

Lecture 17:

Learning Scene Representations

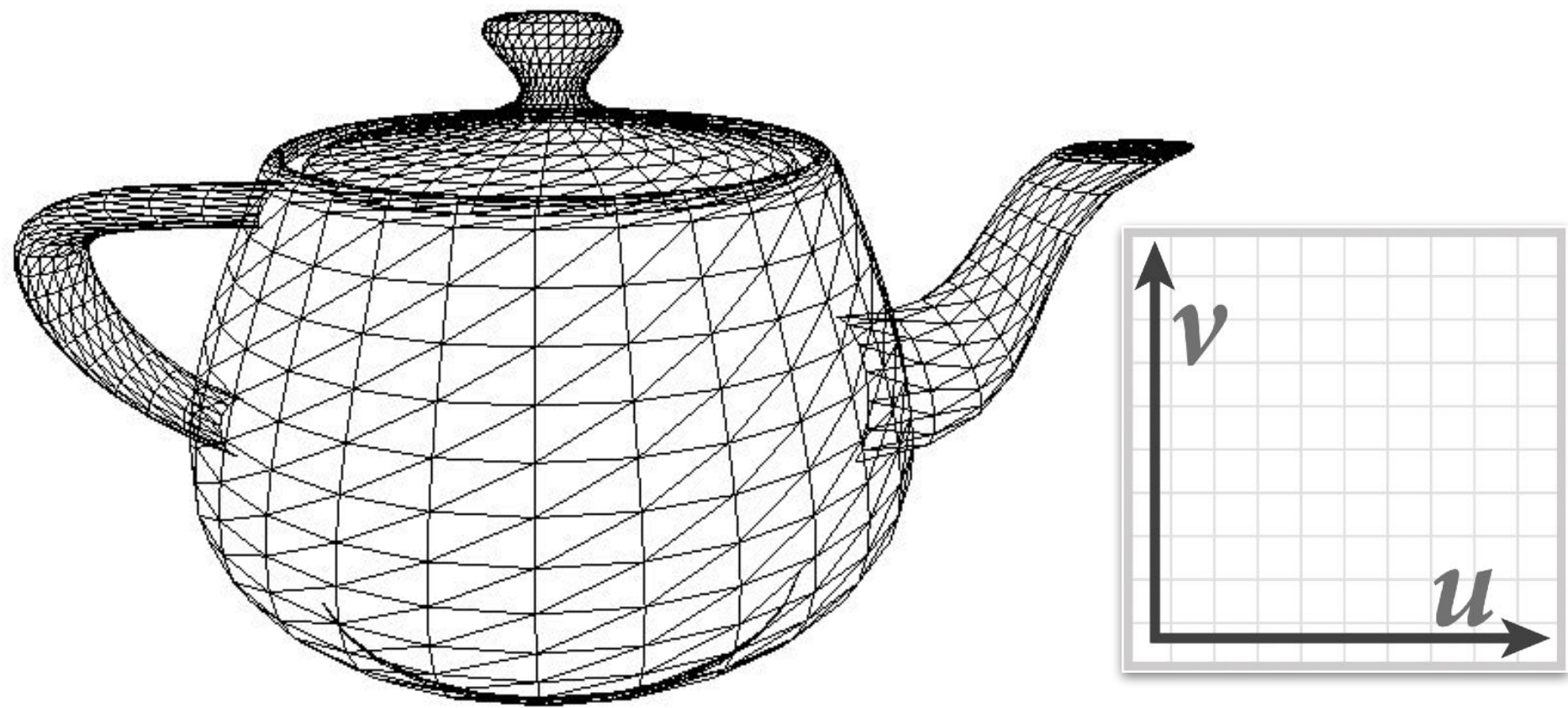
Visual Computing Systems
Stanford CS348K, Spring 2024

**Computer science in a nutshell:
Choose the right representation for the task at hand**

Today's task

- **Recovering a 3D scene representation from a collection of photos**
- **Why?**
 - **So we can render them from novel viewpoints**
 - **So we can perform editing**
 - **Geometric edits vs. material edits vs. lighting edits**
 - **To aid interpretation of their contents**

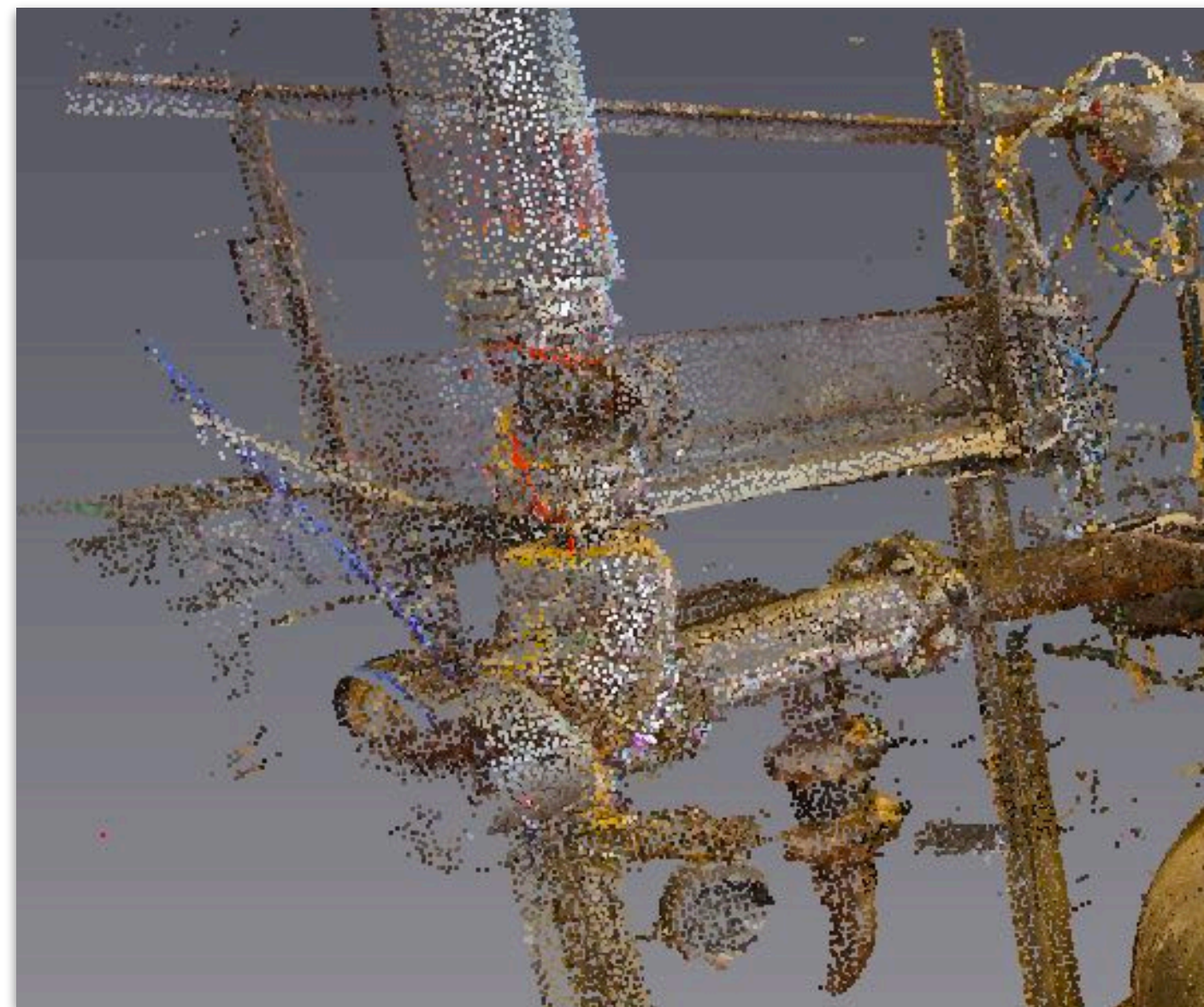
Many scene representations



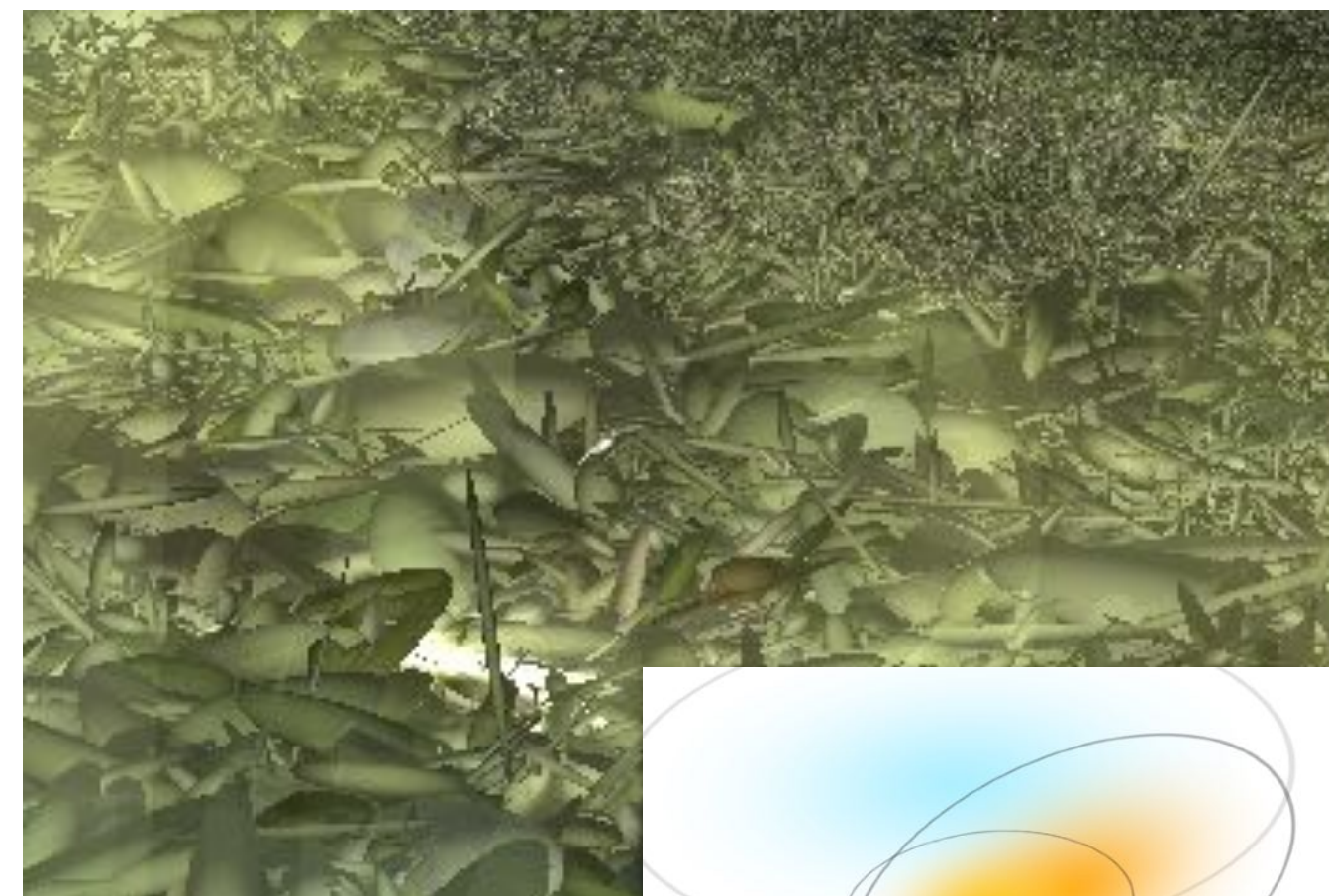
3D triangle mesh + texture map



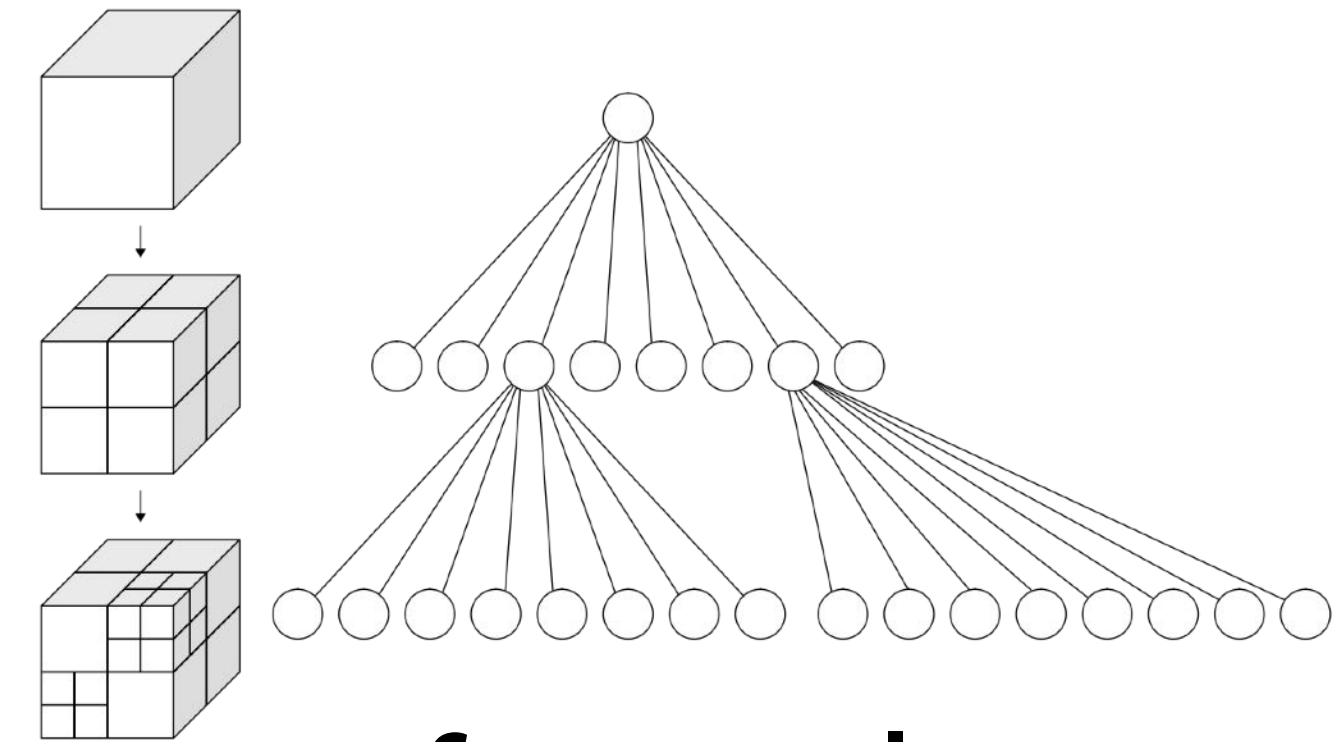
3D volume (voxels)



Point cloud (list of points)

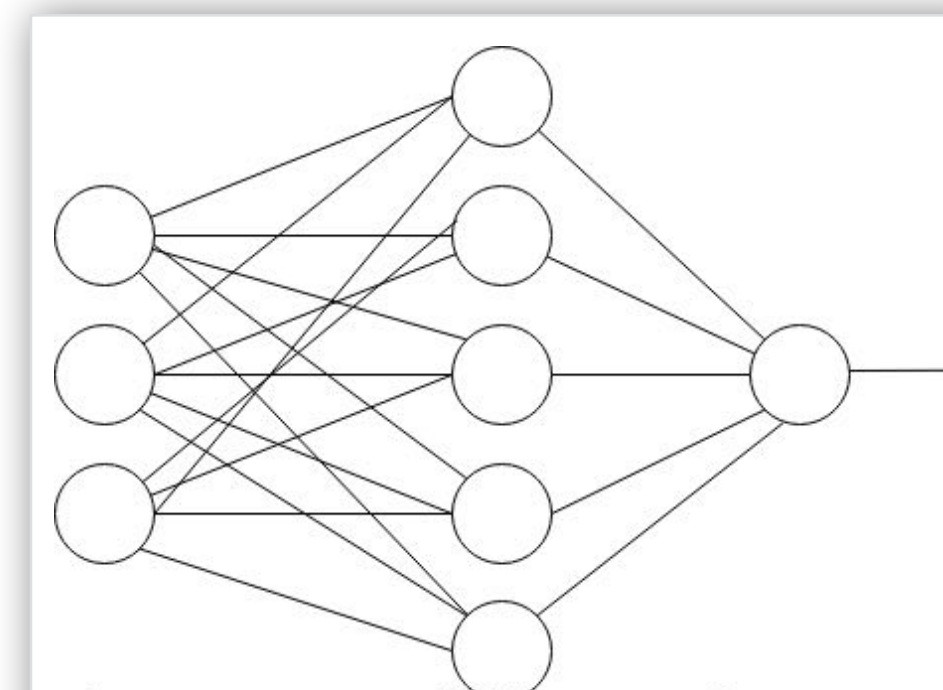


Oriented 3D Gaussians



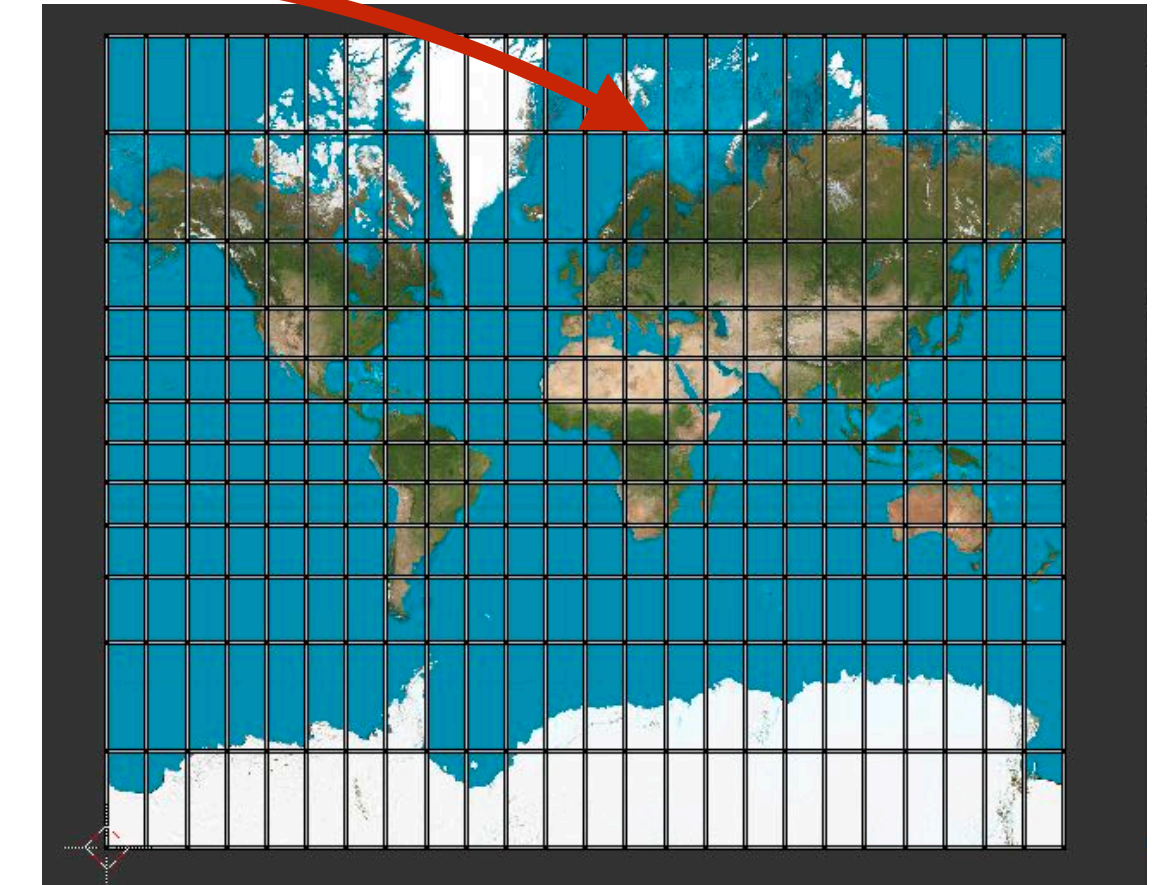
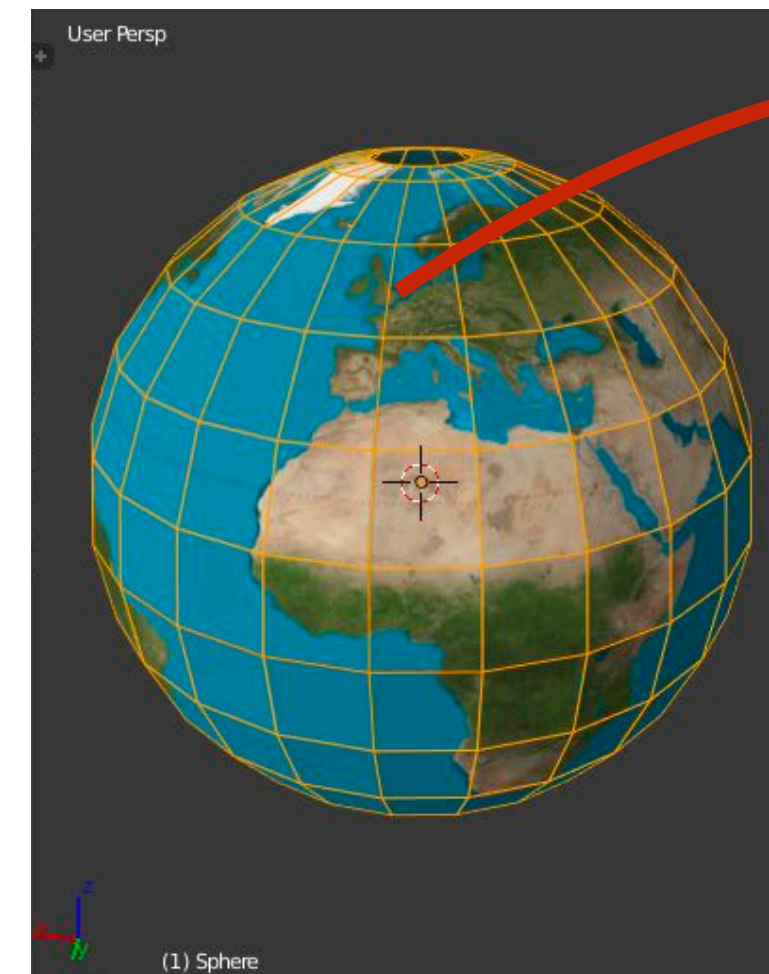
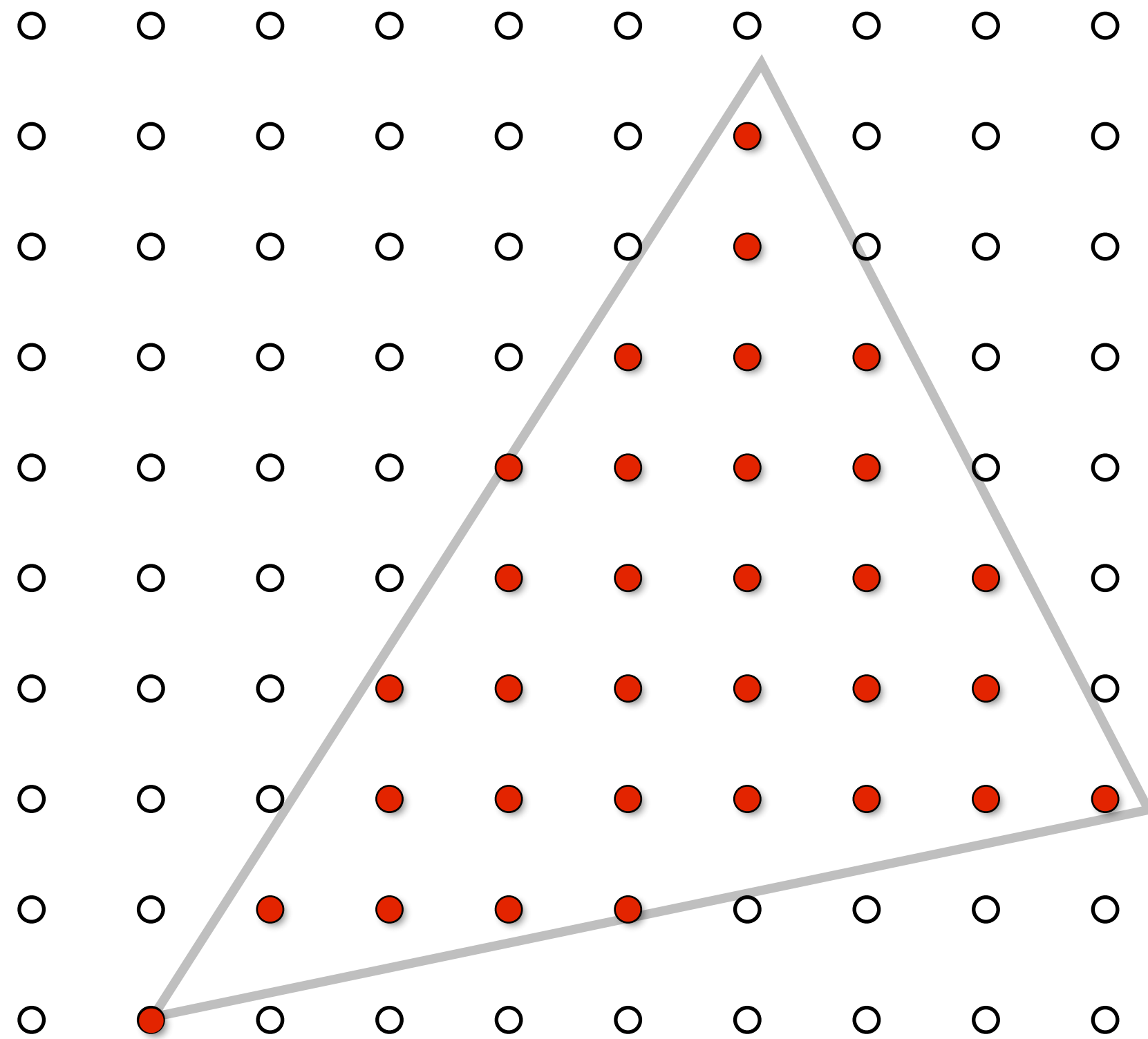
Sparse voxels

DNN (MLP)

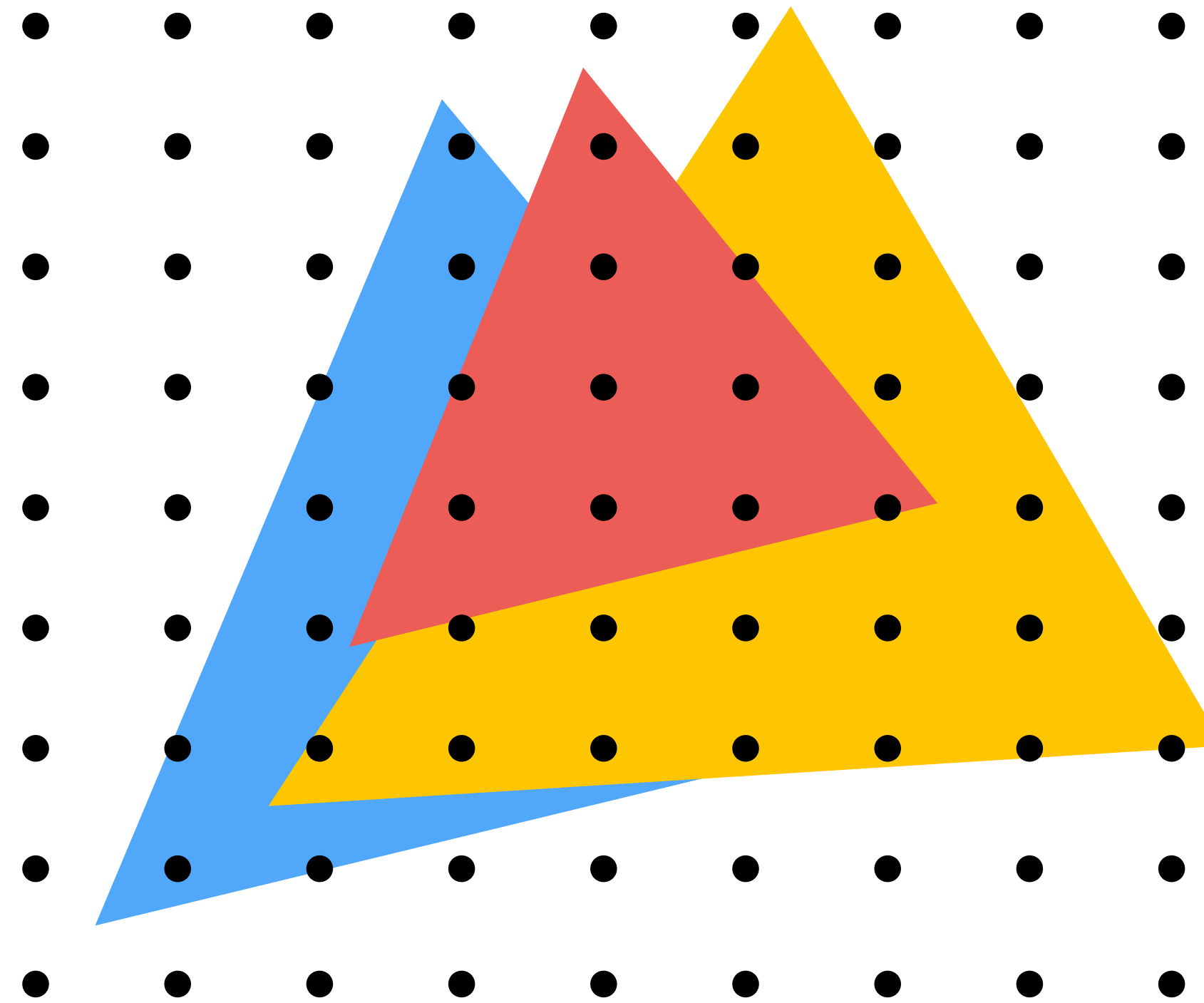


Rendering triangles

Given camera position and 3D position of vertices: (1) project vertices onto screen (2) color pixels within 2D triangle



Example: rendering three opaque triangles



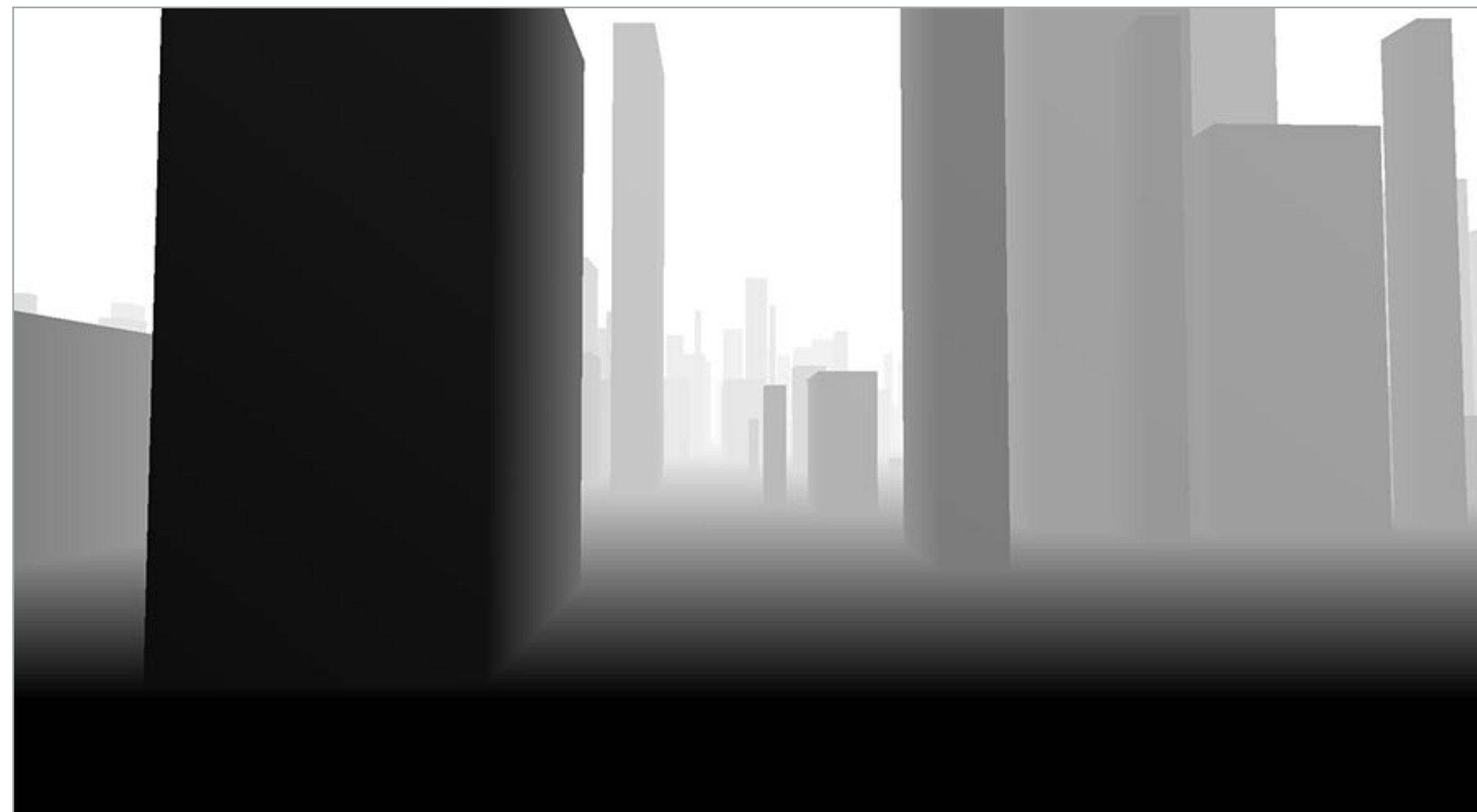
Depth buffer (aka “Z buffer”)

Color buffer:
(stores color per sample... e.g., RGB)



Depth buffer:
(stores depth per sample)

Stores depth of closest surface drawn so far
black = close depth
white = far depth



Occlusion using the depth buffer (opaque surfaces)

```
bool pass_depth_test(d1, d2) {  
    return d1 < d2;  
}
```

```
depth_test(tri_d, tri_color, x, y) {  
  
    if (pass_depth_test(tri_d, depth_buffer[x][y]) {  
  
        // if triangle is closest object seen so far at this  
        // sample point. Update depth and color buffers.  
  
        depth_buffer[x][y] = tri_d;    // update depth_buffer  
        color[x][y] = tri_color;      // update color buffer  
    }  
}
```


Basic rasterization algorithm

Sample = 2D point

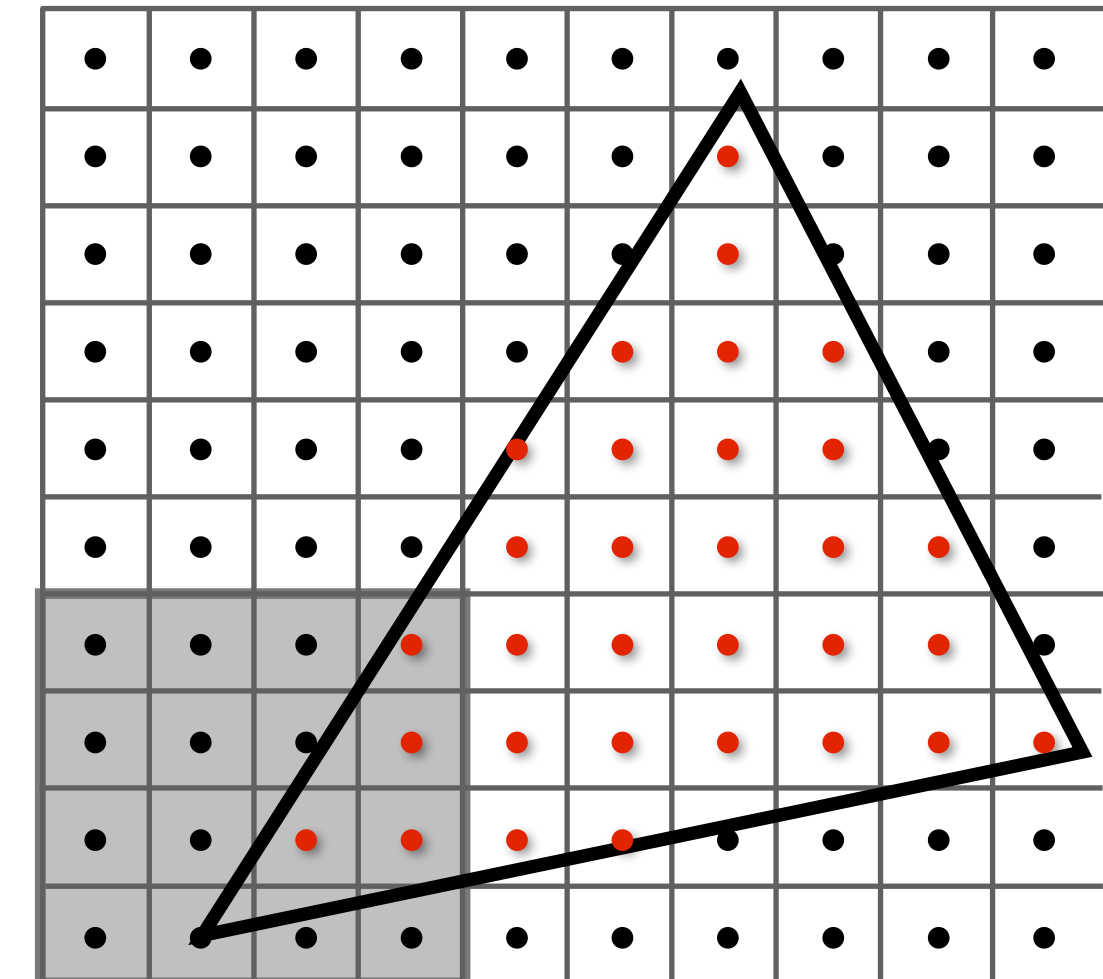
Coverage: 2D triangle/sample tests (does projected triangle cover 2D sample point)

Occlusion: depth buffer

```
initialize z_closest[] to INFINITY // store closest-surface-so-far for all samples
initialize color[] // store scene color for all samples
for each triangle t in scene: // loop 1: over triangles
    t_proj = project_triangle(t)
    for each 2D sample s in frame buffer: // loop 2: over visibility samples
        if (t_proj covers s)
            compute color of triangle at sample
            if (depth of t at s is closer than z_closest[s])
                update z_closest[s] and color[s]
```

“Given a triangle, find the samples it covers”

(finding the samples is relatively easy since they are distributed uniformly on screen)



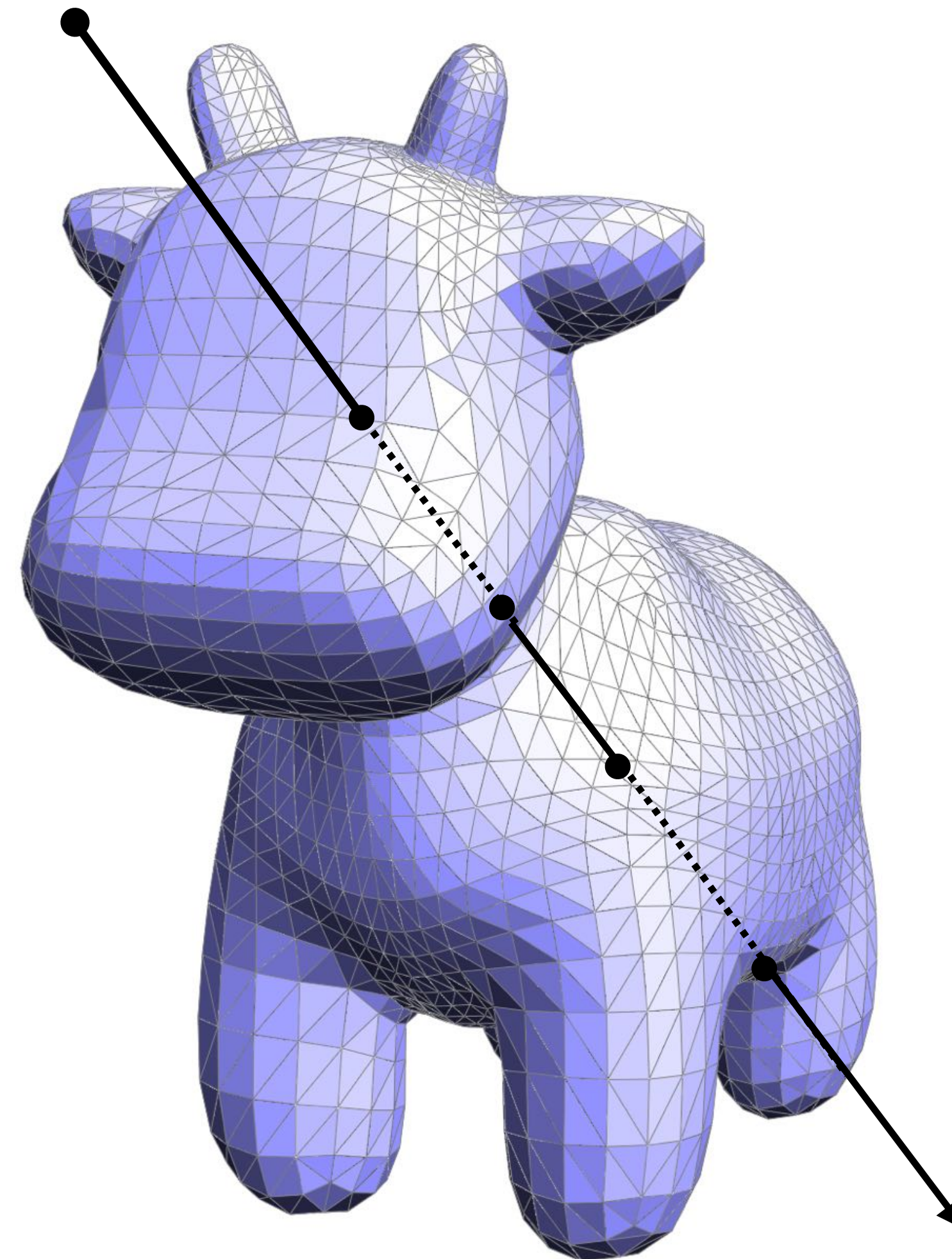
Another way of rendering triangles: ray-scene intersection

Given a scene defined by a set of N primitives and a ray r , find the closest point of intersection of r with the scene

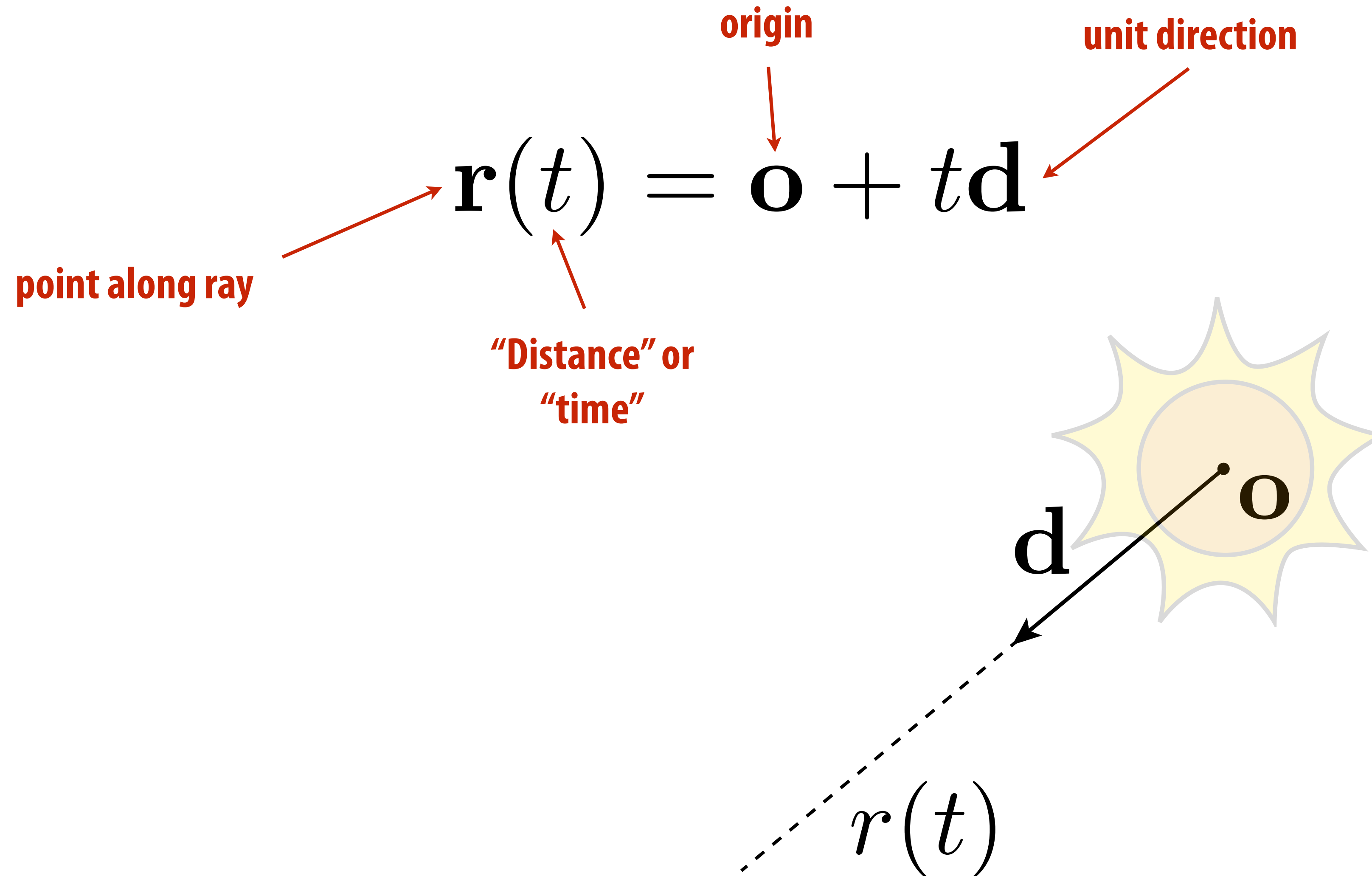
```
t_closest = inf
for each primitive p in scene:
    t = p.intersect(r)
    if t >= 0 && t < t_closest:
        t_closest = t
```

```
// closest hit is:
// r.o + t_closest * r.d
```

(Assume `p.intersect(r)` returns value of t corresponding to the point of intersection with ray r)

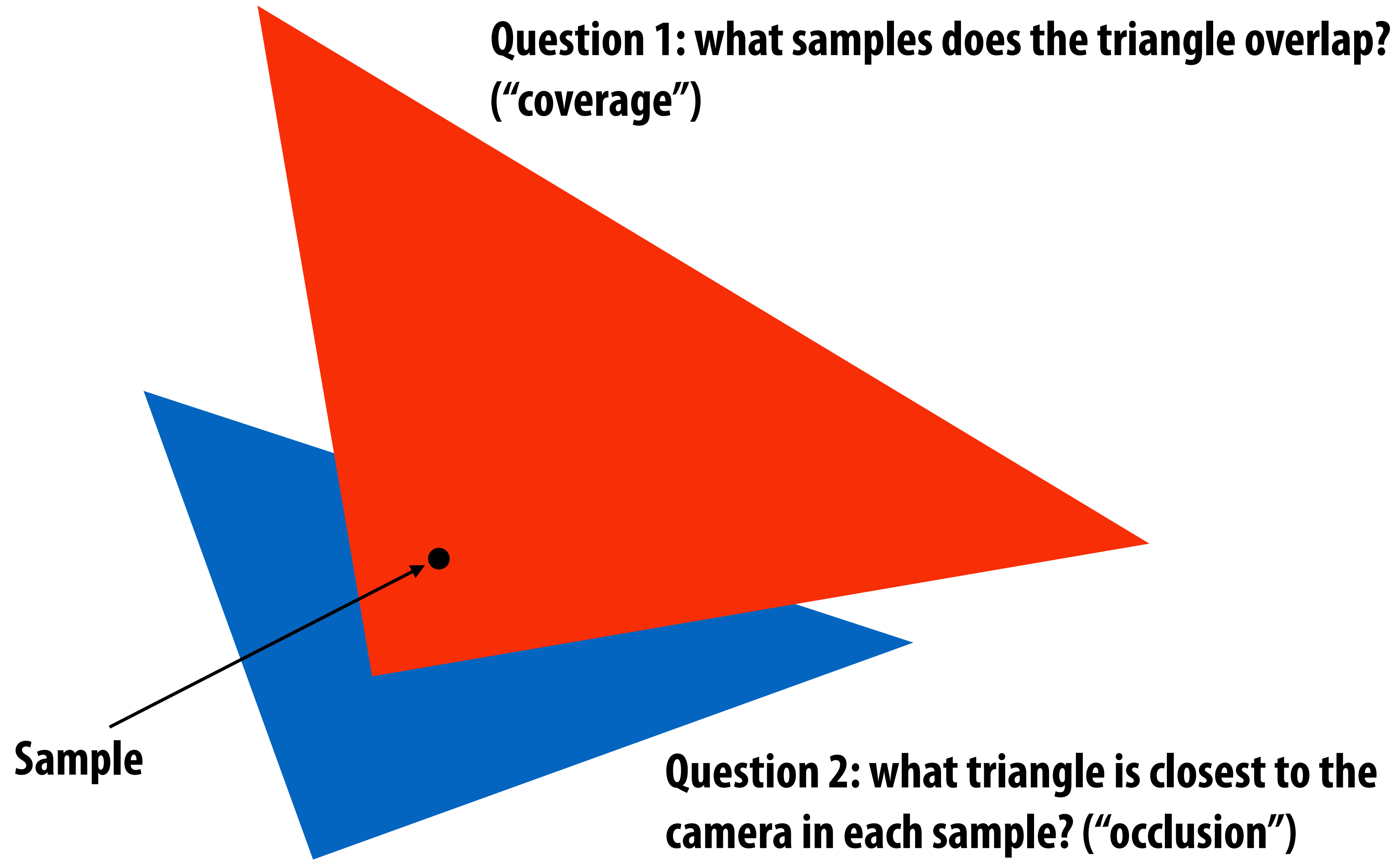


Representing rays



**Rasterization and ray casting are two algorithms for
solving the same problem:
determining surface “visibility” from a virtual camera**

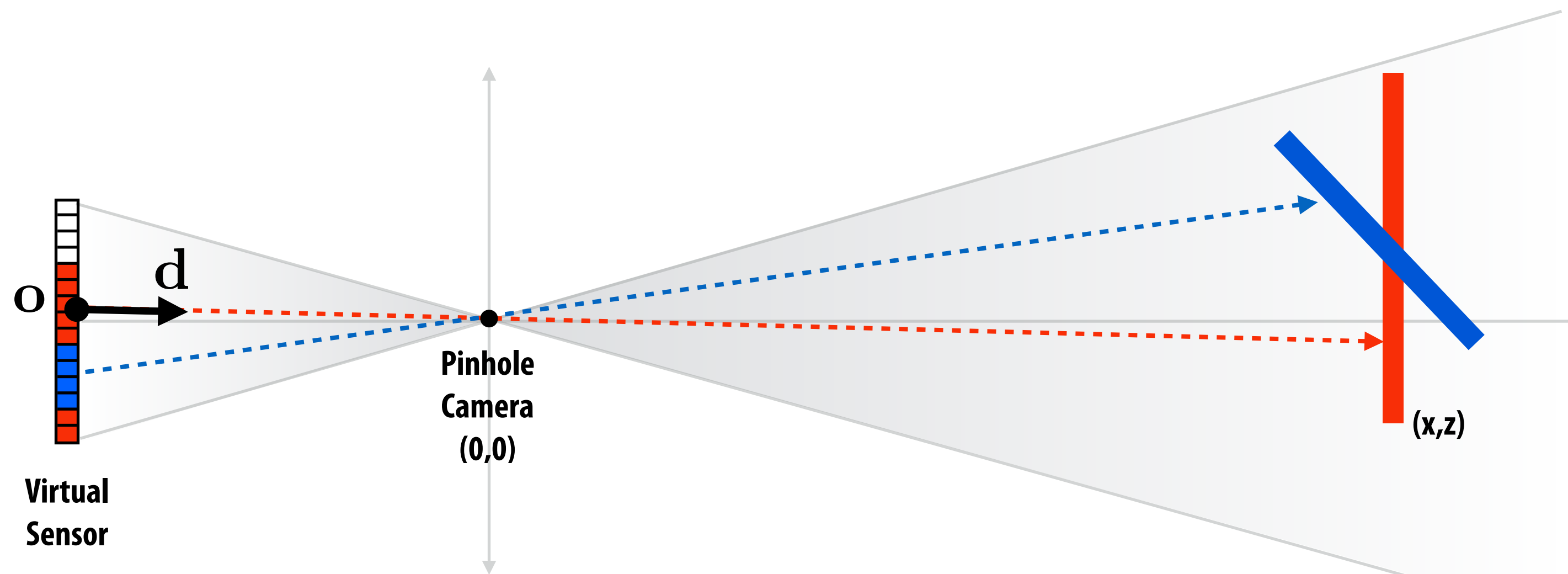
Recall triangle visibility problem:



The visibility problem (described differently)

■ In terms of casting rays from the camera:

- Is a scene primitive hit by a ray originating from a point on the virtual sensor and traveling through the opening of a pinhole camera? (coverage)
- What primitive is the first hit along that ray? (occlusion)



Basic ray casting algorithm

Sample = a ray in 3D

Coverage: 3D ray-triangle intersection tests (does ray “hit” triangle)

Occlusion: closest intersection along ray

```
initialize color[] // store scene color for all samples
for each sample s in frame buffer: // loop 1: over visibility samples (rays)
    r = ray from s on sensor through pinhole aperture
    r.min_t = INFINITY // only store closest-so-far for current ray
    r.tri = NULL;
    for each triangle tri in scene: // loop 2: over triangles
        if (intersects(r, tri)) { // 3D ray-triangle intersection test
            if (intersection distance along ray is closer than r.min_t)
                update r.min_t and r.tri = tri;
        }
    color[s] = compute surface color of triangle r.tri at hit point
```

Compared to rasterization approach: just a reordering of the loops!

“Given a ray, find the closest triangle it hits.”

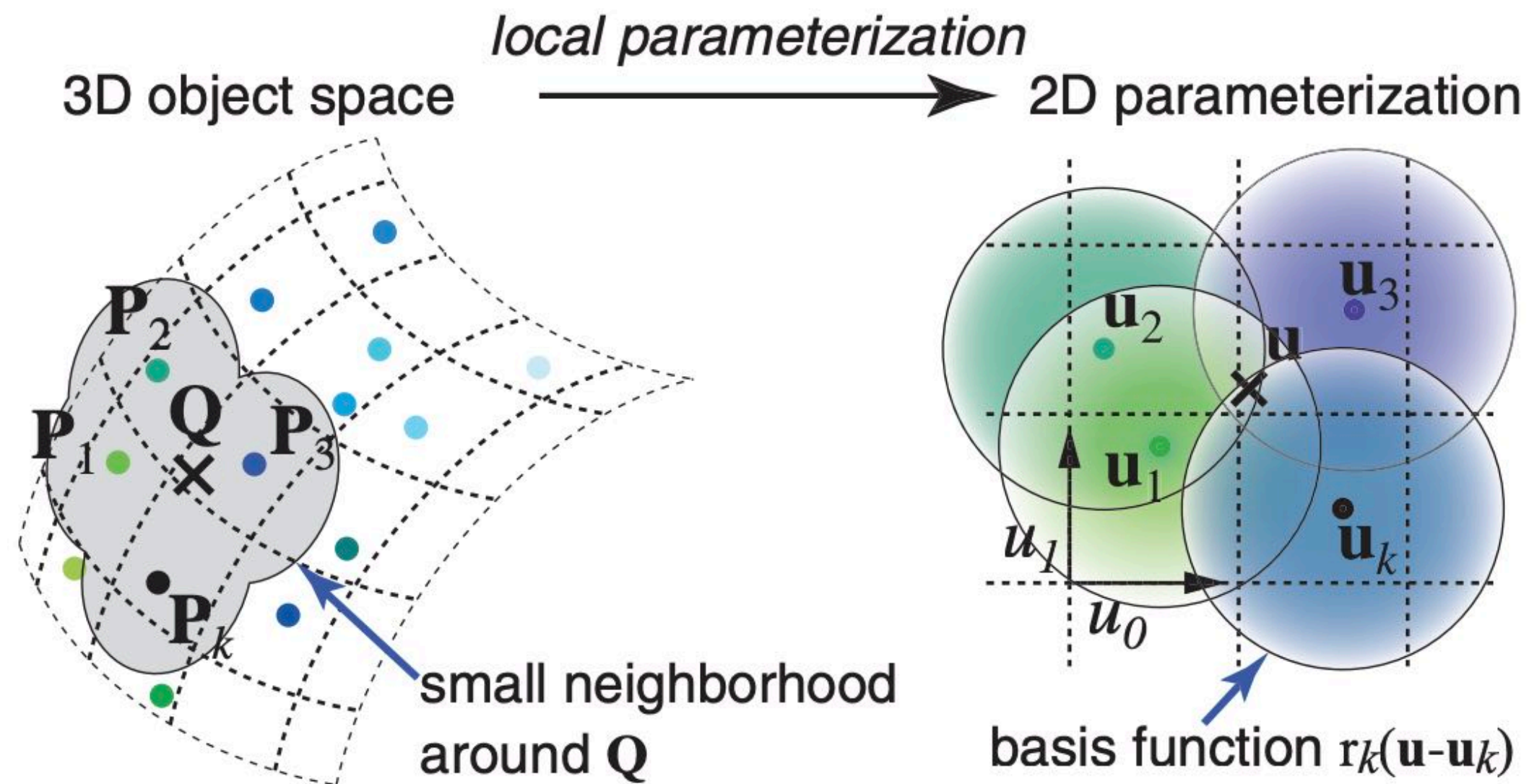
Rendering 3D points

Given camera position and 3D position of point: (1) project point onto screen (2) color pixel if closest

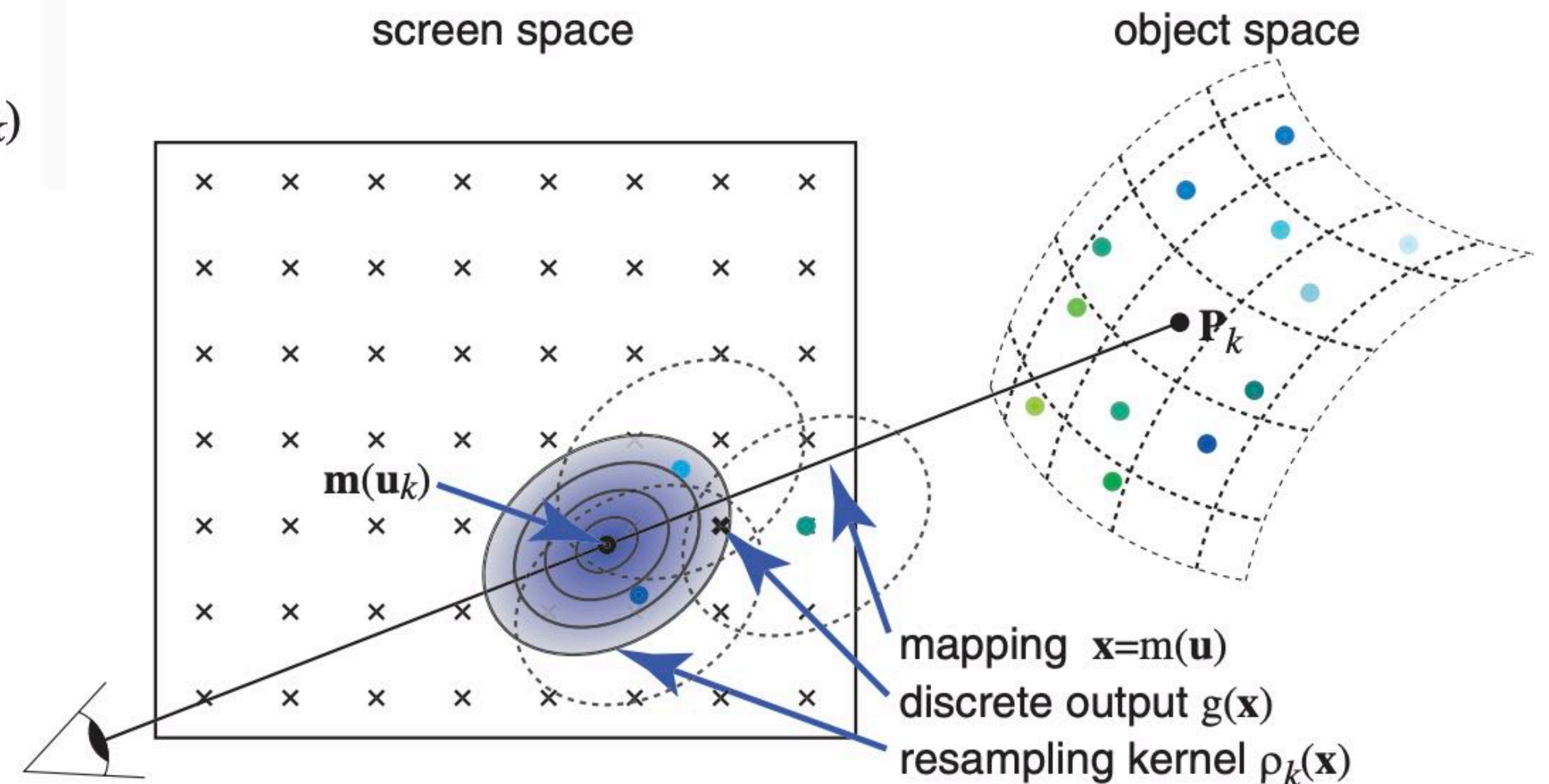


Rendering “splats” / 3D gaussians / surfels

- Treat surface as a collection of “Gaussian blobs” (convolve points with Gaussian filter)



- 3D Gaussians turn into oriented 2D gaussians when projected onto the 2D screen
- Can render the blobs back to front (requires alpha compositing)

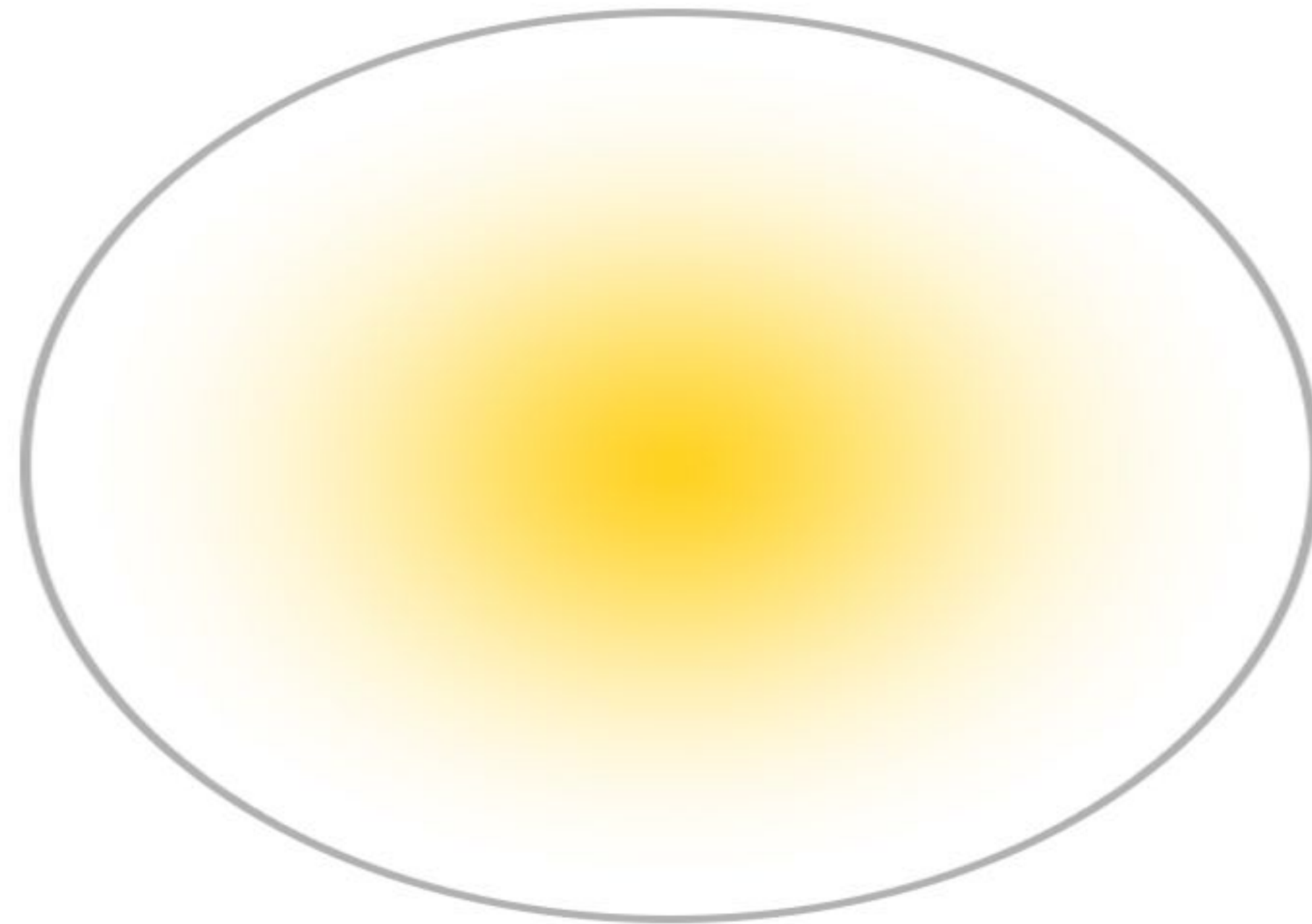
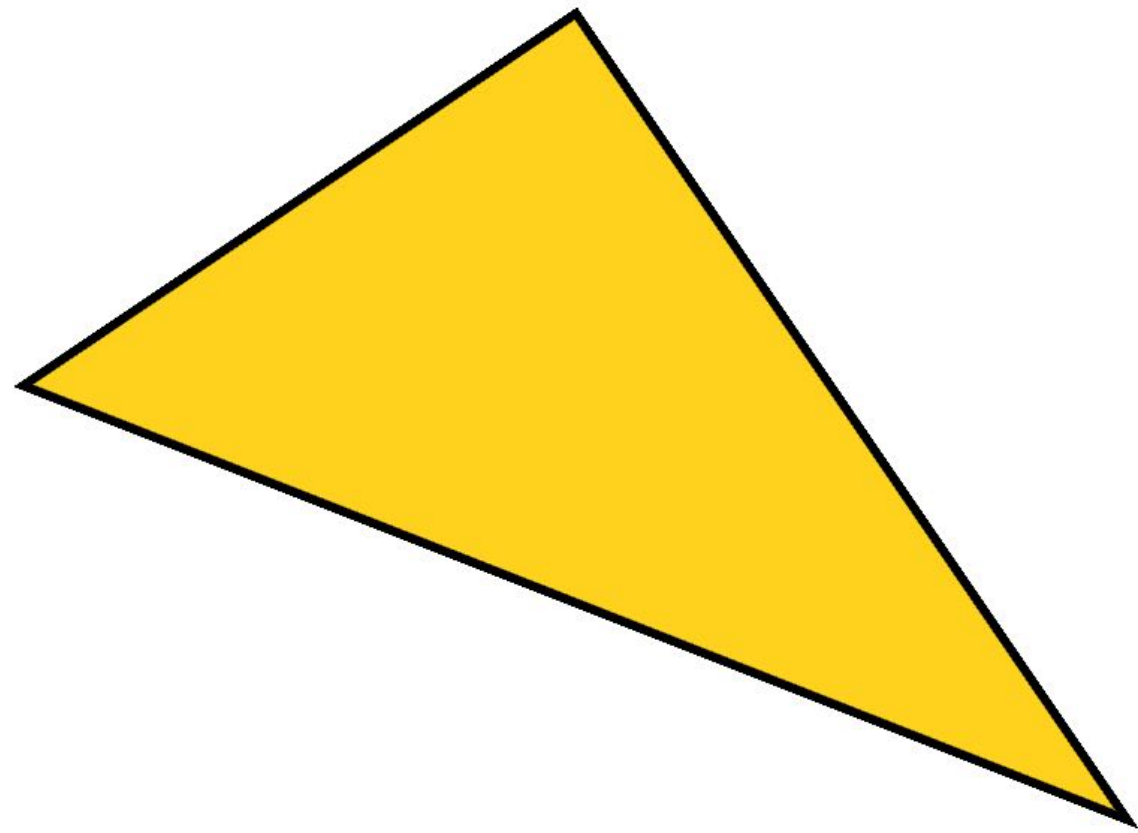






Consider representing a triangle with Gaussians

- Could approximate the triangle with a lot of small 3D gaussians near the triangle's edges
- Not so efficient, eh?



Another representation: regular 3D grid representation

Consider storage requirements:

1024^3 cells

Ignore directional dependency: rgb + 4 bytes/cell
(~4 GB)

Now consider directional dependency on (ϕ, θ)
... much worse

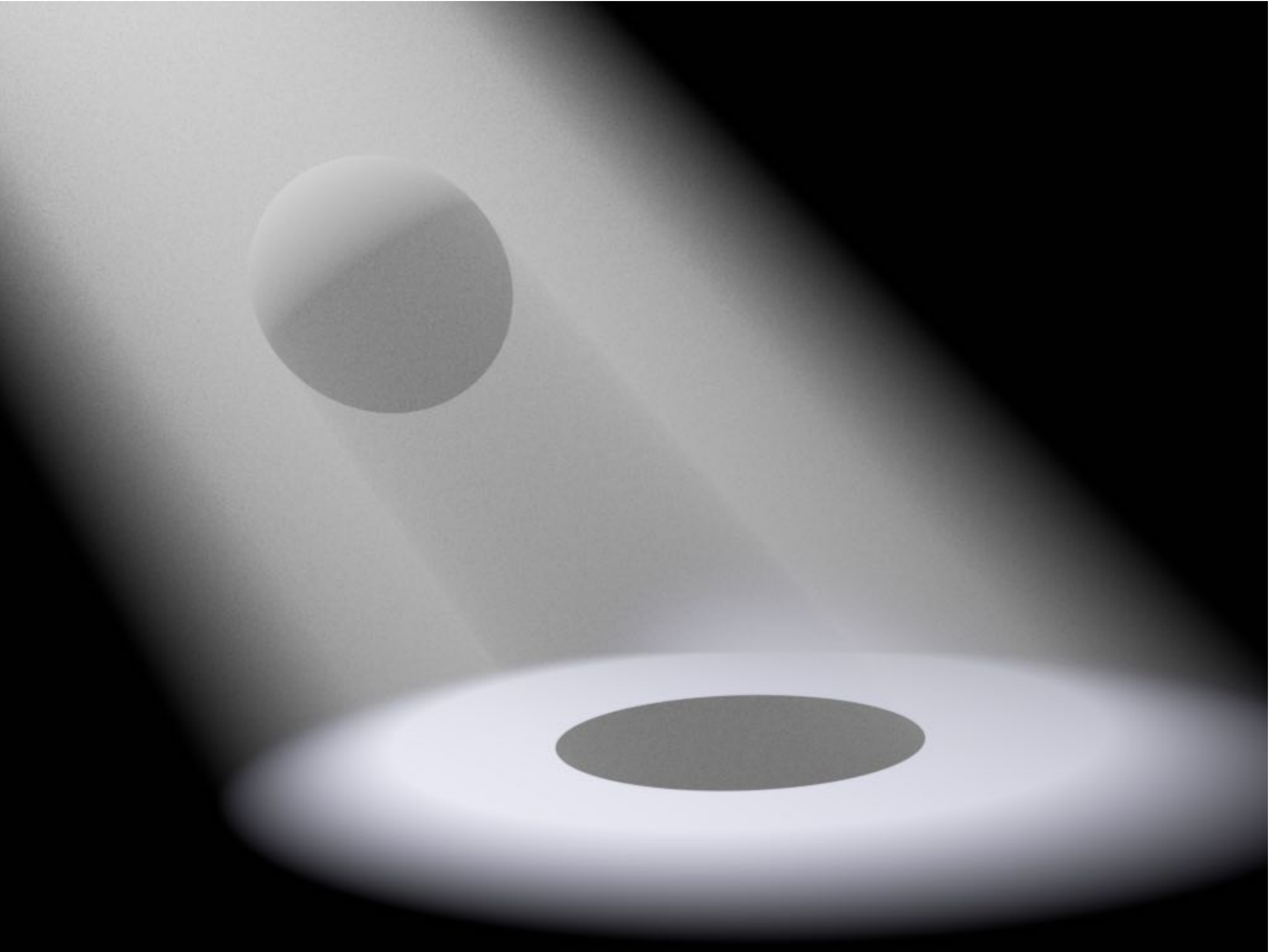


Typical challenge of
dense voxel
representations:
limited resolution



Credit: Voxel Ville NFT (voxelville.io)

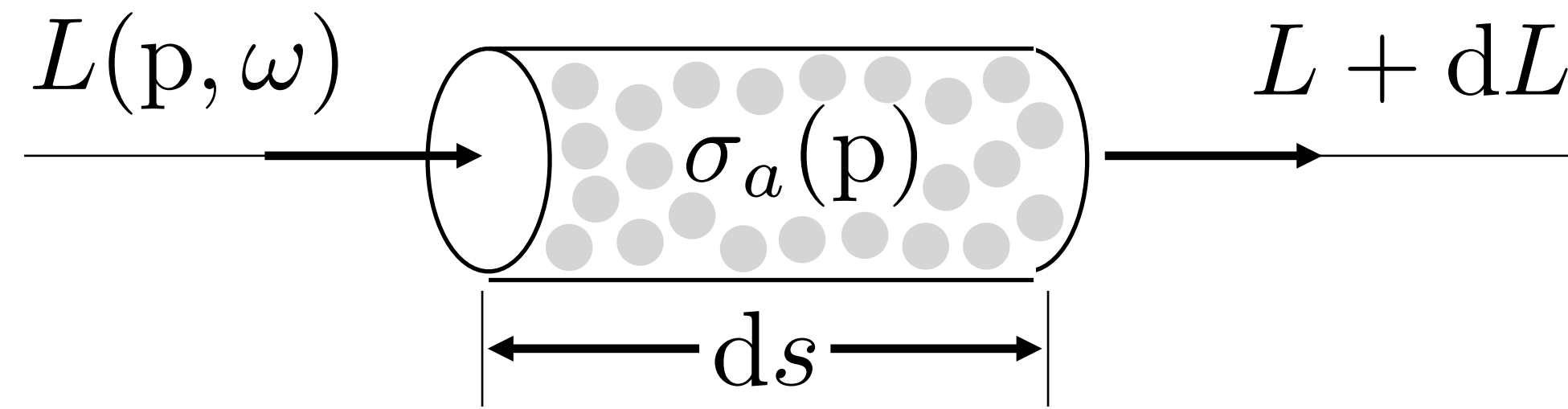
Volumetric effects



Another motivation for non-triangle representations like points/gaussians/volumes: hard to accurately estimate surface triangle mesh in complex real-world situations



Absorption in a volume



$$\mathbf{p} = (x, y, z)$$

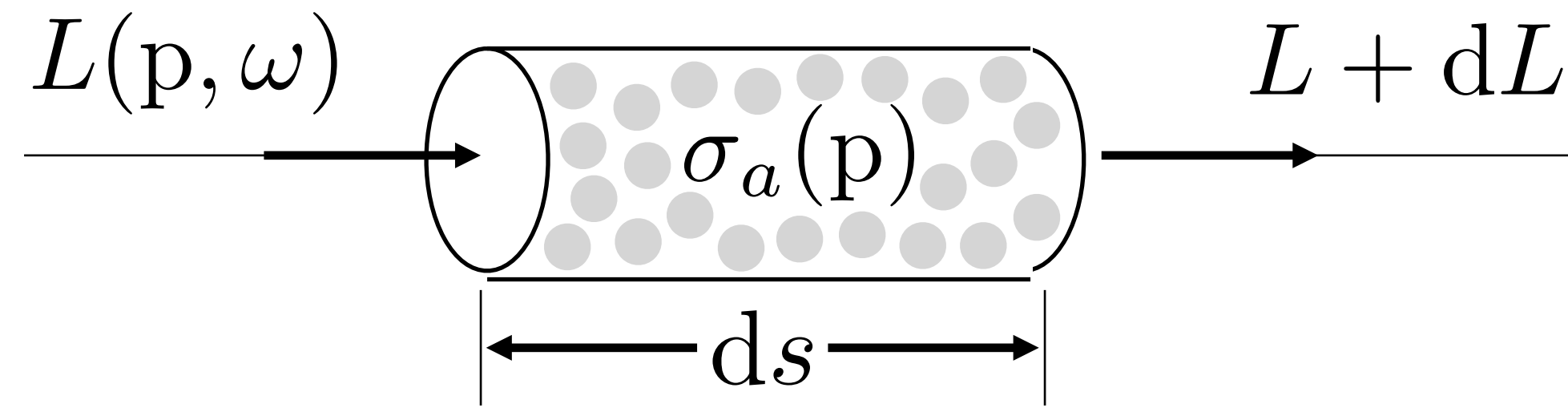
$$\boldsymbol{\omega} = (\phi, \theta)$$

$$dL(p, \omega) = -\sigma_a(p) L(p, \omega) ds$$

$$\frac{dL(p, \omega)}{ds} = -\sigma_a(p) L(p, \omega)$$

- $L(p, \omega)$ radiance along a ray from p in direction ω
- Absorption cross section at point in space: $\sigma_a(p)$
 - Probability of being absorbed per unit length
 - Units: 1/distance

Absorption in a volume



$$\mathbf{p} = (x, y, z)$$

$$\boldsymbol{\omega} = (\phi, \theta)$$

$$\frac{dL(p, \omega)}{L(p, \omega)} = -\sigma_a(p) ds$$

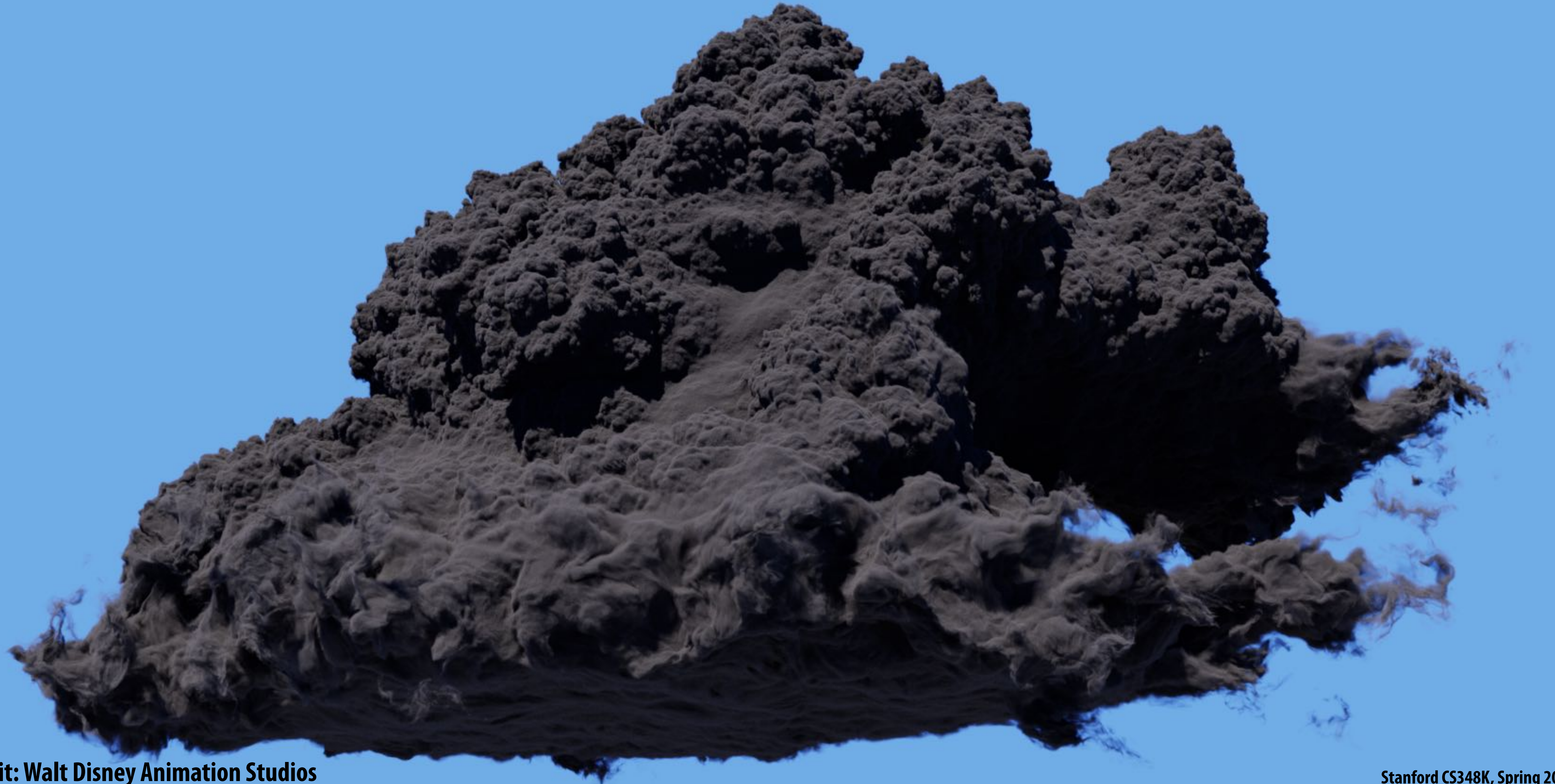
$$L(\mathbf{p} + s\boldsymbol{\omega}, \omega) = e^{-\int_0^s \sigma_a(\mathbf{p} + s'\boldsymbol{\omega}) ds'} L(\mathbf{p}, \omega) = T(s) L(\mathbf{p}, \omega)$$

Transmittance: $T(s) = e^{-\int_0^s \sigma_a(\mathbf{p} + s'\boldsymbol{\omega}, \omega) ds'}$

Absorption: lower density



Absorption: higher density



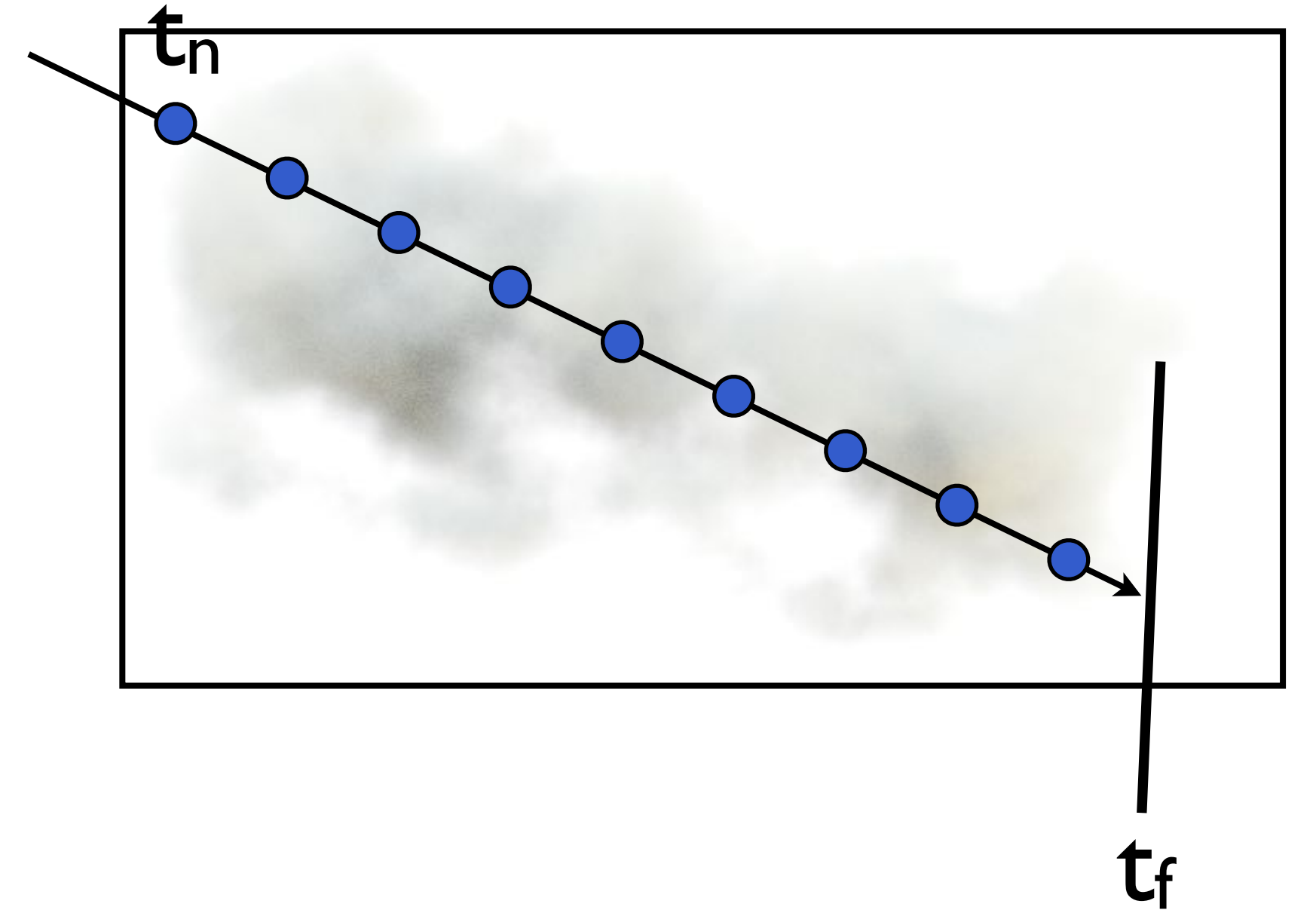
Rendering volumes

$$\sigma(\mathbf{p})$$

$$c(\mathbf{p}, \omega)$$



Volume density and color at all points in space.
e.g., Values stored in a 3D grid



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

Dense 3D volumes = high storage cost

Consider storage requirements:

1024^3 cells

Ignore directional dependency: rgb + 4 bytes/cell
(~4 GB)

Now consider directional dependency on (ϕ, θ)
... much worse



Typical challenge of
dense voxel
representations:
limited resolution

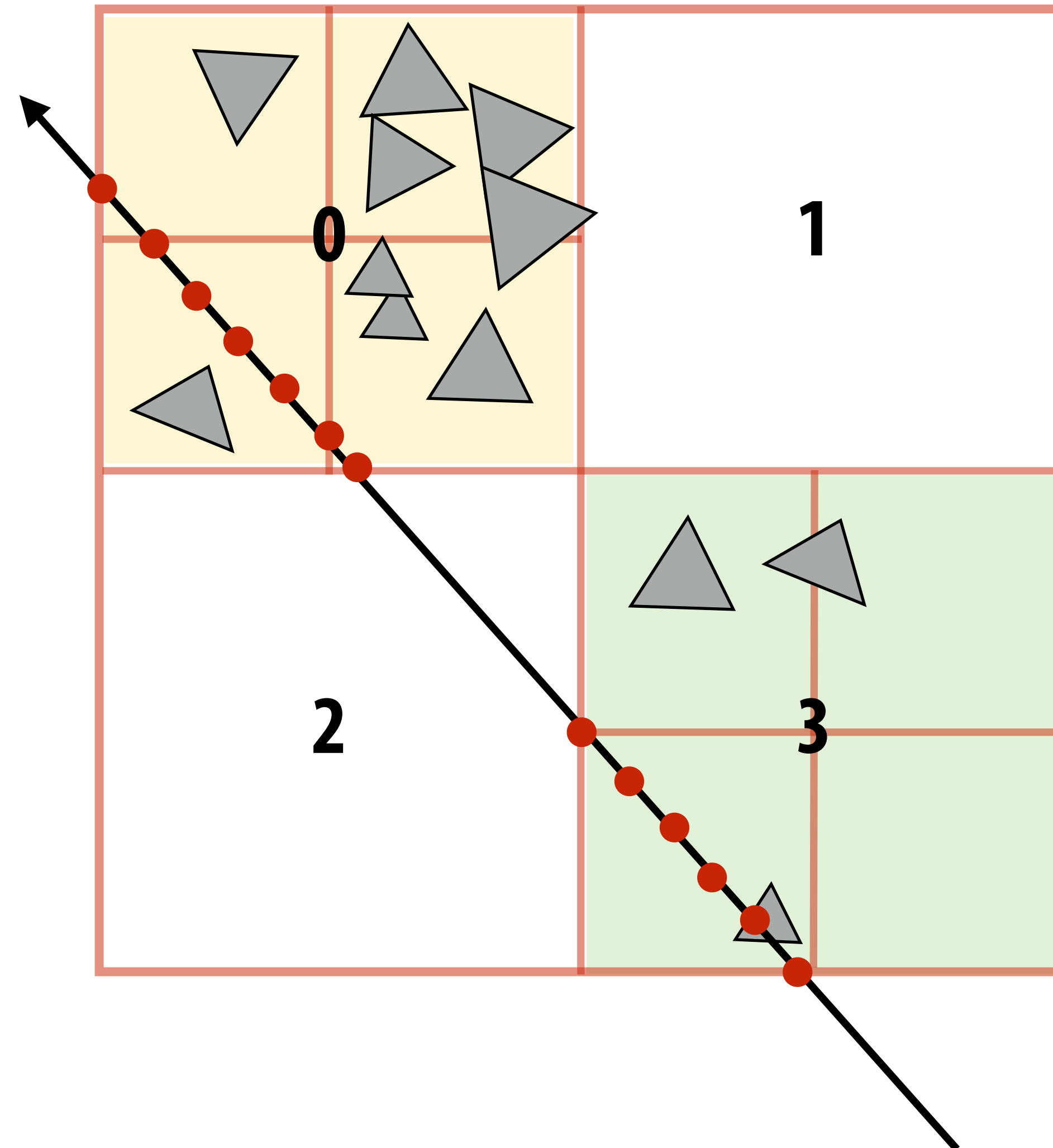
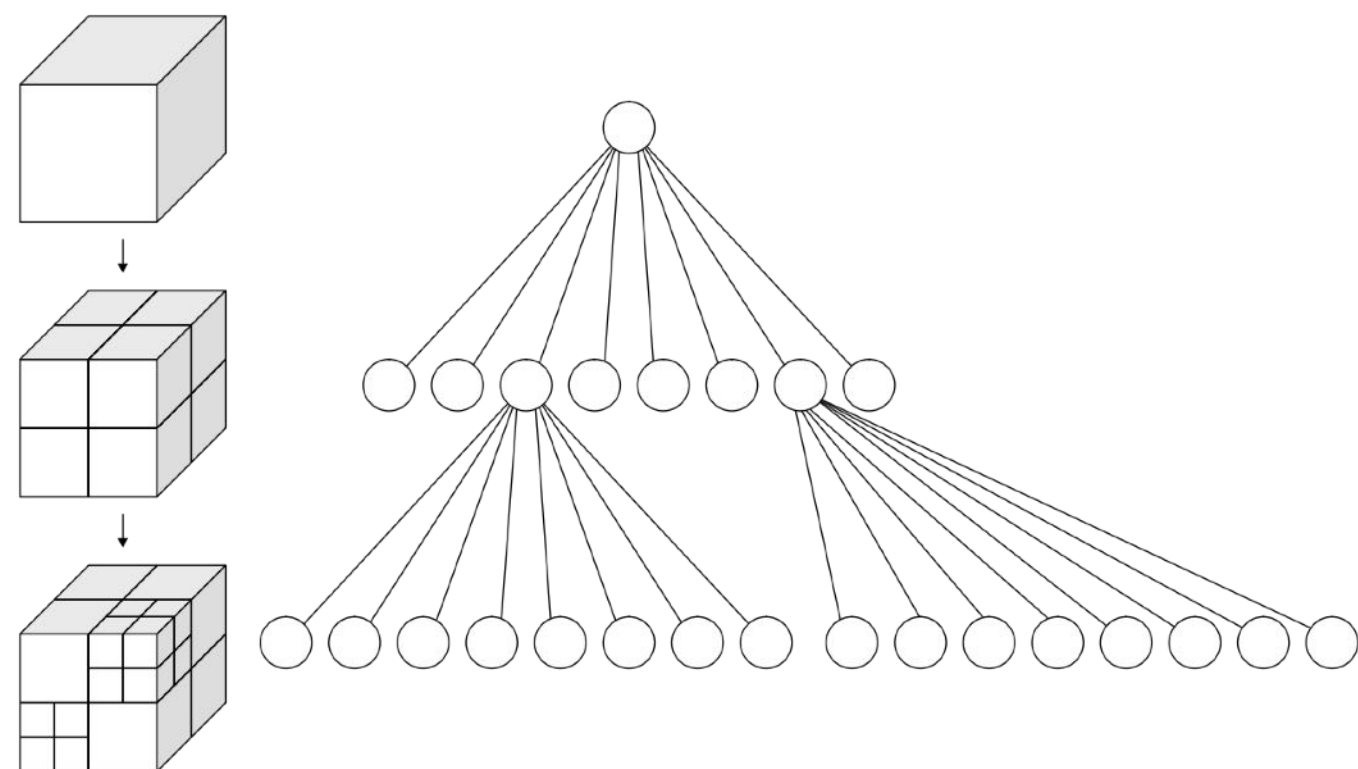
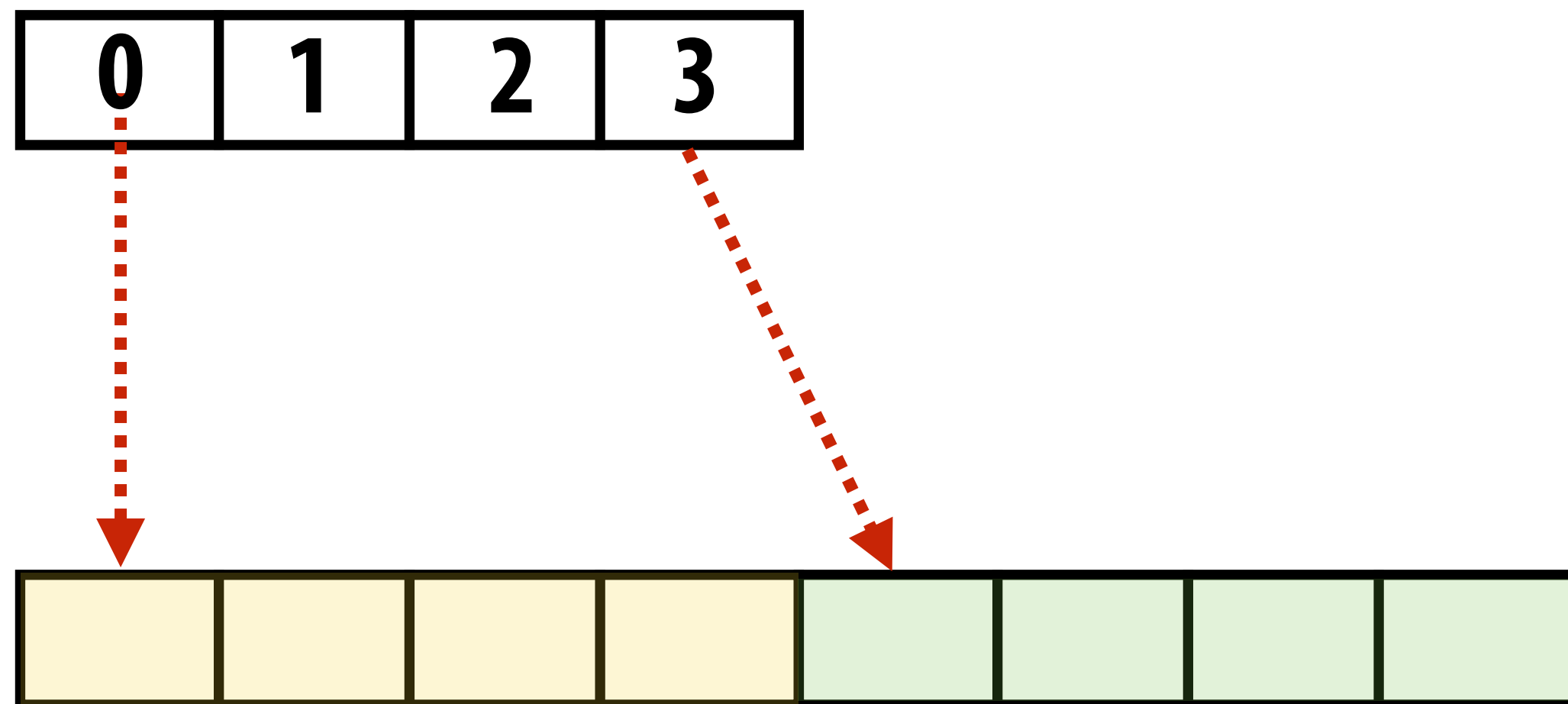


Credit: Voxel Ville NFT (voxelville.io)

Sparse volumes

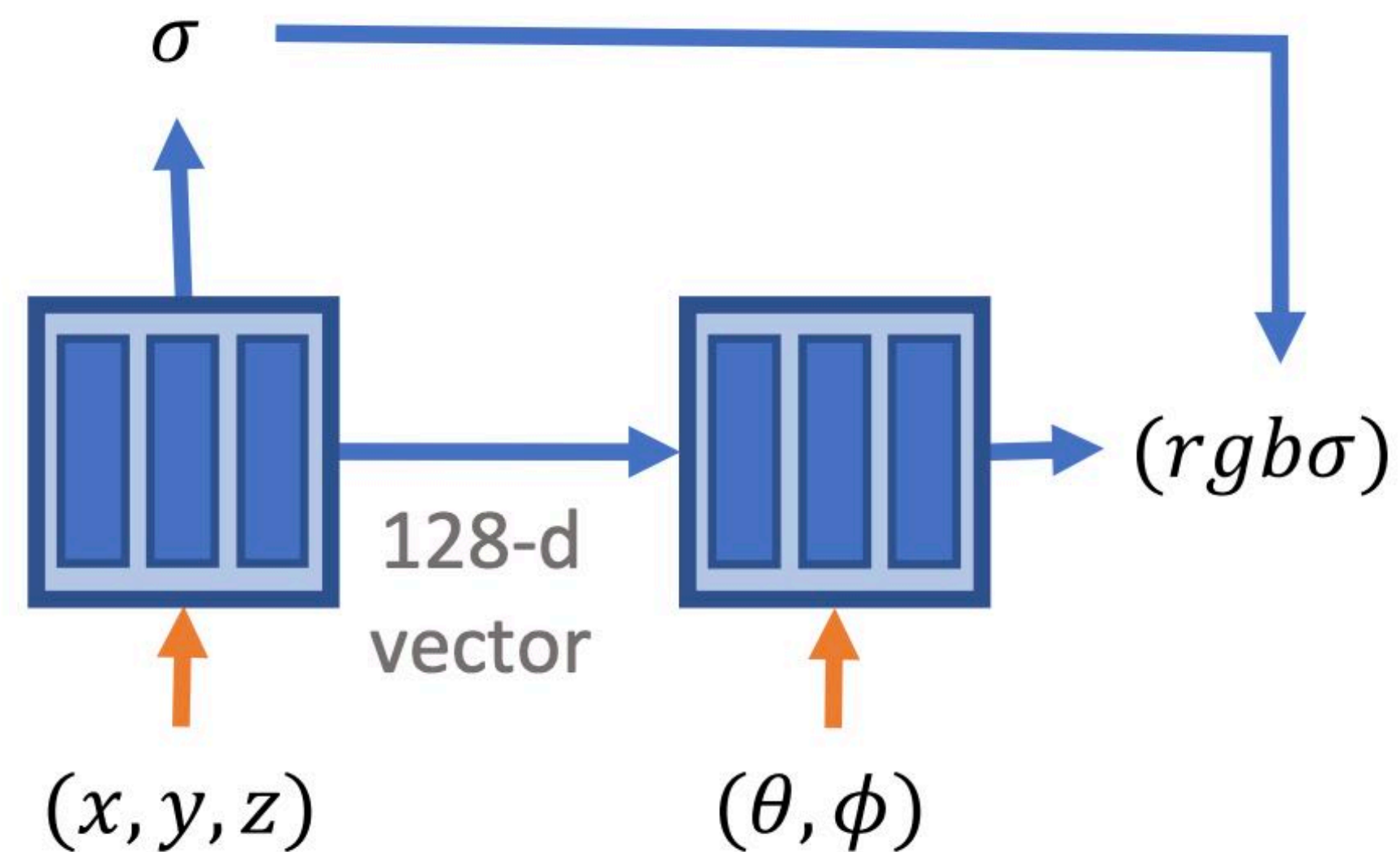
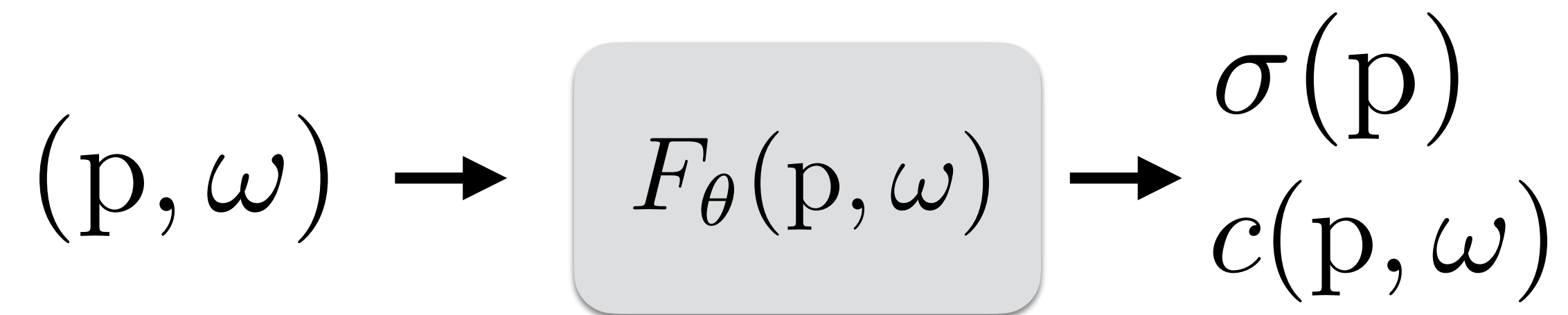
Quad-tree: nodes have 4 children (partitions 2D space)

Octree: nodes have 8 children (partitions 3D space)



A very compressed volume representation

- Use DNN to compress information in a volume
- Why not just learn an approximation to the continuous function that matches observations from different viewpoints?



Many different parameterizations

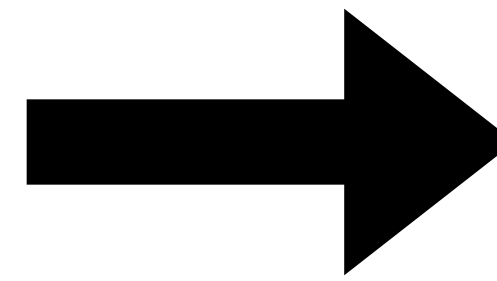
- **Mesh vertex positions + texture values**
- **3D point positions + colors**
- **3D oriented Gaussians + colors**
- **Dense 3D voxels**
- **Sparse 3D voxels**
- **DNN weights**
- **Many combinations not discussed: sparse 3D grid of DNN weights, hash table of DNN weights, etc...**

Reconstruction problem

- **Given many views of a scene for which camera position is known, recover the parameters of a scene representation SO THAT rendering the scene representation from that known view generates the captures image!**
- **Need a differentiable renderer to recover parameters!**

Novel view synthesis problem

Input photos (from a fixed set of views)



Novel views
(camera position different from those in input photos)

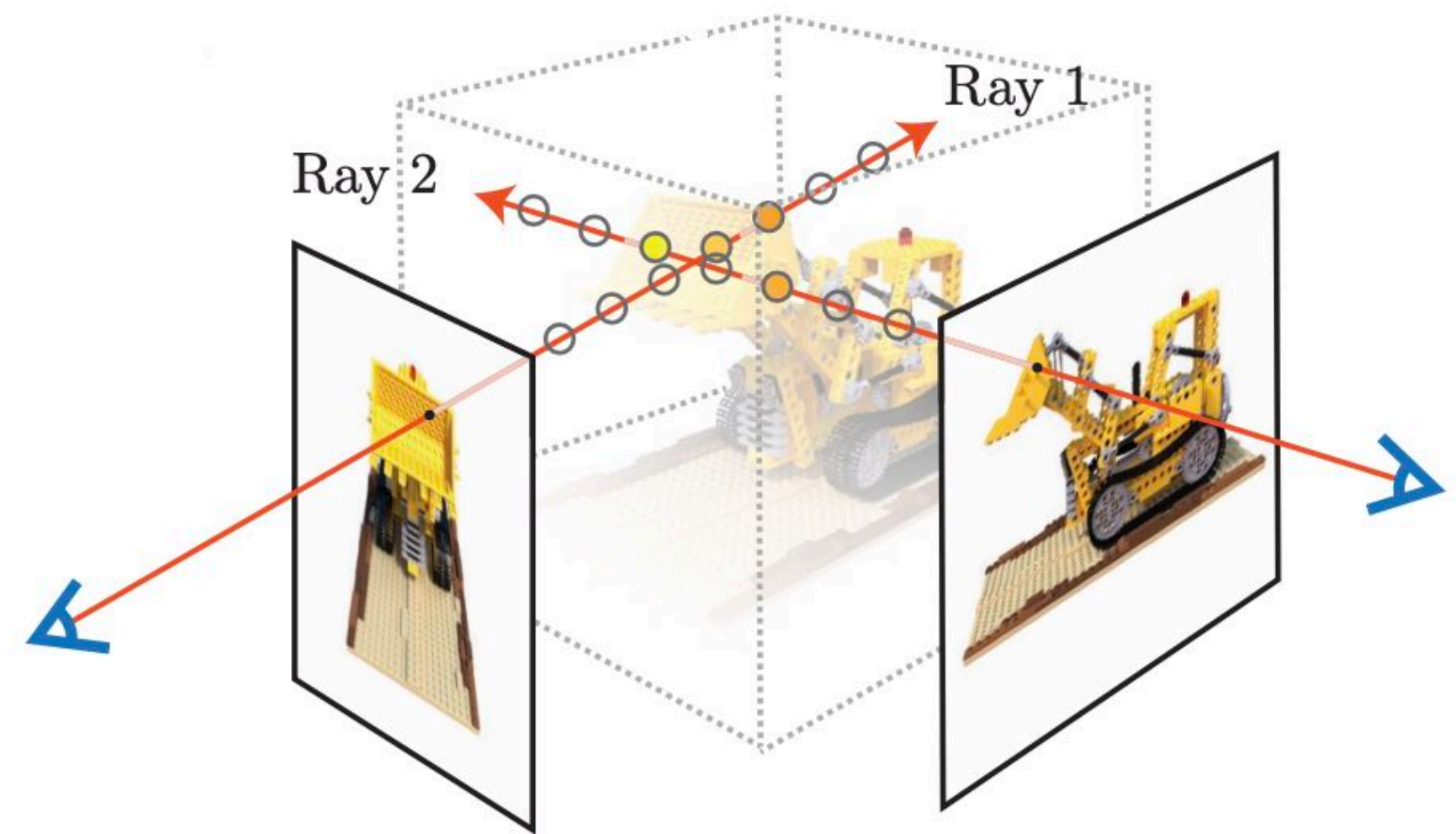
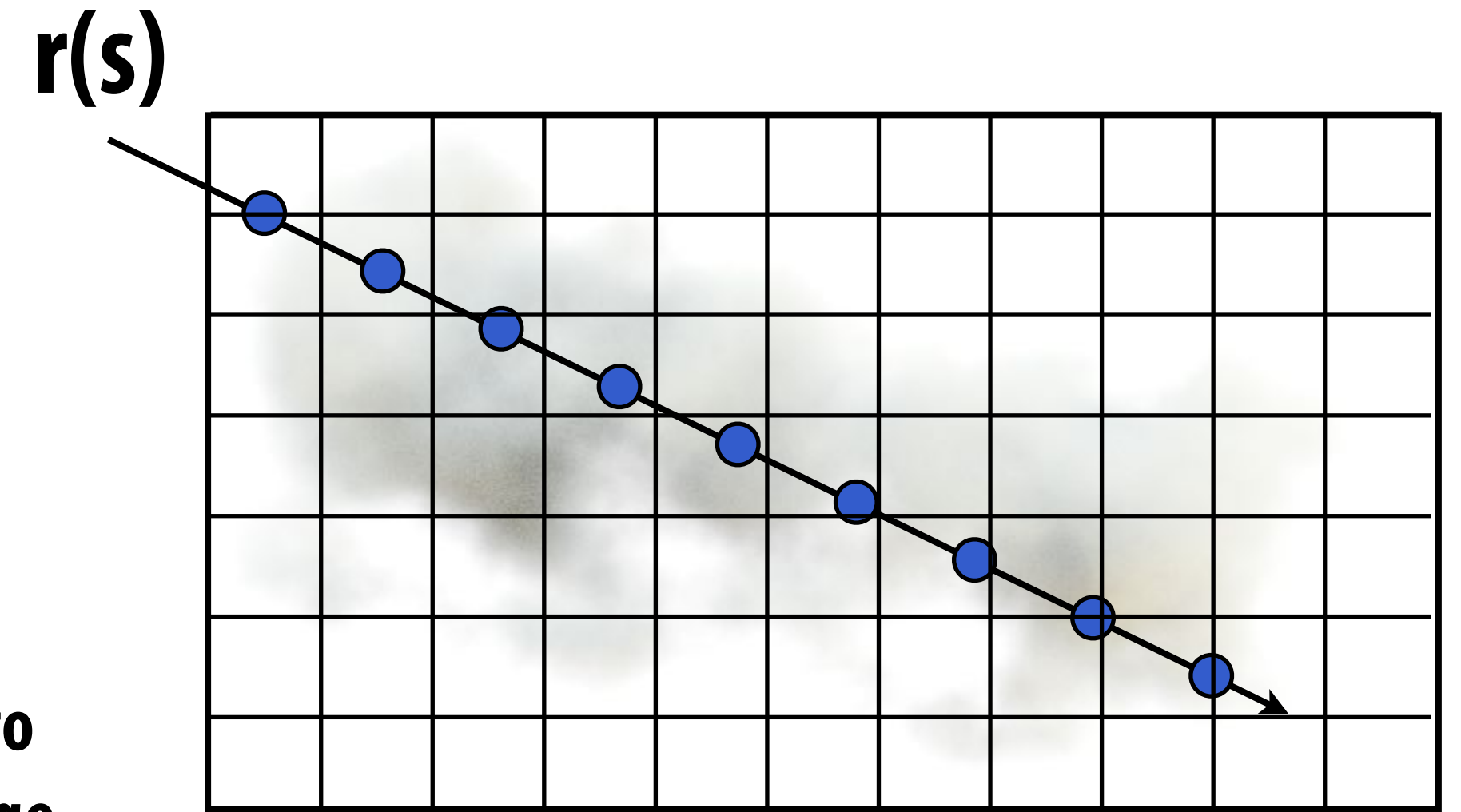


Optimizing volumes

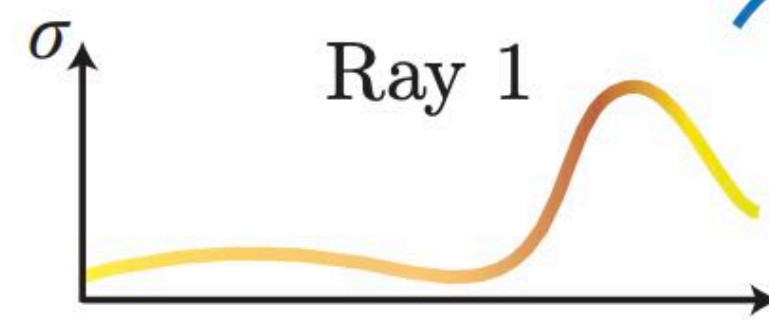
$$C(\mathbf{r}) = \int_{t_n}^{t_f} T(t) \sigma(\mathbf{r}(t)) \mathbf{c}(\mathbf{r}(t), \mathbf{d}) dt, \text{ where } T(t) = \exp\left(-\int_{t_n}^t \sigma(\mathbf{r}(s)) ds\right)$$

Idea: optimize volume values (opacity and color) so that $C(\mathbf{r})$ matches that of photos.

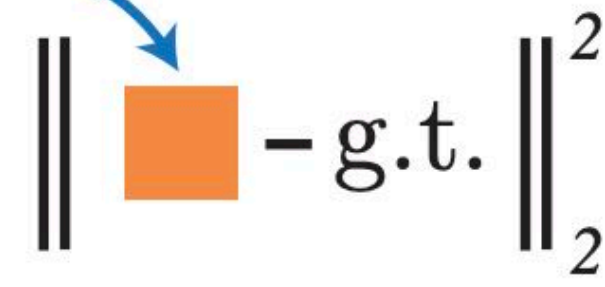
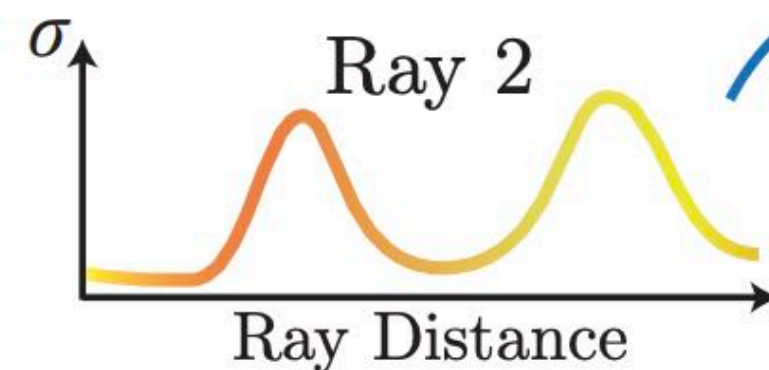
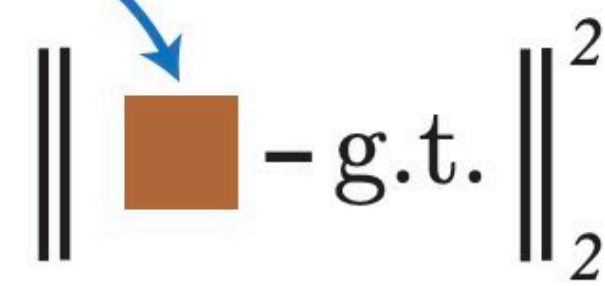
For many rays.... trace through volume... see if the result matches the photo... use error to update volume opacity/color values



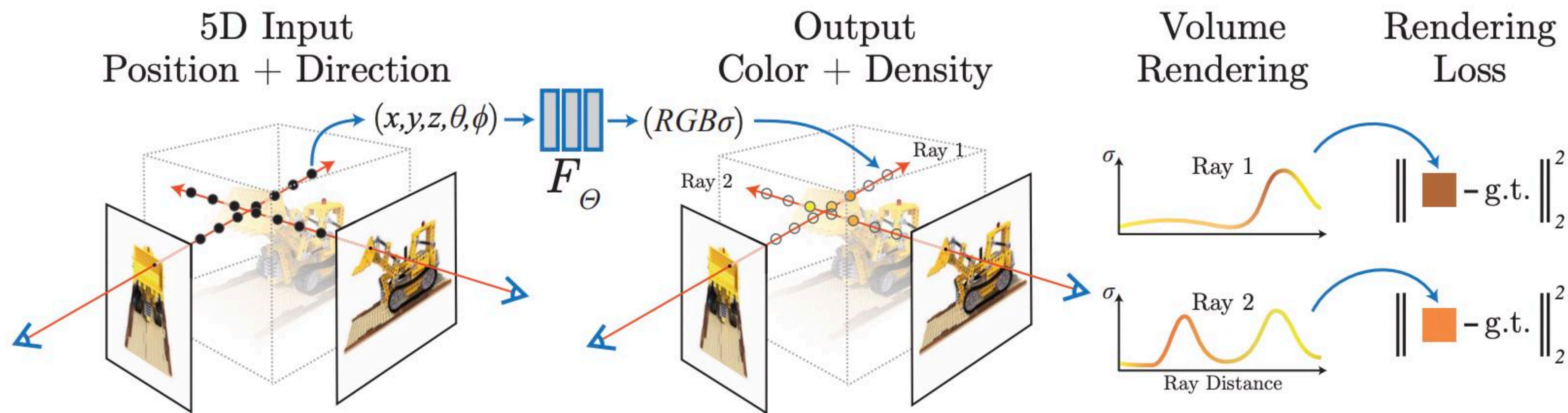
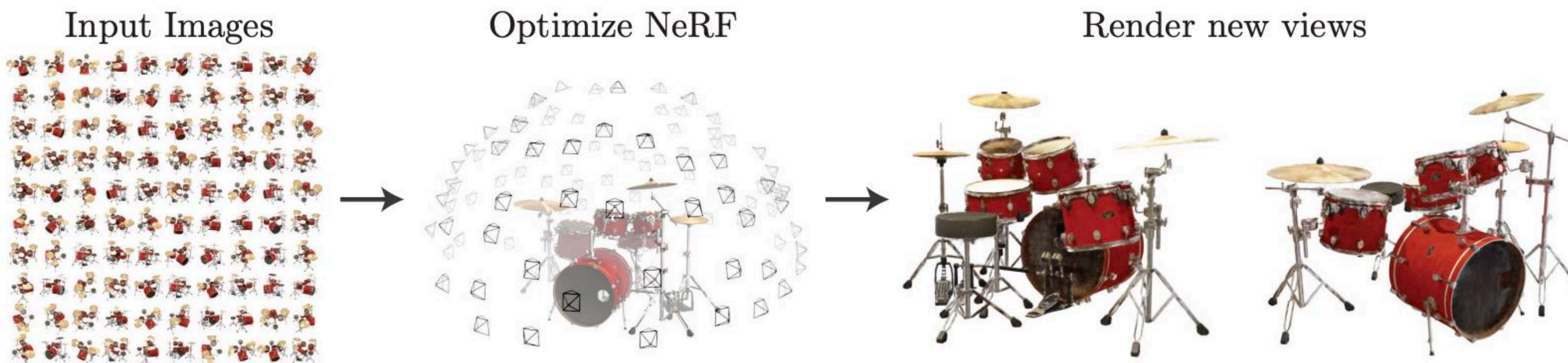
Compute radiance along ray through volume



Compare to actual image



Learning neural radiance fields (NeRF)



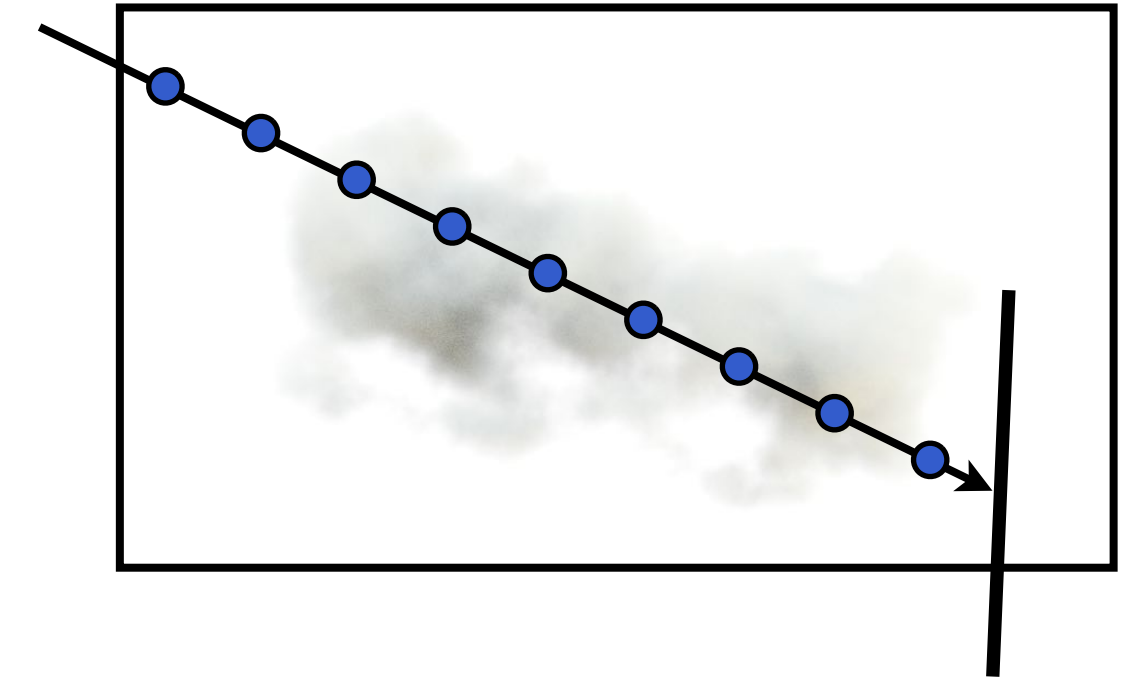
Key idea: differentiable volume renderer to compute $dC/d(\text{color})d(\text{opacity})$

Great visual results!



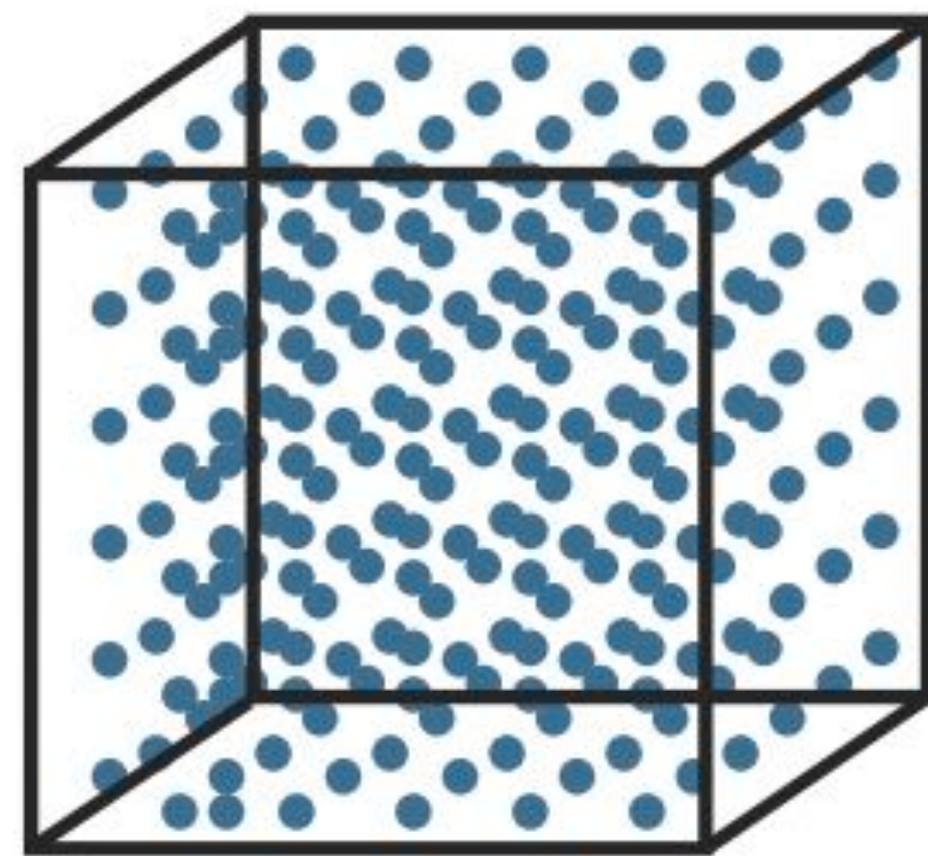
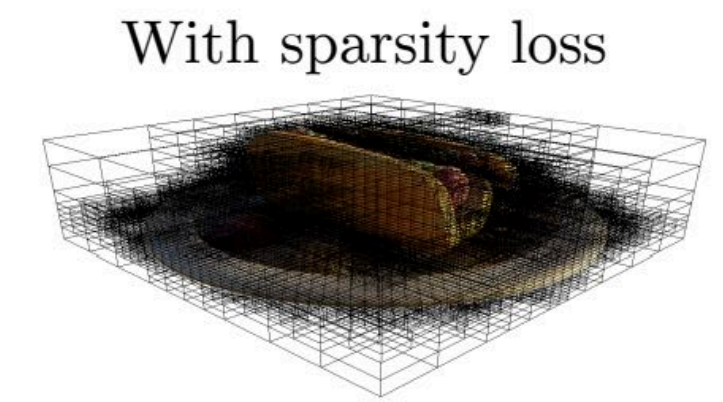
What just happened?

- **Continuous coordinate-based representation vs regular grid: DNN “learns” how to use its weights to produce high-resolution output where needed... given input data**
- **Compact representation: trades-off space for expensive rendering**
 - **Good: a few MBs = effectively very high-resolution dense grid**
 - **Bad: must evaluate DNN every step during ray marching**
 - **And the DNN is a “big” MLP (8-layer x 256)** ← MLP must do real work to associate weights with 5D locations
 - **Bad: must step densely (because we don't know where the surface is)**
- **Compact representation: optimization can learn to interpolate views despite complexity of volume density and radiance function**
 - **Only prior is the separation into positional σ and directional rgb**
 - **Training time: hours to a day to learn a good NeRF**

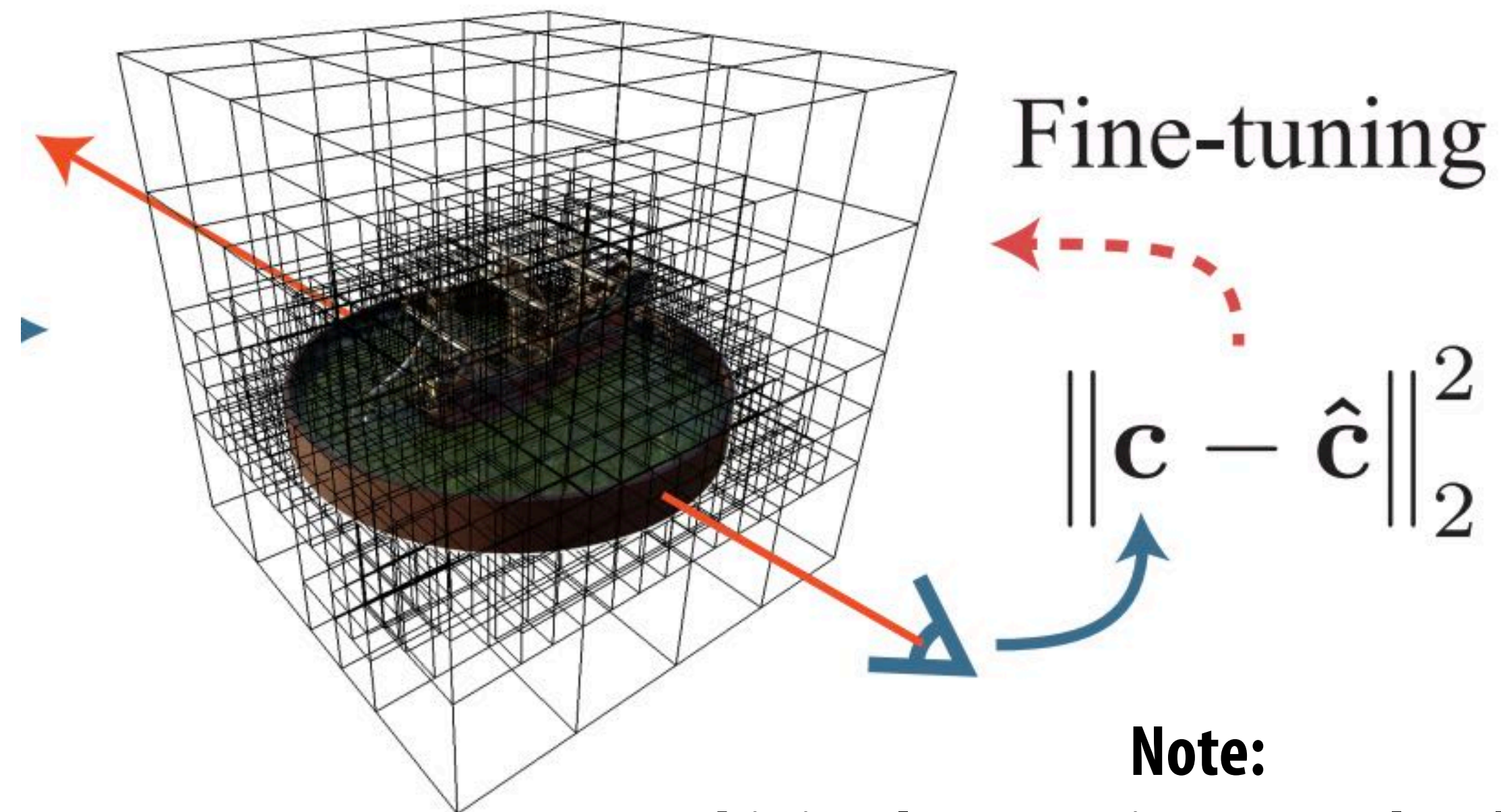


Let's just run optimization for a bit...

- Optimization will push some opacity values to 0
- DNN tells us where the empty space is!
- Then convert dense opacity grid to an octree representation that's more efficient to render from...
- With the octree structure ***fixed***, we can continue to optimize color/density at leaves



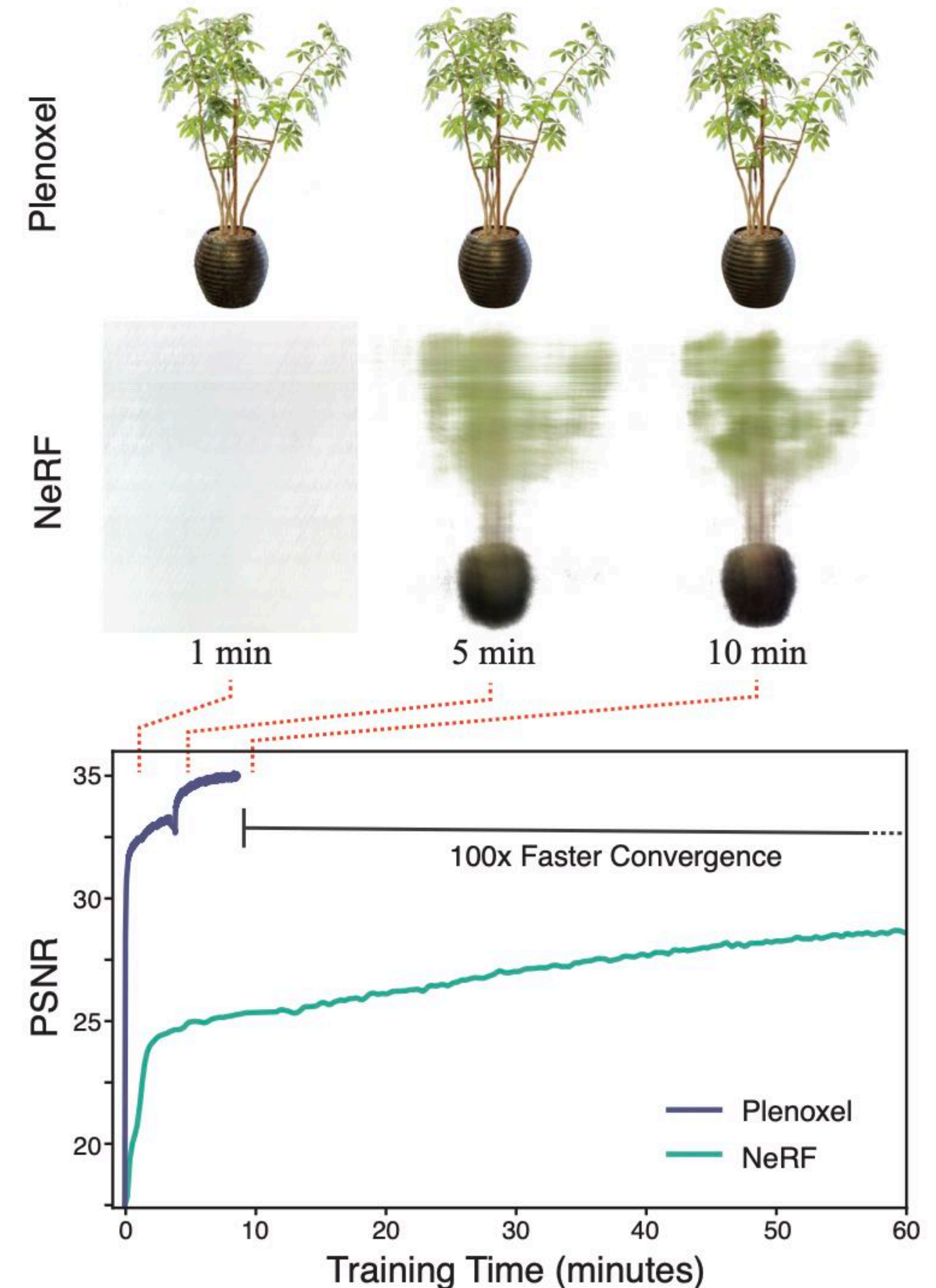
Use the initial MLP to densely sample volume
(Find the empty space that's used to build the octree)



Finally...back to where we began

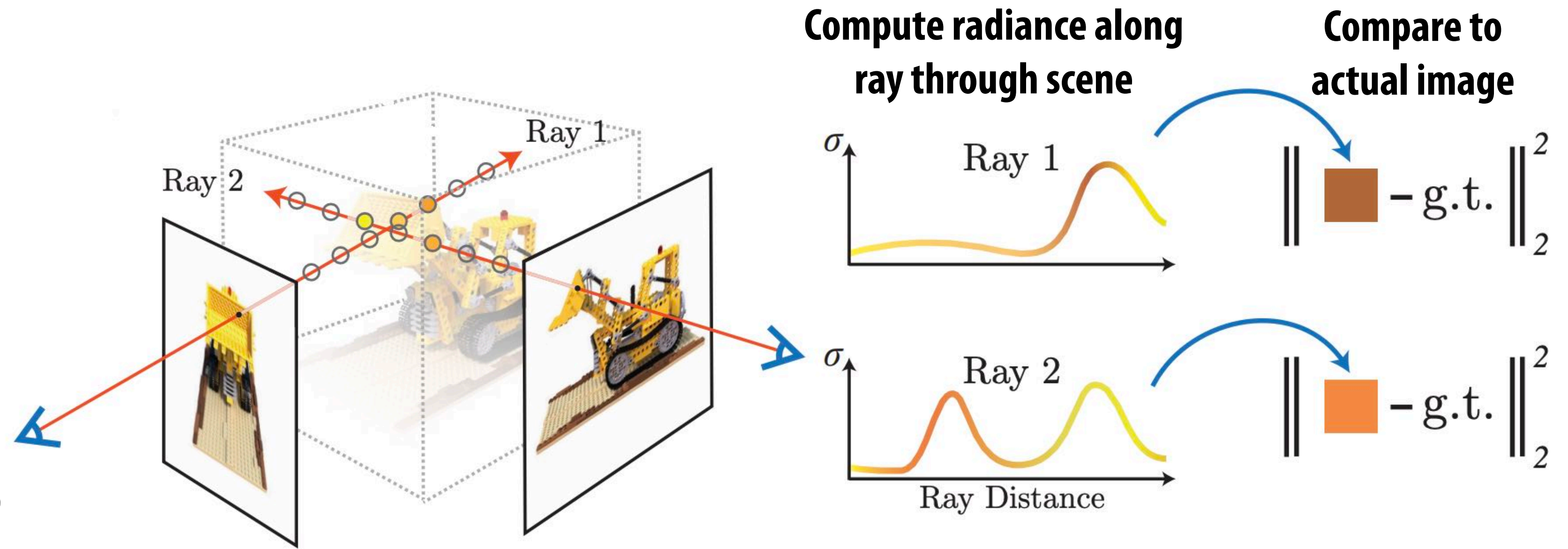
Plenoxels [CVPR 22]

- Start with a dense 3D grid of SH coefficients, optimize those coefficients at low resolution
- Now move to a sparse higher resolution representation (octree)
- Directly optimize for opacities and SH coefficients using differentiable volume rendering
- **No neural networks. Just optimizing the octree representation of baked SH lighting**
- **Takeaway: conventional computer graphics representations are efficient representations to learn/optimize on**



Optimization to produce Gaussians, not voxels

- Earlier in lecture: optimization produces color and opacity at each voxel
- Now: same idea, but optimization chooses color, position, and radius of the Gaussians
 - Now: also need to decide on the number of Gaussians (a bit trickier)



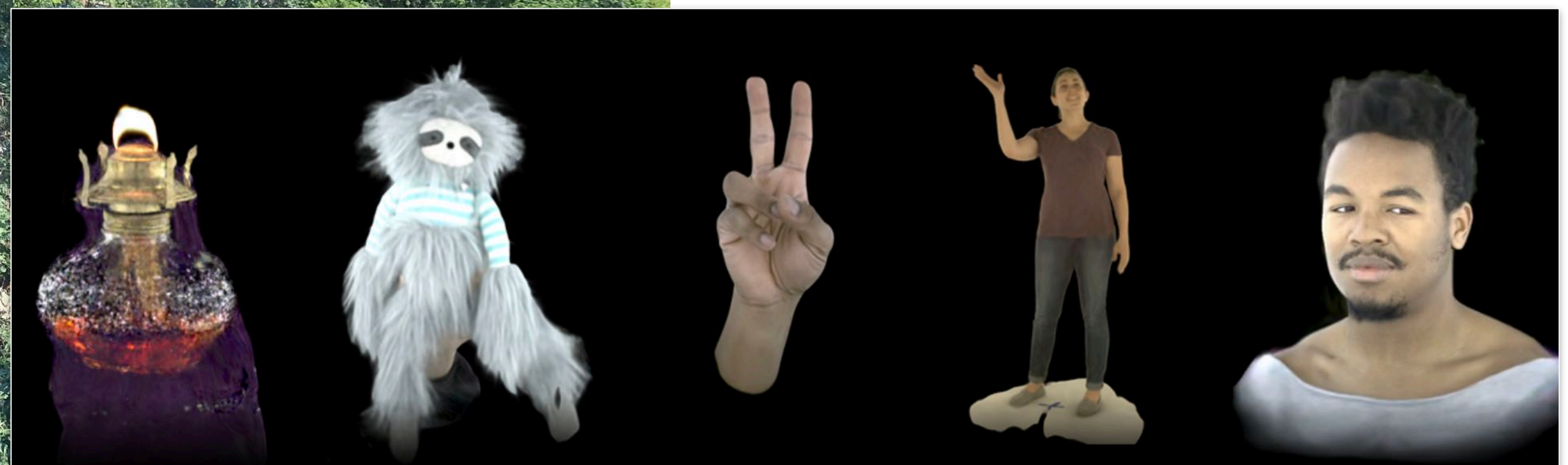
Key idea: differentiable Gaussian splatting rendering to compute $dC/d(\text{color})d(\text{radius})d(\text{location})$

See “3D Gaussian Splatting for Real-Time Radiance Field Rendering” [Kerbl 2023]

Novel View Synthesis



But it's hard to accurately estimate depth or geometry



Discussion: what have we learned?

- **Key idea 1: “unreasonable effectiveness of large-scale optimization”**
 - **High-performance optimization can recover parameter values for complex parameterized models**
 - **Credit: Ren Ng for this perspective**
- **Key idea 2: Many different scene representations can be reconstructed**
 - **Differentiable rendering of these representations is the key technology**
- **There’s a huge “art” to getting optimization to work**
 - **I doubt I could get these things to successfully optimize without a lot of practice and learning myself!**
 - **If I was a early career graphics student, I’d want to become very accomplished in the “art” of getting an optimizer to work for me**
- **Neural representations != preferred representations: neural data compression can be a good thing**
 - **But techniques like Gaussian splatting, sparse voxels, and Plenoxels are strong evidence that even better compact representations are already present (and don’t require resorting to neural representations)**