

CSE 152: Computer Vision

Hao Su

Lecture 19: Final Review



Coverage

- Anything taught this quarter may appear
- >90% on materials from L12 (3D Deep Learning)

Form

- 3 big problems that need calculation
 - Understand your homework well
- A few short Q&A

CSE 152: Computer Vision

Hao Su

Lecture 12: 3D Deep Learning



Credit: Stanford CS231n, L13

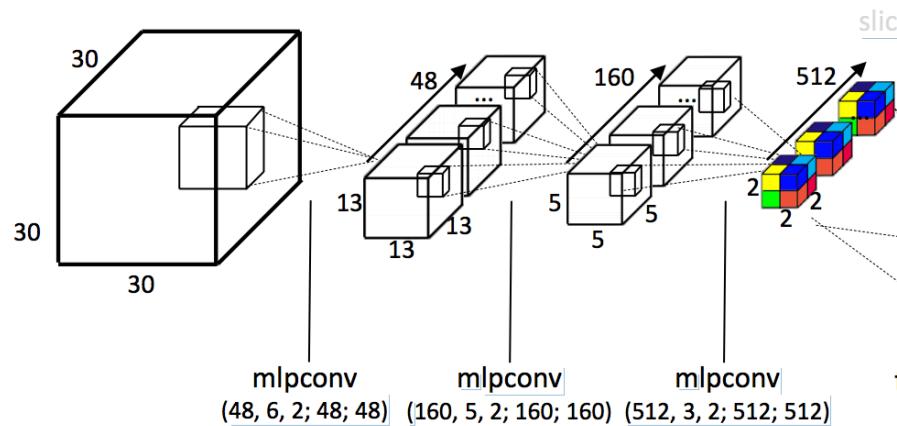
Voxelization

Represent the occupancy of regular 3D grids

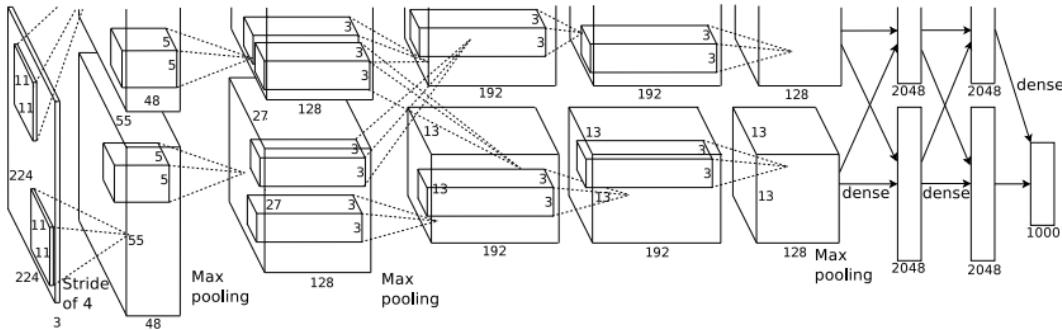


3D CNN on Volumetric Data

3D convolution uses 4D kernels



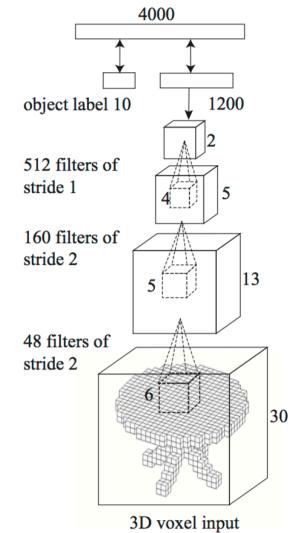
Complexity Issue



AlexNet, 2012

Input resolution: 224x224

$$224 \times 224 = 50176$$



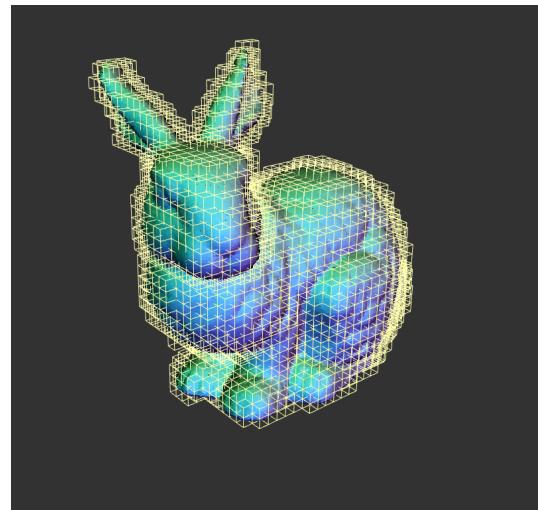
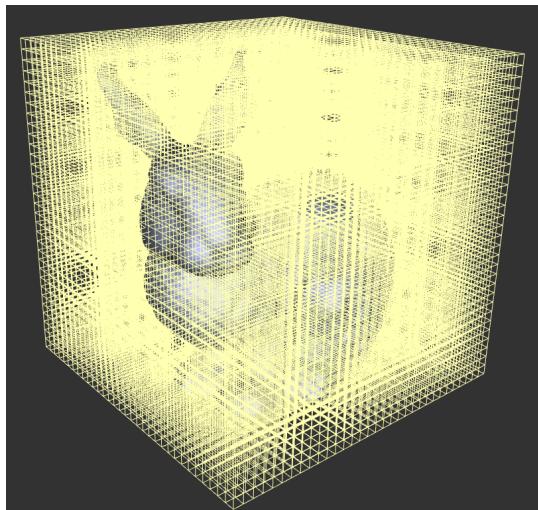
3DShapeNets,
2015

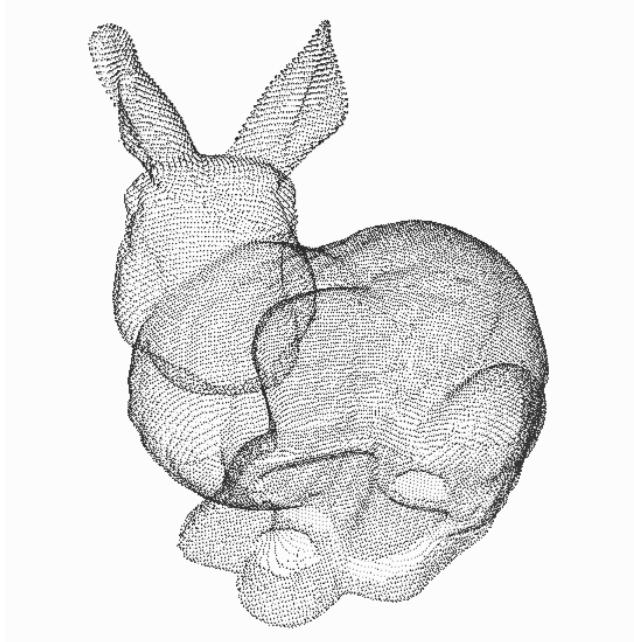
Input resolution: 30x30x30

$$224 \times 224 = 27000$$

Store only the Occupied Grids

- Store the sparse surface signals
- Constrain the computation near the surface

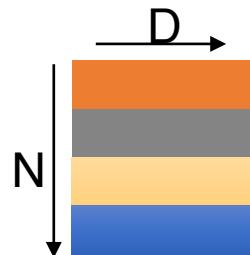




Point cloud
(The most common 3D sensor data)

Properties of a Desired Point Network

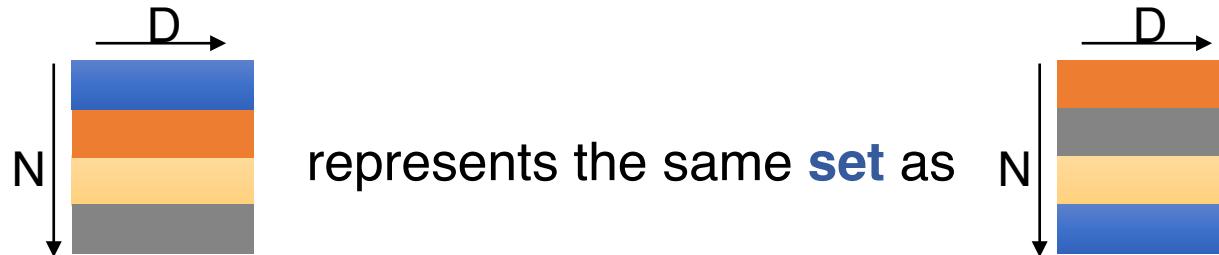
Point cloud: N **orderless** points, each represented by a D dim coordinate



2D array representation

Permutation Invariance

Point cloud: N **orderless** points, each represented by a D dim coordinate

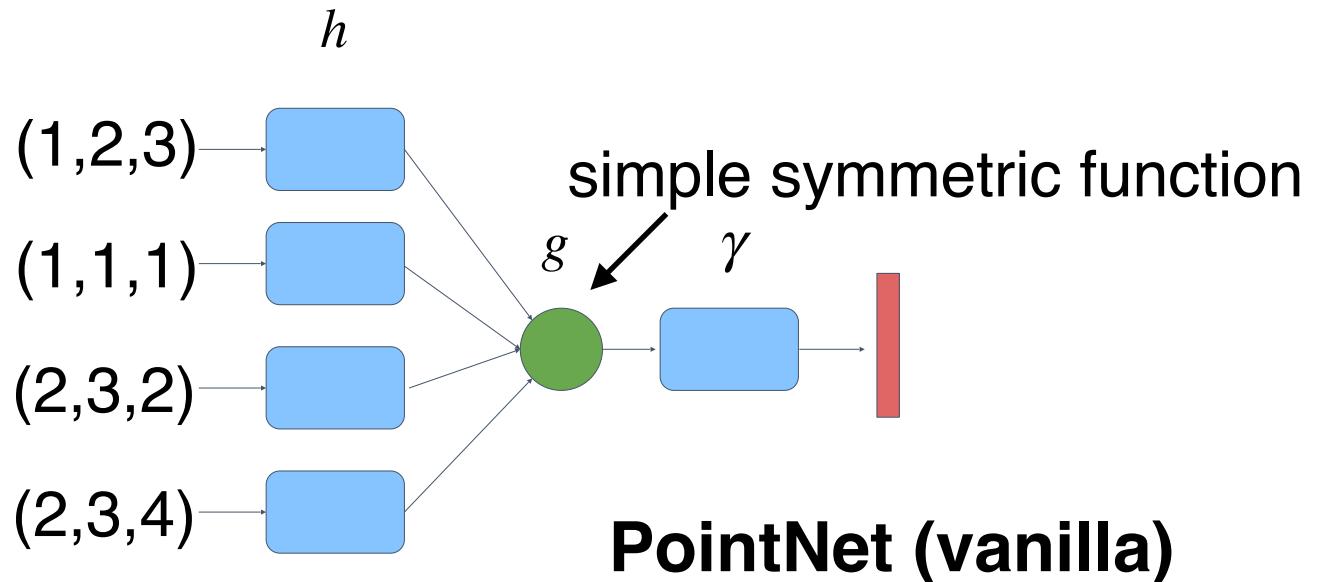


2D array representation

Construct a Symmetric Function

Observe:

$f(x_1, x_2, \dots, x_n) = \gamma \circ g(h(x_1), \dots, h(x_n))$ is symmetric if g is symmetric



CSE 152: Computer Vision

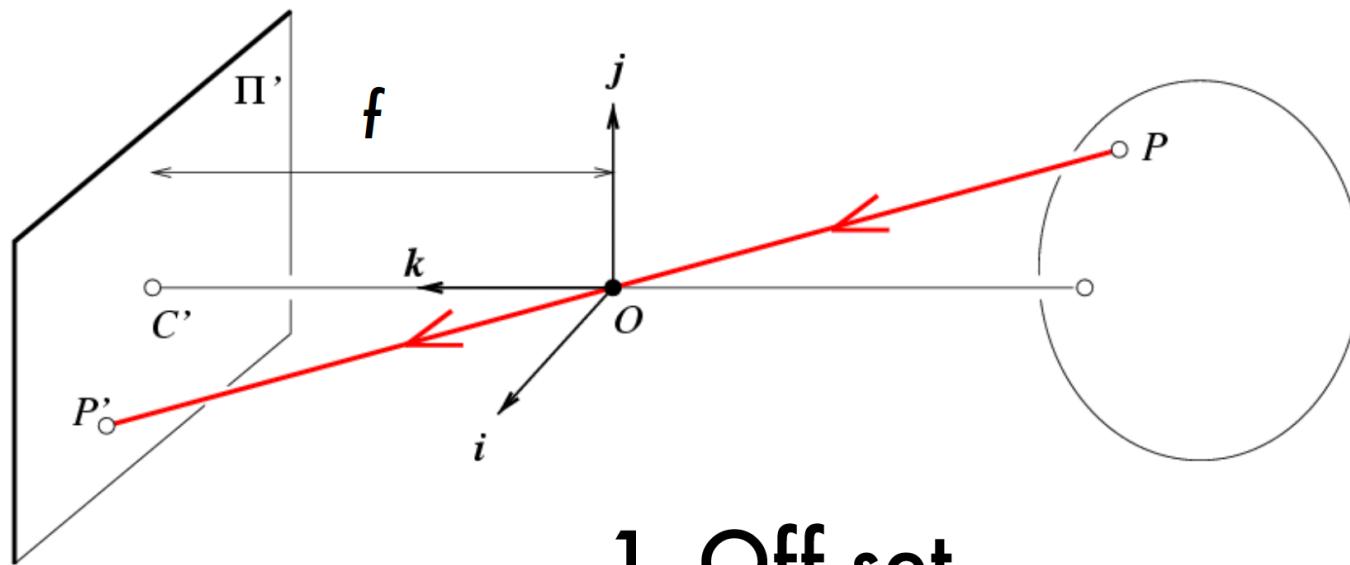
Hao Su

Lecture 13: Camera Models

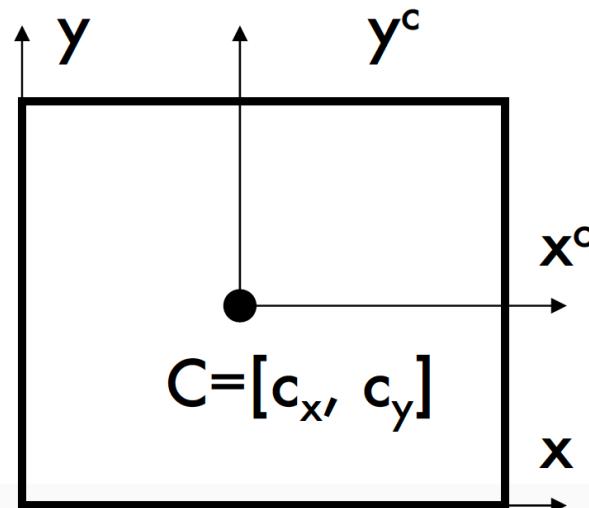


Credit: CS231a, Stanford, Silvio Savarese

Converting to pixels



1. Off set
2. From metric to pixels



$$(x, y, z) \rightarrow \left(f \frac{x}{z} + c_x, f \frac{y}{z} + c_y \right)$$

[Eq. 6]

Units: k, l : pixel/m
 f : m

Non-square pixels
 α, β : pixel

Projective transformation in the homogenous coordinate system

$$P_h' = \begin{bmatrix} \alpha x + c_x z \\ \beta y + c_y z \\ z \end{bmatrix} = \begin{bmatrix} \alpha & 0 & c_x & 0 \\ 0 & \beta & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \quad P_h$$

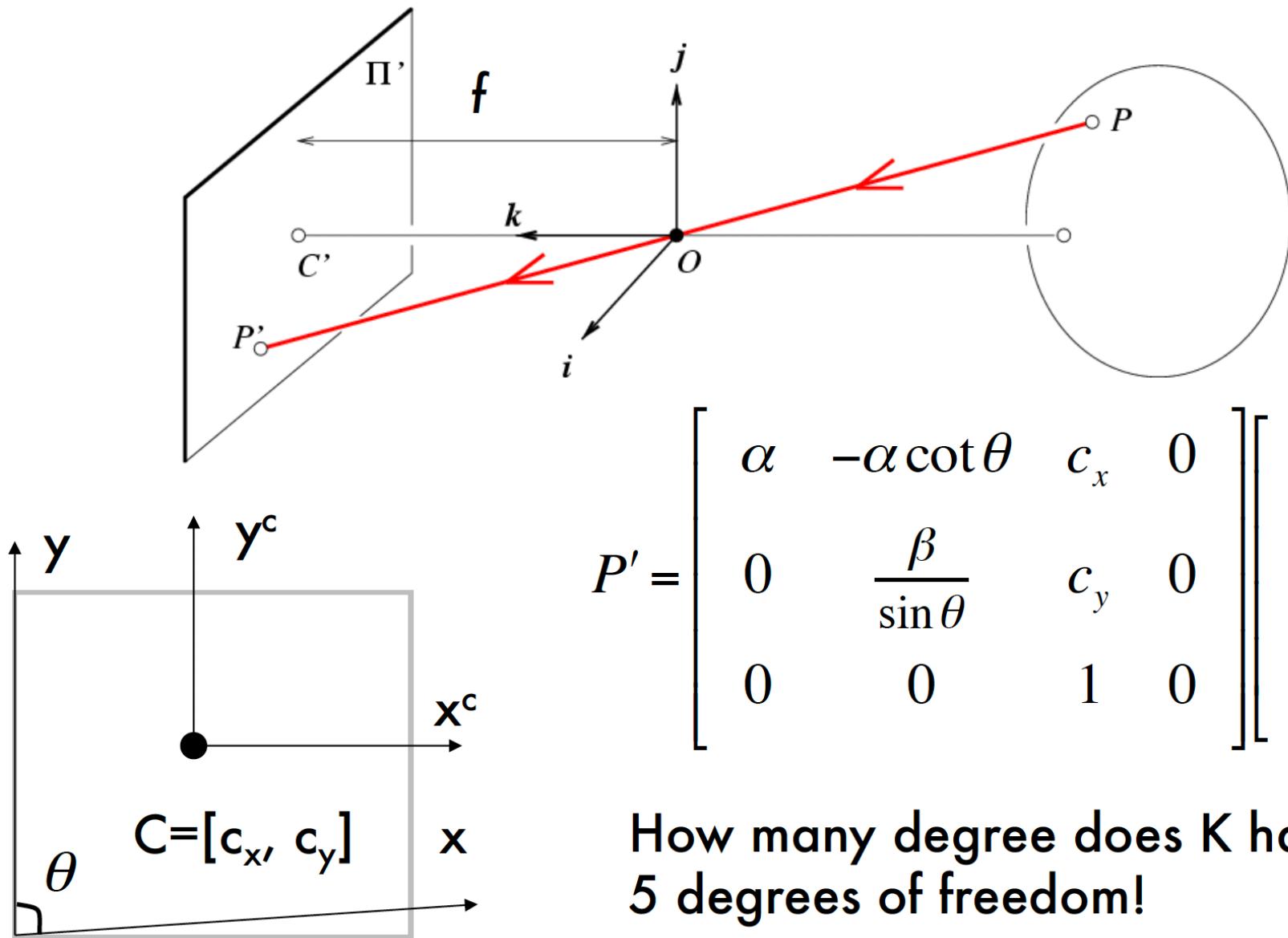
[Eq.8]

Homogenous Euclidian

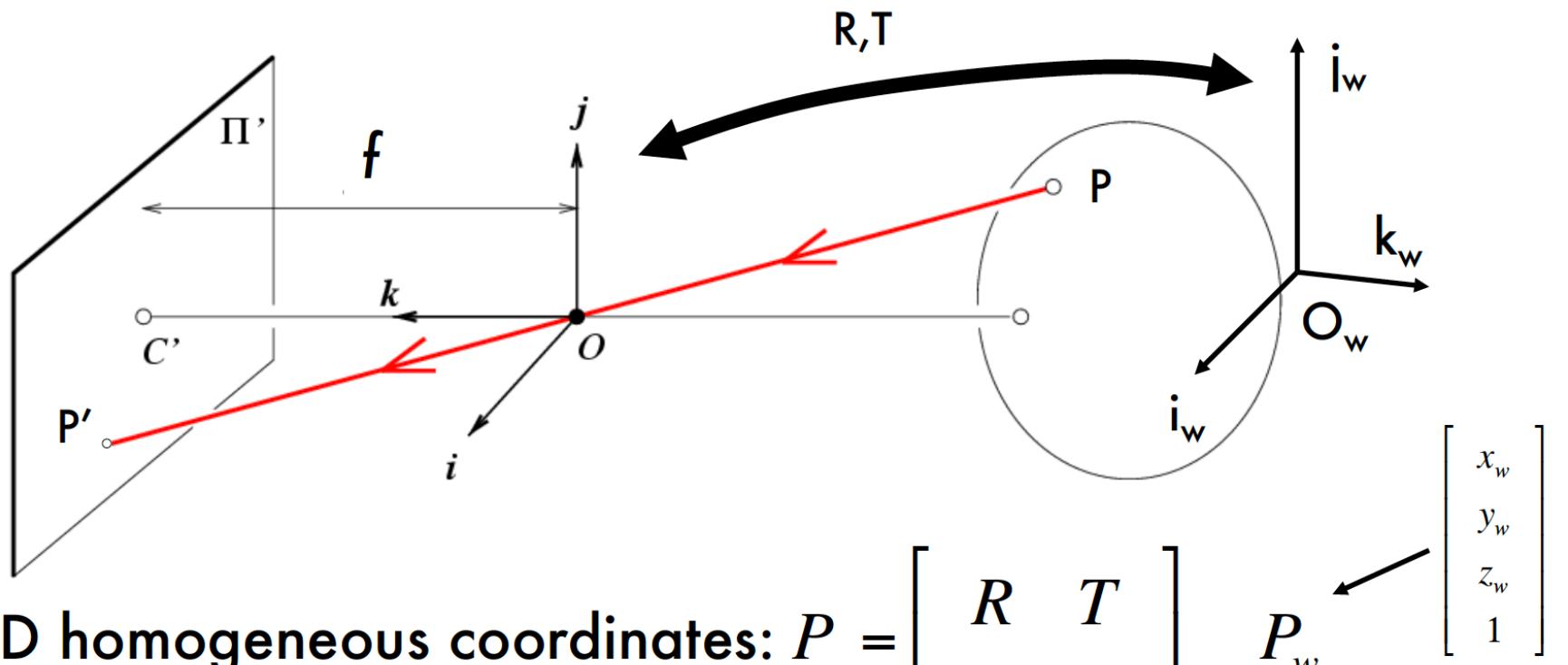
$P_h' \rightarrow P' = (\alpha \frac{x}{z} + c_x, \beta \frac{y}{z} + c_y)$

$$M = \begin{bmatrix} \alpha & 0 & c_x & 0 \\ 0 & \beta & c_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Camera Skewness



World reference system

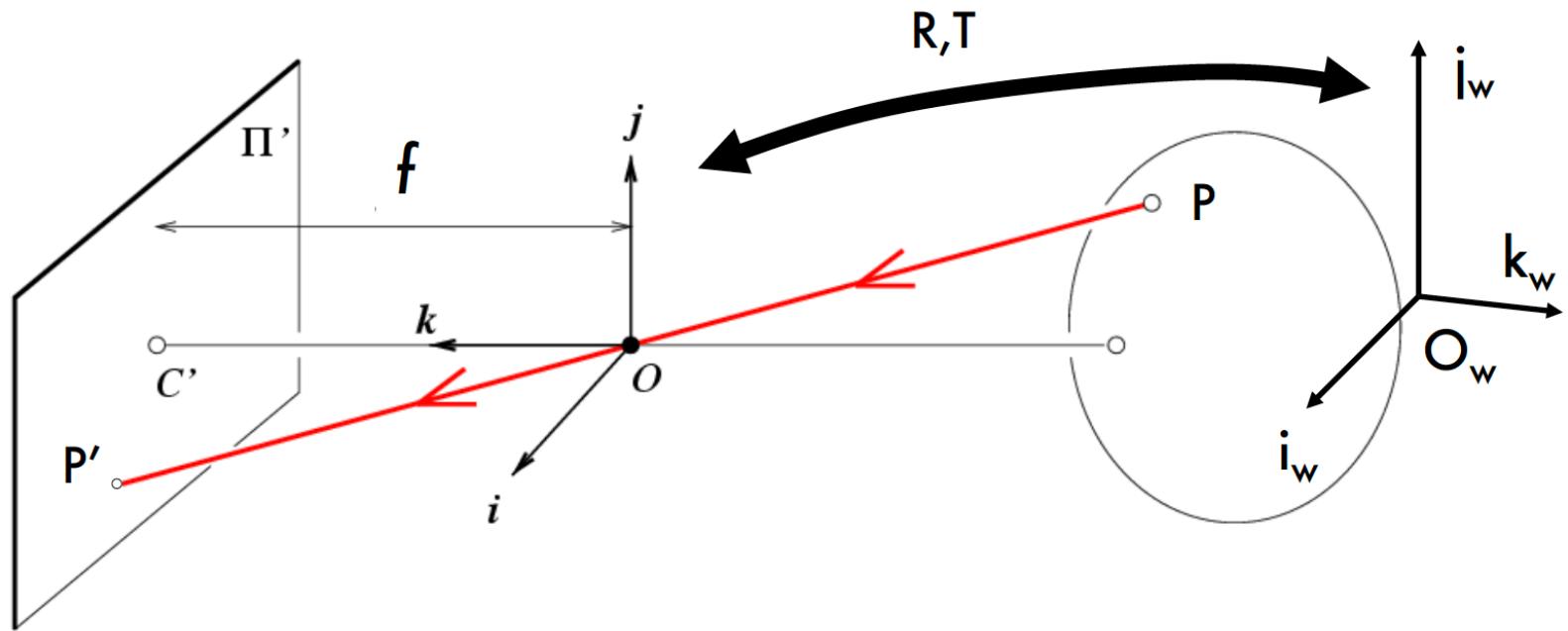


$$P' = K \begin{bmatrix} I & 0 \end{bmatrix} P = K \begin{bmatrix} I & 0 \end{bmatrix} \begin{bmatrix} R & T \\ 0 & 1 \end{bmatrix}_{4 \times 4} P_w = K \boxed{\begin{bmatrix} R & T \end{bmatrix}} P_w$$

intrinsic parameters extrinsic parameters

M [Eq.11]

The projective transformation



$$P'_{3 \times 1} = M_{3 \times 4} P_w = K_{3 \times 3} \begin{bmatrix} R & T \end{bmatrix}_{3 \times 4} P_w_{4 \times 1}$$

How many degrees of freedom?

$$5 + 3 + 3 = 11!$$

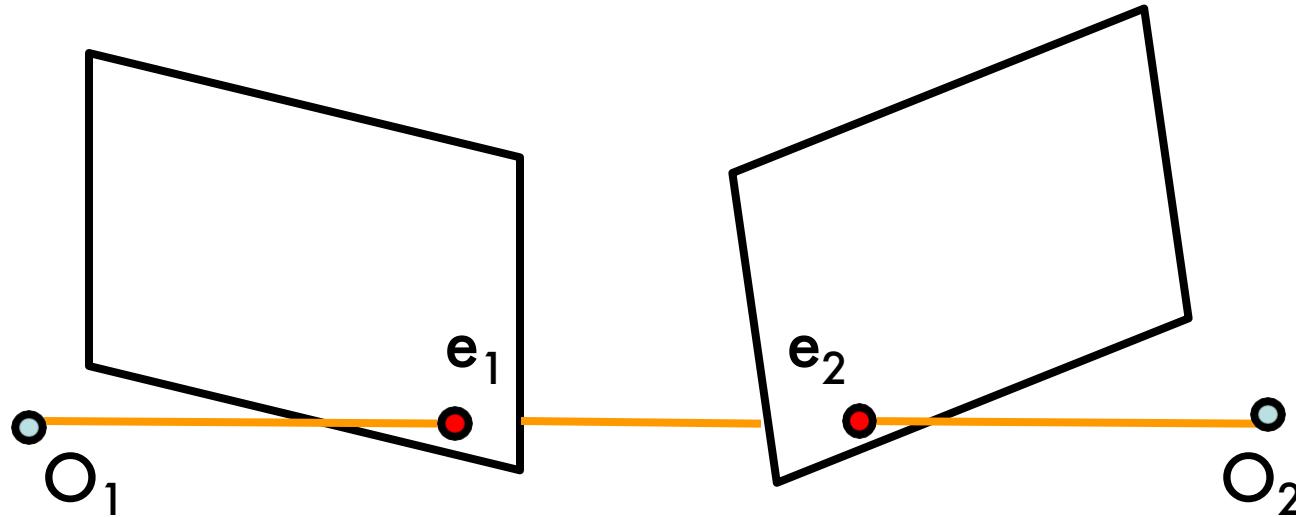
CSE 152: Computer Vision

Hao Su

Lecture 14: Multiview Geometry



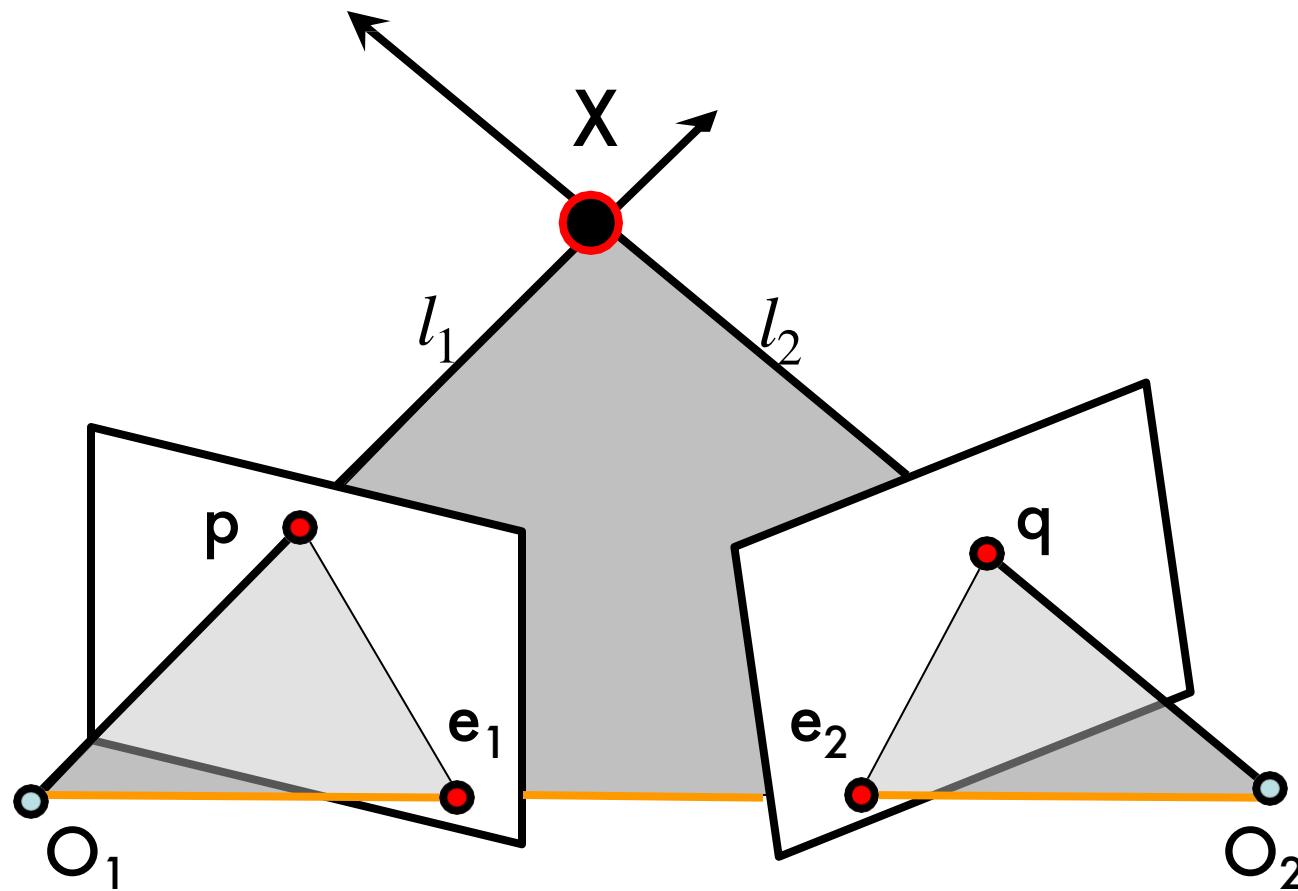
Epipolar Geometry



- Baselines

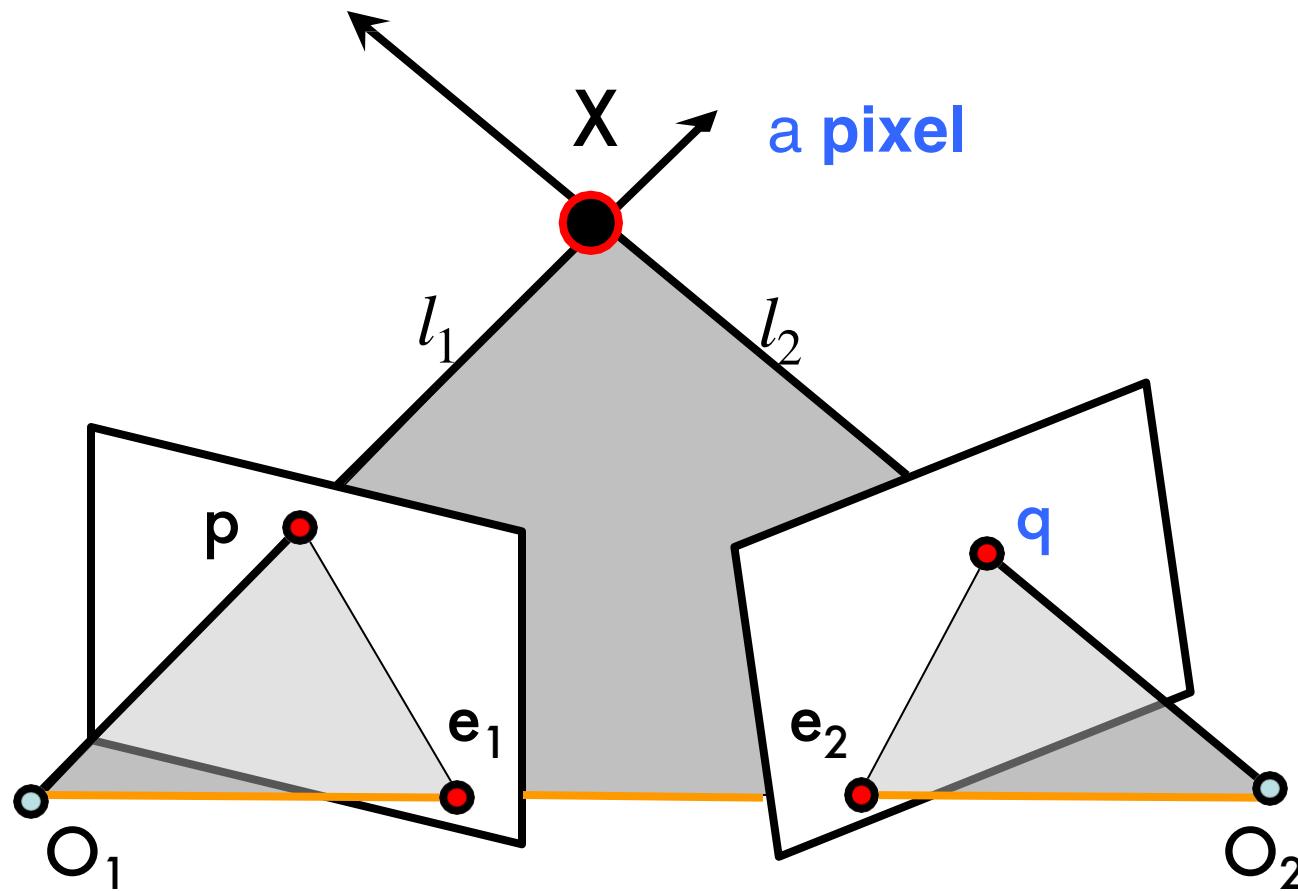
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry



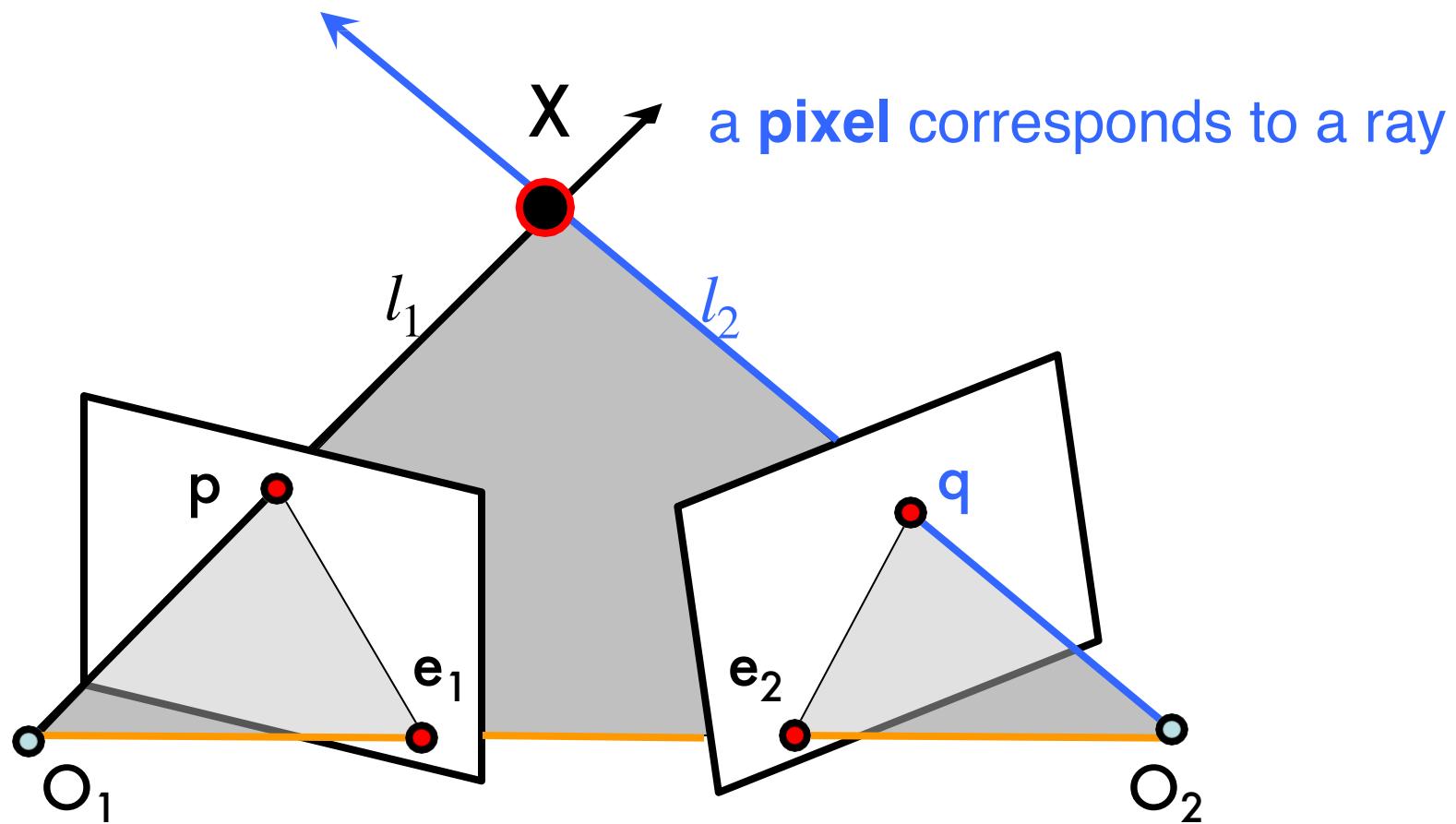
- Baselines
- Epipolar plane
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry



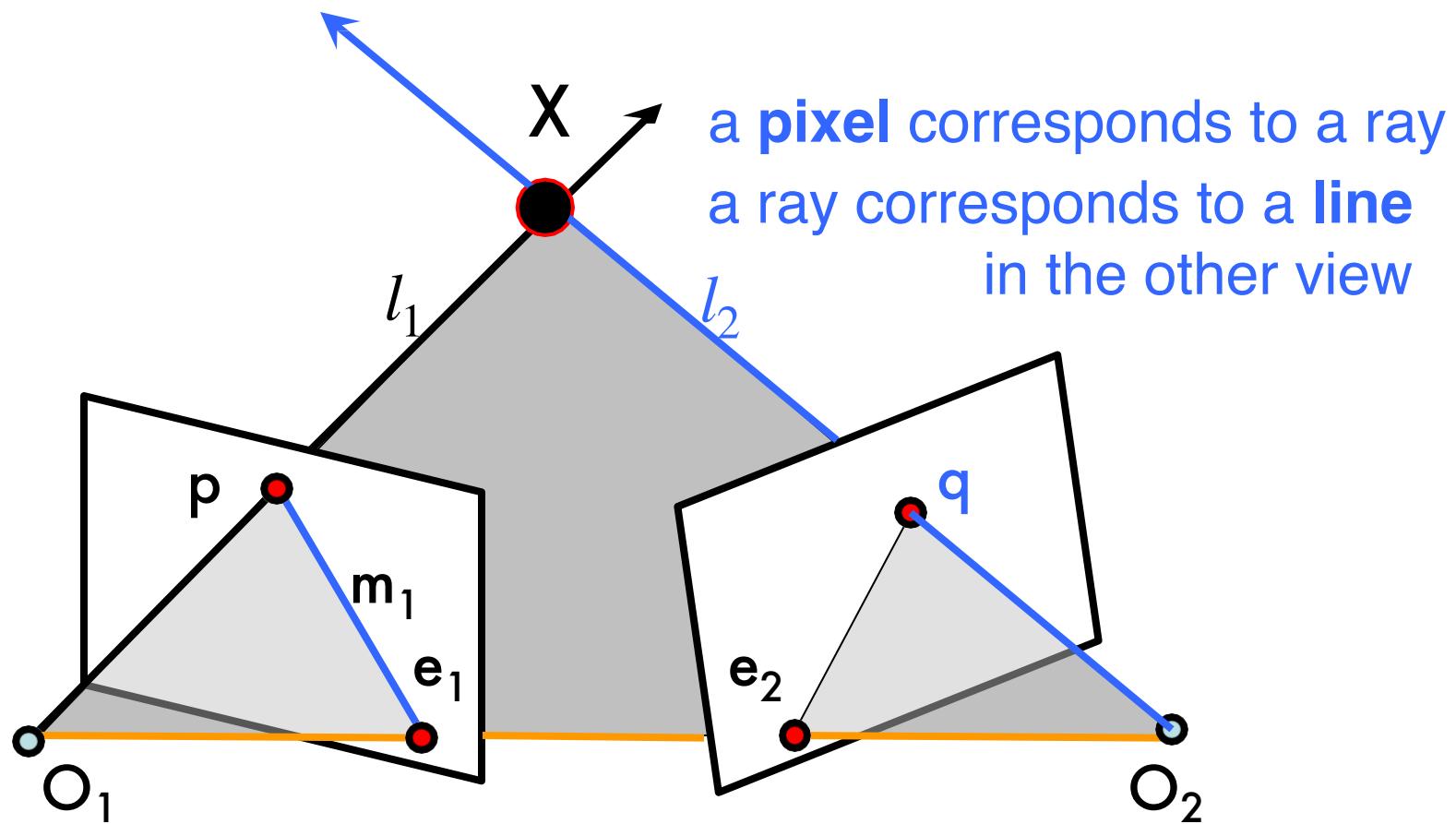
- Baselines
- Epipolar plane
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry



- Baselines
- Epipolar plane
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry

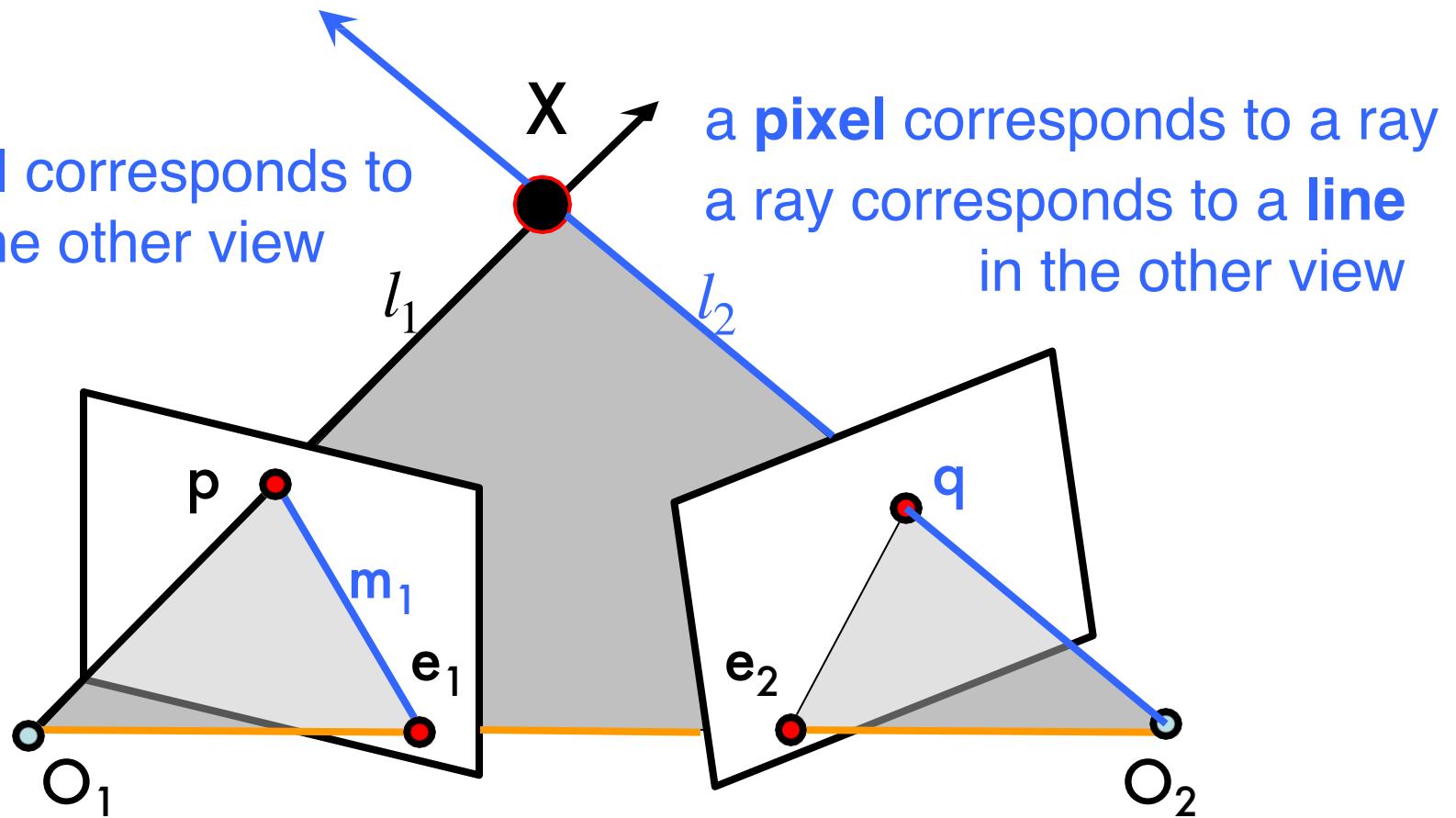


- Baselines
- Epipolar plane

- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry

so, a **pixel** corresponds to
a **line** in the other view



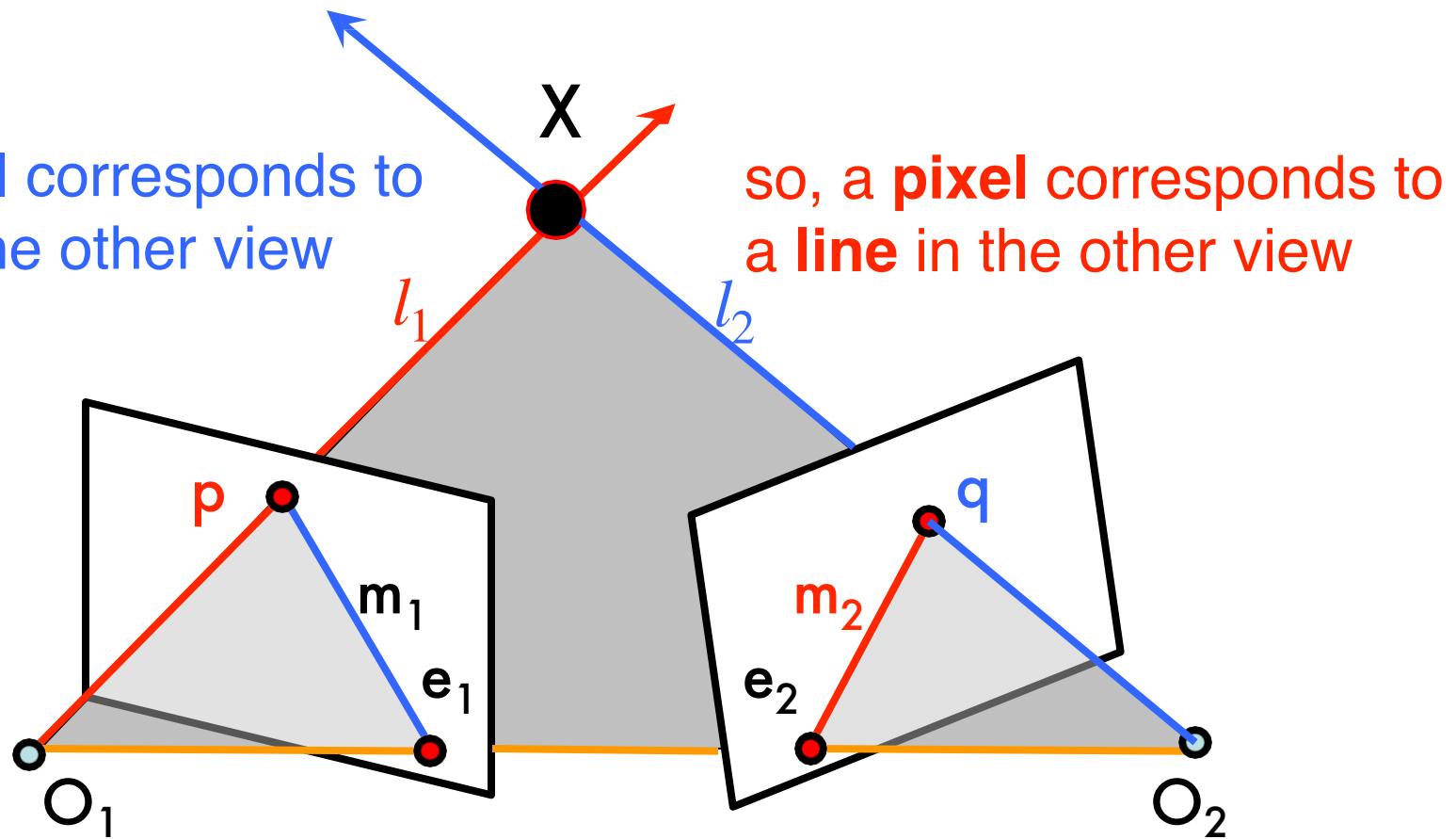
a **pixel** corresponds to a **ray**
a **ray** corresponds to a **line**
in the other view

- Baselines
- Epipolar plane

- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

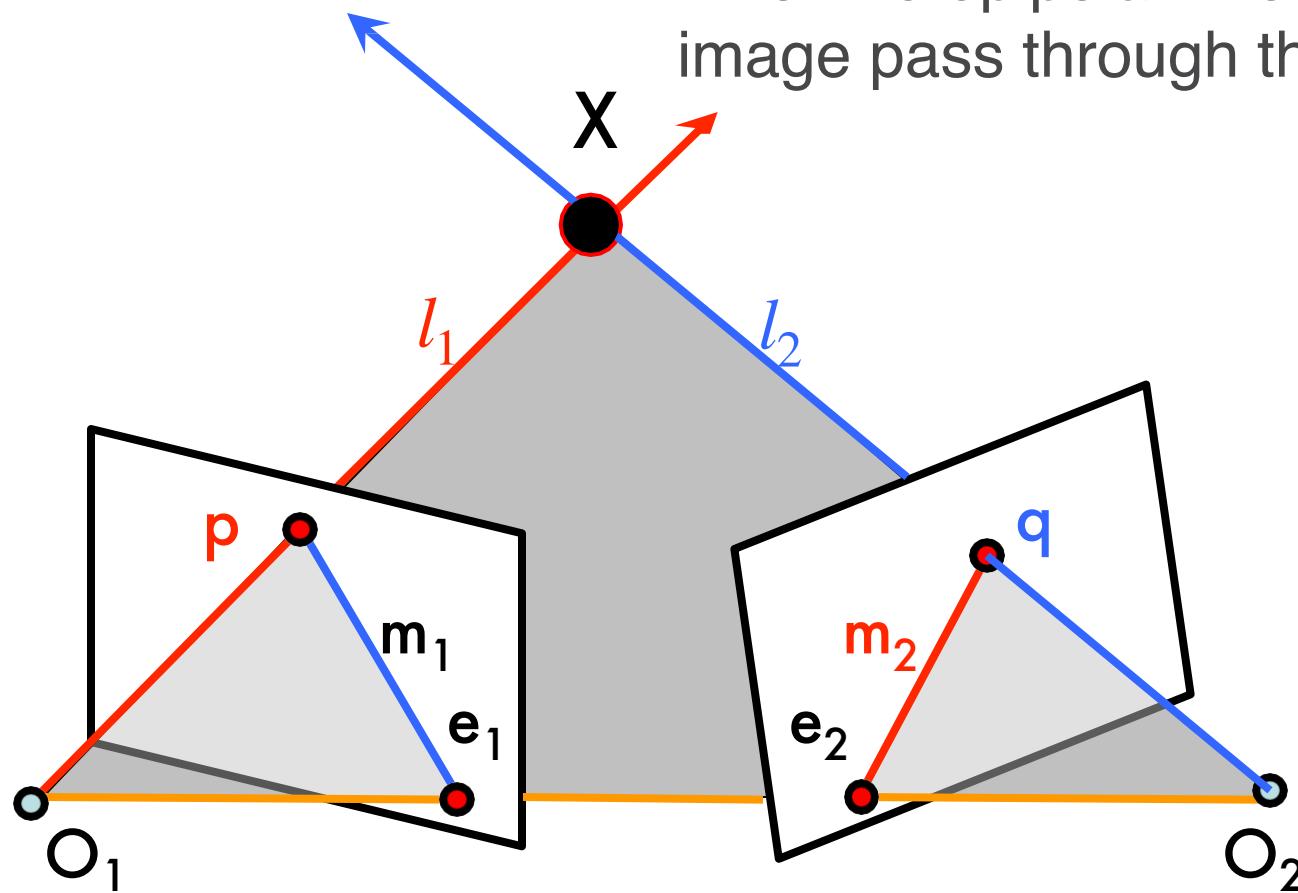
Epipolar Geometry

so, a **pixel** corresponds to
a **line** in the other view



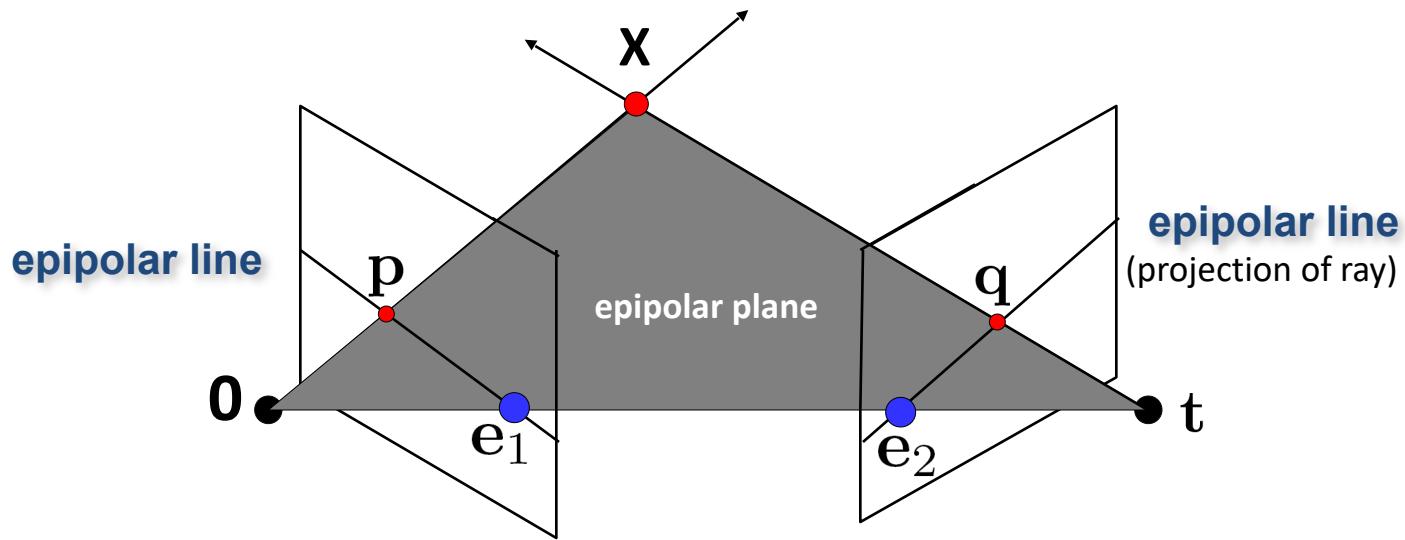
- Baselines
- Epipolar plane
- Epipolar line
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Epipolar Geometry



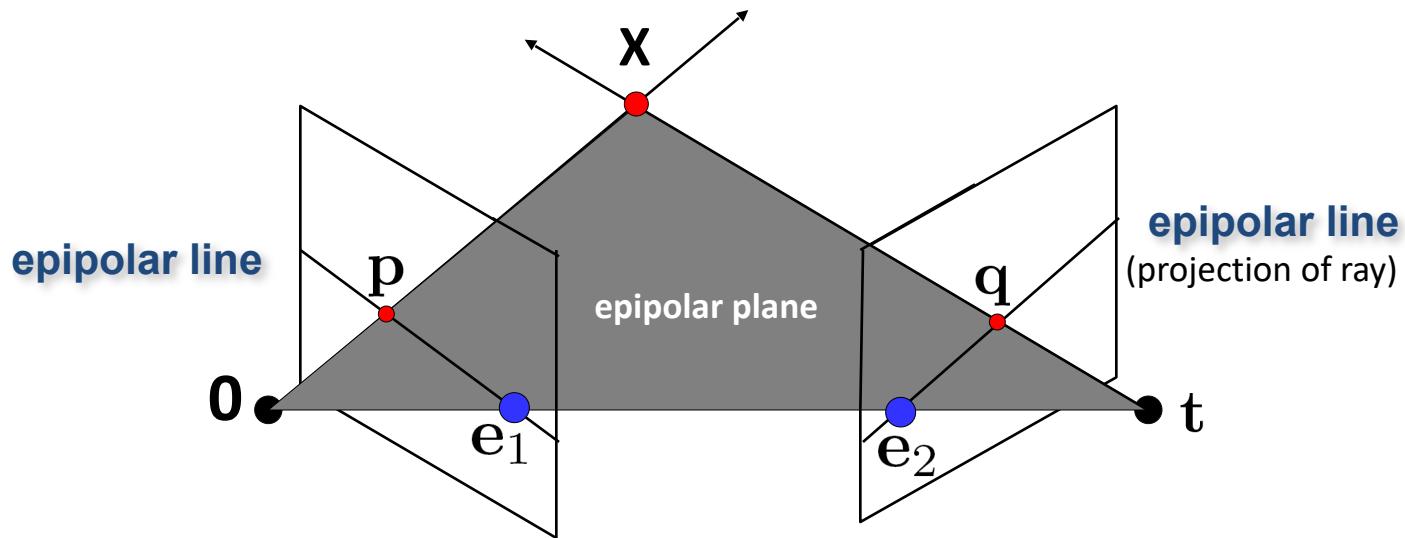
- Baselines
- Epipolar plane
- Epipolar line
- Epipoles: e_1, e_2
 - = intersections of baseline with image planes
 - = projections of the other camera center

Essential Matrix



- Assume p and q in \mathbb{R}^3 are two points on the (virtual) image plane of two cameras
- Denoted by the pinhole frame coordinate in the corresponding cameras

Essential Matrix

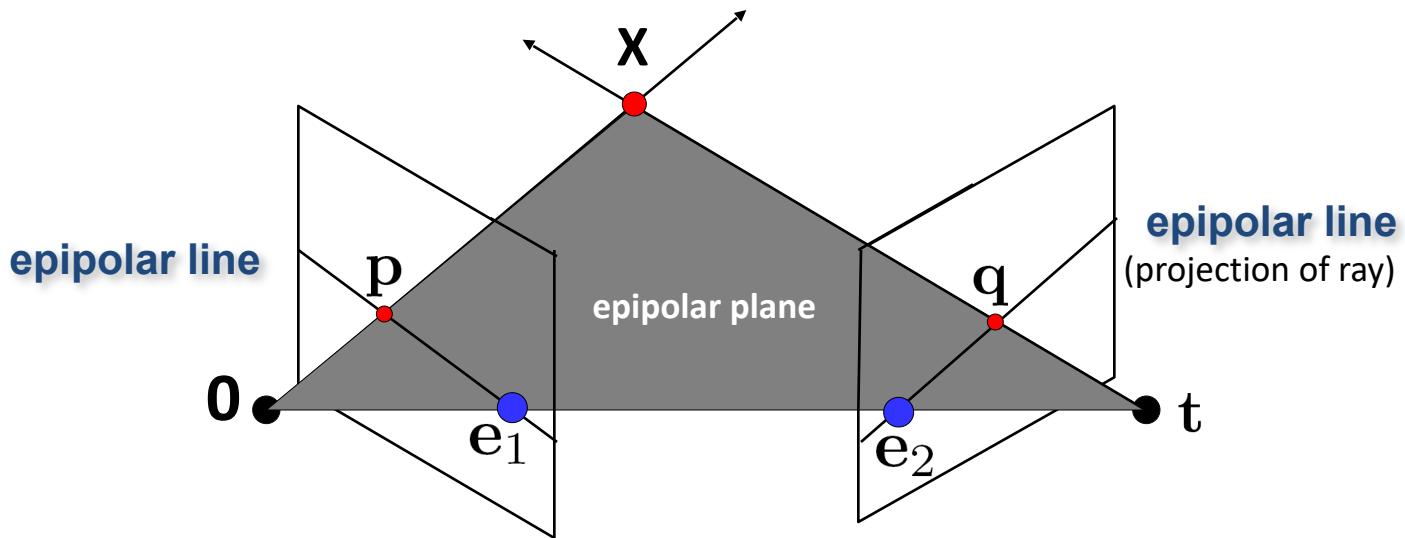


- We have: $(p^1)^T [t]_{\times} R q^2 = 0$
- Define: $\mathbf{E} = [t]_{\times} \mathbf{R}$
- Then, we have:
$$(p^1)^T E q^2 = 0 \xrightarrow{\text{omit superscript}} p^T E q = 0$$

rank(E)=2

Essential matrix

Fundamental Matrix



- Consider intrinsic camera matrices
- Then, \mathbf{p} and \mathbf{q} are in the pinhole frame and pixel counterparts are:

$$\mathbf{p}' = \mathbf{K}_1 \mathbf{p} \quad \mathbf{q}' = \mathbf{K}_2 \mathbf{q}$$

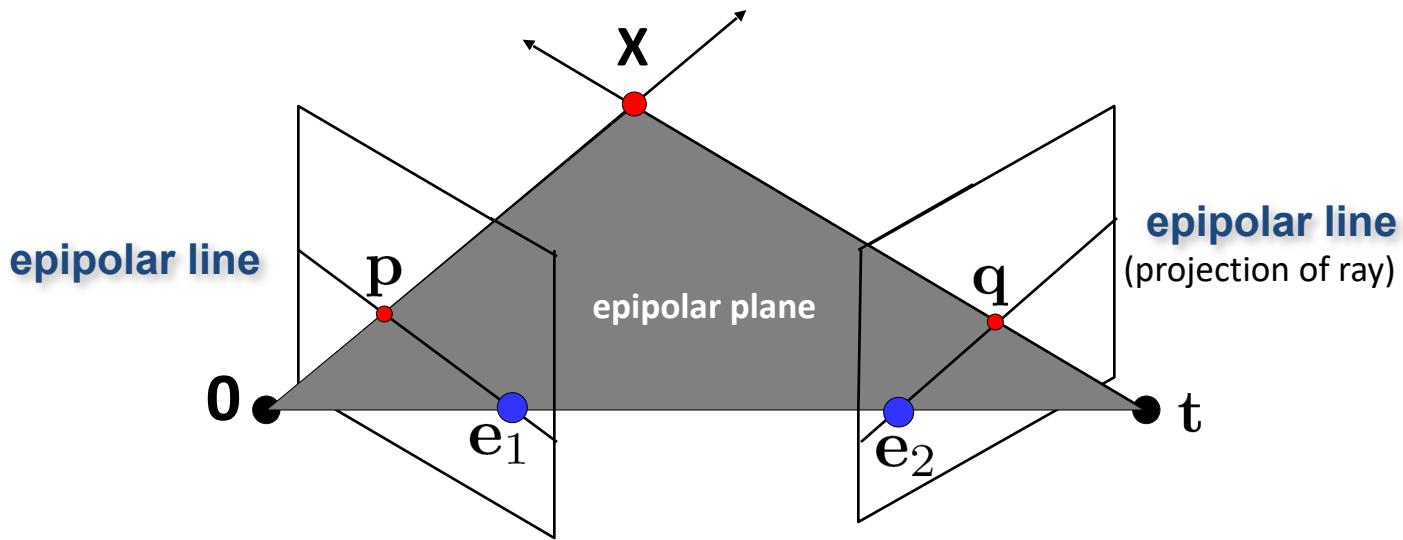
- Recall essential matrix constraint:

$$\mathbf{p}'^T \mathbf{E} \mathbf{q} = 0$$

- Substituting, we have:

$$(\mathbf{K}_1^{-1} \mathbf{p}')^T \mathbf{E} (\mathbf{K}_2^{-1} \mathbf{q}') = 0$$

Fundamental Matrix



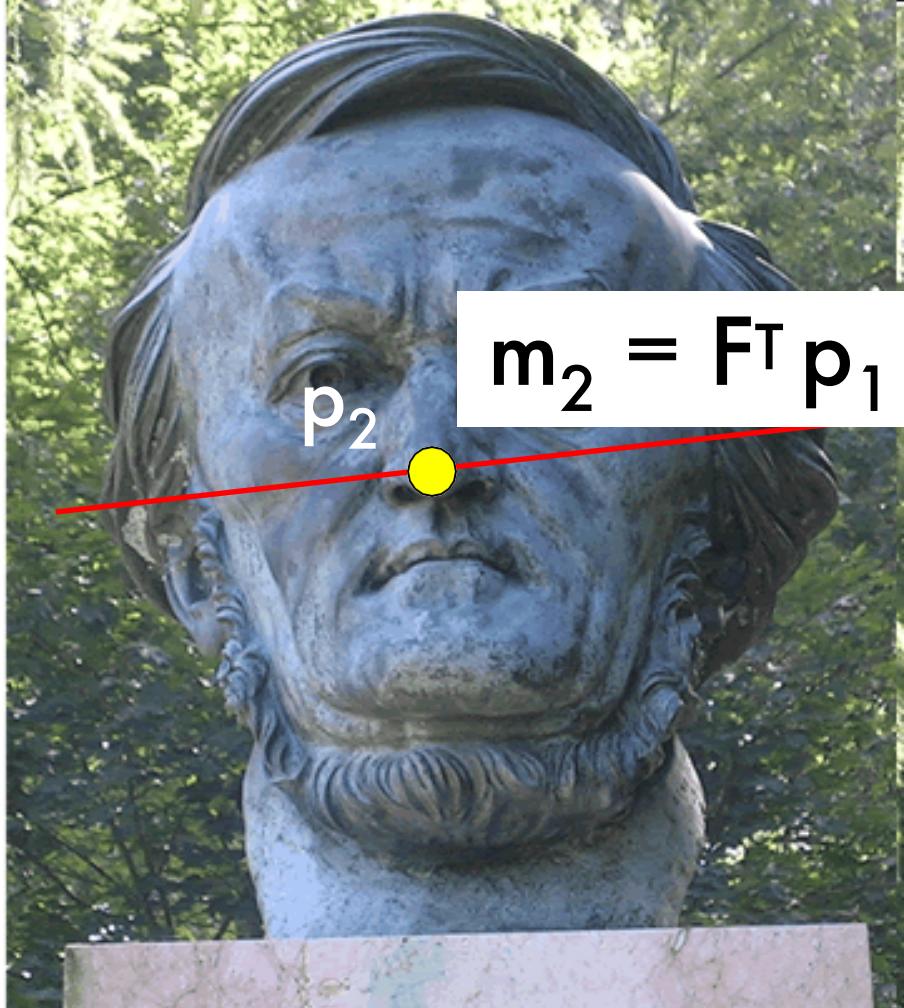
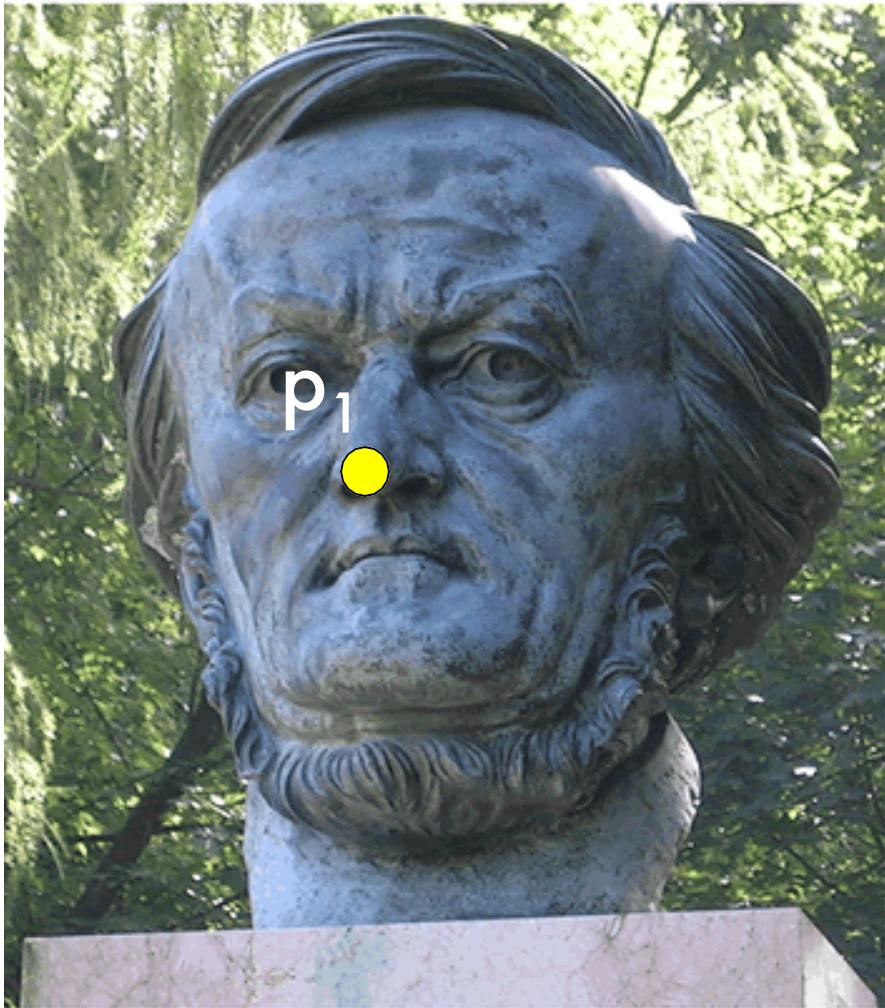
- Essential matrix constraint in pixel space: $(K_1^{-1}p')^T E (K_2^{-1}q') = 0$
- Rearranging: $p'^T K_1^{-T} E K_2^{-1} q' = 0$
- Define: $F = K_1^{-T} E K_2^{-1}$ Fundamental matrix rank(F)=2
- Then, we have: $p'^T F q' = 0$

Epipolar Constraint

$$p_1^T \cdot F p_2 = 0$$

- $w_1 = F p_2$ defines an equation $w_1^T p_1 = 0$
- Note that, p_1 is the corresponding point of p_2 by the derivation of F
- So, $w_1 = F p_2$ defines the epipolar line of p_2

Why F is useful?



- Suppose F is known
- No additional information about the scene and camera is given
- Given a point on left image, we can compute the corresponding epipolar line in the second image

Estimating F

Suppose we have a pair of corresponding points:

$$[\text{Eq. 13}] \quad p^T F p' = 0 \quad \rightarrow$$

$$p = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \quad p' = \begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix}$$

$$(u, v, 1) \begin{pmatrix} F_{11} & F_{12} & F_{13} \\ F_{21} & F_{22} & F_{23} \\ F_{31} & F_{32} & F_{33} \end{pmatrix} \begin{pmatrix} u' \\ v' \\ 1 \end{pmatrix} = 0$$



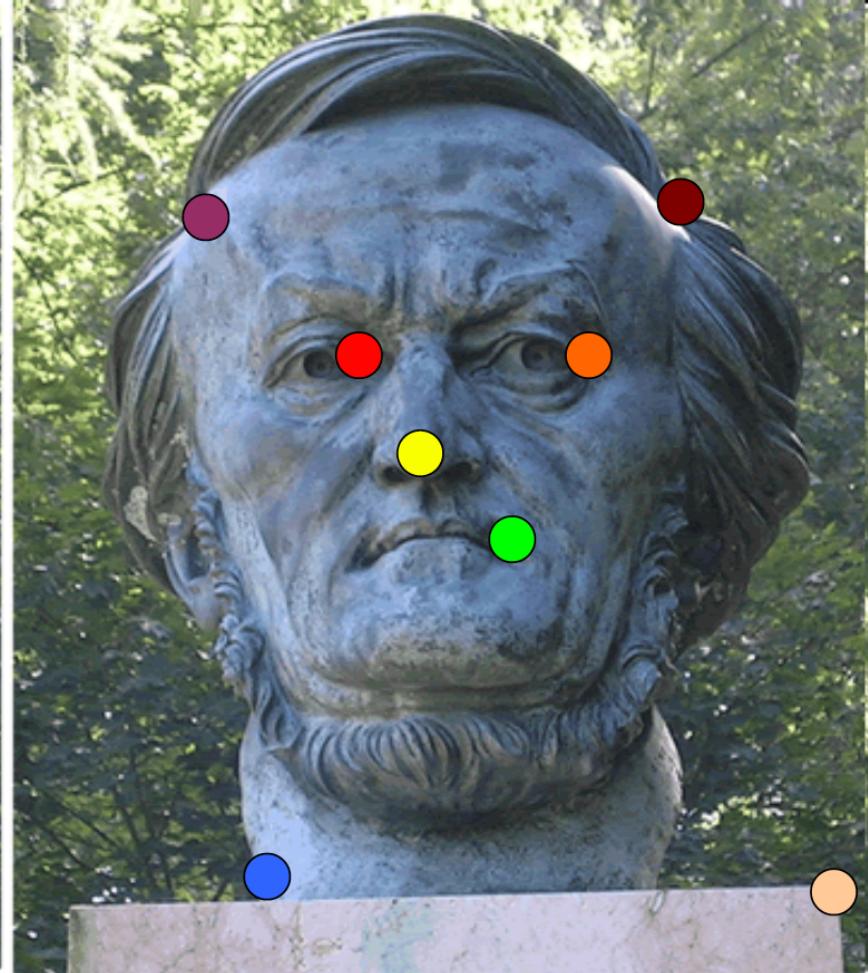
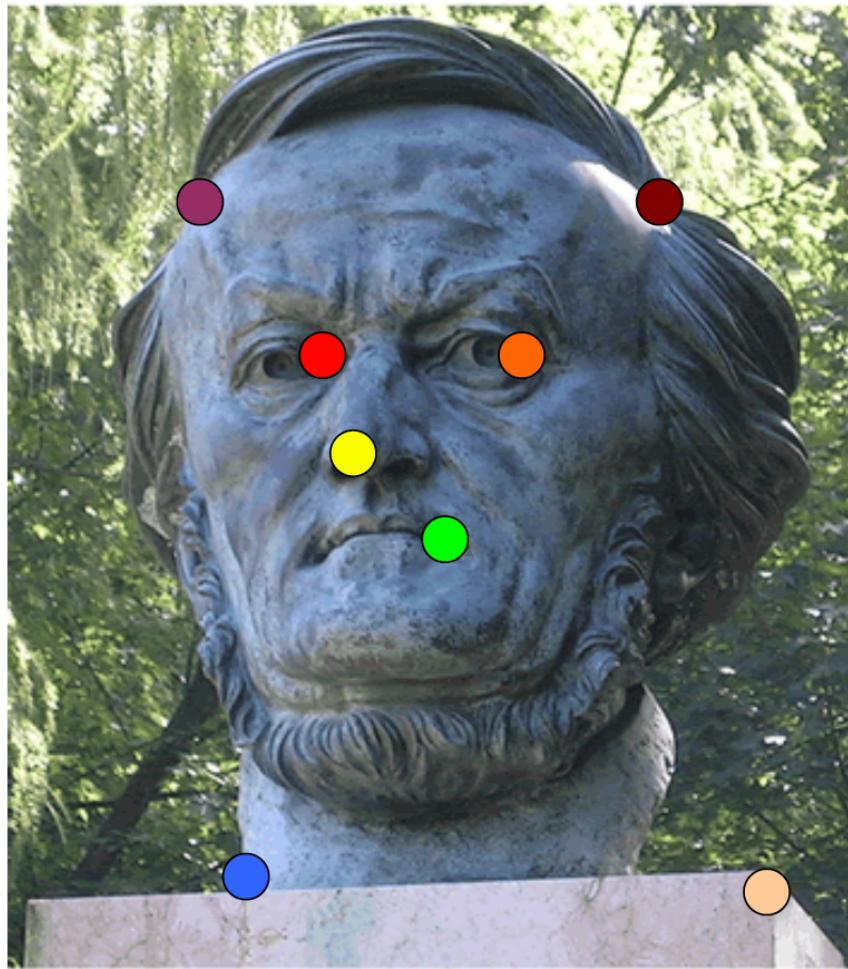
$$(uu', uv', u, vu', vv', v, u', v', 1)$$

$$\begin{pmatrix} F_{11} \\ F_{12} \\ F_{13} \\ F_{21} \\ F_{22} \\ F_{23} \\ F_{31} \\ F_{32} \\ F_{33} \end{pmatrix} = 0$$

[Eq. 14]

Let's take 8 corresponding points

Estimating F

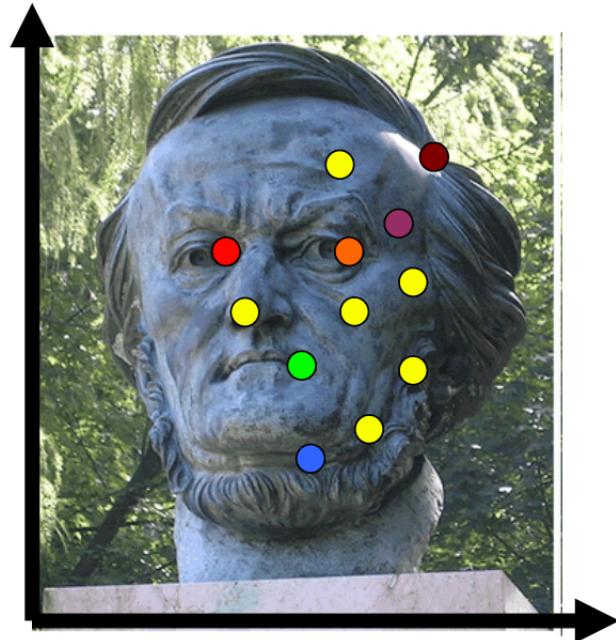


Estimating F

$$\mathbf{W} \begin{pmatrix} u_1 u'_1 & u_1 v'_1 & u_1 & v_1 u'_1 & v_1 v'_1 & v_1 & u'_1 & v'_1 & 1 \\ u_2 u'_2 & u_2 v'_2 & u_2 & v_2 u'_2 & v_2 v'_2 & v_2 & u'_2 & v'_2 & 1 \\ u_3 u'_3 & u_3 v'_3 & u_3 & v_3 u'_3 & v_3 v'_3 & v_3 & u'_3 & v'_3 & 1 \\ u_4 u'_4 & u_4 v'_4 & u_4 & v_4 u'_4 & v_4 v'_4 & v_4 & u'_4 & v'_4 & 1 \\ u_5 u'_5 & u_5 v'_5 & u_5 & v_5 u'_5 & v_5 v'_5 & v_5 & u'_5 & v'_5 & 1 \\ u_6 u'_6 & u_6 v'_6 & u_6 & v_6 u'_6 & v_6 v'_6 & v_6 & u'_6 & v'_6 & 1 \\ u_7 u'_7 & u_7 v'_7 & u_7 & v_7 u'_7 & v_7 v'_7 & v_7 & u'_7 & v'_7 & 1 \\ u_8 u'_8 & u_8 v'_8 & u_8 & v_8 u'_8 & v_8 v'_8 & v_8 & u'_8 & v'_8 & 1 \end{pmatrix} \begin{pmatrix} F_{11} \\ F_{12} \\ F_{13} \\ F_{21} \\ F_{22} \\ F_{23} \\ F_{31} \\ F_{32} \\ F_{33} \end{pmatrix} = 0 \quad [\text{Eqs. 15}]$$

- Homogeneous system $\mathbf{W} \mathbf{f} = 0$
- Rank 8 → A non-zero solution exists (unique)
- If $N > 8$ → Lsq. solution by SVD! → $\hat{\mathbf{F}}$
 $\|\mathbf{f}\| = 1$

Flow of the 8-Point Algorithm



$$Wf = 0, \|f\| = 1$$

↓ Least-square

$$\begin{aligned} & \text{minimize}_{\substack{f \\ s.t.}} && \|Wf\|^2 \\ & && \|f\| = 1 \end{aligned}$$

$$\begin{aligned} & \text{minimize}_{\substack{f \\ s.t.}} && f^T W^T W f \\ & && f^T f = 1 \end{aligned}$$

Do you remember how to solve the problem?
Hint: Check your HW1 (by the SVD of W)

Find F that minimizes

$$\|F - \hat{F}\|$$

Frobenius norm (*)

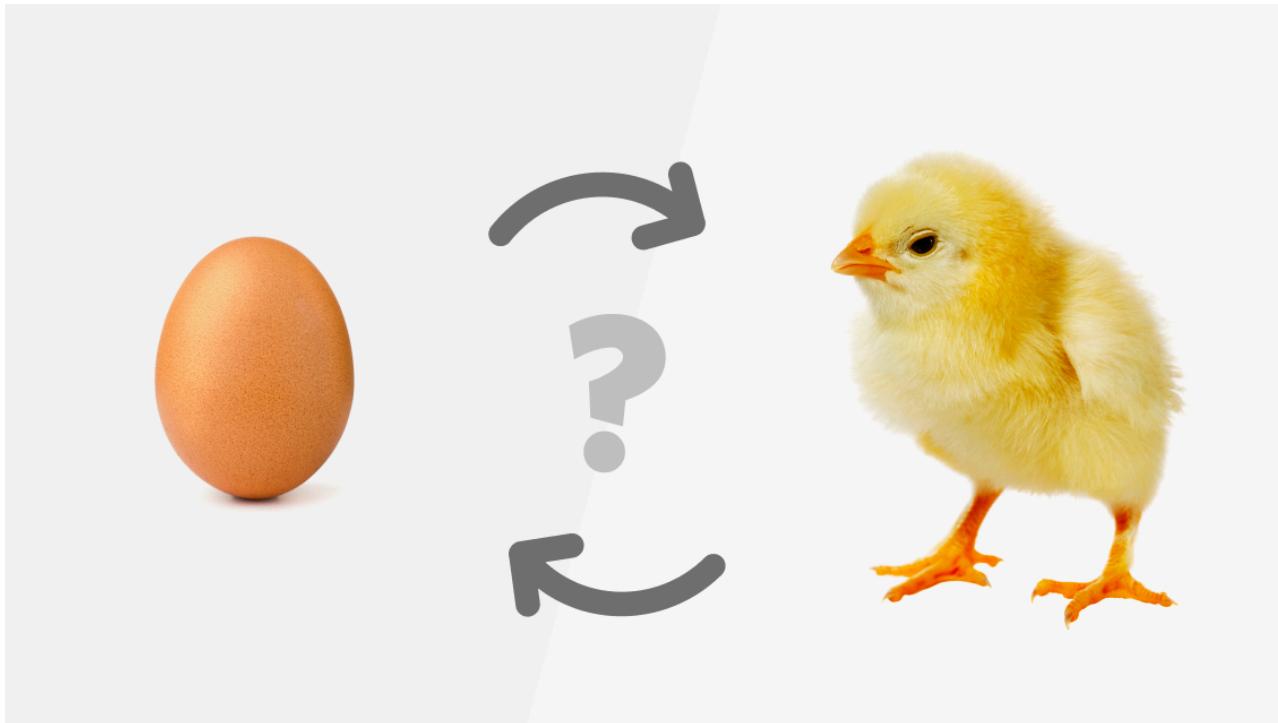
Subject to $\det(F) = 0$

$$F = U \begin{bmatrix} s_1 & 0 & 0 \\ 0 & s_2 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T \quad \text{Where:}$$

$$U \begin{bmatrix} s_1 & 0 & 0 \\ 0 & s_2 & 0 \\ 0 & 0 & s_3 \end{bmatrix} V^T = SVD(\hat{F})$$

Simultaneous Correspondence and F Estimation

- *Simultaneous* Correspondence and F Estimation
 - With F, it is easier to compute correspondence
 - With correspondence, we can estimate F
- A Chicken-or-Egg problem

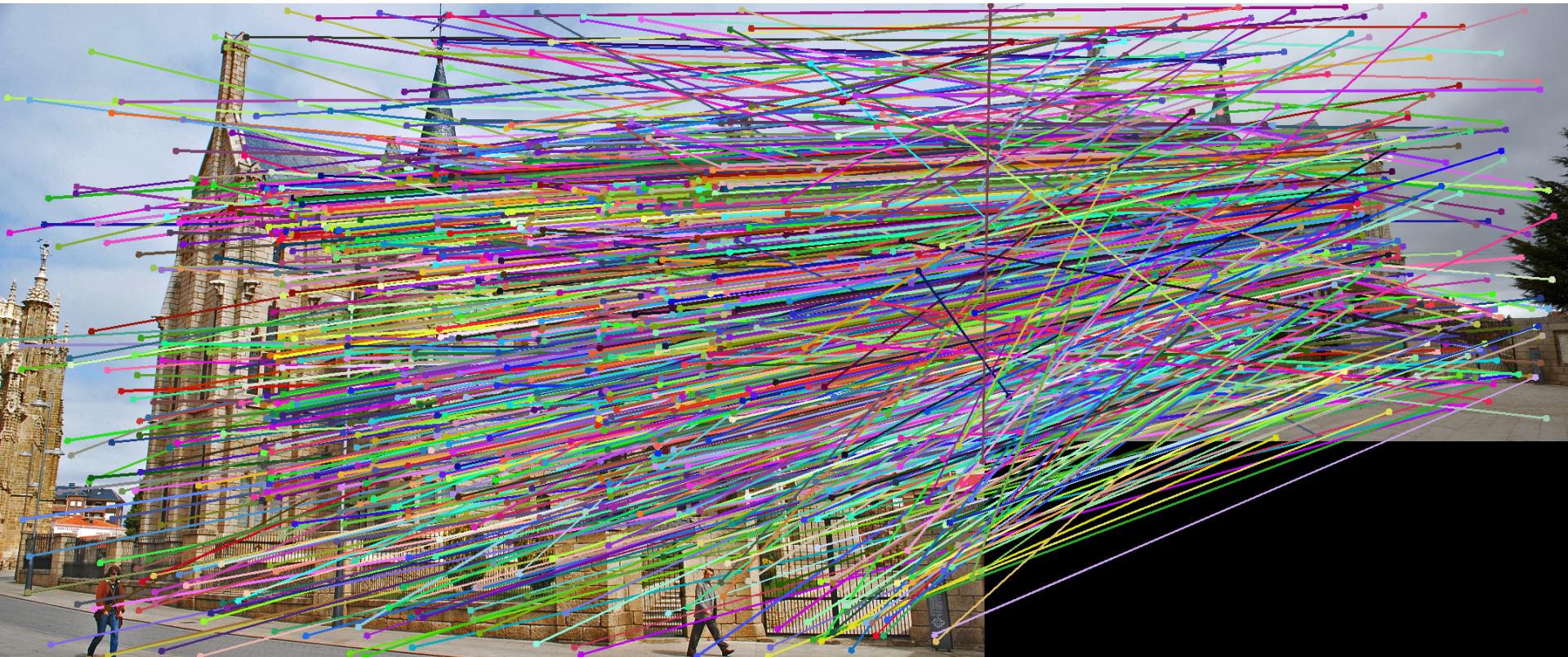


Basic Pipeline: Feature detection

Detect features using, for example, SIFT [Lowe, IJCV 2004] or learning-based keypoint detector



VLFeat's 800 most confident matches among
10,000+ local features.



Basic Pipeline: Repeat the above steps (RANSAC)

for i in range(n):

 randomly choose some pairs

 repeat for m times:

 based on the inliers, estimate F

 based on F, remove pairs with big errors

CSE 152: Computer Vision

Hao Su

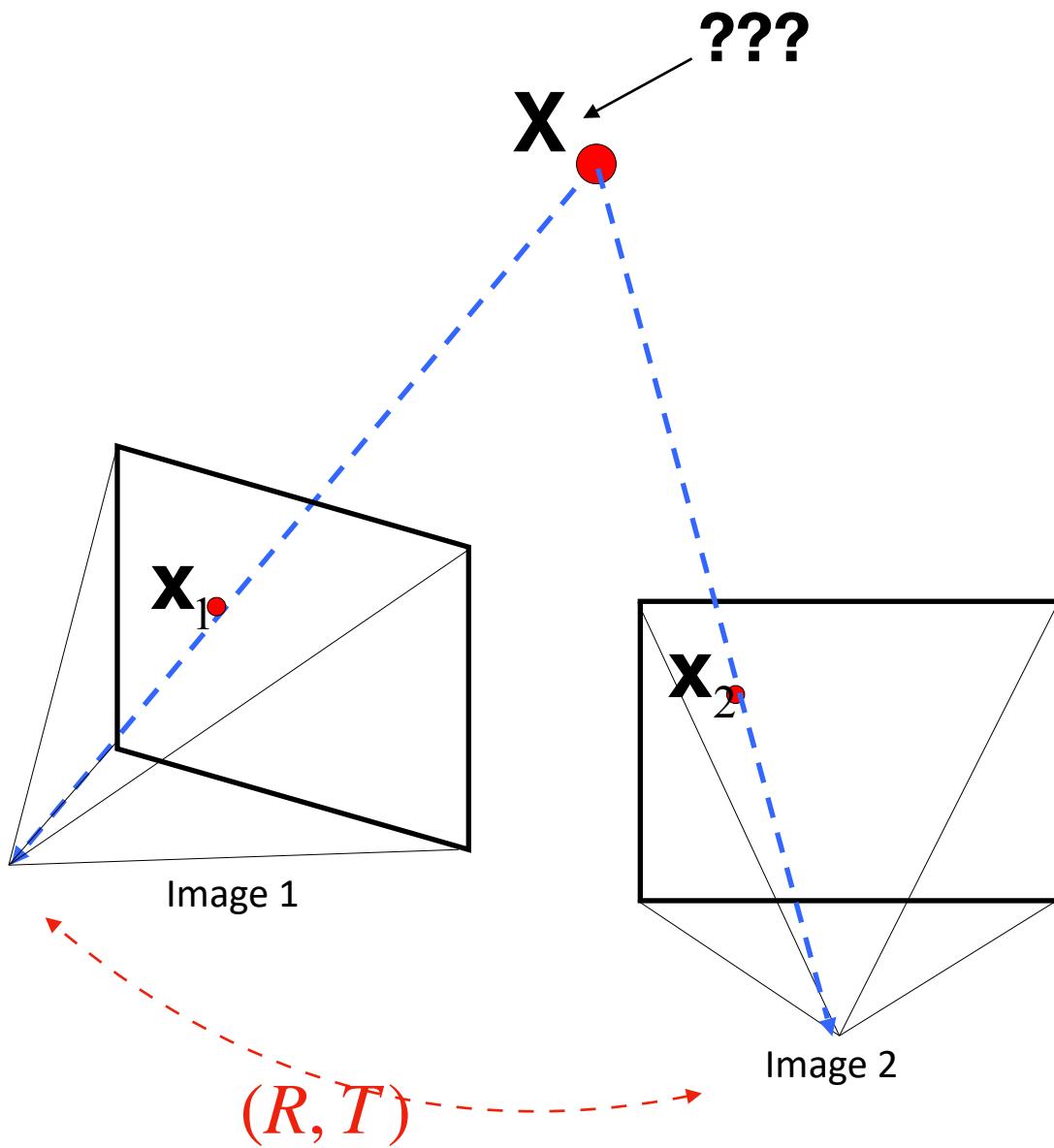
Lecture 16: Stereo Reconstruction



Problem setup

- Known:
 - Two views of the same scene
 - Corresponding points between views
 - Intrinsic camera matrices (K_1, K_2), i.e., camera calibration has been done
 - Fundamental matrix F
- Question: Point coordinates in 3D space

Step I: Estimate (R, T) Between Cameras



$$\mathbf{x}_1 \Leftrightarrow \mathbf{x}_2$$

$$\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$$

Step I: Estimate (R, T) Between Cameras

- Get E from F :

$$F = K_1^{-T} E K_2^{-1}$$

$$E = K_1^T F K_2$$

- Decompose E into skew-symmetric and rotation matrices:

$$\mathbf{E} = [\mathbf{t}]_\times \mathbf{R}$$

Step II: Reprojection Error Minimization

- With K_i , R , t , we can compute the projection matrices for both cameras:

$$P' = M P_w = \boxed{K} \boxed{[R \quad T]} P_w$$

Internal parameters External parameters

- The projective projection equation:

$$p_i = \begin{bmatrix} u_i \\ v_i \end{bmatrix} = \begin{bmatrix} \frac{\mathbf{m}_1 P_i}{\mathbf{m}_3 P_i} \\ \frac{\mathbf{m}_2 P_i}{\mathbf{m}_3 P_i} \end{bmatrix}$$

$M = \begin{bmatrix} \mathbf{m}_1 \\ \mathbf{m}_2 \\ \mathbf{m}_3 \end{bmatrix}$

in pixels

A non-linear transformation, denoted by $p_i = f(P_i)$

Step II: Reprojection Error Minimization

- Minimize sum of squared reprojection errors:

$$\underset{\{P_i\}}{\text{minimize}} \quad \sum_{i=1}^n \sum_{j=1}^2 w_{ij} \left\| f(P_i; R_j, t_j) - \begin{bmatrix} u_{i,j} \\ v_{i,j} \end{bmatrix} \right\|^2$$

indicator variable: *predicted* *observed*
whether point i visible in image j image location image location

- Optimized with non-linear least squares
- LM algorithm (Levenberg-Marquardt) is a popular choice

CSE 152: Computer Vision

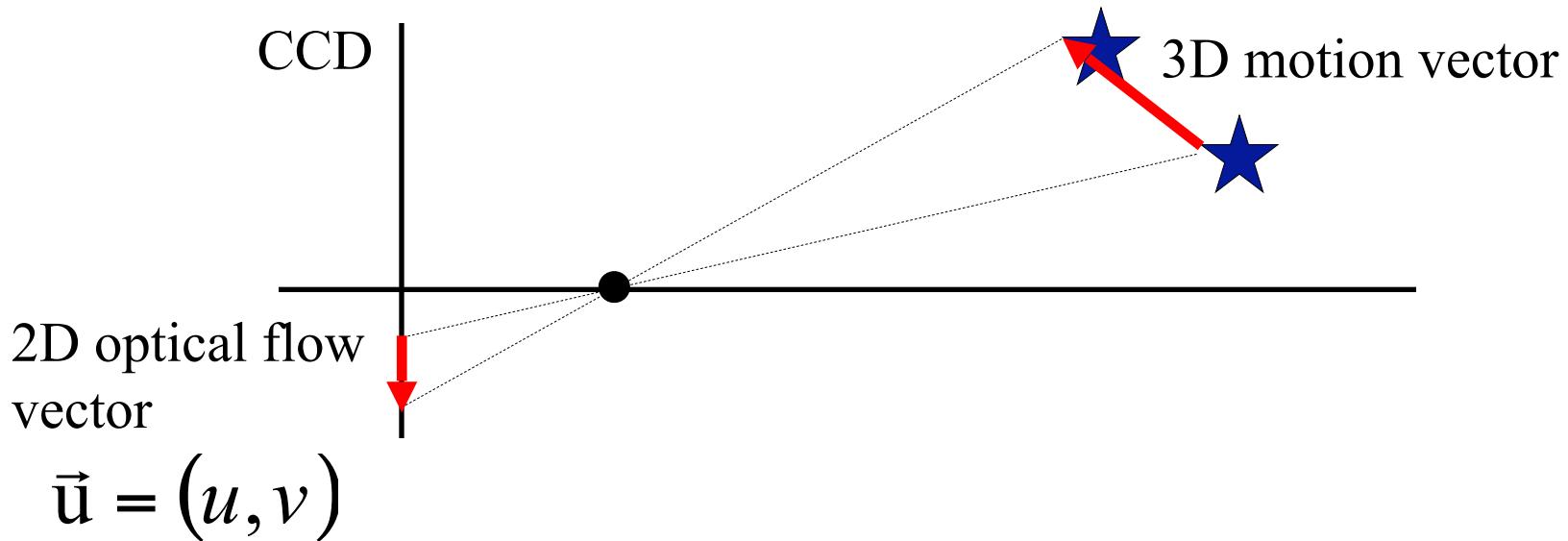
Hao Su

Lecture 17: Motion Estimation



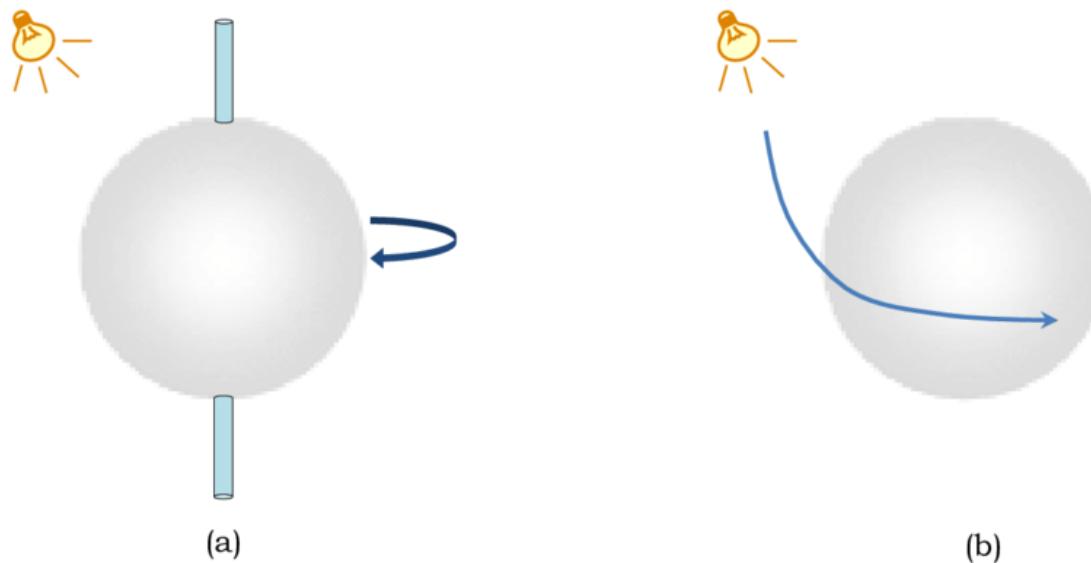
Motion Field & Optical Flow Field

- Motion Field = Real world 3D motion
- Optical Flow Field = Projection of the motion field onto the 2d image

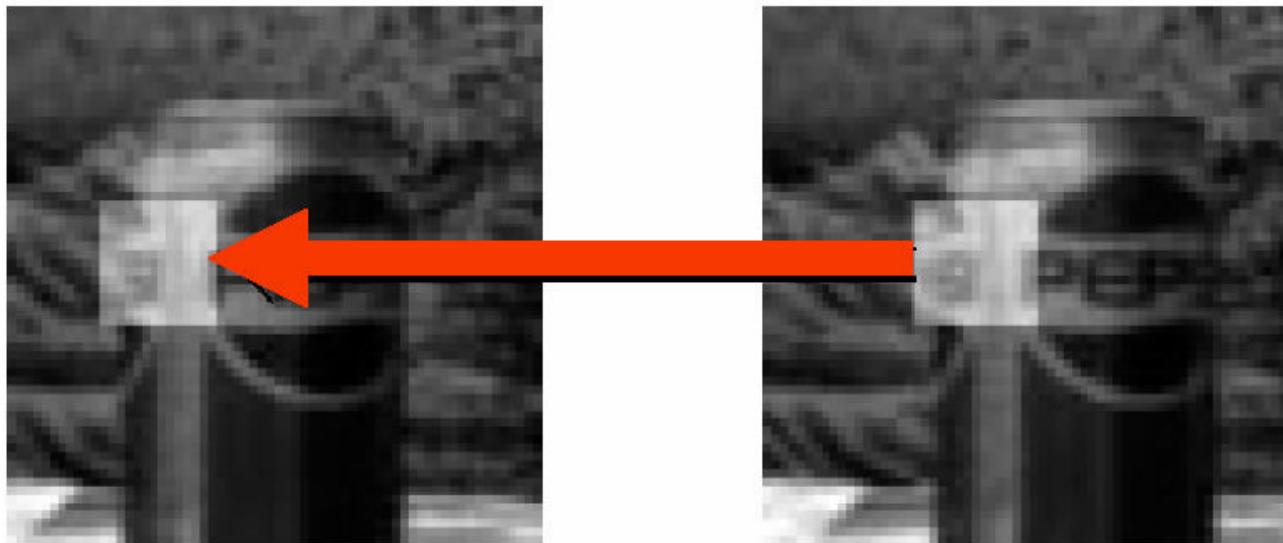


Apparent motion

- Optical flow differs from actual motion field:
 - (a) intensity remains constant, so that no motion is perceived;
 - (b) no object motion exists, however moving light source produces shading changes.



Key Assumptions: Brightness Constancy



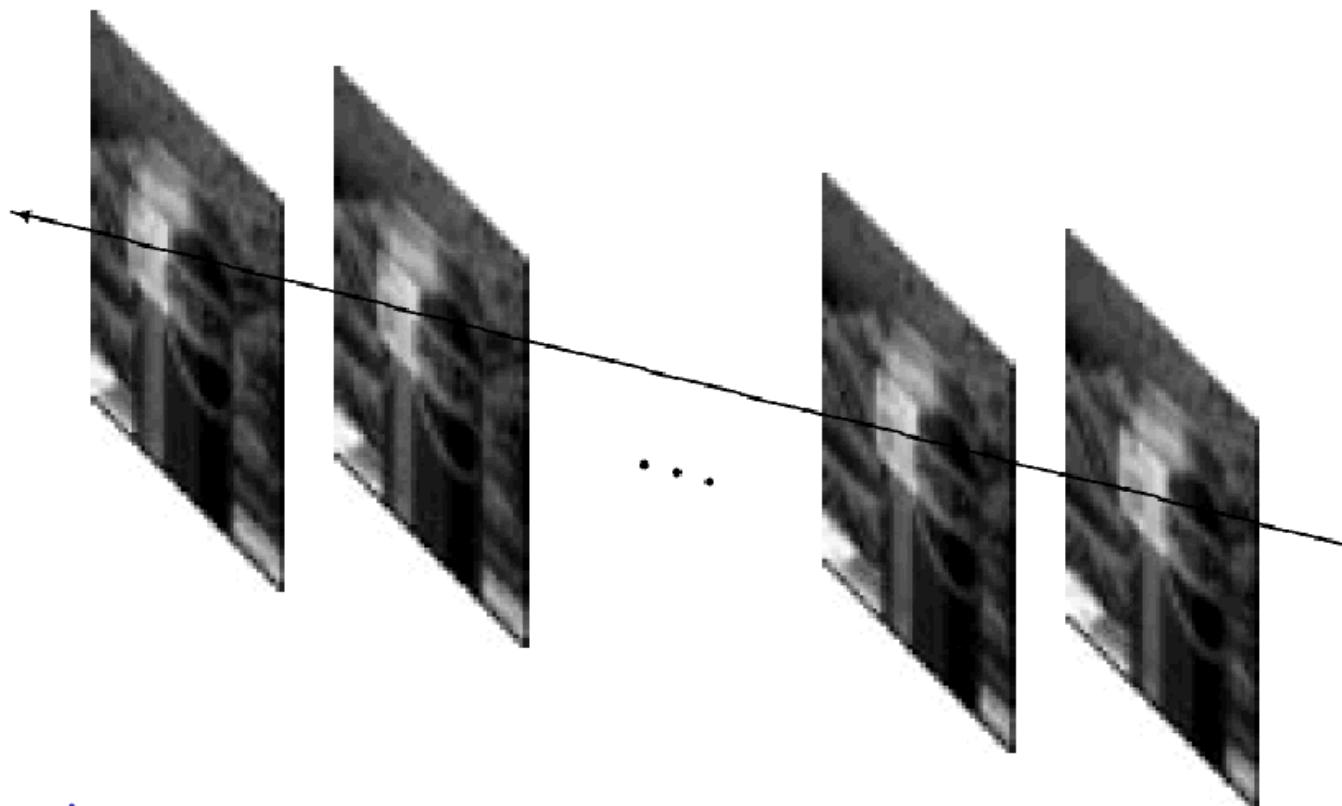
Assumption

Image measurements (e.g. brightness) in a small region remain the same although their location may change.

$$I(x + u, y + v, t + 1) = I(x, y, t)$$

(assumption)

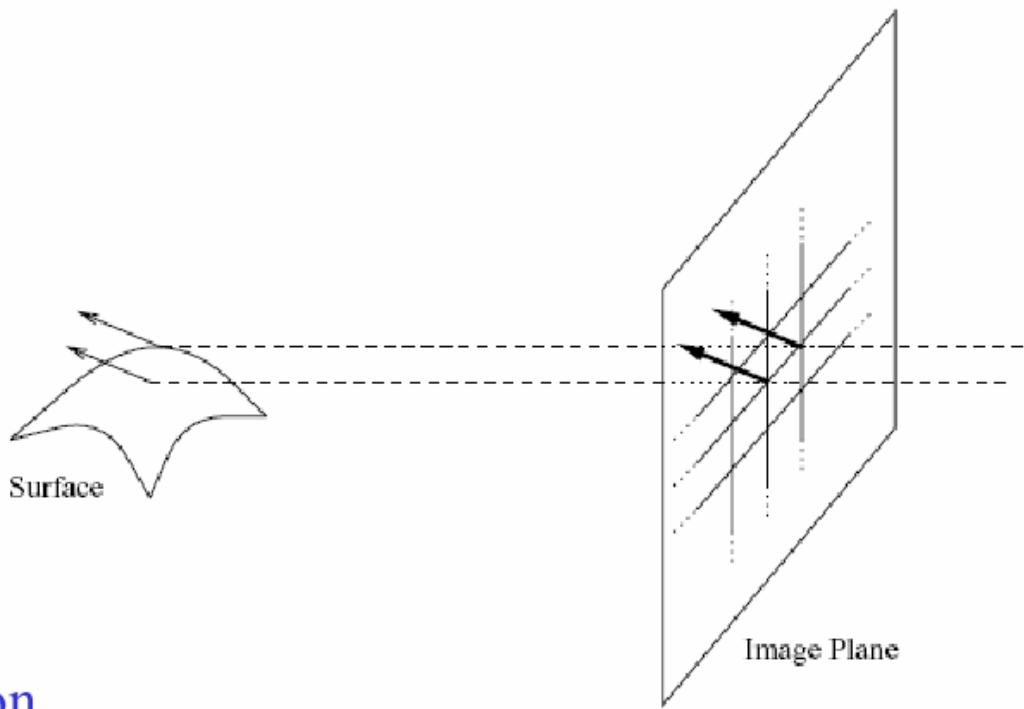
Key Assumptions: Small Motions



Assumption:

The image motion of a surface patch changes gradually over time.

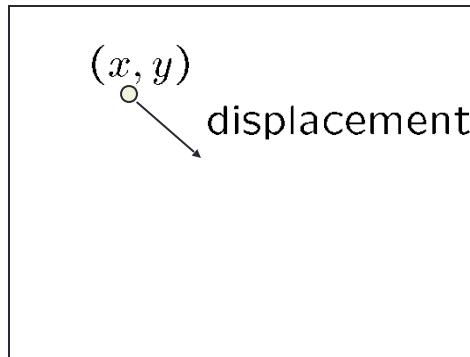
Key Assumptions: Spatial Coherence



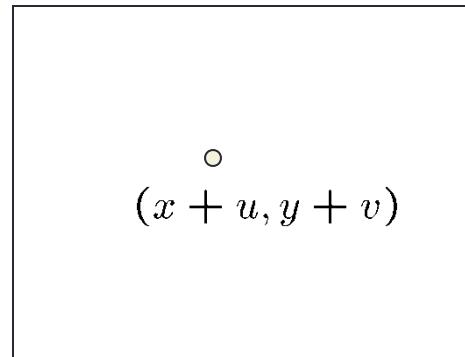
Assumption

- * Neighboring points in the scene typically belong to the same surface and hence typically have similar motions.
- * Since they also project to nearby points in the image, we expect spatial coherence in image flow.

Optical Flow Constraints (grayscale images)



$$I(x, y, t)$$



$$I(x, y, t + 1)$$

- Let's look at these constraints more closely

- Brightness constancy constraint (equation)

$$I(x, y, t) = I(x + u, y + v, t + 1)$$

- Small motion: (u and v are less than 1 pixel, or smoothly varying)
Taylor series expansion of I :

$$I(x + u, y + v, t + 1) = I(x, y, t) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v + \frac{\partial I}{\partial t} + o(1)$$

The Brightness Constancy Constraint

Can we use this equation to recover image motion (u, v) at each pixel?

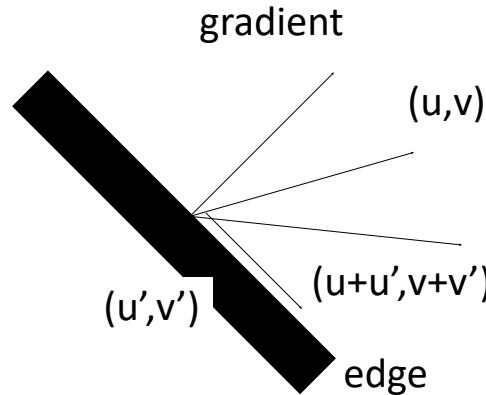
$$0 = I_t + \nabla I \cdot \langle u, v \rangle$$

- How many equations and unknowns per pixel?
 - One equation (this is a scalar equation!), two unknowns (u, v)

The component of the flow perpendicular to the gradient (i.e., parallel to the edge) cannot be measured

If (u, v) satisfies the equation,
so does $(u+u', v+v')$ if

$$\nabla I \cdot [u' \ v']^T = 0$$



Solving the Ambiguity...

B. Lucas and T. Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 674–679, 1981.

- How to get more equations for a pixel?
- **Spatial coherence constraint**

Assume the pixel's neighbors have the same (u, v)

If we use a 5×5 window, that gives us 25 equations per pixel

$$0 = I_t(\mathbf{p}_i) + \nabla I(\mathbf{p}_i) \cdot [u \ v]$$

$$\begin{bmatrix} I_x(\mathbf{p}_1) & I_y(\mathbf{p}_1) \\ I_x(\mathbf{p}_2) & I_y(\mathbf{p}_2) \\ \vdots & \vdots \\ I_x(\mathbf{p}_{25}) & I_y(\mathbf{p}_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(\mathbf{p}_1) \\ I_t(\mathbf{p}_2) \\ \vdots \\ I_t(\mathbf{p}_{25}) \end{bmatrix}$$

Matching Matches Across Images

- Overconstrained linear system

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

$A \quad d = b$
 $25 \times 2 \quad 2 \times 1 \quad 25 \times 1$

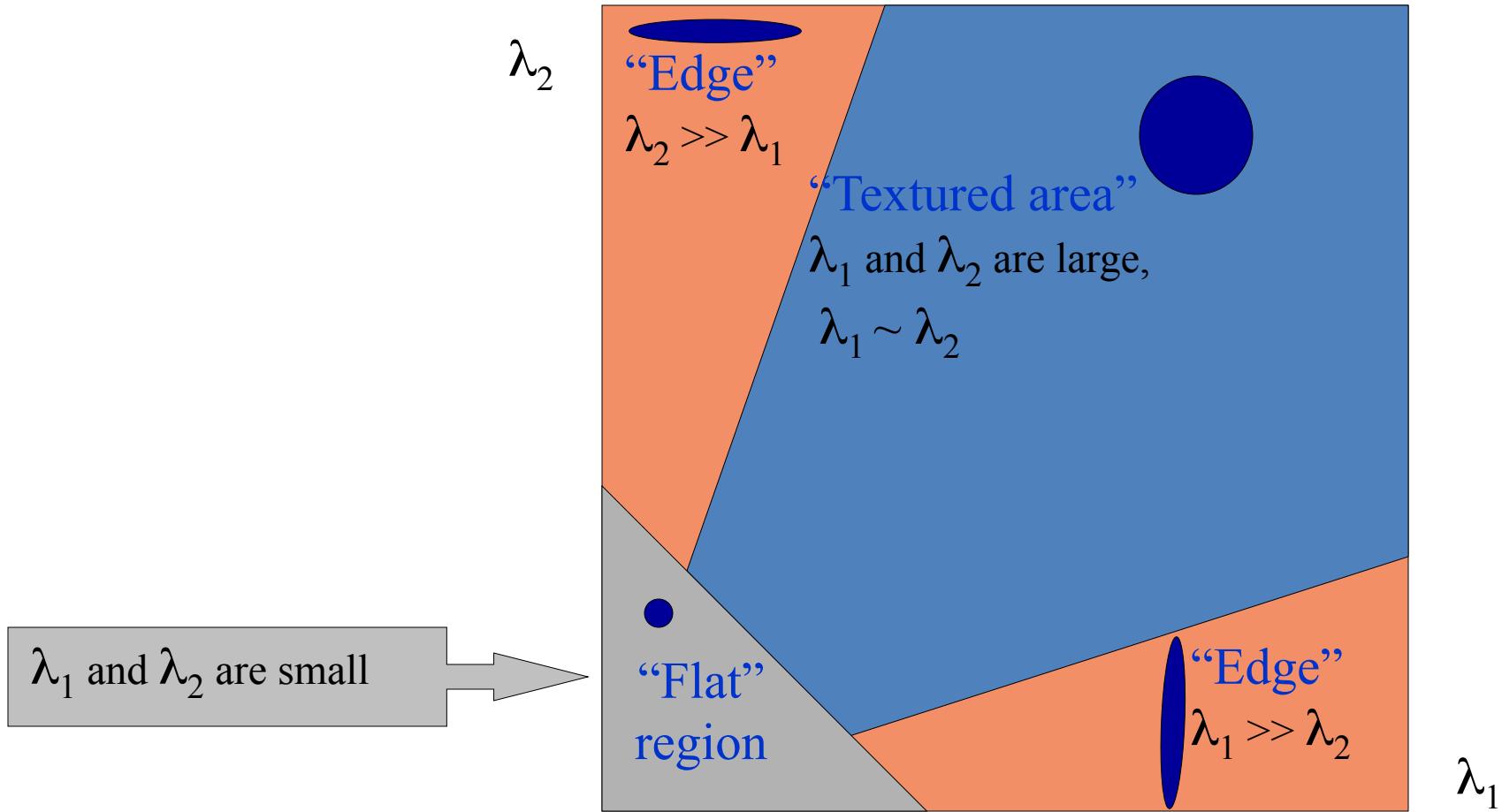
Least squares solution for d given by $(A^T A) d = A^T b$

$$A = \begin{bmatrix} (\nabla I(p_1))^T \\ \vdots \\ (\nabla I(p_n))^T \end{bmatrix} \Rightarrow A^T A = \sum_i \nabla I(p_i) (\nabla I(p_i))^T = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix}$$

The summations are over all pixels in the $K \times K$ window

Interpreting the Eigenvalues

Classification of image points using eigenvalues of $A^T A$



Harris Corner Detector

Cornerness

$$C = \lambda_1 \lambda_2 - \alpha (\lambda_1 + \lambda_2)^2$$

α : constant (0.04 to 0.06)

