



HSE UNIVERSITY



Reconstruction of 3D Scenes

An Overview of Neural Fields for Novel View Synthesis

Kirill Struminsky

Research Fellow @ HSE Centre for Deep Learning and Bayesian Methods

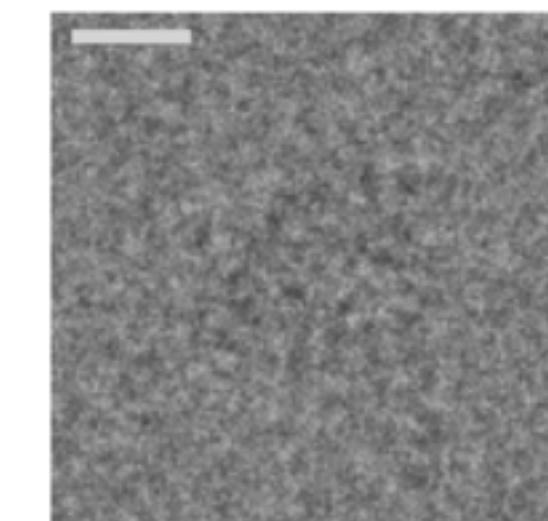
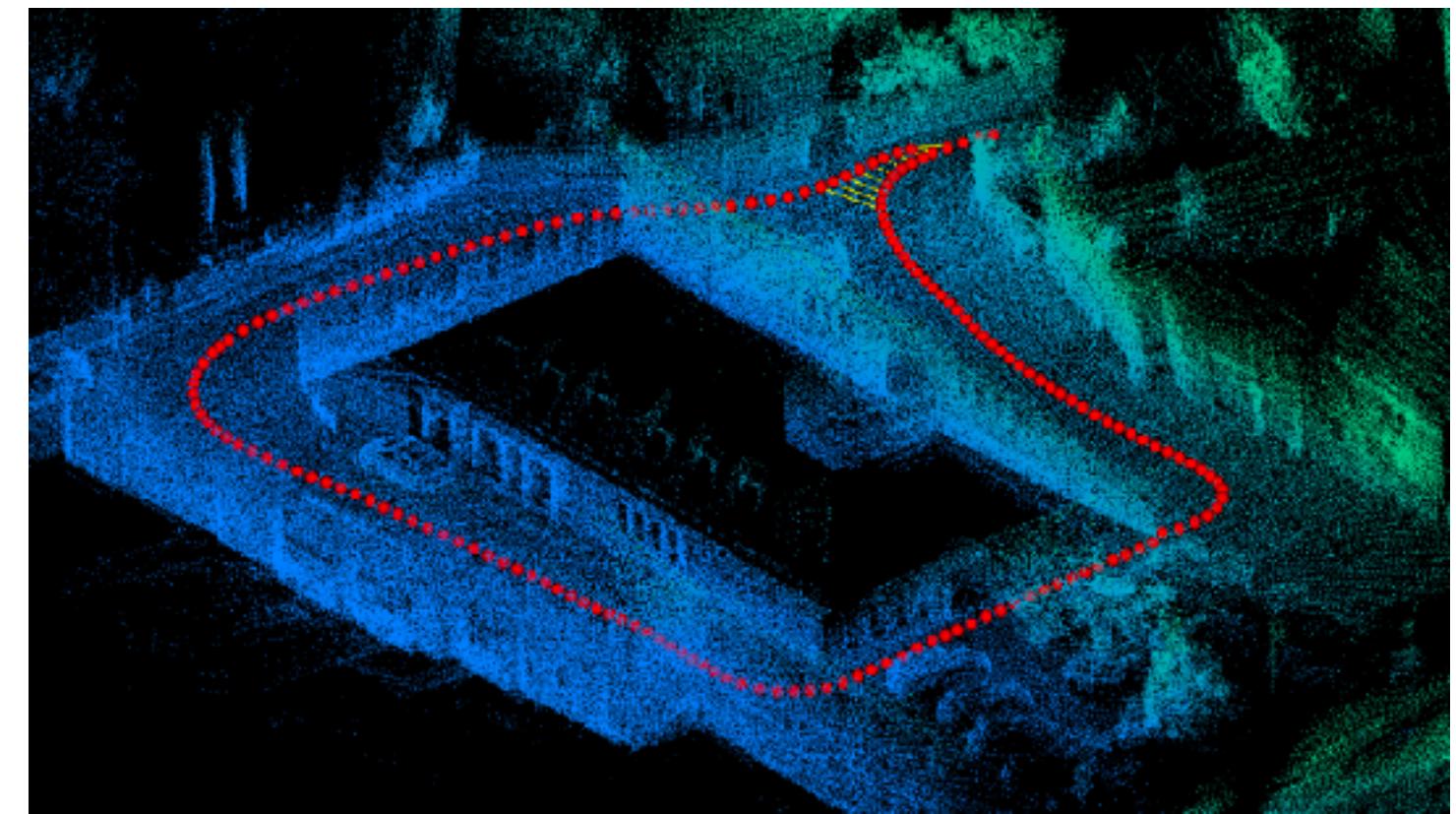
Agenda

- Computer graphics basics
- Neural fields for representing low-dimensional data
- Novel view synthesis

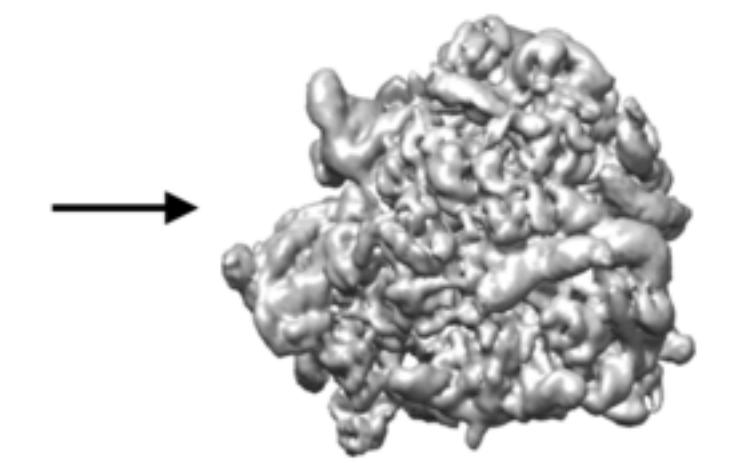
**How to Represent and
Reconstruct the World Around Us?**

Data examples

- Navigation in robotics
 - Mono or stereo cameras
 - Depth sensors, e.g. LiDAR
- Visualisation
 - MRI
 - Electron Microscopy
 - Video conferencing



10^{4-7} projection images



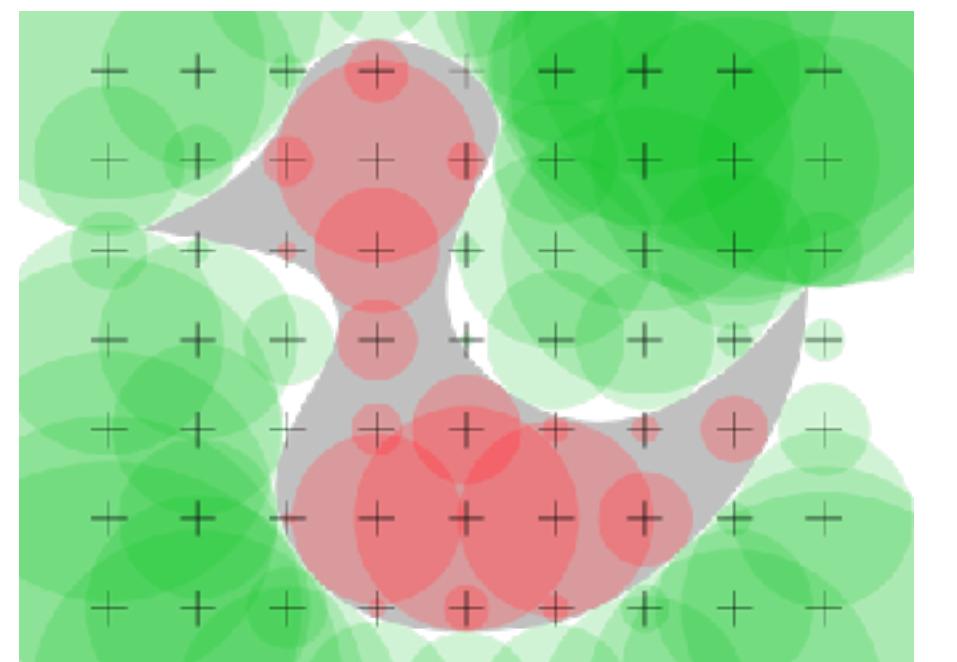
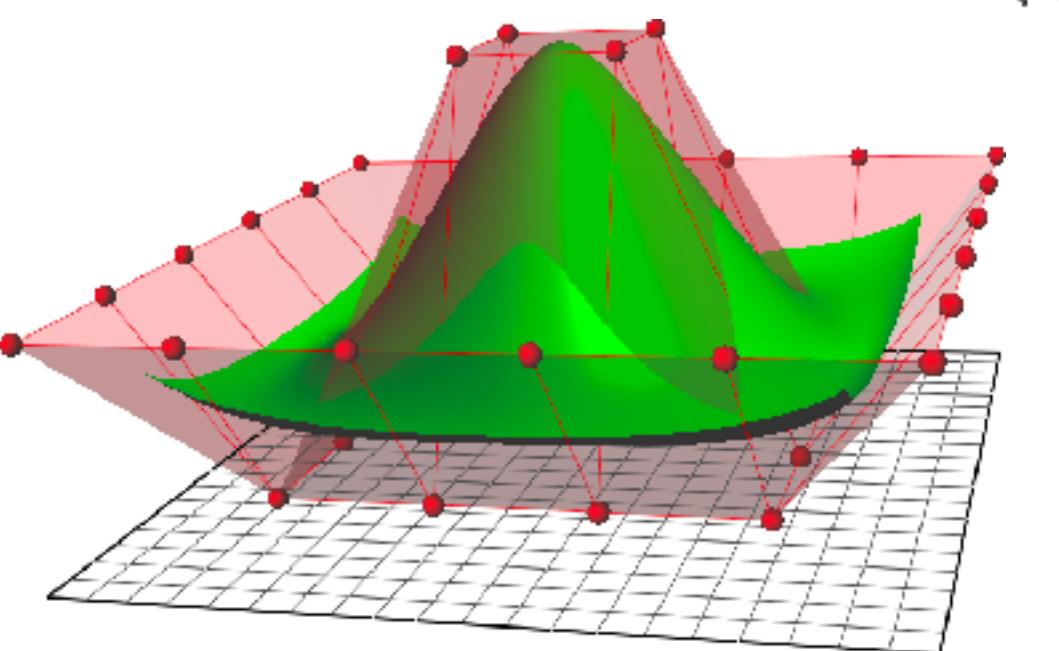
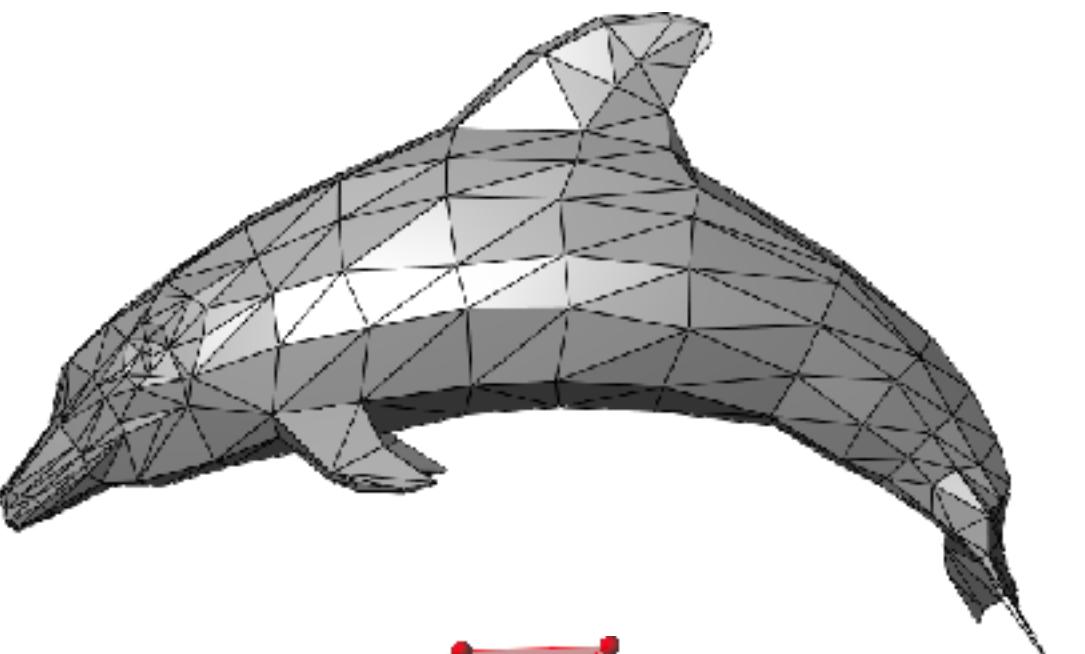
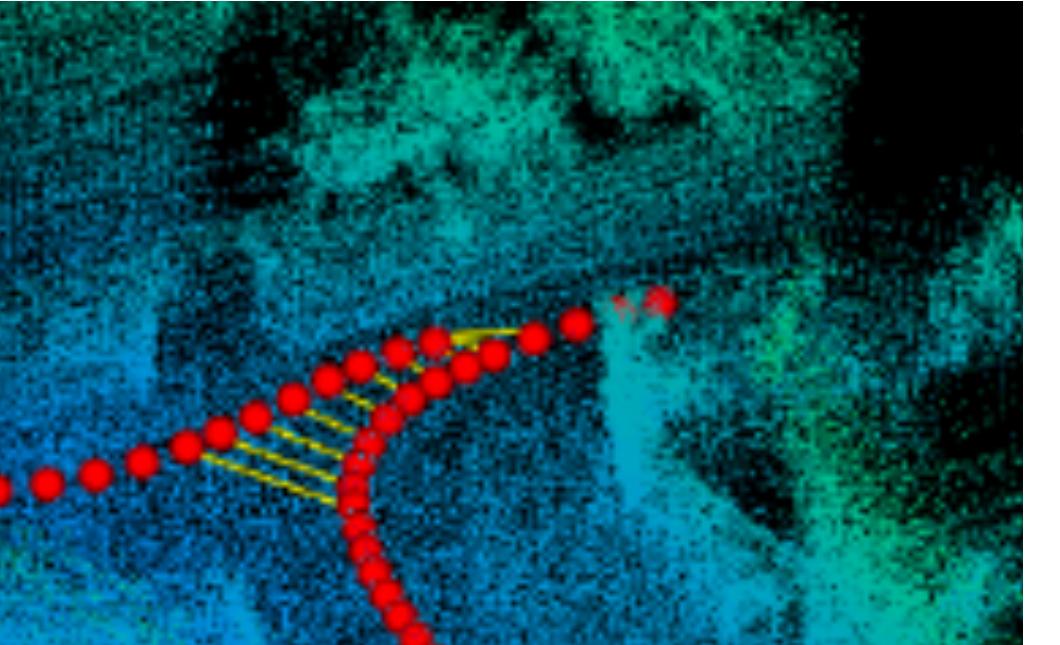
3D electron density



- [https://www.ifp.uni-stuttgart.de/en/research/photogrammetric computer vision/SLAM/](https://www.ifp.uni-stuttgart.de/en/research/photogrammetric_computer_vision/SLAM/)
- <https://www.tesla.com/tesla-gallery>
- https://ru.wikipedia.org/wiki/Беспилотные_автомобили_Яндекса
- Marc Levoy, Efficient Ray Tracing of Volume Data
- Ellen Zhong et al., Reconstructing continuous distributions of 3D protein structure from cryo-EM images
- Project Starline: Feel like you're there, together, <https://www.youtube.com/watch?v=Q13CishCKXY>

Common Representations

- Explicit
 - Point clouds
 - Polygonal meshes
 - Parametric $x = g(t)$
- Implicit
 - Isosurfaces $x : f(x) = 0$



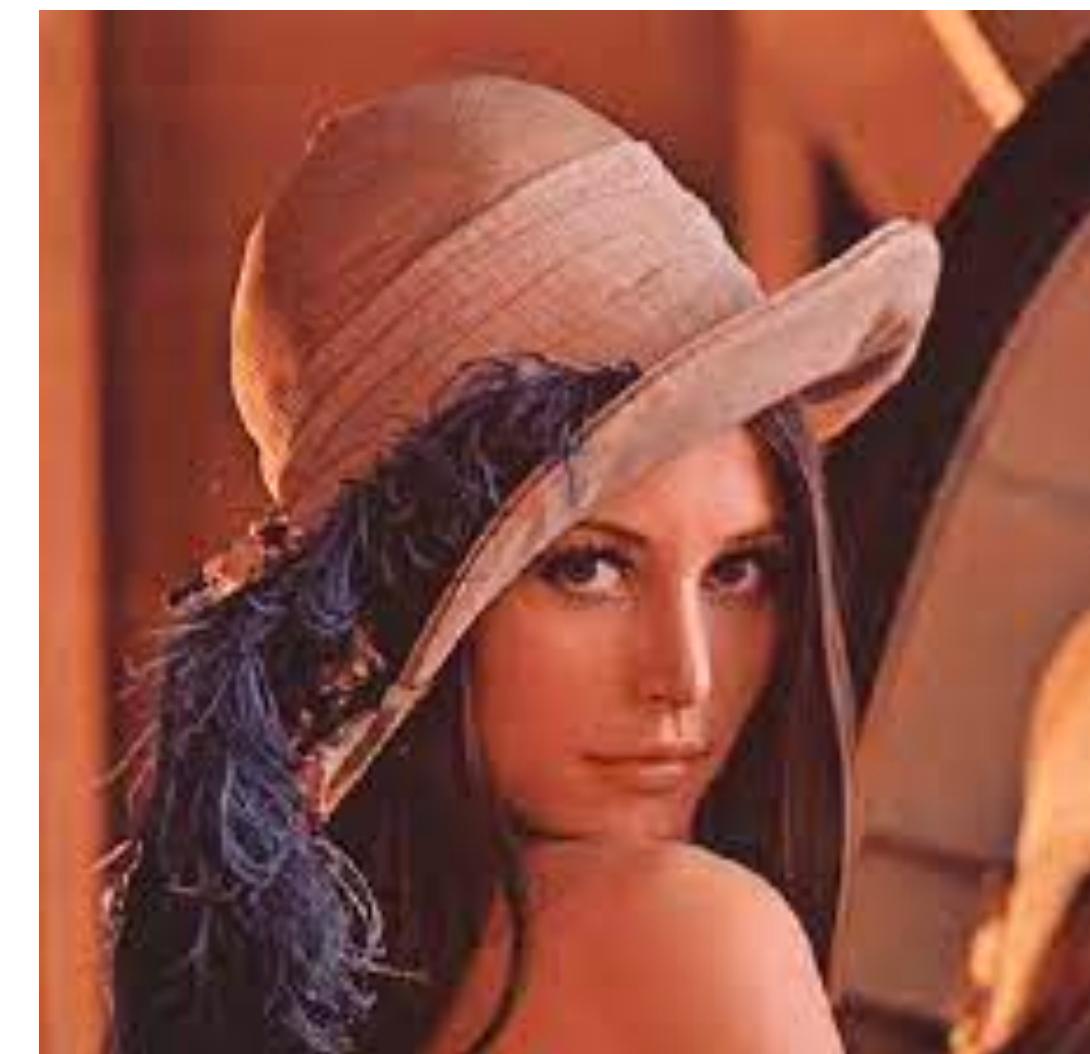
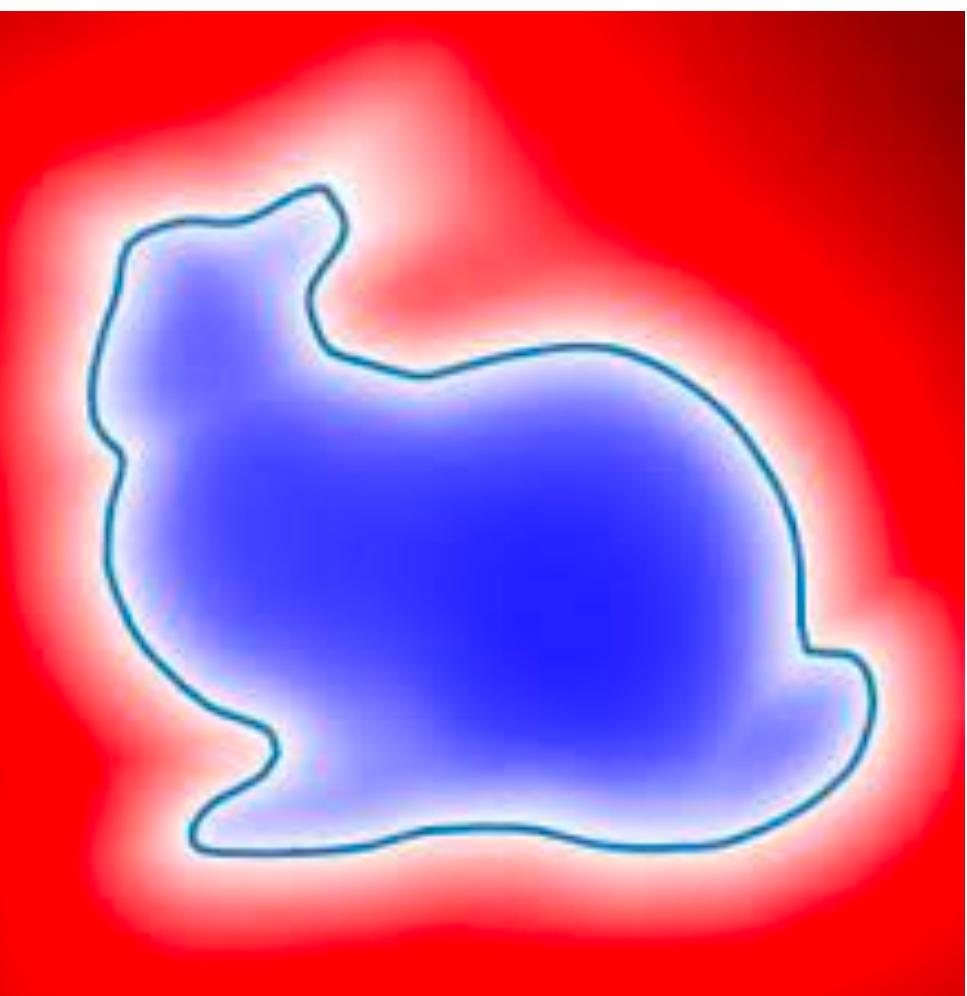
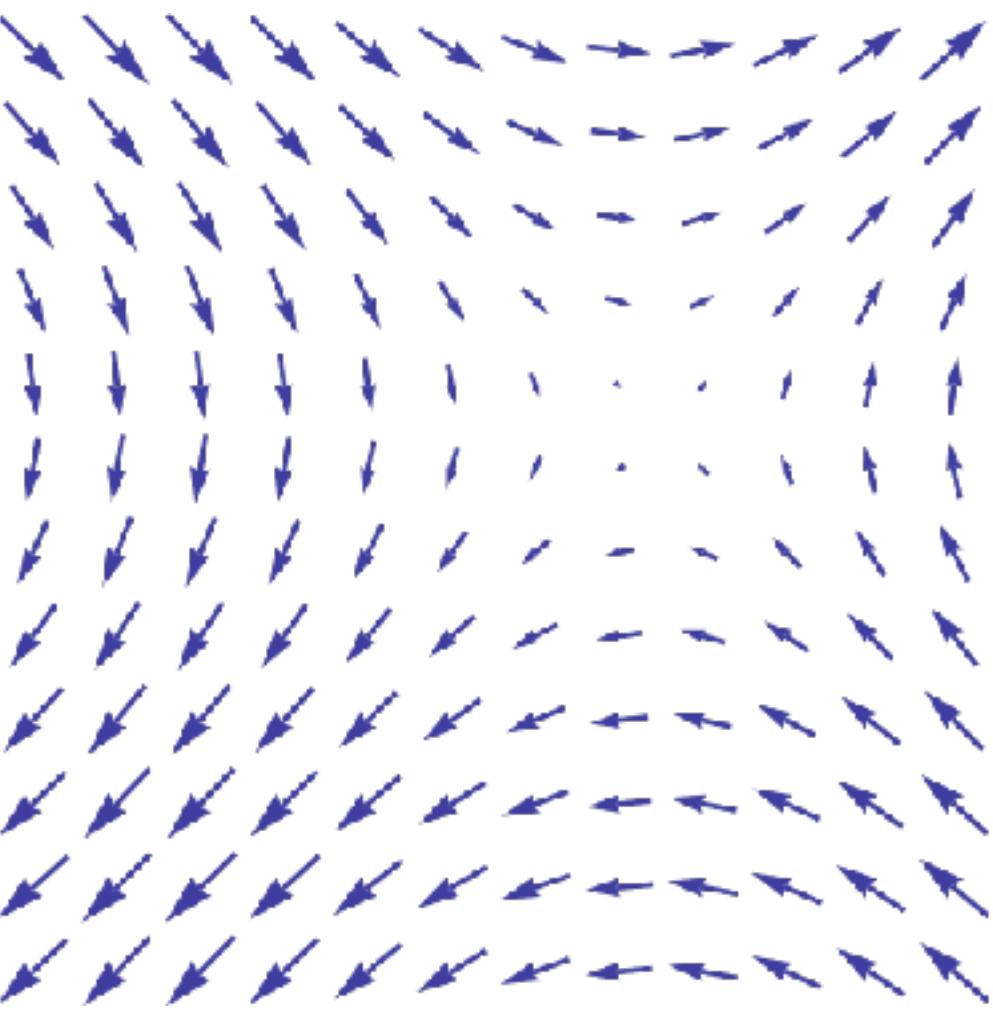
Volumetric Representations

- Surface is a common abstraction in computer graphics
 - Convenience, hardware optimization
 - Sometimes fail to represent real world phenomena
- As an alternative, we will represent scenes with “density” functions
 - Is there something in this point of space?
 - Does light pass through this point in space?
- **Scalar field:** each spatial point stores a scalar value



Vector Fields

- Used in physics and beyond
- Relevant examples
 - Light absorption coefficients
 - Signed Distance Fields
 - Image



<https://en.wikipedia.org/wiki/X-ray>

<https://arxiv.org/abs/1901.05103>

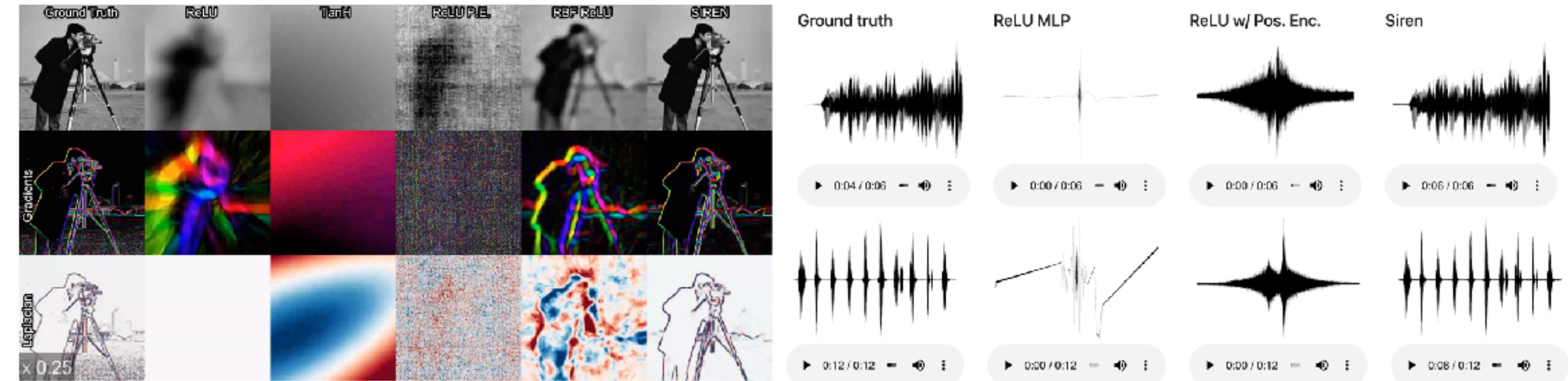
<https://en.wikipedia.org/wiki/Lenna>

Neural Fields

Input: a point in space

Output: vector

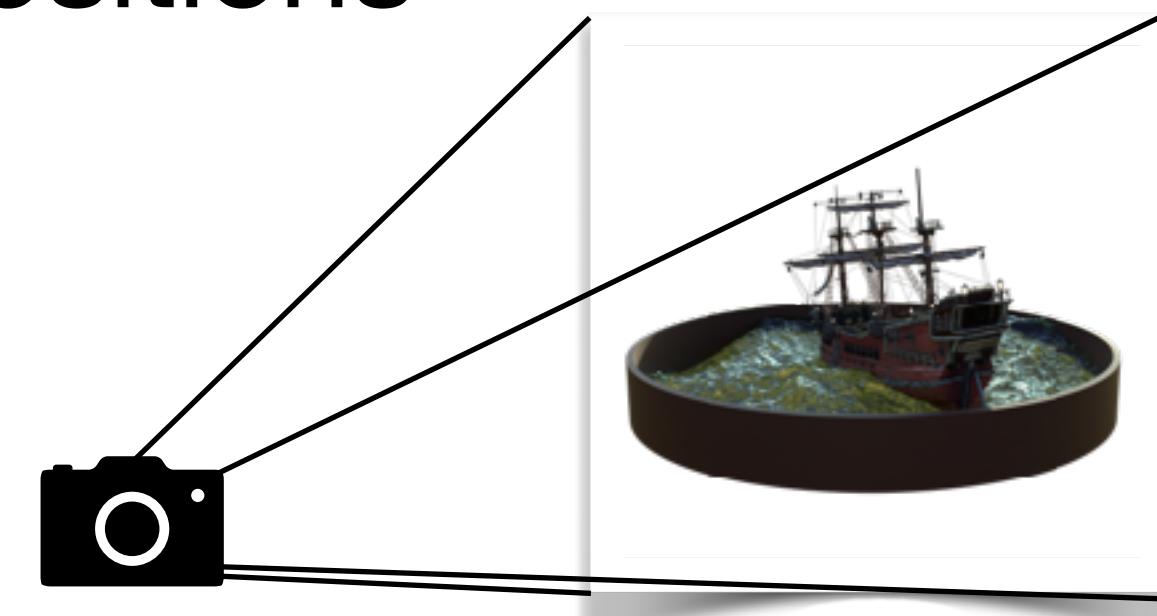
- Architecture: MLP (+ positional embeddings)
- Use cases
 - Photo
 - Video
 - Audio



Neural Radiance Fields

Novel View Synthesis Setup

- Problem setup:
 - A set of pictures taken from a set of pre-determined camera positions
 - Test camera position
- Goal:
 - A picture of a scene taken from the test camera position



Neural Radiance Fields

Representing a Scene with Neural Fields

- Represent a scene with two fields
 - Density: $\sigma(x) : \mathbb{R}^3 \rightarrow \mathbb{R}^+$
 - Radiance: $C(x, d) : \mathbb{R}^3 \times S^2 \rightarrow \mathbb{R}^3$
- Density represents is related to opacity
- Radiance represents color of a point
- Architecture: MLP + positional embeddings



How an Image is Formed?

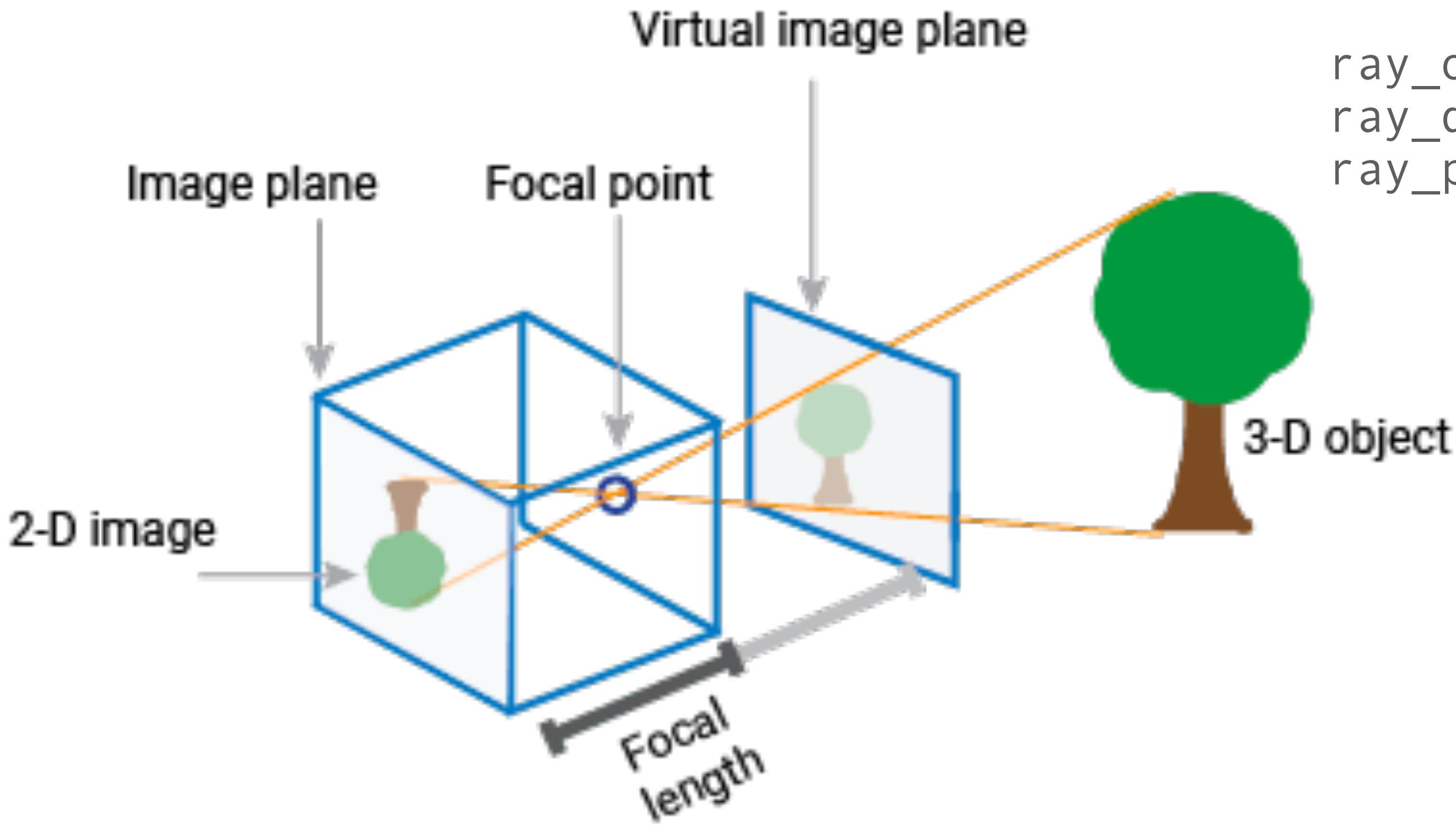
Camera Model

Generate a point on a ray:

$(x, y) \in [-1, 1]^2$ - pixel coordinates

f - focal length

t - distance between the camera and the ray point



```
ray_origin = (0, 0, 0)
```

```
ray_direction = normalize((x, y, f) - ray_origin)
```

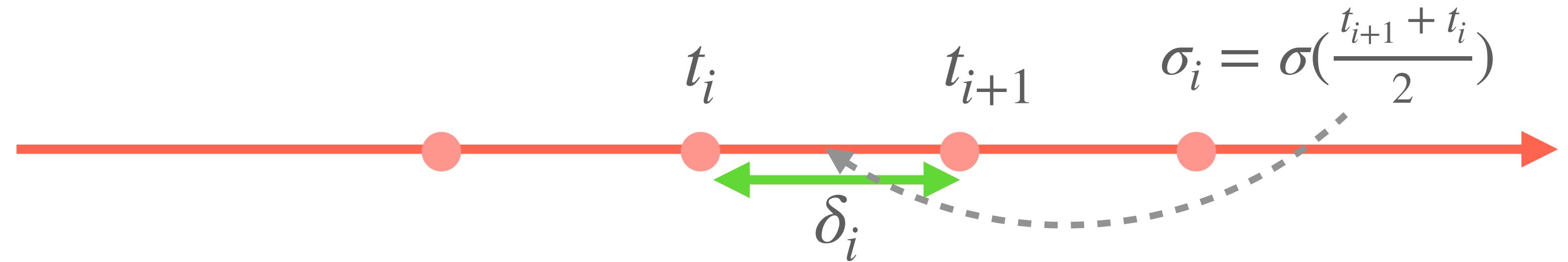
```
ray_point = ray_origin + t * ray_direction
```



Computing Pixel Color

- Density $\sigma(t) \in [0, +\infty)$ is related to opacity at t
- Divide the ray with points t_1, \dots, t_n and define

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



- Expected color along a ray is given by

$$C = \sum_i C(t_i) \cdot \alpha_i \prod_{j < i} (1 - \alpha_j)$$

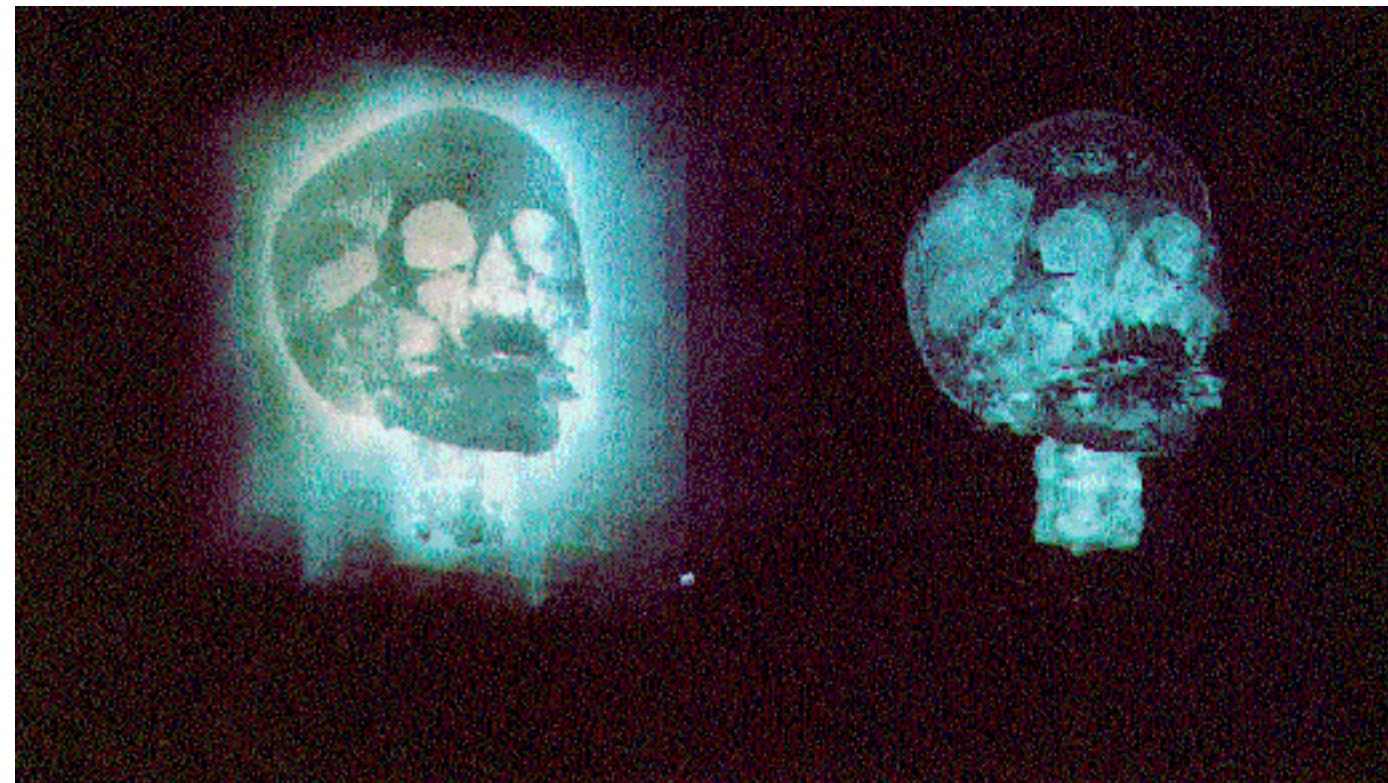


Fig. 11. Costs of rendering Figure 8 using hierarchical enumeration and adaptive termination.

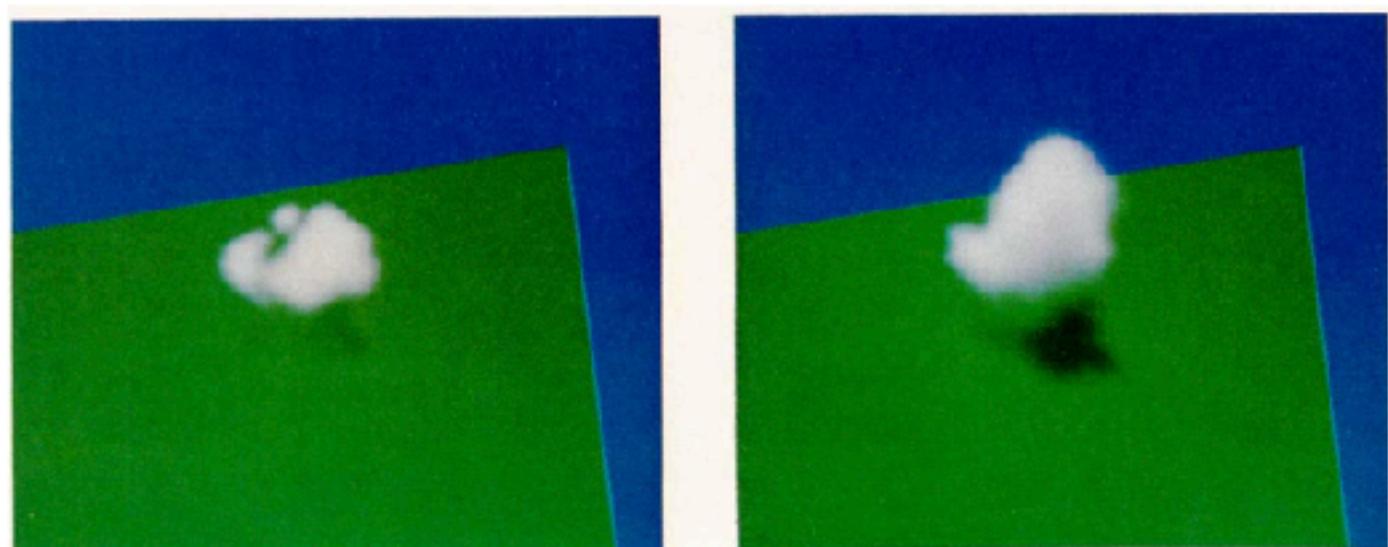


Fig. 5

Fig. 8

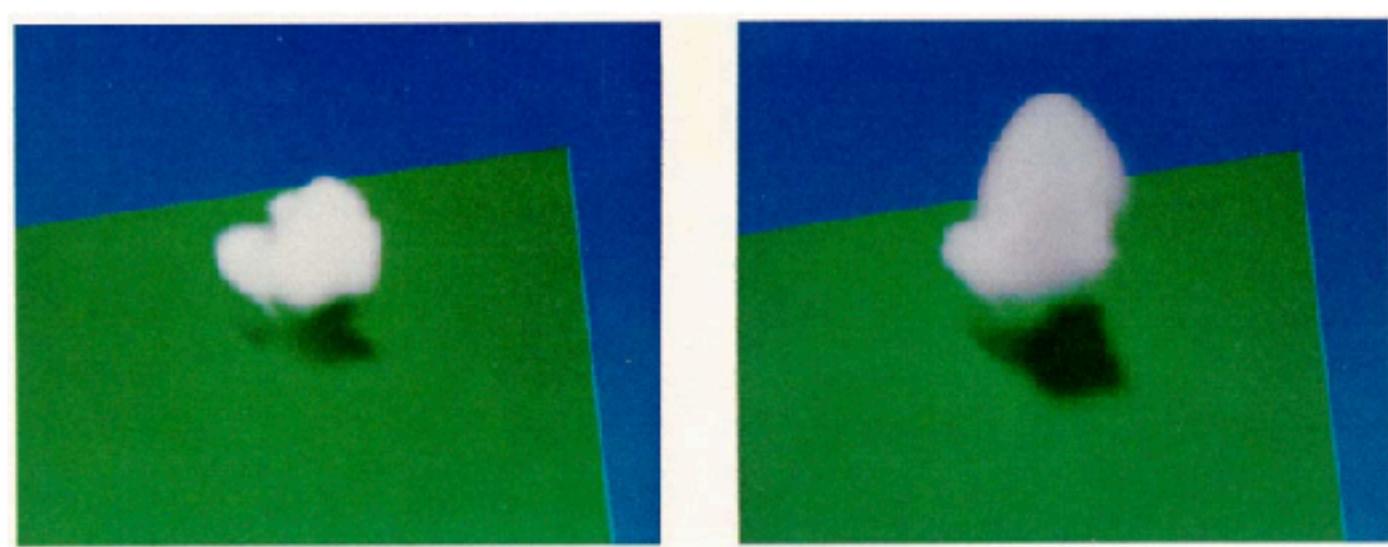


Fig. 6

Fig. 9

Neural Radiance Fields

Optimisation

- Training: minimise the mean squared error $\mathbb{E}_D \|C - C_{gt}\|^2$

$$C = \sum_i C(x_i, d) \alpha_i \prod_{j < i} (1 - \alpha_j)$$

- Optimise w.r.t. parameters of radiance C and density σ

Training Visualisation

- Synthetic data

Coarse iters.: 1
Eps. time: 00:00

Coarse iters.: 1
Eps. time: 00:00

Coarse iters.: 1
Eps. time: 00:00

- Real forward-facing

iters.: 1
Eps. time: 00:00 (>> x5)

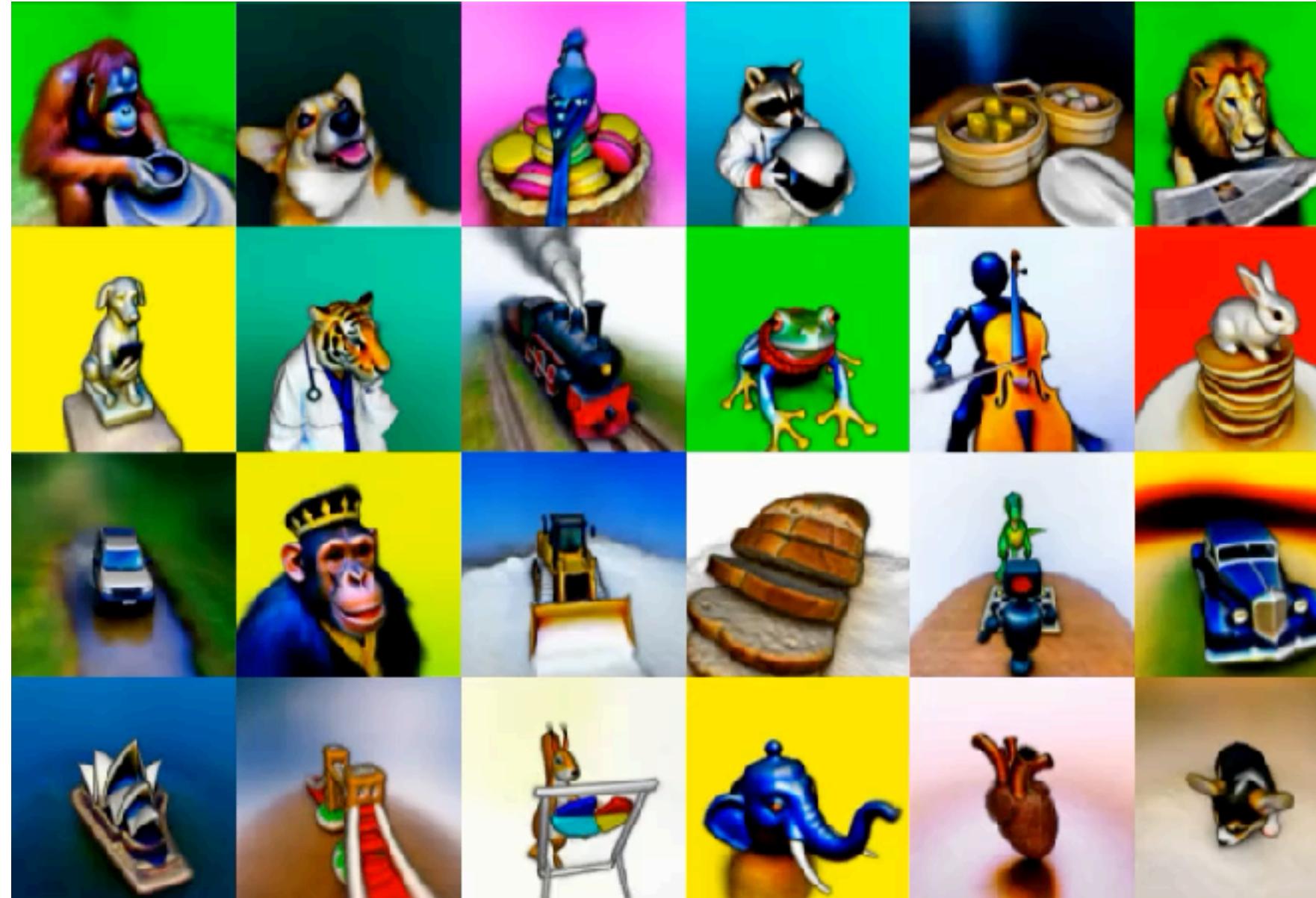
iters.: 1
Eps. time: 00:00 (>> x5)

iters.: 1
Eps. time: 00:00 (>> x5)

Illustrating Model Flexibility

Multimodal Data

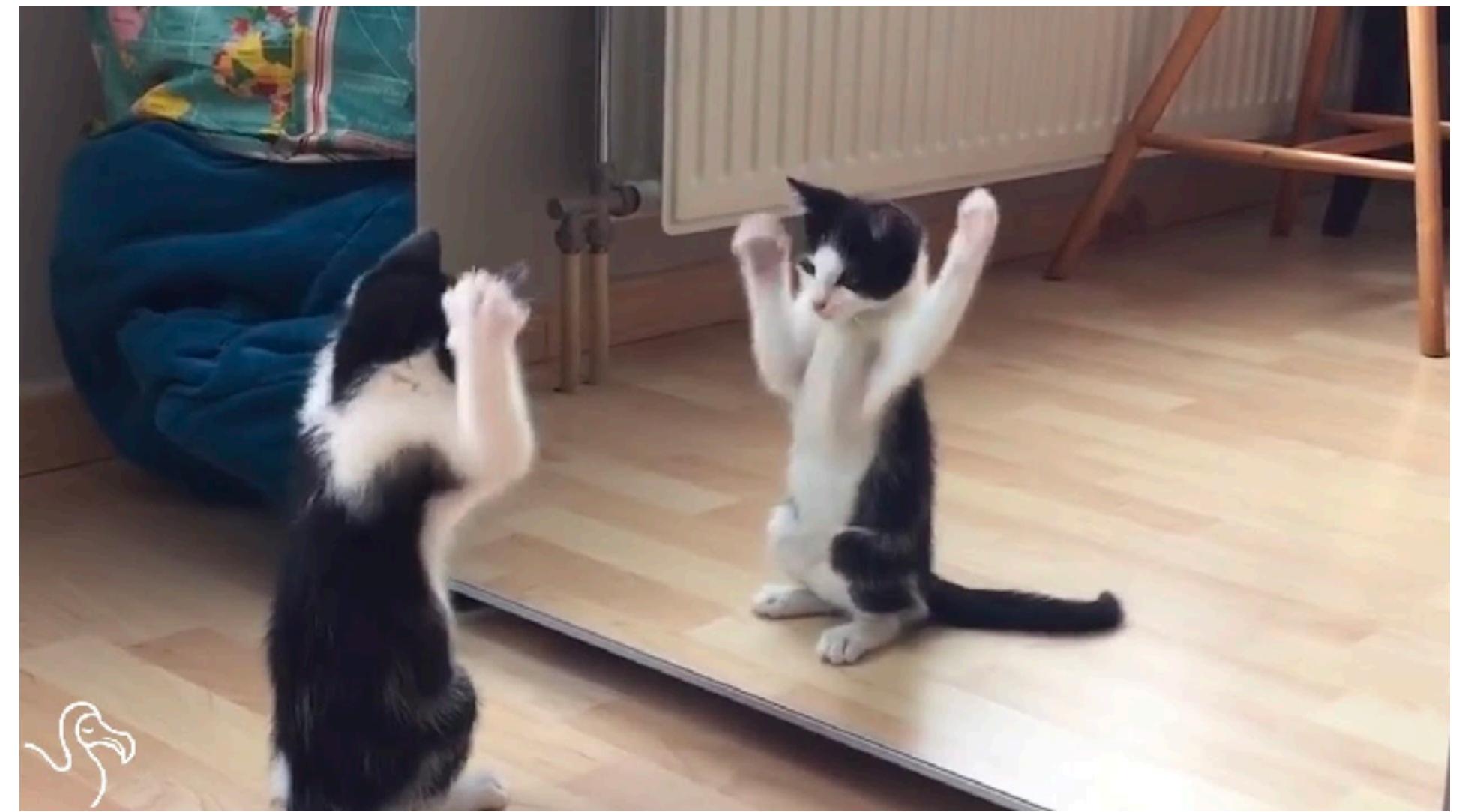
- One can replace MSE reconstruction loss with fancier training signals
 - Leverage text-to-2D diffusion model for scene generation (*second lecture*)
 - CLIP-embeddings for text retrieval and 3D segmentation



Poole B. et al. Dreamfusion: Text-to-3d using 2d diffusion //arXiv preprint arXiv:2209.14988. – 2022.
Kerr J. et al. LERF: Language Embedded Radiance Fields //arXiv preprint arXiv:2303.09553. – 2023.

NeRF pros and cons

- + Simplicity & Flexibility
- + Photorealism
- + Compression (*5mb / scene*)
- Ill-posed problem
- Does not work on real scenes
- Long training time (*48 h. / scene*)
- Slow rendering (*1 min. / frame*)



Improving Rendering Speed

Sparsification

- Each terms in the sum requires running an MLP

$$C = \sum_i C(x_i) \alpha_i \prod_{j < i} (1 - \alpha_j)$$

- If $\sigma(x_i) = 0$, then $1 - \alpha_i = 0$, we can omit i -th term
- Idea: cache $\sigma(\cdot)$, omit i if we know $\sigma(x_i) = 0$
- Outcome: training and rendering 10-30 times faster

Liu L. et al. Neural sparse voxel fields //Advances in Neural Information Processing Systems. – 2020. – T. 33. – C. 15651-15663.

Hedman P. et al. Baking neural radiance fields for real-time view synthesis //Proceedings of the IEEE/CVF International Conference on Computer Vision. – 2021. – C. 5875-5884.

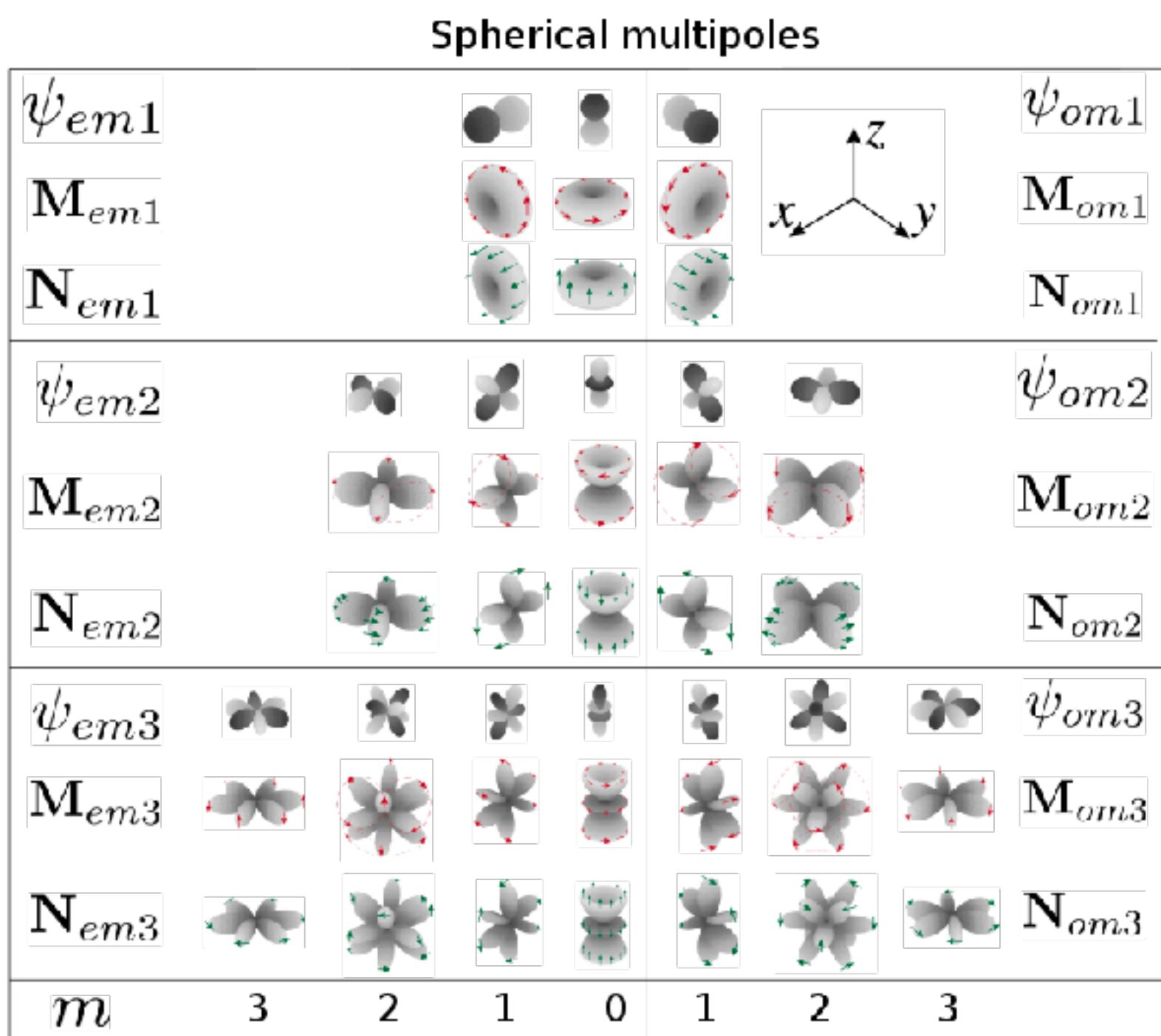
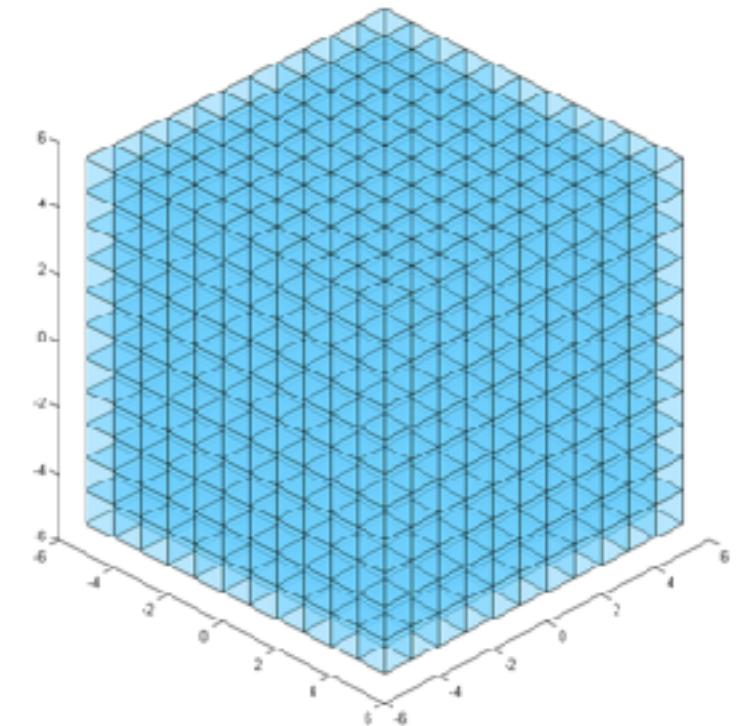
Sun C., Sun M., Chen H. T. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction //Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. – 2022. – C. 5459-5469.

Fridovich-Keil S. et al. Plenoxels: Radiance fields without neural networks //Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. – 2022. – C. 5501-5510.

Li R., Tancik M., Kanazawa A. NerfAcc: A General NeRF Acceleration Toolbox //arXiv preprint arXiv:2210.04847. – 2022.

Volume Grid Representations

- There is no magic in using neural nets
 - We can represent $\sigma(x) : \mathbb{R}^3 \rightarrow \mathbb{R}^+$ with a volumetric grid
 - Radiance $C(x, d) : \mathbb{R}^3 \times S^2 \rightarrow \mathbb{R}^3$ involves 5D grid
- Spherical harmonics - basis for functions on spheres
- Each spatial point x stores basis coordinates $C(x, \cdot)$
- Training time: minutes
- Real-time rendering



Hybrid Approach

- There is no magic in using neural nets
 - We can represent $\sigma(x) : \mathbb{R}^3 \rightarrow \mathbb{R}^+$ with a volumetric grid
 - Radiance $C(x, d) : \mathbb{R}^3 \times S^2 \rightarrow \mathbb{R}^3$ involves 5D grid
- Hybrid representation $C(x, d) = F(G(x), d)$:
 - $G(x)$ - vector field, represented with a voxel grid
 - $F(v, d)$ - tiny MLP
- Leads to higher fidelity compared to spherical harmonics

Instant-NGP

Pros and cons of the above solutions

Is there a solution compromising between the two?

MLP

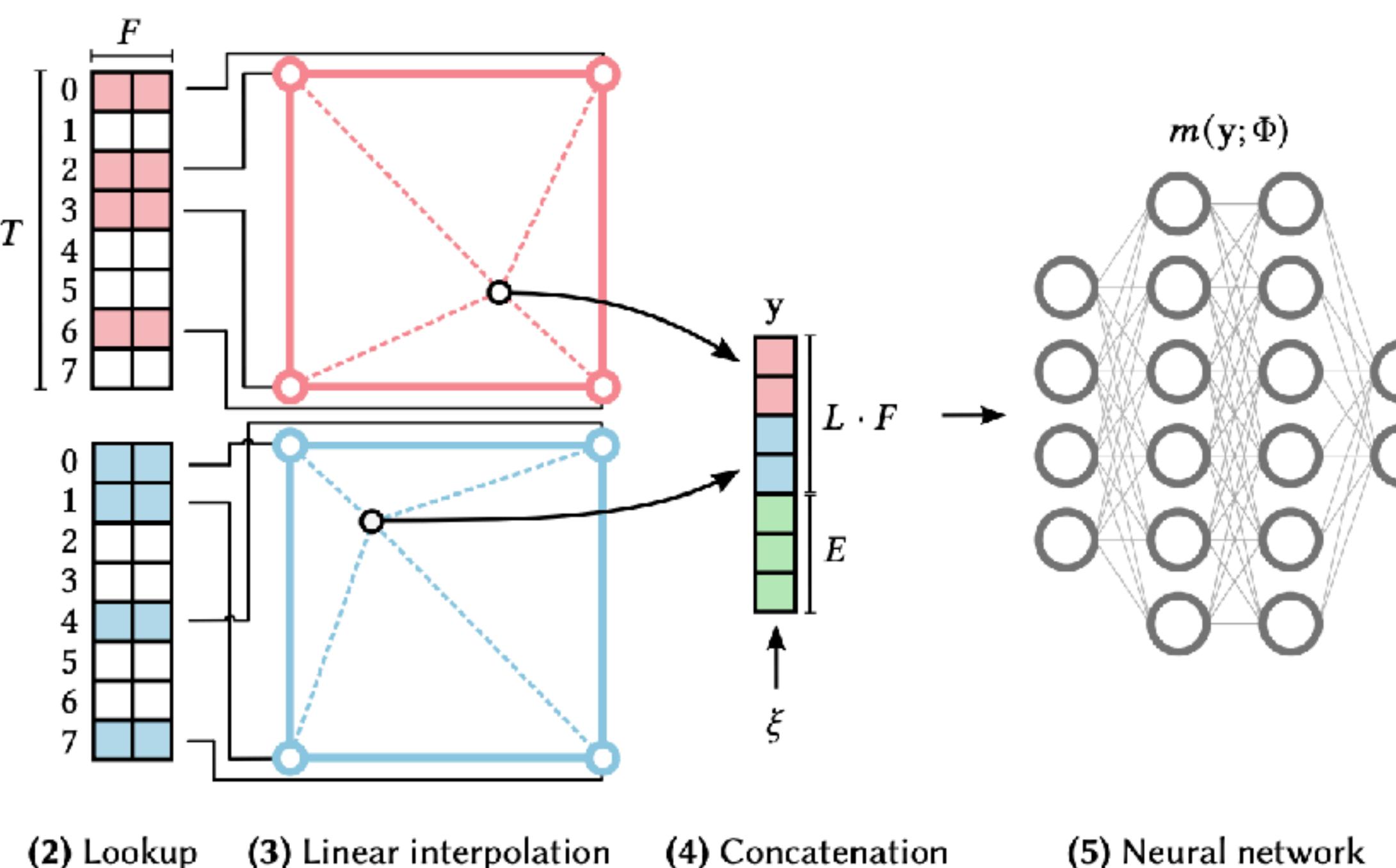
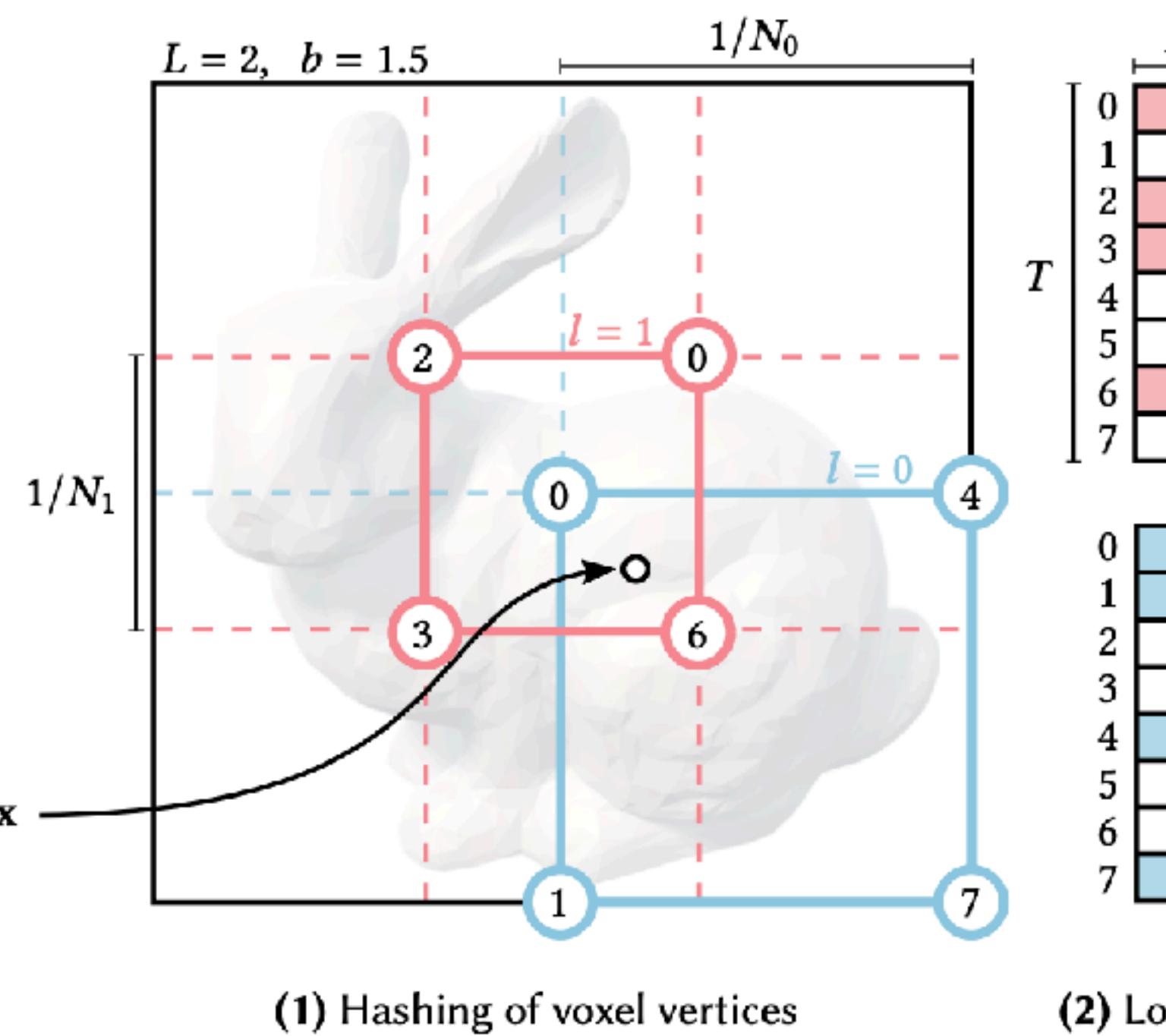
- Size: ~5mb
- Slow
- Trade-off speed for fidelity

Voxel Grid

- Size: ~1Gb
- Fast
- Trade-off memory for fidelity

Instant-NGP

Saving time and memory with HashGrids

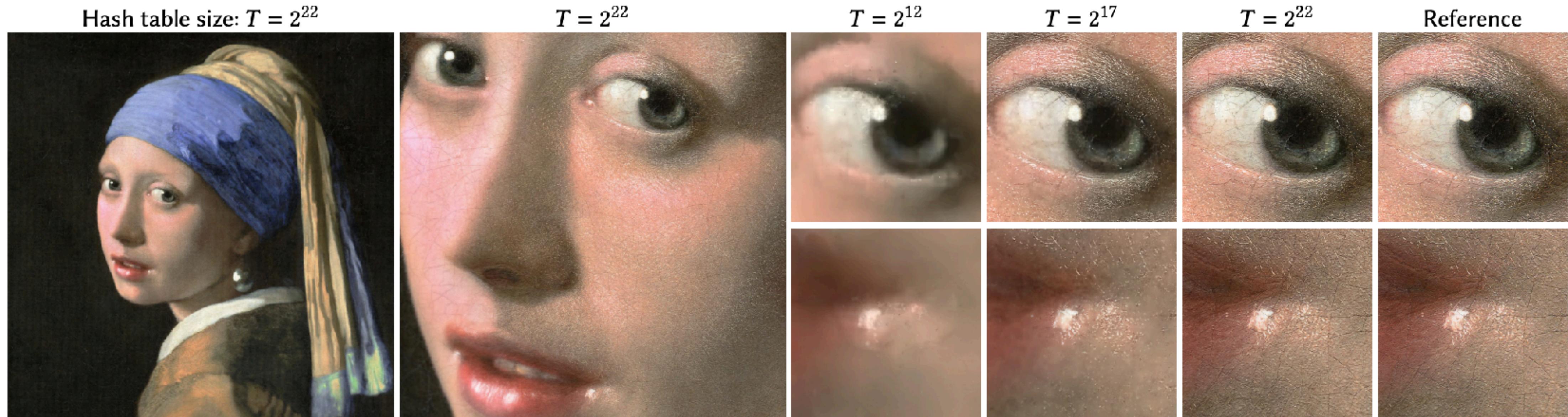


HashGrid

- Size $\sim 100\text{mb}$
- Average speed
- Flexibility factors
 - HashGrid size
 - MLP component

Neural Graphics Primitives

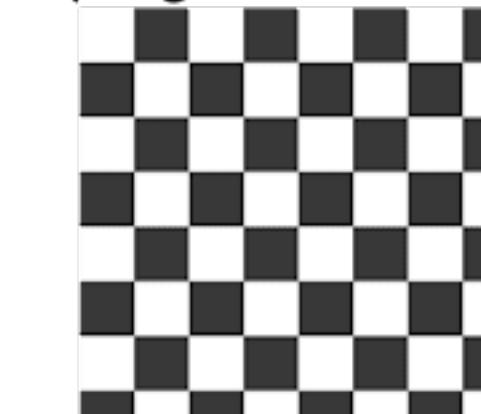
- Network size is 3.4% of the full image size (4×10^8 pixels)
- PSNR 29.8



Improving Reconstruction Fidelity

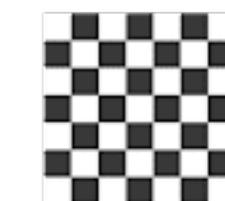
Aliasing

Mip 0
(original texture)



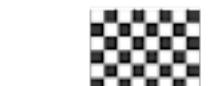
128 x 128

Mip 1



64 x 64

Mip 2



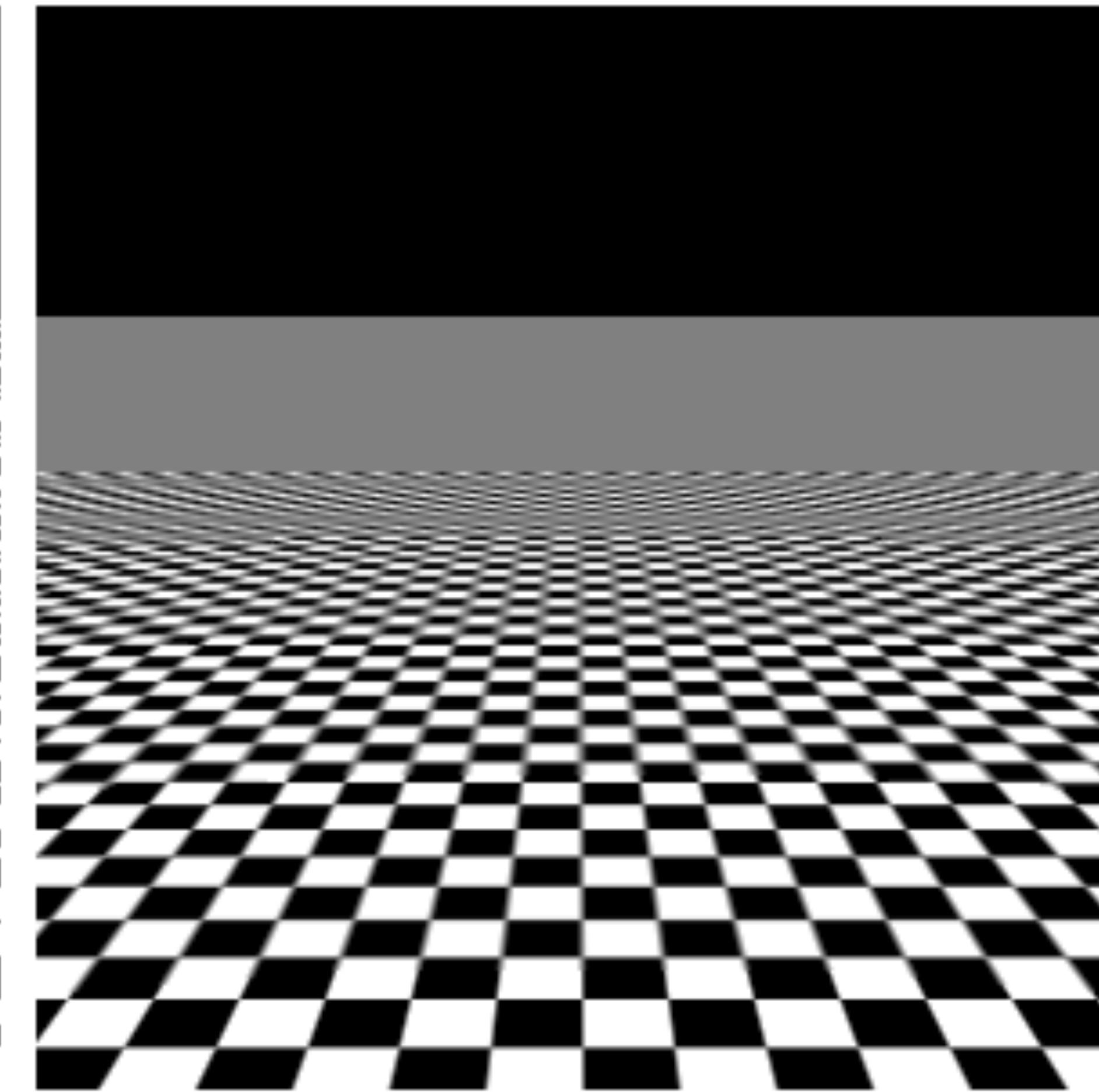
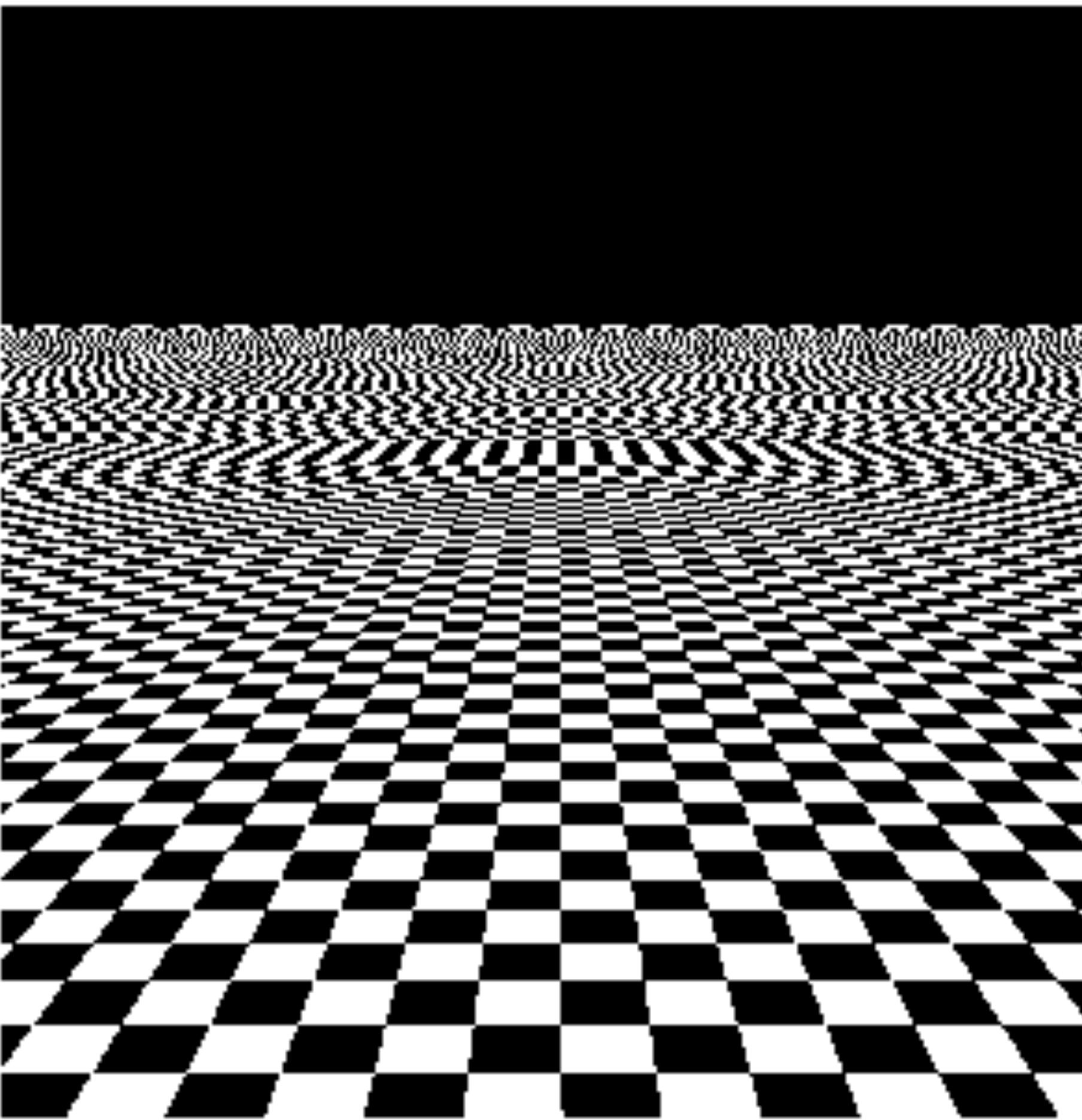
32 x 32

Mip 3



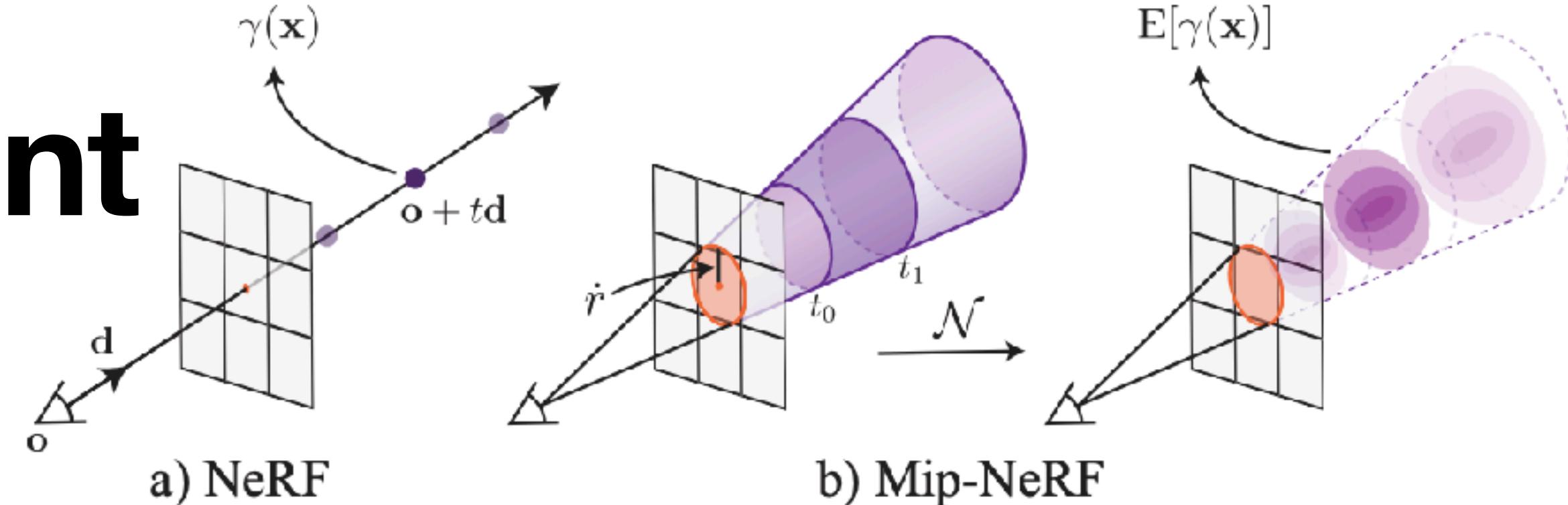
16 x 16

multum in parvo



Taking Scale into Account

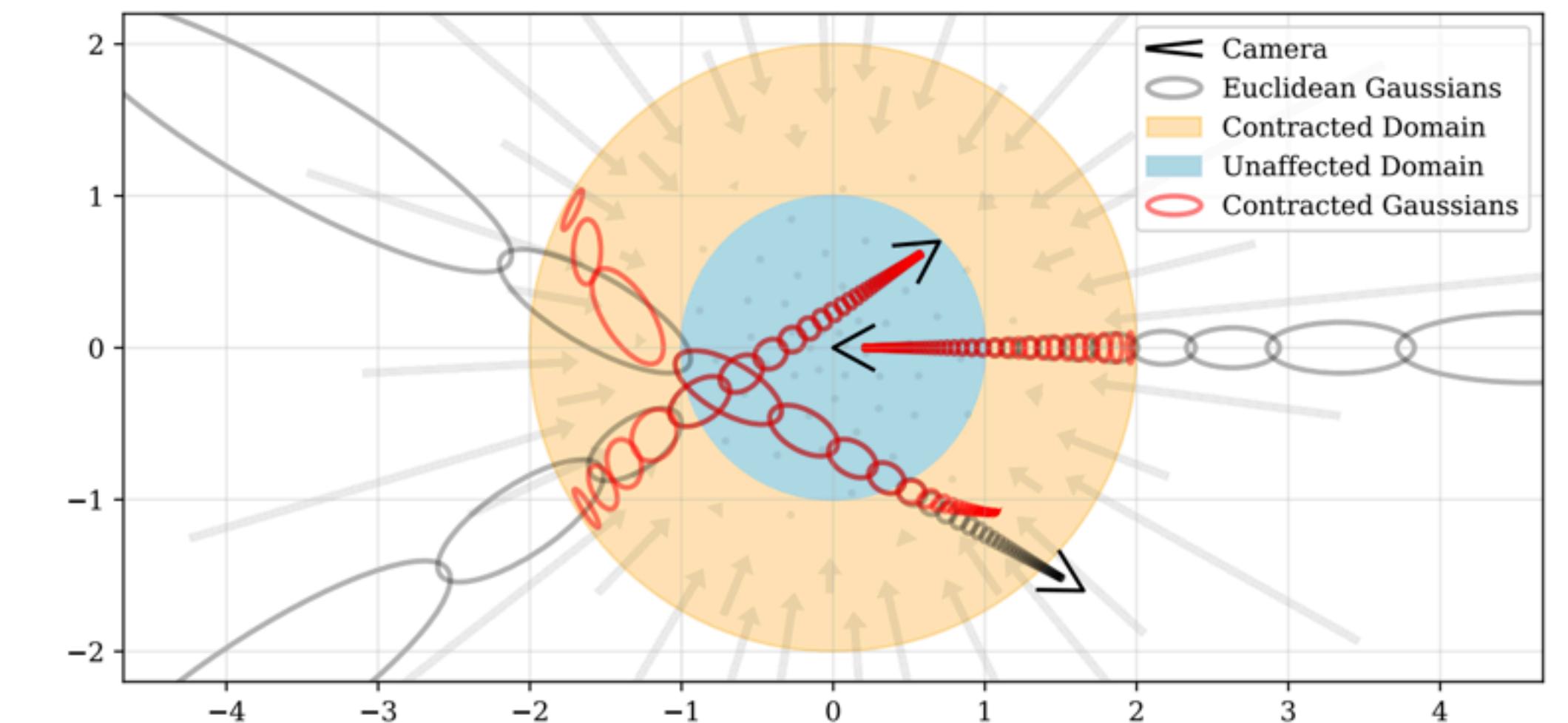
A Parallel Line of Research



- Original NeRF experiments had biased data
 - Camera distances did not change by a significant amount
 - Pixels cover a cone-like region in space
 - Ideal solution - compute the average radiance over the whole cone
 - Mip-NeRF approximates the average radiance/density over a cone segment
 - Idea: compute embedding average rather than output average

Unbounded Scenes

- Original NeRF experiments had biased data
 - All the training scenes were bounded
- NeRF++ models foreground and background with two networks
- Mip-NeRF-360 maps the space into a ball
 - Does not transform scene center
 - The rest is mapped onto a shell



Zip-NeRF

- Merges two research branches
 - Speed: Instant-NGP
 - Fidelity: Mip-NeRF-360
- Mip-NeRF is designed for MLPs
- Zip-NeRF adapts Mip-maps to Instant-NGP



Further References

- Useful tools:
 - Open source solution for training NeRFs <https://docs.nerf.studio/>
 - “There’s an app for that” <https://lumalabs.ai/>
- Hot topics beyond the scope of lecture:
 - NeRFs for video <https://dynibar.github.io/>
 - Better surface reconstruction: <https://research.nvidia.com/labs/dir/neuralangelo/>
 - Gaussian Splatting: recent NeRF alternative <https://gsplat.tech/>