

MVS & Neural Radiance Fields



Video from the original ECCV'20 paper

CS180/280A: Intro to Computer Vision and Computational
Photography
Angjoo Kanazawa and Alexei Efros
UC Berkeley Fall 2023

Logistics

- Project 4 due tonight! Good luck!

Multi-View Stereo

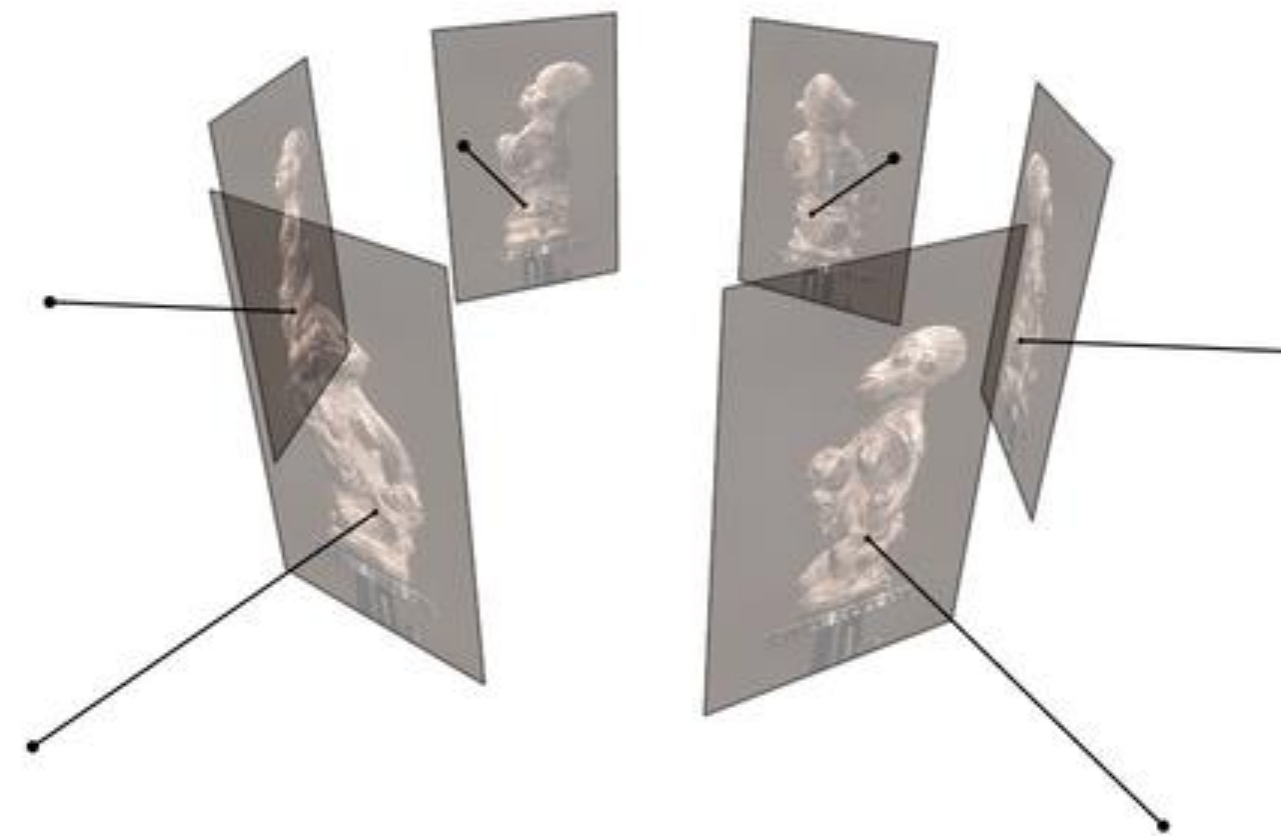
What if we want solid models?



Slide credit: Noah Snavely

Multi-view Stereo (Lots of calibrated images)

- Input: calibrated images from several viewpoints (known camera: intrinsics and extrinsics)
- Output: 3D Model



Figures by Carlos Hernandez

Slide credit: Noah Snavely

In general, conducted in a controlled environment with multi-camera setup that are all calibrated

Whistle in the Form of Female Figure *600 AD - 900 AD*



Details Los Angeles County Museum of Art



Los Angeles County Museum of Art



Sculpture



Mexico

Share

Compare

Saved

Discover

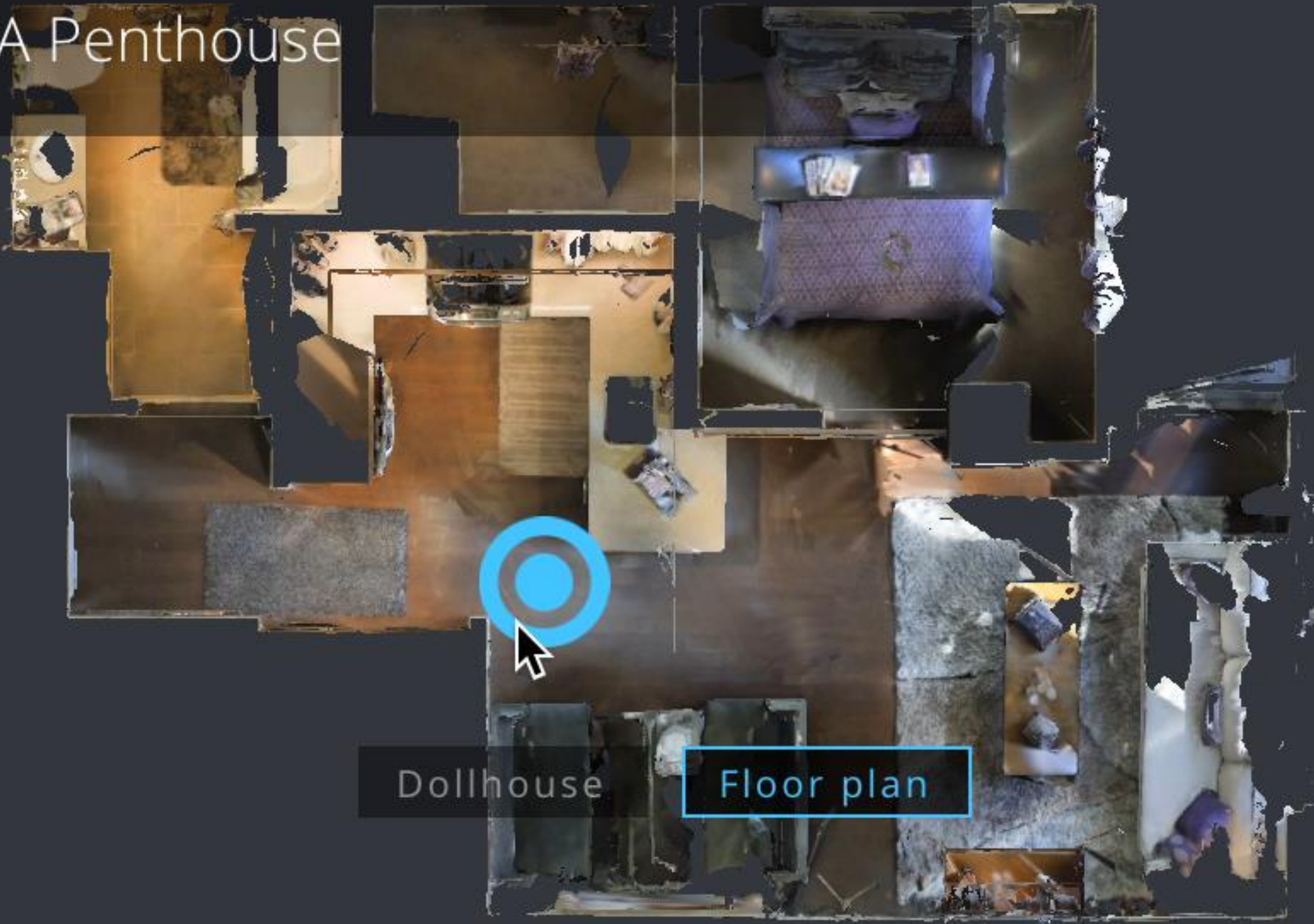
Google

Slide credit: Noah Snavely



< 1BR, 1BA Penthouse

Terms



Dollhouse

Floor plan

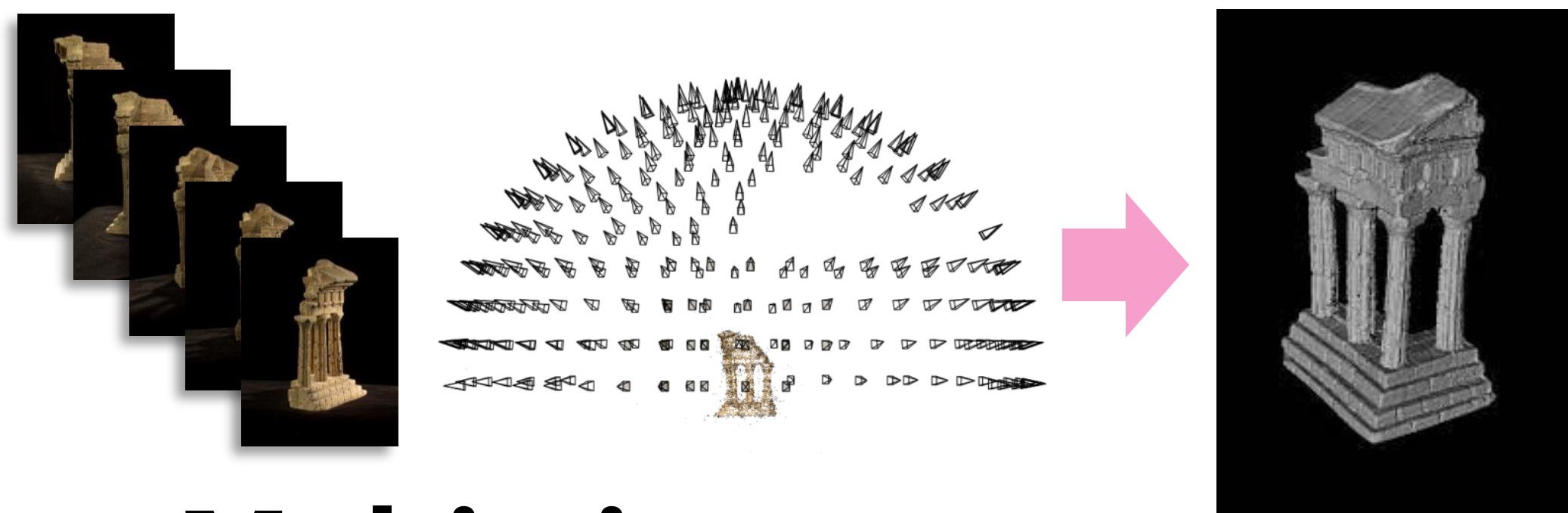


Multi-view Stereo

Problem formulation: given several images of the same object or scene, compute a representation of its 3D shape



Binocular Stereo



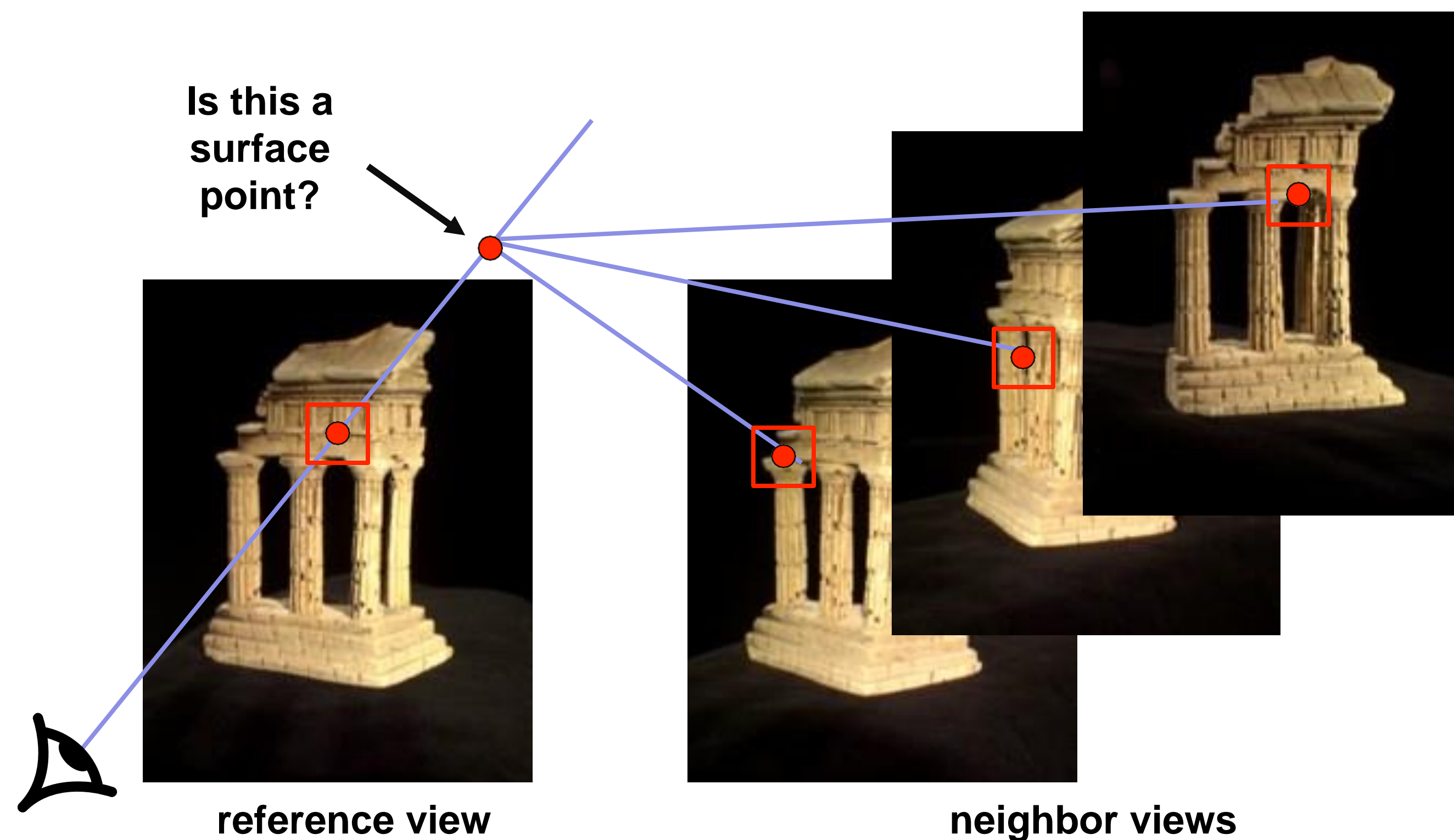
Multi-view stereo

Examples: Panoptic studio



<http://domedb.perception.cs.cmu.edu/>

Multi-view stereo: Basic idea



Source: Y.
Furukawa

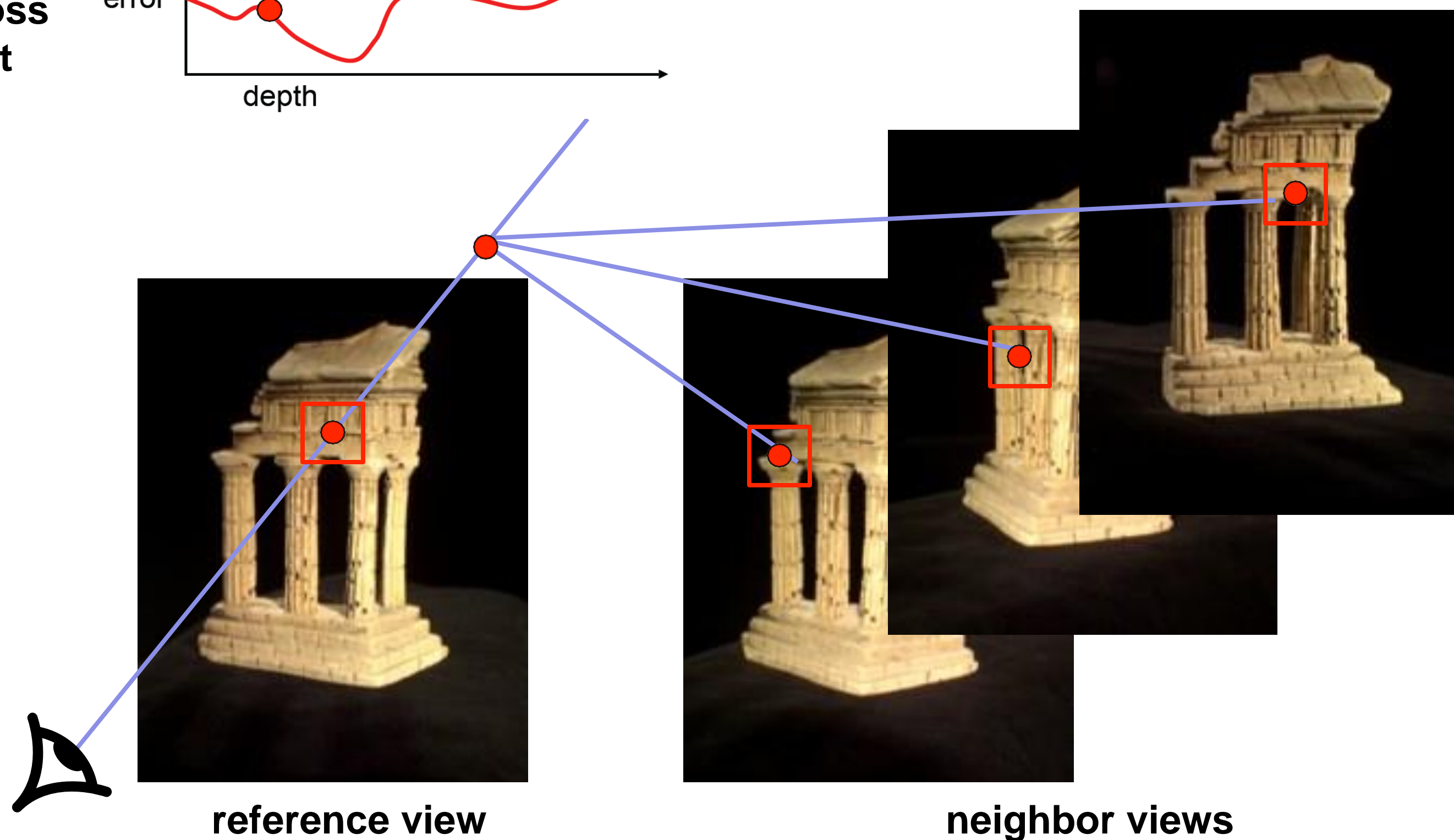
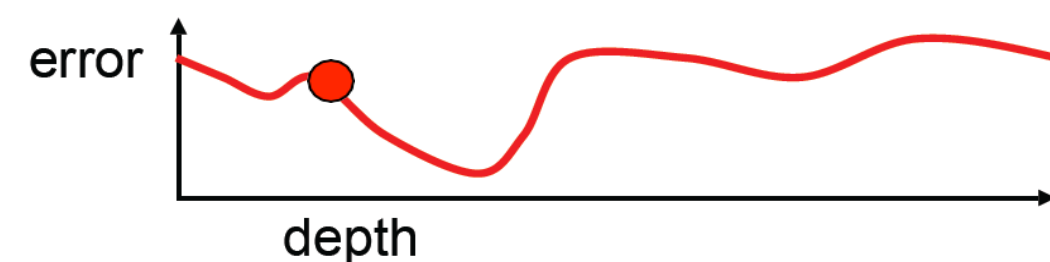
Multi-view stereo: Basic idea

Evaluate the likelihood of geometry at a particular depth for a particular reference patch:



Multi-view stereo: Basic idea

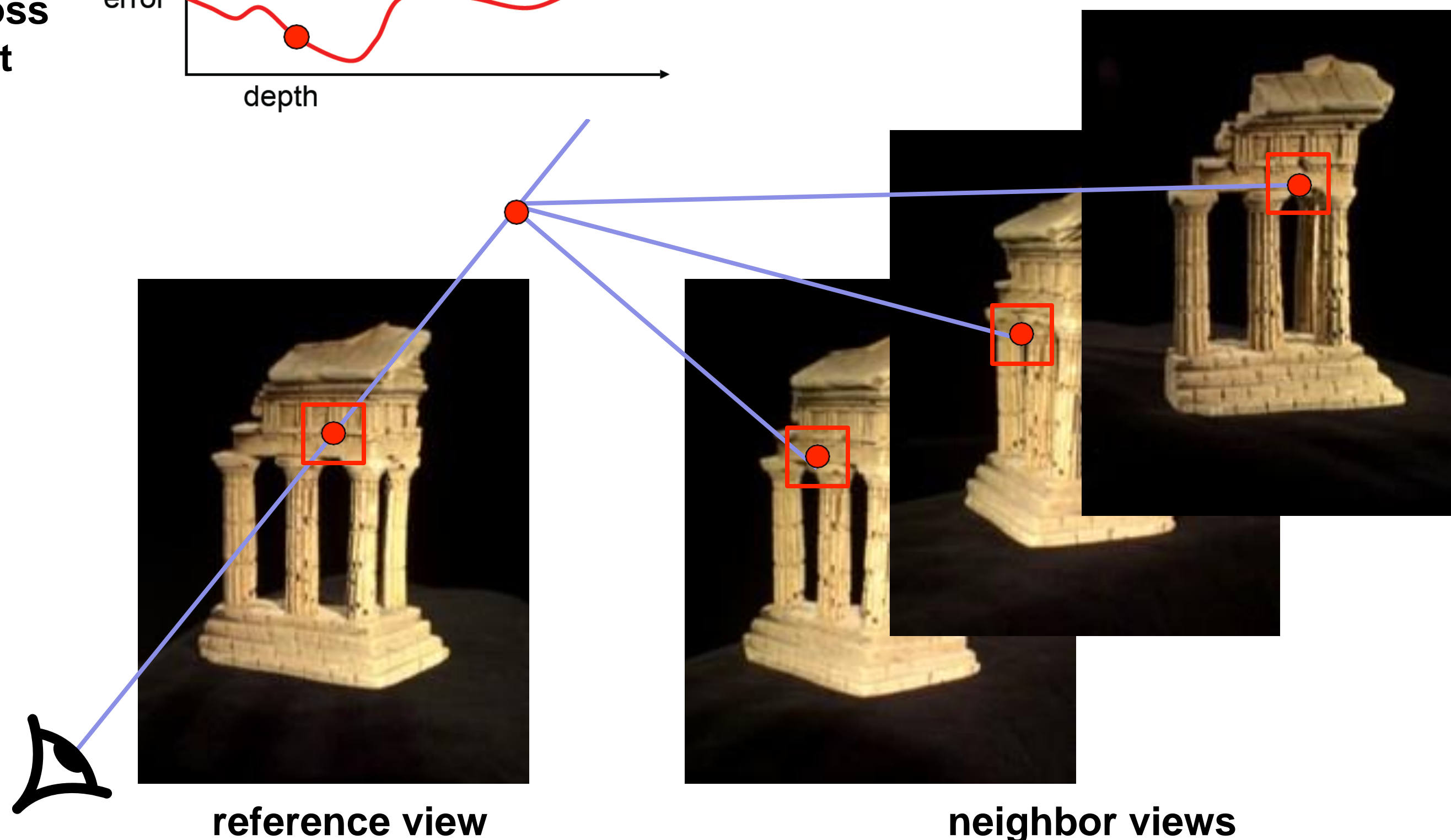
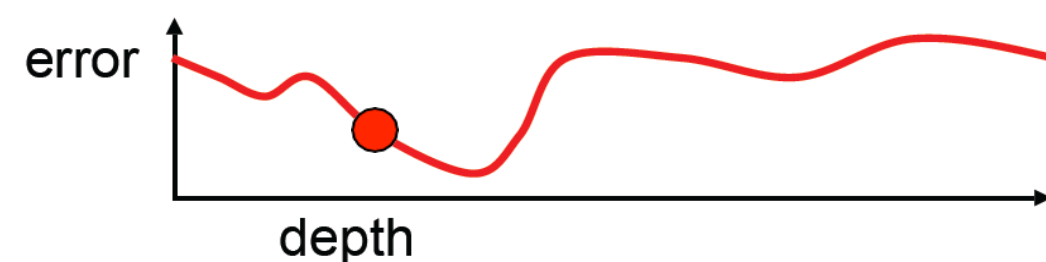
Photometric error across different depths



Source: Y. Furukawa

Multi-view stereo: Basic idea

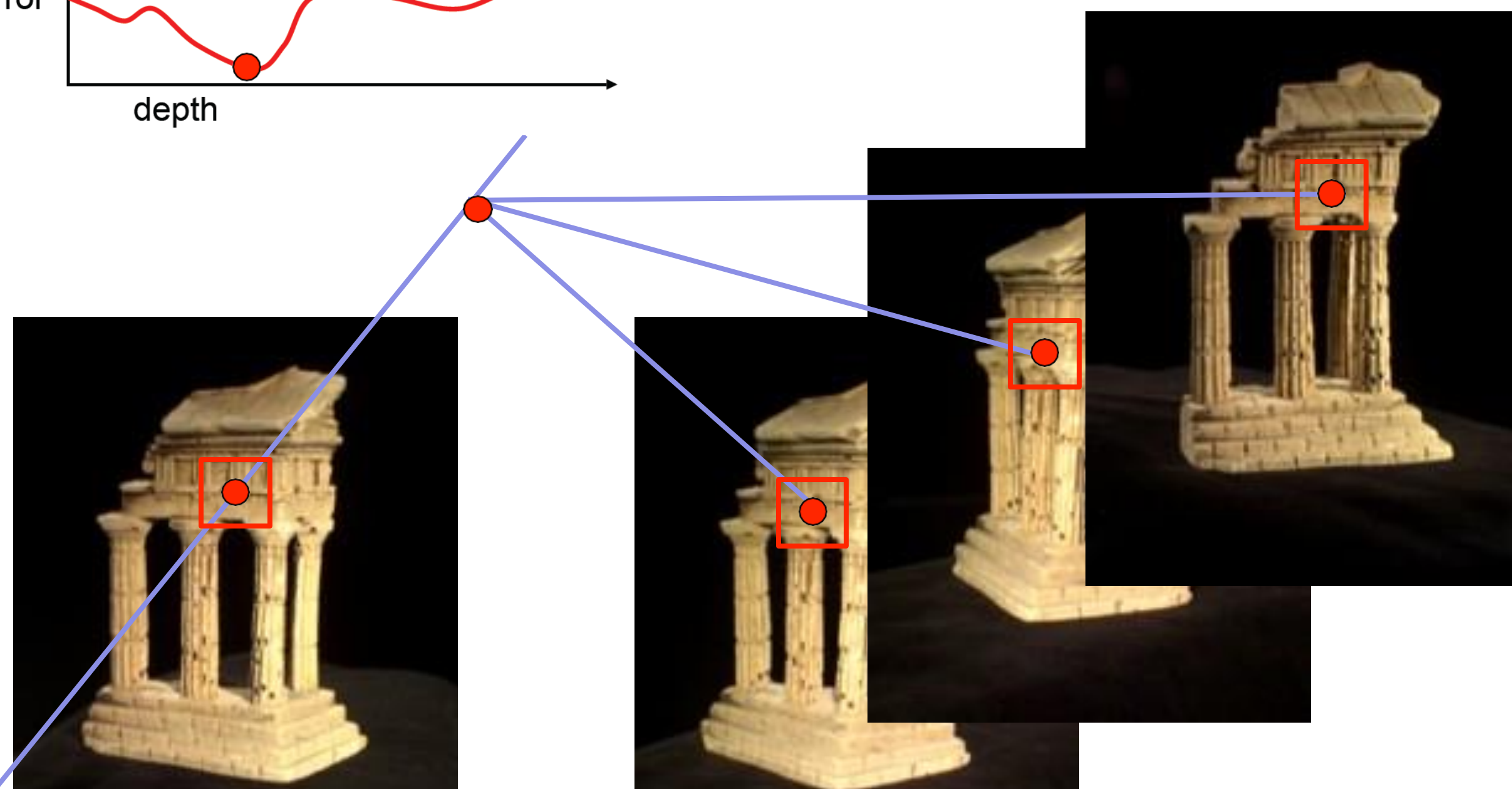
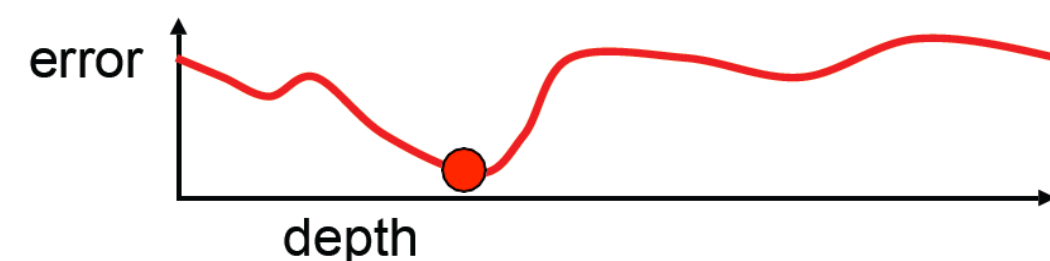
Photometric error across different depths



Source: Y. Furukawa

Multi-view stereo: Basic idea

Photometric error across different depths

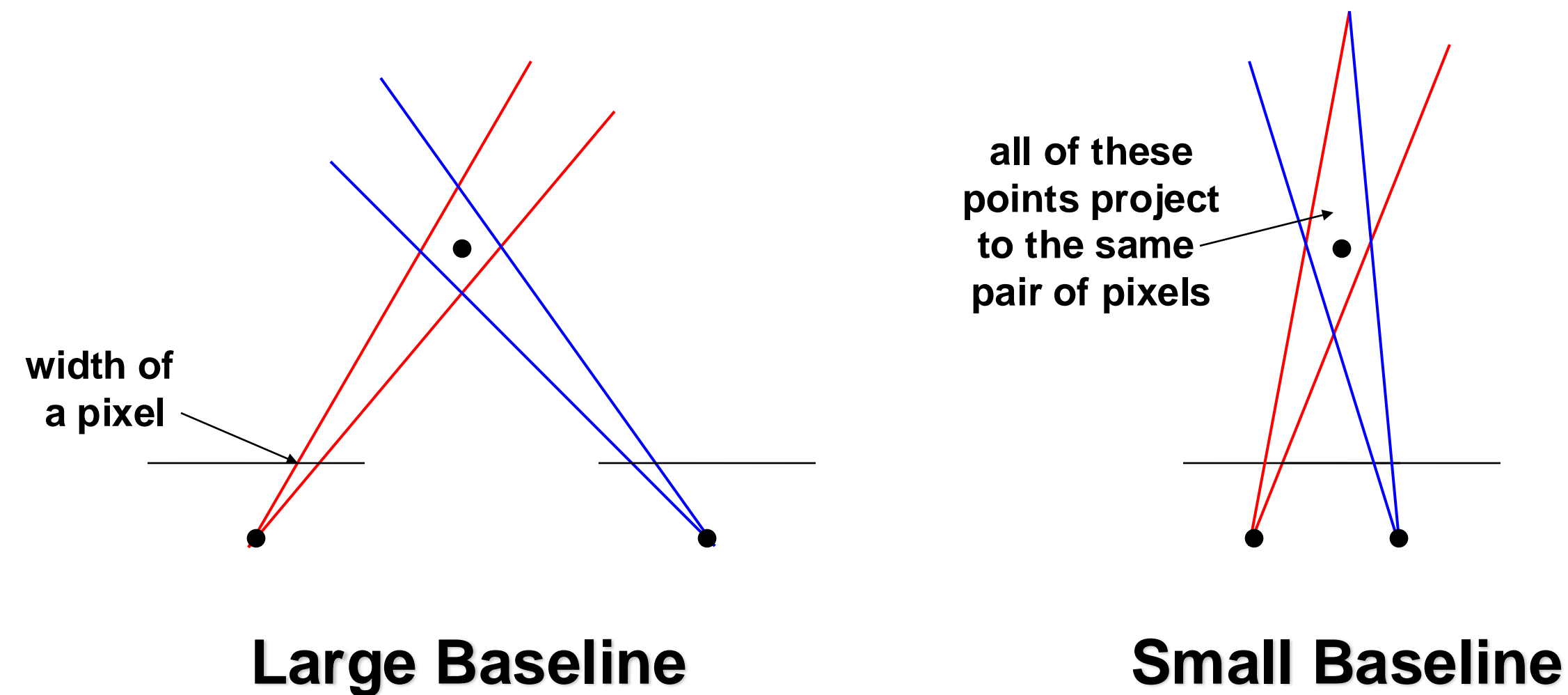


In this manner, solve for a depth map over the whole reference view

Multi-view stereo: advantages over 2 view

- Can match windows using more than 1 other image, giving a **stronger match signal**
- If you have lots of potential images, can **choose the best subset** of images to match per reference image
- Can reconstruct a depth map for each reference frame, and the merge into a **complete 3D model**

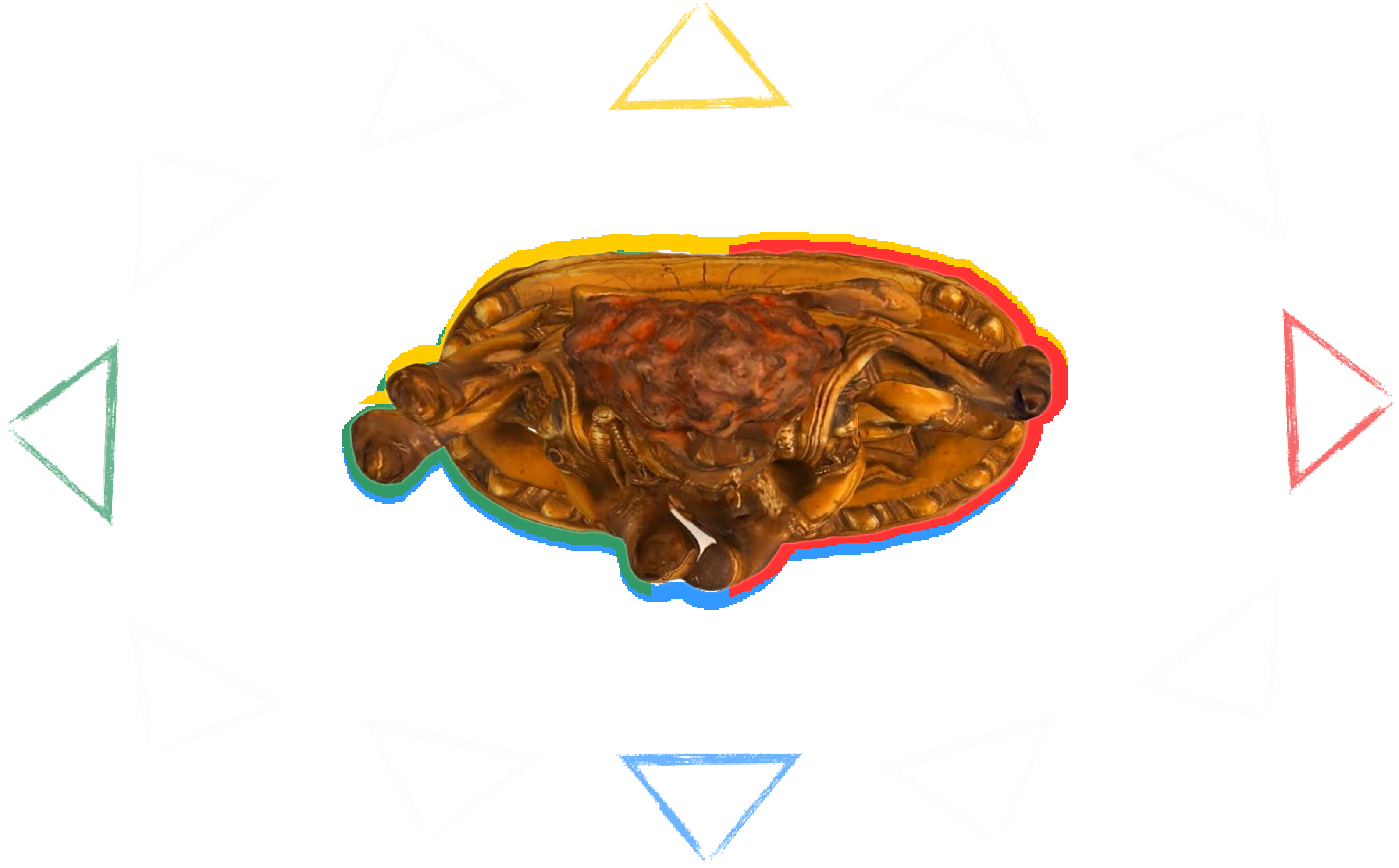
Choosing the baseline



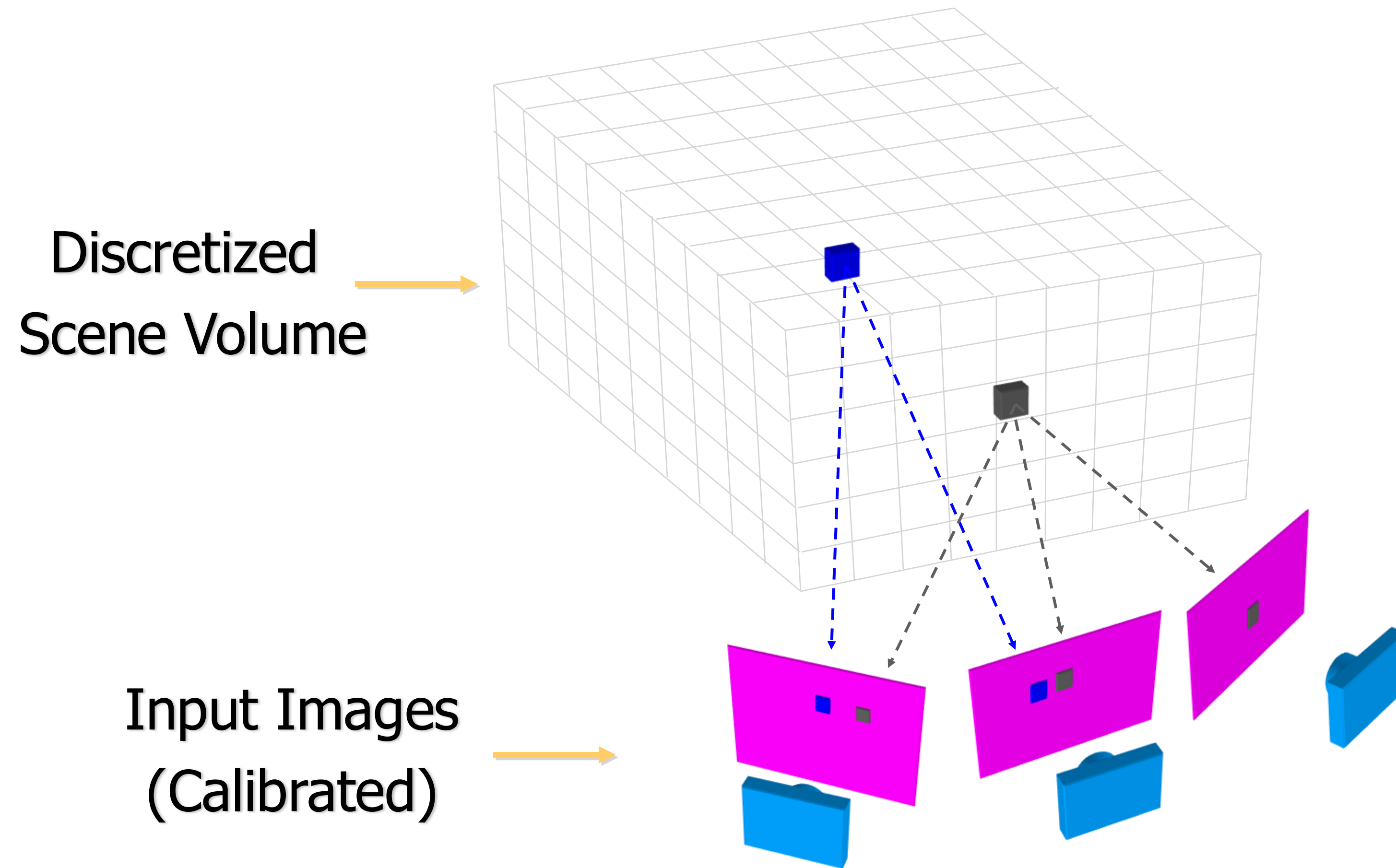
- What's the optimal baseline?
 - Too small: large depth error
 - Too large: difficult search problem





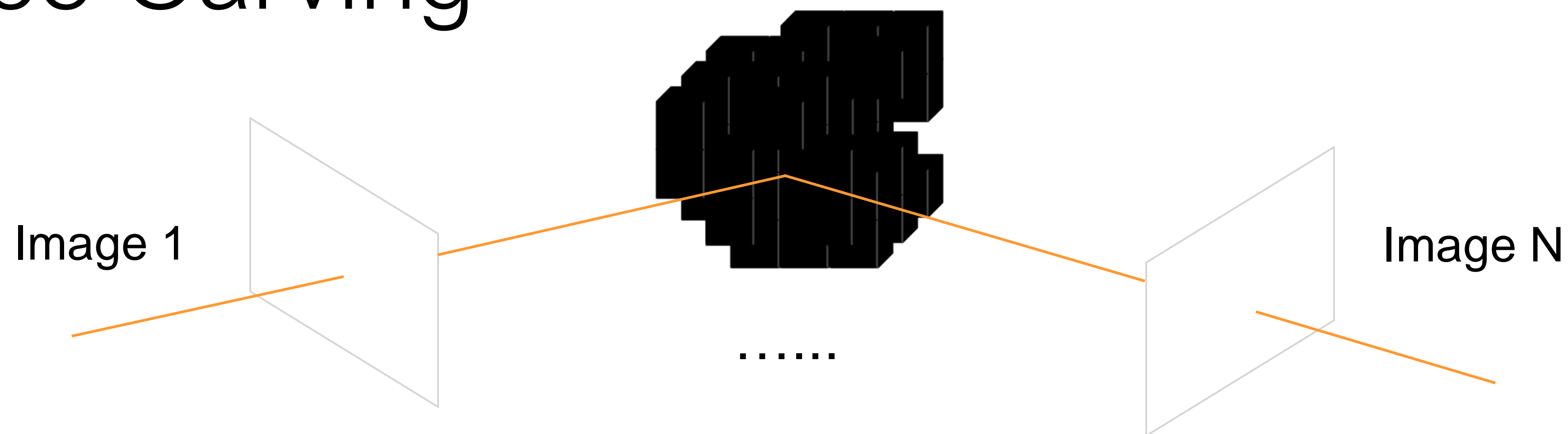


Volumetric stereo



Goal: Assign RGB values to voxels in V
photo-consistent with images

Space Carving



•Space Carving Algorithm

- Initialize to a volume V containing the true scene
- Choose a voxel on the outside of the volume
- Project to visible input images
- Carve if not photo-consistent
- Repeat until convergence

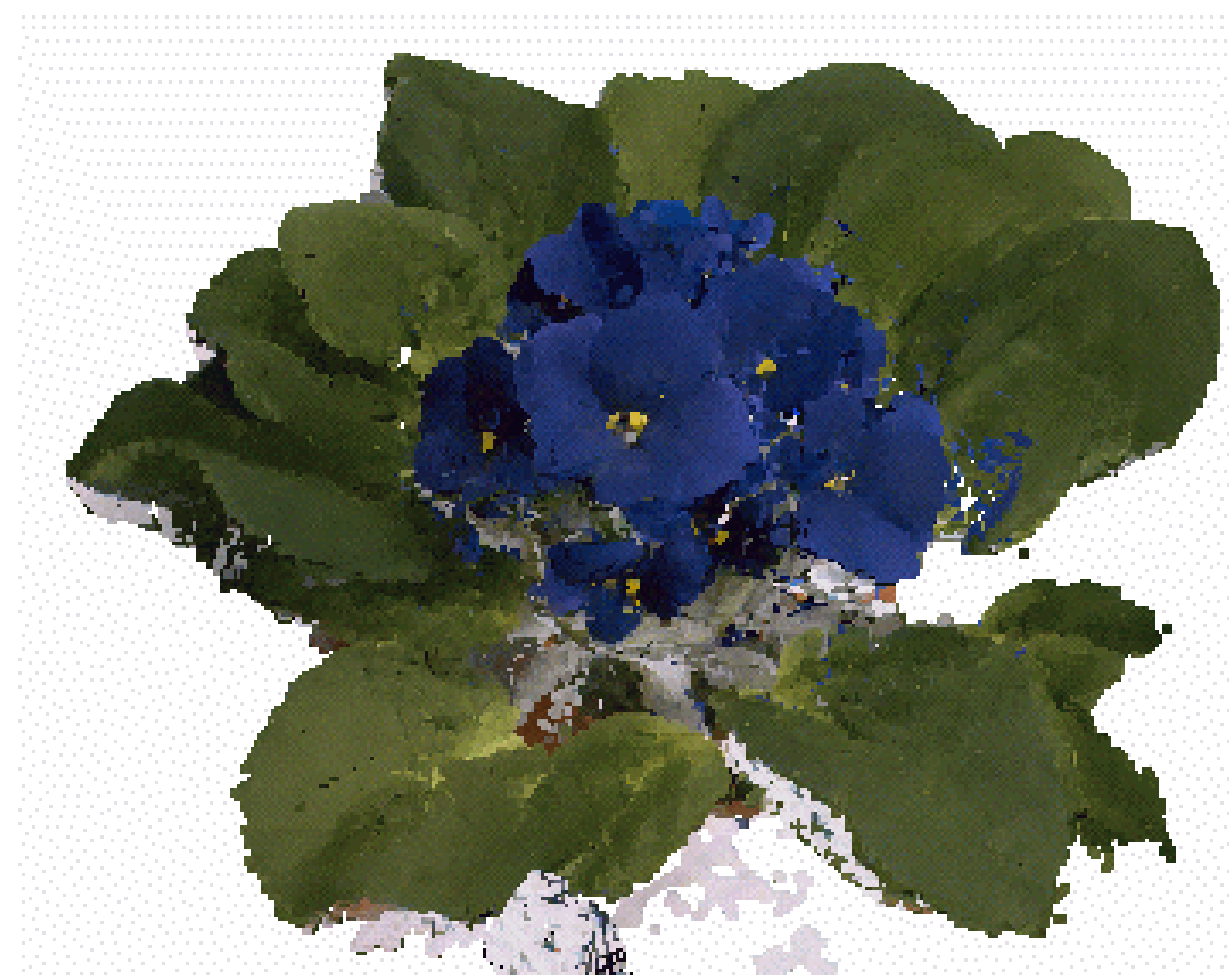
Space Carving Results



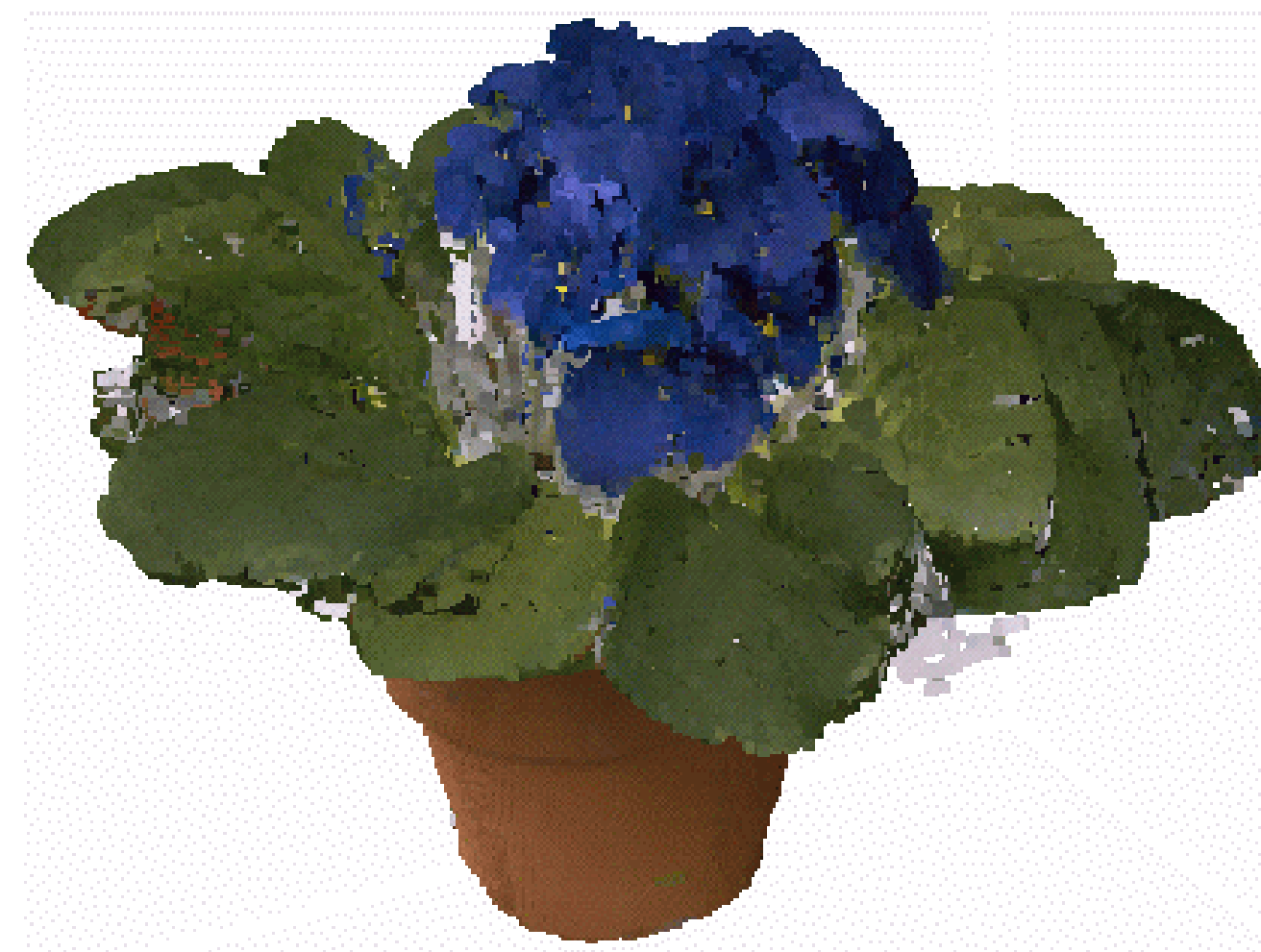
Input Image (1 of 45)



Reconstruction



Reconstruction



Reconstruction

Space Carving Results



**Input Image
(1 of 100)**



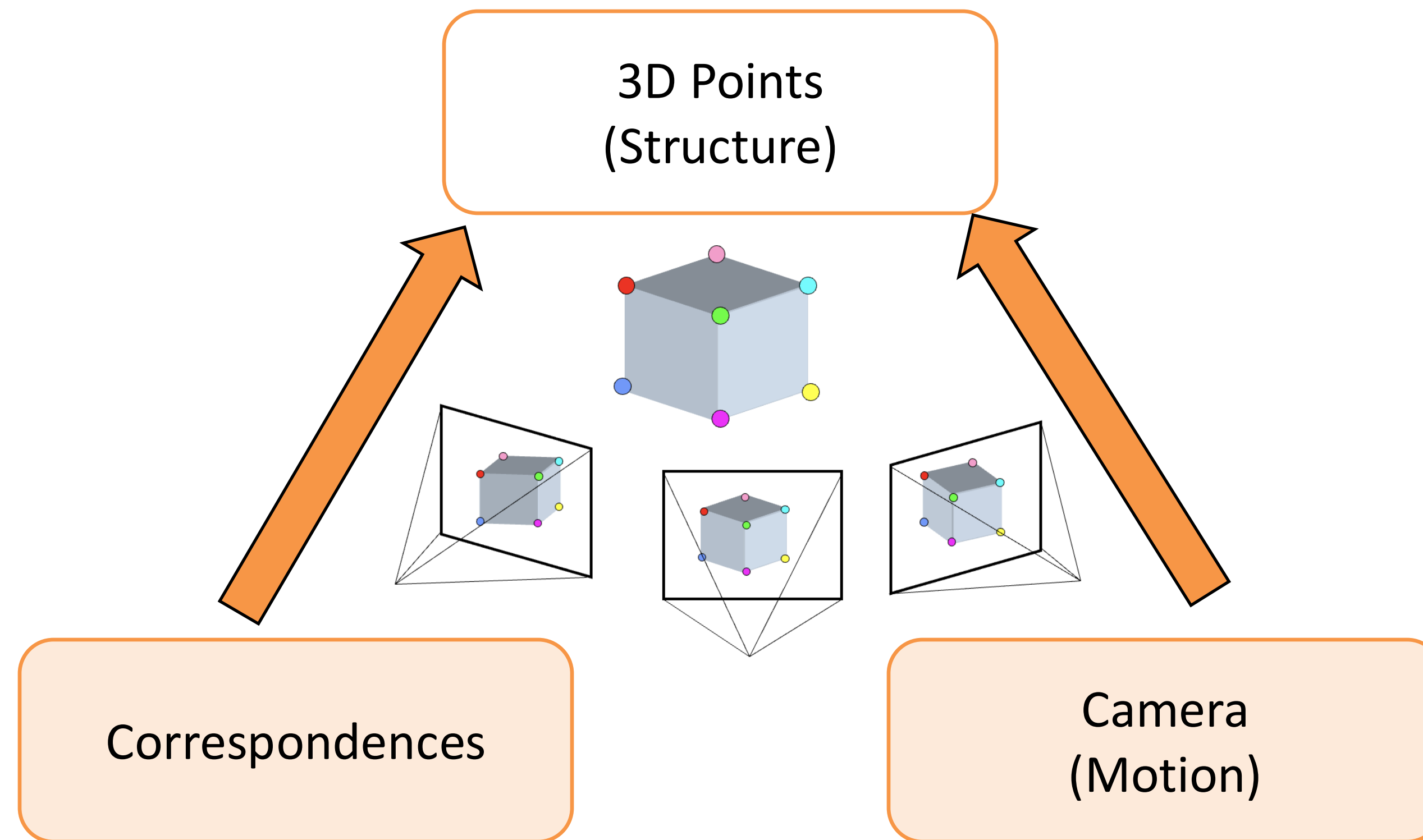
Reconstruction

Tool for you: COLMAP

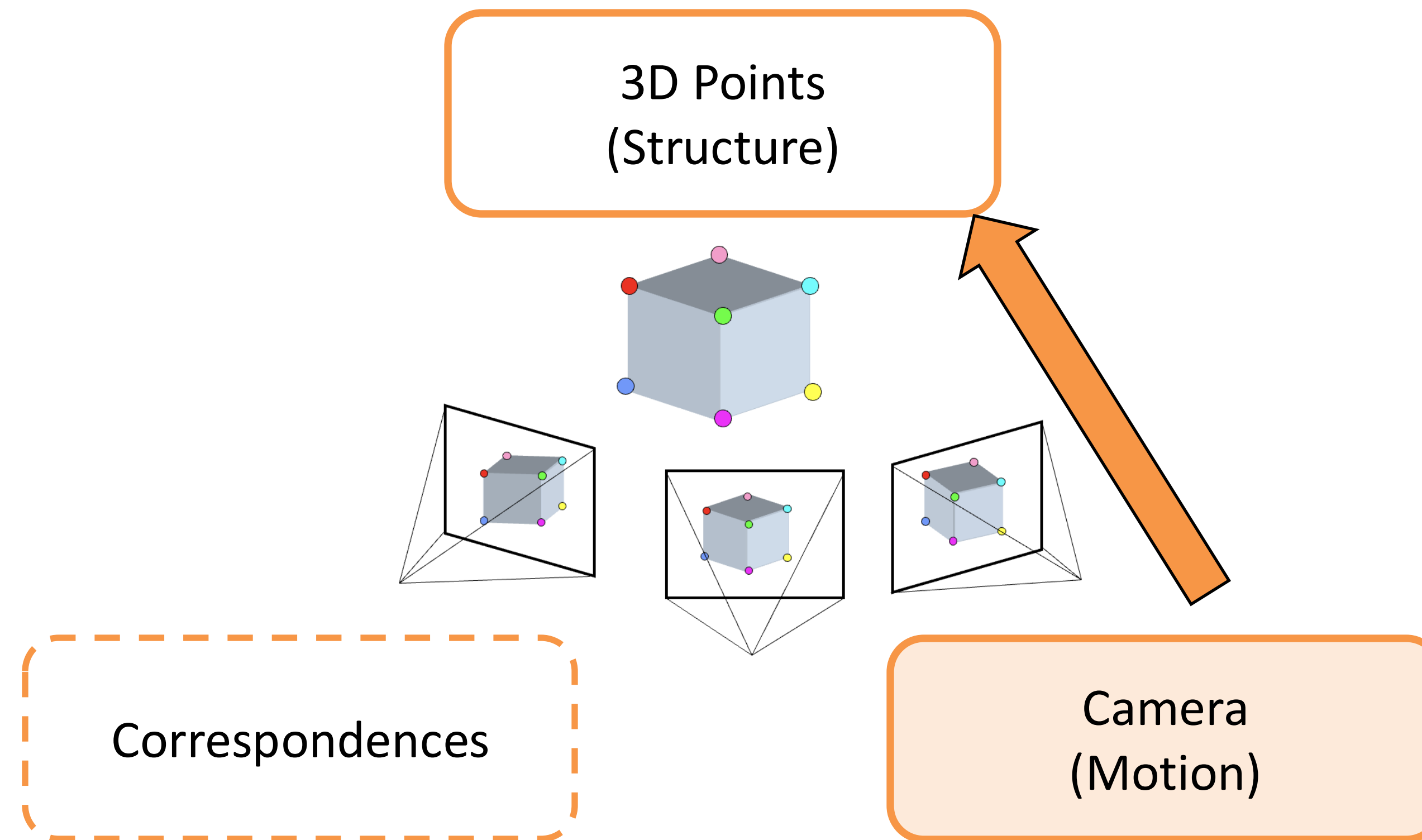
<https://github.com/colmap/colmap>

A general SfM + MVS pipeline

Multi-View Stereo



Volumetric “Neural” Rendering



**Does not use explicit correspondences,
relies on reconstruction loss (Analysis-by-Synthesis)**

Neural Radiance Fields



Video from the original ECCV'20 paper

Capturing Reality



Earliest cave painting (45,500 years old) in Sulawesi, Indonesia

Capturing Reality



Monet's Cathedral series: study of light 1893-1894

Capturing Reality



First self-portrait Cornelius 1839



First Movie - Muybridge 1878

Capturing Reality – in 3D



Capturing Reality – in 3D (MVS – last lecture)



Google Earth 2016~

What is next?

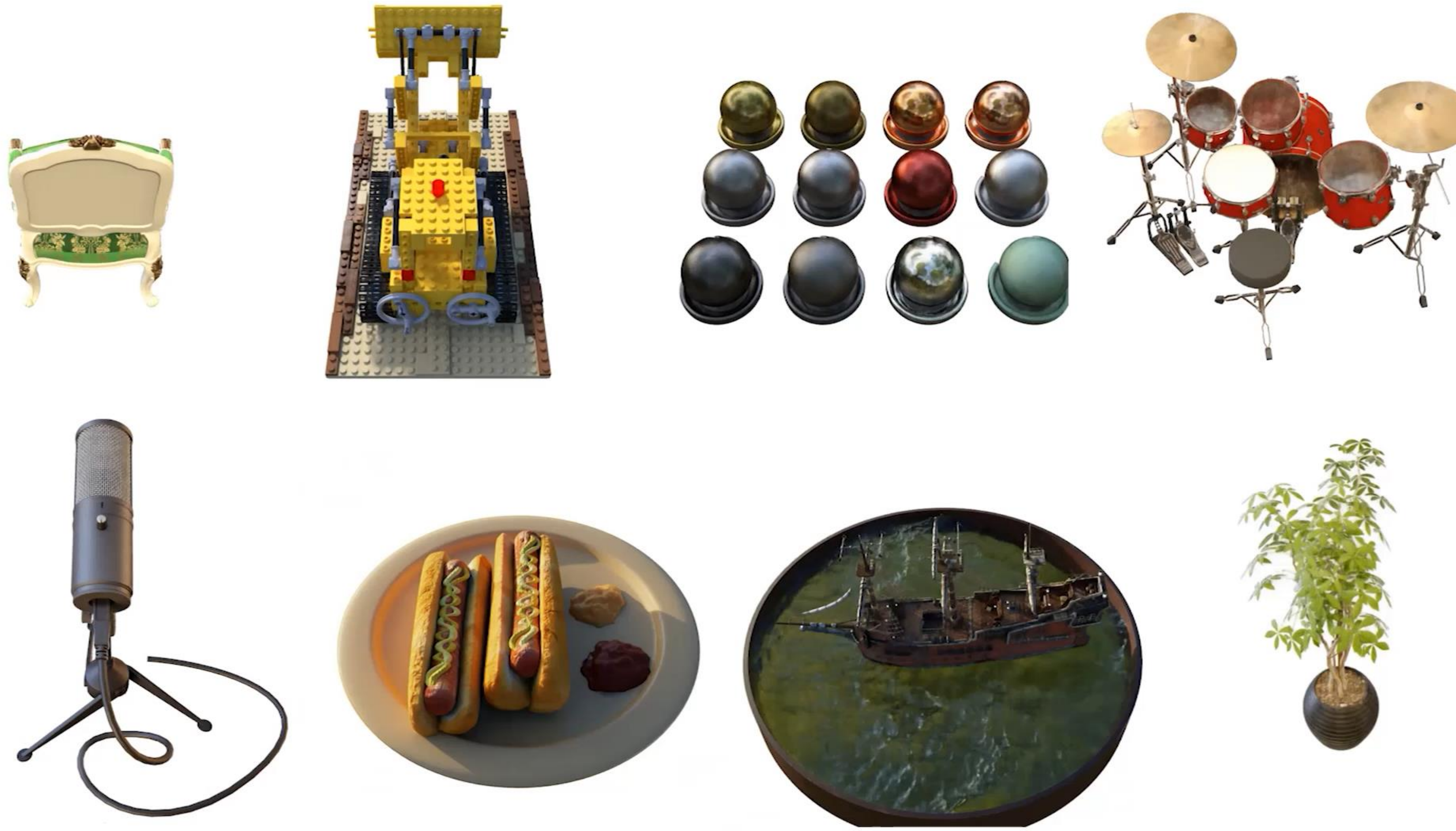
2020: Neural Radiance Field (NeRF)



Mildenhall*, Srinivasan*, Tancik*, Barron, Ramamoorthi, Ng, ECCV 2020

It has been three years

- Original NeRF paper: 4200+ citations in 3 years



Handling Appearance Changes



Real-time Rendering



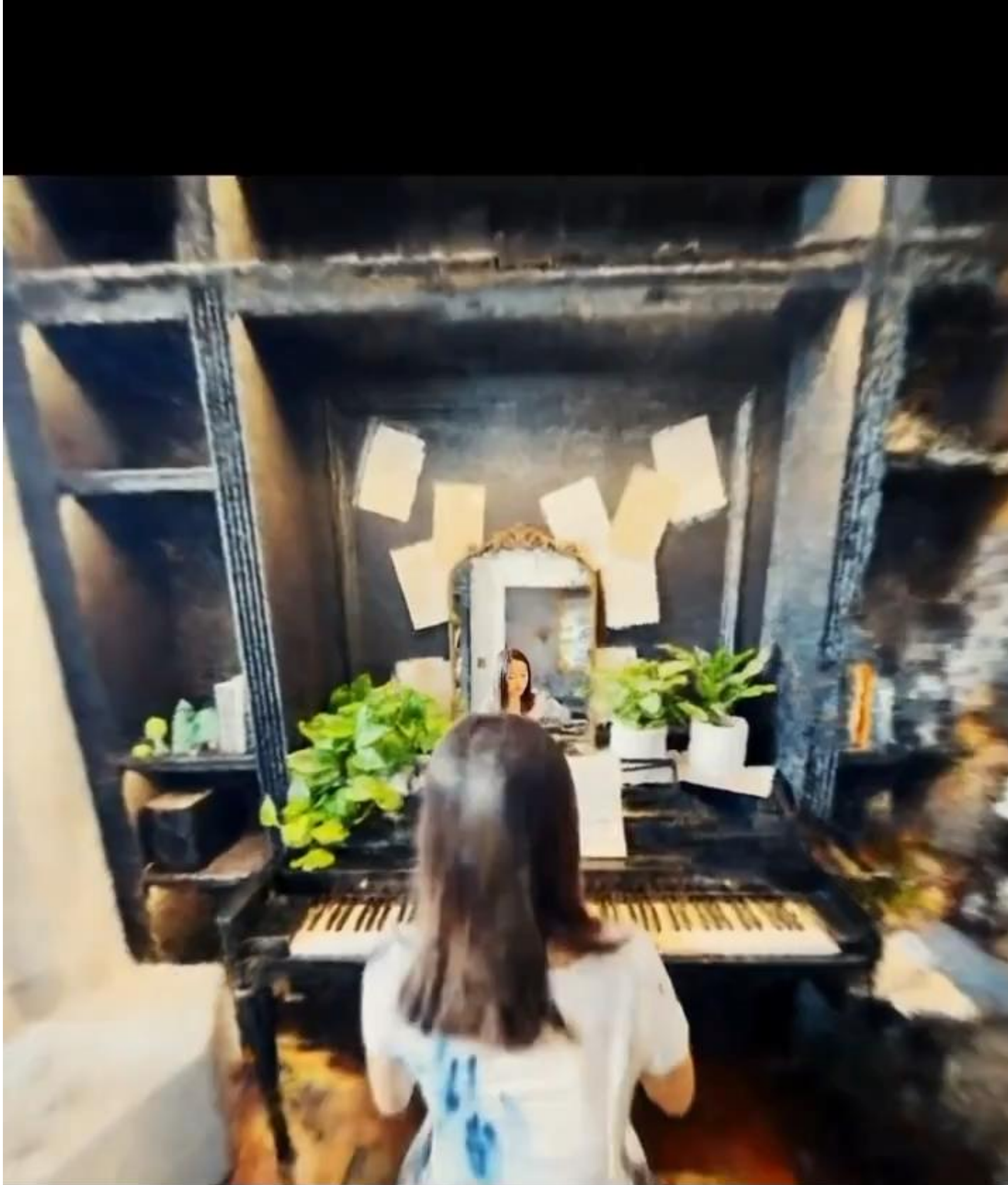
Real-time Inference

INSTANT NEURAL GRAPHICS PRIMITIVES WITH A MULTIRESOLUTION HASH ENCODING

Thomas Müller Alex Evans Christoph Schied Alexander Keller

<https://nvlabs.github.io/instant-ngp>





@karenxcheng, with
InstantNGP [Müller et
al., SIGGRAPH 2022]

VIEWPORT RENDER VIEW

RESUME TRAINING

Show Scene

Show Images

Refresh Page

Resolution: 640x1024px

Time Allocation: 100% spent on viewer

Server Connected

Render Connected

CONTROLS

RENDER

SCENE

LOAD PATH

EXPORT PATH



Height
1080

Width
1920

FOV
50

Seconds
4

FPS
24

ADD CAMERA



Smoothness

0.00

0

1

2

3



CAMERA 0



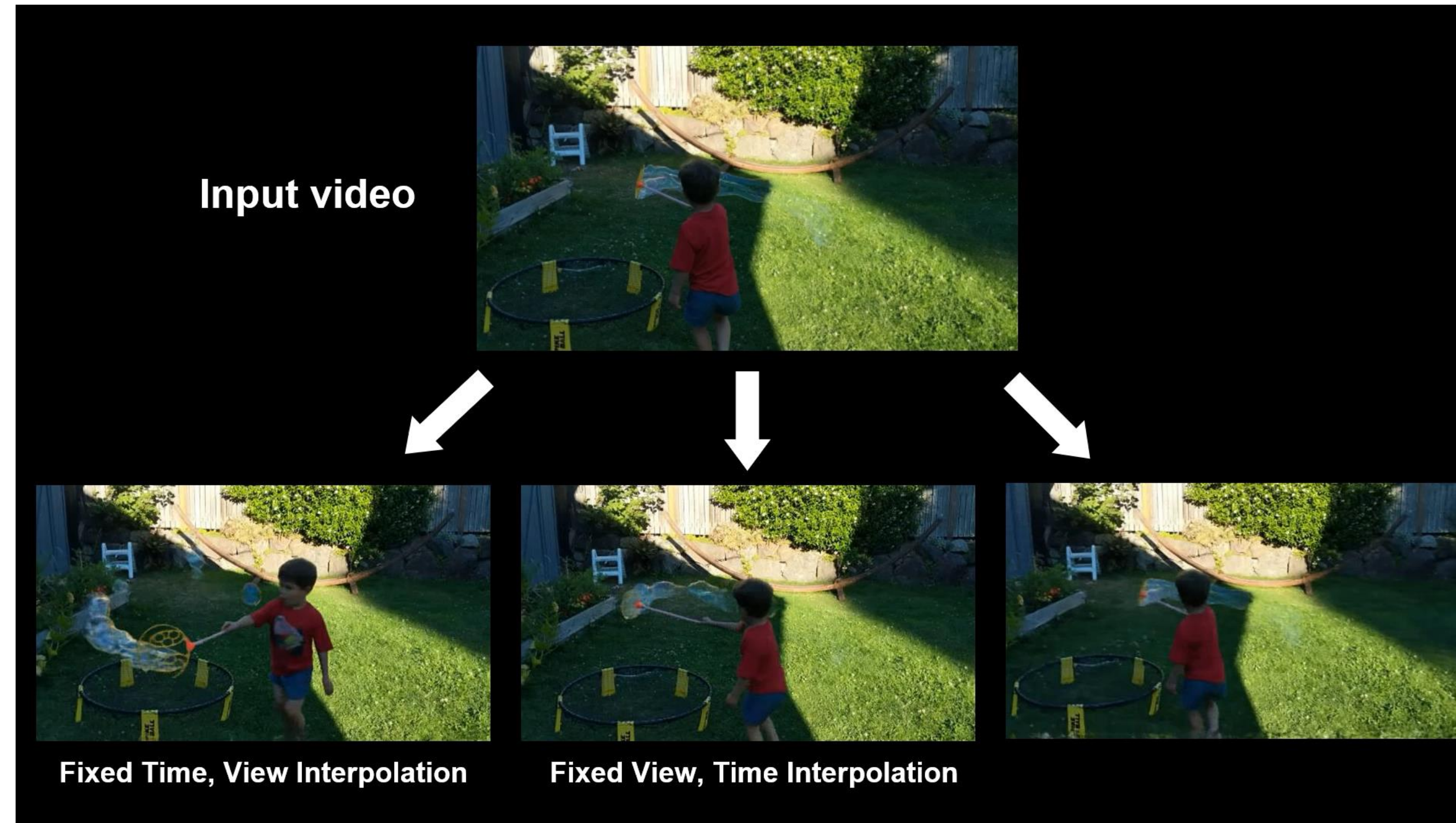
Dynamic NeRFs



[Xian et al., CVPR 2021]



Nerfies [Park et al., ICCV 2021]
HyperNeRF [Park et al., SigAsia 2021]



NSFF [Li et al., CVPR 2021]

Generative 3D Faces



EG3D: Efficient Geometry-aware 3D Generative Adversarial Networks, Chan et al. CVPR 2022



City-Scale NeRFs

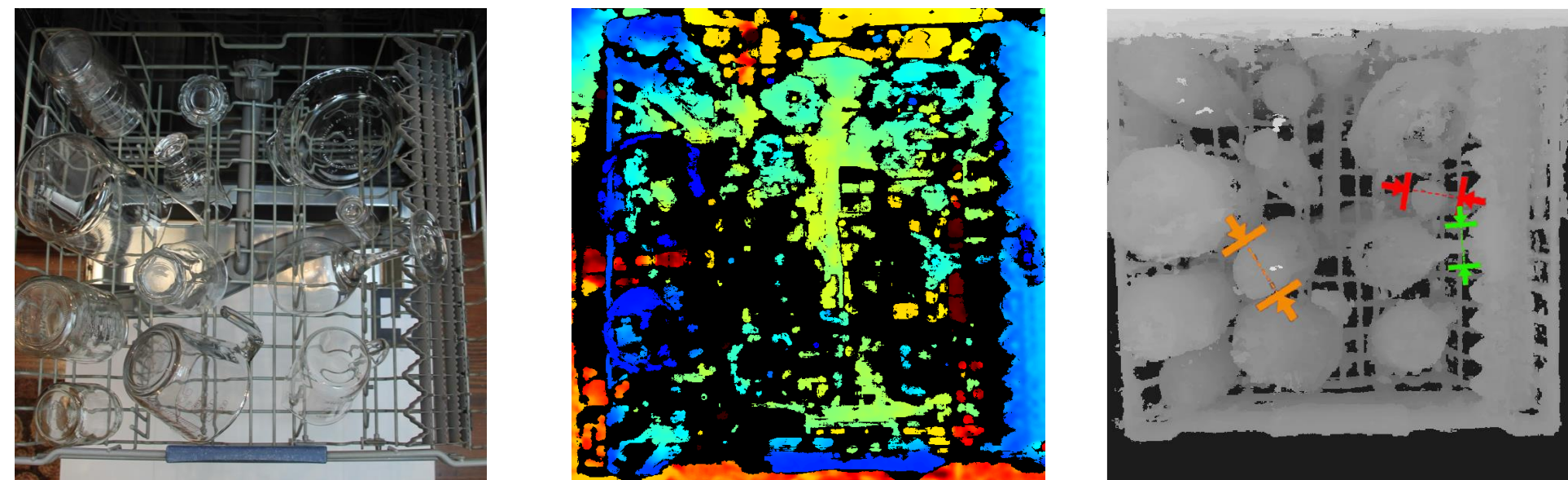
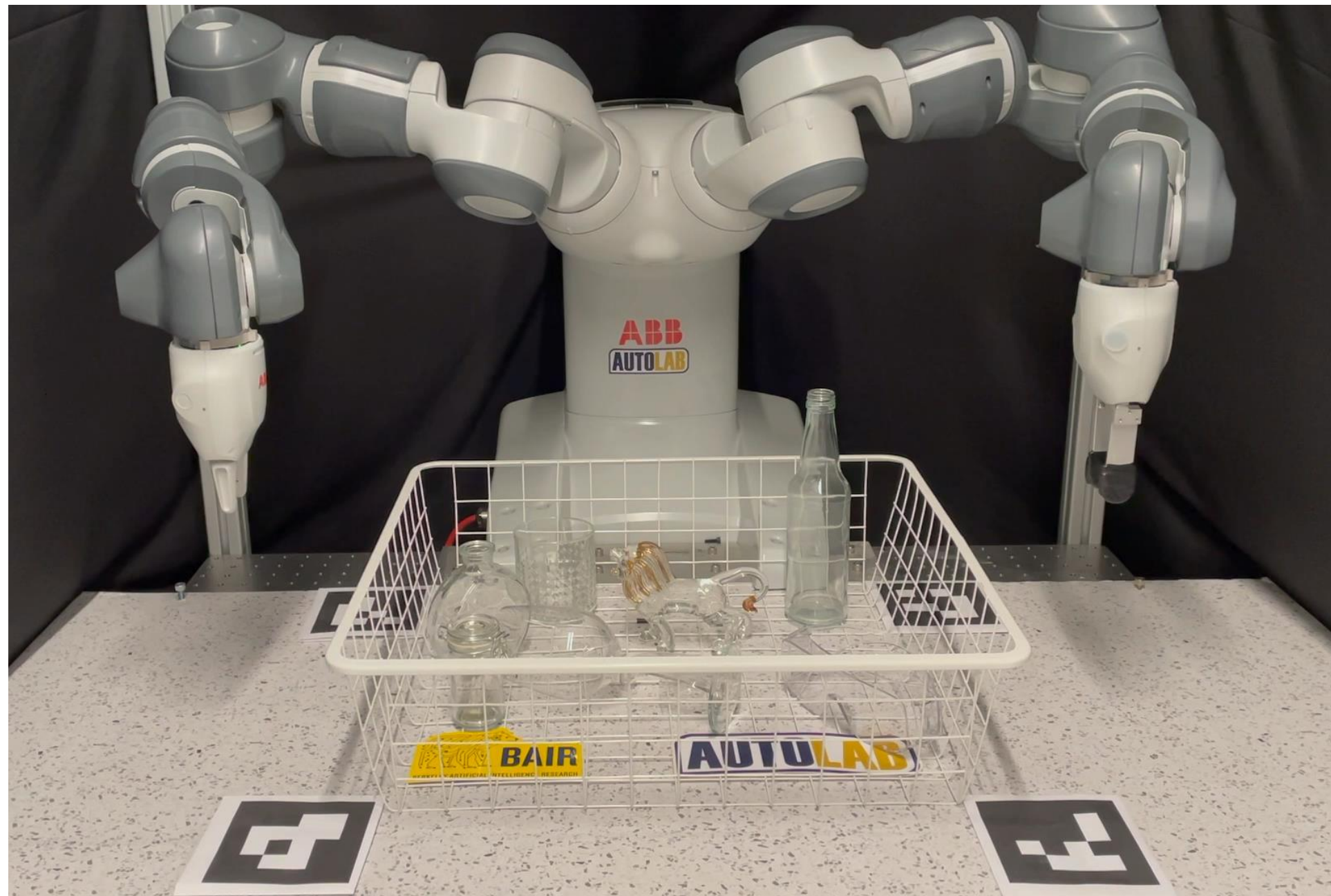


BlockNeRF
[Tancik et al.
CVPR 2022]

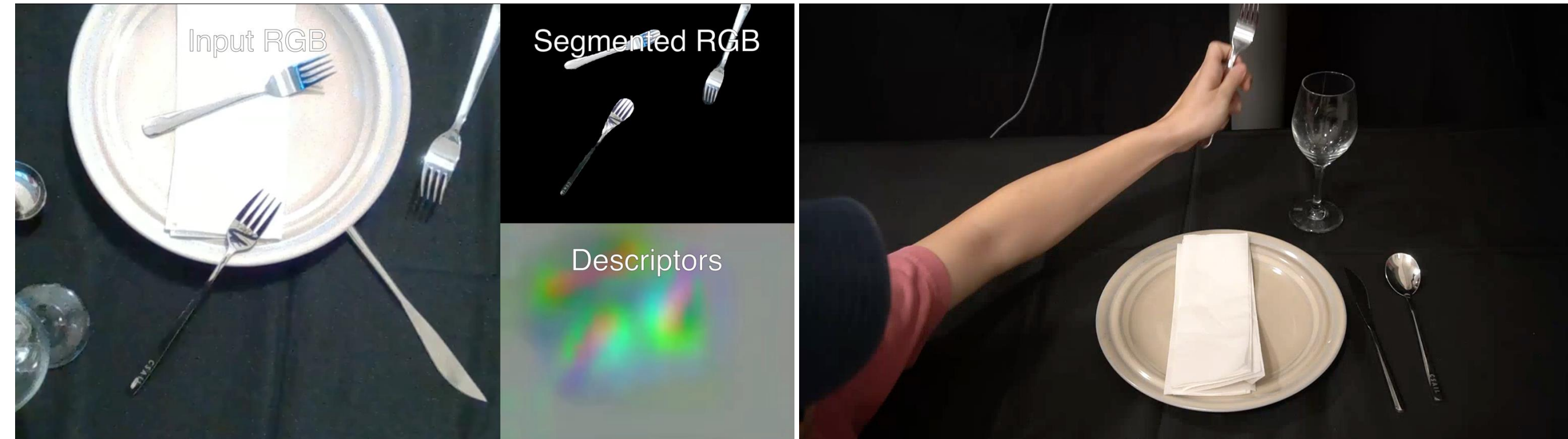


RawNeRF
[Mildenhall et al.
CVPR 2022]

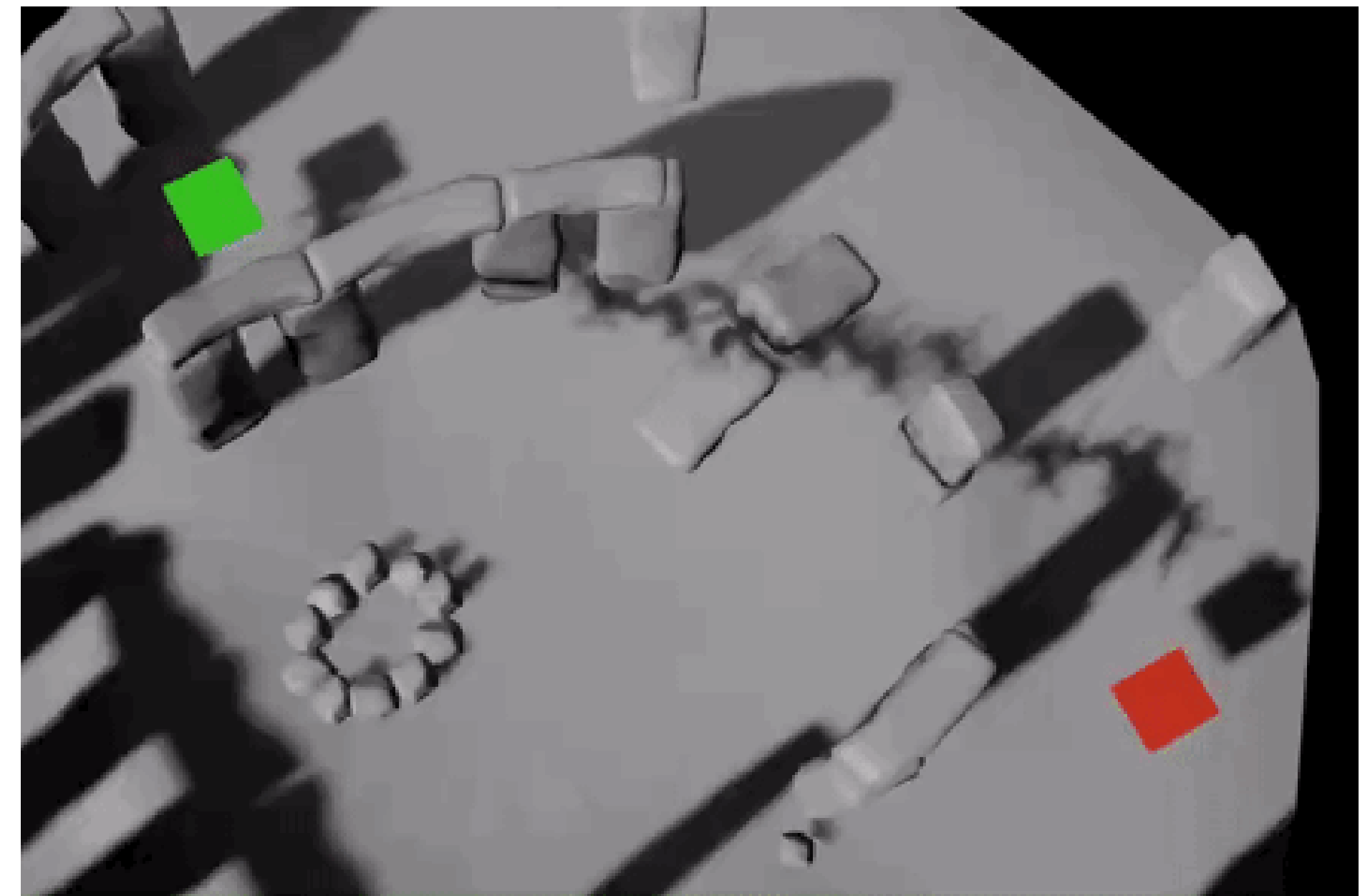
Robotics



Dex-NeRF: Using a Neural Radiance field to Grasp Transparent Objects, [Ichnowski and Avigal et al. CoRL 2021]



NeRF-Supervision: Learning Dense Object Descriptors from Neural Radiance Fields, [Yen-Chen et al. ICRA 2022]



Vision-Only Robot Navigation in a Neural Radiance World [Adamkiewicz and Chen et al. ICRA 2022]

Generating 3D scenes with diffusion models

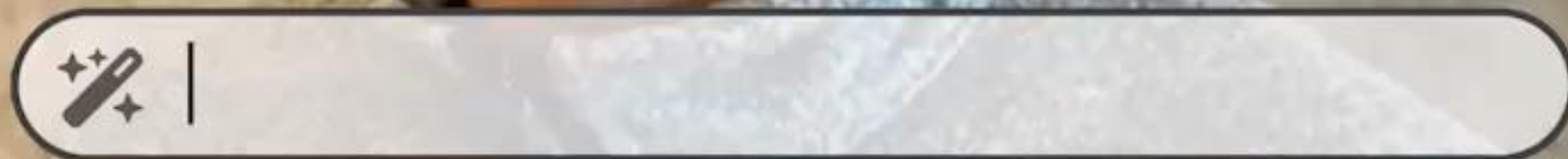


DreamFusion
[Poole et al.
ICLR 2023]

Querying with Language



Editing with Instructions



Goals of the next few lectures

- Visit the fundamentals in Neural Volumetric Rendering by abstracting away recent developments
- Provide first principles + background for you to go and read these papers & play around with the tools
- New Project 5!! Implement these concepts yourself



Capture of UC Berkeley redwoods with

Birds Eye View & Background

Birds Eye View

- What is NeRF?
- How is it different or similar to existing approaches?
- What is its historical context?

Problem Statement

Input:

A set of calibrated Images



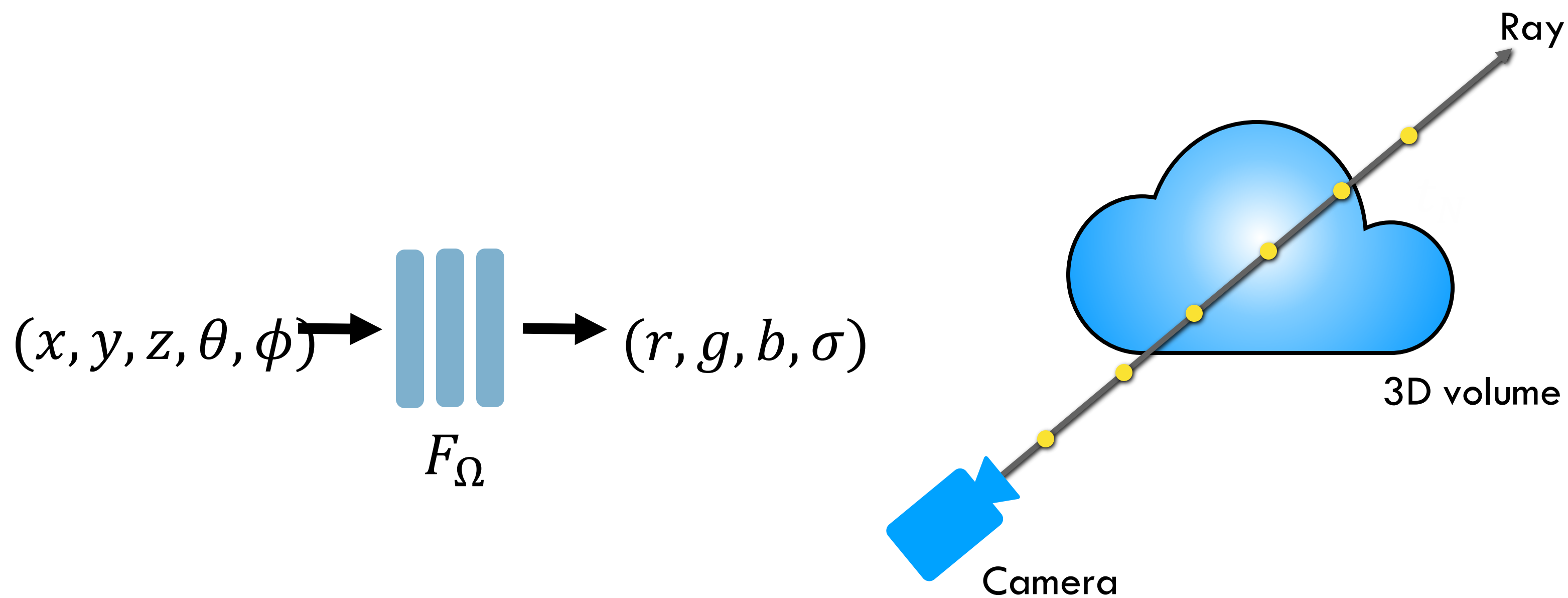
Output:

A 3D scene representation that renders novel views



Three Key Components

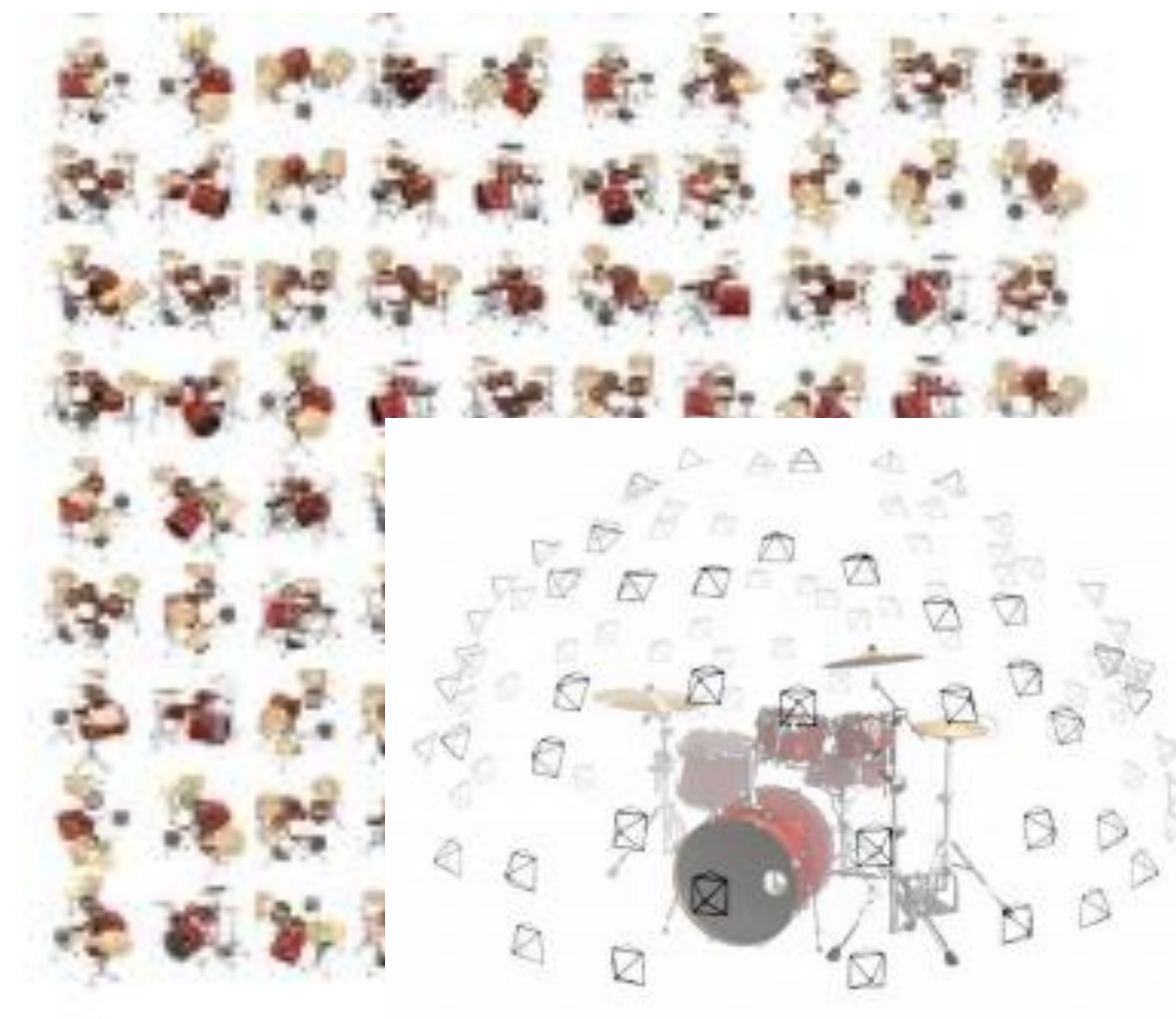
Objective: Synthesize all training views



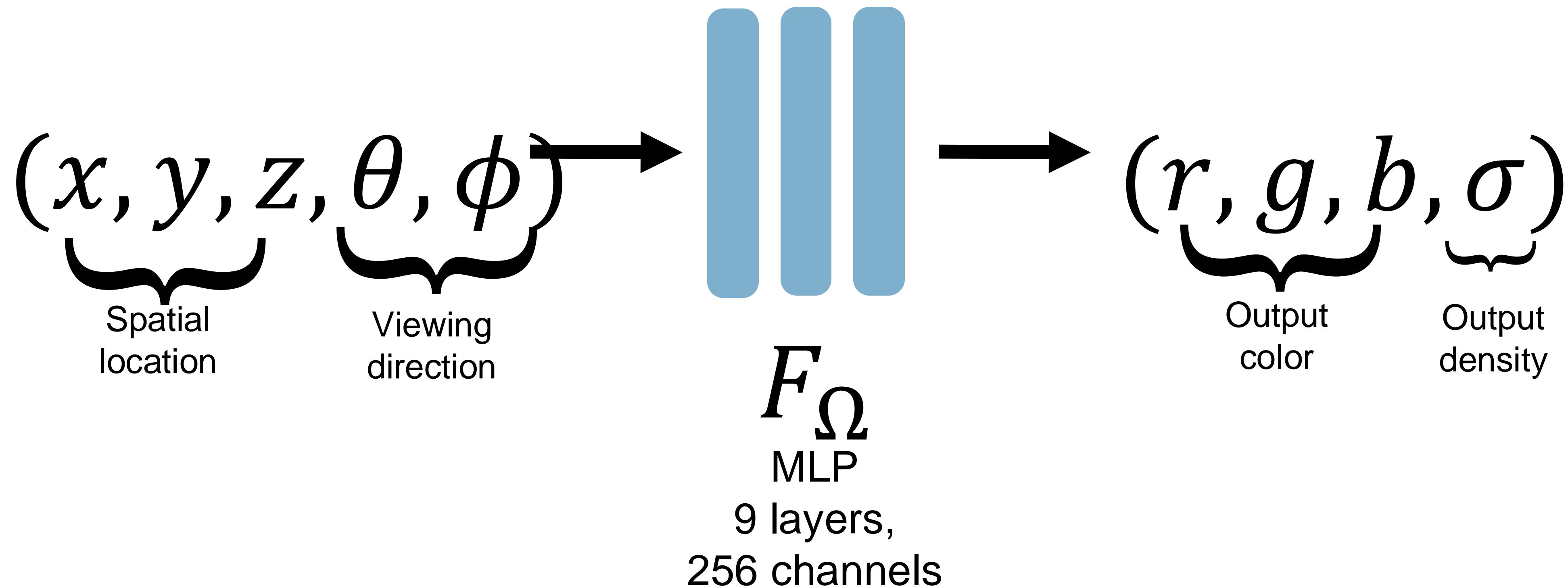
Neural Volumetric 3D Scene Representation

Differentiable Volumetric Rendering Function

Optimization via Analysis-by-Synthesis

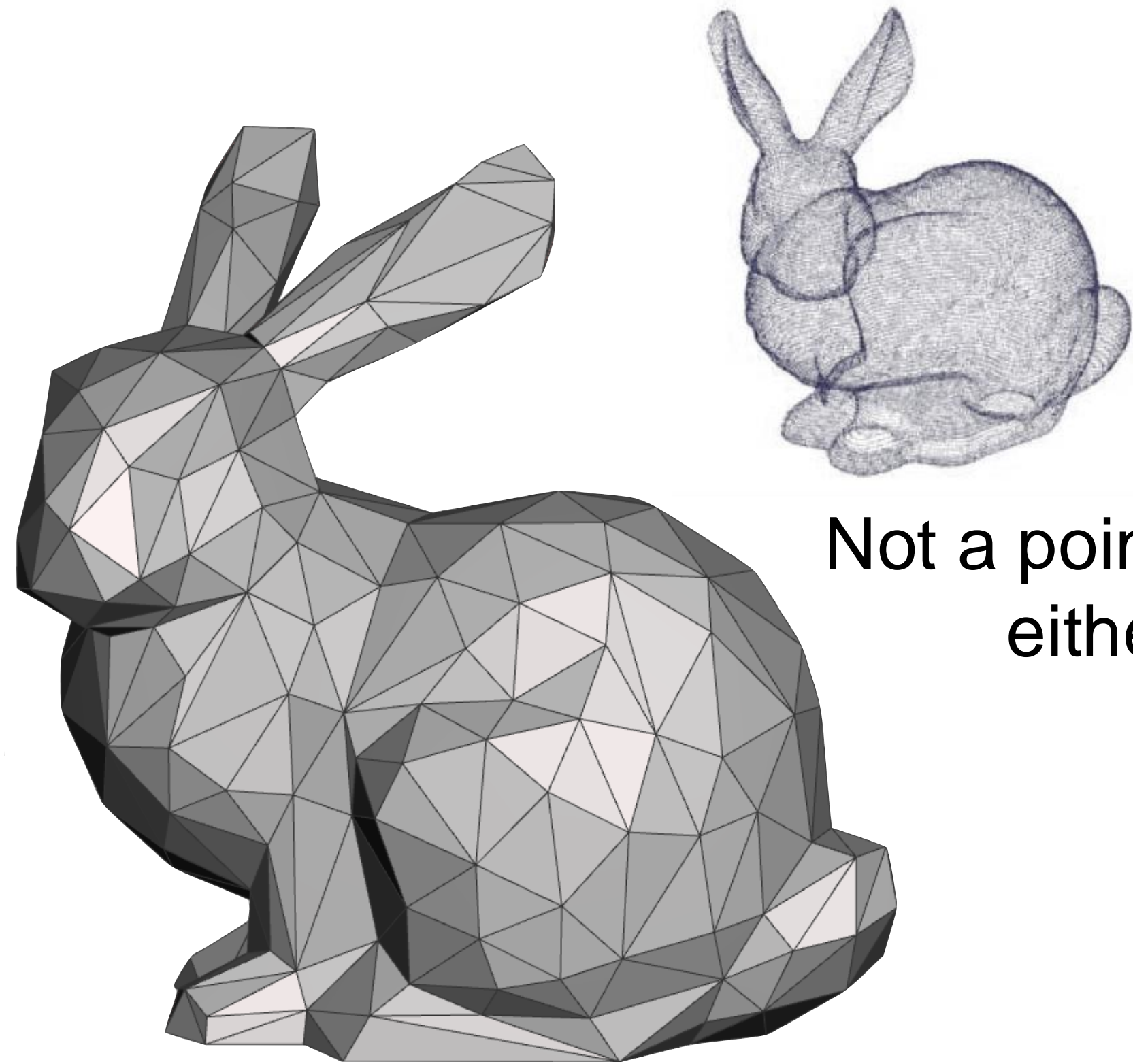


Representing a 3D scene as a continuous 5D function



What kind of a 3D representation is this?

It is not a Mesh



Not a point cloud
either



It is volumetric

It's *continuous* voxels made of shiny transparent cubes

What is the problem that is being solved?



Plenoptic Function

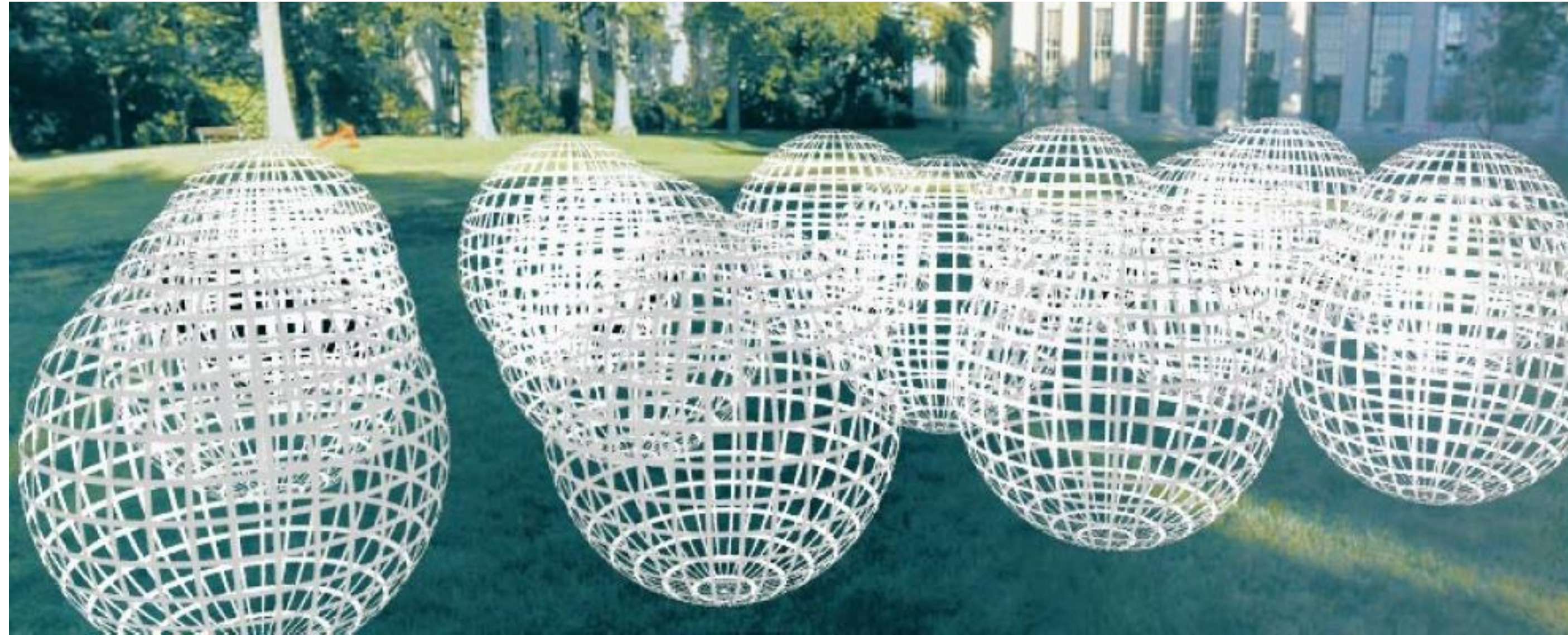


Figure by Leonard McMillan

Q: What is the set of all things that we can ever see?

A: The Plenoptic Function (Adelson & Bergen '91)

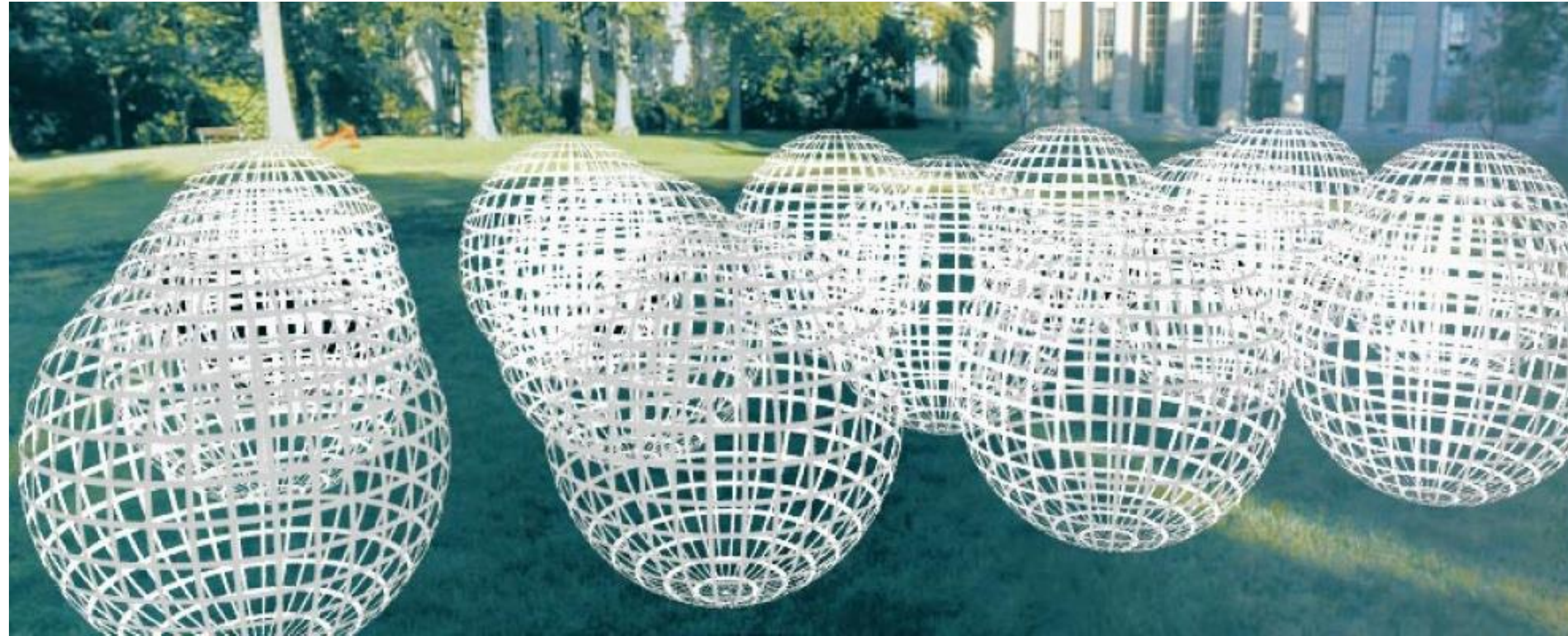
A holographic movie



$$P(\theta, \phi, \lambda, t, V_x, V_y, V_z)$$

- is intensity of light
 - Seen from ANY position and direction
 - Over time
 - As a function of wavelength

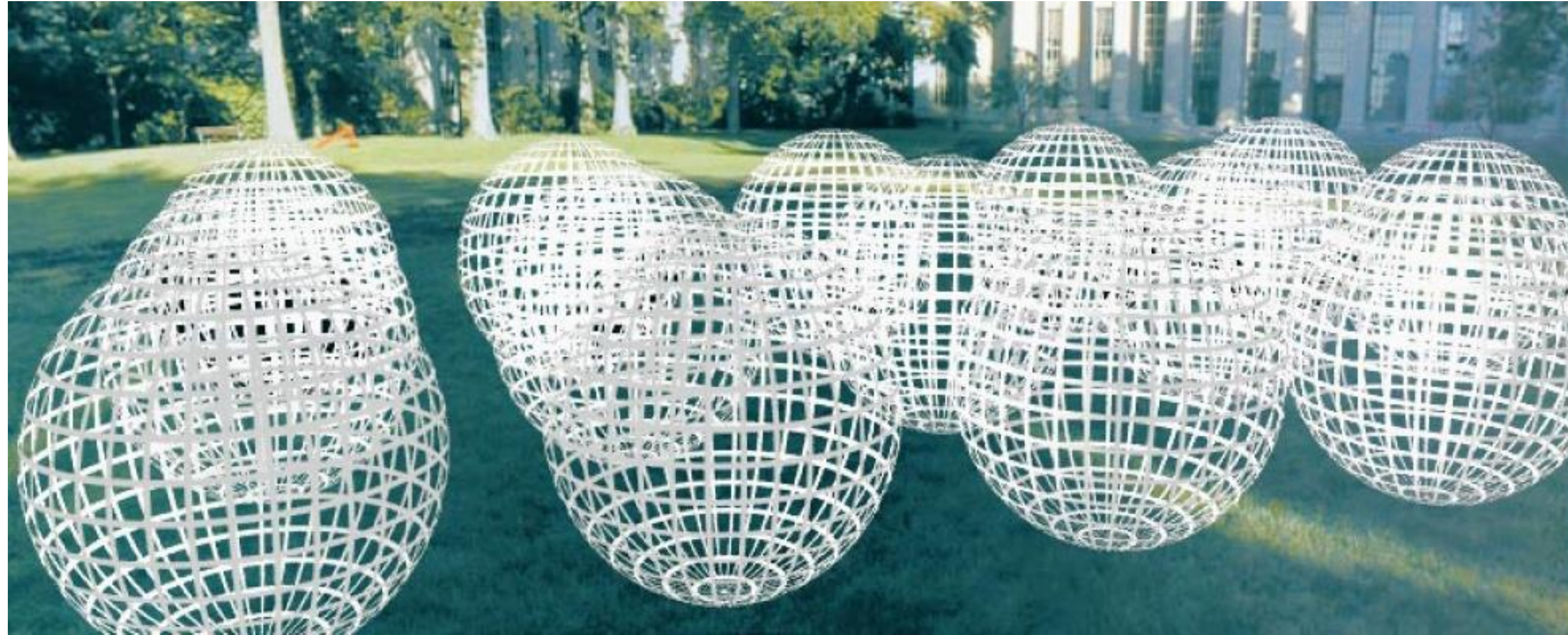
The plenoptic function



$$P(\theta, \phi, \lambda, t, V_x, V_y, V_z)$$

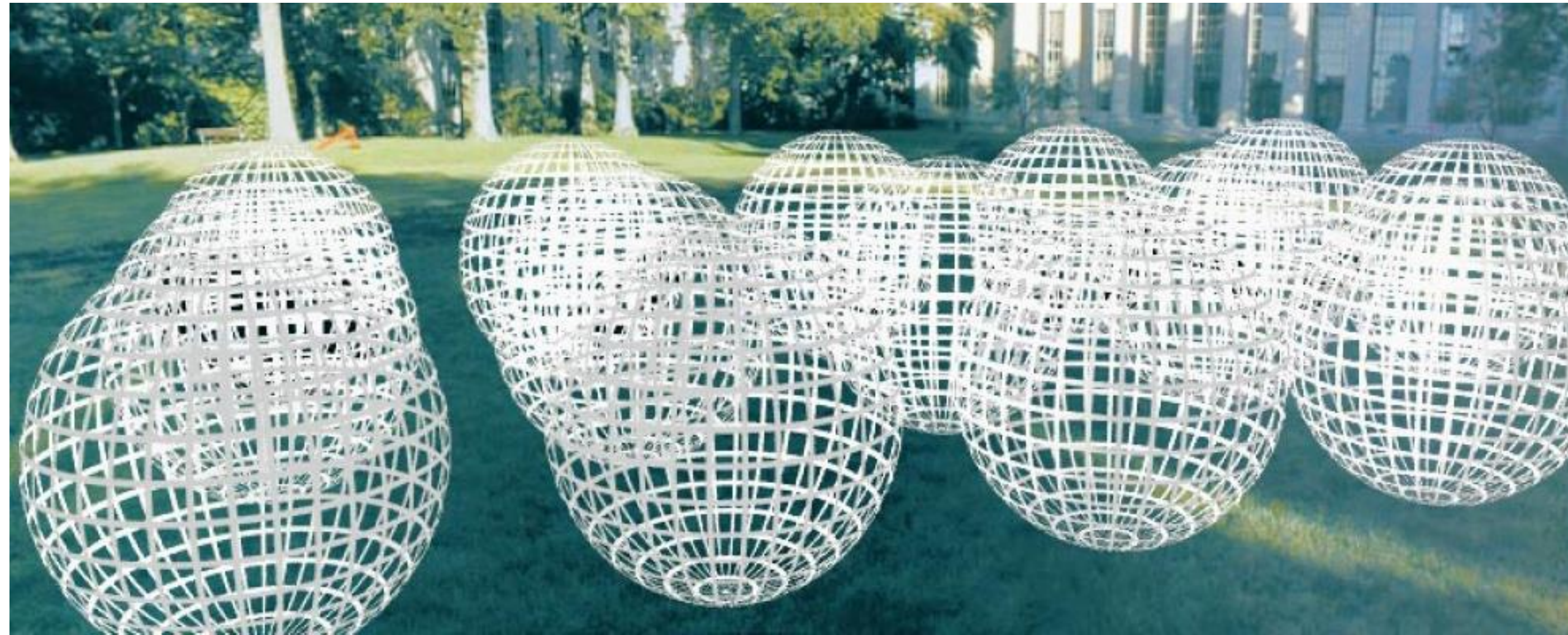
7D function, that can reconstruct every position & direction,
at every moment, at every wavelength
= it recreates the entirety of our visual reality!

Goal: Plenoptic Function from a set of images



- Objective: Recreate the visual reality
- All about recovering photorealistic pixels, not about recording 3D point or surfaces
—Image Based Rendering aka **Novel View Synthesis**

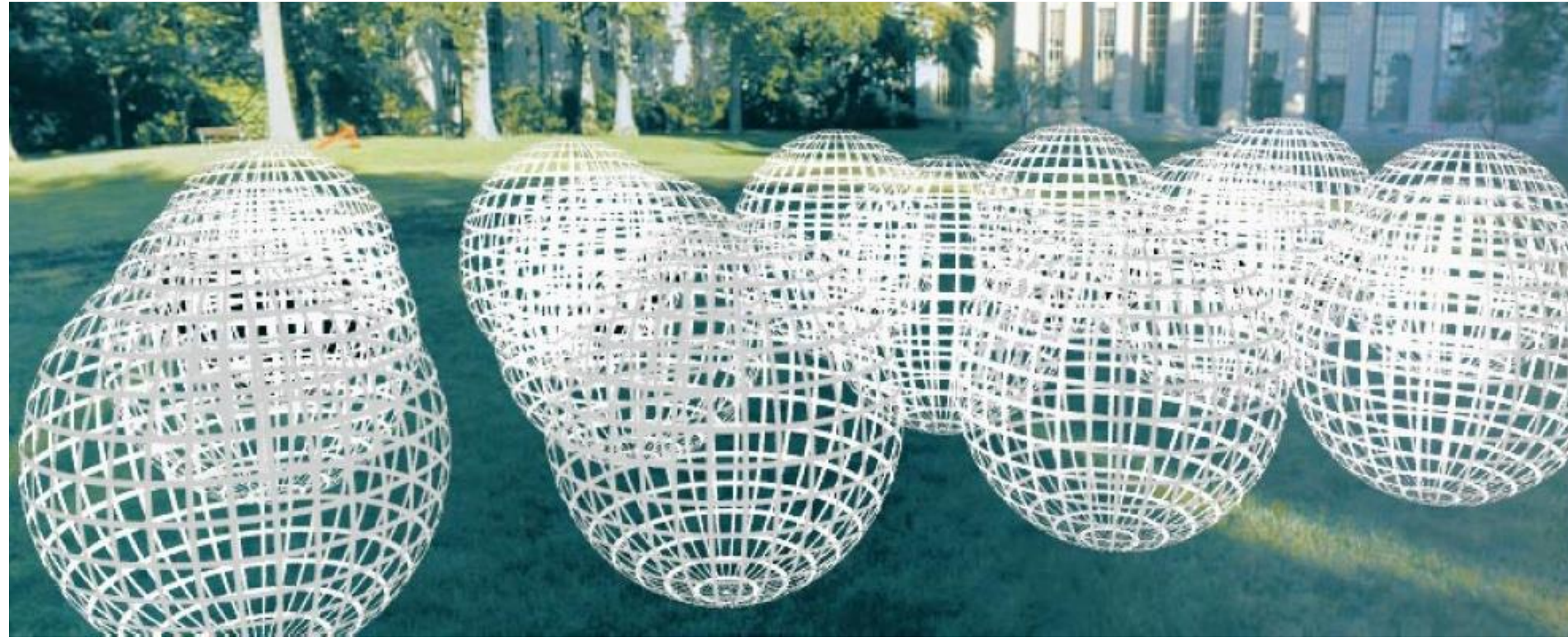
Goal: Plenoptic Function from a set of images



It is a conceptual device

Adelson & Bergen do not discuss how to solve this

Plenoptic Function



Look familiar
😊?

7D function:
2 – direction
1 – wavelength
1 – time
3 – location

$$P(\theta, \phi, \lambda, t, V_x, V_y, V_z) \longrightarrow P(\theta, \phi, V_x, V_y, V_z)$$

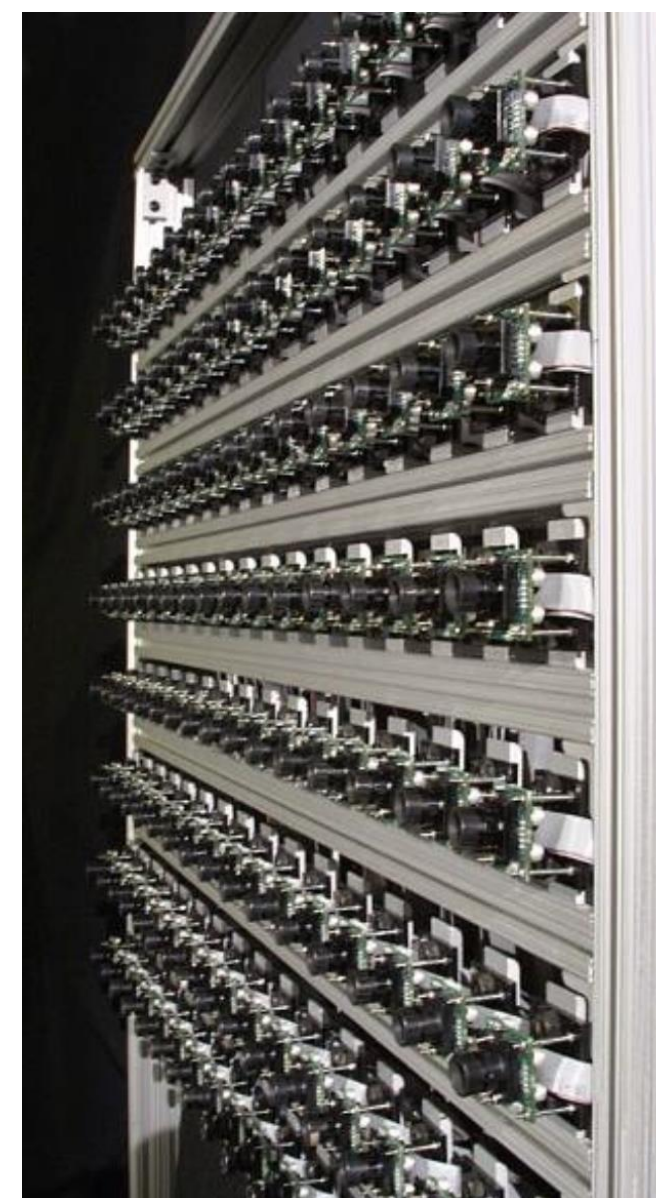
Let's simplify:

1. Remove the time
2. Remove the wavelength & let the function output RGB colors

Lightfield / Lumigraph

- Previous approaches for modeling the Plenoptic Function
- Take a lot of pictures from many views
- Interpolate the rays to render a novel view

Stanford Gantry
128 cameras



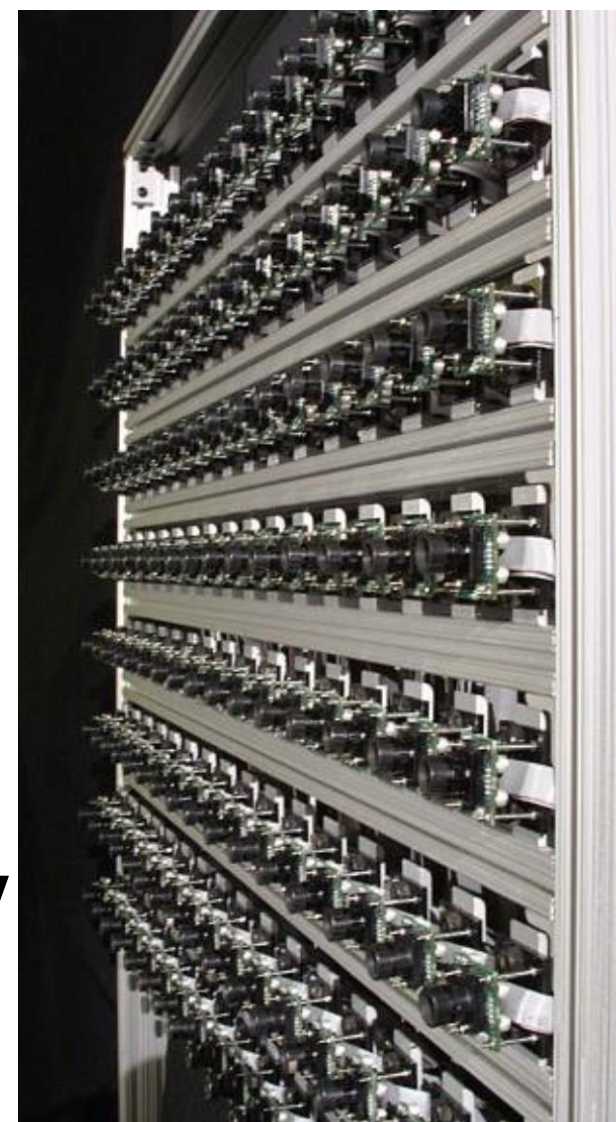
Lytro camera

Lightfield / Lumigraph

- Previous approaches for modeling the Plenoptic Function
- Take a lot of pictures from many views
- Interpolate the rays to render a novel view



Stanford Gantry
128 cameras



Lytro camera

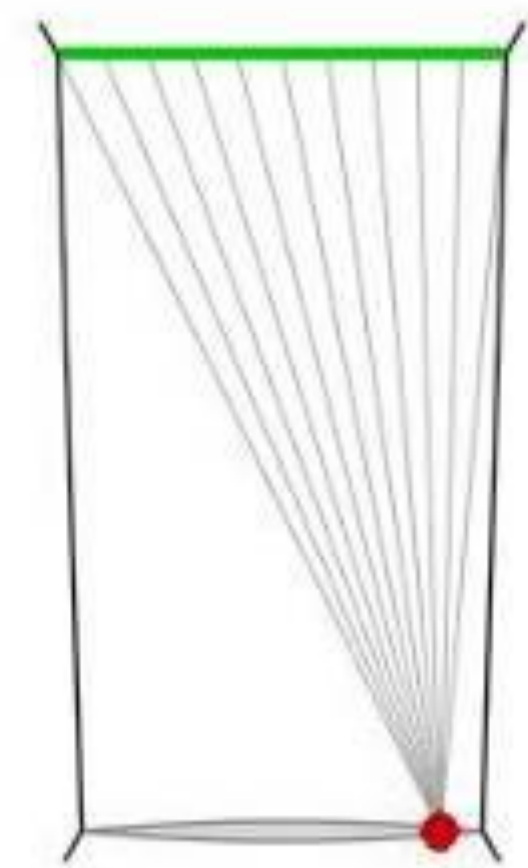


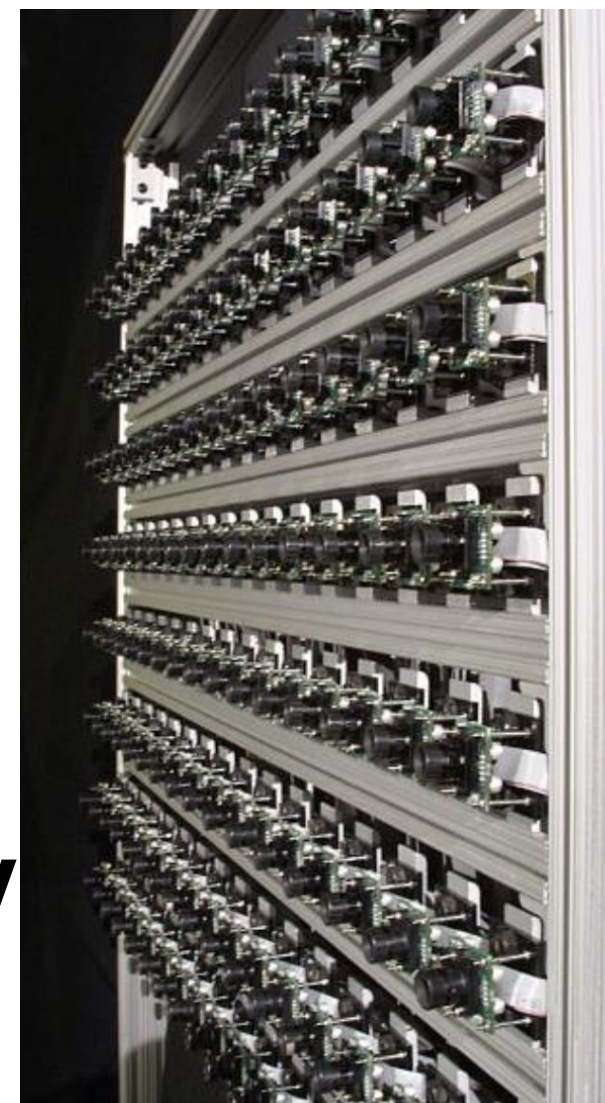
Figure from Marc Levoy

Lightfield / Lumigraph

- Previous approaches for modeling the Plenoptic Function
- Take a lot of pictures from many views
- Interpolate the rays to render a novel view



Stanford Gantry
128 cameras



Lytro camera

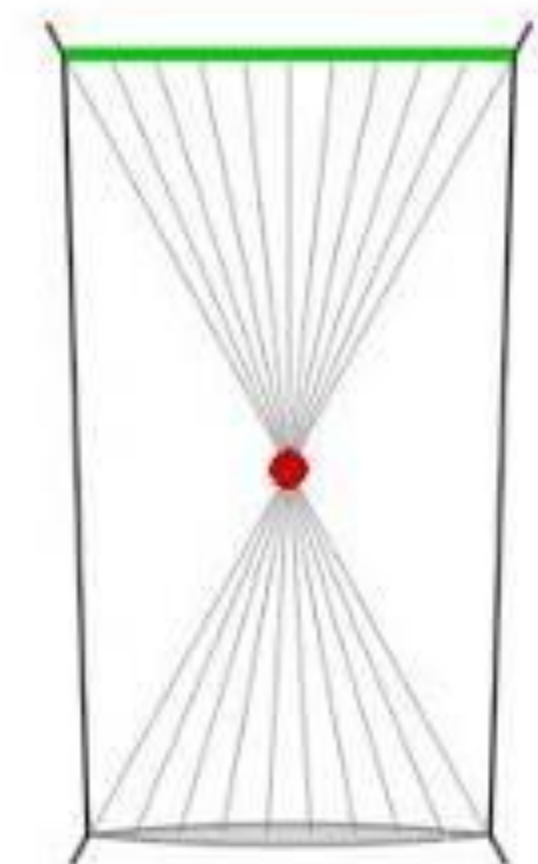
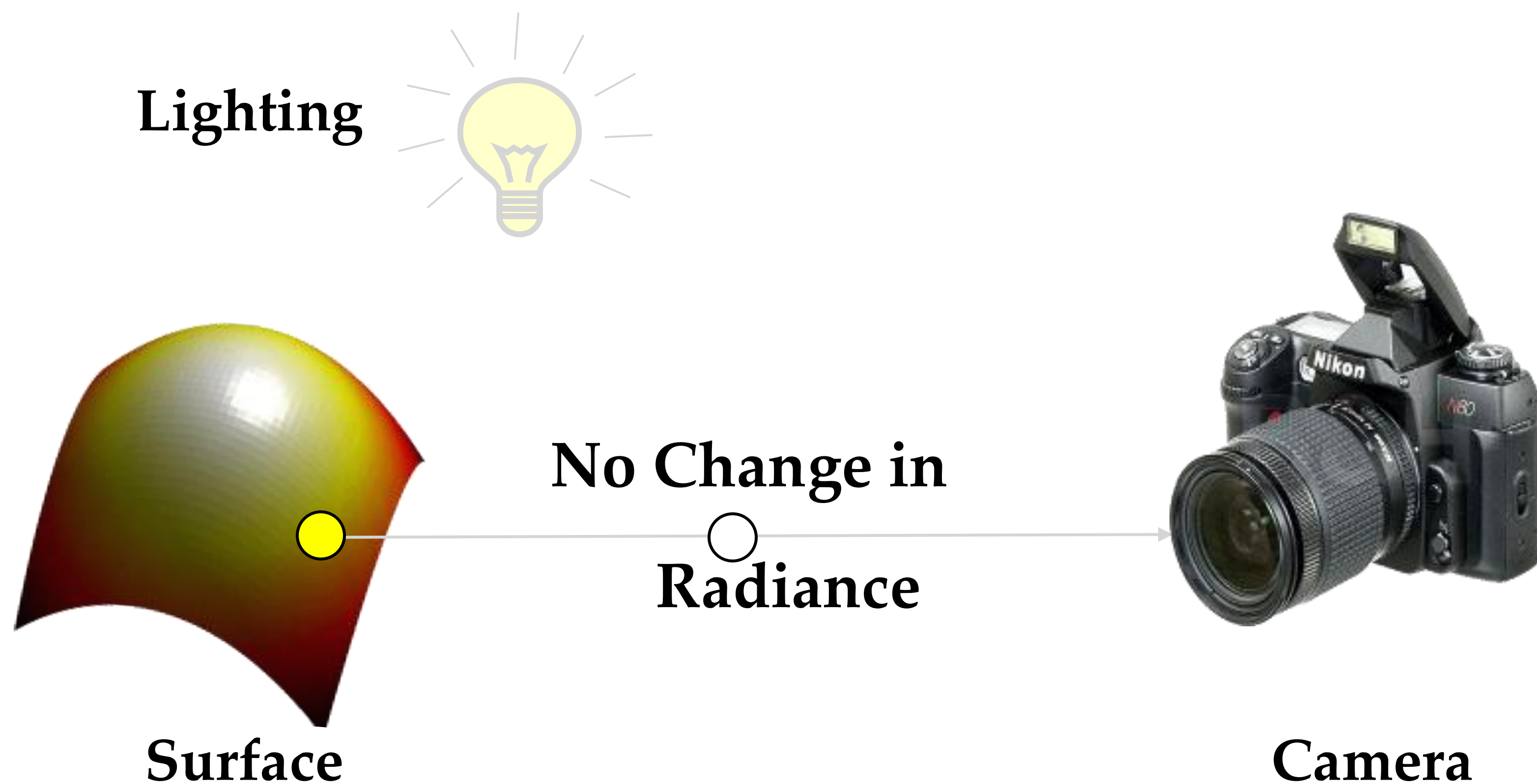


Figure from Marc Levoy

Big Assumption: a ray does not change color

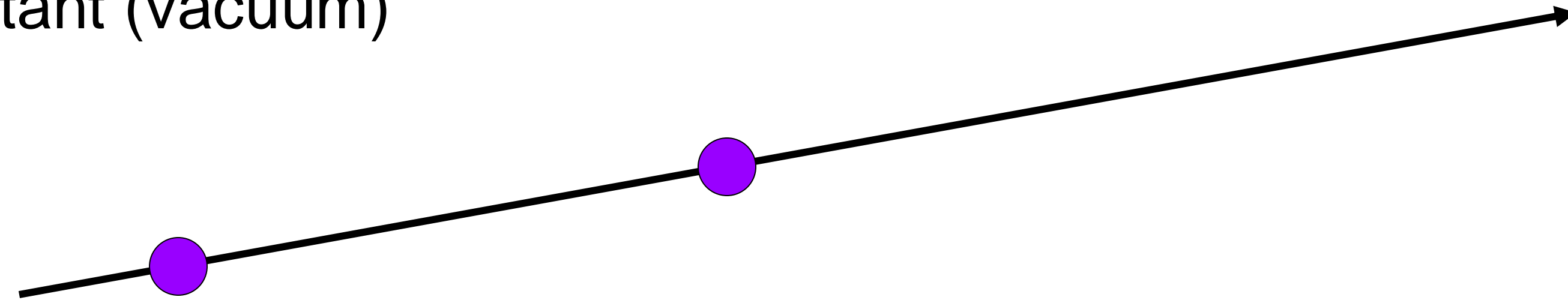


True if there is no occlusion or fog

With this assumption: Ray Reuse

Infinite line

- Assume light is constant (vacuum)



The 5D function

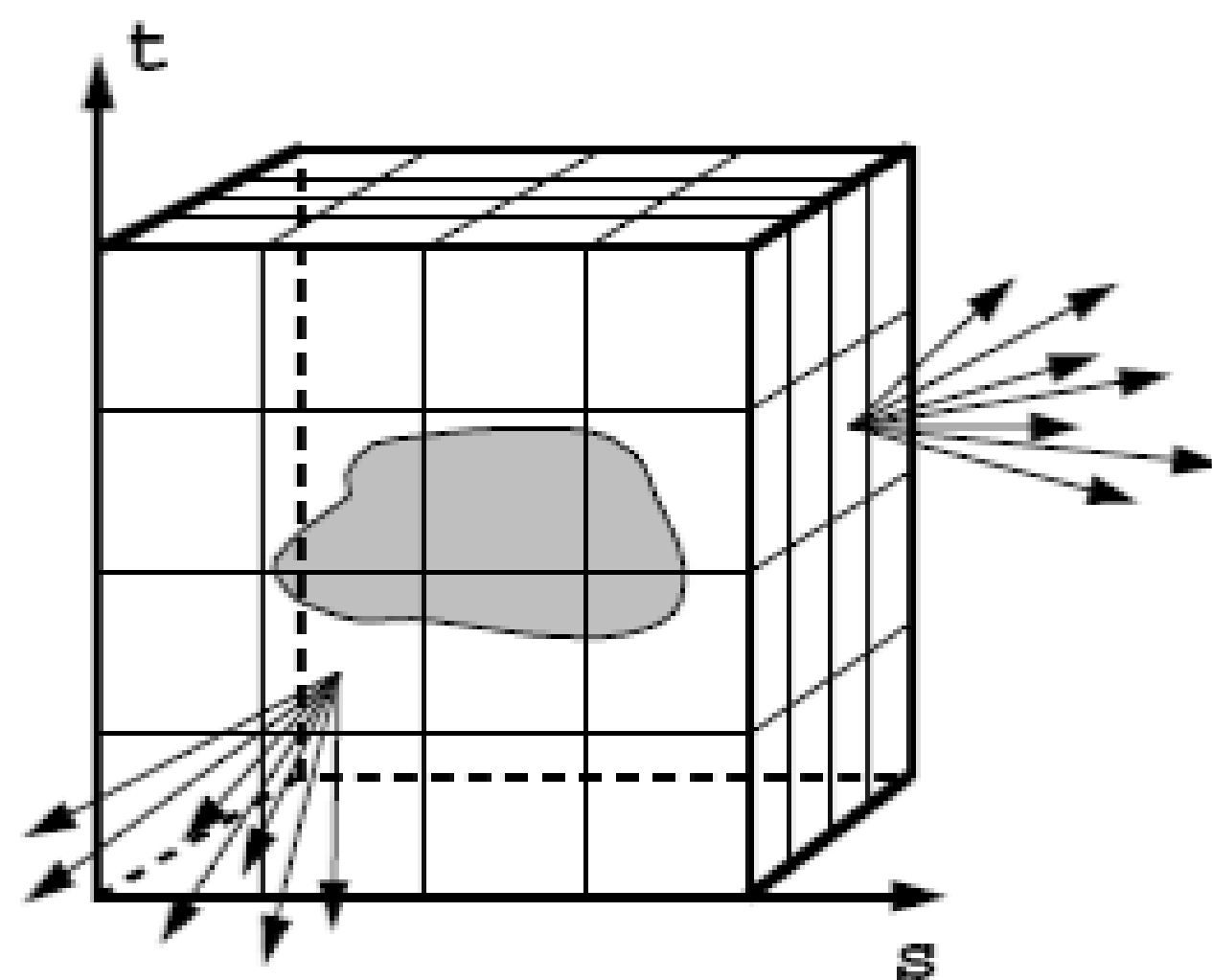
- 3D position
- 2D direction

is now 4D

- 2D direction
- 2D position
- non-dispersive medium

Ray Reuse Assumption

Because of this it only models the
plenoptic surface:

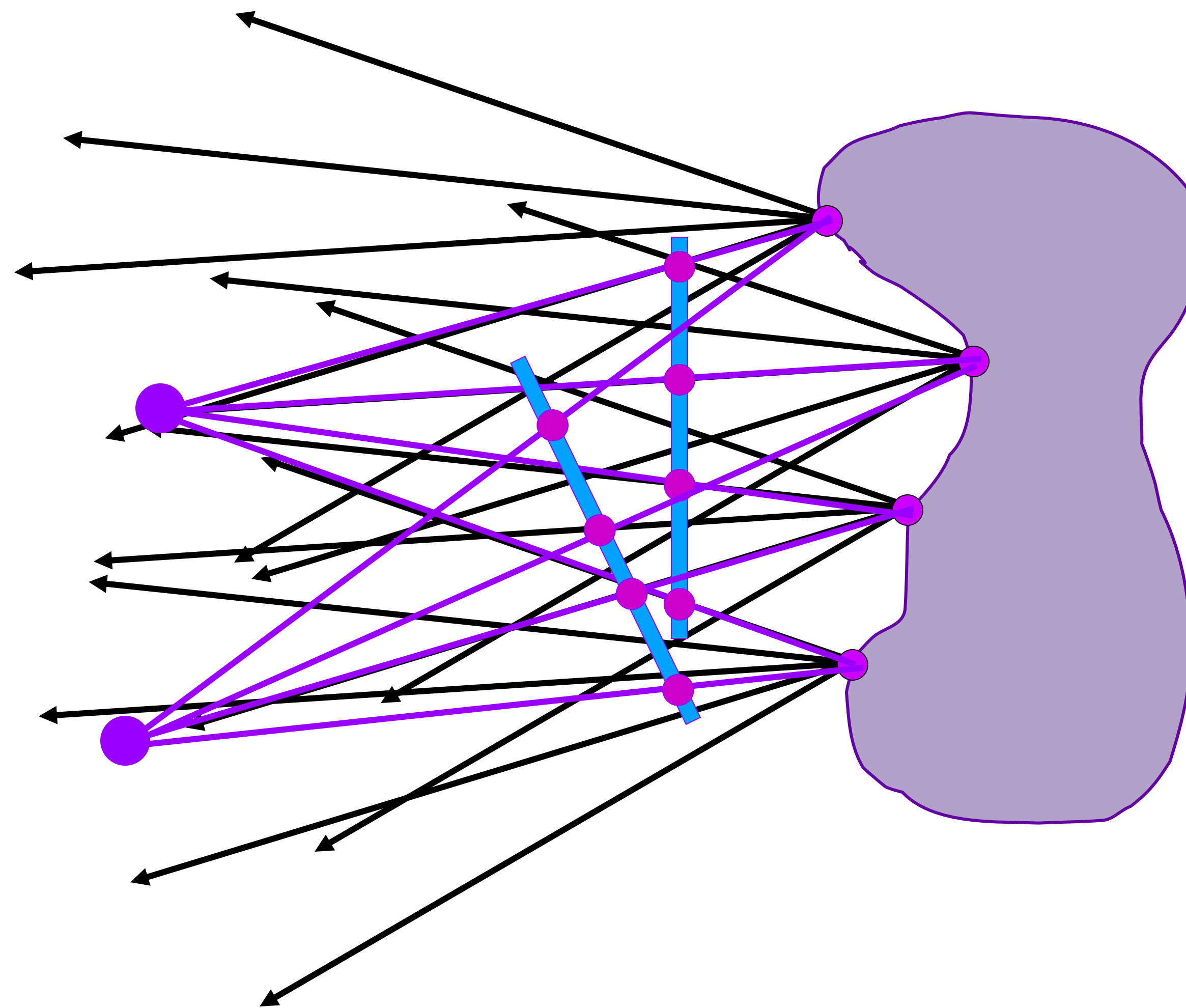


It's like



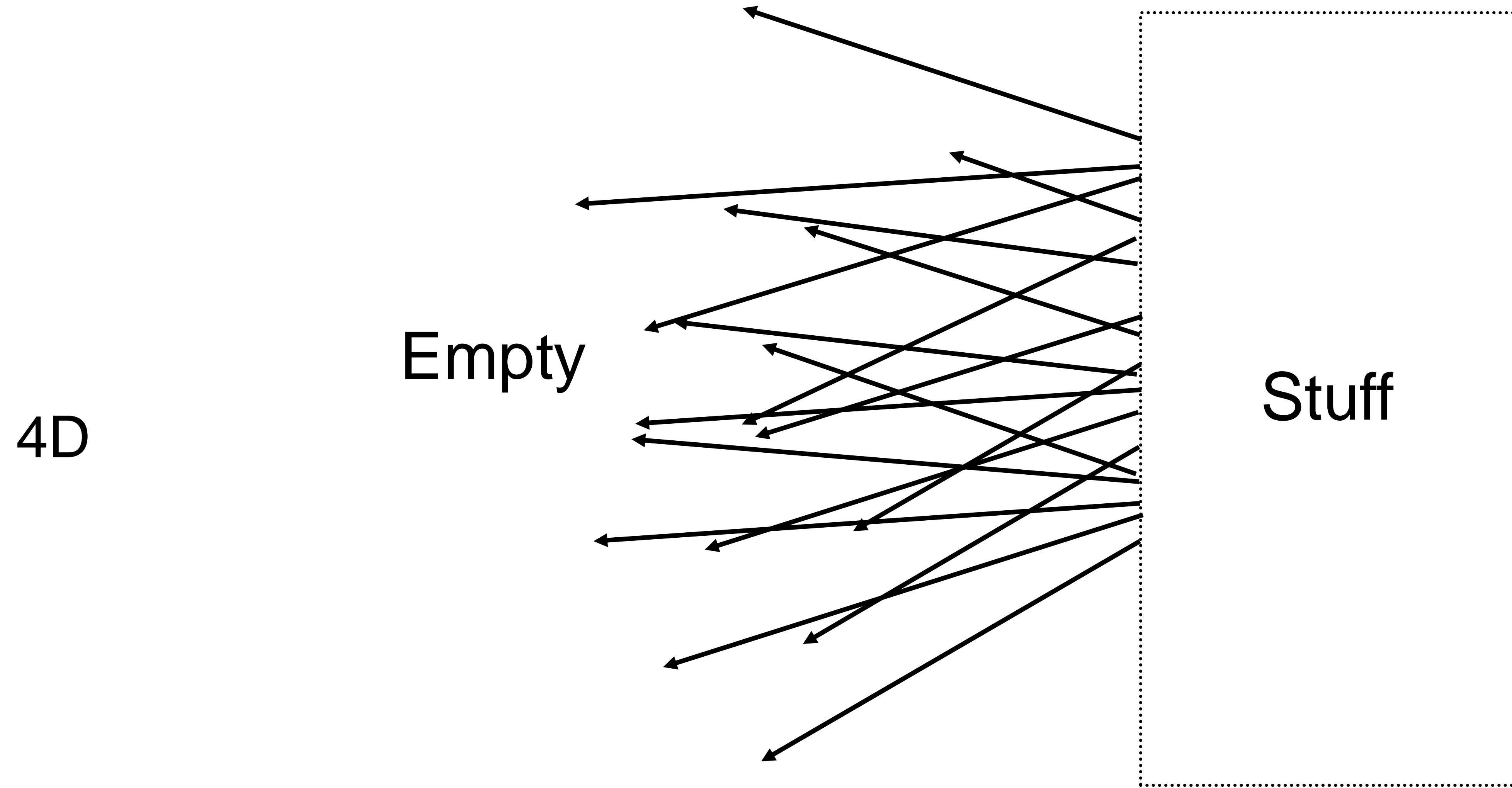
Figure 1: The surface of a cube holds all the radiance information due to the enclosed object.

Synthesizing novel views



Lumigraph / Lightfield

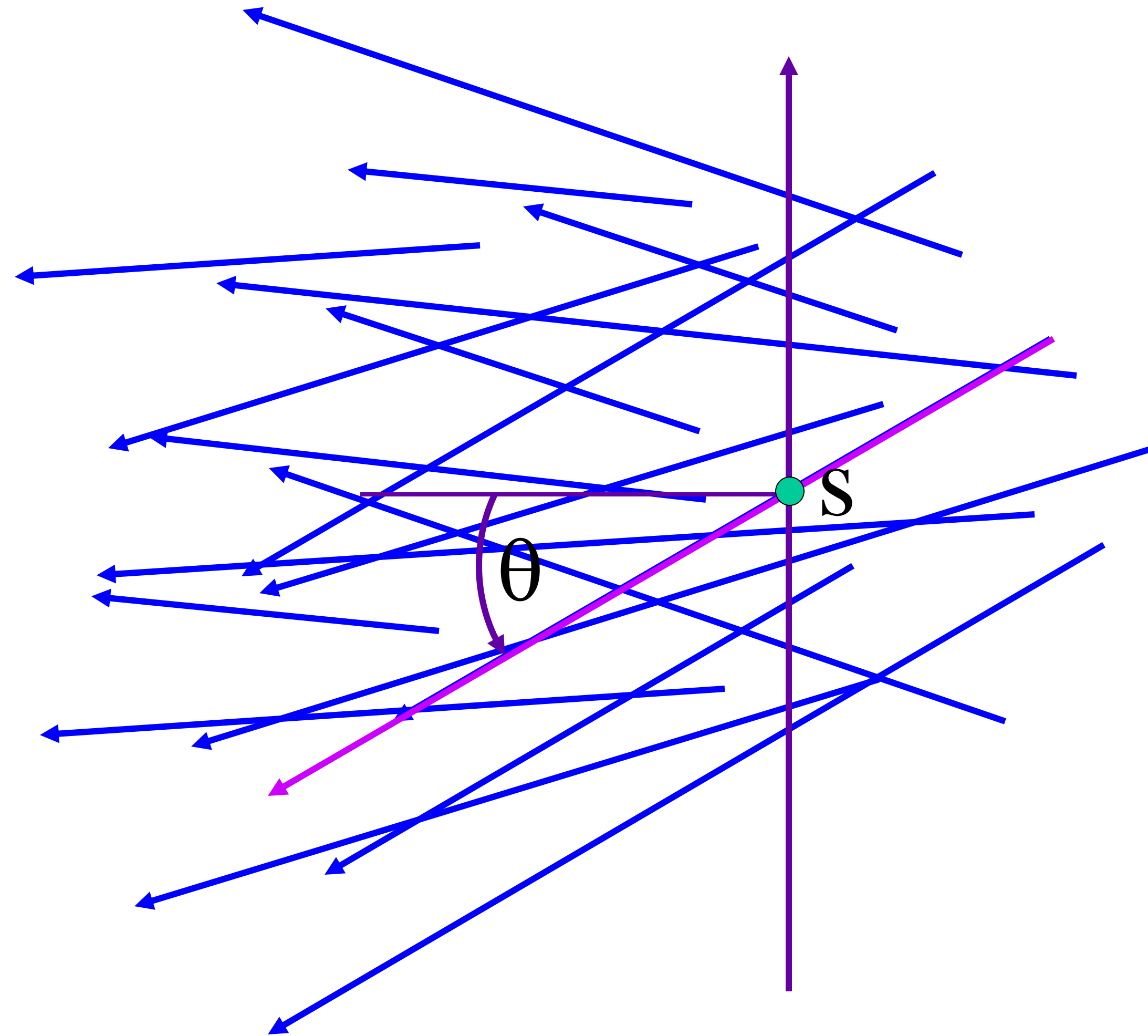
Outside convex space



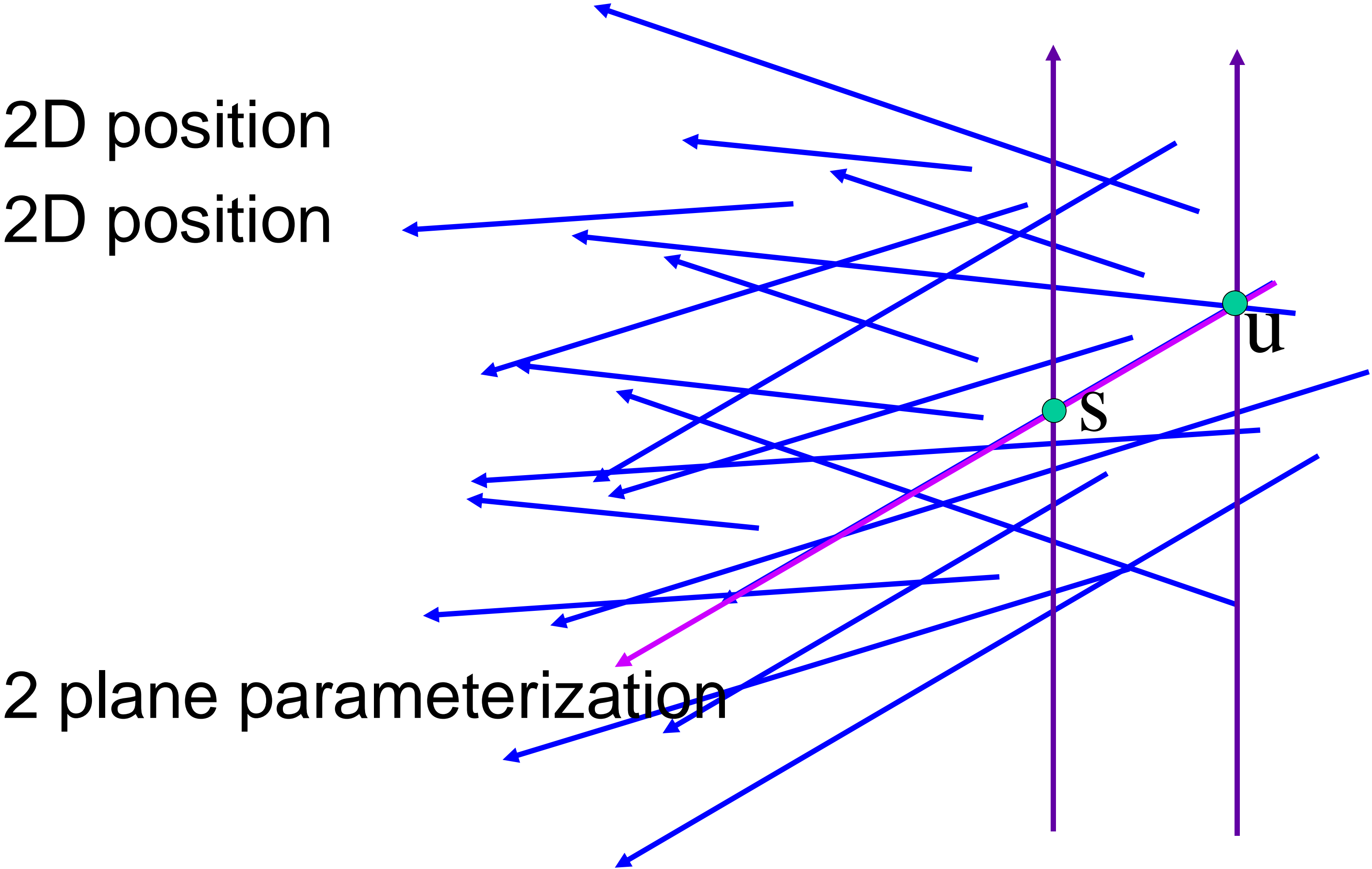
Lumigraph - Organization

2D position

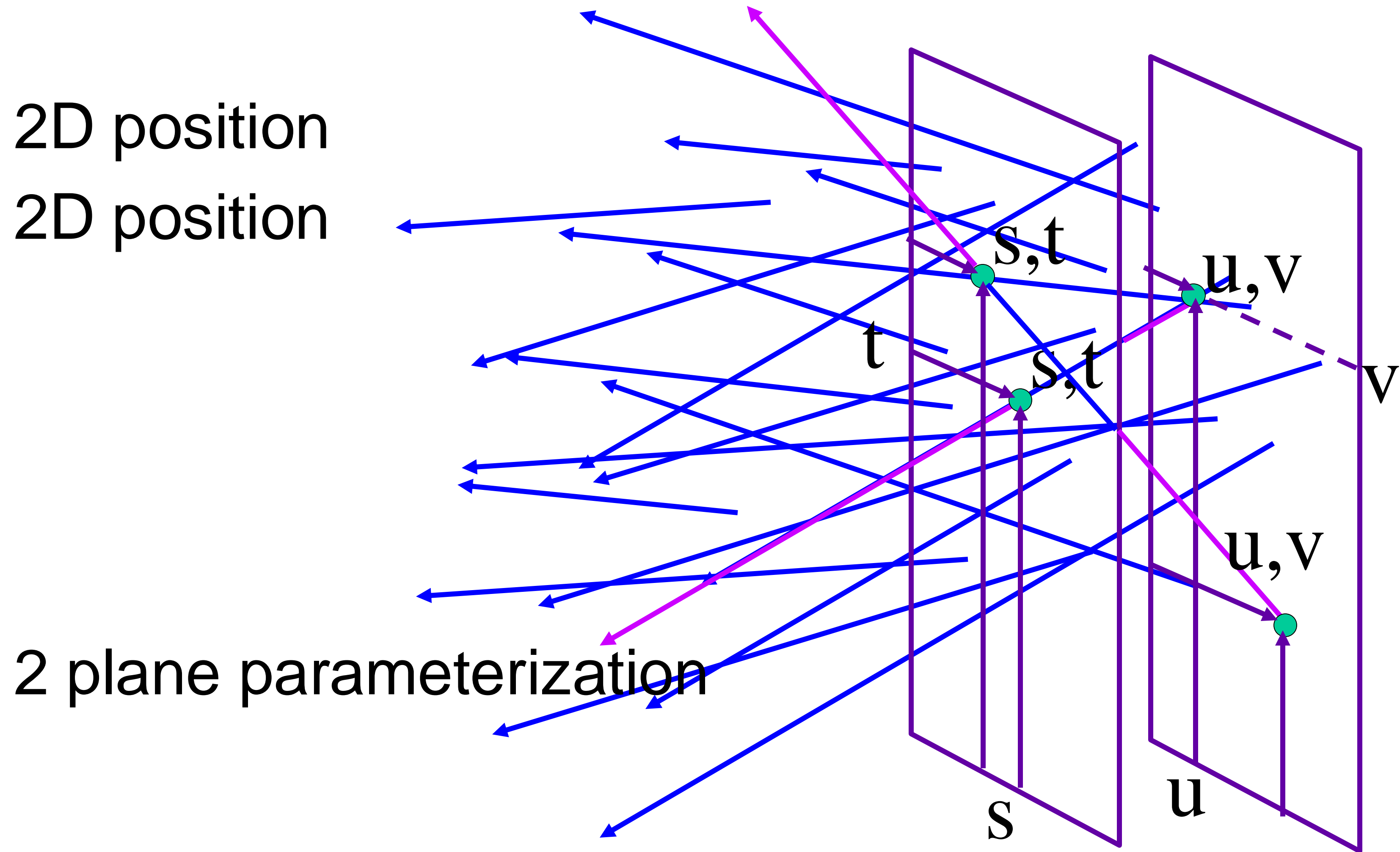
2D direction



Lumigraph - Organization



Lumigraph - Organization

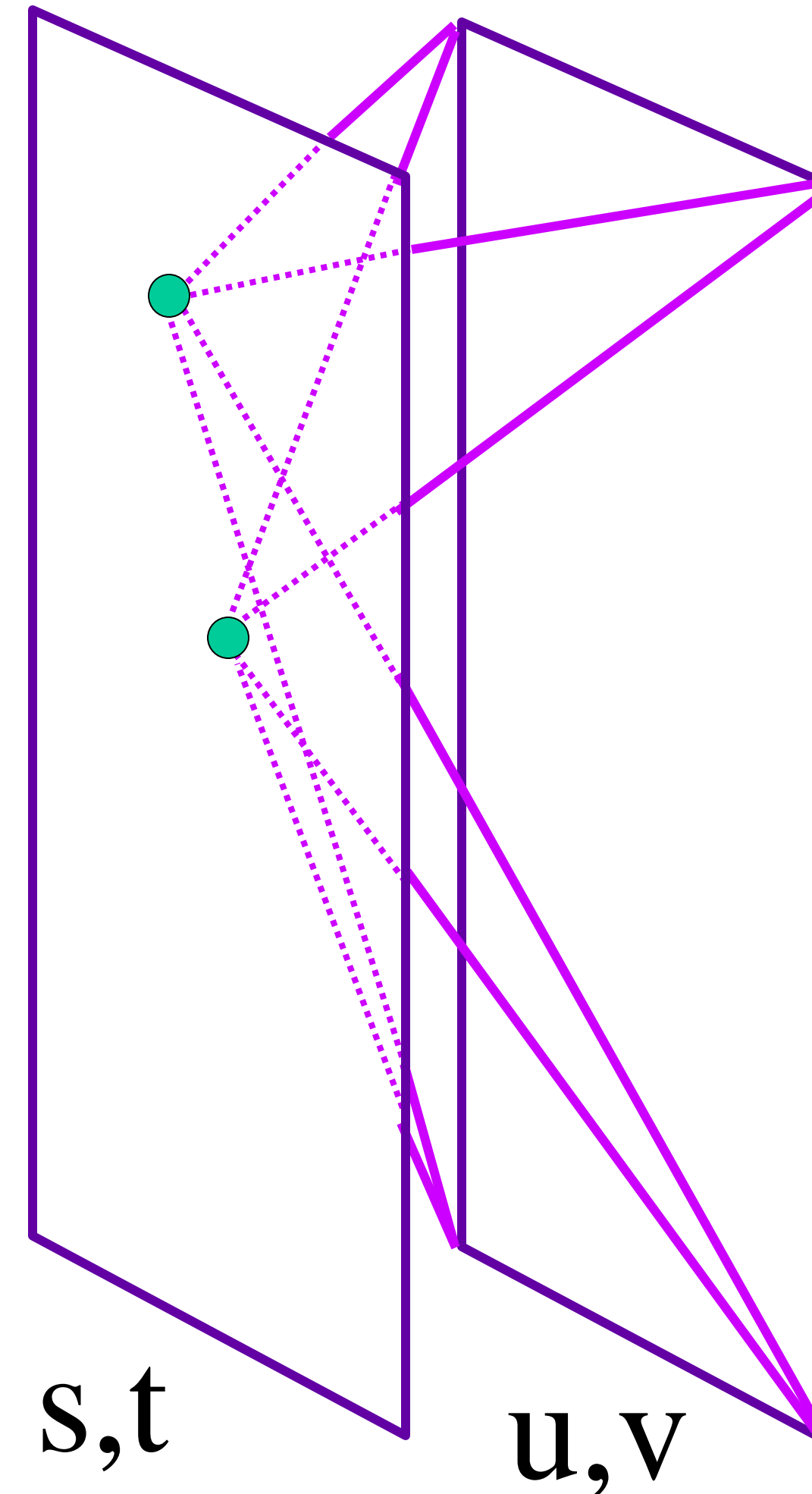


Lumigraph - Organization

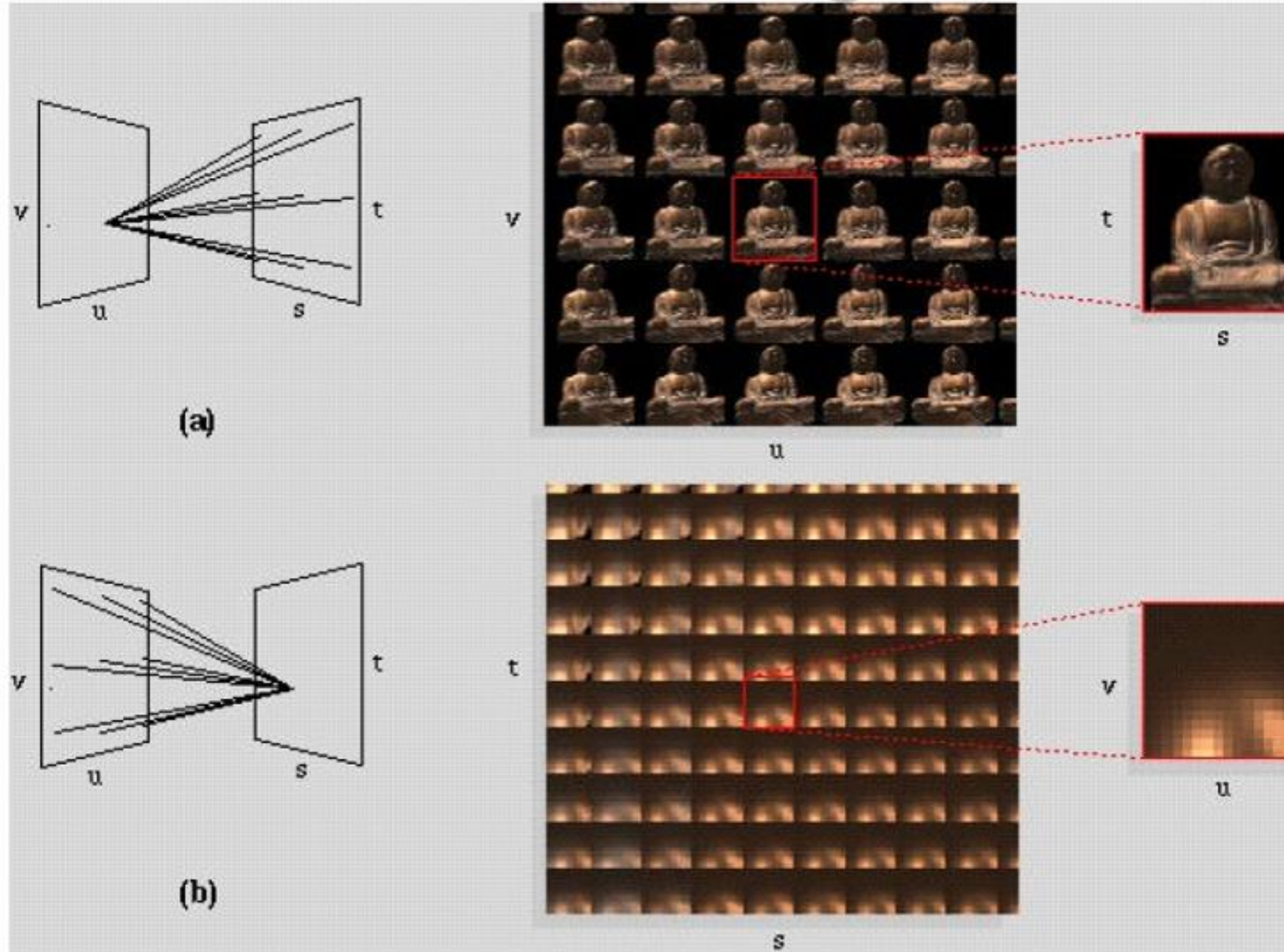
Hold s, t constant

Let u, v vary

An image



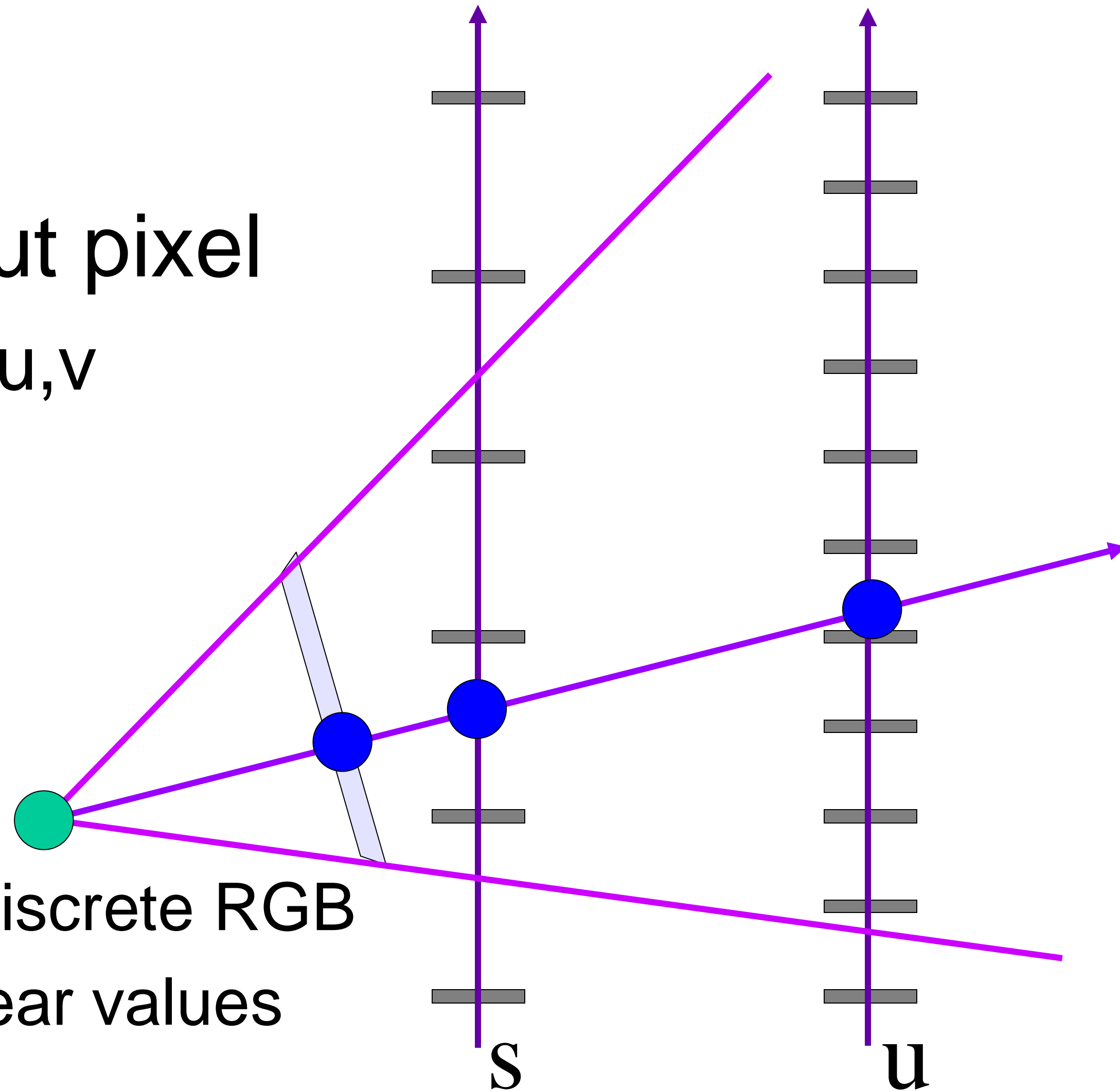
Lumigraph / Lightfield



Novel View Synthesis

- For each output pixel
 - determine s, t, u, v

- either
 - use closest discrete RGB
 - interpolate near values



How NeRF models the Plenoptic Function

$$P(\theta, \phi, V_X, V_Y, V_Z)$$

Look familiar

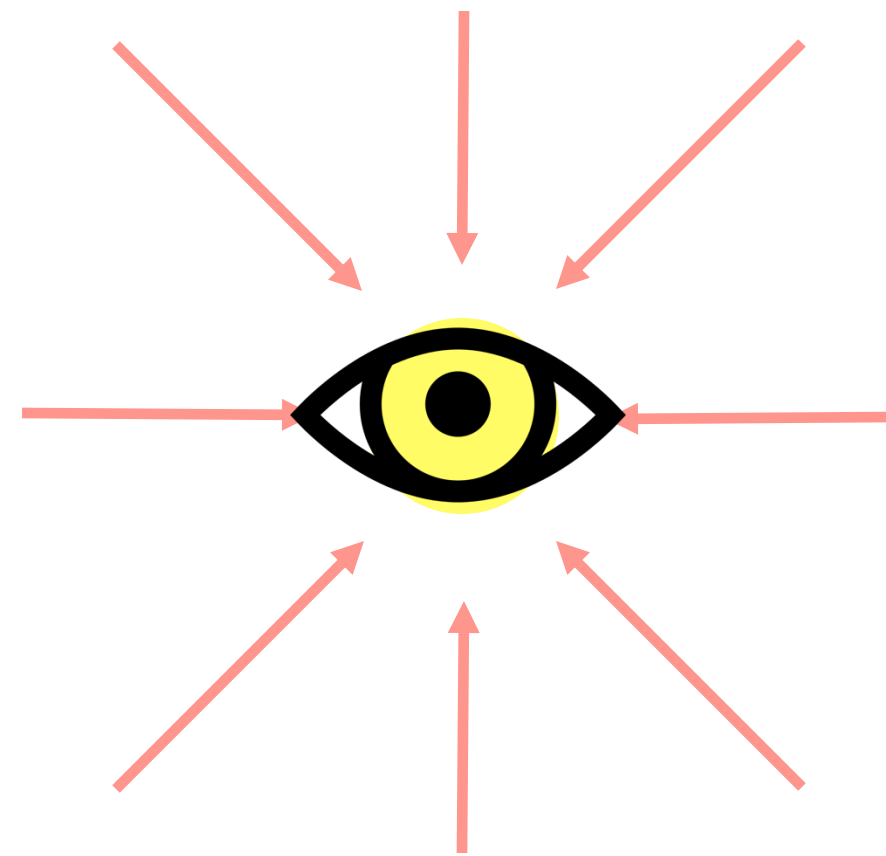


NeRF takes the same input as the Plenoptic Function!

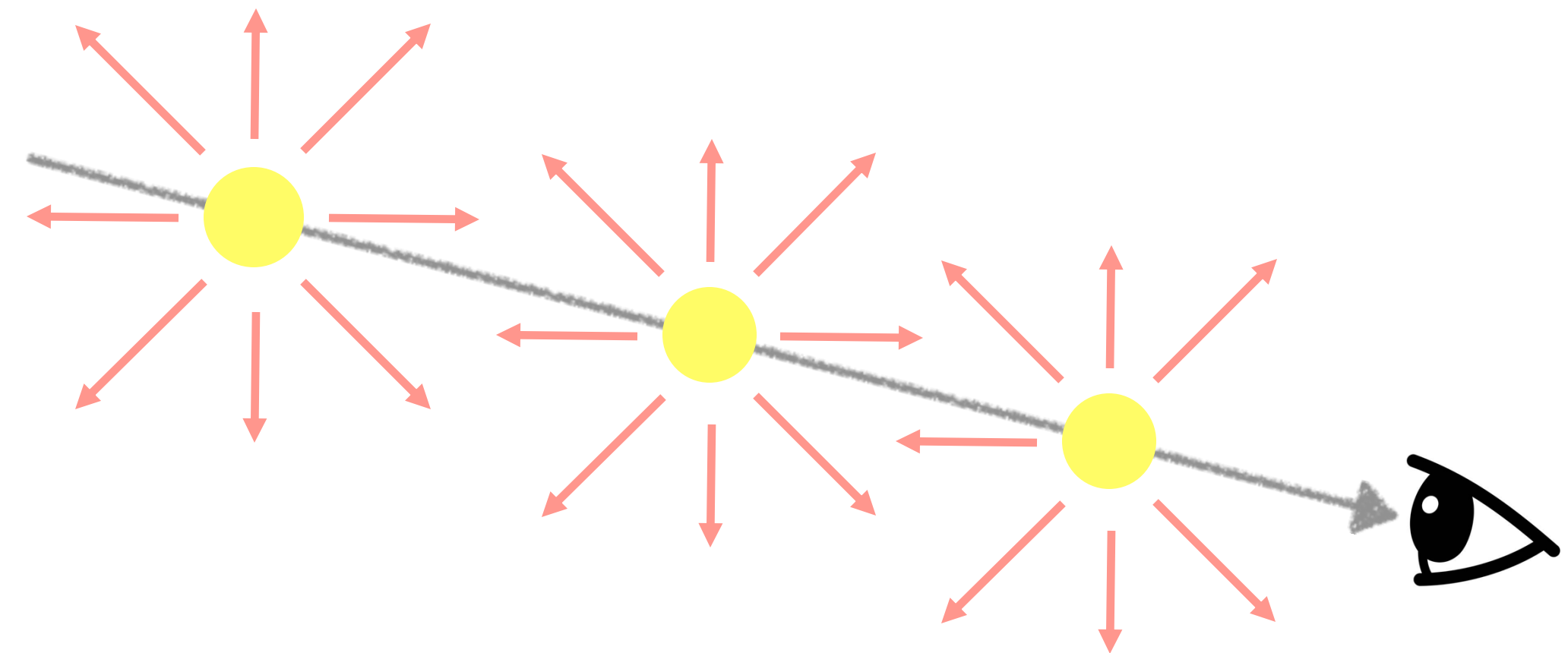
It allows rays to change color. Hence we can fly into the glass bowl (if we had enough observation)



A subtle difference



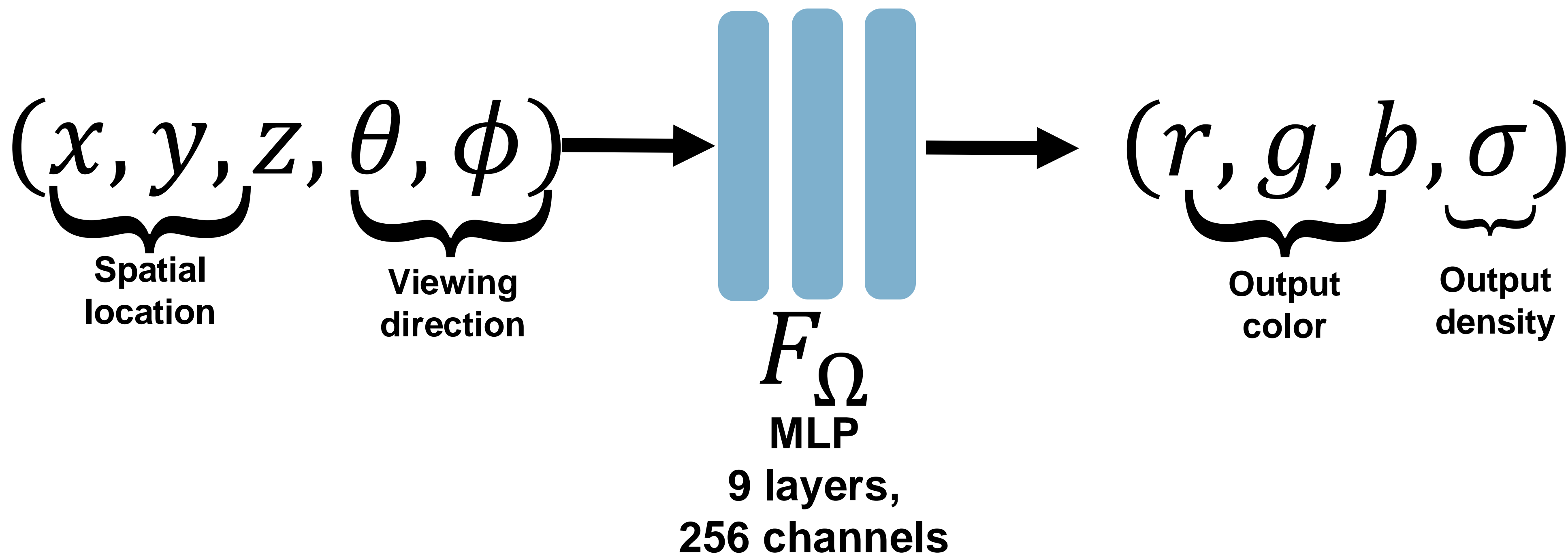
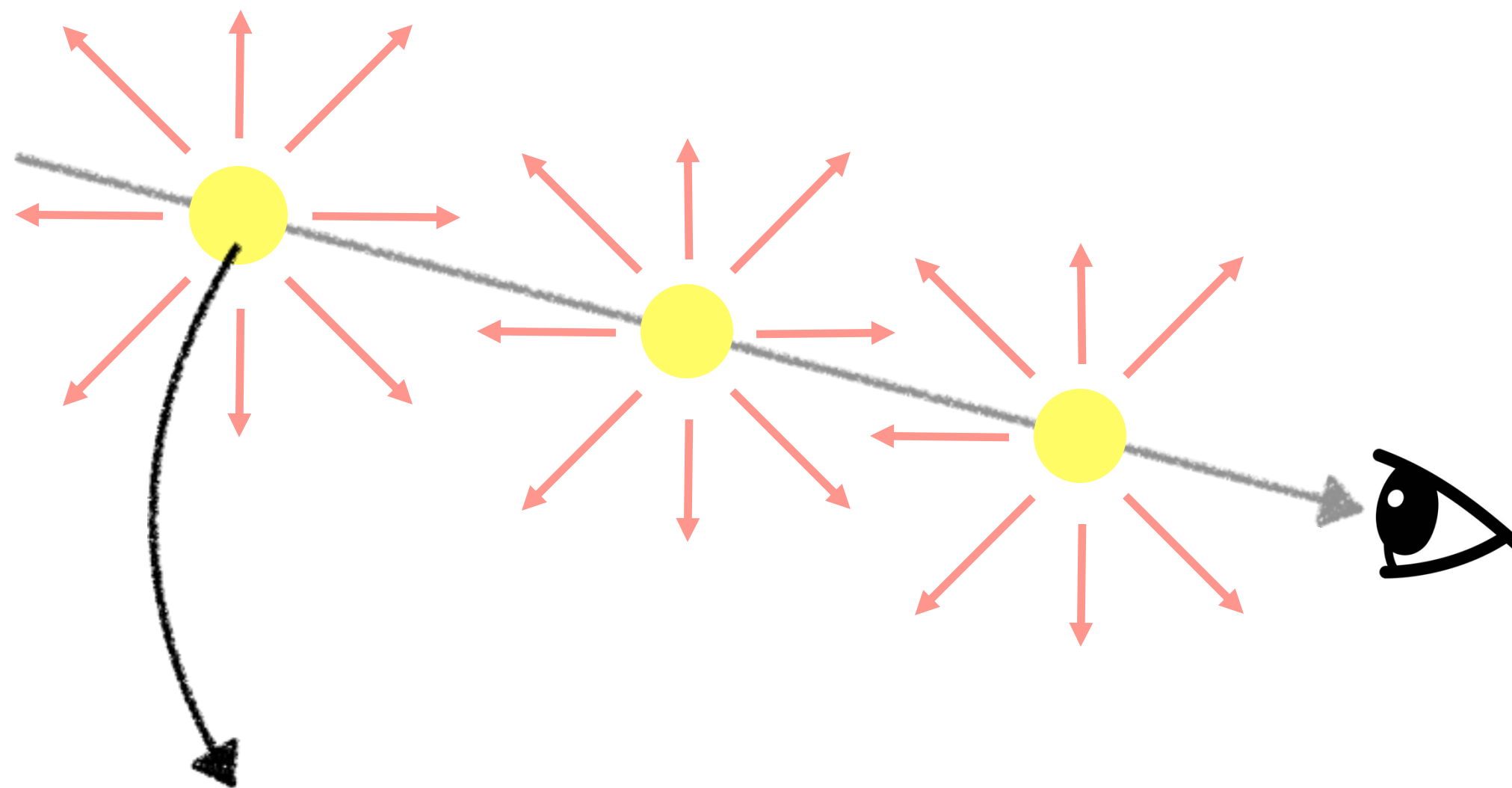
Plenoptic Function



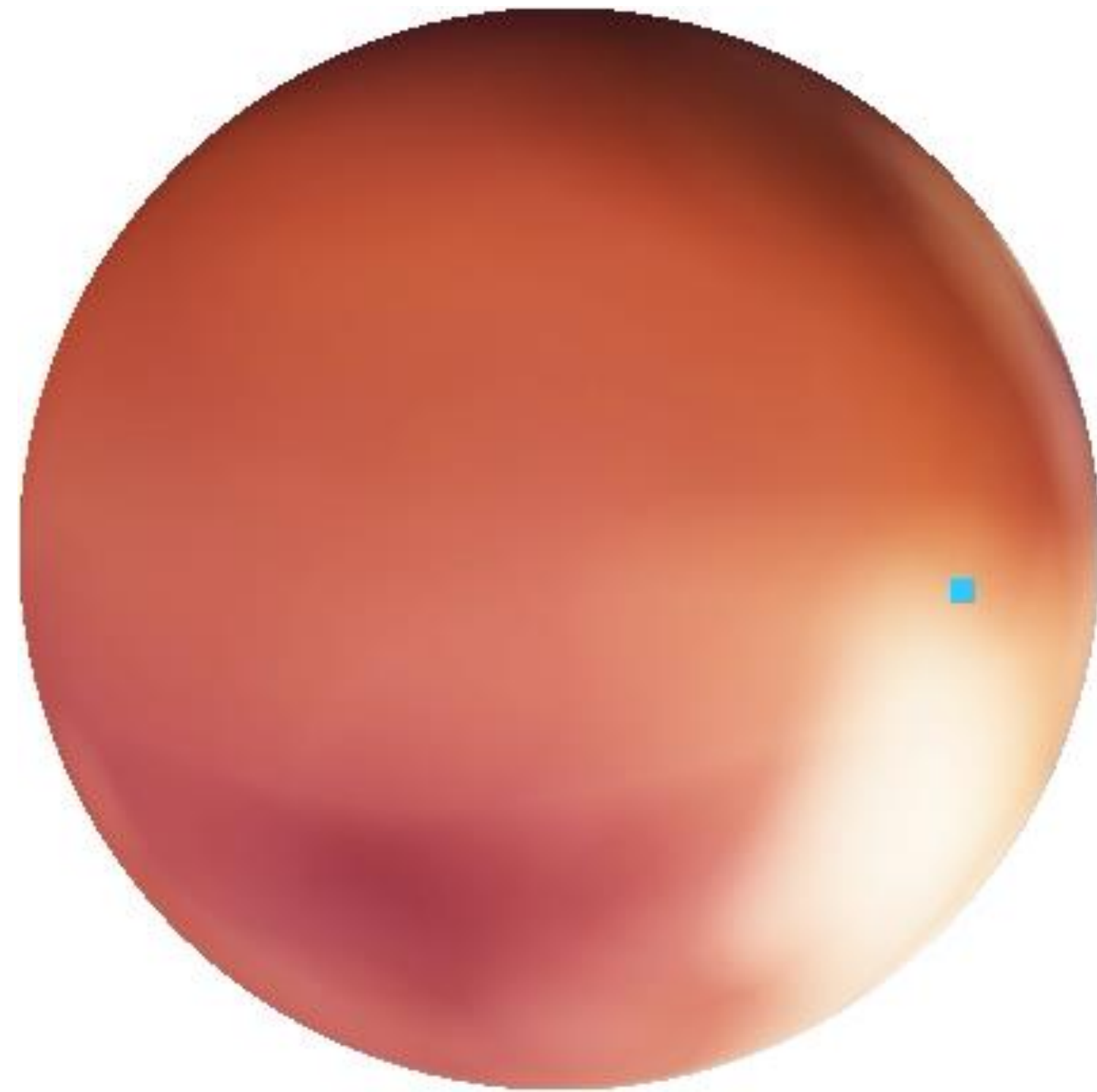
NeRF

So NeRF requires the integration along the viewing ray to compute the Plenoptic Function

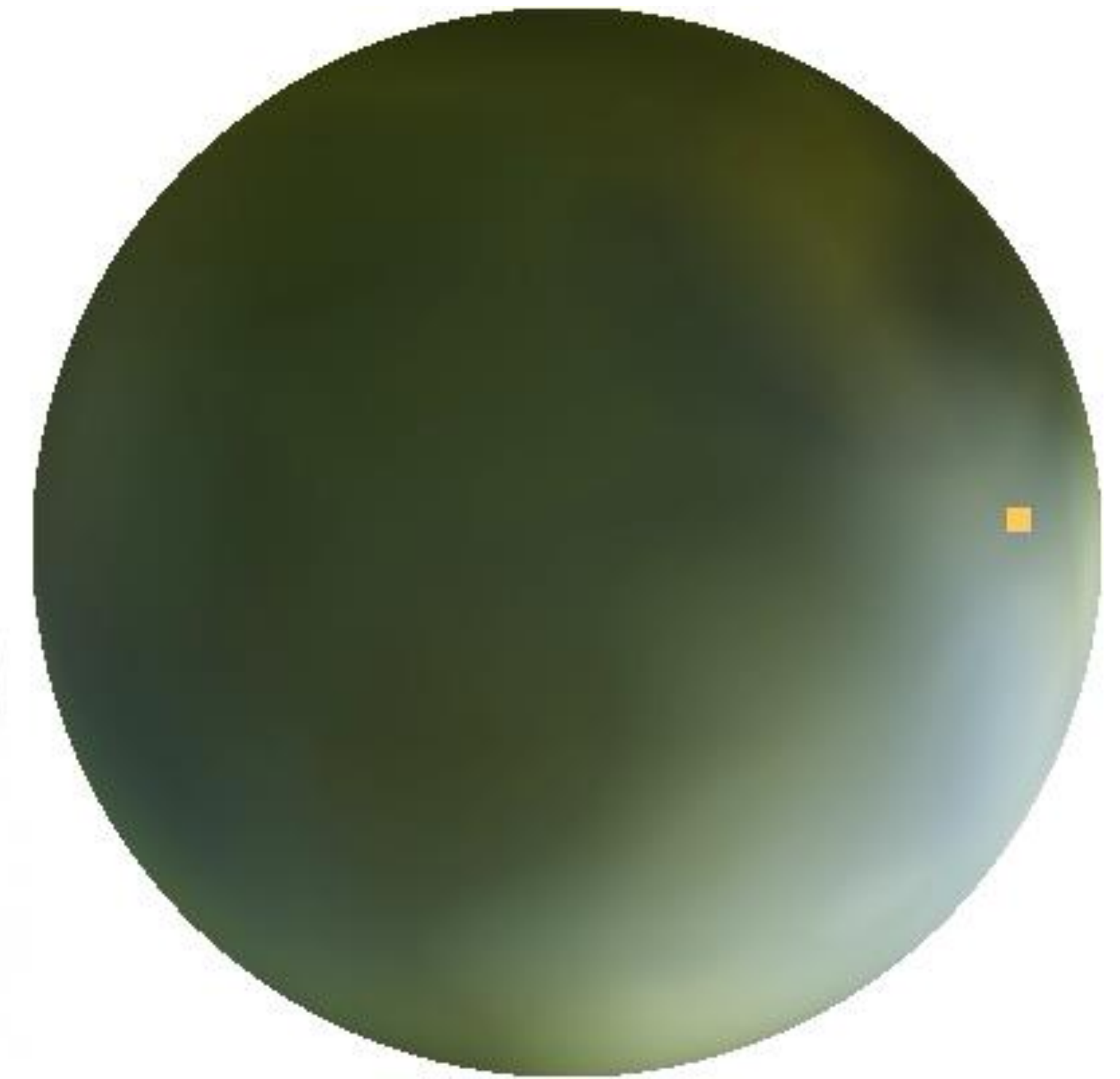
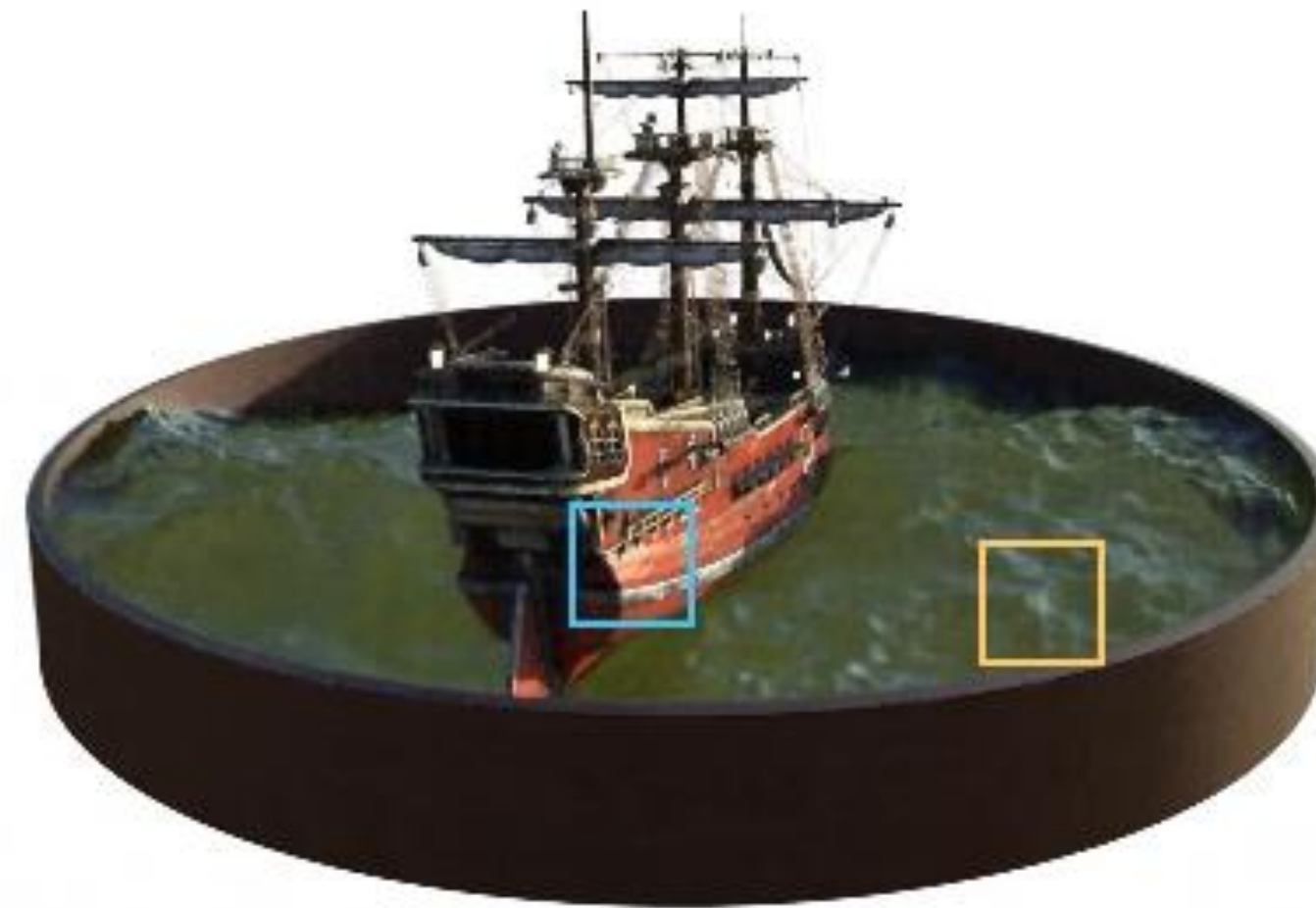
Bottom line: it models the full (5D) plenoptic function!



Visualizing the 2D function on the sphere



Outgoing radiance distribution
for point on side of ship



Outgoing radiance distribution
for point on water's surface

Baking in Light



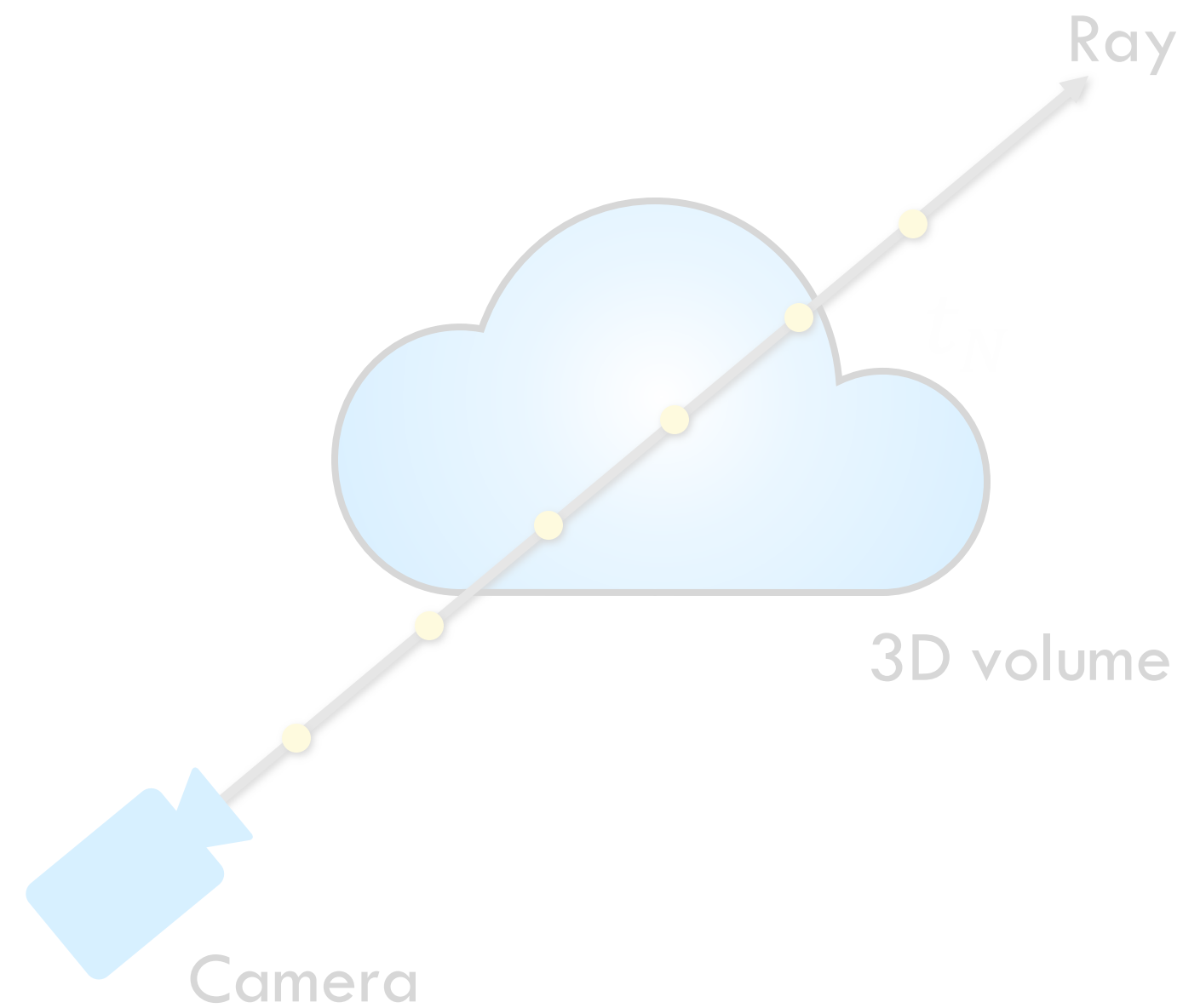
- NeRF can capture non-Lambertian (specular, shiny surfaces) because it models the color in a view-dependent manner
- This is hard to do with meshes unless you model the physical materials & lighting interactions
- But, with Image Based Rendering — All lighting effects are baked in

NeRF in a Slide

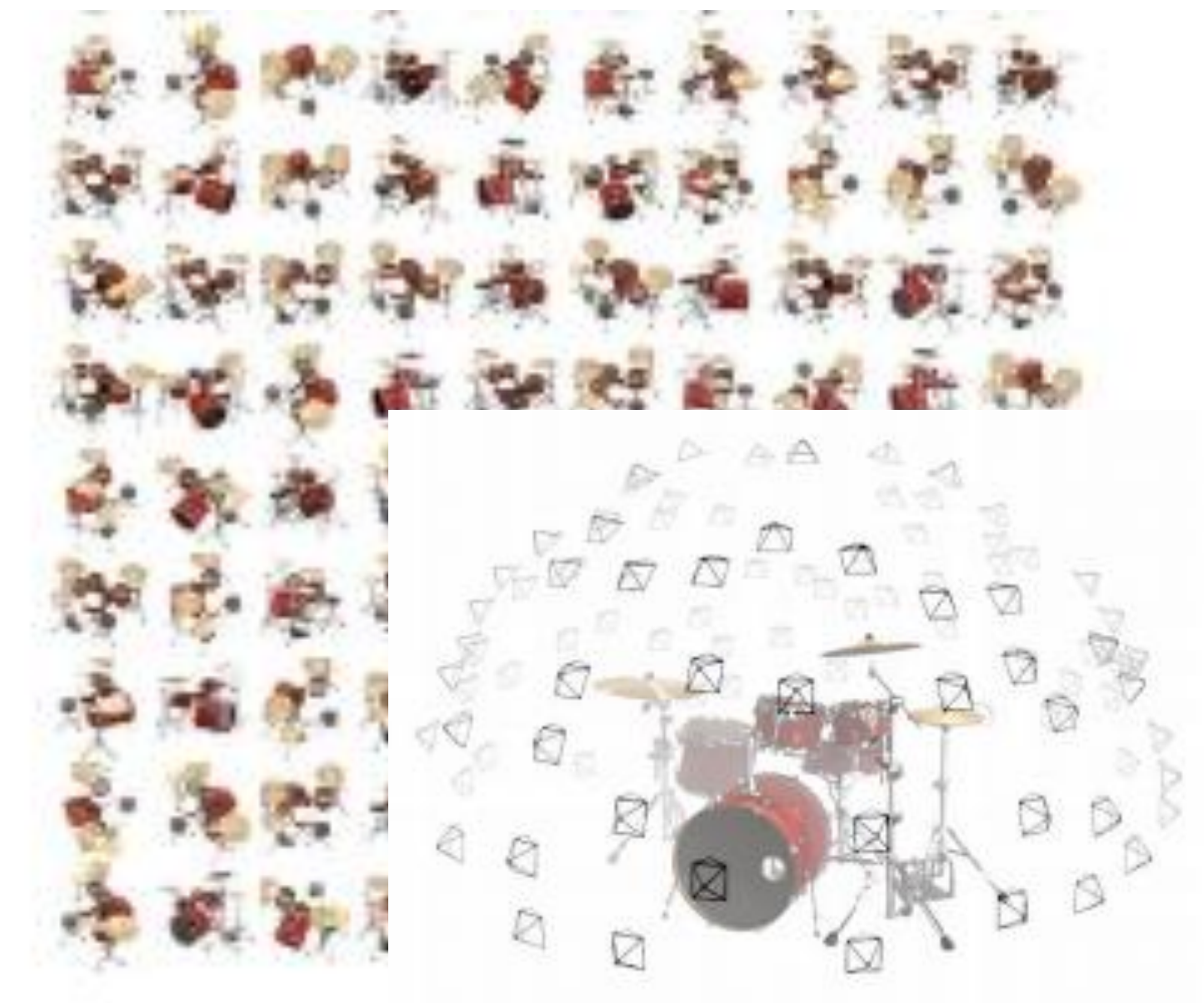
Objective: Reconstruct all training views



Volumetric 3D Scene Representation



Differentiable Volumetric Rendering Function



Optimization via Analysis-by-Synthesis

Unmentioned caveat so far

- Training a NeRF requires a **calibrated** camera!!!!
- Need to know the camera parameters: extrinsic (viewpoint) & intrinsics (focal length, distortion, etc)



How do we get this from images?

Structure from Motion

Or Photogrammetry (1850~)
Long history in Computer Vision

Proc. R. Soc. Lond. B. 203, 405–426 (1979)

Printed in Great Britain

The interpretation of structure from motion

BY S. ULLMAN

*Artificial Intelligence Laboratory, Massachusetts Institute of Technology,
545 Technology Square (Room 808), Cambridge, Massachusetts 02139 U.S.A.*

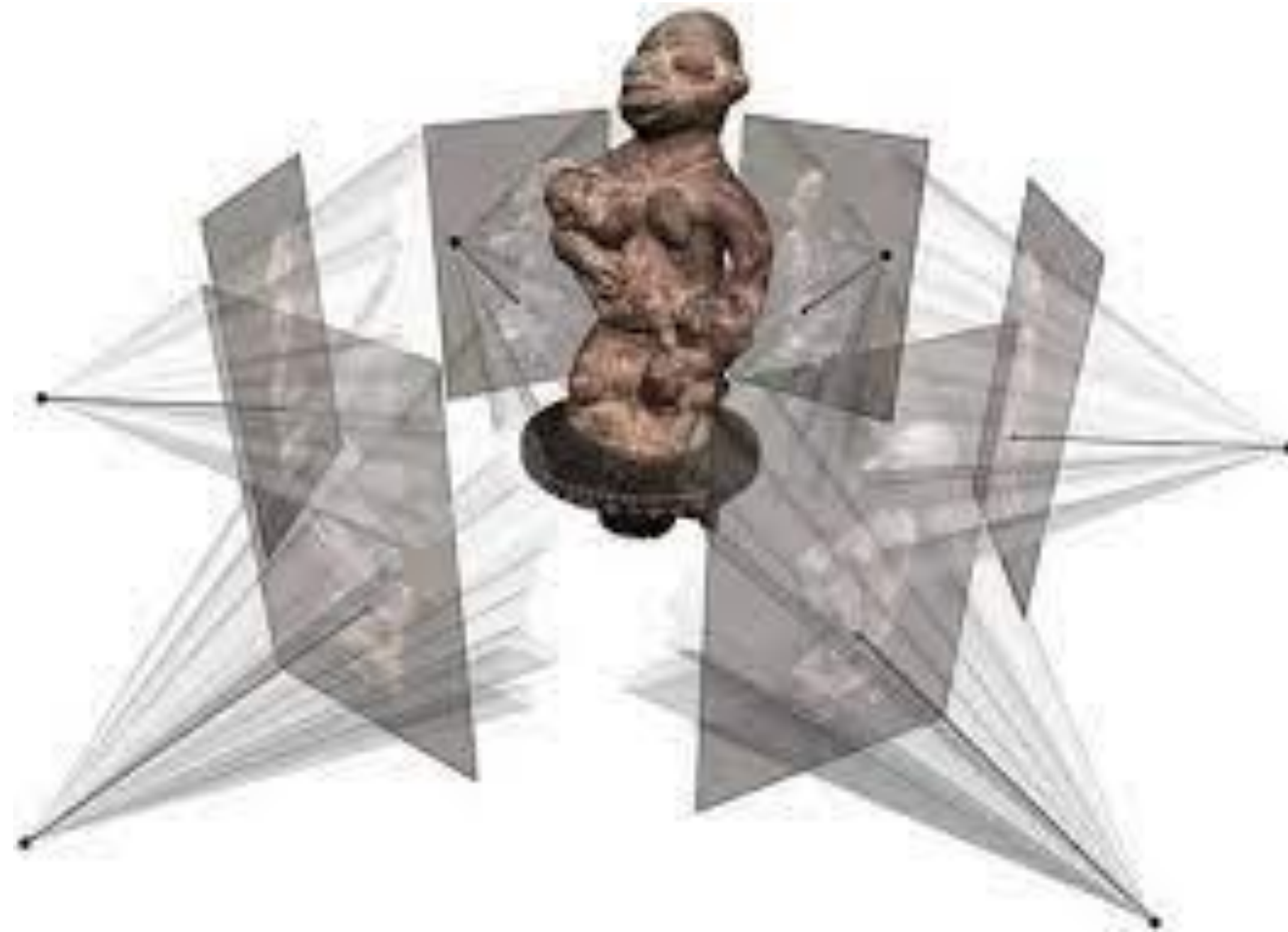
NeRF is AFTER Structure from Motion

- In order to train NeRF you need to run SfM/SLAM on the images to estimate the camera parameters
- In this sense, the problem category is same as that of **Multi-view Stereo**



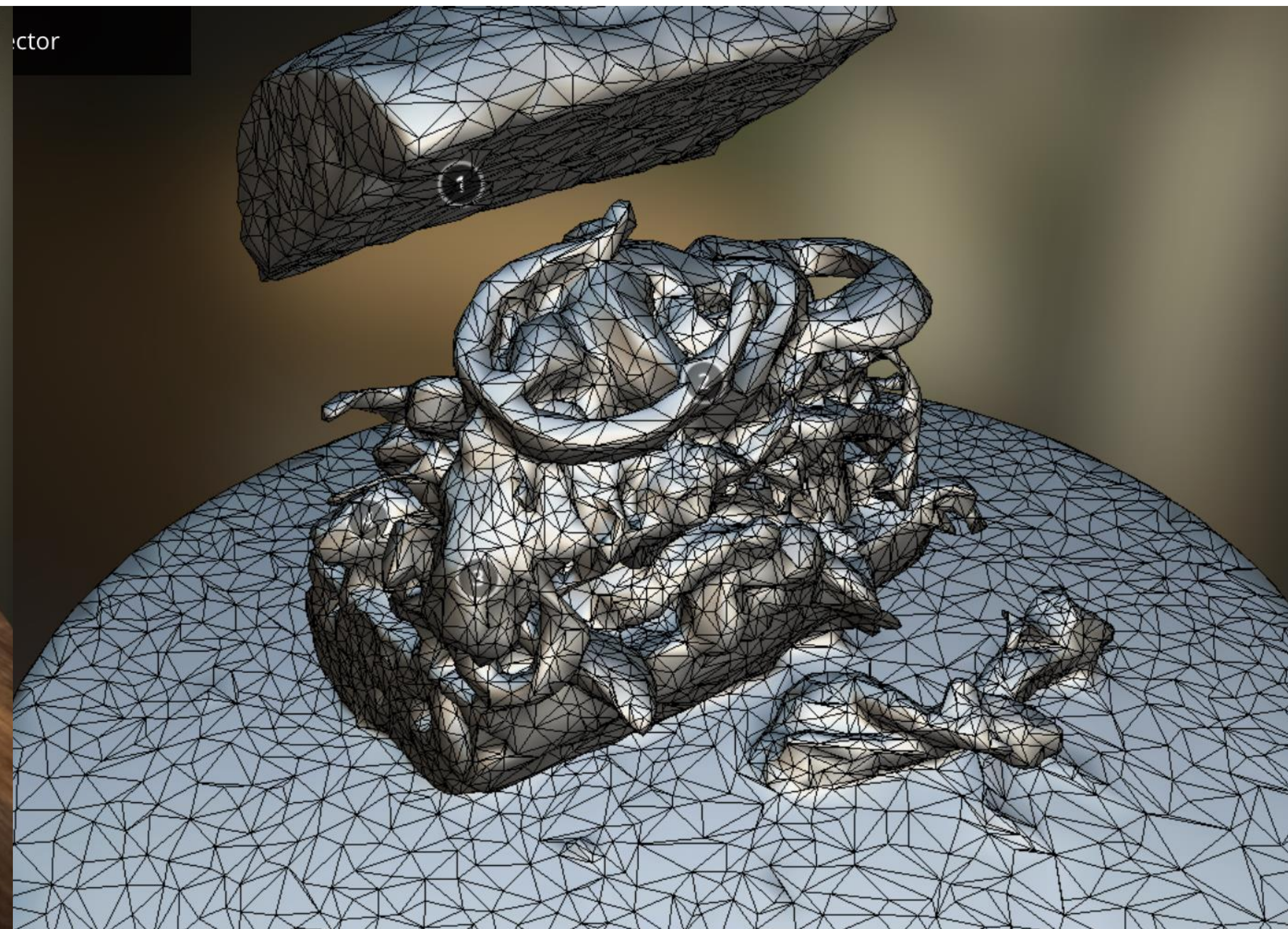
Multi-view Stereo

- Problem: Given calibrated cameras, recover highly detailed 3D **surface** model
- Dense photogrammetry, often the output is textured meshes



Multi-View Stereo

Solutions to MVS is what you see for any existing 3D scanning system, ie sketchfab, or what's in your video game



Multi-View Stereo

Because they often model surfaces, struggles on Thin / Amorphous / Shiny objects





Where NeRF stands

- can do Image Based Rendering well, while also being a 3D representation
- Does not suffer from limitations of surface models
- Easy to optimize from images

Appearance Based
Reconstruction
(Image Based
Rendering)

Physics based
Reconstruction
(3D Surface
Modeling)

NeRFs

Lightfield/Lumigraph
(No 3D representation)

Layered Depth
Images (LDIs)

Multi-Plane
Images (MPIs)

One 3D Surface,
View-Dependent
Texture Mapping

One 3D Surface,
Single Albedo
Texture

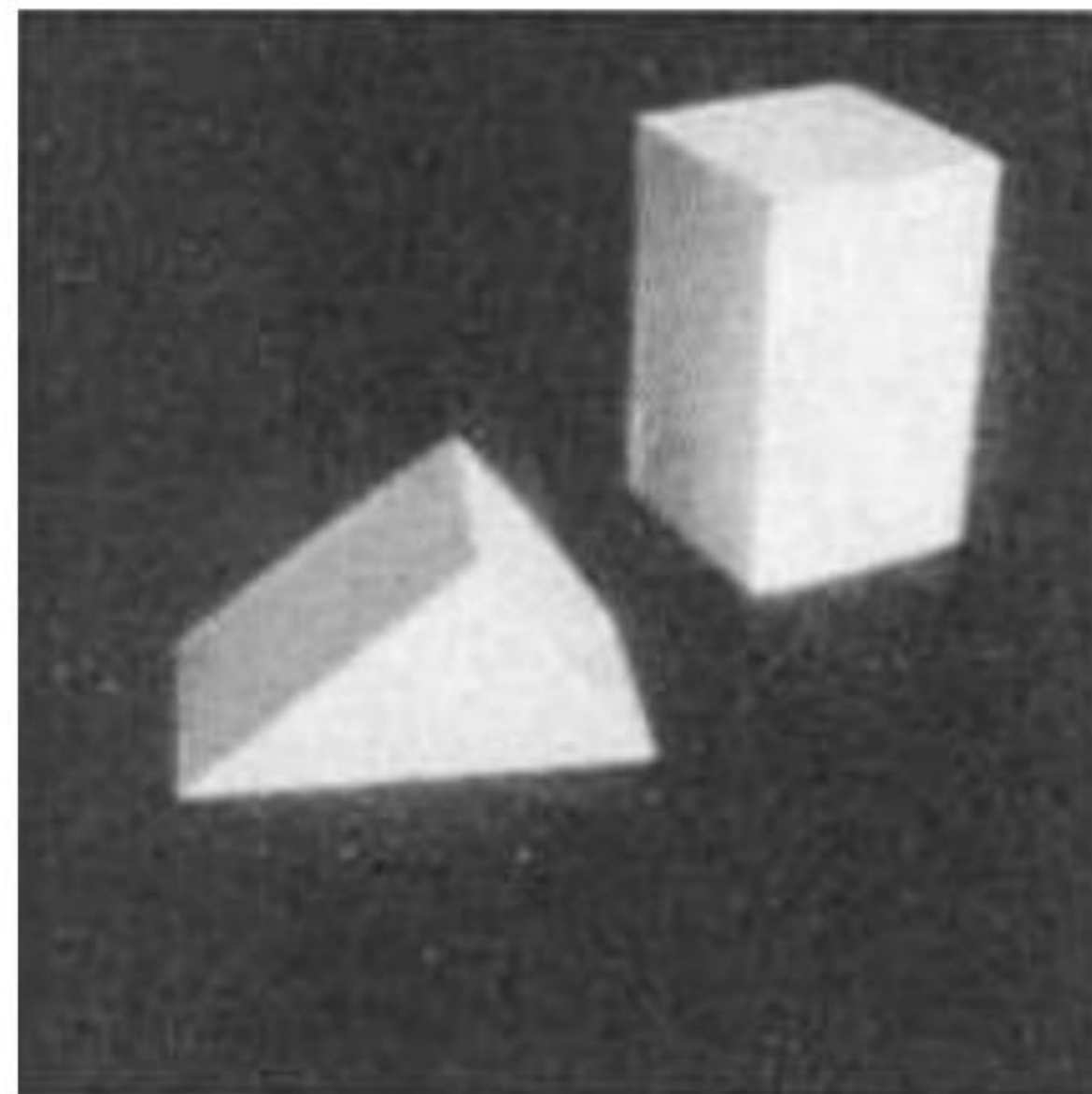
Conventional
Graphics Pipeline

Analysis-by-Synthesis

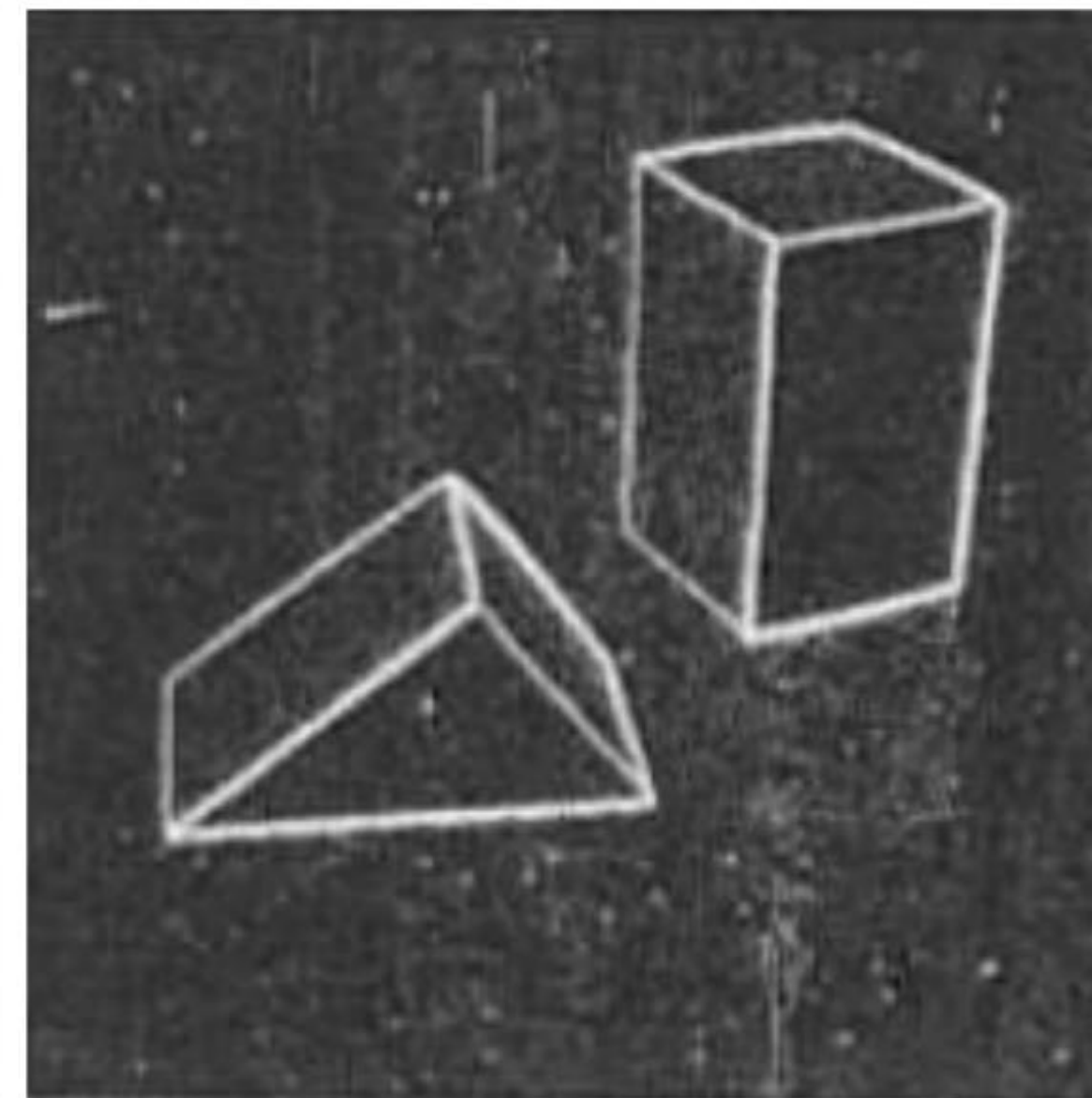


Larry Roberts

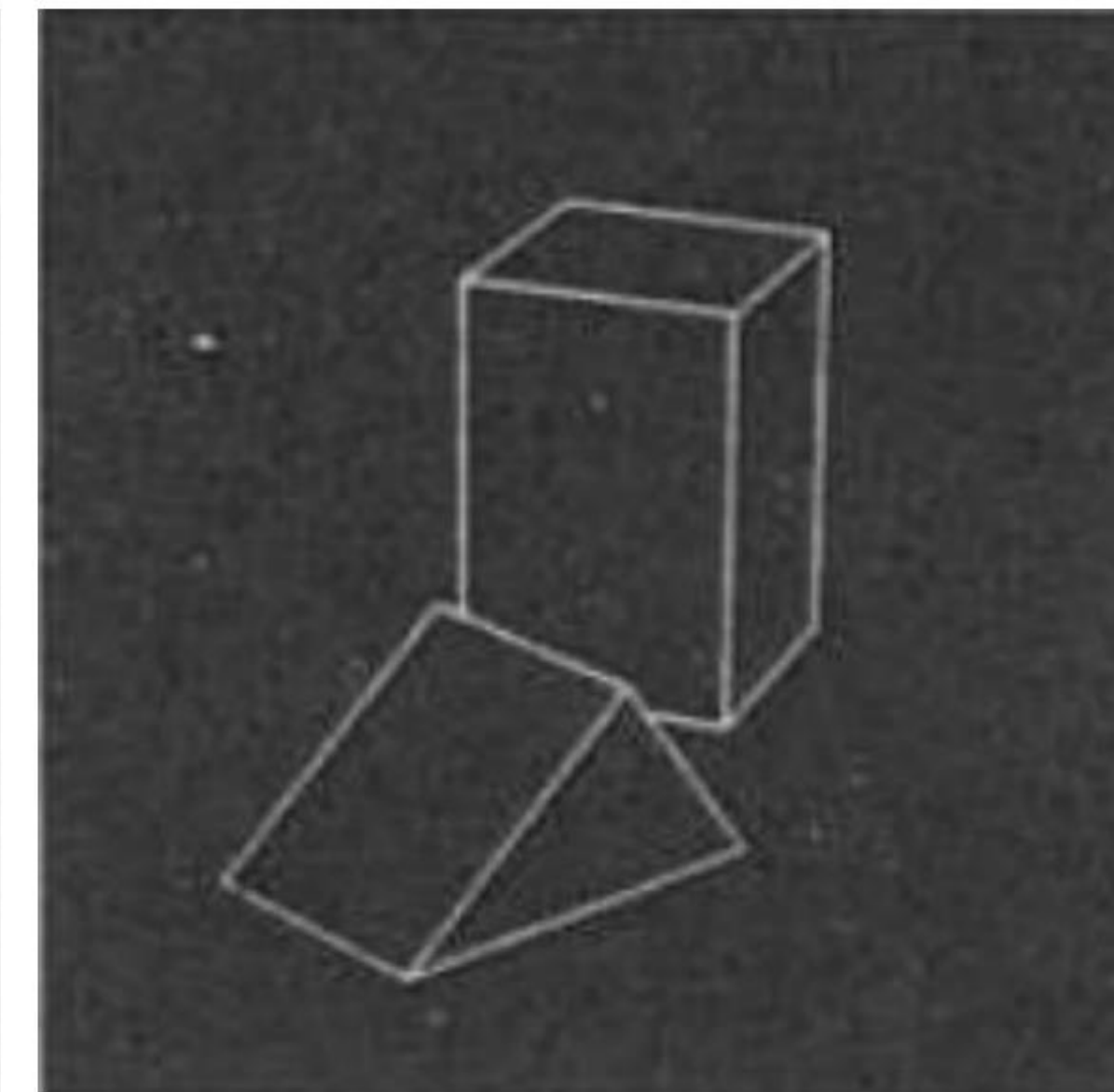
“Father of Computer Vision”



Input image



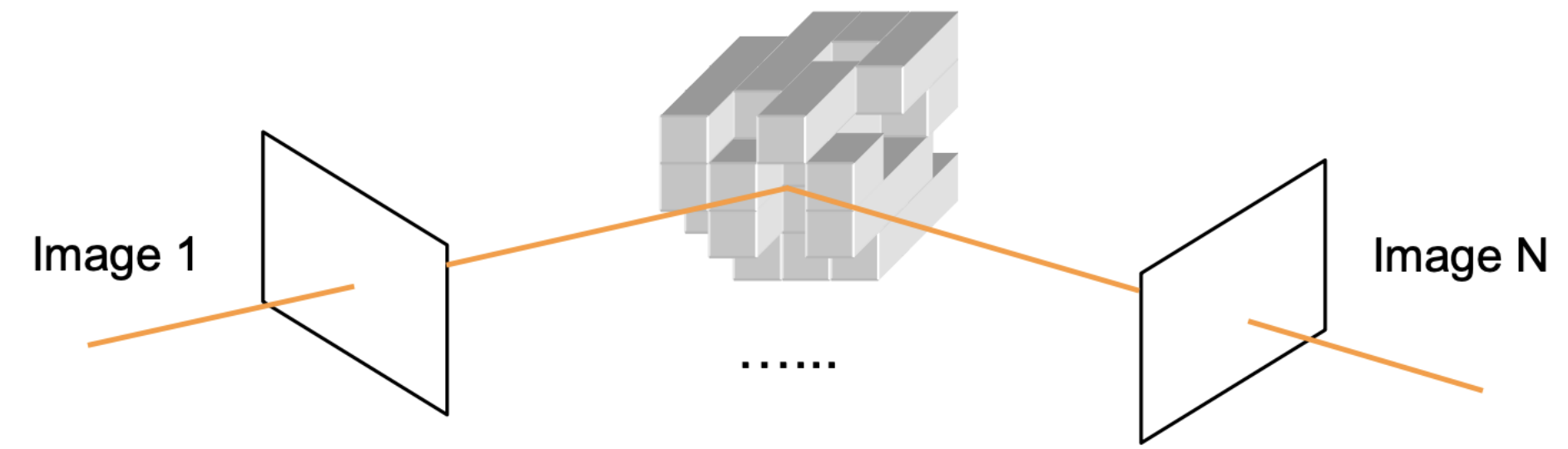
2x2 gradient operator



computed 3D model
rendered from new viewpoint

- History goes way back to the **first** Computer Vision paper!
Roberts: Machine Perception of Three-Dimensional Solids, MIT, 1963

Power of Analysis-by-Synthesis



- Space Carving: A MVS method that used Colored voxels
- But the optimization method was bottom up then.
- Key is optimization via Analysis-by-Synthesis [Plenoxels, Yu et al. 2022]



Input Image (1 of 45)



Reconstruction



Reconstruction



Reconstruction

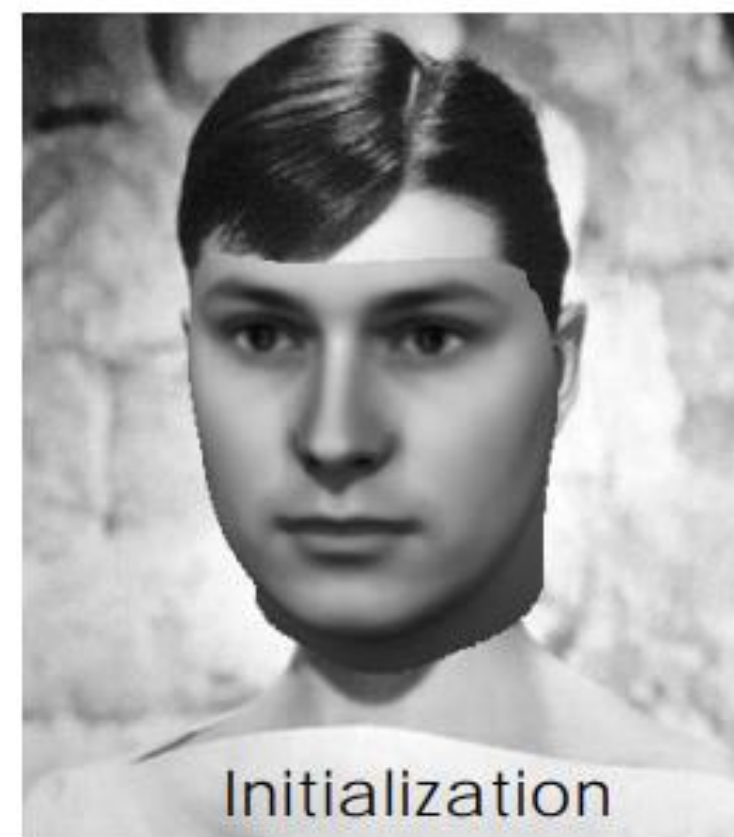
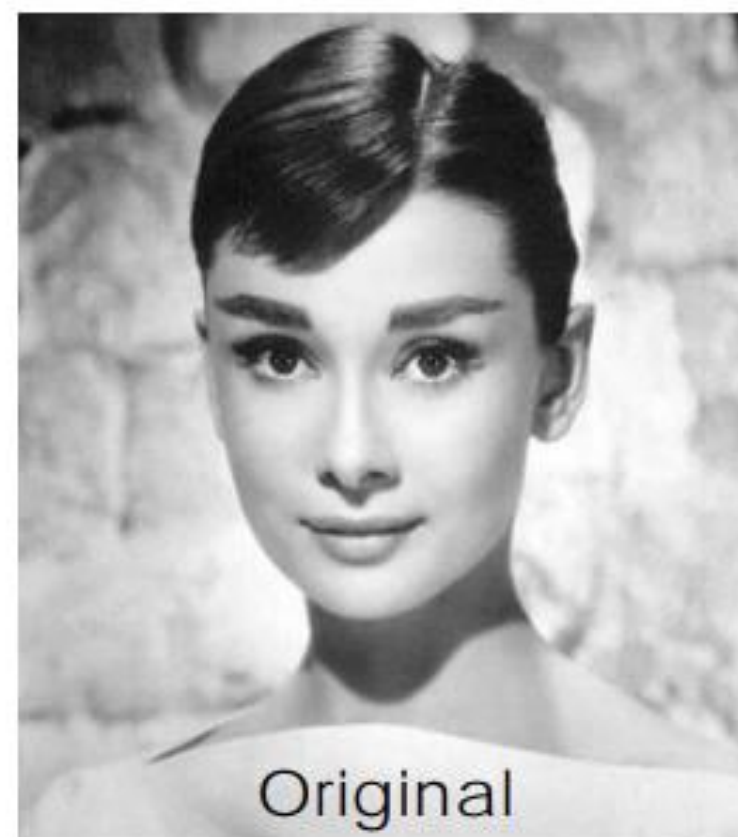


Input Image
(1 of 100)

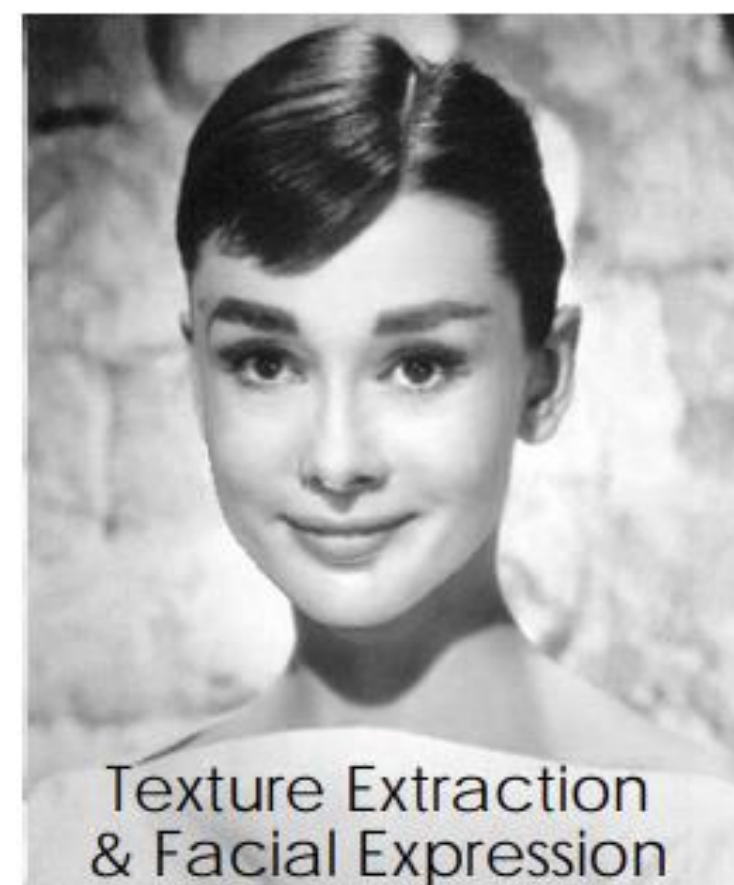


Views of Reconstruction

Analysis-by-Synthesis



3D Reconstruction



Blanz & Vetter 1999

- With custom differentiable renders

Analysis by Synthesis Requires Differentiable Renderers

Next: Deep dive into Volumetric Rendering Function

Where we are

1. Birds Eye View & Background
- 2. Volumetric Rendering Function**
3. Encoding and Representing 3D Volumes
4. Signal Processing Considerations
5. Challenges & Pointers

Volume Rendering

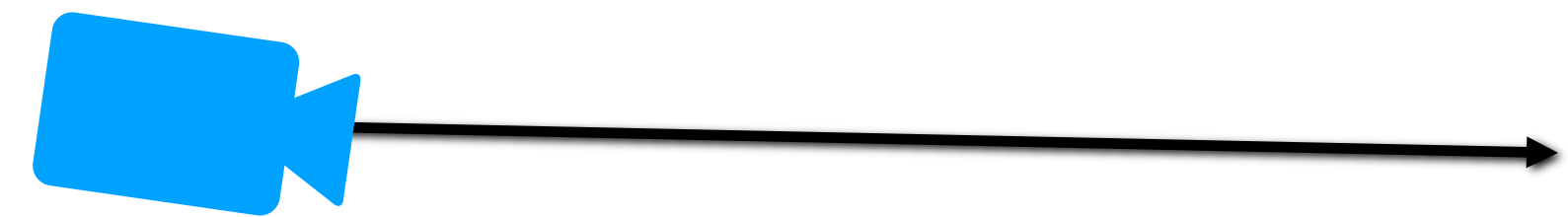
"... in 10 years, all rendering will be volume rendering."

Jim Kajiya at SIGGRAPH '91

Neural Volumetric Rendering

Neural Volumetric Rendering

computing color along rays
through 3D space



What color is this pixel?