

# Two-View Stereo

---



# What do you see in this image?

---



Autostereograms: [www.magiceye.com](http://www.magiceye.com)

# Stereo

---

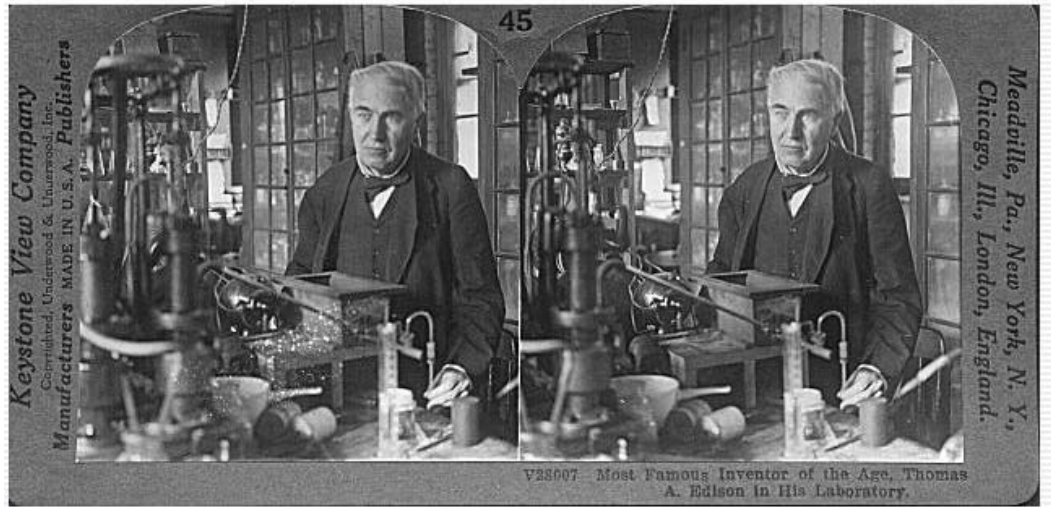
- What cues tell us about scene depth?



# Stereograms

---

- Humans can fuse pairs of images to get a sensation of depth



Stereograms: Invented by Sir Charles Wheatstone, 1838

# Stereograms

---



# Stereograms

---

- Humans can fuse pairs of images to get a sensation of depth



Autostereograms: [www.magiceye.com](http://www.magiceye.com)

# Stereograms

---

- Humans can fuse pairs of images to get a sensation of depth



Autostereograms: [www.magiceye.com](http://www.magiceye.com)

# Problem formulation

---

- Given a calibrated binocular stereo pair, fuse it to produce a depth image

image 1



image 2



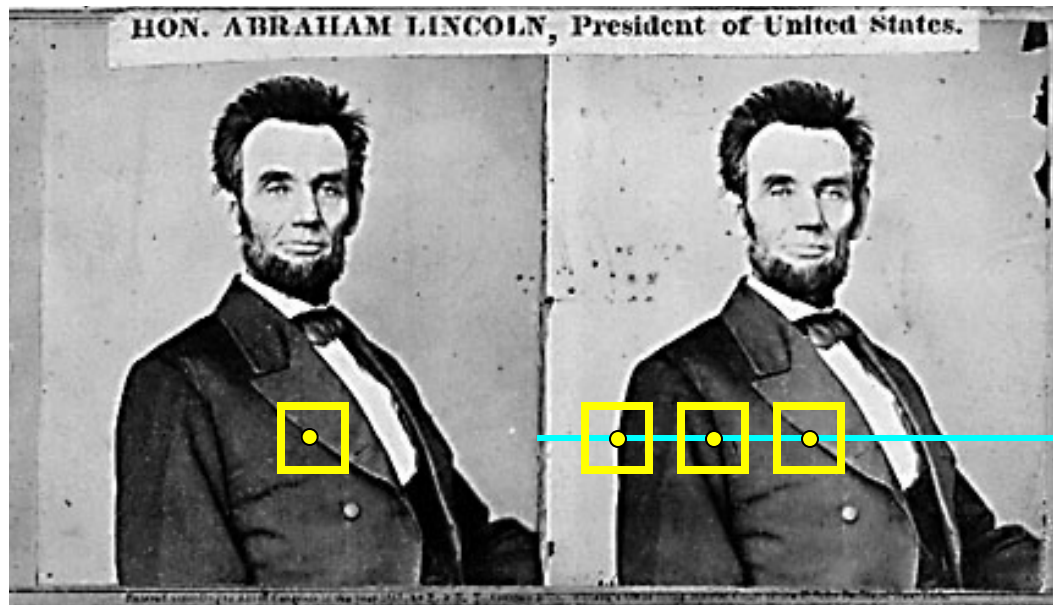
Dense depth map





# Basic stereo matching algorithm

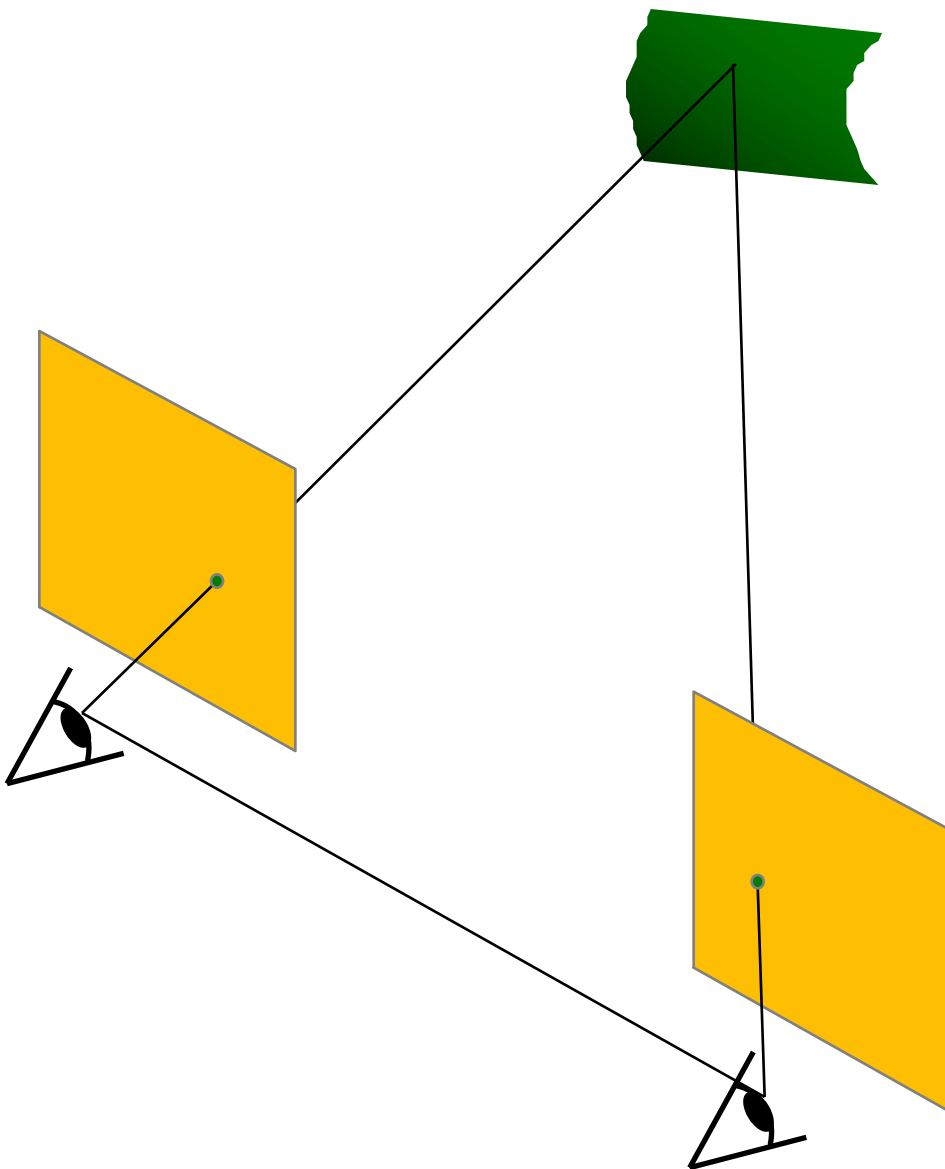
---



- For each pixel in the first image
  - Find corresponding epipolar line in the right image
  - Examine all pixels on the epipolar line and pick the best match
  - Triangulate the matches to get depth information
- Simplest case: epipolar lines are corresponding scanlines
  - When does this happen?

# Simplest Case: Parallel images

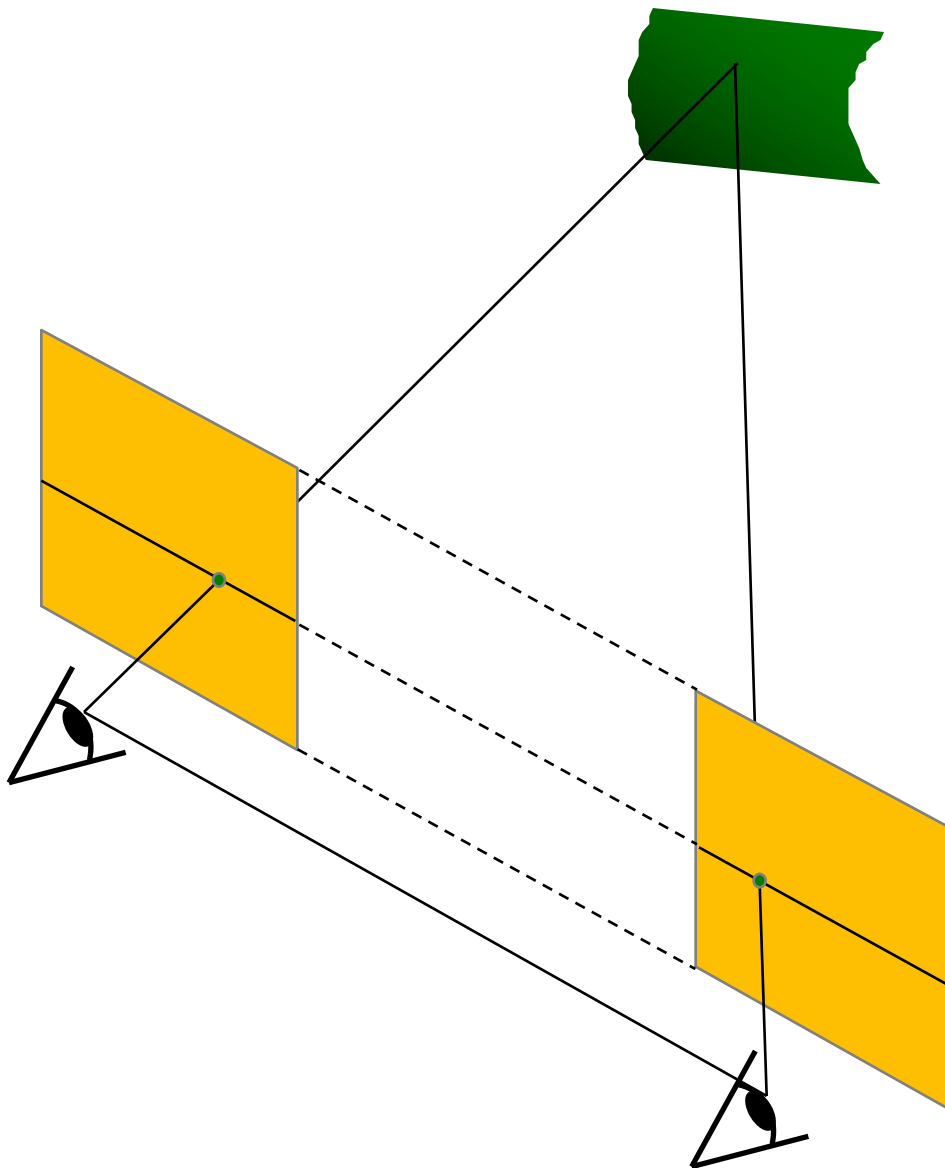
---



- Image planes of cameras are parallel to each other and to the baseline
- Camera centers are at same height
- Focal lengths are the same

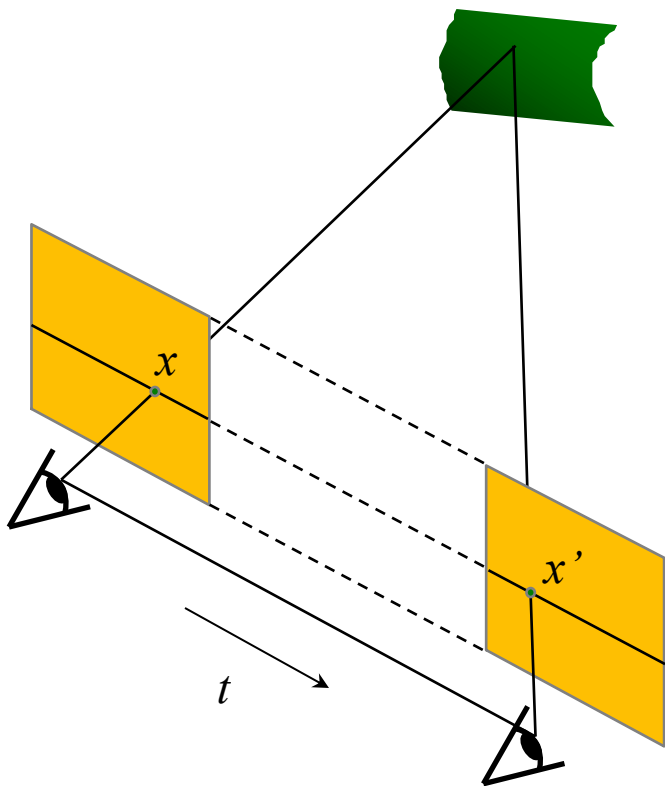
# Simplest Case: Parallel images

---



- Image planes of cameras are parallel to each other and to the baseline
- Camera centers are at same height
- Focal lengths are the same
- Then epipolar lines fall along the horizontal scan lines of the images

# Essential matrix for parallel images



Epipolar constraint:

$$\mathbf{x}'^T \mathbf{E} \mathbf{x} = 0, \quad \mathbf{E} = [\mathbf{t}_\times] \mathbf{R}$$

$$\mathbf{R} = \mathbf{I} \quad \mathbf{t} = (T, 0, 0)$$

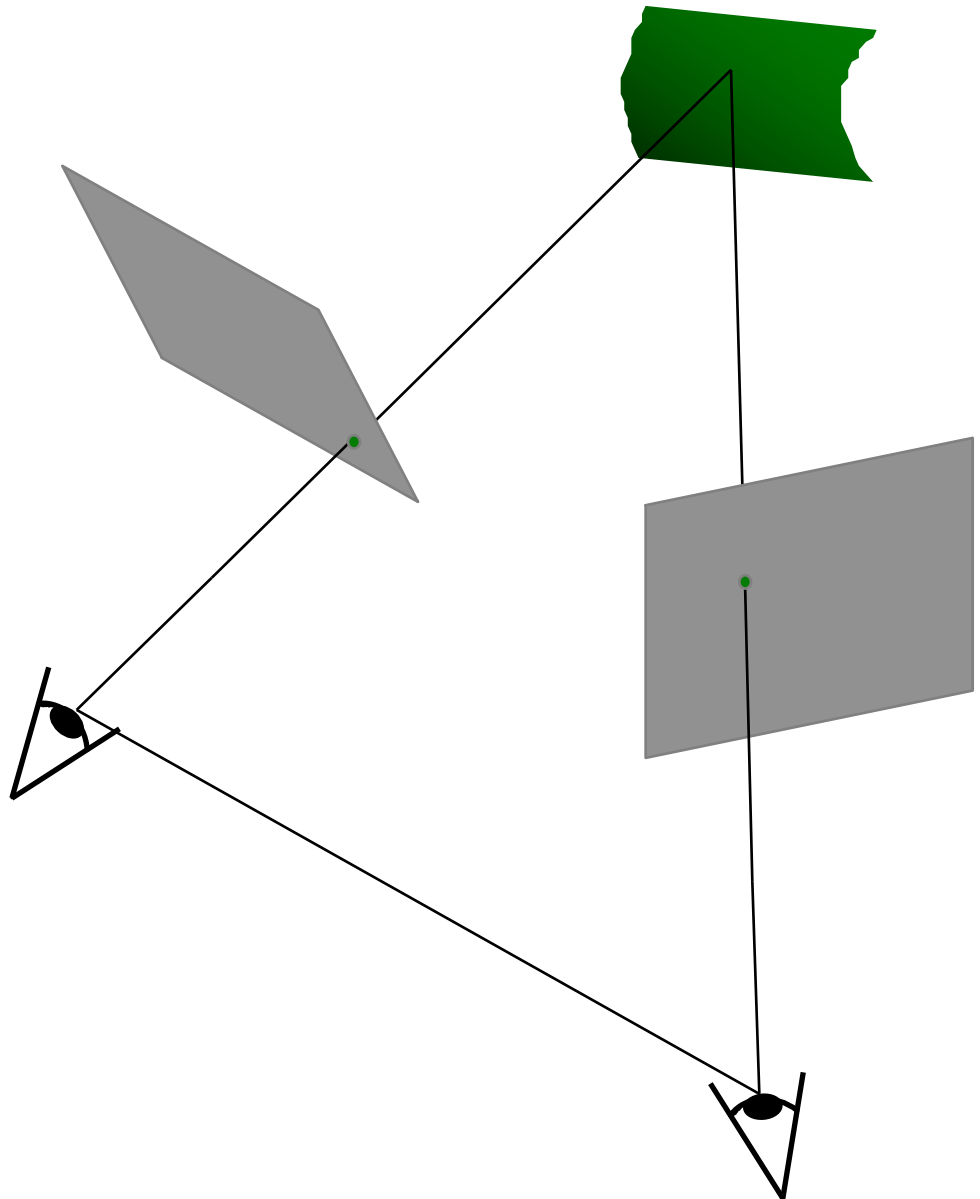
$$\mathbf{E} = [\mathbf{t}_\times] \mathbf{R} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -T \\ 0 & T & 0 \end{bmatrix}$$

$$(u' \quad v' \quad 1) \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & -T \\ 0 & T & 0 \end{bmatrix} \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = 0 \quad (u' \quad v' \quad 1) \begin{pmatrix} 0 \\ -T \\ Tv \end{pmatrix} = 0 \quad Tv' = Tv$$

The y-coordinates of corresponding points are the same!

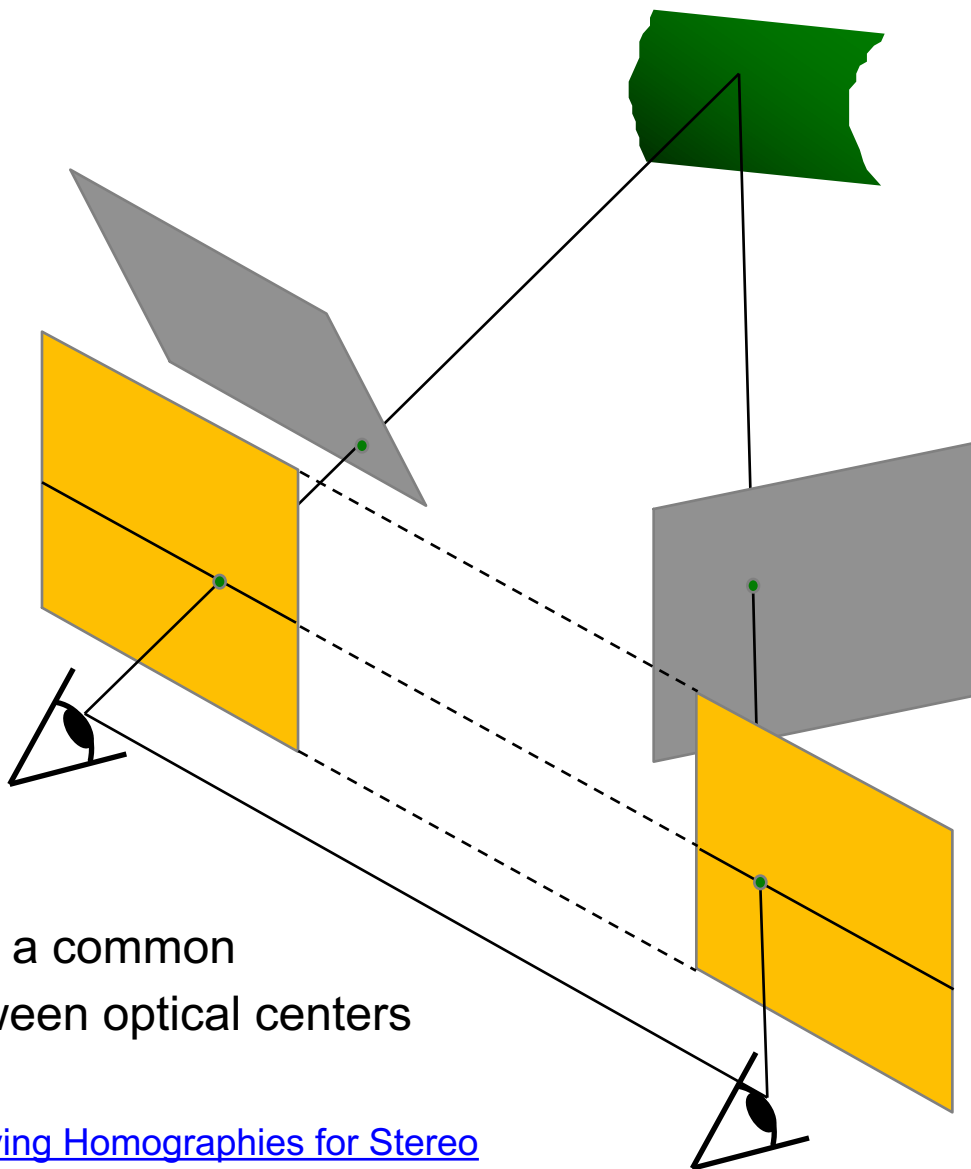
# Stereo image rectification

---



# Stereo image rectification

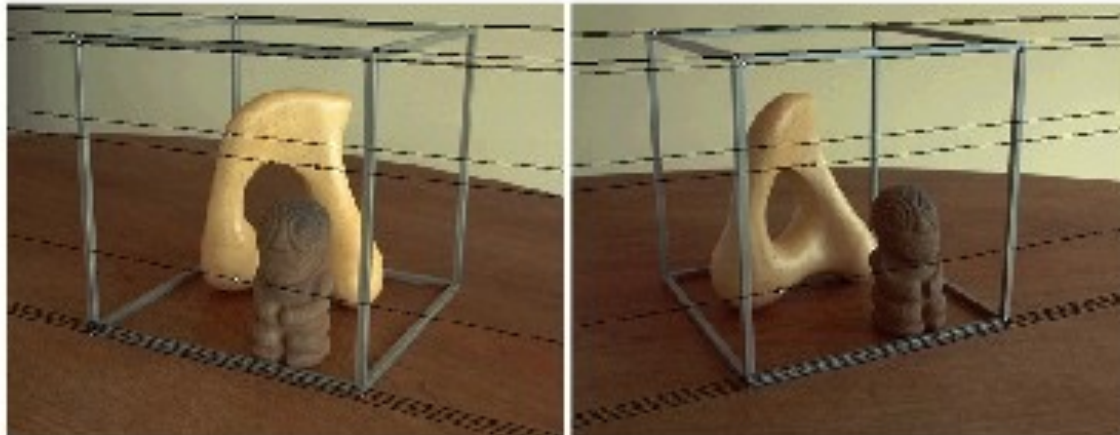
---



- Reproject image planes onto a common plane parallel to the line between optical centers

# Rectification example

---



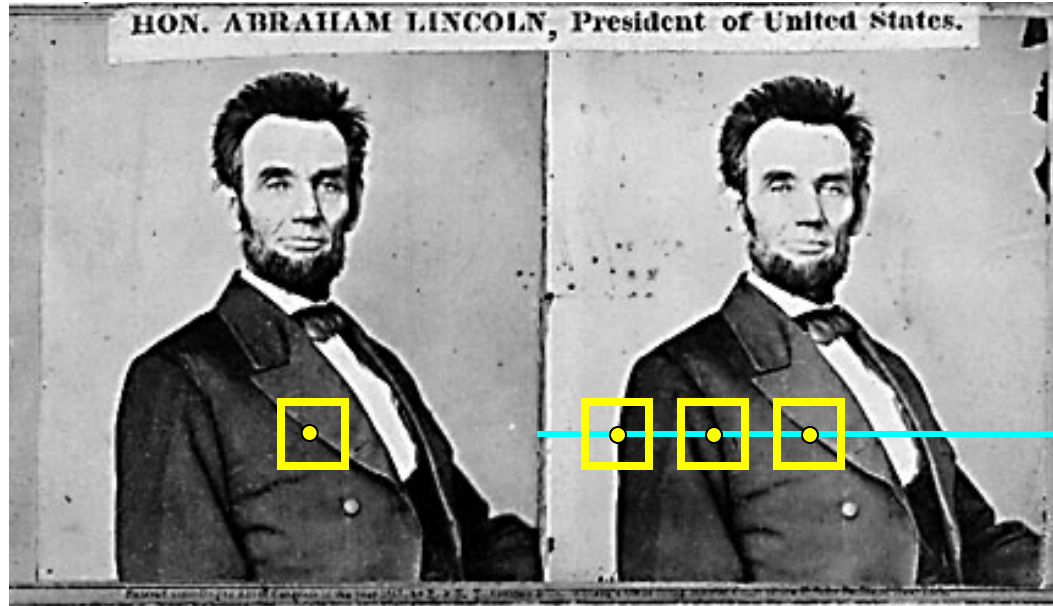
# Another rectification example





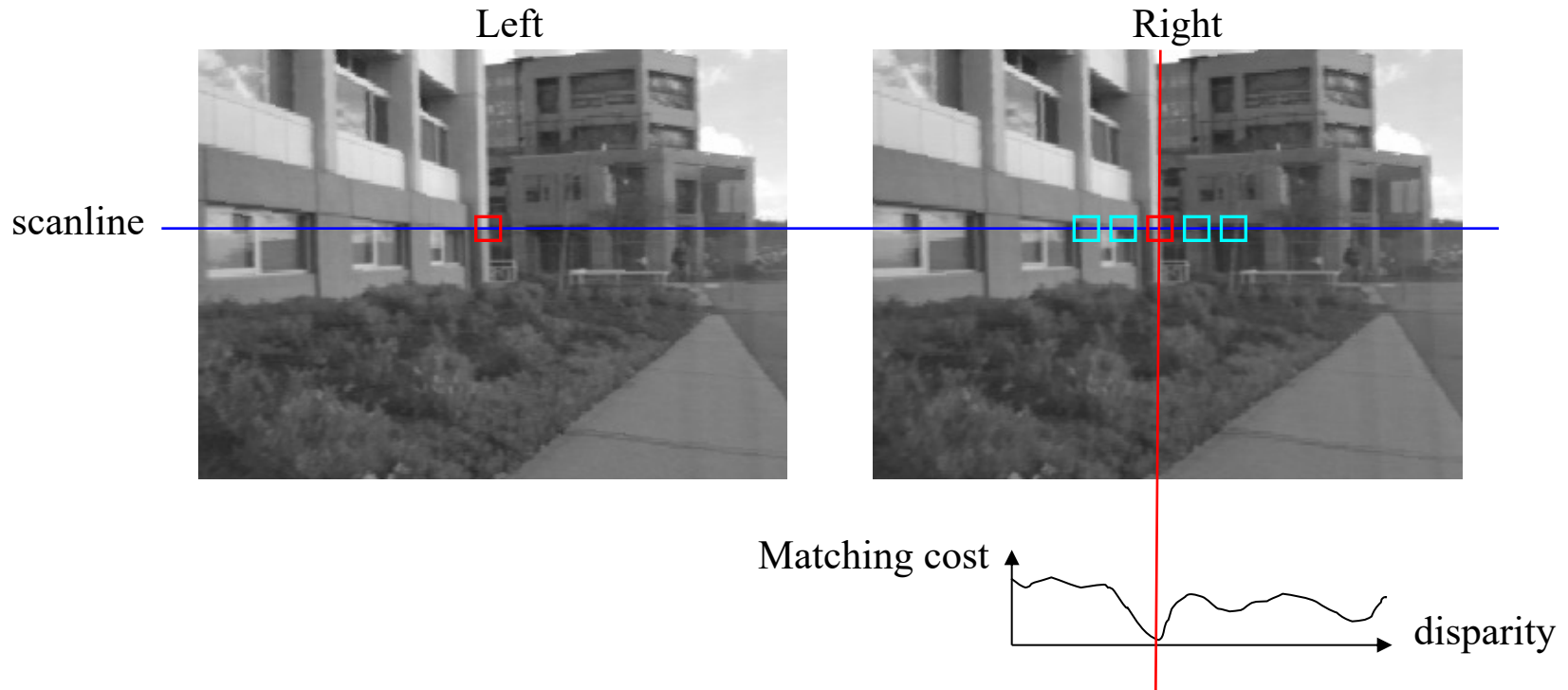
# Basic stereo matching algorithm

---



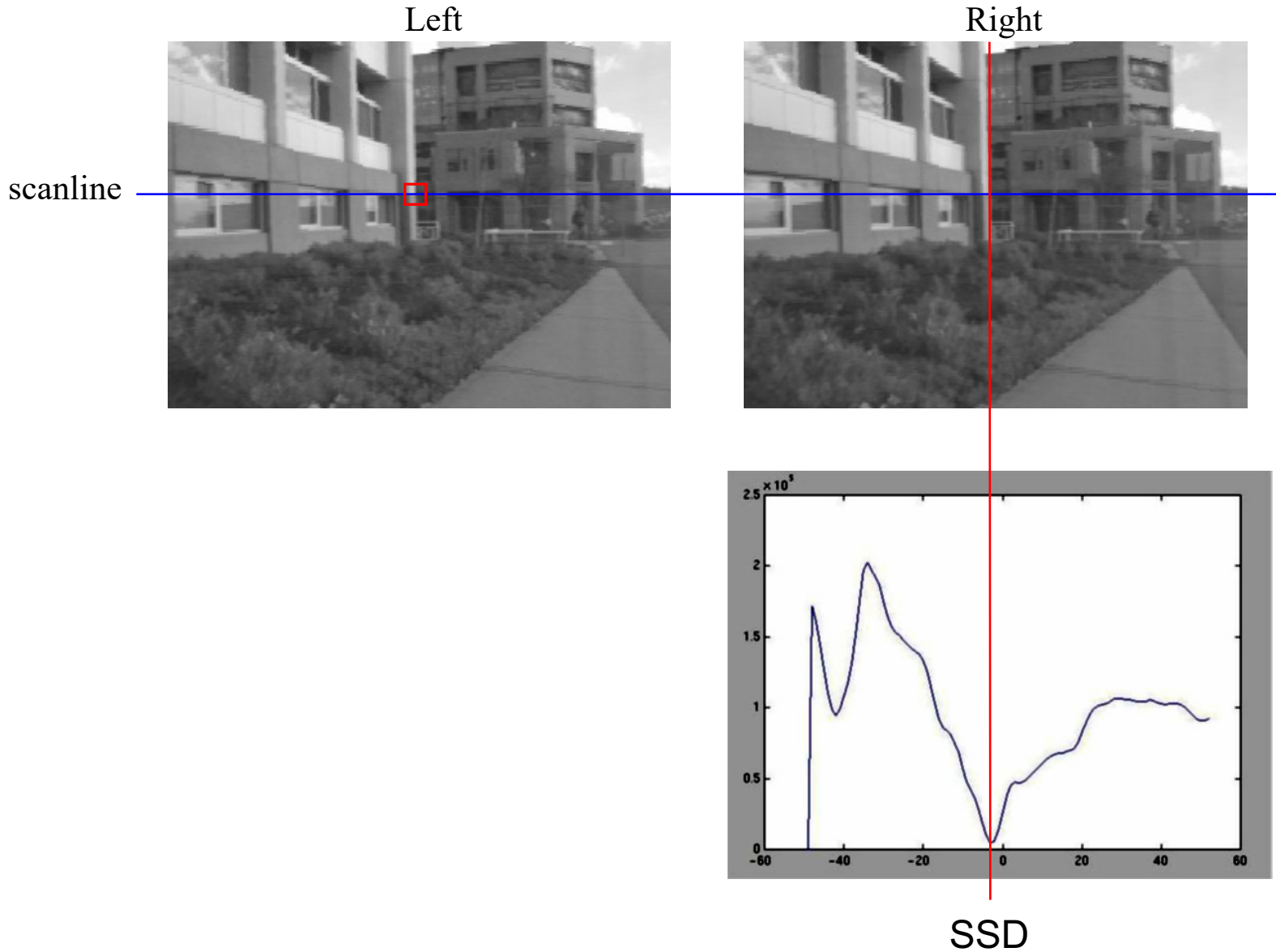
- If necessary, rectify the two stereo images to transform epipolar lines into scanlines
- For each pixel in the first image
  - Find corresponding epipolar line in the right image
  - Examine all pixels on the epipolar line and pick the best match

# Correspondence search

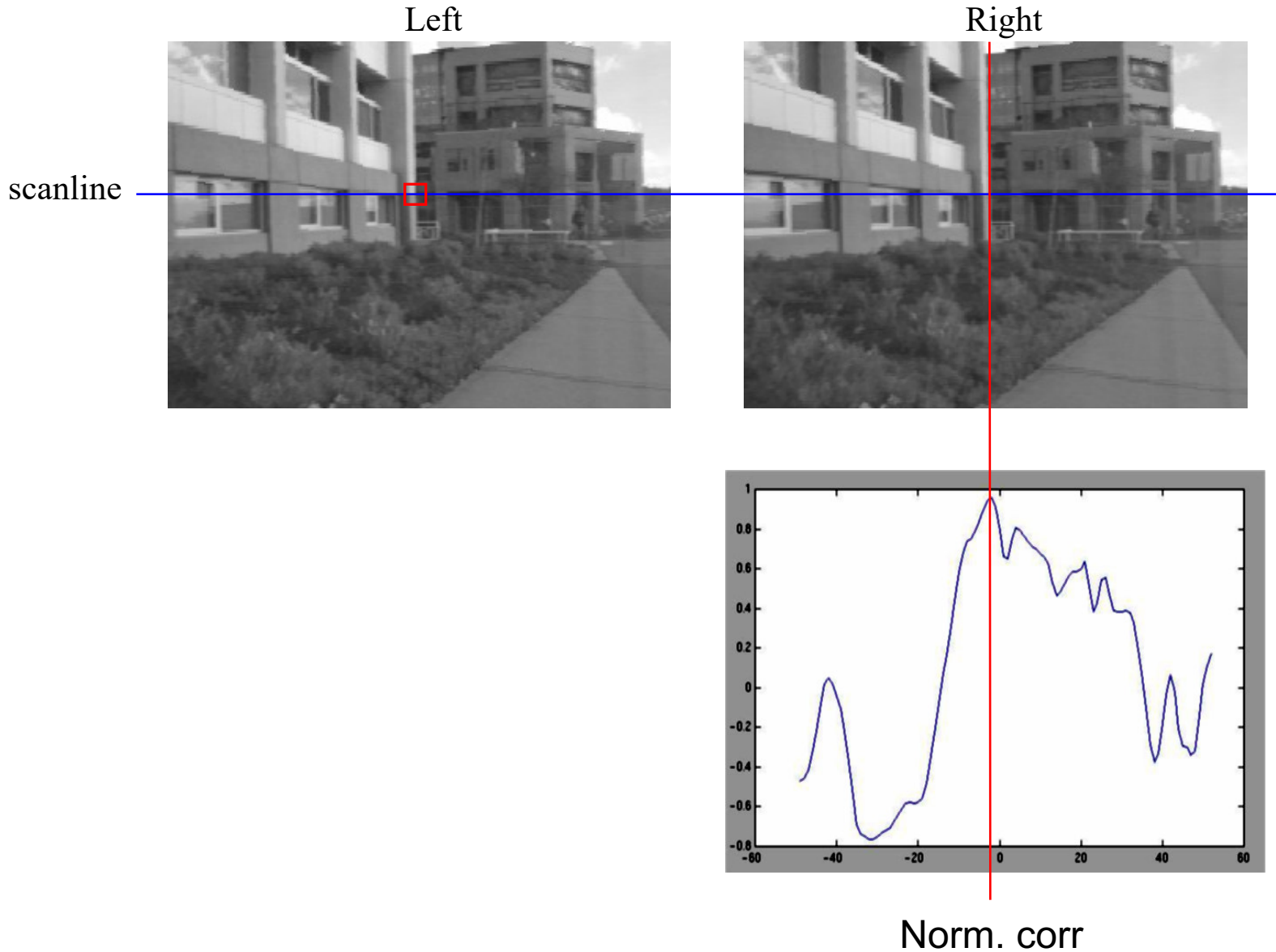


- Slide a window along the right scanline and compare contents of that window with the reference window in the left image
- Matching cost: SSD or normalized correlation

# Correspondence search

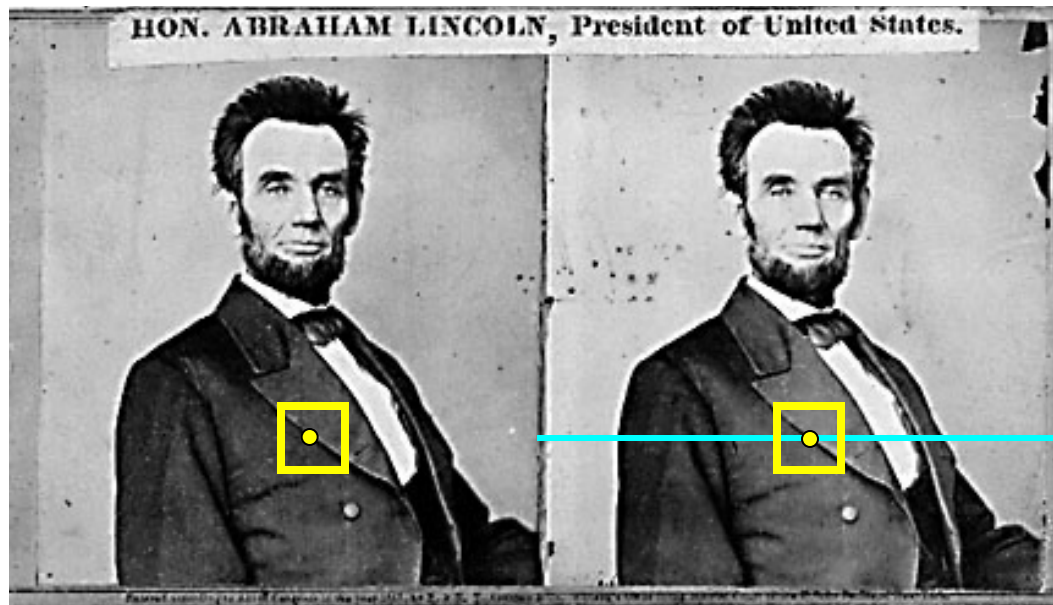


# Correspondence search



# Basic stereo matching algorithm

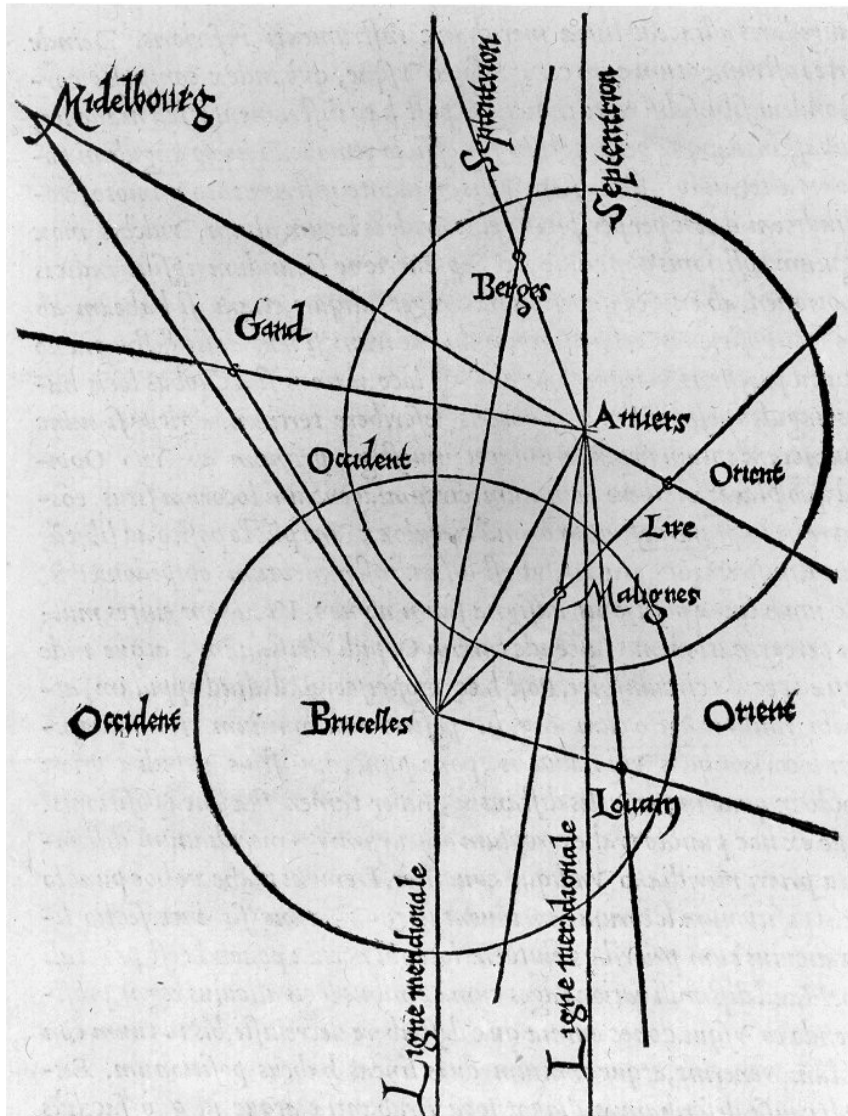
---



- If necessary, rectify the two stereo images to transform epipolar lines into scanlines
- For each pixel  $x$  in the first image
  - Find corresponding epipolar scanline in the right image
  - Examine all pixels on the scanline and pick the best match  $x'$
  - Triangulate the matches to get depth information

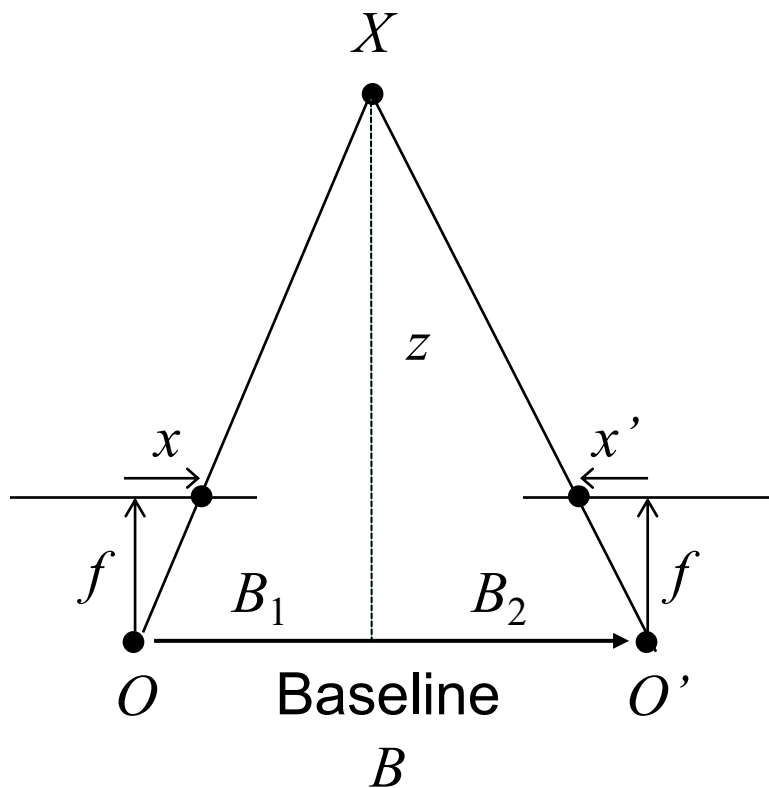
# Triangulation: History

---



From [Wikipedia](#): Gemma Frisius's 1533 diagram introducing the idea of triangulation into the science of surveying. Having established a baseline, e.g. the cities of Brussels and Antwerp, the location of other cities, e.g. Middelburg, Ghent etc., can be found by taking a compass direction from each end of the baseline, and plotting where the two directions cross. This was only a theoretical presentation of the concept — due to topographical restrictions, it is impossible to see Middelburg from either Brussels or Antwerp. Nevertheless, the figure soon became well known all across Europe.

# Depth from disparity



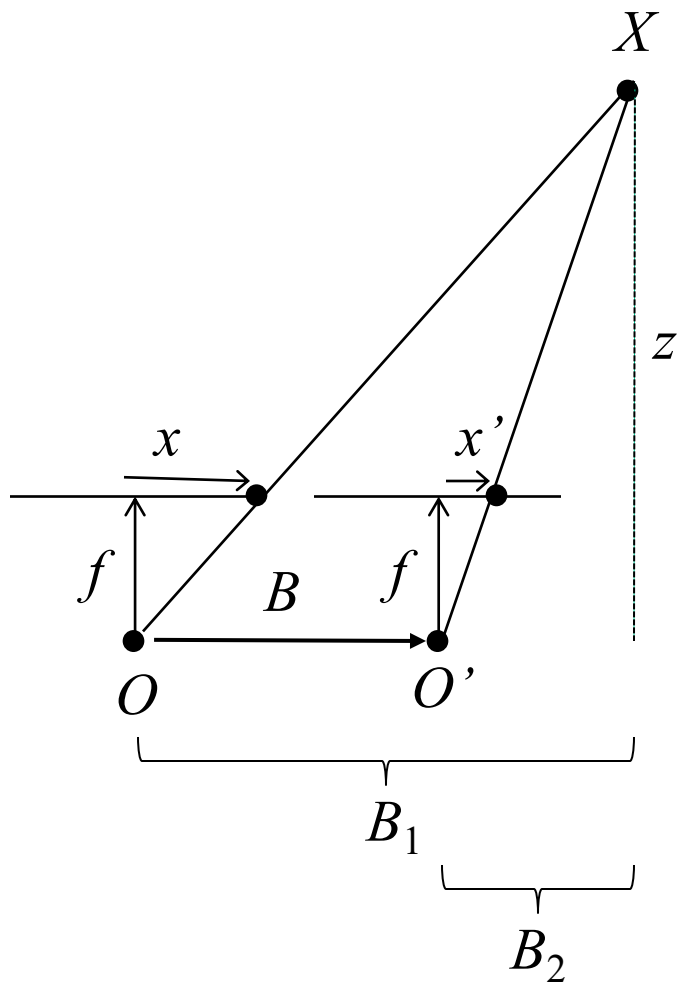
$$\frac{x}{f} = \frac{B_1}{z} \quad \frac{-x'}{f} = \frac{B_2}{z}$$

$$\frac{x - x'}{f} = \frac{B_1 + B_2}{z}$$

$$\text{disparity} = x - x' = \frac{B \cdot f}{z}$$

Disparity is inversely proportional to depth!

# Depth from disparity



$$\frac{x}{f} = \frac{B_1}{z} \quad \frac{x'}{f} = \frac{B_2}{z}$$

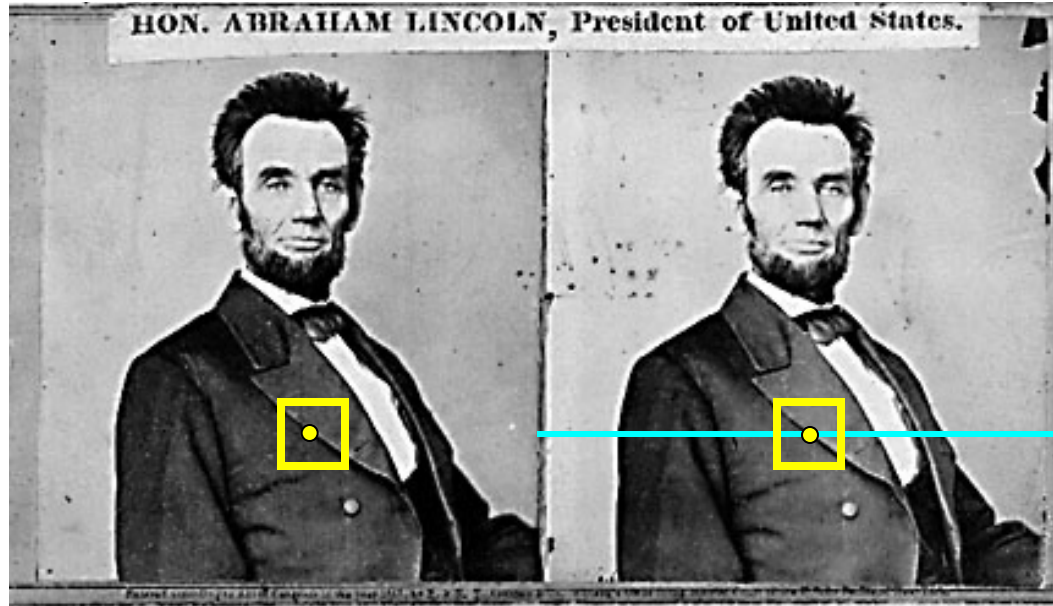
$$\frac{x - x'}{f} = \frac{B_1 - B_2}{z}$$

$$\text{disparity} = x - x' = \frac{B \cdot f}{z}$$



# Basic stereo matching algorithm

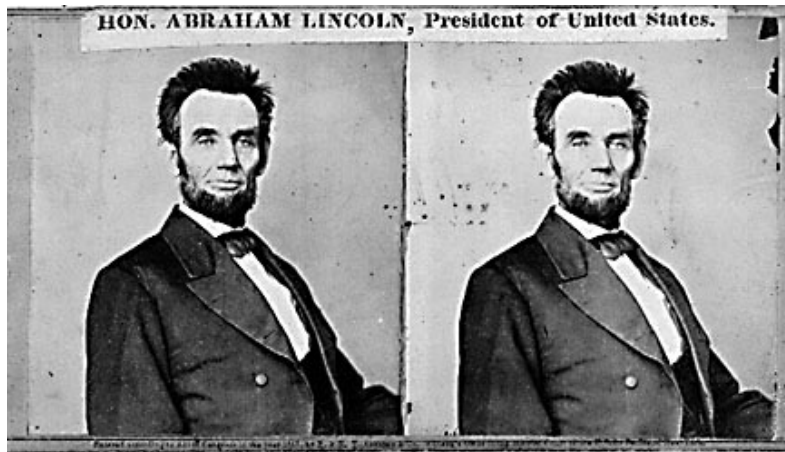
---



- If necessary, rectify the two stereo images to transform epipolar lines into scanlines
- For each pixel  $x$  in the first image
  - Find corresponding epipolar scanline in the right image
  - Examine all pixels on the scanline and pick the best match  $x'$
  - Compute disparity  $x-x'$  and set  $\text{depth}(x) = B \cdot f / (x-x')$

# Failures of correspondence search

---



Textureless surfaces



Occlusions, repetition



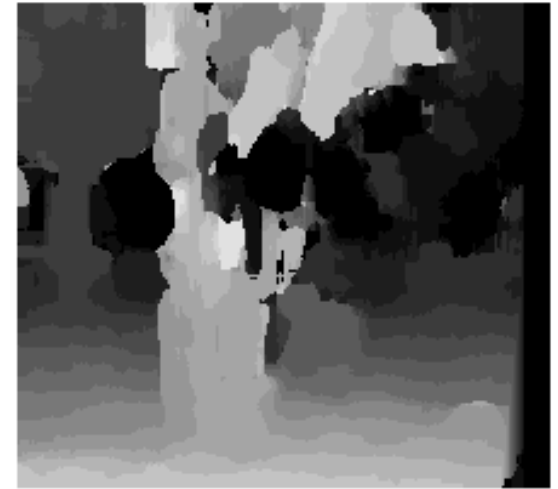
Non-Lambertian surfaces, specularities

# Effect of window size

---



$W = 3$



$W = 20$

- Smaller window
  - + More detail
  - More noise
- Larger window
  - + Smoother disparity maps
  - Less detail

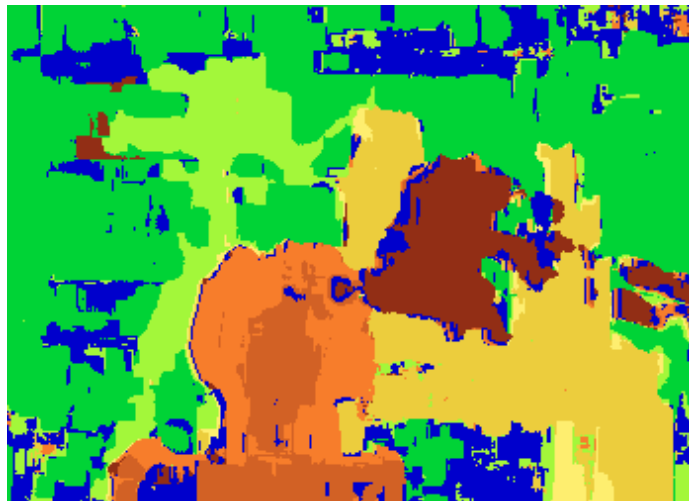
# Results with window search

---

Data



Window-based matching



Ground truth



# Better methods exist...

---



Graph cuts



Ground truth

Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001

For the latest and greatest: <http://www.middlebury.edu/stereo/>

# How can we improve window-based matching?

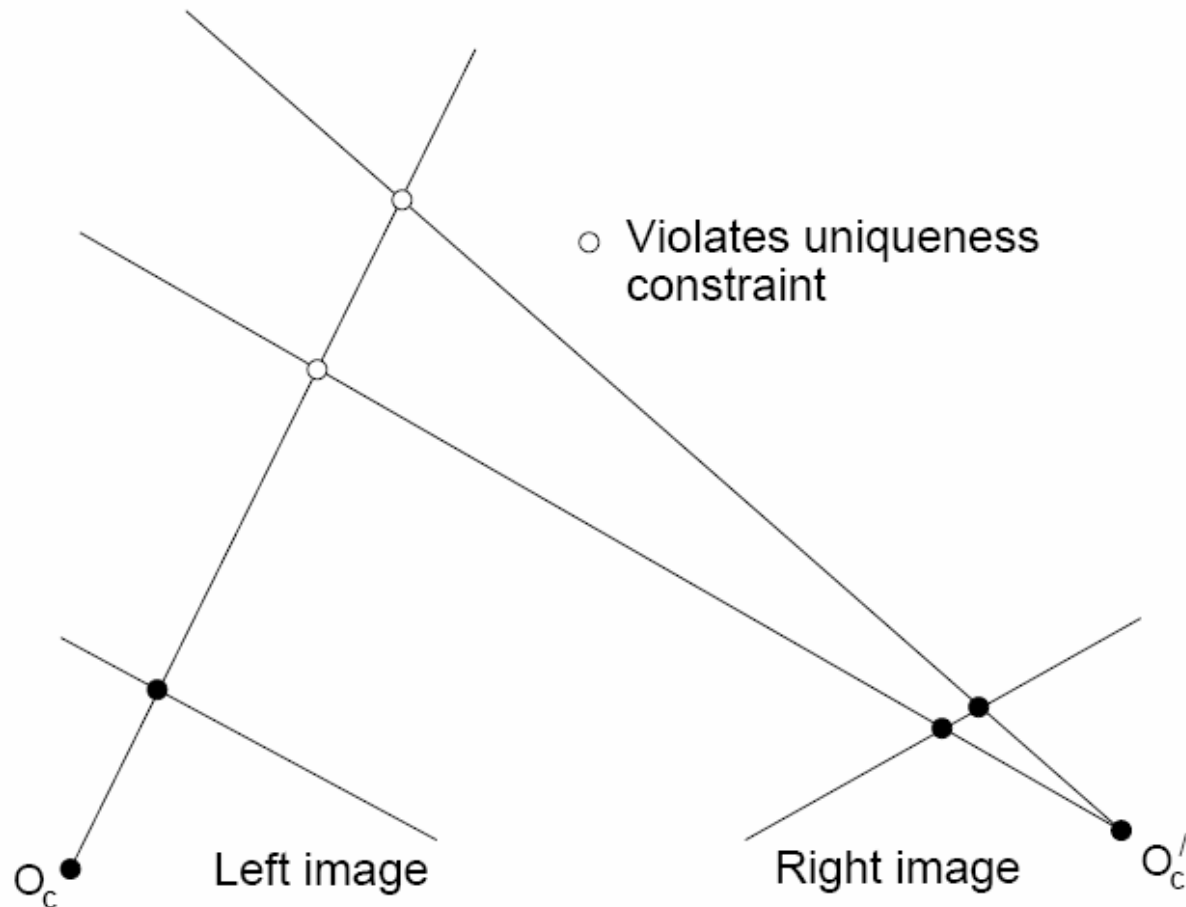
---

- The similarity constraint is **local** (each reference window is matched independently)
- Need to enforce **non-local** correspondence constraints

# Non-local constraints

---

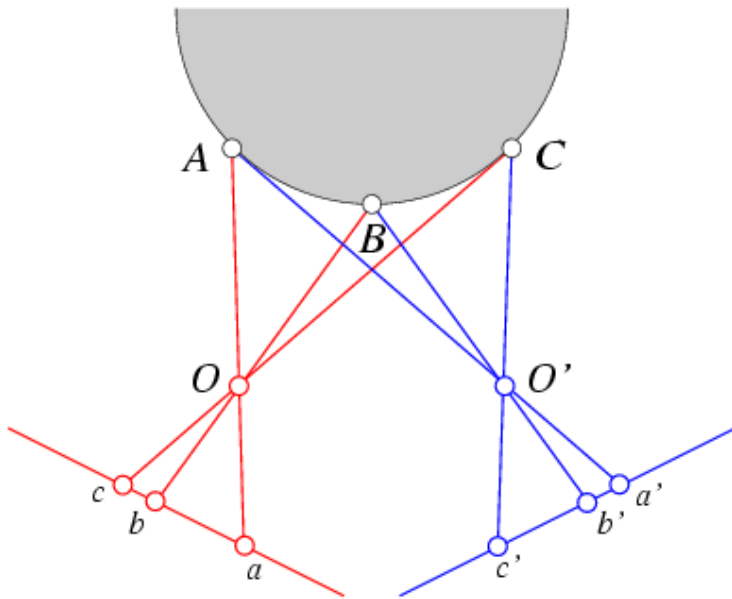
- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image



# Non-local constraints

---

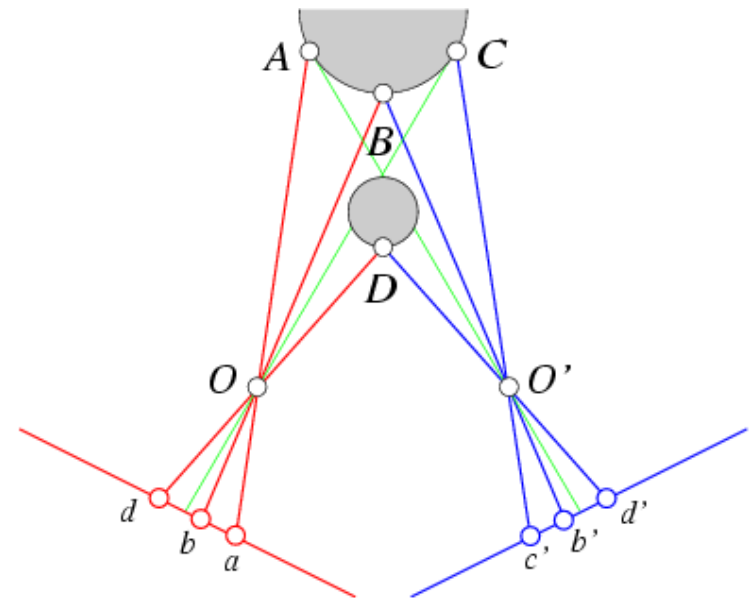
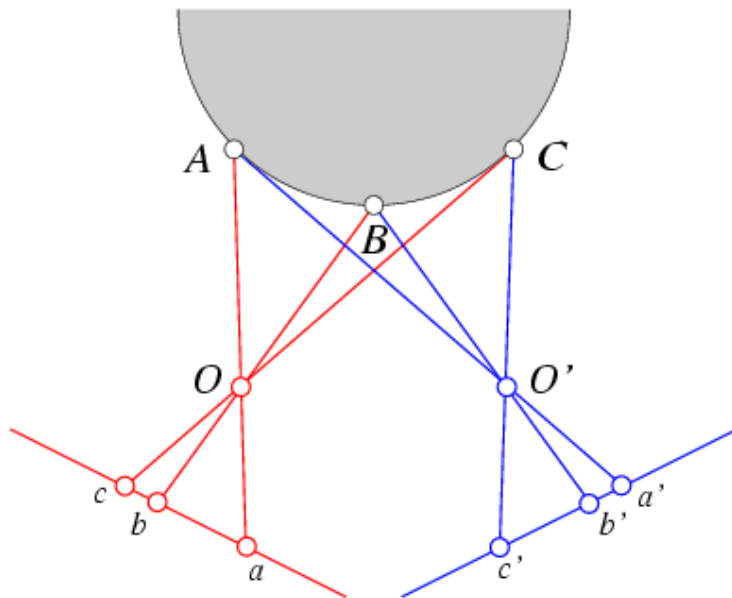
- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views





# Non-local constraints

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views



Ordering constraint doesn't hold

# Non-local constraints

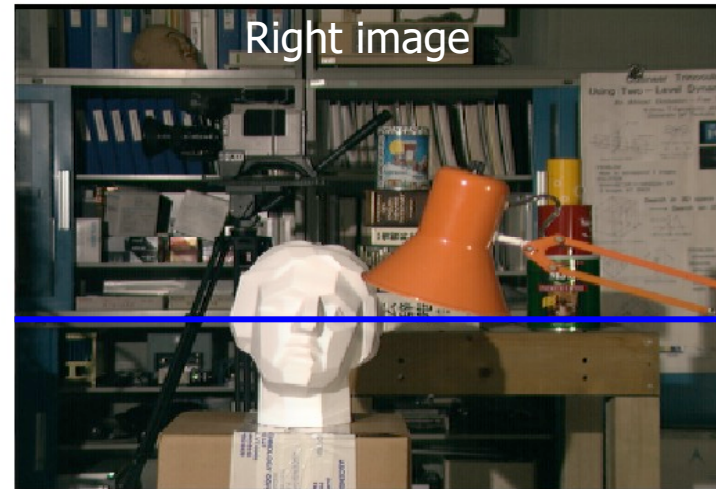
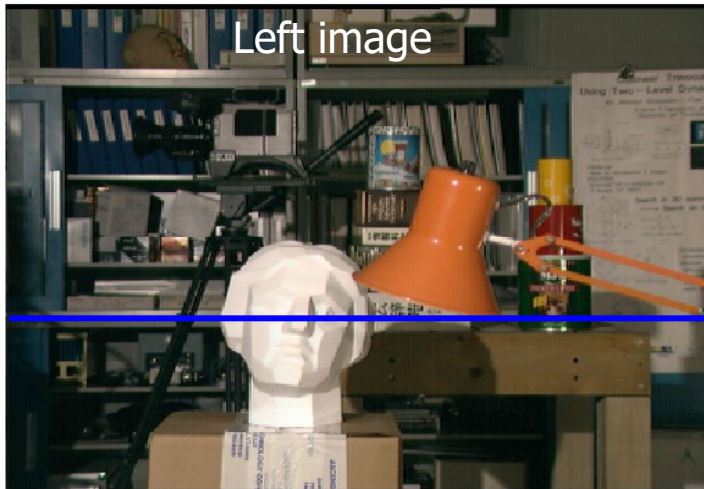
---

- Uniqueness
  - For any point in one image, there should be at most one matching point in the other image
- Ordering
  - Corresponding points should be in the same order in both views
- Smoothness
  - We expect disparity values to change slowly (for the most part)

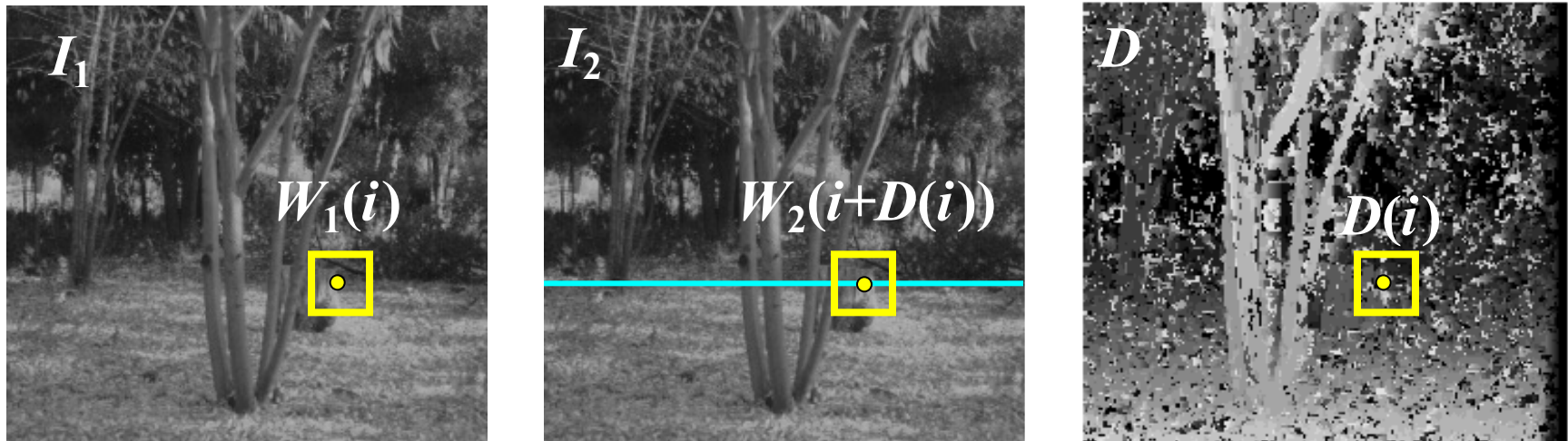
# Scanline stereo

---

- Try to coherently match pixels on the entire scanline
- Different scanlines are still optimized independently



# Stereo matching as global optimization

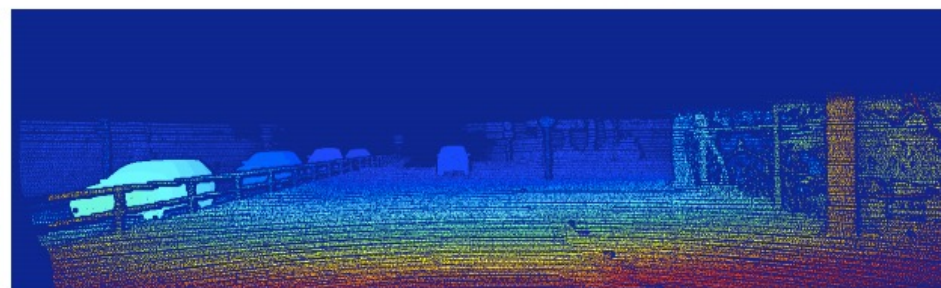
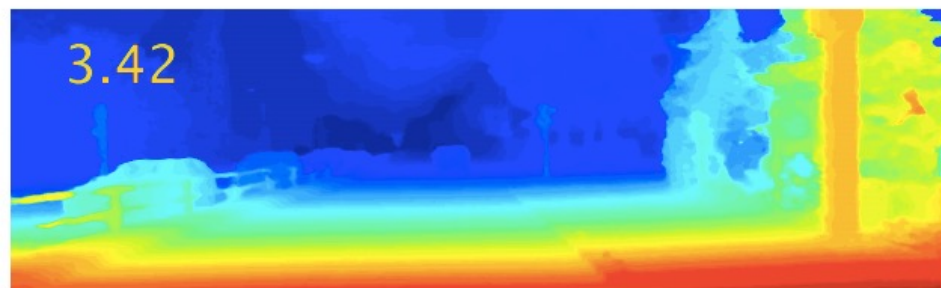
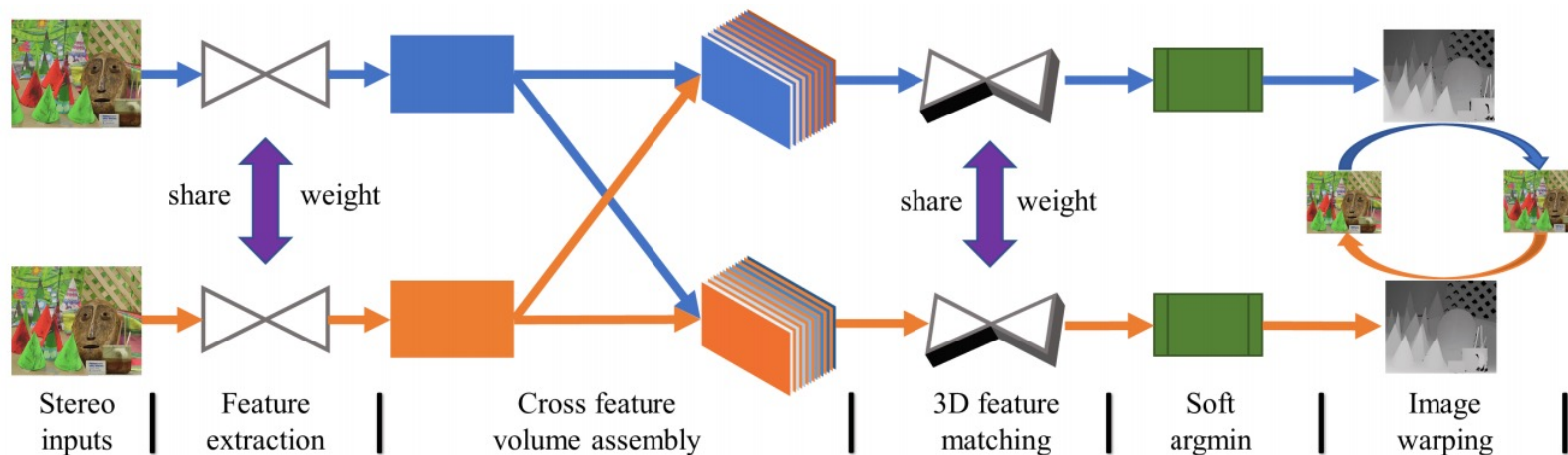


$$E(D) = \underbrace{\sum_i \left( W_1(i) - W_2(i + D(i)) \right)^2}_{\text{data term}} + \lambda \underbrace{\sum_{\text{neighbors } i,j} \rho \left( D(i) - D(j) \right)}_{\text{smoothness term}}$$

- Energy functions of this form can be minimized using *graph cuts*

Y. Boykov, O. Veksler, and R. Zabih, [Fast Approximate Energy Minimization via Graph Cuts](#), PAMI 2001

# Stereo matching as a prediction problem

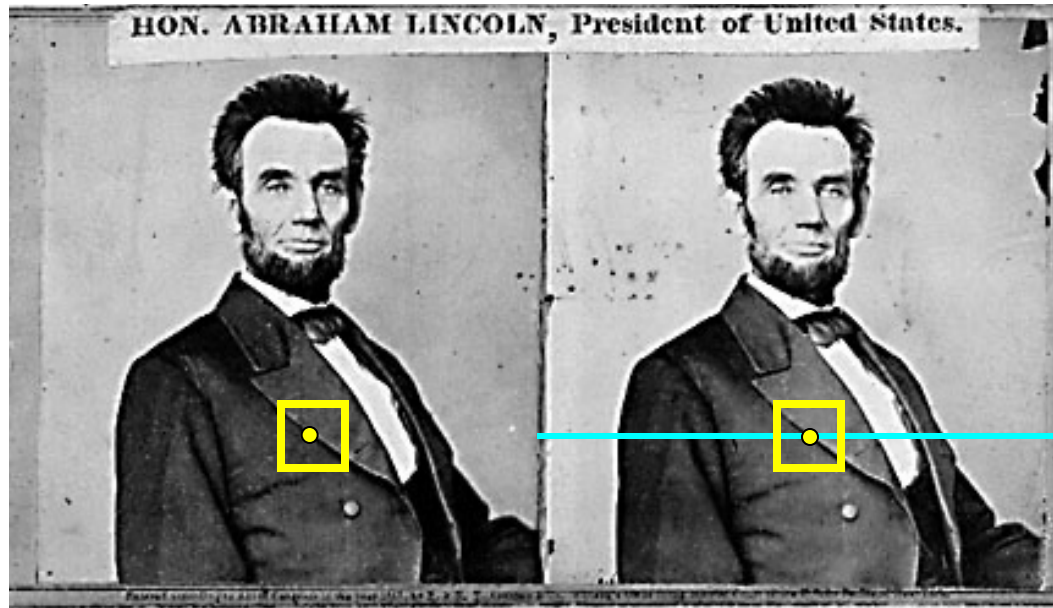


Y. Zhong, Y. Dai, and H. Li, [Self-Supervised Learning for Stereo Matching with Self-Improving Ability](#), arXiv 2017

Slide from L. Lazebnik.

# Review: Basic stereo matching algorithm

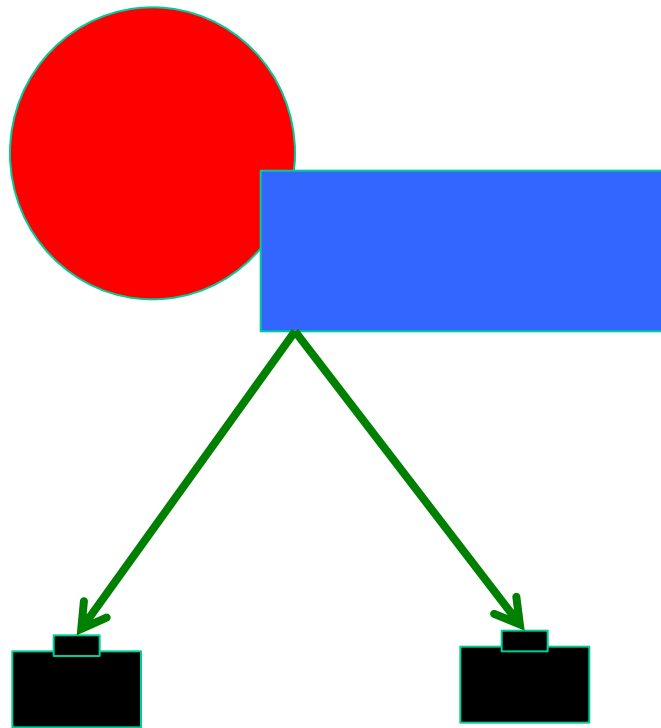
---



- For each pixel  $x$  in the reference image
  - Find corresponding epipolar scanline in the other image
  - Examine all pixels on the scanline and pick the best match  $x'$
  - Compute disparity  $x-x'$  and set  $\text{depth}(x) = B*f/(x-x')$

# Depth from Triangulation

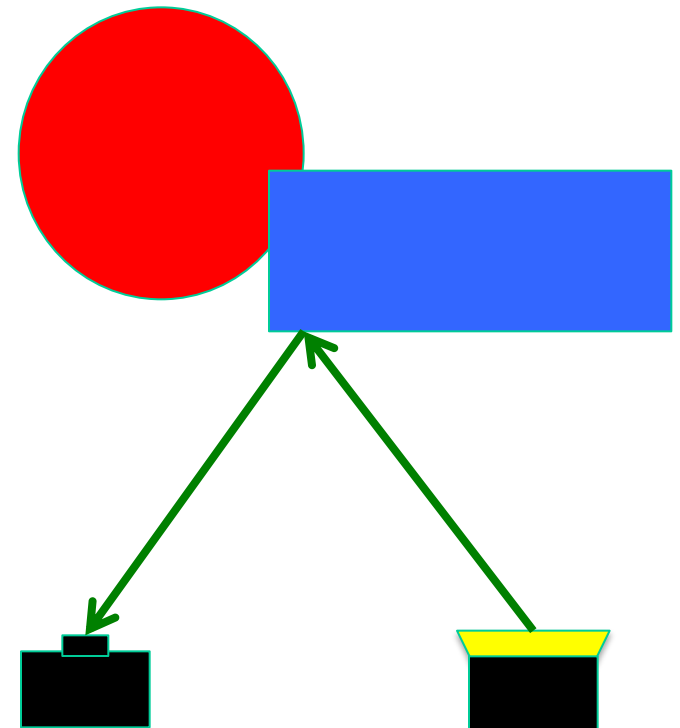
---



Camera 1

Camera 2

Passive Stereopsis



Camera

Projector

Active Stereopsis

Active sensing simplifies the problem of estimating point correspondences

# Kinect: Structured infrared light

---

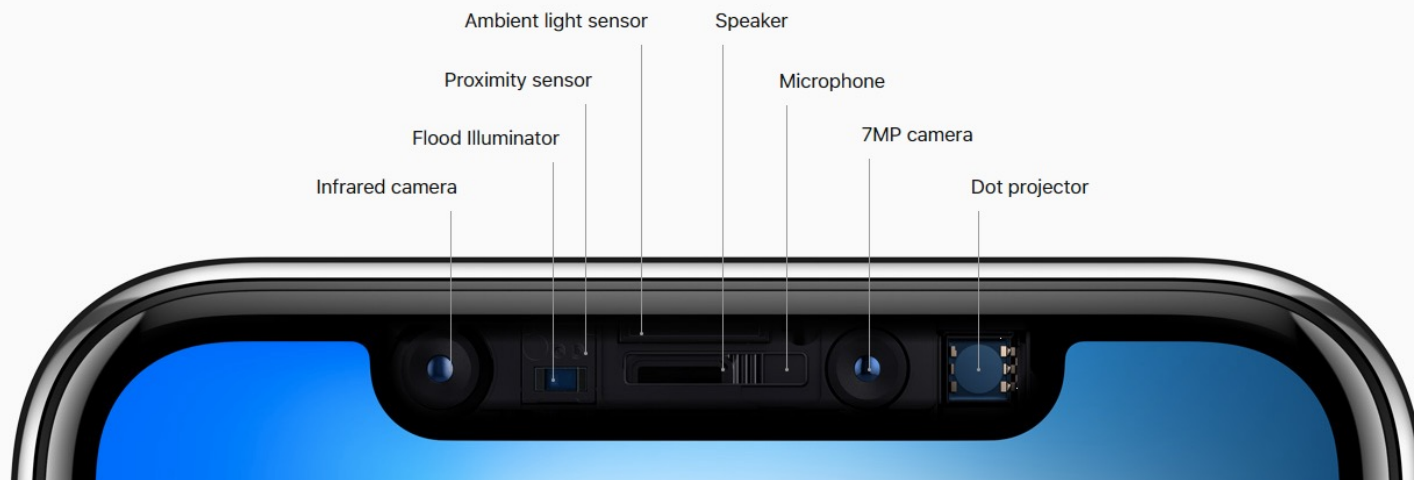


<http://bbzippo.wordpress.com/2010/11/28/kinect-in-infrared/>



# Apple TrueDepth

---

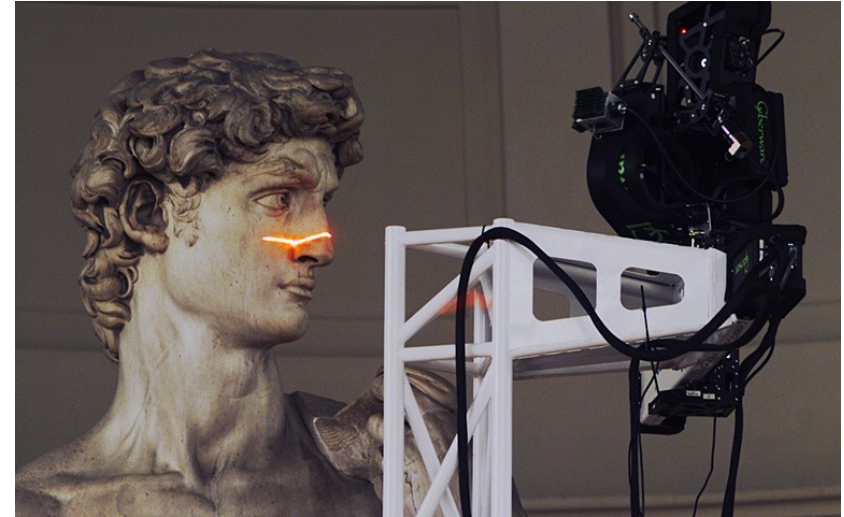
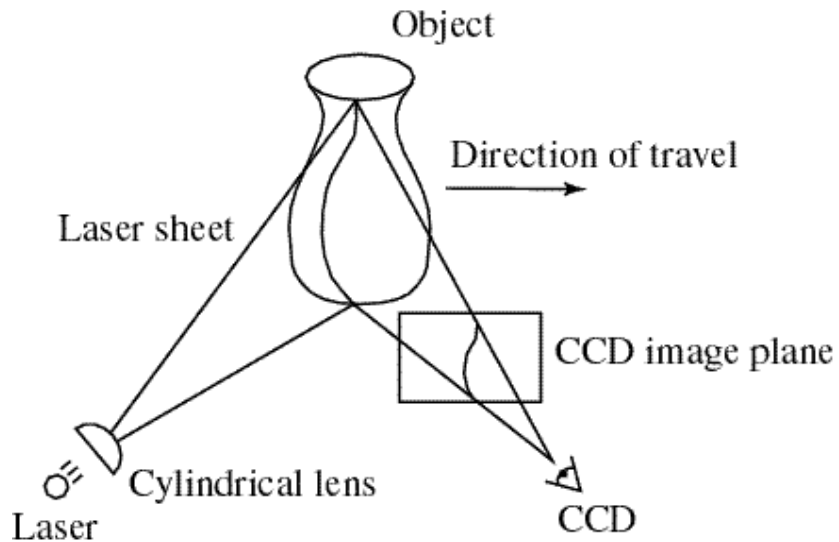


<https://www.cnet.com/news/apple-face-id-truedepth-how-it-works/>



# Laser scanning

---



Digital Michelangelo Project  
Levoy et al.

<http://graphics.stanford.edu/projects/mich/>

## Optical triangulation

- Project a single stripe of laser light
- Scan it across the surface of the object
- This is a very precise version of structured light scanning

# Laser scanned models

---



*The Digital Michelangelo Project*, Levoy et al.

# Laser scanned models

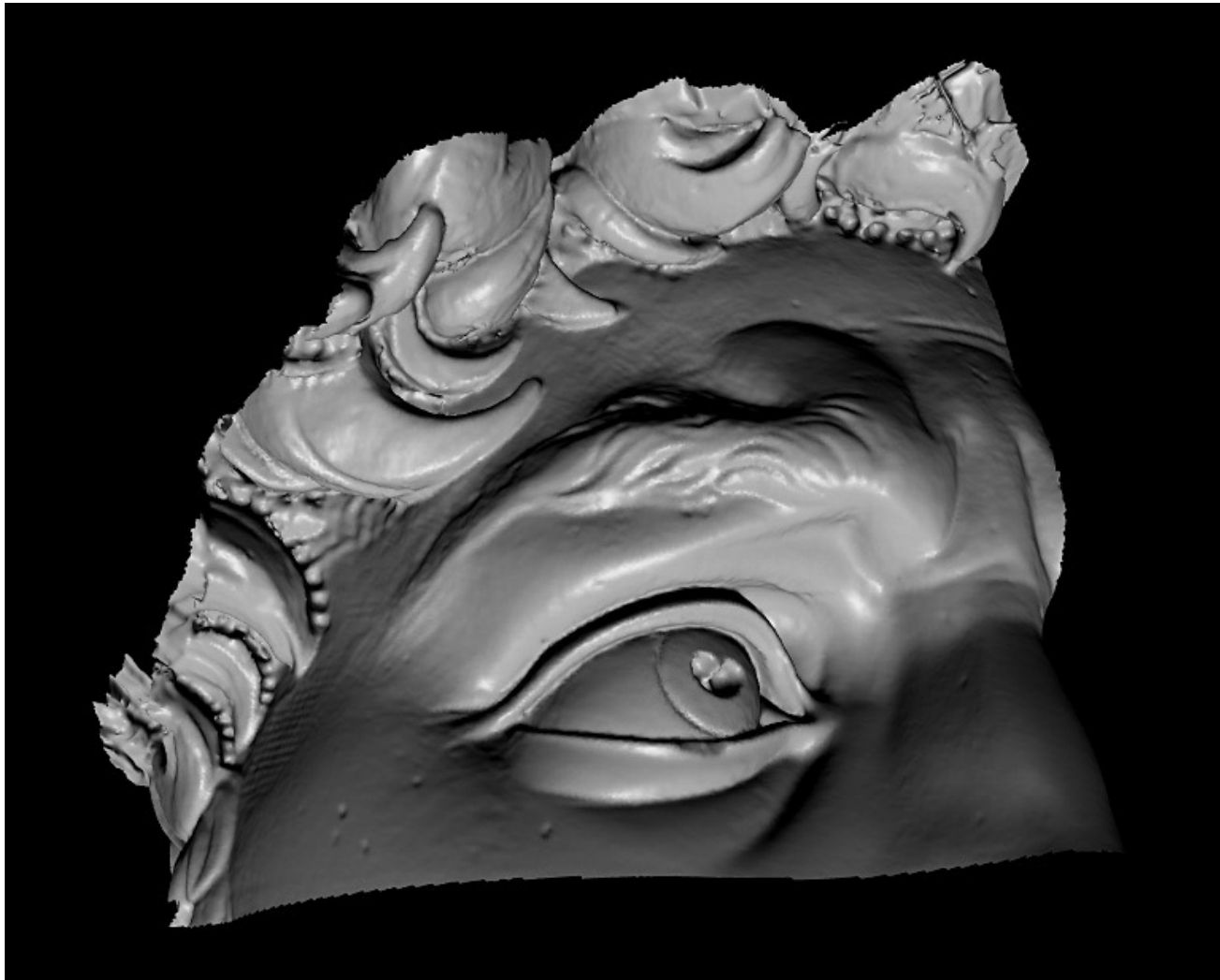
---



*The Digital Michelangelo Project, Levoy et al.*

# Laser scanned models

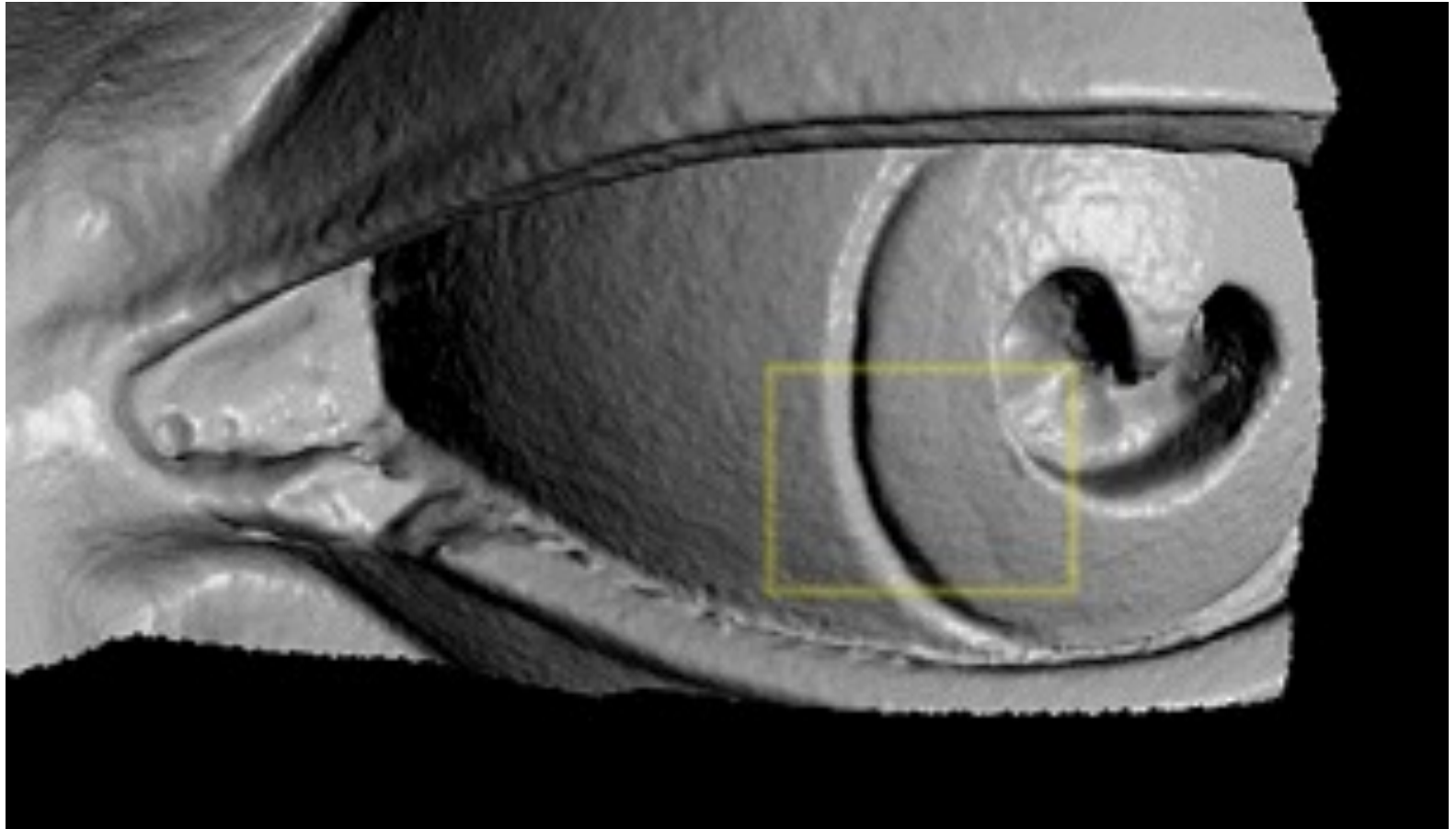
---



*The Digital Michelangelo Project, Levoy et al.*

# Laser scanned models

---

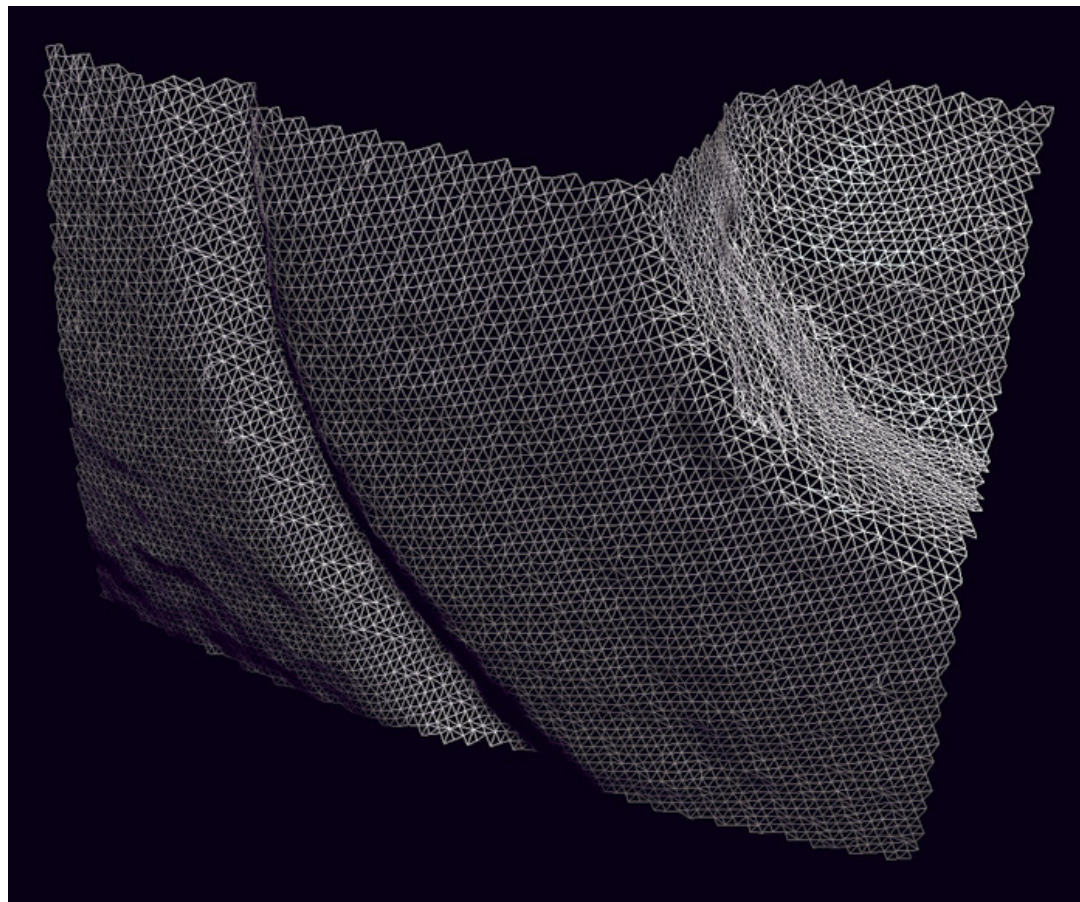


*The Digital Michelangelo Project*, Levoy et al.

# Laser scanned models

---

1.0 mm resolution (56 million triangles)



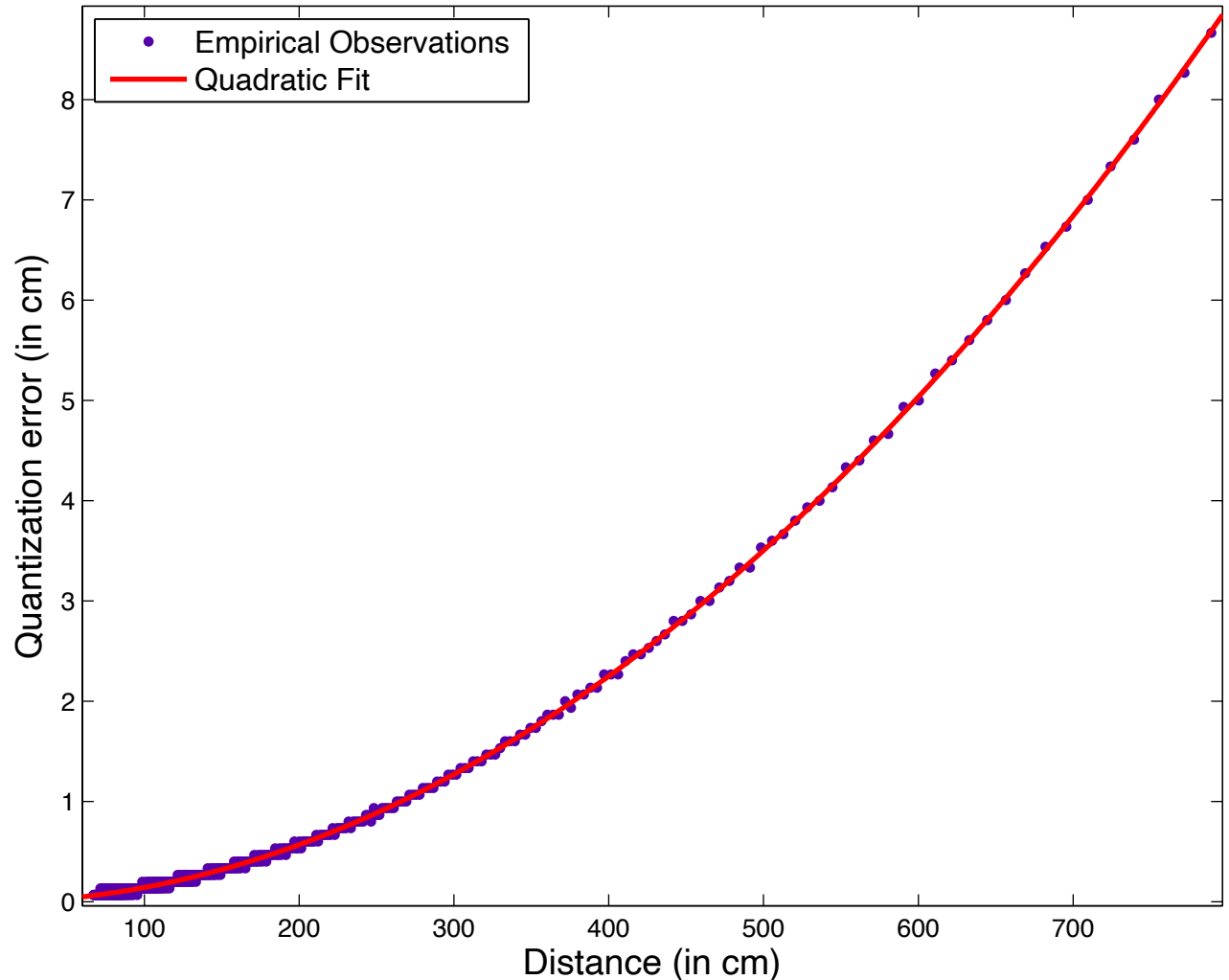
*The Digital Michelangelo Project*, Levoy et al.

# Stereo error(distance)

---

Error in distance estimate increases quadratically with the distance

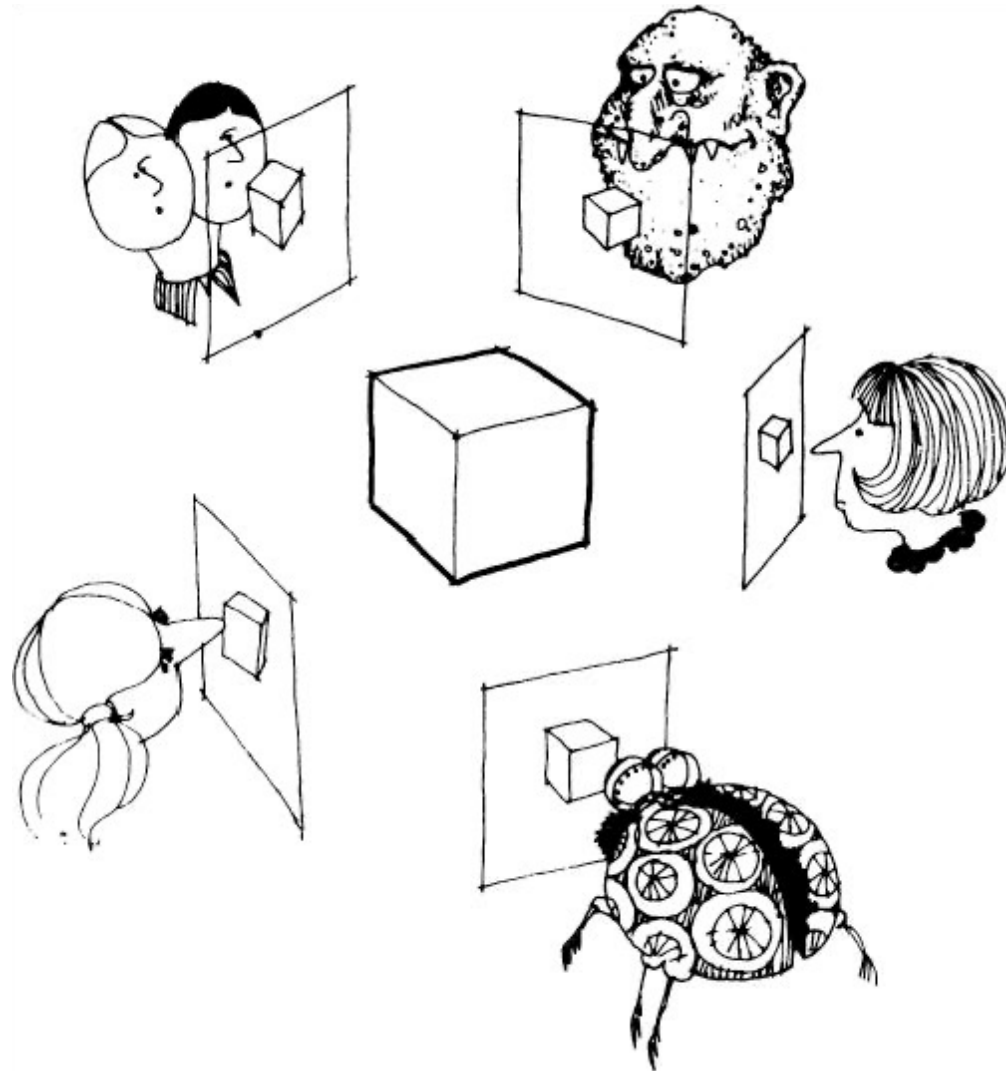
$$\begin{aligned} Z &= \text{distance} \\ d &= \text{disparity} \\ Z &= \frac{C}{d} \\ \delta Z &= \frac{-Z^2}{C} \delta d \\ |\delta Z| &= \frac{Z^2}{C} |\delta d| \\ \text{error} &\propto \text{distance}^2 \end{aligned}$$





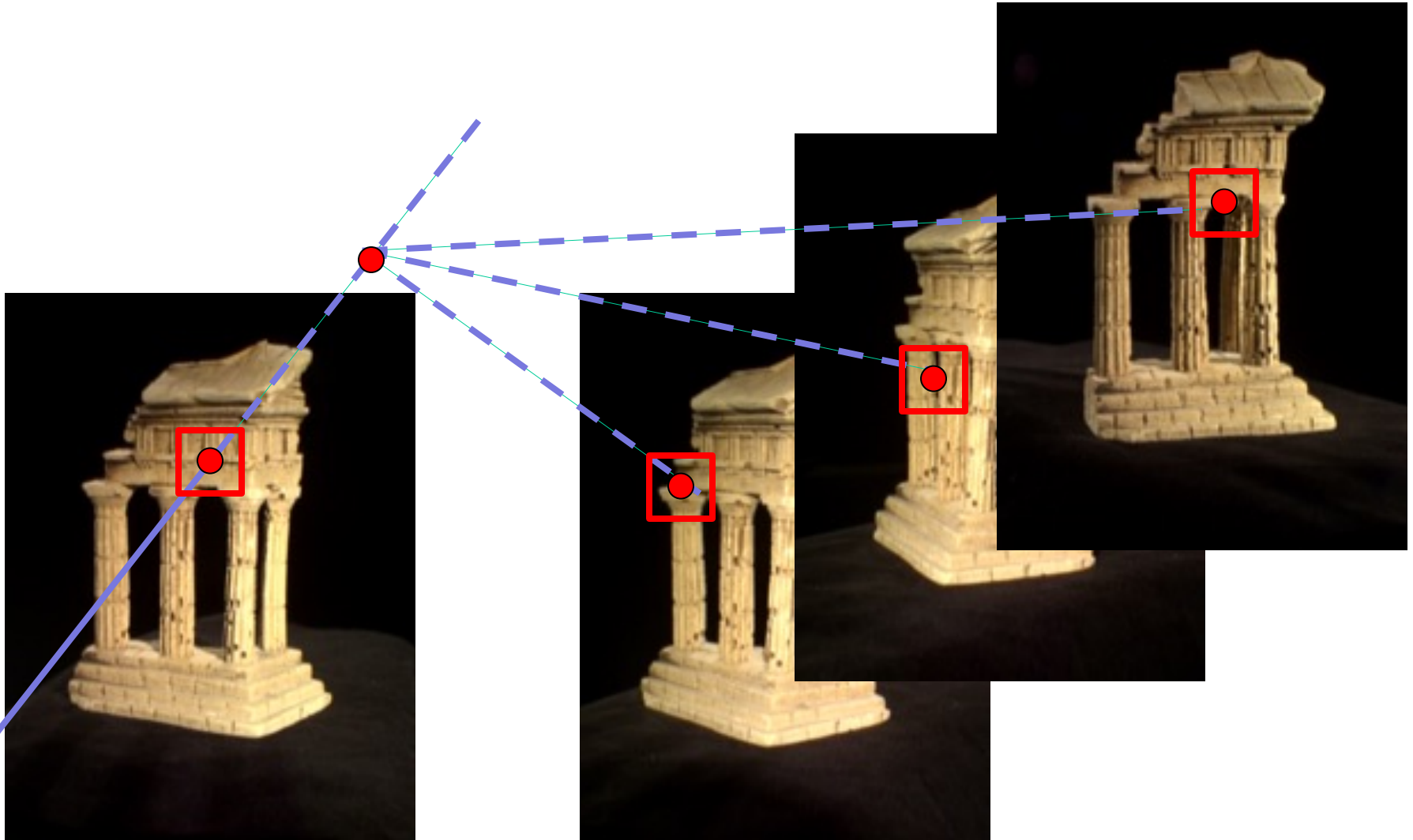
# Multi-view stereo

---

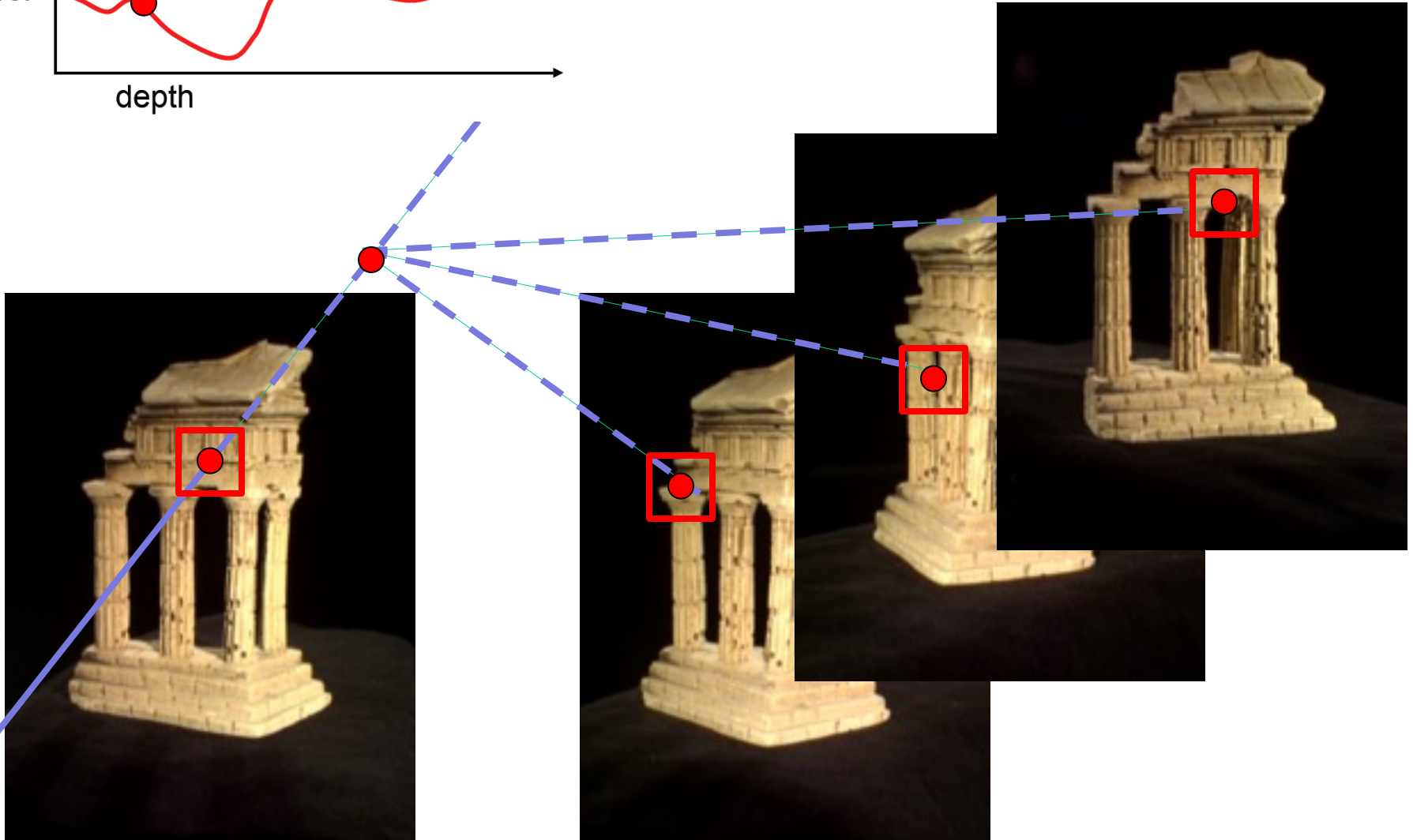
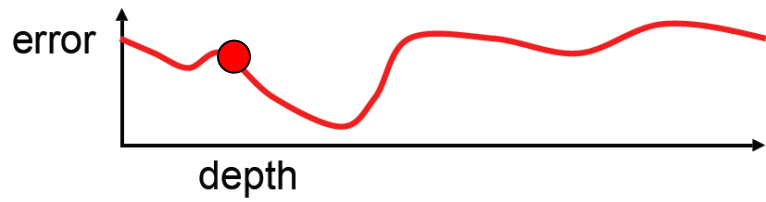


# Multi-view stereo: Basic idea

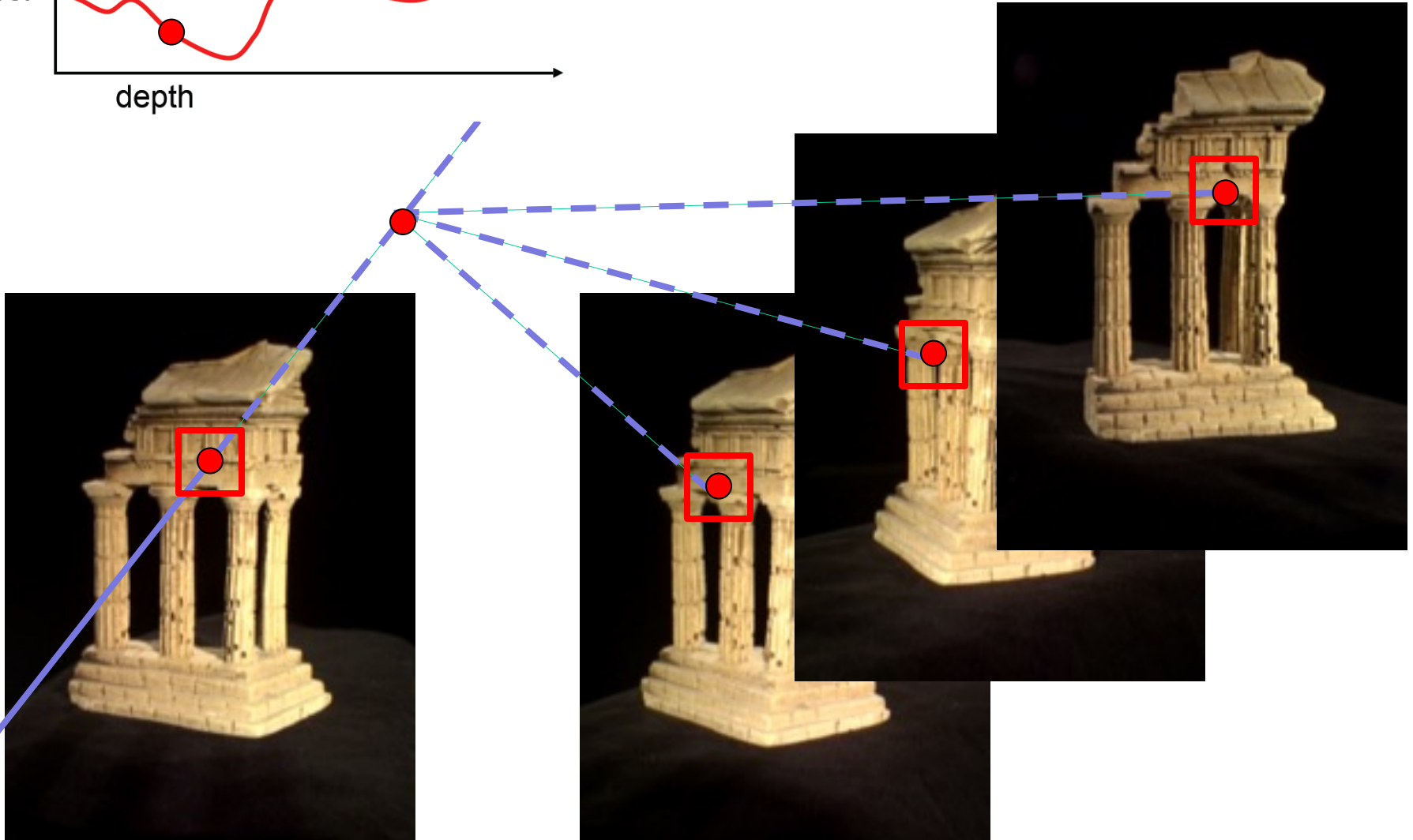
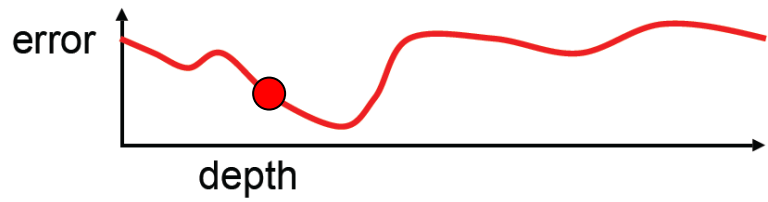
---



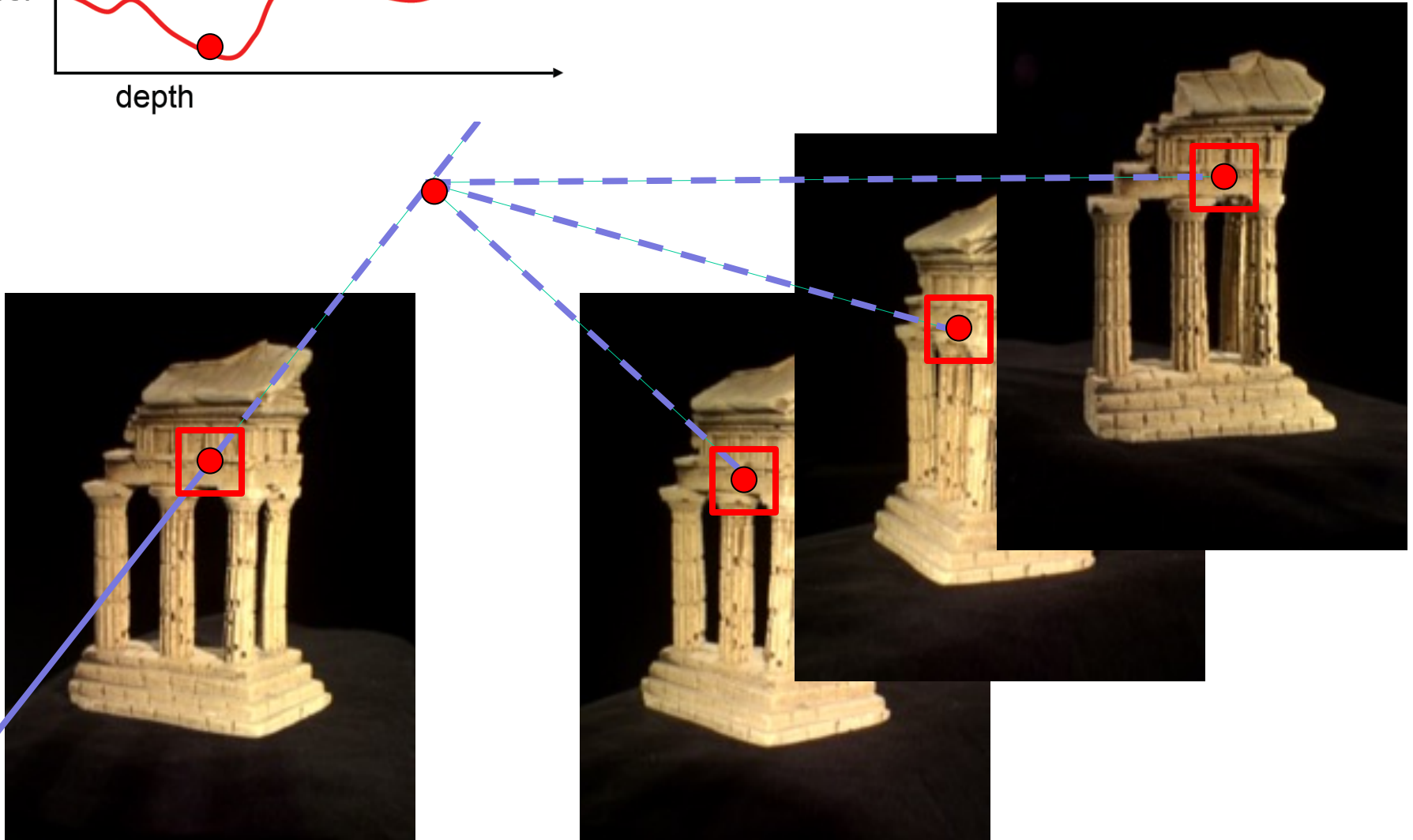
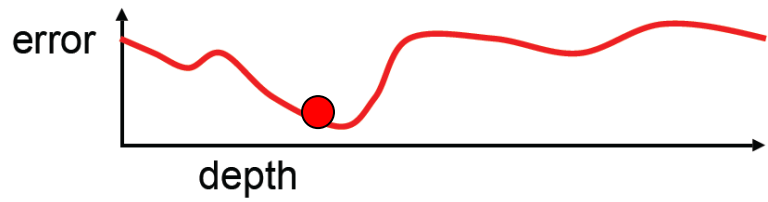
# Multi-view stereo: Basic idea



# Multi-view stereo: Basic idea



# Multi-view stereo: Basic idea



# Towards Internet-Scale Multi-View Stereo



[YouTube video](#), [CMVS software](#)

Y. Furukawa, B. Curless, S. Seitz and R. Szeliski, [Towards Internet-scale Multi-view Stereo](#), CVPR 2010.

# Applications

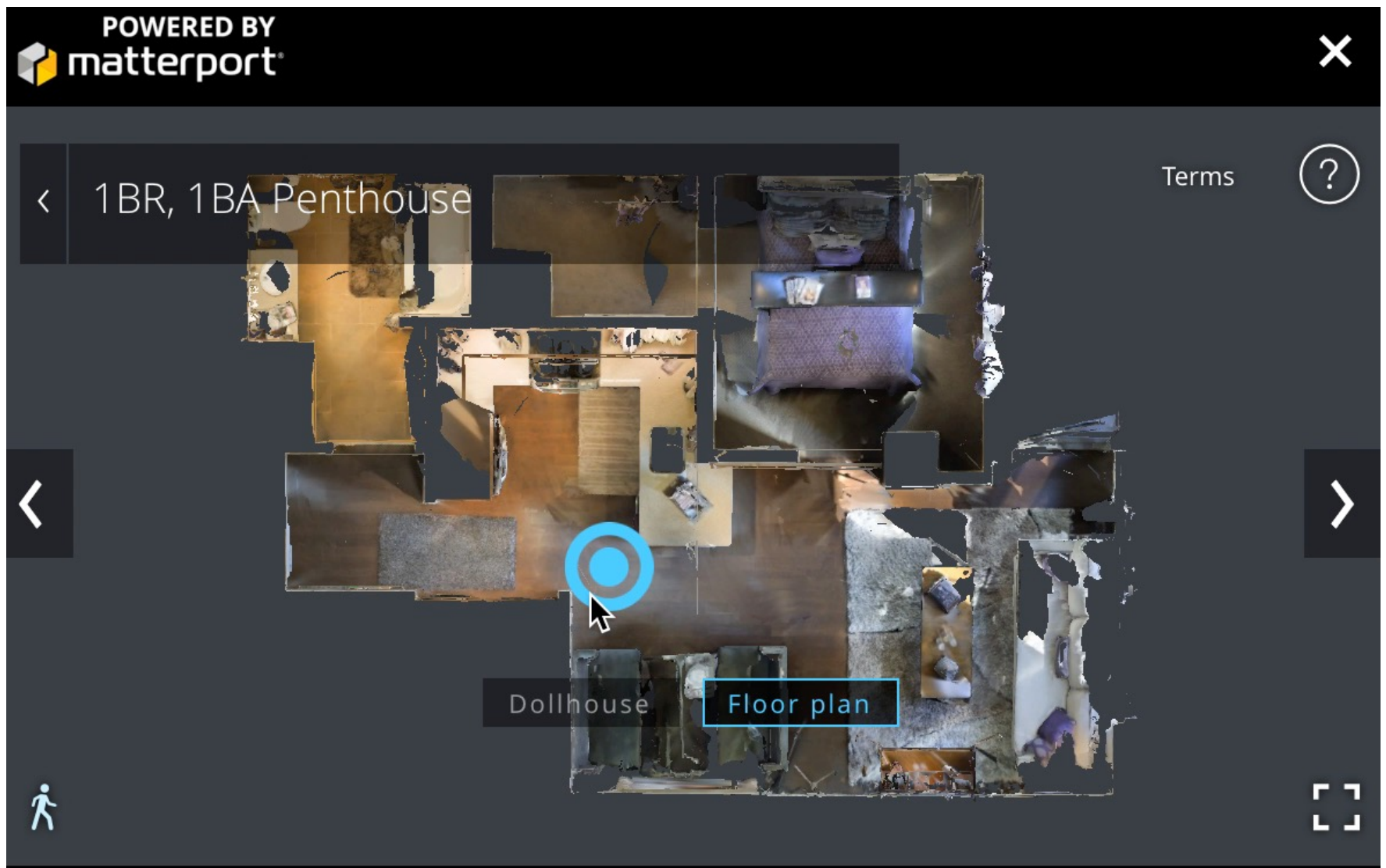
---



Data SIO, NOAA, U.S. Navy, NGA, GEBCO

Google earth

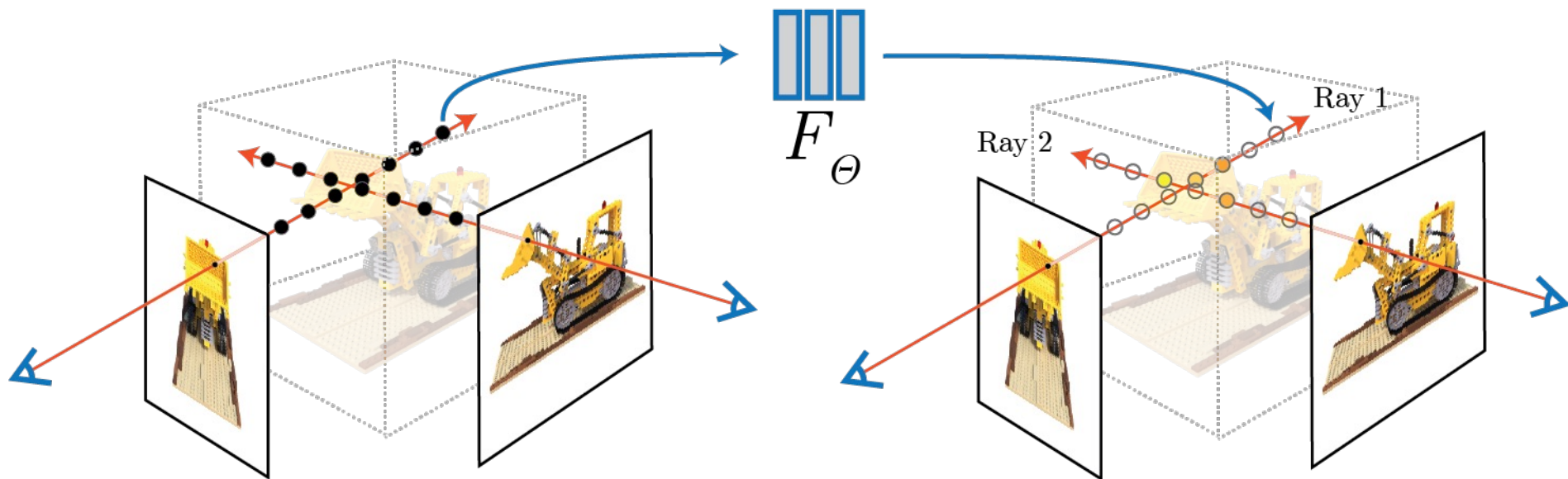
# Applications





# Latest and greatest: NeRF

---



[Representing Scenes as Neural Radiance Fields for View Synthesis](#). ECCV 2020.

Ben Mildenhall, Pratul P. Srinivasan, Matthew Tancik, Jonathan T. Barron, Ravi Ramamoorthi, Ren Ng.