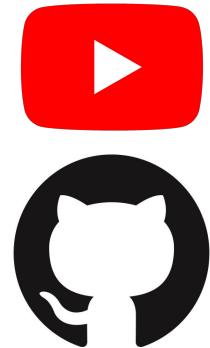




UofT Computer Vision Club

Paper Breakdown Series

NeRFs: Neural Radiance Fields



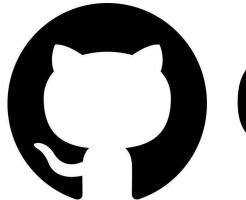
youtube.com/@uoftcv

github.com/ut-cvdp/paper-breakdown

Follow us on social media

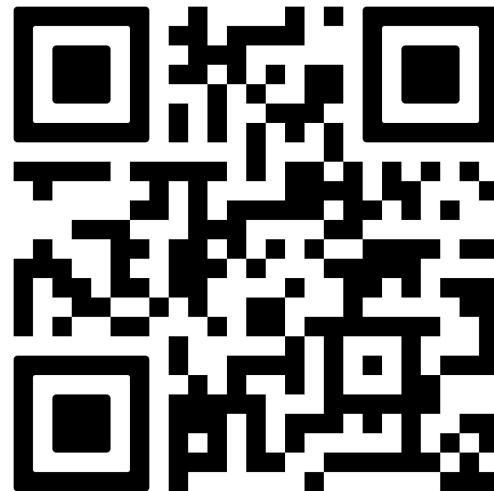


UofT CV is sponsored by



GitHub

Sign up for the GitHub Student Developer Pack:
<https://education.github.com/benefits>



About the Club

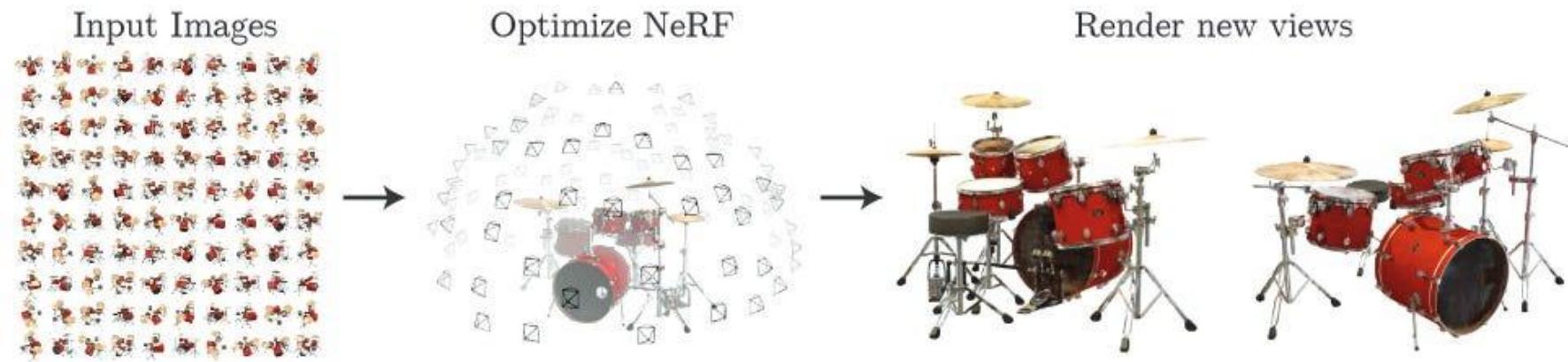
- We are a student club for Computer Vision and Machine Learning affiliated with the University of Toronto.
- We run:
 - paper “discussions” for paper in the past 1 week
 - conference-style research orals for pre-prints
 - Now, paper breakdown series

If you want to help run our club, talk to any of us or drop an email
contact@cvdg.tech.

Vedant

1. Explain problem - show examples as well
2. intuitive recap of previous attempts to solve this
3. what is a neural net - borrow from trace 5d slides
4. what is the 5D input
5. given some images, how do they convert them to these 5d representations
6. volume rendering - eq 1, 2, 3

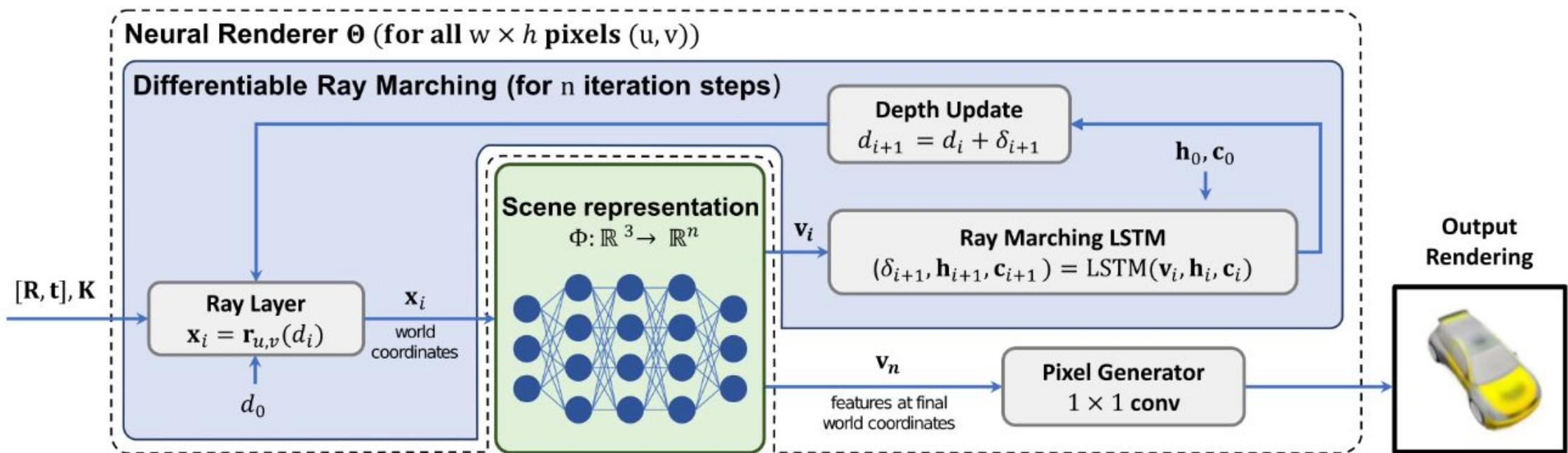
Problem Statement



UNIVERSITY OF
TORONTO

Previous Solutions - SRN

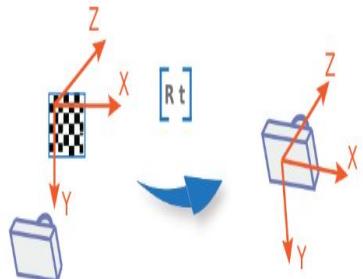
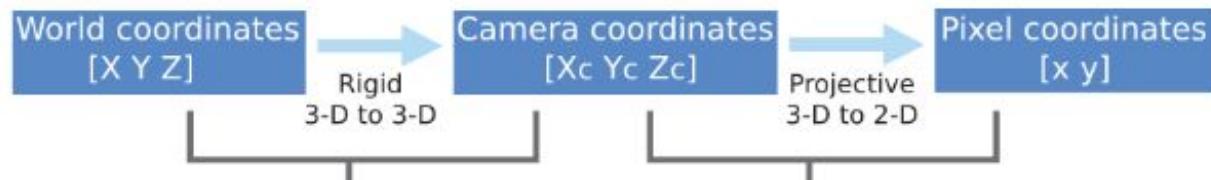
- Take the Implicit and Explicit Camera Matrices as Input.
- The output is a 2D RGB Image which denotes the Scene from a specified viewpoint.
- NERF vs SRNs



Previous Solutions - SRN

Camera Calibration Parameters

The calibration algorithm calculates the camera matrix using the extrinsic and coordinate system to the 3-D camera's coordinate system. The intrinsic param coordinates.



Intrinsic Parameters

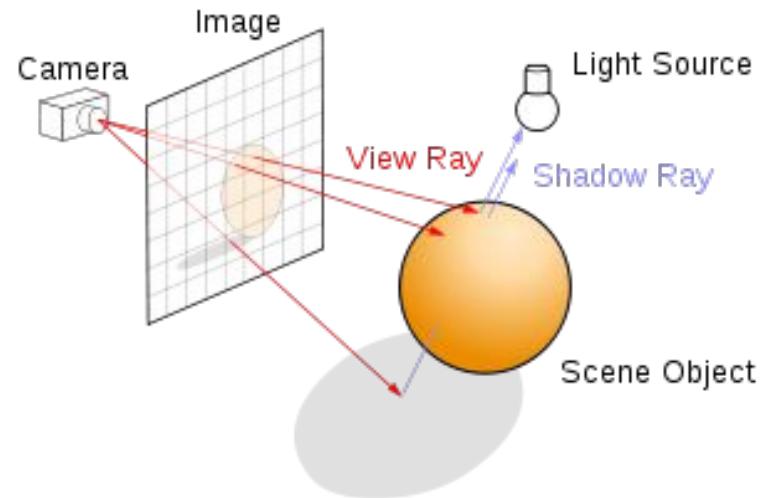
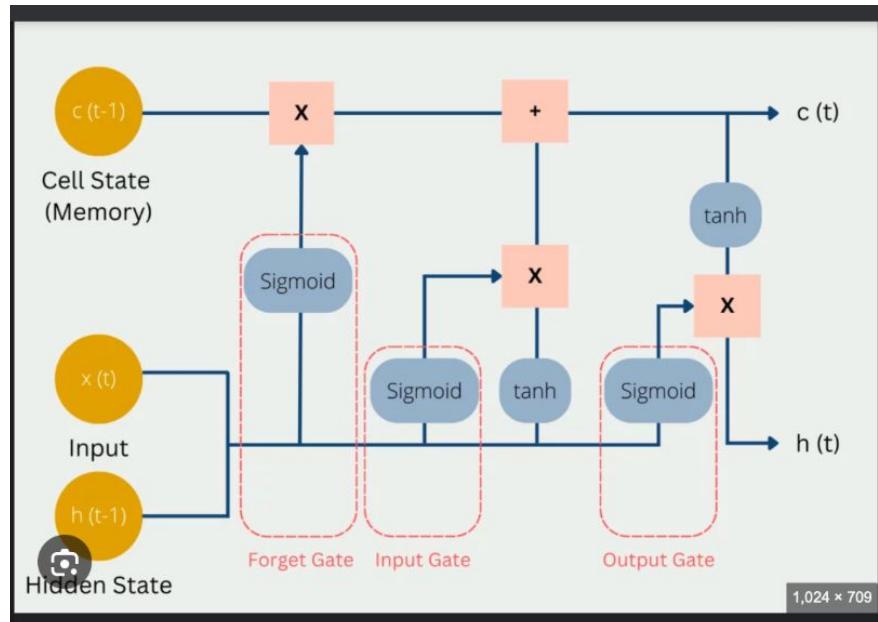
The intrinsic parameters include the focal length, the optical center, also known as the *principal point*, and the skew coefficient.

$$\begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$$



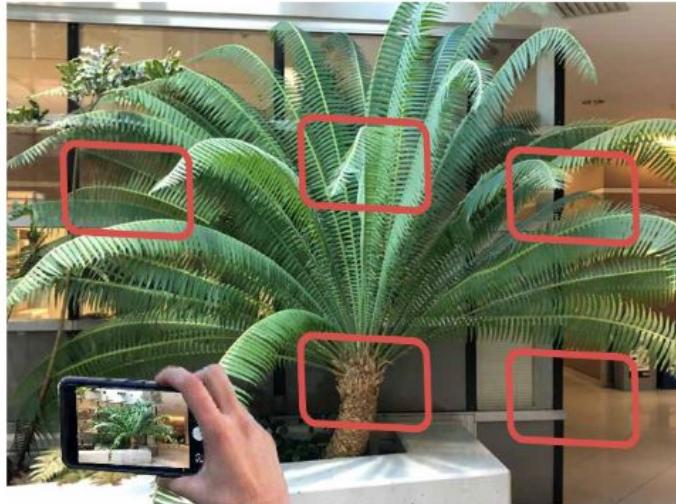
Previous Solutions - SRN

- SRNs employ use of Ray Marching Algorithms



They use the LSTM model architecture in order for Ray Marching ALgorithms.

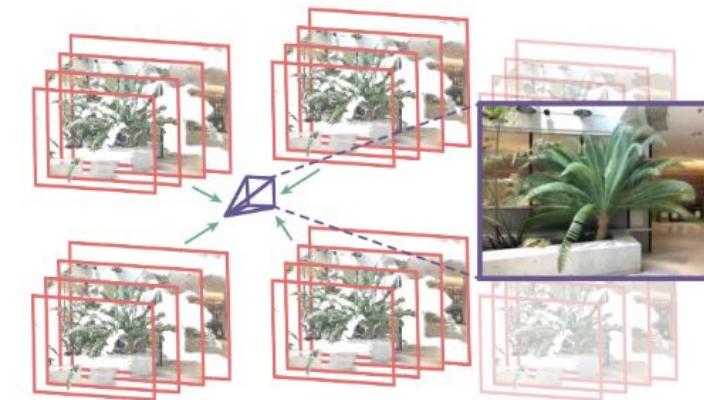
Previous Solutions - LLFF



Fast and easy handheld capture with guideline:
closest object moves at most D pixels between views

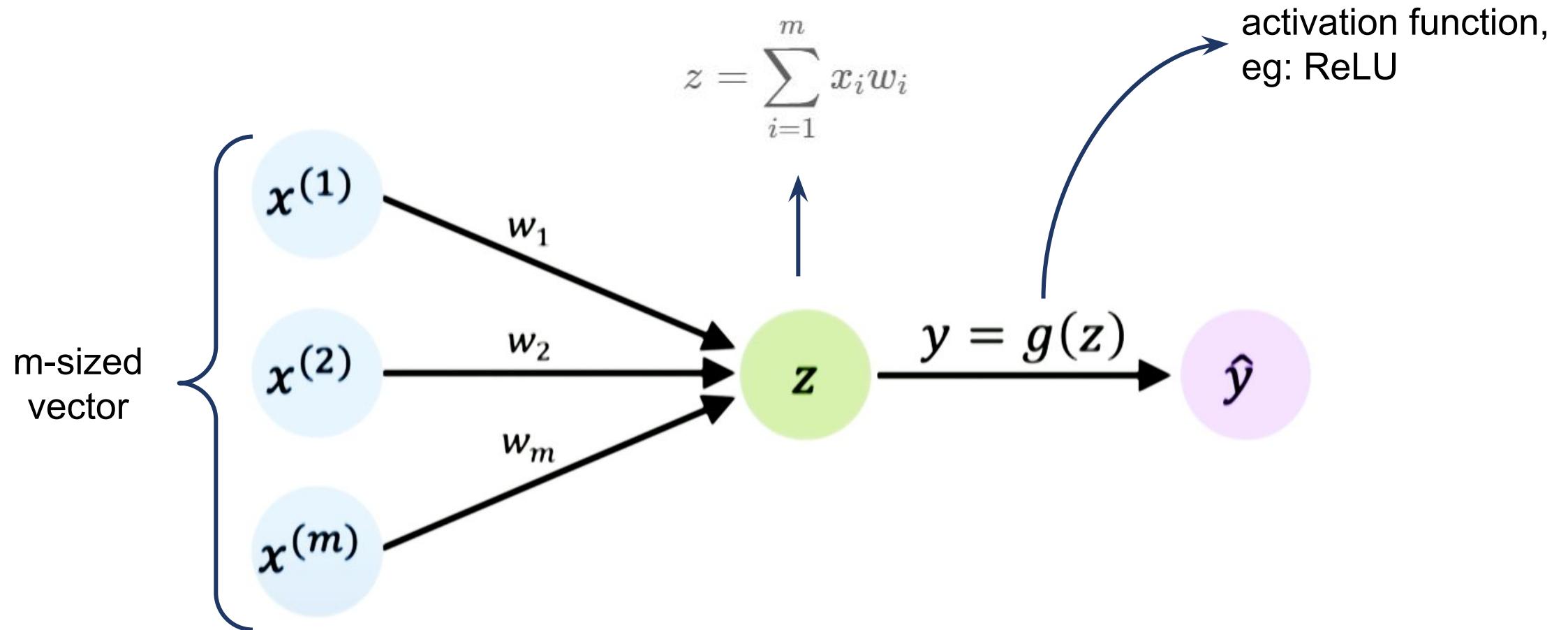


Promote sampled views to local light field
via layered scene representation

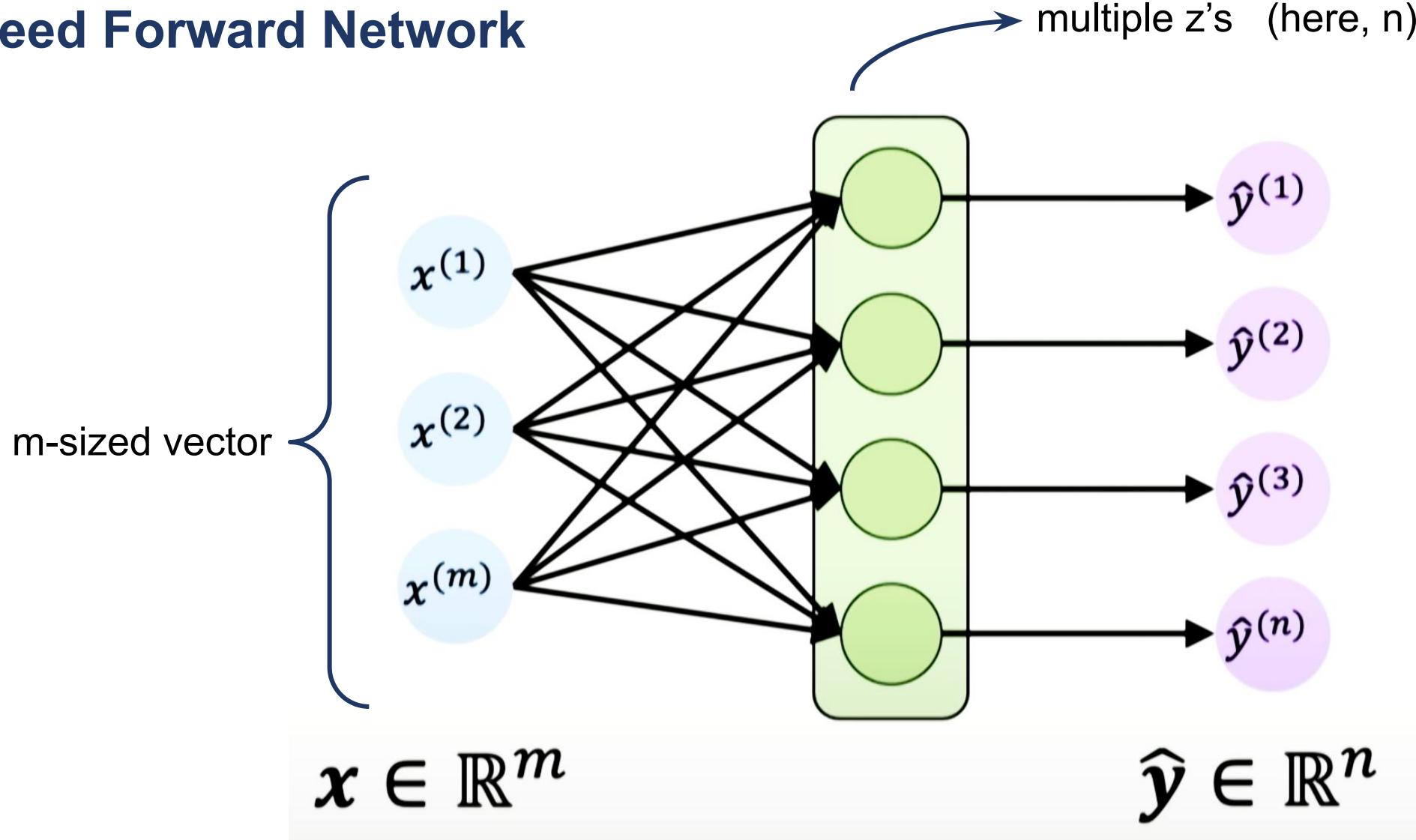


Blend neighboring local light fields
to render novel views

Perceptron

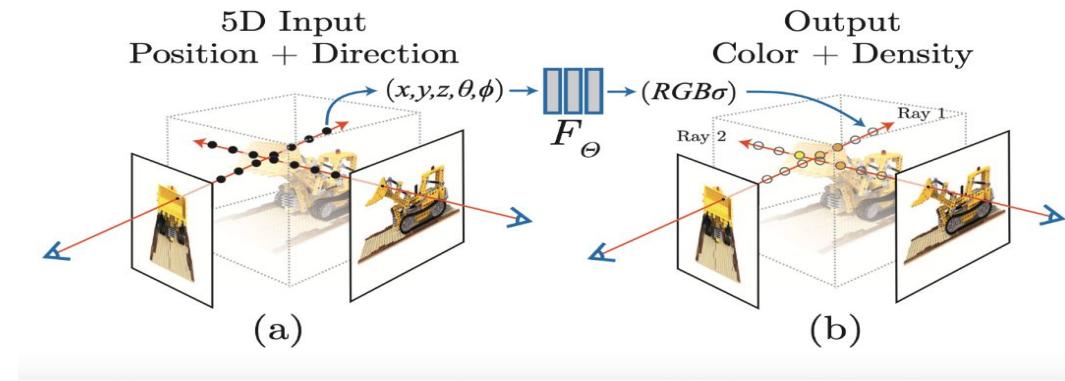


Feed Forward Network

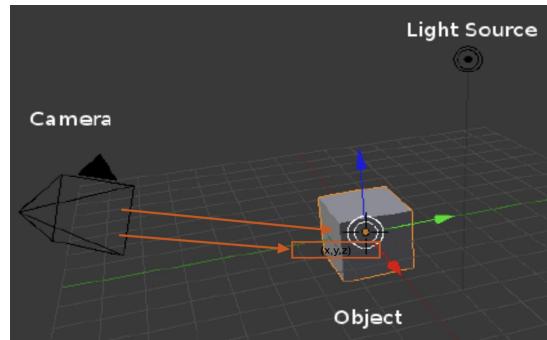


NERFs

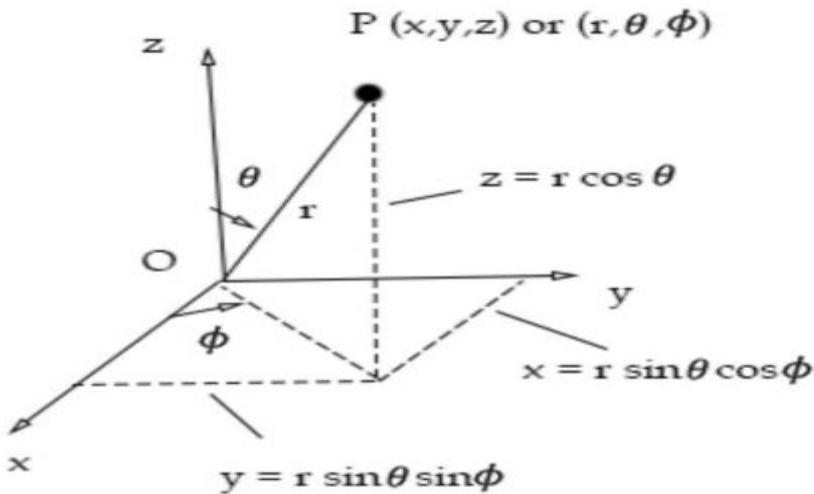
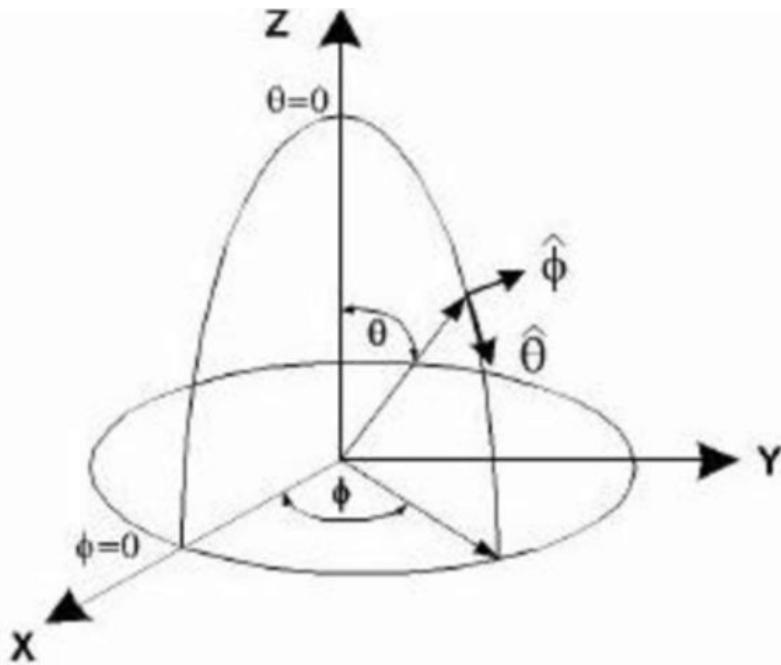
>It is a function which takes 5 coordinates and returns the density and RGB value at the point



>The 3d coordinates of the point are one input.



>The other input is the view direction of the camera angles

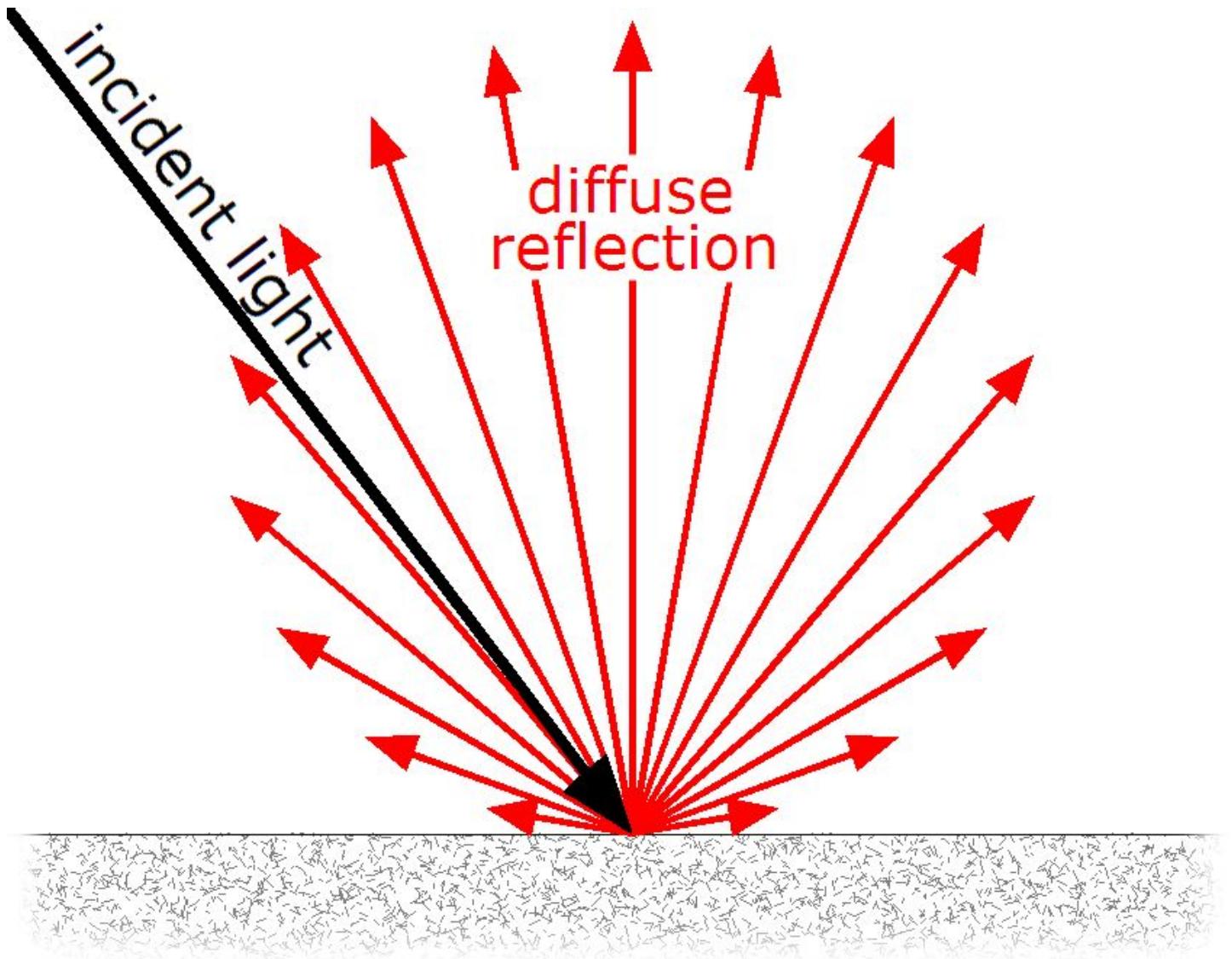


$$\text{cart} \leftrightarrow \text{sph} \quad \left\{ \begin{array}{l} x = \rho \cos \theta \sin \phi, \\ y = \rho \sin \theta \sin \phi, \\ z = \rho \cos \phi, \end{array} \right.$$

$$\left\{ \begin{array}{l} \rho = \sqrt{x^2 + y^2 + z^2}, \\ \theta = \arctan \frac{y}{x}, \\ \phi = \arctan \frac{\sqrt{x^2 + y^2}}{z} \\ \qquad = \arccos \frac{z}{\sqrt{x^2 + y^2 + z^2}}. \end{array} \right.$$



- Path-traced objects with realistic non-Lambertian material



- Inserting virtual objects into real world scenes



Volume Rendering

Utilizing the Output of our model

Once we have the volume density and the directional emitted radiance at a point how are we going to render the image?

$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right).$$

$\sigma(\mathbf{x})$: Colour Density at the position x

$C(\mathbf{r})$: Expected Colour of the ray r

$T(t)$: Accumulated Transmittance along the ray from t_n to t

Calculation of Transmittance

$$\begin{aligned}\mathcal{T}(t + dt) &= \mathcal{T}(t) \cdot (1 - dt \cdot \sigma(t)) \\ \frac{\mathcal{T}(t + dt) - \mathcal{T}(t)}{dt} &\equiv \mathcal{T}'(t) = -\mathcal{T}(t) \cdot \sigma(t)\end{aligned}$$

This is a classical differential equation that can be solved as follows:

$$\mathcal{T}'(t) = -\mathcal{T}(t) \cdot \sigma(t) \tag{3}$$

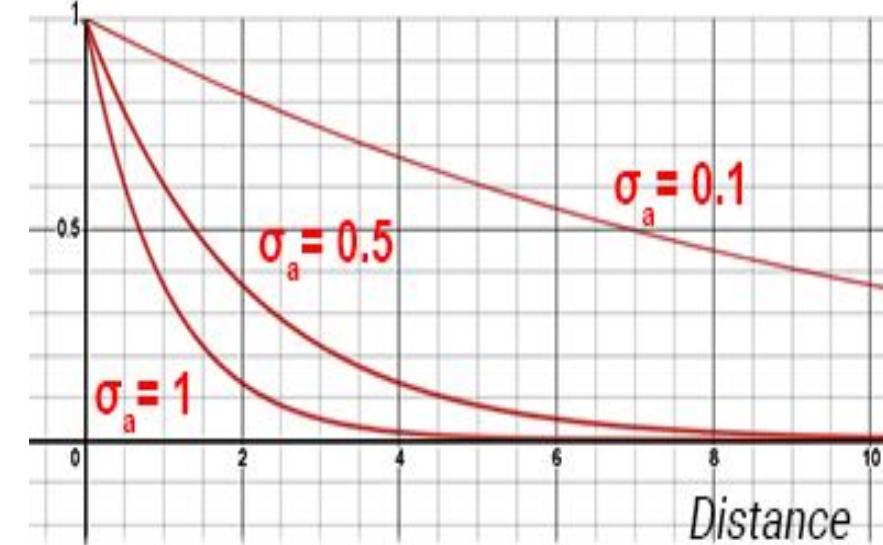
$$\frac{\mathcal{T}'(t)}{\mathcal{T}(t)} = -\sigma(t) \tag{4}$$

$$\int_a^b \frac{\mathcal{T}'(t)}{\mathcal{T}(t)} dt = - \int_a^b \sigma(t) dt \tag{5}$$

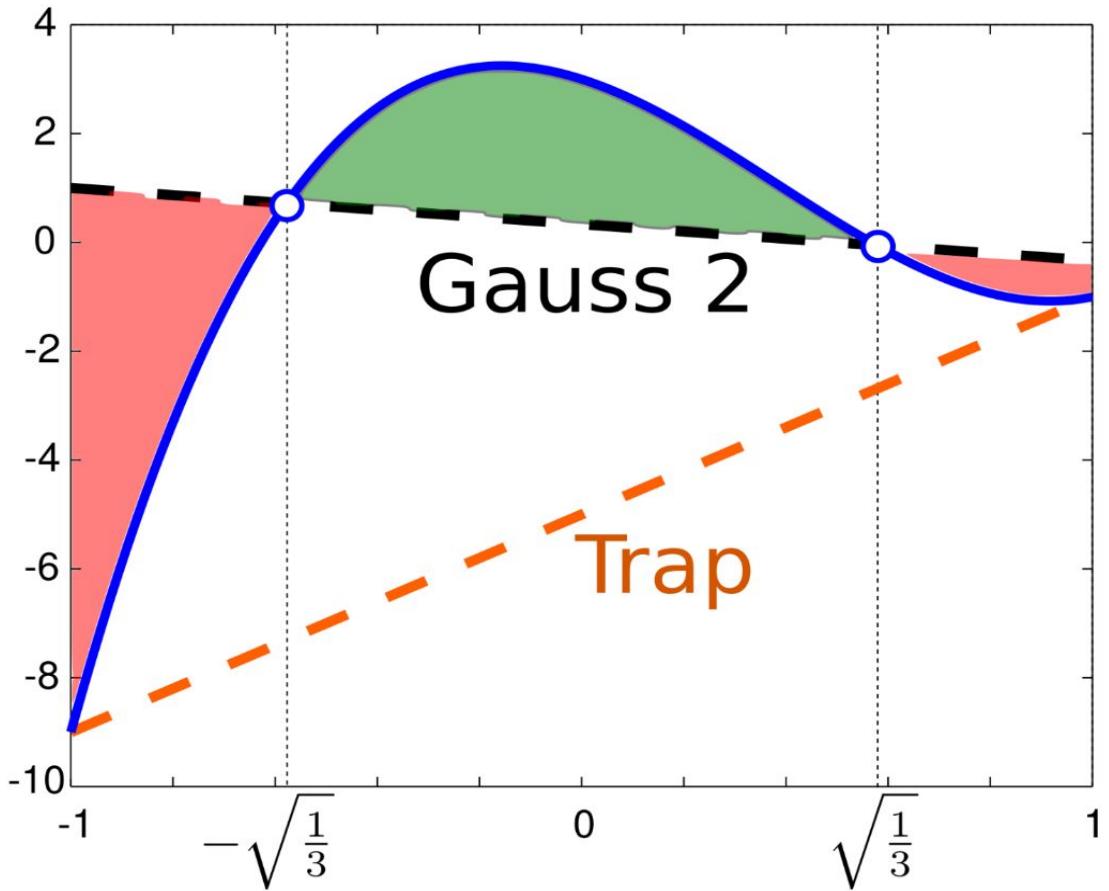
$$\log \mathcal{T}(t)|_a^b = - \int_a^b \sigma(t) dt \tag{6}$$

$$\mathcal{T}(a \rightarrow b) \equiv \frac{\mathcal{T}(b)}{\mathcal{T}(a)} = \exp \left(- \int_a^b \sigma(t) dt \right) \tag{7}$$

Transmittance



Deterministic Quadrature



$$\int_{-1}^1 f(x) dx \approx \sum_{i=1}^n w_i f(x_i),$$

$$w_i = \frac{2}{(1-x_i^2) [P'_n(x_i)]^2}.$$

Stratified Sampling

We partition $[t_n, t_f]$ into N even sized bins and sample from each of the bins.

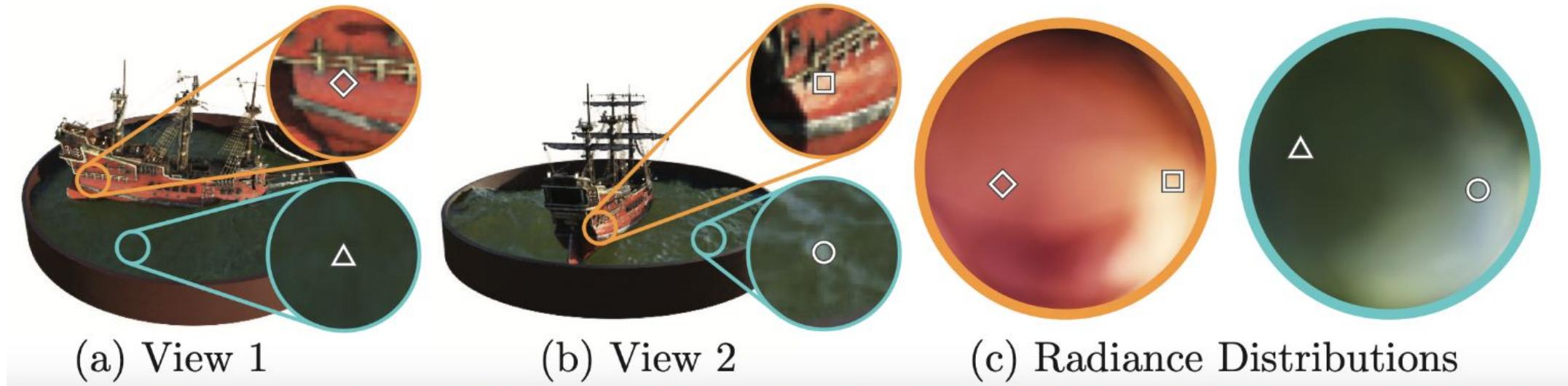
$$t_i \sim \mathcal{U} \left[t_n + \frac{i-1}{N} (t_f - t_n), \; t_n + \frac{i}{N} (t_f - t_n) \right]$$

We transform the integral into a Riemann Sum using the above sampling

$$\hat{C}(\mathbf{r}) = \sum_{i=1}^N T_i (1 - \exp(-\sigma_i \delta_i)) \mathbf{c}_i, \text{ where } T_i = \exp \left(- \sum_{j=1}^{i-1} \sigma_j \delta_j \right)$$

An Example applying Stratified Sampling

View Dependent emitted radial field



Optimization

This is not enough for state-of-the-art yet

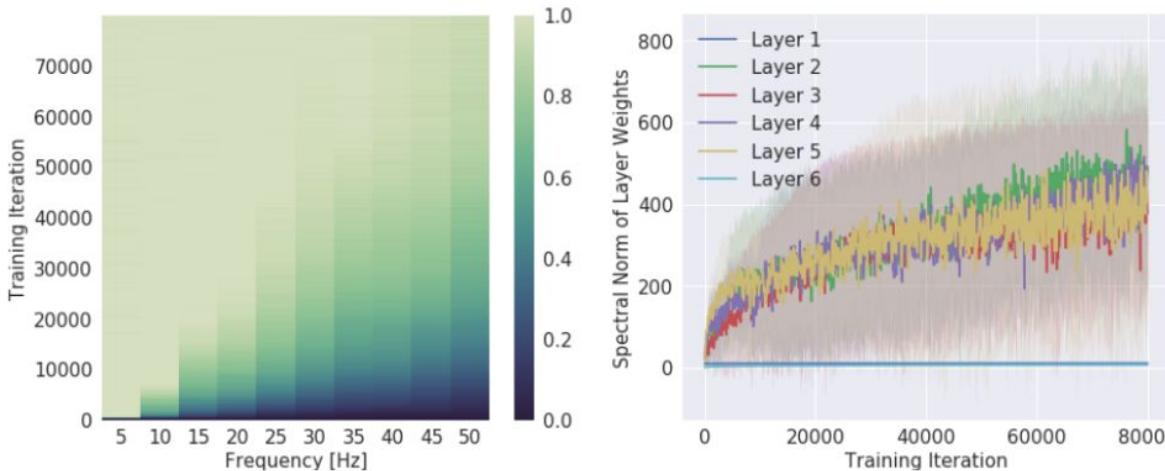
Positional Encoding

- NNs perform poorly on high frequency variations in colour and geometry
- Mapping the inputs to a higher dimensional space using high frequency functions before passing them helps

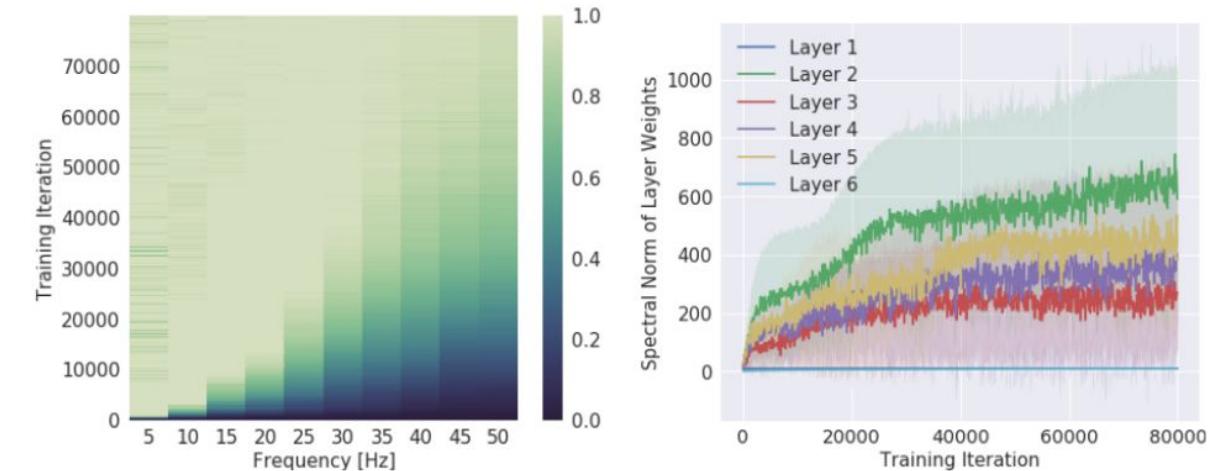
Positional Encoding

On the Spectral Bias of Neural Networks

Nasim Rahaman ^{*1 2} Aristide Baratin ^{*1} Devansh Arpit ¹ Felix Draxler ² Min Lin ¹ Fred A. Hamprecht ²
Yoshua Bengio ¹ Aaron Courville ¹



(a) Equal Amplitudes



(b) Increasing Amplitudes

Positional Encoding

$$F_{\Theta} = F'_{\Theta} \circ \gamma$$

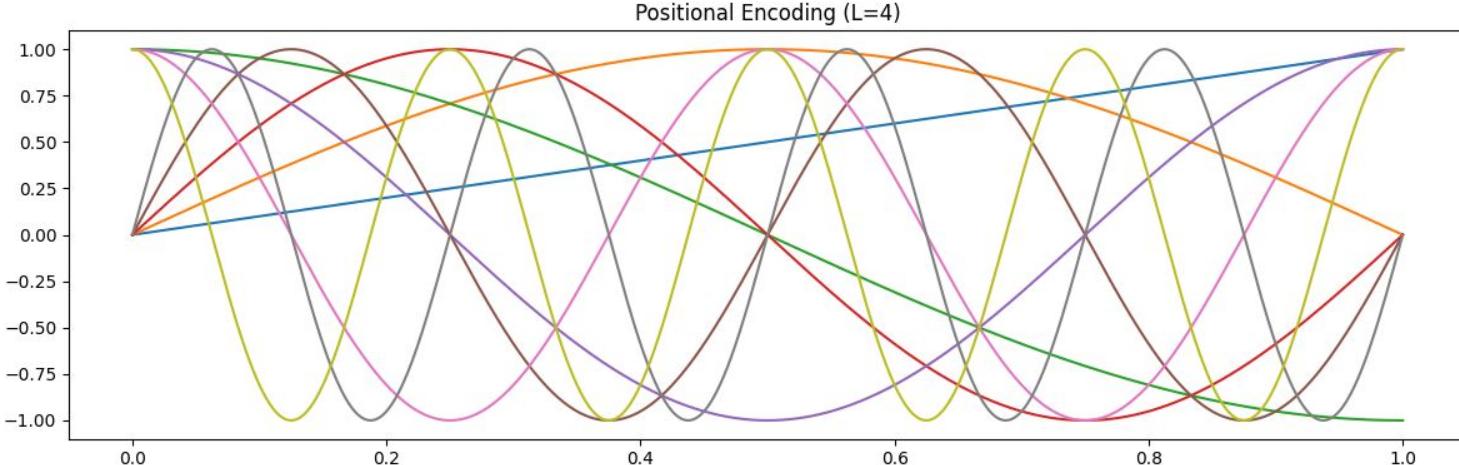
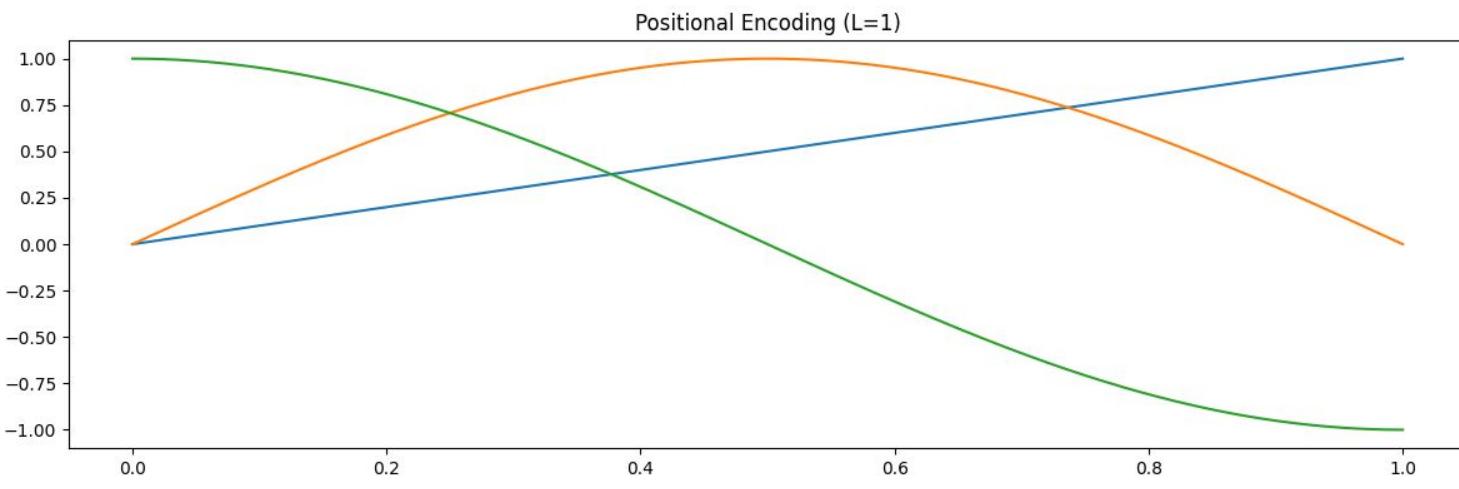
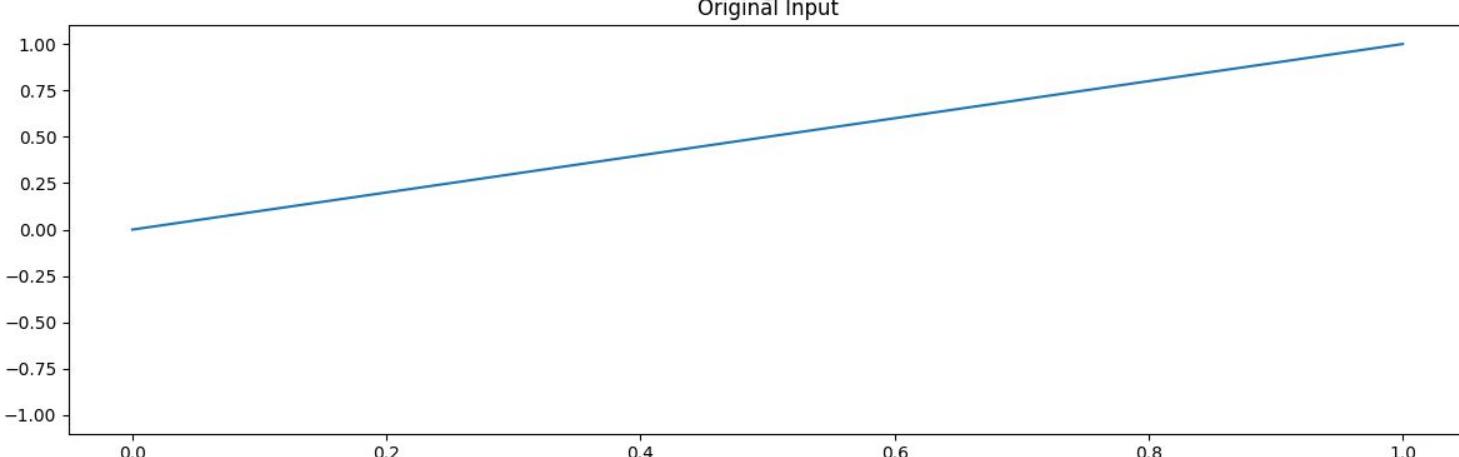
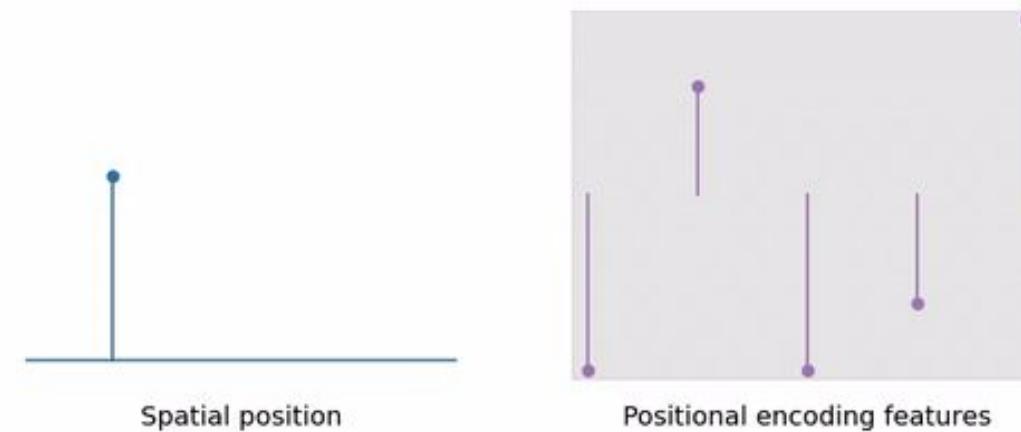
Here γ is a mapping from \mathbb{R} into a higher dimensional space \mathbb{R}^{2L} , and F'_{Θ} is still simply a regular MLP. Formally, the encoding function we use is:

$$\gamma(p) = (\sin(2^0 \pi p), \cos(2^0 \pi p), \dots, \sin(2^{L-1} \pi p), \cos(2^{L-1} \pi p)). \quad (4)$$

Positional Encoding

$$F_{\Theta} = F'_{\Theta} \circ \gamma$$

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



UNIVERSITY OF
TORONTO



Positional Encoding

$$F_{\Theta} = F'_{\Theta} \circ \gamma$$

Here γ is a mapping from \mathbb{R} into a higher dimensional space \mathbb{R}^{2L} , and F'_{Θ} is still simply a regular MLP. Formally, the encoding function we use is:

$$\gamma(p) = (\sin(2^0\pi p), \cos(2^0\pi p), \dots, \sin(2^{L-1}\pi p), \cos(2^{L-1}\pi p)). \quad (4)$$

This function $\gamma(\cdot)$ is applied separately to each of the three coordinate values in \mathbf{x} (which are normalized to lie in $[-1, 1]$) and to the three components of the Cartesian viewing direction unit vector \mathbf{d} (which by construction lie in $[-1, 1]$)

Positional Encoding

$$PE(pos, 2i) = \sin \frac{pos}{10000 \frac{2i}{d_{model}}}$$

$$PE(pos, 2i + 1) = \cos \frac{pos}{10000 \frac{2i}{d_{model}}}$$

Sentence 1

YOUR

PE(0, 0)
PE(0, 1)
PE(0, 2)
...
PE(0, 510)
PE(0, 511)

CAT

PE(1, 0)
PE(1, 1)
PE(1, 2)
...
PE(1, 510)
PE(1, 511)

IS

PE(2, 0)
PE(2, 1)
PE(2, 2)
...
PE(2, 510)
PE(2, 511)

Sentence 2

I

PE(0, 0)
PE(0, 1)
PE(0, 2)
...
PE(0, 510)
PE(0, 511)

LOVE

PE(1, 0)
PE(1, 1)
PE(1, 2)
...
PE(1, 510)
PE(1, 511)

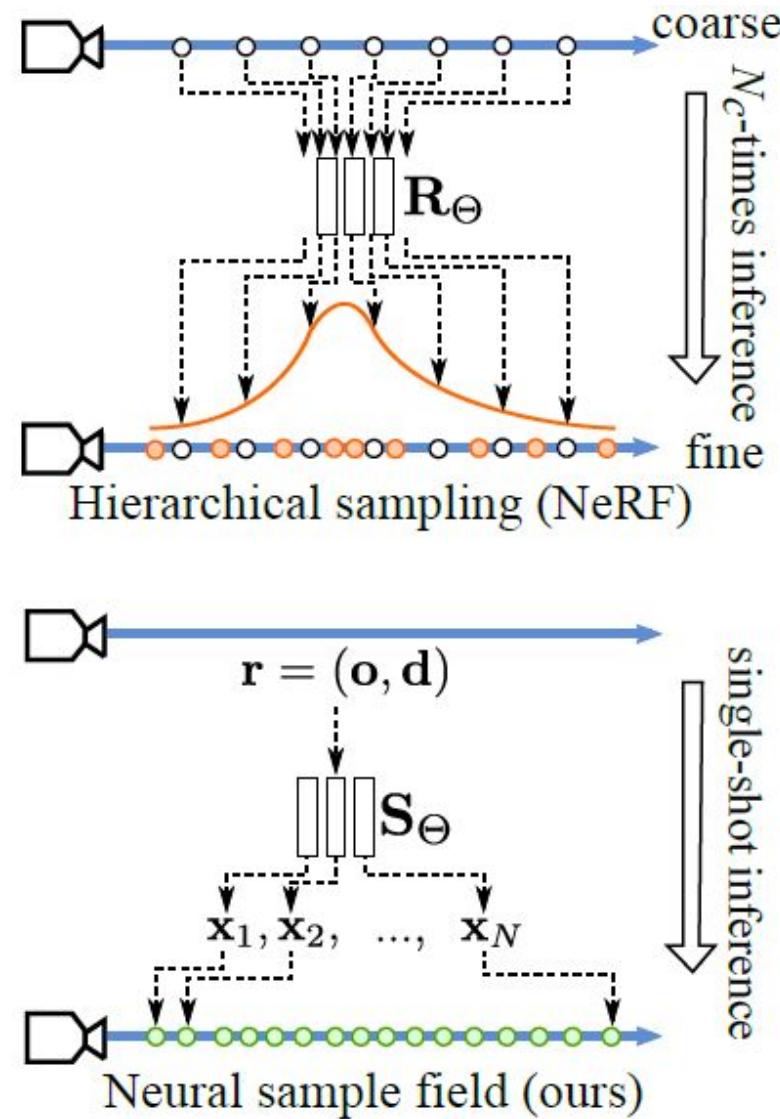
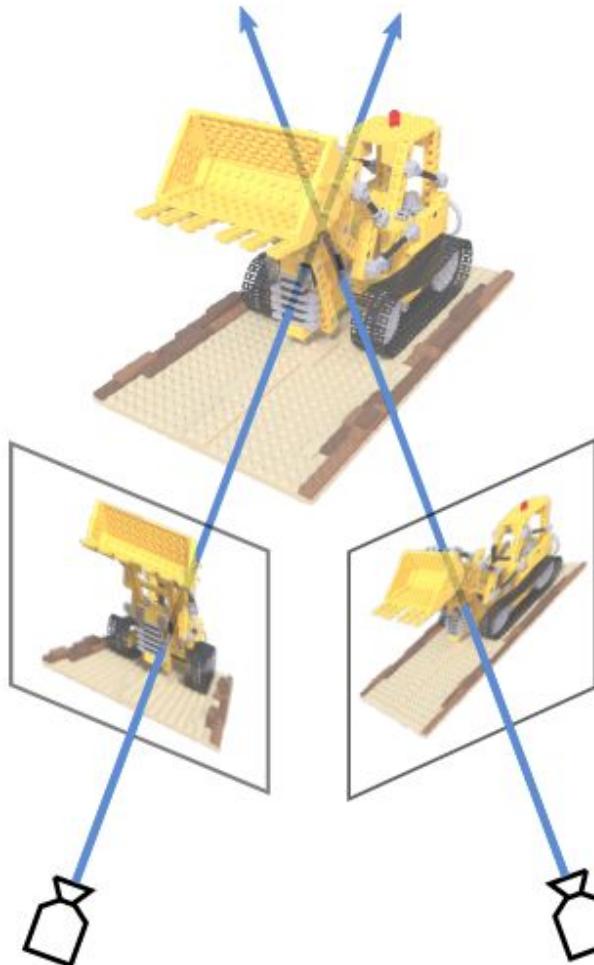
YOU

PE(2, 0)
PE(2, 1)
PE(2, 2)
...
PE(2, 510)
PE(2, 511)

We only need to compute the positional encodings once and then reuse them for every sentence, no matter if it is training or inference.

Hierarchical Volume Sampling

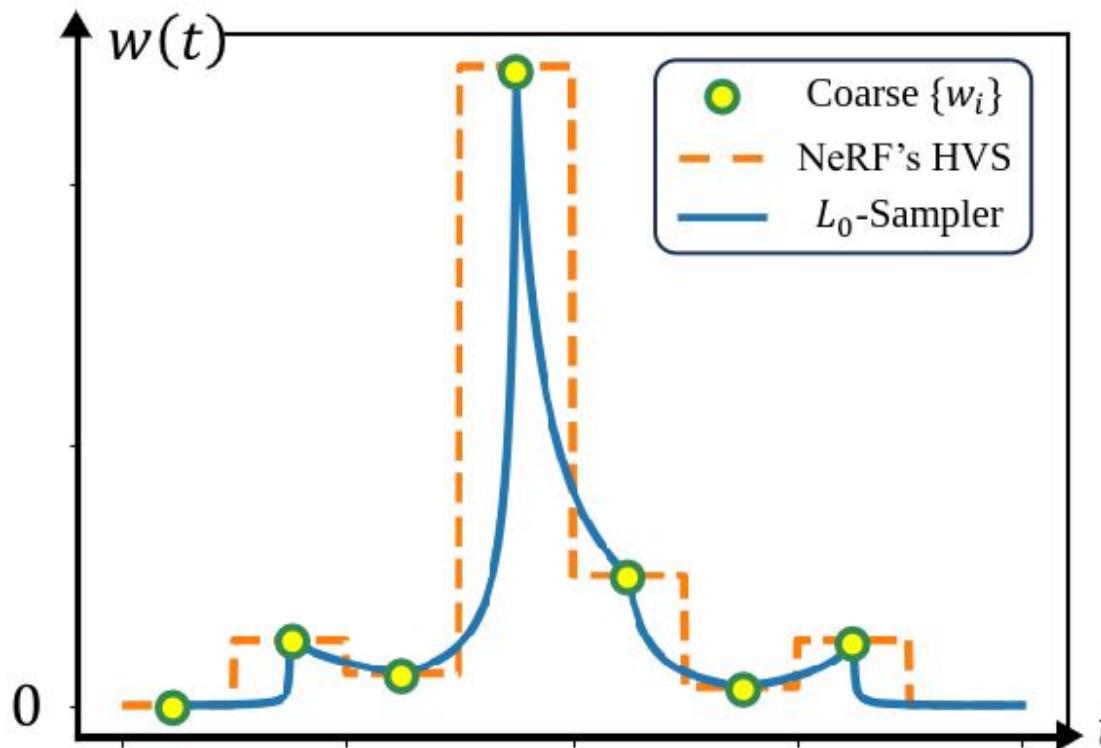
- Using N query points along each camera ray is inefficient: free space and occluded regions that do not contribute to the rendered image are still sampled repeatedly
- Now allocate samples proportionally to their expected effect on the final rendering



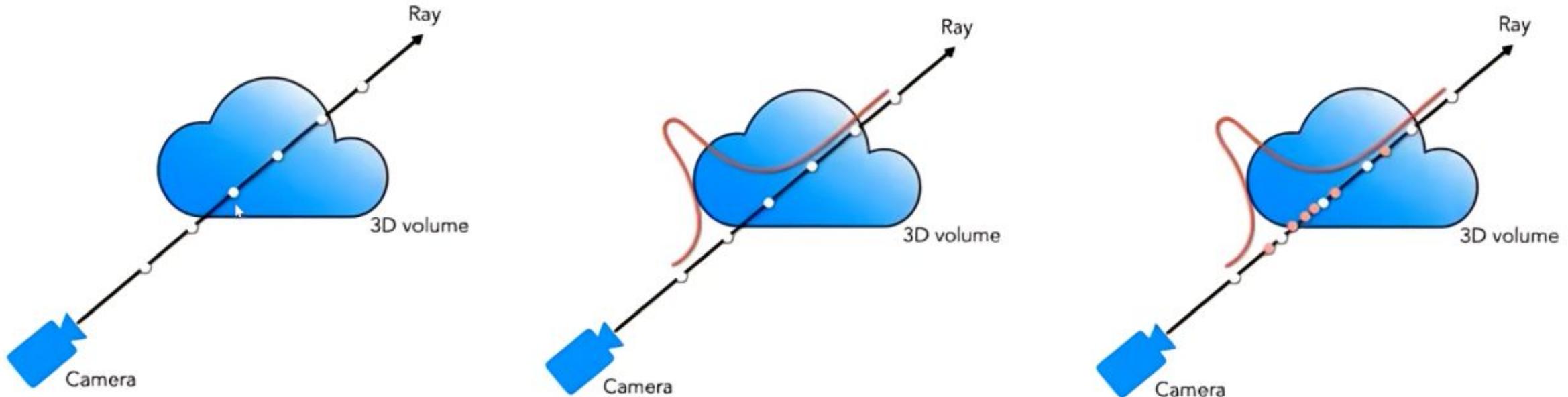
Hierarchical Volume Sampling

$$\hat{C}_c(\mathbf{r}) = \sum_{i=1}^{N_c} w_i c_i, \quad w_i = T_i(1 - \exp(-\sigma_i \delta_i)). \quad (5)$$

Normalizing these weights as $\hat{w}_i = w_i / \sum_{j=1}^{N_c} w_j$ produces a piecewise-constant PDF along the ray. We sample a second set of N_f locations from this distribution



Hierarchical Volume Sampling

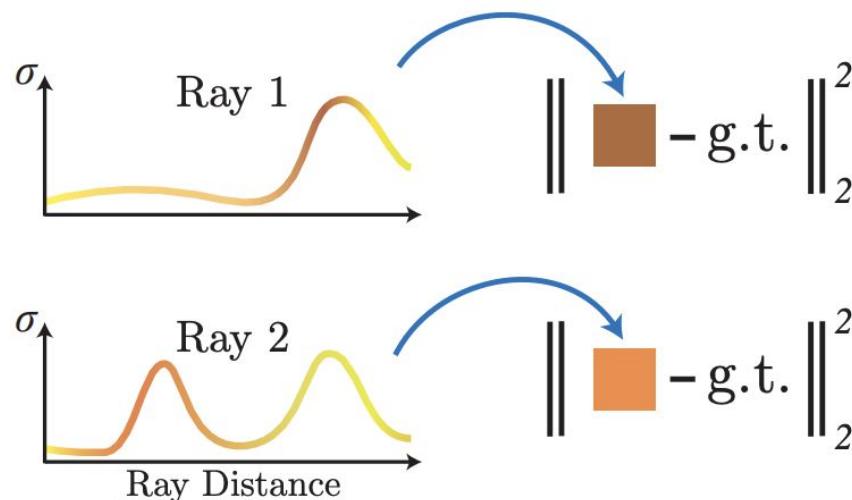


Implementation Details

from both sets of samples. Our loss is simply the total squared error between the rendered and true pixel colors for both the coarse and fine renderings:

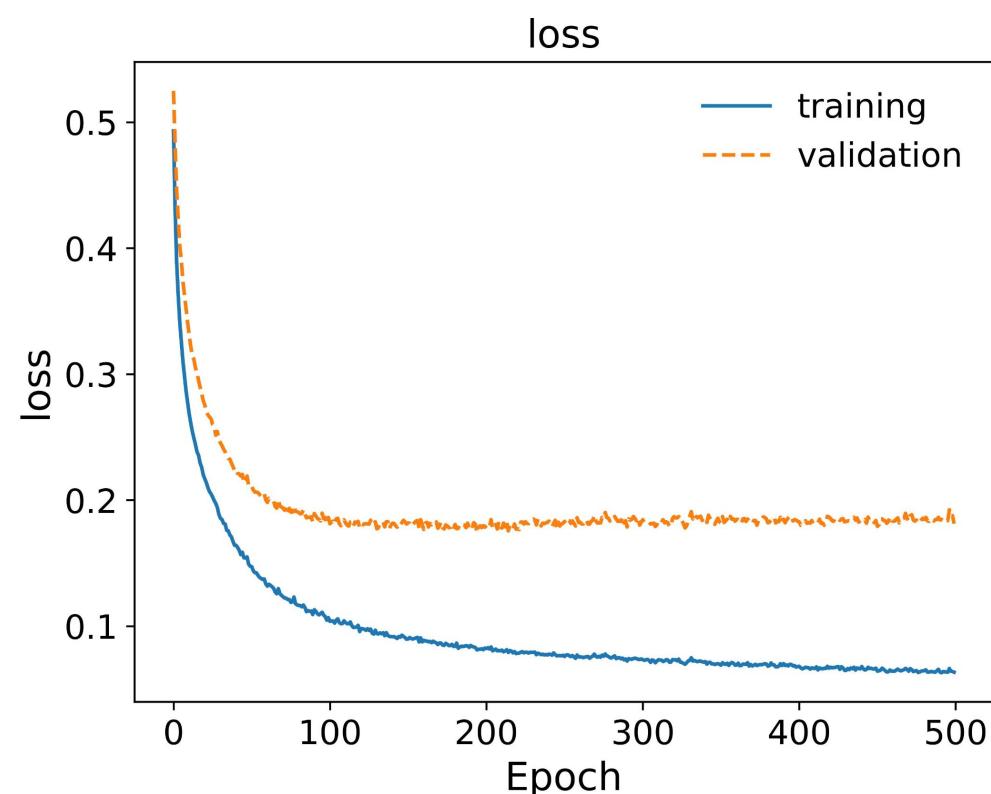
$$\mathcal{L} = \sum_{\mathbf{r} \in \mathcal{R}} \left[\left\| \hat{C}_c(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 + \left\| \hat{C}_f(\mathbf{r}) - C(\mathbf{r}) \right\|_2^2 \right] \quad (6)$$

Volume
Rendering
Rendering
Loss

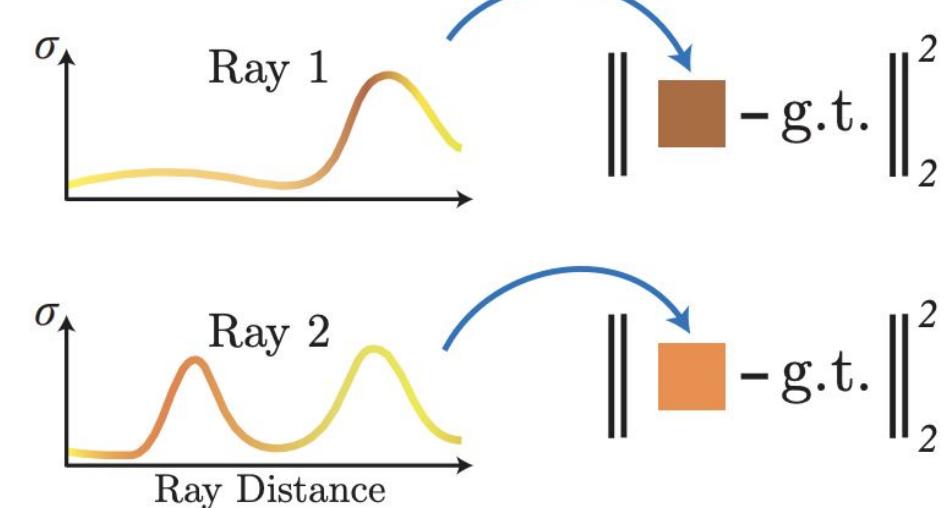


Implementation Details

- Batch Size = 4096, Nc = 64, Nf = 128, Adam
- The optimization for a single scene takes 100-300k iterations to converge on a single NVIDIA V100 GPU (about 1–2 days)



Volume
Rendering Rendering
Loss



Results

Method	Diffuse Synthetic 360° [41]			Realistic Synthetic 360°			Real Forward-Facing [28]		
	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓	PSNR↑	SSIM↑	LPIPS↓
SRN [42]	33.20	0.963	0.073	22.26	0.846	0.170	22.84	0.668	0.378
NV [24]	29.62	0.929	0.099	26.05	0.893	0.160	-	-	-
LLFF [28]	34.38	0.985	0.048	24.88	0.911	0.114	24.13	0.798	0.212
Ours	40.15	0.991	0.023	31.01	0.947	0.081	26.50	0.811	0.250

Table 1: Our method quantitatively outperforms prior work on datasets of both synthetic and real images. We report PSNR/SSIM (higher is better) and LPIPS [50] (lower is better). The DeepVoxels [41] dataset consists of 4 diffuse objects with simple geometry. Our realistic synthetic dataset consists of pathtraced renderings of 8 geometrically complex objects with complex non-Lambertian materials. The real dataset consists of handheld forward-facing captures of 8 real-world scenes (NV cannot be evaluated on this data because it only reconstructs objects inside a bounded volume). Though LLFF achieves slightly better LPIPS, we urge readers to view our supplementary video where our method achieves better multiview consistency and produces fewer artifacts than all baselines.

Limitations and Improvements

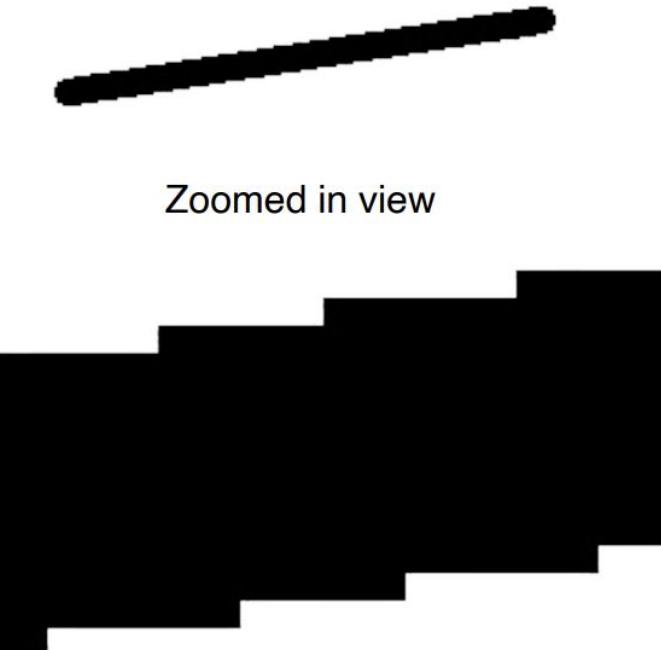
4 years is a long time...

Limitations

- Aliasing problems
- Memory (12+ GB RAM) and compute (2 Days V100) requirement
- Large dataset required (50-100 per scene)

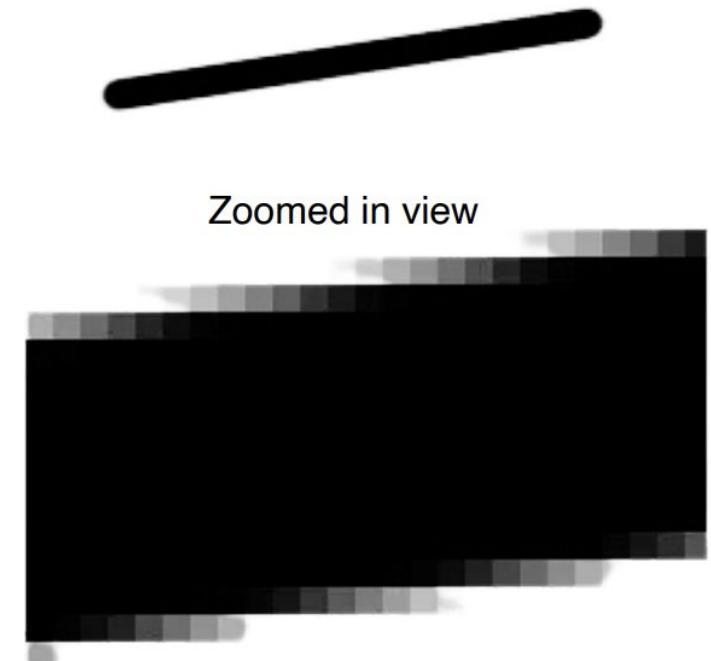


Line with aliasing artifacts



Zoomed in view

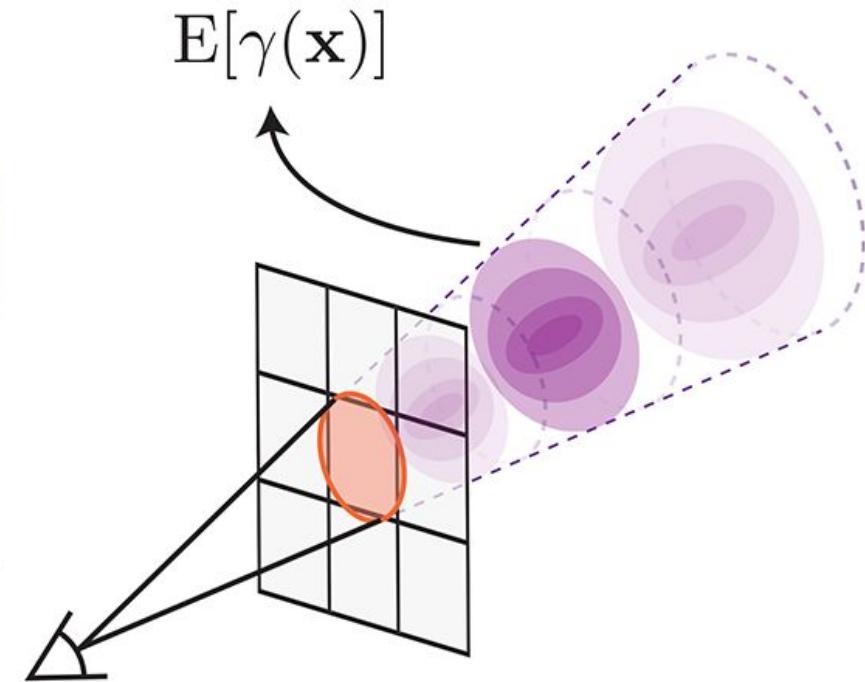
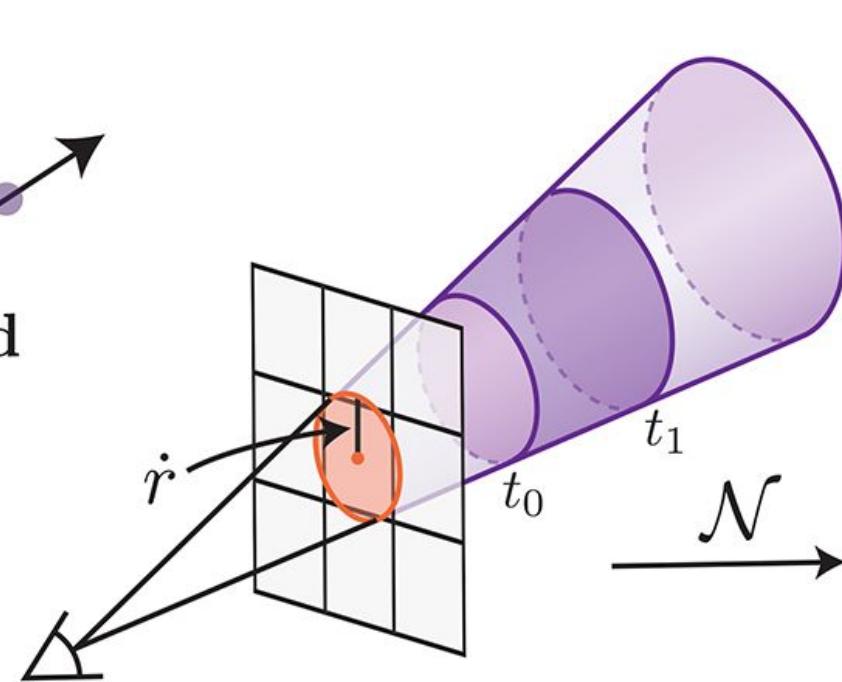
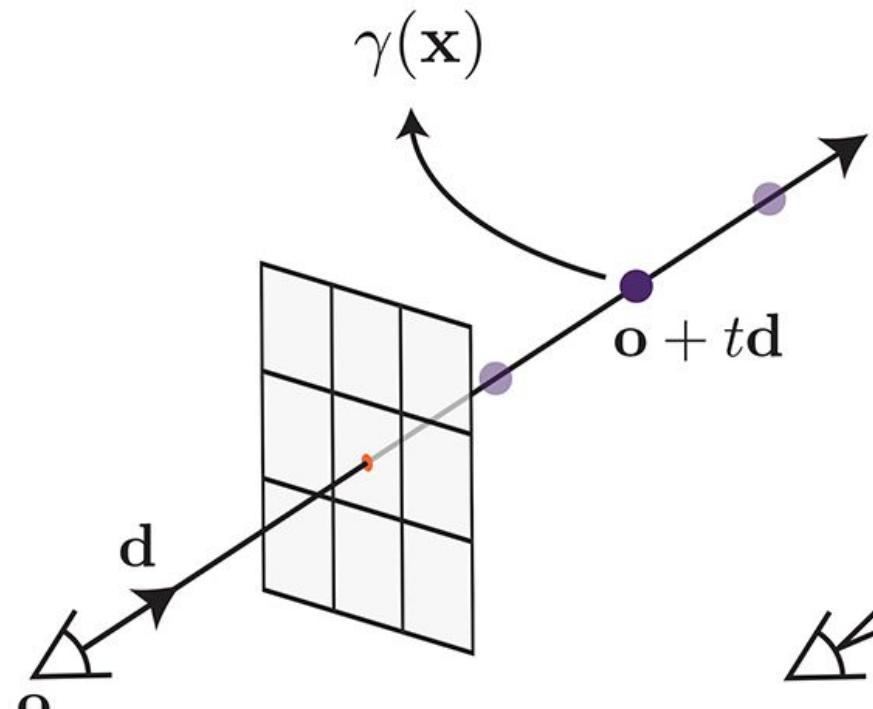
Line with anti-aliasing



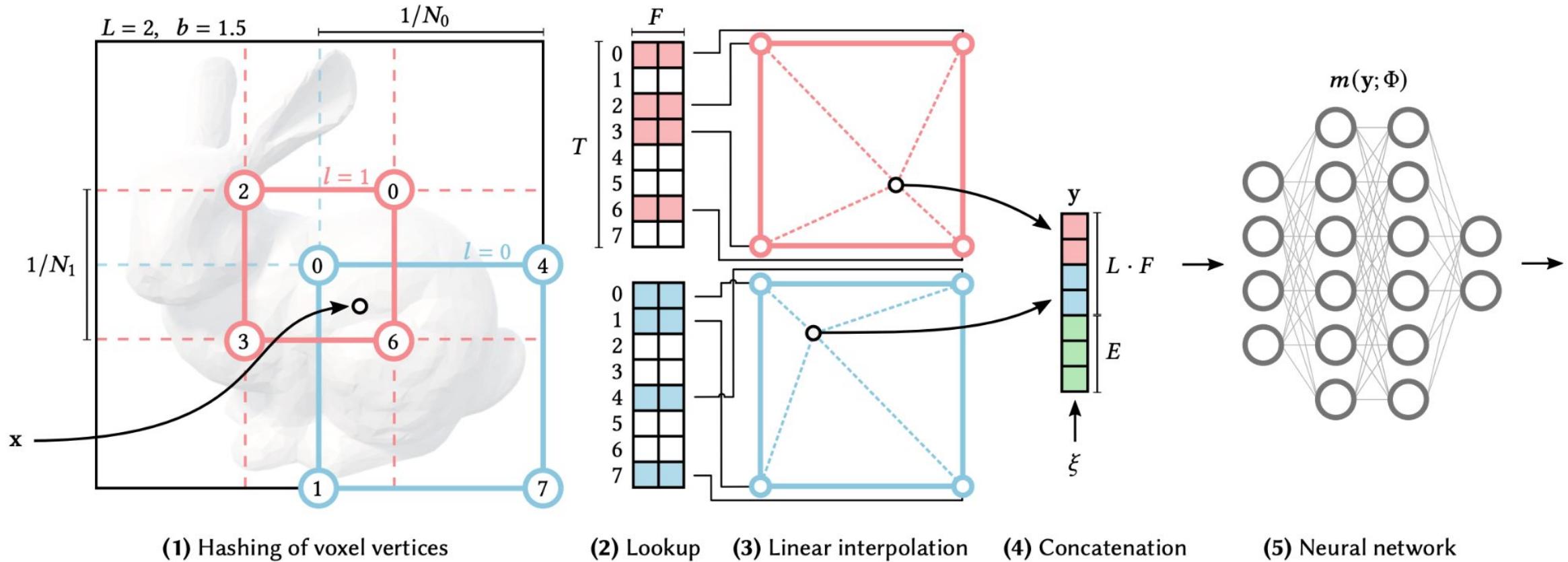
Zoomed in view

Mip-NeRF - Addressing Aliasing and Blurriness (2021)

Google

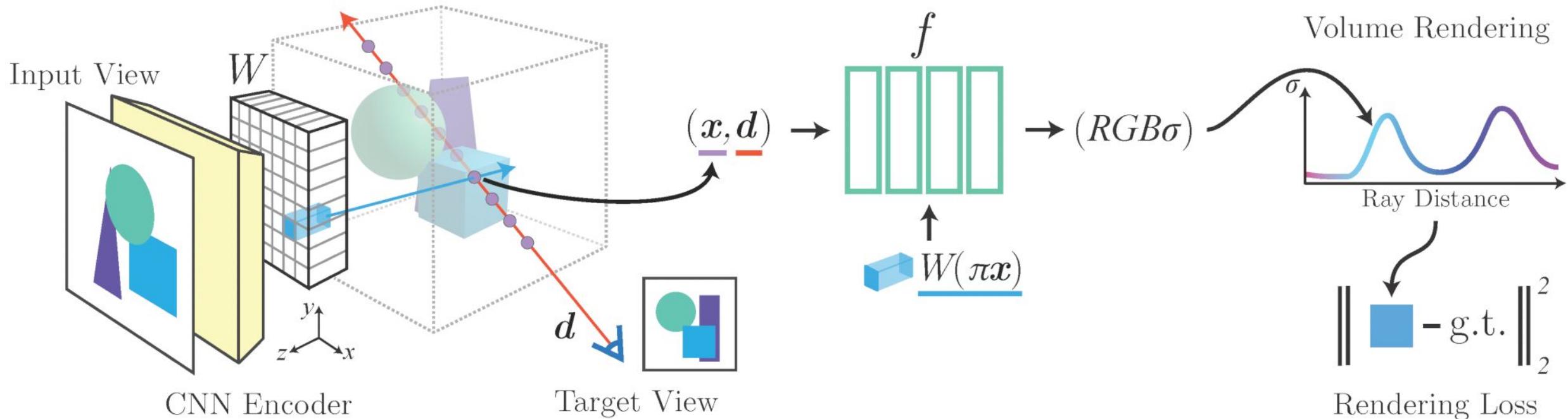


Instant NGP - Addressing Time Constraints (2022)





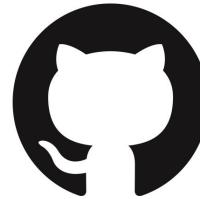
pixelNeRF - Addressing Data Requirements (2021)



Other Works

- <https://www.matthewtancik.com/nerf>
- <https://mobile-nerf.github.io/>
- <https://nerfies.github.io/>
- <https://www.nvidia.com/en-us/research/ai-art-gallery/instant-nerf/>
- <https://meganerf.cmusatyalab.org/>

Slides on

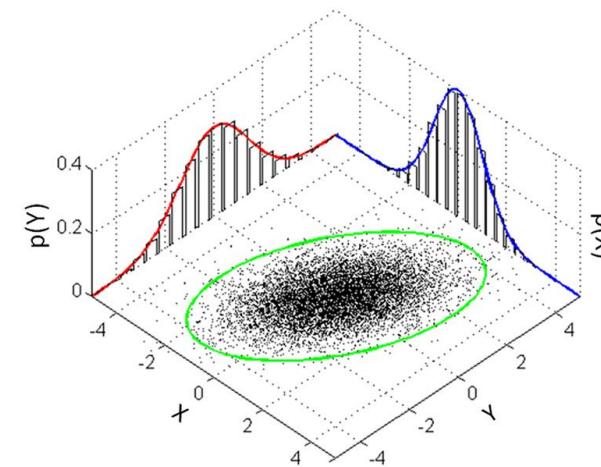


github.com/ut-cvdp/paper-breakdown

Thank You



Up Next: 3D Splatting



UNIVERSITY OF
TORONTO



Follow us on social media

