

GeoDiffuser: Geometry-Based Image Editing with Diffusion Models

Rahul Sajnani^{1,2}, Jeroen Vanbaaar², Jie Min², Kapil Katyal², and Srinath Sridhar^{1,2}

¹ Brown University, Providence, Rhode Island, USA

² Amazon Robotics, North Reading, Massachusetts, USA

ivl.cs.brown.edu/projects/geodiffuser



Fig. 1: We introduce **GeoDiffuser**, a unified method to perform common 2D and 3D image editing tasks like object translation, 3D rotation, object removal, and re-scaling while preserving object style and inpainting disoccluded regions. Our method is a zero-shot optimization-based method that builds on top of a pre-trained diffusion model. We treat image editing as a geometric transformation of parts of the image and bake this directly into a shared attention-based edit optimization. In this figure, the top row shows natural images and the bottom row shows the edit. *vs.*

Abstract. The success of image generative models has enabled us to build methods that can edit images based on text or other user input. However, these methods are bespoke, imprecise, require additional information, or are limited to only 2D image edits. We present GeoDiffuser, a zero-shot optimization-based method that unifies common 2D and 3D image-based object editing capabilities into a single method. Our key insight is to view image editing operations as geometric transformations. We show that these transformations can be directly incorporated into the attention layers in diffusion models to implicitly perform editing operations. Our training-free optimization method uses an objective function that seeks to preserve object style but generate plausible images, for instance with accurate lighting and shadows. It also inpaints disoccluded parts of the image where the object was originally located. Given a natural image and user input, we segment the foreground object [28] and estimate a corresponding transform which is used by our optimization approach for editing. Figure 1 shows that GeoDiffuser can perform common 2D and 3D edits like object translation, 3D rotation, and removal. We present quantitative results, including a perceptual study, that shows how our approach is better than existing methods.

1 Introduction

Image generative models have seen significant progress recently. The most advanced diffusion-based models can now generate high-quality images almost indistinguishable from reality [46,48,50,61]. These models generate images with the desired content and detail by conditioning on text prompts, sometimes in combination with additional information like segmentation masks [64]. They have proliferated in use with many commercial products incorporating them [2–4].

Although realistic image generation is an important capability, in many cases, we may also want to edit generated or existing natural images. While past work relied on computer graphics techniques for image editing [11,27,29,65], recent works have put generative models to use for this problem. In particular, generative models have been shown to enable text-based edits [21,38,57], object stitching [15,53], object removal [48], and interactive edits using user-defined points [39,42,51], 3D transforms [44] or flow [18]. However, these methods have important limitations. Text-based editing methods are imprecise for edits requiring spatial control. Object stitching and removal methods cannot easily be extended to geometric edits. Finally, interactive point-/flow-based methods require additional input such as a text prompt or optical flow.

In this paper, we present **GeoDiffuser**, a method that unifies various image-based object editing capabilities into a single method. We take the view that common user-specified image editing operations can be cast as **geometric transformations** of parts of the image. For instance, 2D object translation or 3D object rotation can be represented as a bijective transformation of the foreground object. However, naively applying this transformation on the image is unlikely to produce plausible edits, for instance, due to mismatched lighting or shadows. To overcome this problem, we use diffusion models, specifically the general editing approach (see Figure 2) enabled by DDIM Inversion [37]. Our key contribution is to bake in the geometric transformation directly within the shared attention layers of a diffusion model to **preserve style** while enabling a wide range of user-specified 2D and 3D edits. Additionally, GeoDiffuser is a zero-shot optimization-based method that operates **without the need for any additional training** and can support any diffusion model with attention layers.

Figure 1 shows common image edits performed by GeoDiffuser on natural images. Without any hyperparameter tuning, our method can perform 2D edits like object translation or removal, or 3D edits like 3D rotation and translation. Given a natural image, we first segment the object of interest [28], and optionally, extract a depth map [60] for 3D edits. For each type of edit, we first compute a geometric transform based on user input and formulate an objective function for optimization. Unlike approaches that first ‘lift’ an object from an image and then stitch the transformed object back into the image [27], we implicitly perform these steps by applying the transform directly to the self- and cross-attention layers. Since attention captures both local and global image interactions, our results exhibit accurate lighting, shadows and reflection while inpainting the disoccluded image regions. Moreover, our objective function incorporates terms to preserve the original style of the transformed object.

We show extensive qualitative results that demonstrate that our method can perform multiple 2D and 3D editing operations using a single approach. To evaluate our method quantitatively, we provide experiments through a perceptual study as well as metrics that measure how well the foreground and background content is preserved during the edit. Results show that our method outperforms existing methods quantitatively while being general enough to perform various kinds of edits. To sum up, our main contributions are:

- A unified image editing approach that formulates common 2D and 3D editing operations as geometric transformations of parts of the image.
- GeoDiffuser, a zero-shot optimization-based approach that incorporates geometric transforms directly within the attention layers of diffusion models enabling realistic edits while preserving object style.
- Qualitative results of numerous 2D and 3D object edits enabled by our method without any hyperparameter tuning (see Figure 1).

2 Related Work

Image editing has been widely studied in computer vision and encapsulates a range of operations, such as object removal and addition [7, 53], style transfer [19, 22, 24, 26], and 2D and 3D transforms [27]. One challenge with this problem is to keep the edit consistent within the *global* context of the image. Traditional methods such as Poisson image editing [45] use gradients of the context to blend edits with existing pixels, while inpainting methods uses boundary and context to fill in pixels [59]. We limit our discussion below to generative model-based and 3D-aware editing methods.

Text-Guided Image Editing: There are several works using generative image models to edit images via changes to the text prompt. The preservation of subject identity in different settings can be achieved by textual inversion along with additional losses to finetune the generative model [49]. *Null-text* inversion is an inversion approach where a null-text embedding is optimized to match an inverted noise trajectory for a given input image along with attention reweighting [37]. Instead of an inversion process, text prompt edits can also be achieved by swap, or re-weighting of cross-attention maps derived from the visual and textual representation [21]. Edits with text prompts can also be achieved by using cross-attention from different prompts to manipulate self-attention maps [10]. Leveraging existing text-to-image models along with [9] gives the ability to generate paired data for finetuning a generative model to achieve text-guided editing results. These methods mostly produce images with style changes or enhancements, or object replacement. [16] leverage prompts and self guidance to perform 2D image edits of scaling and translation. However, it is difficult to guide the diffusion model to perform a specific 3D geometric transform based on a prompt. We extend some of the above approaches to build a method to handle geometric transforms without the need for additional training.

Non-Text-Guided Image Editing: Text-guided edits are mostly limited to appearance and style changes. Non-text-guided edits on the other hand, can

achieve a variety of edits. Point-based editing approaches can perform local image edits. [51] propose a motion supervised latent optimization between the reference and target edit, to guide the denoising to obtain the edit while preserving the object identity. Stroke-based editing can edit larger image regions, or even entire images [35], by projecting strokes onto the image manifold via diffusion. For these methods, edits such as translations are however not possible. ObjectStitch [53] along with inpainting can achieve translation where the denoising diffusion is applied to a target asked region, and guided by the embedding of the object to stitch. However, object style preservation is difficult in this setting. Recent methods [39,40] try to preserve identity and allow for translations while requiring no training. However, these are limited to 2D translations and scaling. An editing approach which first ‘lifts’ the object from a background is proposed in [44]. The background is inpainted and a depth-to-image generative model is used, which performs the denoising conditioned on an input depth. However, this approach needs an additional text prompt while ours does not. Additionally, we support various kinds of edits and not just 3D transforms. [18] is concurrent work that uses flow-guidance for image editing. However, optical flow can be much harder to obtain compared to depth [60]. We present a method that performs 2D and 3D edits using precise geometric transformations while preserving identity and not requiring additional user input.

3D-Aware Editing: Some methods have addressed the 3D editing problem [27] by ‘lifting’ objects into 3D and use 3D meshes and scene illumination to allow for proper blending of the edited object with the existing image context. Other methods use NeRF [13, 20, 58, 62, 63] or works [31, 32] learn over large-scale datasets [12], leverage geometry representations to perform edits but require multi-view images that are difficult to obtain. Edits are also directly applied to generative models, *e.g.*, [43] propose a point-based edit along with motion supervision to guide the neighboring pixels. The authors of [41] propose to represent foreground objects and background as neural feature fields, which can be edited and composited for a final output. The method of [30] addresses limitations of point-based editing in GANs, using template features rather than points for better tracking, and restricting search area around pixels to lines.

3 Background

Denoising Diffusion: We first briefly describe the concept of Denoising Diffusion Probabilistic Models (DDPM) used successfully by diffusion models for image generation [48]. Images can be considered as samples drawn from a data distribution $q(x)$. Denoising diffusion aims to learn a parameterized model $p_\theta(x)$ that approximates $q(x)$ and from which new samples can be drawn. The (forward) diffusion process iteratively adds noise to an input image x_0 , with $t = 0$, according to either a fixed or learned schedule, represented by α_t with $t \in [1, T]$. At each timestep, the latent encoding is performed according to a Gaussian distribution centered around the output of the previous timestep: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha}x_{t-1}, (1 - \alpha_t)\mathbf{I})$. The parameters vary over time such that $p_\theta(x_T) :=$

$\mathcal{N}(\mathbf{0}, \mathbf{I})$. Using the reparameterization trick, the noised version of input x_0 can directly be expressed as: $x_t = \sqrt{\alpha_t}x_0 + \sqrt{1 - \alpha_t}\epsilon_0$.

The reverse process, where noise is gradually removed at each step, can be expressed as the joint distribution $p_\theta(x_{0:T}) = p_\theta(x_T) \prod_{t=1}^T p_\theta^{(t)}(x_{t-1}|x_t)$. Under the assumption of trainable means and fixed variances, a neural network $\hat{\epsilon}_\theta(x_t, t)$ can be trained with the objective of minimizing a variational lower bound to estimate the source noise $\epsilon_0 \sim \mathcal{N}(\epsilon; \mathbf{0}, \mathbf{I})$ that determines x_t from x_0 : $L_\gamma(\epsilon_\theta) := \sum_{t=1}^T \gamma_t \mathbb{E}_{x_0 \sim q(x_0), \epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\|\epsilon_0 - \hat{\epsilon}_\theta(x_t, t)\|_2^2]$. For more details see [34].

Conditioning and Efficiency: This formulation can be extended to the conditional case, *i.e.*, $p_\theta(x_{0:T}|y)$. The condition y could be images, (encoded) text, or something else. The computational bottleneck is the number of denoising timesteps T , however a non-Markovian variant Denoising Diffusion Implicit Models (DDIM) was introduced to reduce the number of timesteps [52]. To further reduce the computational burden, the diffusion process for images can be performed in a lower dimensional latent (feature) space, as proposed by [48]. A perceptually optimized pretrained decoder takes the latent x_1 , and reconstructs the image x_0 . In our work, we use a latent diffusion model together with Classifier-Free Guidance (CFG) [23] for conditioning on text.

Attention: Attention was introduced as an alternative to recurrent networks and large receptive fields in convolution-based neural networks, for capturing local and global context [14, 56]. The scaled dot-product self-attention mechanism adopted in transformers has found widespread application in computer vision applications. The input is a tuple [Query (Q), Key (K), Value (V)], each with learnable parameters via linear layers. An attention layer constructs an attention map $\text{AM}(Q, K) := \text{Softmax}\left(\frac{QK^T}{\sqrt{d}}\right)$ and then computes attention as: $\text{Attention}(Q, K, V) := \text{AM}(Q, K)V$. Here, d is the dimension of the embedding.

In addition to self-attention, the query can be derived from another input source, *e.g.*, another modality, and using the key and values from the first input, the cross-attention between the two inputs can be computed via Section 3 and Section 3. An example of cross-attention is the activation of a word in a sentence with pixels in an image.

The correlation between semantics and pixels for image-text cross-attention can be modified in the denoising diffusion generative image setting to adjust the appearance of a given generated image [21]. In addition, deriving masks from cross-attention to guide self-attention [10] provides the ability to change the appearance of objects while maintaining object identity.

General Editing Framework: Prior works leverage the learned capabilities of diffusion models to perform edits to a given image, rather than a generated one. A general framework (see Figure 2) that is followed in all these works is to first perform an inversion [37, 52] on the image. This inversion provides us with a noise latent that sets a good starting point to regenerate the input image as well as to edit it. Starting from the inverted noise latent, two parallel diffusion processes generate the input image as well as the edited image. The first **reference diffusion process** generates the original input and, in our work, helps preserve un-edited regions of the image. An **edit diffusion process** runs

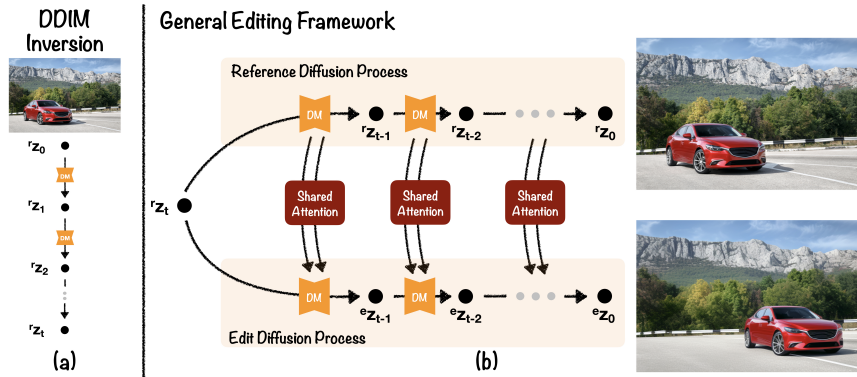


Fig. 2: General image editing framework using diffusion models. **(a) DDIM Inversion:** The process of obtaining noise trajectory $\{r_{z_0}, r_{z_1}, \dots, r_{z_t}\}$ for the reference image [52]. **(b) General Editing Framework:** The Reference Diffusion Process guides the Edit Diffusion Process to achieve the desired edit. In **GeoDiffuser**, we perform *geometric* 2D and 3D edits by transforming the shared attention layers leading to plausible edits that preserves object style, inpainting disoccluded background, and adding details (*e.g.*, the shadow cast by the car).

in parallel that utilizes the attention blocks from the reference process to perform the desired edit. This **shared attention** is a key insight for our proposed work. The editing framework is in sketched Figure 2 (b).

4 GeoDiffuser: Geometry-Based Image Editing

The goal of GeoDiffuser is to enable editing of segmented foreground objects in either natural or generated images. We take the view that common editing operations like 2D translation, 3D object rotation or object removal can be expressed as geometric transformations of parts of the image. Naively applying this transform to segmented foreground objects typically produces poor results w.r.t. image context and does not fill in the disoccluded background. We propose to use diffusion models to realistically edit the image and preserve object style.

Supported Operations: In this paper, we focus on *geometric edits* to an image \mathcal{I} specified by users through sliders that control transformations of foreground objects. In particular, we unify three kinds of edit operations that previously required separate bespoke methods: (1) **2D object edit** operations deal with realistically translating or scaling segmented objects within the image including inpainting the background where the object was originally located. (2) **3D object edit** operations deal with realistically transforming objects based on user-specified 3D rotation, translation or scaling and inpainting any disoccluded background as a result of the edit. Finally, (3) **object removal** refers to the operation of removing the segmented object completely and inpainting the disoccluded background.

In contrast with previous approaches, we formulate edits as an optimization problem based on the shared attention and leverage a pre-trained text-to-image Stable Diffusion model [48] to perform the edit. Notably, our method requires no training and can use any diffusion model with attention. Given an image \mathcal{I} , an object mask M_{obj} , a user-specified 2D or 3D transformation T , our goal is to edit the object in the image and inpaint any disoccluded regions introduced by the edit. To compute T for 3D edits, we use a depth map D obtained from DepthAnything [60] or simply by setting a constant depth of 0.5 m. This enables us to edit in-the-wild natural images without any additional user input.

4.1 Edits via Shared Attention

Each edit operation begins by performing a DDIM inversion [52] on the given image (Figure 2 (a)). Inverting the image provides us with the latent noise trajectory that will guide the edit diffusion process. We then perform the reverse diffusion process along with the geometry-aware attention sharing mechanism as sketched in Figure 2 (b). This attention sharing mechanism along with optimizing for the image latents as well as text embeddings is the key to achieve the desired geometric edit. Figure 3 (a) depicts the process for the shared attention blocks from Figure 2 (b).

Image Inversion: For inversion, we use direct inversion [25] on the image \mathcal{I} with the null prompt "". This inversion provides us with latents $\{^r z_t, ^r z_{t-1} \dots ^r z_0\}$ that preserve the image for the reference denoising process and guide the edit.

2D Edits: GeoDiffuser can perform 2D edits without requiring a depth map. Through a user interface, we can obtain a transformation T corresponding to a desired 2D translation or scaling. We define a 2D signal $S : [0, 1]^2 \rightarrow \mathbb{R}^C$ that stores a per-pixel feature in the image. The signal S can represent the RGB values or even the features of a deep network defined at each coordinate. Given a per-pixel edit \mathcal{F} defined on S , our shared attention mechanism uses \mathcal{F} to transform this signal for the desired edit. In our case, this signal is the Query embedding of the attention layer.

3D Edits: 2D edits are limited as they do not leverage the geometry of objects. We can extend 2D edits to 3D by additionally incorporating depth information D monocular depth estimators [8, 60] or simply a constant billboard depth map. The user specifies a 3D rigid transformation T which can then be used to compute the per-pixel edit \mathcal{F} as

$$\mathcal{F}(S)[u] := S[PTD[u]P^{-1}u].$$

Here, P is the camera intrinsic matrix that is used to project points in the image and u is the coordinate location of the signal. This edit field \mathcal{F} captures the 3D shape of the visible region of the object and provides an estimate of the desired location of the object. Note that if the per-pixel edit field is known, *e.g.*, from optical flow, we do not need a depth map for guidance. However, optical flow is much more challenging to obtain for a single image compared to depth maps.

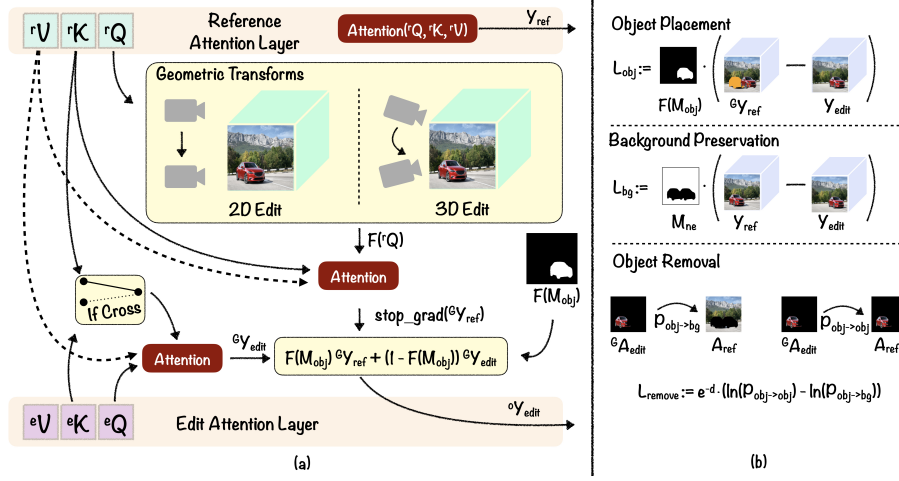


Fig. 3: (a) **GeoDiffuser** attention sharing mechanism that leverages the geometric transformation $\mathcal{F}(\cdot)$ transform the reference attention $\mathcal{G}Y_{ref}$ to guide the edit attention layer. (b) **Optimization Loss Functions** that penalize the latents and text-embeddings to perform the desired geometric edit.

Object Removal: Object removal introduces disocclusions to the background where the object was originally located. We propose an additional loss (see Section 4.2) for the optimization of the diffusion latents to handle such disocclusions. Disocclusions can also occur for 2D and 3D edits, and we consider such edits to be composites of removal and placement operations. Our proposed formulation for latent optimization thus extends to those edits as well.

Shared Attention: A key insight of our work is that we can transform objects by merely applying the edit \mathcal{F} to the query embeddings of the reference attention (Figure 3 (a)). Let ${}^rQ, {}^rK, {}^rV$ be the queries, keys, and values within the diffusion model of the reference denoising process and ${}^eQ, {}^eK, {}^eV$ be the queries, keys, and values of the corresponding attention block in the edit denoising process. The reference attention guidance $\mathcal{G}Y_{ref}$ and edit attention guidance $\mathcal{G}Y_{edit}$ are then given by

$$\mathcal{G}Y_{ref} := \text{Attention}(\mathcal{F}({}^rQ), {}^rK, {}^rV) \quad (1)$$

$$\mathcal{G}Y_{edit} := \begin{cases} \text{Attention}({}^eQ, {}^rK, {}^rV), & \text{if SelfAttention} \\ \text{Attention}({}^eQ, {}^eK, {}^rV), & \text{otherwise} \end{cases} \quad (2)$$

The dot product of the edit query embeddings eQ with the reference key embeddings rK in Eq. 2 provides us with correspondences between the edit and reference denoising process. These correspondences preserve the background and foreground features during the edit. To place the object at the desired location, the edit and reference attention guidance should approximately be the same ($\mathcal{G}Y_{ref} \approx \mathcal{G}Y_{edit}$) for the foreground. Note that they need not be exactly the same

in the case of an ill-defined edit \mathcal{F} . We then transform the output ${}^{\mathcal{O}}Y_{\text{edit}}$ of the edit attention layer by

$${}^{\mathcal{O}}Y_{\text{edit}} := \mathcal{F}(M_{\text{obj}}) \cdot {}^{\mathcal{G}}Y_{\text{ref}} + (1 - \mathcal{F}(M_{\text{obj}})) \cdot {}^{\mathcal{G}}Y_{\text{edit}}, \quad (3)$$

where $\mathcal{F}(M_{\text{obj}})$ refers to the foreground mask after applying the transformation \mathcal{F} . In other words, Eq. 3 aim to preserve identity for the object in the edit at its target location, while simultaneously preserve identity and consistency for the remaining pixels (or background).

4.2 Optimization

GeoDiffuser is a zero-shot optimization-based method that operates **without the need for any additional training**. We achieve this via optimization of the latents for edit guidance. The shared attention guidance provides us with a proxy of where the foreground object must be placed after the edit. We formulate an optimization procedure for the latents, in order to fill in the disocclusions and penalize the deviation of the edit attention guidance from the reference attention guidance. The loss functions used to penalize the diffusion latents in the optimization (shown in Fig. 3 (b)) are explained in detail next.

Background Preservation Loss: Performing shared attention guidance along with optimization could result in the un-edited regions of the image to also be changed. We introduce a background preservation loss to prevent this. Let the mask M_{ne} represent the non-editable region of the image. We define the background preservation loss as

$$\mathcal{L}_{bg} := \text{mean}(M_{\text{ne}} \cdot \|{}^{\mathcal{G}}Y_{\text{edit}} - Y_{\text{ref}}\|_1). \quad (4)$$

Here, $Y_{\text{ref}} = \text{Attention}(Q_{\text{ref}}, K_{\text{ref}}, V_{\text{ref}})$ is the attention block output for the reference de-noising process. The reference attention preserves the style of the image and constrains the optimization towards preserving the background.

Object Preservation Loss: Occasionally, the optimization changes the foreground region of the image. This causes loss of detail in the foreground. To prevent this, we penalize the deviation between the edit guidance and the reference guidance within the transformed foreground mask by

$$\mathcal{L}_{obj} := \text{mean}(\mathcal{F}(M_{\text{obj}}) \cdot \|{}^{\mathcal{G}}Y_{\text{edit}} - {}^{\mathcal{G}}Y_{\text{ref}}\|_1). \quad (5)$$

Note, we don't use this loss for object removal.

Inpainting Loss: To inpaint the disoccluded regions of the image, we maximize the difference between the edit guidance attention map ${}^{\mathcal{G}}A_{\text{edit}} := \text{AM}({}^eQ, {}^rK)$ and the reference guidance attention map $A_{\text{ref}} := \text{AM}({}^rQ, {}^rK)$. Let $\rho_{\text{obj} \rightarrow \text{bg}}$ represent the maximum normalized correlation score for each row in the foreground mask of the attention map ${}^{\mathcal{G}}A_{\text{edit}}$ to each row in the background mask of the reference attention map A_{ref} . We can similarly compute $\rho_{\text{obj} \rightarrow \text{obj}}$ that provides us with the maximum foreground to foreground normalized correlation (see Figure 3 (b)). Our goal is to reduce $\rho_{\text{obj} \rightarrow \text{obj}}$ and increase $\rho_{\text{obj} \rightarrow \text{bg}}$. We want to inject the disoccluded region with features from the background and ensure that the

diffusion process doesn’t in-paint the same features. We penalize for this using

$$\mathcal{L}_{remove} := \text{mean} \left(e^{-d_{\text{obj} \rightarrow \text{bg}}} (\ln(\rho_{\text{obj} \rightarrow \text{obj}}) - \ln(\rho_{\text{obj} \rightarrow \text{bg}})) \right). \quad (6)$$

Here, $d_{\text{obj} \rightarrow \text{bg}}$ is the coordinate distance between the locations of the attention map. The loss weighted by coordinate distance ensures that the foreground region inpaints the region using features within its vicinity. The negative log forces the object to background correlation $\rho_{\text{obj} \rightarrow \text{bg}}$ to increase and also reduces object-object correlation forcing the inpainted region to not be filled by the same object.

Smoothness Constraint: We additionally penalize the edit attention guidance ${}^G Y_{\text{edit}}$ for smoothness by penalizing the absolute value of its gradients using \mathcal{L}_s .

In our experiments, we found that the inpainting loss is hard to optimize and changes every image differently. To combat this, we devise an adaptive optimization scheme that increases the weight of the removal loss if the loss is more than -1.8 and reduce the loss weight if the removal loss is lower than -6 .

4.3 Implementation

The shared attention along with the loss functions defined above, enable performing geometry image edits as a reverse diffusion process by optimizing the latents and text embeddings. To make the optimization faster, we optimize every alternate step for the initial 32 diffusion steps. We set an initial learning rate of 1.5 and linearly decay it to 0. We share attention across all blocks within the UNet till step 45. All our experiments are performed on an Nvidia A40 with a run time of 25 seconds (for removal) up to 50 seconds (for 2D and 3D edits) on image resolution of 512. We use [47] for projecting, splatting, and rendering in our attention sharing mechanism.

5 Results & Experiments

In this section, we present visual examples of our editing results and quantitative results of a perceptual study and other visual metrics of editing quality.

Dataset: To measure the efficacy of our method we collected a dataset of real images from Adobe Stock images [1] to ensure we exclude generative AI images. We collect 70 images corresponding to the prompts *dog, car, cat, bear, mug, lamp, boat, plane, living-room, peaceful scenery*. We also test on real in-the-wild images from [17] and generated images from [44]. For many images in our dataset, we show multiple 2D and 3D edits demonstrating the general editing capabilities of GeoDiffuser.

Baselines: Since there is no extant method that performs all types of edits that we support, we compare each edit type to a different baseline. For the object removal, we compare with a state-of-the-art off-the-shelf **LaMa** image in-painting model [54], dilating object masks to make LaMa work better. For the 3D edit operations, we design a baseline based on **Zero123-XL** [32]. For this baseline, we first use [54] to in-paint the region of the removed foreground object first. We

then Zero123-XL to predict the novel view of the transformed object and composite it to the in-painted background image using Laplacian pyramid blending. We are unable to compare with recent/concurrent work such as DiffusionHandles [44], Motion Guidance [18] and ObjectStitch [53] since no public code is available. Moreover, our tests with [36] on real images produced poor results and therefore we exclude them. For 2D edits, we compare with Dragon Diffusion [39,40]. Since this method requires a text prompt while our method does not, we manually added text descriptions to our dataset. We test our method against Zero123-XL for geometry editing and LaMa for inpainting with a perceptual study. We additionally benchmark our method against Zero123-XL and Dragon Diffusion using community accepted metrics of Mean Distance, Clip Similarity Score, and Warp Error.

5.1 Quantitative Evaluation

Perceptual Study: To evaluate whether participants are satisfied with our editing results compared to other work, we conducted a perceptual study. This was setup as a forced choice questionnaire where participants had to select one of two options as containing the best edit result. Of the two randomly presented options, one was ours and the other was a baseline. In total we presented 70 images (30 for removal, 40 for other transforms) from our dataset. The questions were divided in three categories: edit realism (ER), edit adherence (EA), and removal edit realism (RRE). For removal, we generated results with LaMa [54], and for the remaining two categories, results were generated with Lama followed by Zero123-XL [32]. Each participant answered 12 ER questions, 12 EA questions and 6 RRE questions, for a total of 30 visual questions. In total 53 users participated in the study for which they received no compensation. Please see the supplementary document for more details.

Figure 4 shows the participant preference rate for each division of the study. For RRE, out of the 318 choices, participants preferred our method in 94.06% of the time, which shows that GeoDiffuser is better able to inpaint the disoccluded background regions, especially removing shadows (see Figure 7). For ER, our method was preferred 86.48% out of 636 cases. This demonstrates that GeoDiffuser preserves object style better than other methods, especially in transforming shadows and reflections. Finally, for EA we included 16 2D and 24 3D edits. Our method was preferred 88.48% out of 636 cases. This demonstrates that our method more faithfully performs the intended edit, even challenging ones such as 3D rotation. Whereas the baseline is only capable of performing edits from a more narrow range.

Metrics: In addition to the perceptual study, we also provide metrics to evaluate and compare edit adherence as well as style preservation quantitatively. We performed a total of 102 edits on our dataset comprising of 36 2D edits and 66 3D edits. Previously used metrics such as FID and Image Fidelity (IF) [39, 51] are not suitable for geometric edits because there could be large visual difference (e.g., large translation) and they do not measure disocclusion inpainting quality. Therefore, we use three other metrics to better evaluate methods: (1) The **Mean**

	MD ↓	Warp Error ↓	Clip Similarity ↑
2D & 3D Edits			
Zero123-XL + Lama [31, 54]	19.628	0.148	0.950
GeoDiffuser (Ours)	6.420	0.0930	0.966
2D Edits			
Zero123-XL + Lama [31, 54]	20.29	0.134	0.929
Dragon Diffusion [39]	40.38	0.160	0.965
GeoDiffuser (Ours)	5.65	0.098	0.963

Table 1: Our method adheres to the desired edit having the least **Mean Distance** and **Warp Error**. While DragonDiffusion performs slightly better on CLIP similarity, this metric does not fully measure small edits.

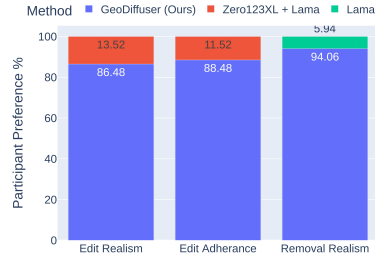


Fig. 4: Results from perceptual study show that participants prefer our edits over [31] and [54] in a majority of the cases.

Distance (MD) metric computes interest points on the foreground of the image using SIFT [33] and finds correspondences between the input and edited image using DiFT [55]. We then measure the distance between the correspondence estimated by DiFT and the edit specified by the user. This metric measures how well each approach transforms the foreground object. (2) the **Warp Error** metric forward warps the foreground region of the input image to the edited image and compute the absolute difference between their pixels for the transformed foreground. This metric measures how well each approach adheres to the edit. (3) the **CLIP Similarity** metric computes the CLIP image embedding [58] of the input and edited image and measures the cosine similarity. A higher cosine similarity indicates better preservation of global context. A good approach should preserve the global context of the image as well as adhere to the edit.

Table 1 shows quantitative comparison of our method with the baselines. GeoDiffuser outperforms baselines with an average **MD** of 5.65 for 2D edits and 6.42 for all edits. We also have the lowest warp loss of 0.098 for 2D edits and 0.093 for all edits. Dragon Diffusion does not perform well in these tasks since their method fails to inpaint disocclusions or preserve the foreground. Zero123-XL baseline performs better but since it is not trained on real-world images, it does not preserve the foreground object well resulting in incorrect DiFT correspondences. All method seem to preserve the context of the scene with a clip score above 0.92. For 2D edits, Dragon Diffusion has the highest clip similarity score of 0.965 however, the edits are performed on images with higher resolution compared to ours which might explain the difference.

Ablations: We perform a visual ablation of our design choices – please see the supplementary material for more details and experiments. Figures 5 and 6 shows the importance of the attention sharing mechanism and adaptive optimization. We can see a degradation in style preservation of the edit when we don’t perform geometric attention sharing till step 45. Without the adaptive optimization, we need image specific tuning for loss weights which is not scalable.



Fig. 5: Ablation of adaptive optimization. Without adaptive optimization, the same losses successfully inpaint some images while others fail (middle row). With our adaptive optimization, the same loss function works well for any image.

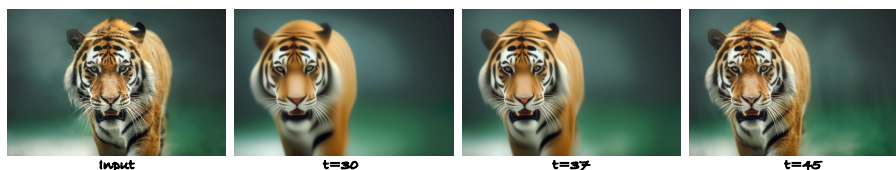


Fig. 6: Ablation on the number of steps for which we use our shared attention mechanism. Increasing the steps t better preserves object style (translation edit shown).

Qualitative Results: We show more qualitative results of 2D and 3D edits performed by GeoDiffuser in Figure 7. Notice how our method not only removes/transforms objects but also the object’s reflection and shadows.

6 Conclusion

GeoDiffuser is a unified method to perform common 2D and 3D object edits on images. Our approach is a zero-shot optimization-based method that uses diffusion models to achieve these edits. The key insight is to formulate image editing as a geometric transformation and incorporate it directly within the shared attention layers in a diffusion model-based editing framework. Results show that our single approach can handle a wide variety of image editing operations, and produces better results compared to previous work.

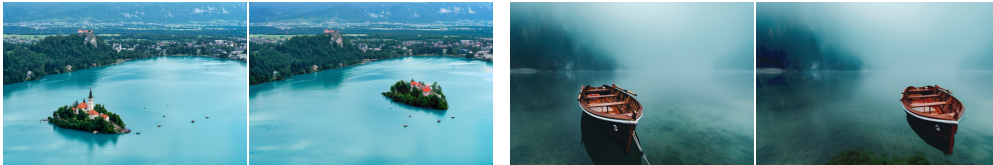
Limitations & Future Work: While we can handle background disocclusions, we cannot yet handle foreground object disocclusions resulting from large 3D motions. Our method also occasionally generates artifacts due to downsampled attention masks. We plan to address these limitations in future work.

Acknowledgements: Part of this work was done during Rahul’s internship at Amazon. This work was additionally supported by NSF grant CNS #2038897, ONR grant N00014-22-1-259, and an AWS Cloud Credits award.

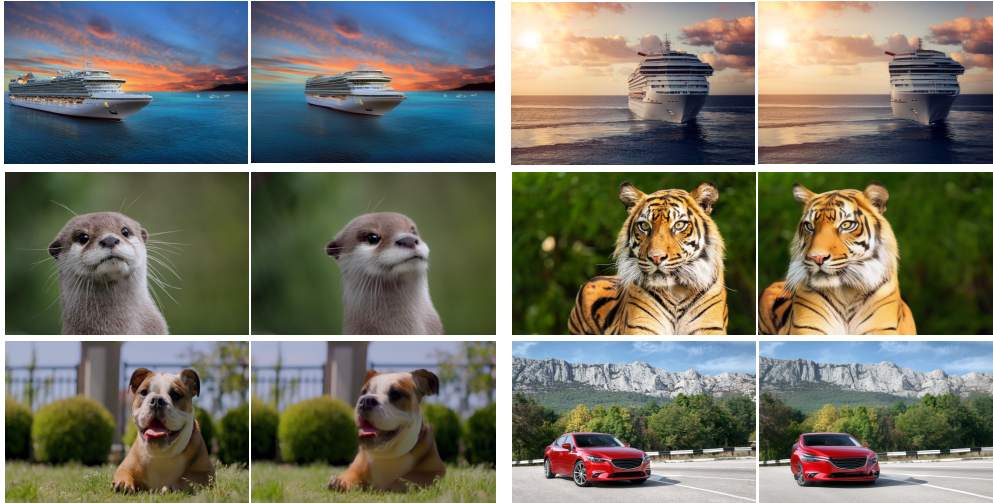
Removal



2D Edits



3D Edits



Scaling



Fig. 7: Qualitative results showing all variations of 2D and 3D edits performed by **GeoDiffuser** on natural images. Notice how our method not only removes/transforms objects but also the object's reflection and shadows (car, couch, boat). For 3D edits, our method produces plausible results for rotations as high as 30° . For scaling, we can perform both uniform and non-uniform scaling operations.

References

1. Adobe stock. <http://web.archive.org/web/20080207010024/http://www.808multimedia.com/winnt/kernel.htm>
2. Amazon titan image generator, multimodal embeddings, and text models are now available in amazon bedrock | aws news blog. <https://aws.amazon.com/blogs/aws/amazon-titan-image-generator-multimodal-embeddings-and-text-models-are-now-available-in-amazon-bedrock/>, (Accessed on 03/04/2024)
3. Dall-e 2. <https://openai.com/dall-e-2>, (Accessed on 03/04/2024)
4. Gemini - chat to supercharge your ideas. <https://gemini.google.com/>, (Accessed on 03/04/2024)
5. Qualtrics. <https://www.qualtrics.com>
6. Abid, A., Abdalla, A., Abid, A., Khan, D., Alfozan, A., Zou, J.: Gradio: Hassle-free sharing and testing of ml models in the wild. arXiv preprint arXiv:1906.02569 (2019)
7. Avrahami, O., Lischinski, D., Fried, O.: Blended diffusion for text-driven editing of natural images. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18208–18218 (2022)
8. Bhat, S.F., Birkel, R., Wofk, D., Wonka, P., Müller, M.: Zoedepth: Zero-shot transfer by combining relative and metric depth. arXiv preprint arXiv:2302.12288 (2023)
9. Brooks, T., Holynski, A., Efros, A.A.: Instructpix2pix: Learning to follow image editing instructions. In: CVPR (2023)
10. Cao, M., Wang, X., Qi, Z., Shan, Y., Qie, X., Zheng, Y.: Masactrl: Tuning-free mutual self-attention control for consistent image synthesis and editing (2023)
11. Chen, T., Zhu, Z., Shamir, A., Hu, S.M., Cohen-Or, D.: 3-sweep: Extracting editable objects from a single photo. ACM Transactions on graphics (TOG) **32**(6), 1–10 (2013)
12. Deitke, M., Schwenk, D., Salvador, J., Weihs, L., Michel, O., Vanderbilt, E., Schmidt, L., Ehsani, K., Kembhavi, A., Farhadi, A.: Objaverse: A universe of annotated 3d objects. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13142–13153 (2023)
13. Dong, J., Wang, Y.X.: Vica-nerf: View-consistency-aware 3d editing of neural radiance fields. Advances in Neural Information Processing Systems **36** (2024)
14. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. In: 9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021. OpenReview.net (2021), <https://openreview.net/forum?id=YicbFdNTTy>
15. Dwibedi, D., Misra, I., Hebert, M.: Cut, paste and learn: Surprisingly easy synthesis for instance detection. In: Proceedings of the IEEE international conference on computer vision. pp. 1301–1310 (2017)
16. Epstein, D., Jabri, A., Poole, B., Efros, A., Holynski, A.: Diffusion self-guidance for controllable image generation. Advances in Neural Information Processing Systems **36** (2024)
17. Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The kitti dataset. The International Journal of Robotics Research **32**(11), 1231–1237 (2013)
18. Geng, D., Owens, A.: Motion guidance: Diffusion-based image editing with differentiable motion estimators. arXiv preprint arXiv:2401.18085 (2024)

19. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. *Advances in neural information processing systems* **27** (2014)
20. Haque, A., Tancik, M., Efros, A.A., Holynski, A., Kanazawa, A.: Instruct-nerf2nerf: Editing 3d scenes with instructions. *arXiv preprint arXiv:2303.12789* (2023)
21. Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., Cohen-Or, D.: Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626* (2022)
22. Hertz, A., Voynov, A., Fruchter, S., Cohen-Or, D.: Style aligned image generation via shared attention. *arXiv preprint arXiv:2312.02133* (2023)
23. Ho, J., Salimans, T.: Classifier-free diffusion guidance (2022)
24. Jing, Y., Yang, Y., Feng, Z., Ye, J., Yu, Y., Song, M.: Neural style transfer: A review. *IEEE transactions on visualization and computer graphics* **26**(11), 3365–3385 (2019)
25. Ju, X., Zeng, A., Bian, Y., Liu, S., Xu, Q.: Direct inversion: Boosting diffusion-based editing with 3 lines of code. *arXiv preprint arXiv:2310.01506* (2023)
26. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 4401–4410 (2019)
27. Kholgade, N., Simon, T., Efros, A., Sheikh, Y.: 3d object manipulation in a single photograph using stock 3d models. *ACM Transactions on graphics (TOG)* **33**(4), 1–12 (2014)
28. Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.Y., et al.: Segment anything. *arXiv preprint arXiv:2304.02643* (2023)
29. Lalonde, J.F., Hoiem, D., Efros, A.A., Rother, C., Winn, J., Criminisi, A.: Photo clip art. *ACM transactions on graphics (TOG)* **26**(3), 3–es (2007)
30. Ling, P., Chen, L., Zhang, P., Chen, H., Jin, Y.: Freedrag: Point tracking is not you need for interactive point-based image editing. In: *arXiv preprint arXiv:2307.04684* (2023)
31. Liu, M., Shi, R., Chen, L., Zhang, Z., Xu, C., Wei, X., Chen, H., Zeng, C., Gu, J., Su, H.: One-2-3-45++: Fast single image to 3d objects with consistent multi-view generation and 3d diffusion. *arXiv preprint arXiv:2311.07885* (2023)
32. Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 9298–9309 (2023)
33. Lowe, G.: Sift-the scale invariant feature transform. *Int. J* **2**(91-110), 2 (2004)
34. Luo, C.: Understanding diffusion models: A unified perspective (2022)
35. Meng, C., He, Y., Song, Y., Song, J., Wu, J., Zhu, J.Y., Ermon, S.: SDEdit: Guided image synthesis and editing with stochastic differential equations. In: *International Conference on Learning Representations* (2022), https://openreview.net/forum?id=aBsCjcPu_tE
36. Michel, O., Bhattad, A., VanderBilt, E., Krishna, R., Kembhavi, A., Gupta, T.: Object 3dit: Language-guided 3d-aware image editing. *Advances in Neural Information Processing Systems* **36** (2024)
37. Mokady, R., Hertz, A., Aberman, K., Pritch, Y., Cohen-Or, D.: Null-text inversion for editing real images using guided diffusion models (2022)
38. Mokady, R., Hertz, A., Aberman, K., Pritch, Y., Cohen-Or, D.: Null-text inversion for editing real images using guided diffusion models. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 6038–6047 (2023)

39. Mou, C., Wang, X., Song, J., Shan, Y., Zhang, J.: Dragondiffusion: Enabling drag-style manipulation on diffusion models (2023)
40. Mou, C., Wang, X., Song, J., Shan, Y., Zhang, J.: Diffeditor: Boosting accuracy and flexibility on diffusion-based image editing. arXiv preprint arXiv:2402.02583 (2024)
41. Niemeyer, M., Geiger, A.: Giraffe: Representing scenes as compositional generative neural feature fields. In: Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR) (2021)
42. Pan, X., Tewari, A., Leimkühler, T., Liu, L., Meka, A., Theobalt, C.: Drag your gan: Interactive point-based manipulation on the generative image manifold. In: ACM SIGGRAPH 2023 Conference Proceedings. pp. 1–11 (2023)
43. Pan, X., Tewari, A., Leimkühler, T., Liu, L., Meka, A., Theobalt, C.: Drag your gan: Interactive point-based manipulation on the generative image manifold. In: ACM SIGGRAPH 2023 Conference Proceedings (2023)
44. Pandey, K., Guerrero, P., Gadelha, M., Hold-Geoffroy, Y., Singh, K., Mitra, N.: Diffusion handles: Enabling 3d edits for diffusion models by lifting activations to 3d (2023)
45. Pérez, P., Gangnet, M., Blake, A.: Poisson image editing. *ACM Trans. Graph.* **22**(3), 313–318 (jul 2003). <https://doi.org/10.1145/882262.882269>, <https://doi.org/10.1145/882262.882269>
46. Ramesh, A., Dhariwal, P., Nichol, A., Chu, C., Chen, M.: Hierarchical text-conditional image generation with clip latents. arXiv preprint arXiv:2204.06125 **1**(2), 3 (2022)
47. Ravi, N., Reizenstein, J., Novotny, D., Gordon, T., Lo, W.Y., Johnson, J., Gkioxari, G.: Accelerating 3d deep learning with pytorch3d. arXiv preprint arXiv:2007.08501 (2020)
48. Rombach, R., Blattmann, A., Lorenz, D., Esser, P., Ommer, B.: High-resolution image synthesis with latent diffusion models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR). pp. 10684–10695 (June 2022)
49. Ruiz, N., Li, Y., Jampani, V., Pritch, Y., Rubinstein, M., Aberman, K.: Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation (2022)
50. Saharia, C., Chan, W., Saxena, S., Li, L., Whang, J., Denton, E.L., Ghasemipour, K., Gontijo Lopes, R., Karagol Ayan, B., Salimans, T., et al.: Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems* **35**, 36479–36494 (2022)
51. Shi, Y., Xue, C., Pan, J., Zhang, W., Tan, V.Y., Bai, S.: Dragdiffusion: Harnessing diffusion models for interactive point-based image editing. arXiv preprint arXiv:2306.14435 (2023)
52. Song, J., Meng, C., Ermon, S.: Denoising diffusion implicit models. In: International Conference on Learning Representations (2021), <https://openreview.net/forum?id=St1giarCHLP>
53. Song, Y., Zhang, Z., Lin, Z., Cohen, S., Price, B., Zhang, J., Kim, S.Y., Aliaga, D.: Objectstitch: Object compositing with diffusion model. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 18310–18319 (2023)
54. Suvorov, R., Logacheva, E., Mashikhin, A., Remizova, A., Ashukha, A., Silvestrov, A., Kong, N., Goka, H., Park, K., Lempitsky, V.: Resolution-robust large mask inpainting with fourier convolutions. In: Proceedings of the IEEE/CVF winter conference on applications of computer vision. pp. 2149–2159 (2022)

55. Tang, L., Jia, M., Wang, Q., Phoo, C.P., Hariharan, B.: Emergent correspondence from image diffusion. arXiv preprint arXiv:2306.03881 (2023)
56. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L.u., Polosukhin, I.: Attention is all you need. In: Guyon, I., Luxburg, U.V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 30. Curran Associates, Inc. (2017), https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf
57. Vinker, Y., Voynov, A., Cohen-Or, D., Shamir, A.: Concept decomposition for visual exploration and inspiration. *ACM Transactions on Graphics (TOG)* **42**(6), 1–13 (2023)
58. Wang, C., Chai, M., He, M., Chen, D., Liao, J.: Clip-nerf: Text-and-image driven manipulation of neural radiance fields. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 3835–3844 (2022)
59. Yang, C., Lu, X., Lin, Z., Shechtman, E., Wang, O., Li, H.: High-resolution image inpainting using multi-scale neural patch synthesis. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (July 2017)
60. Yang, L., Kang, B., Huang, Z., Xu, X., Feng, J., Zhao, H.: Depth anything: Unleashing the power of large-scale unlabeled data. arXiv preprint arXiv:2401.10891 (2024)
61. Yu, L., Shi, B., Pasunuru, R., Muller, B., Golovneva, O., Wang, T., Babu, A., Tang, B., Karrer, B., Sheynin, S., et al.: Scaling autoregressive multi-modal models: Pretraining and instruction tuning. arXiv preprint arXiv:2309.02591 (2023)
62. Yu, L., Xiang, W., Han, K.: Edit-diffnerf: Editing 3d neural radiance fields using 2d diffusion model. arXiv preprint arXiv:2306.09551 (2023)
63. Yuan, Y.J., Sun, Y.T., Lai, Y.K., Ma, Y., Jia, R., Gao, L.: Nerf-editing: geometry editing of neural radiance fields. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 18353–18364 (2022)
64. Zhang, L., Rao, A., Agrawala, M.: Adding conditional control to text-to-image diffusion models. In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*. pp. 3836–3847 (2023)
65. Zheng, Y., Chen, X., Cheng, M.M., Zhou, K., Hu, S.M., Mitra, N.J.: Interactive images: Cuboid proxies for smart image manipulation. *ACM Trans. Graph.* **31**(4), 99–1 (2012)

GeoDiffuser (Supplement)

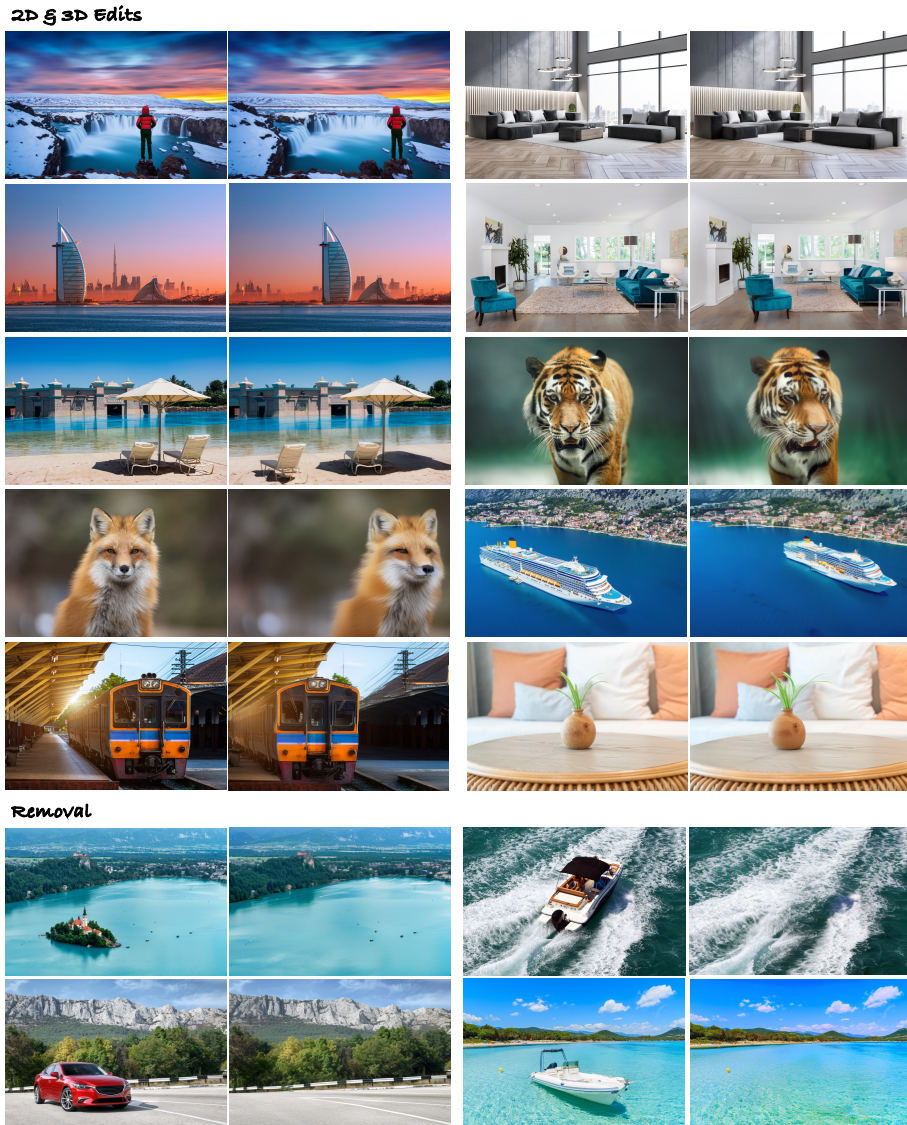


Fig. 1: We display more qualitative results of our method. Each example has the input image in the left and the result of the edit in the right.

1 Qualitative Results

We present some more qualitative results in Figure 1. We also compare our method against prior works in Figure 2 for 2D edits.

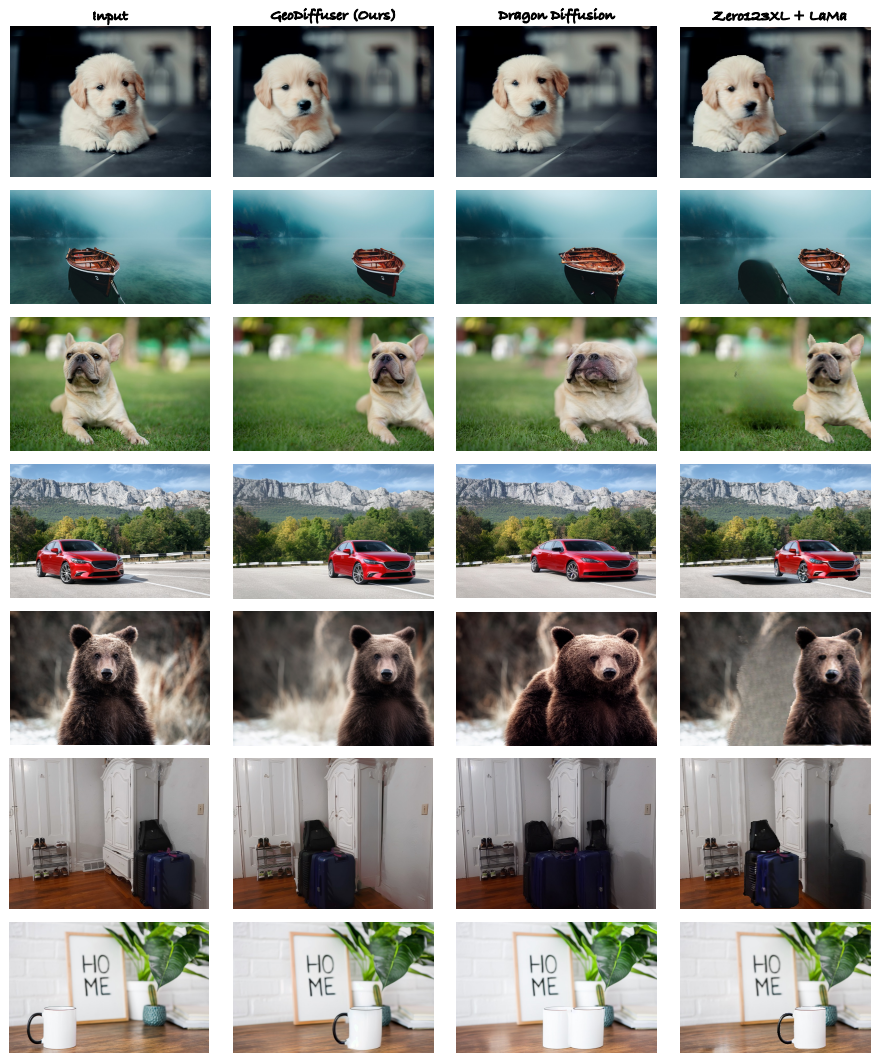


Fig. 2: We perform the same edit using prior works and compare with our work. We show 2D edits here as Dragon Diffusion can not perform 3D edits. Note that Dragon Diffusion requires prompts along with the edit and our method does not.

2 Ablations

We also compute metrics over all the edits by changing the number of shared attention steps and adaptive optimization in Table 1.

	MD ↓	Warp Error ↓	Clip Similarity ↑
Number of Timesteps for Geometric Attention Sharing			
t=30	8.893	0.0965	0.875
t=37	6.994	0.0934	0.928
GeoDiffuser (t = 45)	6.420	0.0930	0.966
Adaptive Optimization			
w/o Adaptive Optimization	10.089	0.0931	0.967
GeoDiffuser (with Adaptive Optimization)	6.420	0.0930	0.966

Table 1: Increasing the number of time steps for shared attention and adaptive optimization both improve the mean distance, warp error, and clip similarity score.

3 Failure Cases

Figure 3 displays examples where our method does not perform well.

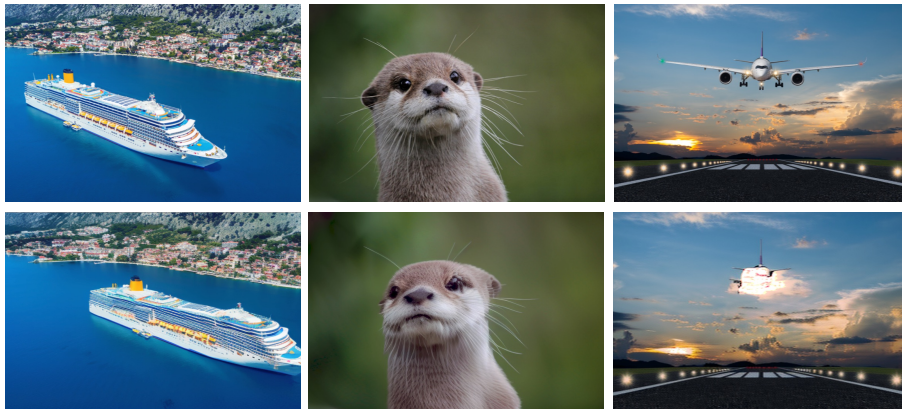


Fig. 3: Each example presents the input image at the top followed by the edited image at the bottom. As our geometric edits are performed in a lower dimensional latent space, we face aliasing and interpolation artefacts as shown in the yellow regions of the ship (left). Occasionally our optimization results in sub-optimal solutions for foreground (middle) and background dis-occlusions (right).

4 Perceptual Study

Our perceptual study was conducted using Qualtrics [5]. We first conducted a pilot study having 2 images per division type with 3 users to ensure that all

questions are clear. These users did not participate in the final study. After getting feedback from the pilot study we conducted the full study. Each participant completed the study within 10 minutes. They were allowed to click and enlarge images for better inspection. We randomized the order of options presented in the study to avoid biases.

5 Object Removal

We detail the object removal loss in Algorithm 1.

Algorithm 1 Object Removal Loss Algorithm

Require: ${}^r Q, {}^r K, {}^e Q, {}^e K$

Ensure: $\mathcal{L}_{remove} := \text{RemovalLoss}({}^r Q, {}^r K, {}^e Q, {}^e K)$

if SelfAttentionBlock **then**

${}^G A_{edit} := \text{AM}({}^e Q, {}^r K)$ ▷ Shared Attention Map

else if CrossAttentionBlock **then**

${}^G A_{edit} := \text{AM}({}^e Q, {}^e K)$ ▷ Shared Attention Map

end if

${}^G A_{ref} := \text{AM}({}^r Q, {}^r K)$

$\rho_{obj \rightarrow bg}, u_{bg} := \text{torch_max}(\text{torch_bmm}({}^G A_{edit}, {}^G A_{ref}) \odot M_{bg}, -1)$ ▷ Foreground to background correlation

$\rho_{obj \rightarrow obj}, _ := \text{torch_max}(\text{torch_bmm}({}^G A_{edit}, {}^G A_{ref}) \odot M_{obj}, -1)$ ▷ Foreground to foreground correlation

$d_{obj \rightarrow bg} := \text{NormalizedCoordinateDistance}(u_{bg})$ ▷ Coordinate distance to the background location having maximum correlation

$\mathcal{L}_{remove} := \text{mean}(e^{-d_{obj \rightarrow bg}} (\ln(\rho_{obj \rightarrow obj}) - \ln(\rho_{obj \rightarrow bg})))$

6 User Interface

See Figures 4 and 5 that display the user interface used to perform edits using GeoDiffuser. We develop this user interface using Gradio [6].

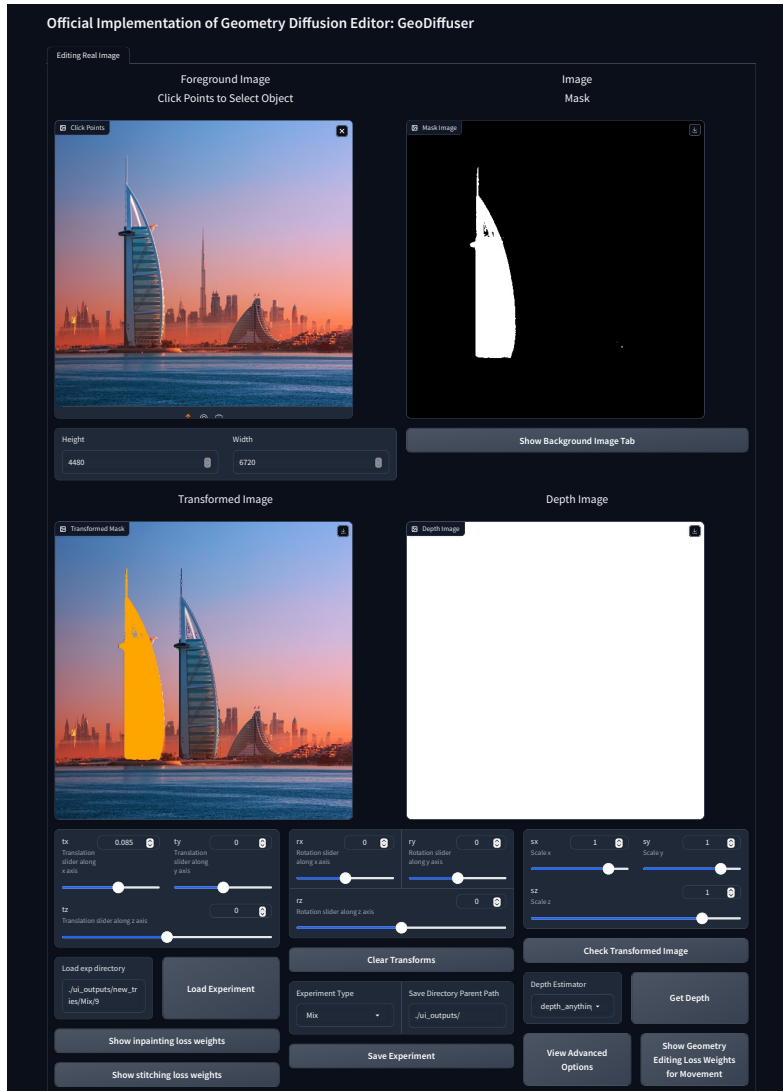


Fig. 4: GeoDiffuser UI that allows users to edit images in the wild. We provide options for users to choose a monocular depth model for geometric editing. The transformed image represents the edit that the user wishes to perform. Here, the orange mask displays the region that needs to be inpainted.

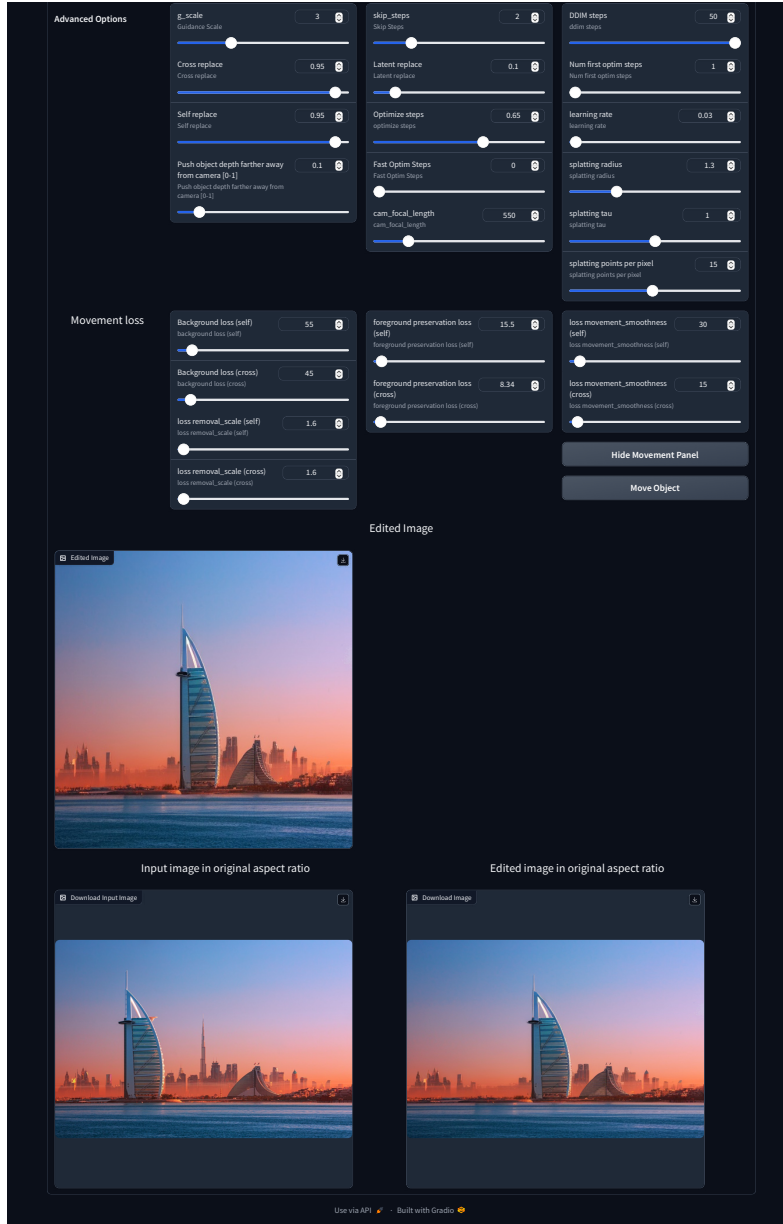


Fig. 5: GeoDiffuser UI also provides options for varying parameters for editing. The edited image in the bottom displays the image after the edit is complete.