

MeshUp: Multi-Target Mesh Deformation via Blended Score Distillation

Hyunwoo Kim
University of Chicago

Itai Lang
University of Chicago

Noam Aigerman
University of Montreal

Thibault Groueix
Adobe Research

Vladimir G. Kim
Adobe Research

Rana Hanocka
University of Chicago

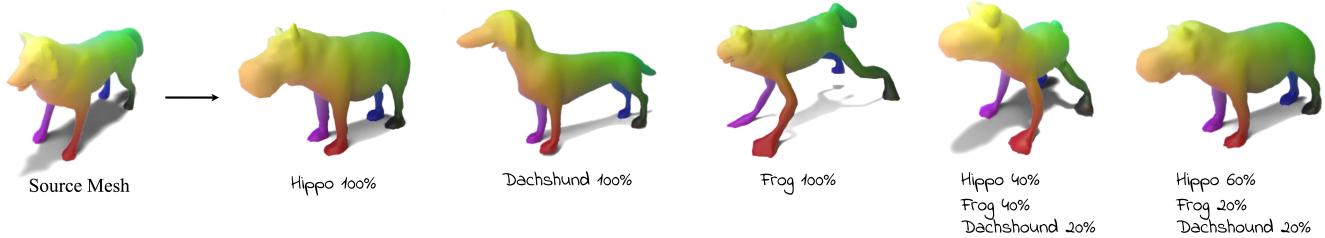


Figure 1. MeshUp is capable of deforming a source mesh into various concepts and into their weighted blends. The target objectives can be text prompts, images, or even mesh. Users can also input a set of *control vertices* to explicitly define where on the mesh the particular concepts should be expressed (Figure 8). The colors on the mesh visualize the point-wise correspondence between the source and the deformed mesh.

Abstract

We propose *MeshUp*, a technique that deforms a 3D mesh towards multiple target concepts, and intuitively controls the region where each concept is expressed. Conveniently, the concepts can be defined as either text queries, e.g., “a dog” and “a turtle,” or inspirational images, and the local regions can be selected as any number of vertices on the mesh. We can effectively control the influence of the concepts and mix them together using a novel score distillation approach, referred to as the *Blended Score Distillation* (*BSD*). *BSD* operates on each attention layer of the denoising U-Net of a diffusion model as it extracts and injects the per-objective activations into a unified denoising pipeline from which the deformation gradients are calculated. To localize the expression of these activations, we create a probabilistic Region of Interest (*ROI*) map on the surface of the mesh, and turn it into 3D-consistent masks that we use to control the expression of these activations. We demonstrate the effectiveness of *BSD* empirically and show that it can deform various meshes towards multiple objectives.

1. Introduction

Deforming mesh is a central task in geometry processing [10, 13, 38, 49, 52, 53]. In particular, it maintains valuable predefined attributes, such as artist-generated tessella-

tion, UV map, textures, and motion functions. Deforming a mesh, however, still remains a task that requires significant expertise, making it difficult for non-experts to creatively manipulate 3D models without knowing their low-level attributes. Addressing this challenge requires an intuitive, high-level control over 3D shapes in a way that can induce any non-expert users’ creative workflows. In this work, we explore the use of diffusion to enable a user-friendly deformation-based 3D content generation.

In addition to the ease of use, creative workflows in generative tasks are also inspired from their ability to synthesize novel imagery—namely, by combining a range of diverse concepts [5]. Some cognitive theories even suggest that the ability to synthesize novel combinations of known concepts and exploring these conceptual ideas is essential to human creativity [41]. While most methods that achieve 3D content generation optimize an implicit representation defined over 3D space [7, 32, 42], these representations are often inappropriate for mesh-specific tasks and cannot reuse any of the attributes defined over an artist-generated mesh. On the other hand, deformation-based approaches such as [17] lack the tools to enable a high-level, creative workflow for users to create novel conceptual imagery (e.g., “a creature with a bear’s head and a frog’s legs”, or mixing across multiple targets) and achieve precise control over their expressions.

Motivated by this observation, we propose MeshUp, a

novel approach that deforms a source mesh towards multiple target concepts defined using a variety of inputs (texts, images, and even meshes), and localizes the region where these concepts are manifested. Given as input various types of user-defined “concepts,” their respective weights, and optionally a set of vertex points on the mesh, our method deforms a mesh to appropriately conform to a localized, weighted mixture of these concepts.

In order to create a mixture of various concepts, we blend the activation maps by running the denoising process for each target and injecting the corresponding maps into a unified denoising U-Net, a method we call Blended Score Distillation (BSD). We then estimate the gradients from the diffusion using Score Distillation Sampling (SDS), a method that enables the inference step of a diffusion model to be performed in a stochastic manner [42, 56], and optimize the mesh deformation parameters, which we represent as Jacobians [1]. For fine-grained control over user-specified local regions, our framework additionally takes as input a set of selected vertices, each for a corresponding concept. Then for these concepts, we create a probability map over the mesh surface by extracting the self-attention maps from a diffusion process run on a batch of multi-viewpoint renderings, and reversely mapping them back to the surface. We then rasterize this probability map to create an attention mask that we use to control the region of deformation within our BSD pipeline (see Figure 6).

We leverage this novel pipeline to build a comprehensive creative modeling tool for concept mixing. The key features of our tool are (1) the support for mesh deformation towards multiple targets, (2) the capability to control both the strength and the region of their expression, (3) the ability to use either text, images, or other meshes as inputs.

2. Related Work

Image Editing Using Diffusion. Following the success of text-to-image generative Diffusion Models [22, 23, 31, 40, 43, 45, 51], many diffusion-based image editing models [4, 6, 9, 11, 20, 21, 29, 39, 61, 66] have been developed. These methods allow introducing custom concepts [16, 44], or enable fine-grained control of which regions and aspects of the image change [6, 11, 20] by weighting, modifying, and transferring the attention weights and activation of the diffusion networks.

Text-to-3D. These pretrained 2D diffusion generative techniques have also been used to enable 3D generation. This is usually accomplished by optimizing a 3D representation so that its rendering matches the desired text prompt [2, 7, 8, 12, 30, 32, 33, 36, 42, 46–48, 54, 57, 58, 60, 65, 67]. These methods often rely on implicit fields as a 3D representation (e.g., NeRFs [37]), which limits their editability, and often requires additional mesh conversion to support stan-

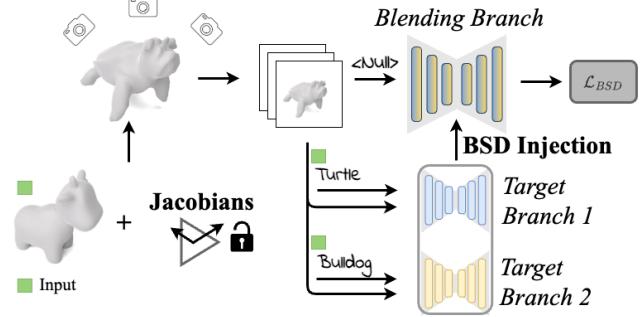


Figure 2. **Overview of Concept Blending.** MeshUp takes as input a 3D mesh and several target objectives, such as the text “Sea turtle” and “Bulldog.” We deform the source mesh by optimizing the per-triangle Jacobians of the mesh. At each iteration, we render the mesh and apply the same random noise for each target objective. Then, we pass the noised renderings and the text input through the U-Net of a pretrained text-to-image model and store the activations associated with each objective (the *Target Branch*). In the *Blending Branch*, we feed the noised rendering of the mesh to the U-Net, but condition it on the null-text embedding. We blend and inject the activations stored at each target branch into the blending branch. The gradients from the blending branch are then back-propagated via Score Distillation Sampling (SDS). After running this process iteratively, the mesh is deformed into a blend of “Sea turtle” and “Bulldog.”

dard graphics pipelines. While some techniques allow editing these implicit fields [3, 28, 46, 63], it is harder to provide local surface control, preserve correspondences (or use them to define continuous interpolations), with these models. A mesh can be extracted as a post-process [42] using marching cubes [34] and even further fine-tuned to match the desired prompt [32, 59], but these meshes would not be consistent with one another, and automatic methods do not produce artist-quality tessellations or UV mappings, necessary for a production-ready asset. In this work, we instead use deformation of a single reference shape guided by multiple concepts (e.g., textual prompts), which enables retaining necessary characteristics of the artist-created asset and enables to create a continuous semantic space interpolating between the concepts.

Mesh Deformation. Traditional mesh deformation is typically based on optimization of correspondences between vertices, faces, or other predictors that derive from these properties. [52, 53] use energy minimizing functions to give users control over the deformation space, while [15, 25] use skinning-based methods that interpolates the coordinate space with respect to the user handles. [50] uses optimal transportation to approximate correspondence across shapes. ARAP [52] and Laplacian surface editing [53] use a variational formulation to regularize the deformation in a way that preserves details and prevents drifting of the geom-

etry. However, these methods do not contain any semantics in their deformation and do not perform concept mixing. Deforming a template mesh to various concepts has been explored even before the advances in neural networks [64], but these techniques required user annotations for rigging meshes via handles and assigning semantic labels.

Several data-driven techniques have been used to learn deformations [1, 14, 18, 19, 35, 55, 62]. A recent class of works leverages text prompts as user inputs for driving a deformation towards an arbitrary textual prompt [17, 26, 36]. These methods use various deformation representations, and we opt to leverage Jacobians since they produce smooth and large-scale global deformations. We also observe that the CLIP objective lacks a full understanding of object details and that Diffusion-based objective, such as SDS [42] provides better guidance. The main goal of this work is to extend these techniques to multi-target deformation, and provide tools to mix, edit, and explore the space of concept combinations.

3. Method

The primary goal of our method is, given N sets of texts or image inputs that define the target "concepts," and their associated weights, w , to deform a source mesh into a shape that represents an effective mixture of these concepts, and control the "strength" of their expression using the weights. To that end, our method runs multiple diffusion pipelines in parallel and mixes their activation matrices within a unified pipeline to yield a single gradient direction that respects the appropriate weighted mixture of the target concepts.

For a framework that deforms a mesh into a specific target, two major design choices should be considered: the objective function (loss) and the representation of the mesh (the parameters to be optimized). For the objective function, we choose the Score Distillation Sampling (SDS) approach [42, 57], a prevalent generative technique that allows the diffusion inference process to be performed in a stochastic manner and thus enables our deformation process to be performed with viewpoint-consistency. While a straightforward application of this objective to mesh deformation would be to directly optimize the vertex positions, this method often leads to sharp artifacts [26] or restricts deformations to only local adjustments [36]. On the other hand, Jacobian-based deformation has been proposed for smooth, continuous, and global deformations, but it has only been used with an L_2 supervision [1] and CLIP similarity loss [17]. Using the SDS objective to supervise the optimization of the Jacobians offers a robust deformation framework with a powerful diffusion-based objective.

In this section, we first overview how one might approach a single target deformation using a combination of Jacobian-based mesh deformation and SDS guidance. We then extend this concept to multi-target deformation via

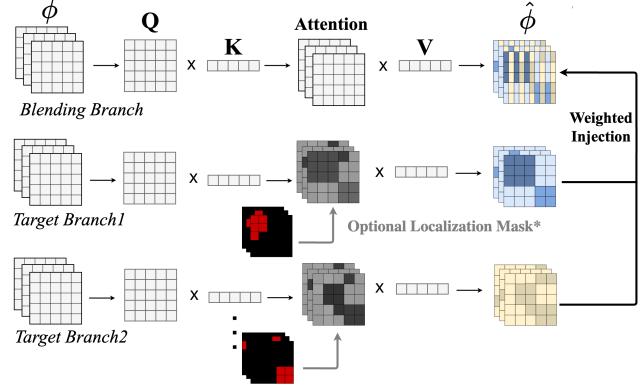


Figure 3. Overview of Blended Score Distillation (BSD). For each attention layer in the denoising U-Net, we inject the activation maps from *Target Branch1* and *Target Branch2* to the *Blending Branch* (*top*), blending the feature representations for each target. *Optional Localization Mask** (*bottom*) indicates the additional mask that we optionally apply over the cross-attention maps for localization control. The mask identifies local regions described by the selected control vertices and different weights are assigned to each of these regions. For more details, please see Figure 6 and Localization Control part of Section 3.

our novel Blended Score Distillation and explain how we achieve local control over the deformations.

Jacobian-Based Mesh Deformation. Our mesh deformation is represented by a per-face Jacobian matrix $J_i \in \mathbf{R}^{3 \times 3}$, where the deformation of a mesh (vertex positions) is computed by optimizing the following least squares problem (*i.e.*, Poisson Equation):

$$\gamma^* = \min_{\gamma} \sum_i t_i \|\nabla_i(\gamma) - J_i\|_2^2, \quad (1)$$

where γ^* is the deformation map embedding the mesh such that its Jacobians $\nabla_i(\gamma)$ are as close as possible to the target Jacobians $\{J_i\}$, the parameters we optimize, and $\{t_i\}$ are the triangle areas. Similar to previous works, we use a differentiable Poisson solver layer [1] to compute the deformation map, and a differentiable renderer [27] to connect this representation to image-based losses [17].

SDS Guidance for a Single-Target Mesh Deformation. To stochastically optimize any arbitrary parameters with respect to a pre-trained 2D diffusion model, [42] proposed the Score Distillation Sampling (SDS) process, where given a rendered image \mathbf{z} and a text condition y , the objective is to minimize the L_2 loss between a sampled noise $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ added to the image, and the noise ϵ_ω predicted by a denoising unet ω at some timestep t , sampled from a uniform distribution $t \sim U(0, 1)$:

$$\mathcal{L}_{\text{Diff}}(\omega, \mathbf{z}, y, \epsilon, t) = w(t) \|\epsilon_\omega(\mathbf{z}_t, y, t) - \epsilon\|_2^2, \quad (2)$$

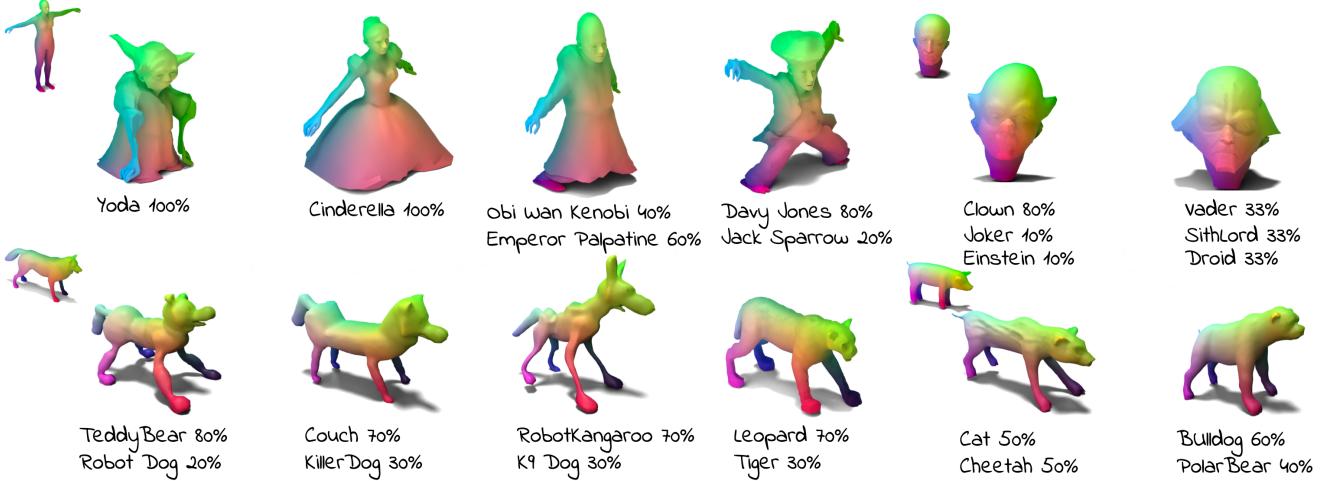


Figure 4. **Results Gallery.** We present a diverse set of 1-way, 2-way, and 3-way blending results of MeshUp. MeshUp can operate on various kinds of source shapes like human body, face, or animals, and can deform them into a blend of multiple concepts.



Figure 5. **Interpolation Between Two Objectives.** We show that we can vary the ratio between two objectives (e.g. going from hippo 100% on the left to Hippo 70%-Dachshund 30%, Hippo 30%-Dachshund 70% and finally Dachshund 100% on the right), effectively interpolating between the shape of the two targets.

where $w(t)$ is a weighting term used in the pretrained diffusion model [42], and \mathbf{z}_t is the rendered image noised at the timestep t . In practice, to compute the gradient of the optimizable parameters efficiently with respect to the loss $\mathcal{L}_{\text{Diff}}$, it has been shown that the gradients through the U-Net of the diffusion model can be omitted [42, 57]. Since we aim to minimize the loss $\mathcal{L}_{\text{Diff}}$ by optimizing each jacobian J_i , we can estimate the gradient of the loss with respect to each jacobian as follows:

$$\nabla_{J_i} \mathcal{L}_{\text{SDS}}(\omega, \mathbf{z}, y, \epsilon, t) = w(t) (\epsilon_\omega(\mathbf{z}_t, y, t) - \epsilon) \frac{\partial \mathbf{z}_t(J_i)}{\partial J_i}, \quad (3)$$

Using this SDS gradient, one can deform a mesh to a single target prompt. A detailed derivation could be found in the supplementary materials.

Following [17], we also find it beneficial to regularize

the deformation by adding a Jacobian regularization loss

$$\mathcal{L}_I = \alpha \sum_{i=1} \|J_i - I\|_2, \quad (4)$$

where α is a hyperparameter determining the regularization strength. This loss penalizes the Jacobians against the identity matrix (which represents the identity deformation) to effectively restrict the magnitude of the deformation. Next, we describe how we extend this framework to multiple targets.

Multi-target Guidance via BSD. Our multi-target architecture is composed of several parallel diffusion branches: one that takes a null text prompt as input (the blending branch), and others with a user-specified target input prompt (the target branches) (see Figure 2). These branches also take the same batch of mesh renderings as input images.

For clarity, let j denote the index for the j^{th} target-branch, each associated with a target "concept." The j^{th} branch would take as input its associated target text, y_j , and a weight w_j that controls the degree to which y_j should be expressed. The key observation is that the activation matrices (ϕ^j) we get at the end of each attention layer represent the "weighted feature space" of each concept, defined over the space of the patch of the input renderings. To blend the features over this space, we perform a weighted interpolation of the activations across the patch dimension, and inject them into the corresponding patch location in the blending branch. Formally, we inject the activation for a single concept as follows:

$$\phi^{\text{blend}} \leftarrow w_j \dot{\phi}^j + (1 - w_j) \dot{\phi}^{\text{blend}} \quad (5)$$

where ϕ^{blend} and $\dot{\phi}^j$ are the activation matrices in the blending branch and target branch j , respectively. To blend two

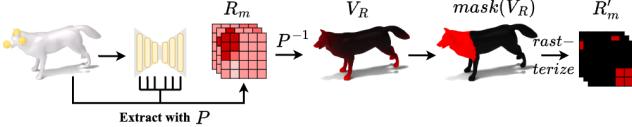


Figure 6. **Overview of Mask Extraction for Localized Control.**

For Localized control, we first extract the self-attention maps that correspond to each *control vertex*. Using the inverse vertex-to-pixel map, P^{-1} , we then project the self-attention maps onto the mesh vertices and create a 3D attention map that describes the Region of Interest (ROI) per target concept (V_R). We normalize and threshold this ROI map to create a 3D mask ($mask(\hat{V}_R)$), and rasterize the map from the same viewpoints as the mesh renderings to generate R'_m , the localization masks to be applied to the cross attention layers of the BSD pipeline.

concepts from the i^{th} and j^{th} target branch, we would inject:

$$\phi^{blend} \leftarrow w_i \phi^i + w_j \phi^j + (1 - w_i - w_j) \phi^{blend} \quad (6)$$

The denoising U-Net from the blending branch utilizes the blended activations ϕ^{blend} to predict the noise added to the image, and the gradients are backpropagated using Equation (3) to update the Jacobians.

Localized Control. Notably, the BSD pipeline is designed in a way that can incorporate a more fine-grained control over the location where each concept is manifested. Specifically, we select a set of *control vertices* as additional inputs, and impose a novel localization constraint over our concept-mixing pipeline by leveraging the self-attention maps extracted from these vertices. We will first go over how we can achieve local control for a single target, then extend this concept to enable localized blending of multiple concepts.

We first begin by mapping the 3D vertex positions to their corresponding pixel locations in a set of rendered images by using a mapping function r that takes as input v , the vertex positions, and c , the camera parameters, to find a vertex-to-pixel mapping P :

$$P = r(v, c). \quad (7)$$

Next, we perform a denoising iteration on these renderings, and using the map P , we extract all the self-attention maps corresponding to the selected control vertices. We then average these maps across the attention layers to form a probabilistic region of interest (ROI) for each rendering, which we henceforth denote as R_m (the ROI map for the m^{th} rendering). We then use the inverse pixel map, P^{-1} to map R_m back to its corresponding vertex positions on the mesh surface:

$$V_R = \sum_m P^{-1}(R_m), \quad (8)$$

where V_R is the 3D probabilistic ROI defined over the mesh vertices. We iteratively update V_R during the BSD

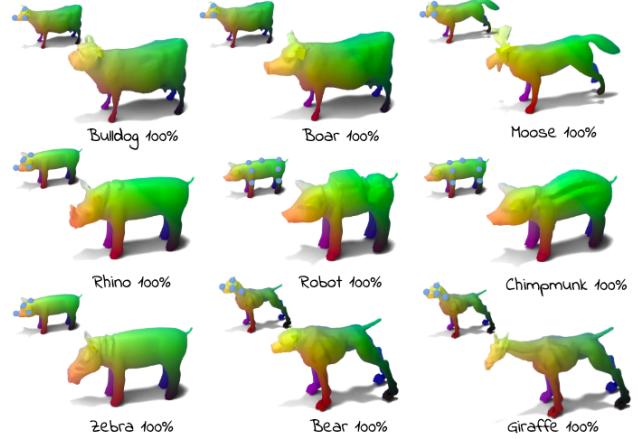


Figure 7. **Local Deformation for a Single Target** We show local deformation results for single text targets. We visualize the source mesh over the top-left corner of each result, and the selected *control vertices* as blue dots. Note how the deformation is constrained to the selected region.

optimization process. We can then get a 3D-consistent 2D ROI map, R'_m , by normalizing V_R with $\hat{V}_R = \frac{V_R - \min(V_R)}{\max(V_R) - \min(V_R)}$, thresholding it at $th = 0.8$ to create a binary mask, $mask_{\hat{V}_R}$.

For single-target deformation, we first deform the entire source mesh by regularly updating the jacobians, and at the end of all iterations, we use this binary mask to manually assign any jacobians that falls out of this mask region to the identity matrix:

$$mask_{\hat{V}_R}(J_i) = I. \quad (9)$$

By solving the poison equation 1 after such assignment, we effectively get a mesh that smoothly deforms to the target only within the region specified by the 3D consistent mask. As we visualize in Figure 7, our method’s significant capability to deform the specified region while preserving its smooth connectivity to the preserved region offers our work to be used as a geometry-editing tool, where given a pre-defined source mesh, the users can select and partially edit specific regions of the mesh using text prompts.

Localized Control for Multiple Concept Blending. To “blend” guidance from a variety of these localized objectives, we first rasterize each of the 3D-mask $mask_{\hat{V}_R}(J_i)$, back to the 2D rendering space,

$$R'_m = Rast(mask_{\hat{V}_R}, v, c). \quad (10)$$

We then use R'_m to mask-out the cross-attention maps, eliminating any association between the target and the unwanted regions of the mesh. Using BSD to mix activation from these masked attention maps yield a guidance score that respects both the weighted blend of multiple targets, as well as their associated local regions, as noticeable in Figure 8

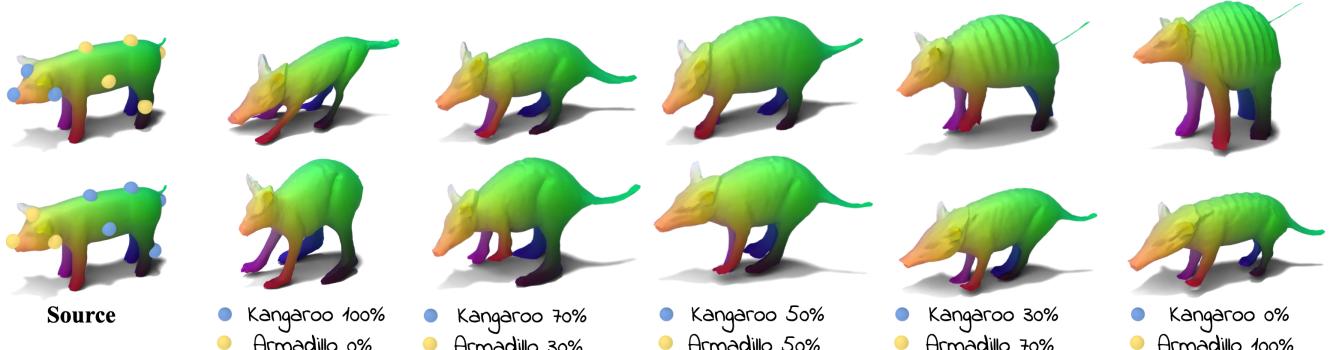


Figure 8. Multiple Local Deformation Control. We show deformation results for different selection of control points. Both columns show deformation results combined with various BSD weights of 1.0, 0.7, 0.5, and 0.3, 0.0, respectively, for the targets “Kangaroo” and “Armadillo.” (*Top row*) the points around the head region (*blue dots*) are assigned to the target “Kangaroo,” and points around the body to “Armadillo” (*yellow dots*). (*Bottom row*) flips the assignment. The figure demonstrates how the deformation results vary according to the assignment of selected control points.

Since mask R'_m is rasterized from a unified 3D ROI map V_R , it is consistent across multiple viewpoints, and thus for the various renderings. Additionally, because V_R is continuously accumulated as the sum of the attention probabilities projected from multiple R_m s, it is guaranteed that the influence from a single attention map is minimized, preventing any particular viewpoints from adding significant variations to the ROI map. We show an ablation of this method in the supplements.

The rasterized R'_m is then used as a binary mask in our usual BSD pipeline to be applied over the cross-attention maps of a desired concept, constraining the area over which the concept can be manifested. Additionally, since self-attention maps extracted from real, non-inverted renderings can be less informative, we optionally fine-tune and overfit LoRA weights to precisely predict the noise from a large batch of multi-viewpoint renderings using the objective from [44]. We supply further details about this, as well as the localization method in the supplements.

Image Targets with Textual Inversion. Text prompts might often be insufficient to describe the desired target and images could be more descriptive in some settings. We leverage textual inversion [16], which converts an image target into a prompt encoding, and use the encoded prompt in place of the target prompt y of the target branch in our BSD framework.

4. Experiments

In this section, we first show multi-target deformation results driven by text or image targets. Additionally, we demonstrate deformations with local control and mesh targets. Finally, we also describe how our method can be used as a regularization term that controls the strength of a deformation. We provide comparisons with various baselines,

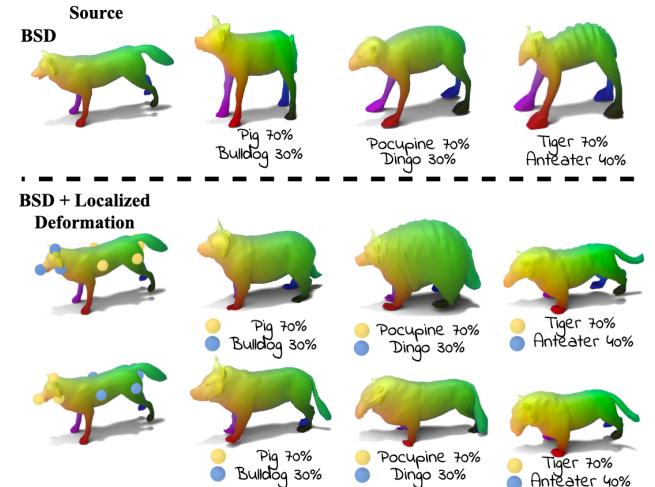


Figure 9. Evaluation of BSD with and without Local Deformation. We evaluate the deformation results with and without the Localized Deformation method. The *top row* shows results using just the naive BSD (our regular multi-target deformation), while the *middle* and *bottom row* shows the results using our localization method. We visualize the selected *control vertices* as blue and yellow dots on the mesh. Note how the results using our Localized Deformation method respect the assigned *control points*, in addition to the mixture of multiple targets.

show an experiment that uses our method to perform key-point interpolation between concepts, and show a qualitative user study of our method in the supplementary material.

4.1. Concept Mixing Results

Multi-Target Results. We demonstrate various multi-target concept mixing results in Figures 1, 4, and 5. Our method successfully mixes diverse concepts (animals, faces, fantasy creatures, and vehicles) with various weights.

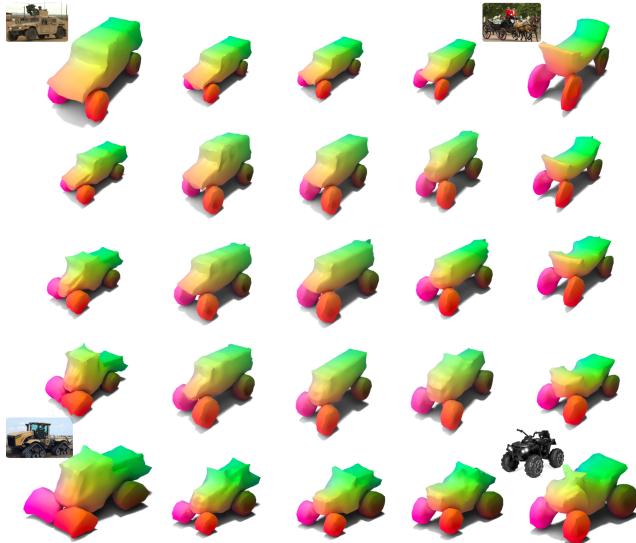


Figure 10. Four-way Blending with Image Targets. We use textual inversion [16] to condition the Text-to-image model with inspiration images represented by their inverted textual token. MeshUp supports as many targets as desired. We demonstrate a four-way blending with four target image concepts. The closer the shape is to the target image, the higher the corresponding blending weight of that target.

The figures illustrate how the same concepts can be mixed with different weights, enabling the user to control which features emerge more prominently. For example, in Figure 5, we can clearly see how with a high weight for the hippo shape, the fat body and the rounded face is prominent. On the other hand, dachshund’s long body and facial features are dominant for examples with higher weights on dachshund.

Localization Control Results. In Figure 7 and Figure 8, we provide examples of localization control, where users can indicate (by selecting the control vertices, visualized in blue and yellow dots) which part of the model should be affected by each target. Note how each of the target features emerges in the user-specified region. This method offers a high level of control over how both mixed/unmixed concepts manifest in the deformed mesh. We also demonstrate how the local deformation is affected by changes in the assignment of weight (w) in Figure 8. We observe that if we give a different emphasis to different parts of the source mesh via the selection of control vertex, BSD conditions the emphasis accordingly on various scales, creating a more versatile space for user control.

Image Targets and the Concept Space. We show the ability of our system to take image concepts as inputs in Figure 10. We find this feature to be especially useful for some concepts that have significant shape variations (e.g., “trucks”), and for those that are difficult to engineer a pre-

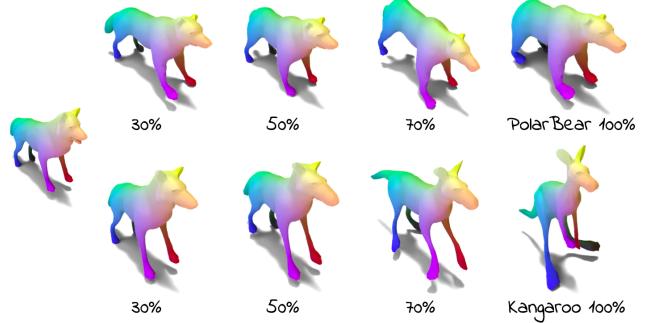


Figure 11. Self-Blending Deformations. We use BSD to inject varied strengths of activations, each from a single target to the blending branch, effectively controlling for the degrees to which the target is expressed in the resulting deformation.



Figure 12. Interpolating Mesh Using Self-Blending Deformation. We show how the Self-Blending capability of MeshUp can be used to interpolate the shapes of two meshes, the **Source** and the **Target**, by using dreambooth to learn the shape of the **Target**, and deforming the **Source** using various weights. Note how the muscular features of the **Target** mesh gradually emerge as we increase the blending weight from 30% to 70%.



Figure 13. Texture Transfer. We show how the texture map initially defined over the source mesh gets transferred without distortion to meshes deformed using our method.

cise prompt for. In this figure, we also illustrate how one can generate a continuous blending space spanning as many as four concepts by sampling different relative weights for each one of them.

Regularizing Mesh Deformation via Self-Blending. In Figure 11 we demonstrate that our BSD pipeline can take a single target objective, and be used to control the strength of a single-target deformation by using various weights, w_j . We use a modified classifier free guidance to achieve this (we provide its details and ablation in the supplements).

Using Mesh as Targets. In Figure 12, we further expand our model’s self-blending capability to interpolate the shapes of two meshes, by gradually deforming the source

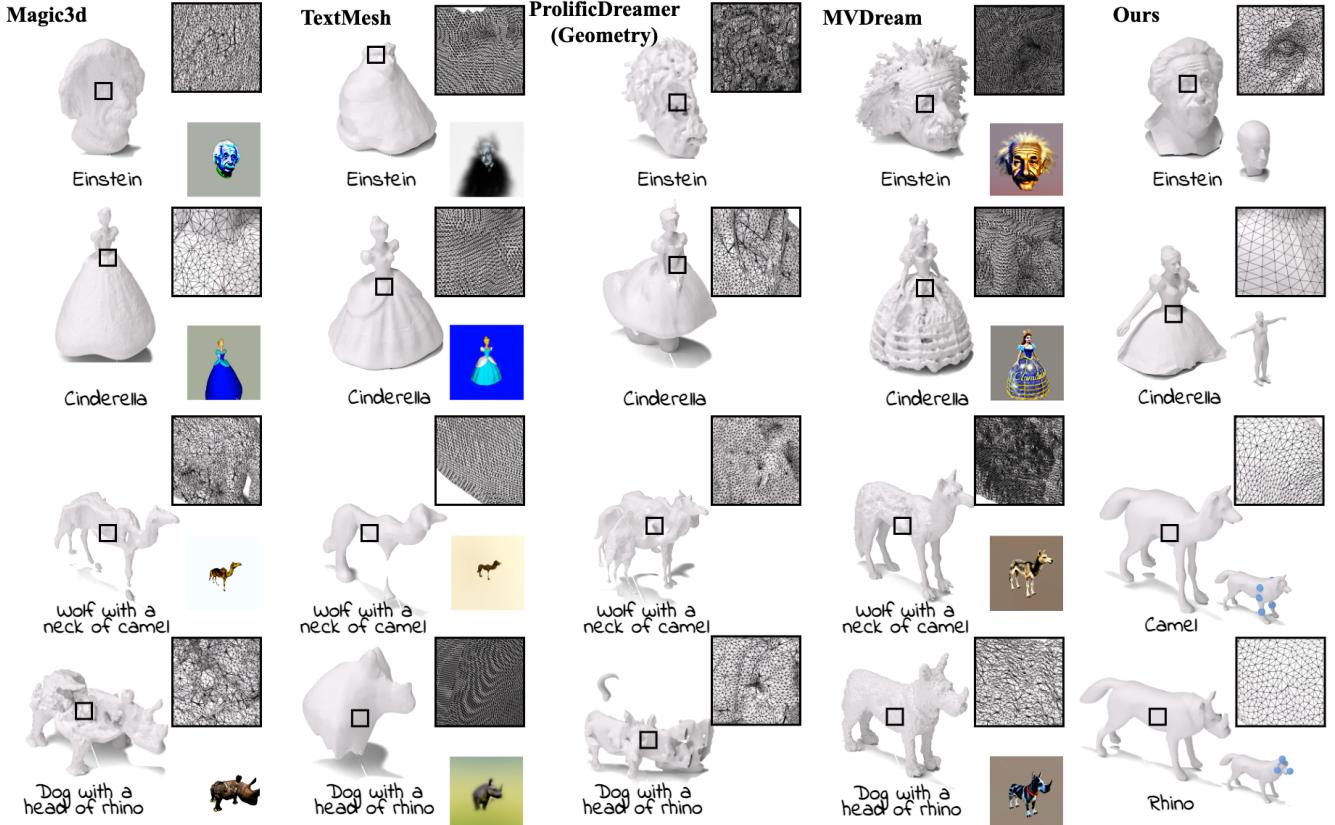


Figure 14. Comparison with other methods. We compare the mesh quality obtained with MeshUp to Magic3D, TextMesh, ProlificDreamer, and MVDream. For Magic3D, TextMesh, and MVDream, we visualize the textured, implicit shape representation on the bottom right of each figure (for ProlificDreamer, we only show the result after the geometry refinement stage). For our first two results, we deform the source mesh (visualized on the bottom right of each result) into the specified targets. For the last two results, we use our localized control method to confine the deformation to the region specified by the control vertices (visualized as blue dots over the source mesh).

mesh into a **Target** mesh. To achieve this, we utilize the multi-viewpoint renderings of the target, and batch 48 renderings per-iteration to fine-tune the UNet of the diffusion model using the objective from DreamBooth [44]. To avoid memory overload, we fine-tune the LoRA weights [24] instead of the whole model. Using the fine-tuned weights with the associated token as the objective, we deform the **Source** using various weights, w_j . Please refer to [44] for details of the training procedure.

Texture Transfer. We demonstrate the utility of deforming from a source shape using our method, as opposed to generating new 3D shapes from scratch. In Figure 13 we show how the texture map defined over the source seamlessly transfers over to other meshes deformed using our method. We can extend this property to transfer other attributes such as motion functions, and we show this example in the supplementary video.

Comparison with Other Methods. Finally, we compare the quality of our mesh outputs to those extracted from Magic3D, TextMesh, ProlificDreamer (geometry refinement stage), MVDream. Not only does our method yield a geometry of much better detail and quality, but the tessellations (visualized on the right side of each figure)

are also superior, a crucial advantage for any mesh-based graphic applications. We also show in the last two *bottom rows* that our localized control method significantly outperforms other methods that use text description to depict the localized deformation results we can achieve using MeshUp. More details of the comparison, including the specific models we used for these experiments, can be found in the supplements.

5. Conclusion and Limitations

In this paper, we propose a versatile framework for mesh deformation that supports creative workflows, enabling deformation via text or image-based concepts, mixing these concepts using various weights, and localizing their expressions.

The deformation of MeshUp is focused on preserving the topology of the initial mesh, and treating it as a shape-prior which prescribes the aesthetic of the deformed results. Thus, our technique would not be suitable for inducing topological changes (such as deforming a sphere to an object with topological holes). We leave the task of topology-modifying mesh modifications to future work. Although we limit our focus to deformation in this paper, another

potential application of our method would be to leverage our technique for generating other mesh parameters, such as textures, materials, and normals.

References

- [1] Noam Aigerman, Kunal Gupta, Vladimir G Kim, Siddhartha Chaudhuri, Jun Saito, and Thibault Groueix. Neural jacobian fields: Learning intrinsic mappings of arbitrary meshes. *SIGGRAPH*, 2022. [2](#), [3](#)
- [2] Mohammadreza Armandpour, Ali Sadeghian, Huangjie Zheng, Amir Sadeghian, and Mingyuan Zhou. Re-imagine the negative prompt algorithm: Transform 2d diffusion into 3d, alleviate janus problem and beyond, 2023. [2](#)
- [3] Chong Bao, Yinda Zhang, Bangbang Yang, Tianxing Fan, Zesong Yang, Hujun Bao, Guofeng Zhang, and Zhaopeng Cui. Sine: Semantic-driven image-based nerf editing with prior-guided editing field, 2023. [2](#)
- [4] Manuel Brack, Felix Friedrich, Dominik Hintersdorf, Lukas Struppek, Patrick Schramowski, and Kristian Kersting. Sega: Instructing text-to-image models using semantic guidance, 2023. [2](#)
- [5] Lindsay Brainard. The curious case of uncurious creation. *Inquiry*, 0(0):1–31, 2023. [1](#)
- [6] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. [2](#)
- [7] Zilong Chen, Feng Wang, and Huaping Liu. Text-to-3d using gaussian splatting, 2023. [1](#), [2](#)
- [8] Xinhua Cheng, Tianyu Yang, Jianan Wang, Yu Li, Lei Zhang, Jian Zhang, and Li Yuan. Progressive3d: Progressively local editing for text-to-3d content creation with complex semantic prompts, 2023. [2](#)
- [9] Jooyoung Choi, Yunjey Choi, Yunji Kim, Junho Kim, and Sungroh Yoon. Custom-edit: Text-guided image editing with customized diffusion models, 2023. [2](#)
- [10] Etienne Corman and Maks Ovsjanikov. Functional characterization of deformation fields. *ACM Transactions on Graphics (TOG)*, 38(1):1–19, 2019. [1](#)
- [11] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance, 2022. [2](#)
- [12] Dale Decatur, Itai Lang, Kfir Aberman, and Rana Hanocka. 3d paintbrush: Local stylization of 3d shapes with cascaded score distillation, 2023. [2](#)
- [13] Ana Dodik, Oded Stein, Vincent Sitzmann, and Justin Solomon. Variational barycentric coordinates. *ACM Transactions on Graphics (TOG)*, 42(6):1–16, 2023. [1](#)
- [14] Marvin Eisenberger, David Novotny, Gael Kerchenbaum, Patrick Labatut, Natalia Neverova, Daniel Cremers, and Andrea Vedaldi. Neuromorph: Unsupervised shape interpolation and correspondence in one go, 2021. [3](#)
- [15] Lawson Fulton, Vismay Modi, David Duvenaud, David I. W. Levin, and Alec Jacobson. Latent-space dynamics for reduced deformable simulation. *Computer Graphics Forum*, 2019. [2](#)
- [16] 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](#), [6](#), [7](#)
- [17] William Gao, Noam Aigerman, Groueix Thibault, Vladimir Kim, and Rana Hanocka. Textdeformer: Geometry manipulation using text guidance. In *ACM Transactions on Graphics (SIGGRAPH)*, 2023. [1](#), [3](#), [4](#)
- [18] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. Deep self-supervised cycle-consistent deformation for few-shot shape segmentation. *SGP*, 2019. [3](#)
- [19] Rana Hanocka, Noa Fish, Zhenhua Wang, Raja Giryes, Shachar Fleishman, and Daniel Cohen-Or. Alignet: Partial-shape agnostic alignment via unsupervised learning. *ACM Trans. Graph.*, 2018. [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](#)
- [21] Amir Hertz, Kfir Aberman, and Daniel Cohen-Or. Delta denoising score, 2023. [2](#)
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. [2](#)
- [23] Susung Hong, Gyuseong Lee, Wooseok Jang, and Seungryong Kim. Improving sample quality of diffusion models using self-attention guidance, 2023. [2](#)
- [24] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. [8](#)
- [25] Alec Jacobson, Zhigang Deng, Ladislav Kavan, and JP Lewis. Skinning: Real-time shape deformation. In *ACM SIGGRAPH 2014 Courses*, 2014. [2](#)
- [26] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. CLIP-mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 Conference Papers*. ACM, 2022. [3](#)
- [27] Samuli Laine, Janne Hellsten, Tero Karras, Yeongho Seol, Jaakko Lehtinen, and Timo Aila. Modular primitives for high-performance differentiable rendering, 2020. [3](#)
- [28] Verica Lazova, Vladimir Guzov, Kyle Olszewski, Sergey Tulyakov, and Gerard Pons-Moll. Control-nerf: Editable feature volumes for scene rendering and manipulation, 2022. [2](#)
- [29] Dongxu Li, Junnan Li, and Steven C. H. Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing, 2023. [2](#)
- [30] Weiyu Li, Rui Chen, Xuelin Chen, and Ping Tan. Sweet-dreamer: Aligning geometric priors in 2d diffusion for consistent text-to-3d, 2023. [2](#)
- [31] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. Glien: Open-set grounded text-to-image generation, 2023. [2](#)
- [32] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3d: High-resolution text-to-3d content creation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2023. [1](#), [2](#)

- [33] Ruoshi Liu, Rundi Wu, Basile Van Hoorick, Pavel Tokmakov, Sergey Zakharov, and Carl Vondrick. Zero-1-to-3: Zero-shot one image to 3d object, 2023. 2
- [34] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *SIGGRAPH*, pages 163–169. ACM, 1987. 2
- [35] Arman Maesumi, Paul Guerrero, Noam Aigerman, Vladimir G. Kim, Matthew Fisher, Siddhartha Chaudhuri, and Daniel Ritchie. Explorable mesh deformation subspaces from unstructured 3d generative models. *SIGGRAPH Asia (Conference track)*, 2023. 3
- [36] Oscar Michel, Roi Bar-On, Richard Liu, Sagie Benaim, and Rana Hanocka. Text2mesh: Text-driven neural stylization for meshes. *arXiv preprint arXiv:2112.03221*, 2021. 2, 3
- [37] Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, and Ren Ng. Nerf: Representing scenes as neural radiance fields for view synthesis. In *ECCV*, 2020. 2
- [38] Niloy J Mitra, Simon Flöry, Maks Ovsjanikov, Natasha Gelfand, Leonidas J Guibas, and Helmut Pottmann. Dynamic geometry registration. In *Symposium on geometry processing*, pages 173–182, 2007. 1
- [39] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models, 2022. 2
- [40] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models, 2023. 2
- [41] Elliot Samuel Paul and Scott Barry Kaufman, editors. *The Philosophy of Creativity*. Oxford University Press, New York, 2014. 1
- [42] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv*, 2022. 1, 2, 3, 4
- [43] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 2
- [44] 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, 2023. 2, 6, 8
- [45] Patrick Schramowski, Manuel Brack, Björn Deiseroth, and Kristian Kersting. Safe latent diffusion: Mitigating inappropriate degeneration in diffusion models, 2023. 2
- [46] Etai Sella, Gal Fiebelman, Peter Hedman, and Hadar Averbuch-Elor. Vox-e: Text-guided voxel editing of 3d objects, 2023. 2
- [47] Junyoung Seo, Wooseok Jang, Min-Seop Kwak, Jaehoon Ko, Hyeonsu Kim, Junho Kim, Jin-Hwa Kim, Jiyoung Lee, and Seungryong Kim. Let 2d diffusion model know 3d-consistency for robust text-to-3d generation, 2023.
- [48] Yichun Shi, Peng Wang, Jianglong Ye, Mai Long, Kejie Li, and Xiao Yang. Mvdream: Multi-view diffusion for 3d generation, 2024. 2
- [49] Justin Solomon, Mirela Ben-Chen, Adrian Butscher, and Leonidas Guibas. As-killing-as-possible vector fields for planar deformation. In *Computer Graphics Forum*, pages 1543–1552. Wiley Online Library, 2011. 1
- [50] Justin Solomon, Fernando de Goes, Gabriel Peyré, Marco Cuturi, Adrian Butscher, Andy Nguyen, Tao Du, and Leonidas Guibas. Convolutional wasserstein distances: efficient optimal transportation on geometric domains. *ACM Trans. Graph.*, 34(4), 2015. 2
- [51] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models, 2022. 2
- [52] Olga Sorkine and Marc Alexa. As-rigid-as-possible surface modeling. In *Symposium on Geometry processing*, pages 109–116. Citeseer, 2007. 1, 2
- [53] Olga Sorkine, Daniel Cohen-Or, Yaron Lipman, Marc Alexa, Christian Rössl, and H-P Seidel. Laplacian surface editing. In *Proceedings of the 2004 Eurographics/ACM SIGGRAPH symposium on Geometry processing*, pages 175–184, 2004. 1, 2
- [54] Junshu Tang, Tengfei Wang, Bo Zhang, Ting Zhang, Ran Yi, Lizhuang Ma, and Dong Chen. Make-it-3d: High-fidelity 3d creation from a single image with diffusion prior, 2023. 2
- [55] Mikaela Angelina Uy, Vladimir G. Kim, Minhyuk Sung, Noam Aigerman, Siddhartha Chaudhuri, and Leonidas Guibas. Joint learning of 3d shape retrieval and deformation. *CVPR*, 2021. 3
- [56] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation, 2022. 2
- [57] Haochen Wang, Xiaodan Du, Jiahao Li, Raymond A. Yeh, and Greg Shakhnarovich. Score jacobian chaining: Lifting pretrained 2d diffusion models for 3d generation. *arXiv preprint arXiv:2212.00774*, 2022. 2, 3, 4
- [58] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation, 2023. 2
- [59] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *arXiv preprint arXiv:2305.16213*, 2023. 2
- [60] Menghua Wu, Hao Zhu, Linjia Huang, Yiyu Zhuang, Yuanxun Lu, and Xun Cao. High-fidelity 3d face generation from natural language descriptions, 2023. 2
- [61] Binxin Yang, Shuyang Gu, Bo Zhang, Ting Zhang, Xuejin Chen, Xiaoyan Sun, Dong Chen, and Fang Wen. Paint by example: Exemplar-based image editing with diffusion models, 2022. 2
- [62] Wang Yifan, Noam Aigerman, Vladimir G Kim, Siddhartha Chaudhuri, and Olga Sorkine-Hornung. Neural cages for detail-preserving 3d deformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 75–83, 2020. 3
- [63] Yu-Jie Yuan, Yang-Tian Sun, Yu-Kun Lai, Yuewen Ma, Rongfei Jia, and Lin Gao. Nerf-editing: Geometry editing of neural radiance fields, 2022. 2
- [64] Mehmet Yumer, Siddhartha Chaudhuri, Jessica Hodgins, and Levent Kara. Semantic shape editing using deformation handles. *ACM Transactions on Graphics*, 34:86:1–86:12, 2015. 3

- [65] Jingbo Zhang, Xiaoyu Li, Ziyu Wan, Can Wang, and Jing Liao. Text2nerf: Text-driven 3d scene generation with neural radiance fields, 2023. [2](#)
- [66] Lvmnin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. [2](#)
- [67] Minda Zhao, Chaoyi Zhao, Xinyue Liang, Lincheng Li, Zeng Zhao, Zhipeng Hu, Changjie Fan, and Xin Yu. Efficientdreamer: High-fidelity and robust 3d creation via orthogonal-view diffusion prior, 2023. [2](#)