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





Diffusion for 3D Shape and Appearance Generation:

- Overview • Input: dataset of optimized neural point clouds Diffusion model is trained jointly on point positions and appearance features $\mathbf{P}_0, \mathbf{F}_0$ Disentangled Generation (here described for appearance-only generation) Appearance-only generation with given point positions \mathbf{P}_0 works by "masking" point positions P100 P0 and generating only features: 1: $\epsilon^{\mathbf{P}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad \mathbf{P}_T = \sqrt{\bar{\alpha}_T} \mathbf{P}_0 + \sqrt{1 - \bar{\alpha}_T} \epsilon^{\mathbf{P}} \quad \mathbf{F}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 2: for t = T, ..., 1 do $(\epsilon_{\theta}^{\mathbf{P}}, \epsilon_{\theta}^{\mathbf{F}}) = T_{\theta}((\mathbf{P}_t, \mathbf{F}_t), t)$ 4: $\mathbf{P}_{t-1} = \sqrt{\bar{\alpha}_{t-1}}\mathbf{P}_0 + \sqrt{1-\bar{\alpha}_{t-1}}\epsilon^{\mathbf{P}}$ 5: $\mathbf{F}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{F}_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha_t}}} \epsilon_{\theta}^{\mathbf{F}} \right) + \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t \epsilon \quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 6: **end for** \rightarrow Intuitively like masked image inpainting, similar to RePaint Disentangled Generation Results (on ShapeNet SRN Cars, SRN Chairs, PhotoShape Chairs) **Appearance-only generation:** Re-sample coarse shape only _____ Re-sample local appearance only I I I I I I I Neural Point Cloud Representation and Point-NeRF Rendering **Represent 3D shape and appearance of object by** $\mathbf{I}_i \subset \mathbb{K}$ a neural point cloud: I I I I I I I I $\mathcal{P} = (\mathbf{P}, \mathbf{F})$ $\to M$ neural points with \blacksquare positions $\mathbf{P} \in \mathbb{R}^{M \times 3}$ and \blacksquare fea- $\in \mathbb{R}^3$ tures $\mathbf{F} \in \mathbb{R}^{M imes D}$ MLP \longrightarrow MLP $G_{\psi} \rightarrow \mathbf{c}$ Can be rendered to images with a Point-NeRF: Like NeRF, but compute colors ${f c}$ and densities σ from neighbouring neural points using MLPs F_{ϕ} and H_{γ} and G_{ψ} **Point-NeRF Autodecoder Optimization** (on N objects from same category) **Optimization objective: Initialization and Regularization:** • Crucial to avoid many-to-one mappings • Optimize point features F and between features and renderings from same category **MLP** parameters ϕ, ψ, γ • We use: **Zero initialization**, **TV regular**-• Objective: Minimize image reization, KL regularization

This work \rightarrow apply diffusion to a neural point cloud 3D representation:

P1000 P900	P800	P700	P600	P500 ·	P400	P300	P200	P1

- Goal: learn distribution \rightarrow generate 3d shape and appearance of new objects • Previous works use different 3D representations: latent codes (Functa), triplanes (SSD-NeRF), voxel grids (DiffRF) • Points have: **3D position + higher-dimensional feature** • Can be rendered to images with Point-NeRF • Figure shows diffusion on neural point cloud + Point-NeRF rendering of the final generated neural point cloud Point positions represent shape, features represent appearance \rightarrow This enables disentangled generation: • Dataset of N objects • **Object**: consists of a





Given:

- neural point cloud $\mathcal{P}_i =$ $(\mathbf{P}_{i}, \mathbf{F}_{i})$ and K images with camera parameters

- construction error for all images of all objects

Philipp Schröppel

Christopher Wewer

Jan Eric Lenssen Eddy IIg

schroepp@cs.uni-freiburg.de

Diffusion Model Training

Thomas Brox





	Estimate nois					
	Compute point positions via forward process					
.)	Dpdate point features via reverse process					

Shape-only generation:

and the second					and the S	
		P	Contra to		Contraction of the second seco	
			Contraction of the second seco		Contraction of the second seco	0000
- and						0
- The second sec				JA -	6	
U				5		
		A				
	JANK AND			R		A
		2				
		P		- a		
	×				Y	

Comparison to disentanglement-capable approaches:

Model		ShapeN	let SR	N	PhotoShape		
	Cars		Chairs		Chairs		
	FID↓	$KID/10^{-3}\downarrow$	FID↓	$KID/10^{-3}\downarrow$	FID↓	$KID/10^{-3}\downarrow$	
GRAF	40.95	19.15	37.19	17.85	34.49	17.13	
D3D	62.34	41.60	45.73	24.33	59.80	36.07	
NPCD (Ours)	28.38	17.62	9.87	3.62	14.45	5.40	

 \rightarrow our approach outperforms previous approaches

Initialization: a) Random init.

Regularization:



Setting

a) Initialization Random initialization Zero initialization

b) Regularization No regularization TV regularization KL regularization TV+KL regularization

c) Model size 40M parameters 300M parameters

- for the application of diffusion models

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





Evaluation of Unconditional Generation Quality

Comparison to 3D diffusion approaches:

Model PhotoShap FID↓ KI[Shape Chairs KID/10 ⁻³	Model	$\begin{vmatrix} SRN Cars \\ FID \downarrow KID/10^{-3} \downarrow \end{vmatrix}$			
DiffRF	15.95	7.93	Functa SSDNeRF	80.3 11.08	- 3.47		
NPCD (Ours)	14.45	5.40	NPCD (Ours)	28.38	17.62		

 \rightarrow Our approach performs better than Functa and DiffRF and worse than SSDNeRF, but enables disentangled generation

Analyses on Initialization and Regularization



Many-to-one mappings:

 \rightarrow how far are features that represent the same appearance away from each other?

 \rightarrow We measure this by computing per-point mean cosine similarities between optimized neural point features of 10 training examples for 100 different seeds:

Init.	Reg.	Cosine sim.
Rand.	×	0.0306
Zero	×	0.7695
Zero	TV	0.9355
Zero	KL	0.9480
Zero	TV,KL	0.9470

 \rightarrow zero initialization and regularizations reduce many-to-one mappings (as shown by Figures on the left and Table above) and improve generation quality (as shown by Table below)

Init.	Reg.	ShapeNet SRN Cars					ShapeNet SRN Chairs				
	C	PSNR ↑	FIDrec↓	KIDrec↓	FID↓	KID↓	PSNR ↑	FIDrec↓	KIDrec↓	FID↓	KID↓
Rand.	×	29.24	37.24	25.37	125.51	97.82	-	_	_	_	-
Zero	×	31.32	18.96	11.17	53.55	35.37	34.91	10.37	4.85	39.19	23.64
Zero	×	31.32	18.96	11.17	53.55	35.37	34.91	10.37	4.85	39.19	23.64
Zero	TV	29.72	22.42	13.71	45.90	28.70	32.38	14.10	6.70	32.87	17.49
Zero	KL	30.02	24.93	15.60	55.01	35.86	34.20	8.37	3.17	18.13	8.17
Zero	TV,KL	29.70	26.12	16.44	43.92	26.53	33.62	8.58	3.34	17.17	7.44
Zero	TV,KL	29.70	26.12	16.44	43.92	26.53	33.62	8.58	3.34	17.17	7.44
Zero	TV,KL	29.70	26.12	16.44	28.38	17.62	33.62	8.58	3.34	9.87	3.62

Summary

 \rightarrow We apply diffusion for 3D shape and appearance generation to a neural point cloud representation

 \rightarrow This enables disentangled generation, as point positions represent coarse shape and features represent appearance

 \rightarrow We provide analyses and insights into many-to-one mappings in auto-decoded latent spaces and how to tame them

Project page and code: https://neural-point-cloud-diffusion.github.io/

Acknowledgements

