

Imagic: Text-Based Real Image Editing with Diffusion Models

Bahjat Kawar*^{1,2}

Huiwen Chang¹

¹Google Research

Shiran Zada*¹

Tali Dekel^{1,3}

²Technion

Oran Lang¹

Inbar Mosseri¹

³Weizmann Institute of Science

Omer Tov¹

Michal Irani^{1,3}

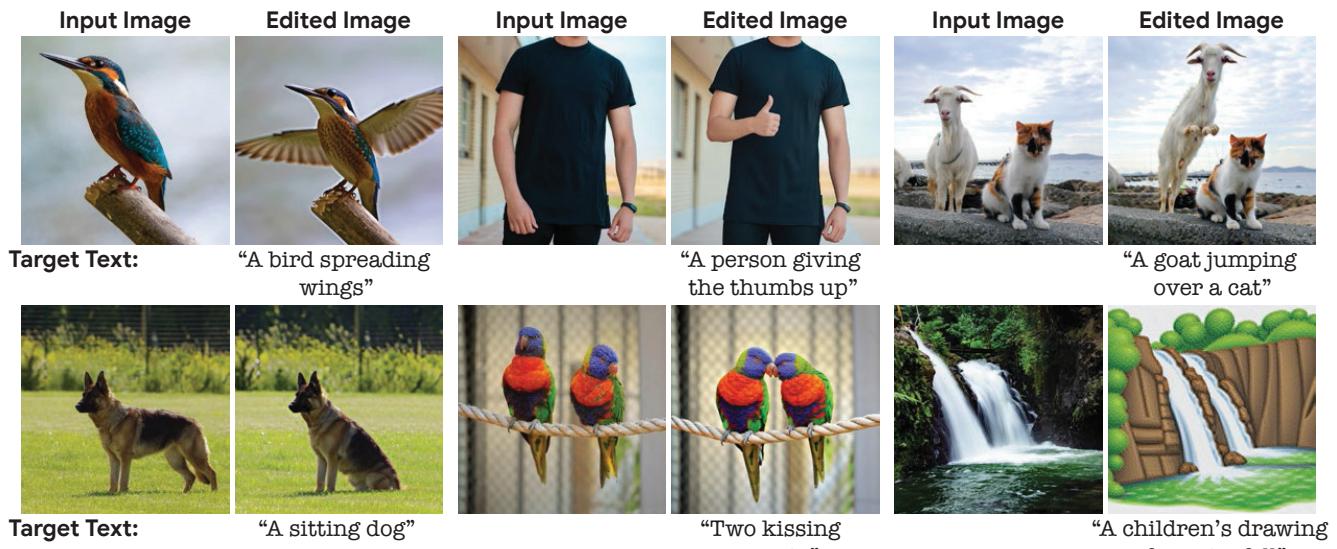


Figure 1. *Imagic – Editing a single real image.* Our method can perform various text-based semantic edits on a single real input image, including highly complex non-rigid changes such as posture changes and editing multiple objects. Here, we show pairs of 1024×1024 input (real) images, and edited outputs with their respective target texts.

Abstract

Text-conditioned image editing has recently attracted considerable interest. However, most methods are currently limited to one of the following: specific editing types (e.g., object overlay, style transfer), synthetically generated images, or requiring multiple input images of a common object. In this paper we demonstrate, for the very first time, the ability to apply complex (e.g., non-rigid) text-based semantic edits to a single real image. For example, we can change the posture and composition of one or multiple objects inside an image, while preserving its original characteristics. Our method can make a standing dog sit down, cause a bird to spread its wings, etc. – each within its single high-resolution user-provided natural image. Contrary to previous work, our proposed method requires only a single input image and a target text (the desired edit). It operates on real images, and

does not require any additional inputs (such as image masks or additional views of the object). Our method, called *Imagic*, leverages a pre-trained text-to-image diffusion model for this task. It produces a text embedding that aligns with both the input image and the target text, while fine-tuning the diffusion model to capture the image-specific appearance. We demonstrate the quality and versatility of *Imagic* on numerous inputs from various domains, showcasing a plethora of high quality complex semantic image edits, all within a single unified framework. To better assess performance, we introduce *TEdBench*, a highly challenging image editing benchmark. We conduct a user study, whose findings show that human raters prefer *Imagic* to previous leading editing methods on *TEdBench*.

1. Introduction

Applying non-trivial semantic edits to real photos has long been an interesting task in image processing [41]. It has attracted considerable interest in recent years, enabled by the considerable advancements of deep learning-based systems. Image editing becomes especially impres-

* Equal contribution.

The first author performed this work as an intern at Google Research.
 Project page: <https://imagic-editing.github.io/>.

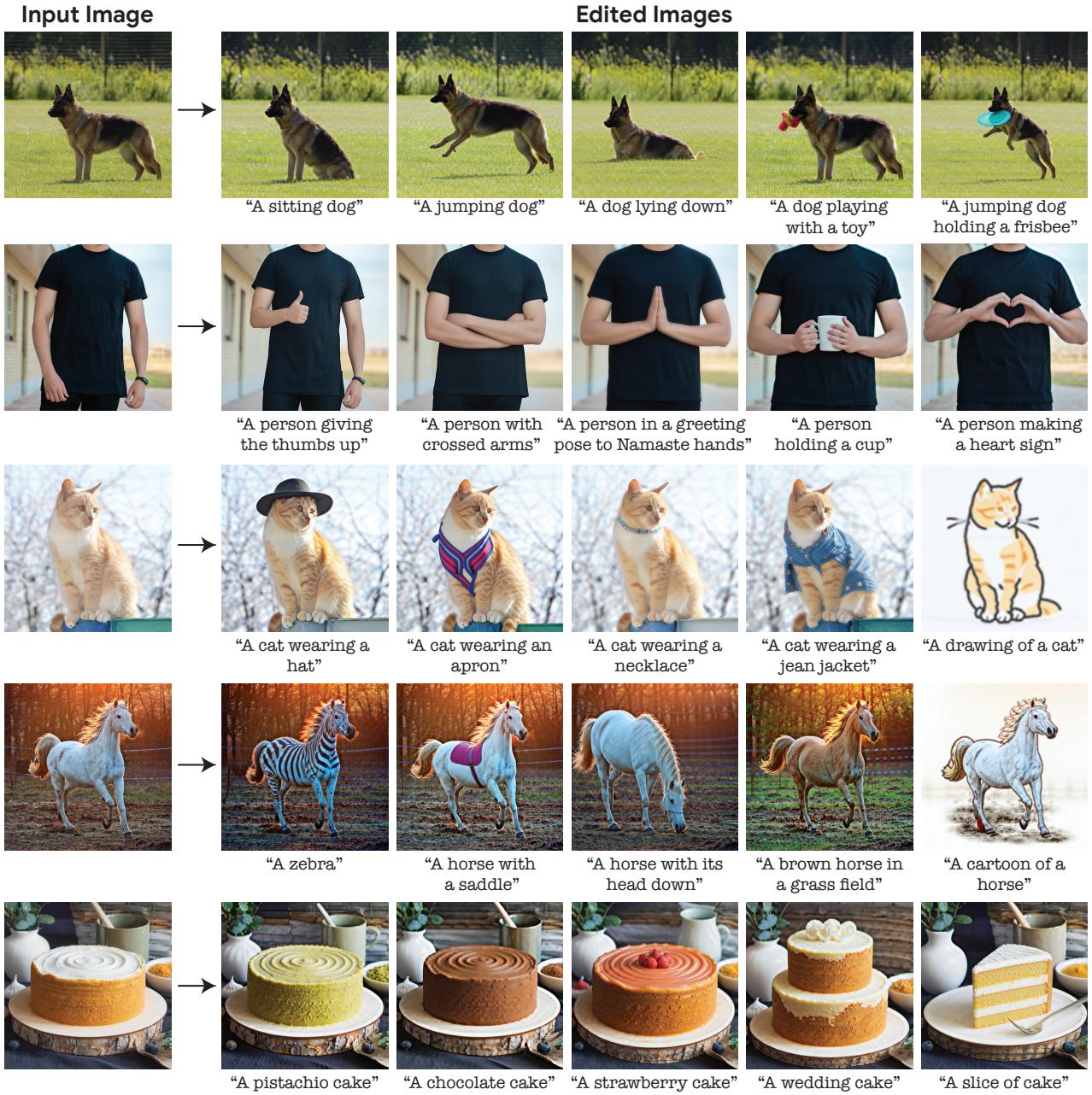


Figure 2. **Different target texts applied to the same image.** *Imagic edits the same image differently depending on the input text.*

sive when the desired edit is described by a simple natural language text prompt, since this aligns well with human communication. Many methods were developed for text-based image editing, showing promising results and continually improving [8, 10, 33]. However, the current leading methods suffer from, to varying degrees, several drawbacks: (i) they are limited to a specific set of edits such as painting over the image, adding an object, or transferring style [6, 33]; (ii) they can operate only on images from a specific domain or synthetically generated images [20, 43]; or (iii) they require auxiliary inputs in addition to the in-

put image, such as image masks indicating the desired edit location, multiple images of the same subject, or a text describing the original image [6, 17, 39, 47, 51].

In this paper, we propose a semantic image editing method that mitigates all the above problems. Given only an input image to be edited and a single text prompt describing the target edit, our method can perform sophisticated non-rigid edits on real high-resolution images. The resulting image outputs align well with the target text, while preserving the overall background, structure, and composition of the original image. For example, we can make two parrots kiss

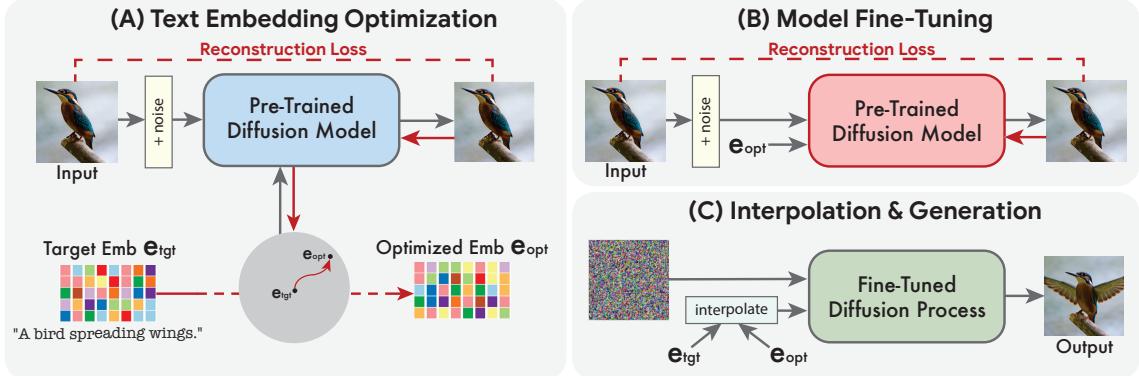


Figure 3. **Schematic description of *Imagic*.** Given a real image and a target text prompt: (A) We encode the target text and get the initial text embedding e_{tgt} , then optimize it to reconstruct the input image, obtaining e_{opt} ; (B) We then fine-tune the generative model to improve fidelity to the input image while fixing e_{opt} ; (C) Finally, we interpolate e_{opt} with e_{tgt} to generate the final editing result.

or make a person give the thumbs up, as demonstrated in Figure 1. Our method, which we call *Imagic*, provides the first demonstration of text-based semantic editing that applies such sophisticated manipulations to a single real high-resolution image, including editing multiple objects. In addition, *Imagic* can also perform a wide variety of edits, including style changes, color changes, and object additions.

To achieve this feat, we take advantage of the recent success of text-to-image diffusion models [47, 50, 53]. Diffusion models are powerful state-of-the-art generative models, capable of high quality image synthesis [16, 22]. When conditioned on natural language text prompts, they are able to generate images that align well with the requested text. We adapt them in our work to edit real images instead of synthesizing new ones. We do so in a simple 3-step process, as depicted in Figure 3: We first optimize a text embedding so that it results in images similar to the input image. Then, we fine-tune the pre-trained generative diffusion model (conditioned on the optimized embedding) to better reconstruct the input image. Finally, we linearly interpolate between the target text embedding and the optimized one, resulting in a representation that combines both the input image and the target text. This representation is then passed to the generative diffusion process with the fine-tuned model, which outputs our final edited image.

We conduct several experiments and apply our method on numerous images from various domains. Our method outputs high quality images that both resemble the input image to a high degree, and align well with the target text. These results showcase the generality, versatility, and quality of *Imagic*. We additionally conduct an ablation study, highlighting the effect of each element of our method. When compared to recent approaches suggested in the literature, *Imagic* exhibits significantly better editing quality and faithfulness to the original image, especially when tasked with sophisticated non-rigid edits. This is further supported by a human perceptual evaluation study, where raters strongly prefer *Imagic* over other methods on a novel

benchmark called *TEdBench* – Textual Editing Benchmark. We summarize our main contributions as follows:

1. We present *Imagic*, the first text-based semantic image editing technique that allows for complex non-rigid edits on a single real input image, while preserving its overall structure and composition.
2. We demonstrate a semantically meaningful linear interpolation between two text embedding sequences, uncovering strong compositional capabilities of text-to-image diffusion models.
3. We introduce *TEdBench* – a novel and challenging complex image editing benchmark, which enables comparisons of different text-based image editing methods.

2. Related Work

Following recent advancements in image synthesis quality [26–29], many works utilized the latent space of pre-trained generative adversarial networks (GANs) to perform a variety of image manipulations [3, 19, 36, 43, 56, 57]. Multiple techniques for applying such manipulations on real images were suggested, including optimization-based methods [1, 2, 25], encoder-based methods [4, 48, 64], and methods adjusting the model per input [5, 9, 15, 49]. In addition to GAN-based methods, some techniques utilize other deep learning-based systems for image editing [8, 12].

More recently, diffusion models were utilized for similar image manipulation tasks, showcasing remarkable results. SDEdit [38] adds intermediate noise to an image (possibly augmented by user-provided brush strokes), then denoises it using a diffusion process conditioned on the desired edit, which is limited to global edits. DDIB [62] encodes an input image using DDIM inversion with a source class (or text), and decodes it back conditioned on the target class (or text) to obtain an edited version. DiffusionCLIP [33] utilizes language-vision model gradients, DDIM inversion [59], and model fine-tuning to edit images using a domain-specific diffusion model. It was also suggested to edit images by

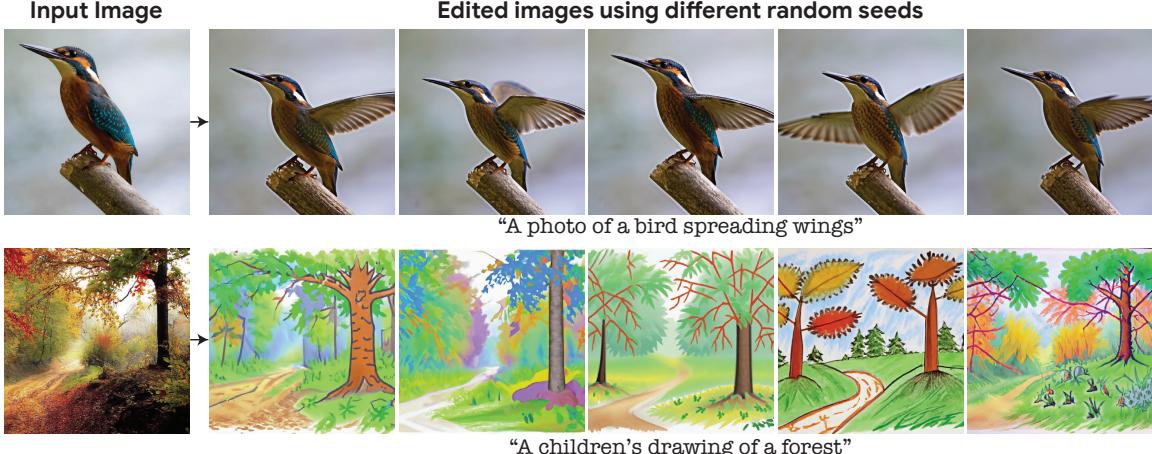


Figure 4. **Multiple edit options.** *Imagic utilizes a probabilistic model, enabling it to generate multiple options with different random seeds.*

synthesizing data in user-provided masks, while keeping the rest of the image intact [6, 14, 39]. Liu et al. [37] guide a diffusion process with a text and an image, synthesising images similar to the given one, and aligned with the given text. Hertz et al. [20] alter a text-to-image diffusion process by manipulating cross-attention layers, providing more fine-grained control over generated images, and can edit real images in cases where DDIM inversion provides meaningful attention maps. Textual Inversion [17] and DreamBooth [51] synthesize novel views of a given subject given 3–5 images of the subject and a target text (rather than edit a single image), with DreamBooth requiring additional generated images for fine-tuning the models. In this work, we provide the first text-based semantic image editing tool that operates on a single real image, maintains high fidelity to it, and applies non-rigid edits given a single free-form natural language text prompt.

3. Imagic: Diffusion-Based Real Image Editing

3.1. Preliminaries

Diffusion models [22, 58, 60, 66] are a family of generative models that has recently gained traction, as they advanced the state-of-the-art in image generation [16, 31, 61, 65], and have been deployed in various downstream applications such as image restoration [30, 52], adversarial purification [11, 40], image compression [63], image classification [69], and others [13, 18, 32, 44, 55, 67].

The core premise of these models is to initialize with a randomly sampled noise image $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$, then iteratively refine it in a controlled fashion, until it is synthesized into a photorealistic image \mathbf{x}_0 . Each intermediate sample \mathbf{x}_t (for $t \in \{0, \dots, T\}$) satisfies

$$\mathbf{x}_t = \sqrt{\alpha_t} \mathbf{x}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_t, \quad (1)$$

with $0 = \alpha_T < \alpha_{T-1} < \dots < \alpha_1 < \alpha_0 = 1$ being hyperparameters of the diffusion schedule, and $\boldsymbol{\epsilon}_t \sim \mathcal{N}(0, \mathbf{I})$. Each refinement step consists of an application of a neural

network $f_\theta(\mathbf{x}_t, t)$ on the current sample \mathbf{x}_t , followed by a random Gaussian noise perturbation, obtaining \mathbf{x}_{t-1} . The network is trained for a simple denoising objective, aiming for $f_\theta(\mathbf{x}_t, t) \approx \boldsymbol{\epsilon}_t$ [22, 58]. This leads to a learned image distribution with high fidelity to the target distribution, enabling stellar generative performance.

This method can be generalized for learning conditional distributions – by conditioning the denoising network on an auxiliary input \mathbf{y} , the network $f_\theta(\mathbf{x}_t, t, \mathbf{y})$ and its resulting diffusion process can faithfully sample from a data distribution conditioned on \mathbf{y} . The conditioning input \mathbf{y} can be a low-resolution version of the desired image [54] or a class label [23]. Furthermore, \mathbf{y} can also be on a text sequence describing the desired image [7, 47, 50, 53]. By incorporating knowledge from large language models (LLMs) [46] or hybrid vision-language models [45], these *text-to-image diffusion models* have unlocked a new capability – users can generate realistic high-resolution images using only a text prompt describing the desired scene. In all these methods, a low-resolution image is first synthesized using a generative diffusion process, and then it is transformed into a high-resolution one using additional auxiliary models.

3.2. Our Method

Given an input image \mathbf{x} and a target text which describes the desired edit, our goal is to edit the image in a way that satisfies the given text, while preserving a maximal amount of detail from \mathbf{x} (e.g., small details in the background and the identity of the object within the image). To achieve this feat, we utilize the text embedding layer of the diffusion model to perform semantic manipulations. Similar to GAN-based approaches [43, 49, 64], we begin by finding meaningful representation which, when fed through the generative process, yields images similar to the input image. We then fine-tune the generative model to better reconstruct the input image and finally manipulate the latent representation to obtain the edit result.

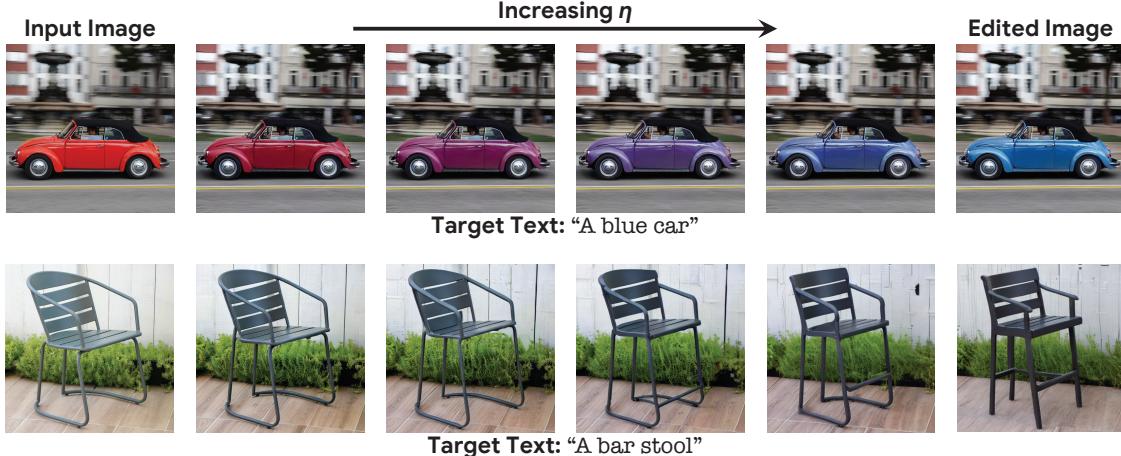


Figure 5. Smooth interpolation. We can smoothly interpolate between the optimized text embedding and the target text embedding, resulting in a gradual editing of the input image toward the required text as η increases (See animated GIFs in supplementary material).

More formally, as depicted in Figure 3, our method consists of 3 stages: (i) we optimize the text embedding to find one that best matches the given image in the vicinity of the target text embedding; (ii) we fine-tune the diffusion models to better match the given image; and (iii) we linearly interpolate between the optimized embedding and the target text embedding, in order to find a point that achieves both fidelity to the input image and target text alignment. We now turn to describe each step in more detail.

Text embedding optimization The target text is first passed through a text encoder [46], which outputs its corresponding text embedding $\mathbf{e}_{tgt} \in \mathbb{R}^{T \times d}$, where T is the number of tokens in the given target text, and d is the token embedding dimension. We then freeze the parameters of the generative diffusion model f_θ , and optimize the target text embedding \mathbf{e}_{tgt} using the denoising diffusion objective [22]:

$$\mathcal{L}(\mathbf{x}, \mathbf{e}, \theta) = \mathbb{E}_{t, \epsilon} \left[\|\epsilon - f_\theta(\mathbf{x}_t, t, \mathbf{e})\|_2^2 \right], \quad (2)$$

where $t \sim Uniform[1, T]$, \mathbf{x}_t is a noisy version of \mathbf{x} (the input image) obtained using $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ and Equation 1, and θ are the pre-trained diffusion model weights. This results in a text embedding that matches our input image as closely as possible. We run this process for relatively few steps, in order to remain close to the initial target text embedding, obtaining \mathbf{e}_{opt} . This proximity enables meaningful linear interpolation in the embedding space, which does not exhibit linear behavior for distant embeddings.

Model fine-tuning Note that the obtained optimized embedding \mathbf{e}_{opt} does not necessarily lead to the input image \mathbf{x} exactly when passed through the generative diffusion process, as our optimization runs for a small number of steps (see top left image in Figure 7). Therefore, in the second stage of our method, we close this gap by optimizing the model parameters θ using the same loss function presented in Equation 2, while freezing the optimized embedding.

This process shifts the model to fit the input image \mathbf{x} at the point \mathbf{e}_{opt} . In parallel, we fine-tune any auxiliary diffusion models present in the underlying generative method, such as super-resolution models. We fine-tune them with the same reconstruction loss, but conditioned on \mathbf{e}_{tgt} , as they will operate on an edited image. The optimization of these auxiliary models ensures the preservation of high-frequency details from \mathbf{x} that are not present in the base resolution. Empirically, we found that at inference time, inputting \mathbf{e}_{tgt} to the auxiliary models performs better than using \mathbf{e}_{opt} .

Text embedding interpolation Since the generative diffusion model was trained to fully recreate the input image \mathbf{x} at the optimized embedding \mathbf{e}_{opt} , we use it to apply the desired edit by advancing in the direction of the target text embedding \mathbf{e}_{tgt} . More formally, our third stage is a simple linear interpolation between \mathbf{e}_{tgt} and \mathbf{e}_{opt} . For a given hyperparameter $\eta \in [0, 1]$, we obtain

$$\bar{\mathbf{e}} = \eta \cdot \mathbf{e}_{tgt} + (1 - \eta) \cdot \mathbf{e}_{opt}, \quad (3)$$

which is the embedding that represents the desired edited image. We then apply the base generative diffusion process using the fine-tuned model, conditioned on $\bar{\mathbf{e}}$. This results in a low-resolution edited image, which is then super-resolved using the fine-tuned auxiliary models, conditioned on the target text. This generative process outputs our final high-resolution edited image $\bar{\mathbf{x}}$.

3.3. Implementation Details

Our framework is general and can be combined with different generative models. We demonstrate it using two different state-of-the-art text-to-image generative diffusion models: Imagen [53] and Stable Diffusion [50].

Imagen [53] consists of 3 separate text-conditioned diffusion models: (i) a generative diffusion model for 64×64 -pixel images; (ii) a super-resolution (SR) diffusion model turning 64×64 -pixel images into 256×256 ones; and (iii) another SR model transforming 256×256 -pixel images

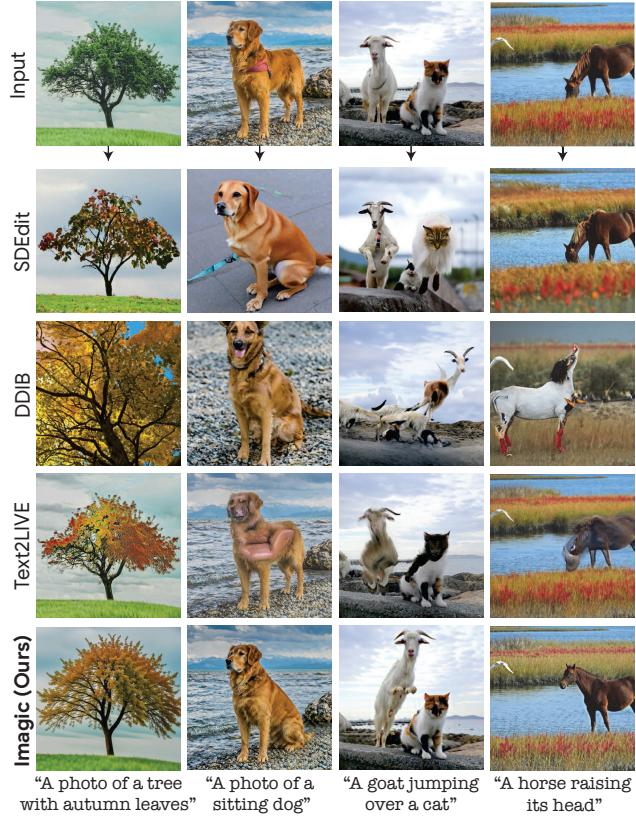


Figure 6. Method comparison. We compare SDEdit [38], DDIB [62], and Text2LIVE [8] to our method. *Imagic* successfully applies the desired edit, while preserving the original image details well.

into the 1024×1024 resolution. By cascading these 3 models [23] and using classifier-free guidance [24], Imagen constitutes a powerful text-guided image generation scheme.

We optimize the text embedding using the 64×64 diffusion model and the Adam [34] optimizer for 100 steps and a fixed learning rate of $1e-3$. We then fine-tune the 64×64 diffusion model by continuing Imagen’s training for 1500 steps for our input image, conditioned on the optimized embedding. In parallel, we also fine-tune the $64 \times 64 \rightarrow 256 \times 256$ SR diffusion model using the target text embedding and the original image for 1500 steps, in order to capture high-frequency details from the original image. We find that fine-tuning the $256 \times 256 \rightarrow 1024 \times 1024$ model adds little to no effect to the results, therefore we opt to use its pre-trained version conditioned on the target text. This entire optimization process takes around 8 minutes per image on two TPUv4 chips.

Afterwards, we interpolate the text embeddings according to Equation 3. Because of the fine-tuning process, using $\eta=0$ will generate the original image, and as η increases, the image will start to align with the target text. To maintain both image fidelity and target text alignment, we choose an intermediate η , usually residing between 0.6 and 0.8 (see Figure 9). We then generate with Imagen [53] with its pro-

vided hyperparameters. We find that using the DDIM [59] sampling scheme generally provides slightly improved results over the more stochastic DDPM scheme.

In addition to Imagen, we also implement our method with the publicly available Stable Diffusion model (based on Latent Diffusion Models [50]). This model applies the diffusion process in the latent space (of size $4 \times 64 \times 64$) of a pre-trained autoencoder, working with 512×512 -pixel images. We apply our method in the latent space as well. We optimize the text embedding for 1000 steps with a learning rate of $2e-3$ using Adam [34]. Then, we fine-tune the diffusion model for 1500 steps with a learning rate of $5e-7$. This process takes 7 minutes on a single Tesla A100 GPU.

4. Experiments

4.1. Qualitative Evaluation

We applied our method on a multitude of real images from various domains, with simple text prompts describing different editing categories such as: style, appearance, color, posture, and composition. We collect high-resolution free-to-use images from Unsplash and Pixabay. After optimization, we generate each edit with 8 random seeds and choose the best result. *Imagic* is able to apply various editing categories on general input images and texts, as we show in Figure 1 and the supplementary material. We experiment with different text prompts for the same image in Figure 2, showing the versatility of *Imagic*. Since the underlying generative diffusion model that we utilize is probabilistic, our method can generate different results for a single image-text pair. We show multiple options for the same edit using different random seeds in Figure 4, slightly tweaking η for each seed. This stochasticity allows the user to choose among these different options, as natural language text prompts can generally be ambiguous and imprecise.

While we use Imagen [53] in most of our experiments, *Imagic* is agnostic to the generative model choice. Thus, we also implement *Imagic* with Stable Diffusion [50]. In Figure 5 (and in the supplementary material) we show that *Imagic* successfully performs complex non-rigid edits also using Stable Diffusion while preserving the image-specific appearance. Furthermore, *Imagic* (using Stable Diffusion) exhibits smooth semantic interpolation properties as η is changed. We hypothesize that this smoothness property is a byproduct of the diffusion process taking place in a semantic latent space, rather than in the image pixel space.

4.2. Comparisons

We compare *Imagic* to the current leading general-purpose techniques that operate on a single input real-world image, and edit it based on a text prompt. Namely, we compare our method to Text2LIVE [8], DDIB [62], and SDEdit [38]. We use Text2LIVE’s default provided hyperparameters. We feed it with a text description of the tar-



Figure 7. **Embedding interpolation.** Varying η with the same seed, using the pre-trained (top) and fine-tuned (bottom) models.

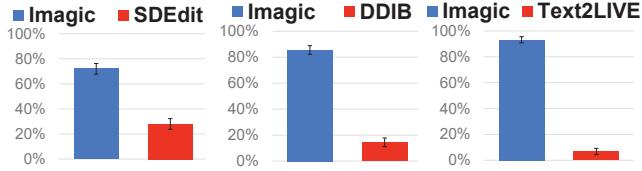


Figure 8. **User study results.** Preference rates (with 95% confidence intervals) for image editing quality of *Imagic* over *SDEdit* [38], *DDIB* [62], and *Text2LIVE* [8].

get object (*e.g.*, “dog”) and one of the desired edit (*e.g.*, “sitting dog”). For *SDEdit* and *DDIB*, we apply their proposed technique with the same *Imagen* [53] model and target text prompt that we use. We keep the diffusion hyperparameters from *Imagen*, and choose the intermediate diffusion timestep for *SDEdit* independently for each image to achieve the best target text alignment without drastically changing the image contents. For *DDIB*, we provide an additional source text.

Figure 6 shows editing results of different methods. For *SDEdit* and *Imagic*, we sample 8 images using different random seeds and display the result with the best alignment to both the target text and the input image. As can be observed, our method maintains high fidelity to the input image while aptly performing the desired edits. When tasked with a complex non-rigid edit such as making a dog sit, our method significantly outperforms previous techniques. *Imagic* constitutes the first demonstration of such sophisticated text-based edits applied on a single real-world image. We verify this claim through a user study in [subsection 4.3](#).

4.3. TEDBench and User Study

Text-based image editing methods are a relatively recent development, and *Imagic* is the first to apply complex non-rigid edits. As such, no standard benchmark exists for evaluating non-rigid text-based image editing. We introduce *TEDBench* (Textual Editing Benchmark), a novel collection of 100 pairs of input images and target texts describing a desired complex non-rigid edit. We hope that future research will benefit from *TEDBench* as a standardized evaluation set for this task.

We quantitatively evaluate *Imagic*’s performance via an

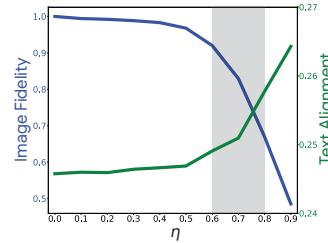


Figure 9. **Editability-fidelity tradeoff.** CLIP score (target text alignment) and $1-LPIPS$ (input image fidelity) as functions of η , averaged over 150 inputs. Edited images tend to match both the input image and text in the highlighted area.

extensive human perceptual evaluation study on *TEdBench*, performed using Amazon Mechanical Turk. Participants were shown an input image and a target text, and were asked to choose the better editing result from one of two options, using the standard practice of Two-Alternative Forced Choice (2AFC) [8, 35, 42]. The options to choose from were our result and a baseline result from one of: *SDEdit* [38], *DDIB* [62], or *Text2LIVE* [8]. In total, we collected 9213 answers, whose results are summarized in [Figure 8](#). As can be seen, evaluators exhibit a strong preference towards our method, with a preference rate of more than 70% across all considered baselines. See supplementary material for more details about the user study and method implementations.

4.4. Ablation Study

Fine-tuning and optimization We generate edited images for different η values using the pre-trained 64×64 diffusion model and our fine-tuned one, in order to gauge the effect of fine-tuning on the output quality. We use the same optimized embedding and random seed, and qualitatively evaluate the results in [Figure 7](#). Without fine-tuning, the scheme does not fully reconstruct the original image at $\eta = 0$, and fails to retain the image’s details as η increases. In contrast, fine-tuning imposes details from the input image beyond just the optimized embedding, allowing our scheme to retain these details for intermediate values of η , thereby enabling semantically meaningful linear interpolation. Thus, we conclude that model fine-tuning is essential for our method’s success. Furthermore, we experiment with the number of text embedding optimization steps in the supplementary material. Our findings suggest that optimizing the text embedding with a smaller number of steps limits our editing capabilities, while optimizing for more than 100 steps yields little to no added value.

Interpolation intensity As can be observed in [Figure 7](#), fine-tuning increases the η value at which the model strays from reconstructing the input image. While the optimal η value may vary per input (as different edits require different intensities), we attempt to identify the region in which the edit is best applied. To that end, we apply our editing scheme with different η values, and calculate the outputs’ CLIP score [21, 45] w.r.t. the target text, and their LPIPS score [68] w.r.t. the input image subtracted from 1. A higher CLIP score indicates better output alignment with the target text, and a higher 1–LPIPS indicates higher fidelity to the input image. We repeat this process for 150 image-text inputs, and show the average results in [Figure 9](#). We observe that for η values smaller than 0.4, outputs are almost identical to the input images. For $\eta \in [0.6, 0.8]$, the images begin to change (according to LPIPS), and align better with the text (as the CLIP score rises). Therefore, we identify this area as the most probable for obtaining satisfactory results. Note that while they provide a good sense of text or image alignment on average, CLIP score and LPIPS are imprecise measures that rely on neural network backbones, and their values noticeably differ for each different input image-text pair. As such, they are not suited for reliably choosing η for each input in an automatic way, nor can they faithfully assess an editing method’s performance.

4.5. Limitations

We identify two main failure cases of our method: In some cases, the desired edit is applied very subtly (if at all), therefore not aligning well with the target text. In other cases, the edit is applied well, but it affects extrinsic image details such as zoom or camera angle. We show examples of these two failure cases in the first and second row of [Figure 10](#), respectively. When the edit is not applied strongly enough, increasing η usually achieves the desired result, but it sometimes leads to a significant loss of original image details (for all tested random seeds) in a handful of cases. As for zoom and camera angle changes, these usually occur before the desired edit takes place, as we progress from a low η value to a large one, which makes circumventing them difficult. We demonstrate this in the supplementary material, and include additional failure cases in *TEdBench* as well.

These limitations can possibly be mitigated by optimizing the text embedding or the diffusion model differently, or by incorporating cross-attention control akin to Hertz et al. [20]. We leave those options for future work. Also, since our method relies on a pre-trained text-to-image diffusion model, it inherits the model’s generative limitations and biases. Therefore, unwanted artifacts are produced when the desired edit involves generating failure cases of the underlying model. For instance, Imagen is known to show sub-standard generative performance on human faces [53]. Additionally, the optimization required by *Imagic* (and other



Figure 10. **Failure cases.** Insufficient consistency with the target text (top), or changes in camera viewing angle (bottom).

editing methods [8]) is slow, and may hinder their direct deployment in user-facing applications.

5. Conclusions and Future Work

We propose a novel image editing method called *Imagic*. Our method accepts a single image and a simple text prompt describing the desired edit, and aims to apply this edit while preserving a maximal amount of details from the image. To that end, we utilize a pre-trained text-to-image diffusion model and use it to find a text embedding that represents the input image. Then, we fine-tune the diffusion model to fit the image better, and finally we linearly interpolate between the embedding representing the image and the target text embedding, obtaining a semantically meaningful mixture of them. This enables our scheme to provide edited images using the interpolated embedding. Contrary to other editing methods, our approach can produce sophisticated non-rigid edits that may alter the pose, geometry, and/or composition of objects within the image as requested, in addition to simpler edits such as style or color. It requires the user to provide only a single image and a simple target text prompt, without the need for additional auxiliary inputs such as image masks.

Our future work may focus on further improving the method’s fidelity to the input image and identity preservation, as well as its sensitivity to random seeds and to the interpolation parameter η . Another intriguing research direction would be the development of an automated method for choosing η for each requested edit.

Societal Impact Our method aims to enable complex editing of real world images using textual descriptions of the target edit. As such, it is prone to societal biases of the underlying text-based generative models, albeit to a lesser extent than purely generative methods since we rely mostly on the input image for editing. However, as with other approaches that use generative models for image editing, such techniques might be used by malicious parties for synthesizing fake imagery to mislead viewers. To mitigate this, further research on the identification of synthetically edited or generated content is needed.

References

- [1] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan: How to embed images into the stylegan latent space? In *Proceedings of the IEEE international conference on computer vision*, pages 4432–4441, 2019. 3
- [2] Rameen Abdal, Yipeng Qin, and Peter Wonka. Image2stylegan++: How to edit the embedded images? In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8296–8305, 2020. 3
- [3] Rameen Abdal, Peihao Zhu, Niloy Mitra, and Peter Wonka. Styleflow: Attribute-conditioned exploration of stylegan-generated images using conditional continuous normalizing flows, 2020. 3
- [4] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021. 3
- [5] Yuval Alaluf, Omer Tov, Ron Mokady, Rinon Gal, and Amit H Bermano. Hyperstyle: Stylegan inversion with hypernetworks for real image editing. *arXiv preprint arXiv:2111.15666*, 2021. 3
- [6] 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. 2, 4
- [7] Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Bryan Catanzaro, et al. eDiff-I: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022. 4
- [8] Omer Bar-Tal, Dolev Ofri-Amar, Rafaail Fridman, Yoni Kassten, and Tali Dekel. Text2LIVE: text-driven layered image and video editing. *arXiv preprint arXiv:2204.02491*, 2022. 2, 3, 6, 7, 8, 15
- [9] David Bau, Hendrik Strobelt, William Peebles, Jonas Wulff, Bolei Zhou, Jun-Yan Zhu, and Antonio Torralba. Semantic photo manipulation with a generative image prior. *arXiv preprint arXiv:2005.07727*, 2020. 3
- [10] Amit H Bermano, Rinon Gal, Yuval Alaluf, Ron Mokady, Yotam Nitzan, Omer Tov, Oren Patashnik, and Daniel Cohen-Or. State-of-the-art in the architecture, methods and applications of stylegan. In *Computer Graphics Forum*, volume 41, pages 591–611. Wiley Online Library, 2022. 2
- [11] Tsachi Blau, Roy Ganz, Bahjat Kawar, Alex Bronstein, and Michael Elad. Threat model-agnostic adversarial defense using diffusion models. *arXiv preprint arXiv:2207.08089*, 2022. 4
- [12] Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11315–11325, 2022. 3
- [13] Wenhui Chen, Hexiang Hu, Chitwan Saharia, and William W Cohen. Re-imagen: Retrieval-augmented text-to-image generator. *arXiv preprint arXiv:2209.14491*, 2022. 4
- [14] 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, 2021. 4
- [15] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. StarGAN v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020. 3
- [16] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. 3, 4
- [17] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 2, 4
- [18] Jin Gao, Jialing Zhang, Xihui Liu, Trevor Darrell, Evan Shelhamer, and Dequan Wang. Back to the source: Diffusion-driven test-time adaptation. *arXiv preprint arXiv:2207.03442*, 2022. 4
- [19] Erik Härkönen, Aaron Hertzmann, Jaakko Lehtinen, and Sylvain Paris. Ganspace: Discovering interpretable gan controls. *arXiv preprint arXiv:2004.02546*, 2020. 3
- [20] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control, 2022. 2, 4, 8
- [21] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. 8
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 3, 4, 5
- [23] Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022. 4, 6
- [24] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 6
- [25] Ali Jahanian, Lucy Chai, and Phillip Isola. On the “steerability” of generative adversarial networks. In *International Conference on Learning Representations*, 2019. 3
- [26] Tero Karras, Miika Aittala, Janne Hellsten, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Training generative adversarial networks with limited data. In *Proc. NeurIPS*, 2020. 3
- [27] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems*, 34, 2021. 3
- [28] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 3
- [29] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020. 3

- [30] Bahjat Kawar, Michael Elad, Stefano Ermon, and Jiaming Song. Denoising diffusion restoration models. In *Advances in Neural Information Processing Systems*, 2022. 4
- [31] Bahjat Kawar, Roy Ganz, and Michael Elad. Enhancing diffusion-based image synthesis with robust classifier guidance. *arXiv preprint arXiv:2208.08664*, 2022. 4
- [32] Bahjat Kawar, Jiaming Song, Stefano Ermon, and Michael Elad. JPEG artifact correction using denoising diffusion restoration models. *arXiv preprint arXiv:2209.11888*, 2022. 4
- [33] Gwanghyun Kim, Taesung Kwon, and Jong Chul Ye. Diffusionclip: Text-guided diffusion models for robust image manipulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2426–2435, 2022. 2, 3
- [34] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun, editors, *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015. 6
- [35] Nicholas Koltkin, Jason Salavon, and Gregory Shakhnarovich. Style transfer by relaxed optimal transport and self-similarity. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10051–10060, 2019. 7, 15
- [36] Oran Lang, Yossi Gandelsman, Michal Yarom, Yoav Wald, Gal Elidan, Avinatan Hassidim, William T. Freeman, Phillip Isola, Amir Globerson, Michal Irani, and Inbar Mosseri. Explaining in style: Training a gan to explain a classifier in stylespace. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 693–702, 2021. 3
- [37] Xihui Liu, Dong Huk Park, Samaneh Azadi, Gong Zhang, Arman Chopikyan, Yuxiao Hu, Humphrey Shi, Anna Rohrbach, and Trevor Darrell. More control for free! image synthesis with semantic diffusion guidance. *arXiv preprint arXiv:2112.05744*, 2021. 4
- [38] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jianjun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: guided image synthesis and editing with stochastic differential equations. In *International Conference on Learning Representations*, 2021. 3, 6, 7, 15
- [39] 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. *arXiv preprint arXiv:2112.10741*, 2021. 2, 4
- [40] Weili Nie, Brandon Guo, Yujia Huang, Chaowei Xiao, Arash Vahdat, and Anima Anandkumar. Diffusion models for adversarial purification. In *International Conference on Machine Learning (ICML)*, 2022. 4
- [41] Byong Mok Oh, Max Chen, Julie Dorsey, and Frédéric Durand. Image-based modeling and photo editing. In *Proceedings of the 28th annual conference on Computer graphics and interactive techniques*, pages 433–442, 2001. 1
- [42] Taesung Park, Jun-Yan Zhu, Oliver Wang, Jingwan Lu, Eli Shechtman, Alexei Efros, and Richard Zhang. Swapping autoencoder for deep image manipulation. *Advances in Neural Information Processing Systems*, 33:7198–7211, 2020. 7, 15
- [43] Or Patashnik, Zongze Wu, Eli Shechtman, Daniel Cohen-Or, and Dani Lischinski. Styleclip: Text-driven manipulation of stylegan imagery. *arXiv preprint arXiv:2103.17249*, 2021. 2, 3, 4
- [44] Vadim Popov, Ivan Vovk, Vladimir Gogoryan, Tasnima Sadekova, and Mikhail Kudinov. Grad-tts: A diffusion probabilistic model for text-to-speech. In *International Conference on Machine Learning*, pages 8599–8608. PMLR, 2021. 4
- [45] 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. 4, 8
- [46] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research*, 21:1–67, 2020. 4, 5, 14
- [47] 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, 3, 4
- [48] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. *arXiv preprint arXiv:2008.00951*, 2020. 3
- [49] Daniel Roich, Ron Mokady, Amit H Bermano, and Daniel Cohen-Or. Pivotal tuning for latent-based editing of real images. *ACM Transactions on Graphics (TOG)*, 42(1):1–13, 2022. 3, 4
- [50] 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. 3, 4, 5, 6
- [51] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. DreamBooth: fine tuning text-to-image diffusion models for subject-driven generation. *arXiv preprint arxiv:2208.12242*, 2022. 2, 4
- [52] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 4
- [53] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamvar Seyed Ghasemipour, Burcu Karagol Ayan, Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, 2022. 3, 4, 5, 6, 7, 8, 14
- [54] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 4

- [55] Hiroshi Sasaki, Chris G Willcocks, and Toby P Breckon. Unit-ddpm: Unpaired image translation with denoising diffusion probabilistic models. *arXiv preprint arXiv:2104.05358*, 2021. 4
- [56] Yujun Shen, Jinjin Gu, Xiaoou Tang, and Bolei Zhou. Interpreting the latent space of gans for semantic face editing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9243–9252, 2020. 3
- [57] Yujun Shen and Bolei Zhou. Closed-form factorization of latent semantics in gans. *arXiv preprint arXiv:2007.06600*, 2020. 3
- [58] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Machine Learning*, pages 2256–2265. PMLR, 2015. 4
- [59] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2020. 3, 6
- [60] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019. 4
- [61] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2020. 4
- [62] Xuan Su, Jiaming Song, Chenlin Meng, and Stefano Ermon. Dual diffusion implicit bridges for image-to-image translation. *arXiv preprint arXiv:2203.08382*, 2022. 3, 6, 7, 15
- [63] Lucas Theis, Tim Salimans, Matthew D Hoffman, and Fabian Mentzer. Lossy compression with Gaussian diffusion. *arXiv preprint arXiv:2206.08889*, 2022. 4
- [64] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *arXiv preprint arXiv:2102.02766*, 2021. 3, 4
- [65] Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021. 4
- [66] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011. 4
- [67] Julia Wolleb, Robin Sandkühler, Florentin Bieder, Philippe Valmaggia, and Philippe C Cattin. Diffusion models for implicit image segmentation ensembles. *arXiv preprint arXiv:2112.03145*, 2021. 4
- [68] 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. 8
- [69] Roland S Zimmermann, Lukas Schott, Yang Song, Benjamin A Dunn, and David A Klindt. Score-based generative classifiers. *arXiv preprint arXiv:2110.00473*, 2021. 4