# CS 554
# Computer Vision

## Homography
## & Image Stitching

**Hamdi Dibeklioğlu**

# Image Stitching

- We are given a bunch of photographs; how do we *stitch* them together?



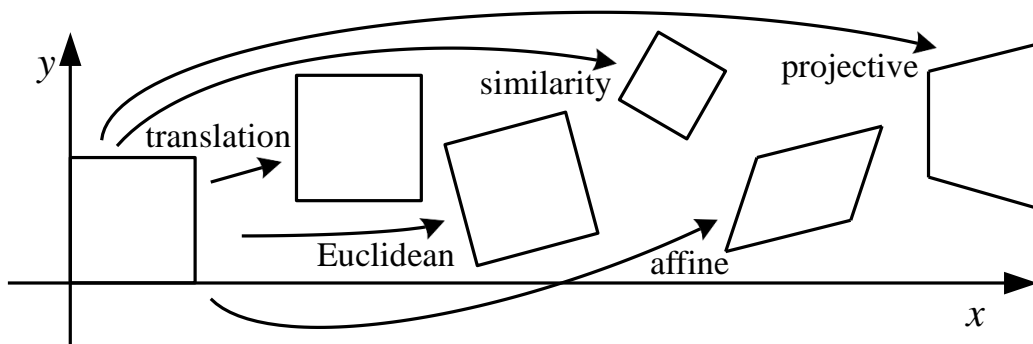- But first, we need to understand *non-linear least squares* and *homographies*

# Non-linear least squares problems

# Non-linear least squares

- Fitting a panography was *linear* in the transformation parameters:

$$E = \sum_i \|f(\mathbf{x}_i; \mathbf{p}) - \mathbf{x}_i'\|^2 = \sum_i \|\mathbf{x}_i + J(\mathbf{x}_i)\mathbf{p} - \mathbf{x}_i'\|^2$$

$$= \sum_i \|J(\mathbf{x}_i)\mathbf{p} - \Delta\mathbf{x}_i\|^2$$

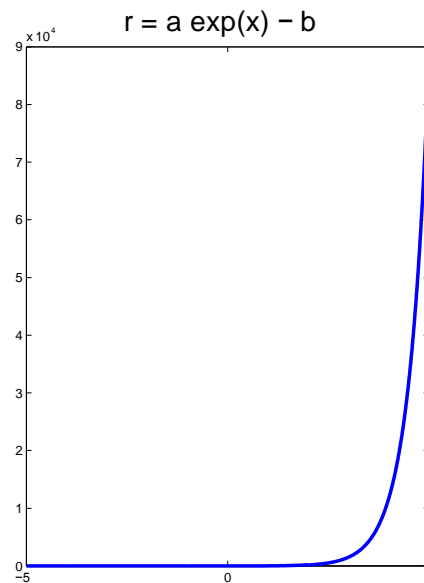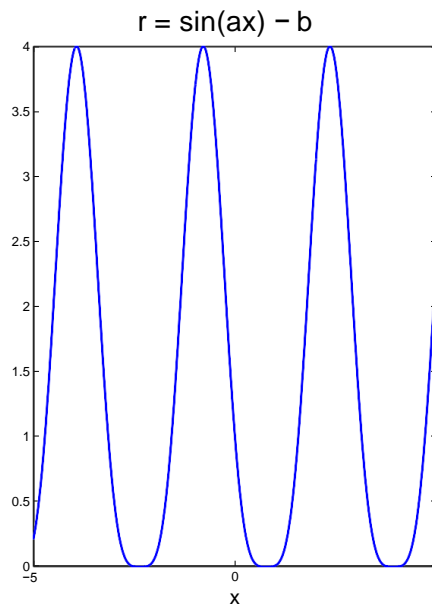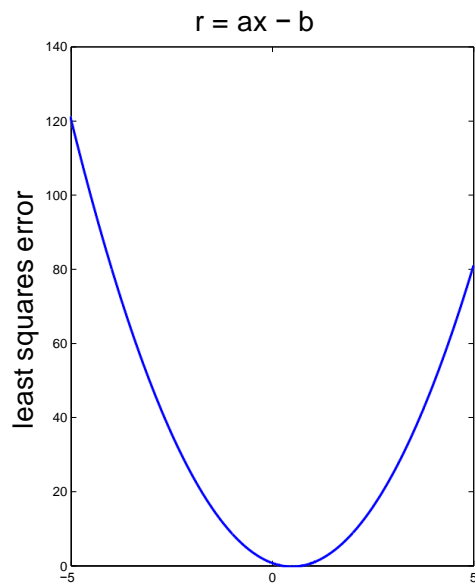- For more complex motion models, the transformation is *non-linear*:

# Non-linear least squares

- Consider the *non-linear least squares problem*:

$$g(\mathbf{x}) = \|f(\mathbf{x}; \mathbf{A}) - \mathbf{b}\|^2$$

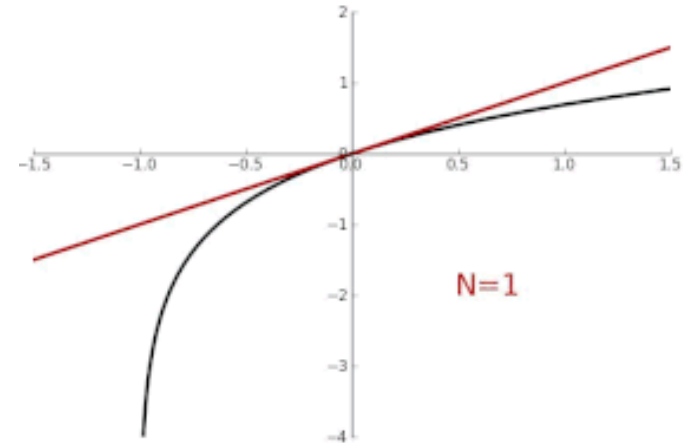- This problem is in general not *convex*; it may have multiple *local minima*:

# Taylor Expansion

- The *Taylor expansion* of the function $f(x)$ around $a$ is given by:

$$\sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x - a)^n$$

- Herein, $f^{(n)}$ denotes the *n*-th derivative



N=1

# Gauss-Newton Method

- Iteratively find parameter updates $\Delta \mathbf{x}$

# Gauss-Newton Method

- Iteratively find parameter updates $\Delta \mathbf{x}$

- Perform a *first-order Taylor expansion* of the residual around the current $\mathbf{x}$

$$\|f(\mathbf{x} - \Delta \mathbf{x}; \mathbf{A}) - \mathbf{b}\|^2 \approx \|f(\mathbf{x}; \mathbf{A}) + J(\mathbf{x})\Delta \mathbf{x} - \mathbf{b}\|^2$$

parameter update
objective

approximation around x

Jacobian

# Gauss-Newton Method

- Iteratively find parameter updates $\Delta \mathbf{x}$

- Perform a *first-order Taylor expansion* of the residual around the current $\mathbf{x}$

$$\|f(\mathbf{x} - \Delta \mathbf{x}; \mathbf{A}) - \mathbf{b}\|^2 \approx \|f(\mathbf{x}; \mathbf{A}) + J(\mathbf{x})\Delta \mathbf{x} - \mathbf{b}\|^2$$

parameter update
objective

approximation around x

Jacobian

- Note that the resulting residual approximation is linear in $\Delta \mathbf{x}$ :

  - The parameter update $\Delta \mathbf{x}$ may be obtained via linear least squares

# Gauss-Newton Method

- Writing down the linear least-squares solution for $\Delta \mathbf{x}$, we obtain:

$$\Delta \mathbf{x} = \left( J(\mathbf{x})^\top J(\mathbf{x}) \right)^{-1} J(\mathbf{x})^\top r(\mathbf{x})$$

↑

"Gauss-Newton approximation to Hessian"

- Gauss-Newton iteratively performs this update: $\mathbf{x} \leftarrow \mathbf{x} - \Delta \mathbf{x}$
- The Taylor expansion just became inaccurate! So iterate the whole process...

# Gauss-Newton Method

- Writing down the linear least-squares solution for $\Delta \mathbf{x}$, we obtain:

$$\Delta \mathbf{x} = \left( J(\mathbf{x})^\top J(\mathbf{x}) \right)^{-1} J(\mathbf{x})^\top r(\mathbf{x})$$

↑

"Gauss-Newton approximation to Hessian"

- Gauss-Newton iteratively performs this update: $\mathbf{x} \leftarrow \mathbf{x} - \Delta \mathbf{x}$

- The Taylor expansion just became inaccurate! So iterate the whole process...

- To implement, you only need to derive *Jacobian*: 
$$J(\mathbf{x}) = \begin{bmatrix} \dfrac{\partial f_1}{\partial x_1} & \cdots & \dfrac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \dfrac{\partial f_m}{\partial x_1} & \cdots & \dfrac{\partial f_m}{\partial x_n} \end{bmatrix}$$

# Newton's Method

- Perform a *second-order Taylor expansion* of $g(\mathbf{x})$ around $\mathbf{x}$:

$$g(\mathbf{x}) \approx \|r(\mathbf{x})\|^2 - 2J(\mathbf{x})r(\mathbf{x})\Delta\mathbf{x} + \left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^2\right](\Delta\mathbf{x})^2$$

with residuals: $r(\mathbf{x}) = f(\mathbf{x}; \mathbf{A}) - \mathbf{b}$

# Newton's Method

- Perform a *second-order Taylor expansion* of $g(\mathbf{x})$ around $\mathbf{x}$ :

$$g(\mathbf{x}) \approx \|r(\mathbf{x})\|^2 - 2J(\mathbf{x})r(\mathbf{x})\Delta\mathbf{x} + \left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^2\right](\Delta\mathbf{x})^2$$

with residuals: $r(\mathbf{x}) = f(\mathbf{x}; \mathbf{A}) - \mathbf{b}$

- This looks a lot like a linear least-squares problem; set gradient to zero:

$$-2J(\mathbf{x})^{\mathrm{T}}r(\mathbf{x}) + 2\left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^{\mathrm{T}}J(\mathbf{x})\right]\Delta\mathbf{x} = 0$$

$$\left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^{\mathrm{T}}J(\mathbf{x})\right]\Delta\mathbf{x} = J(\mathbf{x})^{\mathrm{T}}r(\mathbf{x})$$

# Newton's Method

- Perform a *second-order Taylor expansion* of $g(\mathbf{x})$ around $\mathbf{x}$ :

$$g(\mathbf{x}) \approx \|r(\mathbf{x})\|^2 - 2J(\mathbf{x})r(\mathbf{x})\Delta\mathbf{x} + \left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^2\right] (\Delta\mathbf{x})^2$$

with residuals: $r(\mathbf{x}) = f(\mathbf{x}; \mathbf{A}) - \mathbf{b}$

- This looks a lot like a linear least-squares problem; set gradient to zero:

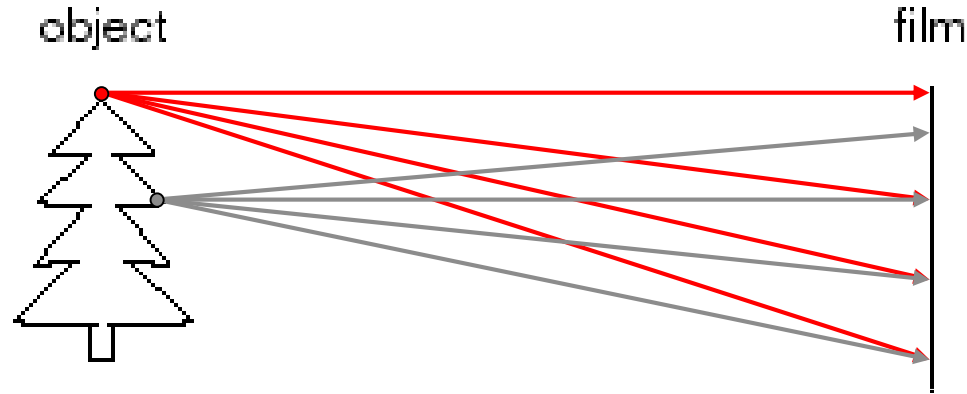$$-2J(\mathbf{x})^\mathrm{T} r(\mathbf{x}) + 2\left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^\mathrm{T} J(\mathbf{x})\right]\Delta\mathbf{x} = 0$$

$$\left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^\mathrm{T} J(\mathbf{x})\right]\Delta\mathbf{x} = J(\mathbf{x})^\mathrm{T} r(\mathbf{x})$$

- Note the similarity of the Newton update with the Gauss-Newton update:

$$\Delta\mathbf{x} = \left[H(\mathbf{x})r(\mathbf{x}) + J(\mathbf{x})^\mathrm{T} J(\mathbf{x})\right]^{-1} J(\mathbf{x})^\mathrm{T} r(\mathbf{x})$$

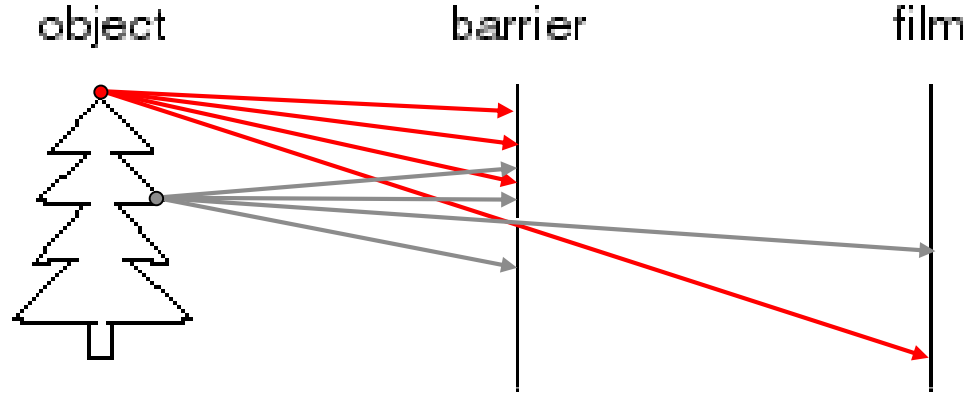# Homography

# Let's design a camera



Idea 1:  put a piece of film in front of an object
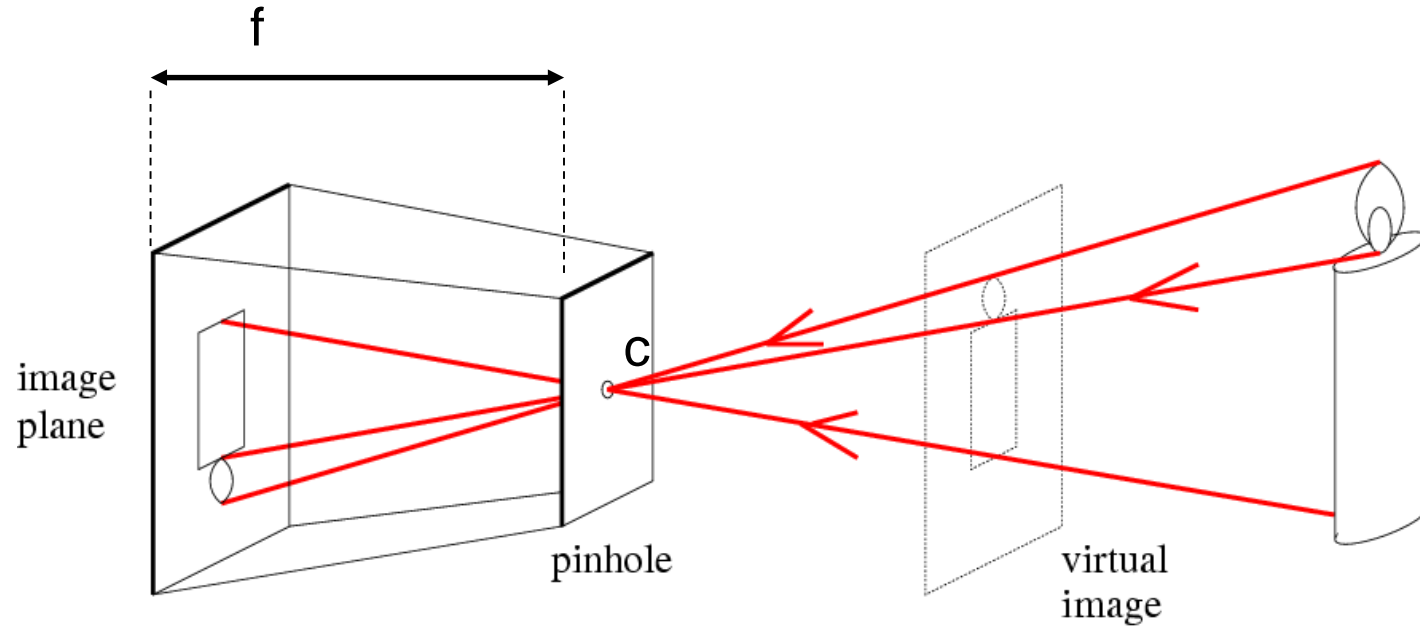
Do we get a reasonable image?

# Let's design a camera
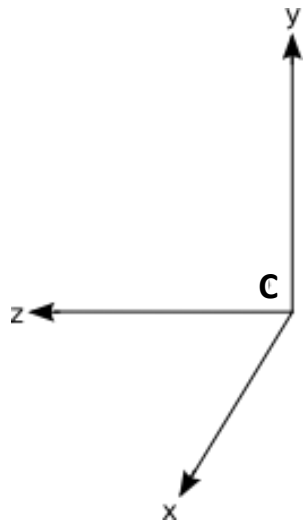


Add a barrier to block off most of the rays
  ◦ This reduces blurring
  ◦ The opening is known as the **aperture**

# Pinhole camera



f = focal length
c = center of projection

# Modeling projection



The coordinate system
- ◦ We will use the pin-hole model as an approximation
- ◦ Put the optical center (**C**enter of projection) at the origin
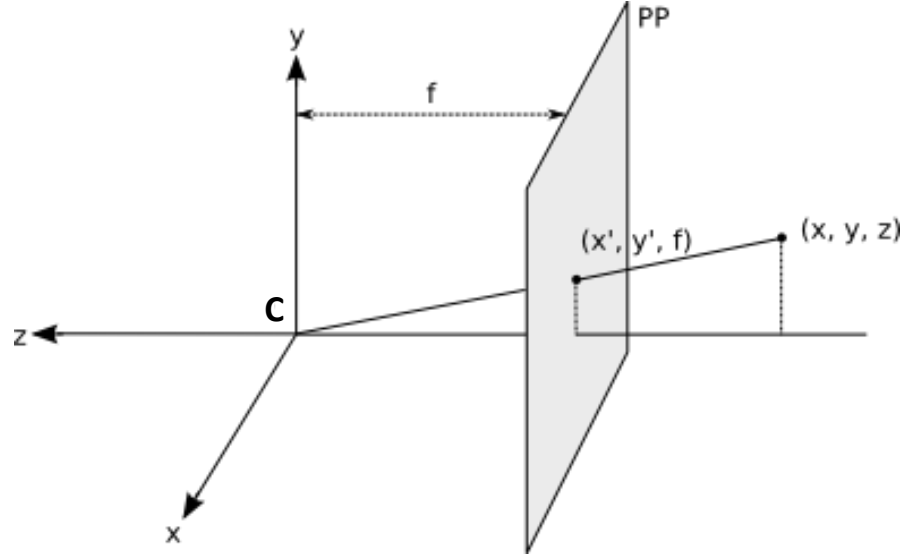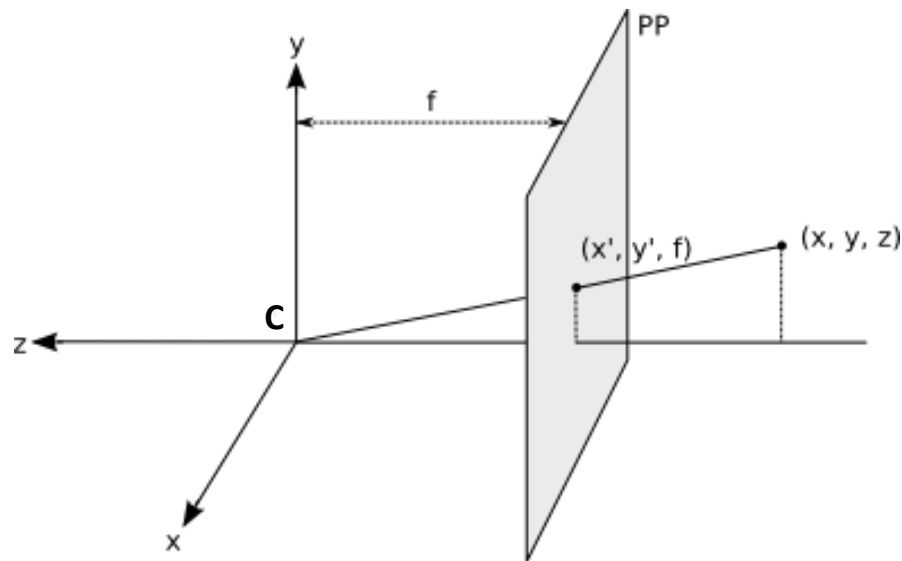- ◦ Where would you put the image (Projection Plane)?

# Modeling projection
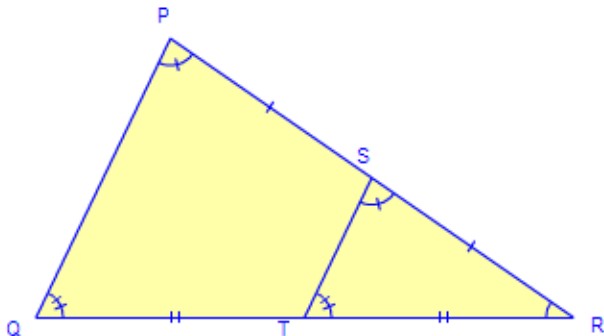


The coordinate system
- ◦ We will use the pin-hole model as an approximation
- ◦ Put the optical center (**C**enter of projection) at the origin
- ◦ Where would you put the image (**P**rojection **P**lane)?
  - ◦ In front of the camera

# Modeling projection



PP

y

f

(x', y', f)

(x, y, z)

C

z

x

## Projection equations

◦ Compute intersection with PP of ray from (x,y,z) to C in terms of *x, y, z,* and *f*
◦ Hint: use similar triangles; ST = PQ * (TR/QR)

$$(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z}, f)$$

- How do you get the 2D projection?

$$(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z})$$

Throw out the z coordinate

P

S

Q        T        R

# Modeling projection

- *Focal length* of the camera influences what is captured on the image plane:



$$(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z}, f)$$

- A large focal length implies small *field of view*, and vice versa

# Homogeneous coordinates

- Is this a linear transformation? $(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z})$

   No; division by z is non-linear

   Why?

   Definition of a linear function:

   $\text{f}(ax) = a\ \text{f}(x)$

   $\text{f}(x + y) = \text{f}(x) + \text{f}(y)$

Lets look at a small example for f = 1/z:

   1+2=3 → f(1) + f(2) = f(1+2) → $\frac{1}{1} + \frac{2}{2} = \frac{3}{3}$ → $1 + 1 = 1$ ???
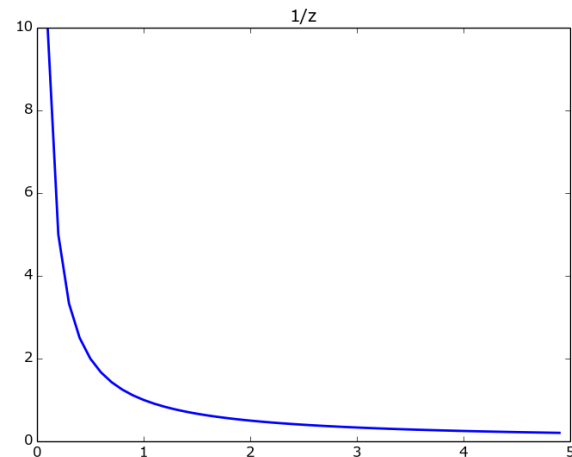
What does 1/z look like?



Trick: add one more coordinate to remember what 1/z was.

   1+2=3 → $(\frac{1}{1}, 1) + (\frac{2}{2}, \frac{1}{2}) = (\frac{3}{3}, \frac{1}{3})$

How would you use this coordinate to fix the problem?

Divide by the extra coordinate: $1 + 2 = 3$

# Homogeneous coordinates

- Is this a linear transformation? $(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z})$

  No; division by z is non-linear

Trick: add one more coordinate:

**2D:**

$$(x, y) \Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

homogeneous image
coordinates

**3D:**

$$(x, y, z) \Rightarrow \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix}$$

homogeneous scene
coordinates

Converting *from* homogeneous coordinates

$$\begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{bmatrix} \Rightarrow (\tilde{x}/\tilde{w}, \tilde{y}/\tilde{w})$$

$$\begin{bmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{z} \\ \tilde{w} \end{bmatrix} \Rightarrow (\tilde{x}/\tilde{w}, \tilde{y}/\tilde{w}, \tilde{z}/\tilde{w})$$

# Homogeneous coordinates

- We are used to describing a location in *Cartesian coordinates*:

$$\mathbf{x} = \begin{bmatrix} x & y \end{bmatrix}^{\mathrm{T}} \qquad\qquad \mathbf{x} = \begin{bmatrix} x & y & z \end{bmatrix}^{\mathrm{T}}$$

- Alternatively, we can describe locations in *homogeneous coordinates*:

$$\tilde{\mathbf{x}} = \begin{bmatrix} \tilde{x} & \tilde{y} & \tilde{w} \end{bmatrix}^{\mathrm{T}} \qquad\qquad \tilde{\mathbf{x}} = \begin{bmatrix} \tilde{x} & \tilde{y} & \tilde{z} & \tilde{w} \end{bmatrix}^{\mathrm{T}}$$

- The corresponding Cartesian coordinates are given by:

$$\mathbf{x} = \begin{bmatrix} \tilde{x}/\tilde{w} & \tilde{y}/\tilde{w} \end{bmatrix}^{\mathrm{T}} \qquad \mathbf{x} = \begin{bmatrix} \tilde{x}/\tilde{w} & \tilde{y}/\tilde{w} & \tilde{z}/\tilde{w} \end{bmatrix}^{\mathrm{T}}$$

- Essentially, you can think of $\tilde{w}$ as a way to deal with object scale (*"disparity"*)

- Homogeneous coordinates are very useful when working with *perspective transformations* (*homographies*)

# Perspective Projection Matrix

Projection is a matrix multiplication using homogeneous coordinates: $(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z})$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/f & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z/f \end{bmatrix} \Rightarrow (x\frac{f}{z}, y\frac{f}{z})$$

divide by the third
coordinate

# Perspective Projection Matrix

Projection is a matrix multiplication using homogeneous coordinates: $(x, y, z) \rightarrow (x\frac{f}{z}, y\frac{f}{z})$

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1/f & 0 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} = \begin{bmatrix} x \\ y \\ z/f \end{bmatrix} \Rightarrow (x\frac{f}{z}, y\frac{f}{z})$$

In practice: split into lots of different coordinate transformations…

$$\begin{pmatrix} \text{2D point (3x1)} \end{pmatrix} = \begin{pmatrix} \text{Camera to pixel coord. trans. matrix (3x3)} \end{pmatrix} \begin{pmatrix} \text{Perspective projection matrix (3x4)} \end{pmatrix} \begin{pmatrix} \text{World to camera coord. trans. matrix (4x4)} \end{pmatrix} \begin{pmatrix} \text{3D point (4x1)} \end{pmatrix}$$

# Camera Matrix

In practice: split into lots of different coordinate transformations…

$$
\begin{bmatrix} \text{2D} \\ \text{point} \\ \text{(3x1)} \end{bmatrix} = \begin{bmatrix} \text{Camera to} \\ \text{pixel coord.} \\ \text{trans. matrix} \\ \text{(3x3)} \end{bmatrix} \begin{bmatrix} \text{Perspective} \\ \text{projection matrix} \\ \text{(3x4)} \end{bmatrix} \begin{bmatrix} \text{World to} \\ \text{camera coord.} \\ \text{trans. matrix} \\ \text{(4x4)} \end{bmatrix} \begin{bmatrix} \text{3D} \\ \text{point} \\ \text{(4x1)} \end{bmatrix}
$$

$$
\tilde{\mathbf{x}} = \mathbf{K} \left[ \mathbf{R} | \mathbf{t} \right] \mathbf{p} = \mathbf{P}\mathbf{p}
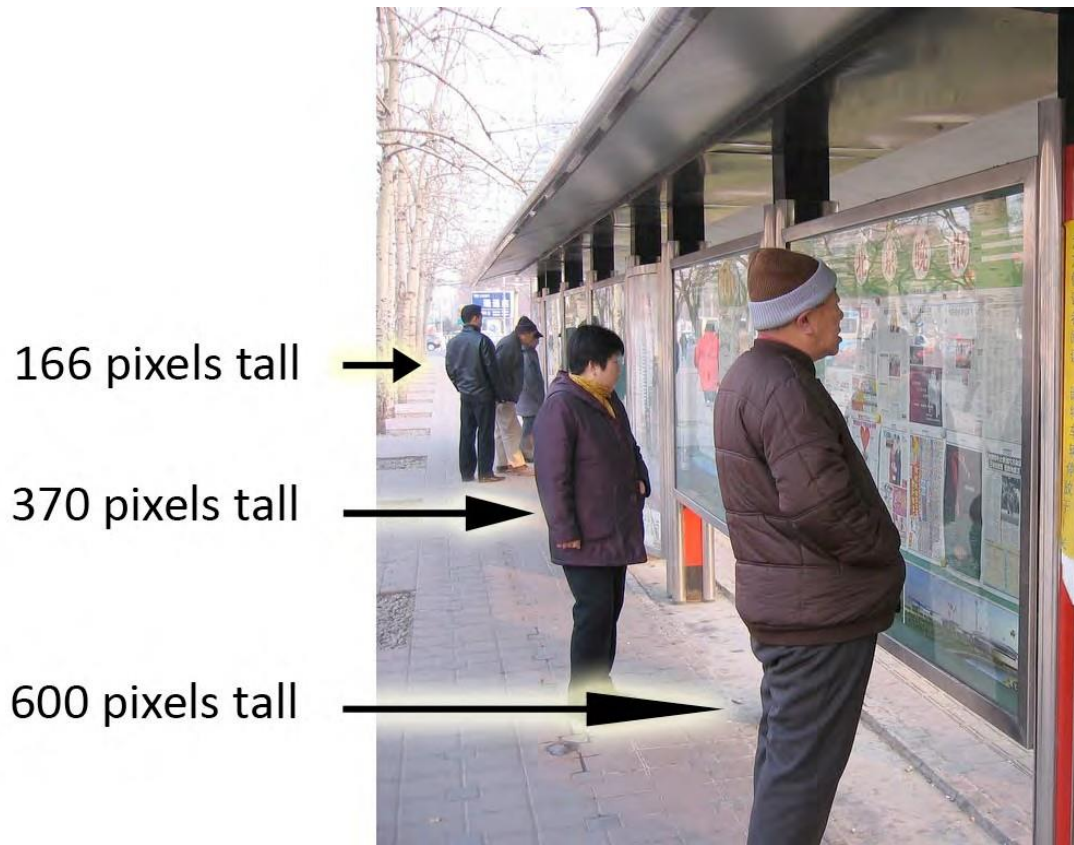$$

camera intrinsics
(calibration matrix)

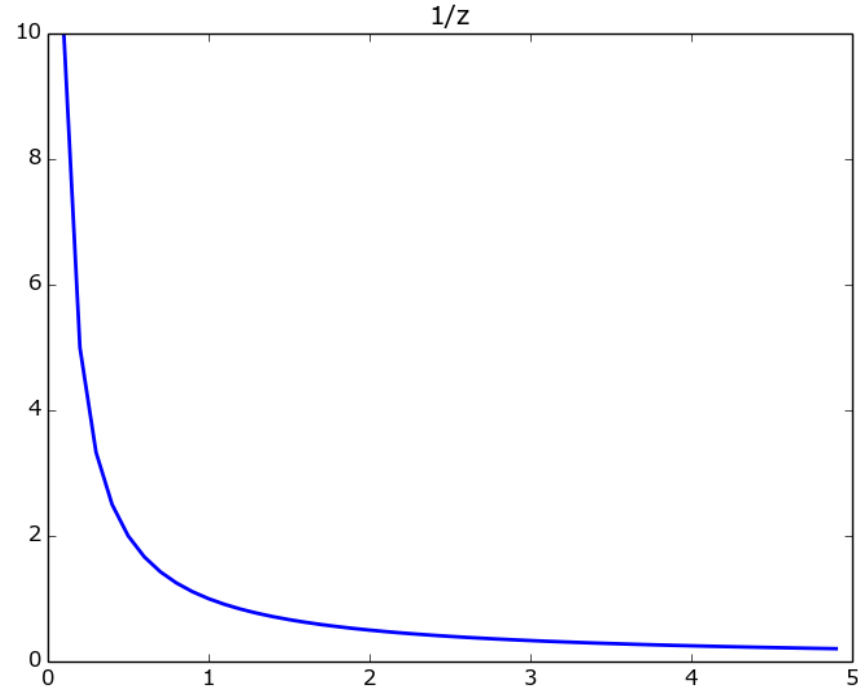camera extrinsics
(rotation + translation)

camera matrix

# How does a camera see a 3D point?

- The camera matrix is an example of a *projective transformation*:



166 pixels tall →
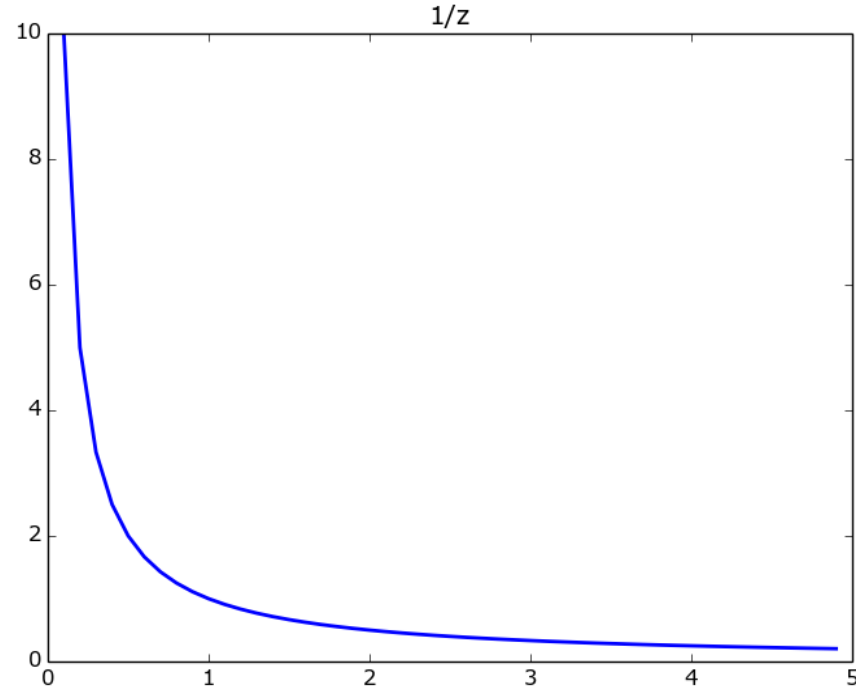
370 pixels tall →

600 pixels tall →

# How does a camera see a 3D point?

What does this curve tell you about perspective effects?

# How does a camera see a 3D point?

What does this curve tell you about perspective effects?



1/z

Play limited role if the object is "far away"

# How does a camera see a place?

- Assume we have two cameras with projection matrices $\tilde{\mathbf{P}}_0$ and $\tilde{\mathbf{P}}_1$

- Where does camera 1 see the point that camera 0 sees at $\tilde{\mathbf{x}}_0$ ?

$$\tilde{\mathbf{x}}_1 \sim \tilde{\mathbf{P}}_1 \tilde{\mathbf{P}}_0^{-1} \tilde{\mathbf{x}}_0 = \mathbf{M}_{10} \tilde{\mathbf{x}}_0$$

# How does a camera see a place?

- Assume we have two cameras with projection matrices $\tilde{\mathbf{P}}_0$ and $\tilde{\mathbf{P}}_1$

- Where does camera 1 see the point that camera 0 sees at $\tilde{\mathbf{x}}_0$ ?

$$\tilde{\mathbf{x}}_1 \sim \tilde{\mathbf{P}}_1 \tilde{\mathbf{P}}_0^{-1} \tilde{\mathbf{x}}_0 = \mathbf{M}_{10} \tilde{\mathbf{x}}_0$$

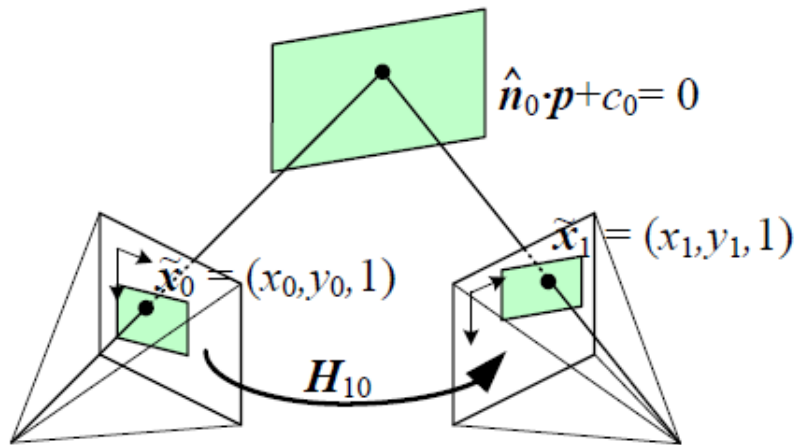- *Disparity* of point is irrelevant, so we can take the 3x3 sub-matrix of $\mathbf{M}_{10}$:

$$\tilde{\mathbf{x}}_1 \sim \tilde{\mathbf{H}}_{10} \tilde{\mathbf{x}}_0$$

- This is known as a *homography* $\tilde{\mathbf{H}}_{10}$ between the two cameras

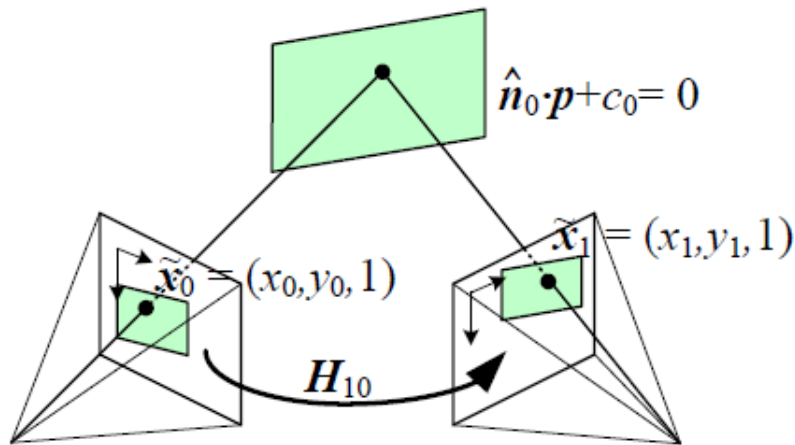* Notation in book is a bit sloppy here: page 56

# How does a camera see a place?

- Illustration of homography between two camera (for point and plane):



$$\hat{n}_0 \cdot p + c_0 = 0$$

$$\tilde{x}_1 = (x_1, y_1, 1)$$

$$\tilde{x}_0 = (x_0, y_0, 1)$$

$$H_{10}$$

$$H_{10} = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & 1 \end{bmatrix}$$

# How does a camera see a place?

- Illustration of homography between two camera (for point and plane):



$$\hat{n}_0 \cdot p + c_0 = 0$$

$$\tilde{x}_1 = (x_1, y_1, 1)$$

$$\tilde{x}_0 = (x_0, y_0, 1)$$

$$H_{10}$$

$$H_{10} = \begin{bmatrix} h_{00} & h_{01} & h_{02} \\ h_{10} & h_{11} & h_{12} \\ h_{20} & h_{21} & 1 \end{bmatrix}$$

- In Cartesian coordinates, the homography is given by:

$$x_1 = \frac{h_{00}x_0 + h_{01}y_0 + h_{02}}{h_{20}x_0 + h_{21}y_0 + 1} \qquad y_1 = \frac{h_{10}x_0 + h_{11}y_0 + h_{12}}{h_{20}x_0 + h_{21}y_0 + 1}$$

# Image stitching

# Image stitching

- Stitching algorithms typically have four main ingredients:

  - Method to determine *correspondences* between images

  - Model describing the set of possible *motions* between images (homography)

  - Algorithm to perform *alignment* of the images (bundle adjustment)

  - Algorithm that *composites* the images after alignment (blending; seam finding)
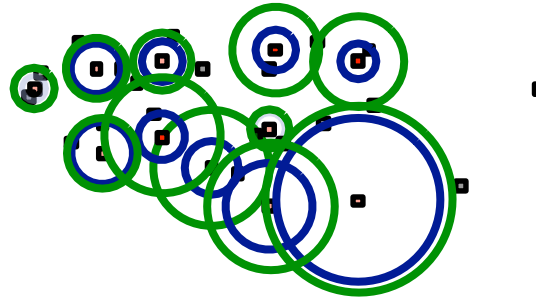
# Determining correspondences
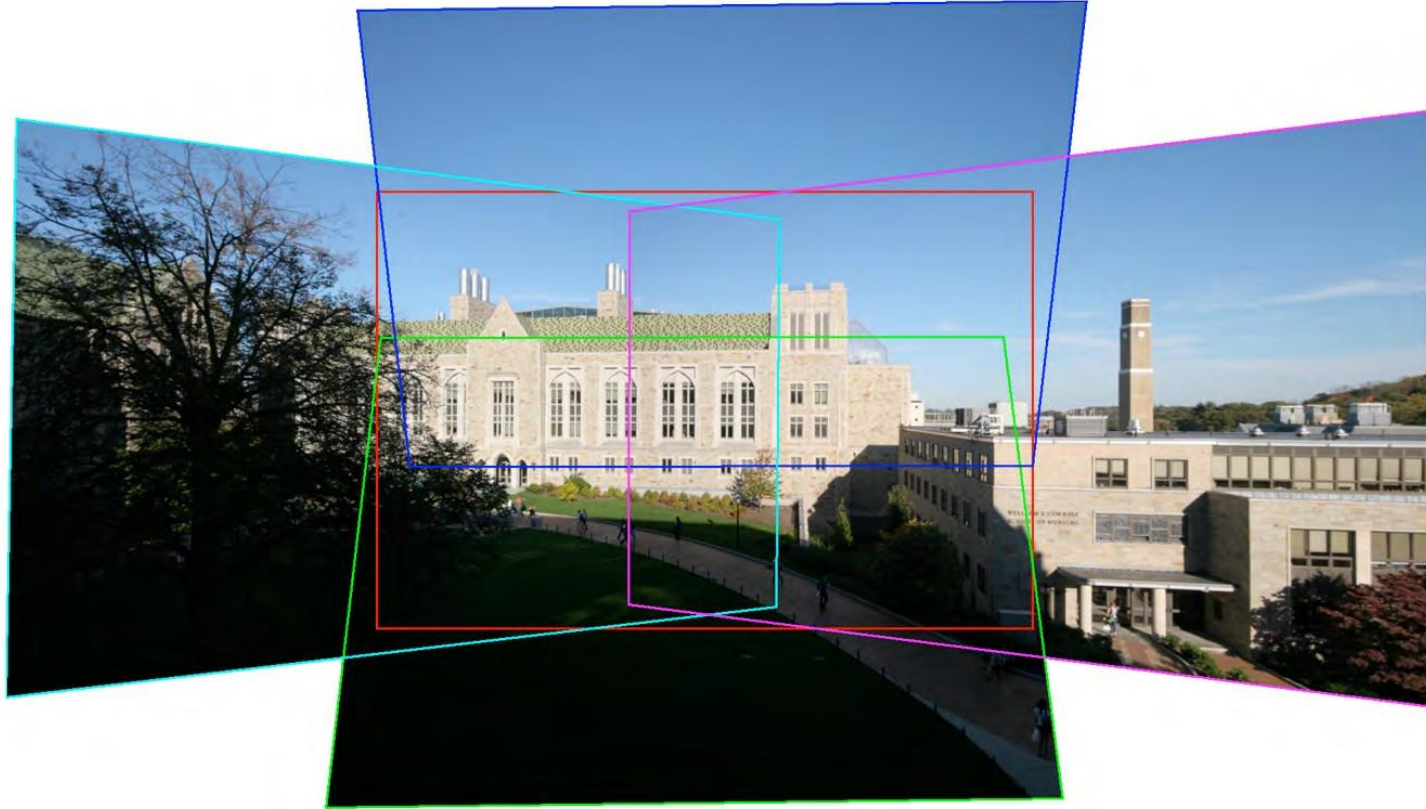
- Feature detection:

- Feature description:

- Feature matching:

# Motion model

- We will consider *homographies* because of their generality:

# Fitting a homography

- Recall the definition of a homography in Cartesian coordinates:

$$x' = f(x, y) = \frac{h_{00}x + h_{01}y + h_{02}}{h_{20}x + h_{21}y + 1} \qquad y' = g(x, y) = \frac{h_{10}x + h_{11}y + h_{12}}{h_{20}x + h_{21}y + 1}$$

- This makes the alignment objective a non-linear least squares problem:

$$\sum_i \left\| \begin{bmatrix} f(x_i, y_i; \mathbf{H}) \\ g(x_i, y_i; \mathbf{H}) \end{bmatrix} - \begin{bmatrix} x_i' \\ y_i' \end{bmatrix} \right\|^2$$

# Fitting a homography

- Recall the definition of a homography in Cartesian coordinates:

$$x' = f(x, y) = \frac{h_{00}x + h_{01}y + h_{02}}{h_{20}x + h_{21}y + 1} \qquad y' = g(x, y) = \frac{h_{10}x + h_{11}y + h_{12}}{h_{20}x + h_{21}y + 1}$$

- This makes the alignment objective a non-linear least squares problem:

$$\sum_i \left\| \begin{bmatrix} f(x_i, y_i; \mathbf{H}) \\ g(x_i, y_i; \mathbf{H}) \end{bmatrix} - \begin{bmatrix} x'_i \\ y'_i \end{bmatrix} \right\|^2$$

- Simply use Gauss-Newton's (or Newton's) method to solve this problem:

  - You may have to use RANSAC to deal with outliers!

# Remark on RANSAC

- The treatment of RANSAC we saw last week was somewhat simplified

- It was not random at all!

# Remark on RANSAC

• The treatment of RANSAC we saw last week was somewhat simplified

• It was not random at all! Full RANSAC actually works as follows:

    1) Select random subset from 50% of "best" matches as initial inliers

    2) Model is fitted to the *hypothetical inliers*

    *3)* Data are tested against the fitted model to determine hypothetical inliers

    4) Return to step 2) until sufficient points are classified as inliers
      (or fixed number of times)

    5) Return to step 1) a fixed number of times

# Fitting a homography

- Consider the following two images and SIFT matches between them:

# Fitting a homography

• Consider the following two images and SIFT matches between them:

# Fitting a homography

- Visualization of the inliers found by RANSAC:

# Fitting a homography

- Visualization of the homography found between the two images:
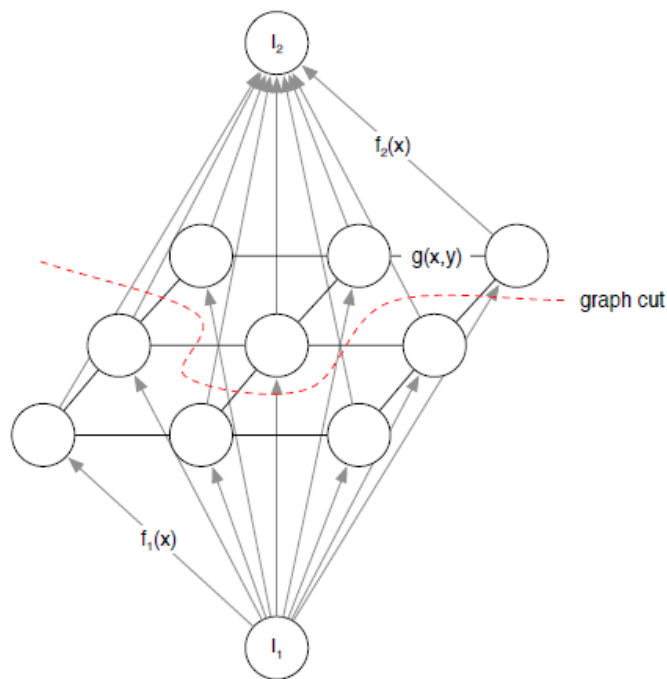
# Alpha blending

- *Alpha blending* averages the images; leads to *"ghosting"* and visible *seams*:

# Finding an optimal seam

• We would like to find a *seam* that minimizes the difference between images

• This can be formulated as a *graph min-cut* problem, as follows:



• When you want pixel $\mathbf{x}$ to be taken from $I_1$, set $f_2(\mathbf{x}) = \infty$ (and vice versa)

• In all other cases: $f_1(\mathbf{x}) = f_2(\mathbf{x}) = 0$

• Set, e.g., $g(x, y) = |I_1(\mathbf{x}) - I_2(\mathbf{y})| + |I_2(\mathbf{x}) - I_1(\mathbf{y})|$

• Efficient graph-cut algorithms exist

# Finding an optimal seam

- Ghosting can be prevented by finding a *seam* at which to switch images:

# Feathering

- Visibility of seam may be reduced by *"blurring"* weights of both images:

# Correcting for exposure differences

• Instead of using seam to select image, we use it to select the *image gradient*:
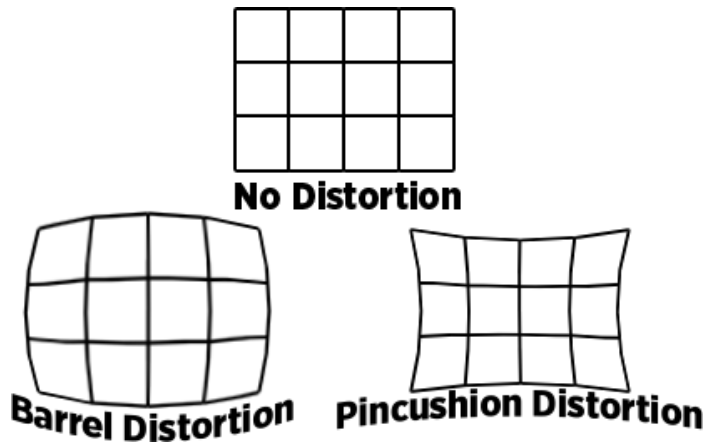
# Correcting for exposure differences

• The stitched image can then be recovered from this gradient:

# Lens distortion

- Homography model may be too simple when lens distortion is present:



- A common way of dealing with this is by incorporating a distortion model in the objective, and also minimizing w.r.t. the parameters of that model
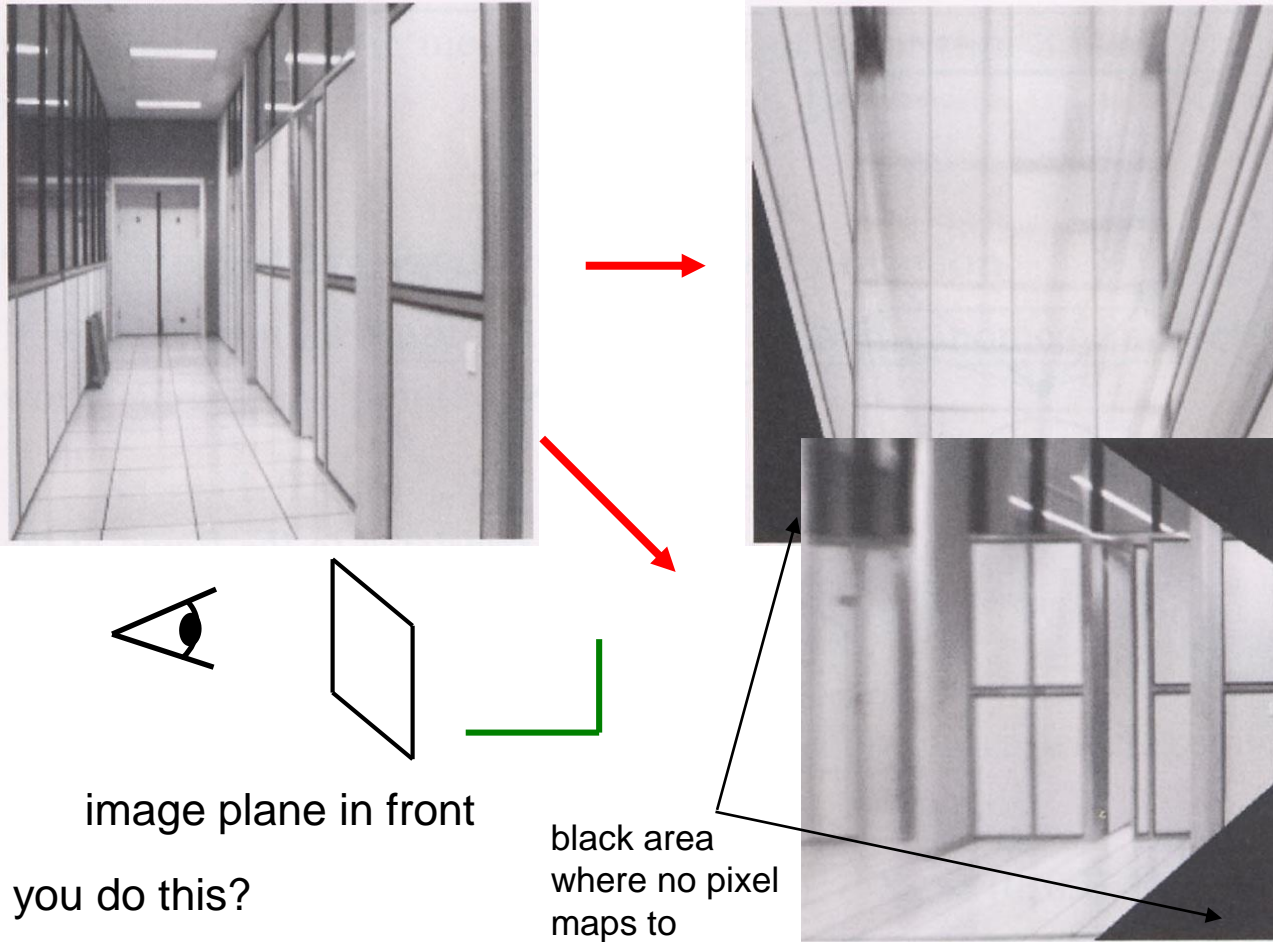
# Other applications of homographies
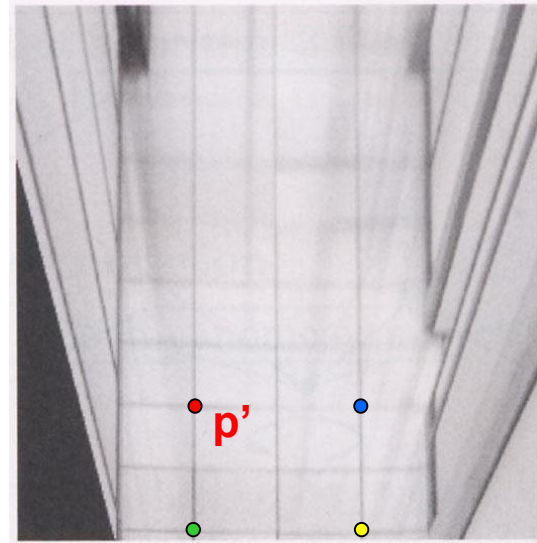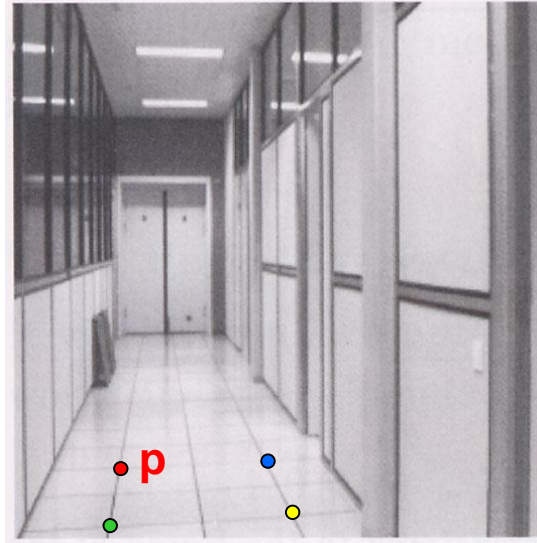
# Other applications of homographies

# Other applications: Image Warping



image plane in front

How would you do this?

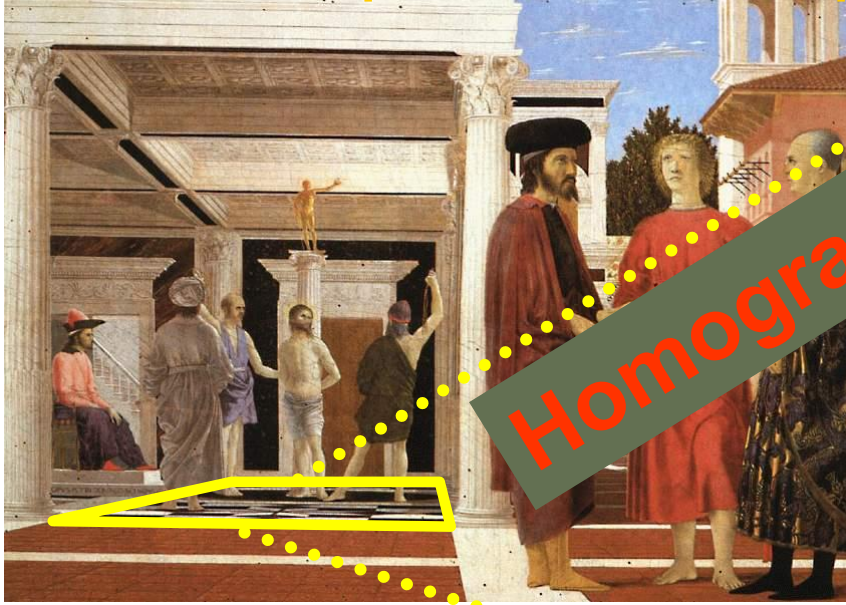black area where no pixel maps to

Source: Steve Seitz

# Other applications: Image rectification

# Other applications: Analyzing patterns/shapes

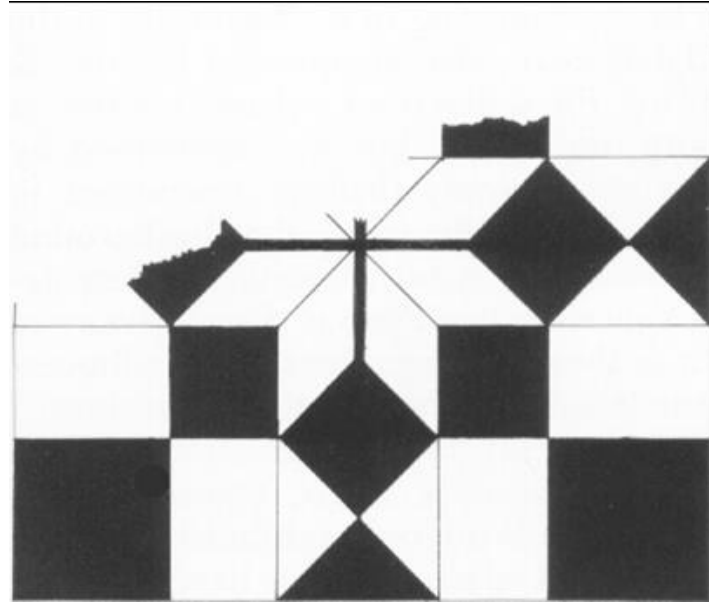**What is the shape of the b/w floor pattern?**

**Homography**



Slide: Criminisi

**The floor (enlarged)**

**Automatically rectified floor**

# Other applications: Analyzing patterns/shapes

**Automatic rectification**



**From Martin Kemp *The Science of Art*** 
***(manual reconstruction)***

# Other applications: Analyzing patterns/shapes



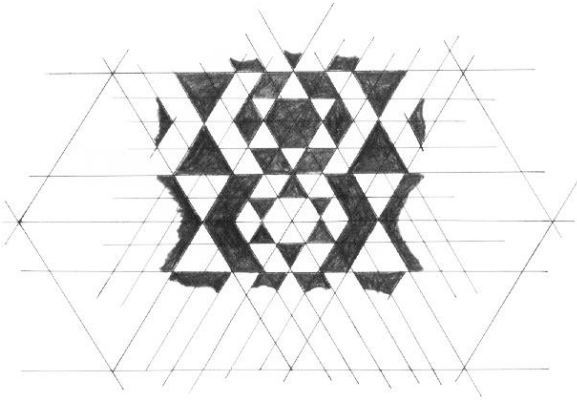*St. Lucy Altarpiece,* D. Veneziano



**Automatically rectified floor**

# Other applications: Analyzing patterns/shapes



**Automatic rectification**

**From Martin Kemp, *The Science of Art* (manual reconstruction)**

Reading material: Section 2 and 9 of Szeliski