3D Vision: Coordinate Spaces



A lot of slides from Noah Snavely +

Shree Nayar's YT series: First principals of Computer Vision

CS180: Intro to Computer Vision and Comp. Photo Angjoo Kanazawa & Alexei Efros, UC Berkeley, Fall 2025

EXCITING NEWS ABOUT CLASS CHOICE AWARDS!

Breaking out of 2D

...now we are ready to break out of 2D







And enter the real world!



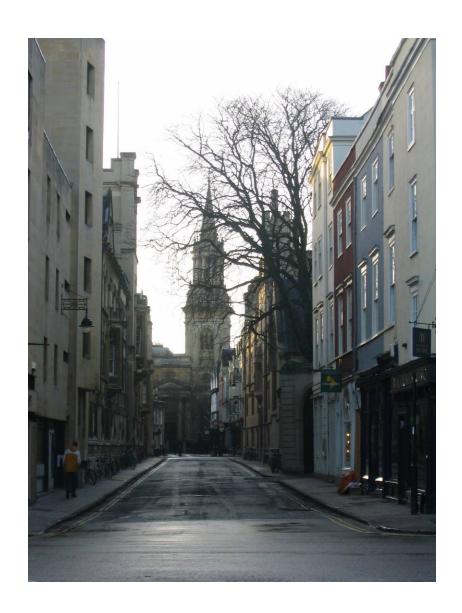
on to 3D...

Enough of images!

We want more of the plenoptic function

We want real 3D scene walk-throughs:

Camera rotation
Camera translation



3D is super cool!



https://rd.nytimes.com/projects/reconstructing-journalistic-scenes-in-3d

3D is super cool!





@capturingreality

@organiccomputer

NeRF in the wild (will get to in few more lectures)



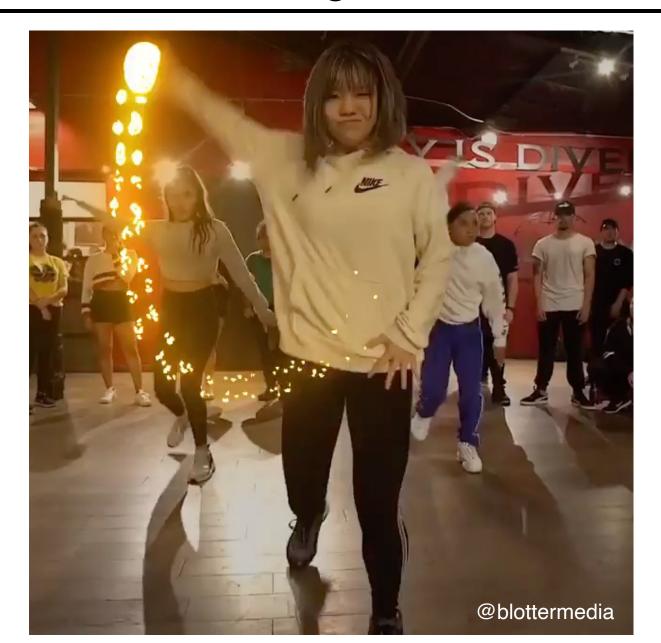
NeRF in the Wild, Martin-Brualla, Radwan et al. CVPR 2021

Not just about 3D reconstruction



[The Chemical Brothers - Wide Open ft. Beck, MV]

3D for video editing



My Research

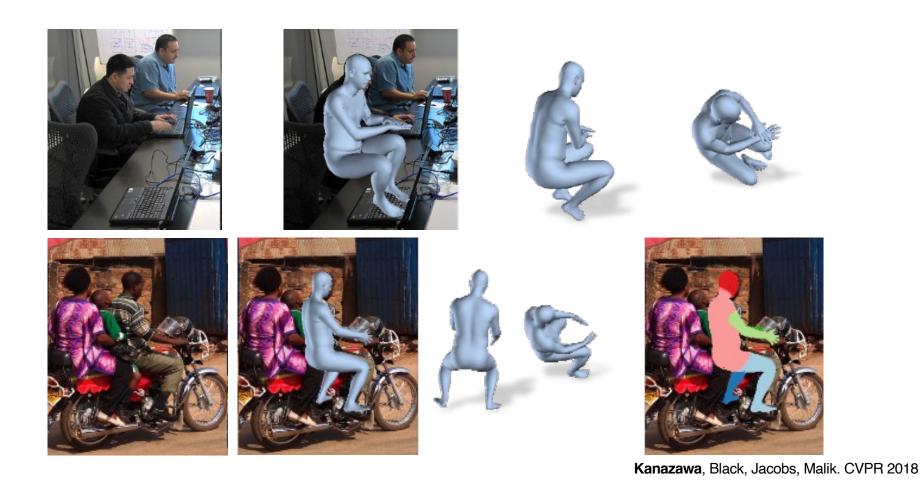
Single-View 3D Human Mesh Recovery



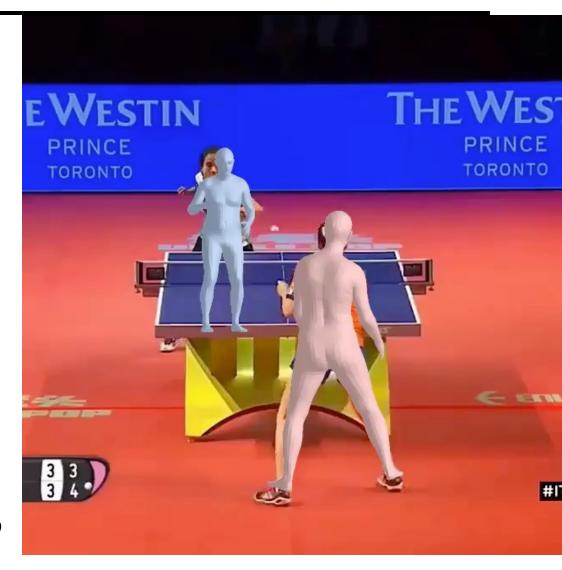


[Bogo*, Kanazawa*, Lassner, Gehler, Romero, Black ECCV '16]

In everyday photos



Or from Video



Kanazawa, Zhang, and Felsen et al. CVPR 2019

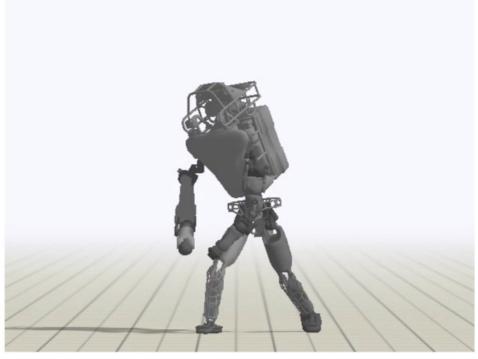
In more detail



Pixel-Aligned Implicit Function for High-Resolution Clothed Human Digitization, Saito, Huang, Natsume, Morishima, **Kanazawa**, Li, ICCV 2019

Teaching robots how to dance from watching YouTube

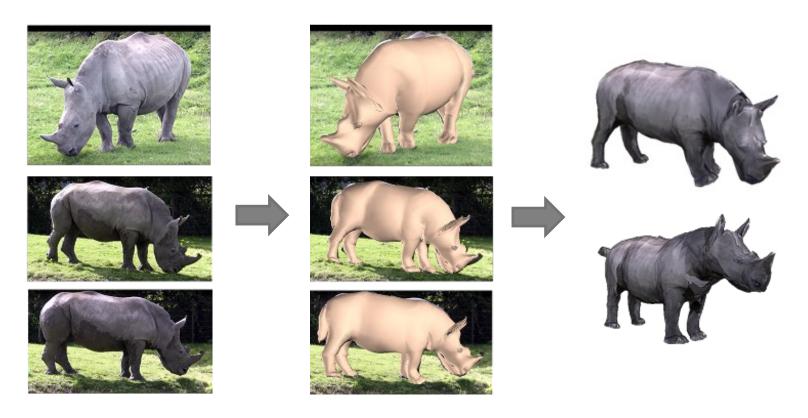




Video Policy

Peng, Kanazawa, Malik, Abbeel, Levine "SFV: Reinforcement Learning of Physical Skills from Videos", SIGGRAPH Asia 2018

Reconstructing Animals with Human Input



Zuffi, Kanazawa, Black, "Lions and Tigers and Bears: Capturing Non-Rigid, 3D, Articulated Shape from Images", CVPR 2018



Flying into an image



Infinite Nature: Perpetual View Generation of Natural Scenes from a Single Image, ICCV 2021

nerfstudio

Matthew Tancik*, Ethan Weber*, Evonne Ng*, Ruilong Li, Brent Yi, Terrance Wang, Alexander Kristoffersen, Jake Austin, Kamyar Salahi, Abhik Ahuja, David McAllister, Angjoo Kanazawa



100+ additional Github contributors



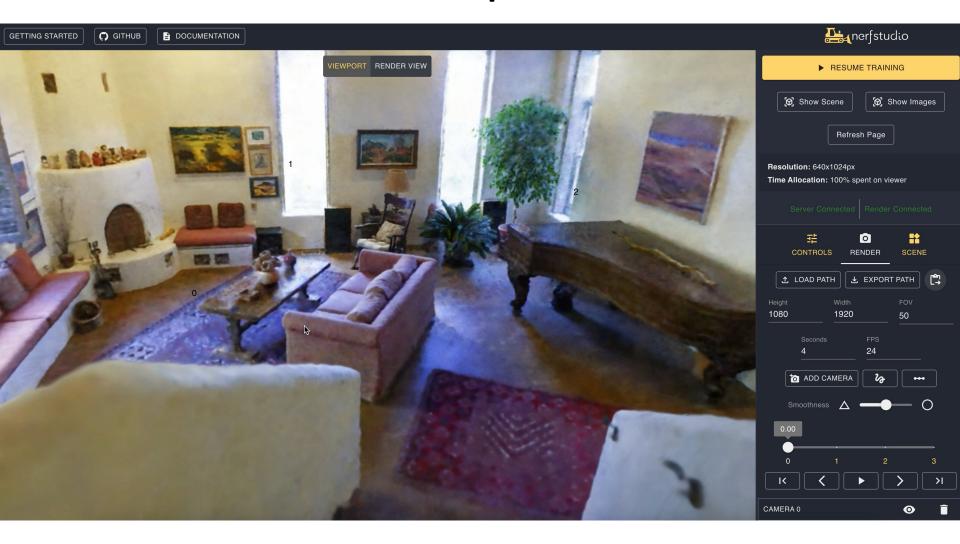


Matt

Ethan

Evonne

3D Capture







so on to 3D...

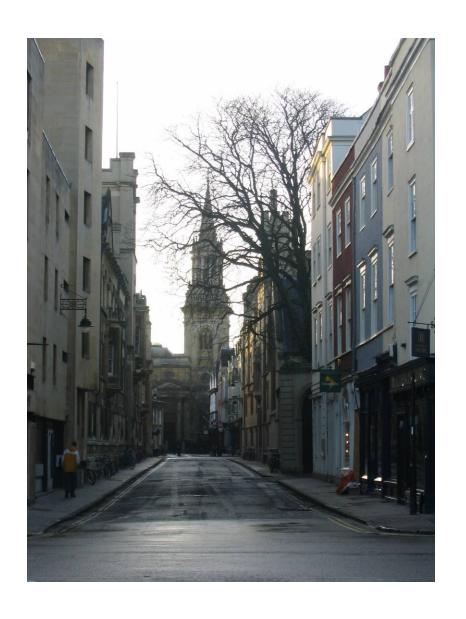
Enough of images!

We want more of the plenoptic function

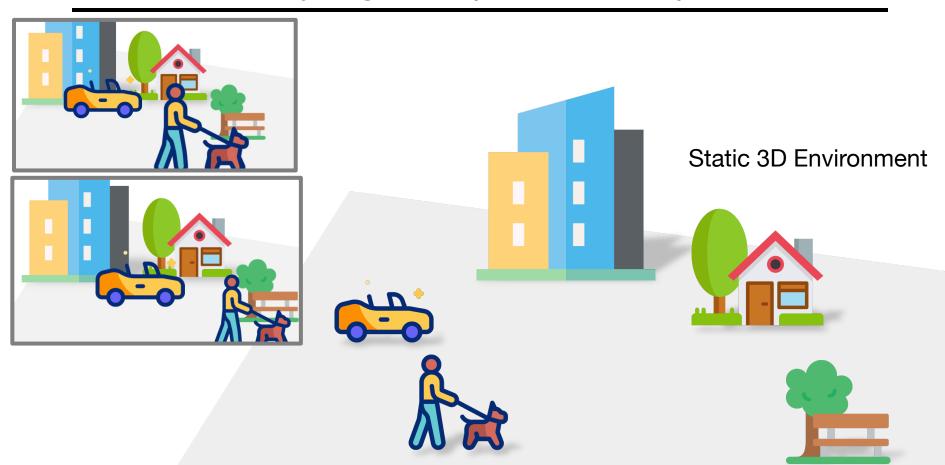
We want real 3D scene walk-throughs:

Camera rotation
Camera translation

Can we do it from a single photograph?



The underlying 4D (3D + time) world

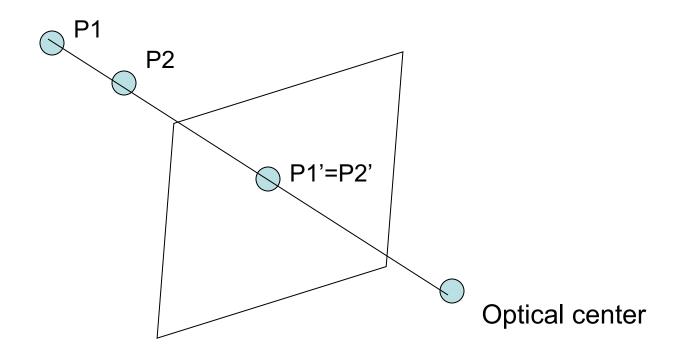


Dynamic 3D Subjects



Do we need multiple views?

 Structure and depth are inherently ambiguous from single views.



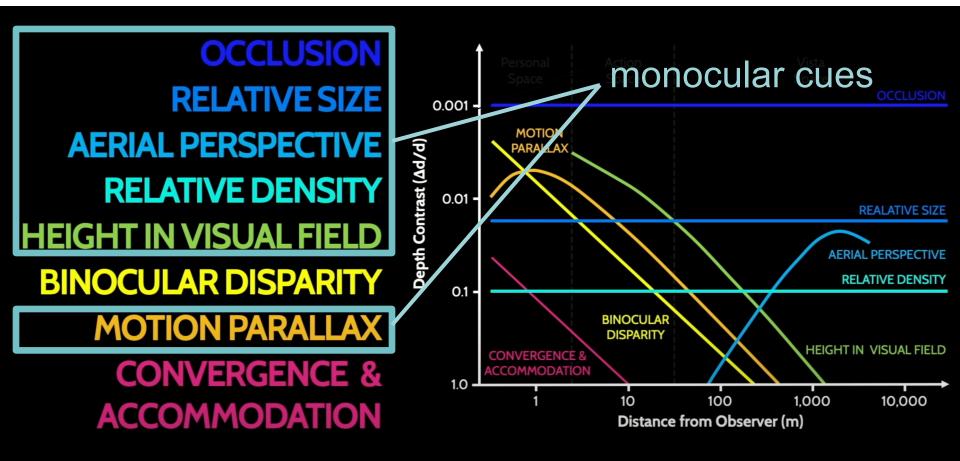
Why multiple views?

Structure and depth are inherently ambiguous from single views.





Human Depth Cues



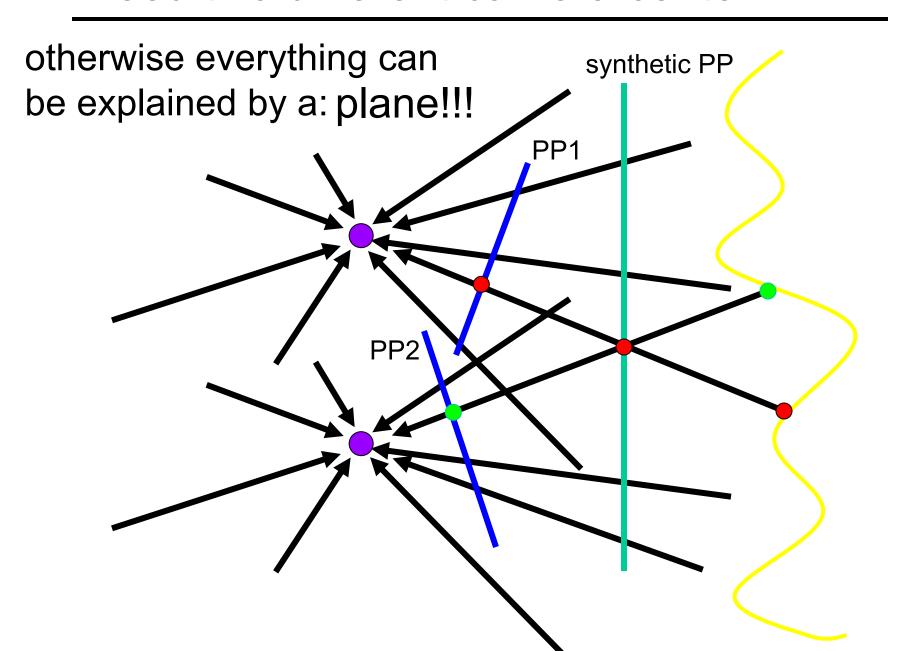
Cutting and Vishton. Perceiving layout and knowing distances. 1995

Geometric Depth Understanding

Ambiguous from a single image

Why?

Need two different camera center



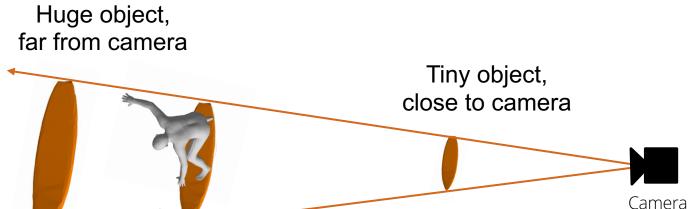
Fundamental Depth Ambiguity in 2D → 3D



Original Image



Same 2D Projection



Infinite Possible 3D Interpretations

2.5D vs 3D

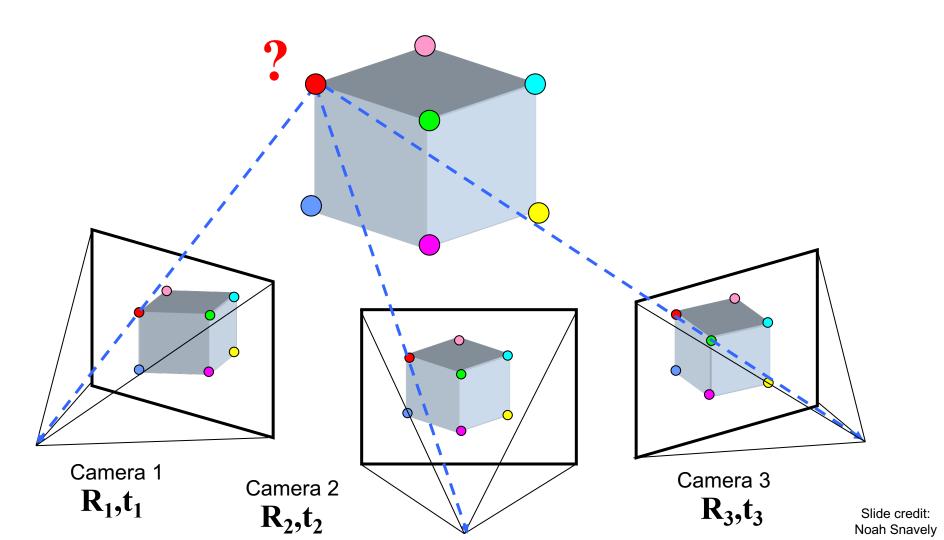
- is 3D = depth from a single image?
- 2.5D = per-pixel depth from a single image





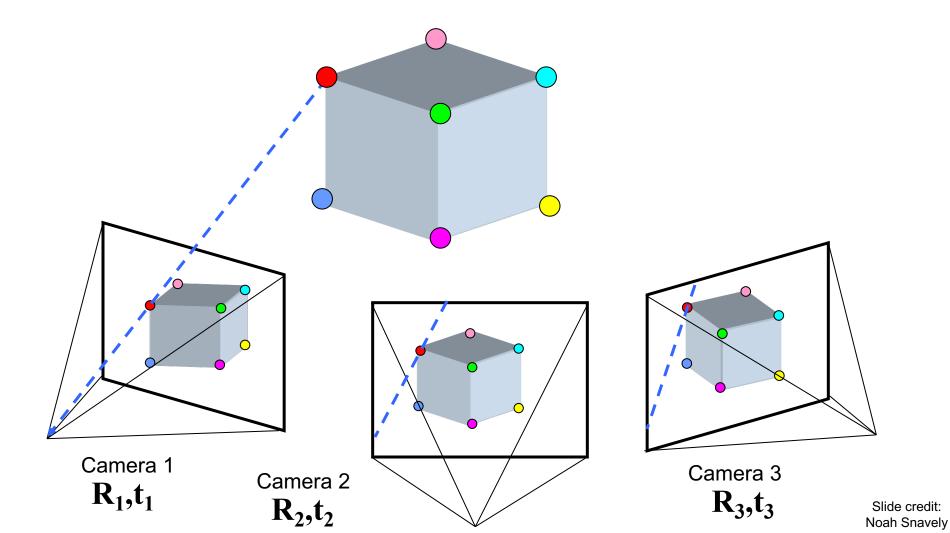
Multi-view geometry problems

• **Structure:** What is the 3D coordinate of a point that can be seen in multiple images?



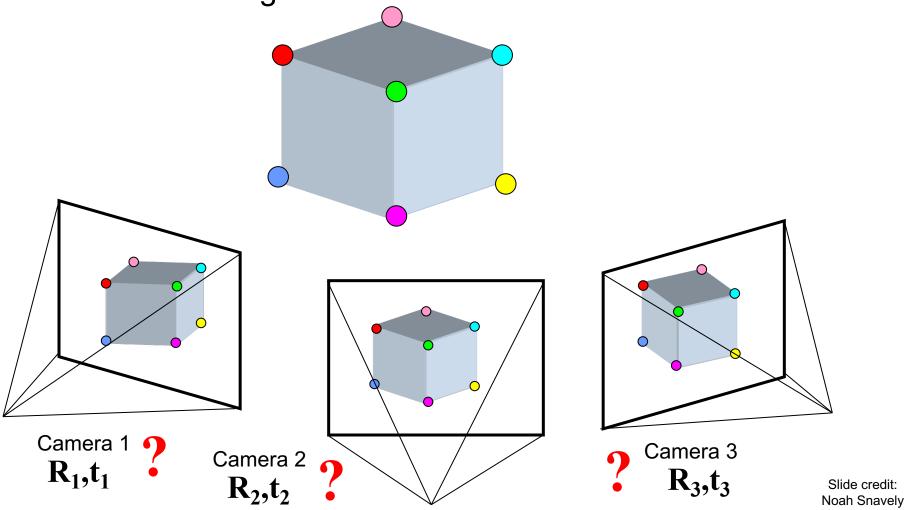
Multi-view geometry problems

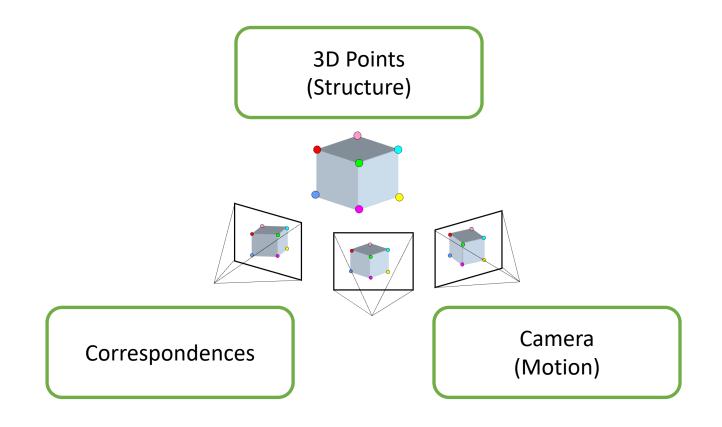
• Correspondence: Given a point in one of the images, where are the corresponding points in the other images?

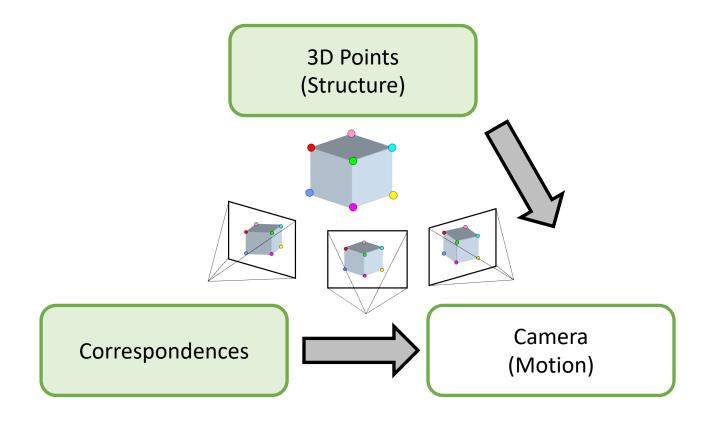


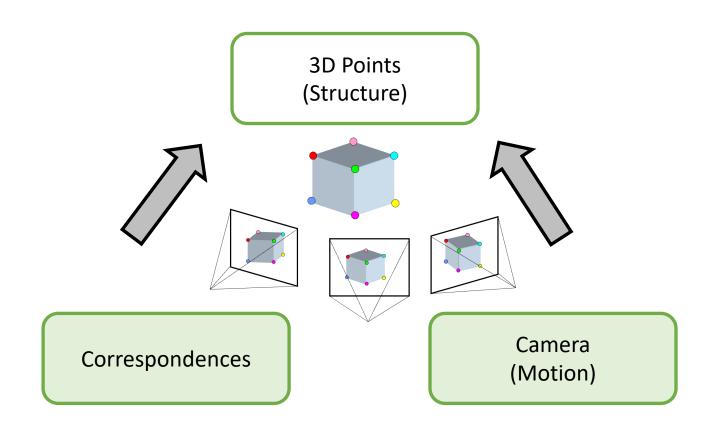
Multi-view geometry problems

 Motion: Given a set of corresponding points in two or more images, what is the relative camera parameters between the images?

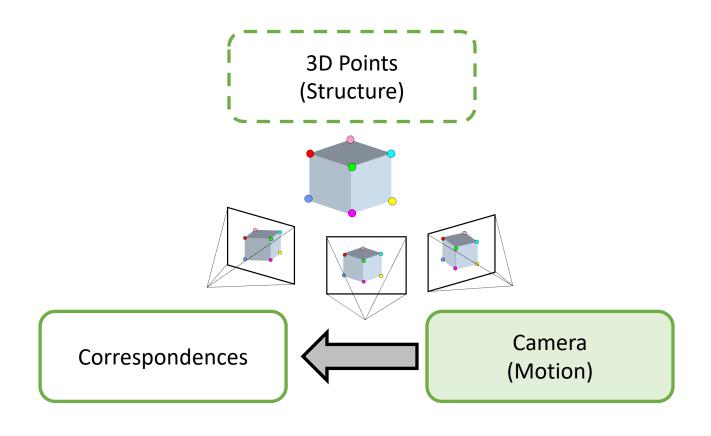






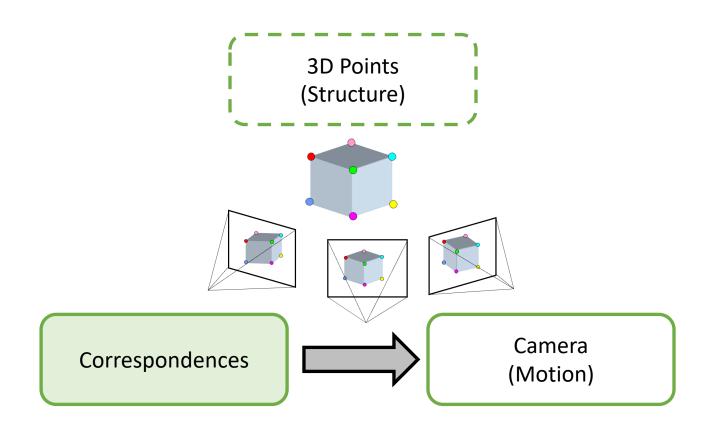








Big picture: 3 key components in 3D





From pixels to the 3D world

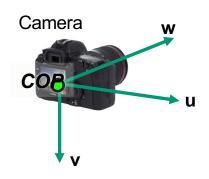
To go from pixels to 3D location in the **world coordinates**, we need to know two things about the camera:

- 1. Position & Orientation of the camera with respect to the world (extrinsics)
- 2. How the camera maps a point in the world to image (intrinsics)

Problem setup

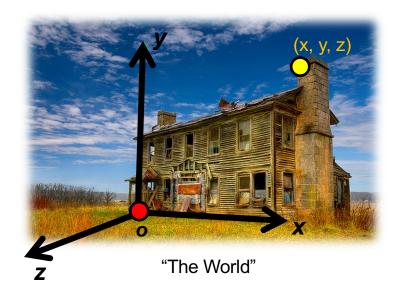
There is a world coordinate frame and camera looking at the world

How can we model the geometry of a camera?



Three important coordinate systems:

- 1. World coordinates
- 2. Camera coordinates
- 3. Image coordinates



Coordinate frames + Transforms

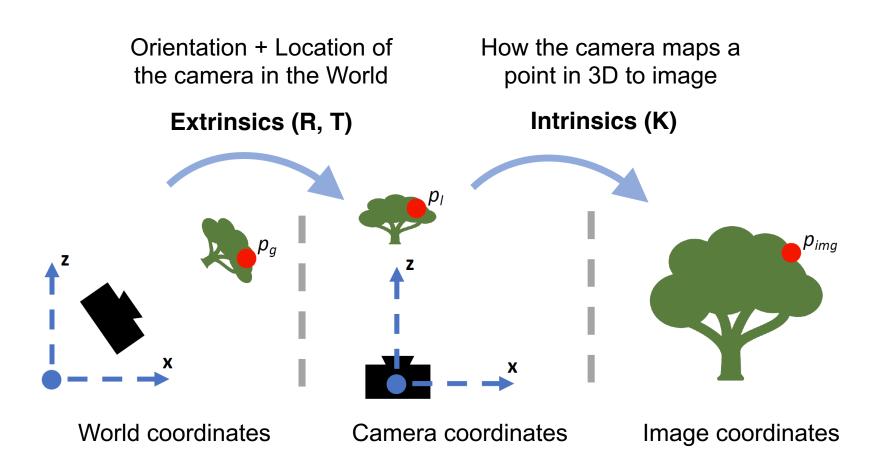


Figure credit: Peter Hedman

Pinhole Camera: Specifics

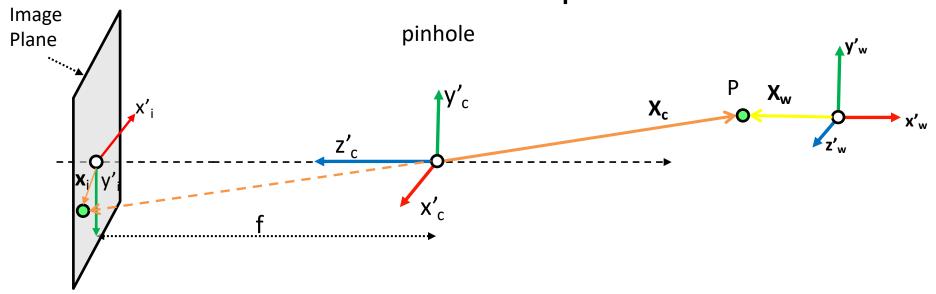


Image Coordinates

Camera Coordinates

World Coordinates

$$\mathbf{x}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}$$

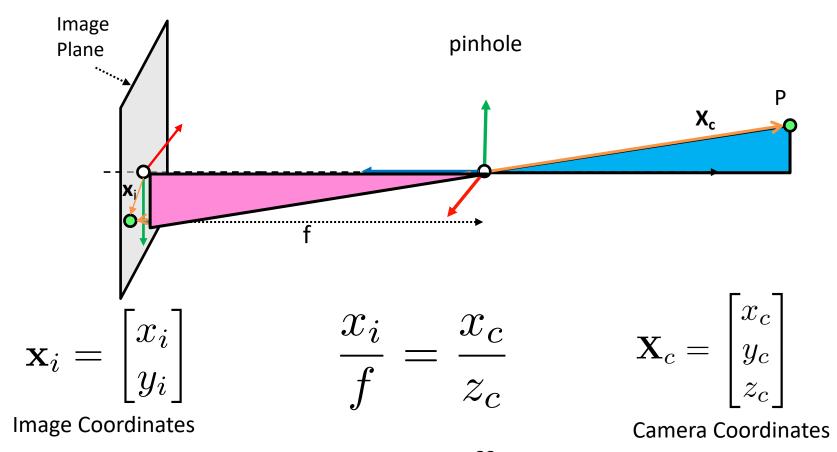
$$\mathbf{X}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

$$\mathbf{X}_w = \begin{bmatrix} x_u \\ y_u \\ z \end{bmatrix}$$

Perspective Projection (3D to 2D)

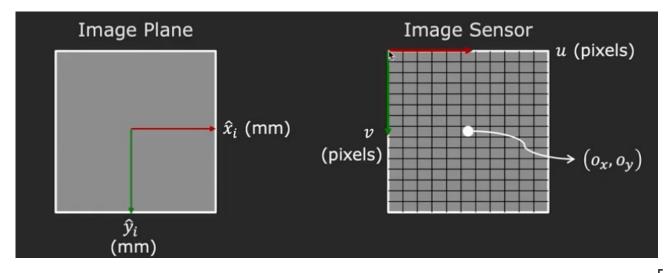
Coordinate
Transformation
(3D to 3D)

Perspective Projection



 $x_i = f \frac{x_c}{z_c}$

Image Plane to Image Sensor Mapping



- 1. Account for pixel density (pixel/mm) & aspect ratio by scalars: $[m_x, m_y] \ m_x x_i, m_y y_i$
- 2. Usually the top left corner is the origin. But in the image plane, the origin is where the optical axis pierces the plane! Need to shift by: (o_x, o_y)

$$u_i = \alpha_x x_i + o_x = \alpha_x f \frac{x_c}{z_c} + o_x$$
 where $[f_x, f_y] = [m_x f, m_y f]$

Pixel Coordinates:

$$u_i = f_x \frac{x_c}{z_c} + o_x$$
 $v_i = f_y \frac{y_c}{z_c} + o_y$

With homogeneous coordinates

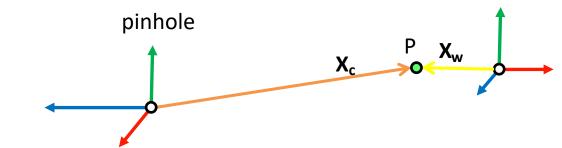
Perspective projection + Transformation to Pixel Coordinates:

$$u_i = f_x \frac{x_c}{z_c} + o_x \quad v_i = f_y \frac{y_c}{z_c} + o_y$$

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} f_x & 0 & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

Intrinsic Matrix

Camera Transformation (3D-to-3D)



Camera Coordinates

$$\mathbf{X}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix} \qquad \qquad \mathbf{X}_w = \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix}$$
Coordinate
Transformation

$$\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix} = \begin{bmatrix} R_{3\times3} & \mathbf{t} \\ \mathbf{0}_{1\times3} & 1 \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$
Extrinsic

Matrix

Putting it all together

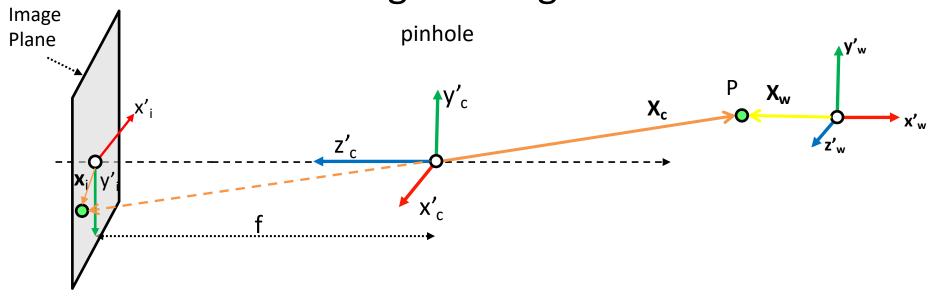


Image Coordinates

Camera Coordinates

World Coordinates

$$\mathbf{x}_i = egin{bmatrix} x_i \ y_i \end{bmatrix}$$
 Perspective Projection $\begin{bmatrix} f_x & 0 & o_x & 0 \ 0 & f_y & o_y & 0 \ 0 & 0 & 1 & 0 \end{bmatrix}$

$$\mathbf{X}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

 $\mathbf{X}_w = \begin{bmatrix} x_w \\ y_w \\ \hat{} \end{bmatrix}$

Coordinate Transformation

$$\begin{bmatrix} R_{3\times3} & \mathbf{t} \\ \mathbf{0}_{1\times3} & 1 \end{bmatrix}$$

Projection Matrix

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} f_x & 0 & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R_{3 \times 3} & \mathbf{t} \\ \mathbf{0}_{1 \times 3} & 1 \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

For completeness, we need to add **skew** (this is 0 unless pixels are shaped like parallelograms – old cameras)

$$K = \begin{bmatrix} f_x & s & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

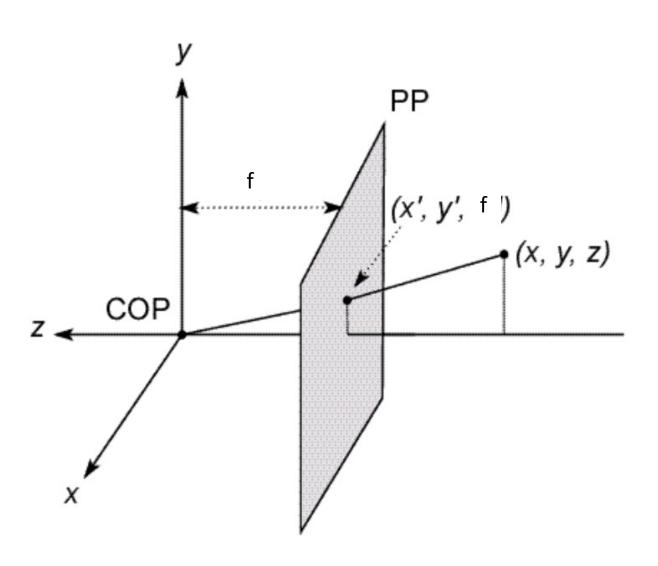
3 x 4 Projection matrix What's the Degrees of Freedom?

Intrinsics: 4 + 1 (skew)

Extrinsic: 3 + 3 = 6

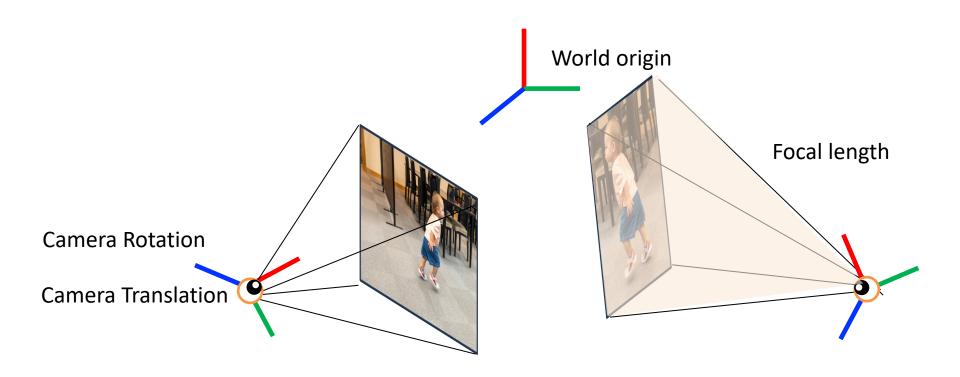
11 unknowns (up to scale)

Another way to draw this

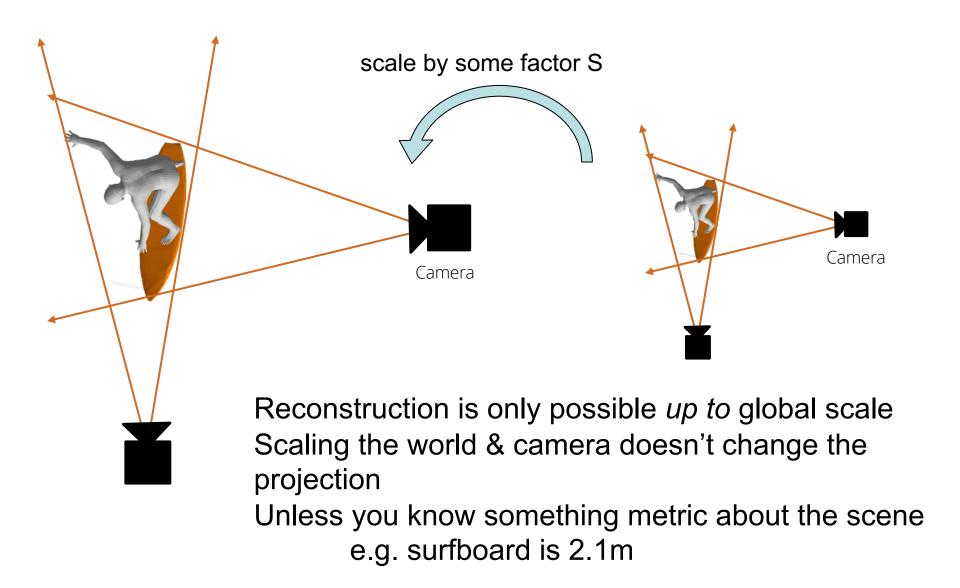


Another picture for the camera

How am I situated in the world (extrinsics) + what is the shape of the ray (intrinsics)



Fundamental Scale Ambiguity



Going from World to Camera

Camera Coordinates

$$\mathbf{X}_c = egin{bmatrix} x_c \ y_c \ z_c \ 1 \end{bmatrix}$$

Extrinsic Matrix:

$$T_{w2c} = \begin{bmatrix} R_{3\times3} & \mathbf{t} \\ \mathbf{0}_{1\times3} & 1 \end{bmatrix}$$

World Coordinates

$$\mathbf{X}_w = egin{bmatrix} x_w \ y_w \ z_w \ 1 \end{bmatrix}$$

$$\mathbf{X}_c = T_{w2c}\mathbf{X}_w$$

Going from Camera to World

Camera Coordinates

$$\mathbf{X}_c = egin{bmatrix} x_c \ y_c \ z_c \ 1 \end{bmatrix}$$

World Coordinates

$$\mathbf{X}_w = egin{bmatrix} x_w \ y_w \ z_w \ 1 \end{bmatrix}$$

Extrinsic Matrix:

$$T_{w2c} = \begin{bmatrix} R_{3\times3} & \mathbf{t} \\ \mathbf{0}_{1\times3} & 1 \end{bmatrix}$$

$$T_{w2c}^{-1}\mathbf{X}_c = \mathbf{X}_w$$

Camera to Image

Image Coordinates

$$\mathbf{x}_i = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix}$$

Intrinsics Matrix:

$$K = \begin{bmatrix} f_x & s & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Camera Coordinates

$$\mathbf{X}_c = egin{bmatrix} x_c \ y_c \ z_c \ 1 \end{bmatrix}$$

$$\mathbf{x}_i = K\mathbf{X}_c$$

Image to Camera?

Image Coordinates

$$\mathbf{x}_i = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \equiv \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix}$$

Camera Coordinates

$$\mathbf{X}_c = egin{bmatrix} x_c \ y_c \ z_c \ 1 \end{bmatrix}$$

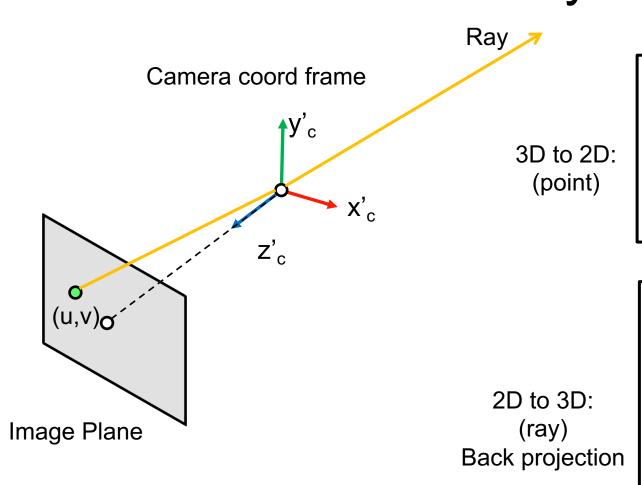
Intrinsics Matrix:

$$K = \begin{bmatrix} f_x & s & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

$$u = f_x \frac{x_c}{z_c} + o_x \qquad \longrightarrow \qquad x = \frac{z}{f_x} (u - o_x)$$

What's the problem?

We don't know the depth! but at the least it will be: on the ray!



$$u = f_x \frac{x_c}{z_c} + o_x$$
$$v = f_y \frac{y_c}{z_c} + o_y$$

$$x = \frac{z}{f_x}(u - o_x)$$

$$y = \frac{z}{f_y}(v - o_y)$$

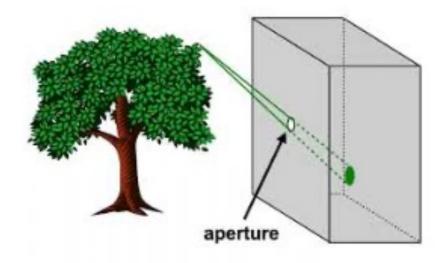
$$z > 0$$

What is your coordinate space?

- In the next project (and in life) always make sure you're in the right coordinate space.
- Q: Which space is the ray defined in?

$$x = \frac{z}{f_x}(u - o_x)$$
 2D to 3D:
$$(\text{ray})$$
 Back projection
$$z > 0$$

Watch these 5 min videos



https://www.youtube.com/watch?v=F5WA26W4JaM https://www.youtube.com/watch?v=g7Pb8mrwcJ0 For the following, assume camera is known through a process called calibration (will get there soon)

We will also only consider 2 views to start with

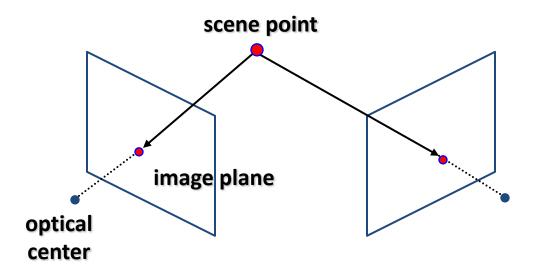
Question

• If we know the camera, can we guess the depth of every point in the image?

A: Yes, with correspondences

Estimating depth with stereo

- Stereo: shape from "motion" between two views
- We'll need to consider:
 - 1. Camera pose ("calibration")
 - 2. Image point correspondences







Stereo vision



Two cameras, simultaneous views

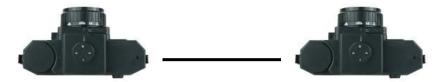


Single moving camera and static scene

Simple Stereo Setup

- Assume parallel optical axes
- Two cameras are calibrated
- Find relative depth

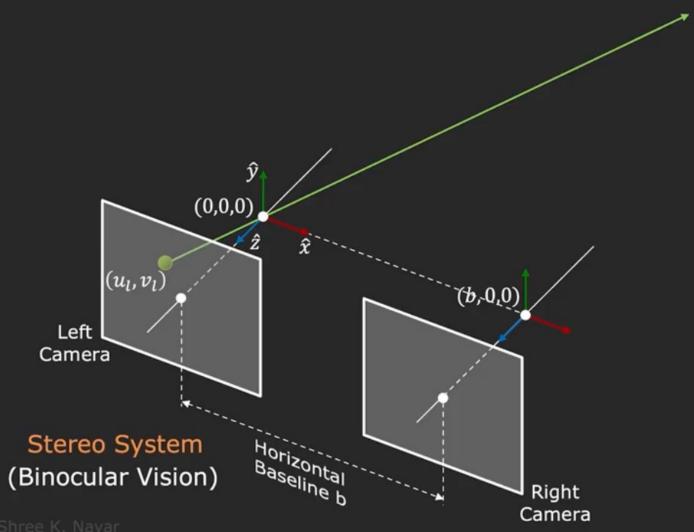




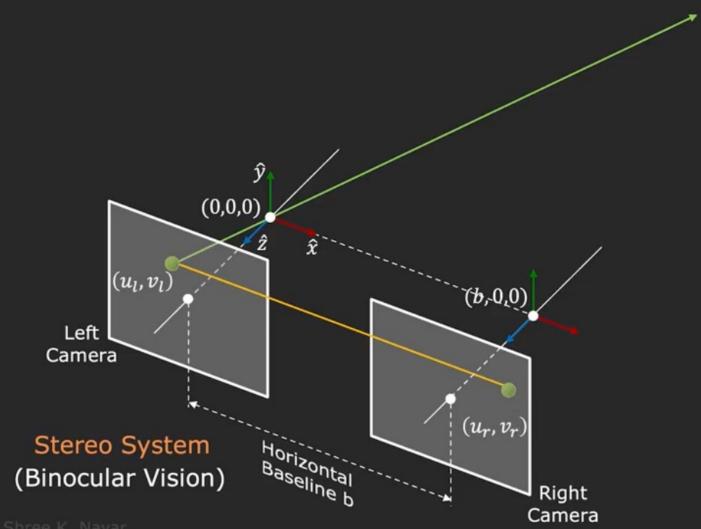
Key Idea: difference in corresponding points to understand shape

Slide credit: Noah Snavely

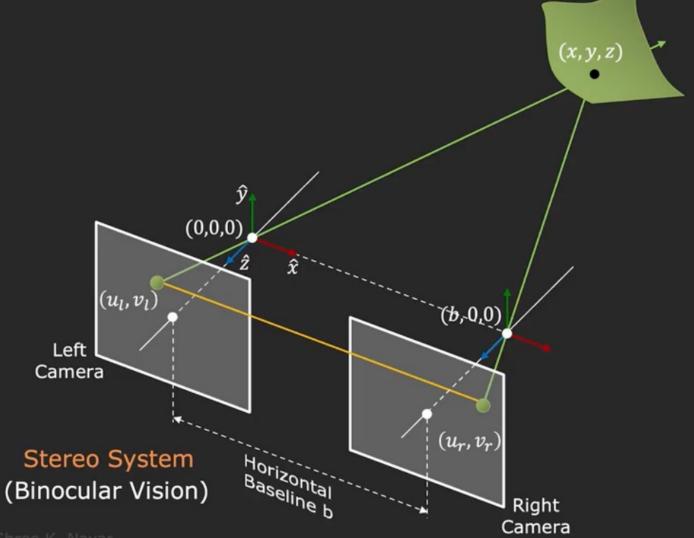
Triangulation using two cameras



With known correspondence



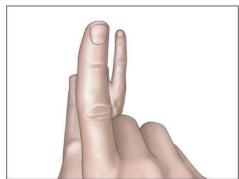
Triangulation using two cameras



We are equipped with binocular vision. Let's try!





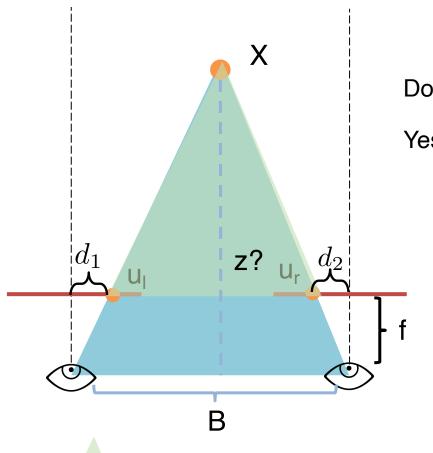


Right retinal image



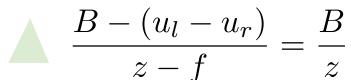
Left retinal image

Solving for Depth in Simple Stereo



Do we have enough to know what is Z?

Yes, similar triangles!



$$z = \frac{fB}{u_l - u_r}$$

disparity (how much corrsp. pixels move)

Base of : $B - (d1 + d_2)$

 $_{\text{coordinates:}}^{\text{in image}} = B - (u_l - u_r)$

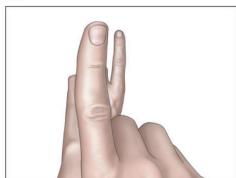




Try with your hands!



(b)

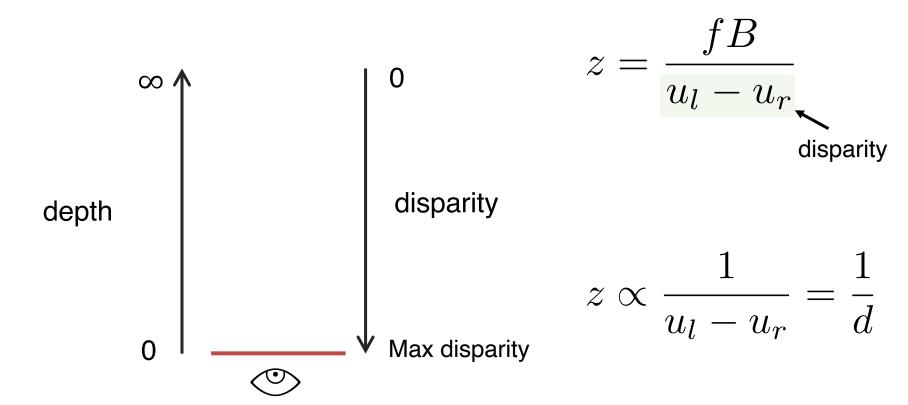


Right retinal image



Left retinal image

Depth is inversely proportional to disparity



what is the disparity of the closer point? what is the disparity of the far away point? Disparity gives you the depth information!

Try again

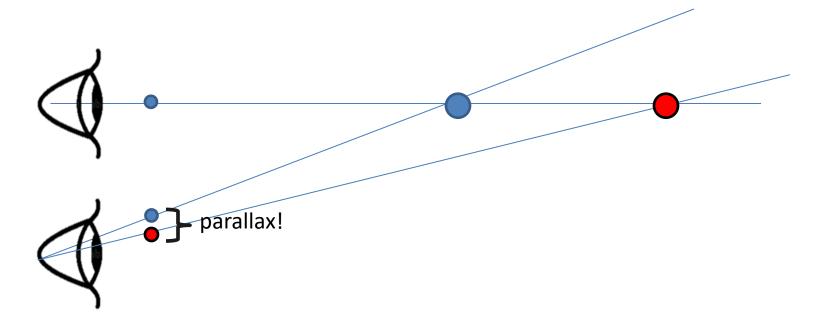
- 1. Setup so your fingers are on the same line of sight from one eye
- 2. Now look in the other eye They move!

Relative displacement is higher as the relative distance grows

== Parallax



Parallax

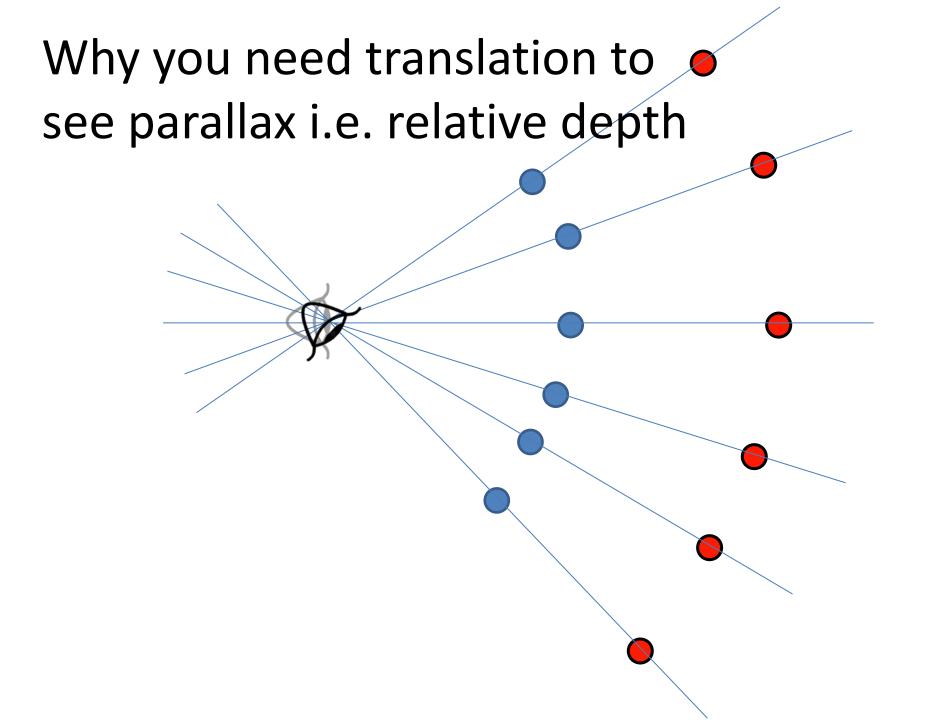


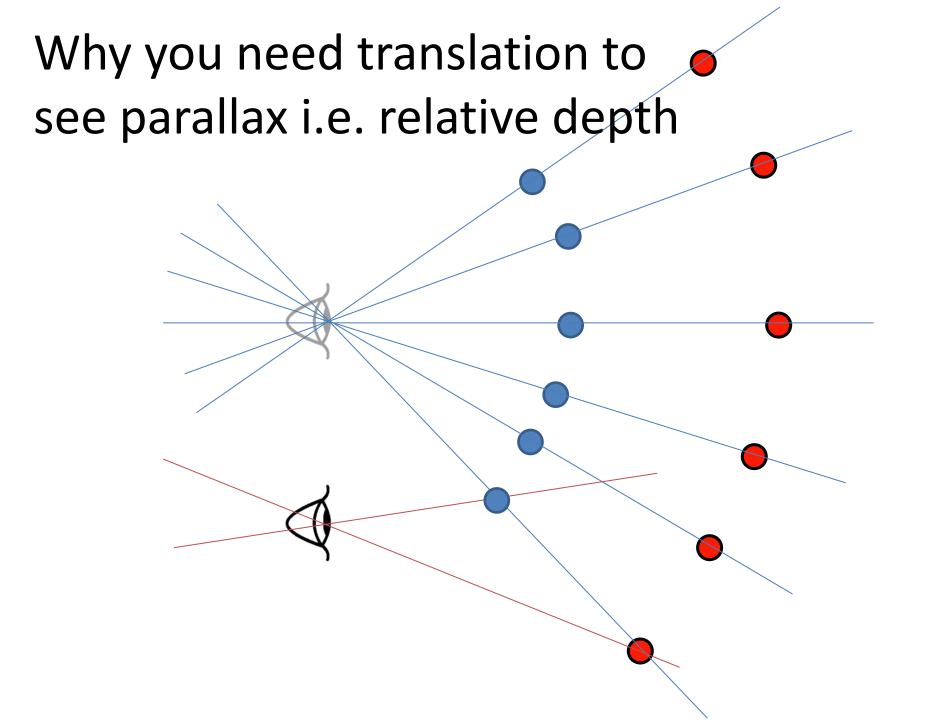
Parallax = from ancient Greek parállaxis

= Para (side by side) + allássō, (to alter)

= Change in position from different view point

Two eyes give you parallax, you can also move to see more parallax = "Motion Parallax"





Stereo Matching: Finding Disparities

Goal: Find the disparity between left and right stereo pairs.



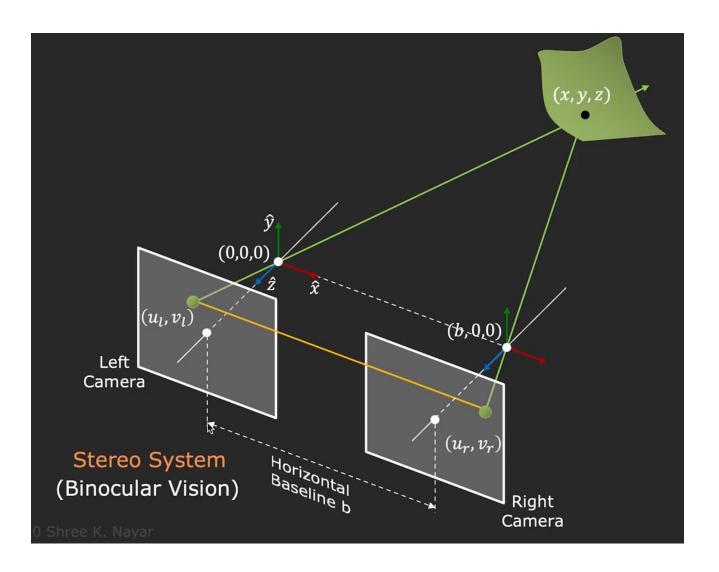
Left/Right Camera Images



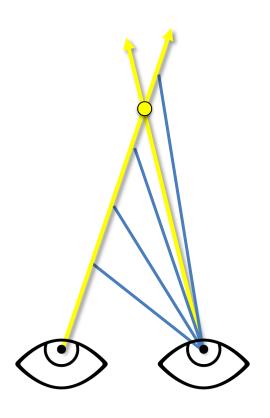
Disparity Map (Ground Truth)

Where is the corresponding point going to be?

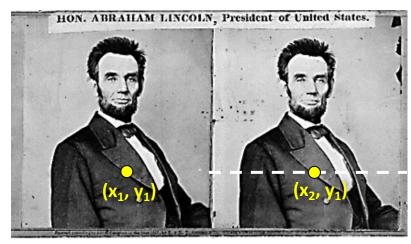
Hint



Epipolar Line



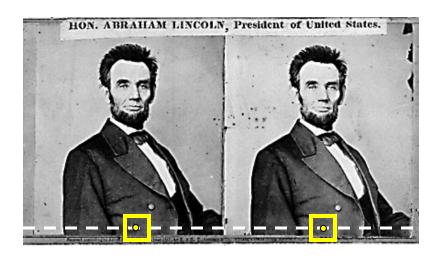
epipolar lines



Two images captured by a purely horizontal translating camera (rectified stereo pair)

 x_1-x_2 = the *disparity* of pixel (x_1, y_1)

Your basic stereo algorithm



For every epipolar line:

For each pixel in the left image

- · compare with every pixel on same epipolar line in right image
- pick pixel with minimum match cost

Improvement: match windows, + clearly lots of matching strategies

Your basic stereo algorithm

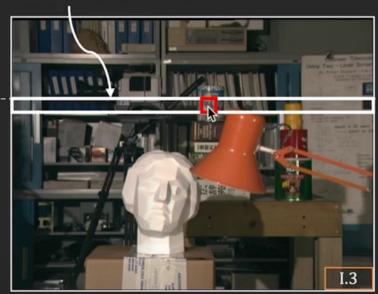
Determine Disparity using Template Matching

Template Window T



Left Camera Image E_l

Search Scan Line L

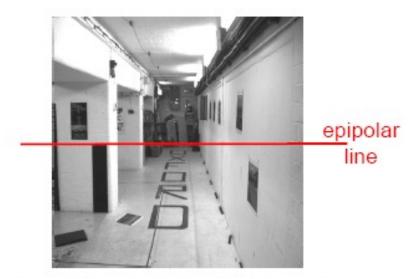


Right Camera Image E_r

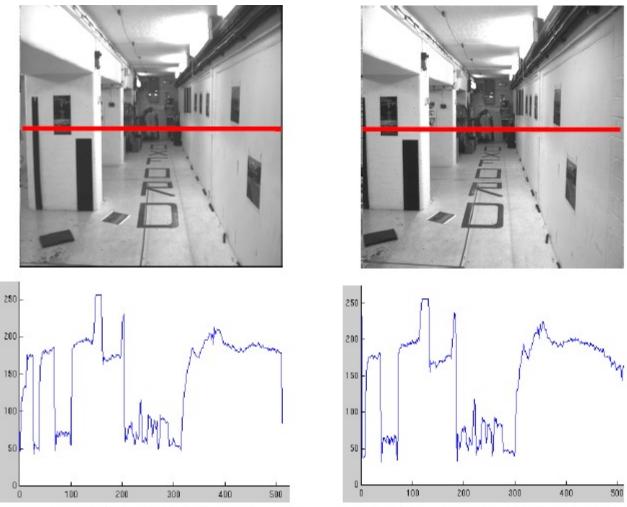
Correspondence problem

Parallel camera example - epipolar lines are corresponding rasters



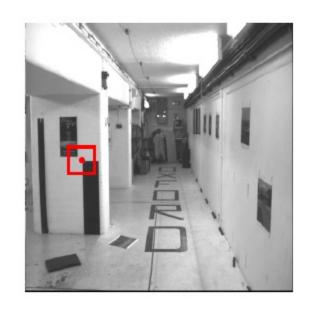


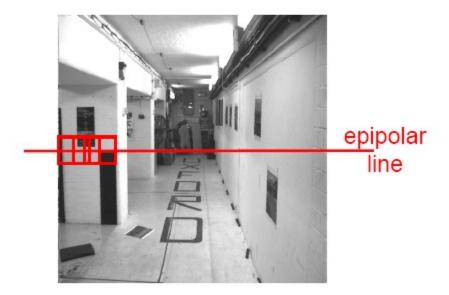
Intensity profiles



Clear correspondence between intensities, but also noise and ambiguity

Correspondence problem





Neighborhood of corresponding points are similar in intensity patterns.

Normalized cross correlation

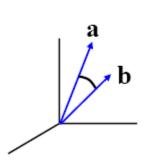
subtract mean: $A \leftarrow A - < A >, B \leftarrow B - < B >$

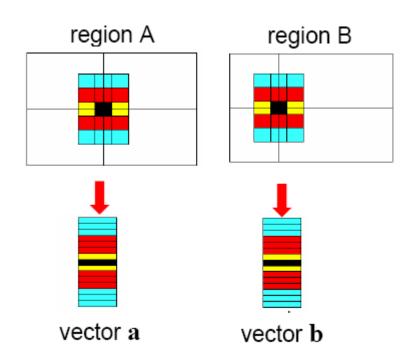
$$NCC = \frac{\sum_{i} \sum_{j} A(i,j) B(i,j)}{\sqrt{\sum_{i} \sum_{j} A(i,j)^{2}} \sqrt{\sum_{i} \sum_{j} B(i,j)^{2}}}$$

Write regions as vectors

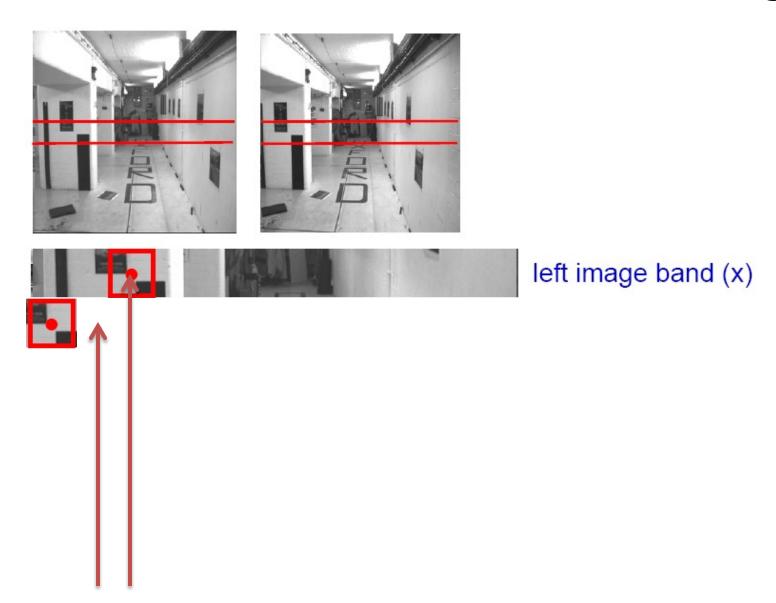
$$\mathtt{A} \to \mathtt{a}, \ \mathtt{B} \to \mathtt{b}$$

$$NCC = \frac{a.b}{|a||b|}$$
$$-1 \le NCC \le 1$$

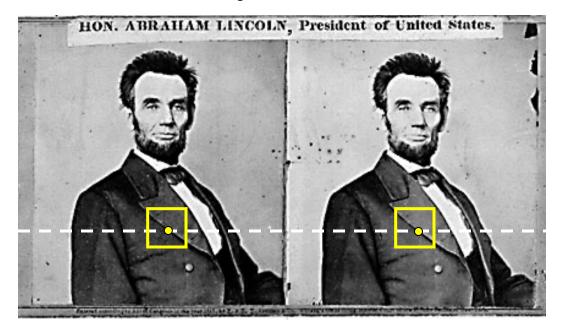




Correlation-based window matching



Dense correspondence search



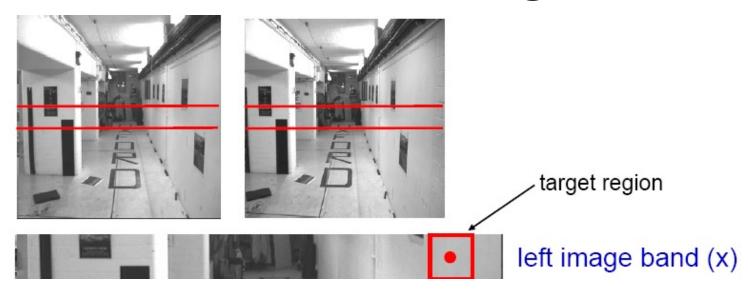
For each epipolar line

For each pixel / window in the left image

- compare with every pixel / window on same epipolar line in right image
- pick position with minimum match cost (e.g., SSD, correlation)

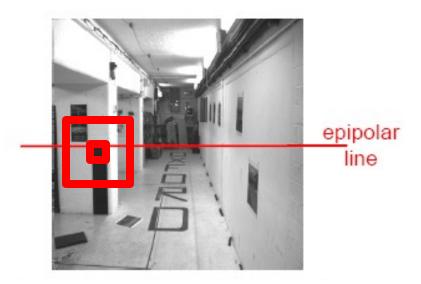
Adapted from Li Zhang Grauman

Textureless regions



Effect of window size





Source: Andrew Zisserman Grauman

Effect of window size







W = 3

W = 20

Want window large enough to have sufficient intensity variation, yet small enough to contain only pixels with about the same disparity.

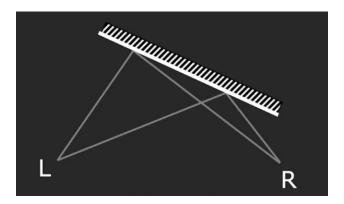
Issues with Stereo

Surface must have non-repetitive texture



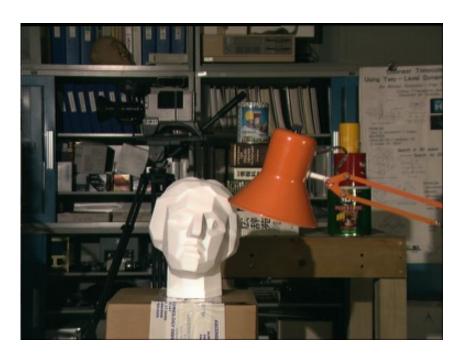


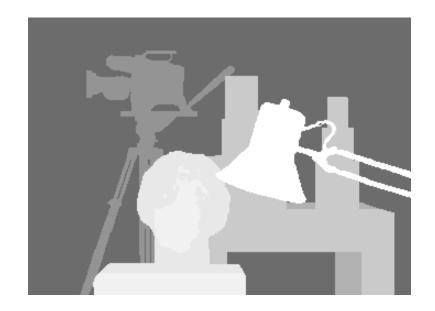
Foreshortening effect makes matching a challenge



Stereo Results

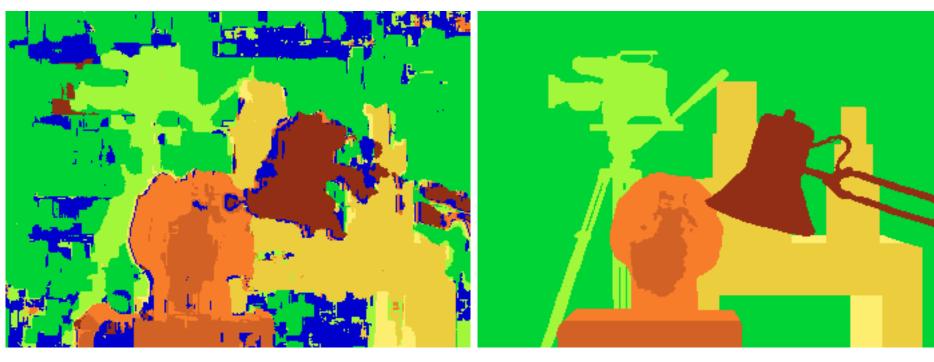
Data from University of Tsukuba





Scene Ground truth

Results with Window Search



Window-based matching (best window size)

Ground truth

Better methods exist...



Energy Minimization

Boykov et al., <u>Fast Approximate Energy Minimization via Graph Cuts</u>, International Conference on Computer Vision, September 1999.

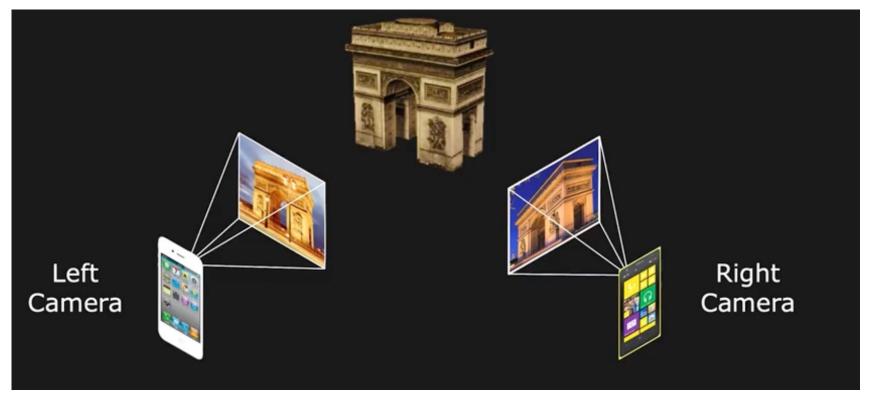
Ground truth

Summary

- With a simple stereo system, how much pixels move, or "disparity" give information about the depth
- Correspondences to measure the pixel disparity

Next: Uncalibrated Stereo

From two arbitrary views



Assume intrinsics are known (fx, fy, ox, oy)