

Virtual Humans – Winter 23/24

Lecture 10_1 – Humans and NeRF

Prof. Dr.-Ing. Gerard Pons-Moll

University of Tübingen / MPI-Informatics

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



Goal: Novel View Synthesis



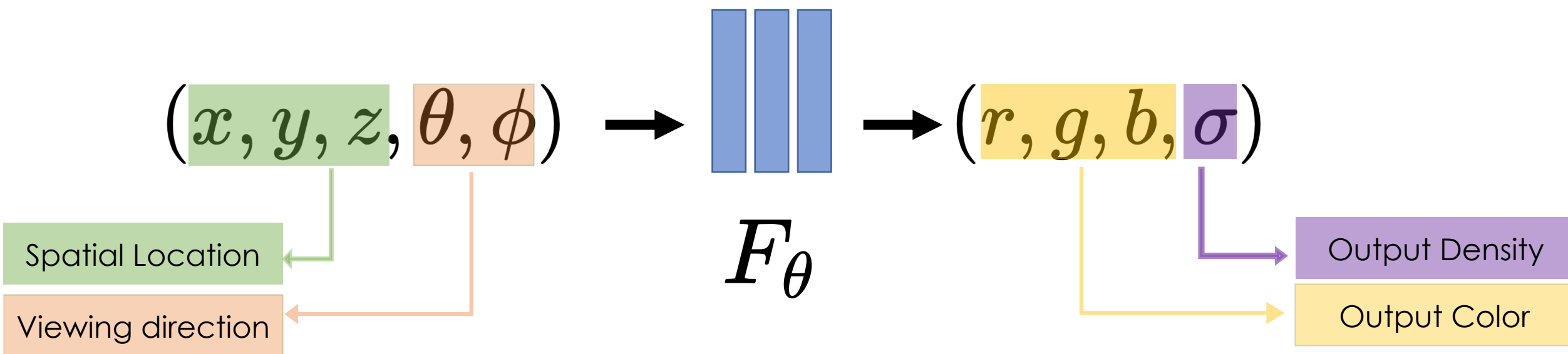
Input:
Sparsely sampled
images of the scene

Learn scene
representation

Novel view synthesis
of the scene

Novel View Synthesis

Learning radiance field representation of scene:



Volume Rendering

Given color and density (r, g, b, σ) , we calculate the color of every camera ray using:

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t)) dt$$

Camera ray: $\mathbf{r}(t) = \mathbf{o} + t\mathbf{d}$

Near and far bounds

Color (r,g,b) at $\mathbf{r}(t)$

Volume density: Probability of a ray terminating at an infinitesimal particle at location $\mathbf{r}(t)$

Accumulated transmittance along ray :
 $T(t) = \exp\left(-\int_{t_n}^{t_f} \sigma(\mathbf{r}(s)) ds\right)$

Volume Rendering in NeRF

Given color and density (r, g, b, σ) , we calculate the color of every camera ray using:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i$$

Uniform N samples

$$t_i \sim \mathcal{U}\left[t_n + \frac{i-1}{N}(t_f - t_n), t_n + \frac{i}{N}(t_f - t_n)\right]$$

Accumulated transmittance along ray

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

Distance between adjacent samples

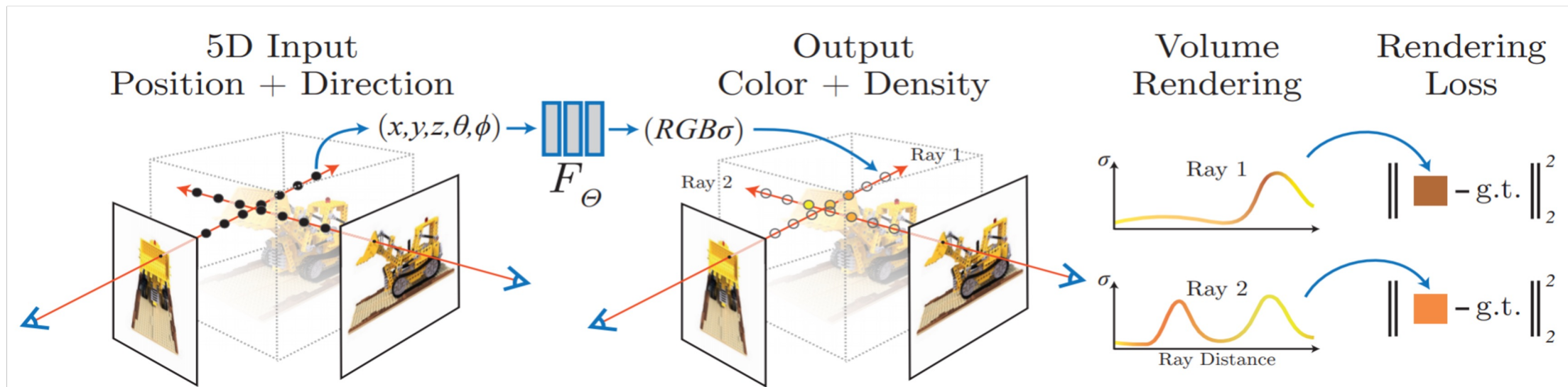
$$\delta_i = t_{i+1} - t_i$$

Alpha (in traditional alpha compositing)

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$

Color (r,g,b) at t_i

Neural Radiance Fields



Training NeRF:

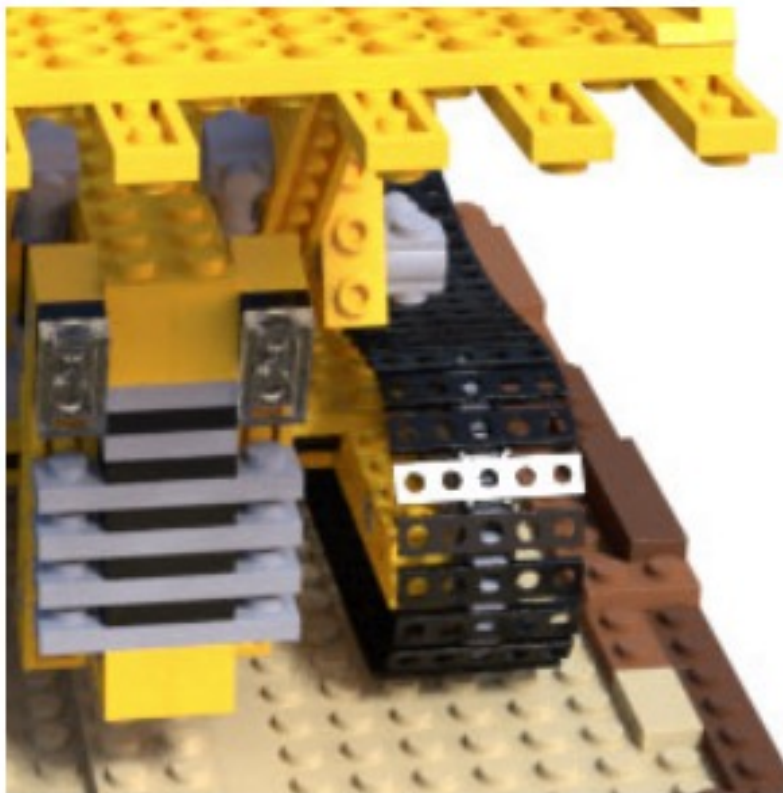
- 1) **March camera rays** through the scene to generate a **sampled set of 3D points**
- 2) Use those **points and** their corresponding **2D viewing directions** as input to the neural network to produce an output **set of colors and densities**
- 3) Use classical volume rendering techniques to **accumulate those colors and densities** into a 2D image
- 4) **Minimize** error between **rendered color and GT color**

Neural Radiance Fields

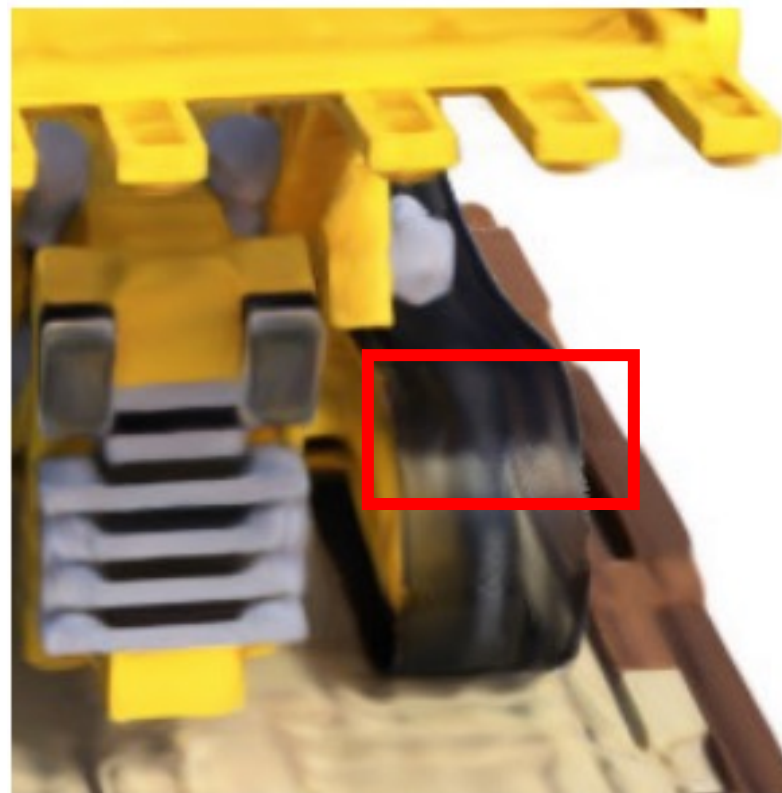


Novel view synthesis using NeRF

Generated results are blurry



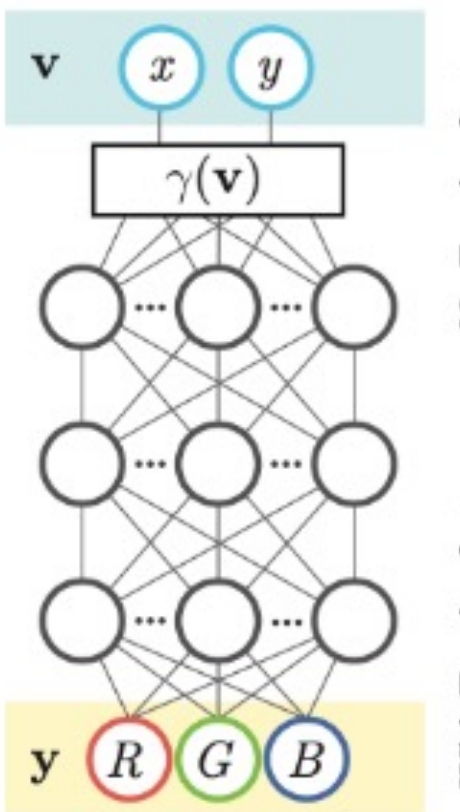
Ground Truth



No Positional Encoding

Why blurry results?

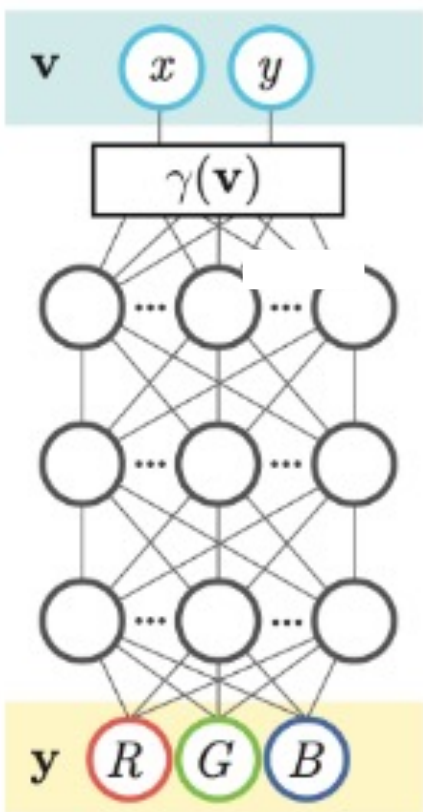
Can coordinate based MLP learn high-frequency details?



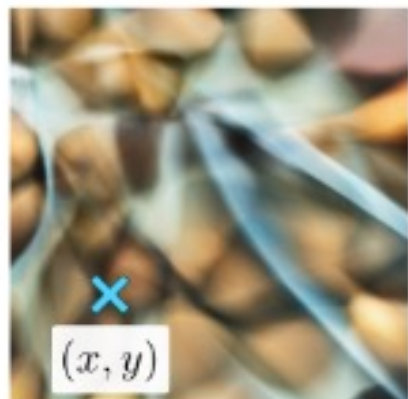
(a) Coordinate-based MLP

Why blurry results?

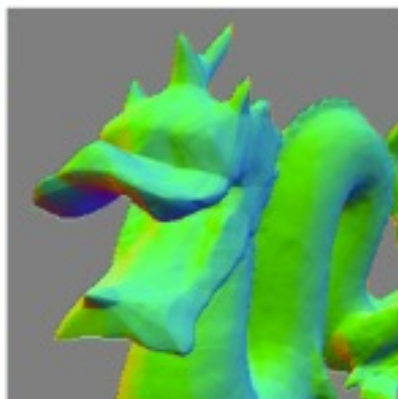
Coordinate based neural network fail to learn high frequency details for all kind of data including RGB image, 3D shape, density , etc



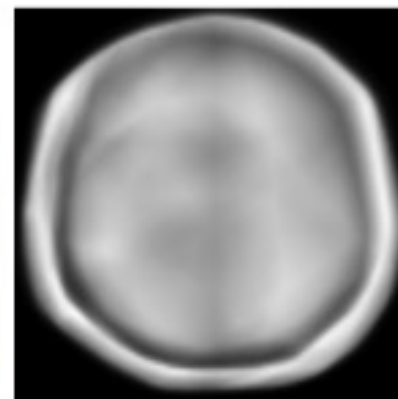
(a) Coordinate-based MLP



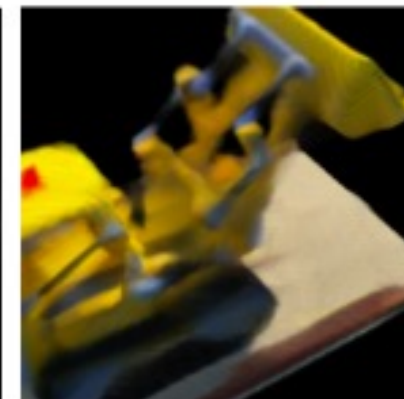
RGB



3D Shape



Density

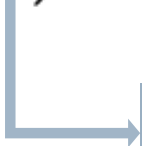


Radiance field
(density, color)

Solution:

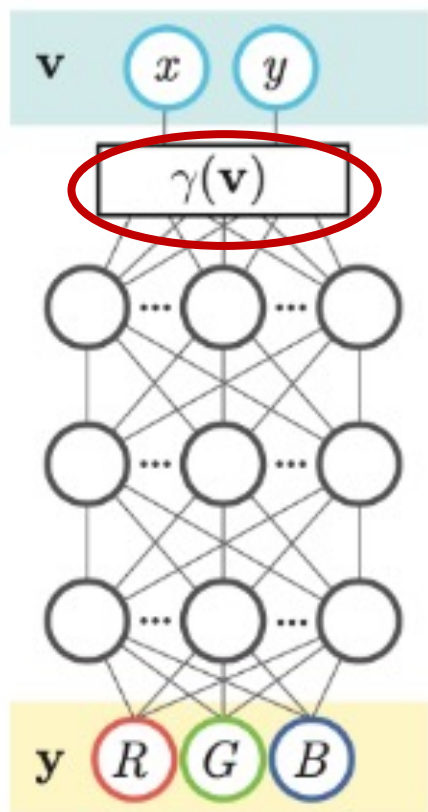
- In naive setting, the bandwidth of the Neural Tangent Kernel limits the spectrum of the recovered/learned function.
- Using a Fourier feature mapping transforms the neural kernel into a stationary kernel in our low-dimensional problem domains and increase the spectrum.

$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p))$$



Coordinate input e.g. pixel location for images, 3D point for NeRF

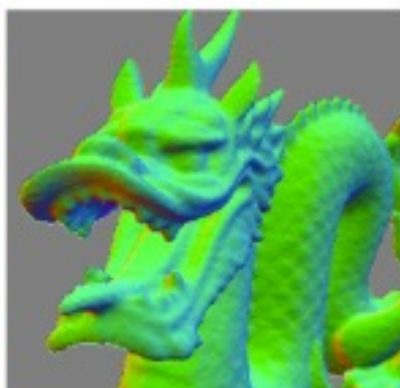
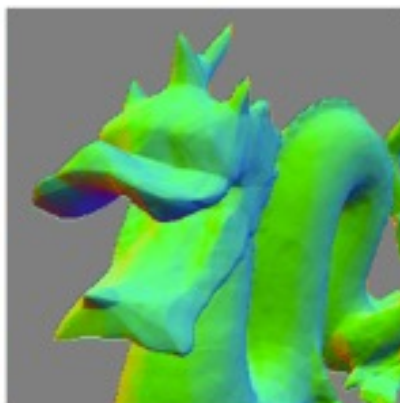
Fourier Features in Coordinate MLPs



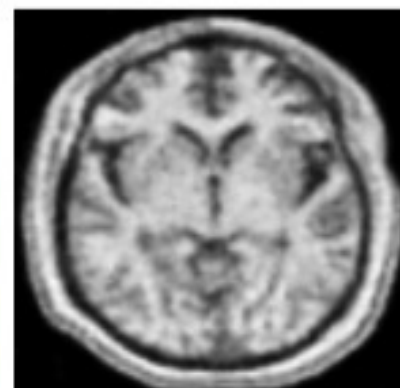
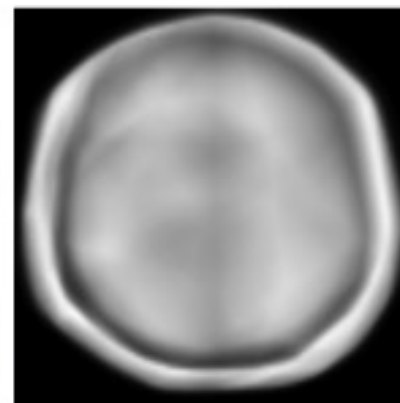
(a) Coordinate-based MLP



RGB



3D Shape



Density



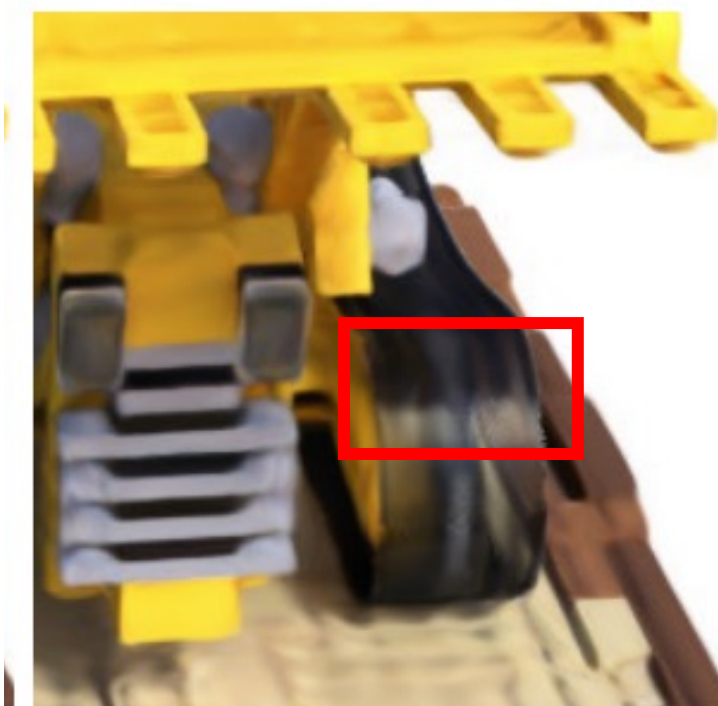
Radiance field
(density, color)

Positional Encoding in NeRF

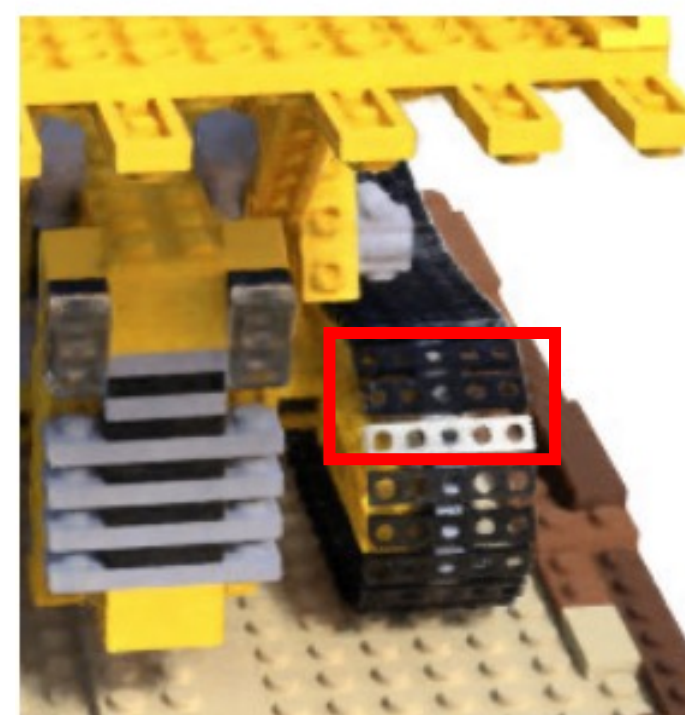
With fourier features/positional encoding, NeRF learns high frequency details.



Ground Truth



No Positional Encoding



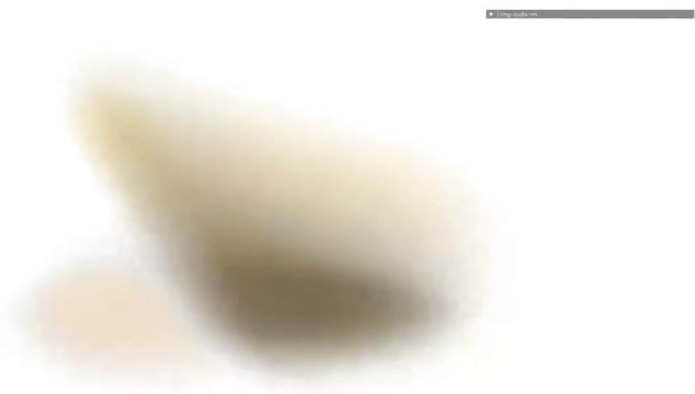
Complete Model

Tannick et al. NeurIPS 20

Positional Encoding in NeRF

With fourier features/positional encoding, NeRF learns high frequency details.

Without positional encoding



Position + View dir.



Neural network
10 layers
128 neurons



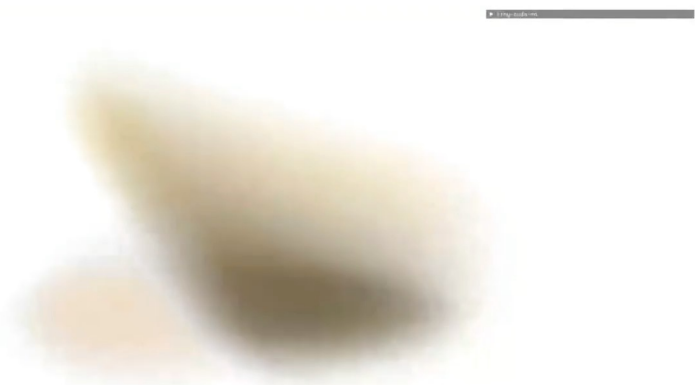
Density & Color



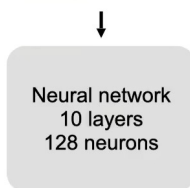
Positional Encoding in NeRF

With fourier features/positional encoding, NeRF learns high frequency details.

Without positional encoding



Position + View dir.



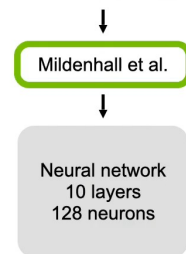
Density & Color



With positional encoding



Position + View dir.

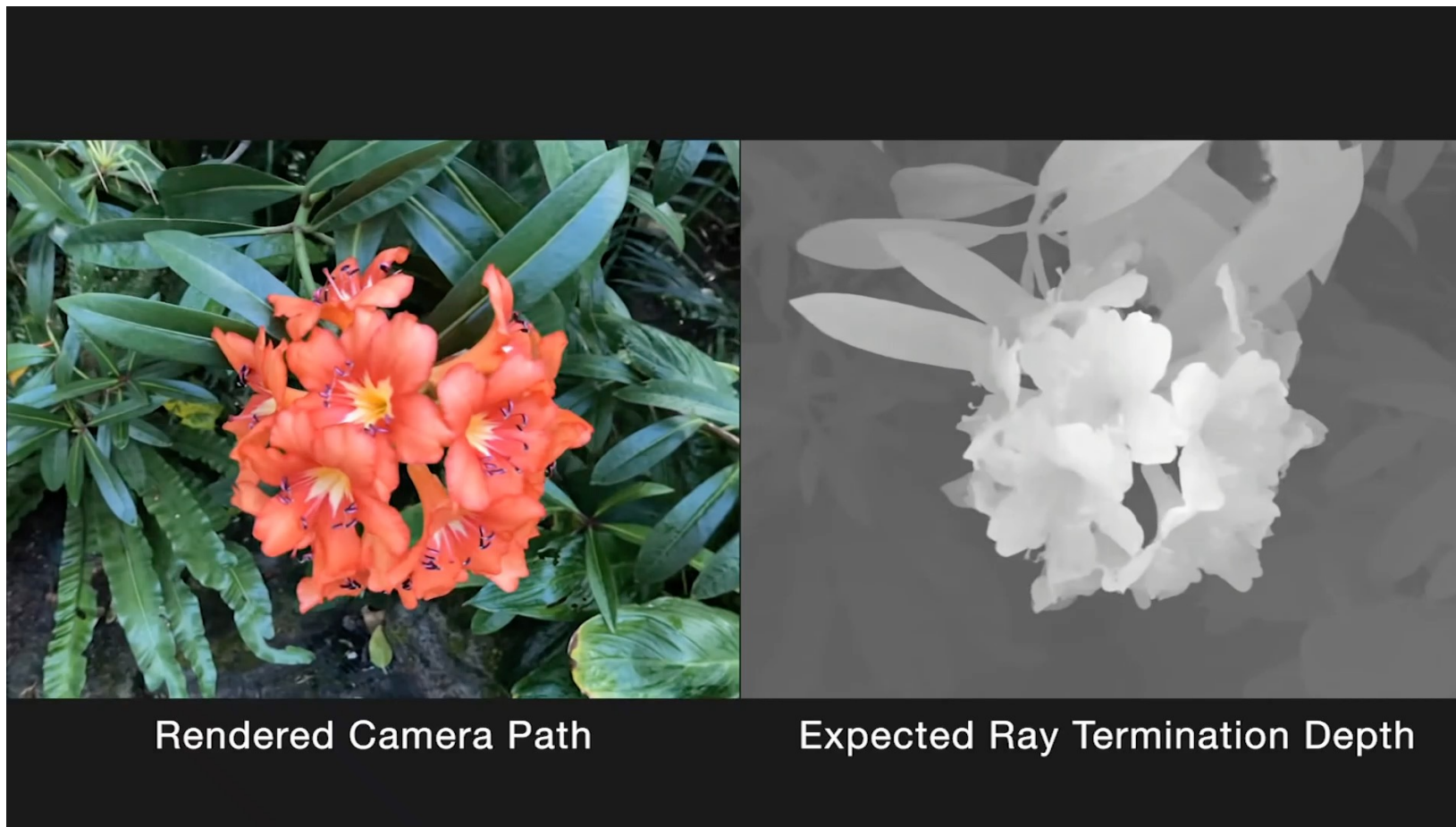


Density & Color



Geometry in NeRF

Scene geometry can be approximated using threshold



Rendered Camera Path

Expected Ray Termination Depth

Limitations of NeRF

1. Scene specific and only static scene can be modeled.

GT image of a dynamic scene

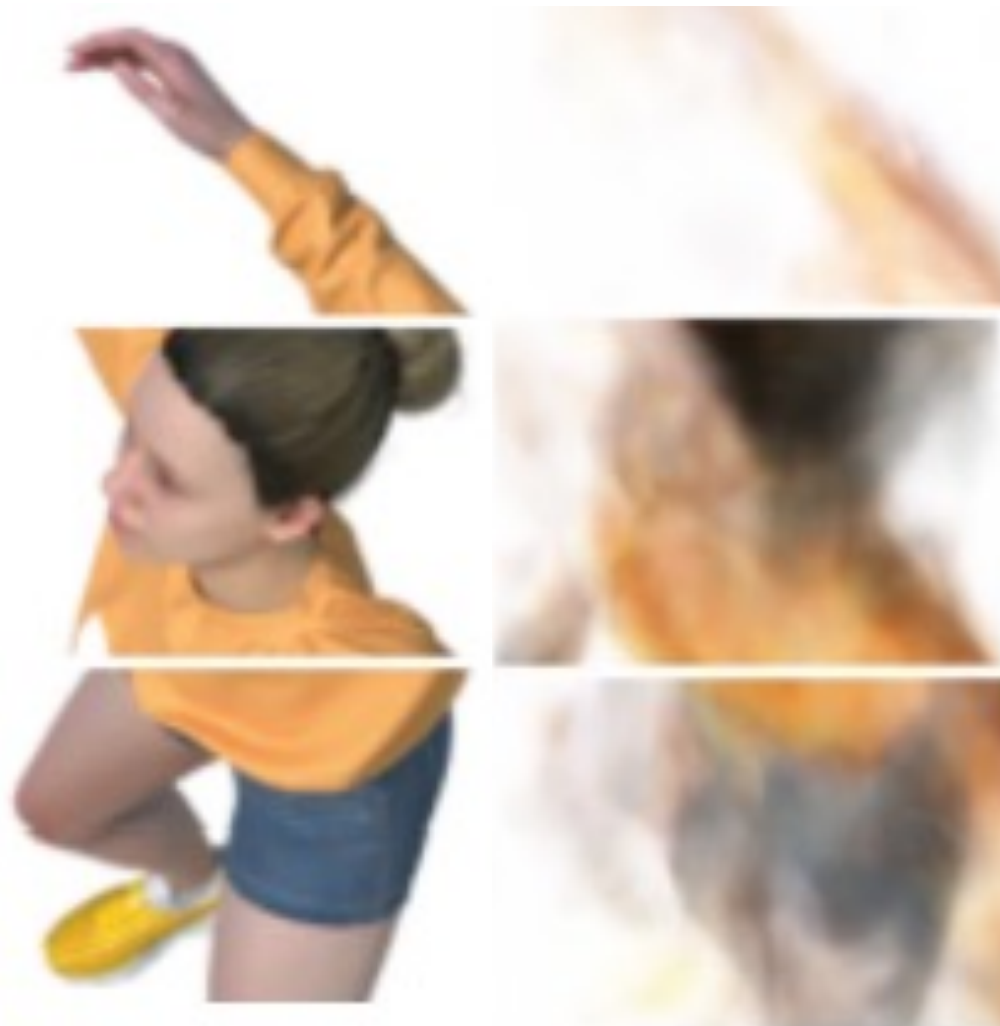


Image generated with NeRF

Limitations of NeRF

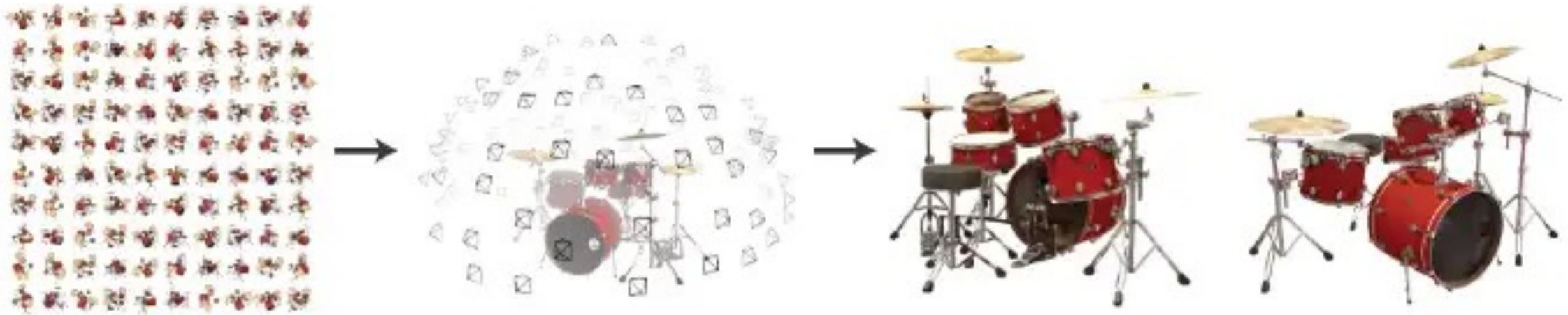
2. No editing and control

- Learned scene cannot be modified.
- Scene is memorized within the network

Limitations of NeRF

3. Generalization

- Scene specific models.
- Large number of images are needed



Limitations of NeRF

4. Expensive training:

- Training is slow(10 hours-up to few days)
- Inference is also not real time

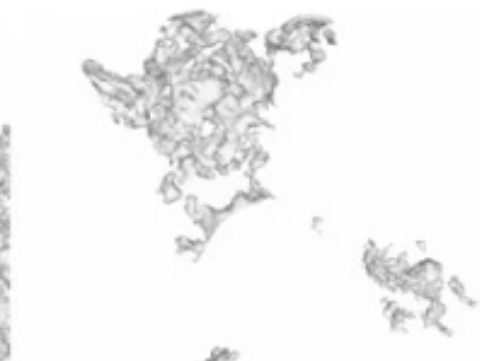
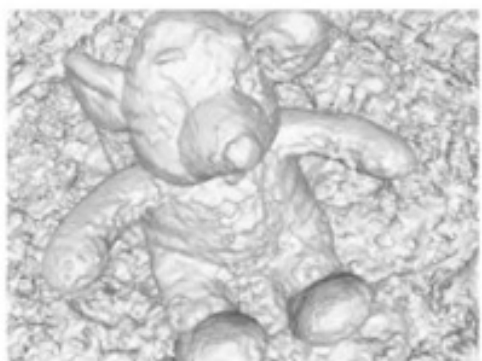
Limitations of NeRF

5. Surface extracted is not accurate and depends on threshold.



GT image

NeRF (density threshold $\sigma = 50$)



$\sigma = 1$

$\sigma = 10$

$\sigma = 50$

$\sigma = 100$

$\sigma = 500$

Limitations of NeRF

1. Scene specific and only static scene can be modeled.

GT image of a dynamic scene



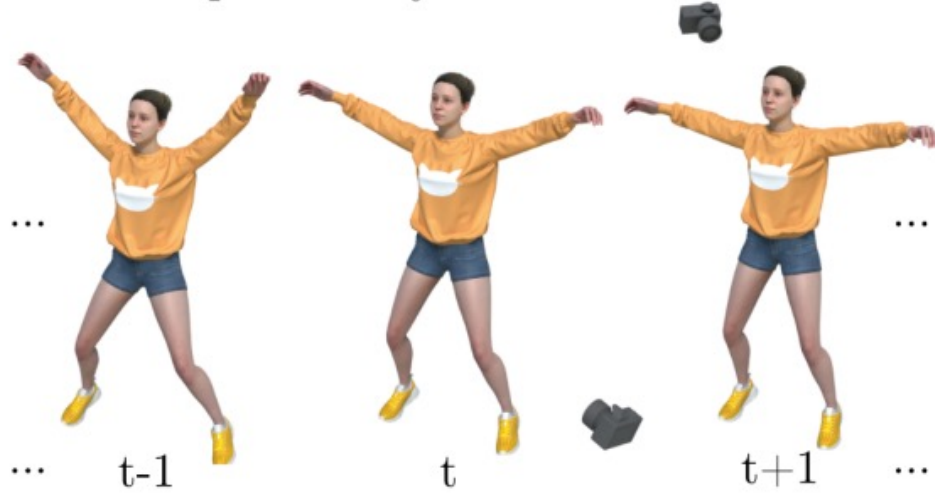
Image generated with NeRF

What about dynamic scenes?

Sparse Monocular RGB



Optimize Dynamic Scene

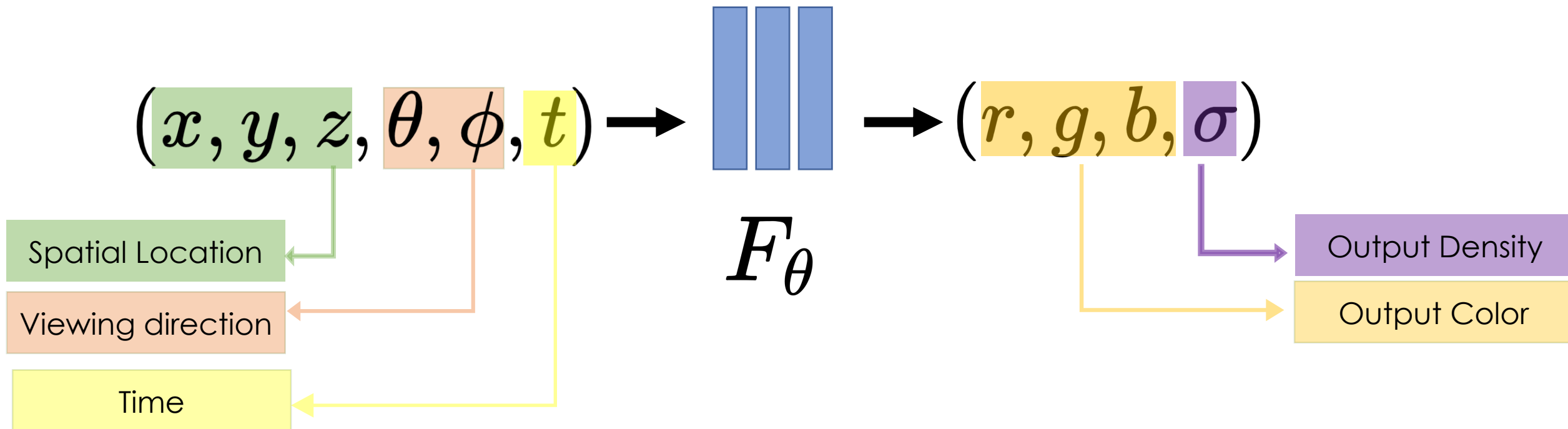


Render New Views at Arbitrary Time



What about dynamic scenes?

Learn radiance field given 3d point, viewing direction and time



Can we learn this mapping directly using NeRF?

Learn radiance field given 3d point, viewing direction and time



Ground truth

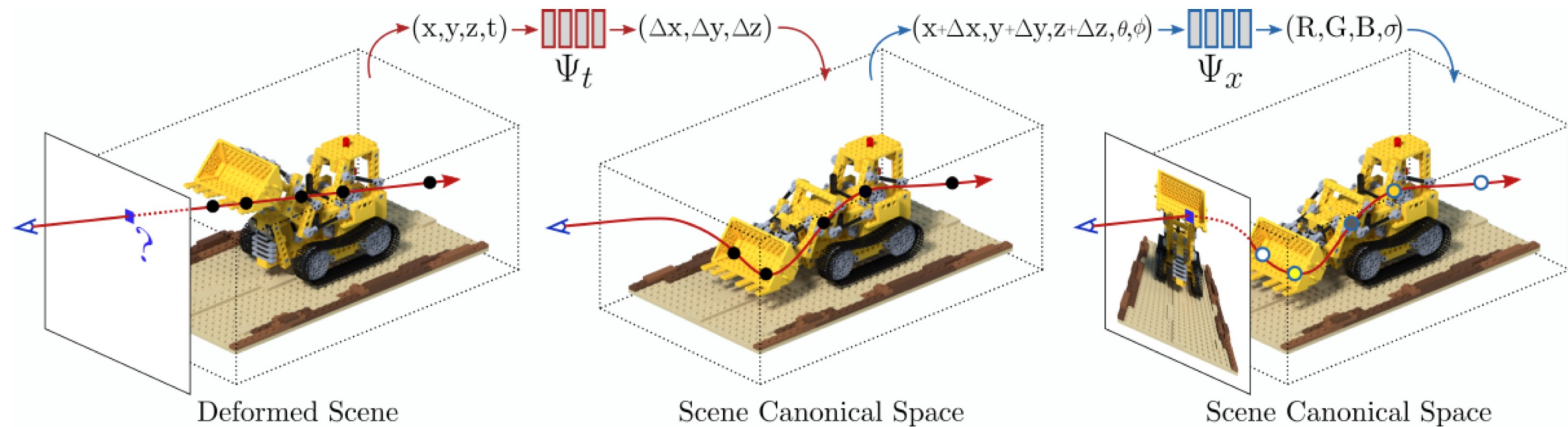


NeRF + time

Pumarola et al. CVPR'21

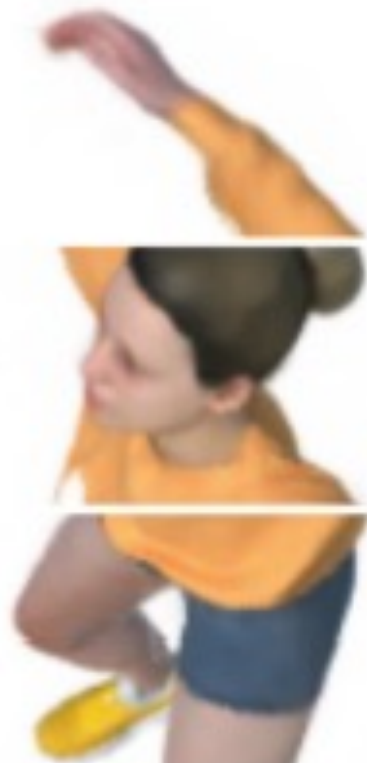
Proposed solution: D-NeRF

Learn canonical shape and radiance field in the canonical shape



Proposed solution: D-NeRF

Learn canonical shape and radiance field in the canonical shape



Ground truth

D-NeRF

NeRF + time

Pumarola et al. CVPR'21

Synthesis Results



D-NeRF



Closest Input View



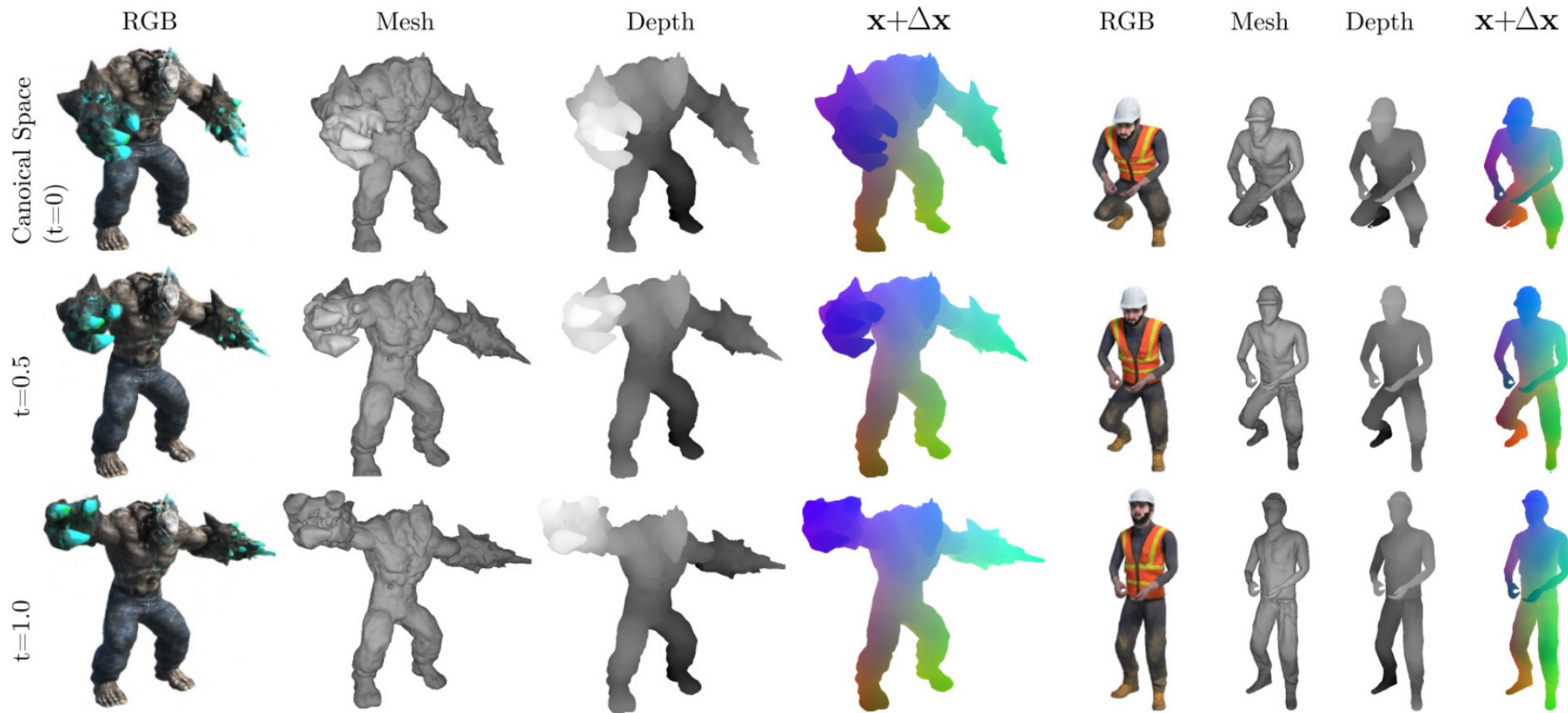
Closest Input Time

Time & View Conditioning:



Pumarola et al. CVPR'21

D-NeRF: Visualization of learned scene representation



Pumarola et al. CVPR'21

Conclusion: Dynamic Scenes with D-NeRF

- Disentangle time dependent deformation from neural rendering network.

Conclusion: Dynamic Scenes with D-NeRF

- Disentangle time dependent deformation from neural rendering network.
- Correspondence between canonical shape and deformed shape is defined by



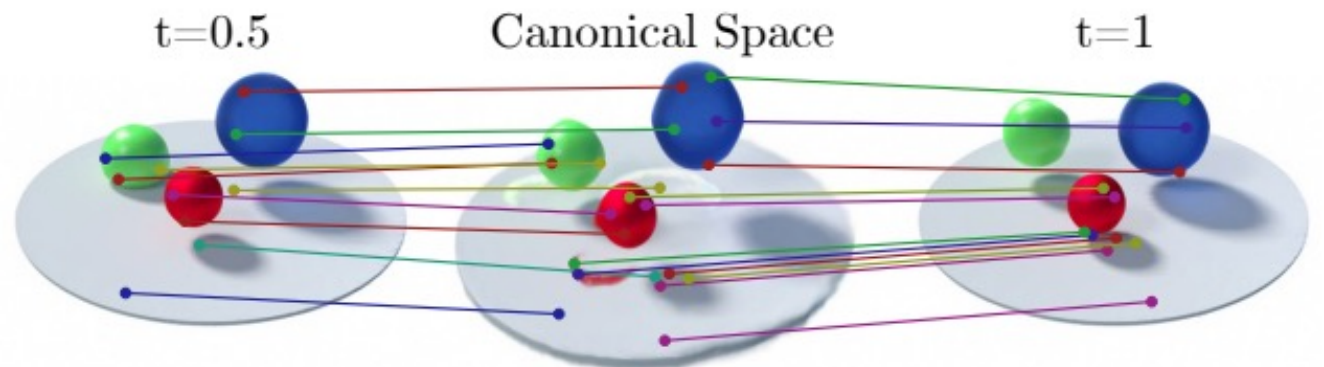
D-NeRF Canonical Mapping
(color-coded as $\mathbf{x} + \Delta\mathbf{x}$)

Conclusion: Dynamic Scenes with D-NeRF

- Disentangle time dependent deformation from neural rendering network.
- Correspondence between canonical shape and deformed shape is defined by
- Time varying shading effects are modeled.



D-NeRF Canonical Mapping
(color-coded as $\mathbf{x} + \Delta\mathbf{x}$)



Limitations of NeRF:

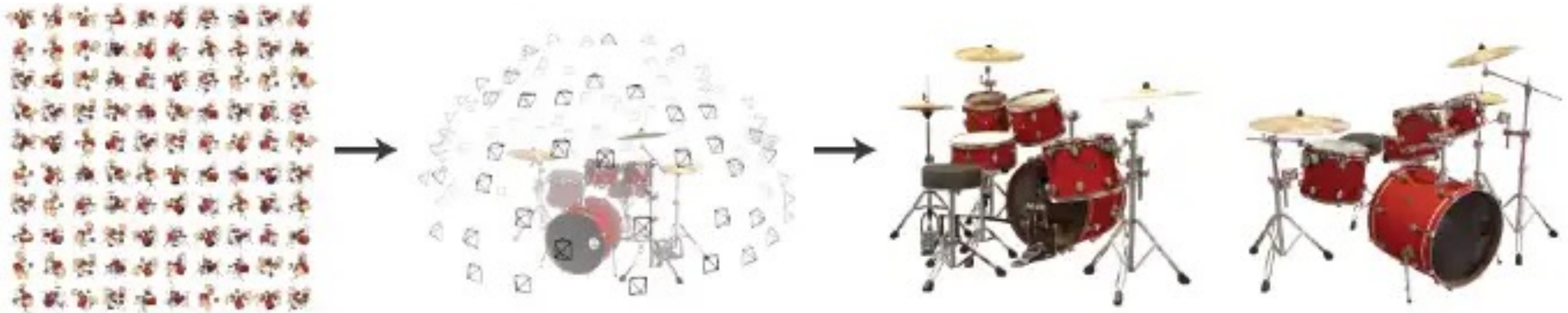
2. No editing and control

- Learned scene cannot be modified.
- Scene is memorized within the network

Limitations of NeRF

3. Generalization

- Scene specific models.
- Large number of images are needed



Control-NeRF

Control NeRF for scene manipulation and rendering



a) Original scenes

T-rex inserted inside the garden scene



b) Inserting objects from one scene into another

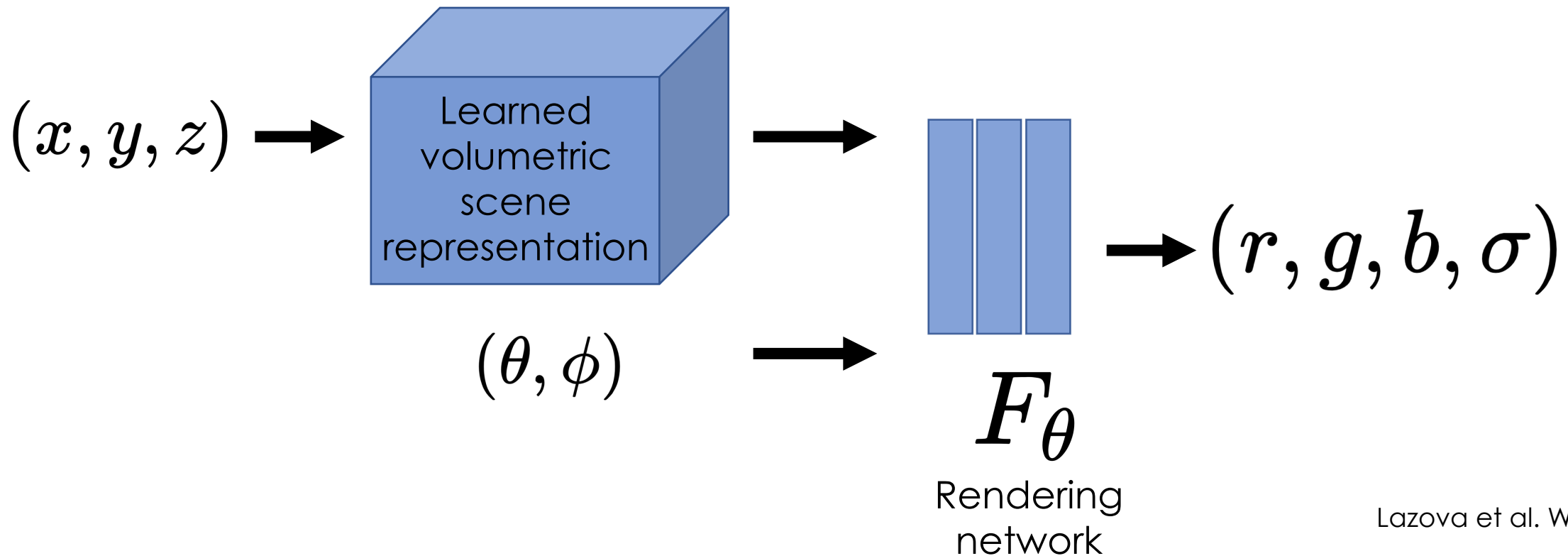
Second T-rex added to the scene



c) Copying and moving objects within the scene

Control-NeRF

- **Prior Work:** Scene is memorized within the neural network, which makes compositing of scenes and editing hard.
- **Key Idea:** Decouple scene representation from neural rendering network.

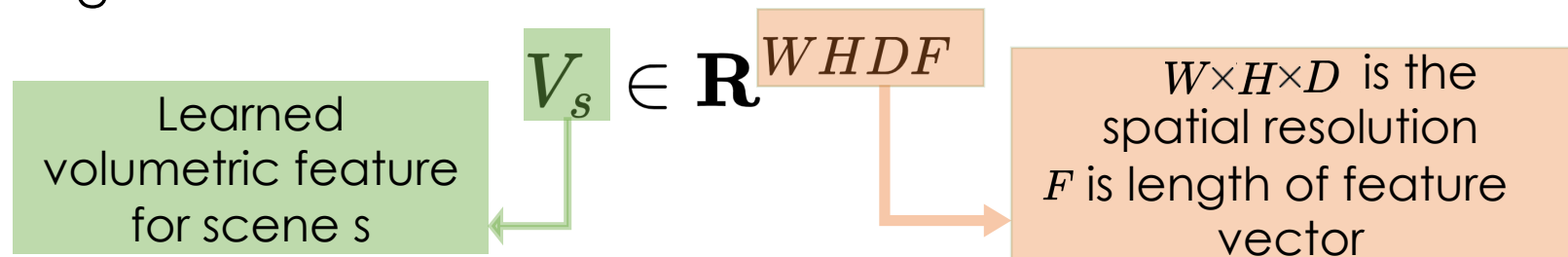


Control-NeRF

1. Scene representation:

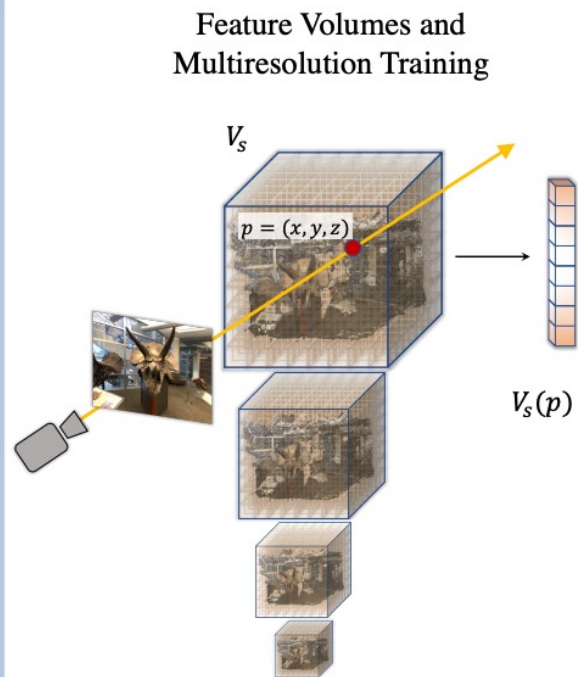


- Given a set of input images $\mathcal{I} = \{I_s^i\}_{i=1}^{N_s}$ from M training scenes
- Where $s \in \mathcal{S}$ is set of training scenes and $M = |\mathcal{S}|$
- Scene representation network learns a volumetric feature
- where $W \times H \times D$ is the spatial resolution of grid which feature vector of length F .



Control-NeRF

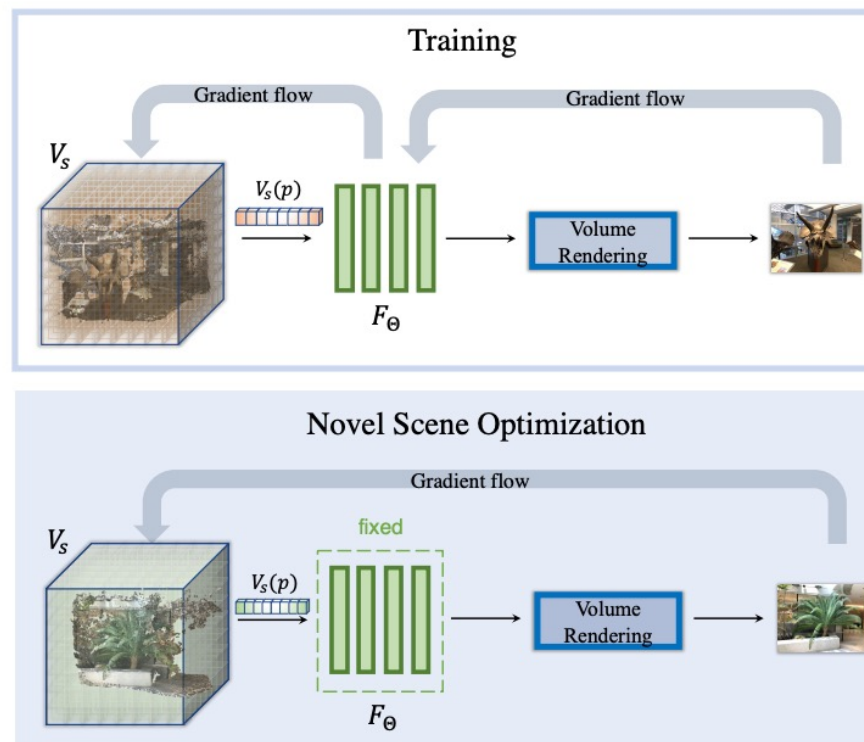
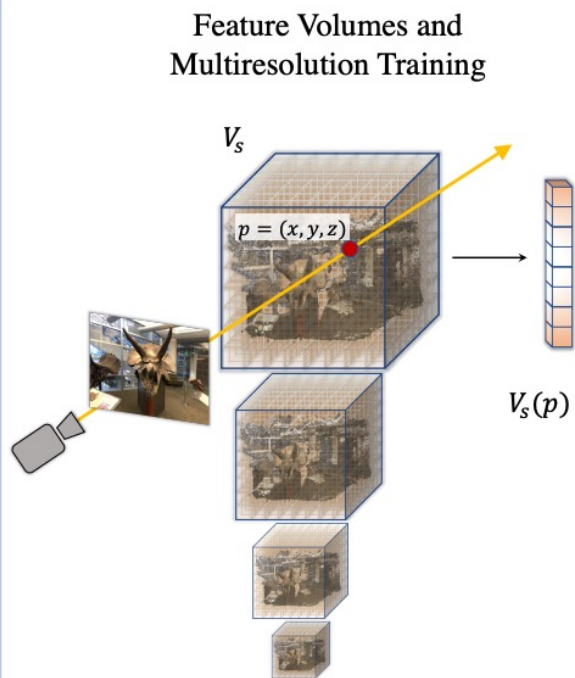
2. Neural rendering network with feature volumes



$V_s(p)$: Volumetric feature at query point p

Control-NeRF

Training and Inference



Control-NeRF: Training

1. Multi-resolution Volume training

- Hierarchical training process is used to compute the volumes in a coarse-to-fine manner.
 - Train low resolution(16^3) volume till convergence.
 - Upsample the learnt feature volume and train till convergence
-
- Improved training time.
 - High-quality image synthesis and manipulation.

Control-NeRF: Training

2. **Multi Scene training**

- Efficient training strategy: Sample one training scene and train for N iterations, before saving the volume grid

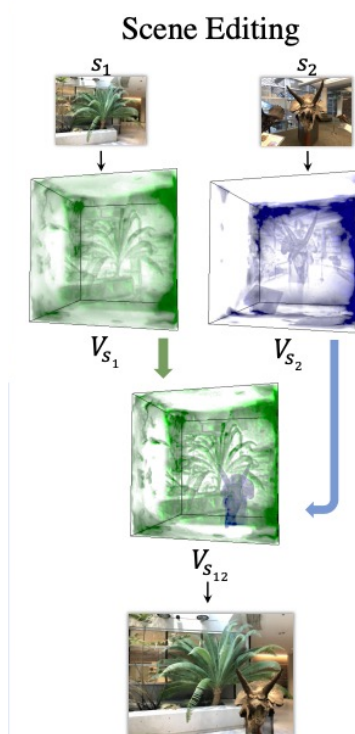
Control-NeRF: Training

3. **Generalization to Novel Scenes**

- Fix neural rendering network and learn feature volume for novel scene.
- Given sufficient training scenes, the learnt radiance function can be applied to optimize for novel scenes more efficiently.

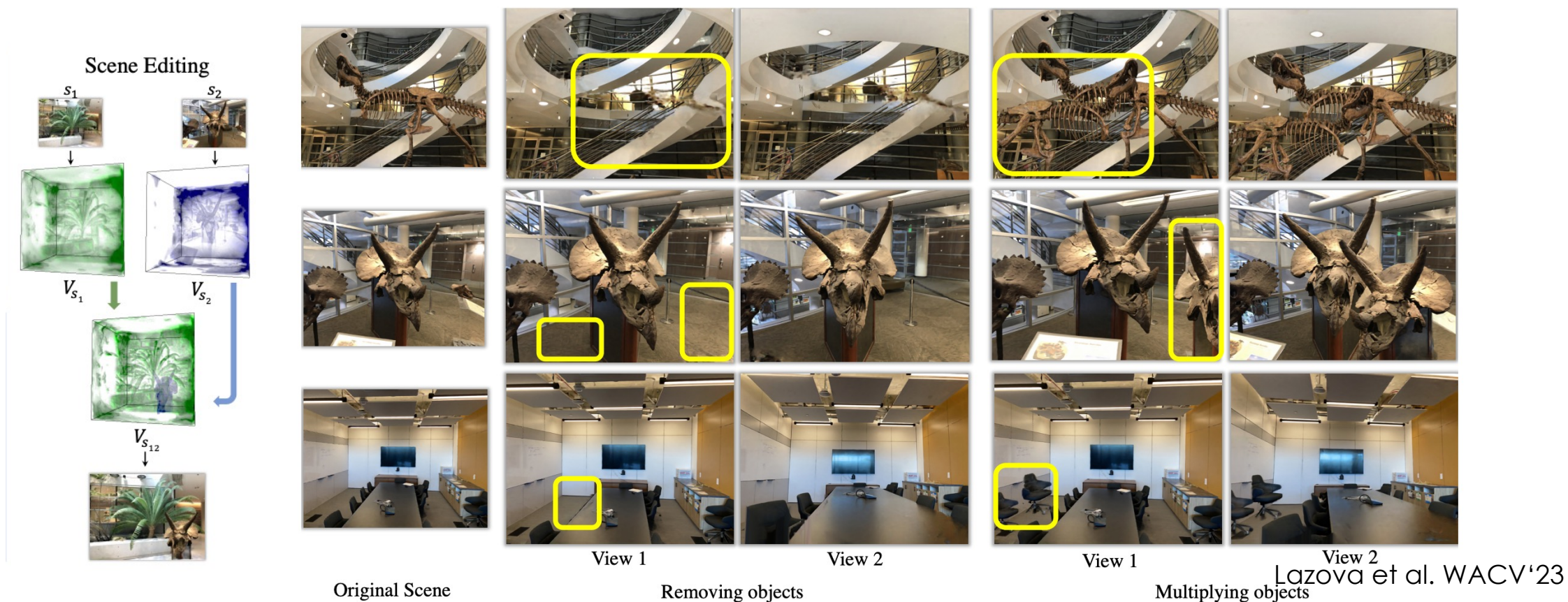
Control NeRF: Scene editing and manipulation

- Scene editing and compositing with Control-NeRF:



Control NeRF: Scene editing and manipulation

- Scene editing and compositing with Control-NeRF:



Goal: Scene Editing



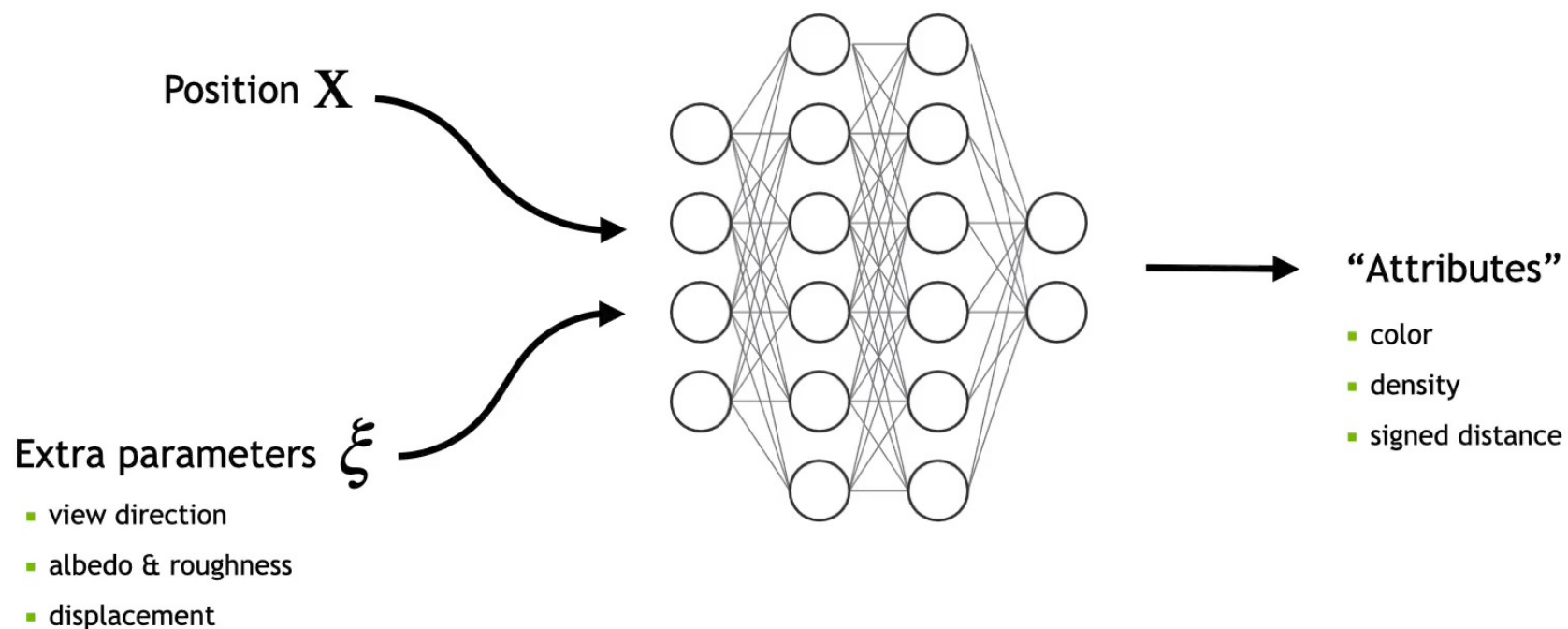
Limitations of NeRF:

4. Expensive training:

- Training is slow(10 hours-upto few days)
- Inference is also not real time

Neural Graphics Primitive

- An object (shape and appearance) represented by queries to a neural network. e.g. images, SDF, NeRF



Making NeRF(NGPs) faster

1. Smaller neural network

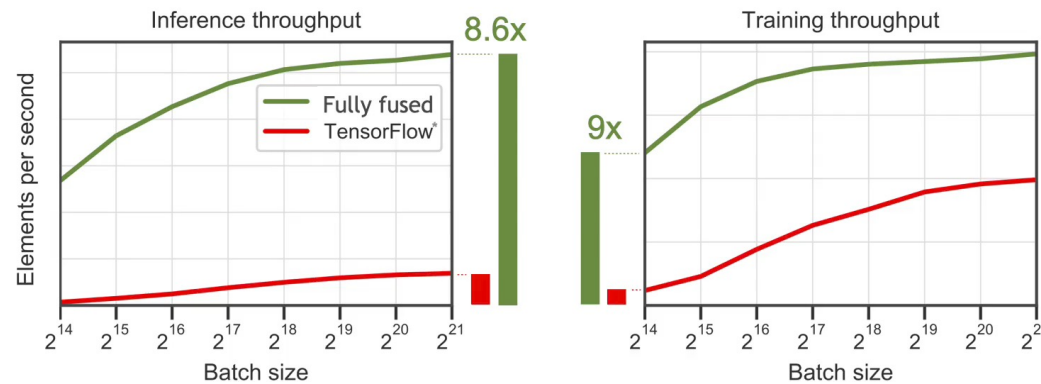
- Standard MLP with L layer, M neurons each, ReLU activations and no biases .For a constant batch size, the cost is:
 - Compute: $\mathcal{O}(M^2)$
 - Memory: $\mathcal{O}(M)$

Making NeRF(NGPs) faster

1. Smaller neural network

- Standard MLP with L layer, M neurons each, ReLU activations and no biases .For a constant batch size, the cost is:
 - Compute: $\mathcal{O}(M^2)$
 - Memory: $\mathcal{O}(M)$

Entire neural network implemented as single (CUDA) kernel



* TensorFlow Version 2.5.0 with XLA enabled



Positional Encoding in NeRF

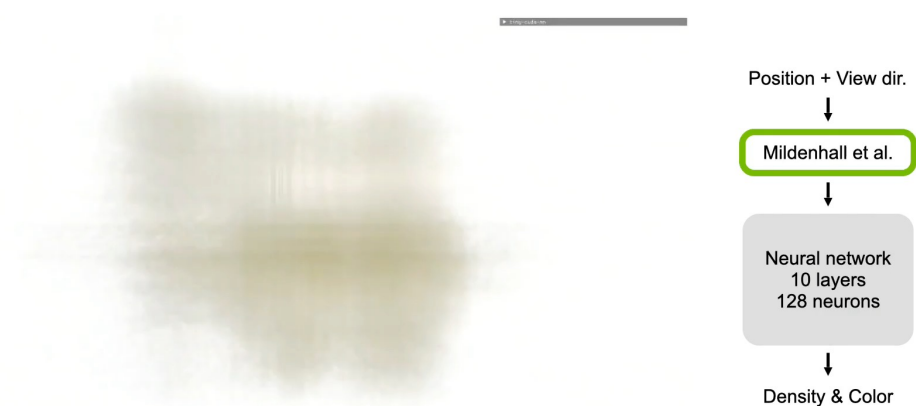
2. Input encoding

Without positional encoding



 NVIDIA

With positional encoding

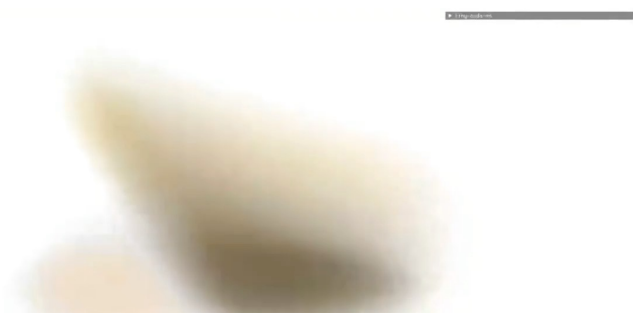


 NVIDIA

Positional Encoding in NeRF

2. Input encoding

Without positional encoding



Position + View dir.

Neural network
10 layers
128 neurons

With positional encoding



Position + View dir.

Mildenhall et al.

Neural network
10 layers
128 neurons

Density & Color

Is it possible to further improve the results with input encoding?

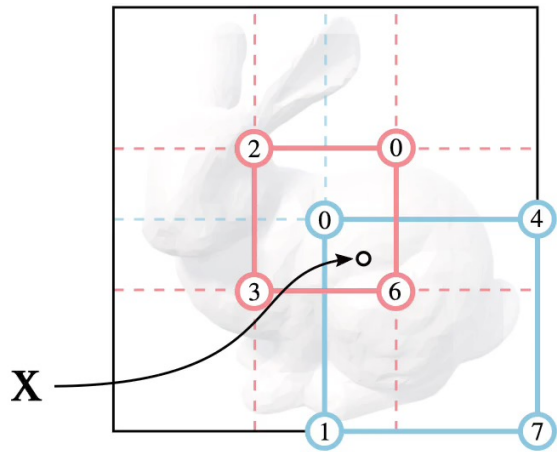


Making NeRF(NGPs) faster

2. Input encoding: Multiresolution Hash Encoding

Making NeRF(NGPs) faster

2. Input encoding: **Multiresolution** Hash Encoding

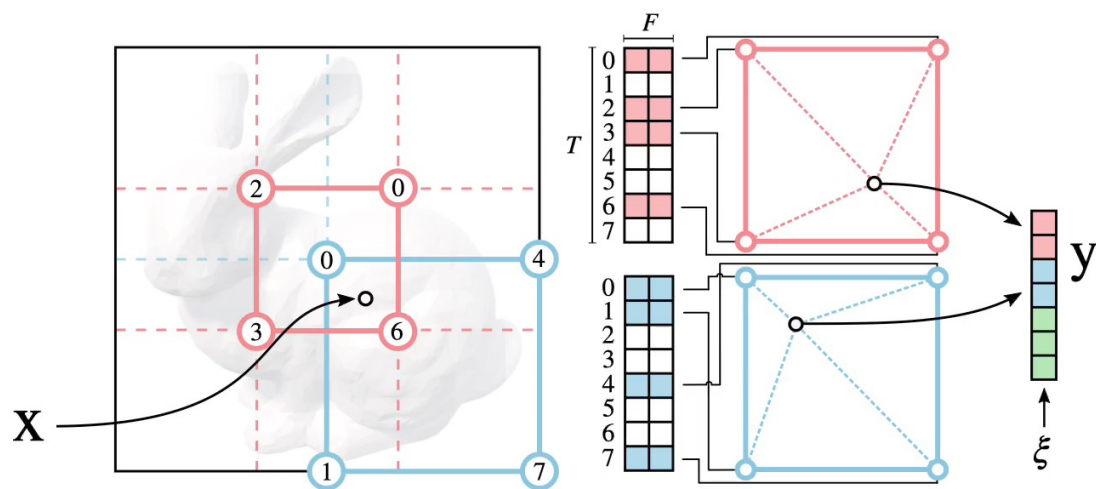


Multiresolution grids:

- Automatic level of details

Making NeRF(NGPs) faster

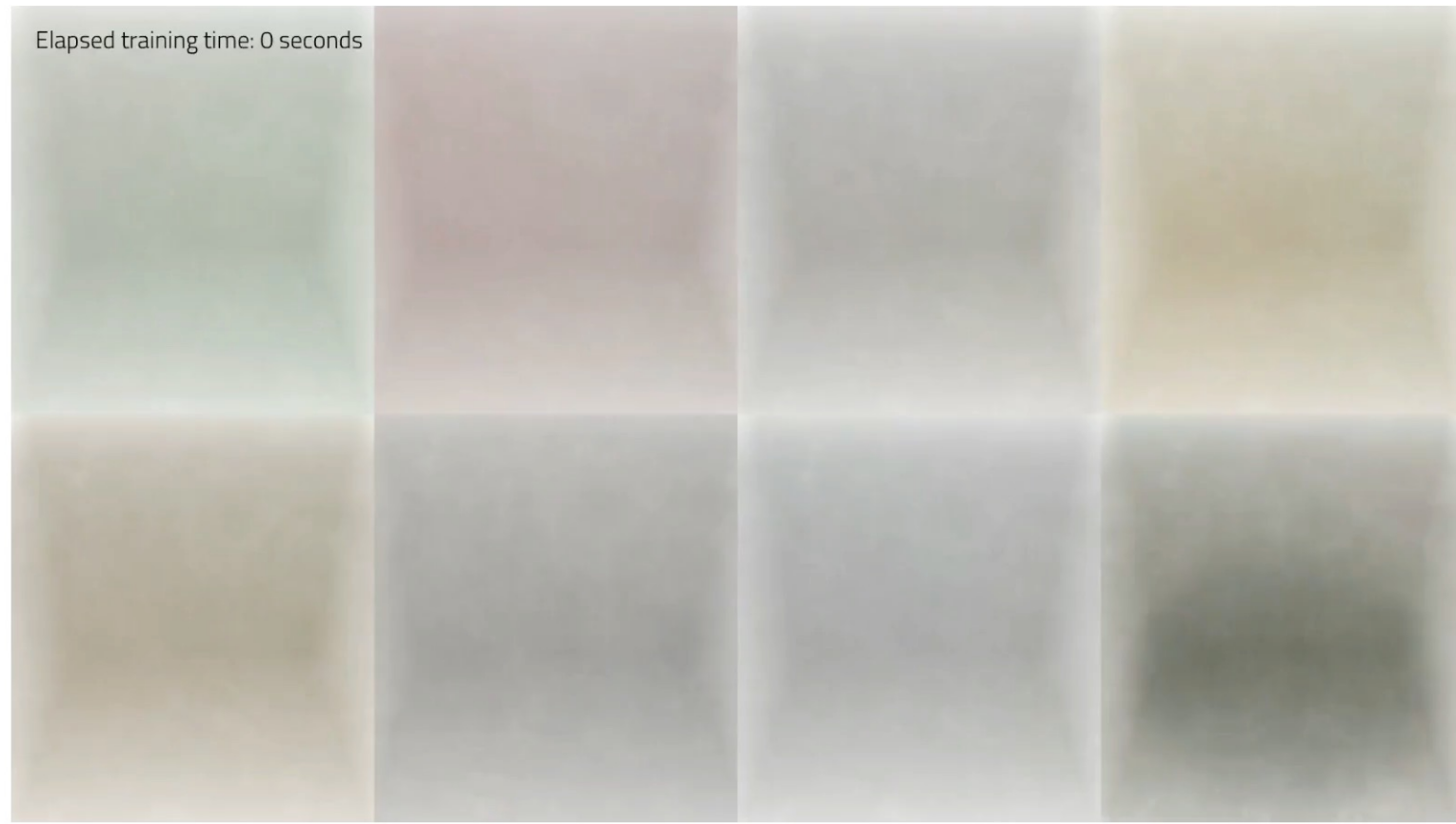
2. Input encoding: Multiresolution Hash Encoding



- Linear interpolation for continuous query.

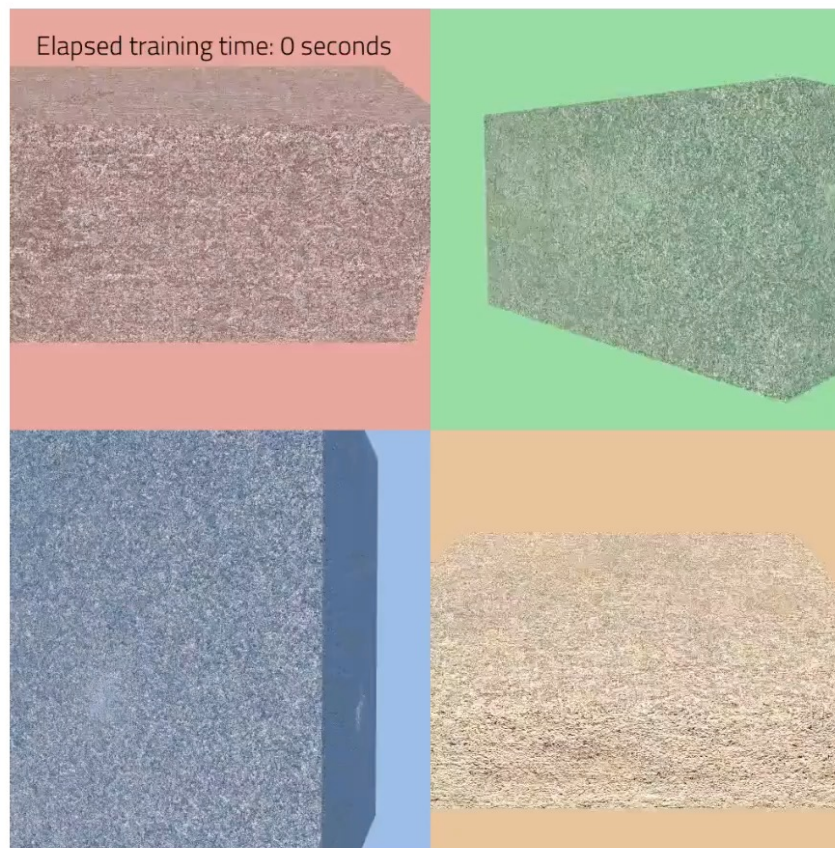
Instant NeRF training

INSTANT NERF TRAINING



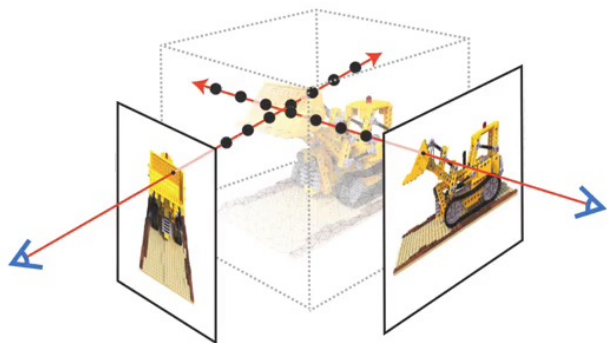
Instant SDF training

INSTANT SIGNED DISTANCE FUNCTION TRAINING



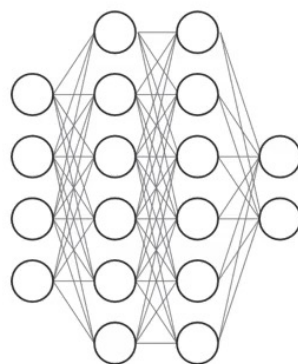
Conclusion: Instant-NGP

Rendering/training algorithm



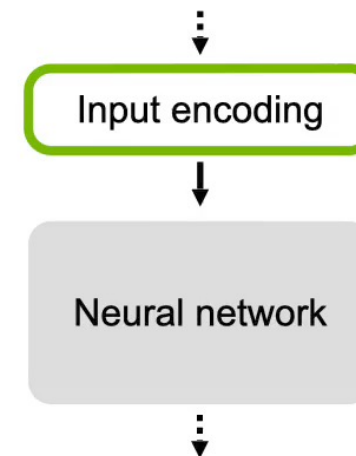
- Task-specific GPU implementation
- **10-100x faster than naïve tensor-based approach**

Small neural network



- Fully fused implementation
- **5-10x faster than TensorFlow**

“Good” input encoding



- Multiresolution hash encoding
- **Better speed-vs-quality tradeoff than prior work**
- **Task agnostic**

Limitations of NeRF

5. Surface extracted is not accurate and depends on threshold.



GT image

NeRF (density threshold $\sigma = 50$)



$\sigma = 1$

$\sigma = 10$

$\sigma = 50$

$\sigma = 100$

$\sigma = 500$

Surface Rendering v/s Volume Rendering

Surface Rendering (Implicit Surfaces)

- + High quality geometry
- + Clear surface definition
- Mask supervision required
- Texture mapping is blurry



Surface Rendering v/s Volume Rendering

Surface Rendering (Implicit Surfaces)

- + High quality geometry
- + Clear surface definition
- Mask supervision required
- Texture mapping is blurry

Volume Rendering (Radiance fields)

- Surface is approximated
- + Without mask supervision
- + High quality novel views and sharp textures/colors



Best of both worlds

Surface Rendering (Implicit Surfaces)

- + High quality geometry
- + Clear surface definition

~~Mask supervision required~~
~~Texture mapping is blurry~~



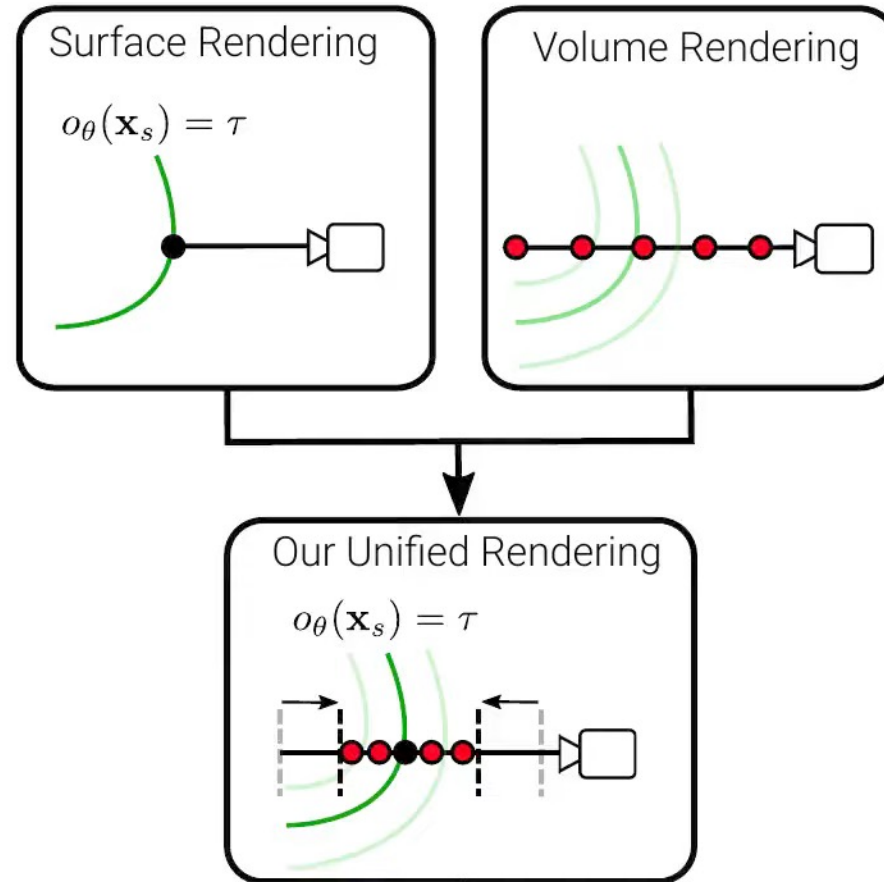
Volume Rendering (Radiance fields)

~~Surface is approximated~~

- + Without mask supervision
- + High quality novel views and sharp textures/colors



Unifying Implicit Surfaces and Radiance Fields



Unifying Implicit Surfaces and Radiance Fields

NeRF Volume rendering:

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i$$

NeRF Volume rendering with density

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j=1}^i (1 - \alpha_j) \mathbf{c}_i$$

$$\alpha_i = 1 - \exp(-\sigma_i \delta_i)$$

Unifying Implicit Surfaces and Radiance Fields

NeRF Volume rendering with density

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j=1}^i (1 - \alpha_j) \mathbf{c}_i$$

Unifying Implicit Surfaces and Radiance Fields

NeRF Volume rendering with density

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j=1}^i (1 - \alpha_j) \mathbf{c}_i$$

Key Idea: For solid objects, $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$, corresponds to an occupancy field O_i at i^{th} sample.

$$\hat{C}_v(\mathbf{r}) = \sum_{i=1}^N o_i \prod_{j=1}^i (1 - o_j) \mathbf{c}_i$$

Unifying Implicit Surfaces and Radiance Fields

NeRF Volume rendering with density

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N \alpha_i \prod_{j=1}^i (1 - \alpha_j) \mathbf{c}_i$$

Key Idea: For solid objects, $\alpha_i = 1 - \exp(-\sigma_i \delta_i)$, corresponds to an occupancy field O_i at i^{th} sample.

$$\hat{C}_v(\mathbf{r}) = \sum_{i=1}^N o_i \prod_{j=1}^i (1 - o_j) \mathbf{c}_i$$

Given occupancy of surface, we can now render the same scene with surface rendering.

Unifying Implicit Surfaces and Radiance Fields

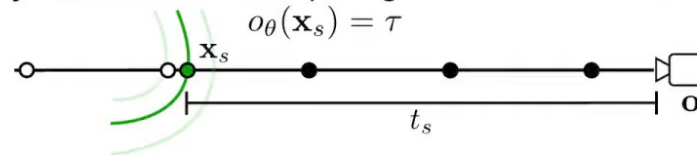
Key Idea:

- Volume rendering in early stage:
 - Optimization without mask
- Surface rendering in later stage:
 - Level-set surfaces

Unifying Implicit Surfaces and Radiance Fields

Rendering Procedure

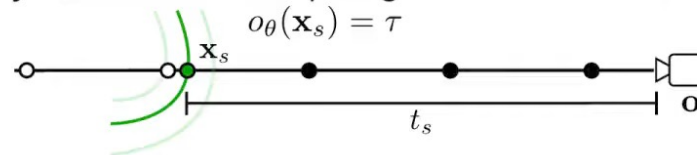
- Find surface along a ray: uniform sampling + iterative secant method



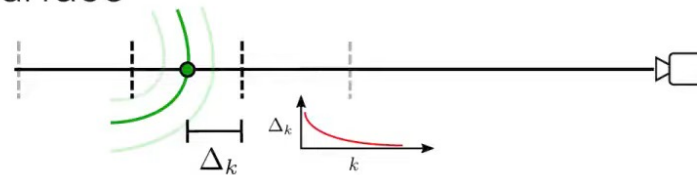
Unifying Implicit Surfaces and Radiance Fields

Rendering Procedure

- ▶ Find surface along a ray: uniform sampling + iterative secant method



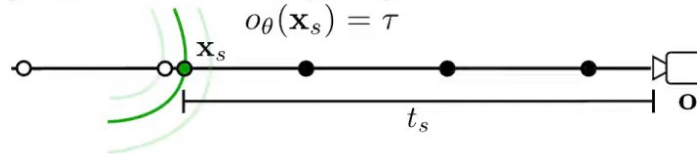
- ▶ Define interval at the surface



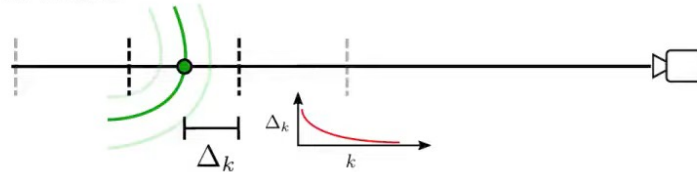
Unifying Implicit Surfaces and Radiance Fields

Rendering Procedure

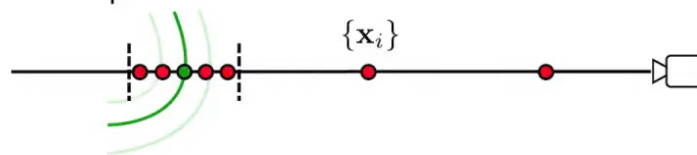
- ▶ Find surface along a ray: uniform sampling + iterative secant method



- ▶ Define interval at the surface



- ▶ Volume rendering with occupancies

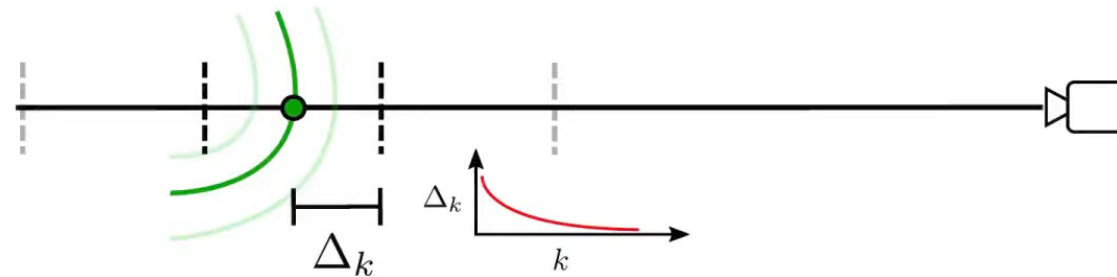


Unifying Implicit Surfaces and Radiance Fields

Rendering Procedure

Transition from Volume rendering to Surface rendering

- ▶ Exponential decay of interval Δ_k wrt. iterations k



Volume Rendering
Rendering

Δk

decreases

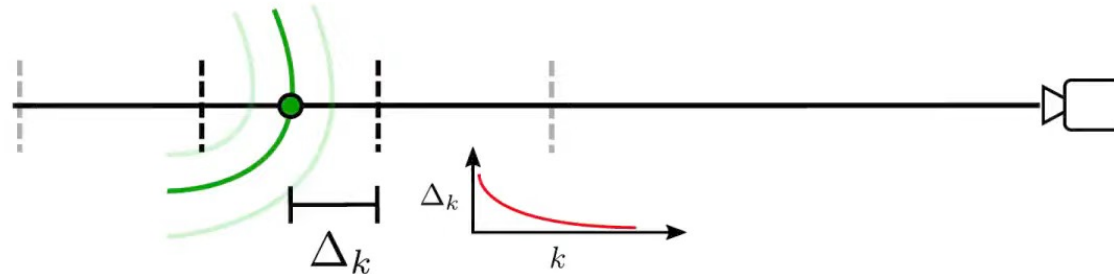
Surface

Unifying Implicit Surfaces and Radiance Fields

Rendering Procedure

Transition from Volume rendering to Surface rendering

- ▶ Exponential decay of interval Δ_k wrt. iterations k



Theorem

Volume and surface rendering become equivalent when reducing the interval and increasing the number of samples.

$$\lim_{\substack{\Delta \rightarrow 0 \\ N \rightarrow \infty}} \hat{C}_v(\mathbf{r}) = \hat{C}_s(\mathbf{r})$$

Unifying Implicit Surfaces and Radiance Fields

Comparison on the DTU MVS dataset



Ours

More on NeRF

Dynamic NeRFs:

- TöRF Attal et al., NeurIPS 2021
- NSFF, Li et al., CVPR 2021
-

NeRF from few images:

- PixelNeRF, Yu et al., CVPR 2021
-

NeRF + implicit surfaces:

- NeUS, Wang et al., NeurIPS 2021
- VolSDF, Yariv et al., NeurIPS 2021
-

- Many more

Checkout NeuralFields: <https://neuralfields.cs.brown.edu/>