

CS479: Machine Learning for 3D Data

Neural Rendering

LECTURE 8
MINHYUK SUNG

Spring 2025
KAIST

Course Road Map

Representations

- Point clouds
- Implicit representation
- **Multi-view images to 3D**
- Hybrid representations
- Meshes
- ~~CAD~~
- Representation Conversion

Applications

- 3D perception (Encoding)
- **Reconstruction (Decoding)**
- Manipulation
- ~~Generation~~

Novel View Synthesis



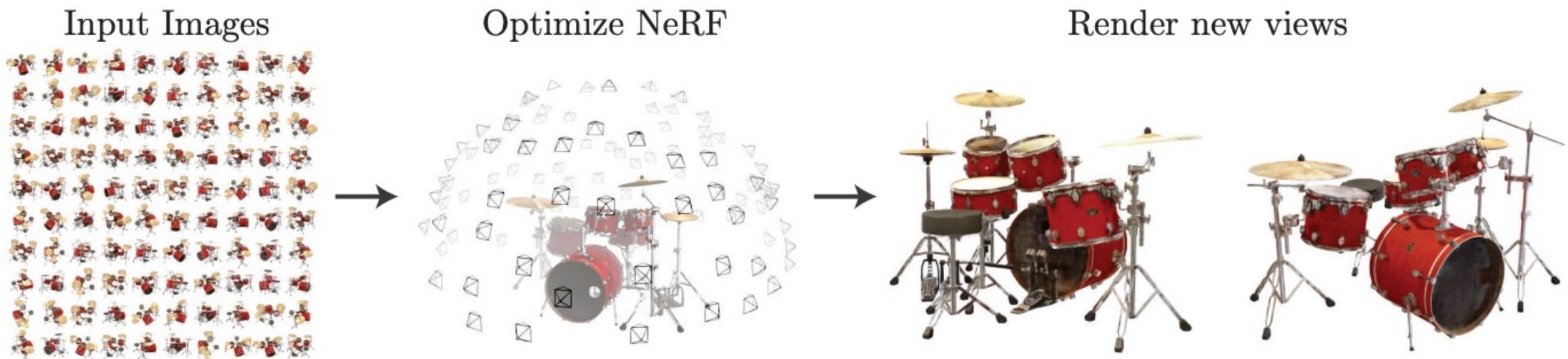
Input: Sparsely sampled images of scene.

Output: New views of the same scene.

Mildenhall et al., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV 2020.

Novel View Synthesis

- 3D reconstruction + Novel view image synthesis.
- The focus is on synthesizing novel view images.
- Only images are given. **No supervision for the 3D shape.**



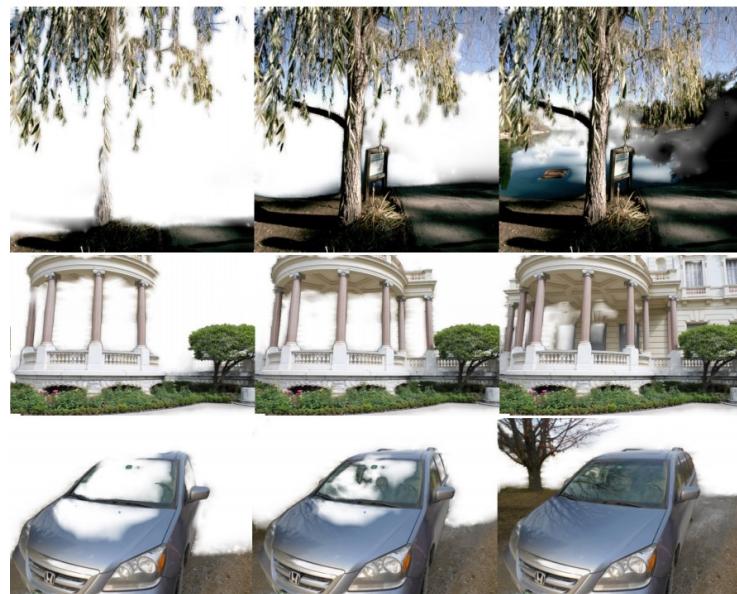
Mildenhall et al., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV 2020.



Neural Novel View Synthesis – Literature

Soft 3D

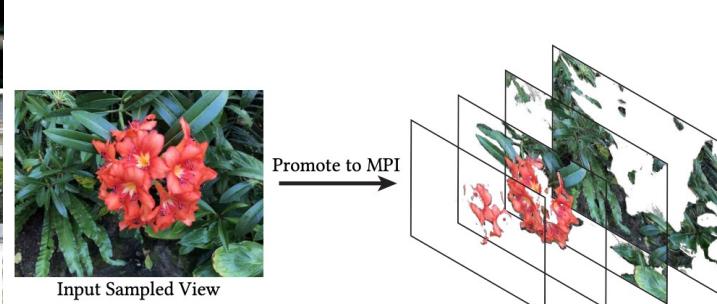
Culmination of non-deep stereo matching techniques.
(Penner & Zhang 2017)



Multi-plane Image Methods

Typical deep learning pipelines - images go into a 3D CNN, big RGBA 3D volume comes out.

(Zhou et al. 2018)
(Srinivasan et al. 2019)
(Mildenhall et al. 2019)
(Flynn et al. 2019)
(Tucker & Snavely 2020)



Neural Volume

Direct gradient descent to optimize an RGBA volume, regularized by a 3D CNN
(Lombardi et al. 2019)



Neural Novel View Synthesis

Training

- **Input:** A set of **multi-view images** of an object or scene, along with their corresponding **camera matrices**.
- **Output:** An implicit 3D representation of the object or scene
- **Loss Function:** L2 loss between the **rendered** image and the corresponding ground truth image for each view

Neural Novel View Synthesis

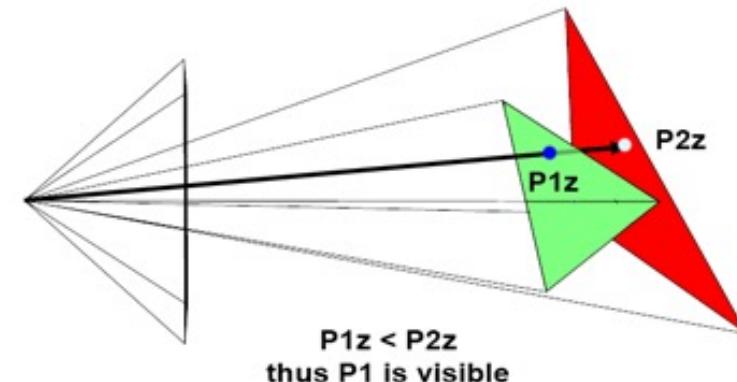
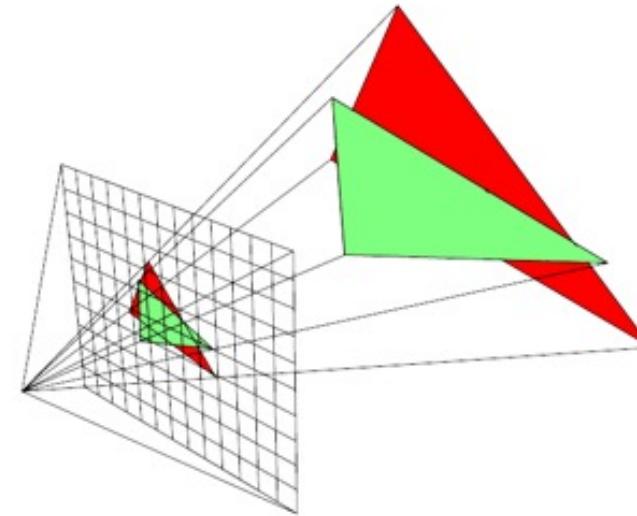
Inference

- **Input:** A camera matrix (not necessarily seen during training).
- **Output:** A 2D image **rendered** from the learned implicit 3D representation corresponding to the given viewpoint.

What is the rendering?

Rasterization – Object-Centric

1. Project each mesh face onto the 2D screen.
2. Interpolate vertex data across pixels within the projected face.
3. Skip pixels already processed with a closer depth; otherwise, write them.

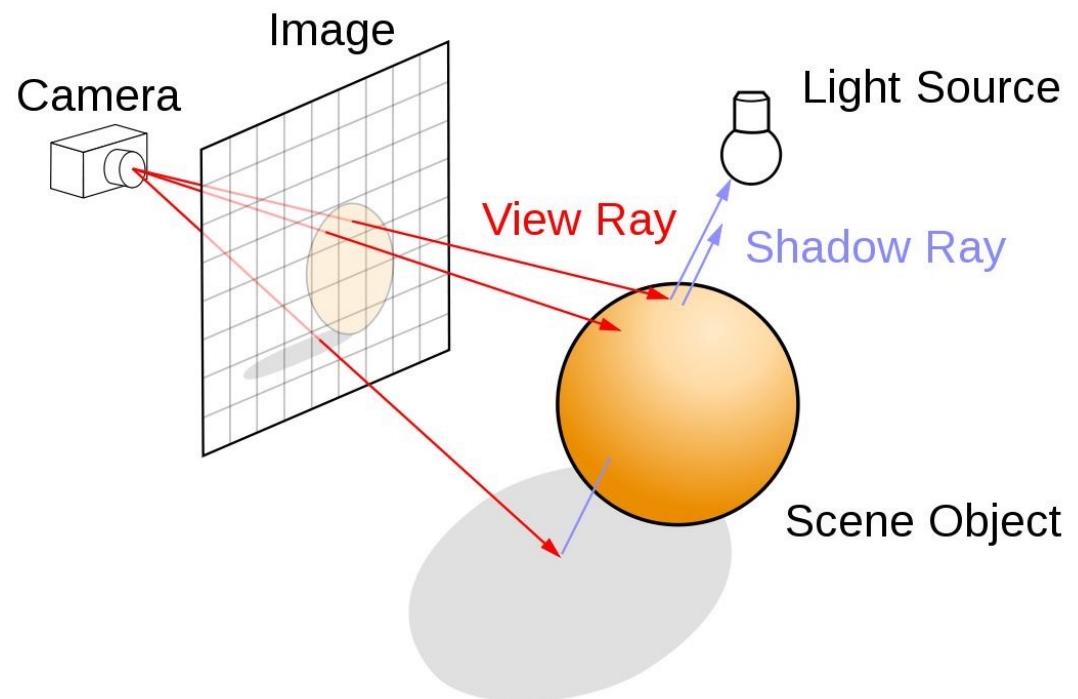


© www.scratchapixel.com

<https://www.scratchapixel.com/>

Ray Casting – Image-Centric

1. Shoot a **ray** from each pixel.
2. Find the first **intersection** point with the 3D object.
3. Retrieve the **color** information.

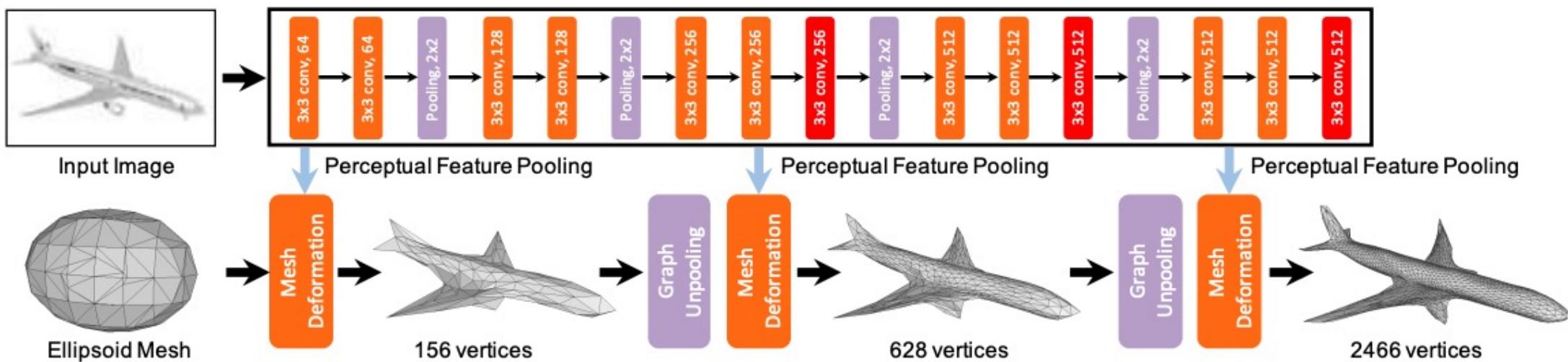


<https://developer.nvidia.com/discover/ray-tracing>

Multi-View to 3D Mesh

Can we reconstruct a 3D object/scene into a **mesh**?

Q. What if we deform a template mesh to align it with images from each viewpoint via rendering?



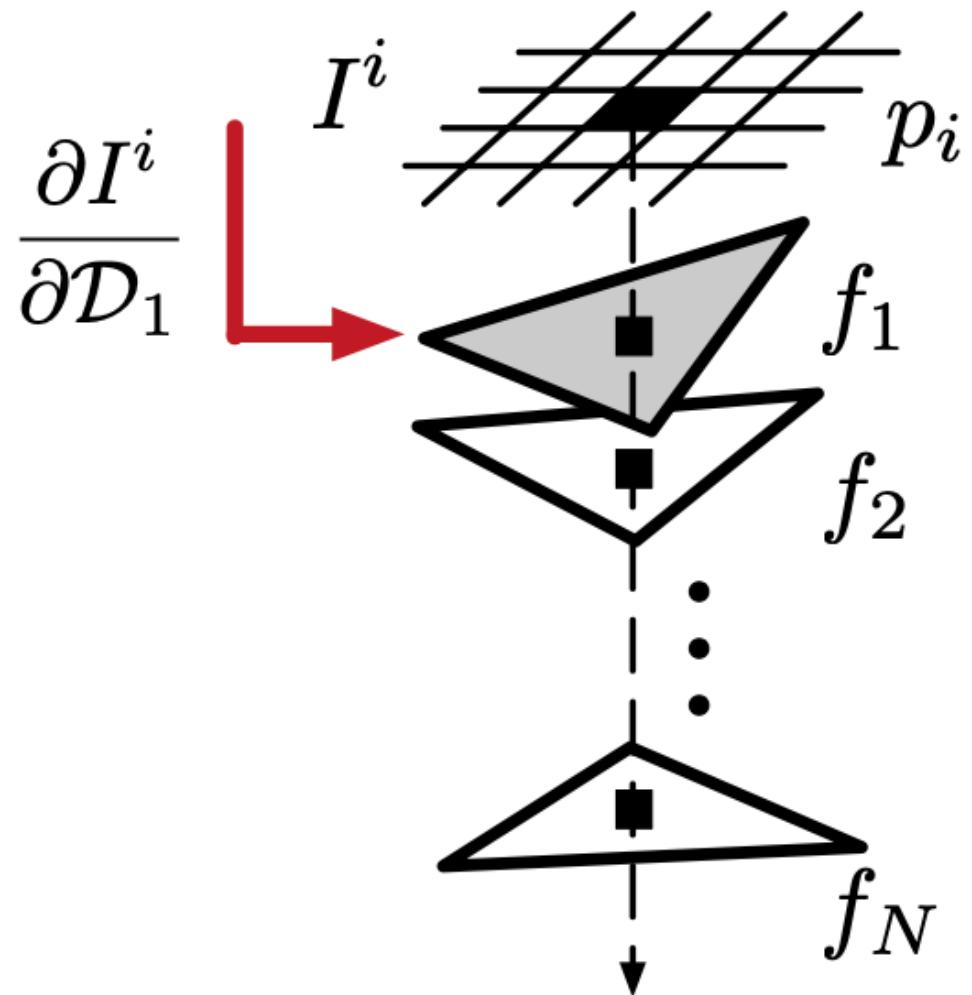
Multi-View to 3D Mesh: Limitations

- Limited to generating shapes that have the **same topology** as the template shape.
- It is very **challenging** to generate a mesh with **arbitrary topology** directly.



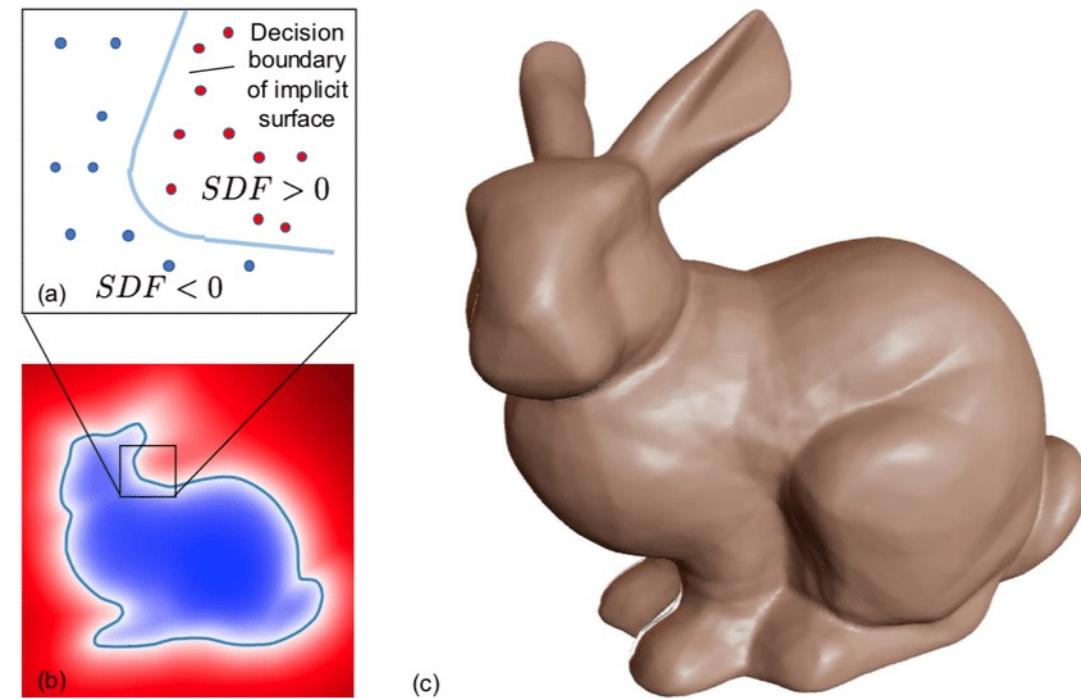
Multi-View to 3D Mesh: Limitations

- It is also **very slow** to update the mesh since only the **first intersection point** is affected in the gradient descent.



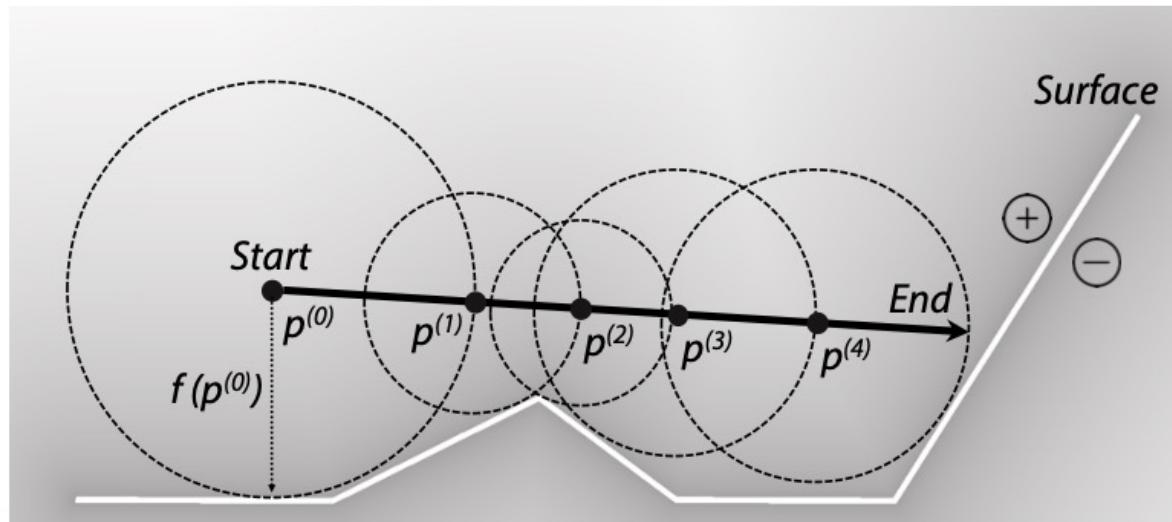
Rendering with Implicit Functions

- **Implicit representations** (such as signed distance functions) are known to be the best for 3D reconstruction.
- But how can we **render** a signed distance function?



Sphere Tracing

- A technique to render an implicit surface with **ray casting**.
- Iterate advancing the ray by the **current signed distance** until the distance approaches (almost) zero.

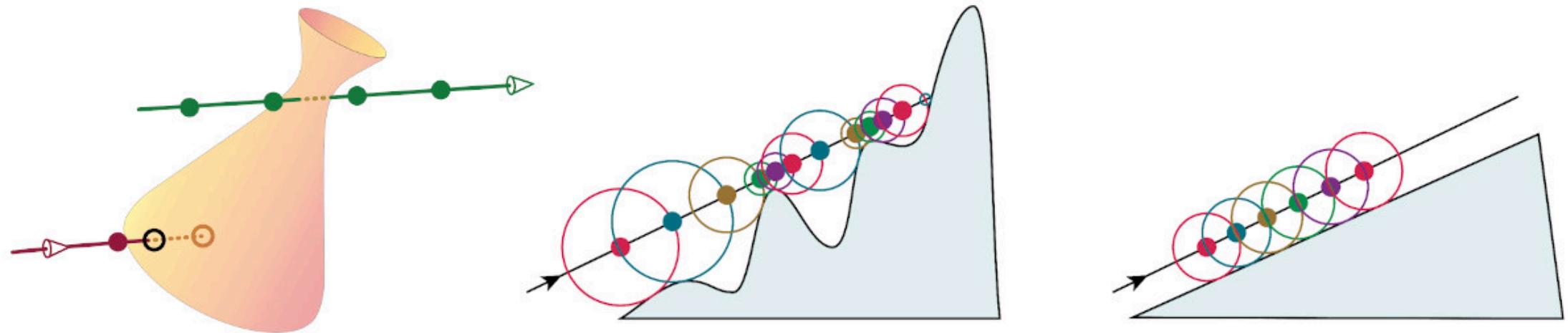


Algorithm 1 Naive sphere tracing algorithm for a camera ray $L : \mathbf{c} + d\tilde{\mathbf{v}}$ over a signed distance fields $f : \mathbb{N}^3 \rightarrow \mathbb{R}$.

- 1: Initialize $n = 0$, $d^{(0)} = 0$, $\mathbf{p}^{(0)} = \mathbf{c}$.
 - 2: **while** not converged **do**:
 - 3: Take the corresponding SDF value $b^{(n)} = f(\mathbf{p}^{(n)})$ of the location $\mathbf{p}^{(n)}$ and make update: $d^{(n+1)} = d^{(n)} + b^{(n)}$.
 - 4: $\mathbf{p}^{(n+1)} = \mathbf{c} + d^{(n+1)}\tilde{\mathbf{v}}$, $n = n + 1$.
 - 5: Check convergence.
 - 6: **end while**
-

Limitations of Sphere Tracing

- Not differentiable.
- Computationally expensive.
- Challenging to achieve convergence.

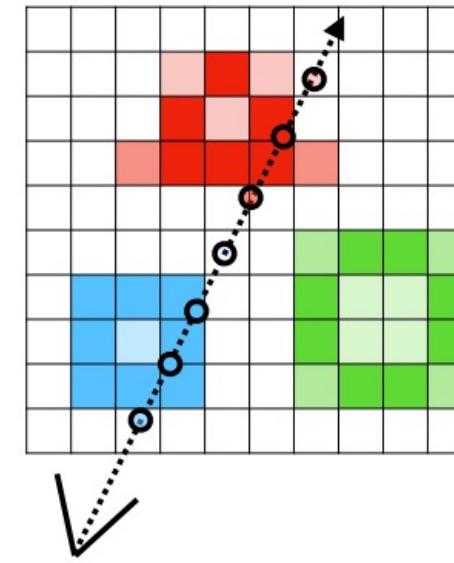
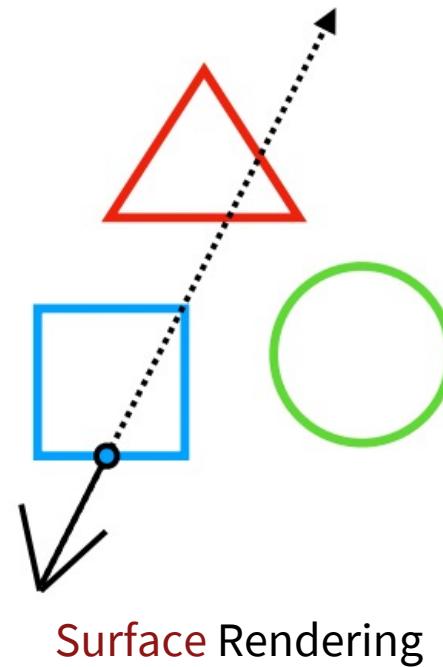


<https://tovacinni.github.io/sdf-explorer/>

Solution – Volume Rendering

Violate the physical constraint!

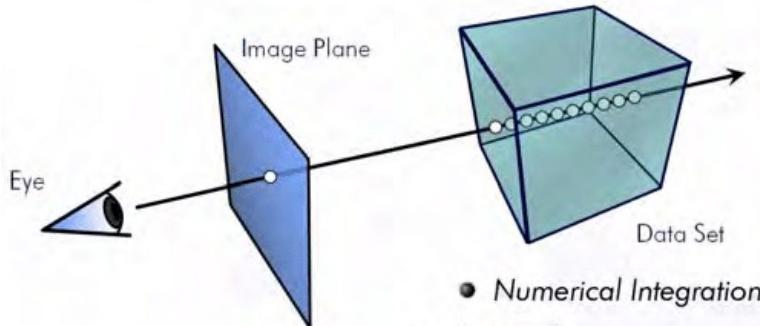
Relax the condition that the shape is a **solid** object, and consider the shape as a **volume**.



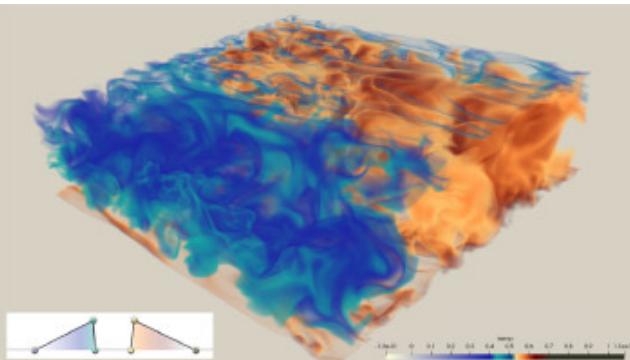
Benjamin Mildenhall, Neural Scene Representations for View Synthesis, PhD Dissertation.

Volume Rendering

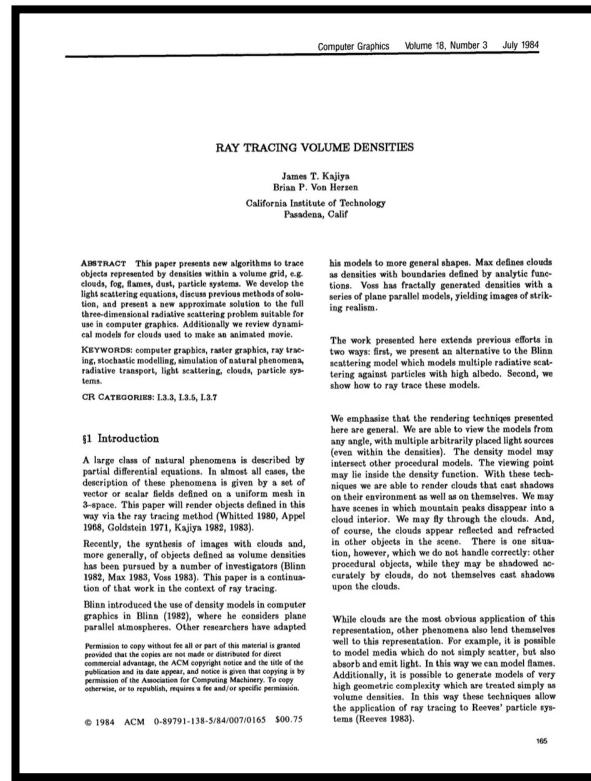
Has been widely studied and utilized for **fluid** and **smoke** simulation and rendering.



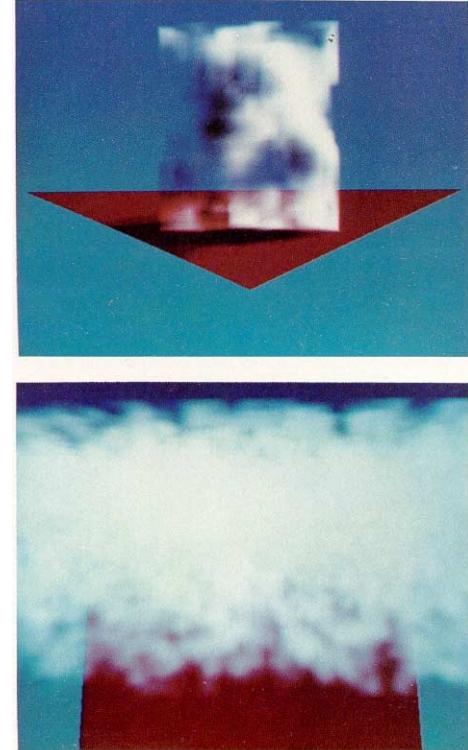
<https://nuguziii.github.io/cg/CG-001/>



<https://www.sciencedirect.com/science/article/pii/S0167819119301280>



James T. Kajiya, Ray Tracing Volume Densities, SIGGRAPH 1984.



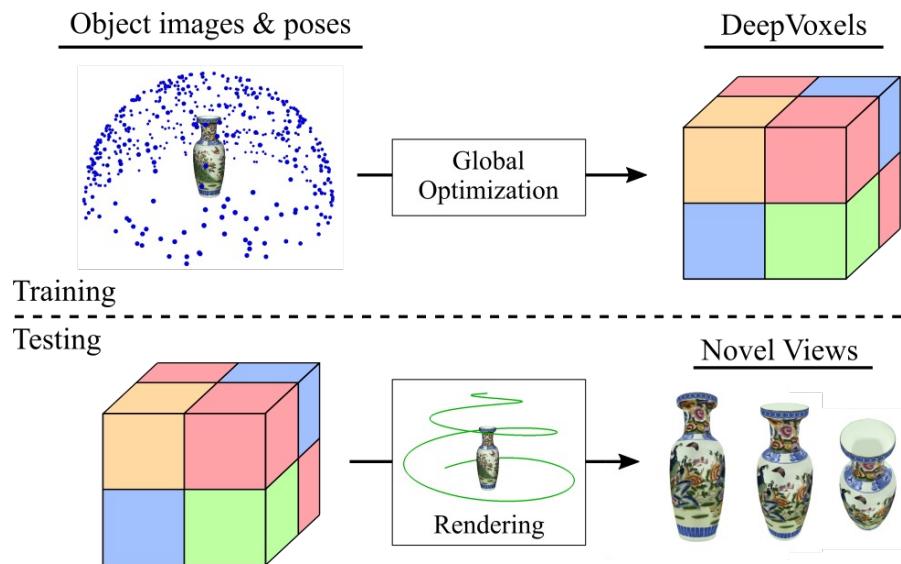
DeepVoxels

Sitzmann et al., DeepVoxels: Learning Persistent 3D Feature Embeddings, CVPR 2019.

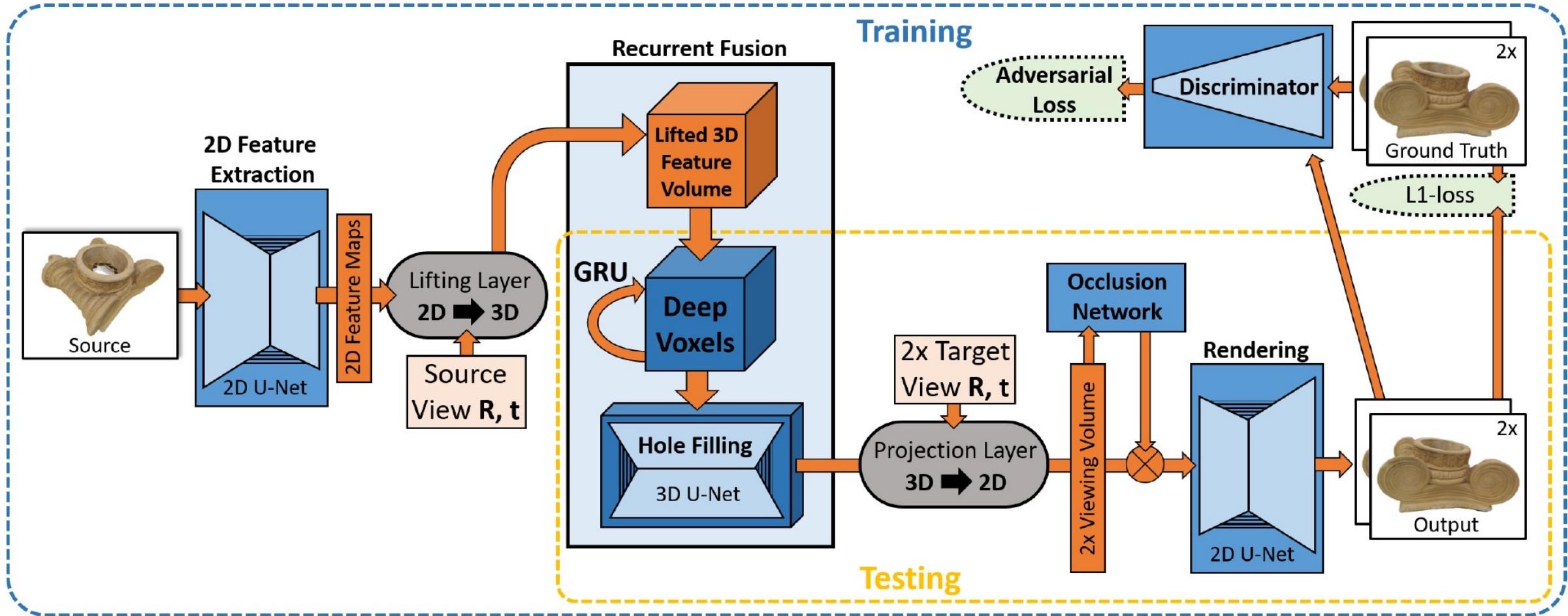
DeepVoxels

Encode the geometry information in **3D voxels**.

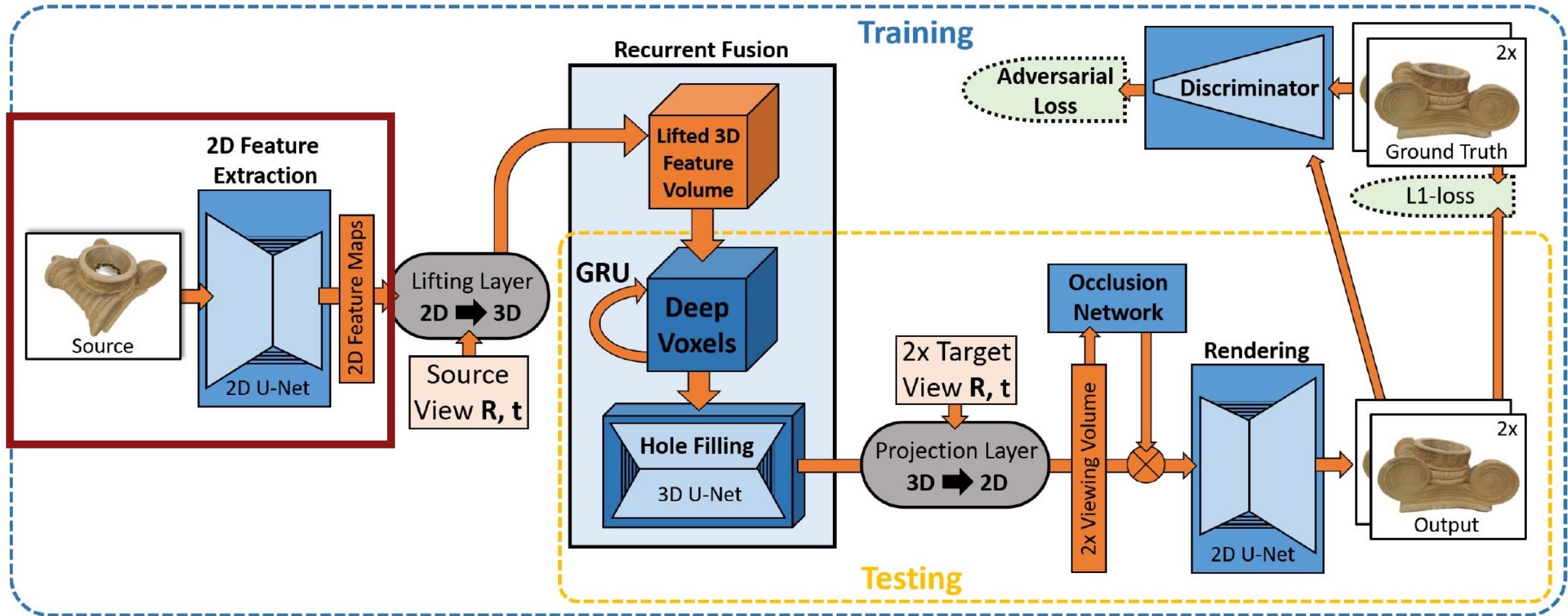
Input: N images with camera extrinsic/intrinsic parameters.
No supervision of the 3D shape.



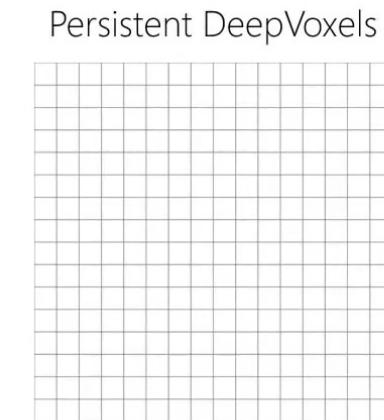
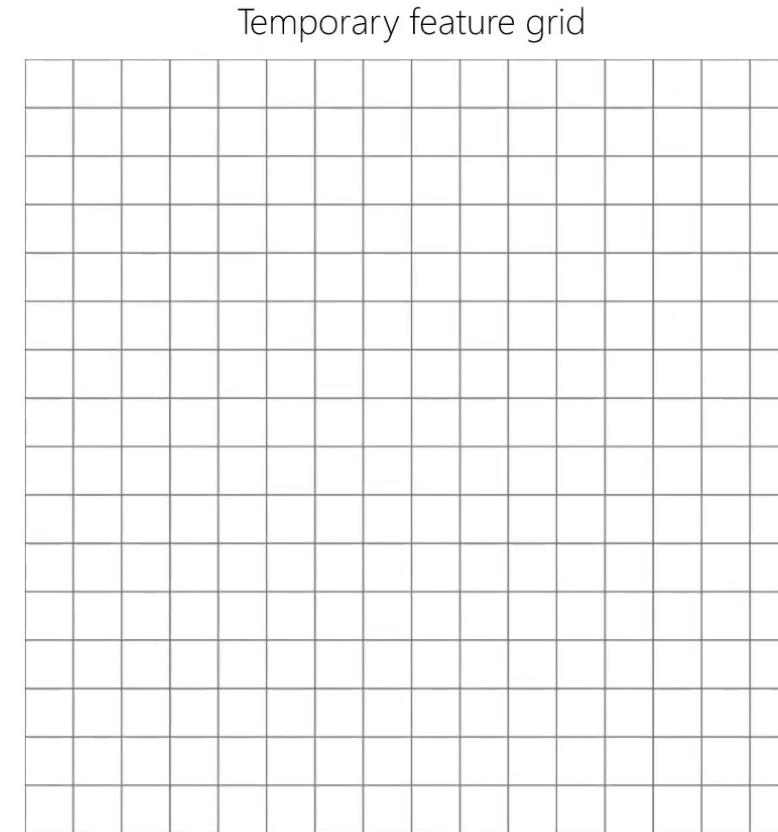
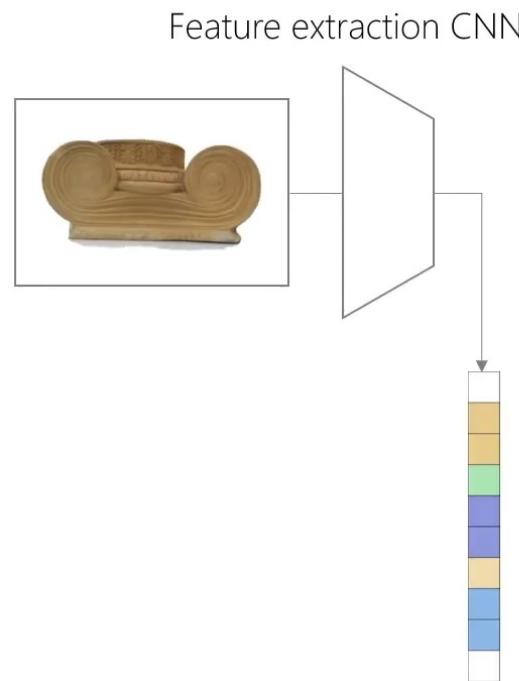
Framework



2D Feature Extraction

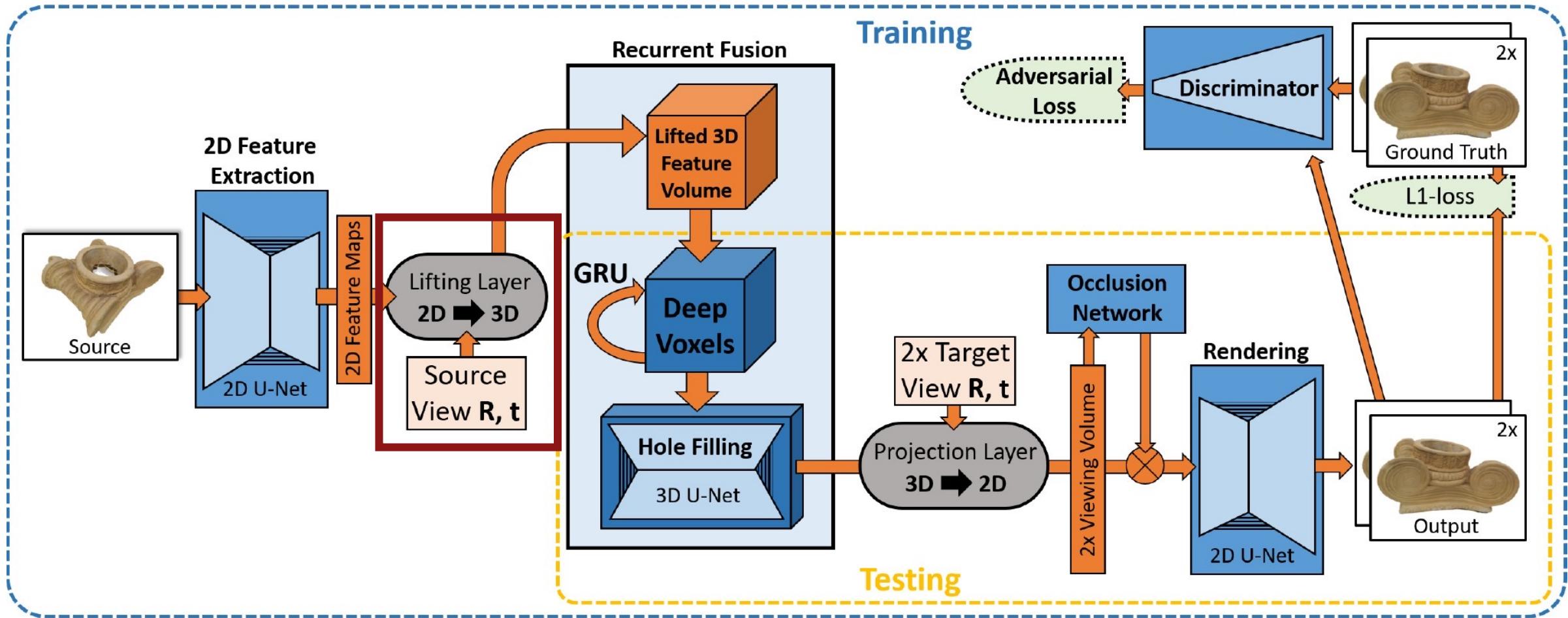


2D Feature Extraction

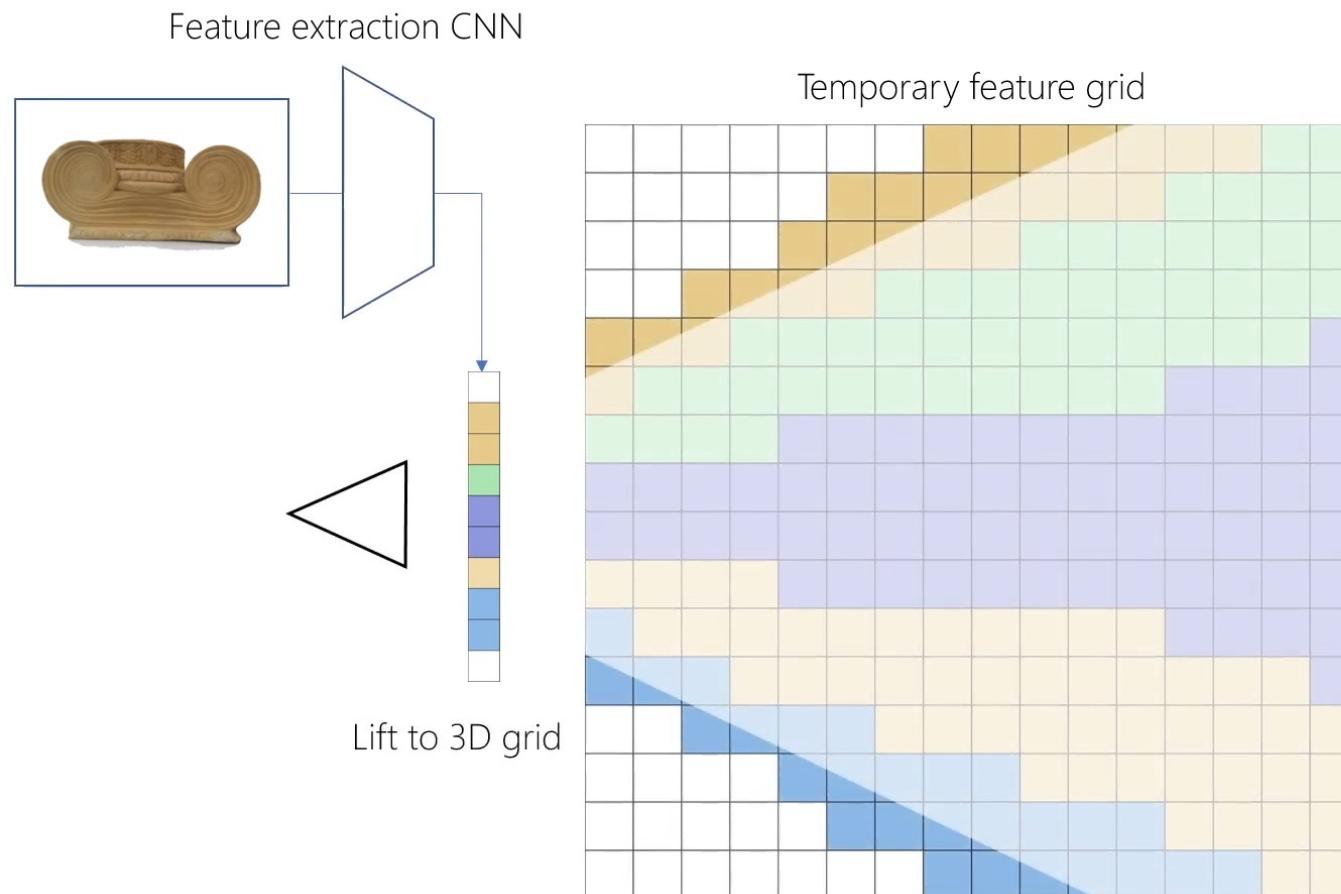


*1D-to-2D Diagram
Image: 1D
Reconstruction: 2D*

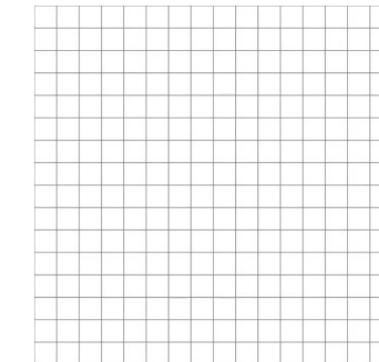
2D-to-3D Lifting



2D-to-3D Lifting

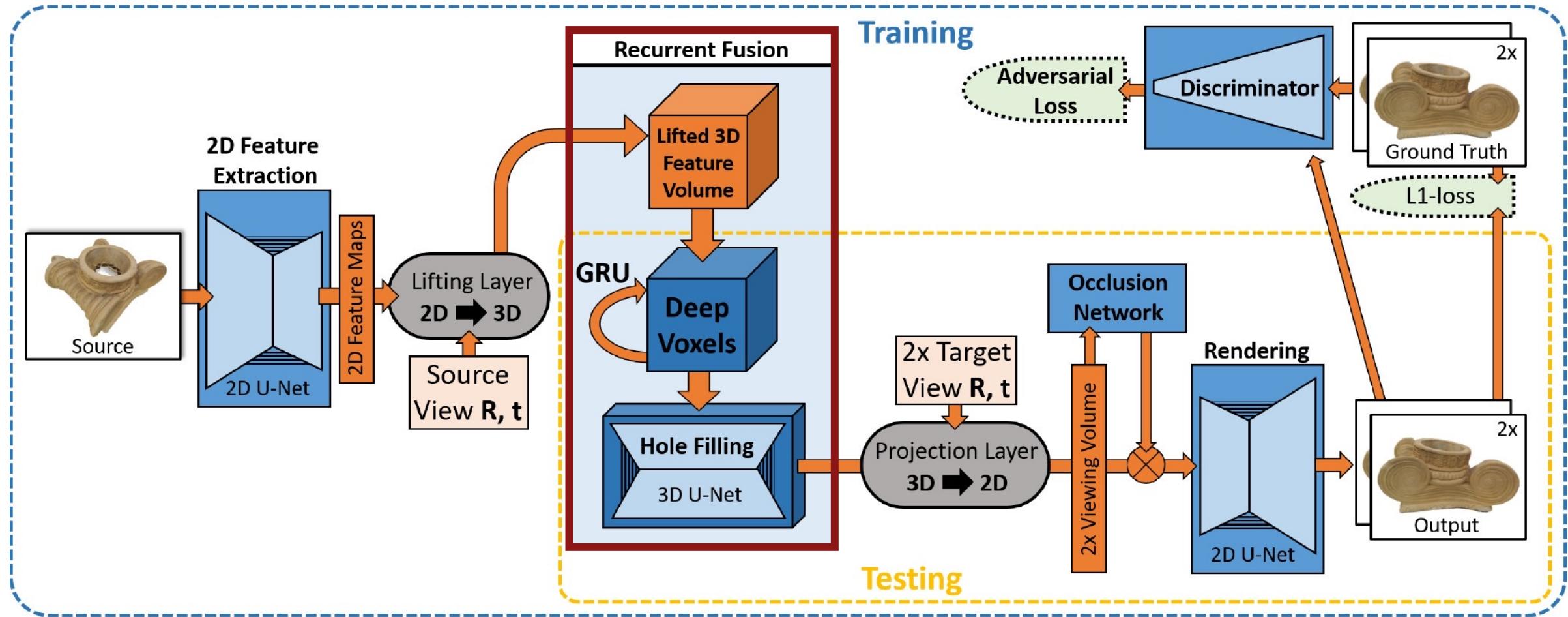


Persistent DeepVoxels

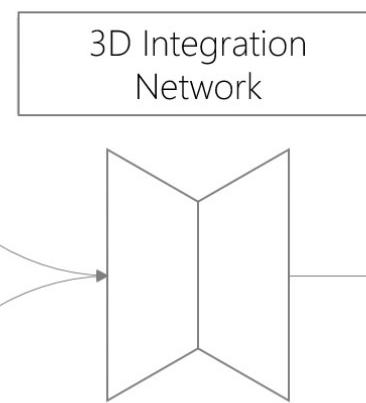
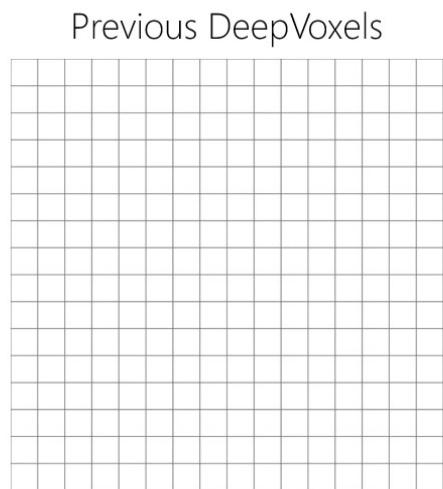


*1D-to-2D Diagram
Image: 1D
Reconstruction: 2D*

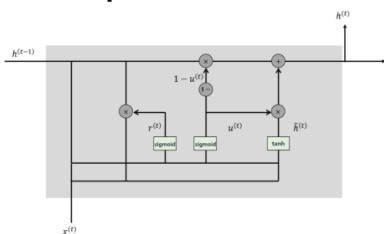
Deep Voxels Fusion



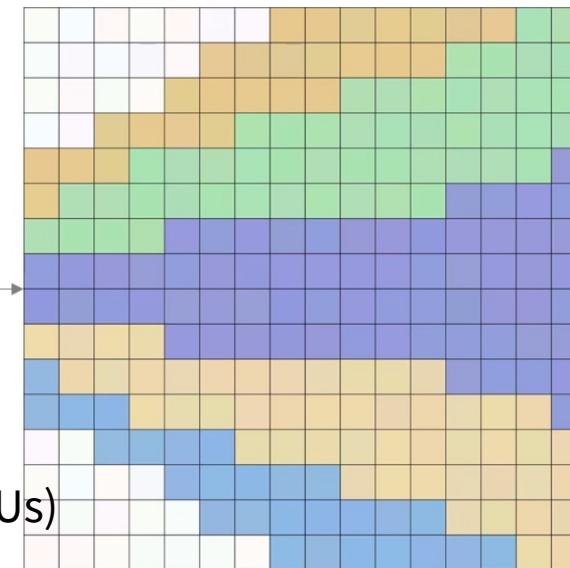
Deep Voxels Fusion



Use Gated Recurrent Units (GRUs)
per voxel

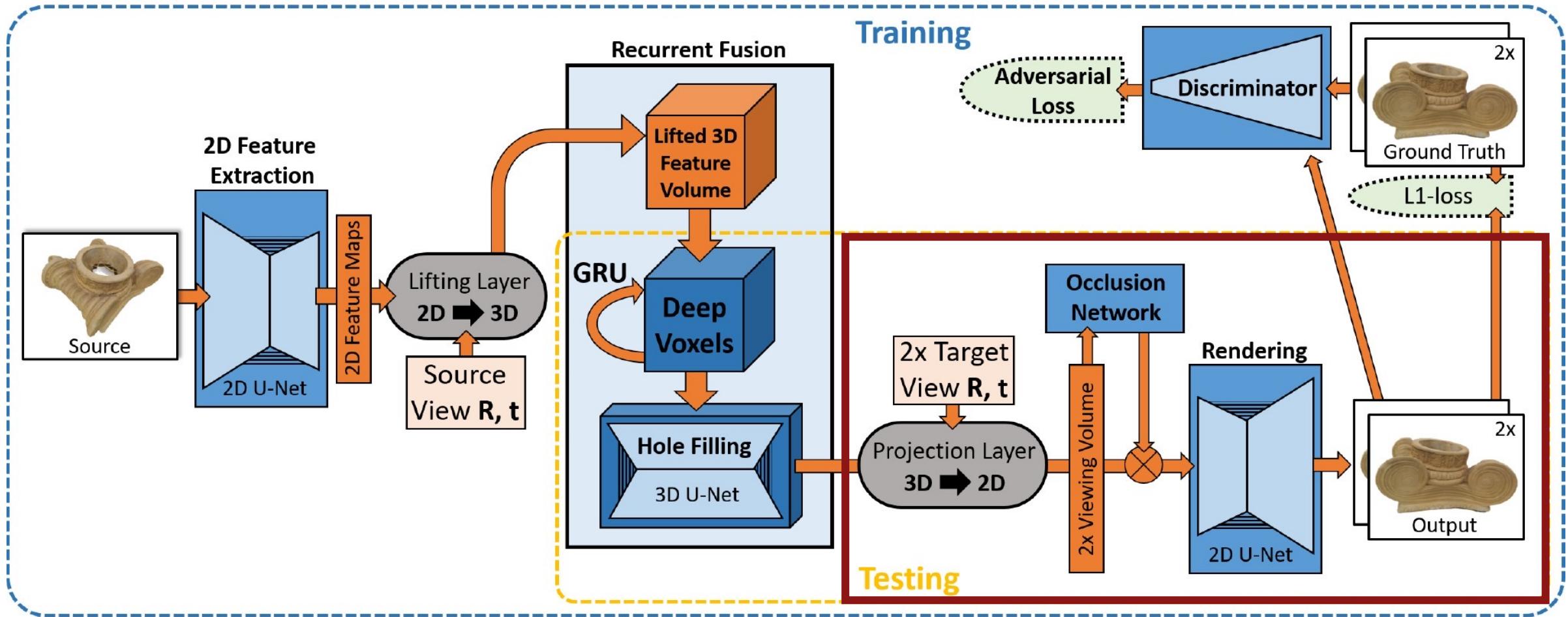


Persistent DeepVoxels
Updated DeepVoxels

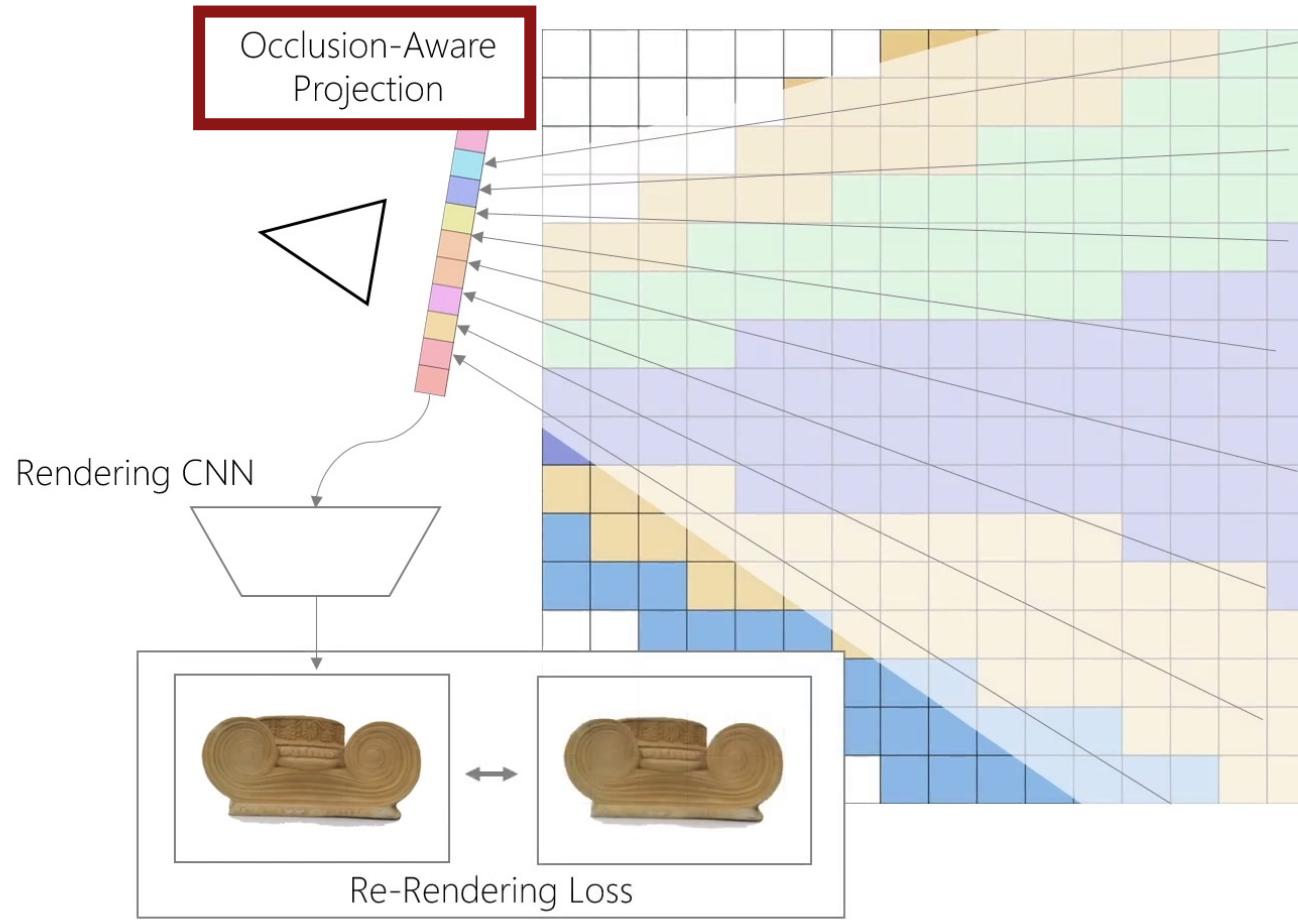


1D-to-2D Diagram
Image: 1D
Reconstruction: 2D

Rendering

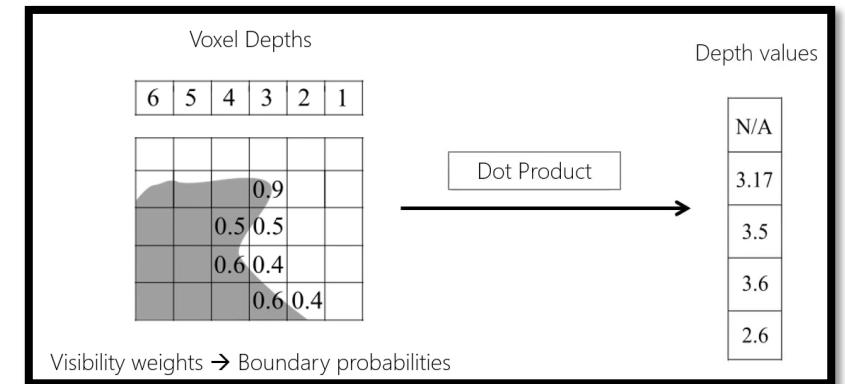


Rendering

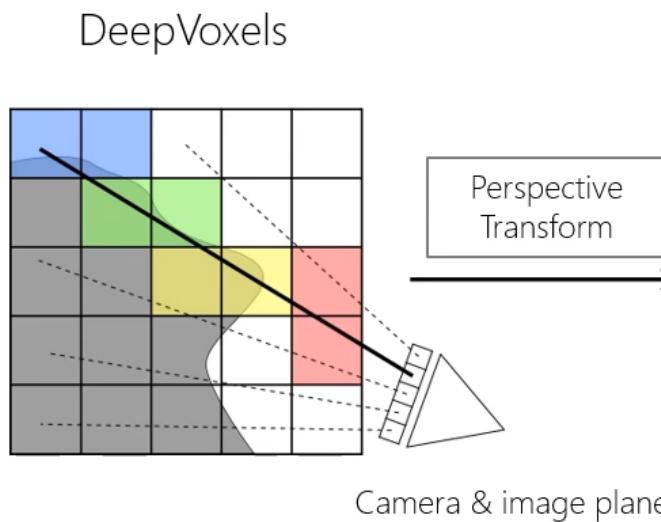


1D-to-2D Diagram
Image: 1D
Reconstruction: 2D

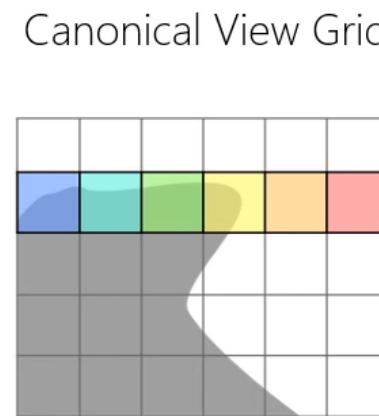
Occlusion-Aware Projection



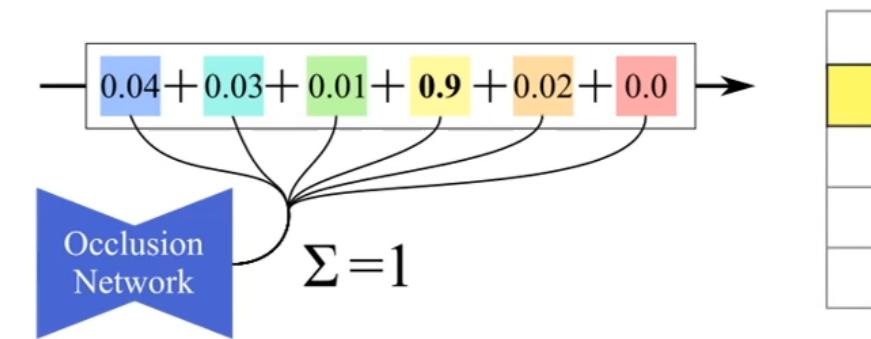
Depth prediction is possible!



Add camera-to-voxel distance to the feature of each voxel.



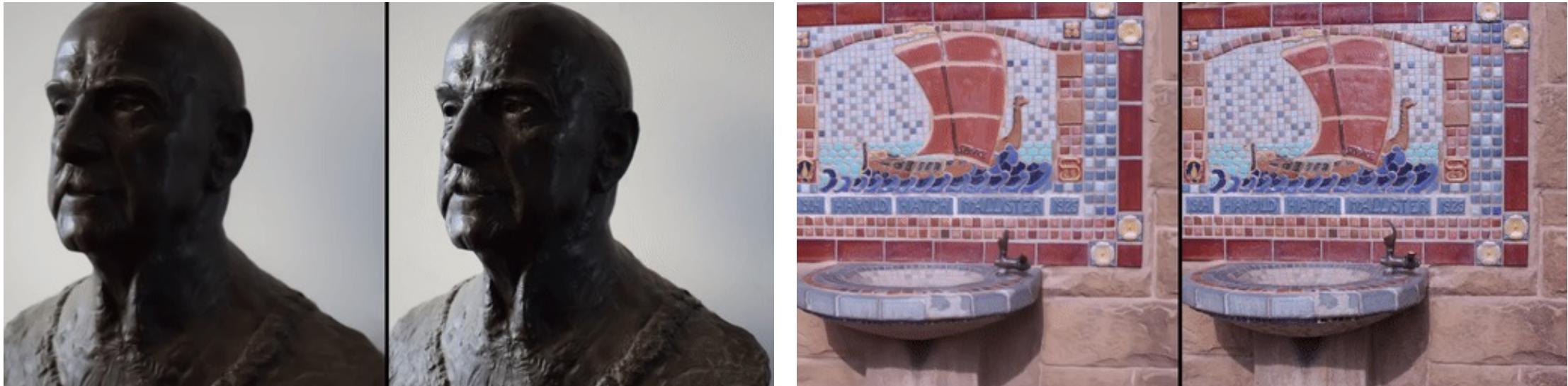
Visibility Reasoning



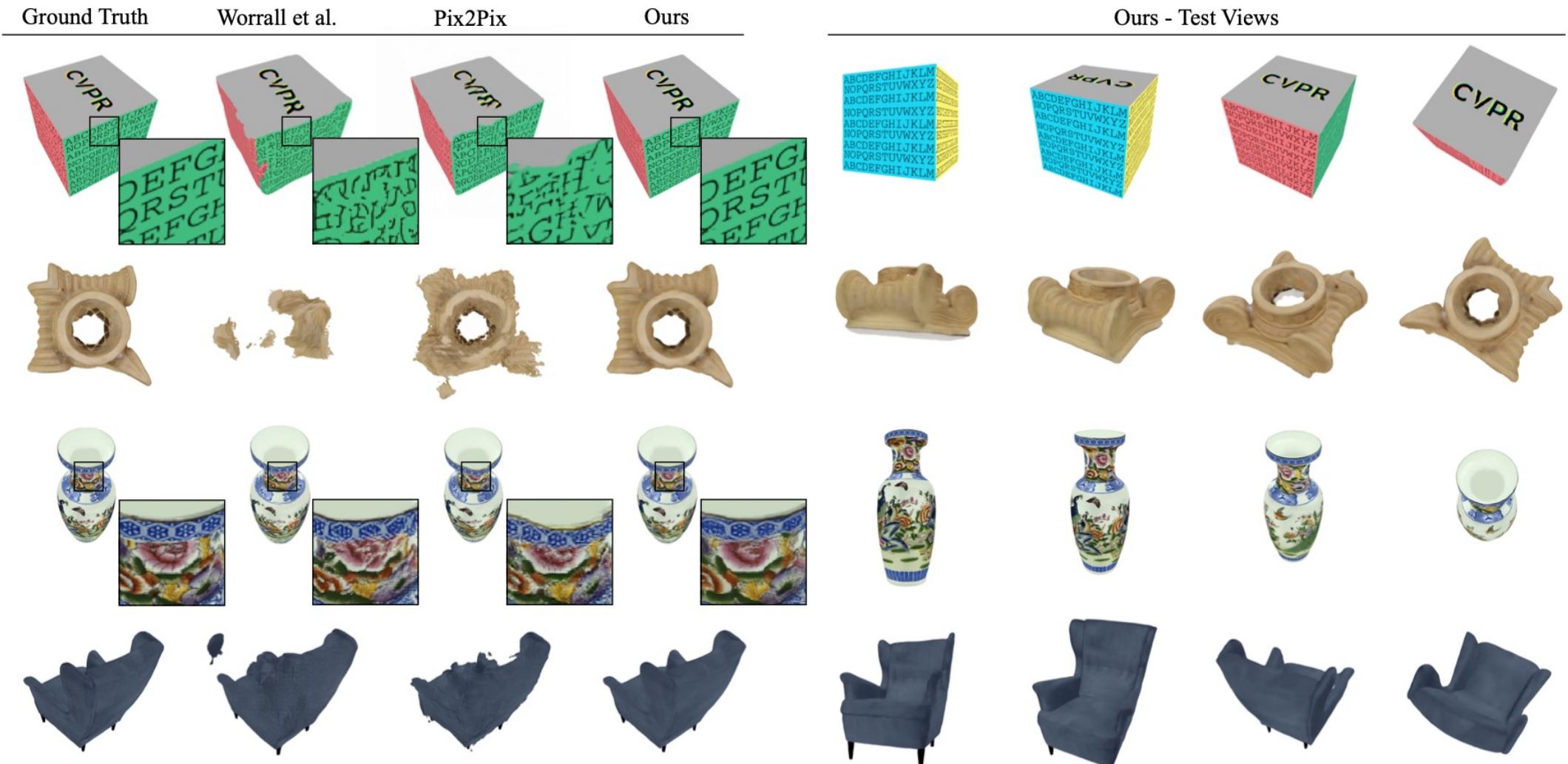
Run a single-layer 3D convolution network to predict visibility per voxel

1D-to-2D Diagram
Image: 1D
Reconstruction: 2D

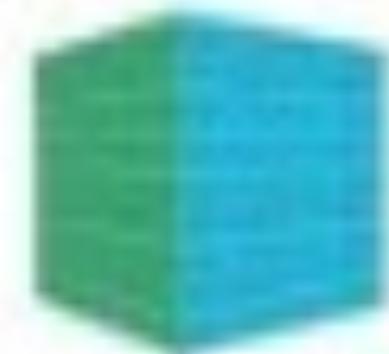
Novel View Synthesis Results



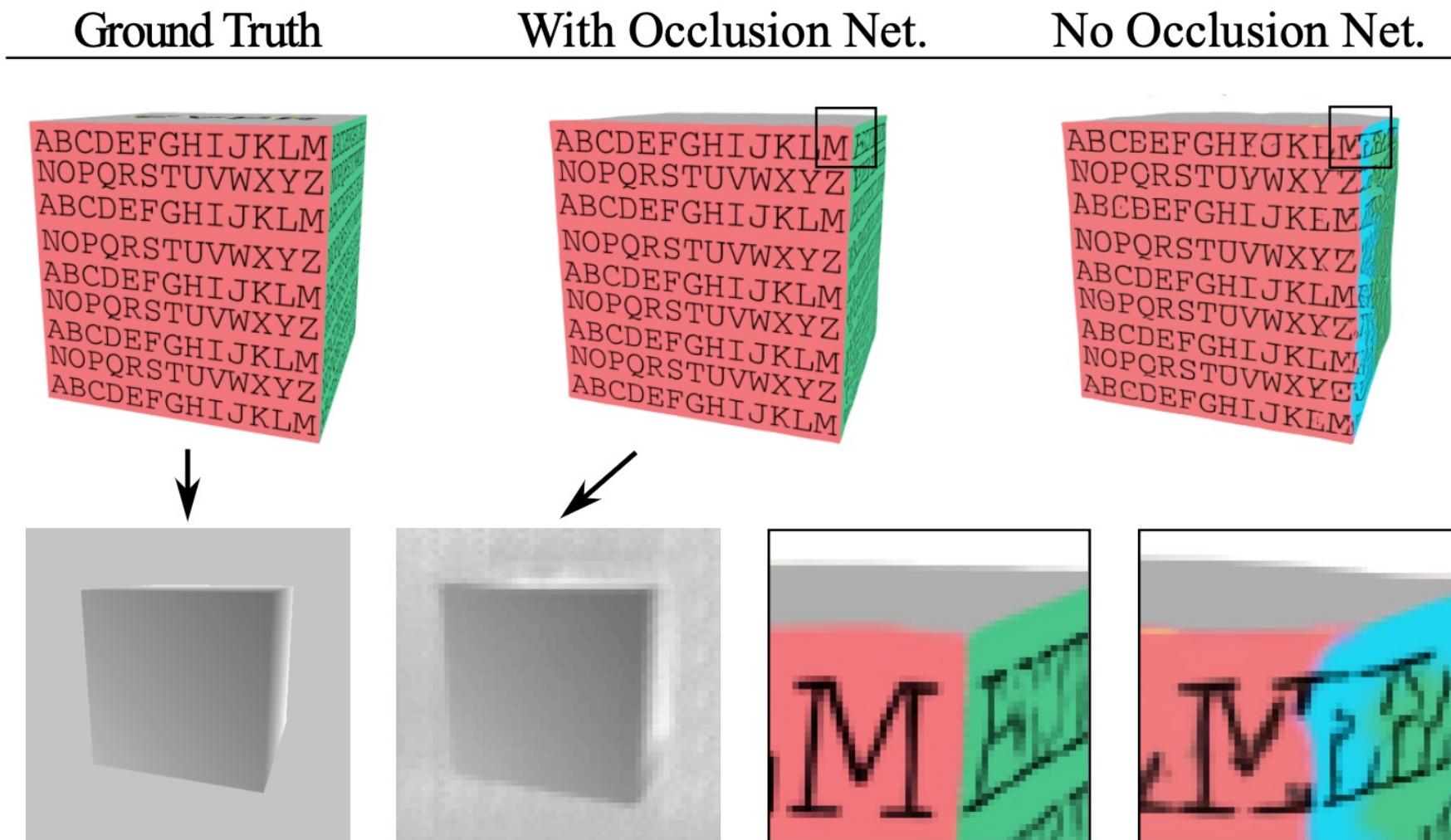
Comparison with Baselines



DeepVoxels: Learning Persistent 3D Feature Embeddings



Ablation Study – Occlusion Reasoning



Deep Voxels – Summary

- Encoding appearance information into the **latent features** of **3D voxels**.
- No 3D geometry supervision is needed.
- Occlusion reasoning leads to better results and generalization to novel viewpoints.

DeepVoxels – Limitations

- Inherits the limitations of using 3D voxels:
 - Memory-heavy
 - Cannot increase the resolution
 - Occlusion-aware projection is **not physically correct.**
- Q.** Why? Any example?

Neural Radiance Fields (NeRF)

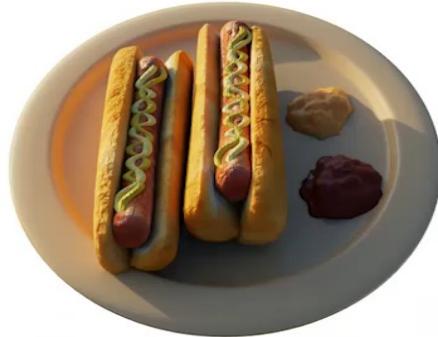
Mildenhall et al., NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV 2020.

Main Differences from DeepVoxels

- Voxels → Implicit function.
- Occlusion-aware projection → Classic volume rendering.
- Latent features → **View-dependent** colors.



Sample Results



Radiance Field

Definition: The **radiance** (luminance) is the power emitted, reflected, transmitted or received by a surface, **per unit solid angle**, per unit projected area.

Roughly speaking, view-dependent color information.

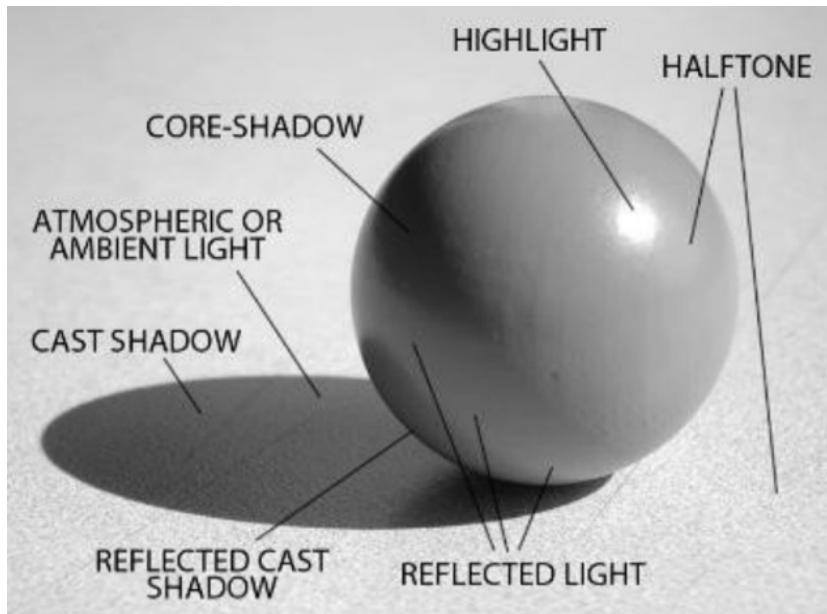
Q. Why?



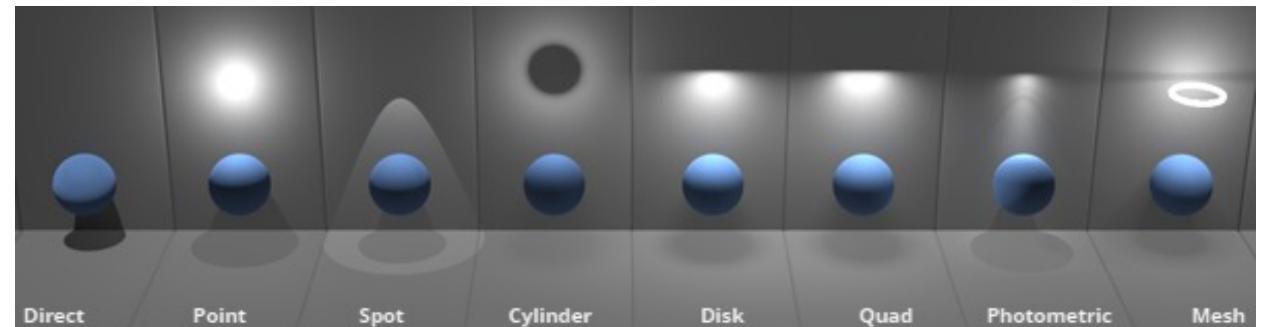
Image Credit: Ren Ng

Radiance Field

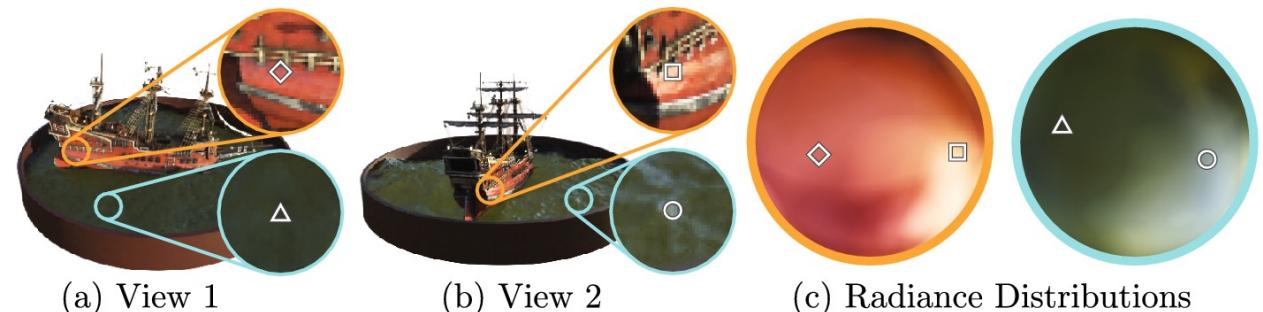
Color can be changed depending on the lighting condition, view direction, material property, and environment.



https://courses.byui.edu/art110_new/art110/glossary/Images/Light.png



<https://docs.arnoldrenderer.com/display/A5AFMUG/Lights>

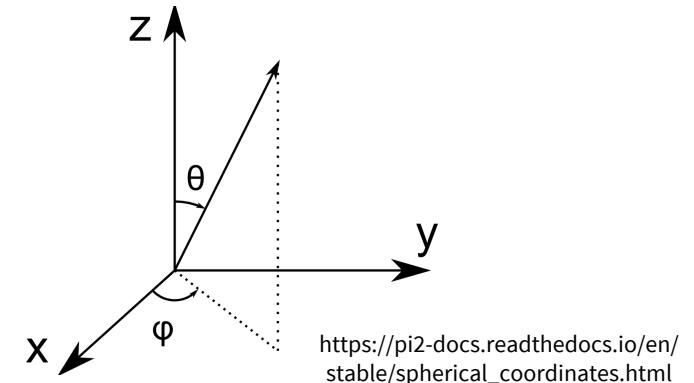


(c) Radiance Distributions

Implicit Functions

Learn two implicit functions using neural networks:

- **Radiance map $c(x, d)$:** $\mathbb{R}^5 \rightarrow \mathbb{R}^3$
 - Input: XYZ 3D location $x \in \mathbb{R}^3$ and viewing direction $d \in \mathbb{R}^2$
 - Output: RGB color
- **Opacity (Density) map $\sigma(x)$:** $\mathbb{R}^3 \rightarrow \mathbb{R}_{\geq 0}$
 - Input: XYZ 3D location $x \in \mathbb{R}^3$
 - Output: Opacity



https://pi2-docs.readthedocs.io/en/stable/spherical_coordinates.html

Volume Rendering

The color of the ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is computed as:

$$C(\mathbf{r}) = \int_0^{t_{max}} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt,$$

where $T(t) = \exp\left(-\int_0^t \sigma(\mathbf{r}(s)) ds\right)$

Transparency:

How likely the ray is to NOT
be terminated until time t .

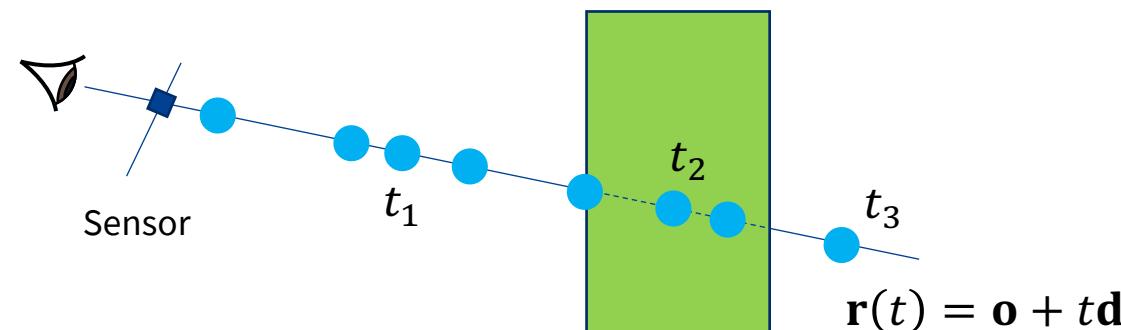
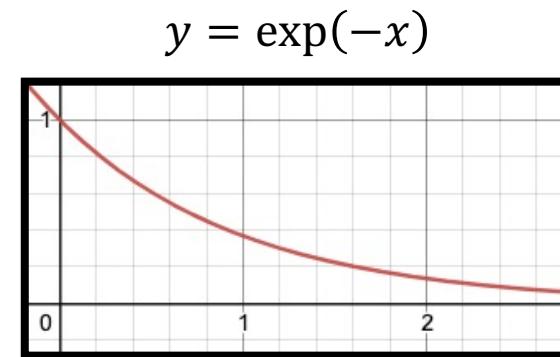


Image credit: In-Young Cho

Volume Rendering

- The transparency function $T(t)$ is defined as a **non-increasing** function.
- It aligns with the **physical concept** that transparency cannot be increased over time along a ray.

$$T(t) = \exp\left(-\int_0^t \sigma(\mathbf{r}(s))ds\right)$$



Volume Rendering

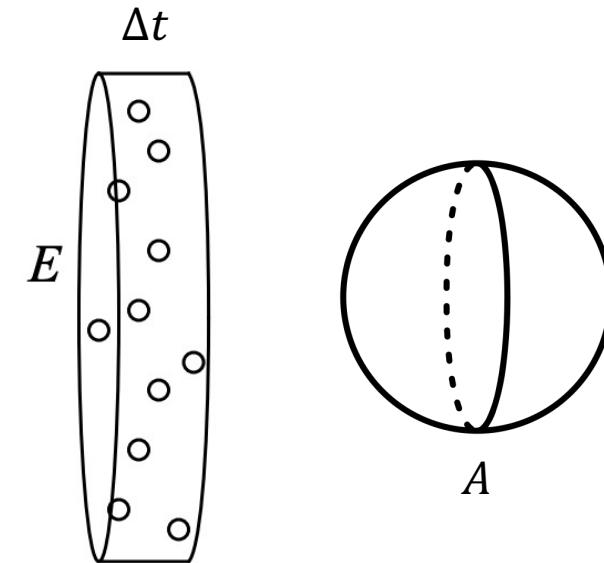
- $\rho(t)$: Particle density.
- $\rho(t)E\Delta t$: Number of particles.
- A : Projected particle area.
- $\frac{A\rho(t)E\Delta t}{E} = A\rho(t)\Delta t$: Ratio of the occluded area.
- **$\sigma(t) := A\rho(t)$: Opacity (extinction coefficient).**
- $T(t)$: Transparency

$$T(t + \Delta t) = T(t)(1 - \sigma(t)\Delta t) = T(t) - T(t)\sigma(t)\Delta t$$

$$\Delta T(t) = T(t + \Delta t) - T(t) = -T(t)\sigma(t)\Delta t$$

Solve ODE with the initial condition $T(0) = 1$:

$$T(t) = \exp\left(-\int_0^t \sigma(t)dt\right)$$



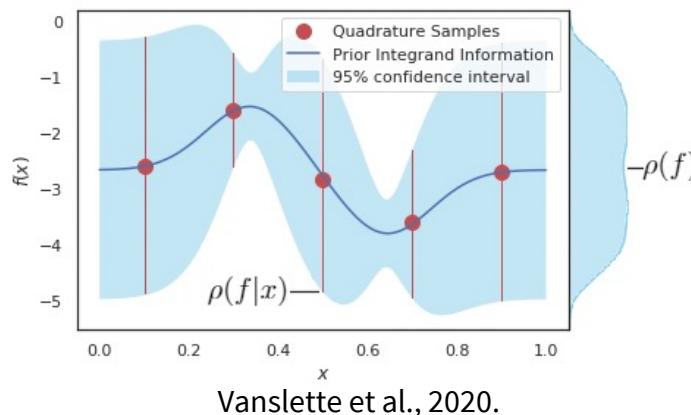
Volume Rendering – Discretization

Use a **stratified sampling**.

Partition the ray range into **evenly-spaced bins** and draw one sample uniformly at random from within each bin.

$$T_i = \exp\left(-\int_0^{t_i} \sigma(s)ds\right) = \exp\left(-\sum_{j=1}^i \sigma_j \Delta t\right)$$

*Transparency until the end
of the i -th bin.*



*Assume that σ_j and c_i are constant
within each bin.*

Vanslette et al., 2020.

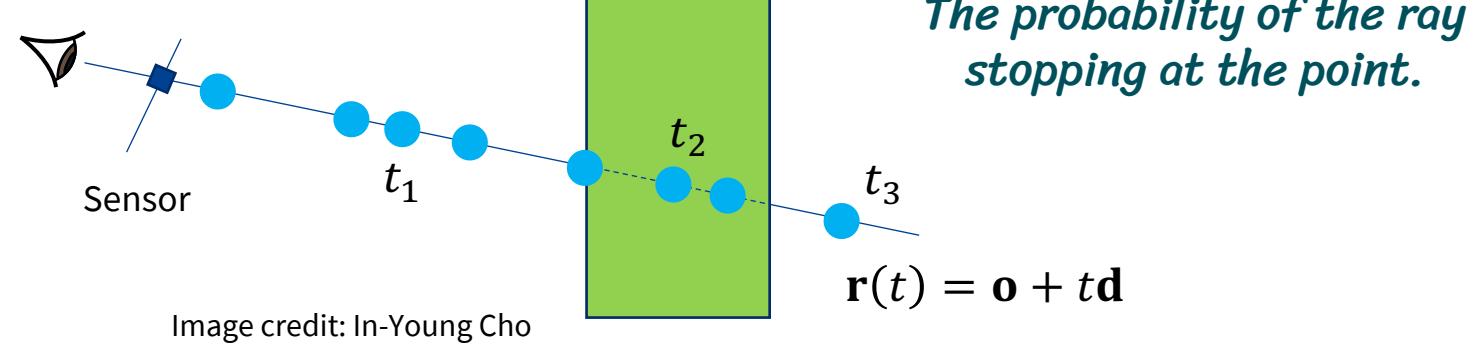
Volume Rendering – Discretization

The color of the ray $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$ is computed as:

$$C(\mathbf{r}) = \int_0^{t_{max}} T(t)\sigma(\mathbf{r}(t))\mathbf{c}(\mathbf{r}(t), \mathbf{d})dt$$

$$\approx \sum_{i=1}^N \int_{t_i}^{t_{i+1}} T(t)\sigma_i \mathbf{c}_i dt = \sum_{i=1}^N \mathbf{c}_i \int_{t_i}^{t_{i+1}} T(t)\sigma_i dt$$

w_i



Volume Rendering – Discretization

Approximation via Quadrature

Integral for each bin

$$\begin{aligned}
 C_i &= c_i \int_{t_i}^{t_{i+1}} \sigma_i T(t) dt = c_i \int_{t_i}^{t_{i+1}} \sigma_i \exp\left(-\int_0^t \sigma(s) ds\right) dt \\
 &= c_i \int_{t_i}^{t_{i+1}} \sigma_i \exp\left(-\int_0^{t_i} \sigma(s) ds\right) \exp\left(-\int_{t_i}^t \sigma(s) ds\right) dt \\
 &= c_i T_i \int_{t_i}^{t_{i+1}} \sigma_i \exp\left(-\sigma_i \int_{t_i}^t ds\right) dt \\
 &= c_i T_i \int_{t_i}^{t_{i+1}} \sigma_i \exp(-\sigma_i(t - t_i)) dt \\
 &= c_i [T_i(1 - \exp(-\sigma_i \Delta t))]^{w_i}
 \end{aligned}$$

For a more precise approximation, compute $T(t)$ based on σ_j and c_i instead of assuming that $T(t)$ is constant within each bin.

The probability of the ray stopping at the point.

Volume Rendering – Discretization

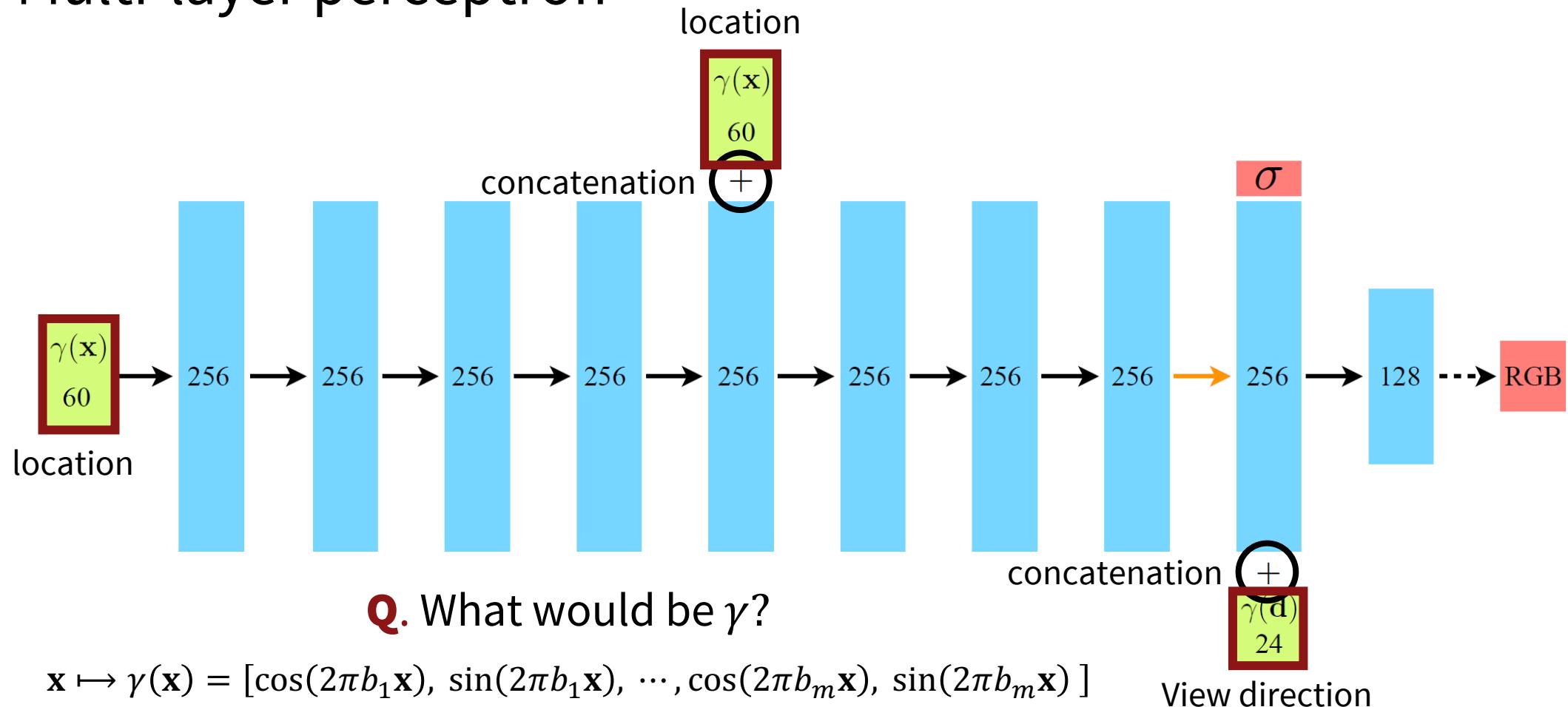
Approximation via Quadrature

$$C(\mathbf{r}) = \sum_{i=1}^N \mathbf{c}_i w_i = \sum_{i=1}^N \mathbf{c}_i T_i (1 - \exp(-\sigma_i \Delta t)),$$

where $T_i = \exp\left(-\sum_{j=1}^{i-1} \sigma_j \Delta t\right)$

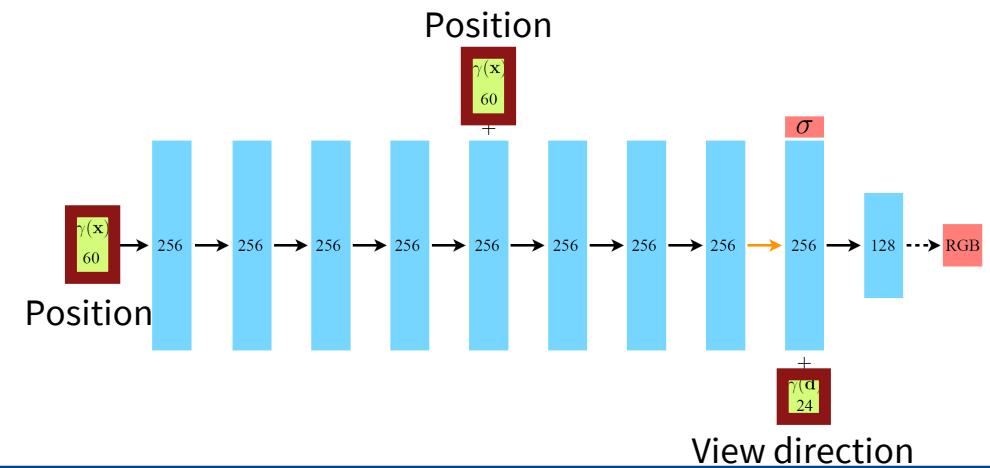
Network Structure

Multi-layer perceptron



Network Structure

- **Position**
 - Mapped to a **higher-dimensional** space through positional encoding and processed by a **deeper** network.
 - Want to learn a **high-frequency** (rapidly fluctuating) **opacity** function.
- **View Direction**
 - Mapped to a **lower-dimensional** space through positional encoding and processed by a **shallower** network.
 - Want to learn a **low-frequency** (smooth) **radiance** function w.r.t. the view direction.

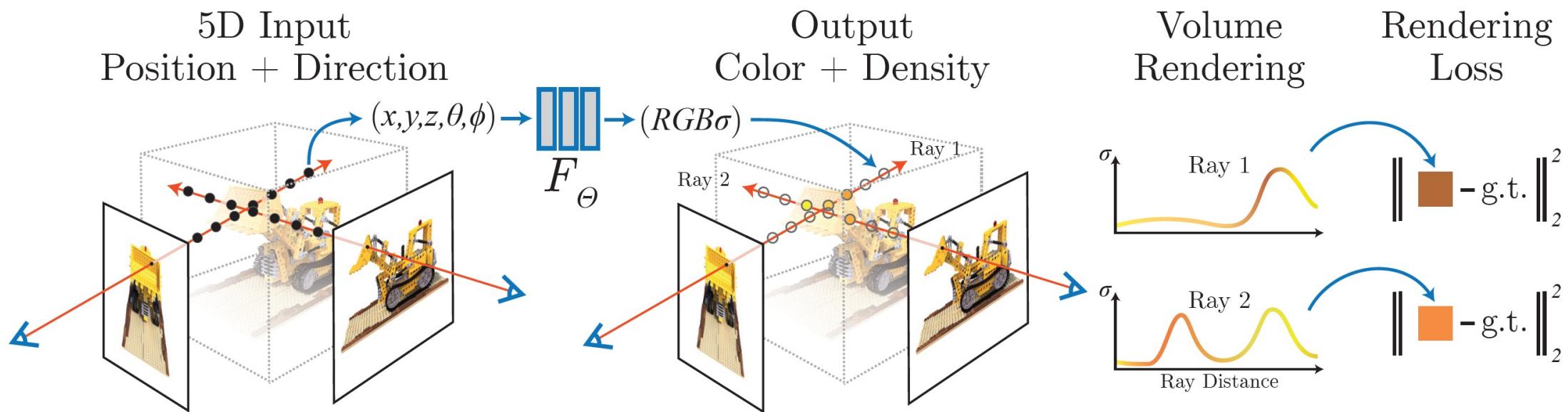


Hierarchical Sampling

- **Two networks:** one for **coarse** sampling, and the other for **fine** sampling.
- First, sample N_c points with the stratified sampling and compute the probability w_i at each point using the **coarse** network.
- Second, sample N_f more points **adaptively** with the probability w_i and compute the final pixel color using the **fine** network and all the $N_c + N_f$ points.

Training

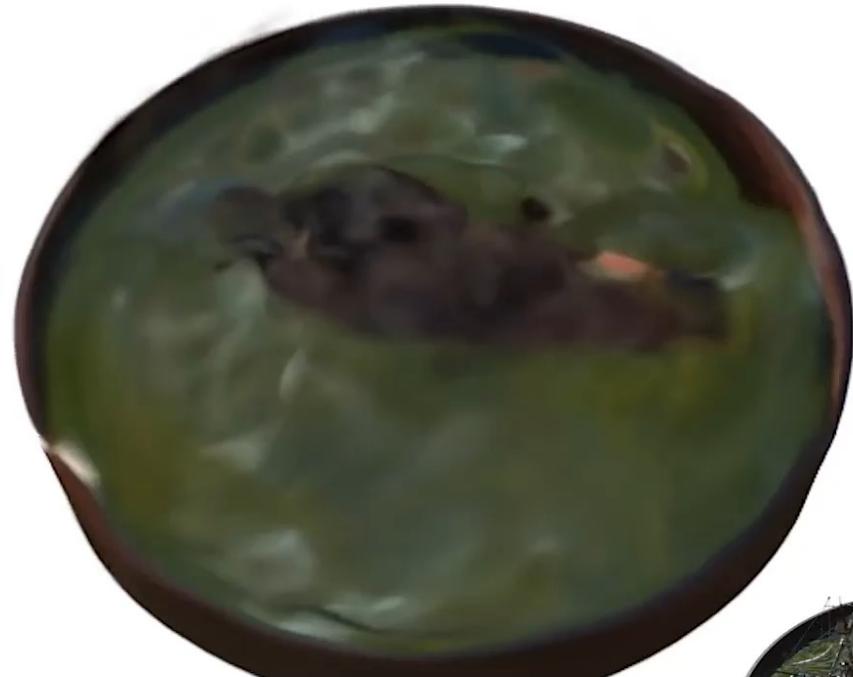
1. Batch pixels.
2. Sample points along rays.
3. Predict color and density at each point.
4. Integrate the outputs w/ volume rendering and compute loss.



Comparisons

Can reconstruct **high-frequency details** while using
much smaller memory (< 5MB MLP parameters)!

SRN [Sitzmann 2019]



NeRF



Voxel-based!

Nearest Input

Geometry Visualization



NeRF – Limitations (1/2)

- Slow speed.
 - Training: 1-2 days for a NVIDIA V100 GPU.
 - Inference: A few seconds for a FHD image.
- *KiloNeRF, FastNeRF, Plenoxels, HashNeRF...*
- Not scalable to large scenes.
- *Neural Scene Graph Rendering, NeRF-OSR, CitiNeRF, Block-NeRF...*

NeRF – Limitations (2/2)

- Requires accurate **camera** information
→ *NeRF in the Wild, BARF, SCNeRF, GNeRF, iNeRF, NeRF--...*
- Cannot handle **dynamic** scenes / moving objects.
→ *Nerfies, HyperNeRF, NeRFlow, D-NeRF, AD-NeRF, NR-NeRF...*
- Cannot change the lighting.
→ *NeRF-W, NeRD, NeRV...*
- Often fails to **precisely reconstruct** the shape.
→ *VolSDF...*
- Cannot generalize to **unseen** scenes – one scene, one model.
→ *GIRAFFE, Pi-GAN, StyleNeRF, EG3D...*

NeRF – More

- Improving output quality
→ *MIP-NeRF, MVSNeRF, DietNeRF, UNISURF, NerfingMFS...*
- Conditional
→ *CRF, GSN, GANcraft, CodeNeRF...*
- Composition
→ *EditNeRF...*
- Articulated
→ *NARF, AnimatableNeRF...*

References:

<https://dellaert.github.io/NeRF/>
<https://dellaert.github.io/NeRF21/>