



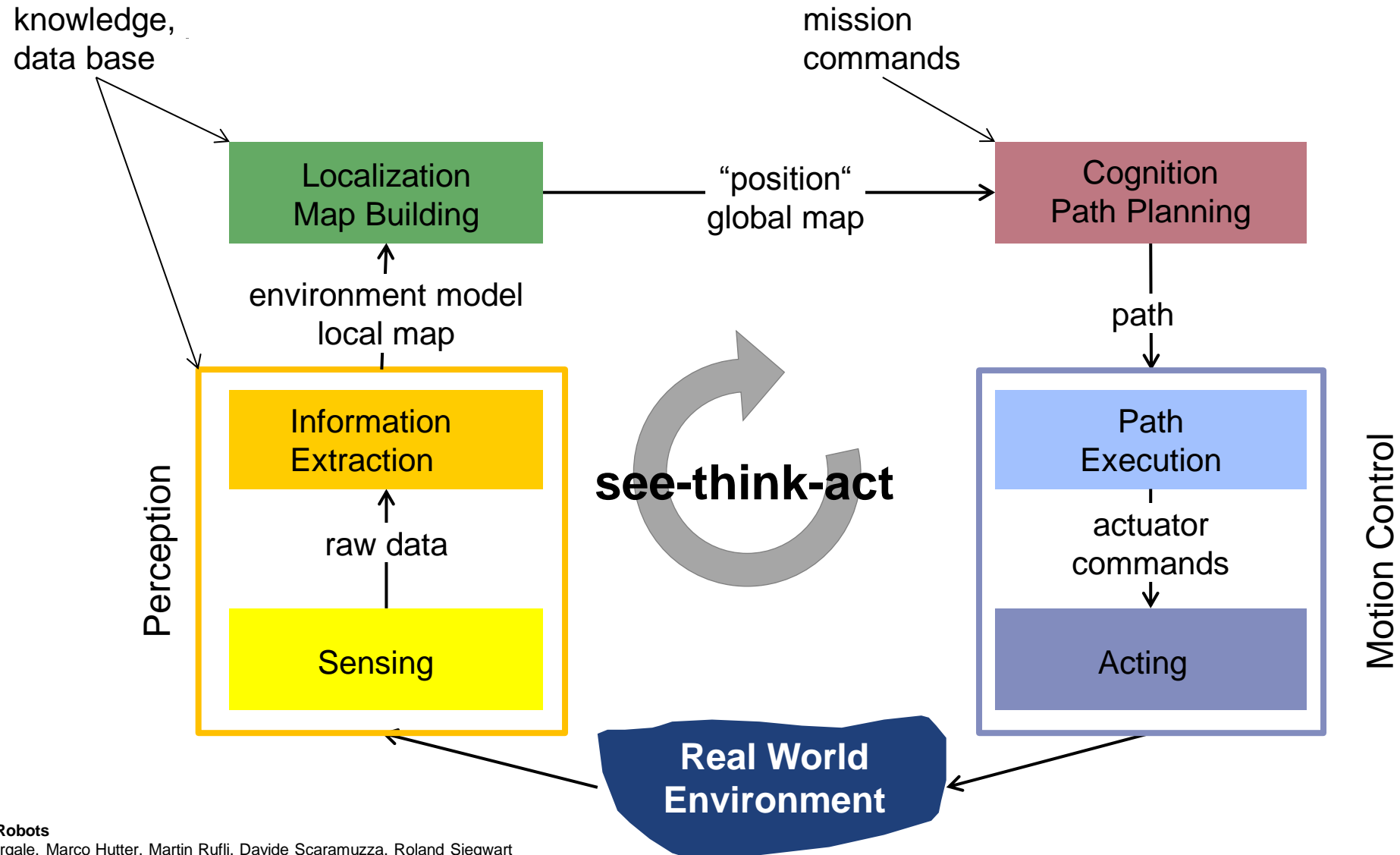
Perception II: Pinhole camera and Stereo Vision

Autonomous Mobile Robots

Davide Scaramuzza

Margarita Chli, Paul Furgale, Marco Hutter, Roland Siegwart

Mobile Robot Control Scheme

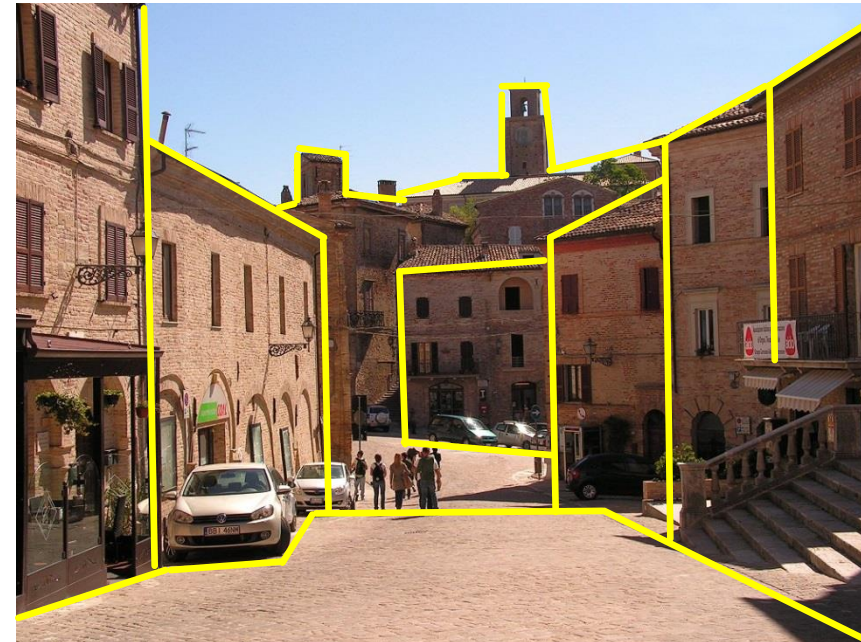


Computer vision | definition

- Automatic extraction of “meaningful” information from images and videos



Semantic information



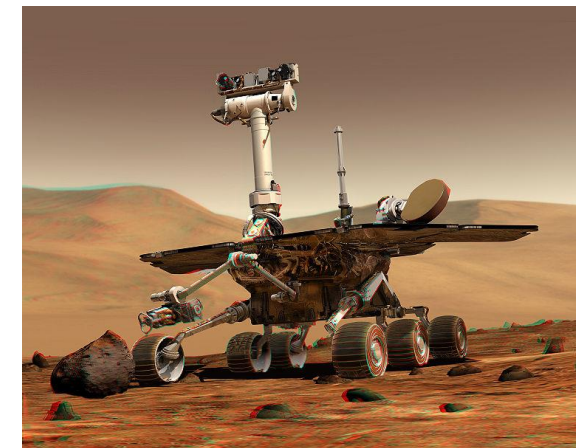
Geometric information

Computer vision | applications

- 3D reconstruction and modeling
- Recognition
- Motion capture
- Augmented reality:
- Video games and tele-operation
- Robot navigation and automotive
- Medical imaging



Google Earth, Microsoft's Bing Maps



Mars rover Spirit used cameras
for visual odometry

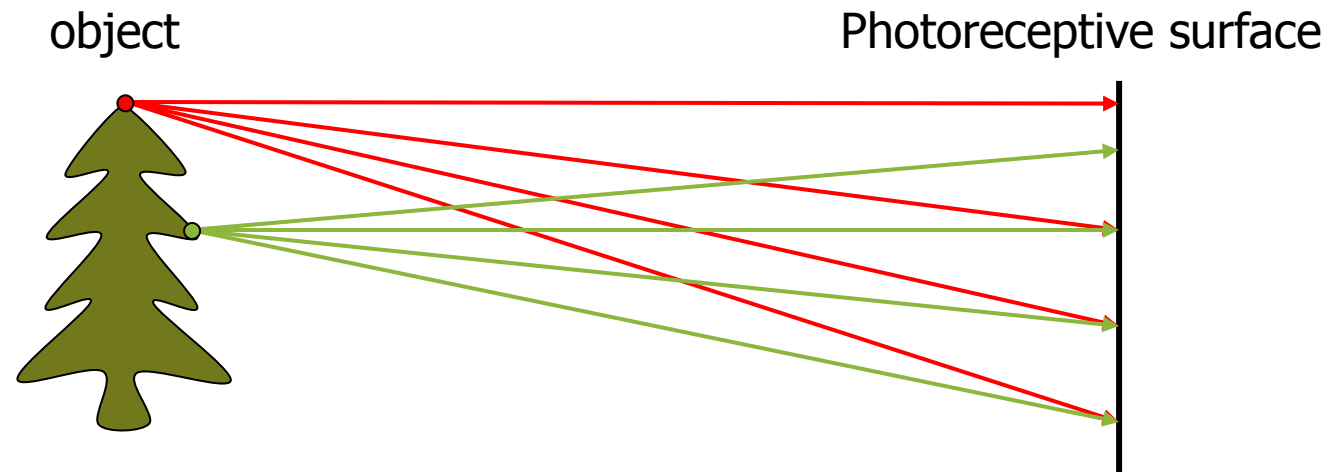
The camera



Sony Cybershot WX1

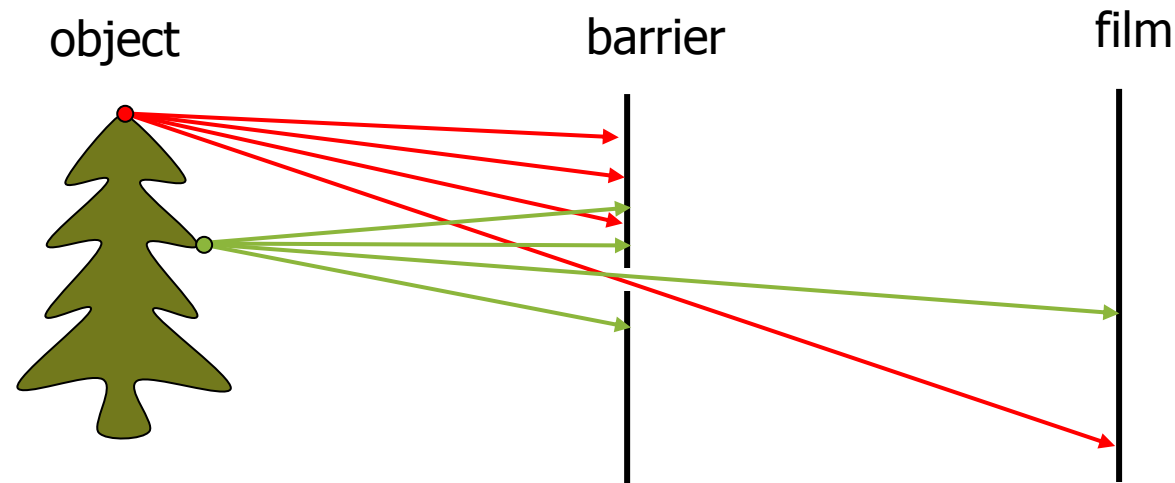
The camera | image formation

- If we place a piece of film in front of an object, do we get a reasonable image?



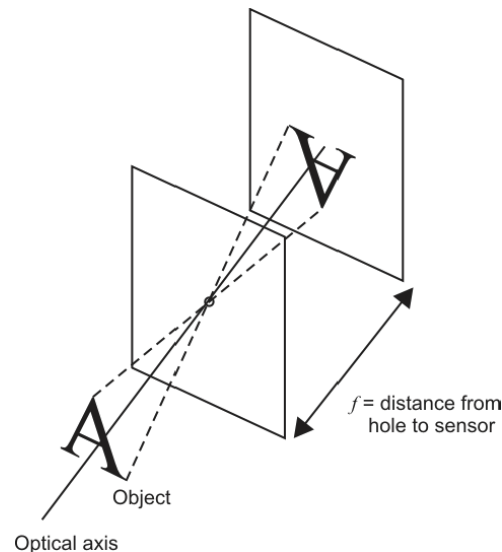
The camera | image formation

- If we place a piece of film in front of an object, do we get a reasonable image?
- Add a barrier to block off most of the rays
 - This reduces blurring
 - The opening is known as the **aperture**



The camera | camera obscura (pinhole camera)

- Pinhole model:
 - Captures **beam of rays** – all rays through a single point
 - The point is called **Center of Projection** or **Optical Center**
 - An “inverted” image is formed on the **Image Plane**
- We will use the pinhole camera model to describe how the image is formed



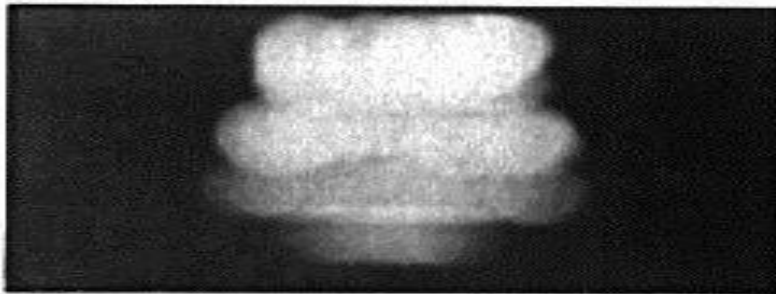
Gemma-Frisius (1508–1555)

Home-made pinhole camera



What can we do
to reduce the blur?

Shrinking the aperture



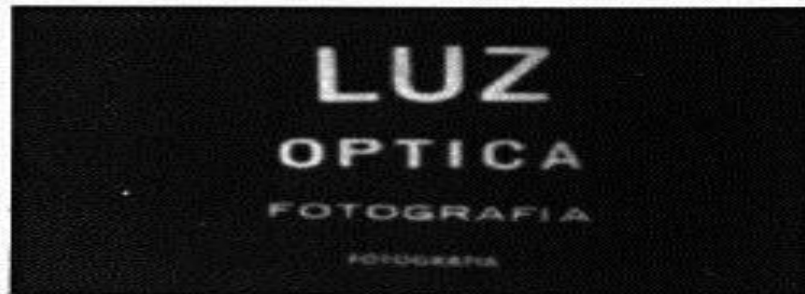
2 mm



1 mm



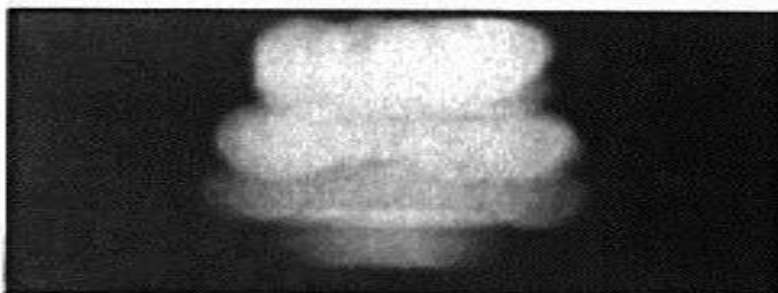
0.6 mm



0.35 mm

Why not make the aperture as small as possible?

Shrinking the aperture



2 mm



1 mm



0.6 mm



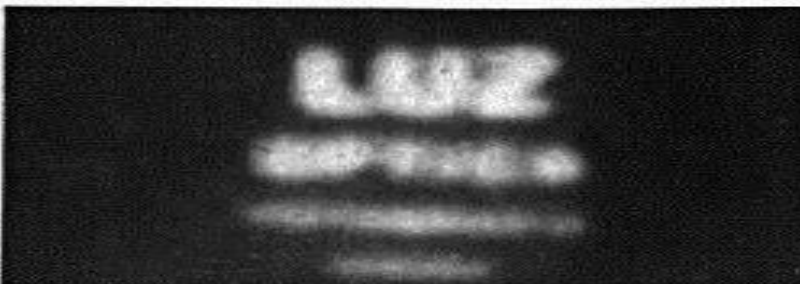
0.35 mm

Why not make the aperture as small as possible?

- Less light gets through (must increase the exposure)
- Diffraction effects...



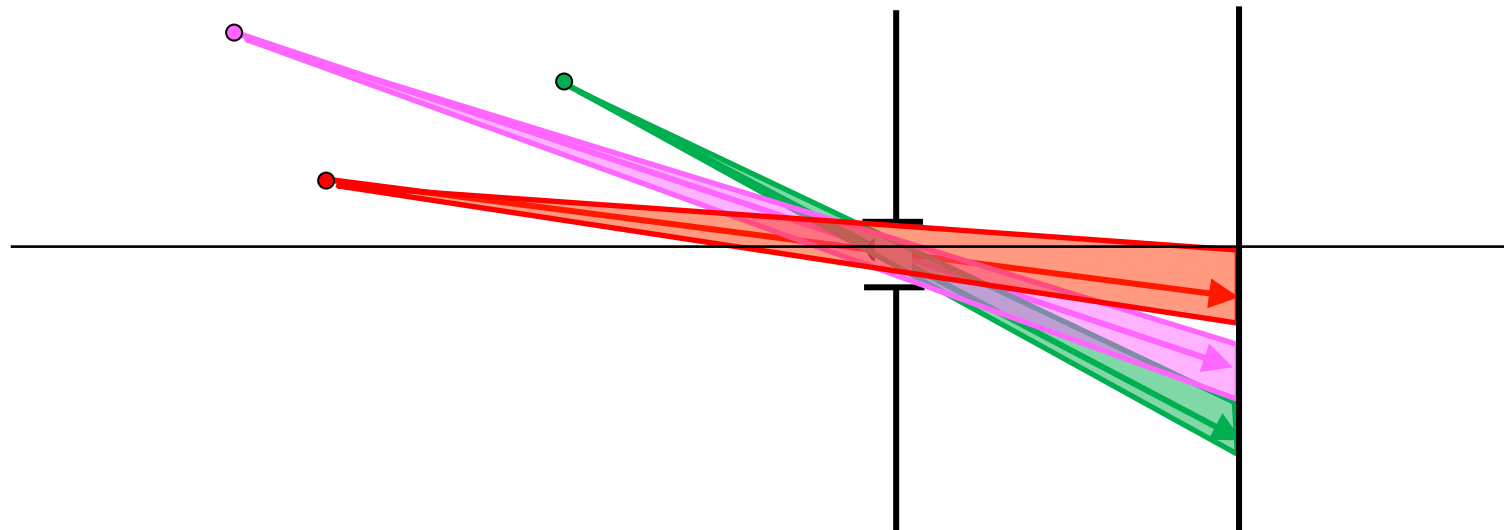
0.15 mm



0.07 mm

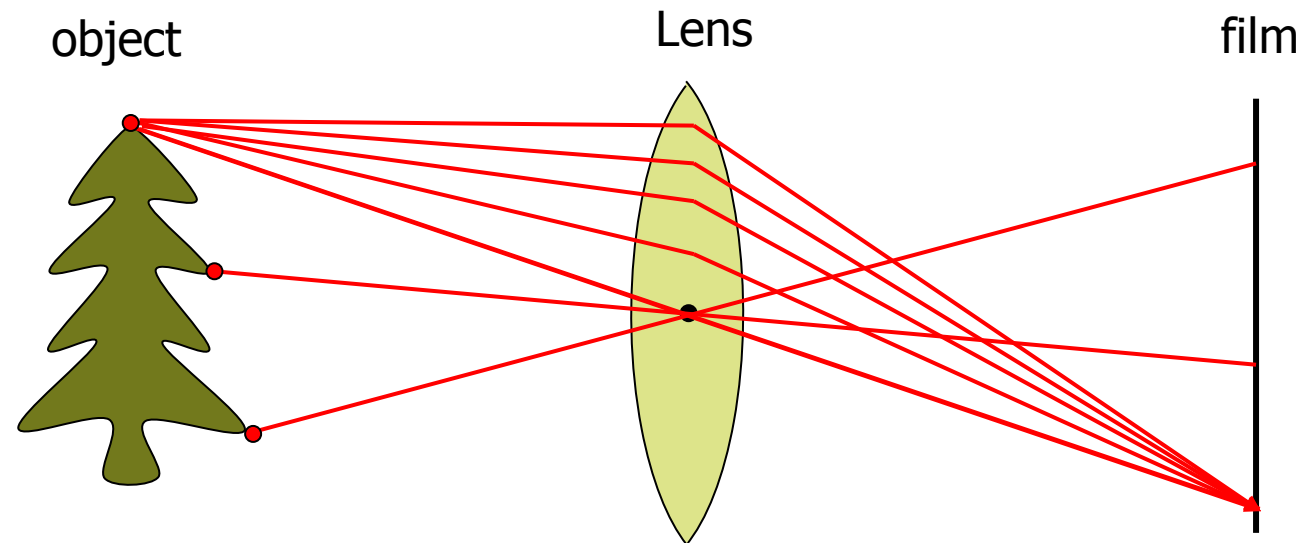
The camera | why use a lens?

- The ideal pinhole: only one ray of light reaches each point on the film
 - \Rightarrow image can be very dim; gives rise to diffraction effects
- Making the pinhole bigger (i.e. aperture) makes the image blurry



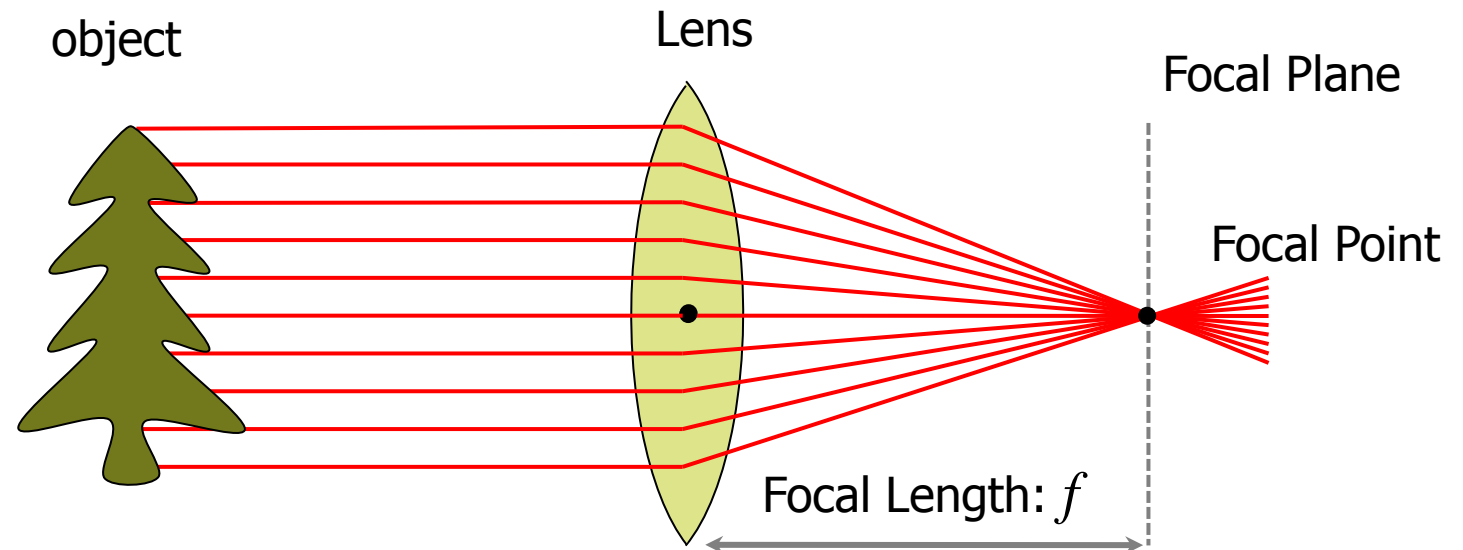
The camera | why use a lens?

- A lens focuses light onto the film
- Rays passing through the **optical center** are not deviated



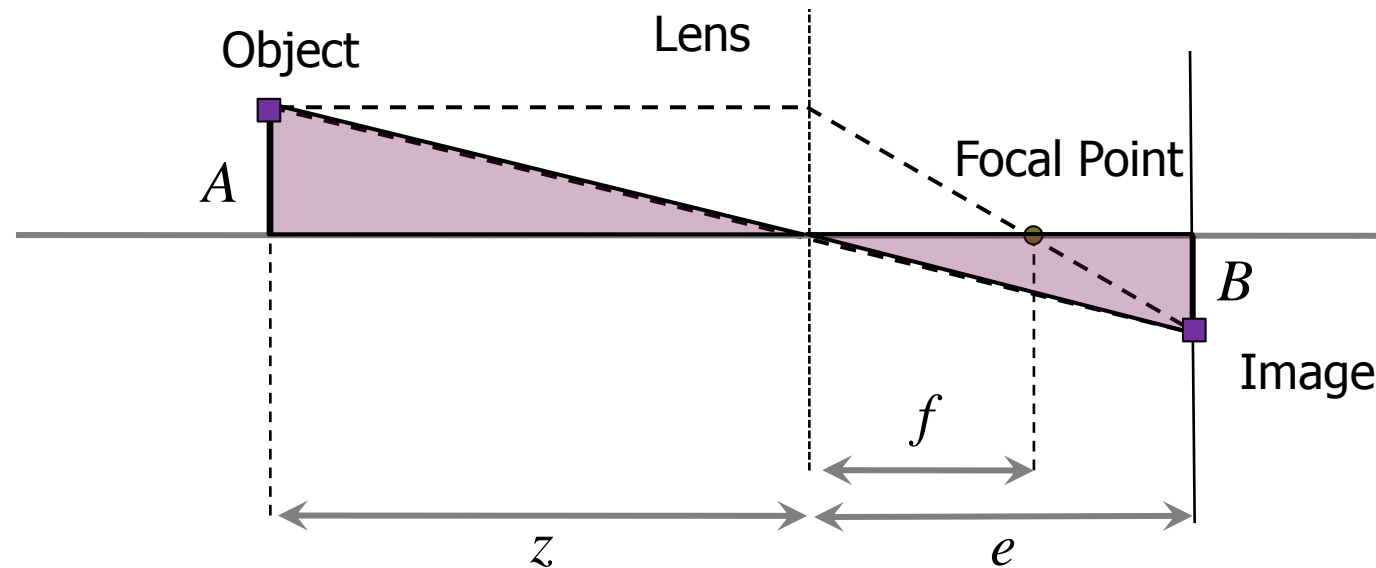
The camera | why use a lens?

- A lens focuses light onto the film
- Rays passing through the **optical center** are not deviated
- All rays parallel to the **optical axis** converge at the **focal point**



The camera | this lens equation

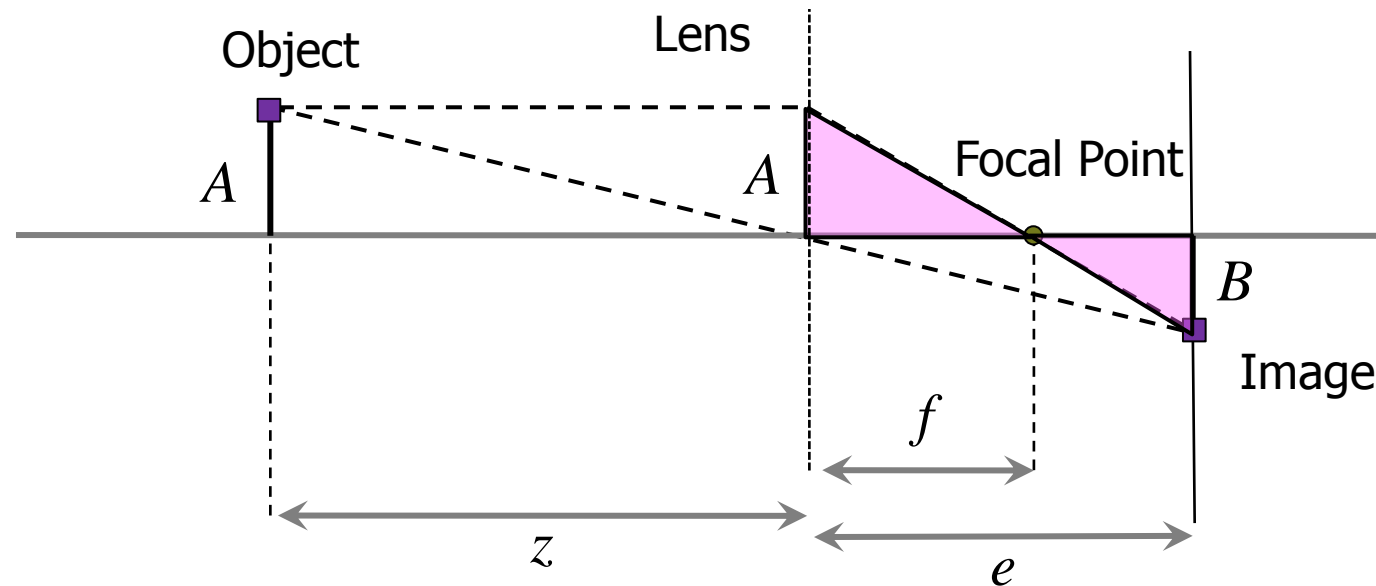
- What is the relationship between f , z , and e ?



- Similar Triangles: $\frac{B}{A} = \frac{e}{z}$

The camera | this lens equation

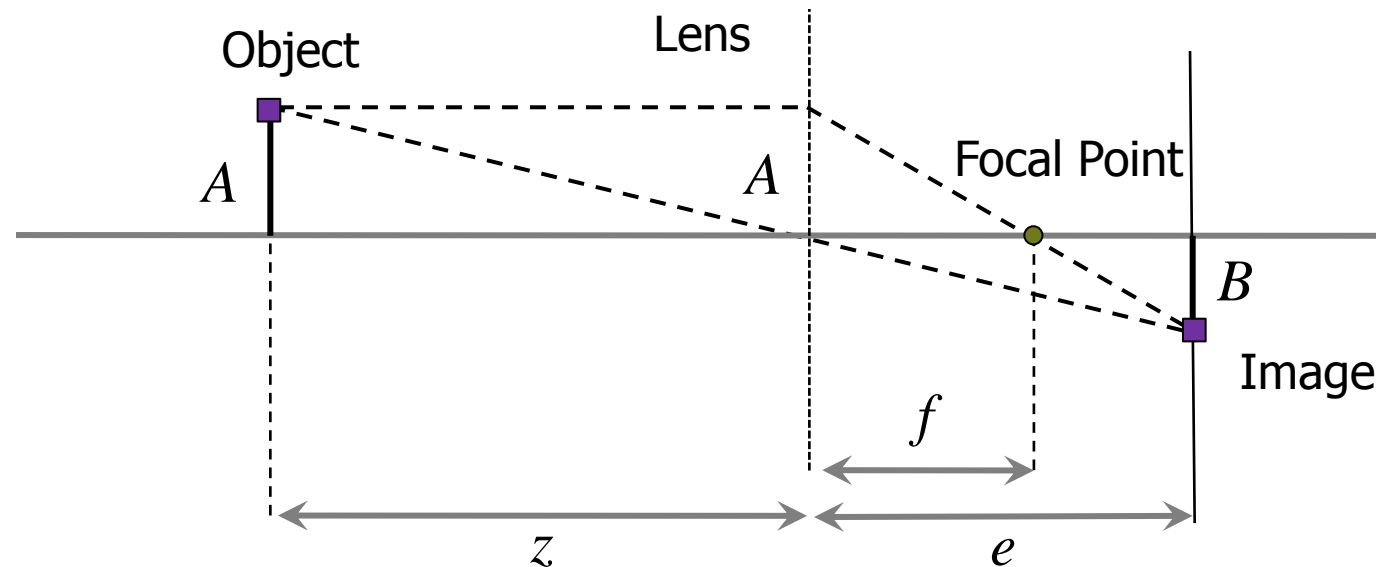
- What is the relationship between f , z , and e ?



- Similar Triangles: $\frac{B}{A} = \frac{e}{z}$
 - $\frac{B}{A} = \frac{e-f}{f} = \frac{e}{f} - 1$
 - $\frac{e}{f} - 1 = \frac{e}{z}$
- “Thin lens equation”

The camera | pinhole approximation

- What happens if $z \gg f$?

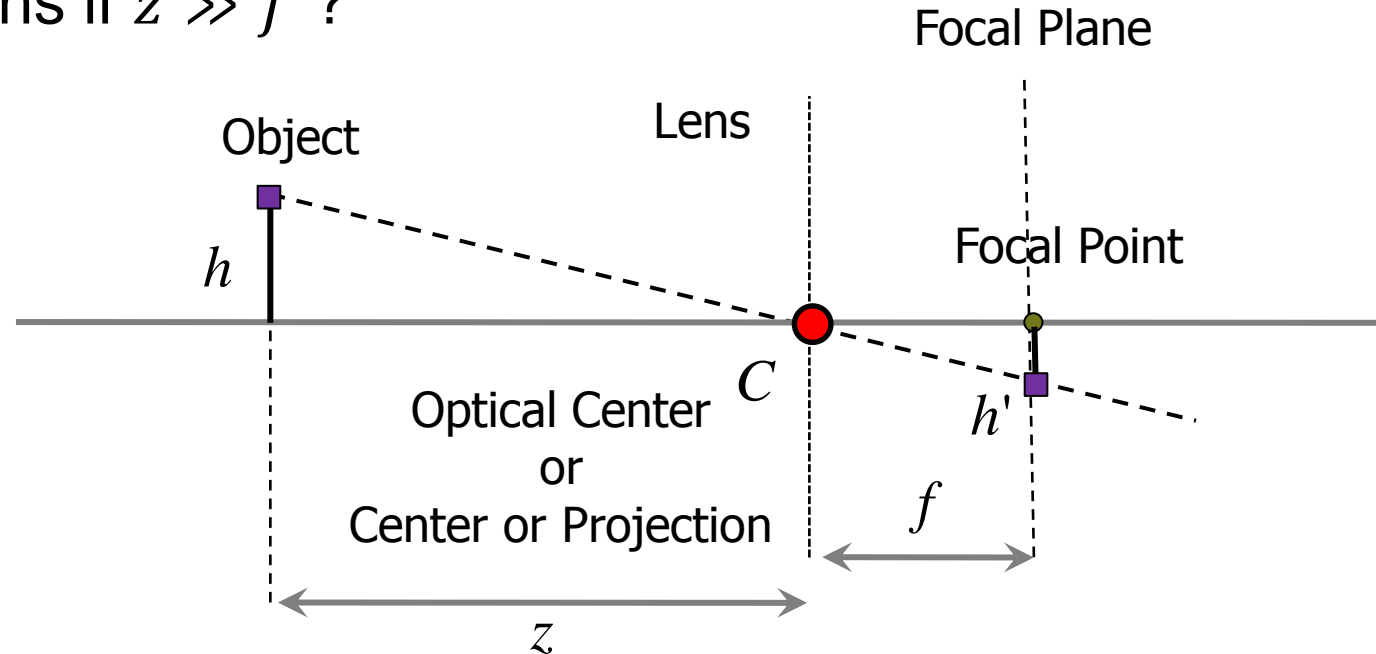


- We need to adjust the image plane such that objects at infinity are in focus

$$\frac{1}{f} = \underbrace{\frac{1}{z}}_{\cong 0} + \frac{1}{e} \Rightarrow \frac{1}{f} \approx \frac{1}{e} \Rightarrow f \approx e$$

The camera | pinhole approximation

- What happens if $z \gg f$?

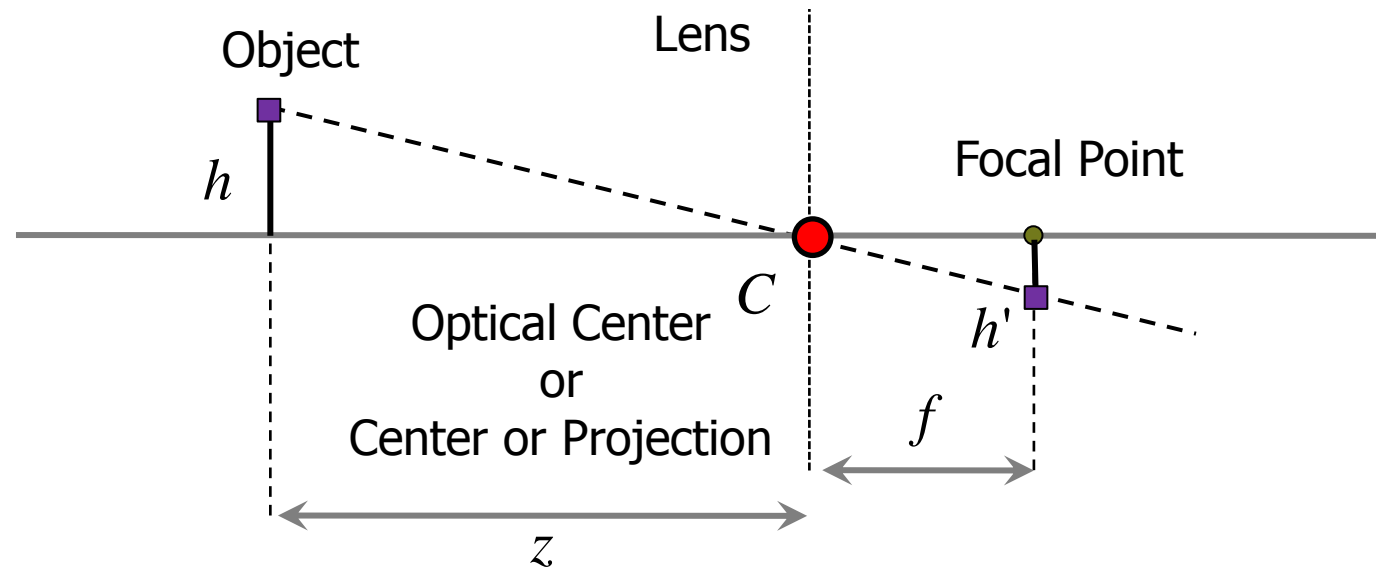


- We need to adjust the image plane such that objects at infinity are in focus

$$\frac{1}{f} = \underbrace{\frac{1}{z}}_{\cong 0} + \frac{1}{e} \Rightarrow \frac{1}{f} \approx \frac{1}{e} \Rightarrow f \approx e$$

The camera | pinhole approximation

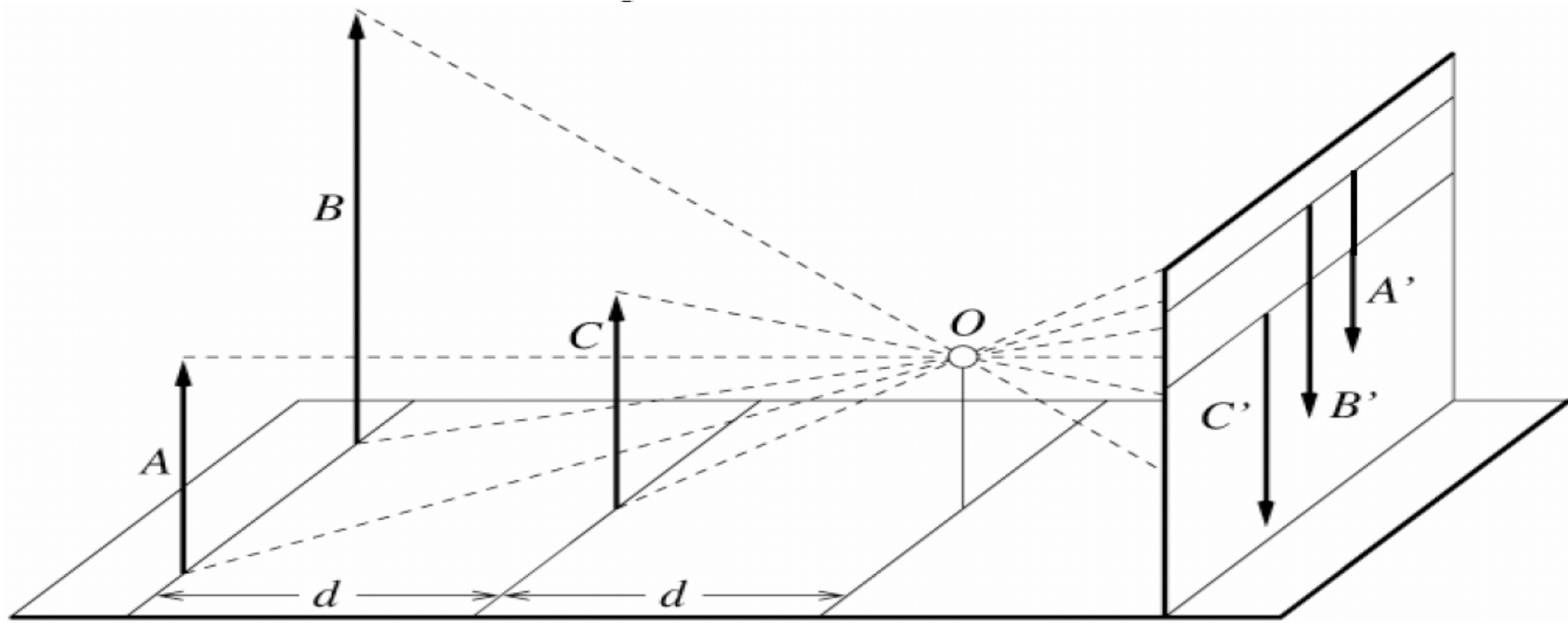
- What happens if $z \gg f$?



$$\frac{h'}{h} = \frac{f}{z} \Rightarrow h' = \frac{f}{z} h$$

Perspective effects

- Far away objects appear smaller



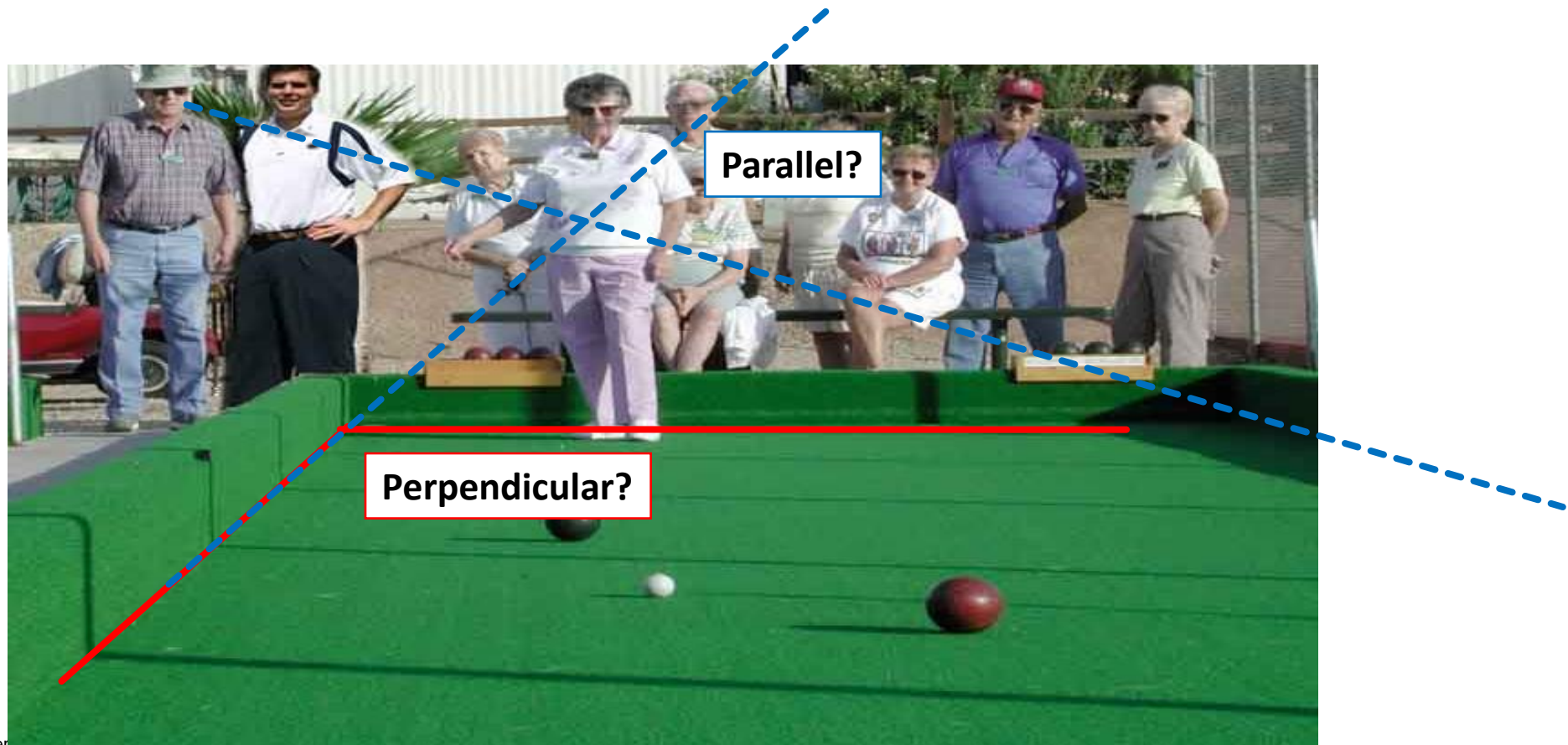
Perspective effects



Projective Geometry

What is lost?

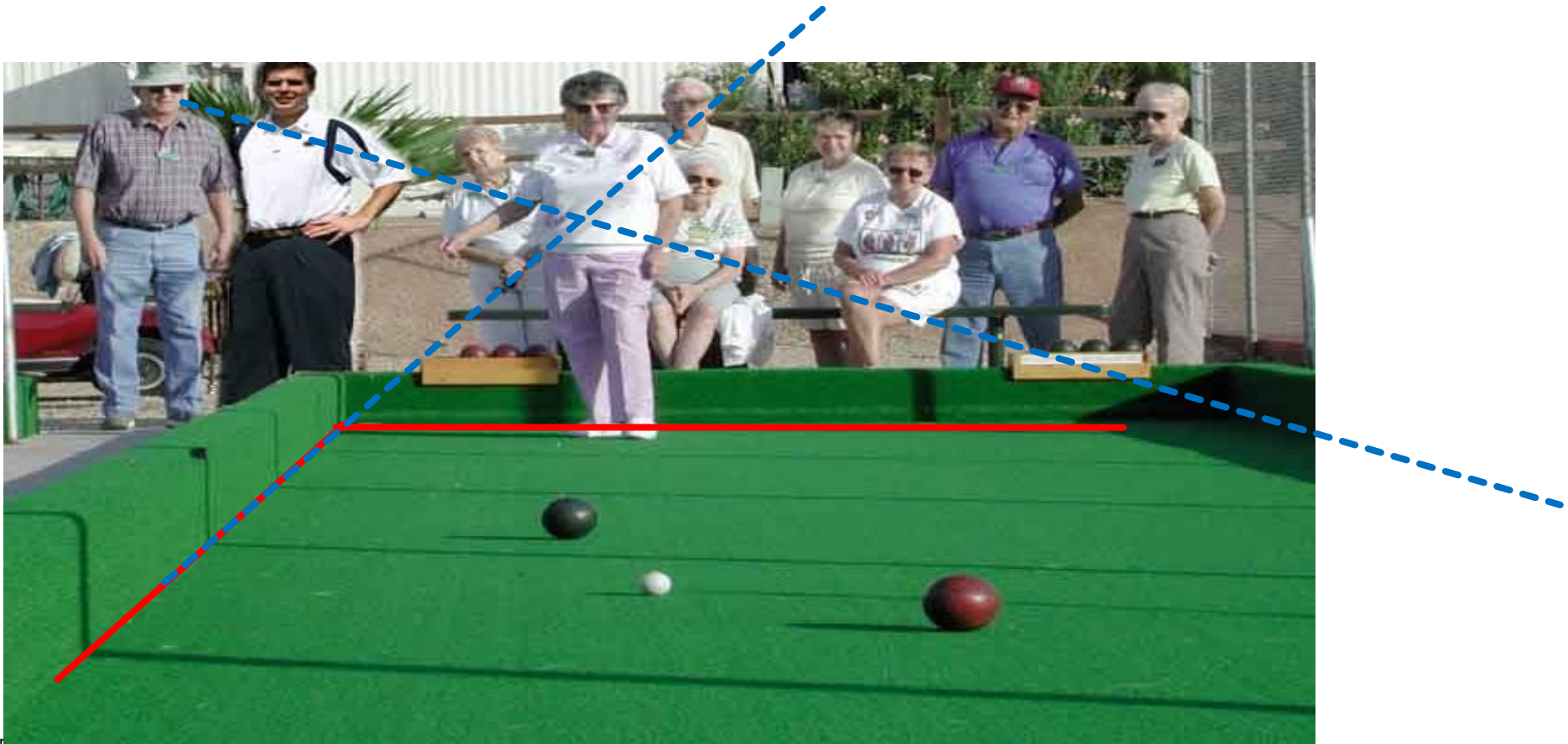
- Length
- Angles



Projective Geometry

What is preserved?

- Straight lines are still straight

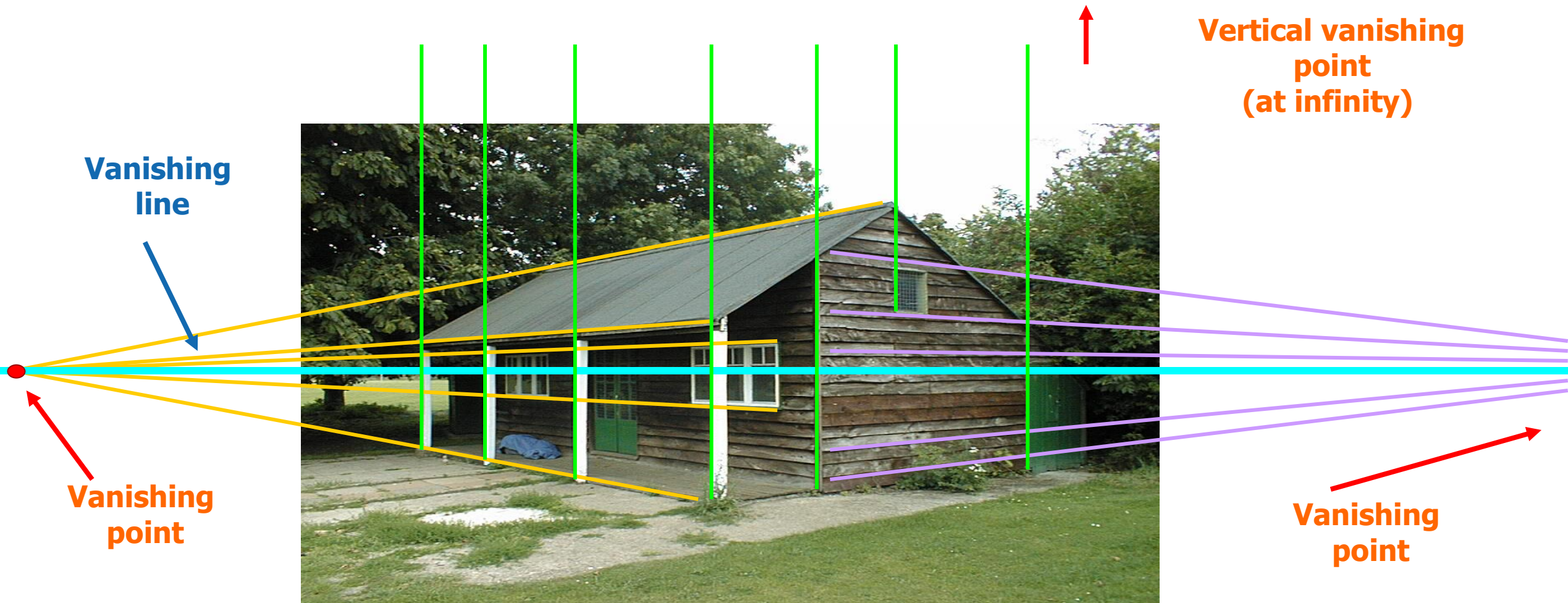


Vanishing points and lines

- Parallel lines in the world intersect in the image at a “vanishing point”



Vanishing points and lines



Perspective and art

- Use of correct perspective projection indicated in 1st century B.C. frescoes
- Skill resurfaces in Renaissance: artists develop systematic methods to determine perspective projection (around 1480-1515)



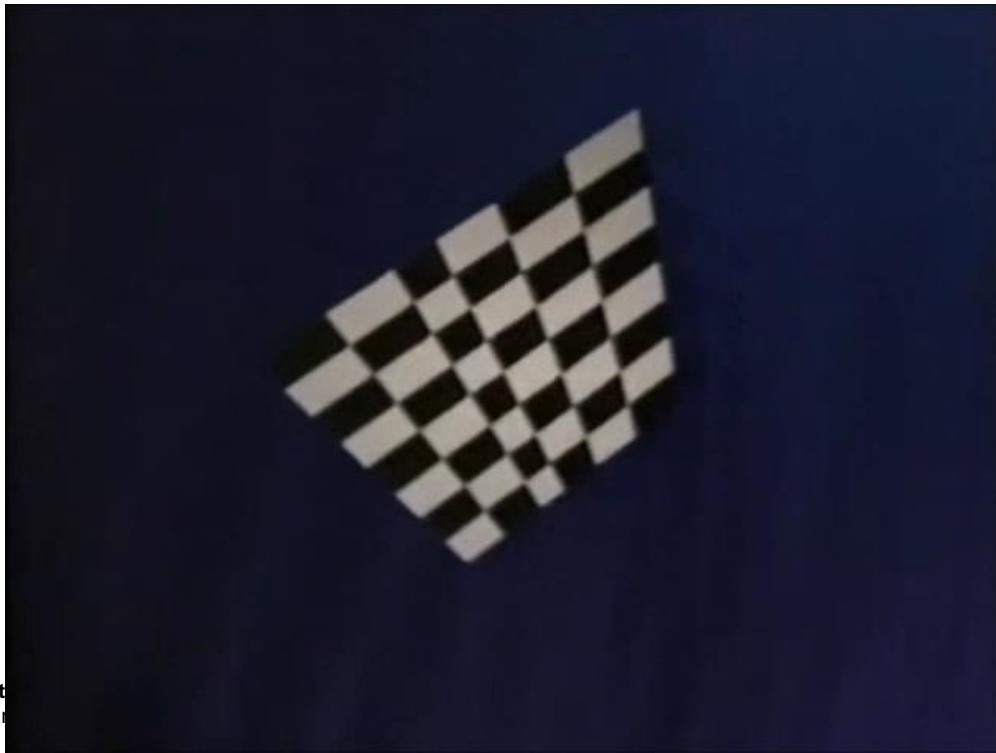
Raphael



Durer, 1525

Playing with Perspective

- Perspective gives us very strong depth cues
⇒ hence we can perceive a 3D scene by viewing its 2D representation (i.e. image)
- An example where perception of 3D scenes is misleading:



“Ames room”

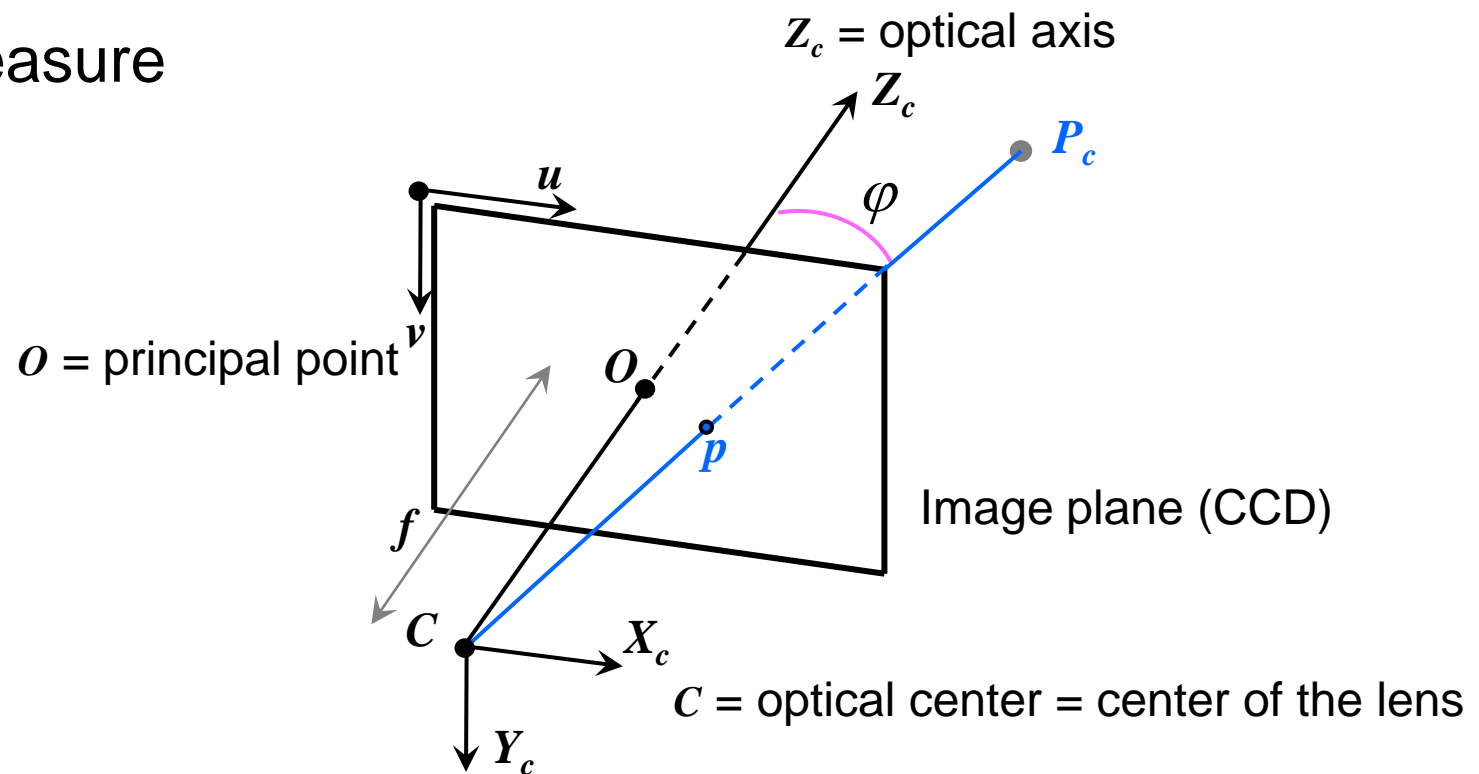
A clip from "The computer that ate Hollywood" documentary. Dr. Vilayanur S. Ramachandran.

Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm

The camera | perspective camera

- For convenience, the image plane is usually represented in front of C such that the image preserves the same orientation (i.e. not flipped)
- A camera does not measure distances but angles!



Perspective projection| from scene points to pixels

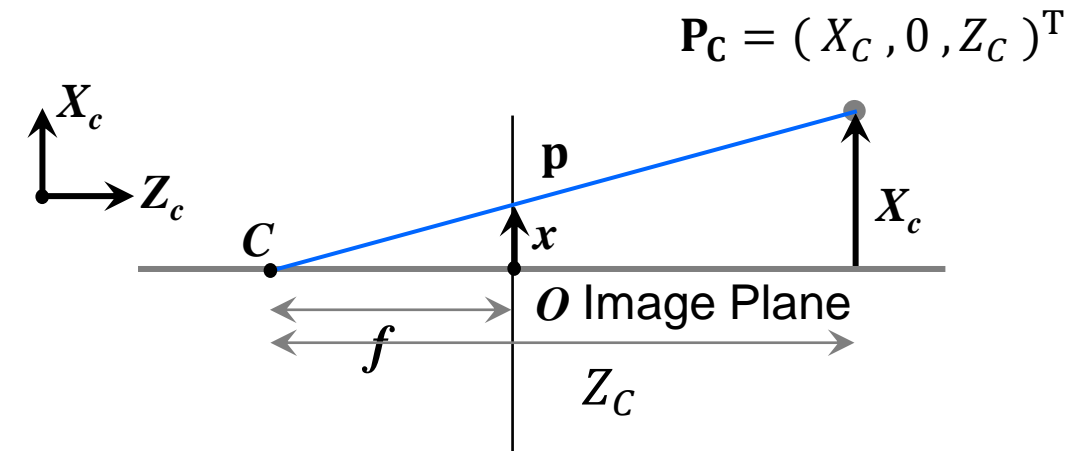
- The Camera point $\mathbf{P}_c = (X_c, 0, Z_c)^T$ projects to $\mathbf{p} = (x, y)$ onto the image plane

- From similar triangles:

$$\frac{x}{f} = \frac{X_c}{Z_c} \Rightarrow x = \frac{fX_c}{Z_c}$$

- Similarly, in the general case:

$$\frac{y}{f} = \frac{Y_c}{Z_c} \Rightarrow y = \frac{fY_c}{Z_c}$$

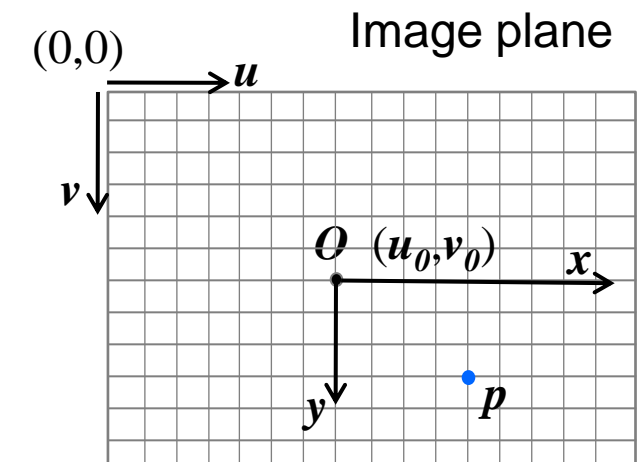


Perspective projection| from scene points to pixels

- To convert \mathbf{p} , from the local image plane coordinates (x, y) to the pixel coordinates (u, v) , we need to account for:
 - The pixel coordinates of the camera optical center $O = (u_0, v_0)$
 - Scale factor k for the pixel-size

$$u = u_0 + kx \Rightarrow u_0 + k \frac{fX_C}{Z_C}$$

$$v = v_0 + ky \Rightarrow v_0 + k \frac{fY_C}{Z_C}$$



- Use Homogeneous Coordinates for linear mapping from 3D to 2D, by introducing an extra element (scale):

$$p = \begin{pmatrix} u \\ v \end{pmatrix} \Rightarrow \tilde{p} = \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$

Perspective projection| from scene points to pixels

$$u = u_0 + kx \Rightarrow u_0 + k \frac{fX_c}{Z_c}$$

$$v = v_0 + ky \Rightarrow v_0 + k \frac{fY_c}{Z_c}$$

- Expressed in matrix form and homogeneous coordinates:

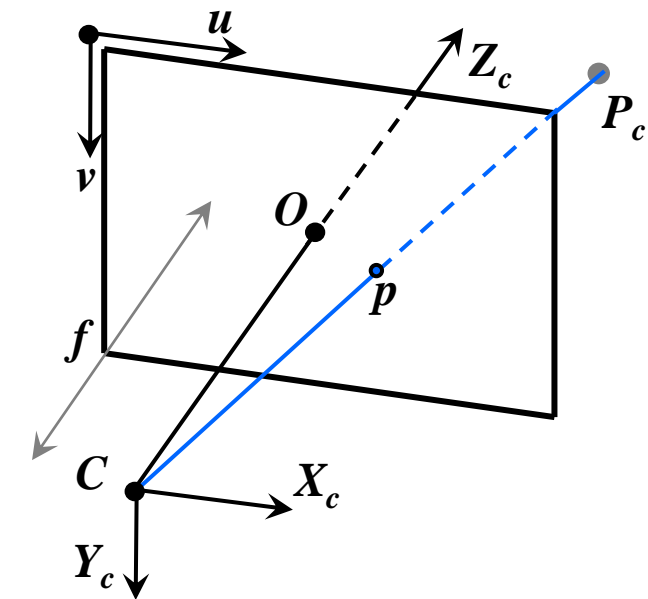
$$\begin{bmatrix} \lambda u \\ \lambda v \\ \lambda \end{bmatrix} = \begin{bmatrix} kf & 0 & u_0 \\ 0 & kf & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

- Or alternatively

$$\begin{bmatrix} \lambda u \\ \lambda v \\ \lambda \end{bmatrix} = \begin{bmatrix} \alpha & 0 & u_0 \\ 0 & \alpha & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = K \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

Focal length in pixels

Intrinsic parameters matrix



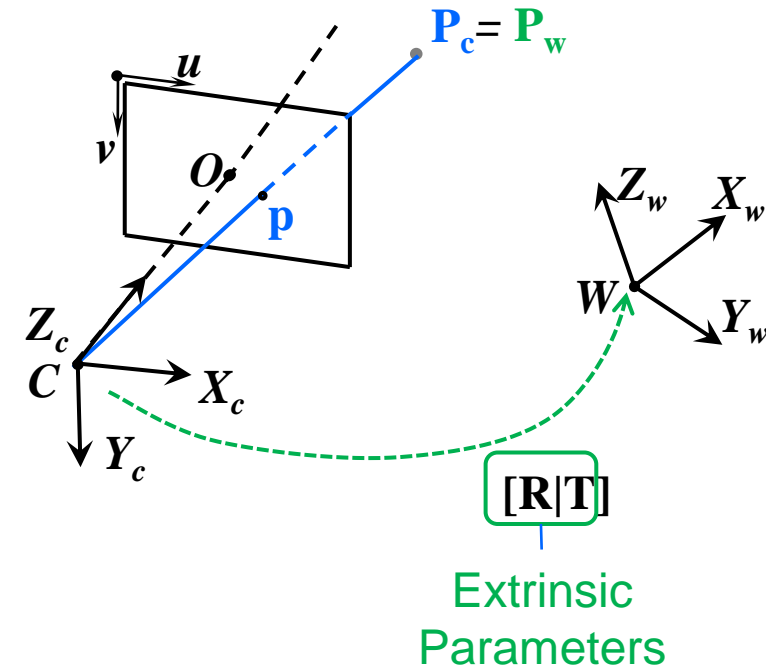
Perspective projection| from scene points to pixels

$$\begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \end{bmatrix} + \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} = \begin{bmatrix} R & | & T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

Perspective Projection Matrix

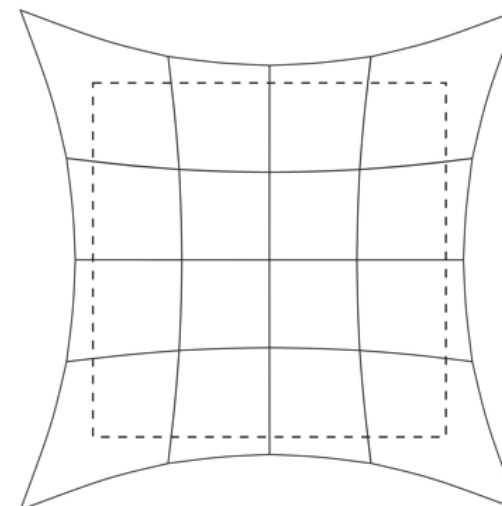
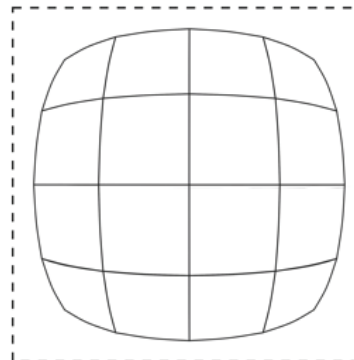
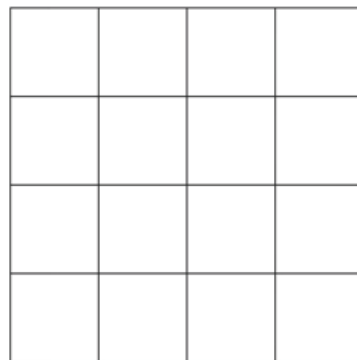
$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \boxed{[R|T]} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$



Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision

Perspective projection| radial distortion



No distortion



Barrel distortion



Pincushion

Perspective projection| radial distortion

- The standard model of radial distortion is a transformation from the ideal coordinates (u, v) (i.e., undistorted) to the real observable coordinates (distorted) (u_d, v_d)
- The amount of distortion of the coordinates of the observed image is a nonlinear function of their radial distance. For most lenses, a simple quadratic model of distortion produces good results

where

$$\begin{bmatrix} u_d \\ v_d \end{bmatrix} = (1 + k_1 r^2) \begin{bmatrix} u - u_0 \\ v - v_0 \end{bmatrix} + \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}$$

$$r^2 = (u - u_0)^2 + (v - v_0)^2$$

Summary: Perspective projection equations

- To recap, a 3D world point $P = (X_w, Y_w, Z_w)$ projects into the image point $p = (u, v)$

$$\lambda p = \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K [R | T] \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad \text{where} \quad K = \begin{bmatrix} \alpha & 0 & u_0 \\ 0 & \alpha & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

and λ is the depth ($\lambda = Z_w$) of the scene point

- If we want to take into account for the radial distortion, then the distorted coordinates (u_d, v_d) (in pixels) can be obtained as

$$\begin{bmatrix} u_d \\ v_d \end{bmatrix} = (1 + k_1 r^2) \begin{bmatrix} u - u_0 \\ v - v_0 \end{bmatrix} + \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}$$

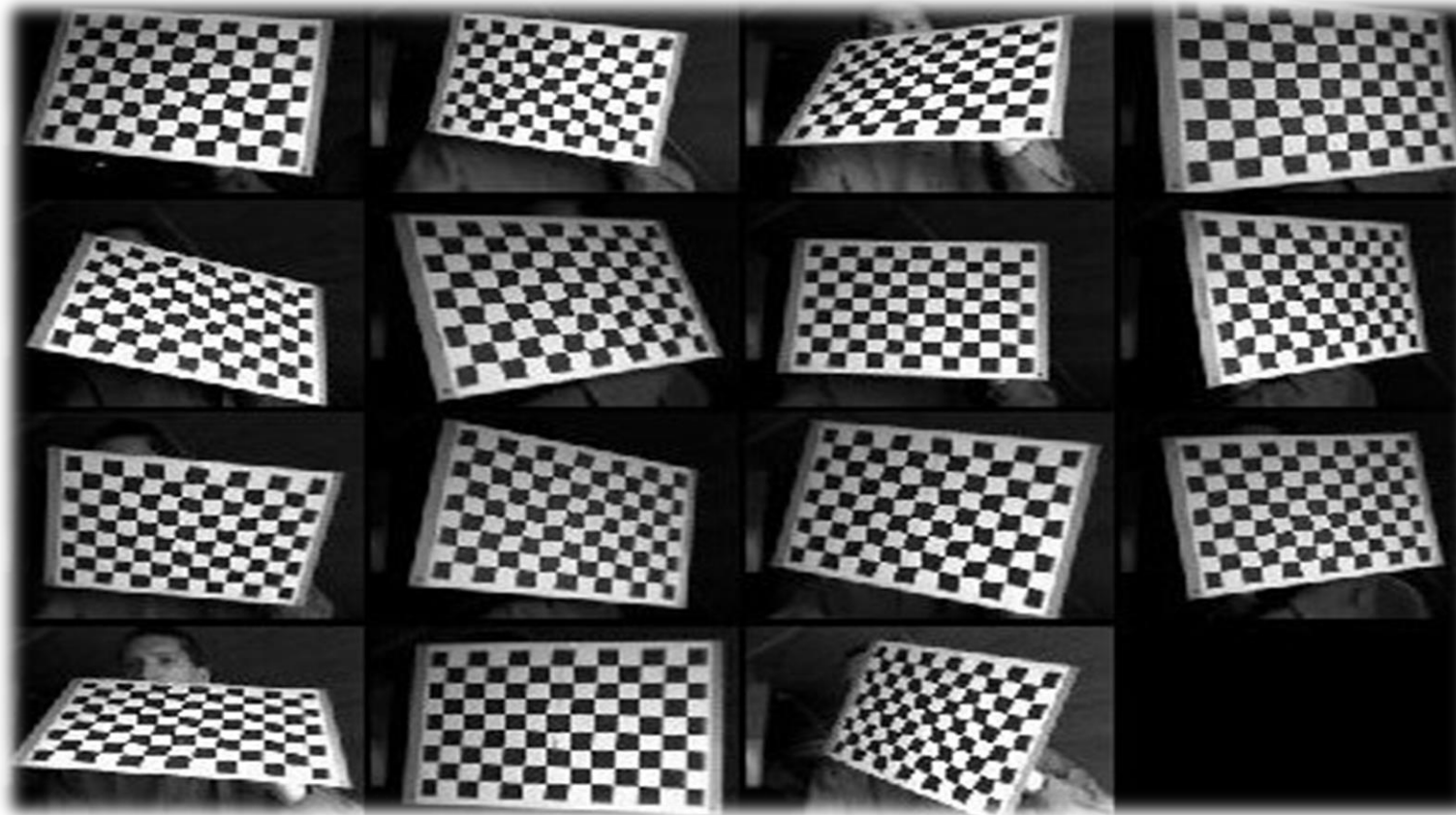
where $r^2 = (u - u_0)^2 + (v - v_0)^2$

Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision

Camera Calibration

- Procedure to determine the *intrinsic parameters* of a camera

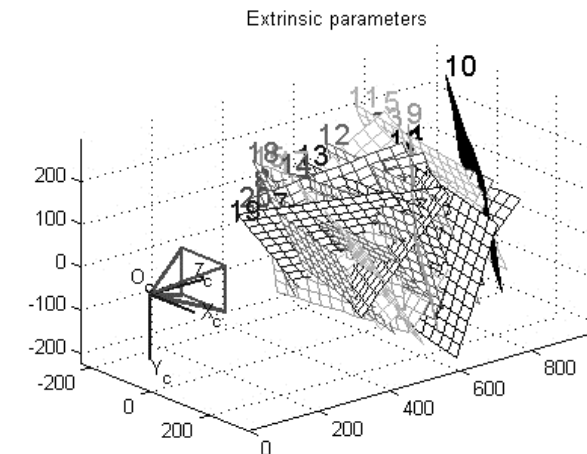
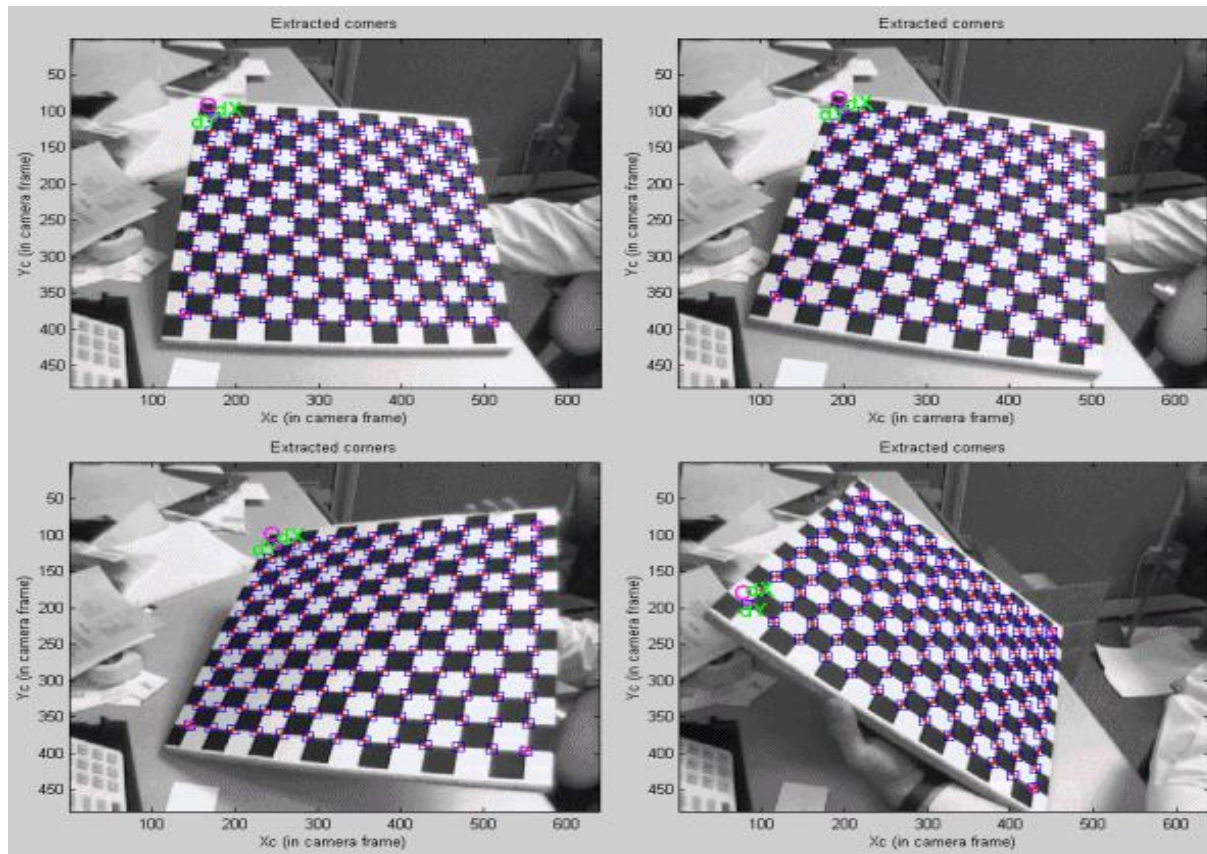


Camera Calibration

- Use camera model to interpret the projection from world to image plane
- Using known correspondences of $p \Leftrightarrow P$, we can compute the unknown parameters K , R , T by applying the perspective projection equation
- ... so associate known, physical distances in the world to pixel-distances in image

Projection Matrix

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \boxed{K[R|T]} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$



Camera Calibration (Direct Linear Transform (DLT) algorithm)

- We know that : $\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R|T] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$
- So there are 11 values to estimate:
(the overall scale doesn't matter, so
e.g. m_{34} could be set to 1)
- Each observed point gives us a pair of equations:

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$u_i = \frac{\lambda u_i}{\lambda} = \frac{m_{11}X_i + m_{12}Y_i + m_{13}Z_i + m_{14}}{m_{31} + m_{32} + m_{33} + m_{34}}$$

$$v_i = \frac{\lambda v_i}{\lambda} = \frac{m_{21}X_i + m_{22}Y_i + m_{23}Z_i + m_{24}}{m_{31} + m_{32} + m_{33} + m_{34}}$$

- To estimate 11 unknowns, we need **at least ?** points to calibrate the camera \Rightarrow solved using linear least squares

Camera Calibration (Direct Linear Transform (DLT) algorithm)

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} = K[R | T] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

- **what we obtained:** the 3x4 projection matrix,
what we need: its decomposition into the camera calibration matrix K , and the rotation R and position T of the camera.
- Use QR factorization to decompose the 3x3 submatrix ($m_{11:33}$) into the product of an upper triangular matrix K and a rotation matrix R (orthogonal matrix)

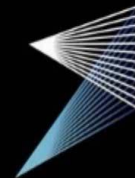
- The translation T can subsequently be obtained by:

$$T = K^{-1} \begin{bmatrix} m_{14} \\ m_{24} \\ m_{34} \end{bmatrix}$$

DLT algorithm applied to multi-robot mutual localization

A Monocular Pose Estimation System based on Infrared LEDs

Karl Schwabe, Matthias Faessler, Elias Mueggler
and Davide Scaramuzza



ROBOTICS &
PERCEPTION
GROUP

rpg.ifi.uzh.ch



University of
Zurich^{UZH}
Department of Informatics

robotics+ Swiss National
Centre of
Competence
in Research

In this case, the camera has been pre-calibrated (i.e., K is known). Can you think of how the DLT algorithm could be modified so that only R and T need to be determined and not K ?

Outline of this lecture

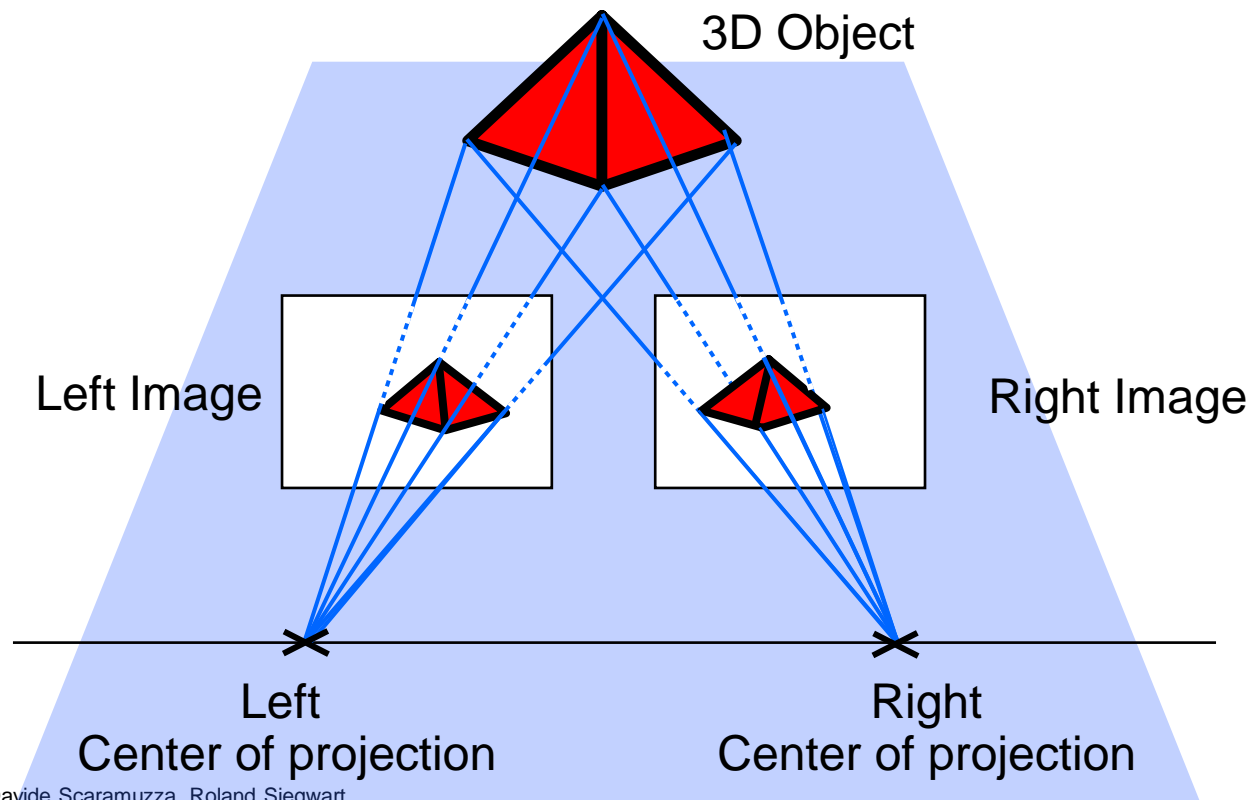
- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision

Stereo Vision versus Structure from Motion

- **Stereo vision:**
is the process of obtaining **depth information** from a pair of images coming from two cameras that look at the same scene from different but **known** positions
- **Structure from Motion:**
is the process of obtaining **depth and motion information** from a pair of images coming from the same camera that looks at the same scene from different positions

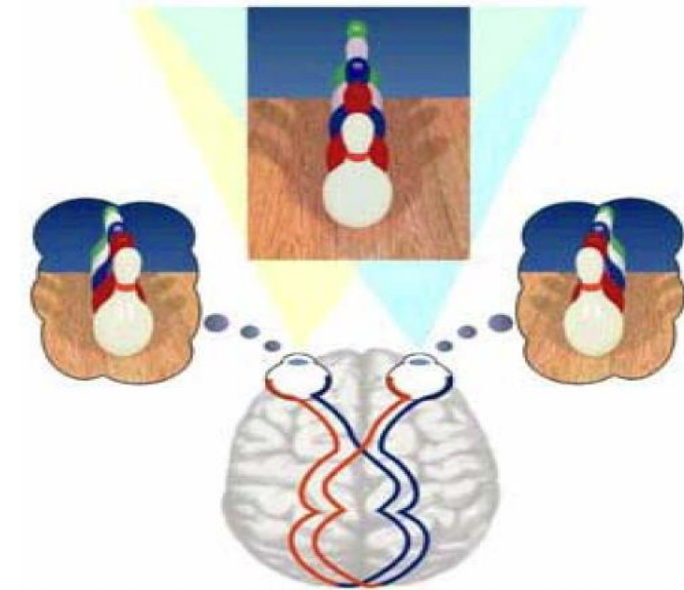
Depth from Stereo

- From a single camera, we can only deduct the **ray** on which each image point lies
- With a stereo camera (binocular), we can solve for the intersection of the rays and recover the 3D structure



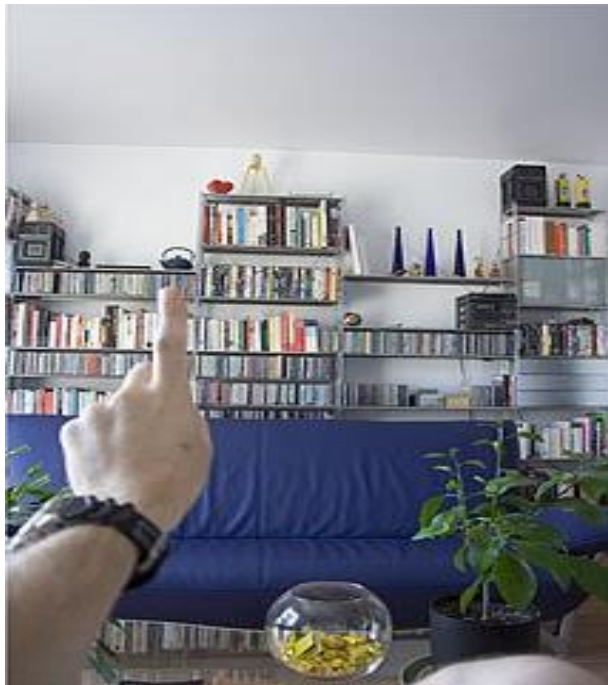
The “human” binocular system

- **Stereopsis:** the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial distortion is also removed. This process is called «**rectification**»



The “human” binocular system

- **Stereopsis:** the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial distortion is also removed. This process is called «**rectification**»



Make a simple test:

1. Fix an object
2. Open and close alternatively the left and right eyes.
 - The horizontal displacement is called **disparity**
 - The smaller the disparity, the farther the object

The “human” binocular system

- **Stereopsis:** the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial distortion is also removed. This process is called «**rectification**»



Make a simple test:

1. Fix an object
2. Open and close alternatively the left and right eyes.
 - The horizontal displacement is called **disparity**
 - The smaller the disparity, the farther the object

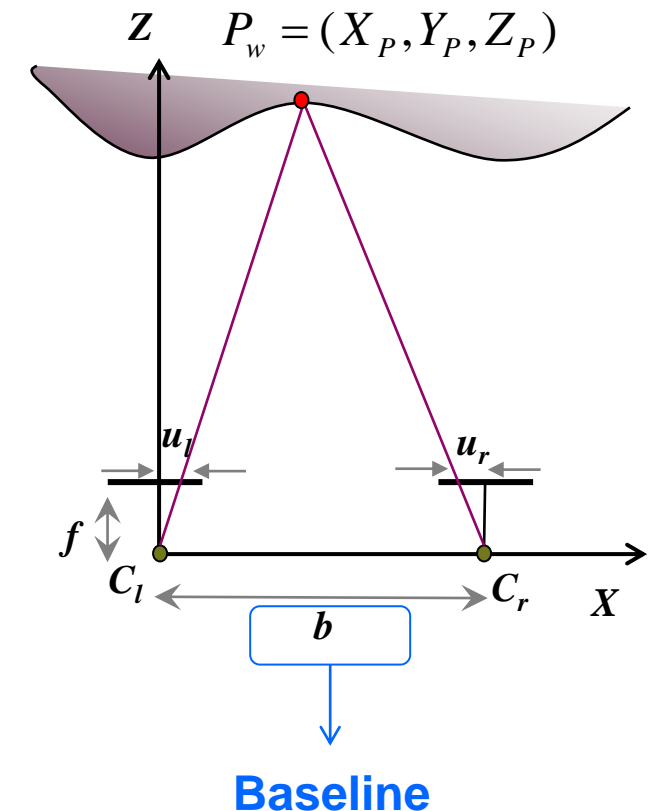
Stereo Vision | simplified case

- An ideal, simplified case assumes that both cameras are **identical** and **aligned** with the x-axis
- Can we find an expression for the depth Z_P of point P_W ?
- From similar triangles:

$$\frac{f}{Z_P} = \frac{u_l}{X_P} \quad \Rightarrow \quad Z_P = \frac{bf}{u_l - u_r}$$

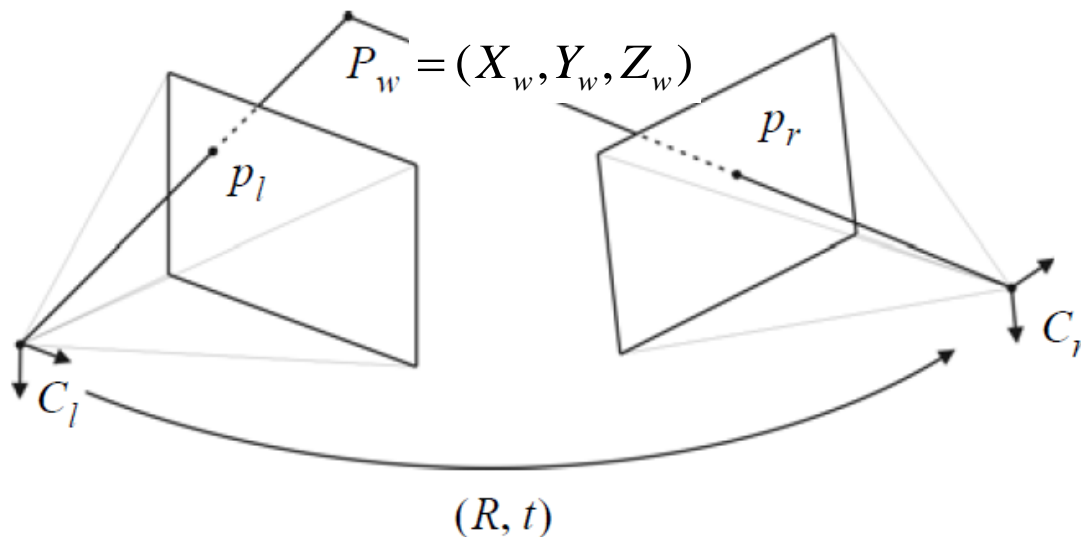
\downarrow
Disparity

- **Disparity** is the difference in image location of the projection of a 3D point in two image planes
- **Baseline** is the distance between the two cameras



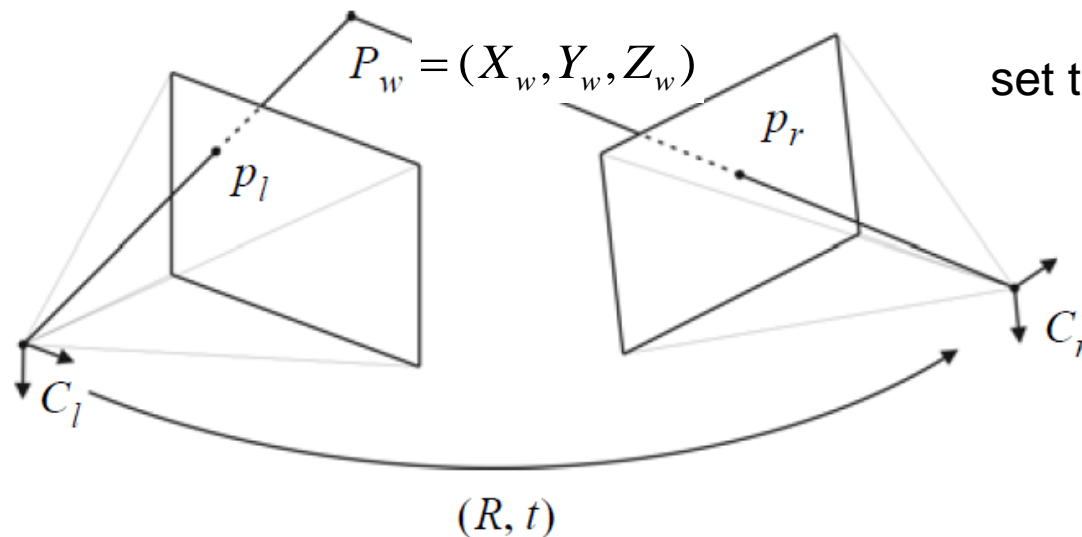
Stereo Vision | general case

- Two identical cameras do not exist in nature!
- Aligning both cameras on a horizontal axis is very difficult
- In order to use a stereo camera, we need to know the intrinsic extrinsic parameters of each camera, that is, the relative pose between the cameras (rotation, translation) \Rightarrow We can solve for this through camera calibration



Stereo Vision | general case

- To estimate the 3D position of P_W we can construct the system of equations of the left and right camera
- Triangulation is the problem of determining the 3D position of a point given a set of corresponding image locations and known camera poses.



Left camera:
set the world frame to coincide
with the left camera frame

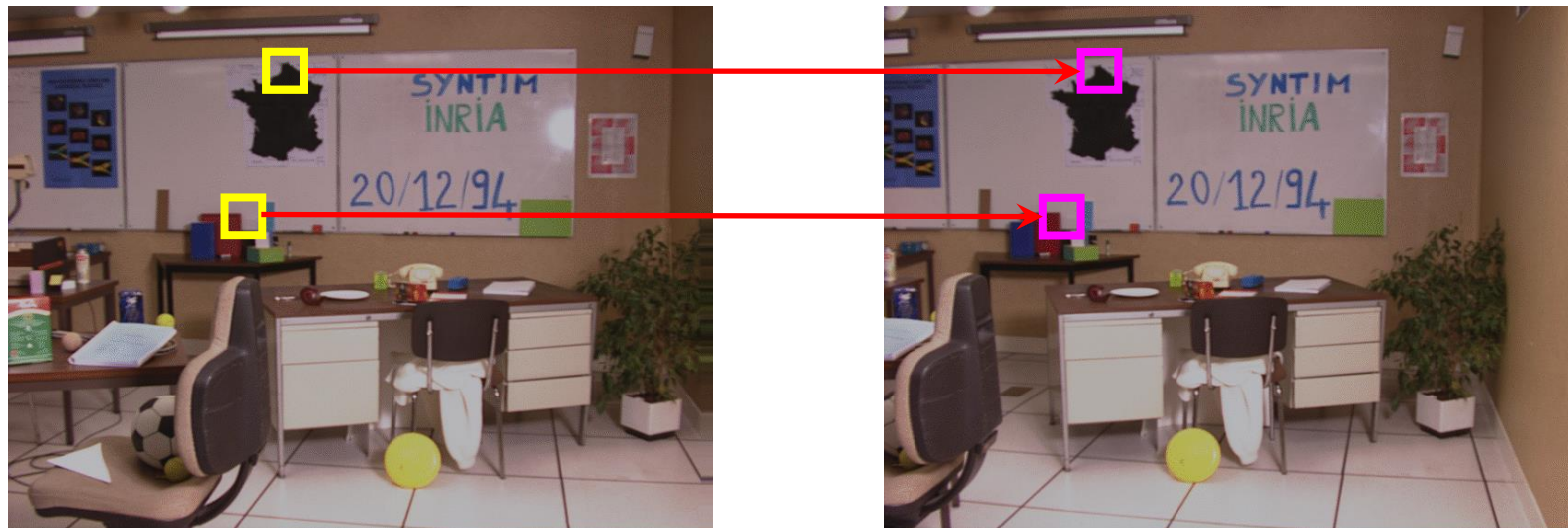
$$\tilde{p}_l = \lambda_l \begin{bmatrix} u_l \\ v_l \\ 1 \end{bmatrix} = K_l \begin{bmatrix} X_w \\ Y_w \\ Z_w \end{bmatrix}$$

Right camera:

$$\tilde{p}_r = \lambda_r \begin{bmatrix} u_r \\ v_r \\ 1 \end{bmatrix} = K_r R \begin{bmatrix} X_w \\ Y_w \\ Z_w \end{bmatrix} + T$$

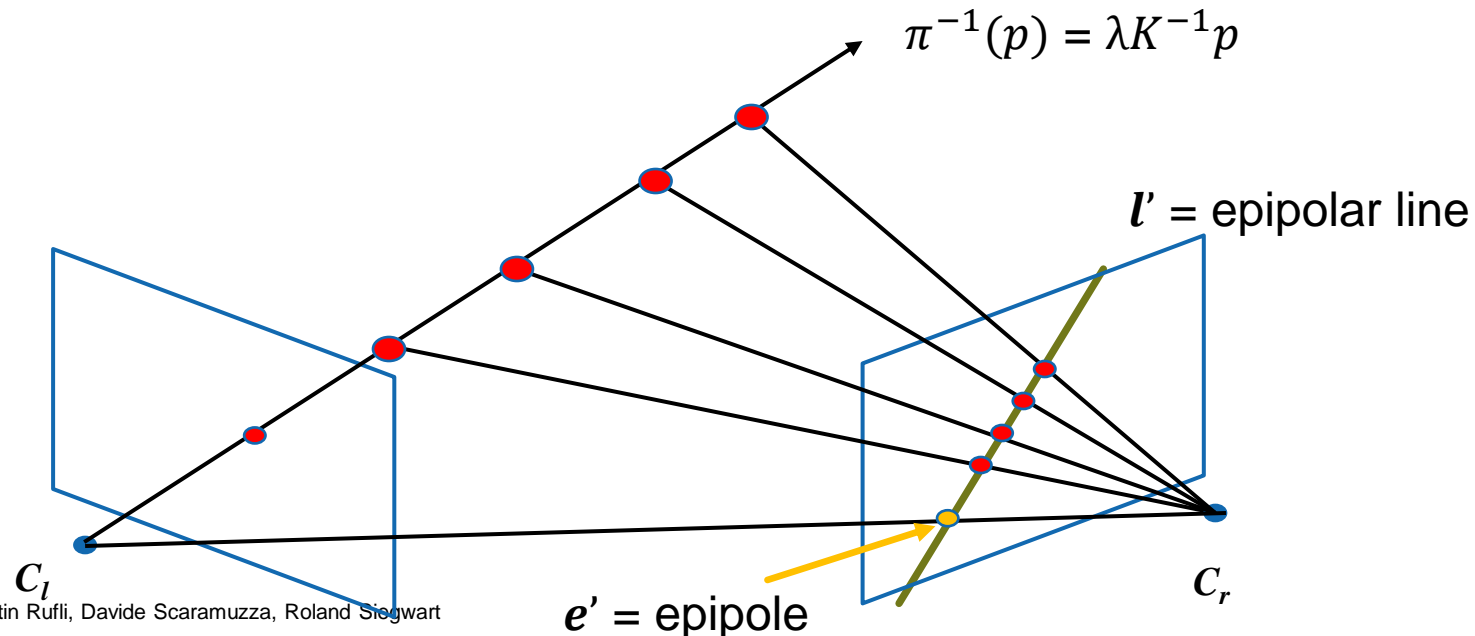
Correspondence Search | the problem

- Goal: identify corresponding points in the left and right images, which are the reprojection of the same 3D scene point
 - Typical similarity measures: Normalized Cross-Correlation (NCC) , Sum of Squared Differences (SSD), Sum of Absolute Differences (SAD), Census Transform
 - Exhaustive image search can be computationally very expensive! Can we make the correspondence search in 1D?



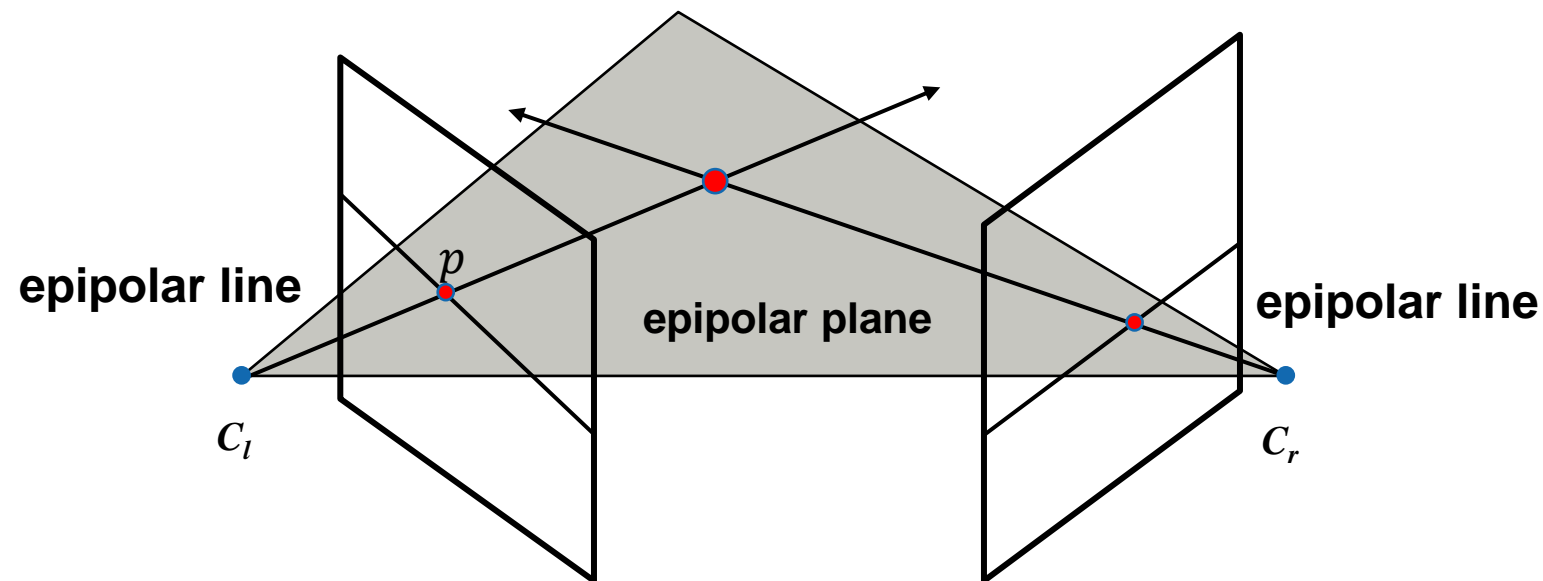
Correspondence Problem

- **Exhaustive** image search can be computationally very expensive!
- Can we make the correspondence search in 1D?
- Potential matches for p have to lie on the corresponding epipolar line l'
 - The **epipolar line** is the projection of the infinite ray $\pi^{-1}(p)$ corresponding to p in the other camera image
 - The **epipole** e' is the projection of the optical center in in the other camera image



Correspondence Search | the epipolar constraint

- The epipolar plane is defined by the image point p and the optical centers
- Impose the epipolar constraint to aid matching: search for a correspondence along the epipolar line



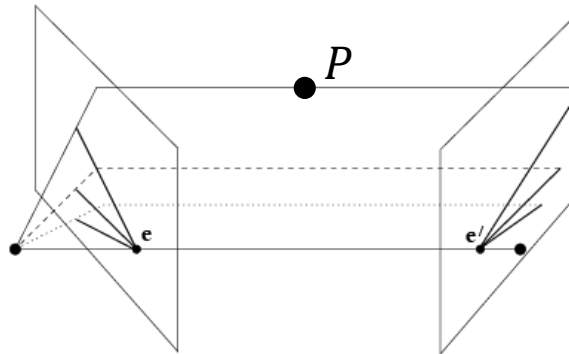
Correspondence Search | the epipolar constraint

- Thanks to the epipolar constraint, corresponding points can be searched for, along epipolar lines \Rightarrow computational cost reduced to 1 dimension!

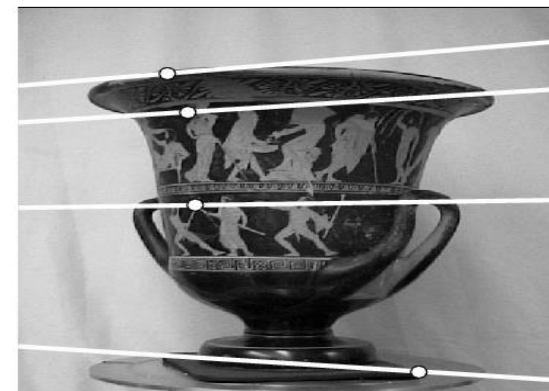


Example: converging cameras

- **Remember:** all the epipolar lines intersect at the epipole
- As the position of the 3D point varies, the epipolar lines “rotate” about the baseline

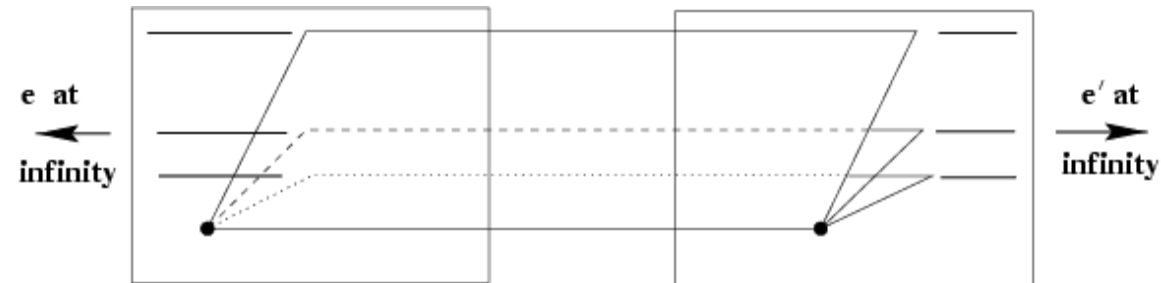


Left image



Right image

Example: horizontally aligned cameras



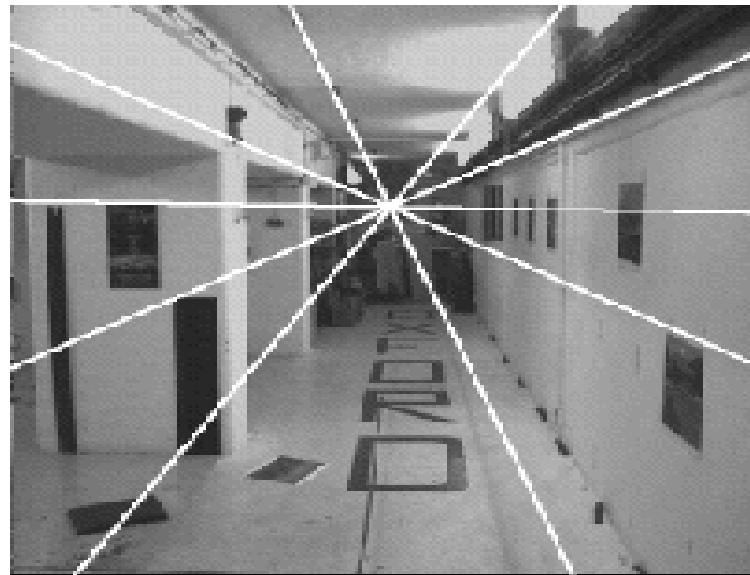
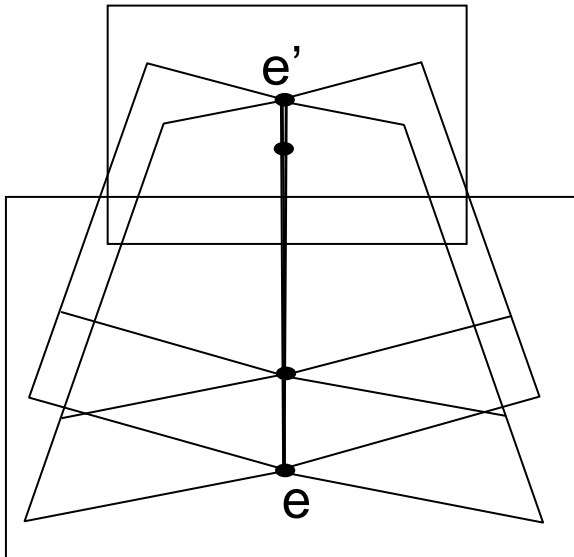
Left image



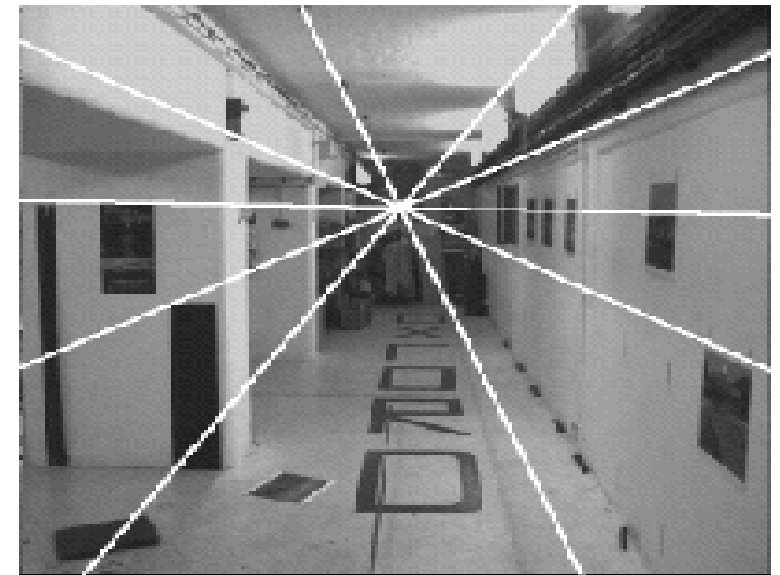
Right image

Example: forward motion (parallel to the optical axis)

- Epipole has the **same coordinates** in both images
- Points move along lines radiating from e: “Focus of expansion”



Left image



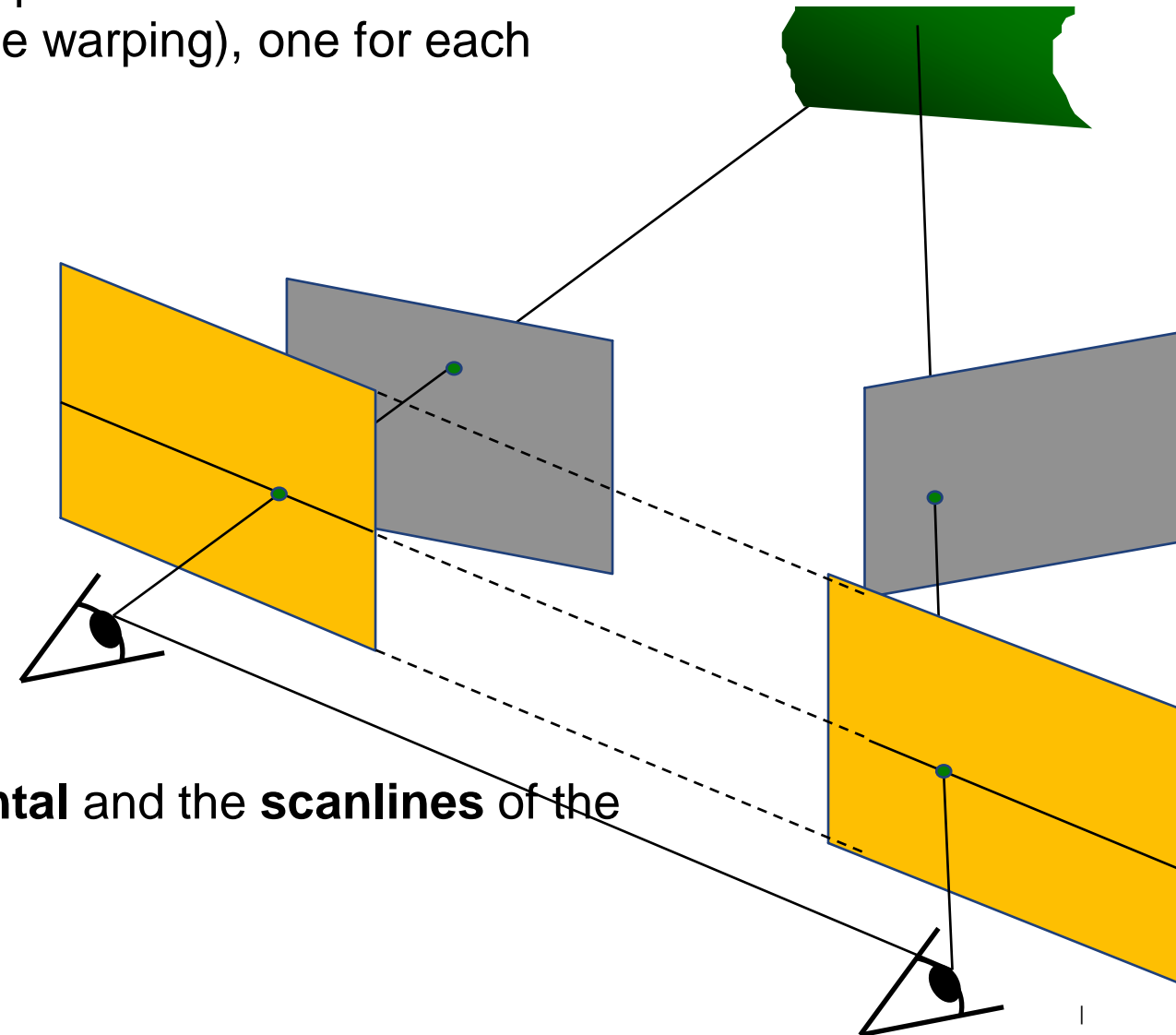
Right image

Stereo Rectification

- Even in commercial stereo cameras the left and right image are never perfectly aligned
- In practice, it is convenient if image scanlines are the epipolar lines
- Stereo rectification warps the left and right images into new “rectified” images, whose epipolar lines are aligned to the baseline

Stereo Rectification

- Reprojects image planes onto a common plane parallel to the baseline
- It works by computing two homographies (image warping), one for each input image reprojection



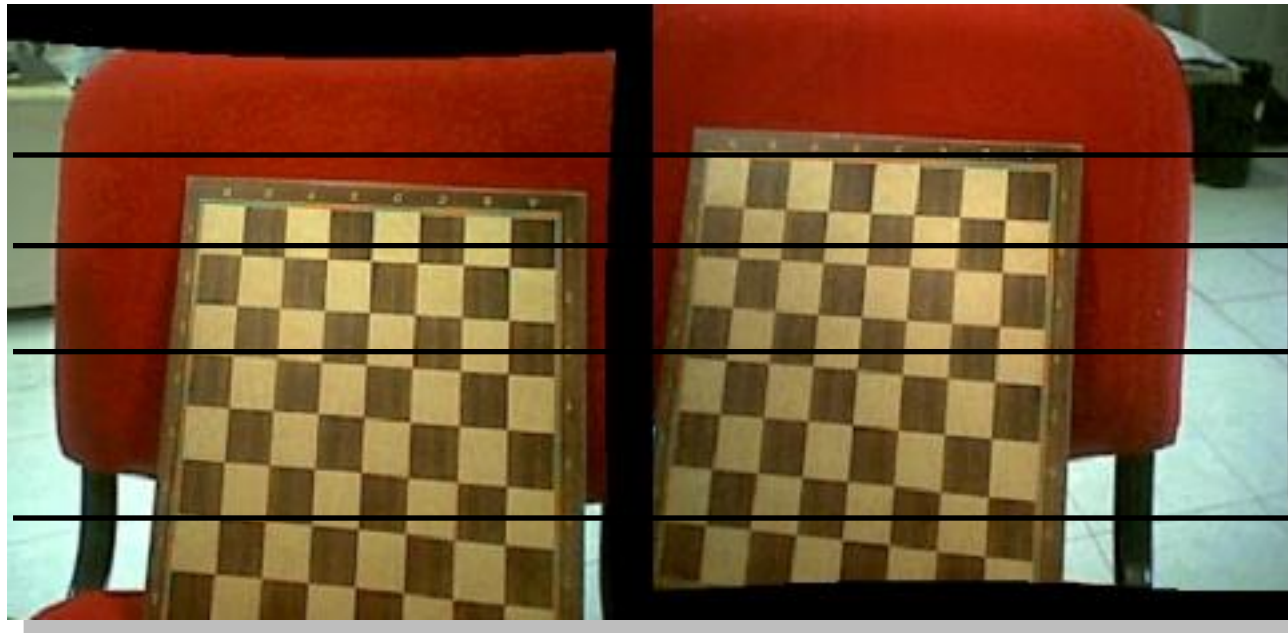
- As a result, the new **epipolar lines** are **horizontal** and the **scanlines** of the left and right image **are aligned**

Epipolar Rectification - Example

- First, remove radial distortion

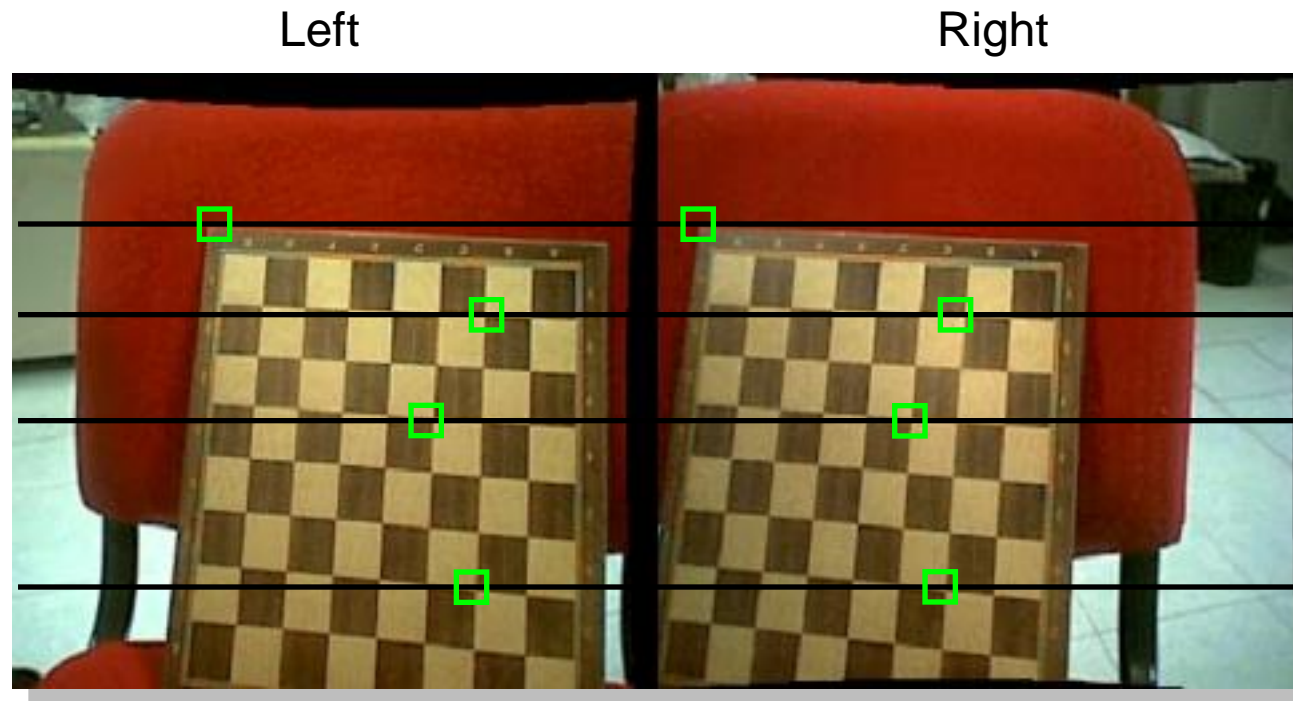
Left

Right

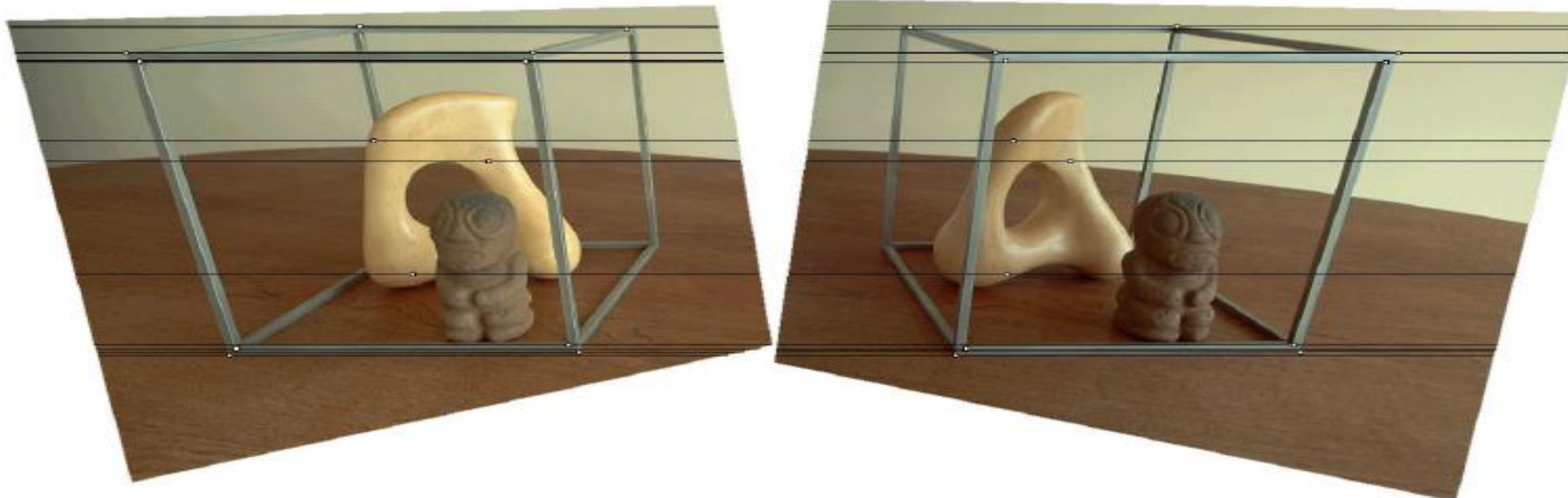


Epipolar Rectification - Example

- First, remove radial distortion
- Then, compute homographies (warping) and rectify



Stereo Rectification: example



Stereo Vision | disparity map

- The disparity map holds the disparity value at every pixel:
 - Identify correspondent points of all image pixels in the original images
 - Compute the disparity ($u_l - u_r$) for each pair of correspondences
- Usually visualized in gray-scale images
- Close objects experience bigger disparity; thus, they appear brighter in disparity map



Left image



Right image



Disparity Map

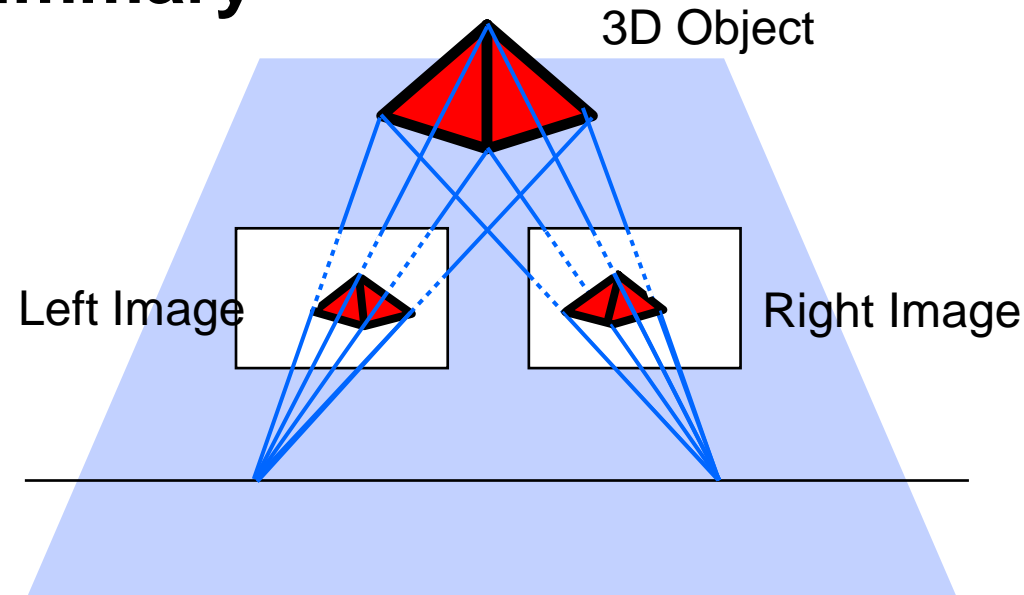
Stereo Vision | disparity map

- The disparity map holds the disparity value at every pixel:
 - Identify correspondent points of all image pixels in the original images
 - Compute the disparity ($u_l - u_r$) for each pair of correspondences
- Usually visualized in gray-scale images
- Close objects experience bigger disparity; thus, they appear brighter in disparity map
- From the disparity, we can compute the depth Z as:

$$Z = \frac{bf}{u_l - u_r}$$



Stereo Vision - summary



1. Stereo camera calibration \Rightarrow compute camera relative pose
2. Epipolar rectification \Rightarrow align images & epipolar lines
3. Search for correspondences
4. Output: compute stereo triangulation or disparity map

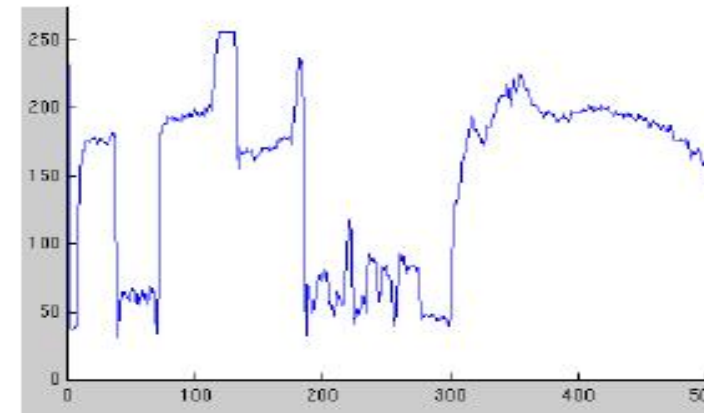
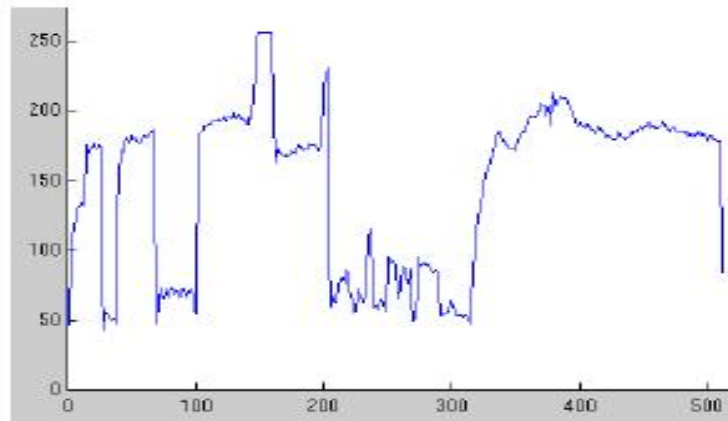
Correspondence problem

- Now that the left and right images are rectified, the correspondence search can be done along the same scanlines



Correspondence problem

- If we look at the intensity profiles of two corresponding scanlines, there is a clear correspondence between intensities but also noise and ambiguities

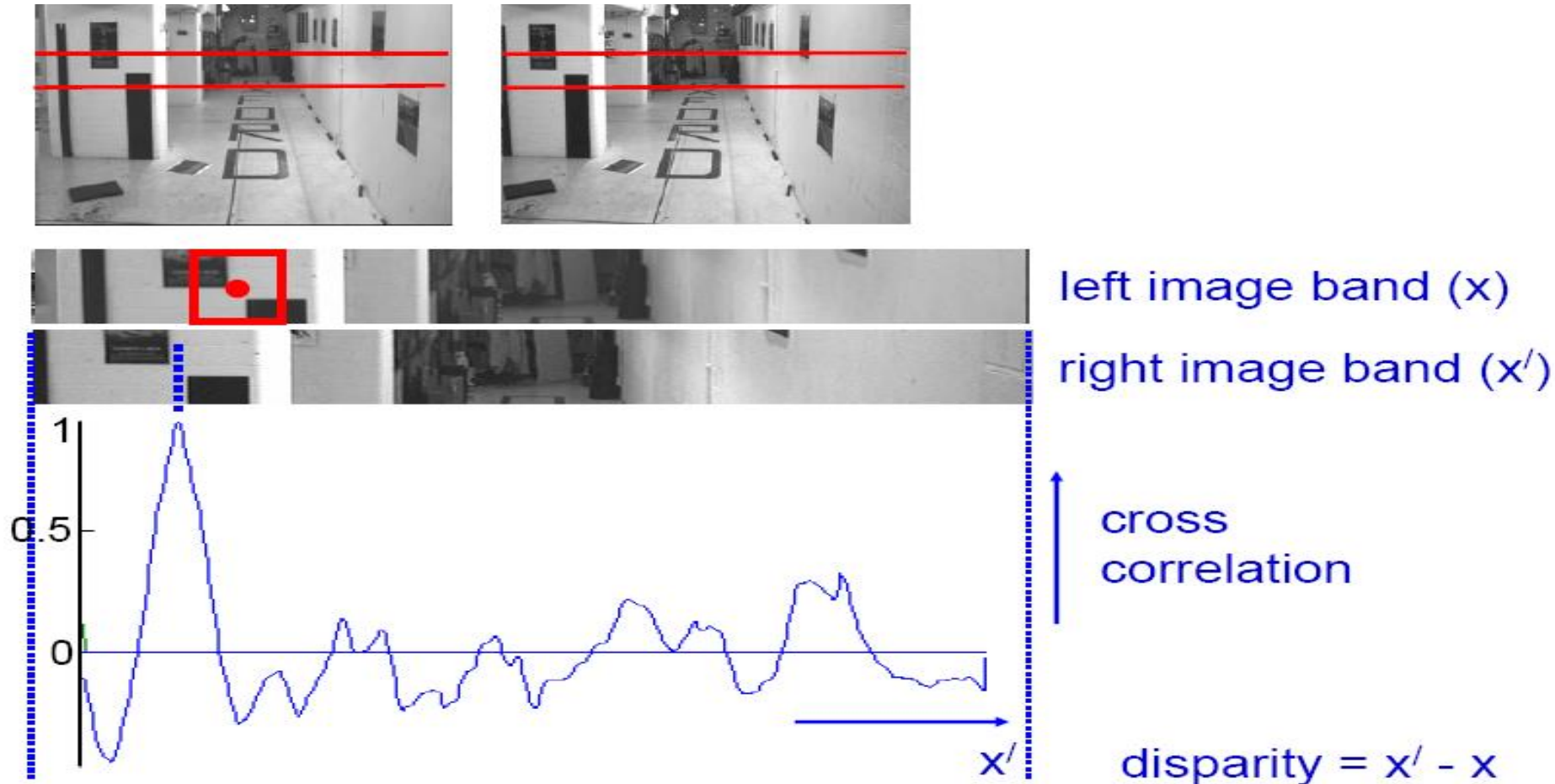


Correspondence problem

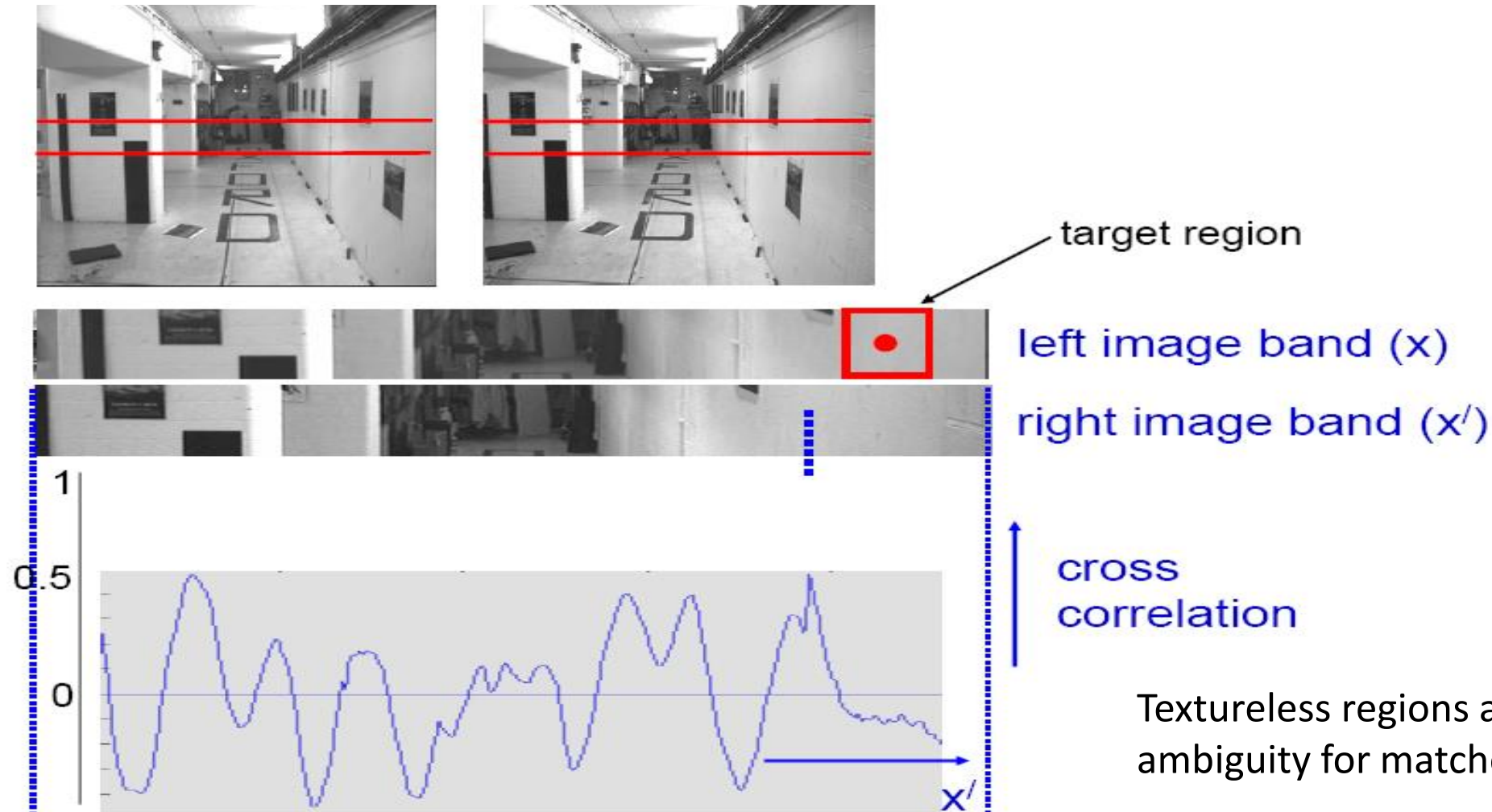
- To average noise effects, use a window around the point of interest
- Neighborhood of corresponding points are similar in intensity patterns
- **Similarity measures:**
 - Zero-Normalized Cross-Correlation (**ZNCC**)
 - Sum of Squared Differences (**SSD**),
 - Sum of Squared Differences (**SAD**)
 - **Census Transform** (Census descriptor plus Hamming distance)



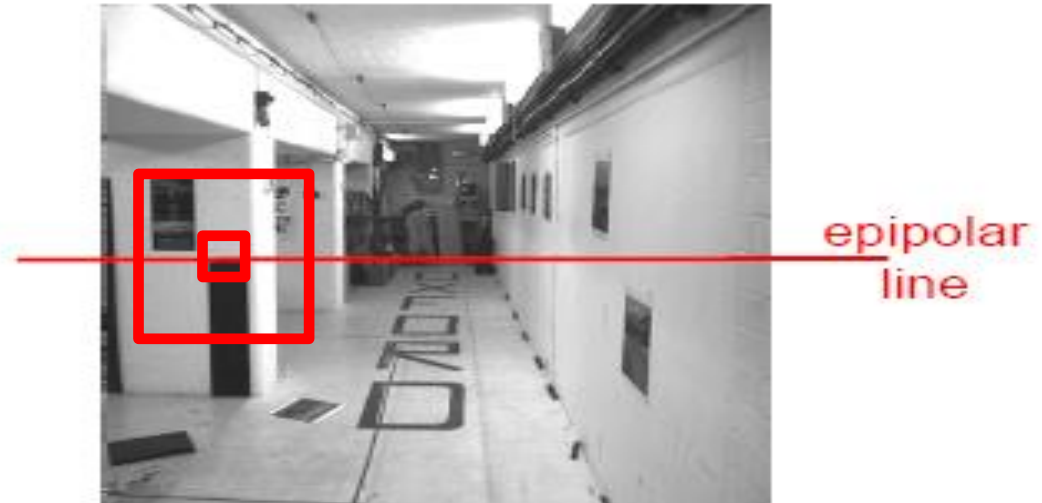
Correlation-based window matching



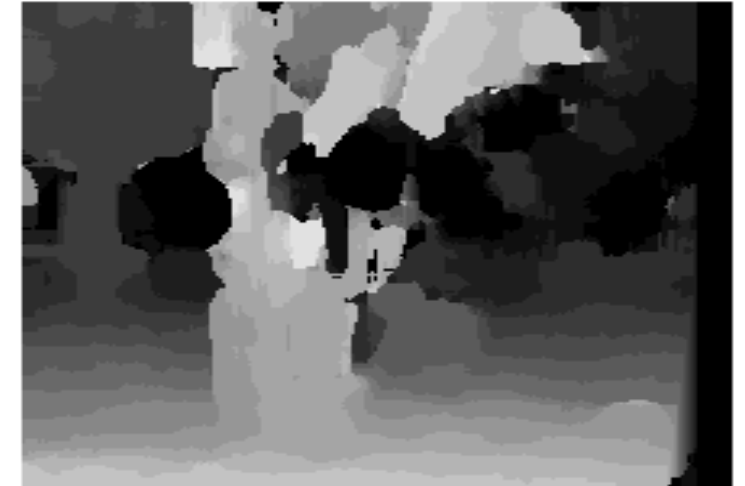
Correspondence Problems: Textureless regions (the aperture problem)



Solution: increase window size

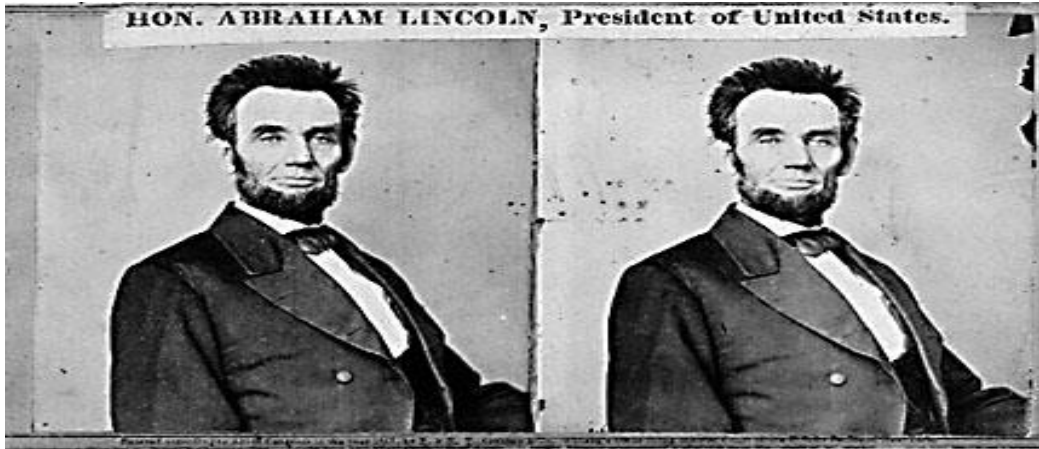


Effects of window size W

 $W = 3$  $W = 20$

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

Failures of correspondence search



Textureless surfaces



Occlusions, repetition

How can we improve window-based matching?

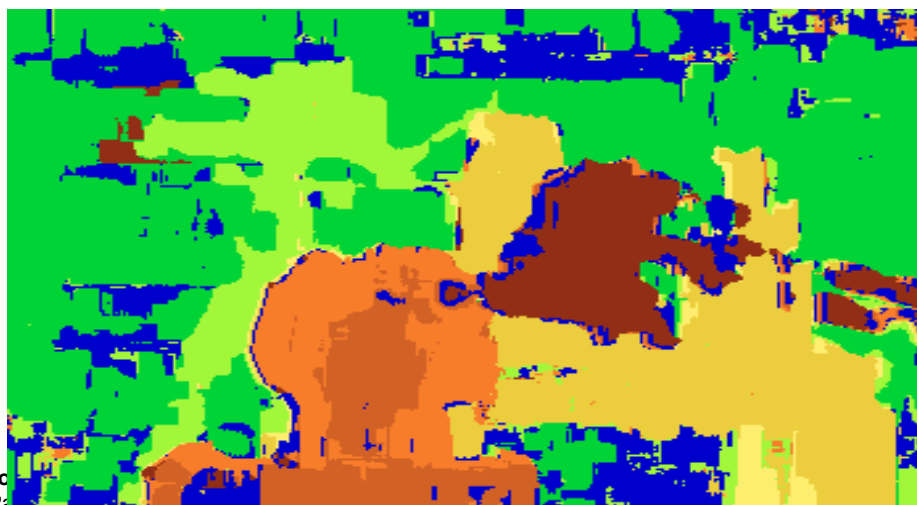
- Beyond the epipolar constraint, there are “soft” constraints to help identify corresponding points
 - Uniqueness
 - Only one match in right image for every point in left image
 - Ordering
 - Points on **same surface** will be in same order in both views
 - Disparity gradient
 - Disparity changes smoothly between points on the same surface

Results with window search

Data



Window-based matching



Ground truth



Better methods exist...



Graph cuts



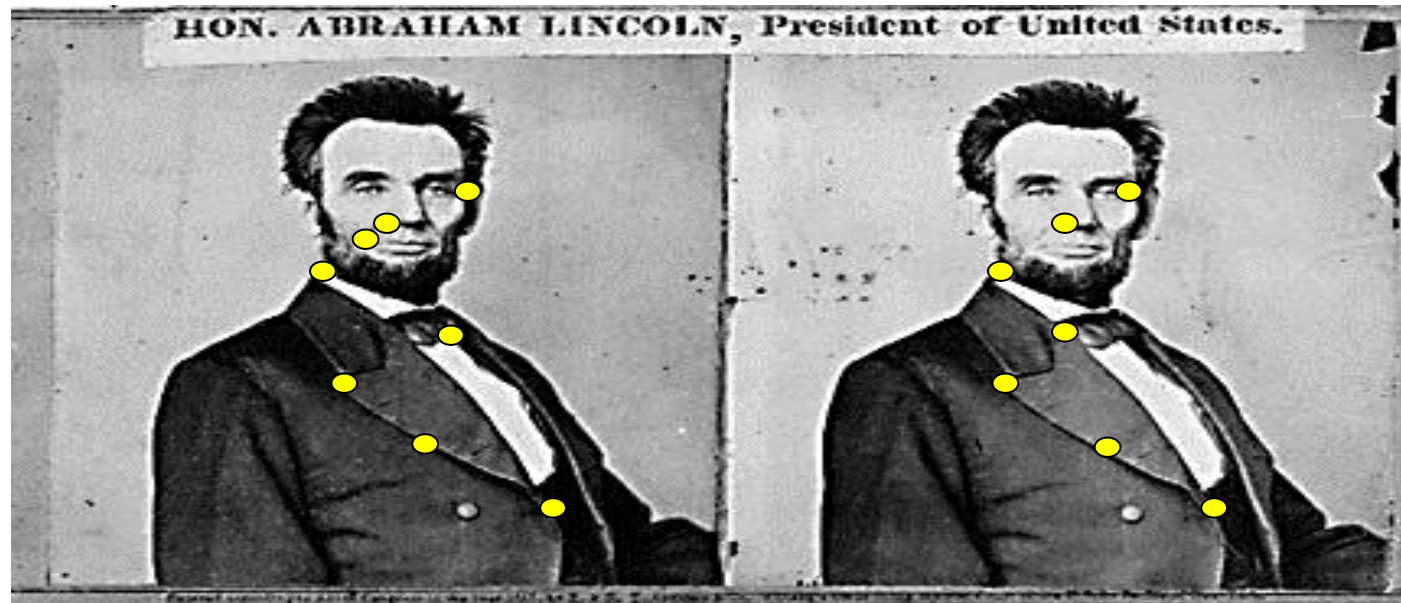
Ground truth

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

For code, datasets, and comparisons all the algorithms: <http://vision.middlebury.edu/stereo/>

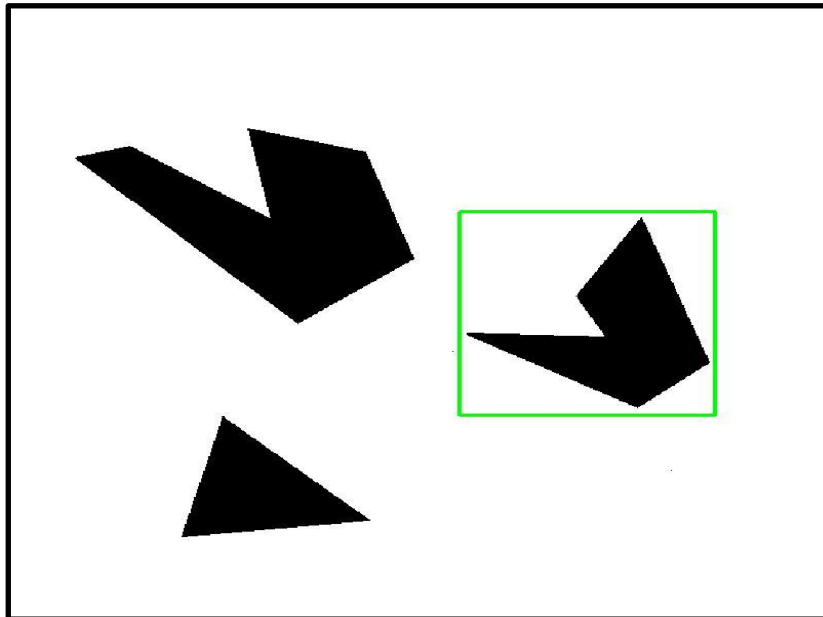
Sparse correspondence search

- Restrict search to sparse set of detected features
- Rather than pixel values (or lists of pixel values) use *feature descriptor* and an associated *similarity metrics*
- Still use epipolar geometry to narrow the search further

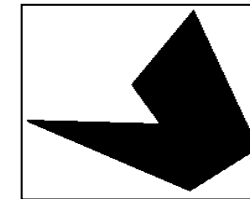


Template matching

- Find locations in an image that are similar to a **template**
- If we look at filters as **templates**, we can use correlation to detect these locations



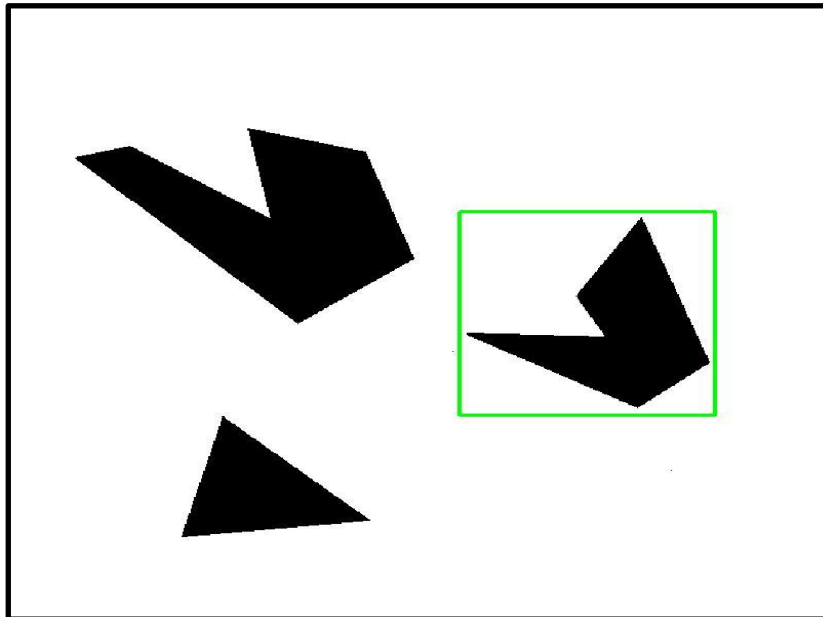
Detected template



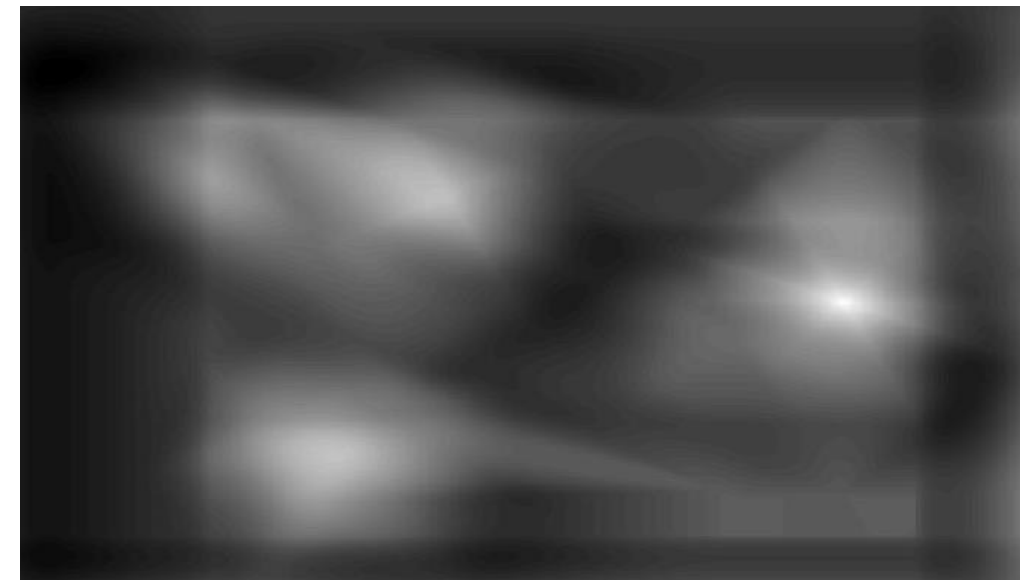
Template

Template matching

- Find locations in an image that are similar to a **template**
- If we look at filters as **templates**, we can use correlation to detect these locations



Detected template



Correlation map

Similarity measures

- Sum of Squared Differences (**SSD**)

$$SSD = \sum_{u=-k}^k \sum_{v=-k}^k (H(u, v) - F(u, v))^2$$

- Sum of Absolute Differences (**SAD**) (used in optical mice)

$$SAD = \sum_{u=-k}^k \sum_{v=-k}^k |H(u, v) - F(u, v)|$$

Similarity measures

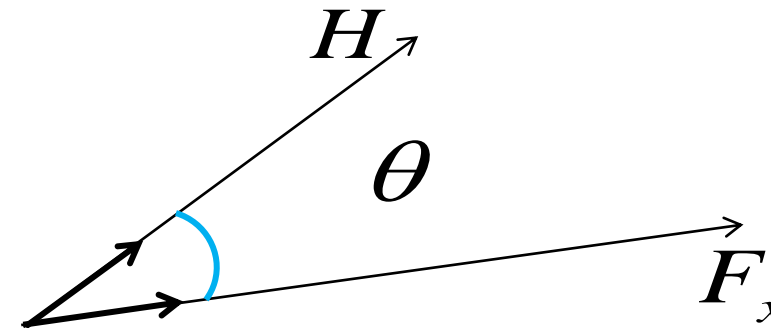
- For ***slight*** invariance to intensity changes, the Zero-mean Normalized Cross Correlation (**ZNCC**) is widely used

$$ZNCC = \frac{\sum_{u=-k}^k \sum_{v=-k}^k (H(u, v) - \mu_H)(F(u, v) - \mu_F)}{\sqrt{\sum_{u=-k}^k \sum_{v=-k}^k (H(u, v) - \mu_H)^2} \sqrt{\sum_{u=-k}^k \sum_{v=-k}^k (F(u, v) - \mu_F)^2}} \left\{ \begin{array}{l} \mu_H = \frac{\sum_{u=-k}^k \sum_{v=-k}^k H(u, v)}{(2N+1)^2} \\ \mu_F = \frac{\sum_{u=-k}^k \sum_{v=-k}^k F(u, v)}{(2N+1)^2} \end{array} \right.$$

Correlation as an inner product

- Considering the filter H and the portion of the image F_x as vectors \Rightarrow their correlation is:

$$\langle H, F_x \rangle = \|H\| \|F_x\| \cos \theta$$

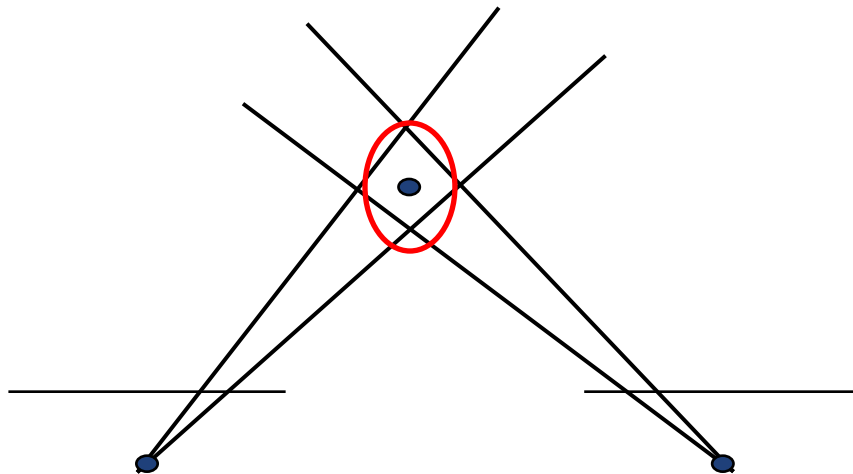


- In **ZNCC** we consider the unit vectors of H and F_x , hence we measure their similarity based on the angle θ . Alternatively, ZNCC maximizes $\cos \theta$

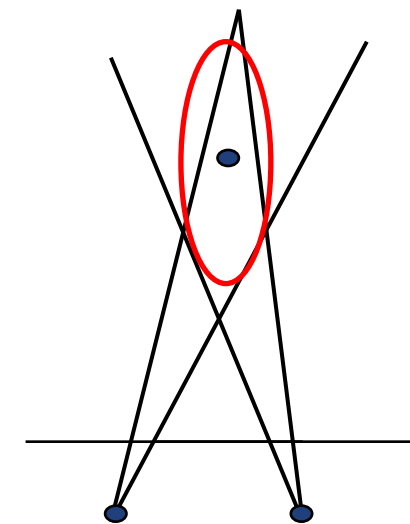
$$\cos \theta = \frac{\langle H, F_x \rangle}{\|H\| \|F_x\|} = \frac{\sum_{u=-k}^k \sum_{v=-k}^k (H(u, v) - \mu_H)(F(u, v) - \mu_F)}{\sqrt{\sum_{u=-k}^k \sum_{v=-k}^k (H(u, v) - \mu_H)^2} \sqrt{\sum_{u=-k}^k \sum_{v=-k}^k (F(u, v) - \mu_F)^2}}$$

Choosing the Baseline

- What's the optimal baseline?
 - **Too small:**
 - Large depth error
 - Can you quantify the error as a function of the disparity?
 - **Too large:**
 - Minimum measurable distance increases
 - Difficult search problem for close objects

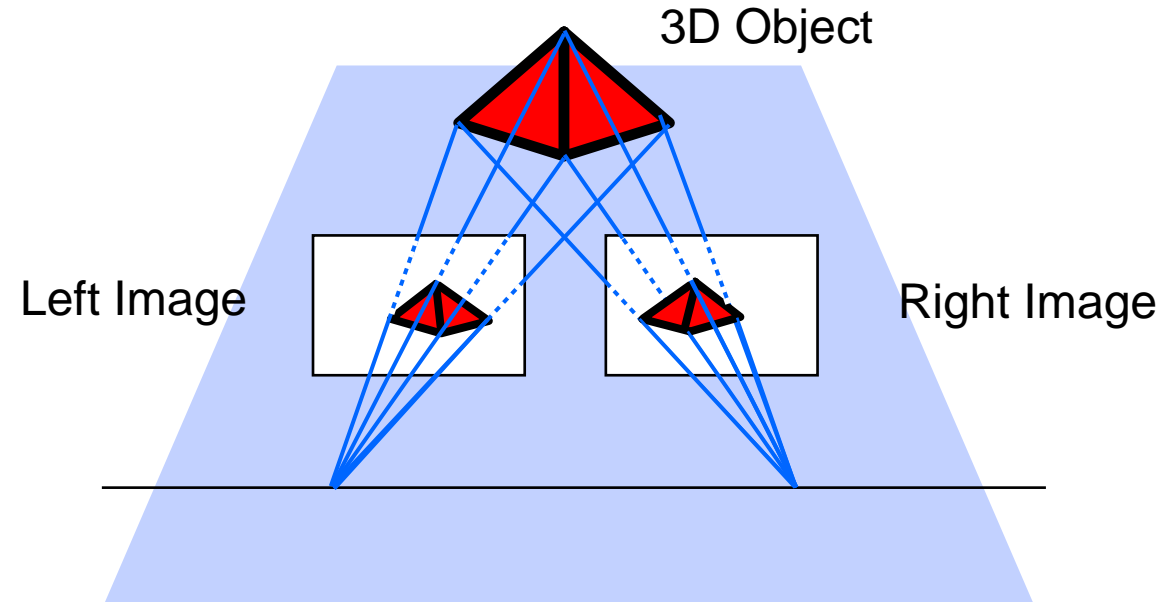


Large Baseline



Small Baseline

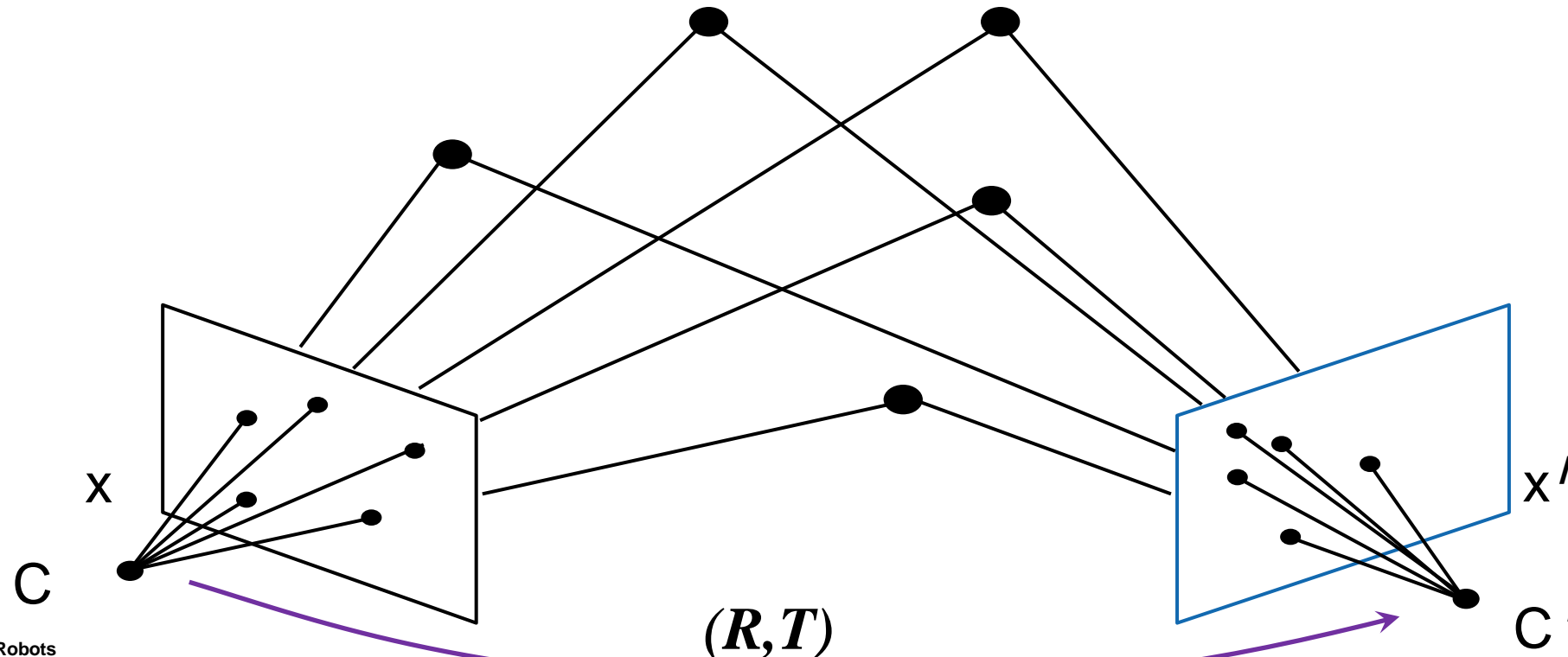
Stereo Vision - summary



1. Stereo camera calibration \Rightarrow compute camera relative pose
2. Epipolar rectification \Rightarrow align images & epipolar lines
3. Search for correspondences
4. Output: compute stereo triangulation or disparity map
5. Consider how baseline & image resolution affect accuracy of depth estimates

SFM: Structure From Motion (watch video segment)

- Given image point correspondences, $\mathbf{x}_i \leftrightarrow \mathbf{x}_i'$, determine \mathbf{R} and \mathbf{T}
- Keep track of point trajectories over multiple views to reconstruct scene structure and motion



Multiple-view structure from motion

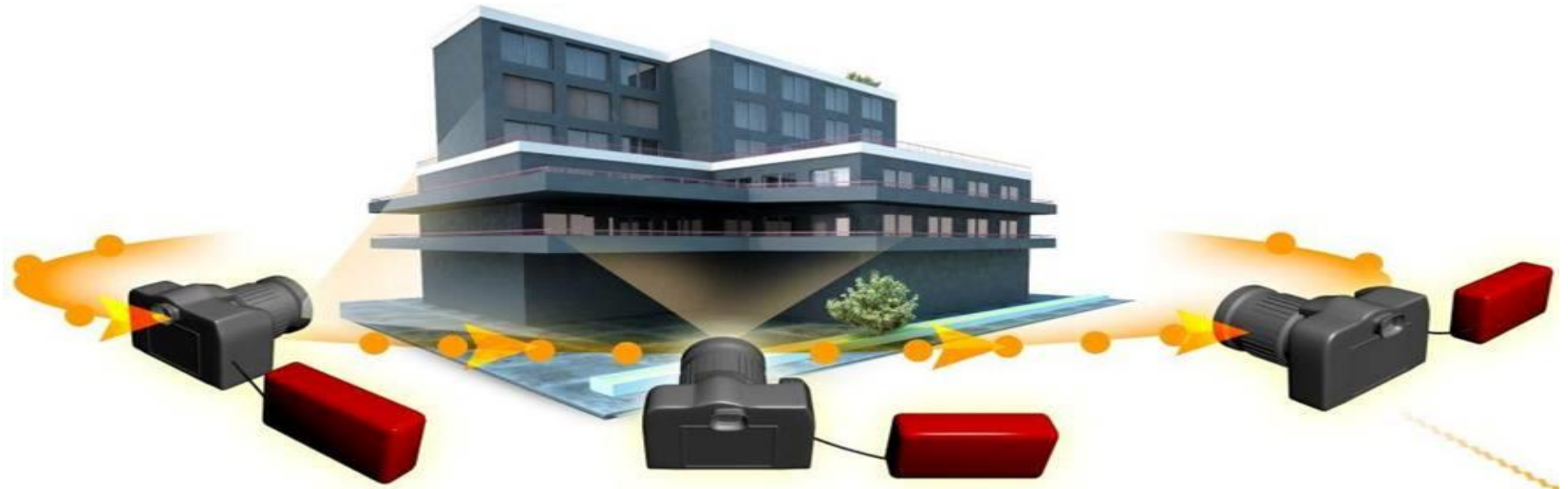


Image courtesy of Nader Salman

Multiple-view structure from motion

- Results of Structure from motion from user images from flickr.com

[Seitz, Szeliski ICCV 2009]



Colosseum, Rome

2,106 images, 819,242 points



San Marco square, Venice

14,079 images, 4,515,157 points