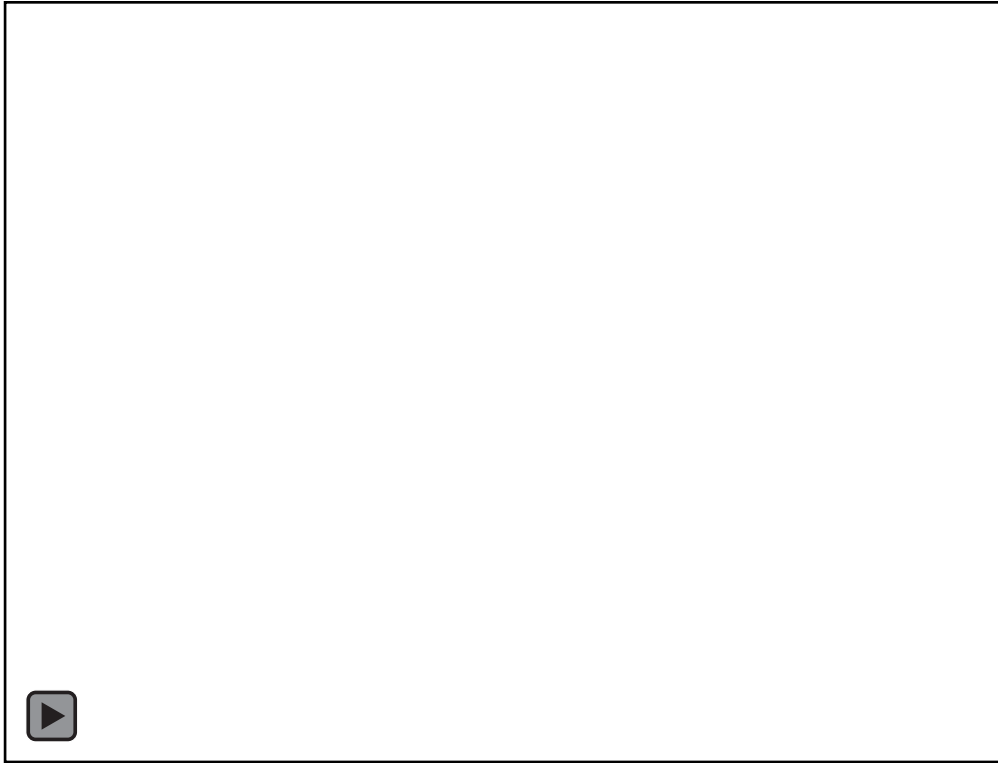


NeRF: Neural Radiance Fields



Computational Photography
Derek Hoiem, University of Illinois

This class

- NeRF: recover 3D model with view-dependent rendering from multiple registered images
- RawNeRF: operate directly on raw images to achieve HDR, denoising, and defocus
- DreamFusion: generate 3D models from text

NeRF: Representing Scenes as Neural Radiance Fields for View Synthesis

Most of following slides from Jon Barron



Ben Mildenhall*



UC Berkeley



Pratul Srinivasan*



UC Berkeley



Matt Tancik*



UC Berkeley



Jon Barron



Google Research



Ravi Ramamoorthi



UC San Diego



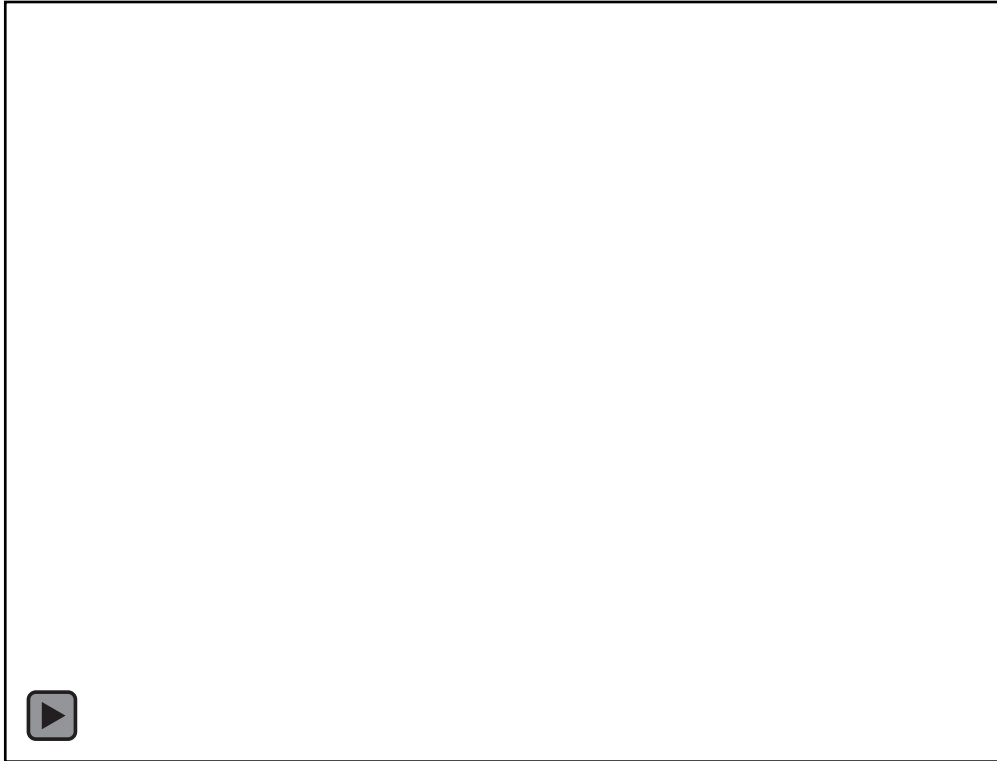
Ren Ng



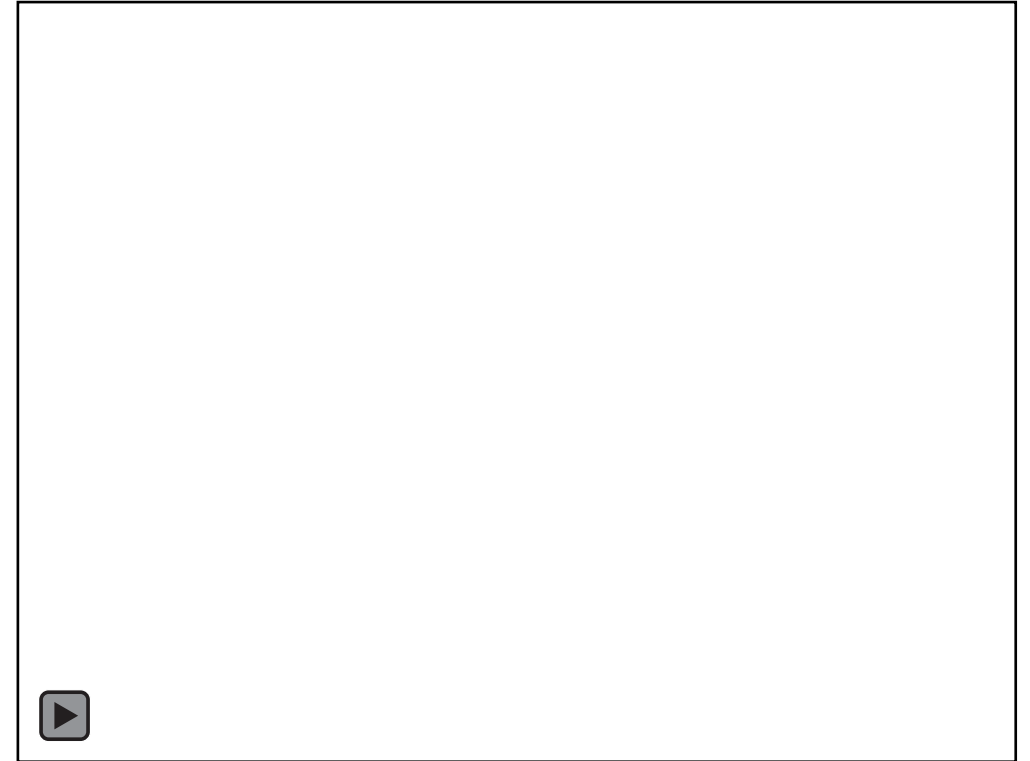
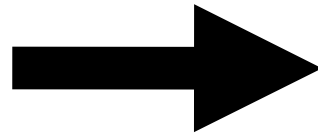
UC Berkeley



Problem: View Interpolation



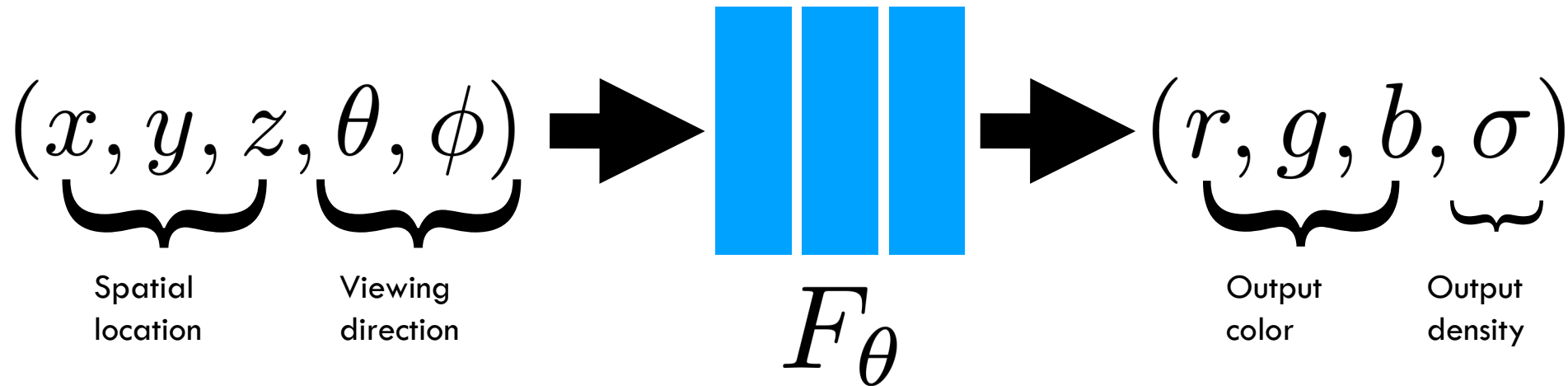
Inputs: sparsely sampled images of scene



Outputs: *new* views of same scene

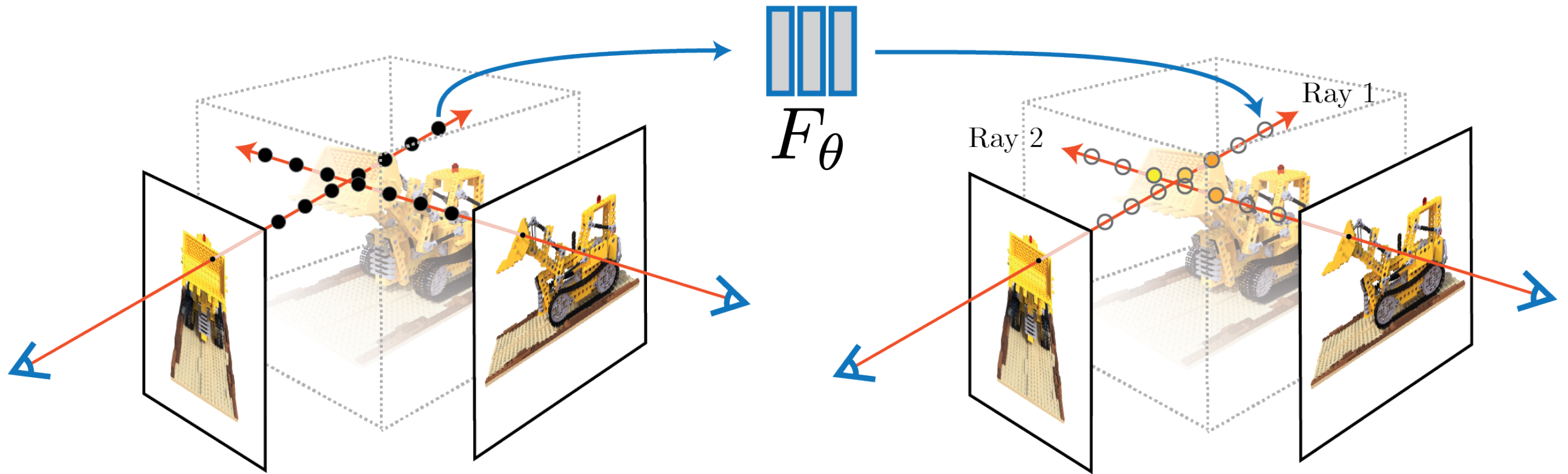
tancik.com/nerf

NeRF (neural radiance fields)



Fully-connected neural
network
9 layers,
256 channels

Volume rendering averages over colors along each ray



Volume rendering is differentiable

Rendering model for ray $r(t) = o + td$:

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

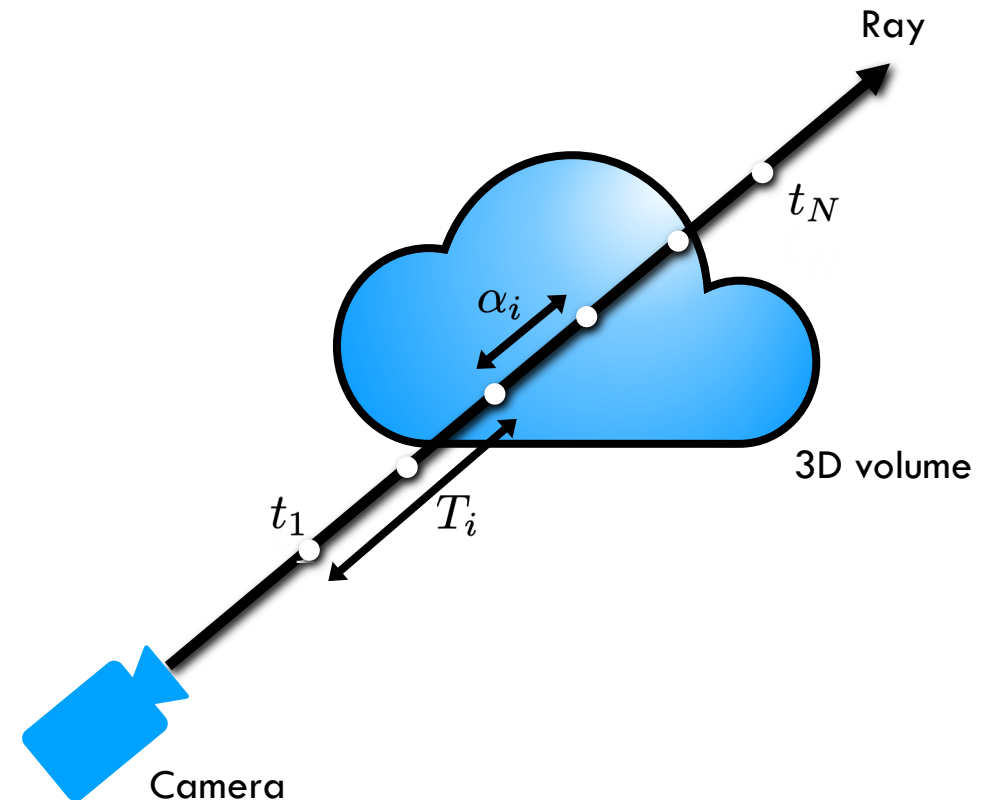
weights colors

How much light is not blocked earlier along ray:

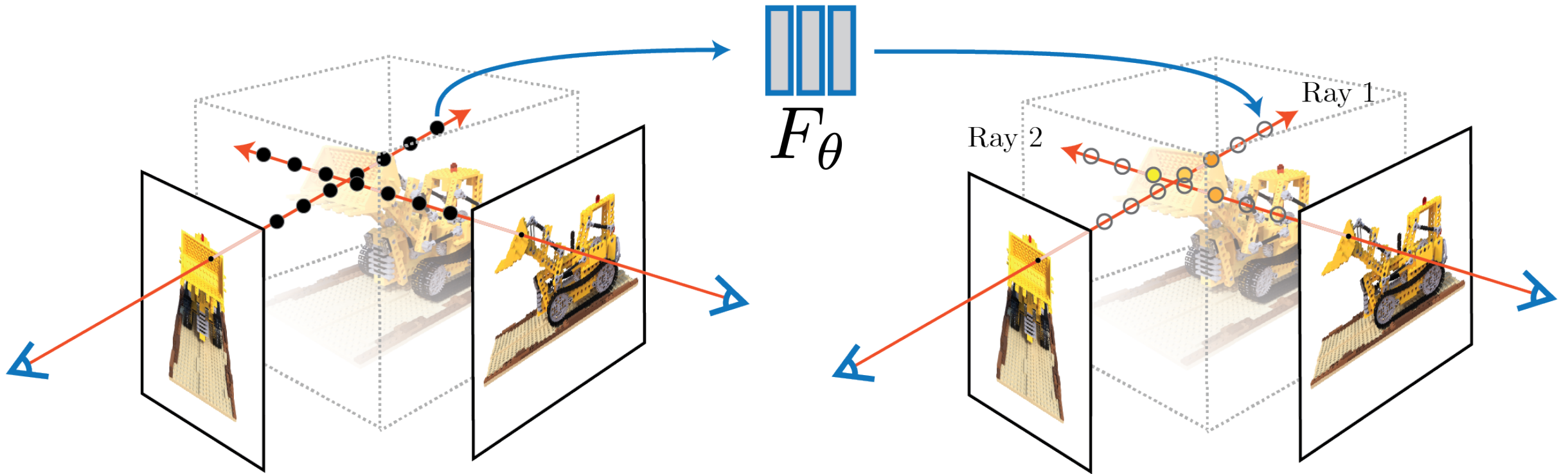
$$T_i = \prod_{j=1}^{i-1} (1 - \alpha_j)$$

How much light is contributed by ray segment i :

$$\alpha_i = 1 - e^{-\sigma_i \delta t_i} \leftarrow \text{Density} * \text{Distance Between Points}$$



Optimize with gradient descent on rendering loss

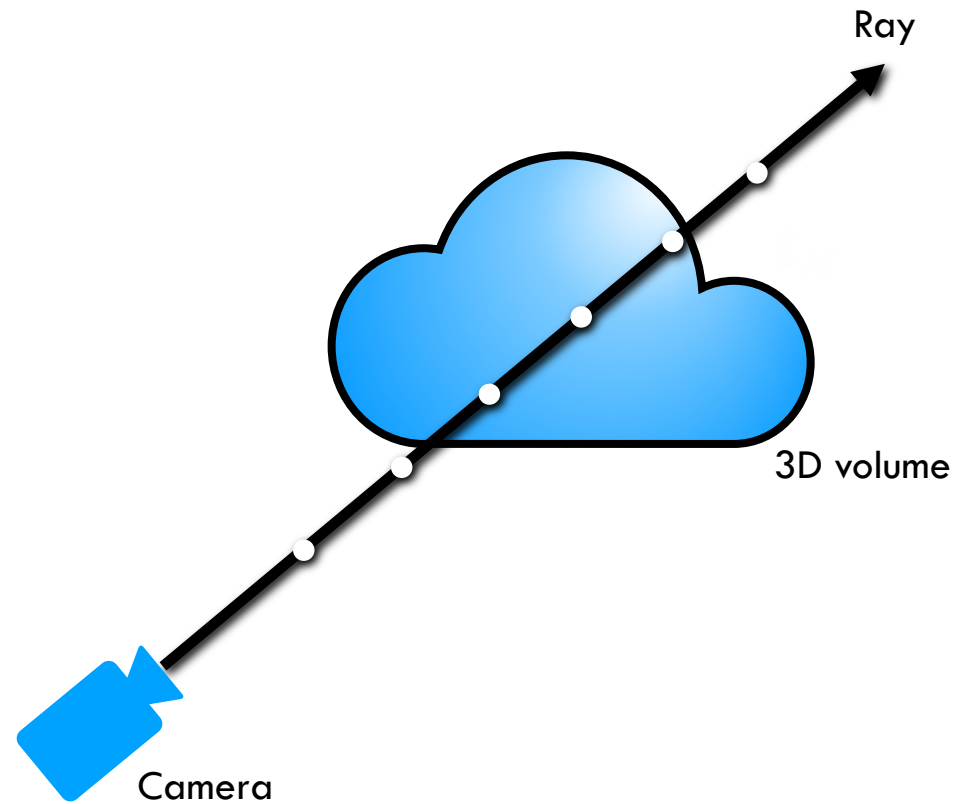


$$\min_{\theta} \sum_i || \text{render}_i(F_\theta) - I_i ||^2$$

Training network to reproduce all input views of the scene



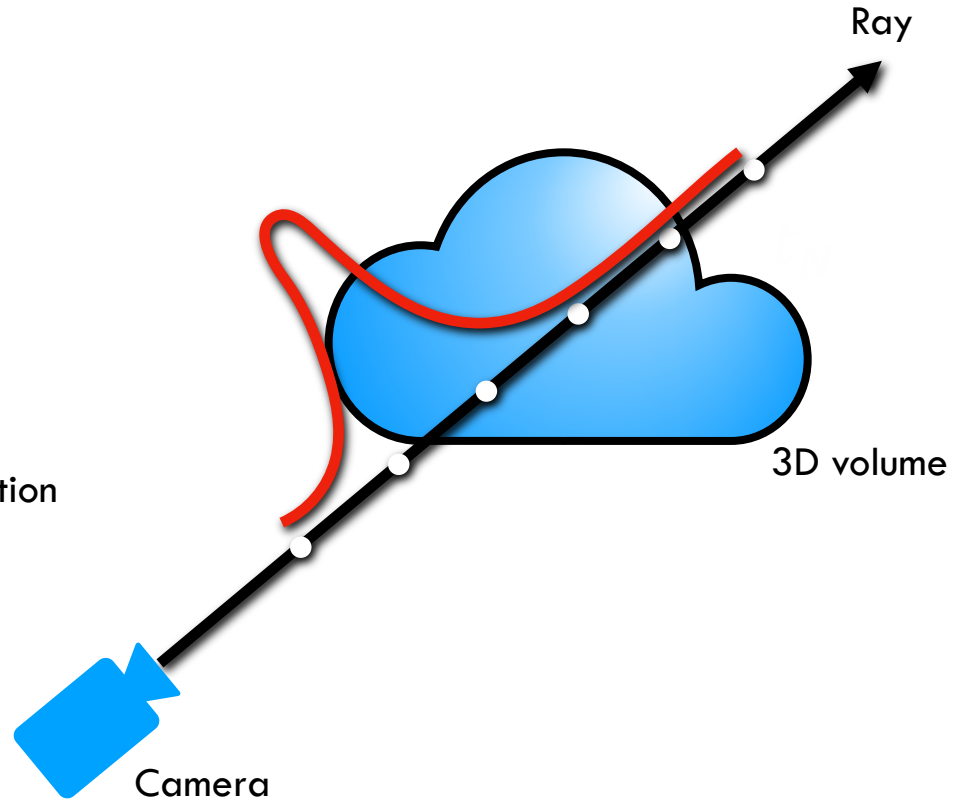
Do two-pass rendering to focus samples near surfaces



Two pass rendering: coarse

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

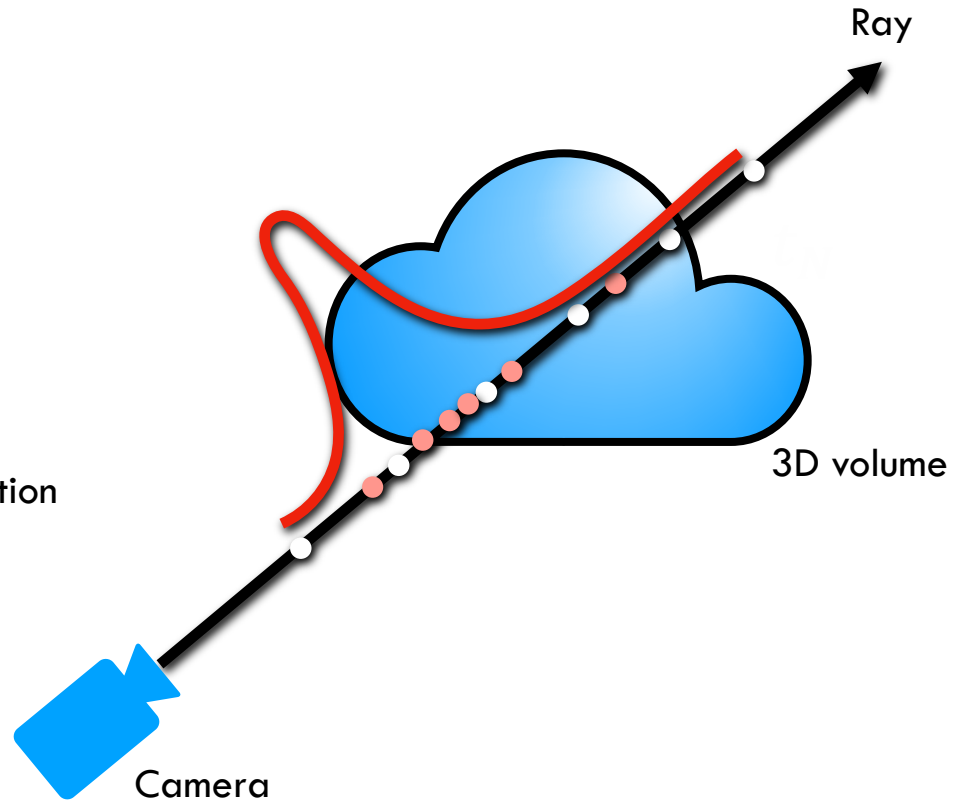
treat weights as probability distribution
for new samples



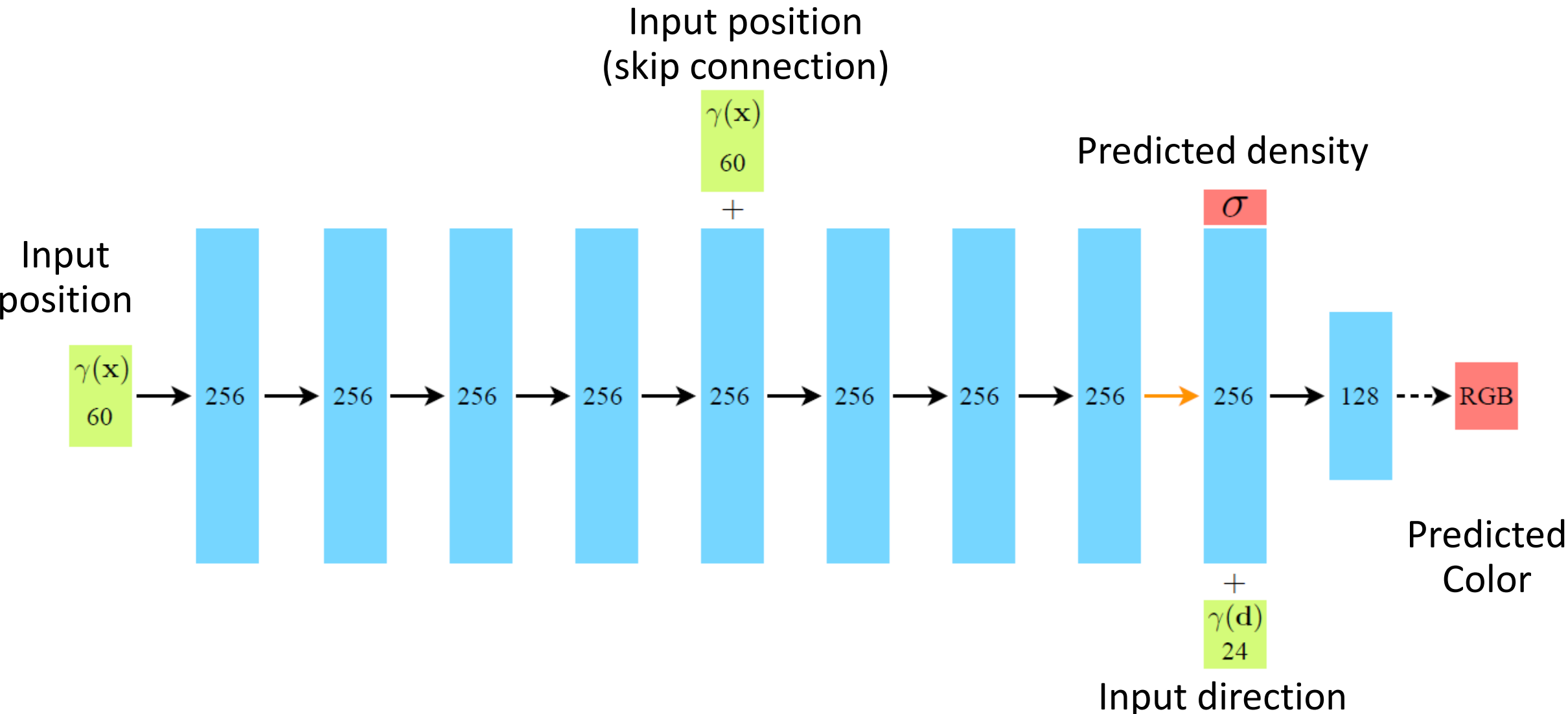
Two pass rendering: fine

$$C \approx \sum_{i=1}^N T_i \alpha_i c_i$$

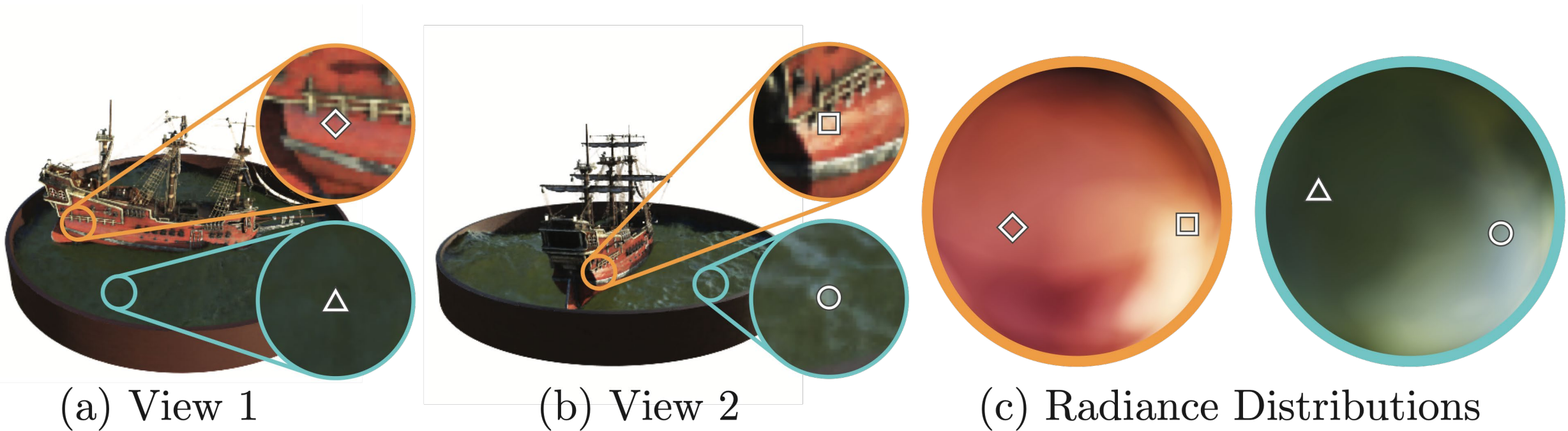
treat weights as probability distribution
for new samples



Network Structure is Dense Multilayer Perceptron



Predicted color depends on the input



Naive implementation produces blurry results



NeRF (Naive)

Naive implementation produces blurry results



NeRF (Naive)



NeRF (with positional encoding)

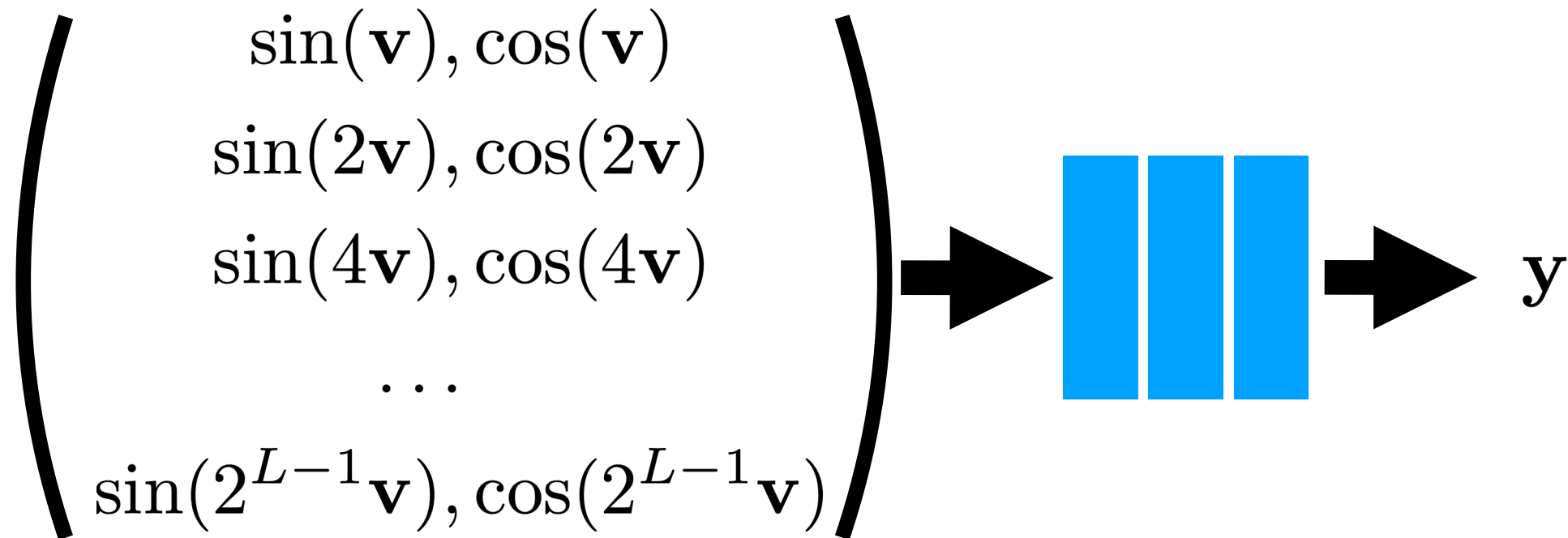
Toy problem: memorizing a 2D image



Toy problem

Ground truth image

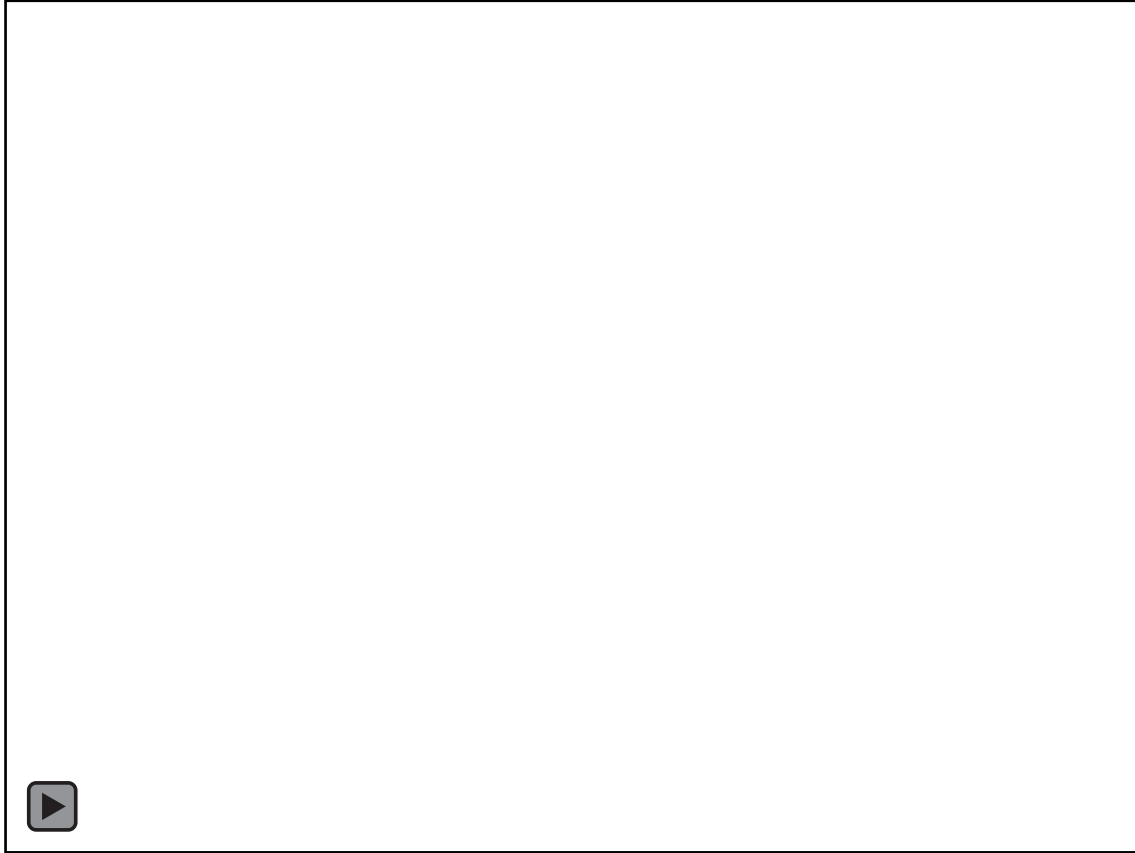




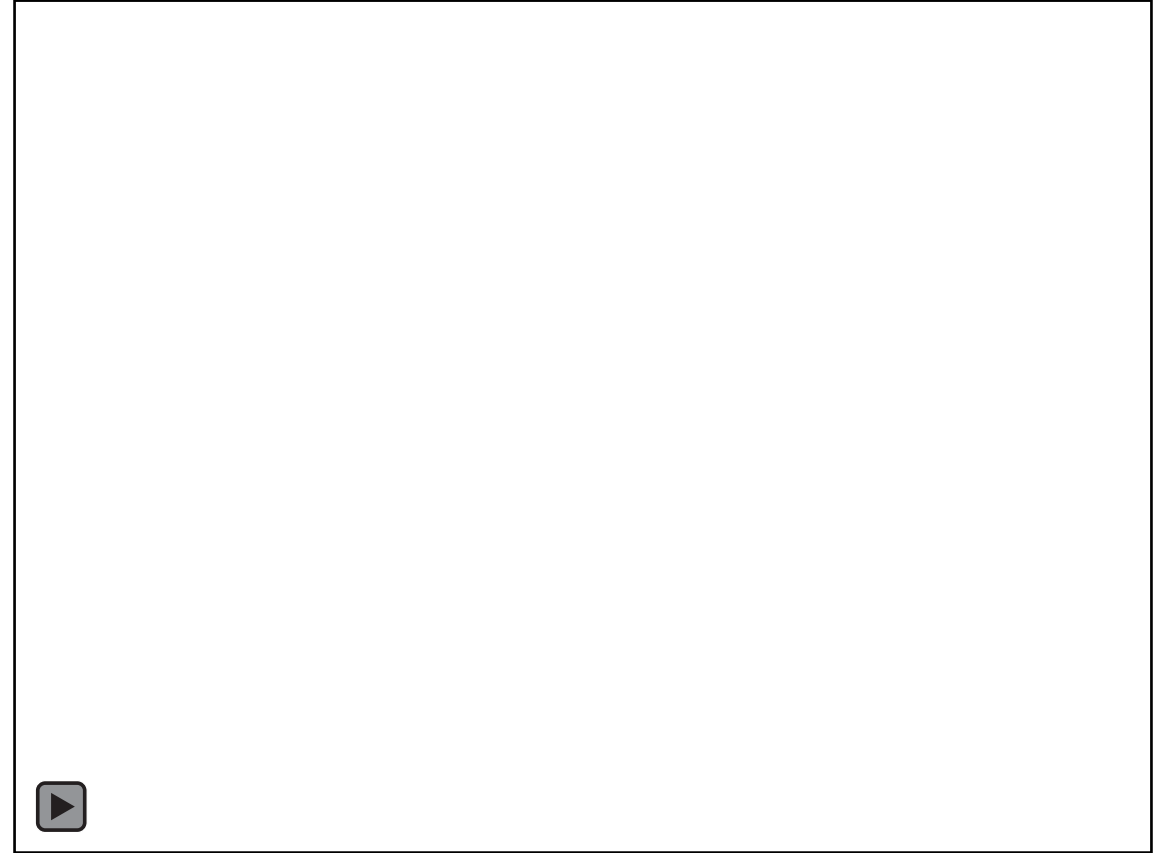
Ground truth image



Positional encoding also directly improves the 3D scene representation

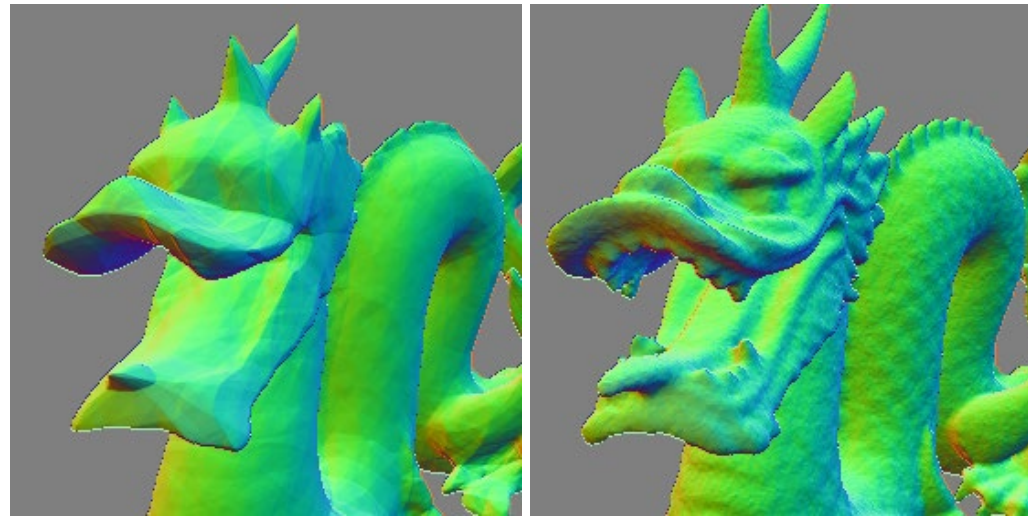


NeRF (Naive)

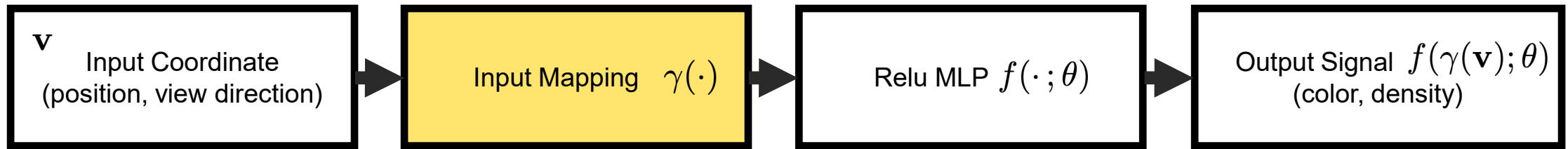


NeRF (with positional encoding)

Fourier Features Let Networks Learn High Frequency Functions in Low Dimensional Domains



Matthew Tancik*, Pratul Srinivasan*, Ben Mildenhall*,
Sara Fridovich-Keil, Nithin Ragahavan, Utkarsh Singhal,
Ravi Ramamoorthi, Jonathan T. Barron, Ren Ng



Positional Encoding [1]: $\gamma(\mathbf{v}) = [\cos(2^0 \mathbf{v}), \sin(2^0 \mathbf{v}), \dots, \cos(2^{L-1} \mathbf{v}), \sin(2^{L-1} \mathbf{v})]$

Random Fourier Features [2]: $\gamma(\mathbf{v}) = [\cos(\mathbf{B}\mathbf{v}), \sin(\mathbf{B}\mathbf{v})]$ $\mathbf{B} \sim \mathcal{N}(0, \sigma^2)$

[1] Vaswani et al.. NeurIPS, 2017
 [2] Rahimi & Recht. NeurIPS, 2007

Neural Tangent Kernel

$$f(\mathbf{x}; \theta) \approx \sum_i (\mathbf{K}^{-1} \mathbf{y})_i k(\mathbf{x}_i, \mathbf{x})$$

Under certain conditions,
neural networks are kernel regression

$$k(\mathbf{x}_i, \mathbf{x}_j) = h_{\text{NTK}}(\langle \mathbf{x}_i, \mathbf{x}_j \rangle)$$

$$h_{\text{NTK}} : \mathbb{R} \rightarrow \mathbb{R}$$

ReLU MLPs correspond to a “dot product” kernel

Dot product of Fourier features

$$\begin{aligned}\langle \gamma(\mathbf{v}_1), \gamma(\mathbf{v}_2) \rangle &= \sum_j (\cos(\mathbf{b}_j^T \mathbf{v}_1) \cos(\mathbf{b}_j^T \mathbf{v}_2) + \sin(\mathbf{b}_j^T \mathbf{v}_1) \sin(\mathbf{b}_j^T \mathbf{v}_2)) \\ &= \sum_j \cos(\mathbf{b}_j^T (\mathbf{v}_1 - \mathbf{v}_2)) \quad (\text{cosine difference trig identity}) \\ &\triangleq h_\gamma(\mathbf{v}_1 - \mathbf{v}_2)\end{aligned}$$

Fourier features turn a dot product of \mathbf{v}_1 and \mathbf{v}_2 into a shift-invariant function of their distance

MLP of Fourier-encoded positions interpolates values between positions

Mapping bandwidth controls underfitting / overfitting



$$\gamma(\mathbf{v}) = [\cos(\mathbf{B}\mathbf{v}), \sin(\mathbf{B}\mathbf{v})] \quad \mathbf{B} \sim \mathcal{N}(0, \sigma^2)$$

Mapping bandwidth controls underfitting / overfitting



$$\gamma(\mathbf{v}) = [\cos(\mathbf{B}\mathbf{v}), \sin(\mathbf{B}\mathbf{v})] \quad \mathbf{B} \sim \mathcal{N}(0, \sigma^2)$$

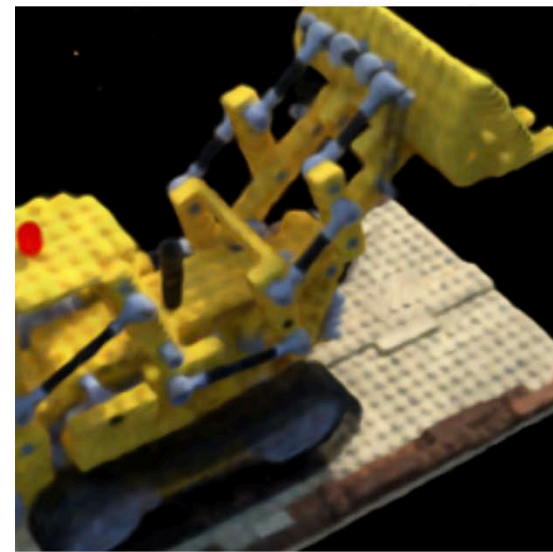
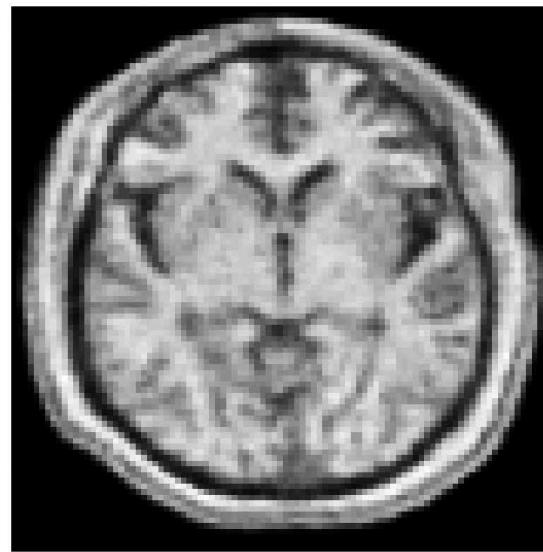
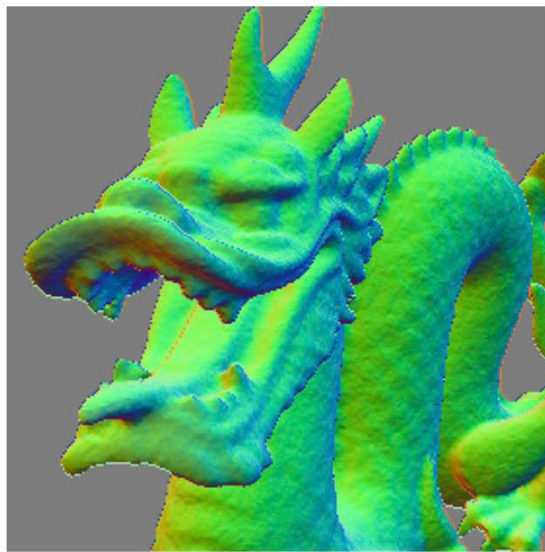
Mapping bandwidth controls underfitting / overfitting



$$\gamma(\mathbf{v}) = [\cos(\mathbf{B}\mathbf{v}), \sin(\mathbf{B}\mathbf{v})] \quad \mathbf{B} \sim \mathcal{N}(0, \sigma^2)$$

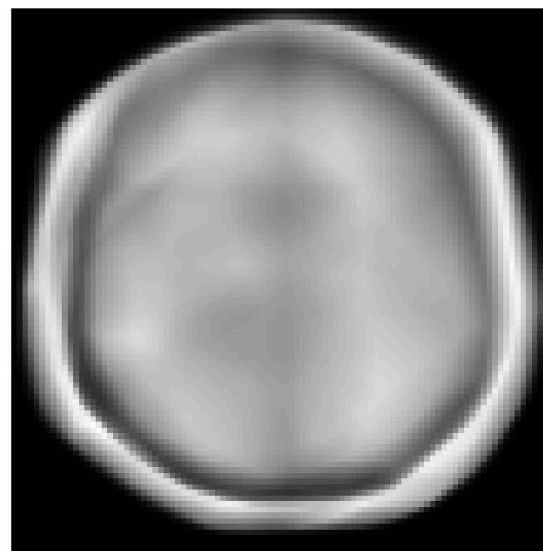
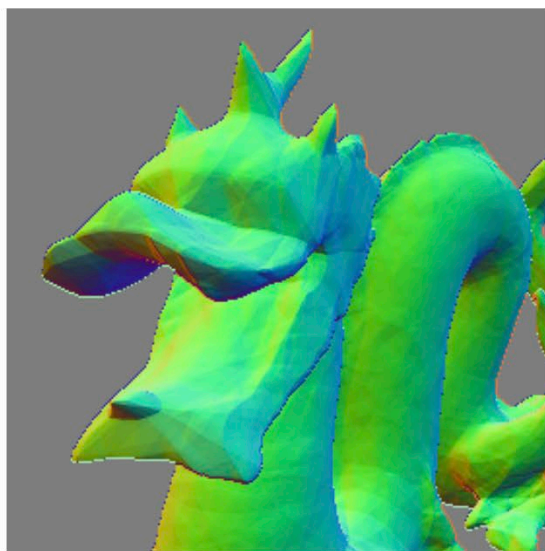
With Fourier features

$$\gamma(\mathbf{v}) = \text{FF}(\mathbf{v})$$



No Fourier features

$$\gamma(\mathbf{v}) = \mathbf{v}$$



(b) Image regression
 $(x, y) \rightarrow \text{RGB}$

(c) 3D shape regression
 $(x, y, z) \rightarrow \text{occupancy}$

(d) MRI reconstruction
 $(x, y, z) \rightarrow \text{density}$

(e) Inverse rendering
 $(x, y, z) \rightarrow \text{RGB, density}$

Very simple!

```
B = SCALE * np.random.normal(shape=(input_dims, NUM_FEATURES))  
x = np.concatenate([np.sin(x @ B), np.cos(x @ B)], axis=-1)  
x = nn.Dense(x, features=256)
```





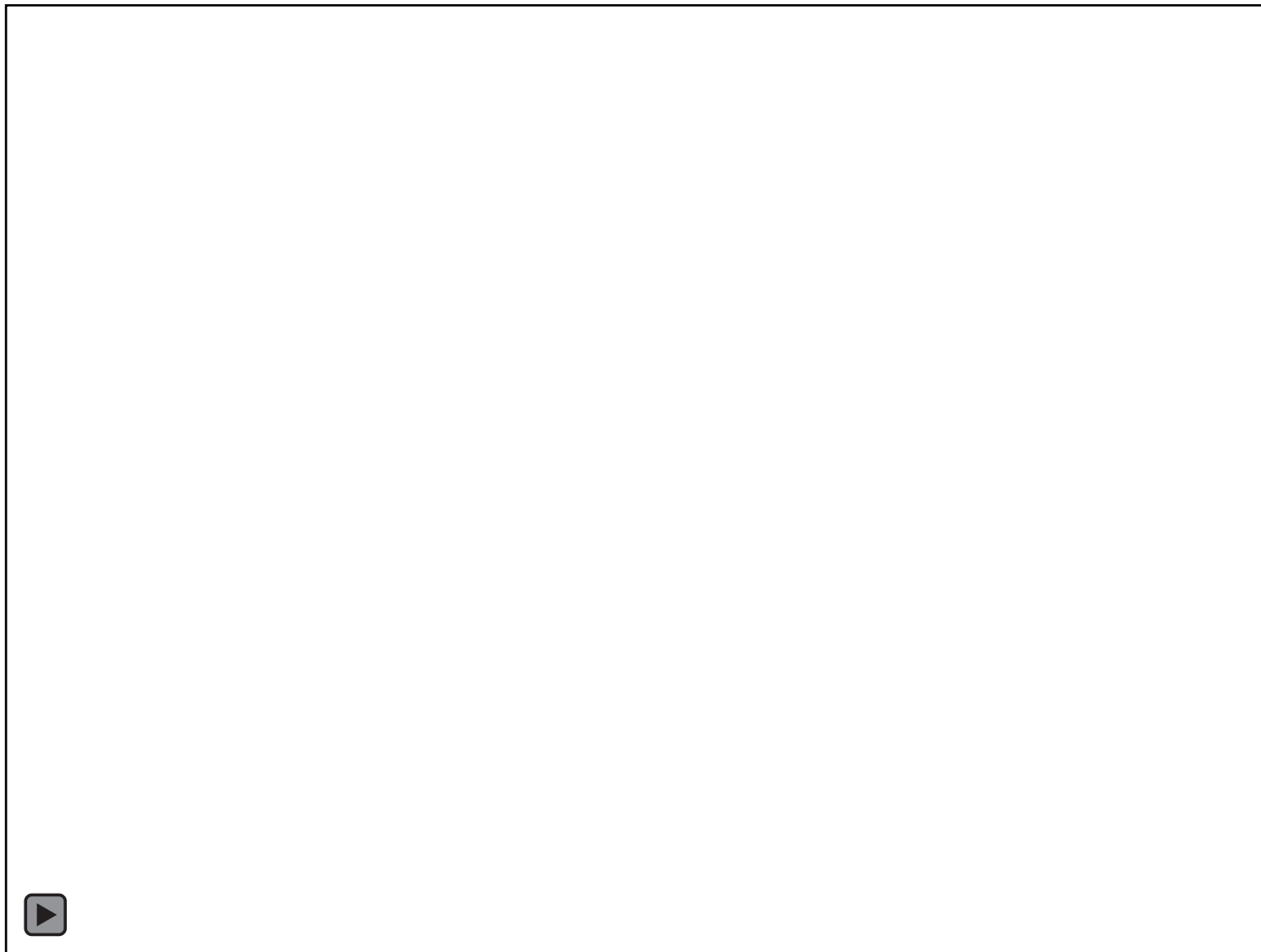
View-Dependent Effects



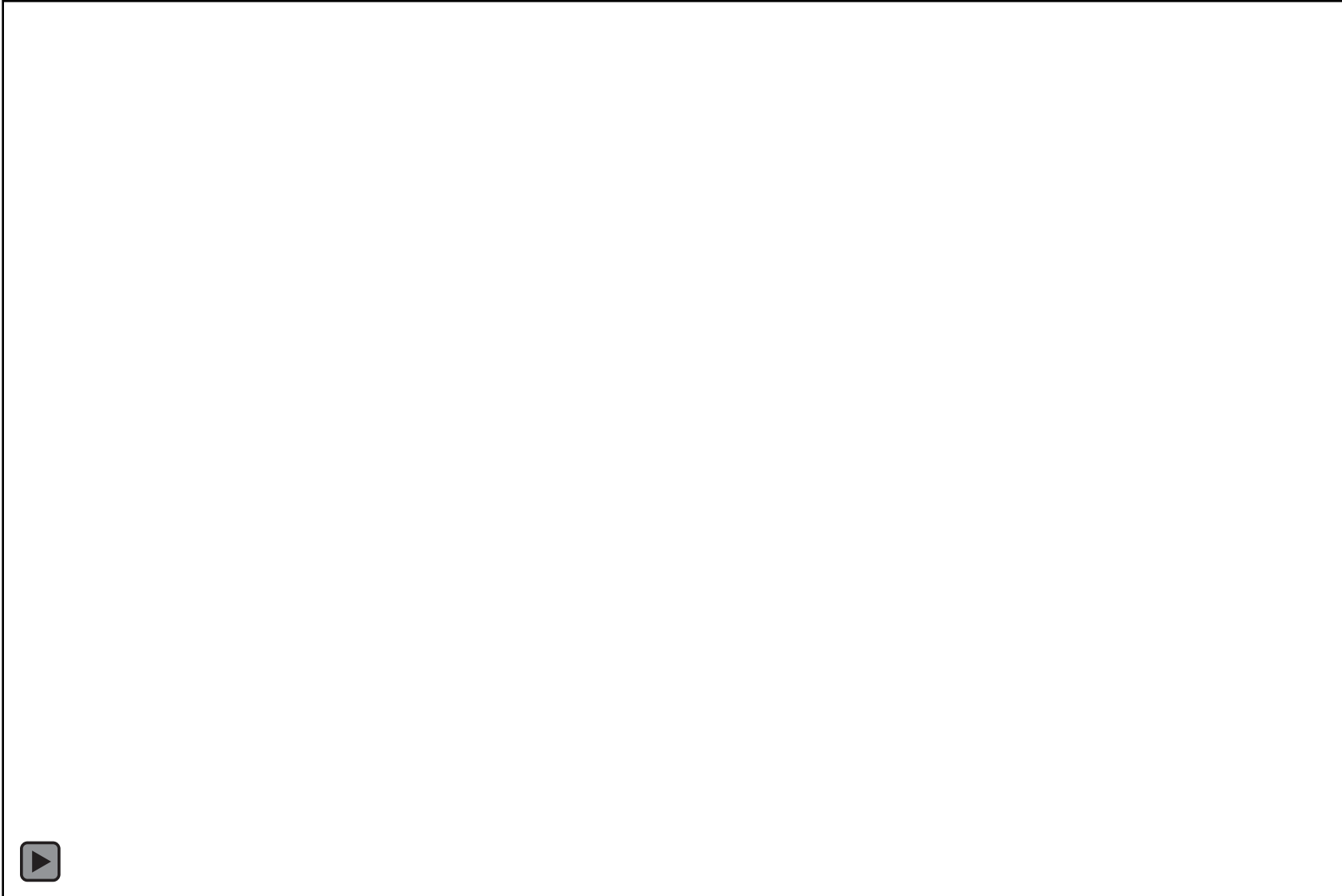
Detailed Geometry & Occlusion



Detailed Geometry & Occlusion



Meshable



MonoSDF (2022)

- Solve for features in multi-resolution grid and pass into MLP w/ positional encoding
- Use per-view surface normal prediction as constraint
- Solve SDF (distance to surface)

MonoSDF: Exploring Monocular Geometric Cues for Neural Implicit Surface Reconstruction

Zehao Yu¹ Songyou Peng^{2,3} Michael Niemeyer^{1,3} Torsten Sattler⁴ Andreas Geiger^{1,3}

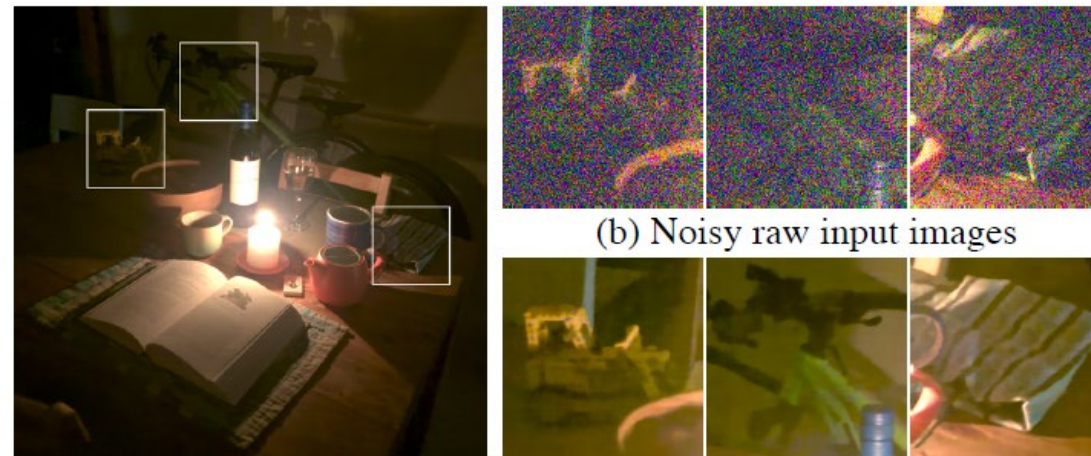


NeRF in the Dark:

High Dynamic Range View Synthesis from Noisy Raw Images

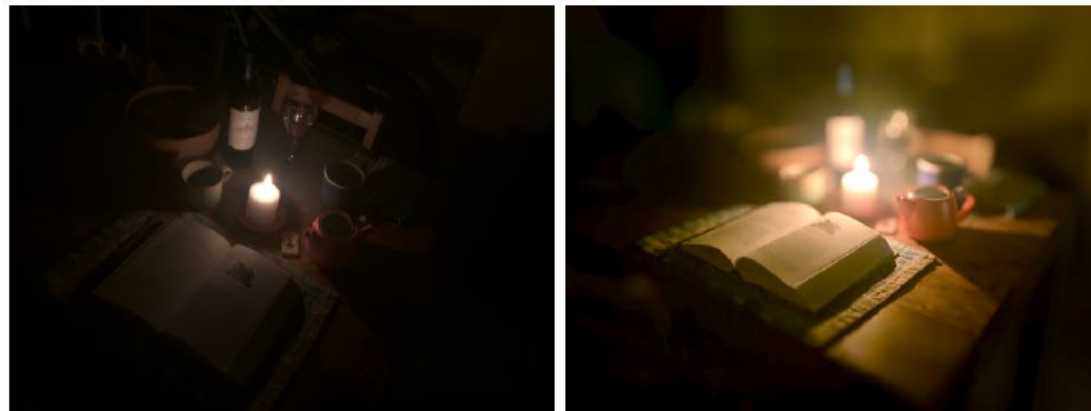
Ben Mildenhall Peter Hedman Ricardo Martin-Brualla Pratul P. Srinivasan Jonathan T. Barron

CVPR 2022



(a) Reconstructed candlelit scene

(c) RawNeRF renderings

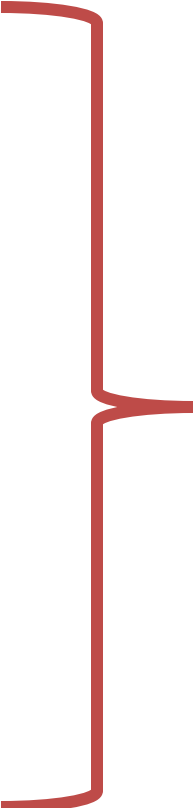


(d) Changing viewpoint, focus, exposure, and tonemapping

<https://www.youtube.com/watch?v=JtBS4KBcKVc>

Measurement noise is complicated by camera processes

1. Raw camera measurement: photons converted to electrical charge, recorded in 10-14 bits
 - Black level subtraction removes noise effects
 - Gaussian measurement noise
2. Color filter demosaicking
 - Color light captured in Bayer pattern
 - Interpolation completes RGB values
3. Color correction and white balance
 - Convert from camera color space to standard RGB
 - White balance to remove tint due to lighting
4. Gamma compression and tone mapping
 - Gamma compression to shift toward darker values with higher human sensitivity
 - Preserve contrast when mapping to 8 bits



Improves viewing quality but complicates noise model and discards information

RawNeRF Method

- COLMAP on jpg images is used for camera registration

- NeRF operates on raw images

- Loss emphasizes darker pixels, since HDR values can get very high

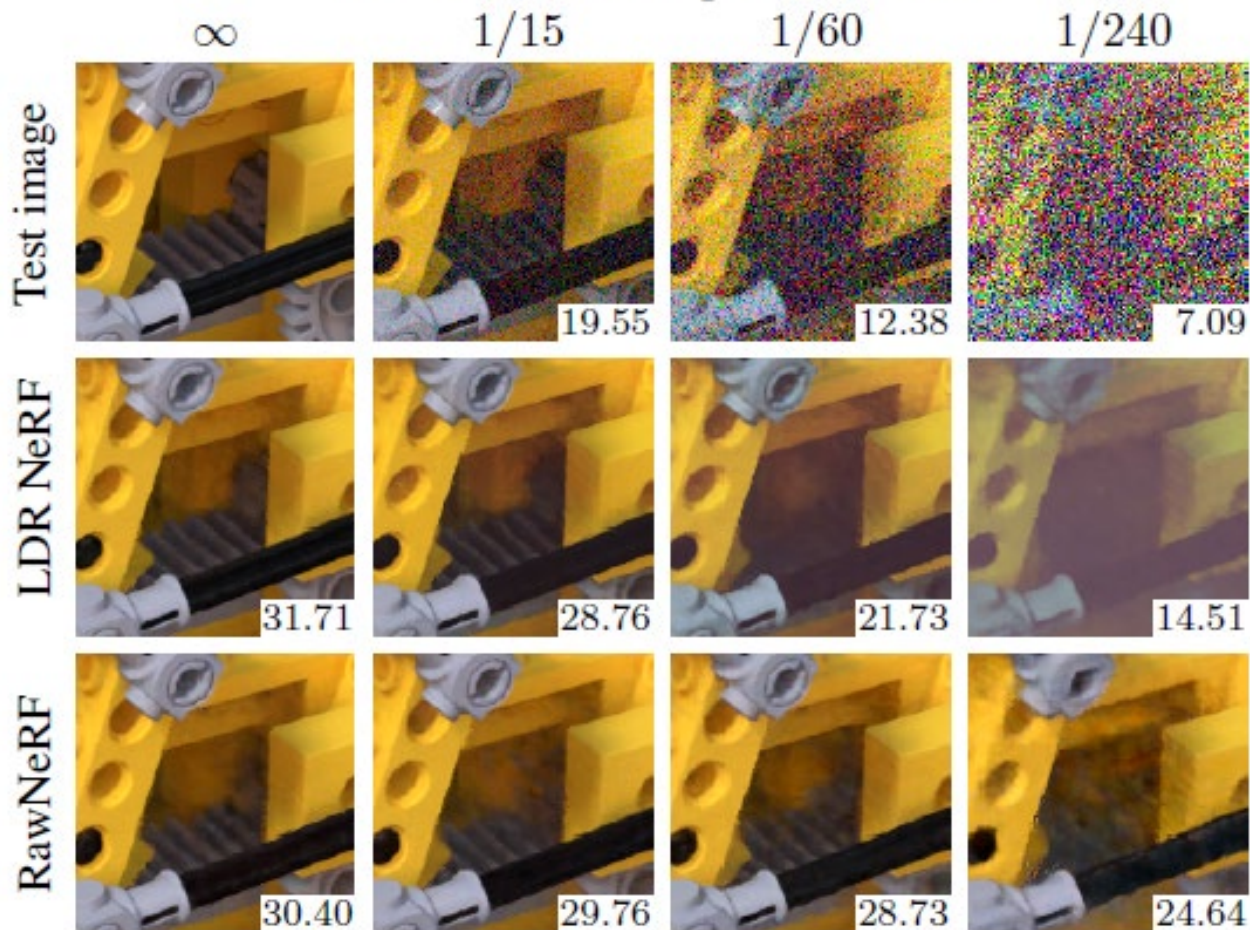
$$L(\hat{y}, y) = \sum_i w_i (\hat{y}_i - y_i)^2 \quad w_i = 1/(\hat{y}_i + \epsilon)$$

- Optimize over exposure, accounting for shutter time

- Use regular NeRF architecture, except final layer is exponential instead of sigmoid

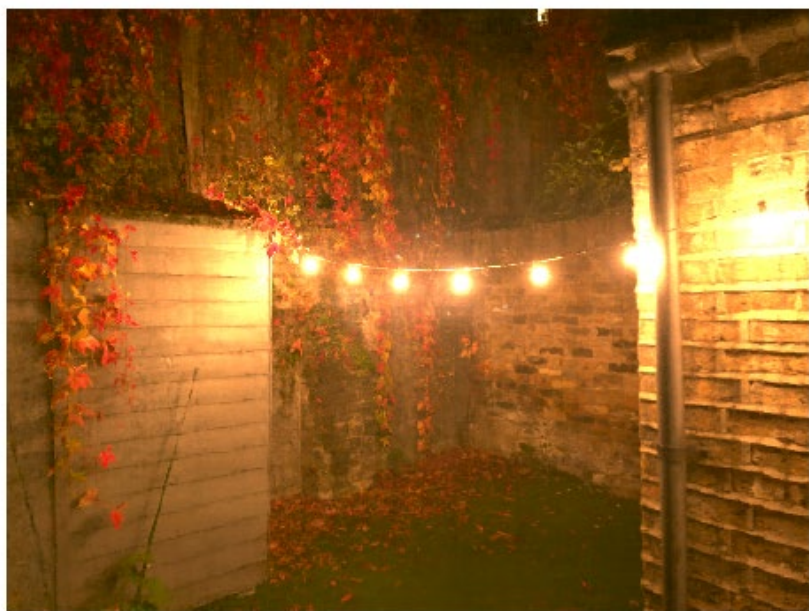
RawNeRF handles noise better

Simulated shutter speed (seconds)

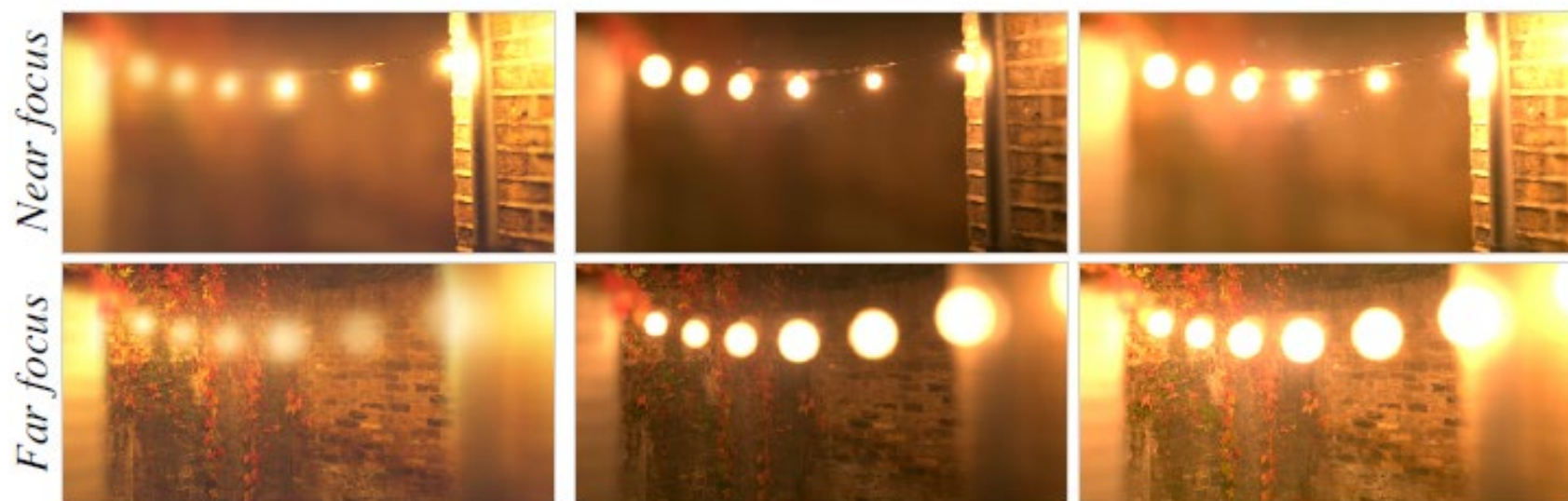


Method	Simulated shutter speed (seconds)						
	∞	1/7	1/15	1/30	1/60	1/120	1/240
Noisy input	-	23.33	19.65	16.03	12.51	9.40	7.18
LDR NeRF	33.16	31.25	29.14	26.10	22.31	18.27	14.87
RawNeRF	32.15	32.11	31.94	31.59	30.94	29.69	27.73

Synthetic re-focusing benefits from HDR and 3D model



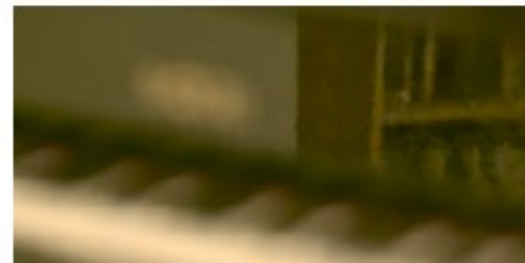
(a) Full RawNeRF output



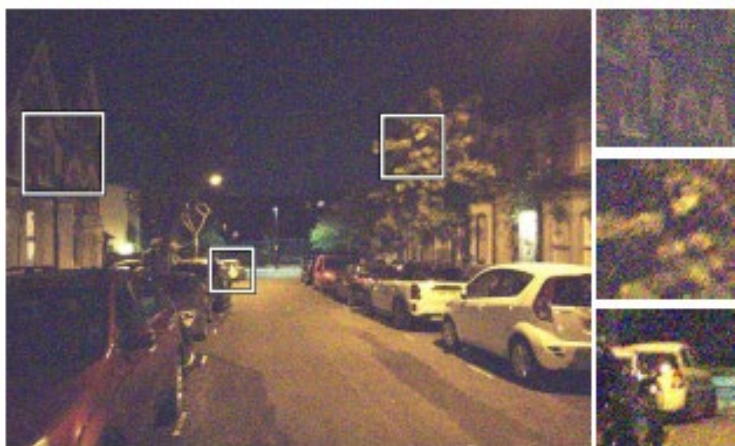
(b) LDR NeRF defocus (c) RawNeRF defocus and exposure variation



(d) Seeing behind objects



(e) Revealing reflections



Noisy test image

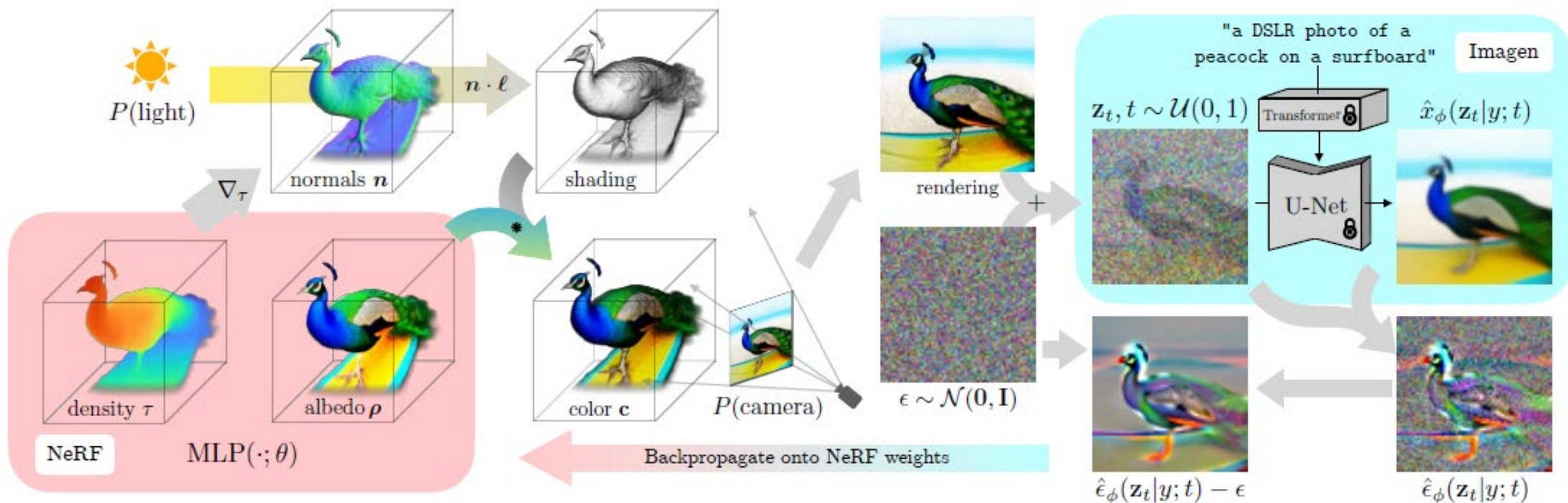
RawNeRF rendering

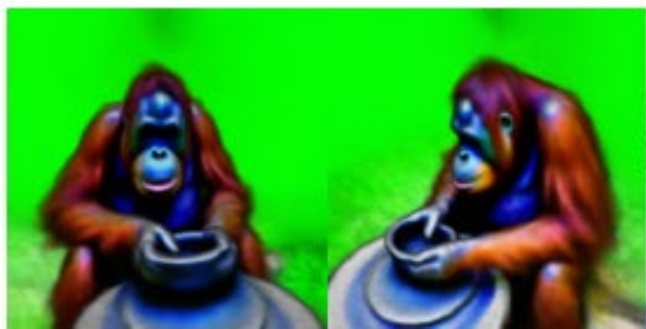
New viewpoint, HDR tonemapping

DREAMFUSION: TEXT-TO-3D USING 2D DIFFUSION

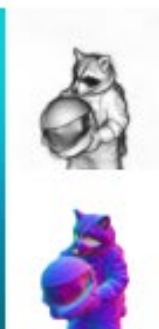
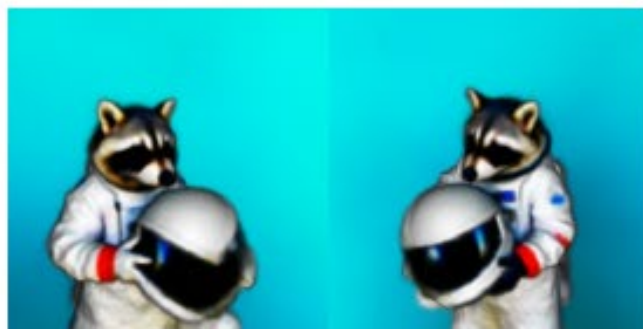
Ben Poole¹, Ajay Jain², Jonathan T. Barron¹, Ben Mildenhall¹

1. **Input:** “a DSLR photo of a peacock on a surfboard, randomly initialized NeRF model, Imagen generator
2. **Generate with Nerf:** albedo and surface normal image, compute shading, render final color
3. **Correct with Imagen:** the generated image (plus noise), subtract noise from delta, and provide as feedback to NeRF

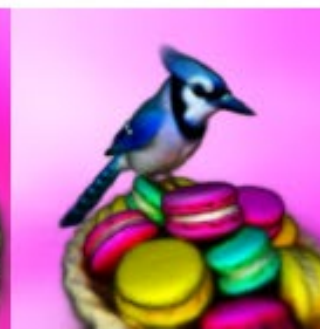




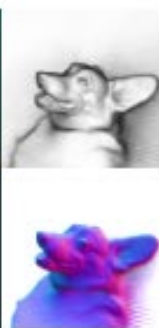
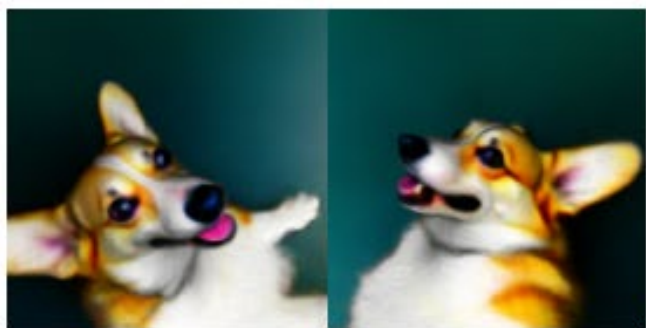
an orangutan making a clay bowl on a throwing wheel*



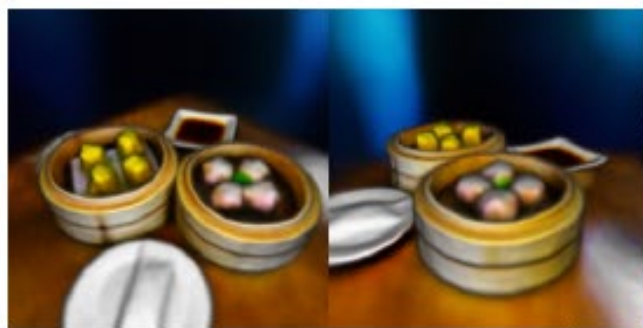
a raccoon astronaut holding his helmet†



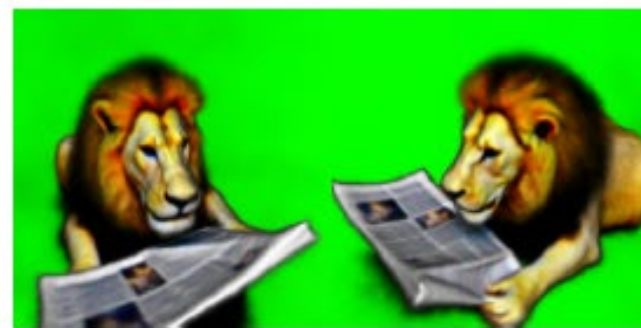
a blue jay standing on a large basket of rainbow macarons*



a corgi taking a selfie*



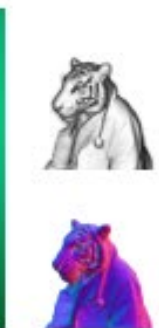
a table with dim sum on it†



a lion reading the newspaper*



Michelangelo style statue of dog reading news on a cellphone



a tiger dressed as a doctor*



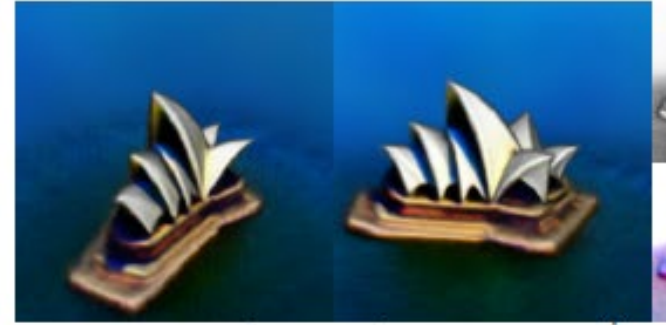
a steam engine train, high resolution*



a frog wearing a sweater*



a humanoid robot playing the cello*



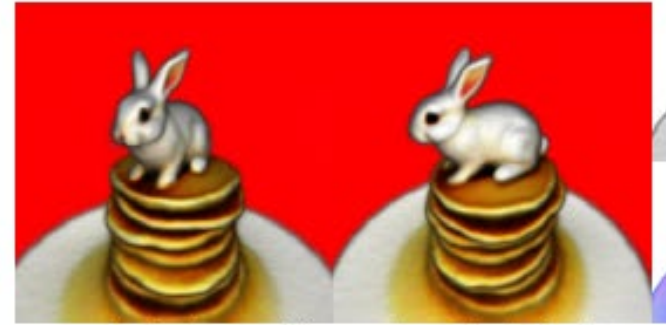
Sydney opera house, aerial view†



an all-utility vehicle driving across a stream†



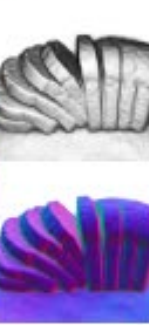
a chimpanzee dressed like Henry VIII king of England*



a baby bunny sitting on top of a stack of pancakes†



a sliced loaf of fresh bread



a bulldozer clearing away a pile of snow*

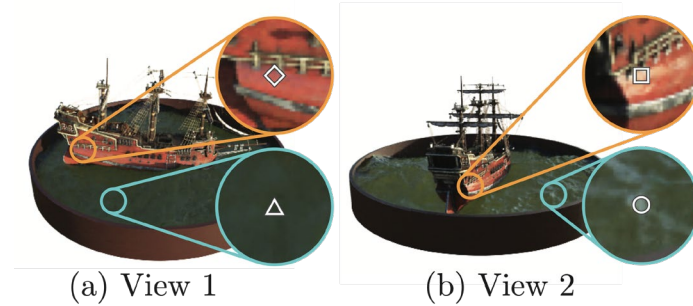


a classic Packard car*



Summary

- NeRF encodes a surface with diffuse and non-diffuse color components by mapping $(x,y,z,direction)$ to $(density, r,g,b)$
- RawNeRF operates directly in raw measurement space, achieving denoising, HDR, and refocus effects
- DreamFusion combines image generation with NeRF to create 3D models from text



a raccoon astronaut holding his helmet†

Coming up

- Enjoy Fall break next week!
- Project 5 due Monday, Nov 28