents. *arXiv preprint arXiv:2204.06125*, 2022. 2, 11
- [34] 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 (CVPR)*, 2022. 2, 5
- [35] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Proceedings of Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pages 234–241, 2015. 3
- [36] Harrison Rosenberg, Shimaa Ahmed, Guruprasad V. Ramesh, Ramya Korlakai Vinayak, and Kassem Fawaz. Unbiased Face Synthesis With Diffusion Models: Are We There Yet? *arXiv preprint arXiv:2309.07277*, 2023. 11
- [37] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv:2205.11487*, 2022. 2, 7, 11
- [38] Dvir Samuel, Rami Ben-Ari, Simon Raviv, Nir Darshan, and Gal Chechik. Generating images of rare concepts using pre-trained diffusion models. *arXiv:2304.14530*, 2023. 16
- [39] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*

*Recognition (CVPR)*, 2023. 9, 11

- [40] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. Laion-5b: An open large-scale dataset for training next generation image-text models. In *Proceedings of NeurIPS Datasets and Benchmarks*, 2022. 11
- [41] Norbert Schwarz. How the questions shape the answers. *American Psychologist*, 1999. 15
- [42] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 2, 3
- [43] Yang Song, Jascha Sohl-Dickstein, Diederik P. Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021. 3
- [44] Dani Valevski, Matan Kalman, Yossi Matias, and Yaniv Leviathan. Unitune: Text-driven image editing by fine tuning an image generation model on a single image. *arXiv:2210.09477*, 2022. 1, 2
- [45] Qiucheng Wu, Yujian Liu, Handong Zhao, Ajinkya Kale, Trung Bui, Tong Yu, Zhe Lin, Yang Zhang, and Shiyu Chang. Uncovering the disentanglement capability in text-to-image diffusion models. *arXiv:1212.08698*, 2022. 2
- [46] Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. *arXiv:2302.04867*, 2023. 2

## Appendix

### A. Broader (Societal) Impact

With LEDITS++, we aim to provide an easy-to-use image editing framework. It lowers the barrier of entry for experienced artists and novices alike, allowing them to unlock the full potential of generative AI in the pursuit of creative expression. Moreover, it puts the user in control for fruitful human-machine collaboration. Crucially, current text-to-image models [29, 33, 37] hold the potential to wield a profound influence on society. When applied in creative and design domains, their dual use offers both promise and peril, as highlighted by prior research [4, 11]. The models are trained on large amounts of data from the web [40], granting them the inherent capacity to generate content that may contravene societal norms, including the creation of inappropriate material like pornography [39], or content that violates law such as child abuse. More alarmingly, the inadvertent generation of inappropriate content is precipitated by spurious correlations within these models. Harmless prompts can lead to the creation of decidedly objectionable content [4, 11]. A prime example of this phenomenon lies in the correlation between specific phrases and the perpetuation of stereotypes, such as the connection between mentions of ethnicity and economic status. For example, an increase of the concept ‘*black person*’ may inadvertently amplify the appearance of the concept ‘*poverty*.’

Conversely, methods like LEDITS++ also possess the potential to mitigate bias and inappropriateness, a prospect highlighted by prior research [10, 11], e.g. through dataset augmentation [36]. Furthermore, established strategies offer means to mitigate the generation of inappropriate content [13, 39] that could be deployed in tandem with LEDITS++. In summary, we advocate for a cautious approach to the utilization of these models, recognizing both the risks and promises they bring to the realm of AI-powered image editing.

### B. Further Examples on LEDITS++ Properties

As discussed in Sec. 4, LEDITS++ versatility benefits from re-sampling to provide variations of edits. The example in Fig. 7 demonstrates the additional control non-deterministic variations provide to the user, which can select the preferred interpretation of the edit instruction.

The precision and versatility of LEDITS++ further benefit from the fact that the magnitude of an editing concept in the output scales monotonically with the edit scale  $s_e$ . In Fig. 8, we can observe the effect of increasing  $s_e$ . Both for positive and negative guidance, the change in scale correlates with the strength of the smile or frown. Consequently, any changes to the input image can be steered intuitively using the edit guidance scale.



Figure 7. LEDITS++ easily produces variations of an edit (different (styles of) sunglasses) by resampling the inversion process. (Best viewed in color)

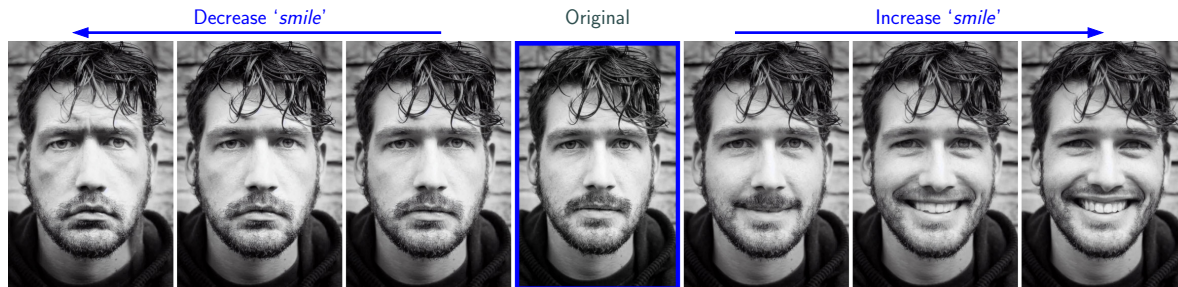


Figure 8. Monotonicity of editing scale with LEDITS++. The original image (middle) is edited with varying scales of the same edit ('smile'). The scale for 'smile' is semantically reflected in the images. (Best viewed in color)

## C. Experimental Details

Subsequently, we provide further details on the experiments presented in the main body of the paper. We first provide information on the reconstruction and runtime experiments (Sec. 4), followed by the masking evaluation (Sec. 5) and multi-conditioning experiments (Sec. 6). Details on the user study are independently described in App. E. All experiments were performed using the respective diffusers<sup>4</sup> implementation (version 0.20.2) with Stable Diffusion 1.5.

### C.1. Properties

First, we go into detail on the reconstruction and runtime experiments presented in Tab. 1.

#### C.1.1 Reconstruction Error

Since the Stable Diffusion VAE already induces errors when reconstructing images, we considered the RMSE over the 64x64 latent image instead. We randomly sampled 100 images from the 2017 COCO validation dataset, which we attempted to reconstruct using the default configuration of each method as described in the respective paper or implementation. For methods that could potentially benefit from a descriptive target prompt of the input image, we considered an empty prompt, COCO caption as prompt, and unconditioned generation (no CFG) and reported the best score. Below, we outline the configuration for each

method.

#### LEDITS++ (Ours) & Edit-friendly DDPM [17]:

Perfect reconstruction for any hyperparameter combination. The error induced by machine precision is inconsequential.

#### Imagic [18]:

1000 embedding learning steps w/ learning rate  $2e-3$  and 1500 model tuning steps w/ learning rate  $5e-7$ . Target prompt is the original image caption. We used 50 generation steps with  $\alpha$  and guidance scale 0.0

#### Pix2Pix-Zero [30]:

Inversion with 100 steps, no CFG,  $\lambda = 20$  for auto correction and KL divergence, and 5 regularization and auto correction steps, respectively. 100 inference steps with cross-attention guidance 0.0 and no CFG. Source and target embeddings are null-vectors.

<sup>4</sup><https://huggingface.co/docs/diffusers>

**DDIM Inversion:** 1000 inversion steps and 50 generation steps. Both without classifier-free guidance

**SDEdit [26]:** 40 diffusion steps (strength 0.8 at 50 default steps) with no CFG.

Notably, the small difference between DDIM and Pix2Pix-Zero only holds for the pure reconstruction of an image. Previous research has shown [28] that for DDIM inversion, the accumulated error increases drastically when using classifier-free guidance. Since CFG is necessary for editing, the "reconstruction" portion during editing becomes worse for DDIM and remains stable for Pix2Pix-Zero.

### C.1.2 Runtime

For runtime measurements, we consider the wallclock runtime on a dedicated NVIDIA A100-SXM4-40GB GPU. As a proxy task, we considered applying an 'oilpainting' style to a photograph. We only measured the inversion and generation loops, discarding any I/O or other processing. For each method, we considered 100 runs with hyperparameters based on the respective paper/official implementation, as outlined below.

**LEDITS++ (Ours):** 20 inversion and 20 generation steps with threshold  $\lambda = 0.1$  and skip for  $t = 0.75T$ .

**Edit-friendly DDPM [17]:** 100 inversion steps and 64 generation steps (i.e. 36 skip steps)

**Imagic [18]:** 1000 text embedding optimization steps, 1500 model finetuning steps, 50 inference steps

**Pix2Pix-Zero [30]:** 100 steps at inversion and inference. 5 regularization steps and auto correction steps for each inversion step.

**DDIM Inversion:** 1000 inversion steps (w/o CFG) and 50 generation steps

**SDEdit [26]:** 40 diffusion steps (strength 0.8 at 50 default steps).

An interesting observation is the fact that LEDITS++ runs faster than SDEdit although both perform 40 diffusion steps overall. However, SDEdit requires 80 total U-Net evaluation (unconditioned and conditioned for each step) step, whereas LEDITS++ only requires 60 (unconditioned at each inversion step and unconditioned + conditioned at each inference step). Performing 2 evaluations of the U-Net is significantly slower than 1 evaluation even if performed as a batch.

## C.2. Implicit Mask Quality

We have already provided detailed information on the experiment in Sec. 5. For further reference, we note that after removing duplicate/ambiguous objects the dataset contains 4983 images and 29307 segmentation objects.

## C.3. Multi-conditioning

The multi-conditioning experiment presented in Sec. 6.1, used the following attributes:

- glasses
- smile
- hat
- wavy hair
- earrings

We used the first 100 images in CelebA that were labeled to not contain any of the five target attributes. Seeds were chosen at random but kept fixed across all experiments. For the LPIPS and CLIP scores, we relied on the default implementation from torchmetrics<sup>5</sup>. Consequently, we used the AlexNet variant with mean reduction and the original ViT-L/14 CLIP checkpoint from OpenAI<sup>6</sup>. For the CLIP scores, we calculated a dedicated score for each of the 3 applied edits and considered the mean for each image.

The hyperparameter variations of each method were run as a grid search over the hyperparameter ranges listed below. Other parameters were kept at their default values. For each method, we ran a grid search over a wider range of parameters to identify reasonable boundaries and subsequently discarded edge values leading to drops in performance.

**LEDITS++ (Ours):** Skip between 0.2 and 0.3, Guidance scale between 10.0 and 15.0, Threshold between 0.7 and 0.9

<sup>5</sup><https://lightning.ai/docs/torchmetrics/stable/>

<sup>6</sup><https://huggingface.co/openai/clip-vit-large-patch14>

<b>Edit-friendly DDPM [17]:</b>	Skip steps between 20 and 40, Guidance scale between 10.0 and 15.0
<b>DDIM Inversion:</b>	Guidance Scale between 1.0 and 15.0
<b>Imagic [18]:</b>	Guidance Scale between 2 and 6, and $\alpha$ between 0.1 and 1.3
<b>Pix2Pix-Zero [30]:</b>	Guidance Scale between 1.0 and 10.0 and cross guidance scale between 0 and 0.15
<b>SDEdit [26]:</b>	Guidance scale between 5.0 and 10.0 and strength between 0.2 and 0.8

## D. TEdBench++

We propose TEdBench++<sup>7</sup>, a revised version of TEdBench [18] which sets a new standard for benchmarking real text-based image editing. It is publicly available, including original images, edit instructions, and edited images with LEDITS++ for benchmarking new methods. Figs. 6 and 11 as well as Tab. 2 demonstrated our generated images with LEDITS++ to improve upon the previous SOTA method, Imagic, setting a new reference for benchmarking. Next to providing better-edited images, we also addressed several inconsistencies in the target texts and missing tasks.

We show several inconsistencies of TEdBench in Fig. 9. First, we corrected ambiguous text prompts such as a *standing* animal that is already standing Fig. 9 (top). This applied to multiple images (horse, cat, bear, etc.). Instead, we propose “an {animal} standing on hind legs” to specify the target text and thus ask for a clear but more challenging edit. Second, we correct for misspellings such as “entrance” which should be “entrance” instead. While this may appear negligible, DMs’ tokenizers may provide completely different tokens in these cases: e.g., one token for the correct word but three tokens for the misspelled word. In Fig. 9, we show the impact of the corrections on the edit success on LEDITS++. Although we use the exact same parameters (seed, etc.), the edit of the original (left) to the middle fails, whereas it is successful for the corrected prompts (right). This way, we provide a higher-quality benchmark.

Further, we added novel tasks to the benchmark, making it more challenging and accounting for a broader range of tasks. In Fig. 6a, we show examples of the new tasks we added: i) multi-editing (adding multiple concepts at the

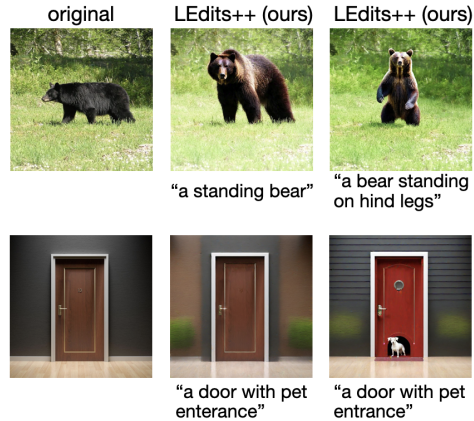


Figure 9. Exemplary inconsistencies in TEdBench and their corrections in TEdBench++. Left is the original image, and in the middle/right are images edited with LEDITS++ for the edit instructions above. All parameters (seed, etc.) are the same. As can be seen, the edit success heavily depends on the clear and correct writing of the words. In the middle column, it does not work (ambiguous (top) and misspelled (bottom)), whereas the edit is successful in the right image (clear and correctly spelled).

same time), ii) object removal (removing an object while staying consistent with the background and overall image composition), iii) style transfer (changing the whole image, i.e. all pixels without changing the overall image composition or object, only their style appearance), and iv) complex replacements (adding and removing multiple concepts at the same time).

With these extensions, we improve on the previous benchmark and propose a higher-quality version. This way, we hope to benefit the research community and set up a new standard for benchmarking text-based real-image editing techniques.

## E. User Study

Next to evaluating with automated metrics such as CLIP and LPIPS scores, we also conducted a study with human evaluators. We focus the user study on TEdBench(++). First, we describe the experimental details for generating the images for the study. Then, we describe the setup of the user study.

**Experimental details** We followed the approach of Kawar *et al.* [18] and generated images for several seeds and hyperparameters and hand-selected the best fitting image (exemplary grid search shown in Fig. 13). Notably, we evaluated only three seeds, whereas Kawar *et al.* evaluated eight seeds. Furthermore, we limited the grid search to a decent but small range for each hyperparameter.

For LEDITS++ (with SD1.5 and SD-XL), we set the number of diffusion steps fixed to 50 steps and

<sup>7</sup><https://figshare.com/s/7adc2b0fe1e0388dd99f>

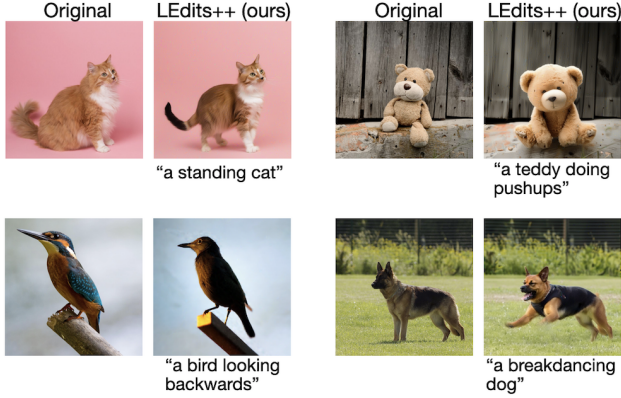


Figure 10. Failure cases of LEdITS++ on TEdBench

grid-searched skip [0.0, 0.1, 0.2, 0.4], masking threshold [0.6, 0.75, 0.9], and guidance scale [10, 15]. As a result, we evaluated 72 images (= 3 seeds  $\times$  4 skips  $\times$  3 thresholds  $\times$  2 scales) per benchmark sample. All other hyperparameters correspond to the default values of the diffusers implementation. Consequently, the generated images with LEdITS++ could be even further improved when evaluating for more hyperparameters, e.g. more seeds (also see open question discussion on seed in App. F).

For Imagic with SD1.5, we relied on 3 seeds and 50 diffusion steps, too. We grid-searched the guidance scale [5.0, 7.5, 10.0] and alpha value [0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1, 1.2, 1.4, 1.6, 1.8, 2.0]. As a result, we evaluated 108 images (= 3 seeds  $\times$  3 scales  $\times$  12 alphas) per benchmark sample. The remaining setup and parameters correspond to the default values of the original Imagic implementation with Imagen [18]. For Imagic with Imagen, we had to rely on their curated outputs of TEdBench since the model is not publicly available.

**User study setup** For the actual user studies, we chose the following setups (cf. Fig. 12) on the platform [thehive.ai](#). Users had to pass a qualifying test, and during the actual labeling, 15% of the tasks a user saw were honeypot tasks (sanity check). Only if the qualifier test was passed error-free and the honeypot accuracy was permanently above 95%, we accepted the given answers to ensure high-quality evaluations. Users could zoom in/out and change several image parameters, such as brightness and contrast to further enable a high-quality assessment.

The first study (see Fig. 12a), which we also describe in the main text in Tab. 2, evaluates the success rate of an editing technique on TEdBench(++). To this end, we asked users to assess the overall success of edits, i.e., if an edit instruction was faithfully realized for a given input image. Fig. 12a shows the setup, in which a user had to choose between two options, whether the edit instruction has been

realized successfully or not. The setup consisted of a general question-and-answer setting for all examples alongside a specific edit text for each image pair. The original image (always left) and the edited image (always right) were shown in the center. Outputs from LEdITS++ and Imagic were interleaved at random. The result is given in Tab. 2 in which LEdITS++ clearly outperforms Imagic for both underlying DMs and both versions of the benchmark benchmarks.

We also conducted a second user study. In this study, we asked the user to assess the image similarity of two methods to a reference image, see Fig. 12b. With this human preference study, we investigate the image-to-image similarity after editing, i.e., if the edited images still look similar to the original one. Participants were shown the input image (middle) and were asked to choose the better editing result from one of the two methods (left and right), using the common practice of Two-Alternative Forced Choice (2AFC). The methods were randomly switched between left and right to avoid confounding factors. To this end, we compared LEdITS++ (with SD-XL) and Imagic (with Imagen) on TEdBench. For this comparison, we considered only images where both methods were labeled successful in the prior study. In that study, users preferred LEdITS++ over Imagic with 60% preference. As outlined previously, this again emphasizes the precision of our method, which preserves the overall image composition and results in high-quality edits. Yet, the preference seems smaller than the results with the LPIPS scores. We found this to be an artifact of imprecise user study design. A preference setup generally suffers from bias such as subjects replacing general questions with more specific ones [41] (e.g. “which is more similar?” might be replaced by “in which did the main object stay the same, regardless of the background?” or “in which are the background and overall image composition better preserved regardless of the main object?”). Hence, it is difficult to draw exact conclusions from this study, but a clear trend is still visible. Moreover, as shown in the main text in Tab. 2, we computed LPIPS scores, which further clarify the results of the user study. Additionally, we broadly discussed the similarity trade-off between object identity/coherence and overall image composition in the limitation sections of the main body (Sec. 7) and appendix (App. F).

## F. Limitations and Further Discussion

In the following, we extend the discussion of the main body with further examples and questions.

**Model Dependency.** In general, we observed the editing success to be dependent on the underlying DM. In Fig. 15, we show that the generation of a *sitting giraffe* depends on the underlying DM. For both editing techniques, the weaker



(a) Image editing or generating a new image?



(b) Coherence Trade-off: compositional robustness vs. object identity

Figure 11. Comparing failure cases of LEDITS++ and Imagic on TEDBench.

SD1.5 variant fails but the more advanced variant succeeds. Upon further investigation, we realized that SD1.5 is incapable at all of generating images of *sitting* giraffes. In Fig. 16, we show exemplary images for the text prompt “a sitting giraffe” (we generated 100 and all showed the same result) and can see that none is actually sitting. In contrast, SD-XL is able to output images of *sitting* giraffes (cf. Fig. 17) and consequently enables LEDITS++ to perform the desired edit. This emphasizes the importance of choosing an apt underlying DM for real image editing and motivates future research to develop more powerful DMs from which editing techniques will benefit, too.

**Failure Cases and Open Questions.** In the following, we want to touch upon open questions and failure cases. In Fig. 10, we show several failure cases of LEDITS++. In the first case, the cat is indeed edited from a sitting to a standing cat. Yet, the identity of the cat has changed, i.e., the shape of the tail and the fur color have changed. We show further examples in Fig. 11. However, defining what makes an edit *acceptable* remains challenging and may differ between users, applications, and context. In general, however, there are two reasons for the discussed limitations. First, LEDITS++ limits its edits to the identified relevant regions. Consequently, the background and image composition will be preserved, but the edit within this region depends on various factors, including hyperparameter strength, underlying DM, and random seed. Second, a lack of descriptiveness of the editing prompt with respect to the specific identity of

an object can lead to changes thereof. Especially if generic terms such as ‘cat’ are used to edit the image. To guarantee the preservation of the object’s identity, methods like Textual Inversion [12], or Break-a-Scene [1] could be employed.

The other three examples in Fig. 10 show failure cases of LEDITS++ beyond changes of object identity, i.e., cases in which the edit instruction is not or not sufficiently realized. There are several reasons for such failures, including a general lack of concept understanding in the underlying DM (discussed above), incorrect masking of the relevant region, wrong hyperparameter choices, or challenging prompts. For example, what exactly is a “breakdancing dog” supposed to look like? We even considered removing this entry from the benchmark but found it very challenging at the same time and, therefore, kept it. Moreover, we found the edit success rate and quality to depend on the used seed. This is in line with current work on the impact of the used seed on the diffusion process [38]. Samuel *et al.* propose a new method to identify fitting seeds, which could be applied to LEDITS++ as well to find satisfying editing seeds.

**Masking and User Interaction** The automatically-inferred implicit masks allow for easy use of LEDITS++ without users tediously providing masks. Nonetheless, user intentions are diverse and cannot always be automatically inferred. Sometimes, individual user masks provide better control over the editing process. Such user masks can be easily integrated into LEDITS++. In



## Question

Was the following edit successful? In other words, has the image changed from left to right according to the text?

'A basket of oranges'

## Answers

Option 1) Yes

Option 2) No



(a) Setup for user study: “was the editing successful?”

## Question

Which image is more similar to the middle image? Left or Right? In other words, which image is closer to the middle?

## Answers

Option 1) left

Option 2) right



(b) Setup for user preference study: “which edited image is closer to the original image?”

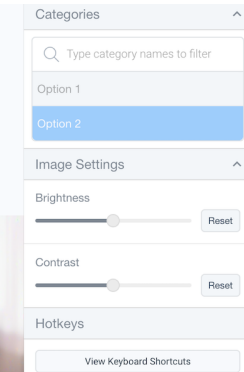
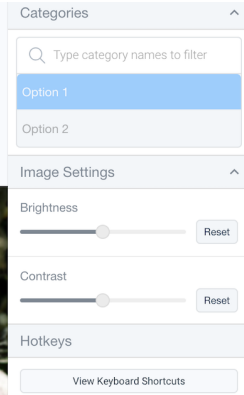


Figure 12. User study setups for both user studies conducted. The first user study evaluates the edit success of an image editing method. The second user study evaluates the user preference between two image editing methods regarding image-to-image similarity.

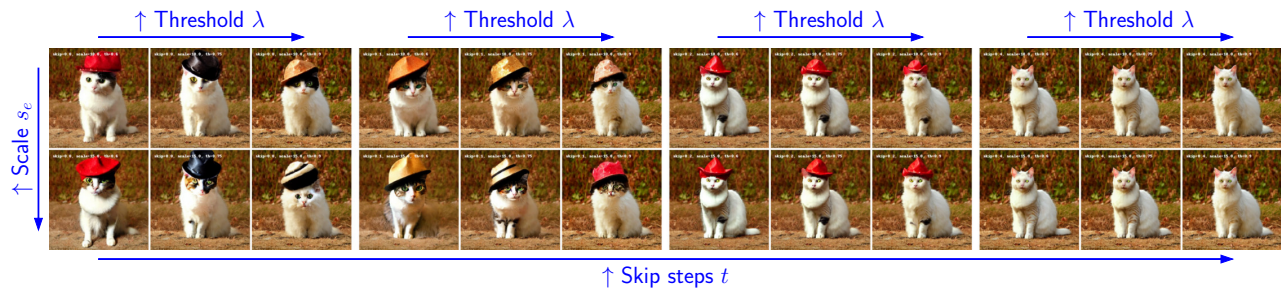


Figure 13. Grid search results for LEDITs++ for TEDBench(++). We show the grid search for the image “cat.jpeg” from the benchmark and the target text “a cat wearing a hat”. From left to right the parameters are increased. We searched three parameters, the skip steps, the guidance scale, and the masking threshold. As can be seen, the stronger the parameters, the more changes are made to the image. On the other hand, too weak parameters do not change the image, i.e. do not realize the target text. This highlights the trade-off between edit success and preservation of the image composition and object identity. All images are generated for the same seed (for TEDBench(++)) we evaluated 3 seeds per benchmark entry, whereas Imagic used 8 seeds).

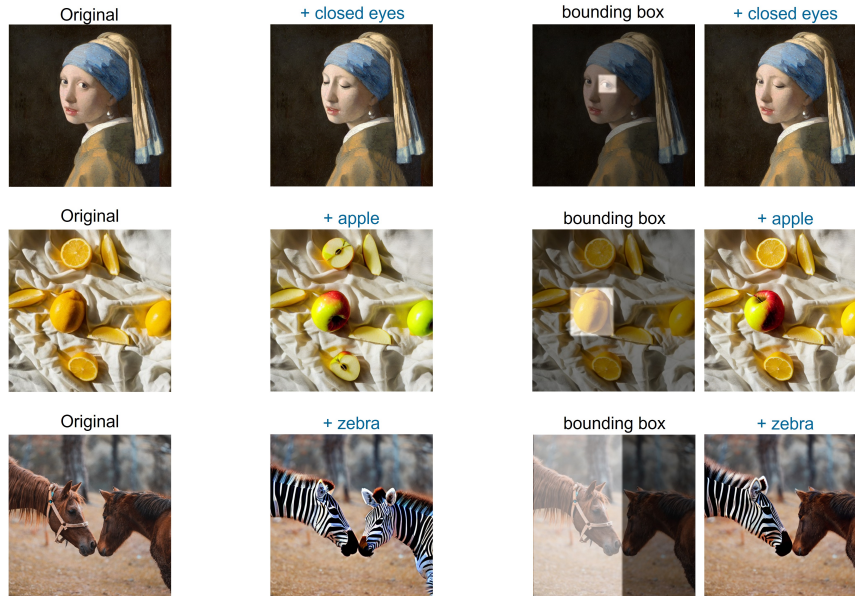


Figure 14. Customized and complex image editing with LEDITS++ for individual user masking. LEDITS++ can be easily extended with user masks to facilitate user preferences. The first column shows the original images and the second and fourth show the edited images with LEDITS++ and the target text above. The third column shows user-provided masks, marking the relevant region for the edit instruction. In the second column, LEDITS++ uses implicit masking as implemented in our default approach, and in the fourth column LEDITS++ uses the explicit user-provided masks from the third column.

Fig. 14, we show customized image editing with individual user masking. LEDITS++ can be easily extended with user masks to facilitate user preferences. Sometimes LEDITS++’s implicit masks do not meet user preferences or it is difficult to textually describe the relevant image region. Next to using dedicated models to obtain image masks, users can simply provide their own masks. In our setup, we did not evaluate this scenario as it drastically increases resources in terms of computation or human labor. Yet, it can be easily integrated into LEDITS++. Fig. 14 shows the original image can be edited well with LEDITS++. Moreover, dedicated user masks help focus on a specific image region. This is specifically helpful for logical and compositional instructions (current models struggle with “left”/“right” etc.), for one specific object if multiple are present (one specific “orange” from several ones), and for both combined. This way, LEDITS++ stays lightweight in its default with implicit masking, but can still and easily handle user masks and thereby implement individual user experience.

## G. Further Results

Subsequently, we present further results, qualitative examples, and visual ablations.

### G.1. Qualitative examples

We show further results in Fig. 21. We remove “cat” and add a diverse set of animals instead. Interestingly, this works for a variety of animals, that share no or only little similarity, such as “flamingo” or “parrot”. Furthermore, the newly occurring background is inpainted semantically sensible, too. Additionally, we show more qualitative examples in Fig. 22.

### G.2. Ablations

In Fig. 13, we show an ablation of LEDITS++ for TED-Bench(++). This grid search illustrates the impact of different hyperparameters on the trade-off between edit success and the preservation of the overall image composition/object identity.

### G.3. Semantic Grounding Ablations

We performed extensive ablations on semantic grounding by re-running LEDITS++ on Sec. 5’s benchmark without any grounding. The results in Fig. 18 show that LEDITS++ will still achieve strong instruction alignment (high CLIP score) without grounding, but semantic masking is key to keeping the generated image similar to the input (low LPIPS score). Moreover, grounding allows for a clearer trade-off between instruction alignment and image similarity in the first place. We believe these ablations foster a deeper understanding of the importance of semantic



Figure 15. Impact of underlying DM. The original image (leftmost) should be edited with the target text “a photo of a sitting giraffe”. The edit success depends on the underlying DM: with SD1.5 it fails whereas it works fine with more advanced DMs (SD-XL and Imagen). This holds for both methods LEDITS++ and Imagic.

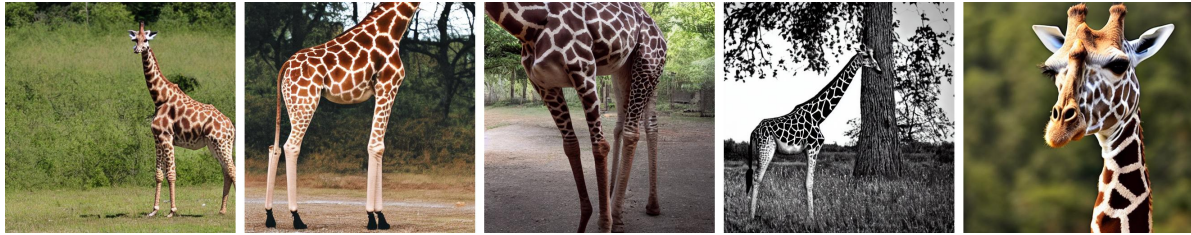


Figure 16. Generated Images with SD1.5 for “a photo of a sitting giraffe”. The model consistently fails to generate a giraffe in that specific pose.



Figure 17. Generated Images with SDXL for “a photo of a sitting giraffe”.

grounding in the LEDITS++ pipeline.

#### G.4. Explanatory Visualization of Monotonicity

The monotonicity of the LEDITS++ guidance scale is an important contribution of the method. Importantly, the inferred masks are mostly isolated from changes to the guidance scale. In the example shown in Fig. 19, we would expect the masks for ‘smiling’ to always target the area around the mouth and eyes. Within these identified regions, the magnitude of applied changes correlates directly with the changing scale, as evident from the provided heatmap. Here, we can also observe that different areas are prioritized depending on the magnitude of the change. The initial focus is clearly on the mouth, with strong changes in the eyes only appearing for larger scales.

#### G.5. Masking and Artifacts

LEDITS++ can faithfully edit reflections/shadows of objects, even with masking (cf. Fig. 1 & 20). This ability

strongly depends on the underlying diffusion model. In the example in Fig. 6b, the underlying diffusion model simply failed to correctly correlate the couple and their shadow/reflection. However, in Fig. 20 where SDXL is applied, the reflection is edited as well.

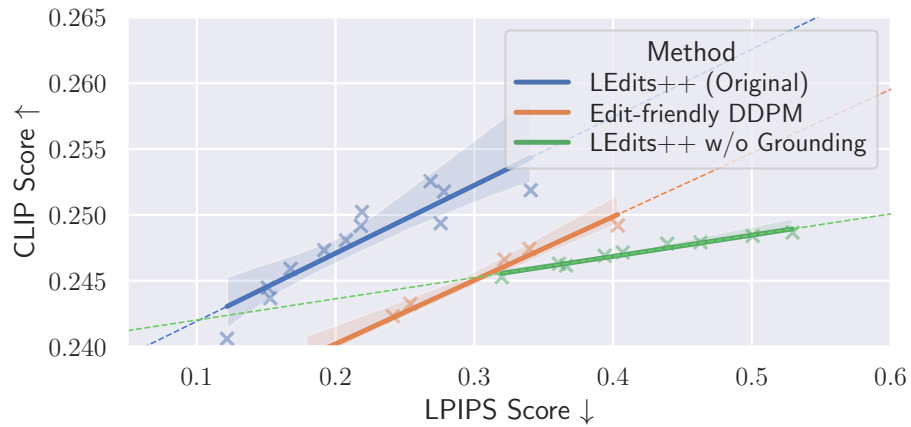


Figure 18. Semantic Grounding Ablations. Semantic grounding is an essential component of LEdits++ that helps preserving overall image composition and realizing concise edit instructions. (Best viewed in color)

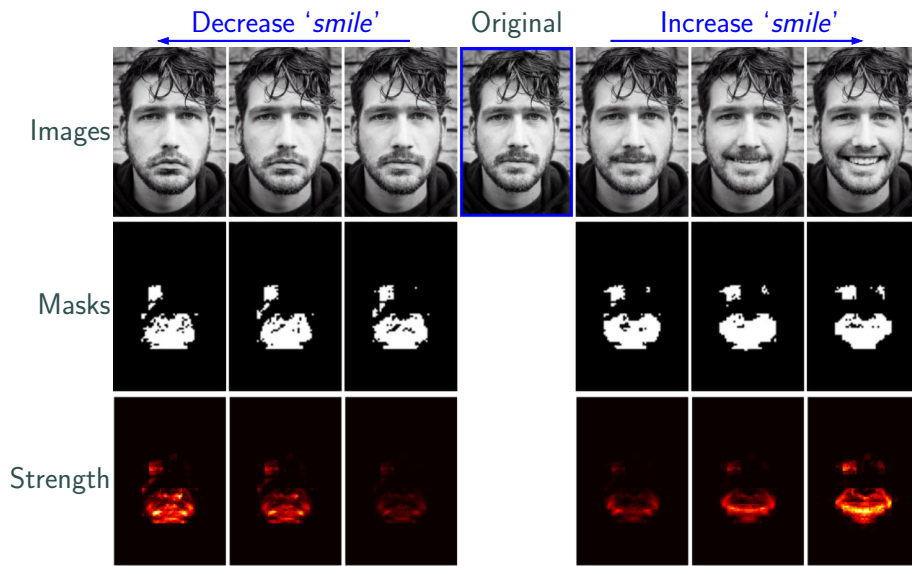


Figure 19. Masking and guidance scale. Within the identified edit regions, the magnitude of applied changes correlates directly with the changing scale. (Best viewed in color)



Figure 20. LEdits++ can easily handle complex edits such as reflections. (Best viewed in color)

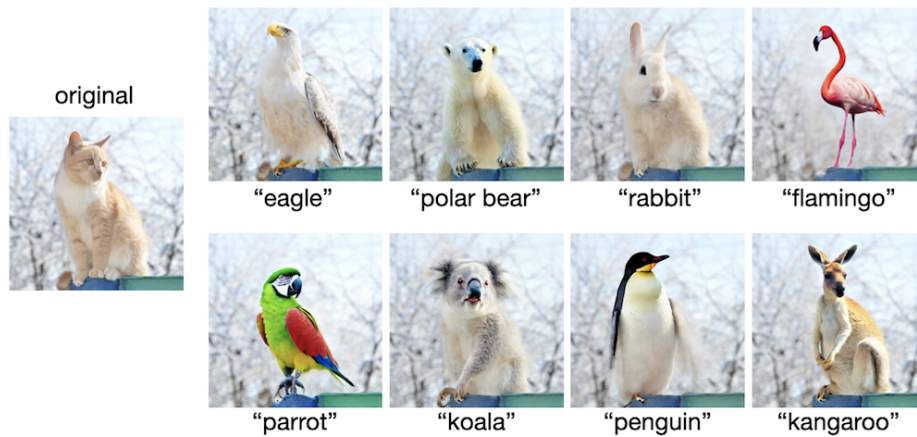


Figure 21. Object replacement with LEDITS++. The leftmost image is the original image and the other images are edited with LEDITS++ and the target text below. We apply diverse replacements of the main object with the overall image composition being preserved. Interestingly, the background is filled and interpolated very well, e.g. for “flamingo” or “parrot”. (Best viewed in color)



Figure 22. Further qualitative examples of image editing with LEDITS++, highlighting its versatility and precision (Best viewed in color)