



浙江大学
ZHEJIANG UNIVERSITY

基于神经表示的三维重建与生成

Fast Inference & Training and 3D-Aware Generative Models

廖依依

VALSE 2022 Tutorial

Outline

- ▶ Real-time rendering of neural radiance fields
- ▶ Fast training of neural scene representation
- ▶ 3D-aware generative model

rendered in real-time on NVIDIA GTX 1080 Ti*



*affordable consumer GPU

Fast Inference

Key Idea: Reduce FLOPs of volume rendering

FLOPs: $H \times W \times K \times L_r$ K : Samples per ray, L_r : FLOPs of color/density retrieval

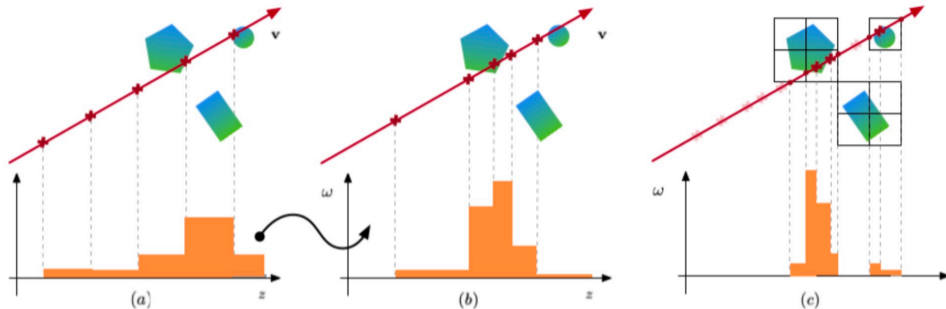
▶ Reducing K

- ▶ Classical techniques: Early ray termination, empty space skipping
- ▶ Adaptive sampling

▶ Reducing L_r

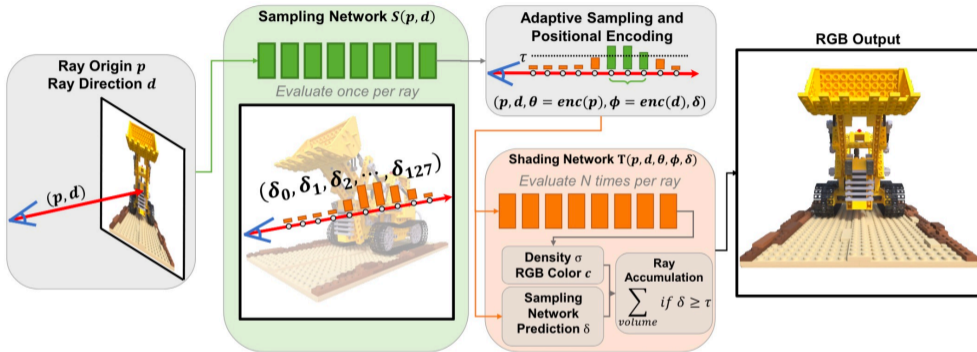
- ▶ Tabulation-based methods
- ▶ Smaller networks

Reducing K : Empty space skipping & Early ray termination



- ▶ ESS: Skipping a sample point \mathbf{x} if \mathbf{x} lies in an unoccupied cell
- ▶ ERT: Skipping $\mathbf{x}_{i+1}, \mathbf{x}_{i+2}$ if transmittance $T_i < \epsilon$

Reducing K : Adaptive sampling



- ▶ Learning where to sample via a **sampling network**
- ▶ Fine-tune shading network to desired sample counts (2,4,8,16) for fast rendering

Reducing L_r : Tabulation-based Methods

Naïve Solution:

1. Train a large MLP $f(\mathbf{x}, \mathbf{d})$
2. save network's output for fast inference

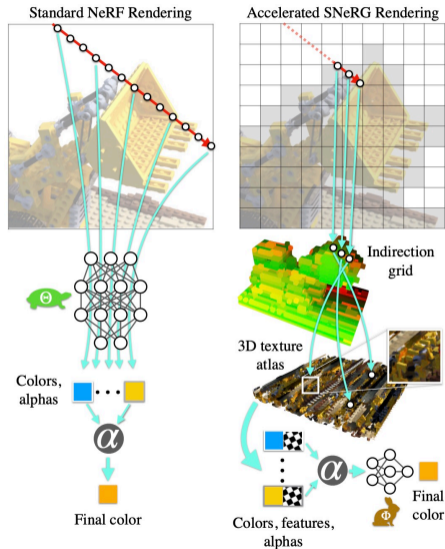
Challenge: Naïve solution requires memory of $O(N^5)$ given 5D input (\mathbf{x}, \mathbf{d})

Key Idea: Modeling **view dependency** \mathbf{d} differently to reduce memory to $O(N^3)$

- ▶ SNeRG, ICCV 2021
- ▶ PlenOctree, ICCV 2021
- ▶ FastNeRF, ICCV 2021

SNeRG

- ▶ Replacing slow MLP evaluations with lookups into **cached sparse 3D grid**
- ▶ Caching **view-independent** colors, features and alphas
- ▶ Model **view dependency** via a small network
- ▶ Combine with ESS + ERT





(a) Frame rendered by our method.



(b) Cross-section of (a).



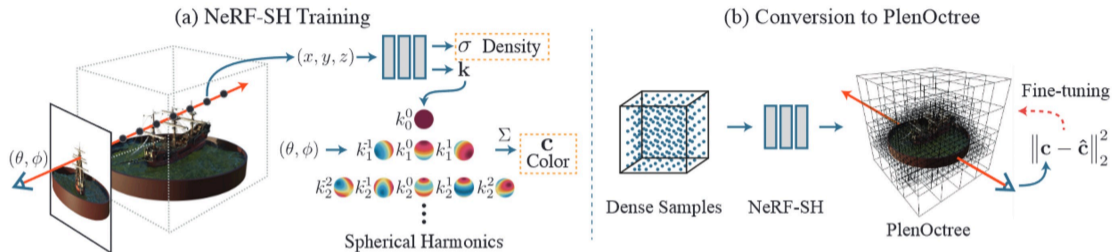
(c) Trained without \mathcal{L}_s .



(d) No \mathcal{L}_s , no visibility culling.

- ▶ $O(N^3)$ still requires **large memory consumption**
- ▶ Culling voxels where the maximum opacity is low, or maximum transmittance $T < \epsilon$ across all views

PlenOctree



- ▶ Modeling view dependency via **spherical harmonics coefficients \mathbf{k}**
- ▶ Use octree structure to skip large empty space, but still **memory costly**

Reducing L_r : Smaller networks

Naïve Solution: Replace the large MLP with a small MLP

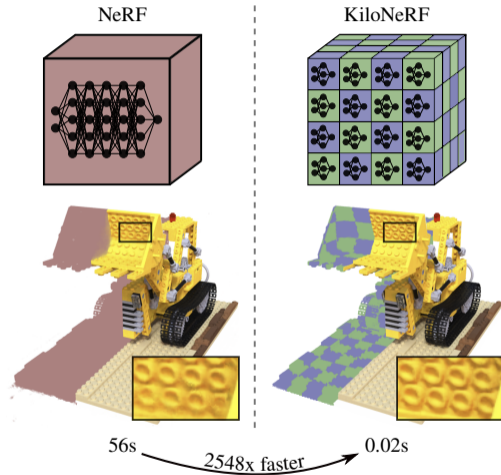
Challenge: Naïve solution leads to degraded image quality

Key Idea: Use a small network to **independently** represent a small region

- ▶ KiloNeRF: ICCV 2021
- ▶ (BlockNeRF, MegaNeRF)

KiloNeRF

- ▶ Replacing **large** MLPs with **small** ones via space partitioning
- ▶ No cache required, more **memory friendly**
- ▶ Combine with ESS + ERT



KiloNeRF

Method	Render time ↓	Speedup ↑
NeRF	56185 ms	–
NeRF + ESS + ERT	788 ms	71
KiloNeRF	22 ms	2548

Table 2: **Speedup Breakdown.** The original NeRF model combined with KiloNeRF’s implementation of ESS and ERT is compared against the full KiloNeRF technique.

- ▶ ESS/ERT and small MLPs both contribute to fast rendering

Fast Inference

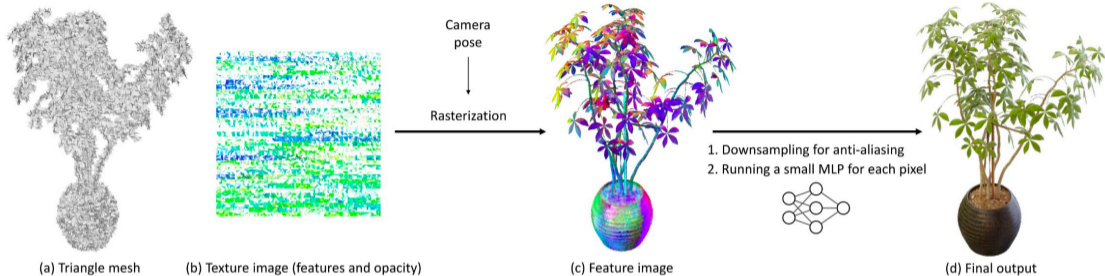
Key Idea: Reduce FLOPs of volume rendering → **image-order rendering**

FLOPs: $H \times W \times K \times L_r$ K : Samples per ray, L_r : FLOPs of color/density retrieval

- ▶ Reducing K
 - ▶ Classical techniques: Early ray termination, empty space skipping
 - ▶ Adaptive sampling
- ▶ Reducing L_r
 - ▶ Tabulation-based methods
 - ▶ Smaller networks

How about **object-order rendering?**

MobileNeRF



- ▶ Exploit fast polygon **rasterization**, using standard GPU rasterization pipeline
- ▶ Avoid semi-transparency, enabling efficient rendering without sorting
- ▶ 10× faster than SNeRG

Discussions

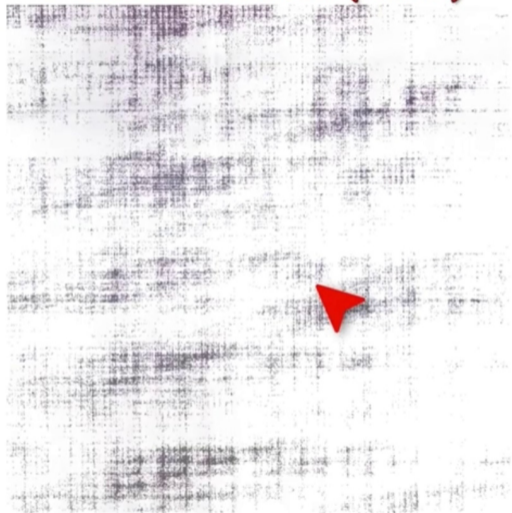
Method	Device	Dataset	FPS	Memory (MB)	PSNR
AdaNeRF*	NVIDIA 3090	DONeRF Dataset (800x800)	20	4	27.5
SNeRG	Laptop with mobile GPU		207.26	4096 (100)	30.38
FastNeRF	NVIDIA 3090		238	16200	29.97
PlenOctree	NVIDIA V100	NeRF Synthetic	167.68	1930	31.71
KiloNeRF	NVIDIA 1080TI	(800x800)	50	50	31
MobileNeRF	NVIDIA 2080TI		744.91	125.75	30.9
	iPhone XS		55.89		

- ▶ Fast rendering is usually achieved by combining **multiple techniques**
- ▶ Image-order: tabulation-based methods are more **efficient** at the cost of **memory**
- ▶ Object-order: leverage standard rendering tools

Outline

- ▶ Real-time rendering of neural radiance fields
- ▶ Fast training of neural scene representation
- ▶ 3D-aware generative model

Previous method (NeRF)



00:04

minutes : seconds

Fast Training

Key Idea: Reduce FLOPs of volume rendering

FLOPs: $H \times W \times K \times L_r$ K : Samples per ray, L_r : FLOPs of color/density retrieval

- ▶ **Reducing K**

- ▶ Empty space skipping

- ▶ **Reducing L_r**

- ▶ Directly optimize $O(N^3)$ voxel grid or optimize $O(N^3)$ local feature + small MLP

Fast Training

Key Idea: Reduce FLOPs of volume rendering

FLOPs: $H \times W \times K \times L_r$ K : Samples per ray, L_r : FLOPs of color/density retrieval

▶ Reducing K

- ▶ Empty space skipping

▶ Reducing L_r

- ▶ Directly optimize $O(N^3)$ voxel grid or optimize $O(N^3)$ local feature + small MLP

Additional benefits of optimizing local tensors?

- ▶ Each voxel or feature tensor is only responsible for a small region
- ▶ Shallow computational graph, allowing for use large learning rate

Fast Training

Challenge: Memory complexity of $O(N^3)$ still **expensive**

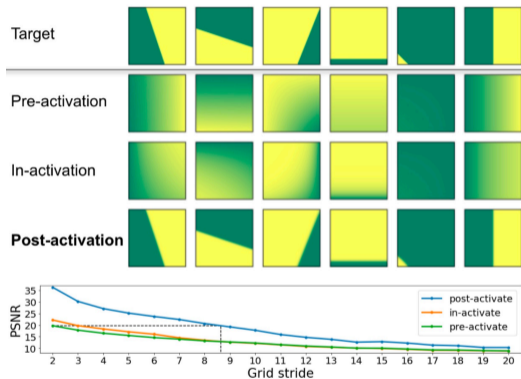
Idea: Reduce required resolution or use compact representation

- ▶ Reduce required resolution via activation design
 - ▶ DVGO
- ▶ Explore sparsity
 - ▶ Plenoxel: Sparse data structure
- ▶ Parameter sharing
 - ▶ Instant NGP: hard parameter-sharing
 - ▶ TensorRF: tensor factorization

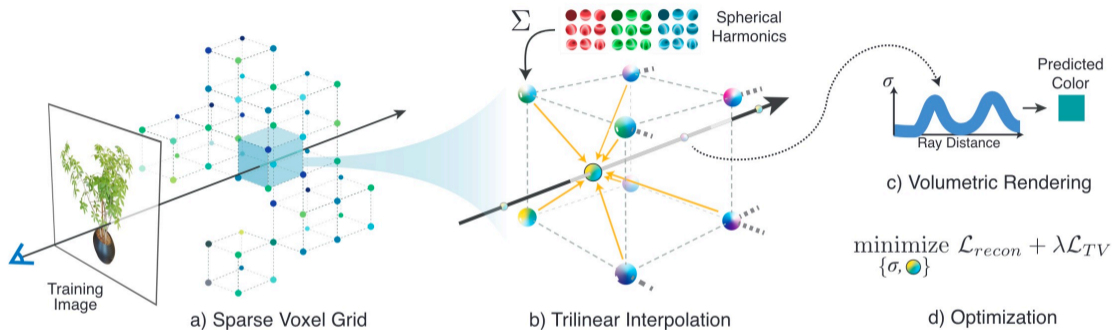
DVGO

- ▶ **Post-activation**: modeling sharp surface within a cell
- ▶ Resolution $160 \times 160 \times 160$ for all NeRF-Synthetic objects

$$\alpha^{(\text{pre})} = \text{interp} \left(\mathbf{x}, \text{alpha} \left(\text{softplus} \left(\mathbf{V}^{(\text{density})} \right) \right) \right),$$
$$\alpha^{(\text{in})} = \text{alpha} \left(\text{interp} \left(\mathbf{x}, \text{softplus} \left(\mathbf{V}^{(\text{density})} \right) \right) \right),$$
$$\alpha^{(\text{post})} = \text{alpha} \left(\text{softplus} \left(\text{interp} \left(\mathbf{x}, \mathbf{V}^{(\text{density})} \right) \right) \right).$$

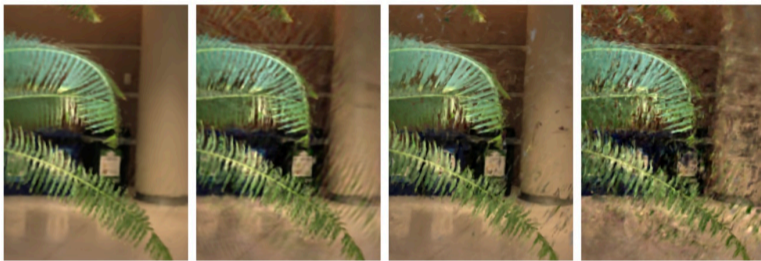


Plenoxel



- ▶ **Sparse data structure** using progressive training
- ▶ Highest resolution $512 \times 512 \times 512$ for NeRF-Synthetic objects

DVGO & Plenoxel



Full

No SH TV

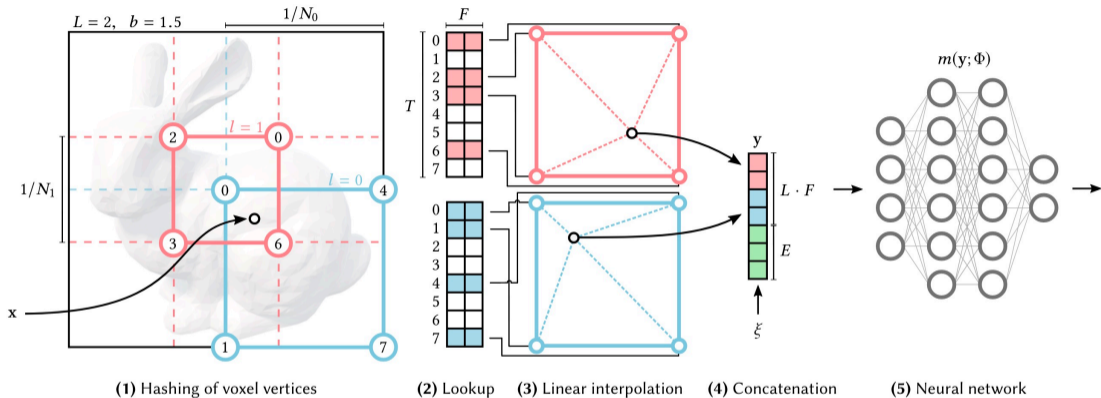
No σ TV

No TV

Regularization: both optimize density voxel grid directly, additional "tricks" required

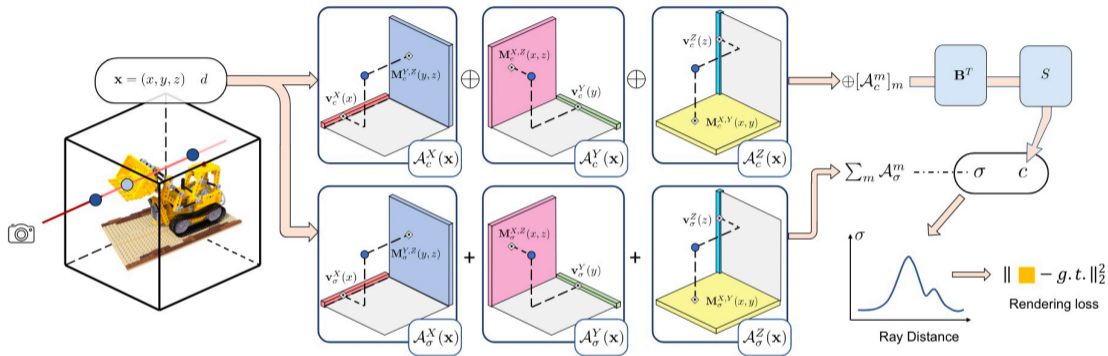
- ▶ DVGO: low density initialization, lower learning rate to voxels visible to fewer views
- ▶ Plenoxel: total variation loss, (sparsity loss for forward-facing scenes)

Instant NGP



- ▶ Multi-resolution feature grid saved in hash table + small MLP
- ▶ Random **Parameter sharing** by not resolving hash confliction

TensorRF



- ▶ Feature grid represented by factorized tensor + small MLP
- ▶ **Parameter sharing** along a certain axis / plane

Discussions

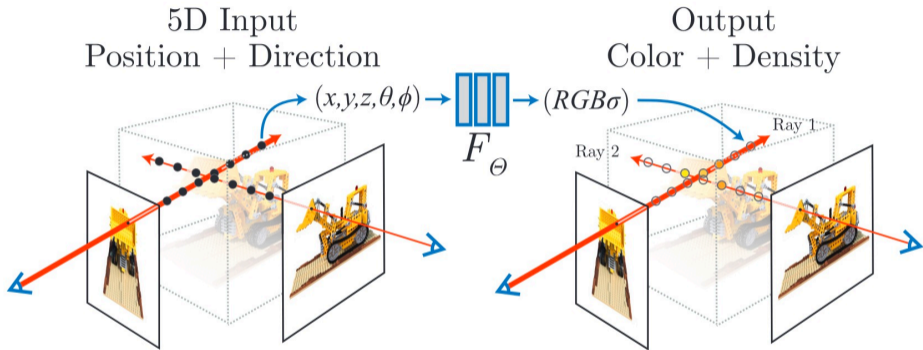
Method	Device	Training Time	Memory (MB)	PSNR	Initial LR
NeRF	NVIDIA V100	1-2 days	6MB	31.01	5e-4
DVGO	NVIDIA RTX2080 Ti	15 mins	612.1	31.95	0.1
Plenoxel	NVIDIA Titan RTX	11 mins	778.1	31.71	30
Instant NGP	NVIDIA RTX 3090	5 mins	50.4	33.18	0.01
Instant NGP's NeRF	NVIDIA RTX 3090	5 mins	–	30.06	–
TensoRF	NVIDIA V100	17 mins	71.8	33.14	0.02

- ▶ Parameter sharing methods are more **memory efficient**
- ▶ Larger **learning rates** of local tensors lead to faster convergence
- ▶ Good **Implementation** is another key to fast training, e.g., Instant NGP's NeRF
- ▶ Empty space skipping is also commonly used

Outline

- ▶ Real-time rendering of neural radiance fields
- ▶ Fast training of neural scene representation
- ▶ 3D-aware generative model

Why 3D-Aware Generative Models?



- ▶ NeRF optimizes the MLP for a **single scene**
- ▶ Not easy to create non-existent scenes

Why 3D-Aware Generative Models?



- ▶ Existing 2D GANs are able to generate high fidelity, novel contents
- ▶ However, there is no **3D controllability**, e.g., control over camera poses

3D-Aware GANs

Naïve Solution: Replace the 2D generator as a conditional radiance field

Challenge: Image-level supervision, **no per-pixel loss**; computationally expensive

Idea: Reducing FLOPs of one forward pass of the generator $H \times W \times K \times L_r$

▶ **Reducing $H \times W$**

- ▶ Patch-based discriminator
- ▶ 2D-CNN based upsampling network

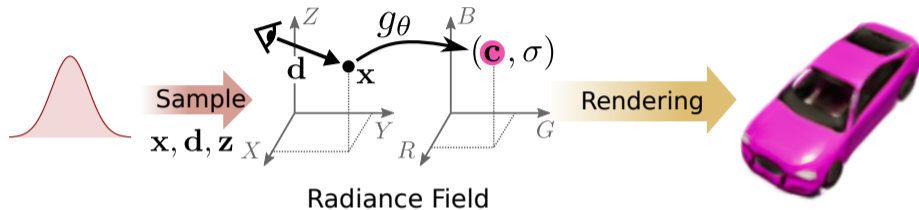
▶ **Reducing K**

- ▶ Sampling on a few isosurfaces
- ▶ Empty space skipping

▶ **Reducing L_r**

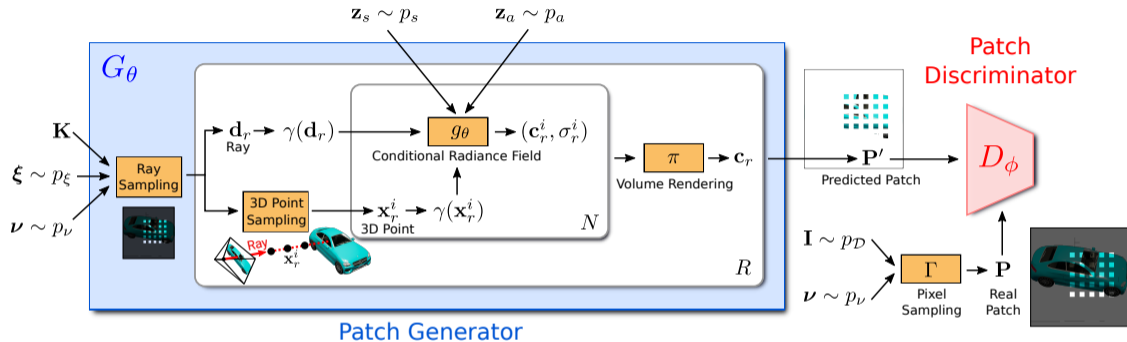
- ▶ Replace (a part of) large MLP with 2D/3D CNN

GRAF: Generative Radiance Fields



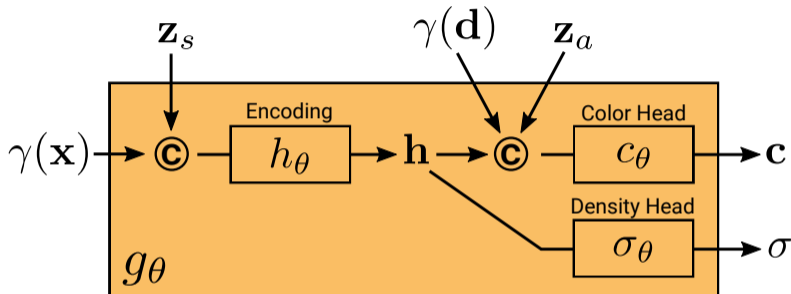
- ▶ **Generative model** for radiance fields
- ▶ Trained from **unstructured** and **unposed** 2D image collections

Reducing $H \times W$: Patch-based discriminator



- ▶ **Conditional** neural radiance fields supervised by **adversarial loss**
- ▶ Sample **camera poses** and **image patches** of size 32×32 pixels

GRAF: Generator



- ▶ \mathbf{z}_s for shape, \mathbf{z}_a for appearance
- ▶ Automatically disentangled \mathbf{z}_s and \mathbf{z}_a

RGB



Depth



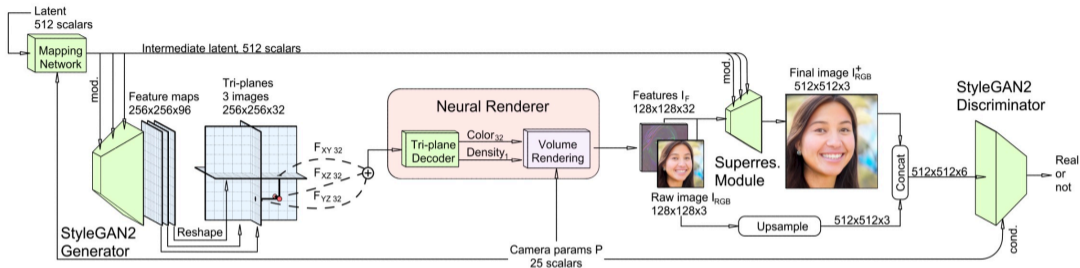
256x256

Reducing $H \times W$: StyleNeRF



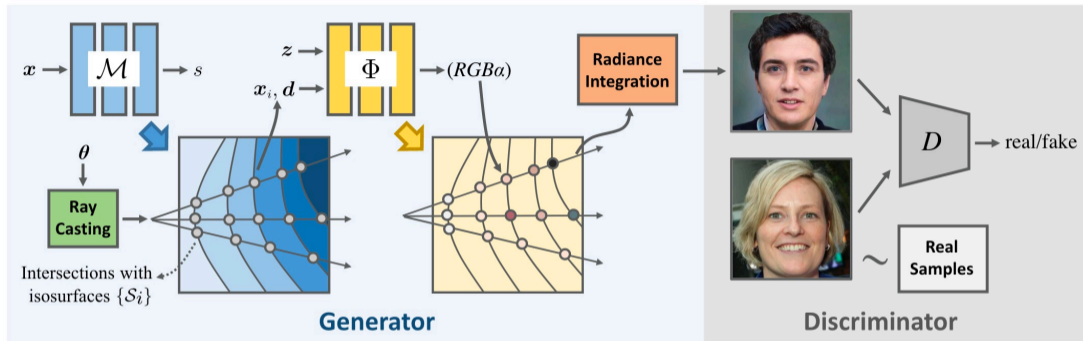
- ▶ Low-res. volume rendering + **2D upsampling neural renderer**
- ▶ High fidelity, harder to preserve multi-view consistency

Reducing $H \times W$ and L_r : EG3D



- ▶ L_r : Replace per-point large MLP with tri-plane 2D generator + small MLP
- ▶ $H \times W$: Uses 2D upsampling neural renderer

Reducing K : GRAM



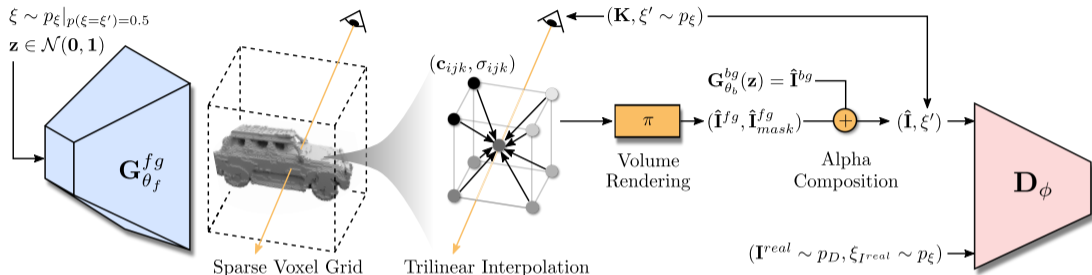
- No 2D CNN, evaluate on 24 or 48 isosurfaces

Reducing K and L_r : VoxGRAF



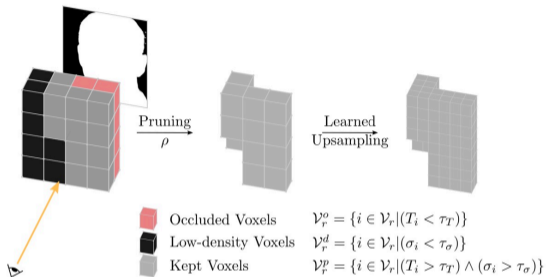
- ▶ Inspired by DVGO and Plenoxel, represent scene as 3D **sparse** voxel grids
- ▶ Fast rendering during inference, one forward pass to generate the voxel grids

VoxGRAF

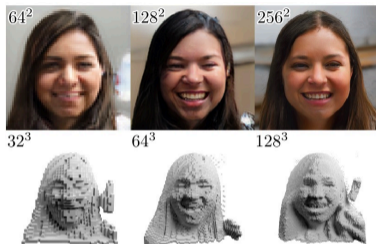


- ▶ Learn sharp surface via regularization
- ▶ Disentangle foreground and background

VoxGRAF



(a) Pruning



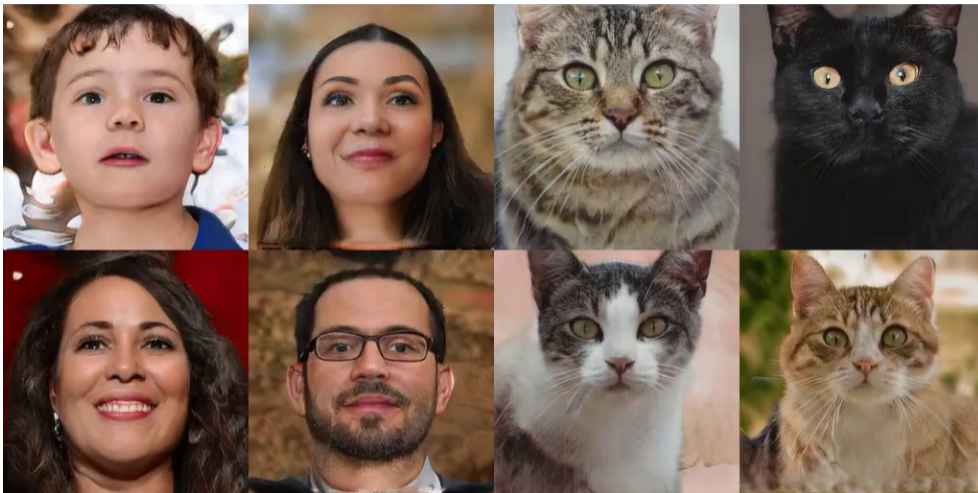
(b) Progressive Growing

- ▶ Density based pruning
- ▶ Progressive Growing for the resolution of voxel grids and 2D image

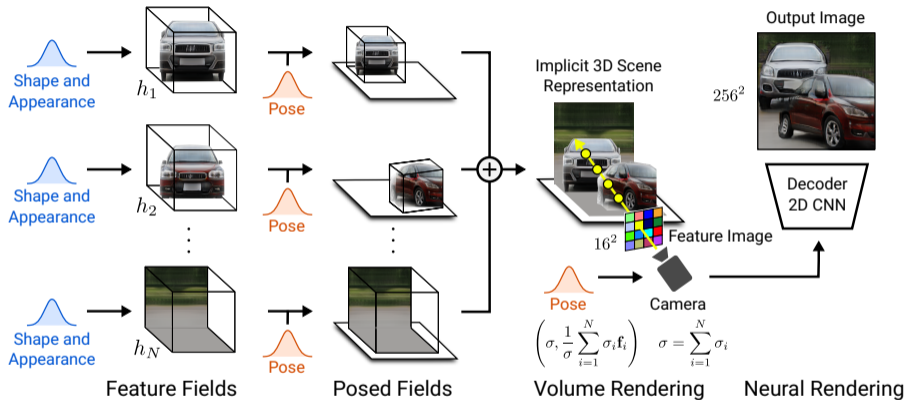
VoxGRAF



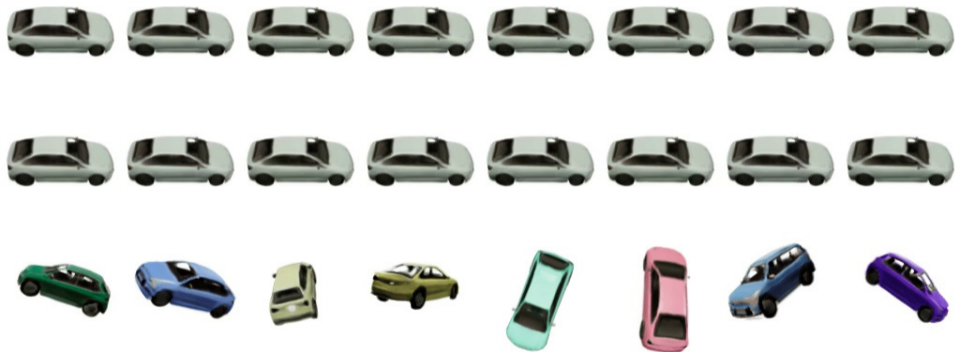
VoxGRAF



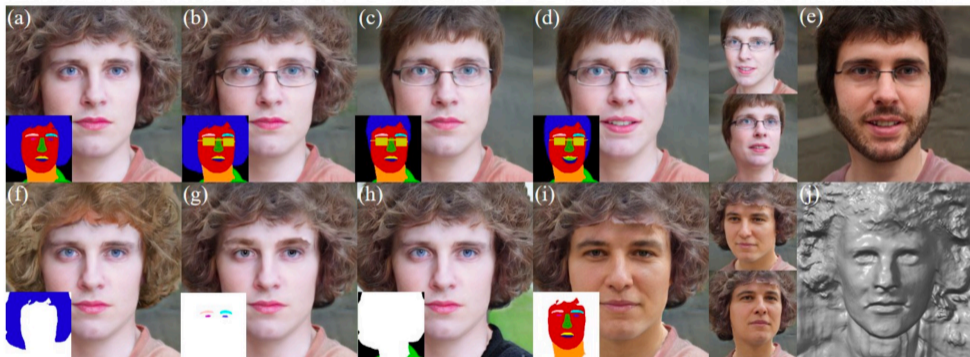
What's more? GIRAFFE: Compositional Generation



What's more? Category-level 6DoF Pose Estimation



What's more? Semantic Editing



Conclusions

- ▶ 3D-aware generative methods share similar ideas to fast inference and training
- ▶ 2D neural renderer leads to better FID
- ▶ Many exciting applications

	FFHQ $R_I = 256^2$	AFHQ $R_I = 256^2$	Carla $R_I = 128^2$
GIRAFFE	31.5	16.1	–
StyleNeRF	8.0	–	–
EG3D	4.8	3.9	–
GRAF	71	121	41
π -GAN	85	47	29.2
GOF	69.2	54.1	29.3
GRAM	17.9	18.5	26.3
VoxGRAF	14.4	9.6	11.3

Thank you!

<https://yiyiliao.github.io/>

https://zju3dv.github.io/inr_tutorial/