

# Neural Point Cloud Diffusion for Disentangled 3D Shape and Appearance Generation

Philipp Schröppel<sup>1</sup> Christopher Wewer<sup>2</sup> Jan Eric Lenssen<sup>2</sup> Eddy Ilg<sup>3</sup> Thomas Brox<sup>1</sup>  
<sup>1</sup>University of Freiburg <sup>2</sup>Max Planck Institute for Informatics <sup>3</sup>Saarland University

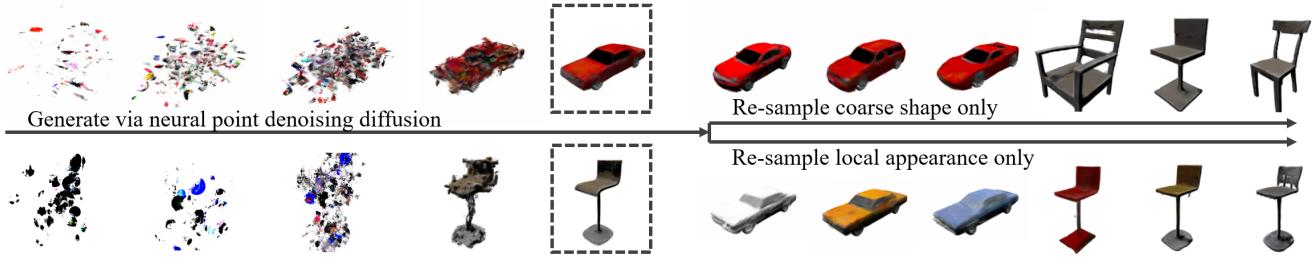


Figure 1. We present a method to model 3D radiance field distributions using neural point denoising diffusion (**left**). Since our representation disentangles coarse object shape from local appearance, we can sample from the individual distributions separately (**right**).

## Abstract

*Controllable generation of 3D assets is important for many practical applications like content creation in movies, games and engineering, as well as in AR/VR. Recently, diffusion models have shown remarkable results in generation quality of 3D objects. However, none of the existing models enable disentangled generation to control the shape and appearance separately. For the first time, we present a suitable representation for 3D diffusion models to enable such disentanglement by introducing a hybrid point cloud and neural radiance field approach. We model a diffusion process over point positions jointly with a high-dimensional feature space for a local density and radiance decoder. While the point positions represent the coarse shape of the object, the point features allow modeling the geometry and appearance details. This disentanglement enables us to sample both independently and therefore to control both separately. Our approach sets a new state of the art in generation compared to previous disentanglement-capable methods by reduced FID scores of 30-90% and is on-par with other non-disentanglement-capable state-of-the-art methods.*

## 1. Introduction

3D assets are used in many practical applications, ranging from engineering to movies and computer games, and will

Code is provided at <https://github.com/lmb-freiburg/neural-point-cloud-diffusion>.

become even more important in virtual spaces and virtual telepresence that will be enabled by AR/VR technology in the future. However, manually creating such assets is a labor-intensive and costly task that requires expert skills. It is even more important that such content cannot just be generated but that the generation can also be controlled to obtain the desired outcome. With the impressive image generation capabilities of diffusion models [9, 13, 31–33, 46], it is appealing to consider such models. Generally, the extension to 3D is still limited and not straightforward [11, 26, 35]. Even more so, none of the current diffusion models allow for disentangled generation of shape and appearance and controlling them separately.

The general challenge for 3D diffusion models lies in selecting the right 3D representation. One track of work explores diffusion models to generate 3D point clouds [4, 22, 26, 45, 47]. While such methods are able to generate the sparse point clouds well, they are not able to model dense geometry or appearance. Another track uses implicit representations [10, 11, 15], triplanes [8, 35] or voxel grids [25] to define the geometry and appearance continuously for each coordinate in a volume that encloses the object. The downside of all of these representations is that they do not provide disentanglement of shape and appearance. The reason for this limitation is the missing invariance of a single factor of variation to changes in others: Global neural field representations, for example, model shape and appearance in joint parameters. Thus, one of these factors cannot be changed independently. Voxel grids or triplanes provide limited invariance of appearance variables to voxel-sized

shifts but fail to be invariant to sub-voxel shifts or more complex, non-rigid sub-voxel deformation.

In contrast, we propose a method that enables individual generation of shape and appearance by introducing a hybrid approach that consists of a neural point cloud hosting a continuous radiance field. The point cloud explicitly disentangles coarse object shape from appearance. Feature vectors model the geometry and appearance of *local parts* [5] and, while the point positions determine *where* a part is, the point features describe *how* the details of a part look like. Notably, the point positions can undergo complex deformations without requiring changes in point features. With this representation, we are able to control the generation of both aspects separately, as illustrated in Fig. 1.

To establish our model, we first train a generalizable Point-NeRF renderer [42] by sharing its weights across many instances of ShapeNet [36] or PhotoShape [29] objects. The obtained neural point clouds then serve as a dataset to train a diffusion model that learns to denoise the point positions and features simultaneously. Different to previous diffusion models, our model operates on high-dimensional latent spaces. In summary, our contributions are:

1. We propose the first approach for object generation that leverages a hybrid approach consisting of a neural point cloud combined with a neural renderer and a diffusion model that operates in a high-dimensional latent space.
2. We identify many-to-one mappings as a crucial obstacle when applying denoising diffusion to autodecoded, high-dimensional latent spaces and present effective regularization schemes to overcome this issue.
3. We show that our approach is capable of successfully disentangling geometry and appearance by allowing to control them separately and that the generation quality our approach significantly outperforms the previous methods GRAF [34] and Disentangled3D [39] by a large margin, while being on-par in generation quality with state-of-the-art methods incapable of disentangling.

## 2. Related work

**Disentangled generation.** Disentangled generation of shape and appearance has been studied in multiple works and is commonly achieved by modeling distinct architecture parts and latent codes [14, 28, 49]. GRAF [27, 34] presented the first generative model for radiance fields that allows to separately control both factors. Disentangled3D (D3D) [39] provides a more explicit disentanglement by leveraging a canonical volume and deformation. In contrast to ours, none of the approaches uses diffusion models or point clouds. We can show that our diffusion model clearly outperforms these previous GAN-based methods.

**Probabilistic diffusion models.** In recent years, diffusion models [37, 38] have emerged as a successful class of generative models. They first define a forward diffusion process in the form of a Markov chain, which gradually transforms the data distribution to a simple known distribution. A model is then trained to reverse this process, which then allows to sample from the learned data distribution. DDPM [13] proposes a diffusion model formulation with various simplifications that enable high quality image synthesis. Follow-up works improved on synthesis quality via refinements of the architecture and sampling procedure to even outperform state-of-the-art generative adversarial networks [9, 16]. In order to scale diffusion models to high resolutions, LDM [32] moves the diffusion process from the image to a latent space with smaller spatial dimensions. In this work, we adopt the DDPM diffusion formulation and apply it to 3D objects on a high-dimensional latent space in the form of neural point clouds.

**NeRF and Point-NeRF.** NeRF [24] represents geometry and appearance as a radiance and density field that can be volumetrically rendered to photorealistic images. Point-NeRF [42] extends NeRF to a parameterization by a point cloud that is obtained from MVSNet [44]. This reduces ambiguity and allows for a much faster reconstruction. The Point-NeRF MLPs are originally trained on a single scene. In this work, we train them jointly on many objects to obtain a generalizable version [41]. In addition, we regularize the neural point cloud features to be optimally suitable for the diffusion model, which we use to generate objects.

**3D diffusion.** There are currently two trends of applying diffusion to 3D: (1) Test-time distillation using large pre-trained image generators [20, 21, 23, 30, 48] and (2) diffusion models on datasets of 3D models [2, 4, 6, 8, 10, 15, 22, 25, 26, 35, 45, 47]. Our approach can be assigned to the second category, and therefore we focus on this direction in the following.

**Diffusion on point clouds.** Unlike alternative 3D representations, such as dense voxel grids or meshes, point clouds are sparse, unlimited to pre-defined topologies, and flexible w.r.t. modifications. Therefore, there are several recent works combining these advantages with the generative power of diffusion models [4, 22, 26, 45, 47]. Except for differences in network architecture and the relation to score matching or diffusion models, first approaches [4, 22, 47] all define the diffusion process directly on 3D point coordinates. In contrast to that, LION [45] applies the idea of LDM [32] to point clouds by denoising latent codes of a hierarchical VAE. However, with only a single dimension, their latent codes are not very expressive and very low dimensional. Following the great success of text-to-image

generation, Point-E [26] trains a transformer-based architecture for generation of RGB point clouds conditioned on complex prompts. Unlike all of these approaches, our method generates point clouds with high-dimensional features encoding detailed shape and appearance.

**Diffusion for 3D object shape and appearance.** Previous approaches for object shape and appearance synthesis opt for other 3D representations. Functa [10] and ShapE [15] generate the weights of implicit (neural) representations such as radiance fields or signed distance functions. DiffRF [25] uses a 3D-UNet to denoise explicit voxel grids storing density and color. The usual training pipeline for these approaches involves first fitting 3D representations to a dataset of multi-view images. Recently, SSDNeRF [8] showed improvements by optimizing the diffusion model and individual NeRFs for each training object in a joint single stage. RenderDiffusion [2] avoids the question about single- or two-stage training by defining the diffusion process again in image space, but it appends the triplane representation from EG3D [6] to the usual diffusion architecture for 3D view conditioning. In contrast to implicit representations, voxel grids and triplanes, our neural point clouds enable the disentanglement of shape and appearance.

### 3. Method

In this section, we describe Neural Point Cloud Diffusion (NCPD), our generative model for 3D shape and appearance via diffusion on neural point clouds. An overview of the method is shown in Fig. 2. At the center of our method is a generalizable neural point representation, which is further described in Sec. 3.1. We discuss characteristics of autodecoder schemes in Sec. 3.2 and provide regularization schemes that enable denoising diffusion on the feature space. Subsequently, in Sec. 3.3 we then present a diffusion model to denoise the neural point positions and features. After the diffusion model is trained, we can sample 3D shape and appearance independently from each other as described in Sec. 3.4.

#### 3.1. Generalizable Point-NeRF

We begin by outlining our representation as an extension of Point-NeRF [42]. A single object is represented by a neural point cloud  $\mathcal{P} = \{(\mathbf{p}_1, \mathbf{f}_1), \dots, (\mathbf{p}_M, \mathbf{f}_M)\} = (\mathbf{P}, \mathbf{F})$  where each 3D point  $i \in \{1, \dots, M\}$  with position  $\mathbf{p}_i \in \mathbb{R}^3$  is associated with a neural feature  $\mathbf{f}_i \in \mathbb{R}^D$ . We denote the full matrix of point positions as  $\mathbf{P} \in \mathbb{R}^{M \times 3}$  and the full matrix of features as  $\mathbf{F} \in \mathbb{R}^{M \times D}$ . In contrast to PointNeRF, we assume the point positions  $\mathbf{P}$  to be given for training objects and  $\mathbf{F}$  to be manually initialized before optimization. We explore different initialization strategies in Sec. 4.6. The neural point representation can be rendered from arbitrary views, as described in the following.

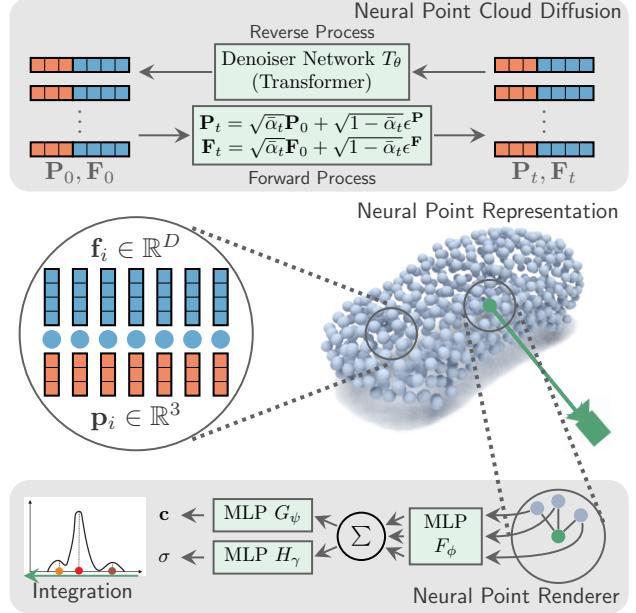


Figure 2. **Overview of neural point cloud diffusion (NCPD).** In the **center** we have a neural point cloud representation, where each point has a position ( $\blacksquare$ ) and an appearance feature ( $\blacksquare$ ). The neural point cloud can be generated with a diffusion model (**top**) and can be rendered via ray integration (**bottom**).

**Volume rendering.** To render a pixel of an image we follow the Point-NeRF [42] procedure. Given camera parameters, we march rays through the scene and sample shading points  $\mathbf{q}$  along the ray. For each shading point, the features  $\mathbf{f}$  of the neighboring points  $\mathbf{p}$  of the neural point cloud are first aggregated to a shading point feature  $\mathbf{f}_q$  via a multi-layer perceptron (MLP)  $F_\phi$  and a weighted combination based on inverse distances:

$$\mathbf{f}_q = \sum_{i=1}^{w_i} \frac{w_i F_\phi(\mathbf{f}_i, \mathbf{q} - \mathbf{p}_i)}{\sum_{i=1}^{w_i}}, \text{ where } w_i = \frac{1}{\|\mathbf{q} - \mathbf{p}_i\|_2}. \quad (1)$$

This feature is then mapped to a color  $\mathbf{c}$  and density  $\sigma$  by separate MLPs  $H_\gamma$  and  $G_\psi$ :

$$\mathbf{c} = G_\psi(\mathbf{f}_q) \quad \sigma = H_\gamma(\mathbf{f}_q) \quad (2)$$

The obtained radiances are then numerically integrated to the pixel color as described in NeRF [24]. Given camera parameters  $\mathbf{v}$ , we denote the full rendering function that renders an image as  $R_{\phi, \psi, \gamma}^\mathbf{v}(\mathbf{P}, \mathbf{F})$ .

**Optimization.** Optimization is done on a dataset of  $N$  objects  $O_1, \dots, O_N$ . Each object  $O_j$  consists of a neural point cloud  $\mathcal{P}_j = (\mathbf{P}_j, \mathbf{F}_j)$  and  $K$  views  $V_{j1}, \dots, V_{jK}$ . Each view  $V_{jk} = (\mathbf{I}_{jk}, \mathbf{v}_{jk})$  consists of a ground truth image  $\mathbf{I}_{jk}$  and corresponding camera parameters  $\mathbf{v}_{jk}$ . The optimization

objective is to jointly find the point features  $\mathbf{F}$  and network parameters  $\phi, \psi, \gamma$  that minimize the image reconstruction error for all views of all objects:

$$\hat{\mathbf{F}}, \hat{\phi}, \hat{\psi}, \hat{\gamma} = \arg \min_{\mathbf{F}, \phi, \psi, \gamma} \sum_{j,k} \mathcal{L} \left( R_{\phi, \psi, \gamma}^{\mathbf{v}_{jk}} (\mathbf{P}_j, \mathbf{F}_j), \mathbf{I}_{jk} \right), \quad (3)$$

with  $\mathcal{L}$  being the mean squared error between rendered and ground truth pixel colors. In contrast to Point-NeRF, we share parameters  $\phi, \psi, \gamma$  of the rendering MLPs over all objects, obtaining a generalizable representation.

### 3.2. Autodecoding for diffusion

The objective given above in Eq. (3) describes an under-constrained optimization problem. We found that, without further regularization, many-to-one mappings between features  $\mathbf{F}$  and renderings  $R_{\theta, \psi, \phi}^{\mathbf{v}} (\mathbf{P}, \mathbf{F})$  emerge. Thus, there are multiple possible point features  $\mathbf{f}_i$  representing the same local appearance information. We support this hypothesis with an empirical verification in Sec. 4.6 and Tab. 4. The many-to-one mappings pose a challenge for the denoising model, as it is trained to produce point estimates of the features  $\mathbf{f}_i$  that are in that case ambiguous.

Most existing latent diffusion methods circumvent this issue by using an autoencoder [32, 45] instead of optimizing representations via backpropagation. Since encoder networks are functions by design, and thus assigning each input value only one output, they do not produce many-to-one mappings between latent representation and output. However, we argue that the autodecoder principle is preferred in many situations, since it does not require an encoder (which is difficult to design for representations like neural points) and often leads to higher quality.

Therefore, we present and analyze a list of strategies to eliminate many-to-one mappings from autodecoder formulations, which are outlined in the following paragraphs.

**Zero initialization.** The first simple, albeit effective strategy is to initialize features  $\mathbf{F}$  with zero instead of randomly sampled values. We can show that this is very effective and encourages convergence to the same minimum.

**Total variation regularization (TV).** Inspired by TV regularization on triplanes [35], we design a TV regularization baseline for our neural point clouds by adding

$$\mathcal{L}_{TV}(\mathbf{F}) = \lambda_{TV} \sum_{i=1}^M \sum_{n \in \mathcal{V}(i)} \frac{\|\mathbf{f}_i - \mathbf{f}_n\|_1}{\|\mathbf{p}_i - \mathbf{p}_n\|_2}, \quad (4)$$

with weighting  $\lambda_{TV}$  to the objective in Eq. (3) for each object, where  $\mathcal{V}(i)$  is a local neighborhood of points around the point with index  $i$ . Intuitively, it encourages that neighboring point features are varying only slightly.

**Variational autodecoder (KL).** We introduce a variational autodecoder by storing vectors of means  $\mu_i$  and isotropic variances  $\Sigma_i$  instead of features  $\mathbf{f}_i$  for each point. For rendering, we obtain the features  $\mathbf{f}_i$  by sampling from the corresponding Gaussians using the reparameterization trick [18]. Additionally, we add a KL divergence loss

$$\mathcal{L}_{KL}(\{\mu_i, \Sigma_i\}_{i=1}^M) = \lambda_{KL} \sum_{i=1}^M KL(\mathcal{N}(\mu_i, \Sigma_i), \mathcal{N}(\mathbf{0}, \mathbf{I}_D)), \quad (5)$$

for each object with weighting  $\lambda_{KL}$ , to minimize the evidence lower bound [18]. Intuitively, the consequences from this regularization are twofold. First, the latent space is regularized to follow a unit Gaussian distribution, reducing many-to-one mappings in the representation. Second, the decoder learns to be more robust to small changes in  $\mathbf{f}$  due to the sampling procedure. Note that our diffusion model regresses  $\mu_i$ , but we write  $\mathbf{f}_i$  in the following for simplicity.

### 3.3. Neural point cloud diffusion

In this section, we describe our diffusion model for neural point cloud representations. As input, we assume a set of optimized representations  $\{\mathcal{P}_j\}_{j=1}^N$  from the first stage. The diffusion model learns the distribution of these representations, which allows us to generate neural point clouds.

**Denoising diffusion background.** Diffusion models [37, 38] learn the distribution  $q(\mathbf{x})$  of data  $\mathbf{x}$  by defining a forward diffusion process in form of a Markov chain with steps  $q(\mathbf{x}_t | \mathbf{x}_{t-1})$  that gradually transform the data distribution into a simple known distribution. A model  $p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t)$  with parameters  $\theta$  is then trained to approximate the steps  $q(\mathbf{x}_{t-1} | \mathbf{x}_t)$  in the Markov chain of the reverse process, which allows to evaluate the likelihood of a given data point, or sample from the learnt distribution.

In case of Gaussian diffusion processes, the forward diffusion process gradually replaces the data with Gaussian noise following a noise schedule  $\beta_1, \dots, \beta_T$ :

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) := \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I}). \quad (6)$$

Further, it is possible to sample  $\mathbf{x}_t$  directly from  $\mathbf{x}_0$ :

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}), \quad (7)$$

with  $\alpha_t := 1 - \beta_t$  and  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ . For small enough step sizes in the noise schedule, the steps in the Markov chain of the reverse process can be approximated with Gaussian distributions:

$$p_\theta(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I}), \quad (8)$$

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \Sigma_\theta(\mathbf{x}_t, t)). \quad (9)$$

The objective hence is to learn  $\mu_\theta(\mathbf{x}_t, t)$  and  $\Sigma_\theta(\mathbf{x}_t, t)$ . DDPM [13] proposes a diffusion model formulation with

various simplifications. Specifically, DDPM suggests to fix  $\Sigma_\theta(\mathbf{x}_t, t)$  and the noise schedule  $\beta_t$  to time-dependent constants. Further, DDPM reparametrizes Eq. (7) to

$$\mathbf{x}_t(\mathbf{x}_0, \epsilon) = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon \text{ with } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (10)$$

and trains the model to directly predict  $\epsilon_\theta(\mathbf{x}_t, t)$ , from which  $\mu_\theta(\mathbf{x}_t, t)$  can be computed.

**Neural point cloud diffusion.** Given the background in DDPM, we turn to describing our neural point cloud diffusion. Conceptually, the neural points clouds  $\mathcal{P} = (\mathbf{P}, \mathbf{F})$  take the place of data points  $\mathbf{x}_0$  in the above DDPM diffusion model formulation.

The distribution of both modalities, point positions  $\mathbf{P}$  and appearance features  $\mathbf{F}$ , is learnt jointly. During training, we sample Gaussian noise  $\epsilon^{\mathbf{P}}$  for all point positions and  $\epsilon^{\mathbf{F}}$  for all point features and use it to obtain the noised neural point cloud  $\mathcal{P}_t = (\mathbf{P}_t, \mathbf{F}_t)$  at a specific timestep  $t$  in the diffusion process via Eq. (10). The denoiser network  $T_\theta((\mathbf{P}_t, \mathbf{F}_t), t) = (\epsilon_\theta^{\mathbf{P}}, \epsilon_\theta^{\mathbf{F}})$  takes the noised neural point cloud and timestep as input and estimates the noise  $\epsilon_\theta^{\mathbf{P}}$  and  $\epsilon_\theta^{\mathbf{F}}$  that was applied to the points and features. The network is optimized by minimizing the average mean squared error on both noise vectors:

$$\mathcal{L}_T = \frac{1}{2}(\text{MSE}(\epsilon^{\mathbf{P}}, \epsilon_\theta^{\mathbf{P}}) + \text{MSE}(\epsilon^{\mathbf{F}}, \epsilon_\theta^{\mathbf{F}})). \quad (11)$$

**Denoiser architecture.** As the architecture for the denoiser network we use a Transformer [26, 40]. As input, the transformer receives  $M + 1$  tokens, one token per point plus one additional token encoding  $t$ . Point tokens are obtained by concatenating the point position and feature of each point in the noisy point cloud and projecting them with a linear layer. Similarly, the  $t$  token is obtained by projection with its own linear layer onto the same dimensionality. After encoding, all tokens are processed by transformer layers, including self-attention and MLPs. Finally, the resulting output tokens corresponding to the  $M$  points are projected back to the dimensions of the input point positions and features and interpreted as noise predictions  $\epsilon_\theta^{\mathbf{P}}$  and  $\epsilon_\theta^{\mathbf{F}}$ .

### 3.4. Disentangled generation

Given a trained NPCD model, we can naively sample from the joint distribution  $p(\mathbf{P}, \mathbf{F})$  of point positions and features by sampling positions and features from a unit Gaussian distribution and using the transformer for iterative denoising (c.f. Fig. 1 for a visualization). To achieve individual generation of shape and appearance, one needs to sample from conditional distributions  $p(\mathbf{P}|\mathbf{F})$  or  $p(\mathbf{F}|\mathbf{P})$  instead.

To sample from  $p(\mathbf{F}|\mathbf{P})$  given point positions  $\mathbf{P}_0$ , we obtain the initial noisy neural point cloud  $\mathcal{P}_T = (\mathbf{P}_T, \mathbf{F}_T)$  by

$$\mathbf{P}_T = \sqrt{\bar{\alpha}_t} \mathbf{P}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon^{\mathbf{P}}, \quad \mathbf{F}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (12)$$

i.e., only sampling  $\mathbf{F}_T$  from noise while obtaining  $\mathbf{P}_T$  as a noisy variant of  $\mathbf{P}_0$ . Then, the denoising transformer is applied. After each reverse process iteration, we reset  $\mathbf{P}_t$  to the values from Eq. (12) for the appropriate  $t$ . Sampling from  $p(\mathbf{P}|\mathbf{F})$  can be done analogously.

Note that the conditional sampling procedure described above is enabled by the neural point representations, which allow to individually modify the disentangled variables for coarse shape and local appearance.

## 4. Experiments

In this section, we provide experimental results for the presented NPCD method. We begin by introducing the experimental setup in Sec. 4.1 and used metrics in Sec. 4.2. Then, we evaluate the main contribution of our method in Sec. 4.3, i.e. the disentangled generation of coarse geometry and appearance. We compare against previous generative approaches that allow disentangled generation, namely GRAF [34] and Disentangled3D [39], and show our superior generation quality. Next, we compare against recent diffusion models without disentangling capabilities in Sec. 4.4. Here, we compare against DiffRF [25], Functa [10] and SSDNeRF [8]. Many existing 3D generative models model only shape but not appearance. Thus, to complement existing comparisons, we also provide a shape-only comparison in Sec. 4.5. Finally, we analyze many-to-one mappings due to auto-decoded features as a problem for diffusion models and propose regularization methods as effective countermeasures in Sec. 4.6.

### 4.1. Datasets and experimental setup

**Data.** We use the cars and chairs categories of the ShapeNet SRN dataset [7, 36]. The cars split contains 2,458 training objects and 704 test objects, while the chairs split contains 4,611 training and 1,317 test objects. We use the original renderings with 50 views per training object and 251 views per test object. The images have a resolution of 128x128 pixels. For all training objects, the poses are sampled randomly from a sphere. For all test objects the poses follow a spiral on the upper hemisphere. We extract point clouds with 30k points from the mesh and subsample them to 512 points with farthest point sampling.

Besides ShapeNet SRN, we also use the PhotoShape Chairs dataset [29]. The dataset contains 15,576 objects and features more realistic textures on top of ShapeNet meshes. We use the same test split as DiffRF [25], which consists of 1,552 objects. From the remaining objects, we randomly select 2,480 objects for training. Further, we use the same renderings as DiffRF [25], which consist of 200 views per object on an Archimedean spiral. We use images at a resolution of 128x128 pixels. We use the same point clouds with 512 points as for SRN chairs.

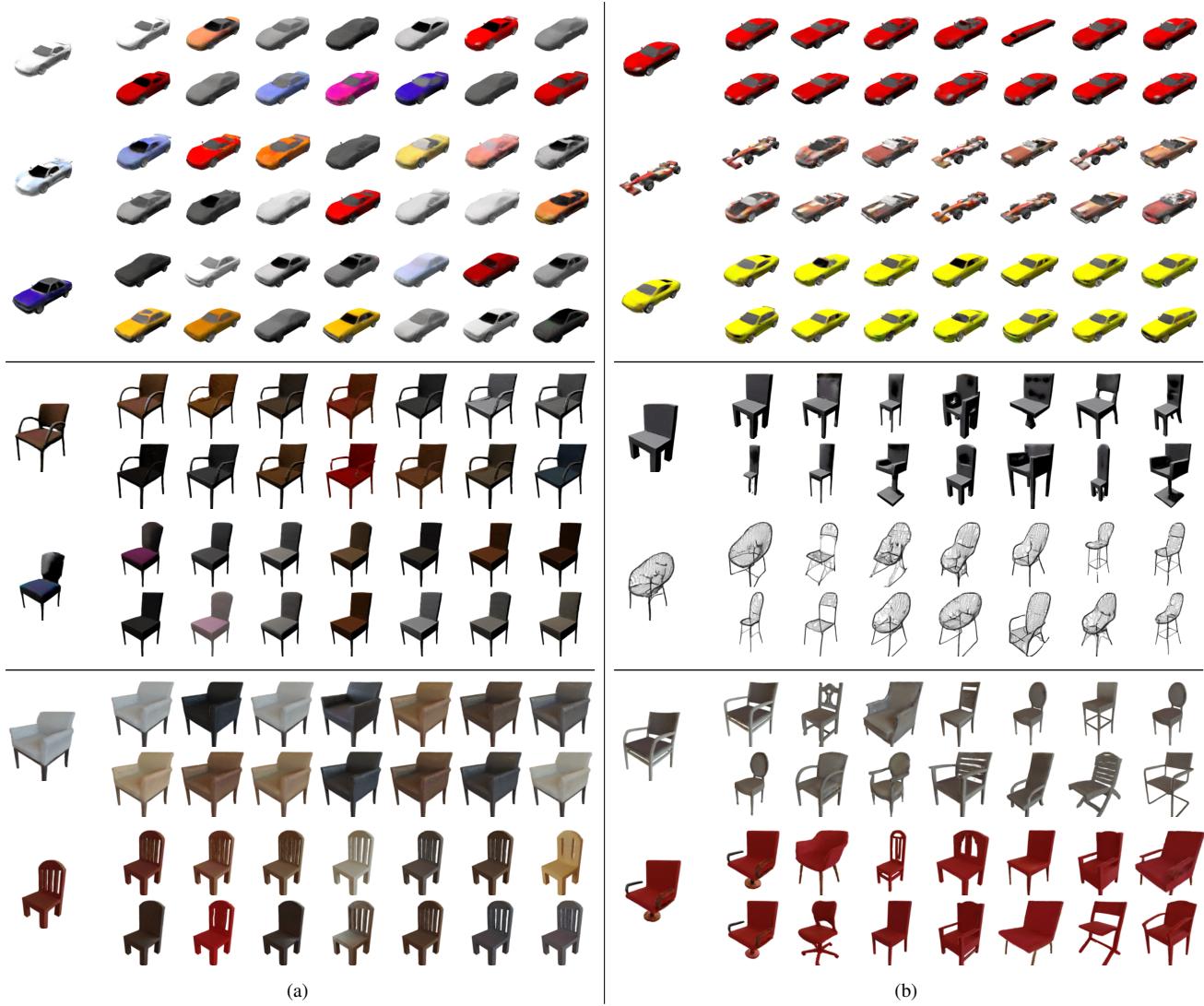


Figure 3. **Qualitative examples of disentangled generation** on SRN cars, SRN chairs, PhotoShape chairs. **(a) Appearance-only generation:** we show a generated object and objects with re-sampled appearance. **(b) Shape-only generation:** we show a generated object and objects with re-sampled coarse shape. We can get diverse samples of local appearance or coarse shape when the respective other is given.

**Training details.** We construct neural point clouds by zero-initializing features for each point. In generalizable Point-NeRF, we optimize the reconstruction loss in Eq. (3) with the TV and KL regularizers from Eq. (4) and Eq. (5). Further details on network architectures and training parameters are provided in the supplementals. For the diffusion model training, we normalize the neural point clouds and use DDPM [13] diffusion model parameters. Further details on the denoiser architecture, diffusion model parameters, and training parameters are provided in the supplementals.

## 4.2. Metrics

To measure the quality and diversity of the generated samples of the diffusion model, we report the FID [12] and

KID [3] metrics. FID and KID compare the appearance and diversity of two image sets by computing features for each image with an Inception model and comparing the feature distributions of the two sets. We use the images of the test set objects as the reference set. For comparability, we follow the evaluation procedures of previous works: on SRN Cars and Chairs, we generate the same number of objects as in the test set and render them from the same poses; on PhotoShape Chairs, we generate 1,000 objects and render them from 10 poses that are sampled randomly from the Archimedean spiral poses. Furthermore, for the shape-only evaluation of our generated point clouds representing the coarse geometry, we employ 1-nearest-neighbor accuracy w.r.t. Chamfer and Earth Mover’s Distance [43]. Last, we

conduct a quantitative analysis by reporting the per-point mean cosine similarities between optimized neural point features of 10 random training examples over 100 different seeds with a fixed renderer. This measures the extent of *many-to-one mappings*, i.e. how far features that represent the same appearance are away from each other.

### 4.3. Disentangled generation

A major advantage of using neural point clouds as 3D representation is their intrinsic disentanglement of shape and appearance: the point positions represent the coarse geometry and the features model local geometry and appearance. As described in Sec. 3.4, this property enables the proposed method to generate both modalities separately, even though a joint distribution is modeled with a single diffusion model. In the following, we present the results of this disentangled 3D shape and appearance generation.

NPCD allows to explicitly control the point positions or point features throughout the diffusion process. As a consequence, we can re-sample appropriate features that fit to the given point positions or vice versa, resulting in generating new samples with one modality fixed. Results for the appearance-only generation are shown in Fig. 3a and for the shape-only generation in Fig. 3b. It can be seen that our method succeeds in generating diverse novel shapes or appearances when one modality is fixed. Note that our method performs an actual recombination and does more than retrieval of objects from the training dataset.

Previous approaches that are able to generate disentangled 3D shape and appearance are GRAF [34] and Disentangled3D (D3D) [39], which are both GAN-based. We provide a quantitative comparison to these approaches regarding generation quality in Tab. 1 and a qualitative comparison in Fig. 4. Our comparisons show that our proposed method is capable of disentangled generation with superior quality than these previous approaches.

Model	ShapeNet SRN				PhotoShape	
	Cars		Chairs		Chairs	
	FID↓	KID/ $10^{-3}$ ↓	FID↓	KID/ $10^{-3}$ ↓	FID↓	KID/ $10^{-3}$ ↓
GRAF [34]	40.95	19.15	37.19	17.85	34.49	17.13
D3D [39]	62.34	41.60	45.73	24.33	59.80	36.07
NPCD (Ours)	<b>28.38</b>	<b>17.62</b>	<b>9.87</b>	<b>3.62</b>	<b>14.45</b>	<b>5.40</b>

Table 1. Comparison to disentanglement-capable approaches. The numbers show that we clearly outperform previous generative models that allow disentangled generation.

### 4.4. 3D diffusion comparison

We compare NPCD with Functa [10], SSDNeRF [8], and DiffRF [25], previous works that generate 3D shape and appearance on medium-scale datasets with 3D diffusion mod-

els. We compare against SSDNeRF and Functa on SRN Cars and against DiffRF on Photoshape Chairs. Quantitative results are provided in Tab. 2. Our proposed method performs better than Functa and DiffRF methods and worse than SSDNeRF regarding the FID and KID metrics. However, none of these methods enable disentangled generation.

Model	PhotoShape Chairs		SRN Cars	
	FID↓	KID/ $10^{-3}$ ↓	FID↓	KID/ $10^{-3}$ ↓
DiffRF	15.95	7.93	-	-
NPCD (Ours)	<b>14.45</b>	<b>5.40</b>	28.38	17.62

Table 2. Comparison to 3D diffusion models for unconditional 3D shape and appearance generation. Our NPCD model achieves better scores than DiffRF and Functa. SSDNeRF performs slightly better. However, none of the other models enable disentangled generation.

### 4.5. Shape-only comparison

Since our neural radiance field is build on top of a coarse point cloud, we evaluate the geometry of NPCD samples by comparing with the state of the art in point cloud generation in Tab. 3. Even though the point cloud defines only the coarse structure beneath a fine radiance field, the quality and diversity of our generated point clouds is comparable to the ones from approaches that are specialized for shape-only generation.

Model	SRN Cars		SRN Chairs	
	CD↓	EMD↓	CD↓	EMD↓
r-GAN [1]	94.46	99.01	83.69	99.70
PointFlow [43]	58.10	56.25	62.84	60.57
SoftFlow [17]	64.77	60.09	59.21	60.05
DPF-Net [19]	62.35	54.48	62.00	58.53
Shape-GF [4]	63.20	56.53	68.96	65.48
PVD [47]	<b>54.55</b>	53.83	56.26	53.32
LION [45]	<b>53.41</b>	<b>51.14</b>	<b>53.70</b>	<b>52.34</b>
NPCD (Ours)	60.23	<b>52.41</b>	60.50	58.84

Table 3. Shape-only comparison. We evaluate the point cloud generation part of our approach individually. Despite being just the coarse structure of a finer radiance field on top, NPCD can compete with the state of the art in point cloud generation.

### 4.6. Analysis

As diffusion on hybrid point clouds and local radiance fields has not been done before, we conduct ablation studies and analyze various novel design choices. Here, we analyze the effects of different initialization strategies, feature dimensionality and regularization methods in the generalizable



Figure 4. **Comparison against previous generative models that allow disentangled generation.**: While we present the first diffusion model allowing disentangled generation, earlier works are GAN-based. It can be seen that our model generates examples in much higher quality, as also evident from the metrics in Tab. 1.

Point-NeRF and diffusion model. We provide a more detailed analysis in the supplementals.

**Neural point cloud initialization** Regarding the initialization of the neural point cloud features, we analyze initialization with features sampled from a Gaussian distribution against a zero initialization. Interestingly, we find that these different initializations strongly affect the feature space of the trained models, which directly translates to differences in generation quality (c.f. supplementals). Tab. 4 indicates that the simple measure of zero initialization is able to largely mitigate the many-to-one mappings in the MLP decoder and provide much more coherent features.

**Neural point cloud regularization** As we assume that the structure of the feature space strongly affects the diffusion model, we analyze the effects of applying KL and TV regularizations to the neural point features during the generalizable Point-NeRF training. Tab. 4 shows that both regularizations further decrease the ambiguity of the latent space over the zero initialization. We can summarize that a combination of TV and KL regularization provides overall the best results (c.f. supplementals). Overall, we found an appropriate initialization and regularization to be key ingredients for successful neural point cloud diffusion.

## 5. Conclusion

We presented a diffusion approach for neural point clouds with high-dimensional features that represent local parts of objects and can be rendered to images. We have shown that, in contrast to all other current 3D diffusion models, the neural point cloud enables our method to effectively

Init.	Reg.	$\lambda$	Cosine sim.
Rand.	$\times$	-	0.0306
Zero	$\times$	-	0.7695
Zero	TV	3.5e-6	0.9355
Zero	KL	1e-6	0.9480
Zero	TV,KL	3e-7,1e-7	0.9470

Table 4. **Auto-decoded feature similarity.** We compute per-point mean cosine similarities between optimized neural point features of 10 training examples for 100 different seeds. Zero initialization and regularization effectively reduce the ambiguity of the auto-decoded latent space. This improves generation quality significantly, as shown in the supplementals.

*disentangle* the coarse shape from the fine geometry and appearance and to control both factors separately. In experiments on ShapeNet and PhotoShape we could show that we clearly outperform previous methods that perform disentangled generation of 3D objects, while being competitive with unconditional generation methods. Further, we provide an extensive analysis and identified many-to-one mappings in auto-decoded latent spaces as the main challenge for the successful training of a latent diffusion model. Therefore, we proposed a suitable initialization and regularization of the neural point features as effective countermeasures.

The research leading to these results is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under the project numbers 401269959 and 417962828.

## References

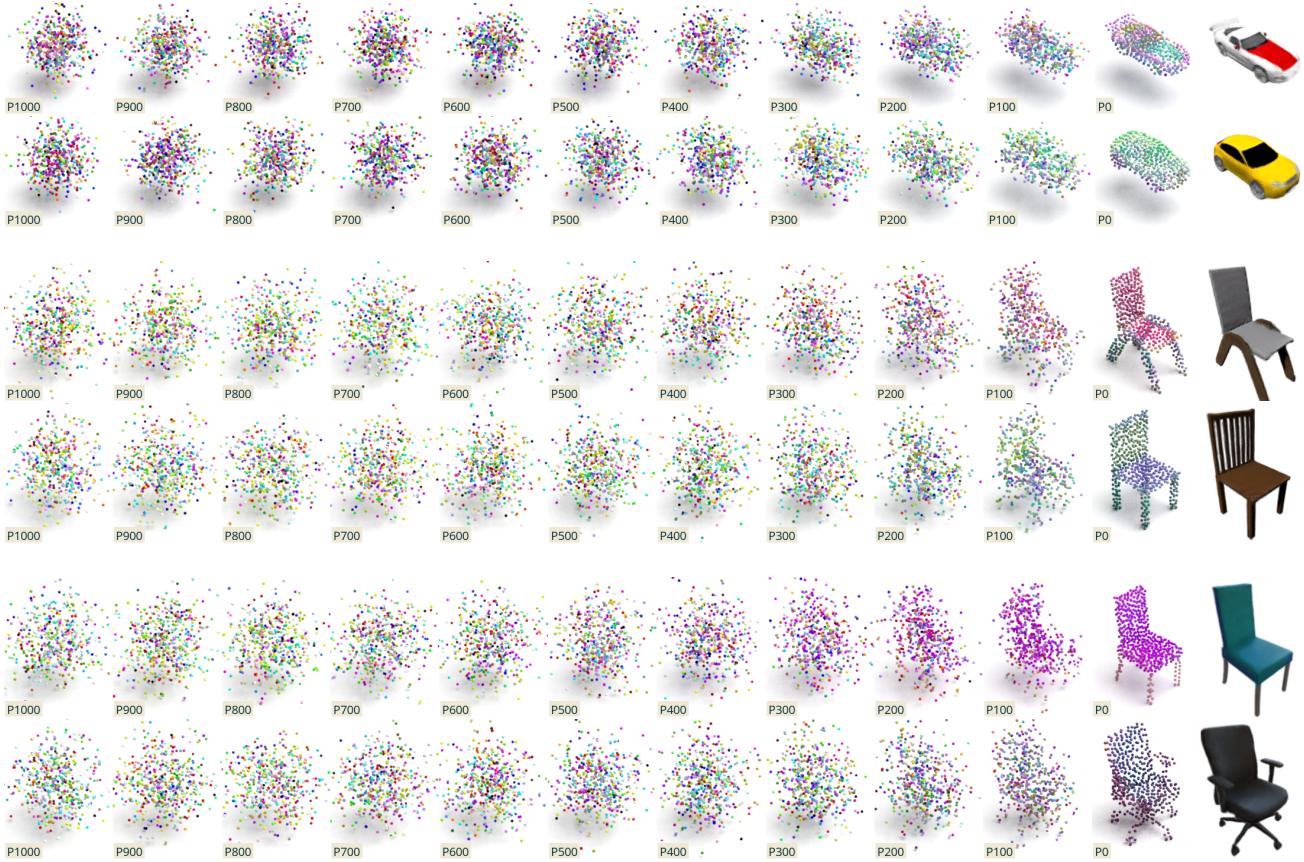
- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3D point clouds. In *ICML*, 2018. 7
- [2] Titas Auciukevičius, Zexiang Xu, Matthew Fisher, Paul Henderson, Hakan Bilen, Niloy J. Mitra, and Paul Guerrero. Renderdiffusion: Image diffusion for 3d reconstruction, inpainting and generation. In *CVPR*, 2023. 2, 3
- [3] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. In *ICLR*, 2018. 6
- [4] Ruojin Cai, Guandao Yang, Hadar Averbuch-Elor, Zekun Hao, Serge Belongie, Noah Snavely, and Bharath Hariharan. Learning gradient fields for shape generation. In *ECCV*, 2020. 1, 2, 7
- [5] Rohan Chabra, Jan E. Lenssen, Eddy Ilg, Tanner Schmidt, Julian Straub, Steven Lovegrove, and Richard Newcombe. Deep local shapes: Learning local SDF priors for detailed 3D reconstruction. In *ECCV*, 2020. 2
- [6] Eric R. Chan, Connor Z. Lin, Matthew A. Chan, Koki Nagano, Boxiao Pan, Shalini De Mello, Orazio Gallo, Leonidas Guibas, Jonathan Tremblay, Sameh Khamis, Tero Karras, and Gordon Wetzstein. Efficient geometry-aware 3D generative adversarial networks. In *CVPR*, 2022. 2, 3
- [7] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An information-rich 3d model repository. Technical Report [arXiv:1512.03012](#), Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015. 5
- [8] Hansheng Chen, Jiatao Gu, Anpei Chen, Wei Tian, Zhuowen Tu, Lingjie Liu, and Hao Su. Single-stage diffusion nerf: A unified approach to 3d generation and reconstruction, 2023. 1, 2, 3, 5, 7
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. In *NeurIPS*, 2021. 1, 2
- [10] Emilien Dupont, Hyunjik Kim, S. M. Ali Eslami, Danilo Jimenez Rezende, and Dan Rosenbaum. From data to functa: Your data point is a function and you can treat it like one. In *ICML*, 2022. 1, 2, 3, 5, 7
- [11] Ziya Erkoç, Fangchang Ma, Qi Shan, Matthias Nießner, and Angela Dai. HyperDiffusion: Generating implicit neural fields with weight-space diffusion, 2023. 1
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *NeurIPS*, 2017. 6
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. 1, 2, 4, 6
- [14] Wonbong Jang and Lourdes Agapito. Codenerf: Disentangled neural radiance fields for object categories. In *ICCV*, 2021. 2
- [15] Heewoo Jun and Alex Nichol. Shap-E: Generating conditional 3d implicit functions. [arXiv:2305.02463](#), 2023. 1, 2, 3
- [16] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. In *NeurIPS*, 2022. 2
- [17] Hyeongju Kim, Hyeonseung Lee, Woo Hyun Kang, Joun Yeop Lee, and Nam Soo Kim. Softflow: Probabilistic framework for normalizing flow on manifolds. In *NeurIPS*, 2020. 7
- [18] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. In *ICLR*, 2014. 4
- [19] Roman Klokov, Edmond Boyer, and Jakob Verbeek. Discrete point flow networks for efficient point cloud generation. In *ECCV*, 2020. 7
- [20] 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 *CVPR*, 2023. 2
- [21] 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. [arXiv preprint arXiv:2303.11328](#), 2023. 2
- [22] Shitong Luo and Wei Hu. Diffusion probabilistic models for 3d point cloud generation. In *CVPR*, 2021. 1, 2
- [23] Luke Melas-Kyriazi, Iro Laina, Christian Rupprecht, and Andrea Vedaldi. Realfusion: 360deg reconstruction of any object from a single image. In *CVPR*, 2023. 2
- [24] 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, 3
- [25] Norman Müller, Yawar Siddiqui, Lorenzo Porzi, Samuel Rota Bulo, Peter Kortschieder, and Matthias Nießner. Diffrrf: Rendering-guided 3d radiance field diffusion. In *CVPR*, 2023. 1, 2, 3, 5, 7
- [26] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-E: A system for generating 3d point clouds from complex prompts. [arXiv:2212.08751](#), 2022. 1, 2, 3, 5
- [27] Michael Niemeyer and Andreas Geiger. Giraffe: Representing scenes as compositional generative neural feature fields. In *CVPR*, 2021. 2
- [28] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *CVPR*, 2020. 2
- [29] Keunhong Park, Konstantinos Rematas, Ali Farhadi, and Steven M. Seitz. Photoshape: Photorealistic materials for large-scale shape collections. *ACM TOG*, 2018. 2, 5
- [30] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. In *ICLR*, 2023. 2
- [31] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, 2021. 1
- [32] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2, 4

- [33] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *NeurIPS*, 2022. 1
- [34] Katja Schwarz, Yiyi Liao, Michael Niemeyer, and Andreas Geiger. Graf: Generative radiance fields for 3d-aware image synthesis. In *NeurIPS*, 2020. 2, 5, 7, 8
- [35] J. Ryan Shue, Eric Ryan Chan, Ryan Po, Zachary Ankner, Jiajun Wu, and Gordon Wetzstein. 3d neural field generation using triplane diffusion. In *CVPR*, 2023. 1, 2, 4
- [36] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *NeurIPS*, 2019. 2, 5
- [37] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 2, 4
- [38] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *NeurIPS*, 2019. 2, 4
- [39] Ayush Tewari, MalliKarjun B R, Xingang Pan, Ohad Fried, Maneesh Agrawala, and Christian Theobalt. Disentangled3d: Learning a 3d generative model with disentangled geometry and appearance from monocular images. In *CVPR*, 2022. 2, 5, 7, 8
- [40] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 5
- [41] Christopher Wewer, Eddy Ilg, Bernt Schiele, and Jan Eric Lenssen. Simmp: Learning self-similarity priors between neural points. In *ICCV*, 2023. 2
- [42] Qiangeng Xu, Zexiang Xu, Julien Philip, Sai Bi, Zhixin Shu, Kalyan Sunkavalli, and Ulrich Neumann. Point-NeRF: Point-based neural radiance fields. In *CVPR*, 2022. 2, 3
- [43] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge Belongie, and Bharath Hariharan. Pointflow: 3d point cloud generation with continuous normalizing flows. In *ICCV*, 2019. 6, 7
- [44] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. In *ECCV*, 2018. 2
- [45] Xiaohui Zeng, Arash Vahdat, Francis Williams, Zan Gojcic, Or Litany, Sanja Fidler, and Karsten Kreis. LION: Latent point diffusion models for 3d shape generation. In *NeurIPS*, 2022. 1, 2, 4, 7
- [46] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 1
- [47] Linqi Zhou, Yilun Du, and Jiajun Wu. 3d shape generation and completion through point-voxel diffusion. In *ICCV*, 2021. 1, 2, 7
- [48] Zhizhuo Zhou and Shubham Tulsiani. Sparsefusion: Distilling view-conditioned diffusion for 3d reconstruction. In *CVPR*, 2023. 2
- [49] Jun-Yan Zhu, Zhoutong Zhang, Chengkai Zhang, Jiajun Wu, Antonio Torralba, Josh Tenenbaum, and Bill Freeman. Vi-
- sual object networks: Image generation with disentangled 3d representations. In *NeurIPS*, 2018. 2

# Neural Point Cloud Diffusion for Disentangled 3D Shape and Appearance Generation – Appendix –

## A1. Visualization of the neural point cloud diffusion process

Figure A1 shows a visualization of the neural point cloud diffusion process for unconditional generation on ShapeNet Cars, ShapeNet Chairs, and PhotoShape Chairs.



**Figure A1. Visualization of the neural point cloud diffusion process.** We generate the shape and appearance of 3D objects on ShapeNet Cars, ShapeNet Chairs, and PhotoShape Chairs with the proposed Neural Point Cloud Diffusion (NPCD) model. We visualize the neural point clouds  $\mathcal{P}_t = (\mathbf{P}_t, \mathbf{F}_t)$  from intermediate timesteps  $t$  of the diffusion process. In total, the diffusion process of NPCD has 1000 timesteps and we visualize every 100th timestep. The features of the neural point clouds are visualized by taking the first three PCA components as RGB color. The last visualized neural point cloud  $\mathcal{P}_0$  represents the final generated 3D object. Additionally, we visualize a Point-NeRF rendering of the final neural point cloud.

## A2. Implementation details

### A2.1. Generalizable Point-NeRF

**Architecture** For the aggregation MLP  $F_\phi$ , we use 4 linear layers with a hidden dimension of 256, each followed by a LeakyReLU, and an output projection linear layer that maps to 256d. For the color MLP  $G_\psi$ , we use 4 linear layers with a hidden dimension of 256, each followed by a LeakyReLU, and an output projection linear layer that maps to 3d. For the density MLP  $H_\gamma$ , we use 1 linear layer with a hidden dimension of 256, followed by a LeakyReLU, and an output projection linear layer that maps to 1d.

**Training parameters** We construct training samples by splitting the available views per object into groups of 5 views per sample. In each iteration, we use 8 samples and 112 pixels of each view in each sample. The image reconstruction loss is hence optimized for an effective batch size of  $8 \cdot 50 \cdot 112 = 44800$  pixels. For the volumetric rendering of each ray, we sample 128 shading points. In all experiments, we train Point-NeRF for 1000 epochs. We use the Adam optimizer with a constant learning rate of  $1e-3$ .

### A2.2. Diffusion model

**Architecture** For the denoiser network of the diffusion model, we use a standard transformer architecture with 24 layers, a feature dimension of 1024d and 16 heads [3]. This architecture has ca. 300M parameters.

**Diffusion model parameters** For the diffusion model, we use the linear noise schedule from DDPM [2] with 1,000 steps and  $\beta$  ranging from 0.0001 to 0.02. We normalize the neural point clouds such that the positions are unit Gaussian distributed and the features are in the range  $[-1, 1]$  (and apply the inverse transform later before rendering the representation). During sampling, we clip the coordinates and features to the respective minimum and maximum values of the training dataset.

**Training parameters** We train the diffusion model for 1M iterations with a batch size of 32. We train on a single GPU with 16 bit and employ flash attention [1]. We use an exponential moving average over the model parameters with a decay of 0.9999. On ShapeNet Cars, we use a constant learning rate of  $7e-5$  and on ShapeNet Chairs and PhotoShape Chairs, we use a constant learning rate of  $4e-5$ . On all datasets, we use a weight decay of 0.01.

### A3. Analysis

As described in the main paper, we conduct ablations studies regarding the effects of different initialization strategies, feature dimensionality and regularization methods in the generalizable Point-NeRF and diffusion model. These parameters affect both, the quality of the reconstructions and the structure of the neural point cloud feature space. Both properties in turn affect the diffusion model that is trained on the resulting neural point clouds. As the reconstructed objects serve as training data for the diffusion model, the reconstruction quality likely is the upper bound of the generation quality. On the other hand, the structure of the feature space affects how well the data distribution can be learned by the diffusion model.

#### A3.1. Setup

We conduct the analyses with the same settings as described in the main paper. The only difference is that, for computational reasons, we conduct the analyses with a smaller transformer denoiser network with 40M parameters. Thus, the numbers of this configuration might differ from the best configuration in the main paper. At the end of the analysis, we compare this 40M parameter model to the 300M parameter model from the main paper for the best initialization strategy, dimensionality, and regularization parameters. We conduct the analyses on ShapeNet Cars and Chairs. In case the results on Cars are very clear, we omit the corresponding experiments on Chairs.

#### A3.2. Neural point cloud initialization

Regarding the initialization of the neural point cloud features, we analyze initialization with features sampled from a Gaussian distribution against a zero initialization. Interestingly, we find that these different initializations strongly affect the feature space of the trained models. To illustrate this, we visualize neural point features of reconstructed objects in Fig. A2.

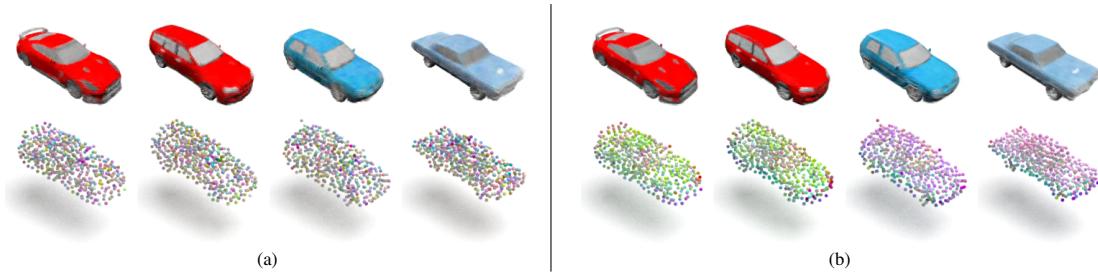


Figure A2. The first row shows the generalizable Point-NeRF renderings of four *reconstructed* training objects. The second row shows a visualization of the features from the point clouds by taking the first three PCA components as RGB colors. **(a) Random initialization:** The features learned starting from a random initialization are distributed randomly across an object and differ across objects with the same appearance. **(b) Zero initialization:** Features learned starting with a zero initialization are coherent within an object and across objects with the same appearance.

This effect is measured quantitatively via the cosine similarities in Tab. 4 of the main paper. As shown in Fig. A3 and Tab. A1a, the more coherent features from the zero initialization, are vital to enable the successful training of a diffusion model.



Figure A3. Generated samples from diffusion models trained on neural point clouds from generalizable Point-NeRFs that were optimized with different initialization strategies: **(a) Random initialization** leads to many artifacts in the appearance of generated samples. **(b) Zero initialization** leads to more diversity and fewer artifacts.

Setting	Dim.	Init.	Reg.	$\lambda$	ShapeNet SRN Cars						ShapeNet SRN Chairs					
					PSNR↑	FIDrec↓	KIDrec↓	FID↓	KID↓	PSNR↑	FIDrec↓	KIDrec↓	FID↓	KID↓		
<b>a) Initialization</b>																
Random initialization	32	Rand.	$\times$	-	29.24	37.24	25.37	125.51	97.82	-	-	-	-	-	-	
Zero initialization	32	Zero	$\times$	-	31.32	18.96	11.17	53.55	35.37	34.91	10.37	4.85	39.19	23.64	-	
<b>b) Dimensionality</b>																
16D features	16	Zero	$\times$	-	30.60	22.56	13.86	56.02	39.30	-	-	-	-	-	-	
32D features	32	Zero	$\times$	-	31.32	18.96	11.17	53.55	35.37	34.91	10.37	4.85	39.19	23.64	-	
128D features	128	Zero	$\times$	-	32.65	19.33	11.60	73.93	52.09	-	-	-	-	-	-	
<b>c) Regularization</b>																
No regularization	32	Zero	$\times$	-	31.32	18.96	11.17	53.55	35.37	34.91	10.37	4.85	39.19	23.64	-	
TV regularization	32	Zero	TV	3.5e-6	29.72	22.42	13.71	45.90	28.70	32.38	14.10	6.70	32.87	17.49	-	
KL regularization	32	Zero	KL	1e-6	30.02	24.93	15.60	55.01	35.86	34.20	8.37	3.17	18.13	8.17	-	
TV+KL regularization	32	Zero	TV, KL	3e-7, 1e-7	29.70	26.12	16.44	43.92	26.53	33.62	8.58	3.34	17.17	7.44	-	
<b>d) Model size</b>																
40M parameters	32	Zero	TV, KL	3e-7, 1e-7	29.70	26.12	16.44	43.92	26.53	33.62	8.58	3.34	17.17	7.44	-	
300M parameters	32	Zero	TV, KL	3e-7, 1e-7	29.70	26.12	16.44	28.38	17.62	33.62	8.58	3.34	9.87	3.62	-	

Table A1. **Analysis of the generalizable Point-NeRF reconstructions and the diffusion model generations regarding:** **a)** feature initialization, **b)** feature dimensionality, **c)** feature regularization ( $\lambda$  is the weight of the regularization loss) and **d)** model size. The PSNR, FIDrec and KIDrec metrics in the gray columns measure the quality of the reconstructions in the generalizable Point-NeRF optimization stage. The FID and KID metrics measure the quality of the diffusion model generations. The reported KID is multiplied with  $10^3$ . Zero initialization clearly outperforms random initialization. 32D features outperform 16D and 128D features. Combined TV+KL regularization outperforms having no regularization or TV or KL regularization alone.

### A3.3. Neural point cloud feature dimension

We analyze 16D, 32D and 128D features for the neural point clouds in Tab. A1b and Fig. A4.

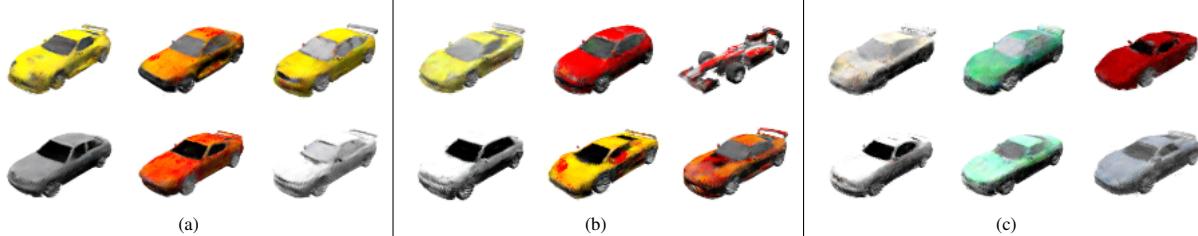


Figure A4. Generated samples from diffusion models trained on neural point clouds with different numbers of feature dimensions: **(a)** 16D features, **(b)** 32D features, **(c)** 128D features. We observe a clear difference between 16 and 32D features, while the difference to 128D is small.

As indicated by the PSNR metrics, higher feature dimensions allow better training data reconstructions in the generalizable Point-NeRF training. However, the FID and KID metrics for the generation performance of the diffusion model decrease for 128D. As the visual quality of 32D and 128D in Fig. A4 is very similar and the FID and KID metrics for 32D are better, we choose to continue with the 32D features.

### A3.4. Neural point cloud regularization

As stated in the main paper, we analyze the effects of applying KL and TV regularizations to the neural point features during the generalizable Point-NeRF training. The cosine similarities in Tab. 4 of the main paper shows that both regularizations further decrease the ambiguity of the latent space over the zero initialization. However, we observe different behaviors w.r.t. quantitative and qualitative results. On the one hand, among TV and KL regularization, TV leads to the better FID and KID scores in Tab. A1c. The qualitative comparison in Fig. A5 on the other hand shows that KL regularization results in

cleaner samples with less artifacts compared to TV regularization. As a consequence, we try a combination of TV and KL regularization, which leads to an equally high visual quality and the best FID and KID scores.



Figure A5. Generated samples from diffusion models trained on neural point clouds from Point-NeRFs that were optimized with different regularization strategies: (a) No regularization (FID 53.55), (b) TV regularization (FID 45.90), (c) KL regularization (FID 55.01), (d) TV+KL regularization (FID 43.92). TV regularization increases performance regarding the FID and KID metrics. KL regularization leads to cleaner qualitative results. The proposed method NPCD uses a combination of both regularizations, which improves quantitative and qualitative results.

### A3.5. Model size

Lastly, we compare the performance of the transformer model with 40M parameters that was used in the analyses, with the performance of the transformer model with 300M parameters, that was used in the main paper. The architectures of these models are as follows:

- 40M parameter model: 12 layers, hidden dimension 512, 8 heads; trained with batch size 64 and learning rate  $1e-4$ .
- 300M parameter model: 24 layers, hidden dimension 1024, 16 heads; trained with batch size 32 and learning rate  $7e-5$  on Cars and  $4e-5$  on Chairs.

The quantitative comparison in Tab. A1d shows that the quality of the generated samples scales with the model size.

#### A4. Qualitative results for unconditional generation on ShapeNet Cars

Figure A6 shows unconditional generations of the proposed NPCD model trained on ShapeNet Cars.

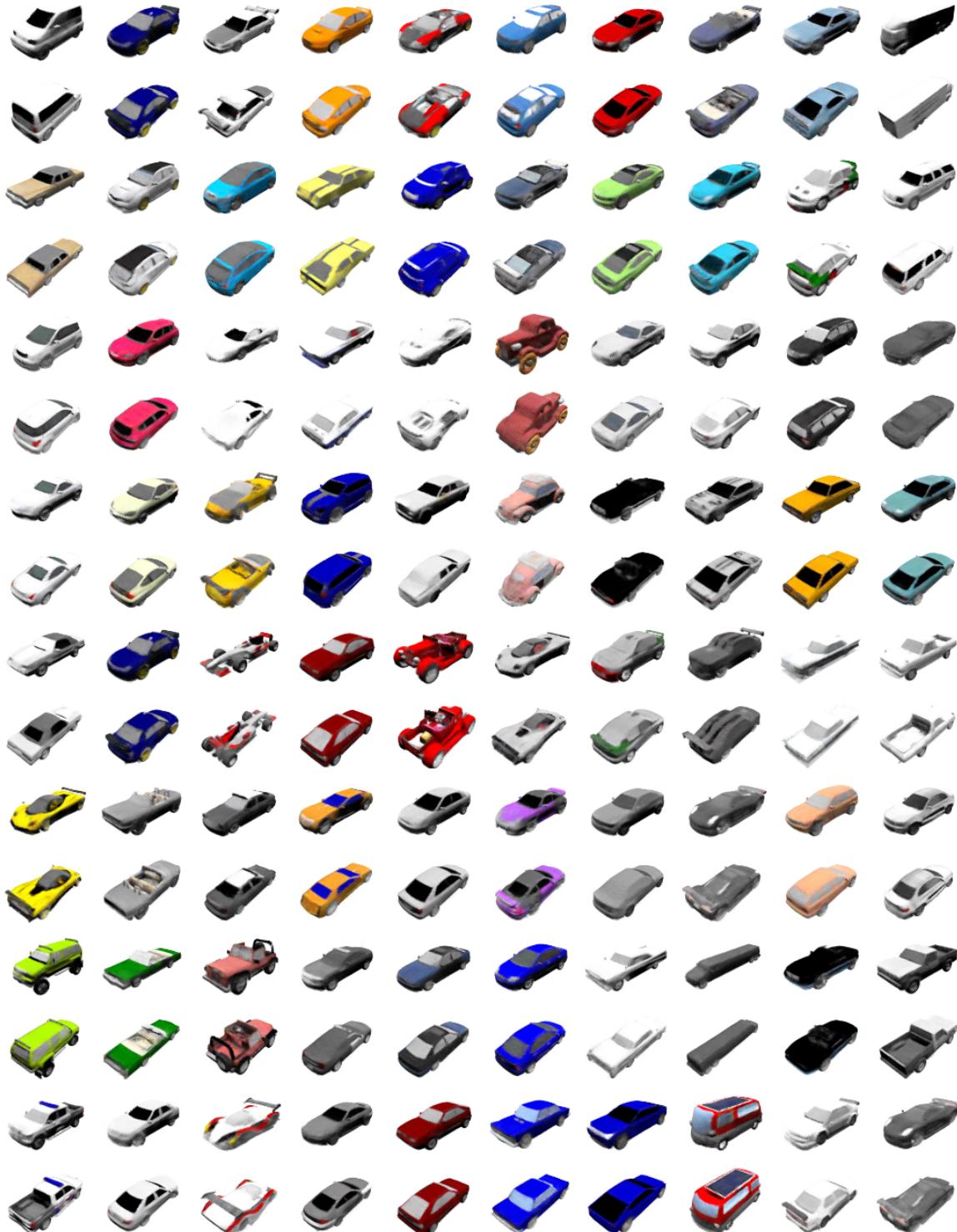


Figure A6. Unconditional generations from the proposed NPCD model trained on ShapeNet Cars. Each generated object is visualized from two different viewpoints. Note that the shown objects are not cherry-picked.

## A5. Qualitative results for unconditional generation on ShapeNet Chairs

Figure A7 shows unconditional generations from the proposed NPCD model trained on ShapeNet Chairs.



Figure A7. Unconditional generations of the proposed NPCD model trained on ShapeNet Chairs. Each generated object is visualized from two different viewpoints. Note that the shown objects are not cherry-picked.

## A6. Baseline details

For our comparisons on ShapeNet Cars, ShapeNet Chairs, and PhotoShape Chairs, we retrain GRAF and Disentangled3D on these datasets. For training, both approaches require a dataset of images, the distribution of camera poses  $p_\xi$  that correspond to the images, and a known camera matrix  $\mathbf{K}$ . Further, the volumetric rendering in both approaches requires given near and far clipping planes. Training is conducted in a GAN framework. The generator renders images for randomly sampled shape and appearance latent codes and randomly sampled camera poses. The discriminator compares generated and real images.

**ShapeNet Cars** For training on ShapeNet Cars, according to the the SRN rendering parameters, we set the radius to 1.3, the field of view to 52 (chosen such that it results in the correct focal length), and sample camera poses from the full hemisphere. To define the near and far planes, we compute the cube that bounds all pointclouds. For the Cars training split, this gives:

- x-coordinate range:  $-0.30903995$  to  $0.30898672$
- y-coordinate range:  $-0.48859358$  to  $0.49043328$
- z-coordinate range:  $-0.27073205$  to  $0.2709335$

The nearest possible point hence has the following distance to the camera:

$$1.3 - \sqrt{0.30903995^2 + 0.49043328^2 + 0.2709335^2} = 1.3 - 0.6398714 = 0.6601285815238953.$$

The furthest point hence has the following distance to the camera:  $1.3 + 0.6398714 = 1.9398714184761048$ .

Based on this, we set the near and far planes to 0.5 and 2.1.

**ShapeNet Chairs** For training on ShapeNet Cars, according to the the SRN rendering parameters, we set the radius to 2.0, the field of view to 52 (chosen such that it results in the correct focal length), and sample camera poses from the full hemisphere. To define the near and far planes, we compute the cube that bounds all pointclouds. For the Chairs training split, this gives:

- x-coordinate range:  $-0.5$  to  $0.5$
- x-coordinate range:  $-0.5$  to  $0.5$
- x-coordinate range:  $-0.5$  to  $0.5$

The nearest possible point hence has the following distance to the camera:  $2.0 - \sqrt{0.5^2 + 0.5^2 + 0.5^2} = 2.0 - 0.87 = 1.13$ .

The furthest point hence has the following distance to the camera:  $2.0 + 0.87 = 2.87$ .

Based on this, we set the near and far planes to 1.0 and 3.0.

**PhotoShape Chairs** For training on PhotoShape Chairs, according to the PhotoShape Chairs rendering parameters, we set the radius to 2.5 and the field of view to 39.6 (chosen such that it results in the correct focal length). As PhotoShape Chairs views were always rendered from the same set of camera poses, we randomly sample poses from this discrete set during training. As the 3D object shapes are the same as on ShapeNet Chairs, we compute the near and far clipping planes comparably to ShapeNet Chairs. Given the different radius, this results in near and far planes of 1.25 and 3.75.

## References

- [1] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *NeurIPS*, 2022. [2](#)
- [2] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, 2020. [2](#)
- [3] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-E: A system for generating 3d point clouds from complex prompts. [arXiv:2212.08751](#), 2022. [2](#)