

# UDiFF: Generating Conditional Unsigned Distance Fields with Optimal Wavelet Diffusion

Junsheng Zhou<sup>1\*</sup>, Weiqi Zhang<sup>1\*</sup>, Baorui Ma<sup>1,2†</sup>, Kanle Shi<sup>3</sup>, Yu-Shen Liu<sup>1†</sup>, Zhizhong Han<sup>4</sup>

School of Software, Tsinghua University, Beijing, China<sup>1</sup>

BAAI, Beijing, China<sup>2</sup>, Kuaishou Technology, Beijing, China<sup>3</sup>

Department of Computer Science, Wayne State University, Detroit, USA<sup>4</sup>

zhoujs21@mails.tsinghua.edu.cn, zwq23@mails.tsinghua.edu.cn, brma@baai.ac.cn

shikanle@kuaishou.com, liuyushen@tsinghua.edu.cn, h312h@wayne.edu



Figure 1. Diverse shapes with and without open surfaces generated by our UDiFF model. **Top-Left:** Conditional generation of clothes with prompts ‘A short-sleeved dress in spiderman style’, ‘A Batman upper with long sleeves’, ‘A superman pant’, ‘A camouflage slip dress’. **Around:** A shape gallery generated by UDiFF conditionally and unconditionally.

## Abstract

Diffusion models have shown remarkable results for image generation, editing and inpainting. Recent works explore diffusion models for 3D shape generation with neural implicit functions, i.e., signed distance function and occupancy function. However, they are limited to shapes with closed surfaces, which prevents them from generating diverse 3D real-world contents containing open surfaces. In this work, we present UDiFF, a 3D diffusion model for un-

signed distance fields (UDFs) which is capable to generate textured 3D shapes with open surfaces from text conditions or unconditionally. Our key idea is to generate UDFs in spatial-frequency domain with an optimal wavelet transformation, which produces a compact representation space for UDF generation. Specifically, instead of selecting an appropriate wavelet transformation which requires expensive manual efforts and still leads to large information loss, we propose a data-driven approach to learn the optimal wavelet transformation for UDFs. We evaluate UDiFF to show our advantages by numerical and visual comparisons with the latest methods on widely used benchmarks. Page: <https://weiqi-zhang.github.io/UDiFF>.

\*Equal contribution. † Corresponding authors. This work was supported by National Key R&D Program of China (2022YFC3800600), the National Natural Science Foundation of China (62272263, 62072268), and in part by Tsinghua-Kuaishou Institute of Future Media Data.

## 1. Introduction

Probabilistic diffusion models [17, 60] have largely revolutionized 2D content generation. Recent advancements, such as DALL-E 2 [55] and Stable Diffusion [56], have been widely used in text-to-image generation, image inpainting, etc. A series of works [37, 59] try to replicate these success in 3D content generation by developing diffusion models for point clouds or voxels, but fails to produce high fidelity results due to the limited resolution in voxels and the discreteness of points. Recent approaches [10, 21, 69] explore diffusion models to generate 3D shapes as neural implicit functions, i.e., signed distance function (SDF) [38, 52] and occupancy function (Occ) [45]. However, they are limited to generate closed shapes since both SDF and Occ model the internal and external relations of 3D locations for representing 3D shapes. This prevents previous 3D implicit diffusion models from generating diverse 3D real-world contents containing open surfaces.

Another challenge in diffusion-based 3D generative models is how to define a compressing transform schema for achieving compact implicit representations which can be learned by diffusion models efficiently. Some works train a variational auto-encoder (VAE) [26] for converting shapes into triplane [15, 57] or single latents [49] for latent diffusion. However, the relative limited 3D data makes it difficult to train a stable VAE. Instead, another series of works (e.g. WaveGen [21]) seek to leverage explicit transform in another domain (e.g. wavelet transform [11]) for direct compression. Nevertheless, they need to select an appropriate wavelet type, which demands extensive manual efforts and can still result in significant information loss during wavelet recovery.

To address these issues, we propose UDiFF, a 3D diffusion model for unsigned distance fields [8, 72] which is capable of generating textured 3D shapes without geometric limits on the surface watertightness (e.g. contain open surfaces). Compared to commonly-used SDF or Occ, UDF has proven to be an advanced representation that supports arbitrary typologies and remain strong generalization. Going beyond pure unconditioned models, we incorporate conditions achieved from CLIP [54] models to UDiFF by introducing conditional cross-attentions. This enables to control 3D generation using the text and image signals. Previous works merely focus on generating geometries which lead to a lack of appearance and prevent them from creating diverse and visual-appealing 3D models, while we get inspiration from Text2Tex [4] to simultaneously generate textures for universal 3D content creation.

Adapting existing SDF-based diffusion models directly to UDF does not work well. The difficulty arises from the significantly greater complexity of UDF compared to SDF, particularly in the context of the non-differential zero-level set. To solve this issue, we introduce UDiFF as a diffusion

model in the spatial-frequency domain based on an optimal wavelet transformation, which produces a compact representation space for UDF generation. Instead of engaging in selecting a suitable wavelet transformation, which is tedious and often results in significant information loss, we employ a data-driven approach to obtain an optimal wavelet filter for representing UDFs. We minimize the unsigned distance errors during a self-reconstruction through the wavelet transformation, especially near the zero-level set of UDFs. This preserves the geometry details during wavelet transformation, which leads to the high-fidelity generation of 3D geometries. We evaluate UDiFF for generating 3D shapes with open surfaces and closed surfaces using conditions or unconditionally under DeepFashion3D [79] and ShapeNet [3] datasets. The experimental results demonstrate that UDiFF achieves promising generation performance compared to the existing state-of-the-art approaches, in both qualitative and quantitative evaluations. Our main contributions can be summarized as follows.

- We propose UDiFF, a 3D diffusion model for unsigned distance fields which is capable of generating real-world textured 3D shapes with open surfaces from text conditions or unconditionally.
- We introduce an optimal wavelet transformation for UDF through data-driven optimization, and justify that the spatial-frequency domain learned through this transformation is a compact domain suitable for UDF generation.
- We evaluate UDiFF for generating 3D shapes with both open and closed surfaces, and show our superiority over the state-of-the-art methods.

## 2. Related Work

With the rapid development of deep learning, the neural networks have shown great potential in 3D applications [19, 24, 32, 33, 41, 43, 63, 64, 66, 70, 73, 75–77]. We mainly focus on learning generative Neural Implicit Functions with networks for generating 3D shapes.

### 2.1. Neural Implicit Representations

Recently, Neural Implicit Functions (NIFs) have shown promising results in surface reconstruction [42, 45, 52], novel view synthesis [46, 48], image super-resolution [1, 58], etc. The NIFs approaches train a neural network to represent shapes and scenes with signed distance functions (SDFs) [9, 52] or binary occupancy [45, 53], where the marching cubes algorithm [36] is then used to extract surfaces from the learned NIFs. OccNet and DeepSDF [45, 52] are the pioneers of NIFs which learn global latent codes for representing 3D shapes with MLP-based decoder to achieve occupancies or signed distances. The subsequent approaches [23, 53] leverage more latent codes to represent detailed local geometries. PCP [40] and OnSurf [39]

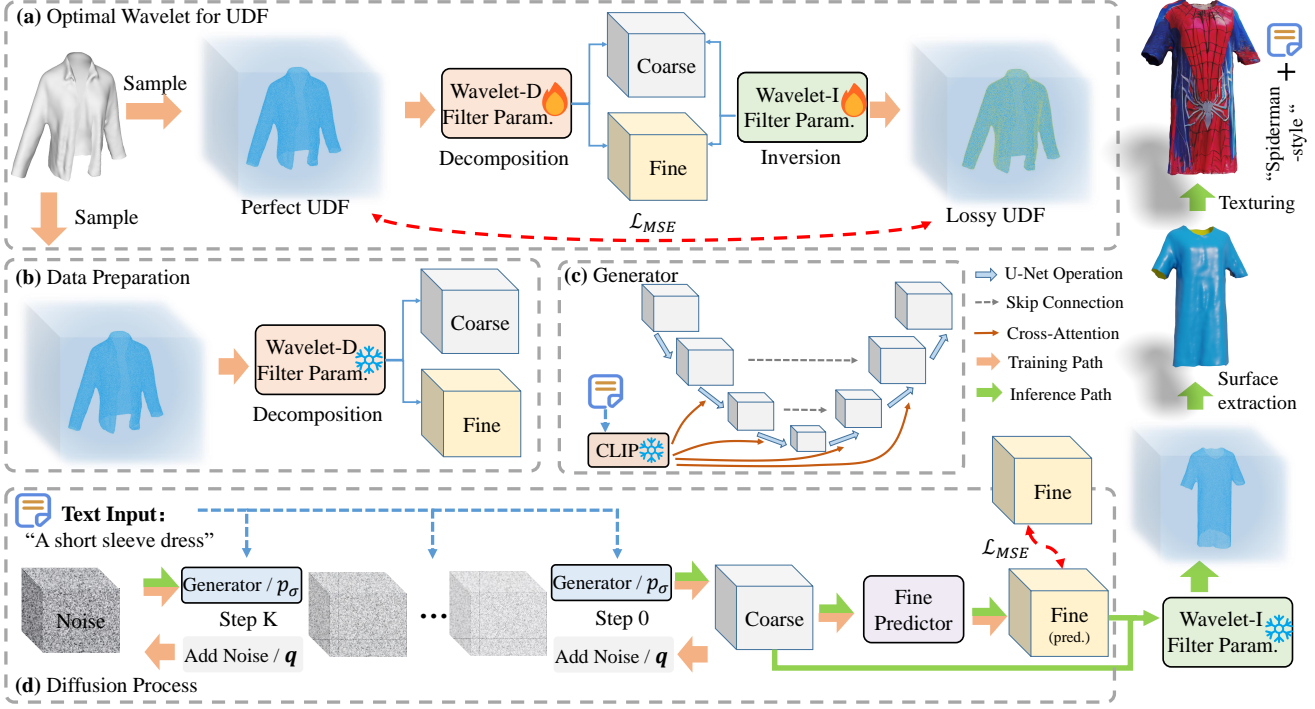


Figure 2. **Overview of UDiFF.** (a) We propose a data-driven approach to attain the optimal wavelet transformation for UDF generation. We optimize wavelet filter parameters through the decomposition and inversion by minimizing errors in UDF self-reconstruction. (b) We fix the learned decomposition wavelet parameters and leverage it to prepare the data as a compact representation of UDFs including pairs of coarse and fine coefficient volumes. (c) is the architecture of the generator in diffusion models, where text conditions are introduced with cross-attentions. (d) The diffusion process of UDiFF. We train the generator to produce coarse coefficient volumes from random noises guided by input texts and train the fine predictor to predict fine coefficient volumes from the coarse ones. Follow the **green** arrows for inference, we start from a random noise and an input text to leverage the trained generator to produce a coarse coefficient volume. The trained fine predictor then predicts the fine coefficient volume. Together with the coarse one, we recover the UDFs with the fixed pre-optimized inversion wavelet filter parameters. Finally, we extract surfaces from UDFs and further texture them with the guiding text.

introduce predictive context priors and on-surface prior to enhance the representation ability of NIFs.

Occupancy and SDFs are mainly suitable to represent closed shapes. Recent works explore the neural unsigned distances (UDFs) [5, 8, 34, 35, 62, 72, 74] to represent shapes and scenes with open surfaces. NDF [8] designs a hierarchical neural network to learn UDFs with ground truth distance supervisions. GIFS [68] learns UDFs and represents shapes with query relationships. CAP-UDF [72] and LevelSetUDF [74] propose consistency-aware constraints and level set projections to stabilize the optimization of UDFs and produce more accurate geometries.

## 2.2. Diffusion-based 3D Generative Models

Generating 3D contents plays the key role in augmented/virtual reality and has been widely explored in the past few years. Earlier works transfer the success of GAN [13], VAE [26] and the flow-based model [25] in image generation to the 3D domain for generating 3D shapes represented as point clouds [2, 22, 31, 50, 67] and voxels [59, 65]. PointDiff [37] introduces the powerful diffusion models for

point cloud generation. Some advanced works [61, 78] combining the voxel and point representations were proposed for more robust 3D generation with diffusion models.

More recently, some approaches [7, 12, 21, 30, 57] try to combine the diffusion models and neural implicit representations for generating high-quality 3D shapes. These methods generate signed distance fields [10, 15, 21, 28, 30] or occupancy fields [69] with diffusion models and extract the meshes from the fields with the marching cubes [36]. For the efficient training of diffusion models, methods like Diffusion-SDF [10] and 3D-LDM [49] train a VAE for converting shapes into latent codes for latent diffusion. But the relative small number of 3D samples for training makes it difficult to train a stable VAE. WaveGen [21] was proposed to explicitly compress SDFs in frequency domain with wavelet transform, but it is limited to the information loss during the wavelet recovery.

The advances in NIFs-based 3D generative models have shown significant improvements in the generation qualities, however, they are limited to generate closed surfaces. This prevents them from generating diverse 3D contents in real

world. In this work, we focus on generating UDFs for open surfaces with textures using a 3D diffusion model.

### 3. Method

**Overview.** The overview of UDiFF is shown in Fig. 2. UDiFF is a 3D generative model which takes texts as conditions and generates general textured 3D shapes with either open or closed surfaces. We will start by introducing the novel approach to obtain an optimal wavelet transform for a compact UDF representation and the data preparation process for training diffusion models in Sec. 3.1. We then present the designed conditional diffusion framework for UDF generation and the generator network in Sec. 3.2. Finally, we extract surfaces from the generated UDF and further add textures on the mesh with the guiding text in Sec. 3.3.

#### 3.1. Optimal Wavelet Transformation for UDFs

One main challenge in diffusion-based 3D generative models is to search for a compact representation space for diffusion model to learn efficiently. WaveGen [21] adopts an explicit wavelet transform on the SDF volumes ( $256^3$ ) to decompose them into coarse coefficient volumes and fine coefficient volumes with much lower resolutions. The naive wavelet transform leads to large information loss since the manually selected wavelet is not capable of representing various shapes as accurate distance functions.

To represent UDFs in a compact way, we follow WaveGen to adopt multi-scale wavelet transform [11, 44] as the compressing schema, keeping only the coefficients at a relative small scale of  $\mathcal{J} = 3$  for efficient shape learning. However, the UDF is significantly more complex and unstable than SDF, particularly in the area of non-differential zero-level sets, where the geometry details that the wavelet compressing does not preserve will severely affect the generation of UDFs. Thus, a suitable wavelet filter with much less information loss but remains compact and efficient for UDFs is vital.

To this end, instead of manually searching for the appropriate wavelet filter which demands costly efforts and is still hard to reduce the information loss, we propose a data-driven approach to learn the optimal wavelet filter parameters for UDFs through learning-based optimization as shown in Fig. 2(a). Specifically, we define a learnable biorthogonal wavelet filter which consists of a decomposition filter  $\phi_\theta^D$  and an inversion filter  $\phi_\delta^I$  with learnable filter parameters  $\theta$  and  $\delta$ . Given a set of shapes  $\{S_i\}_{i=1}^N$ , we first sample the UDF volume  $U_i$  for each shape at a resolution of  $256^3$  and truncate the distance values in  $U_i$  to  $[0, 0.1]$ , and then compress it into a coarse coefficient volume and a fine coefficient volume with the learnable decomposition filter  $\phi_\theta^D$  as:

$$\{C_i, F_i\} = \phi_\theta^D(U_i). \quad (1)$$

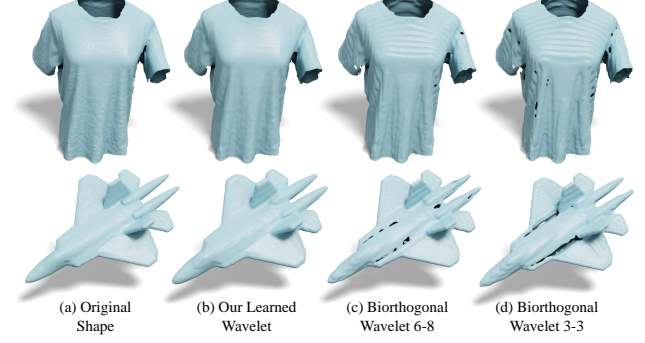


Figure 3. **Comparisons of reconstructions with different wavelet filters.** (a) The input shapes from DeepFashion3D [79] and ShapeNet [3], from where we sample UDFs to prepare compact wavelet representations. (b) The surfaces extracted from the recovered UDF with decomposition and inversion by our learned wavelet filter. (c,d) The surfaces extracted from the recovered UDF with manual chosen wavelet filters.

We then predict the lossy UDF  $\hat{U}$  from  $C_i$  and  $F_i$  with the learnable inversion filter  $\phi_\delta^I$  as:

$$\{\hat{U}_i\} = \phi_\delta^I(C_i, F_i). \quad (2)$$

The target is to optimize the filter parameters  $\theta$  and  $\delta$  by minimizing the information loss during wavelet decomposition and inversion, formulated as:

$$\min_{\theta, \delta} \sum_{i=1}^N \mathcal{L}_{\text{MSE}}(w_i^\gamma \hat{U}_i, w_i^\gamma U_i). \quad (3)$$

where  $w_i^\gamma$  is the weights for enforcing the optimization to focus on the space near the zero-level set of UDF.  $w_i^\gamma$  has the same size as  $U_i$  for weighting each grid in the UDF volume, where we define  $w_i^\gamma$  according to a threshold  $\gamma$  to mask the grids with distances larger than  $\gamma$ .

After data-driven optimization of the wavelet filters  $\phi_\theta^D$  and  $\phi_\delta^I$ , we learn the optimal wavelet transform with much less information loss and can faithfully reconstruct the original UDF while remains compact. We show the comparison on the wavelet filters in Fig. 3, where the surfaces reconstructed from UDF with our learned wavelet filter in Fig. 3(b) are much smoother and more accurate than the reconstructions with common filters like Biorthogonal wavelet 3-3 in Fig. 3(d). Specifically, Biorthogonal wavelet 6-8 in Fig. 3(c) is the carefully chosen filter by WaveGen from a series of wavelet filters, where our learned filter significantly outperforms the manually selected filters in compressing and recovering UDF. The reason is that the filters learned by data-driven optimizing from UDF datasets are much more suitable to specific characters of UDFs, which preserves more geometry details.

With the learned optimal wavelet filter, we then leverage it to represent UDFs as a compact representation for train-



ing diffusion models. As shown in Fig. 2 (b), we fix the parameters for  $\phi_\theta^D$  and produce the paired coarse efficient  $\{C_i\}_{i=1}^M$  volumes and fine efficient volumes  $\{F_i\}_{i=1}^M$  by decomposing  $U_i$  with Eq. (1).

### 3.2. Conditional UDF Diffusion

**Generator Architecture.** We first introduce the network details of diffusion generators for 3D volumes, as shown in Fig. 2(c). The generator shares a similar U-Net architecture as Stable-Diffusion [51, 56], where the 2D convolutions are replaced with 3D ones for handling 3D volumes. Each U-Net operation in Fig. 2(c) contains  $3 \times 3 \times 3$  residual blocks, pooling layers and down/up-sampling layers. For introducing text conditions to diffusion models, we first encode the input texts with frozen CLIP [54] models to produce text embeddings and then fuse them into the volume features with cross-attention layers.

**Learning Diffusion Models.** We develop our 3D generative model UDiFF based on diffusion probabilistic models [17, 60]. The diffusion process is to generate coarse coefficient volumes which represents the general geometry of 3D shapes from random volume noises, as shown in Fig. 2 (d). We define  $\{C_0, C_1, \dots, C_T\}$  as the forward process  $q(C_{0:T})$  which gradually transforms a real data  $C_0$  into Gaussian noise ( $C_T$ ) by adding noises, where  $C_0$  is a sample from the coarse coefficient data  $\{C_i\}_{i=1}^M$ . The diffusion backward process  $p_\sigma(C_{0:T})$  leverages the generator with parameter  $\sigma$  to denoise  $C_T$  into a real data sample. The learning schema is to train the generator to maximize the generation probability of the target, i.e.  $p_\sigma(C_0)$ . We follow DDPM [17] to simplify the optimization target to predict noises  $\epsilon_\sigma$  with the generator, formulated as:

$$\min_{\sigma} \mathbb{E}_{C_0, t, \epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_\sigma(C_t, t)\|^2 \right], \quad (4)$$

where  $t$  is a time step and  $\epsilon$  is a noise volume sampled from the unit Gaussian distribution  $\mathcal{N}$ .

**Condition-Guided 3D Diffusion.** Up to this point, we have covered the generative diffusion process without conditions. For a controllable generation of unsigned distance fields, we further introduce a conditioning mechanism [56] into the diffusion process by cross-attention. Specifically, given an input text  $y$ , we first leverage a frozen CLIP text encoder  $\tau$  to project  $y$  into the condition embedding  $\tau(y)$ . The embedding is then fused into the U-Net layers of generator with cross attention modules implemented as  $\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V$ , where  $Q = W_Q^{(i)} \cdot \varphi_i(C_t)$ ,  $K = W_K^{(i)} \cdot \tau(y)$  and  $V = W_V^{(i)} \cdot \tau(y)$ . Here,  $\varphi_i(C_t)$  is the output of an intermediate layer of the U-Net and  $W_Q^{(i)}$ ,  $W_K^{(i)}$  &  $W_V^{(i)}$  are learnable matrices.

The cross-attention mechanism learns a mapping from the input text condition to the coefficient volumes which

Table 1. **Quantitative comparison of shape generation under DeepFashion3D dataset.** MMD-CD scores and MMD-EMD scores are scaled by  $10^3$  and  $10^2$ , respectively.

| Method             | COV $\uparrow$ |              | MMD $\downarrow$ |              | 1-NNA $\downarrow$ |              |
|--------------------|----------------|--------------|------------------|--------------|--------------------|--------------|
|                    | CD             | EMD          | CD               | EMD          | CD                 | EMD          |
| PointDiff [37]     | 68.67          | 64.56        | <b>11.01</b>     | 15.53        | 83.21              | 87.69        |
| WaveGen [21]       | 62.34          | 51.89        | 15.56            | 17.03        | 92.93              | 94.83        |
| Diffusion-SDF [10] | 67.09          | 62.03        | 14.79            | 16.63        | 88.98              | 92.63        |
| LAS-Diffusion [71] | 67.40          | 56.01        | 14.59            | 16.53        | 88.61              | 91.41        |
| Ours               | <b>69.62</b>   | <b>67.72</b> | 11.60            | <b>14.01</b> | <b>81.83</b>       | <b>82.14</b> |

represent the geometric generations. The optimizing target in Eq. (4) is then modified as:

$$\min_{\sigma} \mathbb{E}_{C_0, y, t, \epsilon \sim \mathcal{N}(0,1)} \left[ \|\epsilon - \epsilon_\sigma(C_t, \tau(y), t)\|^2 \right], \quad (5)$$

**Fine Predictor.** The last module for learning to generate UDFs is the fine predictor  $f$  which predicts fine coefficient volumes from the generated coarse ones. We follow WaveGen [21] to implement  $f$  with the similar U-Net architecture as the generator. We train  $f$  with pairs of coarse and fine coefficient volumes  $\{C_i, F_i\}$  with MSE loss to minimize the differences between  $F_i$  and the prediction  $f(C_i)$ .

### 3.3. Generating Novel 3D Shapes

**Generating UDFs at Inference.** With the learned optimal wavelet filters and the trained conditional diffusion models, we can now generate novel 3D shapes as shown in the green arrows in Fig. 2. Starting from a random volume noise and an input text  $y$ , we leverage the generator to produce a coarse coefficient  $C'$  volume by removing noises iteratively with the guidance of  $y$ . The fine predictor then predicts the fine coefficient volume  $F'$ , together with  $C'$  to generate the UDF  $U'$  by wavelet inversion with the learned filter  $\phi_\delta^I$  as Eq. (2).

**Surface Extraction and Texturing** After generating a novel UDF  $U'$ , we extract the zero-level set of  $U'$  as a surface. The recent works [14, 72] leverage the gradients at UDF as the signals to mesh UDFs, however, the approximated gradients of generated UDF may not be stable enough at the zero-level set, which leads to errors and holes. We therefore adopt DCUDF [18] with double covering to mesh the generated UDF of UDiFF. Please refer to the supplementary for the adaptations to DCUDF. Finally, to create visual-appealing 3D models, we drew inspiration from Text2Tex [4]. This helps to generates textures for the extracted mesh while leveraging the text guidance within a progressive rendering-based texturing framework.

## 4. Experiment

In this section, we evaluate our proposed UDiFF on the task of shape generation. We first demonstrate the performance of UDiFF in generating novel shapes with open surfaces



Figure 4. Visual comparison with state-of-the-arts on the generated shapes under DeepFashion3D dataset. The front and back faces are rendered with different colors for a clear distinguish on open surfaces.

in Sec. 4.1. Next, we conduct experiments on generating shape with closed surfaces in Sec. 4.2. The ablation studies are shown in Sec. 4.3.

#### 4.1. Open-Surface Shape Generation

**Dataset.** For evaluations in the task of open-surface shape generation, we conduct experiments under DeepFashion3D dataset [79]. The DeepFashion3D dataset is a real-captured 3D dataset of open-surface clothes, containing 1,798 models reconstructed from real garments. It covers 10 categories and 563 garment instances. The dataset is randomly split into training and testing sets by the ratio 80% and 20 %. To get the text conditions for training UDiFF, we first render each model from the front facing view to obtain the image representing the model. We then leverage BLIP2 [29] for captioning the images and keep the caption description of the rendered image for each model as the text condition for the model. We further mix the category description of each model provided in the dataset into the text condition as a supplementary.

**Metrics.** For a fair comparison with various methods, we



Figure 5. Conditional generations produced by UDiFF and Shap-E. The front and back faces are rendered with different colors for a clear distinguish on open surfaces.

conduct the quantitative evaluations on the unconditional shape generation. We randomly generate 1,000 shapes with the trained model and uniformly sample 2,048 points on each generated shape. We follow previous works [21, 37] to evaluate the generation quality using Minimum matching distance (MMD), Coverage (COV) and 1-NN classifier accuracy (1-NNA). MMD measures the geometry accuracy of the generated shapes. COV indicates the ability of the generated shapes to cover the shapes in the test set. 1-NNA is designed to measure how well a classifier differentiates the generated shapes from the given shapes in the testing set. Lower is better for MMD, higher is better for COV and the closer to 50 % the better for 1-NNA.

**Baselines.** We compare UDiFF with the state-of-the-art methods in terms of the shape generation quality. PointDiff [37], WaveGen [21], Diffusion-SDF [10] and LAS-Diffusion [71]. PointDiff uses point cloud data for training, where we sample 2,048 points on each model and leverage the official code for training. All the previous implicit-based shape generation methods represent shapes as SDF or Occ, where the watertight meshes are required to generate the SDF/Occ data for training. Therefore, we leverage the commonly-used manifold method [20] for preprocessing the open-surfaces in DeepFashion3D. After that, we follow the official codes of these methods for training unconditional models with the watertight meshes.

**Comparison.** The quantitative comparison is shown in Tab. 1. where UDiFF achieves the best performance compared to the previous state-of-the-art methods. The main reason



Figure 6. Image conditioned generation of UDiFF. (a) Open-surface geometries generated with image guidance. (b) An example of generating textured shapes with image guidance.

is that all the previous implicit-based methods fail to handle the open-surfaces, where the needed manifold preprocessing leads to large bias on the original open-surface shapes. While the proposed UDiFF represents shapes as unsigned distance fields and is able to handle general shapes with or without open surfaces, leading to superior performance compared to other methods.

The visual comparison is shown in Fig. 4, where the proposed UDiFF significantly outperforms the previous works in generating visual-appealing clothes with open surfaces. We render the inside and outside surfaces in different colors for a clear difference on open surfaces. The PointDiff generates the point cloud to represent shapes, which do not require the manifold preprocess. However, it struggles to produce high-fidelity generations due to the discreteness of points.

**Text-conditional Generation.** For evaluations in generation with conditions, we further train a conditional model and generate shapes with the guidance from provided text prompts. We visually compare the generations with those produced by Shap-E under the same texts as shown in Fig. 5. The results demonstrate that UDiFF generates more accurate and high-fidelity predictions from the texts. UDiFF also produces more realistic textures thanks to the powerful Text2Tex [4]. On the contrary, Shap-E struggles to generate correct geometries and textures.

**Image-conditional Generation.** We further justify that UDiFF can receive diverse signals except texts (e.g. images) for conditional generation. This is achieved by leveraging the pre-aligned text and image representations of the CLIP model, where we adopt the frozen CLIP image encoder to achieve the image embeddings to guide UDiFF generation by cross-attention, without requiring extra training on images. We show the image-conditional generations of UDiFF in Fig. 6. The textures on the right of Fig. 6 is achieved with Text2Tex [4] on the text prompt predicted from the image with BLIP2 [29], i.e., ‘A white floral shirt with a long sleeves’.

## 4.2. Closed Shape Generation

**Dataset and metrics.** For the closed shape generation, we follow the common setting of previous methods [21, 28] to conduct generation experiments under the airplane and



Figure 7. Visual comparison with state-of-the-arts on the generated shapes under ShapeNet dataset.

chair classes of the ShapeNet [3] dataset. We randomly generate 2,000 shapes with the trained model and uniformly sample 2,048 points on each generated shape. We follow previous works [21, 37] to evaluate the generation quality using MMD, COV and 1-NNA. We compare our method with all the baselines using their officially provided pre-trained models and codes.

**Comparison.** We compare UDiFF with the state-of-the-art methods including IM-GAN [6], Voxel-GAN [27], PointDiff [37], SPAGHETTI [16], WaveGen [21] and SALAD [28]. We show the quantitative comparison in Tab. 2, where the results are directly borrowed from WaveGen and SALAD for a fair comparison.

The comparison demonstrates UDiFF also has the capability to generate high-fidelity watertight geometries with only closed surfaces. We justify that UDiFF is a general shape generator to produce general shapes with open surfaces and closed surfaces. We achieve the comparable performance with the state-of-the-art method SALAD, and also significantly outperform the baseline WaveGen which also leverages wavelet transformation as the compact representation. The reason is that our proposed approach for learning optimal wavelet filter largely reduces the information loss during transformation, which leads to more accurate and diverse generations. We further show the visual comparison of some generated shapes of different methods in Fig. 7. We can see that the shapes generated by our method are more faithful than IM-GAN and SPAGHETTI



Table 2. **Quantitative comparison of shape generation under ShapeNet dataset.** MMD-CD scores and MMD-EMD scores are scaled by  $10^3$  and  $10^2$ , respectively.

| Method              | COV $\uparrow$ |              | Chair MMD $\downarrow$ |              | 1-NNA $\downarrow$ |              | COV $\uparrow$ |              | Airplane MMD $\downarrow$ |              | 1-NNA $\downarrow$ |              |
|---------------------|----------------|--------------|------------------------|--------------|--------------------|--------------|----------------|--------------|---------------------------|--------------|--------------------|--------------|
|                     | CD             | EMD          | CD                     | EMD          | CD                 | EMD          | CD             | EMD          | CD                        | EMD          | CD                 | EMD          |
| IM-GAN [6]          | <b>56.49</b>   | 54.50        | 11.79                  | 14.52        | 61.98              | 63.45        | 61.55          | 62.79        | 3.320                     | 8.371        | 76.21              | 76.08        |
| Voxel-GAN [27]      | 43.95          | 39.45        | 15.18                  | 17.32        | 80.27              | 81.16        | 38.44          | 39.18        | 5.937                     | 11.69        | 93.14              | 92.77        |
| PointDiff [37]      | 51.47          | 55.97        | 12.79                  | 16.12        | 61.76              | 63.72        | 60.19          | 62.30        | 3.543                     | 9.519        | 74.60              | 72.31        |
| SPAGHETTI [16]      | 49.48          | 50.22        | 14.7                   | 15.85        | 72.34              | 69.46        | 56.86          | 58.83        | 4.260                     | 8.930        | 79.36              | 78.86        |
| SALAD (Global) [28] | 49.71          | 48.75        | 11.71                  | 14.12        | 62.72              | 61.25        | 54.88          | 59.33        | 3.877                     | 8.958        | 82.20              | 80.35        |
| SALAD [28]          | 56.42          | 55.16        | 11.69                  | 14.29        | <b>57.82</b>       | <b>58.41</b> | 63.16          | <b>65.39</b> | 3.636                     | 8.238        | <b>73.92</b>       | <b>71.08</b> |
| WaveGen [21]        | 49.63          | 50.15        | 12.12                  | 14.25        | 65.04              | 62.87        | 60.94          | 59.09        | 3.528                     | 7.964        | 75.77              | 72.93        |
| Ours                | 52.58          | <b>55.99</b> | <b>11.67</b>           | <b>14.04</b> | 65.96              | 63.42        | <b>64.77</b>   | 63.78        | <b>3.151</b>              | <b>7.798</b> | 74.48              | 78.99        |



Figure 8. Text-conditioned generation produced by UDiFF and AutoSDF under ShapeNet dataset.

Table 3. **Ablation studies on the framework design.** MMD-CD scores and MMD-EMD scores are scaled by  $10^3$  and  $10^2$ .

| Method              | COV $\uparrow$ |              | MMD $\downarrow$ |              | 1-NNA $\downarrow$ |              |
|---------------------|----------------|--------------|------------------|--------------|--------------------|--------------|
|                     | CD             | EMD          | CD               | EMD          | CD                 | EMD          |
| W/o learned wavelet | 64.52          | 65.02        | 13.24            | 15.26        | 85.06              | 86.22        |
| W/o fine predictor  | 66.36          | 65.18        | 12.37            | 14.48        | 83.62              | 84.17        |
| Full                | <b>69.62</b>   | <b>67.72</b> | <b>11.60</b>     | <b>14.01</b> | <b>81.83</b>       | <b>82.14</b> |

by producing finer details and cleaner surfaces, and have less bumpy geometries than WaveGen thanks to the optimal wavelet filter to significantly reduce information loss.

**Conditional Generation.** We further train a text-conditional model under the ‘Chair’ category of the ShapeNet dataset. We visually compare the generations produced by AutoSDF [47] and our propose UDiFF under the same texts as shown in Fig. 8. The results demonstrate that UDiFF generates more accurate and high-fidelity predictions from the texts compared to AutoSDF.

### 4.3. Ablation Studies

**Framework Design.** To evaluate the major components in our methods, we conduct ablation studies under the DeepFashion3D dataset [79] and report the performance in Tab. 3. We first justify the effectiveness of the proposed

Table 4. **Ablation studies on the effect of wavelet optimization.** We report the L2 Chamfer Distance scaled by  $10^5$ .

| Method | Haar                      | Biorthogonal3-3           | Biorthogonal6-8 |
|--------|---------------------------|---------------------------|-----------------|
| CD     | 264.8                     | 46.04                     | 42.92           |
| Method | Learnable $\phi_\theta^D$ | Learnable $\phi_\delta^I$ | Both            |
| CD     | 36.12                     | 32.15                     | <b>28.51</b>    |

optimal wavelet transformation by replacing our learned wavelet filter with the previous carefully chosen wavelet filter by WaveGen [21], i.e., Biorthogonal 6-8. The result is shown as ‘W/o learned wavelet’. We then remove the fine predictor of UDiFF to recover the 3D shapes with only the generated coarse coefficients as shown in ‘W/o fine predictor’. The ablation study results demonstrate that effect of designs in UDiFF by significantly improving the generation performance.

**The Effect of Wavelet Optimization.** We further evaluate the effect of our proposed wavelet optimization to achieve optimal wavelet filter. The result is shown in Tab. 4, where we conduct evaluations under the test set of DeepFashion3D [79] and report the L2 Chamfer Distance between the ground truth meshes and the recovered meshes with wavelet filters Haar, Biorthogonal3-3, Biorthogonal6-8 and ours. We show the performance of only optimizing decomposition filter parameters  $\phi_\theta^D$  and fix inversion filter parameters  $\phi_\delta^I$  as ‘Learnable  $\phi_\theta^D$ ’, and only optimize  $\phi_\delta^I$  with fixed  $\phi_\theta^D$  as ‘Learnable  $\phi_\delta^I$ ’. The best performance is achieved with optimizing both  $\phi_\delta^I$  and  $\phi_\theta^D$  as ‘Both’.

## 5. Conclusion

In this work, we present UDiFF, a 3D diffusion model for conditional or unconditional generating textured 3D shapes with open and closed surfaces. We leverage a diffusion model to learn distributions of UDFs in a spatial-frequency space established through an optimal wavelet transformation for UDFs, which is obtained by data-driven optimizations. The evaluations on widely used benchmarks show our superior performance over the latest methods in generating shapes with either open and closed surfaces.



## References

- [1] Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18392–18402, 2023. [2](#)
- [2] Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snaveley, and Bharath Hariharan. Learning gradient fields for shape generation. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 364–381. Springer, 2020. [3](#)
- [3] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3D model repository. *arXiv preprint arXiv:1512.03012*, 2015. [2](#), [4](#), [7](#), [12](#)
- [4] Dave Zhenyu Chen, Yawar Siddiqui, Hsin-Ying Lee, Sergey Tulyakov, and Matthias Nießner. Text2tex: Text-driven texture synthesis via diffusion models. *arXiv preprint arXiv:2303.11396*, 2023. [2](#), [5](#), [7](#), [12](#)
- [5] Weikai Chen, Cheng Lin, Weiyang Li, and Bo Yang. 3psdf: Three-pole signed distance function for learning surfaces with arbitrary topologies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18522–18531, 2022. [3](#)
- [6] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5939–5948, 2019. [7](#), [8](#)
- [7] Yen-Chi Cheng, Hsin-Ying Lee, Sergey Tulyakov, Alexander G Schwing, and Liang-Yan Gui. Sdfusion: Multimodal 3d shape completion, reconstruction, and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4456–4465, 2023. [3](#)
- [8] Julian Chibane, Gerard Pons-Moll, et al. Neural unsigned distance fields for implicit function learning. *Advances in Neural Information Processing Systems*, 33:21638–21652, 2020. [2](#), [3](#)
- [9] Gene Chou, Ilya Chugunov, and Felix Heide. Gensdf: Two-stage learning of generalizable signed distance functions. In *Advances in Neural Information Processing Systems*. [2](#)
- [10] Gene Chou, Yuval Bahat, and Felix Heide. Diffusion-sdf: Conditional generative modeling of signed distance functions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2262–2272, 2023. [2](#), [3](#), [5](#), [6](#)
- [11] Ingrid Daubechies. The wavelet transform, time-frequency localization and signal analysis. *IEEE transactions on information theory*, 36(5):961–1005, 1990. [2](#), [4](#)
- [12] Jun Gao, Tianchang Shen, Zian Wang, Wenzheng Chen, Kangxue Yin, Daiqing Li, Or Litany, Zan Gojcic, and Sanja Fidler. Get3d: A generative model of high quality 3d textured shapes learned from images. *Advances In Neural Information Processing Systems*, 35:31841–31854, 2022. [3](#)
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020. [3](#)
- [14] Benoit Guillard, Federico Stella, and Pascal Fua. MeshUDF: Fast and differentiable meshing of unsigned distance field networks. *European Conference on Computer Vision*, 2022. [5](#), [12](#)
- [15] Anchit Gupta, Wenhan Xiong, Yixin Nie, Ian Jones, and Barlas Oğuz. 3dgen: Triplane latent diffusion for textured mesh generation. *arXiv preprint arXiv:2303.05371*, 2023. [2](#), [3](#)
- [16] Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-Hornung, and Daniel Cohen-Or. Spaghetti: Editing implicit shapes through part aware generation. *ACM Transactions on Graphics (TOG)*, 41(4):1–20, 2022. [7](#), [8](#)
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. [2](#), [5](#)
- [18] Fei Hou, Xuhui Chen, Wencheng Wang, Hong Qin, and Ying He. Robust zero level-set extraction from unsigned distance fields based on double covering. *arXiv preprint arXiv:2310.03431*, 2023. [5](#), [12](#)
- [19] Han Huang, Yulun Wu, Junsheng Zhou, Ge Gao, Ming Gu, and Yu-Shen Liu. Neusurf: On-surface priors for neural surface reconstruction from sparse input views. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. [2](#)
- [20] Jingwei Huang, Hao Su, and Leonidas Guibas. Robust watertight manifold surface generation method for shapenet models. *arXiv preprint arXiv:1802.01698*, 2018. [6](#)
- [21] Ka-Hei Hui, Ruihui Li, Jingyu Hu, and Chi-Wing Fu. Neural wavelet-domain diffusion for 3d shape generation. In *SIGGRAPH Asia 2022 Conference Papers*, pages 1–9, 2022. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [22] Le Hui, Rui Xu, Jin Xie, Jianjun Qian, and Jian Yang. Progressive point cloud deconvolution generation network. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XV 16*, pages 397–413. Springer, 2020. [3](#)
- [23] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, Thomas Funkhouser, et al. Local implicit grid representations for 3D scenes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6001–6010, 2020. [2](#)
- [24] Chuan Jin, Tieru Wu, and Junsheng Zhou. Multi-grid representation with field regularization for self-supervised surface reconstruction from point clouds. *Computers & Graphics*, 2023. [2](#)
- [25] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems*, 31, 2018. [3](#)
- [26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. [2](#), [3](#)
- [27] Marian Kleineberg, Matthias Fey, and Frank Weichert. Adversarial generation of continuous implicit shape representations. *arXiv preprint arXiv:2002.00349*, 2020. [7](#), [8](#)
- [28] Juil Koo, Seungwoo Yoo, Minh Hieu Nguyen, and Minhyuk Sung. Salad: Part-level latent diffusion for 3d shape generation and manipulation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14441–14451, 2023. [3](#), [7](#), [8](#)

- [29] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv preprint arXiv:2301.12597*, 2023. 6, 7
- [30] Muheng Li, Yueqi Duan, Jie Zhou, and Jiwen Lu. Diffusion-sdf: Text-to-shape via voxelized diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12642–12651, 2023. 3
- [31] Ruihui Li, Xianzhi Li, Ka-Hei Hui, and Chi-Wing Fu. Sgan: Sphere-guided 3d shape generation and manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–12, 2021. 3
- [32] Shujuan Li, Junsheng Zhou, Baorui Ma, Yu-Shen Liu, and Zhizhong Han. NeAF: Learning neural angle fields for point normal estimation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2023. 2
- [33] Shujuan Li, Junsheng Zhou, Baorui Ma, Yu-Shen Liu, and Zhizhong Han. Learning continuous implicit field with local distance indicator for arbitrary-scale point cloud upsampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2024. 2
- [34] Yu-Tao Liu, Li Wang, Jie Yang, Weikai Chen, Xiaoxu Meng, Bo Yang, and Lin Gao. Neudf: Leaning neural unsigned distance fields with volume rendering. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 237–247, 2023. 3
- [35] Xiaoxiao Long, Cheng Lin, Lingjie Liu, Yuan Liu, Peng Wang, Christian Theobalt, Taku Komura, and Wenping Wang. Neuraludf: Learning unsigned distance fields for multi-view reconstruction of surfaces with arbitrary topologies. *arXiv preprint arXiv:2211.14173*, 2022. 3
- [36] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3D surface construction algorithm. *ACM Siggraph Computer Graphics*, 21(4):163–169, 1987. 2, 3, 12
- [37] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2837–2845, 2021. 2, 3, 5, 6, 7, 8
- [38] Baorui Ma, Zhizhong Han, Yu-Shen Liu, and Matthias Zwicker. Neural-Pull: Learning signed distance function from point clouds by learning to pull space onto surface. In *International Conference on Machine Learning*, pages 7246–7257. PMLR, 2021. 2
- [39] Baorui Ma, Yu-Shen Liu, and Zhizhong Han. Reconstructing surfaces for sparse point clouds with on-surface priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [40] Baorui Ma, Yu-Shen Liu, Matthias Zwicker, and Zhizhong Han. Surface reconstruction from point clouds by learning predictive context priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [41] Baorui Ma, Haoge Deng, Junsheng Zhou, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang. Geodream: Disentangling 2d and geometric priors for high-fidelity and consistent 3d generation. *arXiv preprint arXiv:2311.17971*, 2023. 2
- [42] Baorui Ma, Yu-Shen Liu, and Zhizhong Han. Learning signed distance functions from noisy 3d point clouds via noise to noise mapping. In *International Conference on Machine Learning (ICML)*, 2023. 2
- [43] Baorui Ma, Junsheng Zhou, Yu-Shen Liu, and Zhizhong Han. Towards better gradient consistency for neural signed distance functions via level set alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17724–17734, 2023. 2
- [44] Stephane G Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *IEEE transactions on pattern analysis and machine intelligence*, 11(7):674–693, 1989. 4
- [45] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3D reconstruction in function space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4460–4470, 2019. 2
- [46] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *European Conference on Computer Vision*, 2020. 2
- [47] Paritosh Mittal, Yen-Chi Cheng, Maneesh Singh, and Shubham Tulsiani. Autosdf: Shape priors for 3d completion, reconstruction and generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 306–315, 2022. 8
- [48] Thomas Müller, Alex Evans, Christoph Schied, and Alexander Keller. Instant neural graphics primitives with a multiresolution hash encoding. *ACM Transactions on Graphics (ToG)*, 41(4):1–15, 2022. 2
- [49] Gimin Nam, Mariem Khelifi, Andrew Rodriguez, Alberto Tono, Linqi Zhou, and Paul Guerrero. 3d-ldm: Neural implicit 3d shape generation with latent diffusion models. *arXiv preprint arXiv:2212.00842*, 2022. 2, 3
- [50] Alex Nichol, Heewoo Jun, Pratul Dharwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv preprint arXiv:2212.08751*, 2022. 3
- [51] Alexander Quinn Nichol and Pratul Dharwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning*, pages 8162–8171. PMLR, 2021. 5
- [52] Jeong Joon Park, Peter Florence, Julian Straub, Richard Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 165–174, 2019. 2
- [53] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision*, pages 523–540. Springer, 2020. 2
- [54] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 2, 5

- [55] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. [2](#)
- [56] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. [2](#), [5](#), [12](#)
- [57] J Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20875–20886, 2023. [2](#), [3](#)
- [58] Vincent Sitzmann, Julien Martel, Alexander Bergman, David Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. *Advances in Neural Information Processing Systems*, 33:7462–7473, 2020. [2](#)
- [59] Edward J Smith and David Meger. Improved adversarial systems for 3d object generation and reconstruction. In *Conference on Robot Learning*, pages 87–96. PMLR, 2017. [2](#), [3](#)
- [60] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. [2](#), [5](#)
- [61] Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, Karsten Kreis, et al. Lion: Latent point diffusion models for 3d shape generation. *Advances in Neural Information Processing Systems*, 35:10021–10039, 2022. [3](#)
- [62] Li Wang, Jie Yang, Weikai Chen, Xiaoxu Meng, Bo Yang, Jintao Li, and Lin Gao. Hsdf: Hybrid sign and distance field for modeling surfaces with arbitrary topologies. In *Advances in Neural Information Processing Systems*. [3](#)
- [63] Xin Wen, Junsheng Zhou, Yu-Shen Liu, Hua Su, Zhen Dong, and Zhizhong Han. 3D shape reconstruction from 2D images with disentangled attribute flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3803–3813, 2022. [2](#)
- [64] Xin Wen, Peng Xiang, Zhizhong Han, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Yu-Shen Liu. PMP-Net++: Point cloud completion by transformer-enhanced multi-step point moving paths. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):852–867, 2023. [2](#)
- [65] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. *Advances in neural information processing systems*, 29, 2016. [3](#)
- [66] Peng Xiang, Xin Wen, Yu-Shen Liu, Yan-Pei Cao, Pengfei Wan, Wen Zheng, and Zhizhong Han. Snowflake point deconvolution for point cloud completion and generation with skip-transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(5):6320–6338, 2023. [2](#)
- [67] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4541–4550, 2019. [3](#)
- [68] Jianglong Ye, Yuntao Chen, Naiyan Wang, and Xiaolong Wang. GIFS: Neural implicit function for general shape representation. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. [3](#)
- [69] Biao Zhang, Jiapeng Tang, Matthias Niessner, and Peter Wonka. 3dshape2vecset: A 3d shape representation for neural fields and generative diffusion models. *arXiv preprint arXiv:2301.11445*, 2023. [2](#), [3](#)
- [70] Wenyuan Zhang, Ruofan Xing, Yunfan Zeng, Yu-Shen Liu, Kanle Shi, and Zhizhong Han. Fast learning radiance fields by shooting much fewer rays. *IEEE Transactions on Image Processing*, 2023. [2](#)
- [71] Xin-Yang Zheng, Hao Pan, Peng-Shuai Wang, Xin Tong, Yang Liu, and Heung-Yeung Shum. Locally attentional sdf diffusion for controllable 3d shape generation. *arXiv preprint arXiv:2305.04461*, 2023. [5](#), [6](#)
- [72] Junsheng Zhou, Baorui Ma, Yu-Shen Liu, Yi Fang, and Zhizhong Han. Learning consistency-aware unsigned distance functions progressively from raw point clouds. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2022. [2](#), [3](#), [5](#), [12](#)
- [73] Junsheng Zhou, Xin Wen, Yu-Shen Liu, Yi Fang, and Zhizhong Han. Self-supervised point cloud representation learning with occlusion auto-encoder. *arXiv e-prints*, pages arXiv–2203, 2022. [2](#)
- [74] Junsheng Zhou, Baorui Ma, Shujuan Li, Yu-Shen Liu, and Zhizhong Han. Learning a more continuous zero level set in unsigned distance fields through level set projection. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2023. [3](#)
- [75] Junsheng Zhou, Baorui Ma, Wenyuan Zhang, Yi Fang, Yu-Shen Liu, and Zhizhong Han. Differentiable registration of images and lidar point clouds with voxelpoint-to-pixel matching. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2023. [2](#)
- [76] Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang. Uni3d: Exploring unified 3d representation at scale. *International Conference on Learning Representations*, 2024.
- [77] Junsheng Zhou, Xin Wen, Baorui Ma, Yu-Shen Liu, Yue Gao, Yi Fang, and Zhizhong Han. 3d-oae: Occlusion auto-encoders for self-supervised learning on point clouds. *IEEE International Conference on Robotics and Automation (ICRA)*, 2024. [2](#)
- [78] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5826–5835, 2021. [3](#)
- [79] Heming Zhu, Yu Cao, Hang Jin, Weikai Chen, Dong Du, Zhangye Wang, Shuguang Cui, and Xiaoguang Han. Deep fashion3d: A dataset and benchmark for 3d garment reconstruction from single images. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part I 16*, pages 512–530. Springer, 2020. [2](#), [4](#), [6](#), [8](#), [12](#)

## Supplementary Material

### A. More Visualizations

In this section, we provide more qualitative illustrations of the generation results produced by UDiFF.

#### Category-conditional generation on DeepFashion3D.

We provide extra open-surface shape generations achieved by the UDiFF model trained under DeepFashion3D [79] dataset with the cloth categories as the conditions. Specifically, we generate 8 categories of cloth shapes, including “long sleeve dress”, “long sleeve upper”, “pants”, “no sleeve dress”, “no sleeve upper”, “dress”, “shot sleeve dress” and “shot sleeve upper”. The visualizations are shown in Fig. 10 and 11, where UDiFF generates diverse and novel shapes correctly corresponds to the text conditions.

**Unconditional generation on ShapeNet.** We further provide more unconditional shape generation results achieved by the UDiFF model trained under single categories of ShapeNet [3] dataset. We generate shapes of the “chair” and “airplane” categories. The visualizations are shown in Fig. 12, where UDiFF generates visual-appealing shapes.

### B. Analysis on Meshing and Texturing

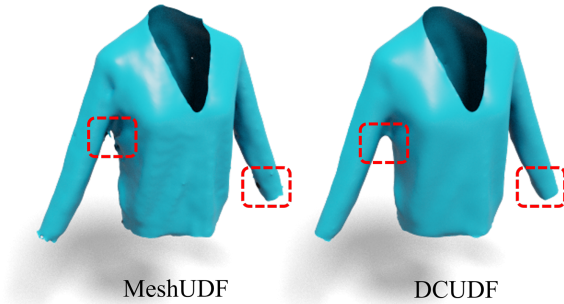


Figure 9. Mesh extraction comparisons between MeshUDF and DCUDF.

**Meshing.** Different from SDFs, UDFs fail to extract surfaces by the marching cubes [36] since UDFs cannot perform inside/outside tests on 3D grids. Recent works [14, 72] leverage the gradients at UDF grids as the signals to mesh UDFs. However, for the generated UDFs, the approximated gradients may not be stable enough at the zero-level set, which leads to errors and holes. The approximated gradient at a grid point  $q$  is defined as the direction from  $q$  to the neighbour grid  $q_n$  where the UDF from  $q$  to  $q_n$  increases rapidly the most. We adopt DCUDF [18] with double covering to mesh the generated UDF of UDiFF, which results in more continuous surfaces. We make an adaption to DCUDF on the double covering operation to replace the time-consuming optimizations with an explicit vertices re-

finement strategy. We move each vertices against the surface normals with a stride of unsigned distances to reach the zero-level sets, and then leverage the min-cut algorithm to achieve the final model. We show the comparison of meshing the generated UDF with MeshUDF [14] and DCUDF [18] in Fig. 9.

**Texturing.** We leverage Text2Tex [4] to generate textures for the extracted meshes. This is achieved with a progressive texture generation process and a texture refinement process. Specifically, we first render the texture-less initial mesh from the preset viewpoint and generate the appearance according to the text prompt with the depth-guided stable-diffusion [56]. We then adjust to the next preset viewpoint and repeat the appearance generation process until the last preset viewpoint where the whole mesh is textured. Finally, we optimize the textures with automatically selected viewpoints for refinement.



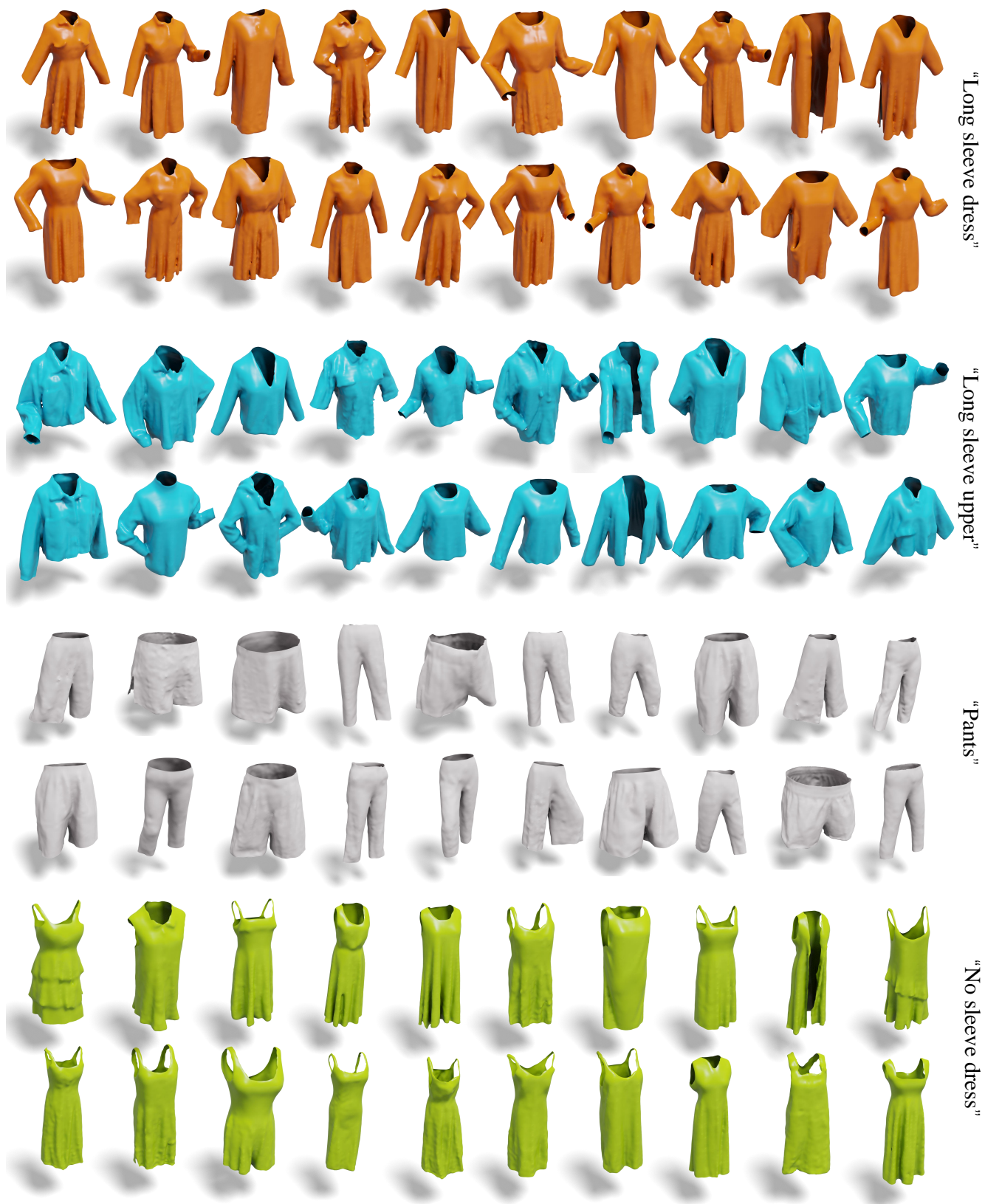


Figure 10. Category conditional generations under DeepFashion3D dataset. Here, we visualize the conditional generations of categories "long sleeve dress", "long sleeve upper", "pants" and "no sleeve dress"

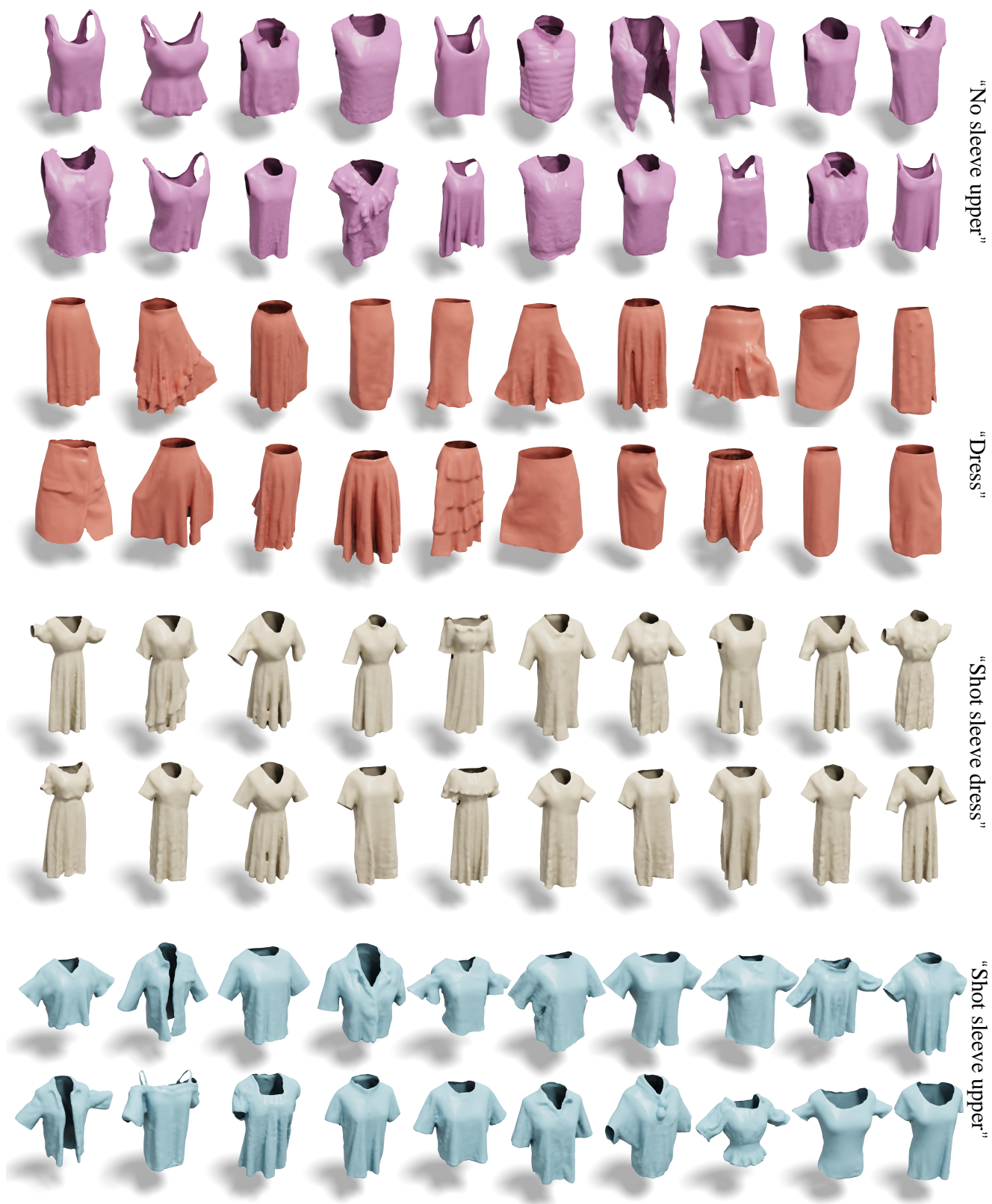


Figure 11. Category conditional generations under DeepFashion3D dataset. Here, we visualize the conditional generations of categories *"no sleeve upper"*, *"dress"*, *"shot sleeve dress"* and *"shot sleeve upper"*

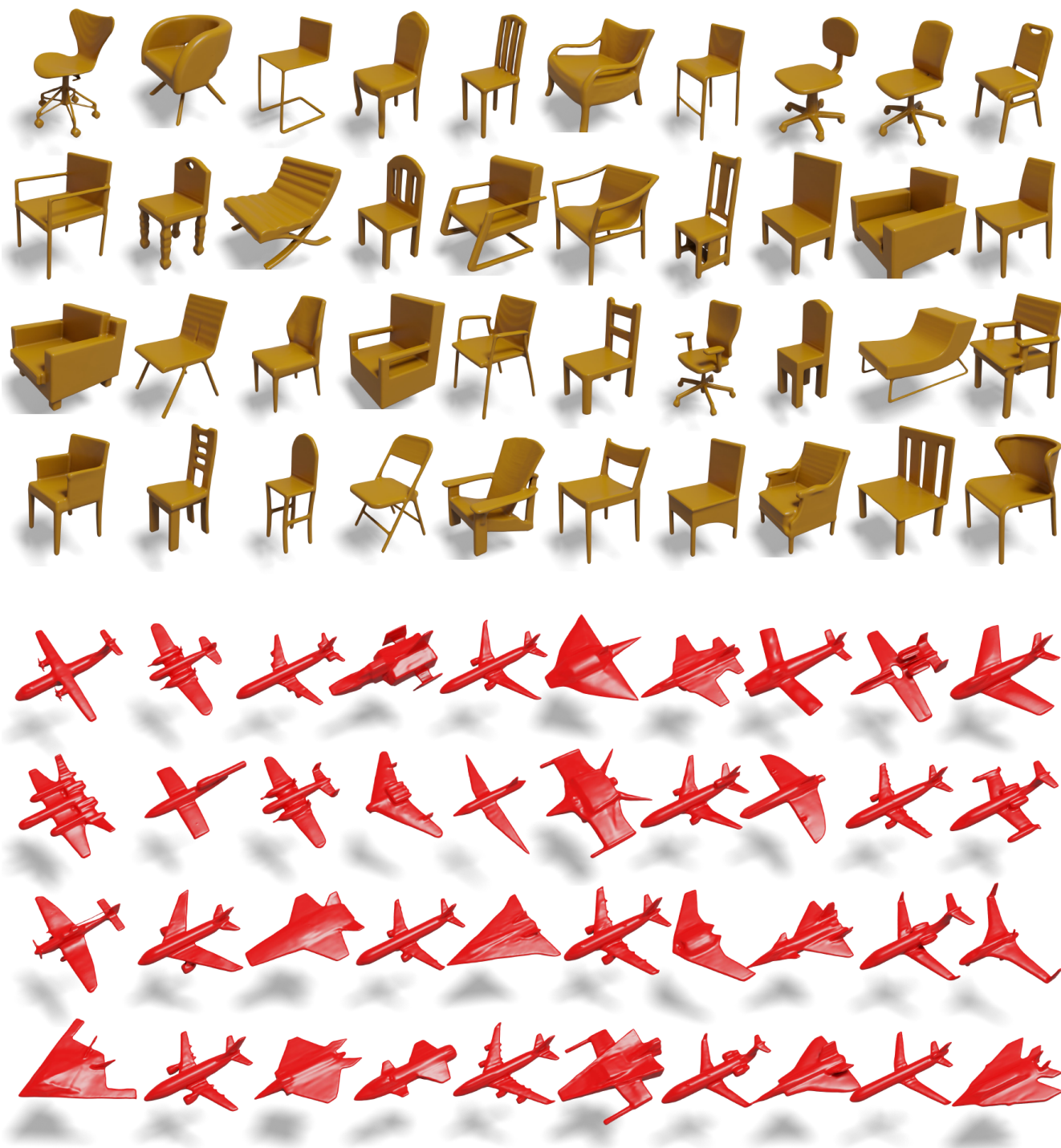


Figure 12. Unconditional generations under the “chair” and “airplane” categories of the ShapeNet dataset.