

— CM50248 — 2017/2018 —

Visual Understanding 1

Dr Christian Richardt

Today's roadmap

1. Optical flow
2. Stereo matching
3. Structure-from-motion + bundle adjustment
4. Introduction to recognition for vision

Optical flow

Why estimate visual motion?

- Visual motion can be annoying
 - Camera instabilities, jitter
 - Measure it; remove it (stabilize)
- Visual motion indicates dynamics in the scene
 - Moving objects, behavior
 - Track objects and analyze trajectories
- Visual motion reveals spatial layout
 - Motion parallax

Classes of techniques

- Feature-based methods
 - Extract visual features (corners, textured areas) and track them
 - Sparse motion fields, but possibly robust tracking
 - Suitable especially when image motion is large (10s of pixels)
- Direct-methods
 - Directly recover image motion from spatio-temporal image brightness variations
 - Global motion parameters directly recovered without an intermediate feature motion calculation
 - Dense motion fields, but more sensitive to appearance variations
 - Suitable for video and when image motion is small (< 10 pixels)

Patch matching (revisited)

- How do we determine correspondences?
 - block matching or SSD (sum squared differences)

$$E(x, y; d) = \sum_{(x', y') \in N(x, y)} [I_L(x' + d, y') - I_R(x', y')]^2$$



The brightness constraint

Brightness constancy equation:

$$J(x, y) \approx I(x + u(x, y), y + v(x, y))$$

Or, equivalently, minimize :

$$E(u, v) = (J(x, y) - I(x + u, y + v))^2$$

Linearizing (assuming small u, v)
using Taylor series expansion:

$$J(x, y) \approx I(x, y) + I_x(x, y) \cdot u(x, y) + I_y(x, y) \cdot v(x, y)$$

Gradient constraint (or the optical flow constraint)

$$E(u, v) = (I_x \cdot u + I_y \cdot v + I_t)^2$$

Minimizing: $\frac{\partial E}{\partial u} = \frac{\partial E}{\partial v} = 0$

$$I_x(I_x u + I_y v + I_t) = 0$$

$$I_y(I_x u + I_y v + I_t) = 0$$

In general $I_x, I_y \neq 0$

Hence, $I_x \cdot u + I_y \cdot v + I_t \approx 0$

Patch translation [Lucas-Kanade]

Assume a single velocity for all pixels within an image patch

$$E(u, v) = \sum_{x, y \in \Omega} (I_x(x, y)u + I_y(x, y)v + I_t)^2$$

Minimizing

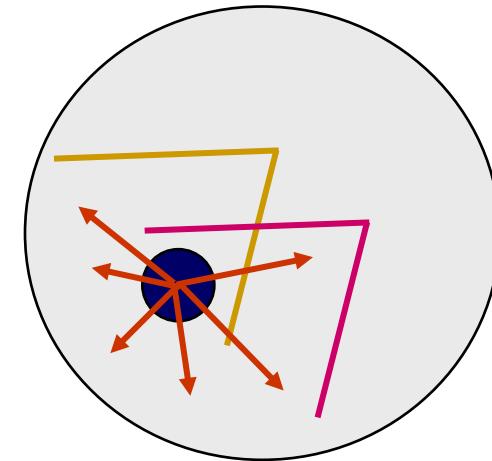
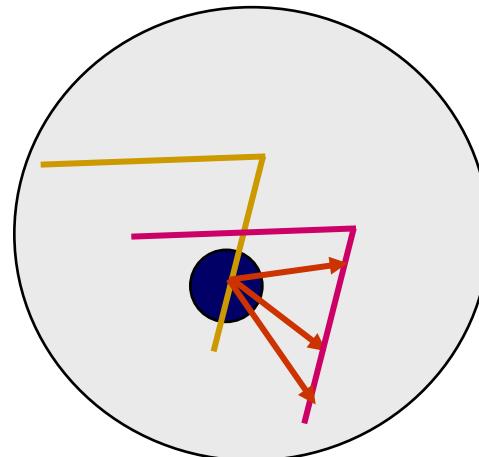
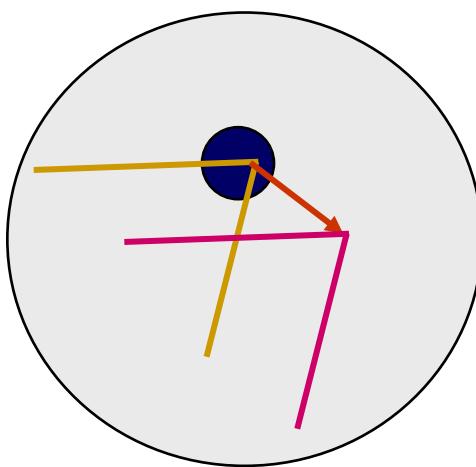
$$\begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{pmatrix} u \\ v \end{pmatrix} = - \begin{pmatrix} \sum I_x I_t \\ \sum I_y I_t \end{pmatrix}$$

$$\left(\sum \nabla I \nabla I^T \right) \vec{U} = - \sum \nabla I I_t$$

LHS: sum of the 2×2 outer product of the gradient vector

Local patch analysis

- How *certain* are the motion estimates?



The aperture problem

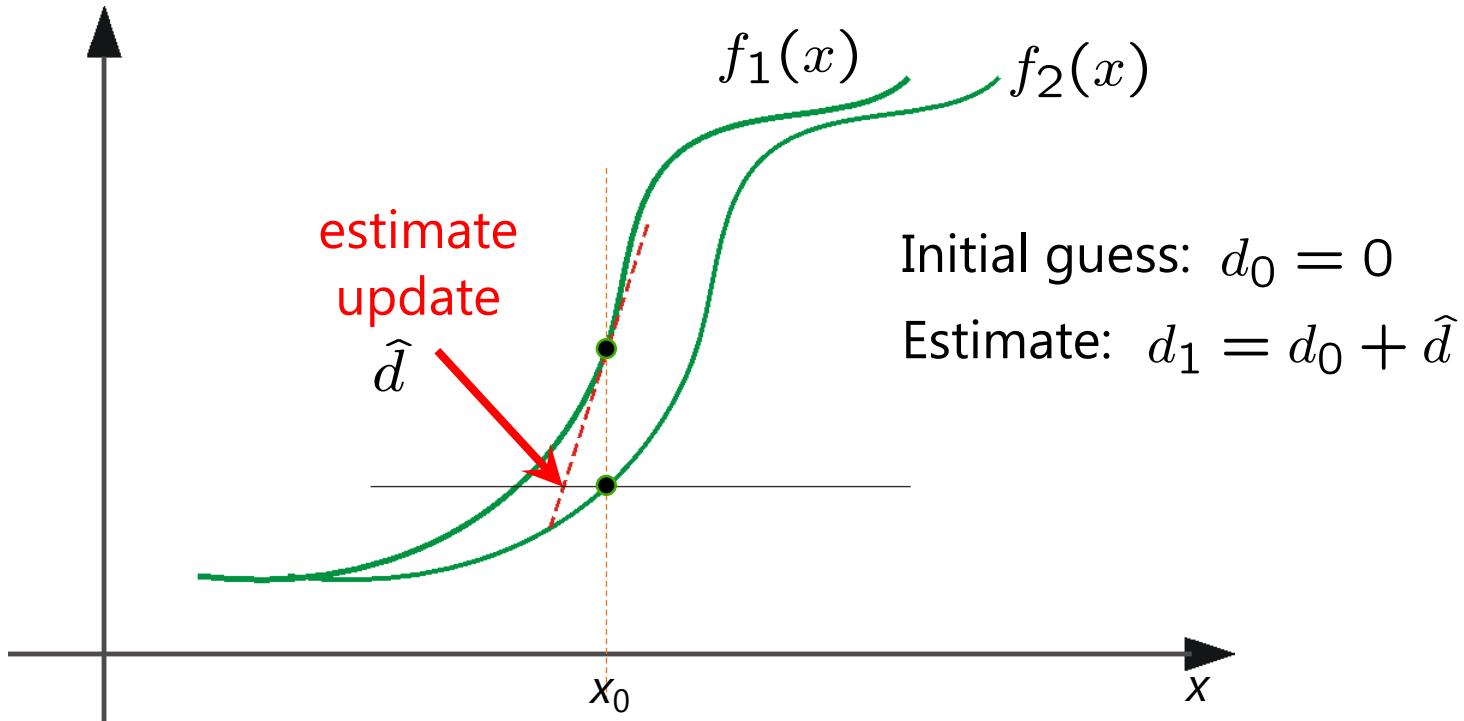
Let $M = \sum (\nabla I)(\nabla I)^T$ and $b = \begin{bmatrix} -\sum I_x I_t \\ -\sum I_y I_t \end{bmatrix}$

- Algorithm: At each pixel compute U by solving $MU=b$
- M is singular if all gradient vectors point in the same direction
 - e.g., along an edge
 - of course, trivially singular if the summation is over a single pixel or there is no texture
 - i.e., only *normal flow* is available (aperture problem)
- Corners and textured areas are OK

Iterative refinement

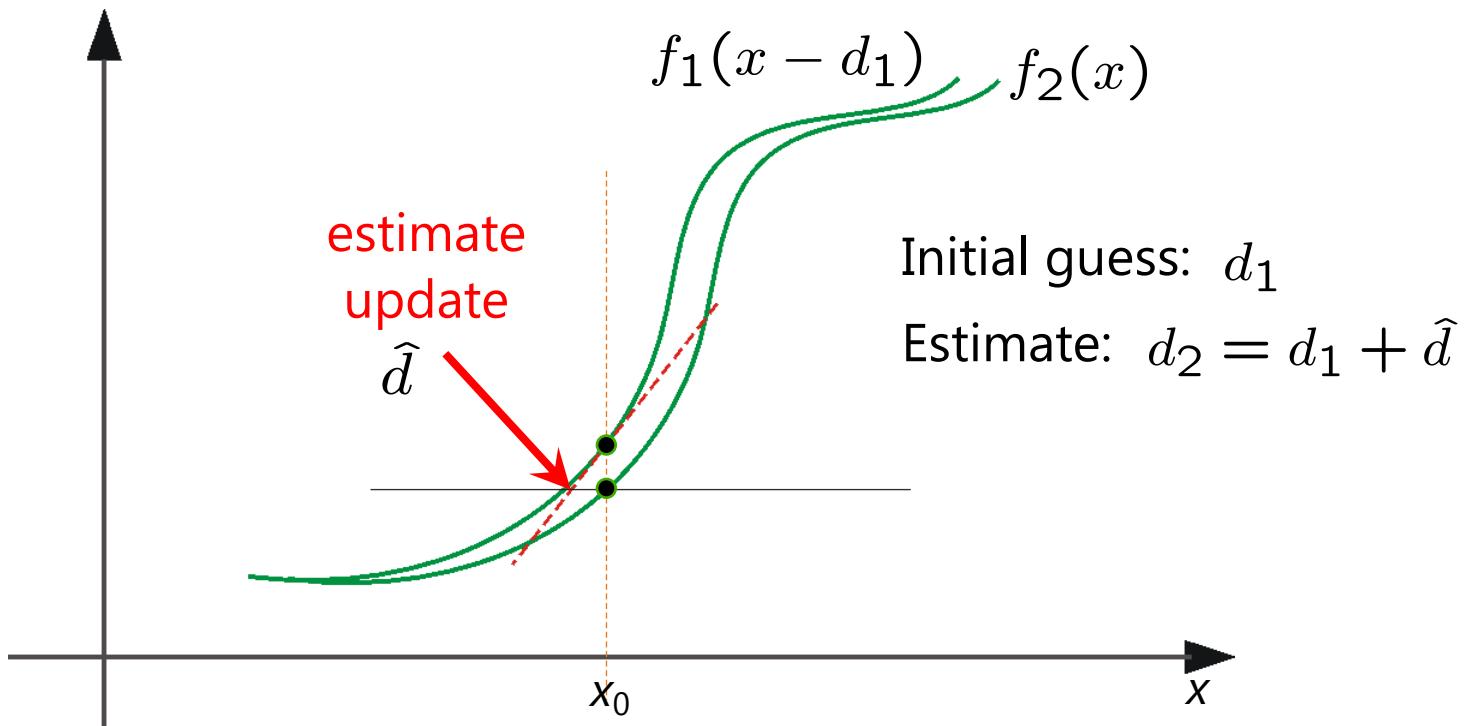
- Estimate velocity at each pixel using one iteration of Lucas–Kanade estimation
- Warp one image toward the other using the estimated flow field
 - (easier said than done)
- Refine estimate by repeating the process

Optical flow: iterative estimation

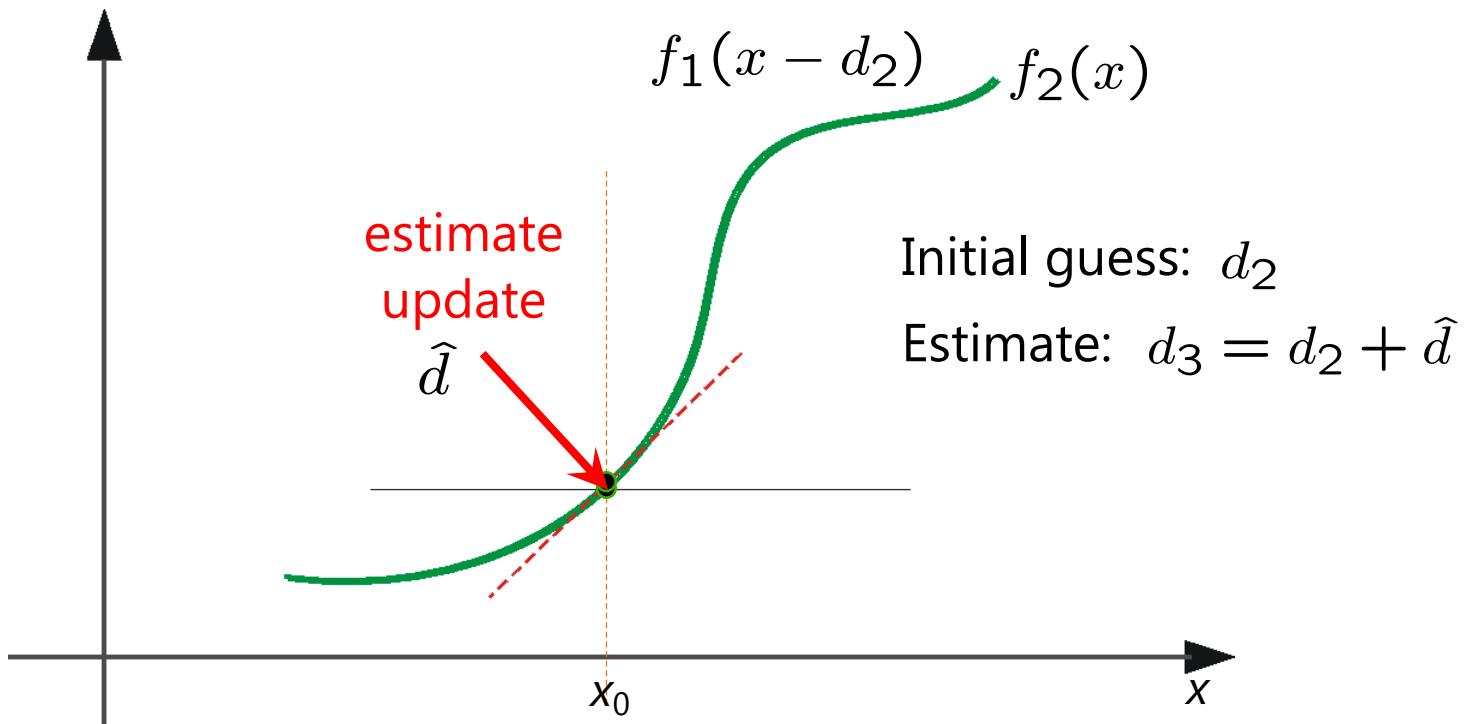


(using d for *displacement* here instead of u)

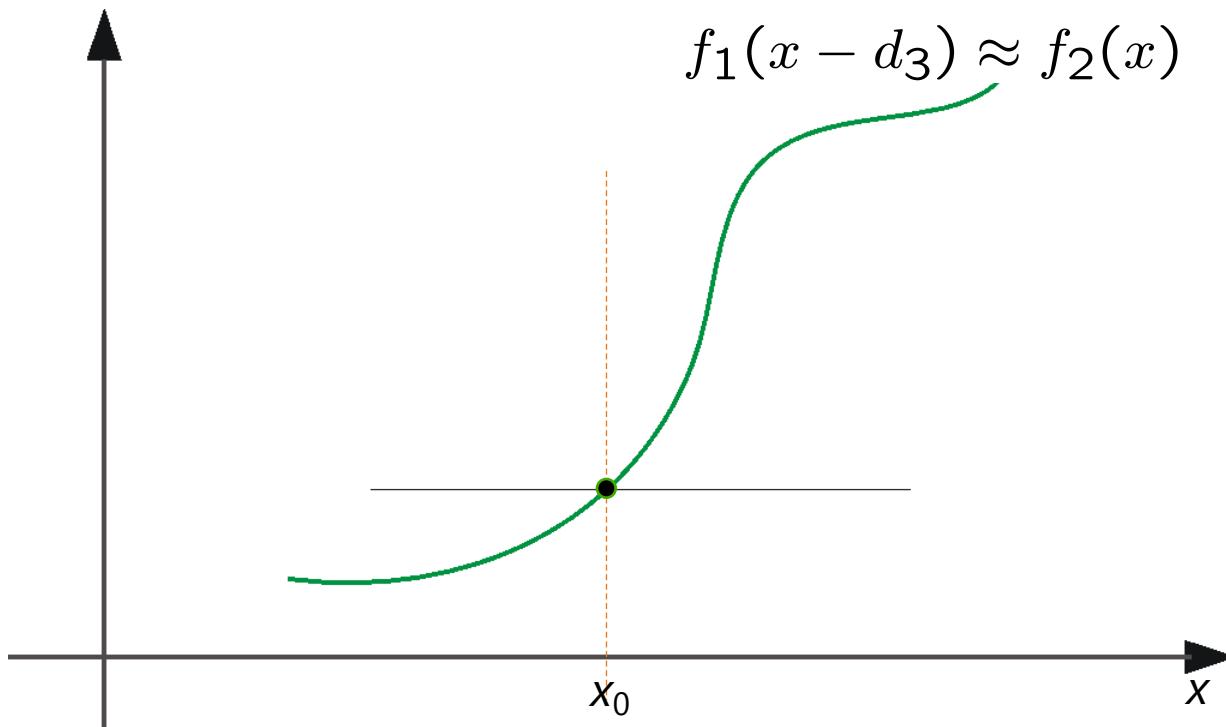
Optical flow: iterative estimation



Optical flow: iterative estimation



Optical flow: iterative estimation



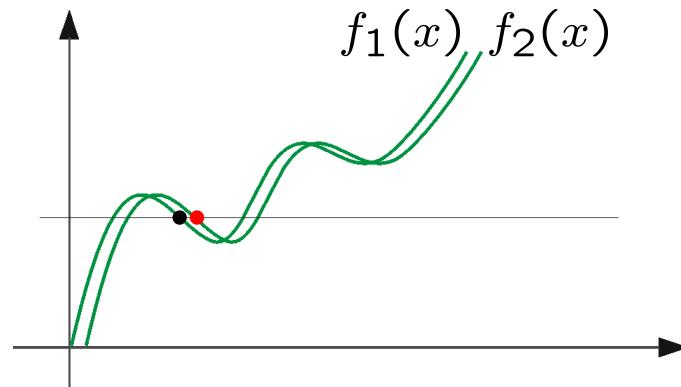
Optical flow: iterative estimation

- Some implementation issues:
 - Warping is not easy (ensure that errors in warping are smaller than the estimate refinement)
 - Warp one image, take derivatives of the other so you don't need to recompute the gradient after each iteration
 - Often useful to low-pass filter the images before motion estimation (for better derivative estimation, and linear approximations to image intensity)

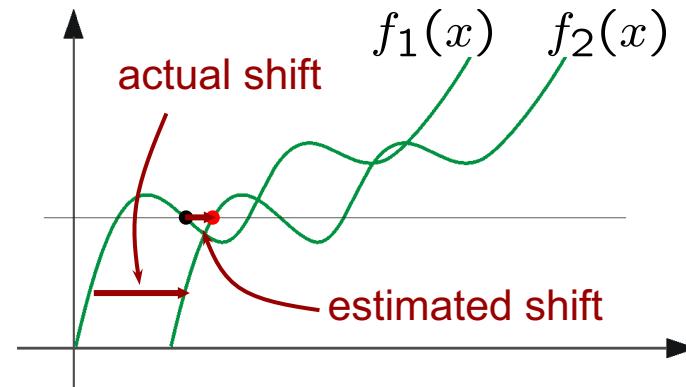
Optical flow: aliasing

Temporal aliasing causes ambiguities in optical flow because images can have many pixels with the same intensity.

How do we know which 'correspondence' is correct?



*nearest match is correct
(no aliasing)*



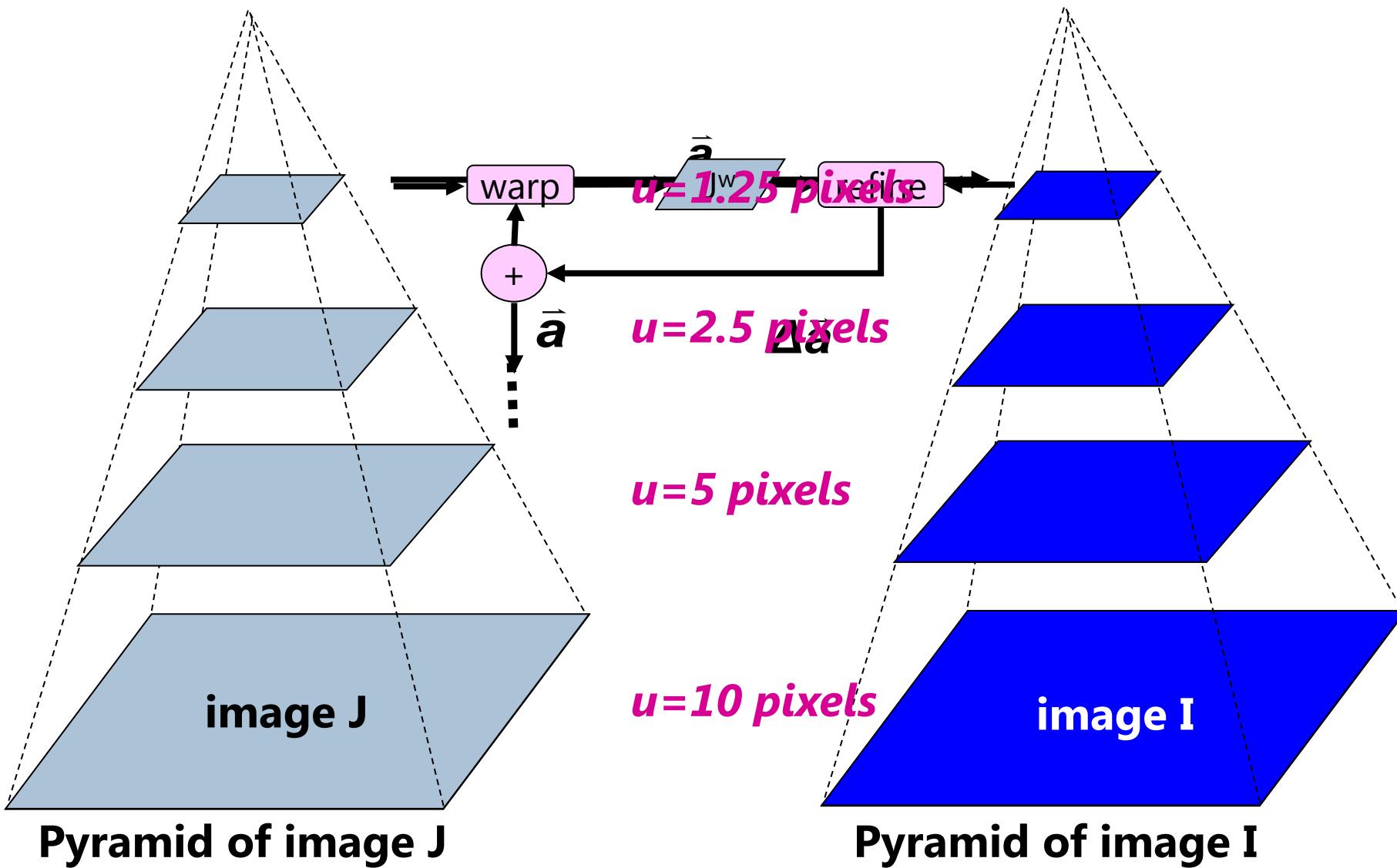
*nearest match is incorrect
(aliasing)*

To overcome aliasing: coarse-to-fine estimation.

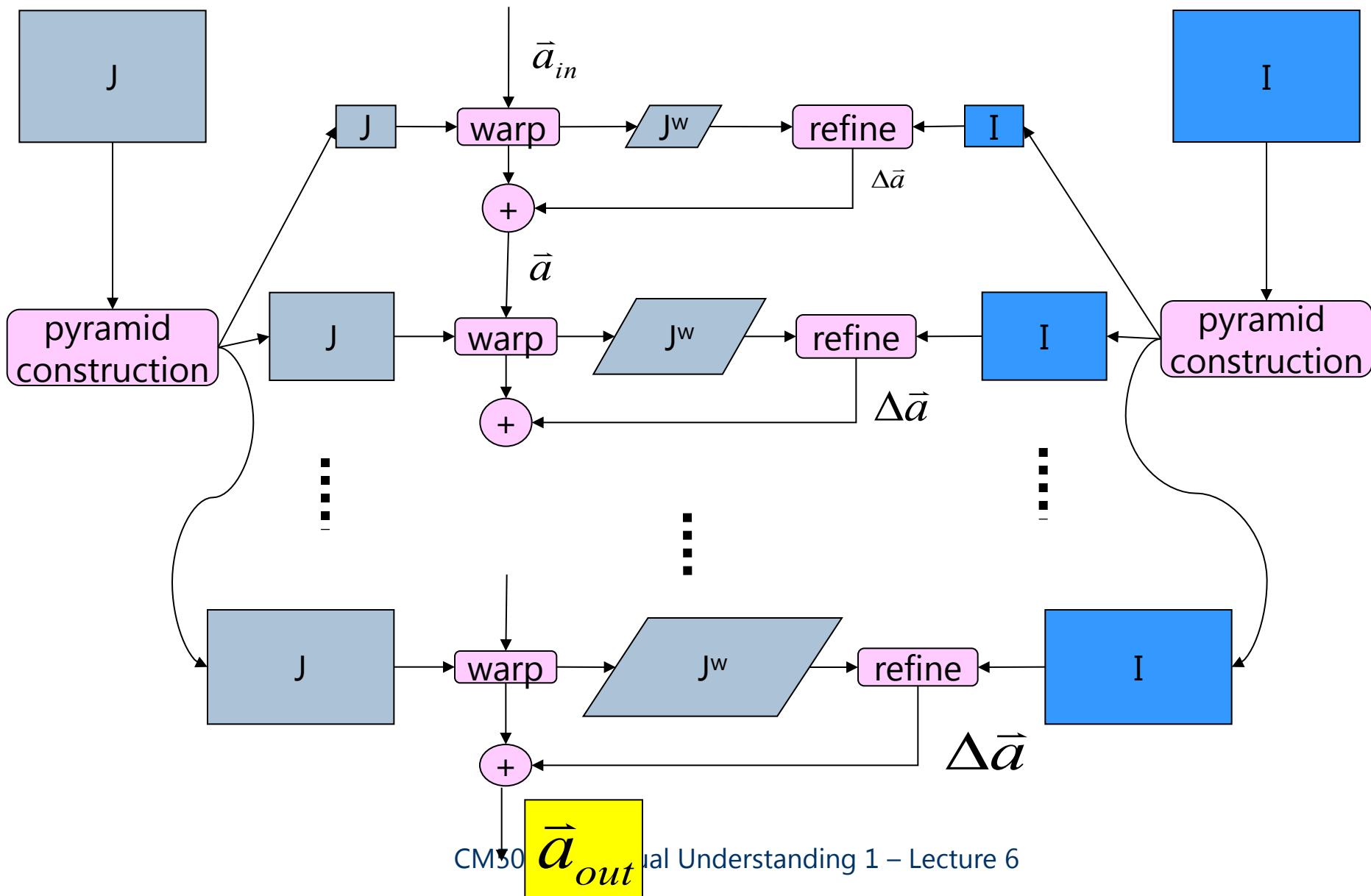
Limits of the gradient method

- Fails when intensity structure in window is poor
- Fails when the displacement is large (typical operating range is motion of 1 pixel)
 - Linearization of brightness is suitable only for small displacements
- Also, brightness is not strictly constant in images
 - actually less problematic than it appears, since we can pre-filter images to make them look similar

Coarse-to-fine estimation



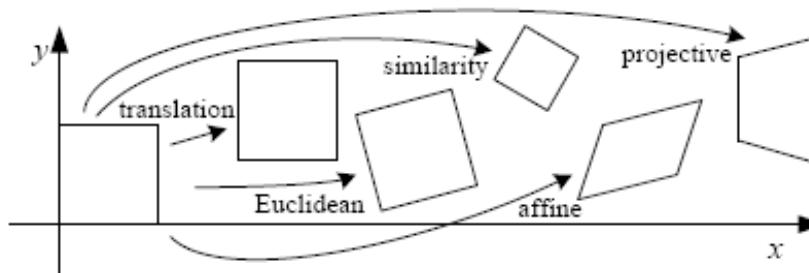
Coarse-to-fine estimation



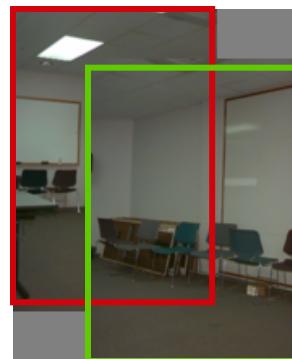
Global (parametric) motion models

- 2D models:
 - Affine
 - Quadratic
 - Planar projective transform (homography)
- 3D models:
 - Instantaneous camera motion models
 - Homography + epipole
 - Plane + Parallax

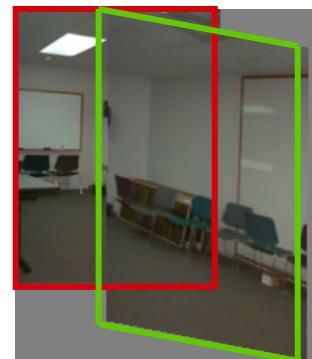
Motion models



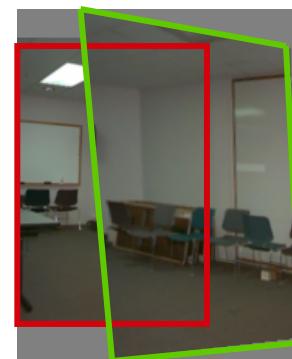
Translation



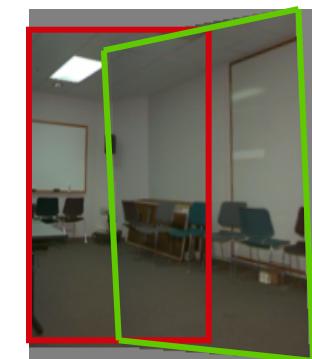
Affine



Perspective



3D rotation



2 unknowns

6 unknowns

8 unknowns

3 unknowns

Example: affine motion

$$u(x, y) = a_1 + a_2x + a_3y$$

$$v(x, y) = a_4 + a_5x + a_6y$$

Substituting into the brightness constancy equation:

$$I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t \approx 0$$

Each pixel provides 1 linear constraint in 6 *global* unknowns

Least Square Minimization (over all pixels):

$$Err(\vec{a}) = \sum [I_x(a_1 + a_2x + a_3y) + I_y(a_4 + a_5x + a_6y) + I_t]^2$$

Other 2D motion models

Quadratic – instantaneous approximation to planar motion

$$\begin{aligned} u &= q_1 + q_2x + q_3y + q_7x^2 + q_8xy \\ v &= q_4 + q_5x + q_6y + q_7xy + q_8y^2 \end{aligned}$$

Projective – exact planar motion

$$x' = \frac{h_1 + h_2x + h_3y}{h_7 + h_8x + h_9y}$$

$$y' = \frac{h_4 + h_5x + h_6y}{h_7 + h_8x + h_9y}$$

and

$$u = x' - x, \quad v = y' - y$$

3D motion models

Instantaneous camera motion:

Global parameters: $\Omega_X, \Omega_Y, \Omega_Z, T_X, T_Y, T_Z$

Local Parameter: $Z(x, y)$

$$u = -xy\Omega_X + (1+x^2)\Omega_Y - y\Omega_Z + (T_X - T_Zx)/Z$$

$$v = -(1+y^2)\Omega_X + xy\Omega_Y - x\Omega_Z + (T_Y - T_Zx)/Z$$

Homography+Epipole

Global parameters: $h_1, \dots, h_9, t_1, t_2, t_3$

Local Parameter: $\gamma(x, y)$

$$x' = \frac{h_1x + h_2y + h_3 + \gamma t_1}{h_7x + h_8y + h_9 + \gamma t_3}$$

$$y' = \frac{h_4x + h_5y + h_6 + \gamma t_1}{h_7x + h_8y + h_9 + \gamma t_3}$$

$$\text{and: } u = x' - x, \quad v = y' - y$$

Residual Planar Parallax Motion

Global parameters: t_1, t_2, t_3

Local Parameter: $\gamma(x, y)$

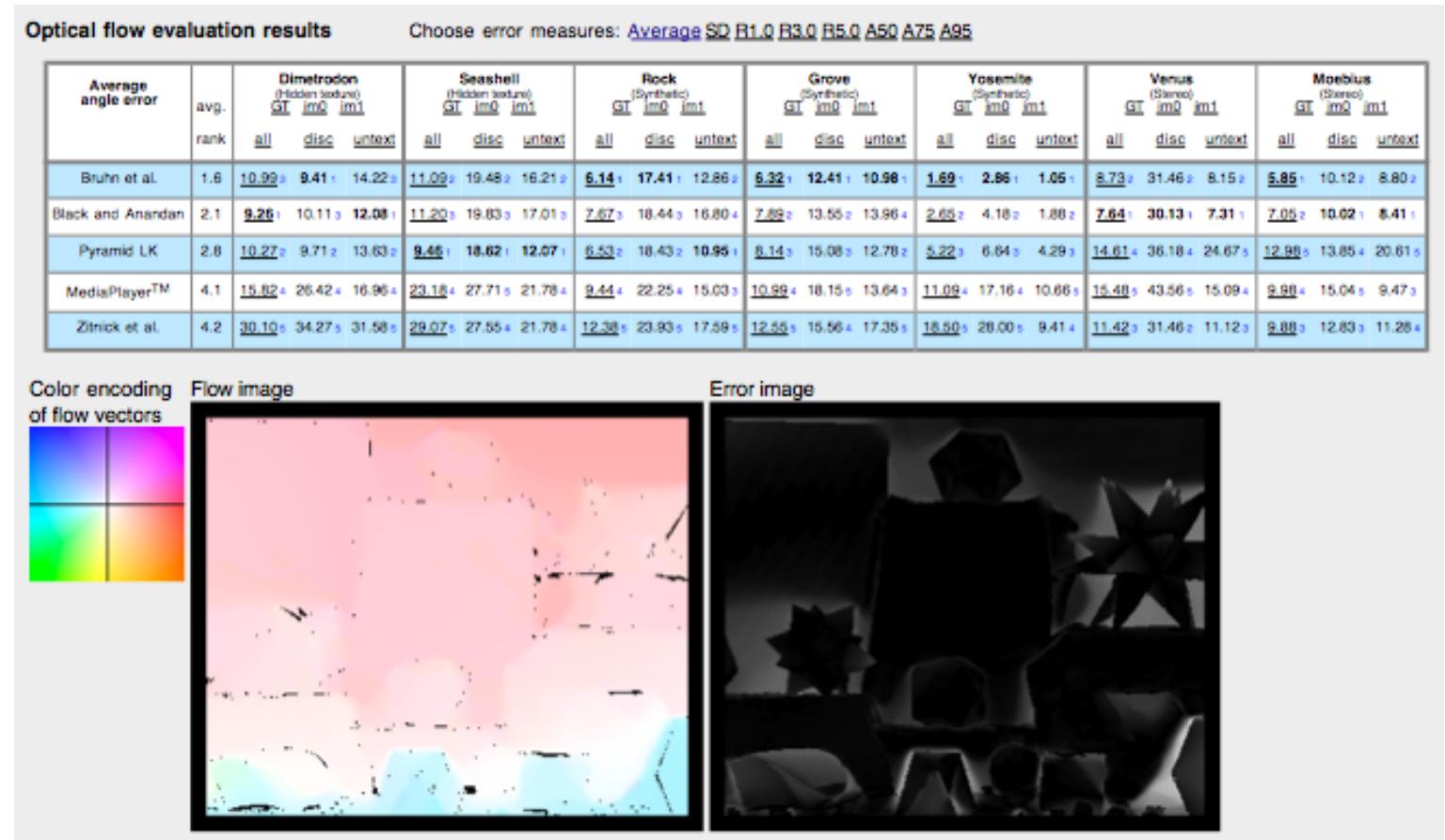
$$u = x^w - x = \frac{\gamma}{1+\gamma t_3} (t_3x - t_1)$$

$$v = y^w - x = \frac{\gamma}{1+\gamma t_3} (t_3y - t_2)$$

How well do these techniques work?

- A database and evaluation methodology for optical flow
- Simon Baker, Daniel Scharstein, JP Lewis, Stefan Roth, Michael Black and Richard Szeliski
- ICCV 2007

<http://vision.middlebury.edu/flow/>



Conclusions about optical flow

- Difficulty:
 - substantially more challenging than previous datasets
- Diversity:
 - Substantial variation in difficulty across the datasets
- Motion GT vs Interpolation:
 - Best algorithms for one are not the best for the other
- Comparison with Stereo:
 - Performance of existing flow algorithms appears weak

Stereo matching

Stereo matching

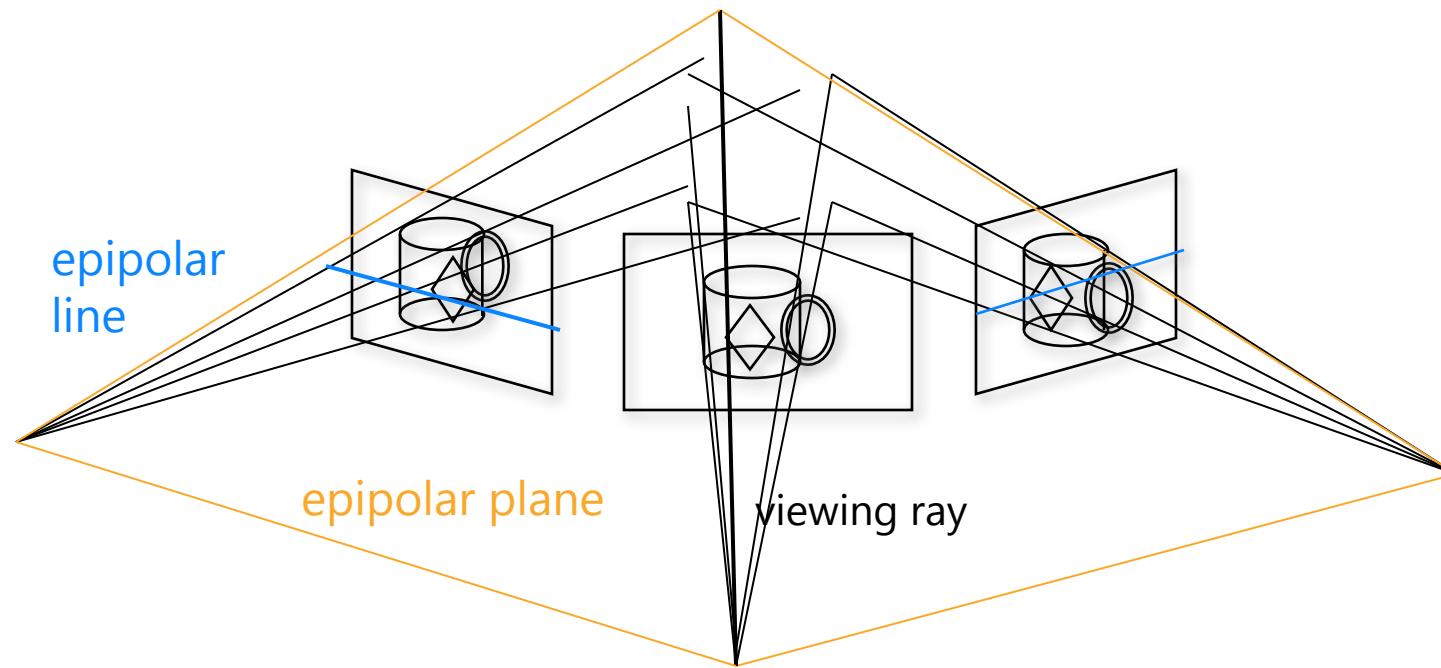
- Given two or more images of the same scene or object, compute a representation of its shape
- What are some possible representations?
 - depth maps
 - volumetric models
 - 3D surface models
 - planar (or offset) layers

Stereo matching

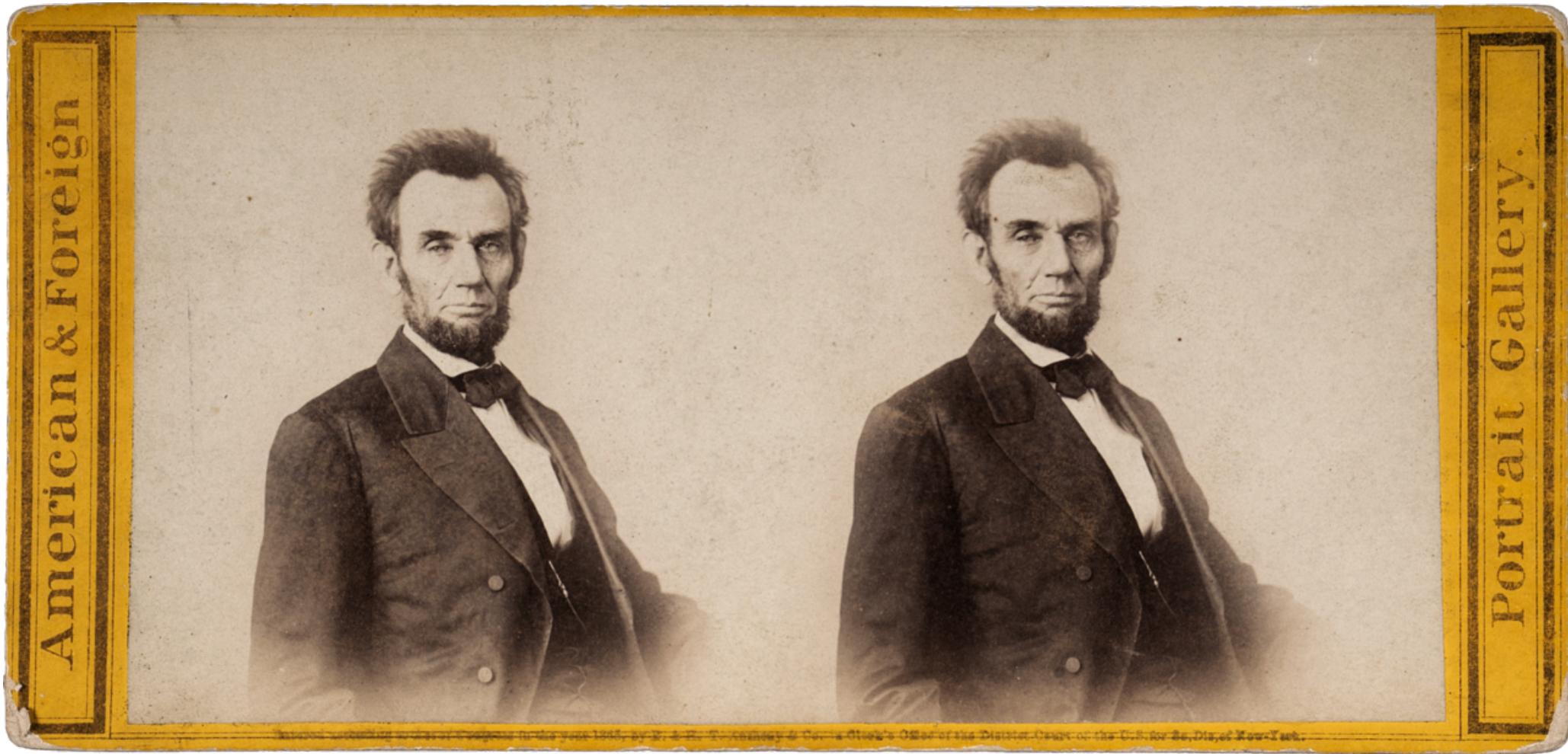
- What are some possible algorithms?
 - match “features” and interpolate
 - match edges and interpolate
 - match all pixels with windows (coarse-fine)
 - use optimization:
 - iterative updating
 - dynamic programming
 - energy minimization (regularization, stochastic)
 - graph algorithms

Stereo: epipolar geometry

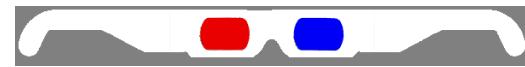
- Match features along epipolar lines



Stereo image pair



Anaglyphs



<http://www.rainbowsymplicity.com/freestuff.html>

(Wikipedia for images)

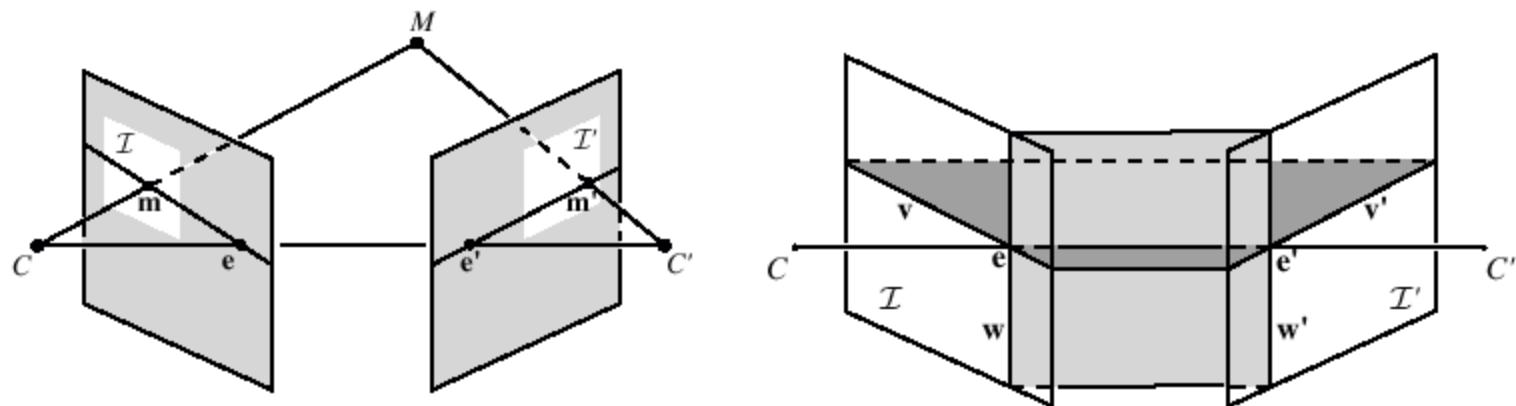
Public Library, Stereoscopic Looking Room, Chicago, by Phillips, 1923

Stereo: epipolar geometry

- for two images (or images with collinear camera centers), can find epipolar lines
- epipolar lines are the projection of the pencil of planes passing through the centers
- **Rectification:**
warping the input images (perspective transformation)
so that epipolar lines are horizontal

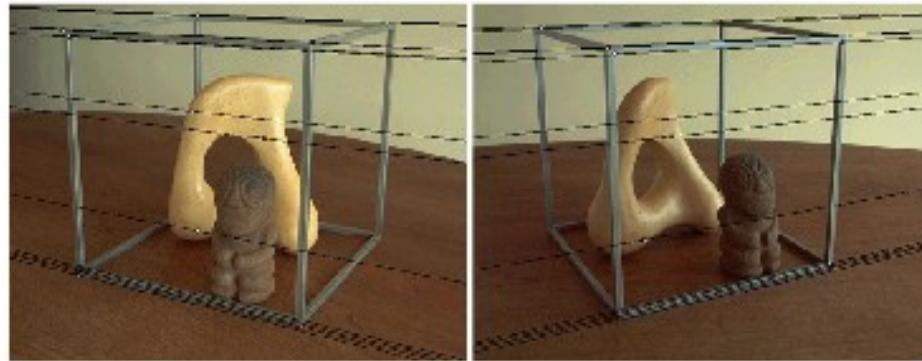
Rectification

- Project each image onto same plane, which is parallel to the epipole
- Resample lines (and shear/stretch) to place lines in correspondence, and minimize distortion

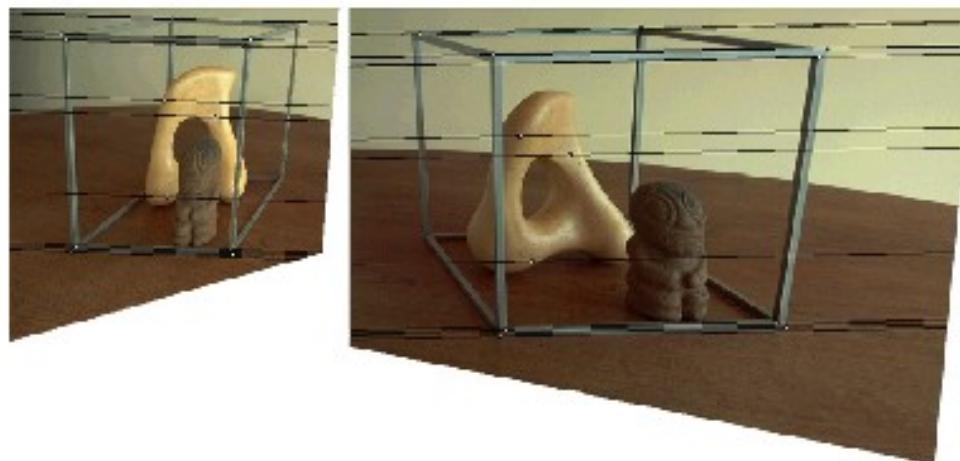


- [Loop and Zhang, CVPR'99]

Rectification



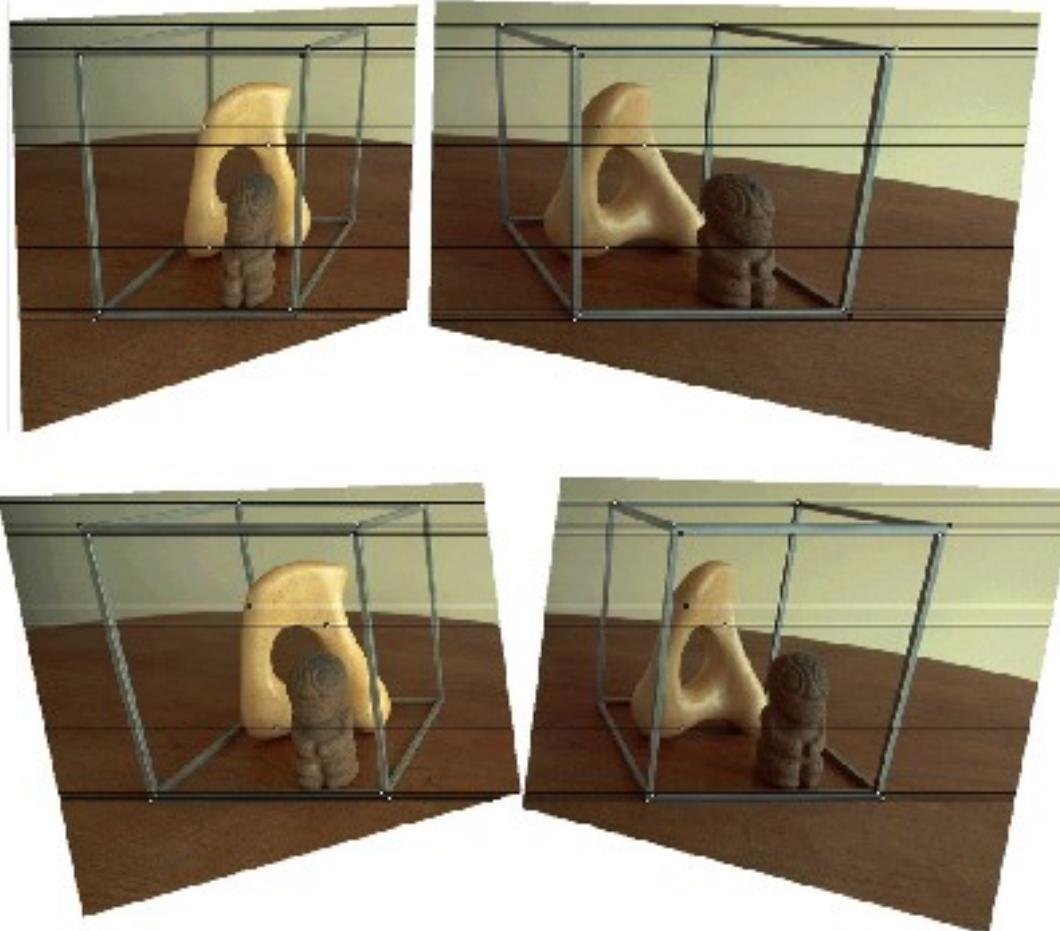
(a) Original image pair overlaid with several epipolar lines.



(b) Image pair transformed by the specialized projective mapping H_p and H'_p . Note that the epipolar lines are now parallel to each other in each image.

BAD!

Rectification



(c) Image pair transformed by the similarity H_r and H'_r . Note that the image pair is now rectified (the epipolar lines are horizontally aligned).

(d) Final image rectification after shearing transform H_s and H'_s . Note that the image pair remains rectified, but the horizontal distortion is reduced.

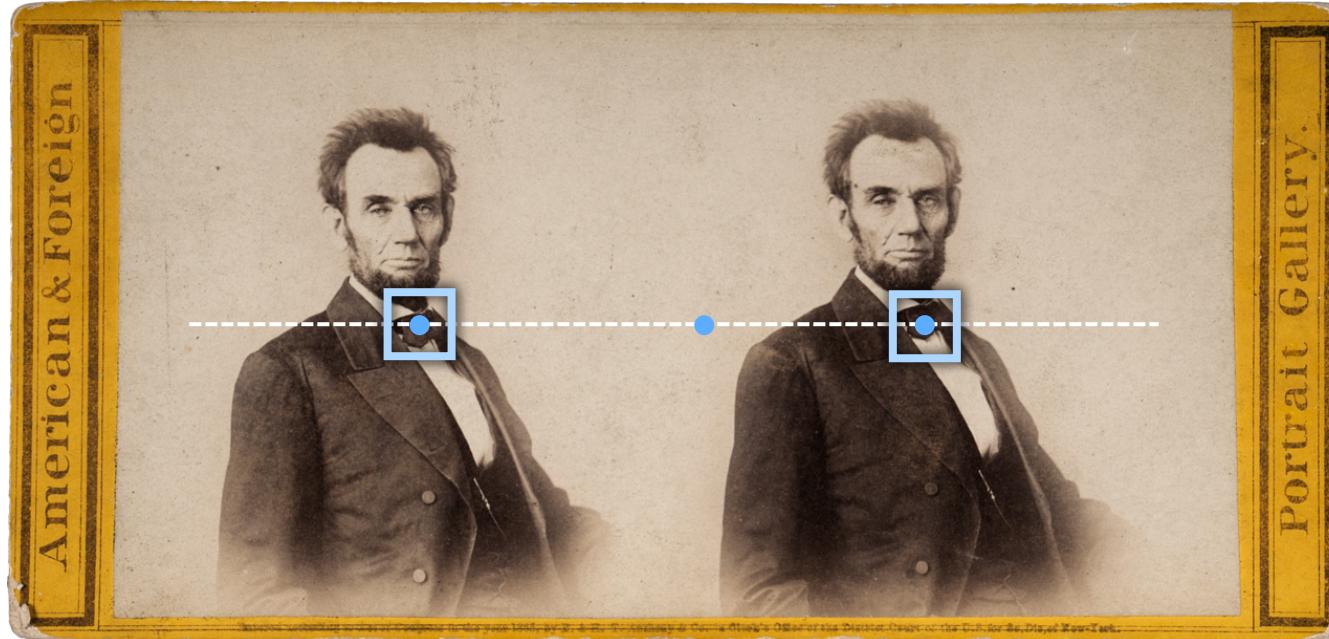
GOOD!

Finding correspondences

- apply feature matching criterion (e.g., correlation or Lucas-Kanade) at all pixels simultaneously
- search only over epipolar lines (many fewer candidate positions)



Your basic stereo algorithm



For each epipolar line

For each pixel in the left image

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

Improvement: match windows

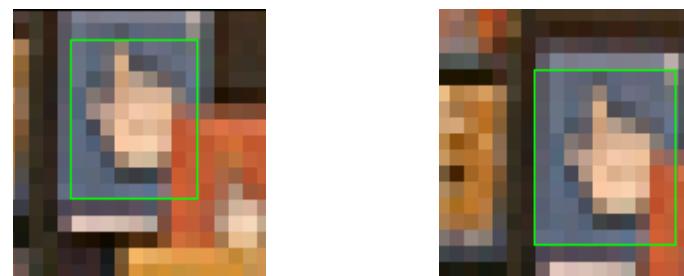
- This should look familiar ...

Image registration

- How do we determine correspondences?
 - block matching or SSD (sum squared differences)

$$E(x, y; d) = \sum_{(x', y') \in N(x, y)} [I_L(x' + d, y') - I_R(x', y')]^2$$

d is the disparity (horizontal motion)



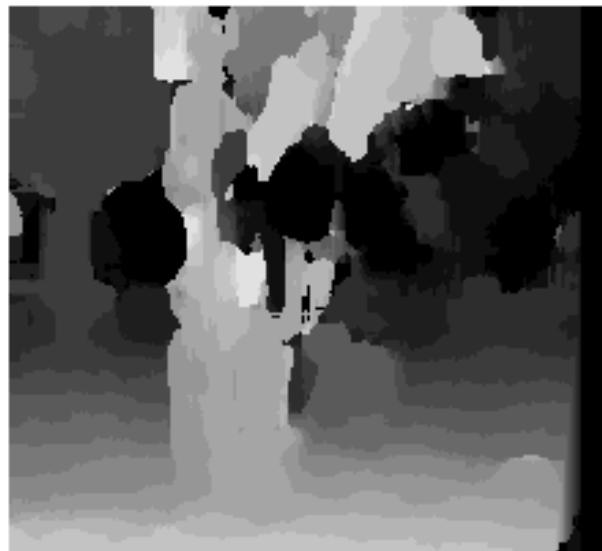
- How big should the neighborhood be?

Neighborhood size

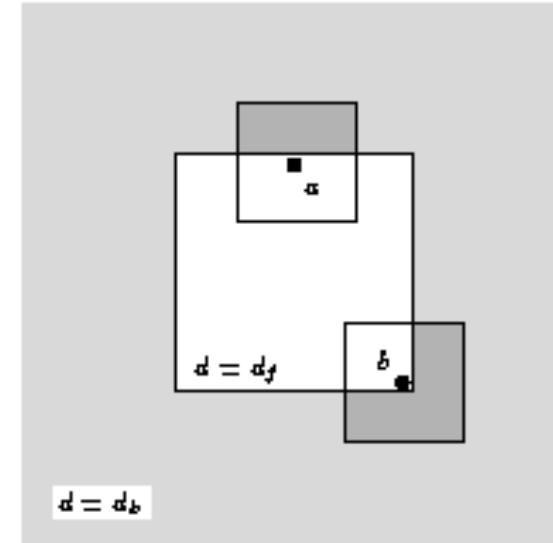
- Smaller neighborhood: more details
- Larger neighborhood: fewer isolated mistakes



$w = 3$



$w = 20$



Matching criteria

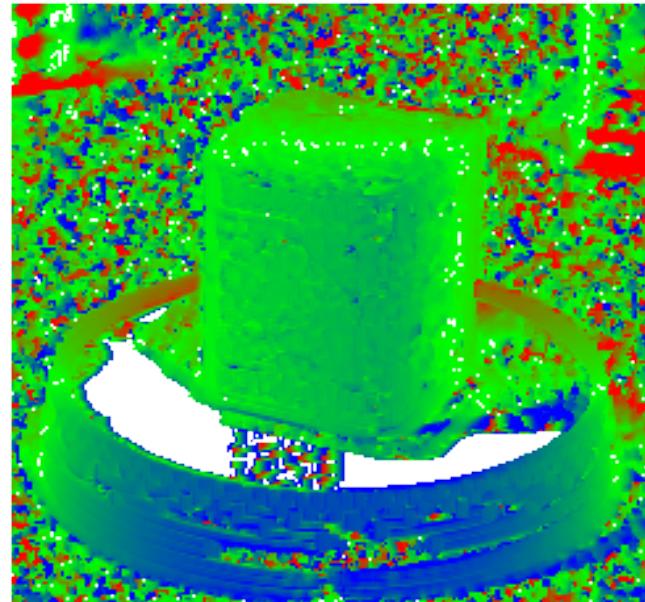
- Raw pixel values (correlation)
- Band-pass filtered images [Jones & Malik 92]
- “Corner” like features [Zhang, ...]
- Edges [many people...]
- Gradients [Seitz 89; Scharstein 94]
- Rank statistics [Zabih & Woodfill 94]

Stereo: certainty modeling

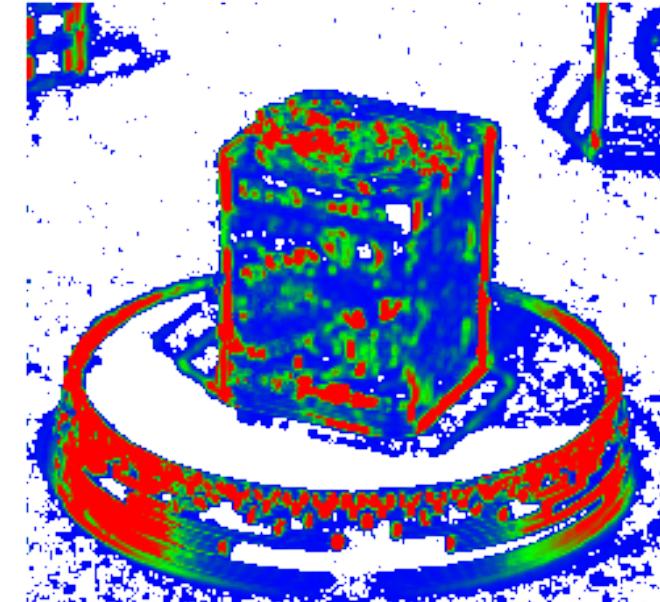
- Compute certainty map from correlations



input



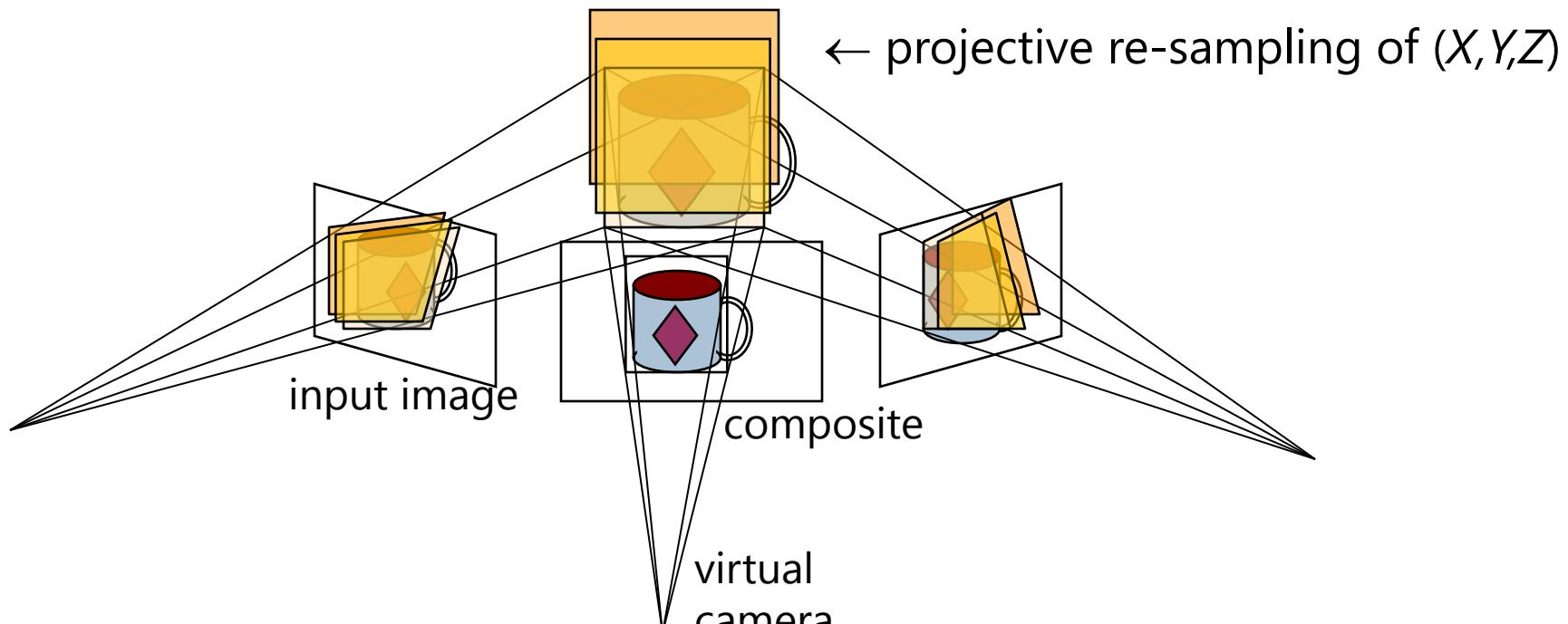
depth map



certainty map

Plane sweep stereo

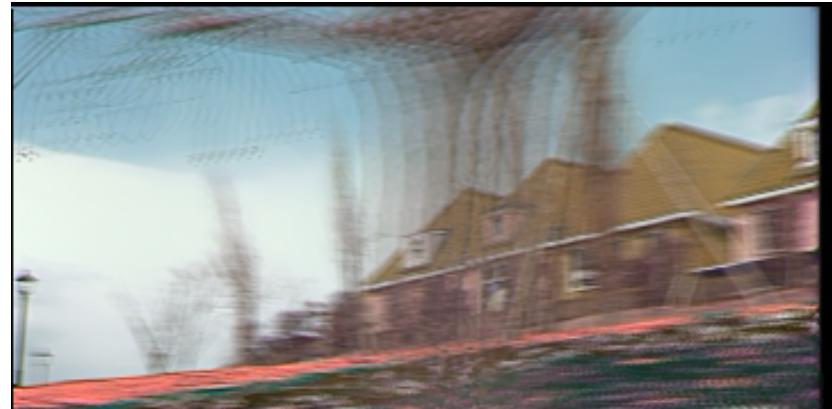
- Sweep family of planes through volume



- each plane defines an image \Rightarrow composite homography

Plane sweep stereo

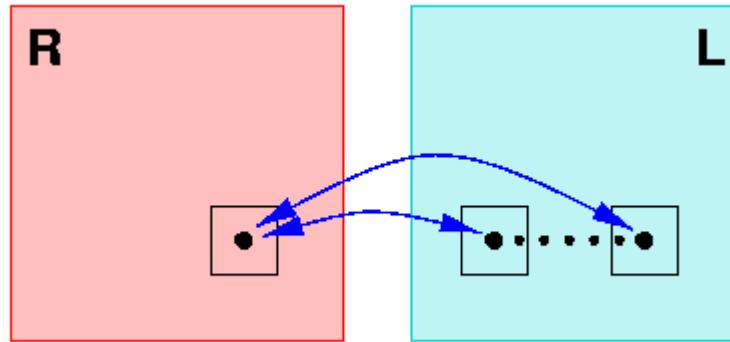
- For each depth plane
 - compute composite (mosaic) image — mean



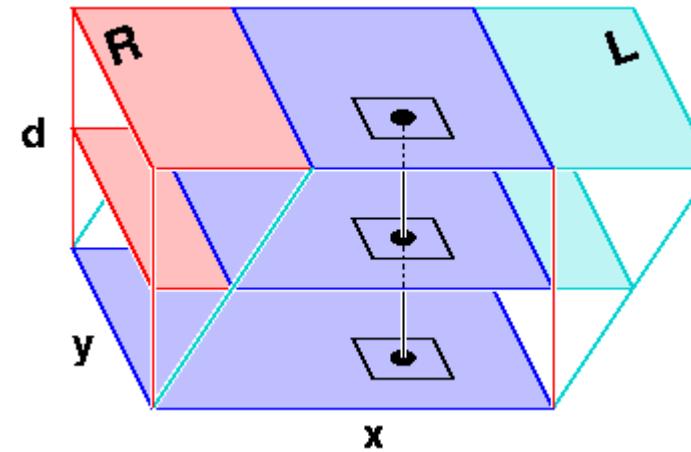
- compute error image — variance
 - convert to confidence and aggregate spatially
- Select winning depth at each pixel

Plane sweep stereo

- Re-order (pixel / disparity) evaluation loops



for every pixel,
for every disparity
compute cost



for every disparity
for every pixel
compute cost

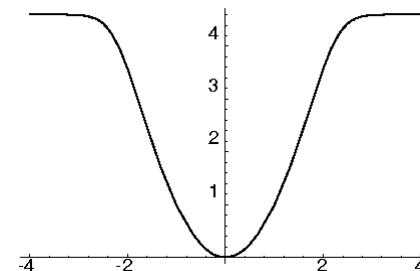
Stereo matching framework

1. For every disparity, compute raw matching costs

$$E_0(x, y; d) = \rho(I_L(x' + d, y') - I_R(x', y'))$$

Why use a robust function?

- occlusions, other outliers



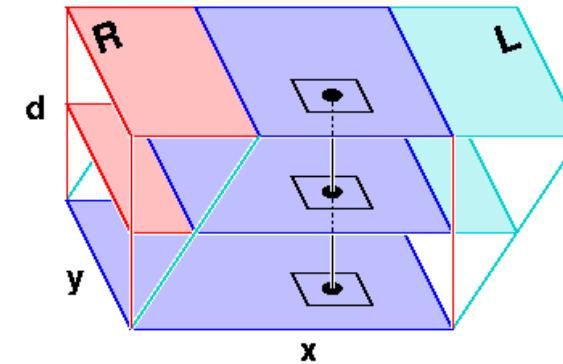
- Can also use alternative match criteria

Stereo matching framework

2. Aggregate costs spatially

$$E(x, y; d) = \sum_{(x', y') \in N(x, y)} E_0(x', y', d)$$

- Here, we are using a box filter (efficient moving average implementation)
- Can also use weighted average, [non-linear] diffusion...

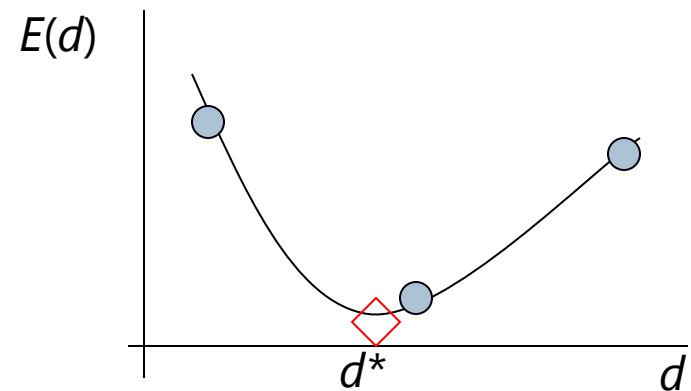


Stereo matching framework

3. Choose winning disparity at each pixel

$$d(x, y) = \arg \min_d E(x, y; d)$$

4. Interpolate to sub-pixel accuracy

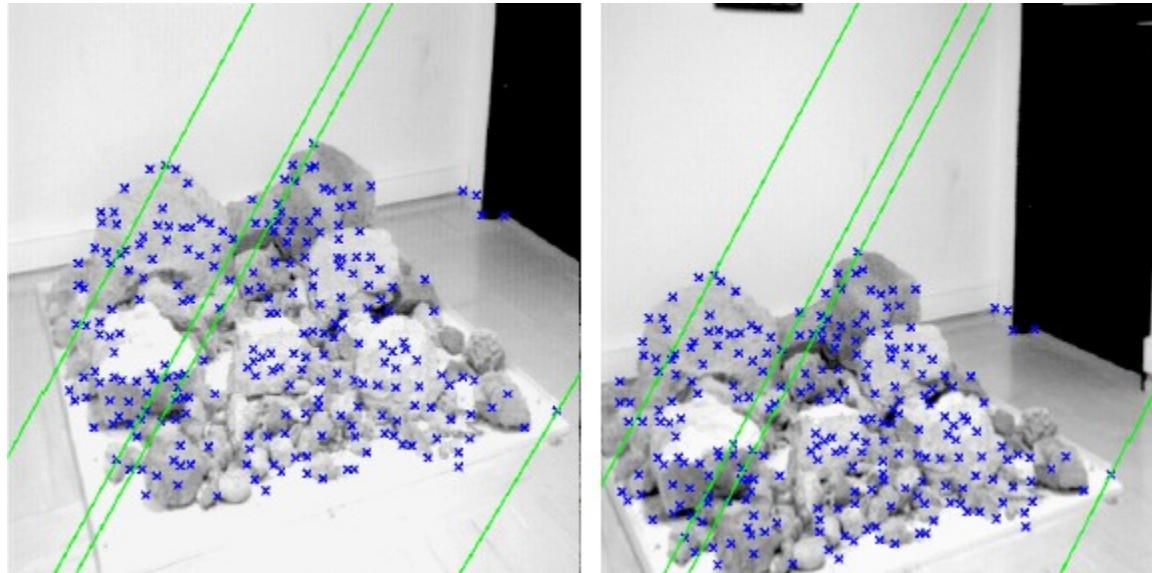


Traditional stereo matching

- Advantages:
 - gives detailed surface estimates
 - fast algorithms based on moving averages
 - sub-pixel disparity estimates and confidence
- Limitations:
 - narrow baseline \Rightarrow noisy estimates
 - fails in textureless areas
 - gets confused near occlusion boundaries

Feature-based stereo

- Match “corner” (interest) points



- Interpolate complete solution

Data interpolation

- Given a sparse set of 3D points, how do we interpolate to a full 3D surface?
 - Scattered data interpolation [Nielson93]
 - triangulate
 - put onto a grid and fill (use pyramid?)
 - place a kernel function over each data point
 - minimize an energy function

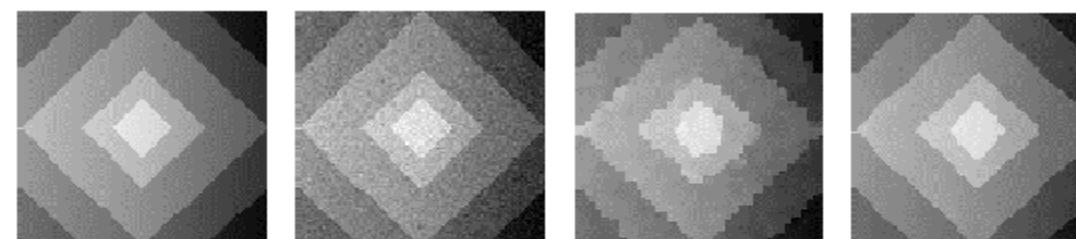
Graph cuts

- Solution technique for general 2D problem

$$E_{\text{total}}(\mathbf{d}) = E_{\text{data}}(\mathbf{d}) + \lambda E_{\text{smoothness}}(\mathbf{d})$$

$$E_{\text{data}}(\mathbf{d}) = \sum_{x,y} f_{x,y}(d_{x,y})$$

$$\begin{aligned} E_{\text{smoothness}}(\mathbf{d}) &= \sum_{x,y} \rho(d_{x,y} - d_{x-1,y}) \\ &\quad + \sum_{x,y} \rho(d_{x,y} - d_{x,y-1}) \end{aligned}$$



(a) original image

(b) observed image

(c) local min w.r.t.
standard moves

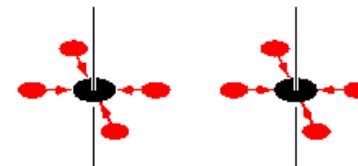
(d) local min w.r.t.
 α -expansion moves

Markov random field (MRF)

- Probability distribution on disparity field $d(x, y)$

$$p_P(d_{x,y}|\mathbf{d}) = p_P(d_{x,y}|\{d_{x',y'}, (x', y') \in \mathcal{N}(x, y)\})$$

$$p_P(\mathbf{d}) = \frac{1}{Z_P} e^{-E_P(\mathbf{d})}$$



$$E_P(\mathbf{d}) = \sum_{x,y} \rho_P(d_{x+1,y} - d_{x,y}) + \rho_P(d_{x,y+1} - d_{x,y})$$

- Enforces *smoothness* or coherence on field

MAP estimate

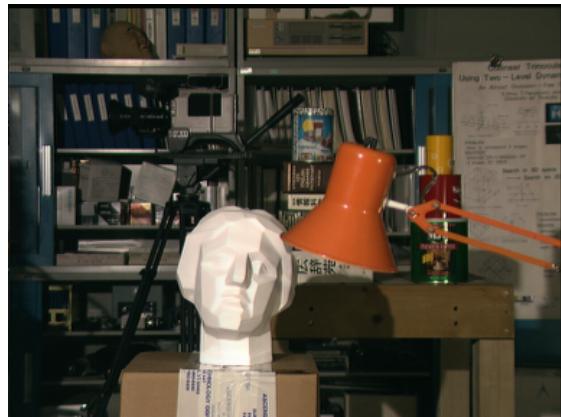
- Maximize posterior likelihood

$$\begin{aligned} E(\mathbf{d}) &= -\log p(\mathbf{d}|I_L, I_R) \\ &= \sum_{x,y} \rho_P(d_{x+1,y} - d_{x,y}) + \rho_P(d_{x,y+1} - d_{x,y}) \\ &\quad + \sum_{x,y} \rho_M(I_L(x + d_{x,y}, y) - I_R(x, y)) \end{aligned}$$

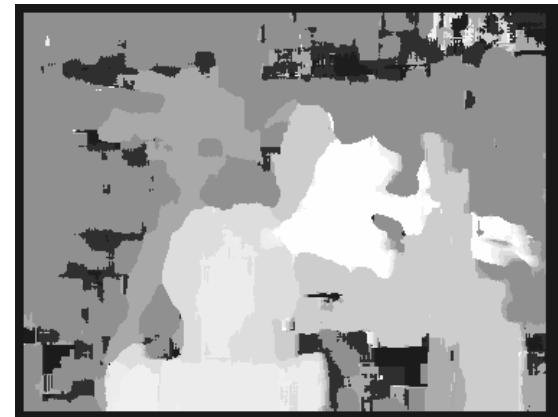
- Equivalent to *regularization* (energy minimization with smoothness constraints)

Depth map results

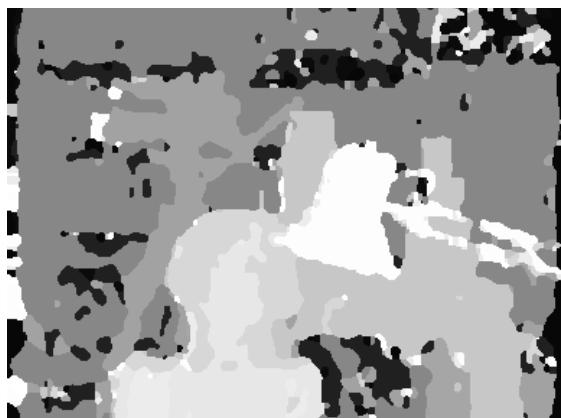
Input image



Sum of abs. diff. (SAD)



Mean field



Graph cuts



Middlebury stereo benchmark

Error Threshold = 1		Sort by nonocc			Sort by all			Sort by disc					
Algorithm	Avg.	Tsukuba ground truth			Venus ground truth			Teddy ground truth			Cones ground truth		
		Rank	nonocc	all	disc	nonocc	all	disc	nonocc	all	disc	nonocc	all
AdaptingBP [17]	2.8	<u>1.11</u> 6	1.37 3	5.79 7	<u>0.10</u> 1	0.21 2	<u>1.44</u> 1	<u>4.22</u> 4	7.06 2	11.8 4	<u>2.48</u> 1	7.92 2	7.32 1
DoubleBP2 [35]	2.9	<u>0.88</u> 1	<u>1.29</u> 1	4.76 1	<u>0.13</u> 3	0.45 5	1.87 5	<u>3.53</u> 2	8.30 3	9.63 1	<u>2.90</u> 3	8.78 8	7.79 2
DoubleBP [15]	4.9	<u>0.88</u> 2	1.29 2	4.76 2	<u>0.14</u> 5	0.60 13	2.00 7	<u>3.55</u> 3	8.71 5	9.70 2	<u>2.90</u> 4	9.24 11	7.80 3
SubPixDoubleBP [30]	5.6	<u>1.24</u> 10	1.76 13	5.98 8	<u>0.12</u> 2	0.46 6	1.74 4	<u>3.45</u> 1	8.38 4	10.0 3	<u>2.93</u> 5	8.73 7	7.91 4
AdaptOvrSeqBP [33]	9.9	<u>1.69</u> 22	2.04 21	5.64 6	<u>0.14</u> 4	0.20 1	1.47 2	<u>7.04</u> 14	11.1 7	16.4 11	<u>3.60</u> 11	8.96 10	8.84 10
SymBP+occ [7]	10.8	<u>0.97</u> 4	1.75 12	5.09 4	<u>0.16</u> 6	0.33 3	2.19 8	<u>6.47</u> 8	10.7 6	17.0 14	<u>4.79</u> 24	10.7 21	10.9 20
PlaneFitBP [32]	10.8	<u>0.97</u> 5	1.83 14	5.26 5	<u>0.17</u> 7	0.51 8	1.71 3	<u>6.65</u> 9	12.1 13	14.7 7	<u>4.17</u> 20	10.7 20	10.6 19
AdaptDispCalib [36]	11.8	<u>1.19</u> 8	1.42 4	6.15 9	<u>0.23</u> 9	0.34 4	2.50 11	<u>7.80</u> 19	13.6 21	17.3 17	<u>3.62</u> 12	9.33 12	9.72 15
Segm+visib [4]	12.2	<u>1.30</u> 15	1.57 5	6.92 18	<u>0.79</u> 21	1.06 18	6.76 22	<u>5.00</u> 5	<u>6.54</u> 1	12.3 5	<u>3.72</u> 13	8.62 6	10.2 17
C-SemiGlob [19]	12.3	<u>2.61</u> 29	3.29 24	9.89 27	<u>0.25</u> 12	0.57 10	3.24 15	<u>5.14</u> 6	11.8 8	13.0 6	<u>2.77</u> 2	8.35 4	8.20 5
SO+borders [29]	12.8	<u>1.29</u> 14	1.71 9	6.83 15	<u>0.25</u> 13	0.53 9	2.26 9	<u>7.02</u> 13	12.2 14	16.3 9	<u>3.90</u> 15	9.85 16	10.2 18
DistinctSM [27]	14.1	<u>1.21</u> 9	1.75 11	6.39 11	<u>0.35</u> 14	0.69 16	2.63 13	<u>7.45</u> 18	13.0 17	18.1 19	<u>3.91</u> 16	9.91 18	8.32 7
CostAggr+occ [39]	14.3	<u>1.38</u> 17	1.96 17	7.14 19	<u>0.44</u> 16	1.13 19	4.87 19	<u>6.80</u> 11	11.9 10	17.3 16	<u>3.60</u> 10	8.57 5	9.36 13
OverSegmBP [26]	14.5	<u>1.69</u> 23	1.97 18	8.47 24	<u>0.51</u> 18	0.68 15	4.69 18	<u>6.74</u> 10	11.9 12	15.8 8	<u>3.19</u> 8	8.81 9	8.89 11
SegmentSupport [28]	15.1	<u>1.25</u> 11	1.62 7	6.68 13	<u>0.25</u> 11	0.64 14	2.59 12	<u>8.43</u> 24	14.2 22	18.2 20	<u>3.77</u> 14	9.87 17	9.77 16
RegionTreeDP [18]	15.7	<u>1.39</u> 19	1.64 8	6.85 16	<u>0.22</u> 8	0.57 10	1.93 6	<u>7.42</u> 17	11.9 11	16.8 13	<u>6.31</u> 30	11.9 27	11.8 23
EnhancedBP [24]	16.6	0.94 3	1.74 10	5.05 3	0.35 15	0.86 17	4.34 17	8.11 22	13.3 19	18.5 22	5.09 27	11.1 23	11.0 21

Traditional stereo

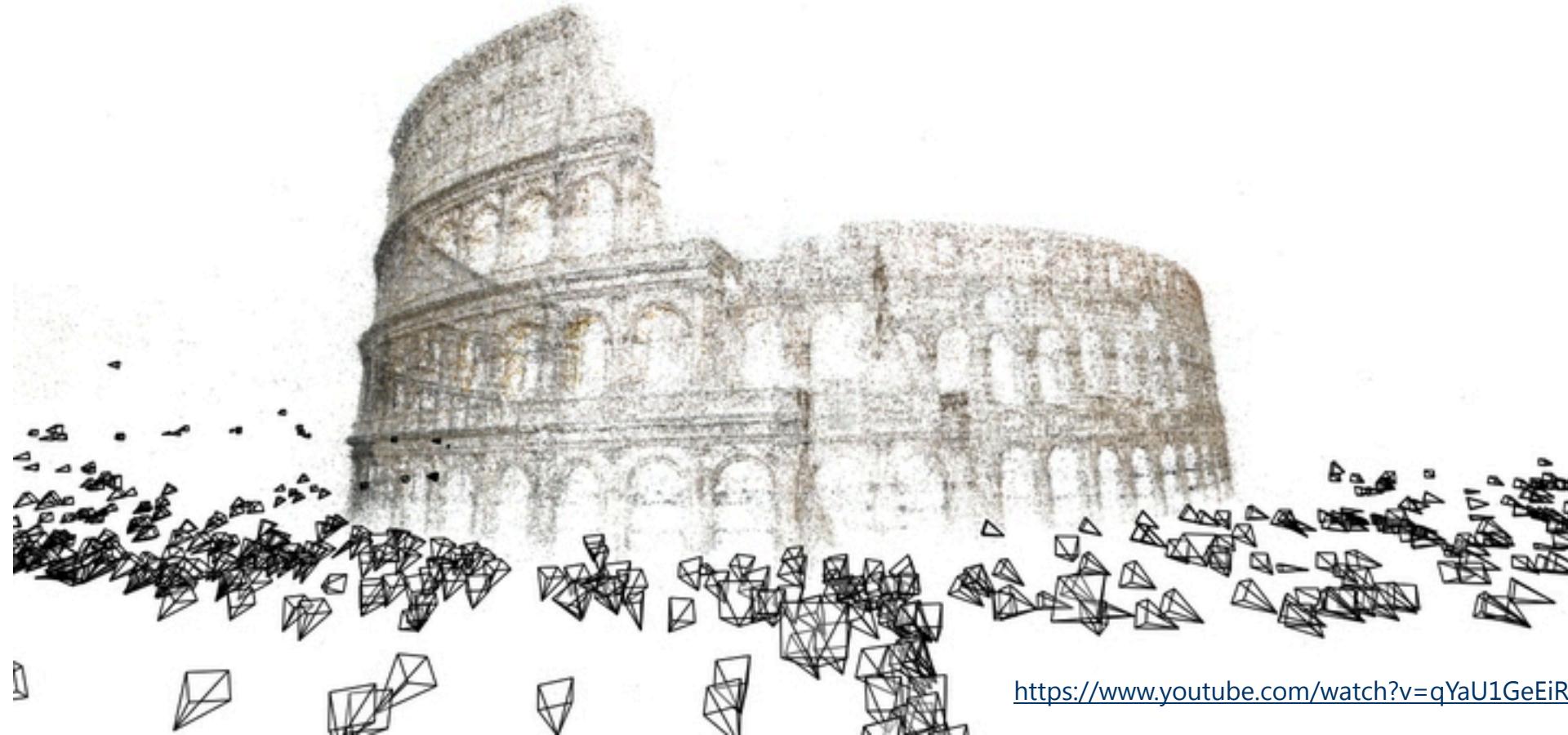
- Advantages:
 - works very well in non-occluded regions
- Disadvantages:
 - restricted to two images (not)
 - gets confused in occluded regions
 - can't handle *mixed* pixels

Slides by Szymon Rusinkiewicz and Jianxiong Xiao

Structure-from-motion & bundle adjustment

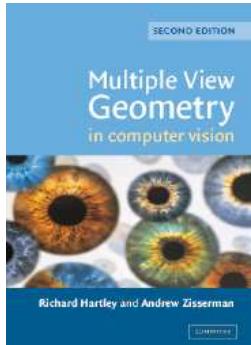


Building Rome in a Day



<https://www.youtube.com/watch?v=qYaU1GeEiR8>

How?



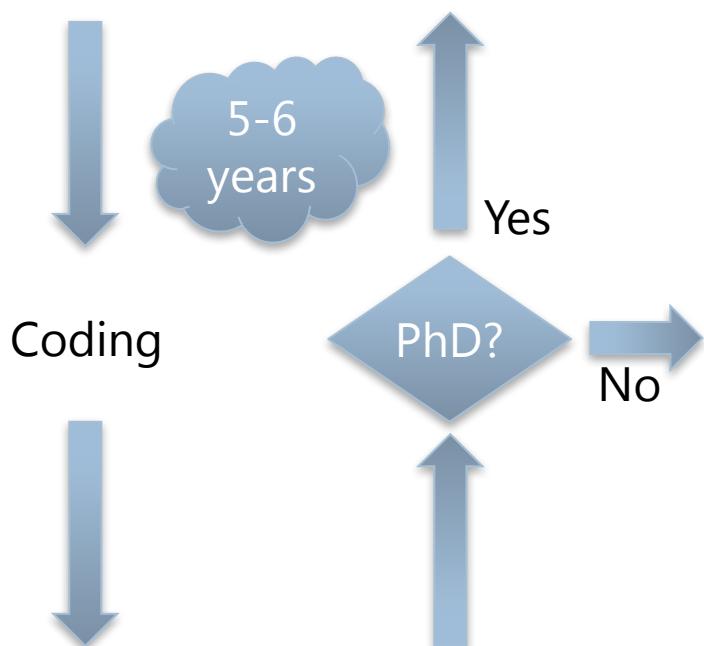
Reading: the MVG bible (need ~2 years)



Coding



Crying: Not Working At All



Steps

Images	→ Points:	Structure from Motion
Points	→ More points: Multiple View Stereo	
Points	→ Meshes:	Model Fitting
+	Meshes → Models:	Texture Mapping
=====	Images → Models:	Image-based Modeling

Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

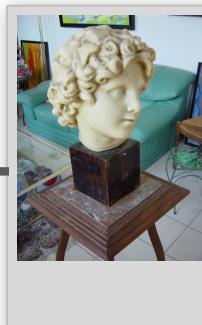
Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

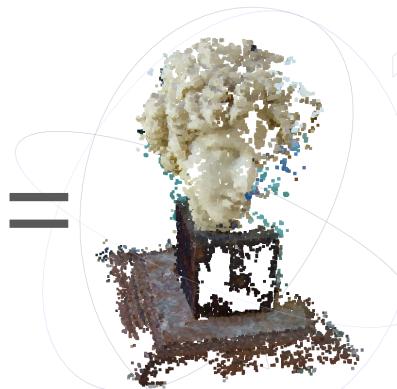
= Images → Models: Image-based Modeling



+



+



Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

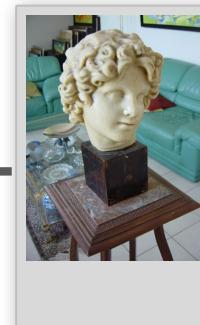
= Images → Models: Image-based Modeling



+



+



=



Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

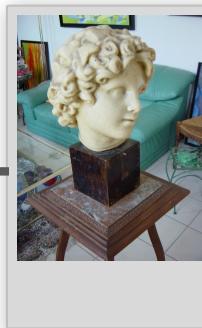
Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

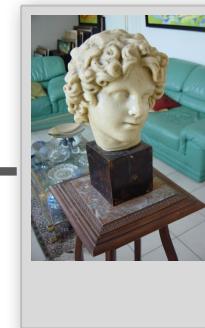
= Images → Models: Image-based Modeling



+



+



+

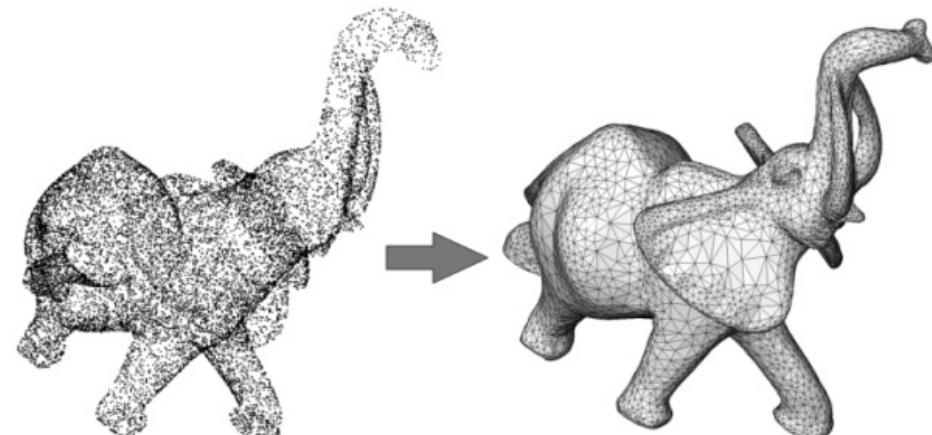


=



Steps

- Images → Points: Structure from Motion
- Points → More points: Multiple View Stereo
- Points → Meshes: Model Fitting
- + Meshes → Models: Texture Mapping
-
- = Images → Models: Image-based Modeling



Steps

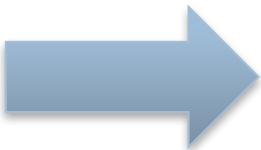
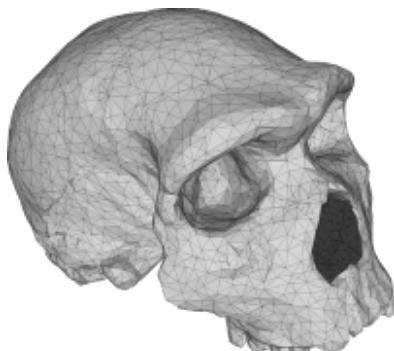
Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

= Images → Models: Image-based Modeling



Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

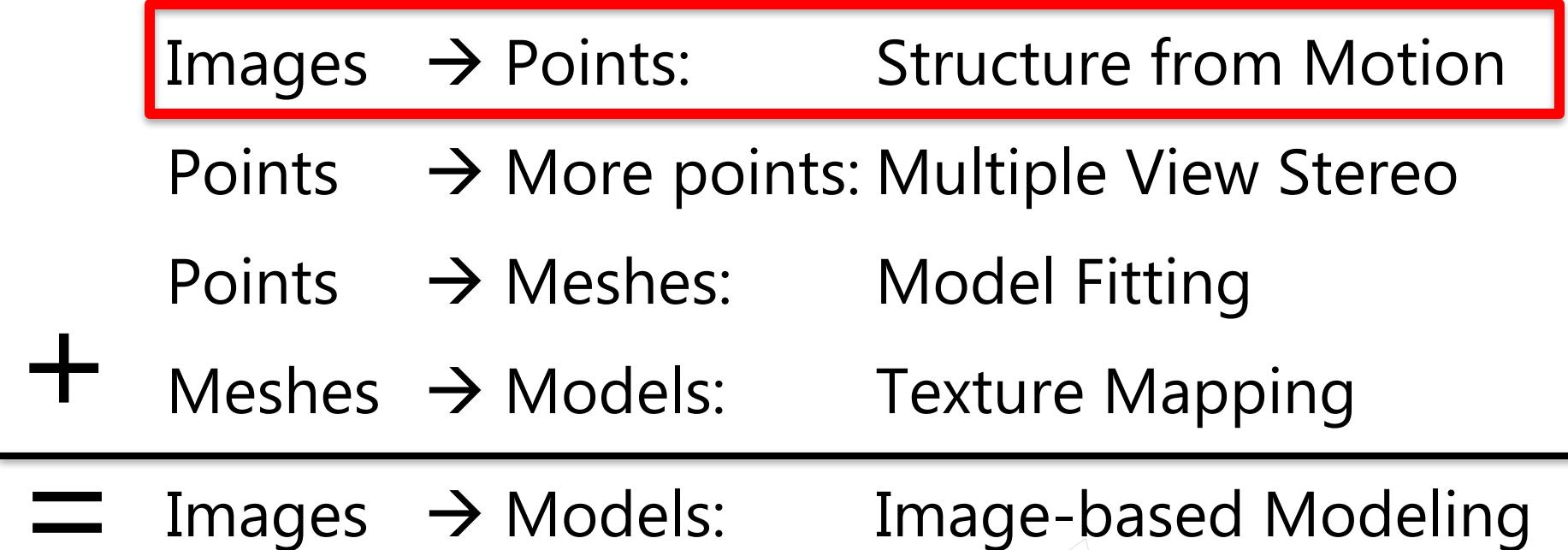
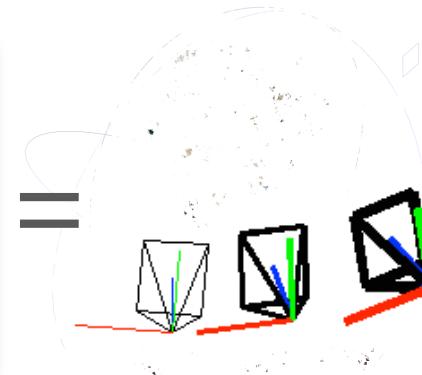
Points → Meshes: Model Fitting

+

Meshes → Models: Texture Mapping

= **Images → Models: Image-based Modeling**

Steps

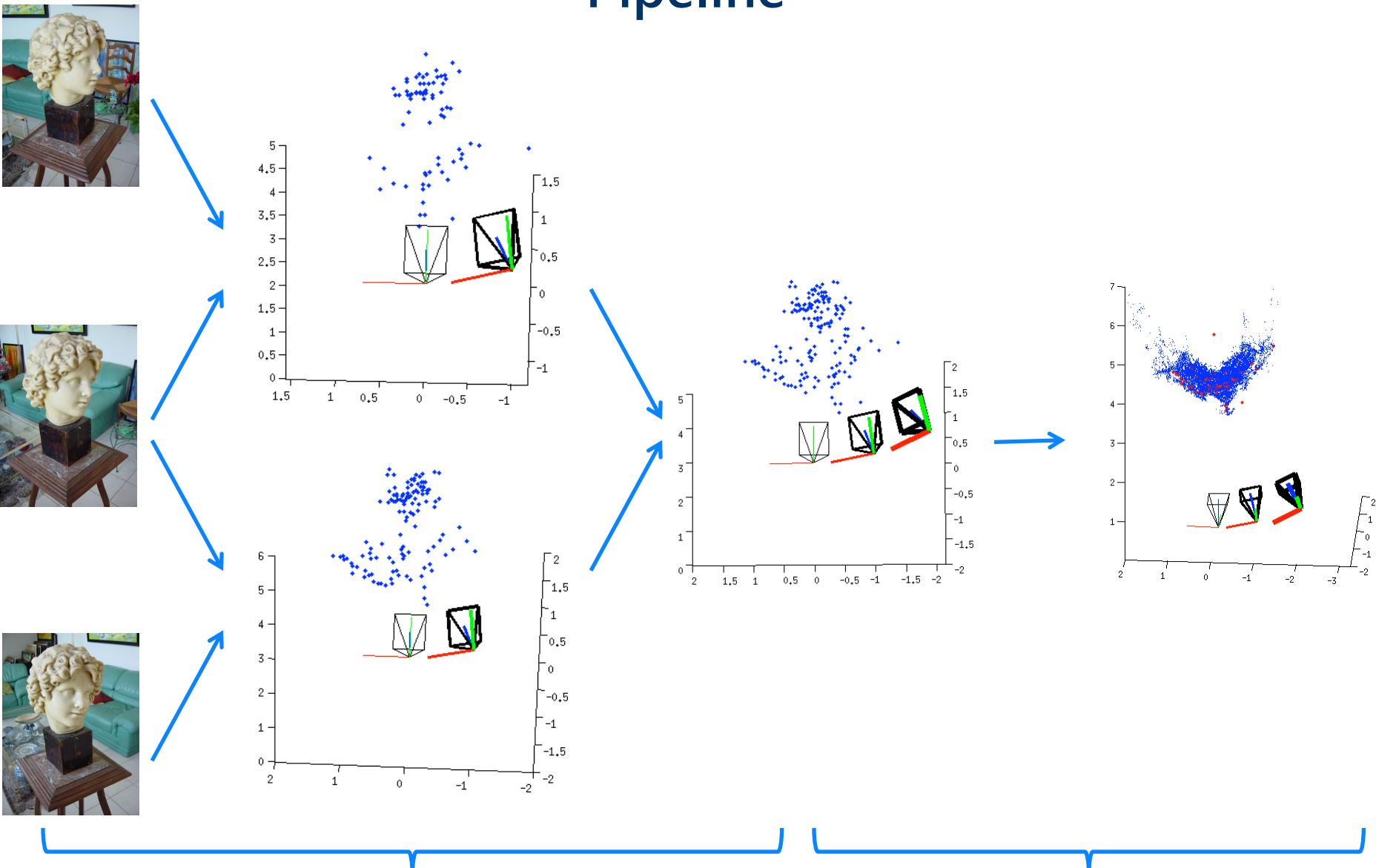
 $+$  $+$ 

Structure from motion

- Structure = 3D point cloud of the scene
- Motion = camera location and orientation
- SFM = get the point cloud from moving cameras
- Structure and motion: joint problems to solve



Pipeline

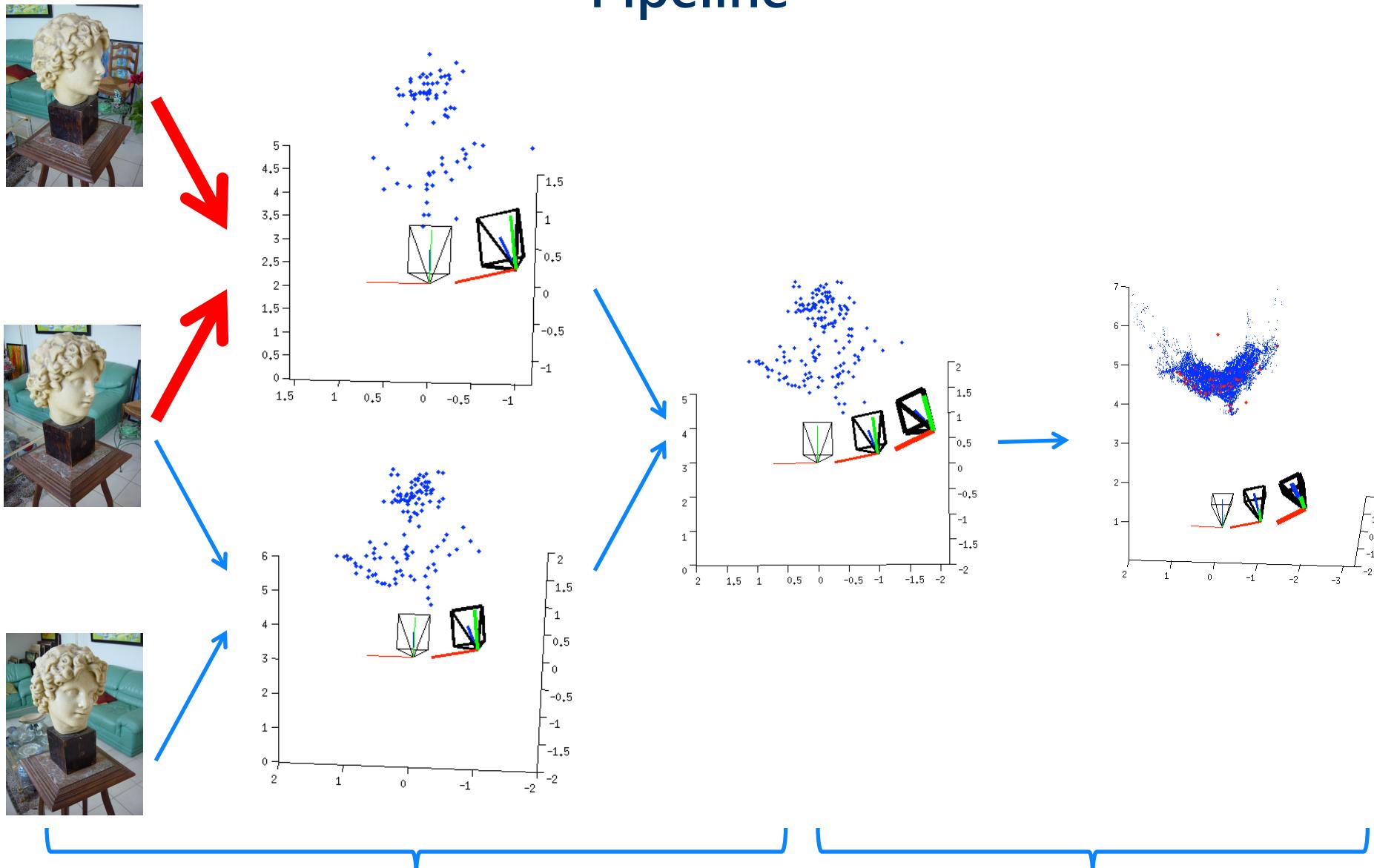


Structure from Motion (SFM)

CM50248 – Visual Understanding 1 – Lecture 6

Multi-view Stereo (MVS)

Pipeline

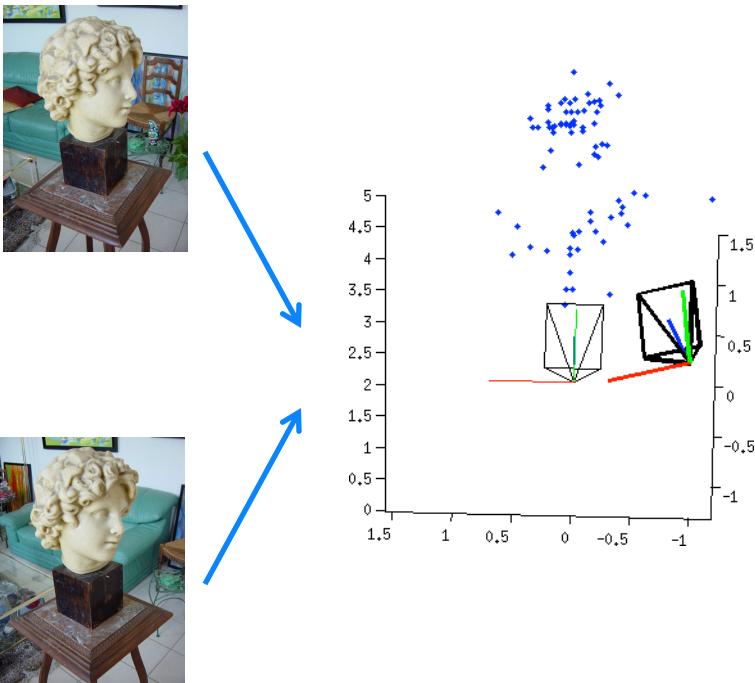


Structure from Motion (SfM)

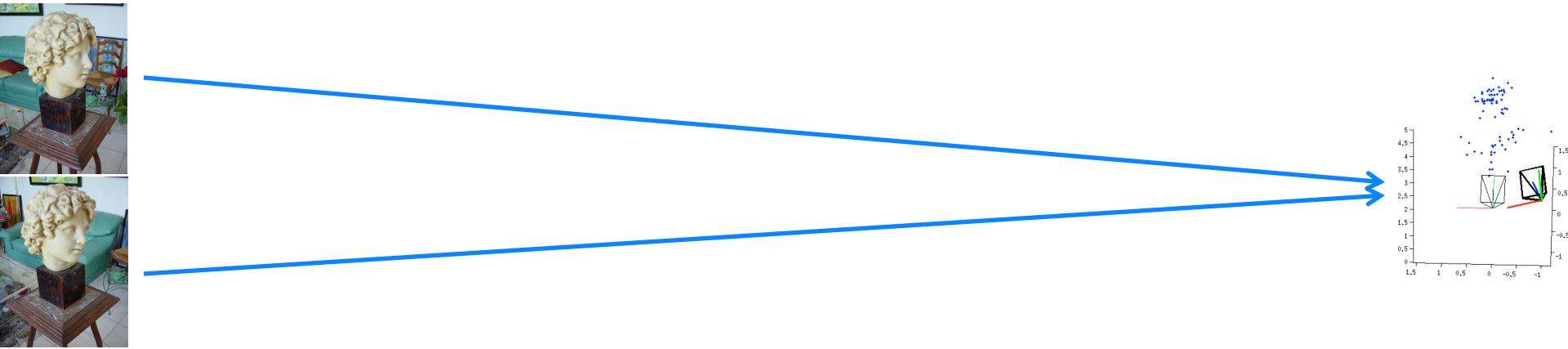
CM50248 – Visual Understanding 1 – Lecture 6

Multi-view Stereo (MVS)

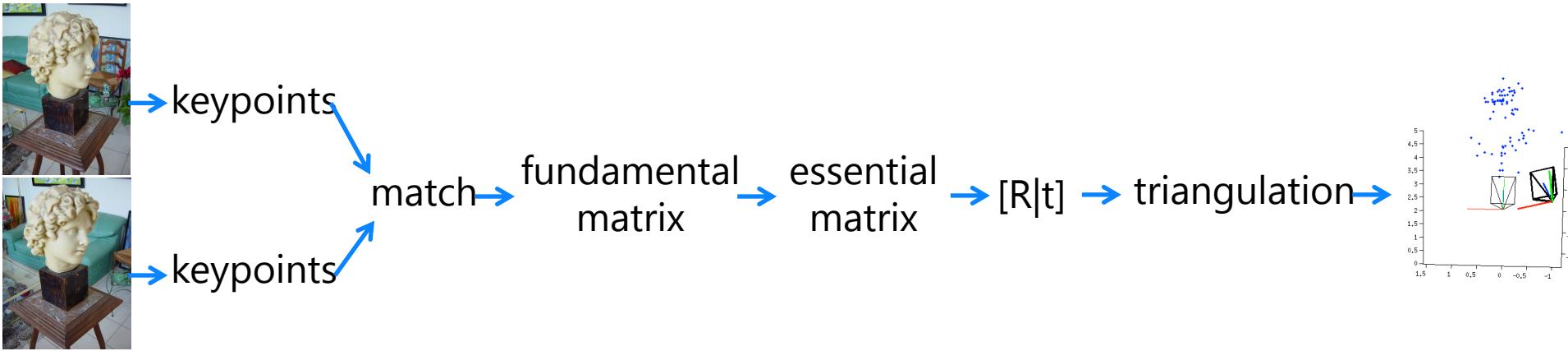
Two-view reconstruction



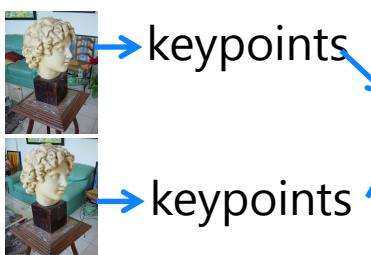
Two-view reconstruction



Two-view reconstruction



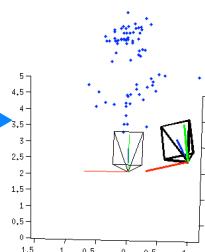
Keypoint detection



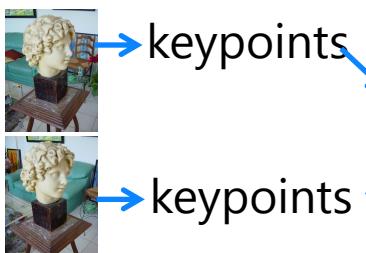
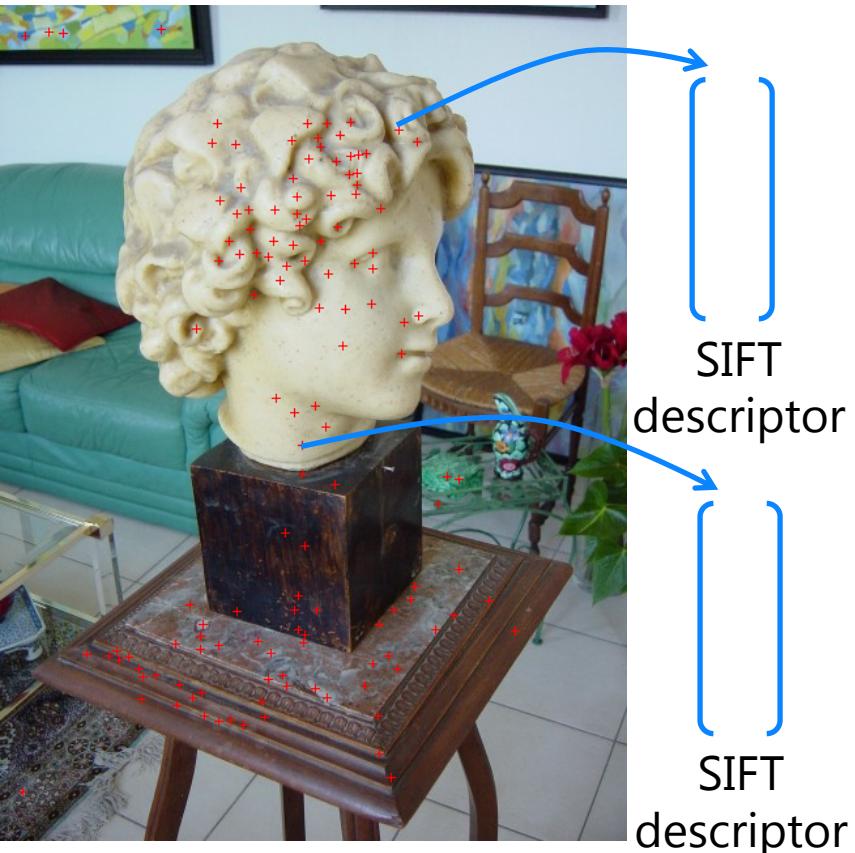
keypoints
match
fundamental matrix
essential matrix
 $[R|t]$ → triangulation

2017-10-20

CM50248 – Visual Understanding 1 – Lecture 6



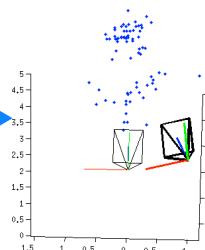
Descriptor for each point



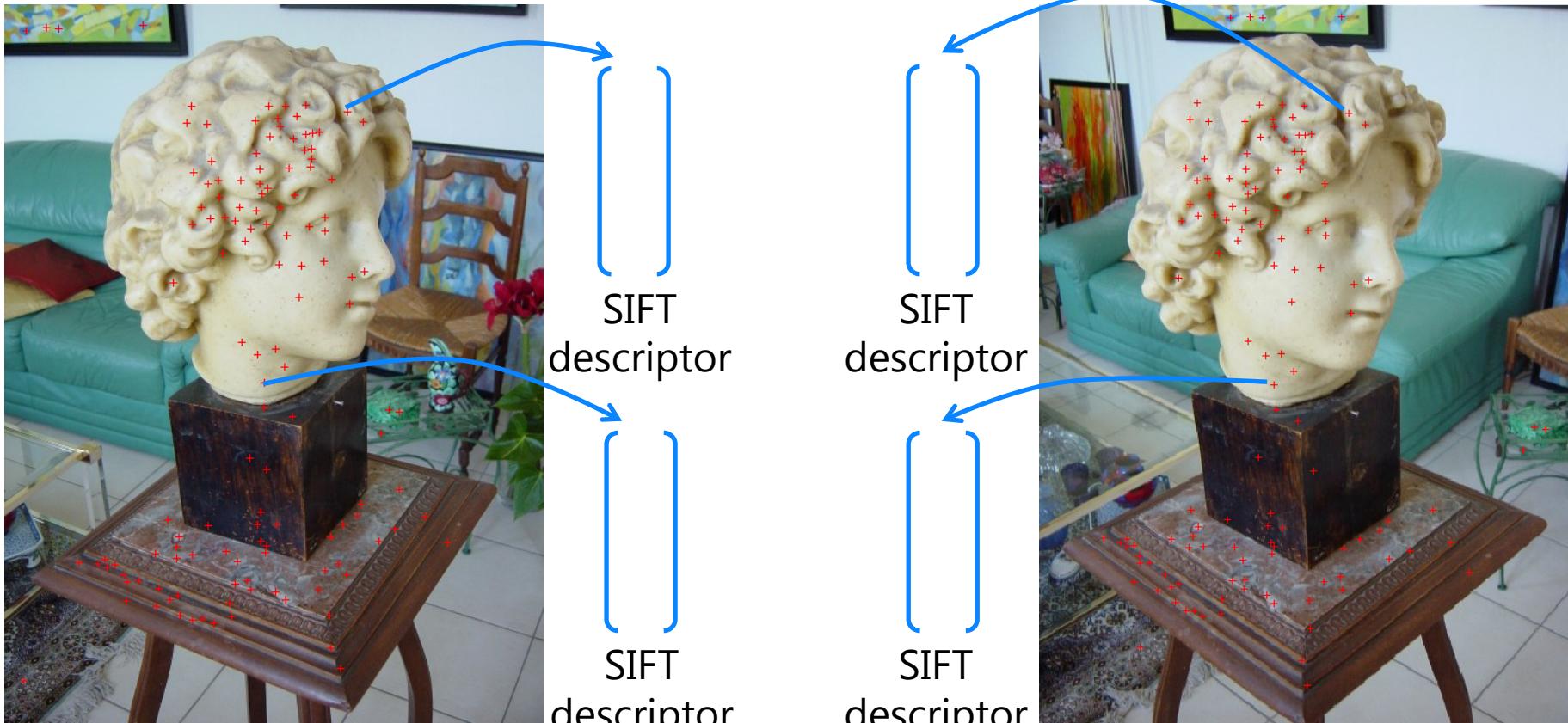
2017-10-20

CM50248 – Visual Understanding 1 – Lecture 6

keypoints → match → fundamental matrix → essential matrix → $[R|t]$ → triangulation



Same for the other images



keypoints

keypoints

match

fundamental
matrix

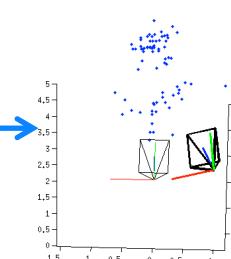
CM50248 – Visual Understanding 1 – Lecture 6

SIFT
descriptor

SIFT
descriptor

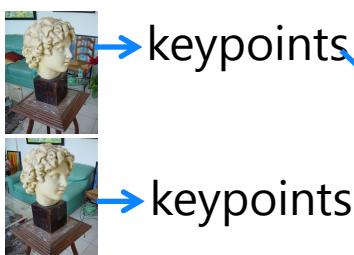
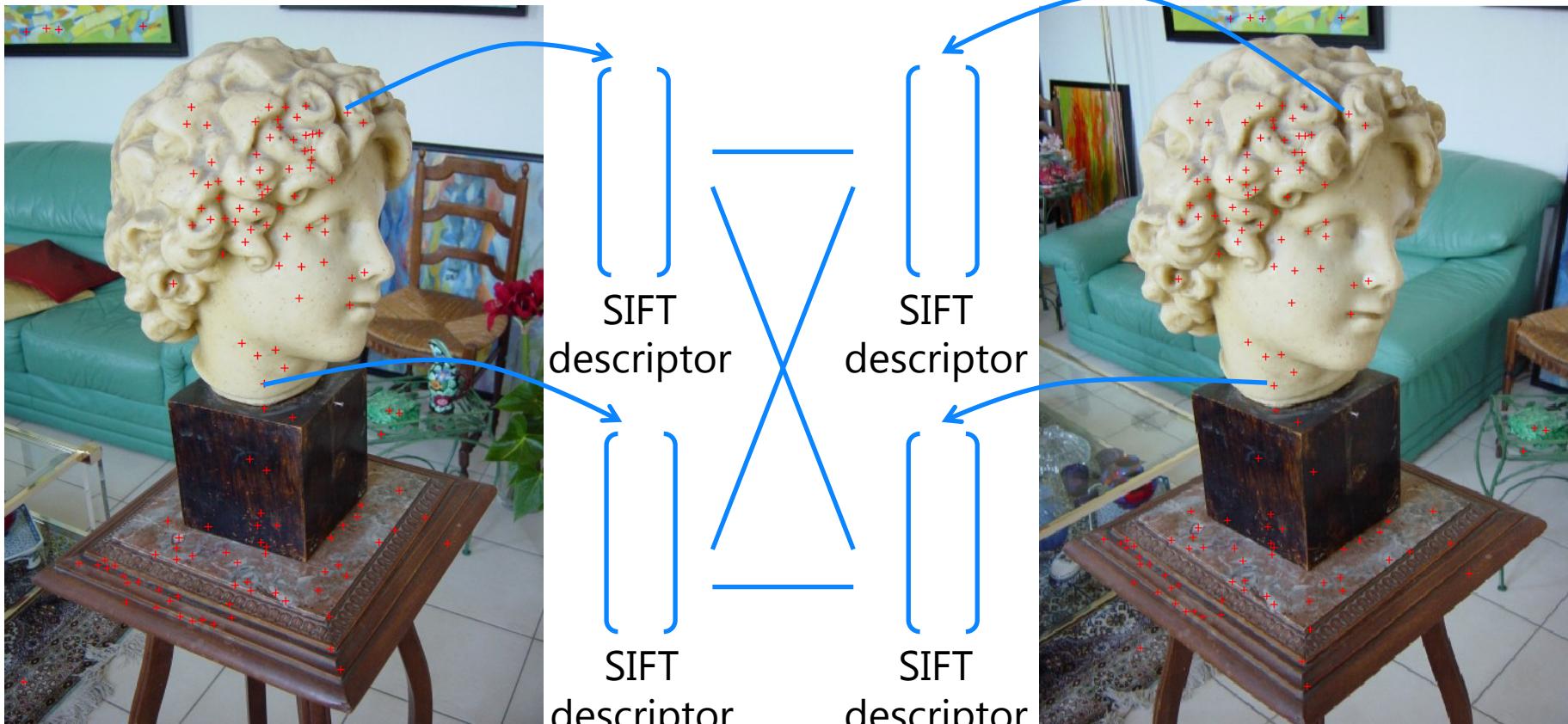
SIFT
descriptor

2017-10-20



79

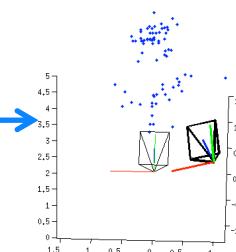
Point match for correspondences



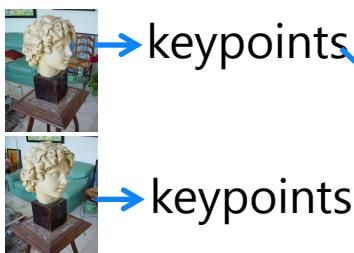
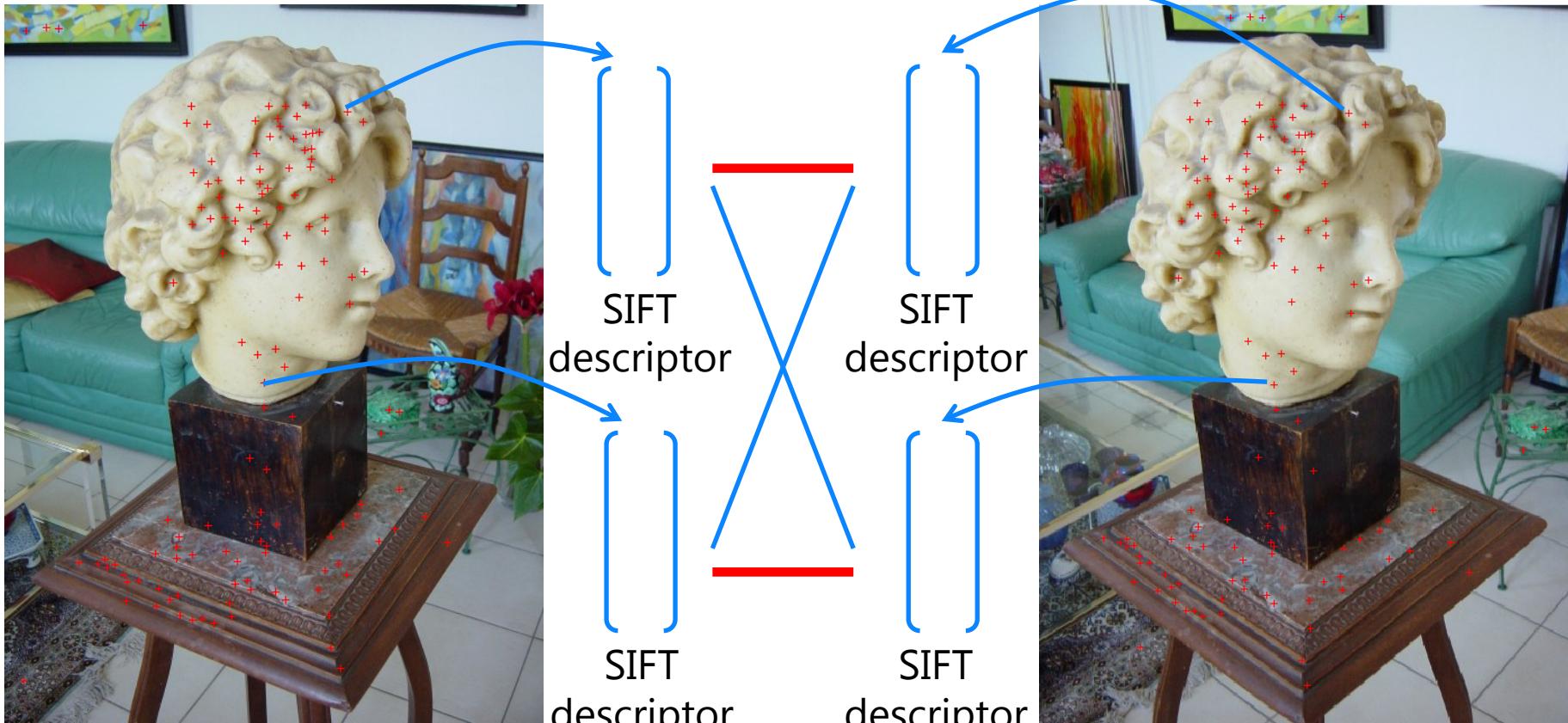
match → fundamental matrix → essential matrix → $[R|t]$ → triangulation

2017-10-20

CM50248 – Visual Understanding 1 – Lecture 6



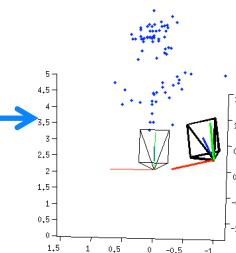
Point match for correspondences



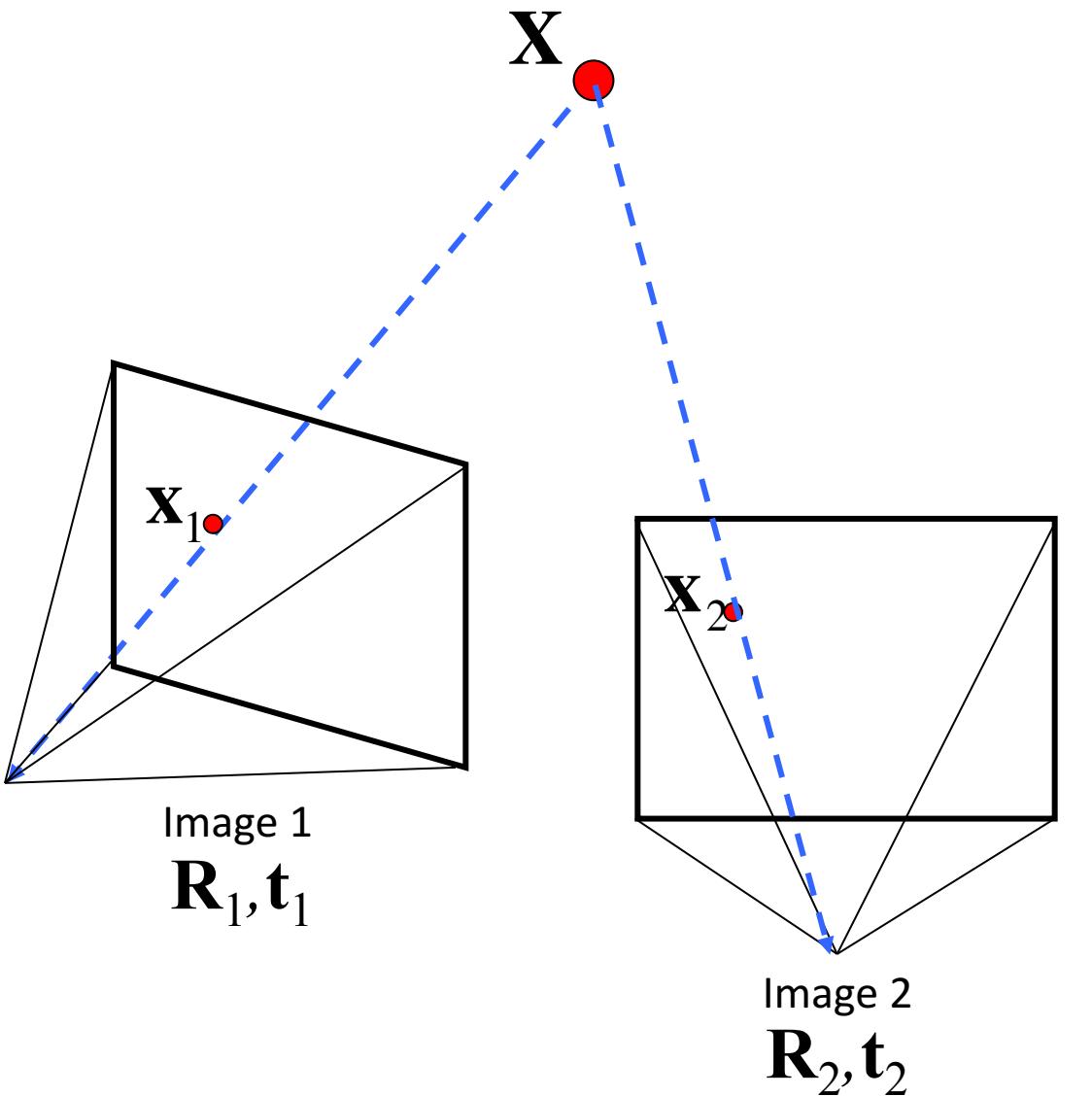
keypoints
match
fundamental matrix
essential matrix
 $[R|t]$ → triangulation

2017-10-20

CM50248 – Visual Understanding 1 – Lecture 6



Fundamental Matrix



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

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

Exercise

What is the difference between a fundamental matrix and a homography?
(Both of them are to explain 2D point to 2D point correspondences.)



Estimating the fundamental matrix

- Given a correspondence
- The basic incidence relation is $\mathbf{x}_1 \leftrightarrow \mathbf{x}_2$

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



$$\begin{bmatrix} x_1x_2, x_1y_2, x_1, y_1x_2, y_1y_2, y_1, x_2, y_2, 1 \end{bmatrix}$$

Need 8 points

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

Estimating the fundamental matrix

$\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2 = 0$ for 8 point correspondences:

$\mathbf{x}_1^1 \leftrightarrow \mathbf{x}_2^1, \mathbf{x}_1^2 \leftrightarrow \mathbf{x}_2^2, \mathbf{x}_1^3 \leftrightarrow \mathbf{x}_2^3, \mathbf{x}_1^4 \leftrightarrow \mathbf{x}_2^4, \mathbf{x}_1^5 \leftrightarrow \mathbf{x}_2^5, \mathbf{x}_1^6 \leftrightarrow \mathbf{x}_2^6, \mathbf{x}_1^7 \leftrightarrow \mathbf{x}_2^7, \mathbf{x}_1^8 \leftrightarrow \mathbf{x}_2^8$

$$\begin{bmatrix} x_1^1 x_2^1 & x_1^1 y_2^1 & x_1^1 & y_1^1 x_2^1 & y_1^1 y_2^1 & y_1^1 & x_2^1 & y_2^1 & 1 \\ x_1^2 x_2^2 & x_1^2 y_2^2 & x_1^2 & y_1^2 x_2^2 & y_1^2 y_2^2 & y_1^2 & x_2^2 & y_2^2 & 1 \\ x_1^3 x_2^3 & x_1^3 y_2^3 & x_1^3 & y_1^3 x_2^3 & y_1^3 y_2^3 & y_1^3 & x_2^3 & y_2^3 & 1 \\ x_1^4 x_2^4 & x_1^4 y_2^4 & x_1^4 & y_1^4 x_2^4 & y_1^4 y_2^4 & y_1^4 & x_2^4 & y_2^4 & 1 \\ x_1^5 x_2^5 & x_1^5 y_2^5 & x_1^5 & y_1^5 x_2^5 & y_1^5 y_2^5 & y_1^5 & x_2^5 & y_2^5 & 1 \\ x_1^6 x_2^6 & x_1^6 y_2^6 & x_1^6 & y_1^6 x_2^6 & y_1^6 y_2^6 & y_1^6 & x_2^6 & y_2^6 & 1 \\ x_1^7 x_2^7 & x_1^7 y_2^7 & x_1^7 & y_1^7 x_2^7 & y_1^7 y_2^7 & y_1^7 & x_2^7 & y_2^7 & 1 \\ x_1^8 x_2^8 & x_1^8 y_2^8 & x_1^8 & y_1^8 x_2^8 & y_1^8 y_2^8 & y_1^8 & x_2^8 & y_2^8 & 1 \end{bmatrix} = 0 \quad \xrightarrow{\text{blue arrow}} \quad \mathbf{A}\mathbf{f} = \mathbf{0}$$

Direct Linear Transformation (DLT)

Algebraic error vs. geometric error

- Algebraic Error

$$\min \| \mathbf{A}\mathbf{f} \|$$

- Geometric Error (better)

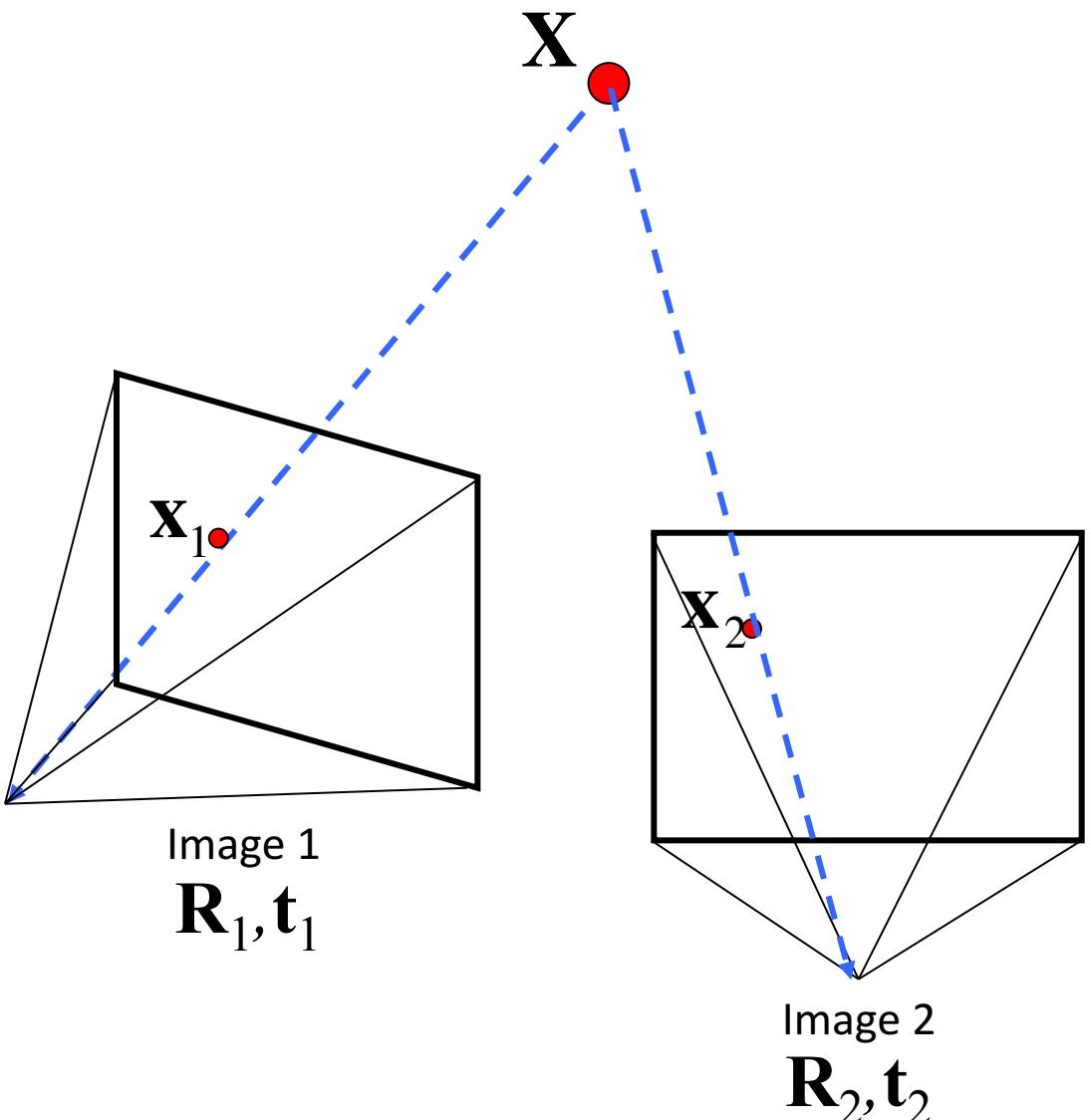
$$\min \sum_j d\left(\mathbf{x}_1^j, \mathbf{F}\mathbf{x}_2^j\right)^2 + d\left(\mathbf{x}_2^j, \mathbf{F}^T\mathbf{x}_1^j\right)^2 \quad \text{Unit: pixel}$$

Solved by (non-linear) least square solver (e.g. Ceres)

RANSAC to estimate the fundamental matrix

- For many times
 - Pick 8 points
 - Compute a solution for \mathbf{F} using these 8 points
 - Count number of inliers that with $\mathbf{x}_1^T \mathbf{F} \mathbf{x}_2$ close to 0
- Pick the one with the largest number of inliers

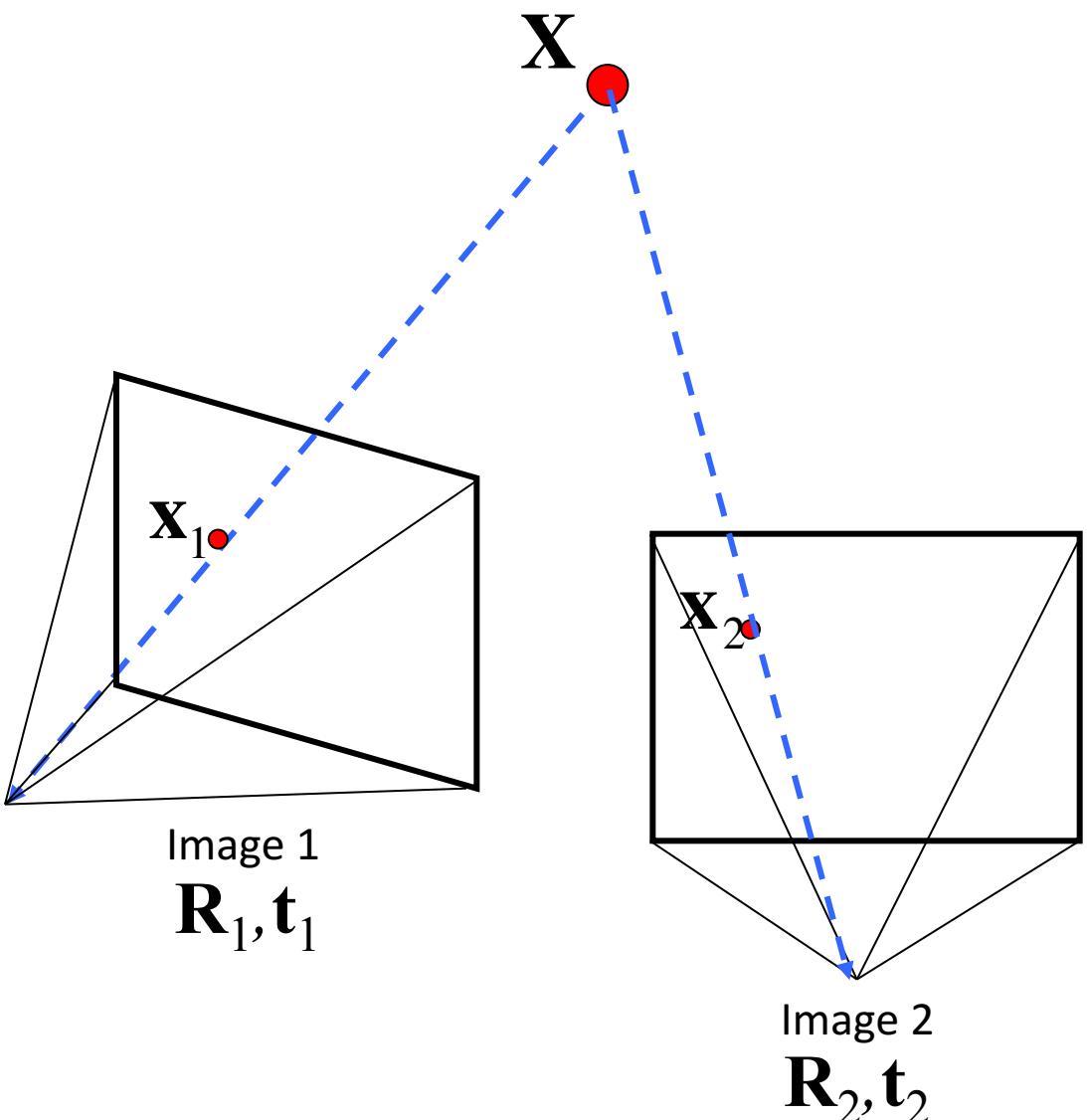
Fundamental matrix → essential matrix



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

$$\mathbf{E} = \mathbf{K}_1^T \mathbf{F} \mathbf{K}_2$$

Essential matrix $\rightarrow [\mathbf{R} | \mathbf{t}]$



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

$$\mathbf{E} = \mathbf{K}_1^T \mathbf{F} \mathbf{K}_2$$

Essential matrix → $[R|t]$

Result 9.19. For a given essential matrix

$$E = U \text{diag}(1, 1, 0) V^T,$$

and the first camera matrix $P_1 = [I|0]$, there are four possible choices for the second camera matrix P_2 :

$$P_2 = [UWV^T | +u_3]$$

$$P_2 = [UWV^T | -u_3]$$

$$P_2 = [UW^T V^T | +u_3]$$

$$P_2 = [UW^T V^T | -u_3]$$

$$W = \begin{bmatrix} 0 & -1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Four Possible Solutions

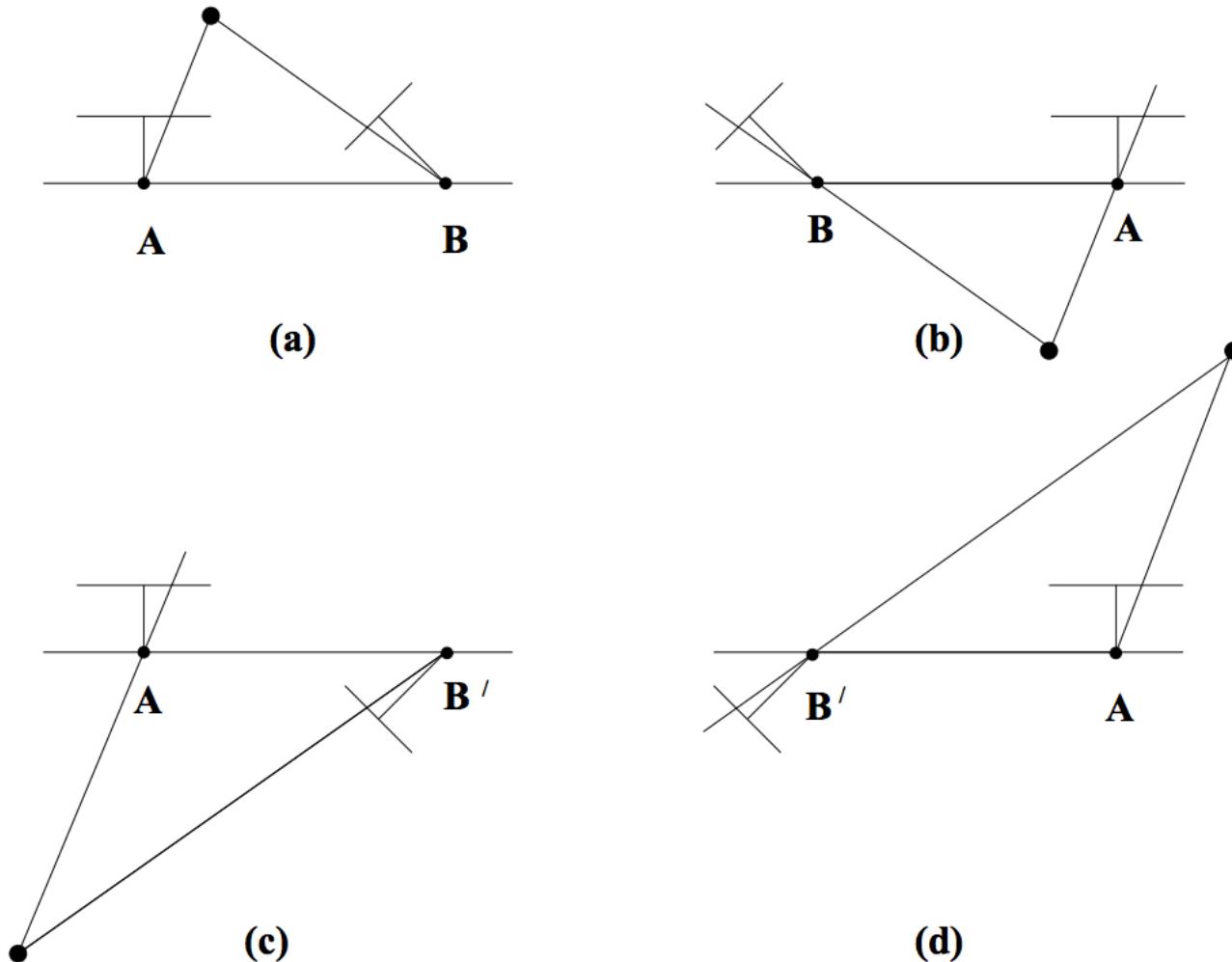
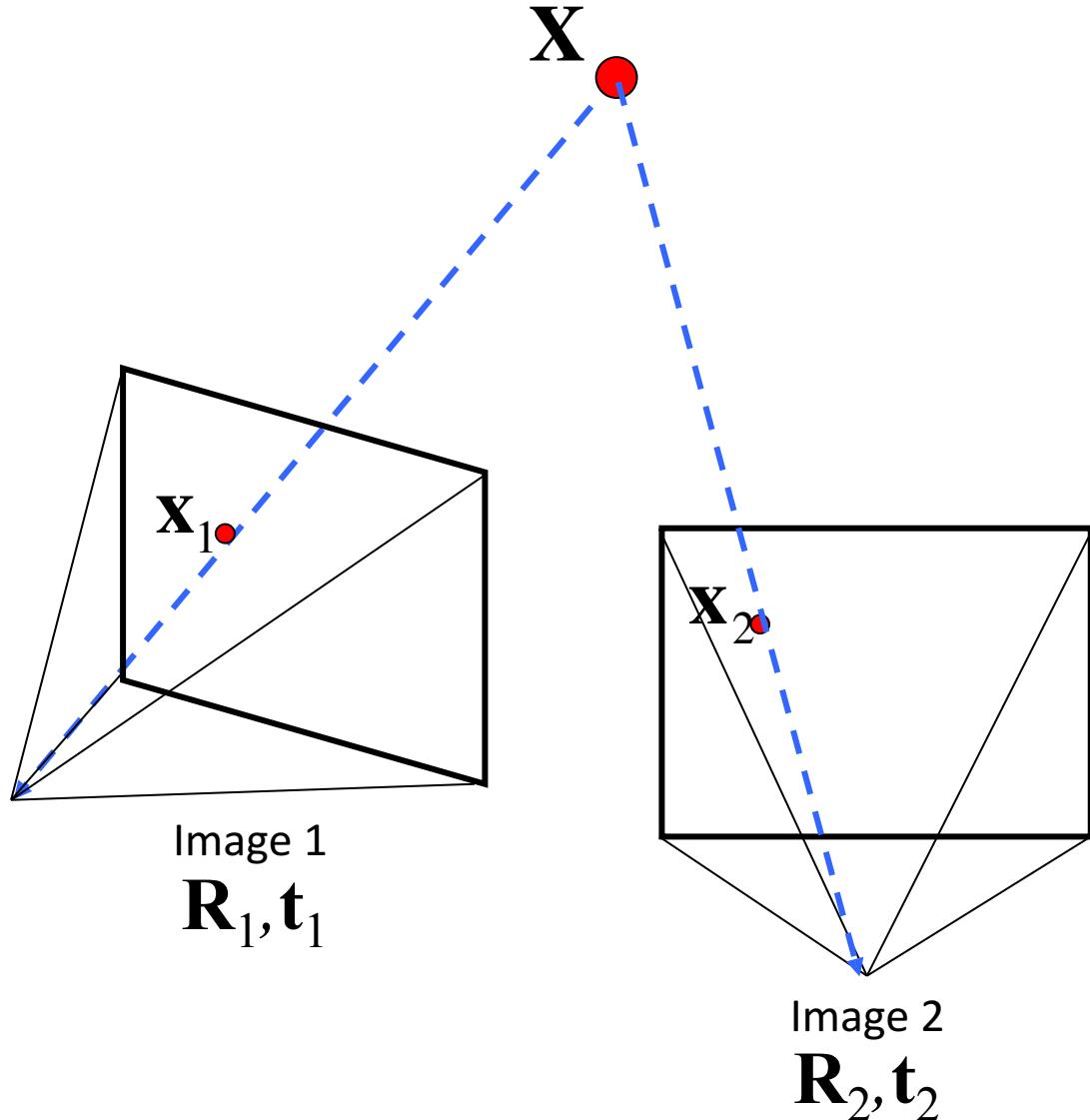


Fig. 9.12. **The four possible solutions for calibrated reconstruction from E.** Between the left and right sides there is a baseline reversal. Between the top and bottom rows camera B rotates 180° about the baseline. Note, only in (a) is the reconstructed point in front of both cameras.

Page 260 of the bible (Multiple View Geometry, 2nd Ed)

CM50248 – Visual Understanding 1 – Lecture 6

Triangulation



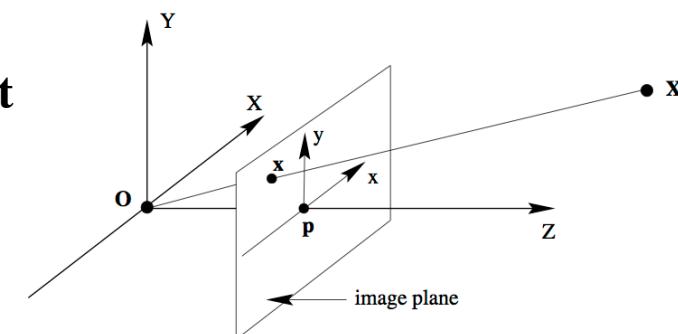
In front of the camera?

- Camera Extrinsic $[\mathbf{R}|\mathbf{t}]$

$$\begin{bmatrix} X_{cam} \\ Y_{cam} \\ Z_{cam} \end{bmatrix} = \mathbf{R} \begin{bmatrix} X_{world} \\ Y_{world} \\ Z_{world} \end{bmatrix} + \mathbf{t} \quad \leftrightarrow \quad \begin{bmatrix} X_{world} \\ Y_{world} \\ Z_{world} \end{bmatrix} = \mathbf{R}^{-1} \left(\begin{bmatrix} X_{cam} \\ Y_{cam} \\ Z_{cam} \end{bmatrix} - \mathbf{t} \right) = \mathbf{R}^T \begin{bmatrix} X_{cam} \\ Y_{cam} \\ Z_{cam} \end{bmatrix} - \mathbf{R}^T \mathbf{t}$$

- Camera Center

$$\begin{bmatrix} X_{cam} \\ Y_{cam} \\ Z_{cam} \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \rightarrow \quad \mathbf{C} = \begin{bmatrix} X_{world} \\ Y_{world} \\ Z_{world} \end{bmatrix} = \mathbf{R}^T \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} - \mathbf{R}^T \mathbf{t} = -\mathbf{R}^T \mathbf{t}$$



- View Direction

$$\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \rightarrow \quad \left(\mathbf{R}^T \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} - \mathbf{R}^T \mathbf{t} \right) - (\mathbf{C}) = \left(\mathbf{R}(3,:)^T - \mathbf{R}^T \mathbf{t} \right) - \left(-\mathbf{R}^T \mathbf{t} \right) = \mathbf{R}(3,:)^T$$

Camera Coordinate System

World Coordinate System

In front of the camera?

- A point \mathbf{X}
- Direction from camera center to point $\mathbf{X} - \mathbf{C}$
- Angle between two vectors

$$\mathbf{A} \cdot \mathbf{B} = \|\mathbf{A}\| \|\mathbf{B}\| \cos \theta$$

- Angle between $\mathbf{X} - \mathbf{C}$ and view virection
- Just need to test $(\mathbf{X} - \mathbf{C}) \cdot \mathbf{R}(3,:)^T > 0$?

Pick the solution

with maximal number of points in front of both cameras

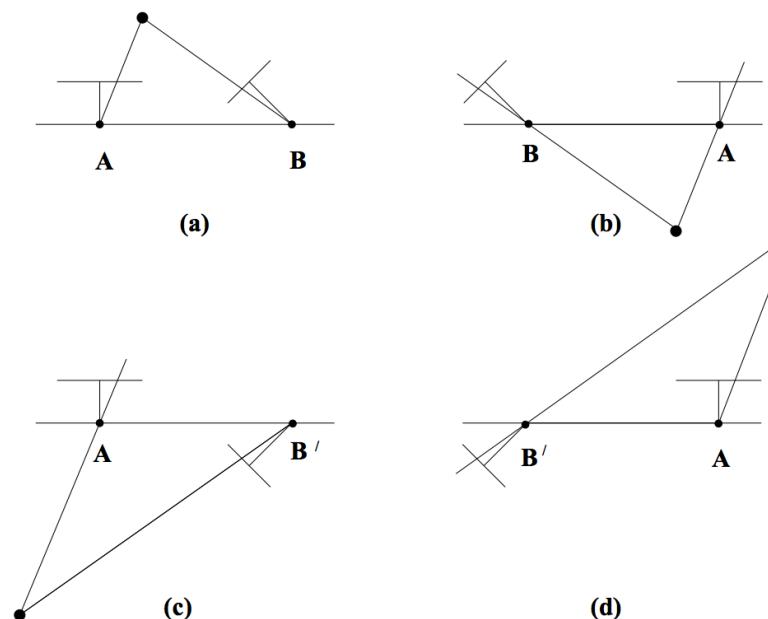
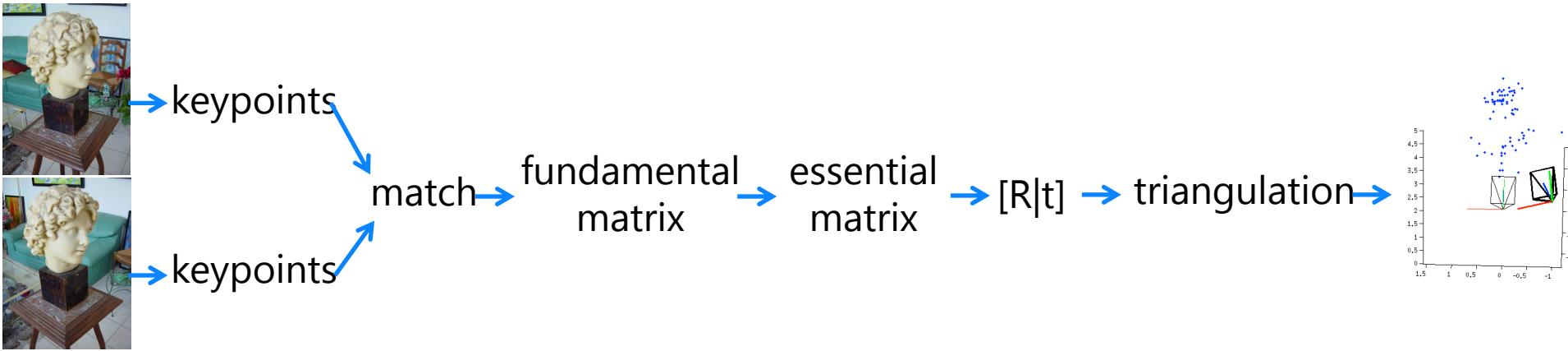
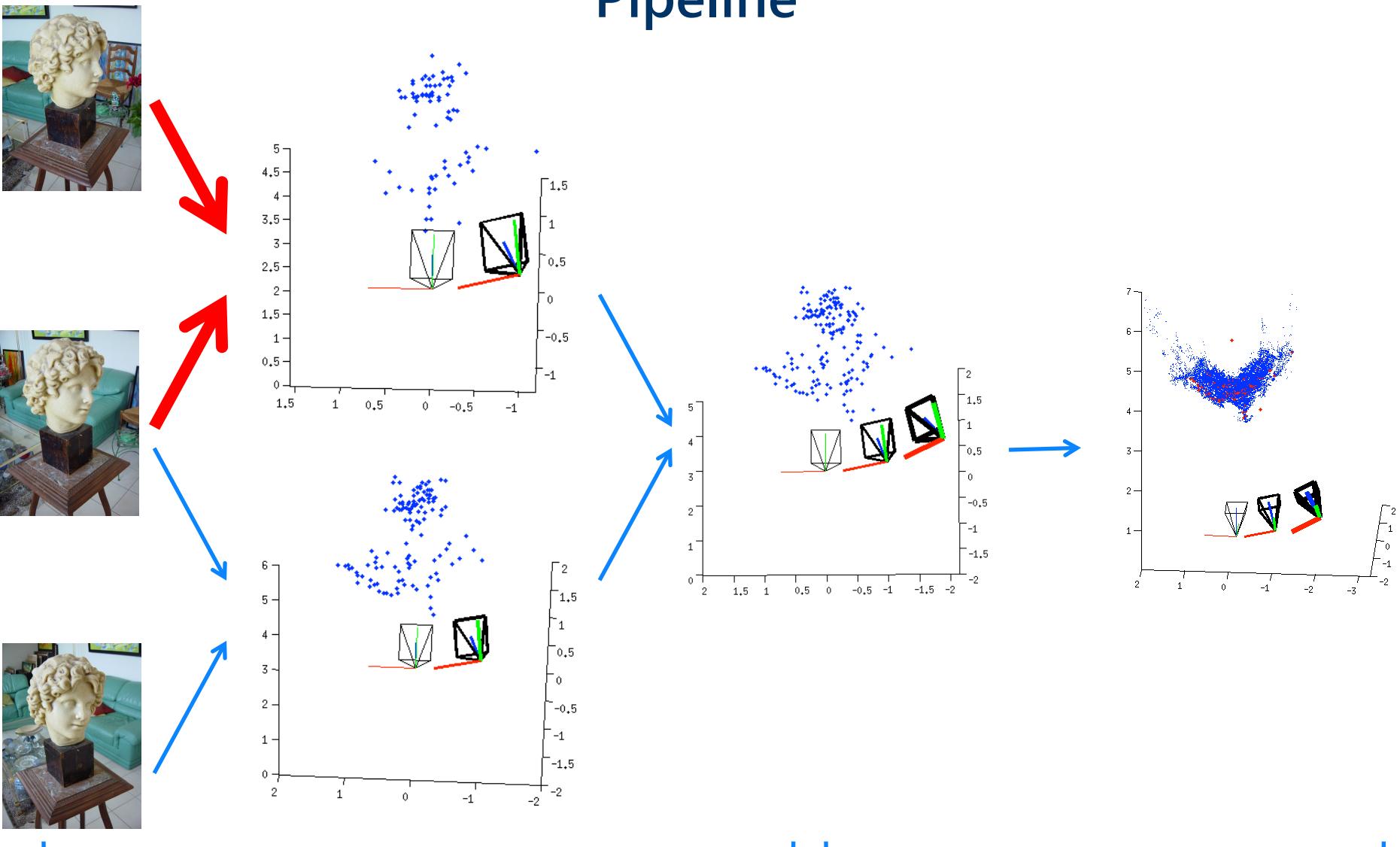


Fig. 9.12. The four possible solutions for calibrated reconstruction from E. Between the left and right sides there is a baseline reversal. Between the top and bottom rows camera B rotates 180° about the baseline. Note, only in (a) is the reconstructed point in front of both cameras.

Two-view reconstruction



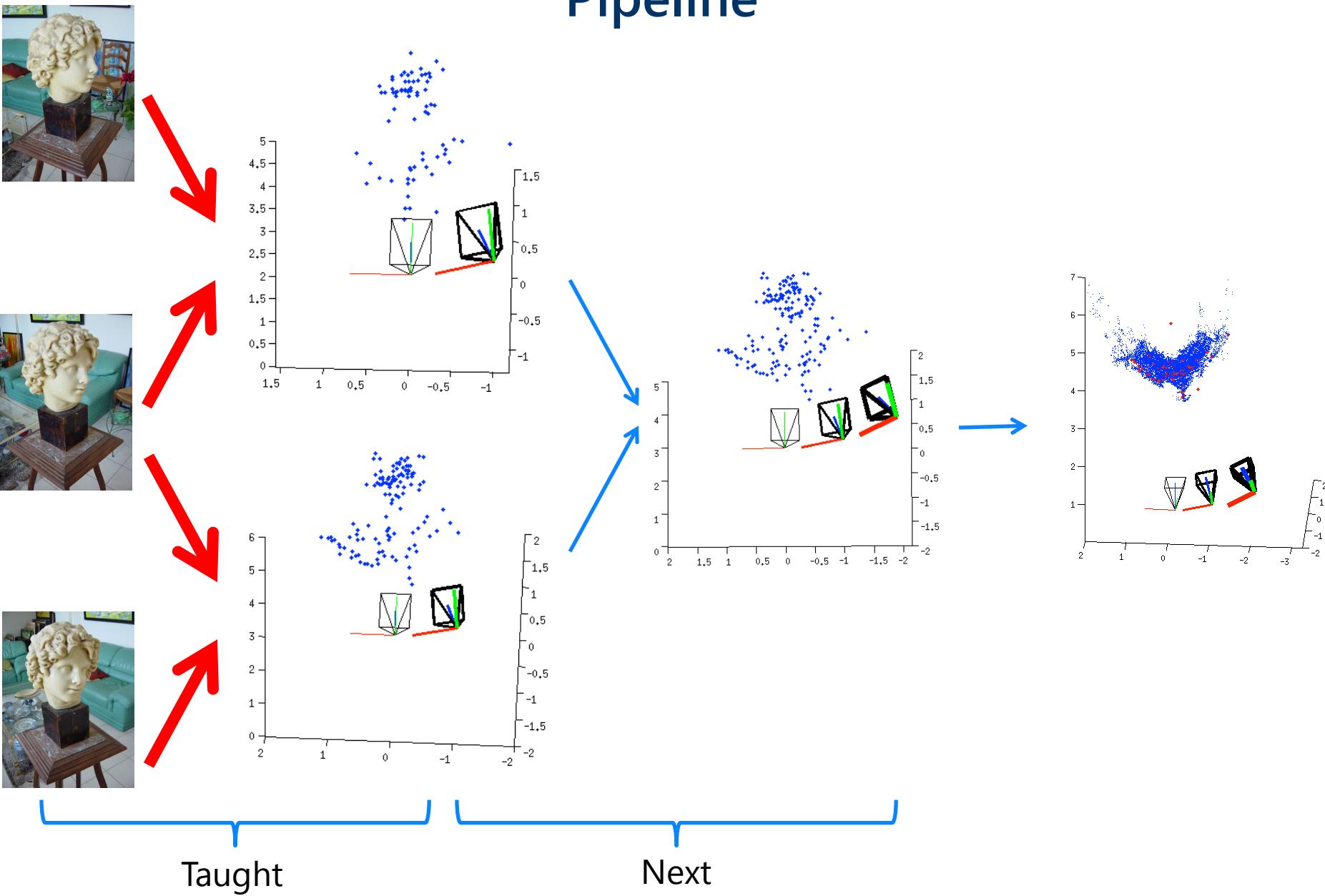
Pipeline



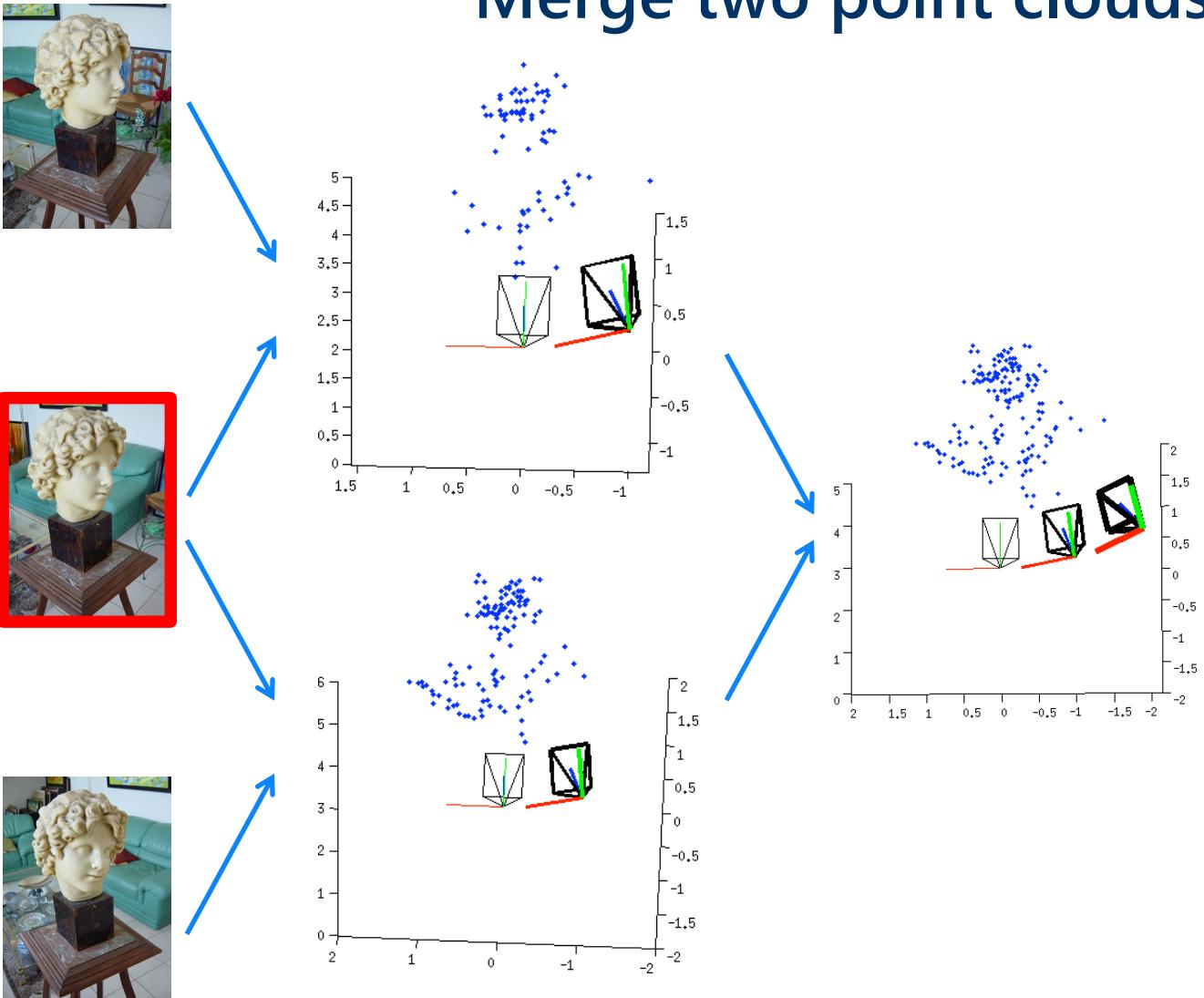
Structure from Motion (SFM)

Multi-view Stereo (MVS)

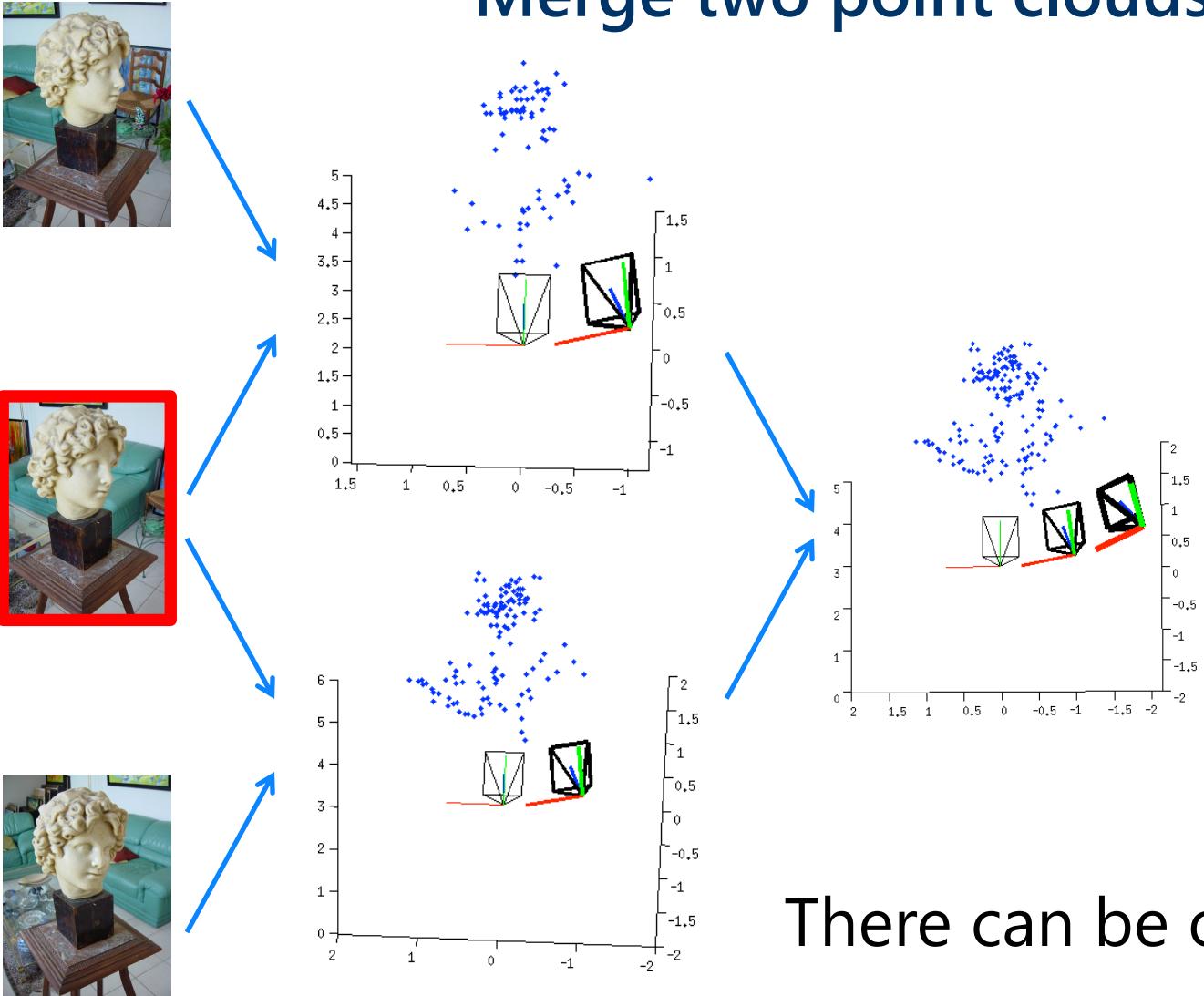
Pipeline



Merge two point clouds



Merge two point clouds



There can be only one $[R_2 | t_2]$

Merge two point clouds

- From the 1st and 2nd images, we have

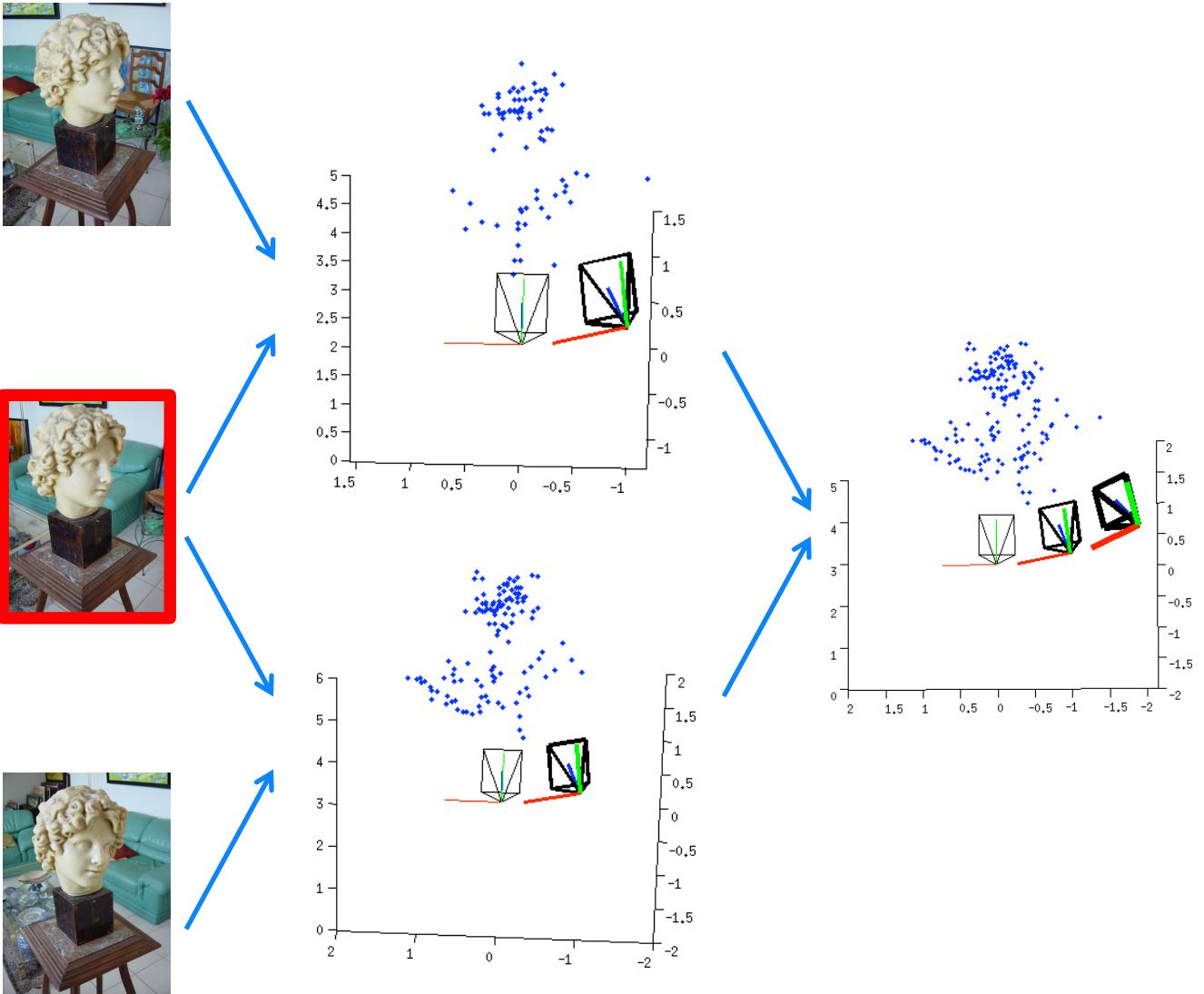
$$[\mathbf{R}_1 | \mathbf{t}_1] \text{ and } [\mathbf{R}_2 | \mathbf{t}_2]$$

- From the 2nd and 3rd images, we have

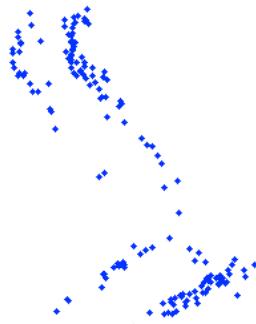
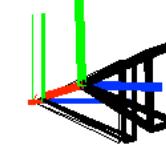
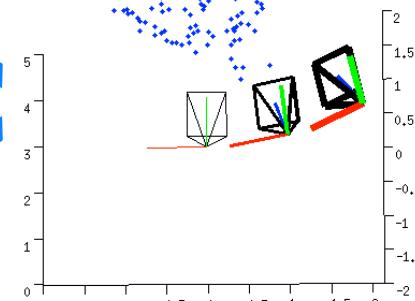
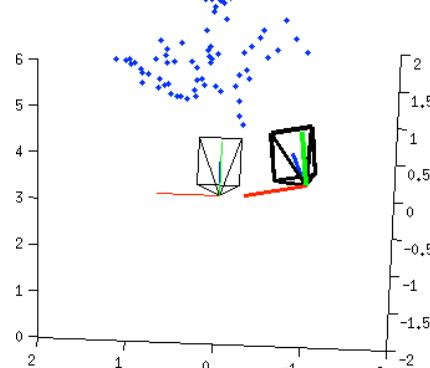
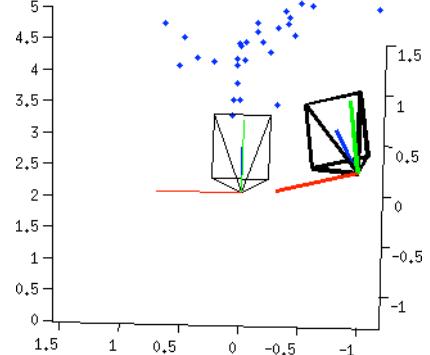
$$[\mathbf{R}_2 | \mathbf{t}_2] \text{ and } [\mathbf{R}_3 | \mathbf{t}_3]$$

- **Exercise:** How to transform the coordinate system of the second point cloud to align with the first point cloud so that there is only one $[\mathbf{R}_2 | \mathbf{t}_2]$?

Merge two point clouds

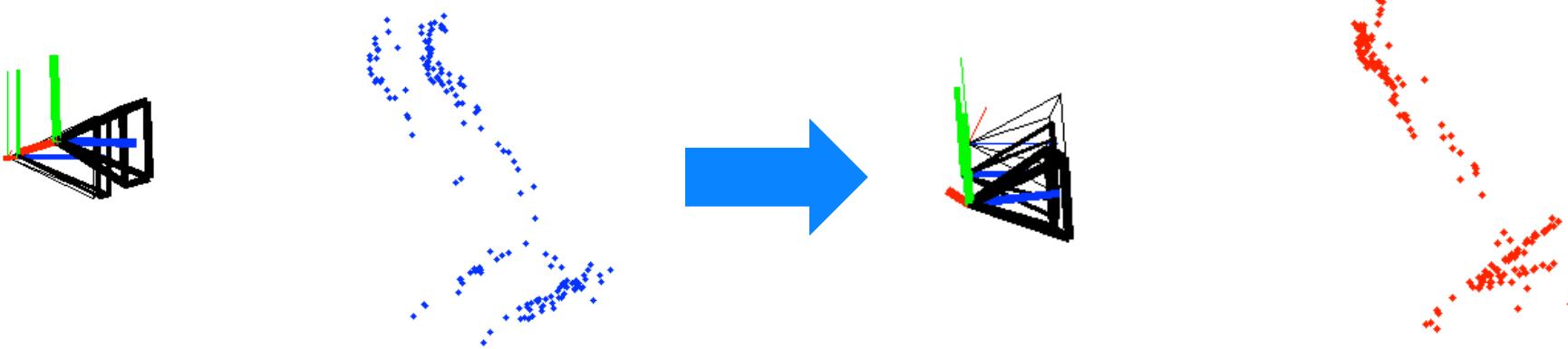


Oops

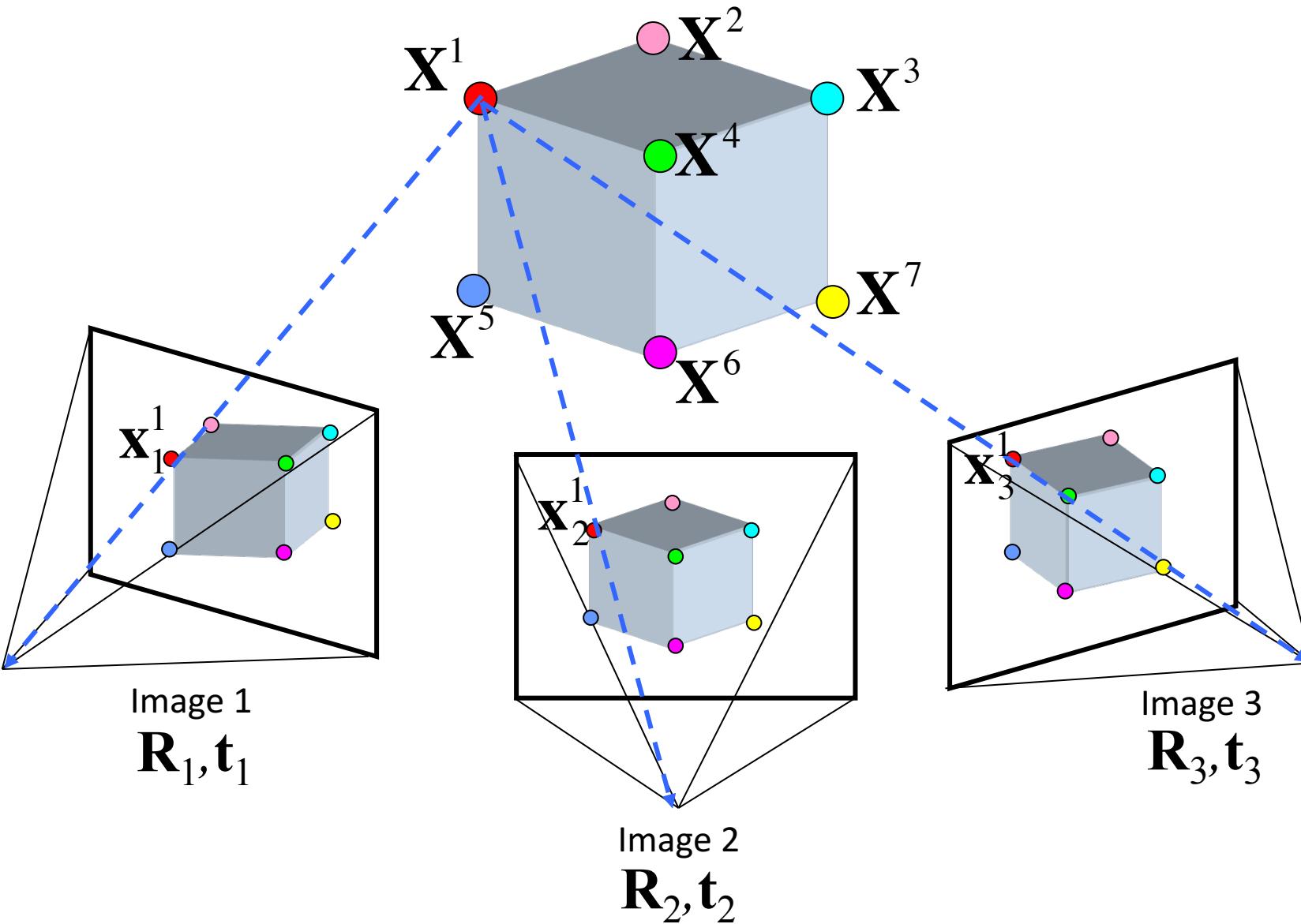


Seen from a different angle

Bundle adjustment



"America's Next Top Model"



“America's Next Top Model”

	Point 1	Point 2	Point 3
Image 1	$\mathbf{x}_1^1 = \mathbf{K}[\mathbf{R}_1 \mathbf{t}_1] \mathbf{X}^1$	$\mathbf{x}_1^2 = \mathbf{K}[\mathbf{R}_1 \mathbf{t}_1] \mathbf{X}^2$	
Image 2	$\mathbf{x}_2^1 = \mathbf{K}[\mathbf{R}_2 \mathbf{t}_2] \mathbf{X}^1$	$\mathbf{x}_2^2 = \mathbf{K}[\mathbf{R}_2 \mathbf{t}_2] \mathbf{X}^2$	$\mathbf{x}_2^3 = \mathbf{K}[\mathbf{R}_2 \mathbf{t}_2] \mathbf{X}^3$
Image 3	$\mathbf{x}_3^1 = \mathbf{K}[\mathbf{R}_3 \mathbf{t}_3] \mathbf{X}^1$		$\mathbf{x}_3^3 = \mathbf{K}[\mathbf{R}_3 \mathbf{t}_3] \mathbf{X}^3$

Rethinking the SFM problem

- Input: observed 2D image position

$$\begin{array}{cc} \tilde{\mathbf{x}}_1^1 & \tilde{\mathbf{x}}_1^2 \\ \tilde{\mathbf{x}}_2^1 & \tilde{\mathbf{x}}_2^2 & \tilde{\mathbf{x}}_2^3 \\ \tilde{\mathbf{x}}_3^1 & & \tilde{\mathbf{x}}_3^3 \end{array}$$

- Output:

Unknown camera parameters (with some guess)

$$[\mathbf{R}_1|\mathbf{t}_1], [\mathbf{R}_2|\mathbf{t}_2], [\mathbf{R}_3|\mathbf{t}_3]$$

Unknown point 3D coordinate (with some guess)

$$\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \dots$$

Bundle adjustment

A valid solution $[\mathbf{R}_1|\mathbf{t}_1], [\mathbf{R}_2|\mathbf{t}_2], [\mathbf{R}_3|\mathbf{t}_3]$ and $\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \dots$

must let

Re-projection $=$

$$\begin{cases} \mathbf{x}_1^1 = \mathbf{K}[\mathbf{R}_1|\mathbf{t}_1]\mathbf{X}^1 & \mathbf{x}_1^2 = \mathbf{K}[\mathbf{R}_1|\mathbf{t}_1]\mathbf{X}^2 \\ \mathbf{x}_2^1 = \mathbf{K}[\mathbf{R}_2|\mathbf{t}_2]\mathbf{X}^1 & \mathbf{x}_2^2 = \mathbf{K}[\mathbf{R}_2|\mathbf{t}_2]\mathbf{X}^2 & \mathbf{x}_2^3 = \mathbf{K}[\mathbf{R}_2|\mathbf{t}_2]\mathbf{X}^3 \\ \mathbf{x}_3^1 = \mathbf{K}[\mathbf{R}_3|\mathbf{t}_3]\mathbf{X}^1 & & \mathbf{x}_3^3 = \mathbf{K}[\mathbf{R}_3|\mathbf{t}_3]\mathbf{X}^3 \end{cases}$$

Observation

$$\begin{cases} \tilde{\mathbf{x}}_1^1 & \tilde{\mathbf{x}}_1^2 \\ \tilde{\mathbf{x}}_2^1 & \tilde{\mathbf{x}}_2^2 & \tilde{\mathbf{x}}_2^3 \\ \tilde{\mathbf{x}}_3^1 & & \tilde{\mathbf{x}}_3^3 \end{cases}$$

Bundle adjustment

A valid solution $[\mathbf{R}_1|\mathbf{t}_1], [\mathbf{R}_2|\mathbf{t}_2], [\mathbf{R}_3|\mathbf{t}_3]$ and $\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \dots$ must let the reprojection close to the observation, i.e. to minimize the reprojection error

$$\min \sum_i \sum_j (\tilde{\mathbf{x}}_i^j - \mathbf{K}[\mathbf{R}_i | \mathbf{t}_i] \mathbf{X}^j)^2$$

Bundle adjustment

A valid solution $[\mathbf{R}_1|\mathbf{t}_1], [\mathbf{R}_2|\mathbf{t}_2], [\mathbf{R}_3|\mathbf{t}_3]$ and $\mathbf{X}^1, \mathbf{X}^2, \mathbf{X}^3, \dots$ must let the reprojection close to the observation, i.e. to minimize the reprojection error

$$\min \sum_i \sum_j (\tilde{\mathbf{x}}_i^j - \mathbf{K}[\mathbf{R}_i | \mathbf{t}_i] \mathbf{X}^j)^2$$

Question: What is the unit of this objective function?

Solving this optimization problem

- Theory:

The Levenberg–Marquardt algorithm

http://en.wikipedia.org/wiki/Levenberg-Marquardt_algorithm

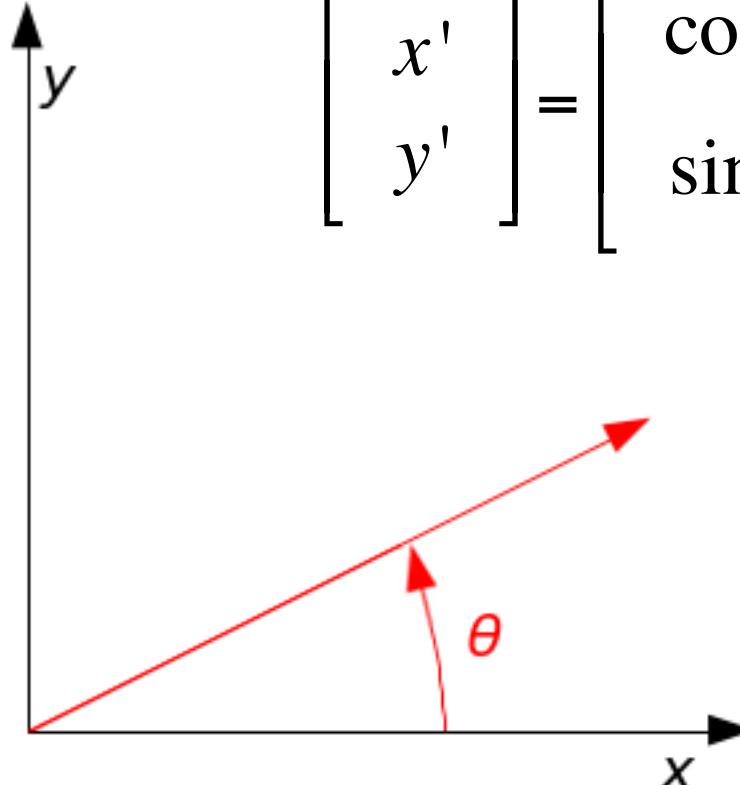
- Practice:

The Ceres-Solver from Google

<http://code.google.com/p/ceres-solver/>

Parameterizing rotation matrix

- 2D Rotation

$$\begin{bmatrix} x' \\ y' \end{bmatrix} = \begin{bmatrix} \cos(\theta) & -\sin(\theta) \\ \sin(\theta) & \cos(\theta) \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}$$


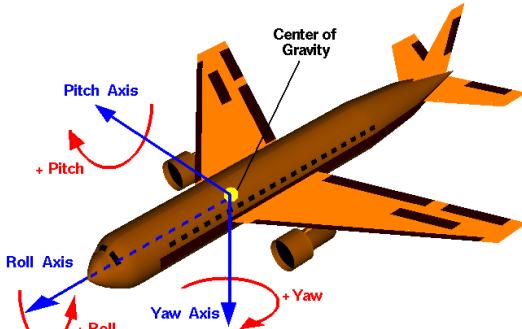
A 2D Cartesian coordinate system with x and y axes. A red vector originates from the origin. The angle between the positive x-axis and the vector is labeled θ . A small red arrow points along the vector to indicate its orientation.

$$\mathbf{R}^T = \mathbf{R}^{-1}, \det \mathbf{R} = 1$$

3D rotation

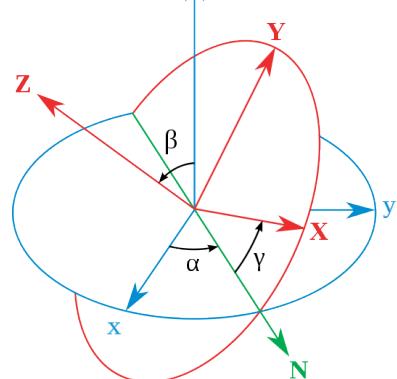
$$\mathbf{R}_x(\theta) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos(\theta) & -\sin(\theta) \\ 0 & \sin(\theta) & \cos(\theta) \end{bmatrix}$$
$$\mathbf{R}_y(\theta) = \begin{bmatrix} \cos(\theta) & 0 & \sin(\theta) \\ 0 & 1 & 0 \\ -\sin(\theta) & 0 & \cos(\theta) \end{bmatrix}$$
$$\mathbf{R}_z(\theta) = \begin{bmatrix} \cos(\theta) & -\sin(\theta) & 0 \\ \sin(\theta) & \cos(\theta) & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Yaw, pitch and roll are α, β, γ



$$\mathbf{R}_x(\gamma)\mathbf{R}_y(\beta)\mathbf{R}_z(\alpha)$$

Euler angles are α, β, γ

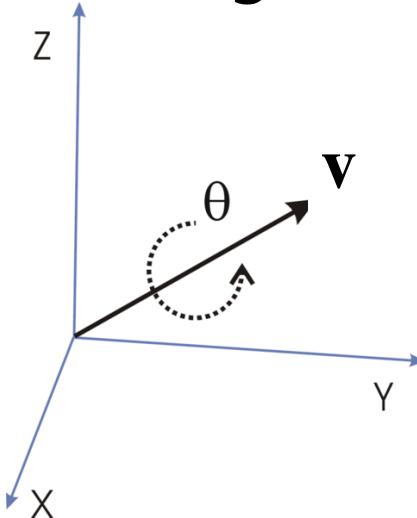


$$\mathbf{R}_z(\gamma)\mathbf{R}_x(\beta)\mathbf{R}_y(\alpha)$$

http://en.wikipedia.org/wiki/Rotation_matrix

3D rotation

Axis-angle representation



Quaternions

$$\mathbf{q} = \left(\mathbf{v} \sin\left(\frac{\theta}{2}\right), \cos\left(\frac{\theta}{2}\right) \right)^T$$

Avoid Gimbal Lock!

Recommendation!

Triplet Representation $\mathbf{v}\theta$ (3 dof) ← Not over-parameterized

Rodrigues' rotation formula

$$\mathbf{k}_{rot} = \mathbf{k} \cos \theta + (\mathbf{v} \times \mathbf{k}) \sin \theta + \mathbf{v} (\mathbf{v} \cdot \mathbf{k}) (1 - \cos \theta)$$

http://en.wikipedia.org/wiki/Axis-angle_representation

http://en.wikipedia.org/wiki/Quaternions_and_spatial_rotation

http://en.wikipedia.org/wiki/Rodrigues%27_rotation_formula

Initialization matters

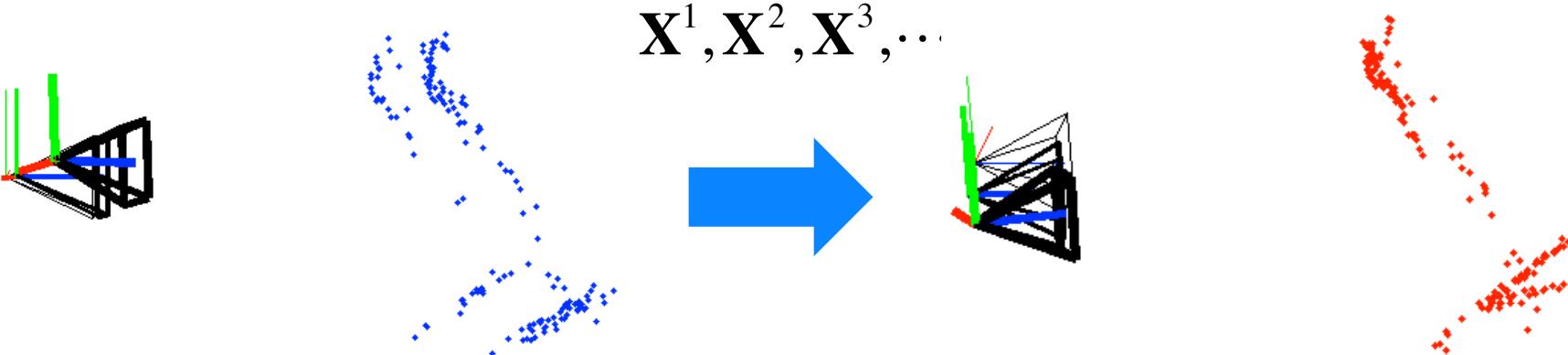
- Input: Observed 2D image position
- Output:

$$\begin{array}{cc} \tilde{\mathbf{x}}_1^1 & \tilde{\mathbf{x}}_1^2 \\ \tilde{\mathbf{x}}_2^1 & \tilde{\mathbf{x}}_2^2 \quad \tilde{\mathbf{x}}_2^3 \\ \tilde{\mathbf{x}}_3^1 & \tilde{\mathbf{x}}_3^3 \end{array}$$

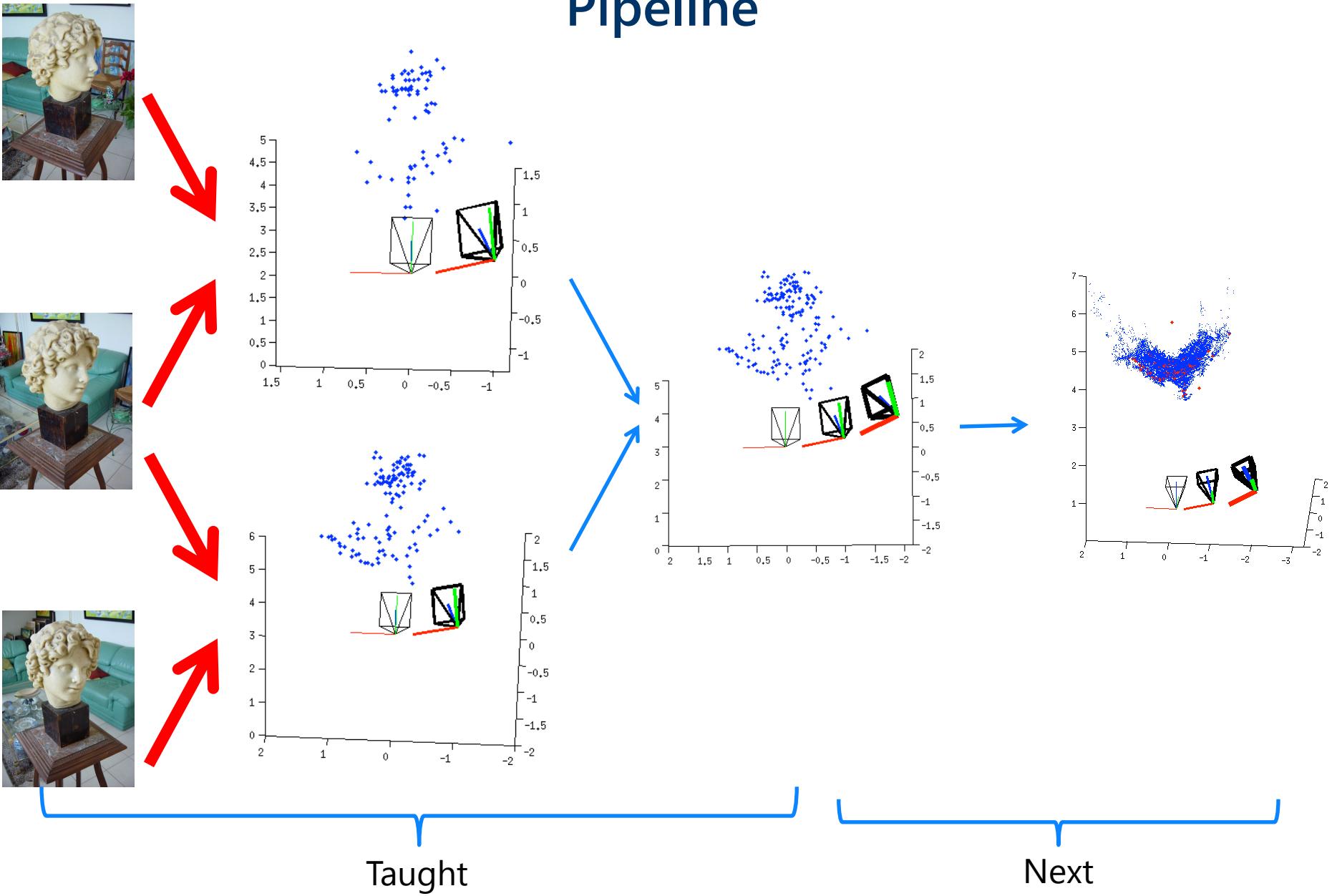
Unknown camera parameters (with some guess)

$$[\mathbf{R}_1|\mathbf{t}_1], [\mathbf{R}_2|\mathbf{t}_2], [\mathbf{R}_3|\mathbf{t}_3]$$

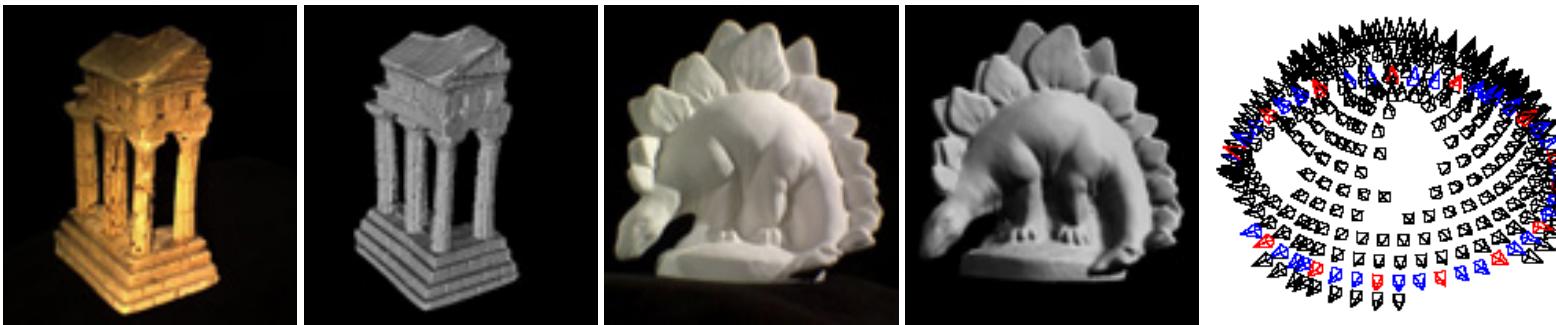
Unknown point 3D coordinate (with some guess)



Pipeline



Multi-view stereo



State-of-the-art:

COLMAP: <https://colmap.github.io/>

Pixelwise View Selection for Unstructured Multi-View Stereo
JL Schönberger, E Zheng, M Pollefeys, J-M Frahm, ECCV 2016.

Benchmark:

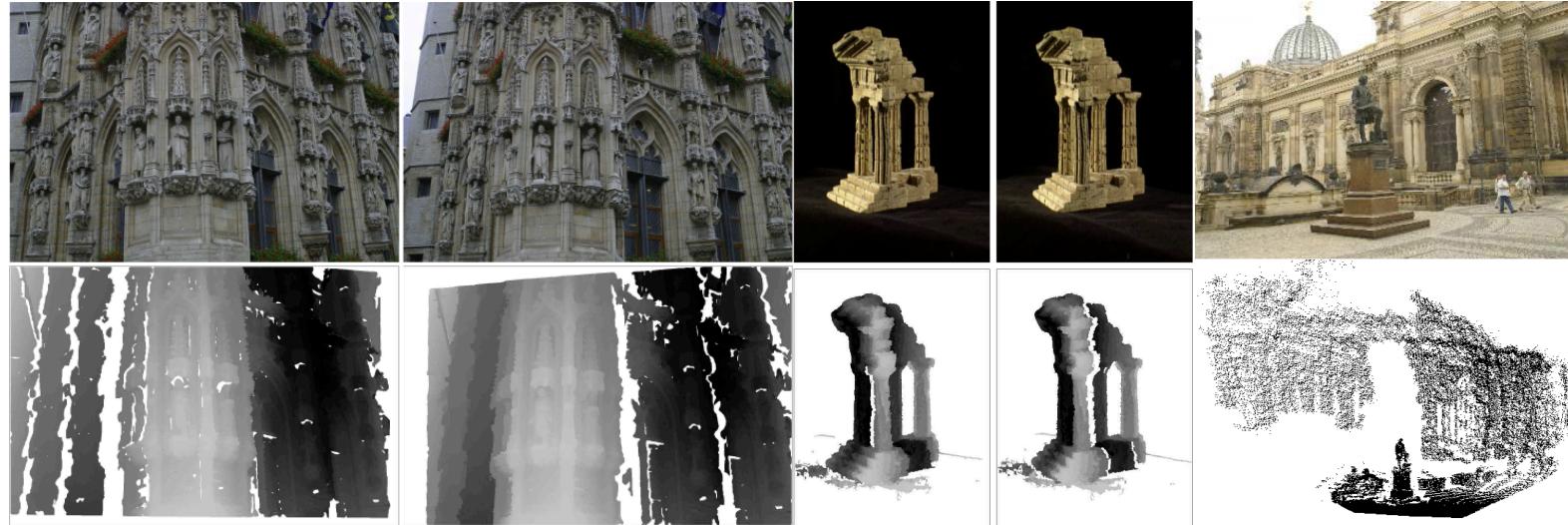
<http://vision.middlebury.edu/mview/>

A Comparison and Evaluation of Multi-View Stereo Reconstruction Algorithms.
SM Seitz, B Curless, J Diebel, D Scharstein, R Szeliski. CVPR 2006.

Baseline:

Multi-view stereo revisited. M Goesele, B Curless, SM Seitz. CVPR 2006.

Key idea: matching propagation

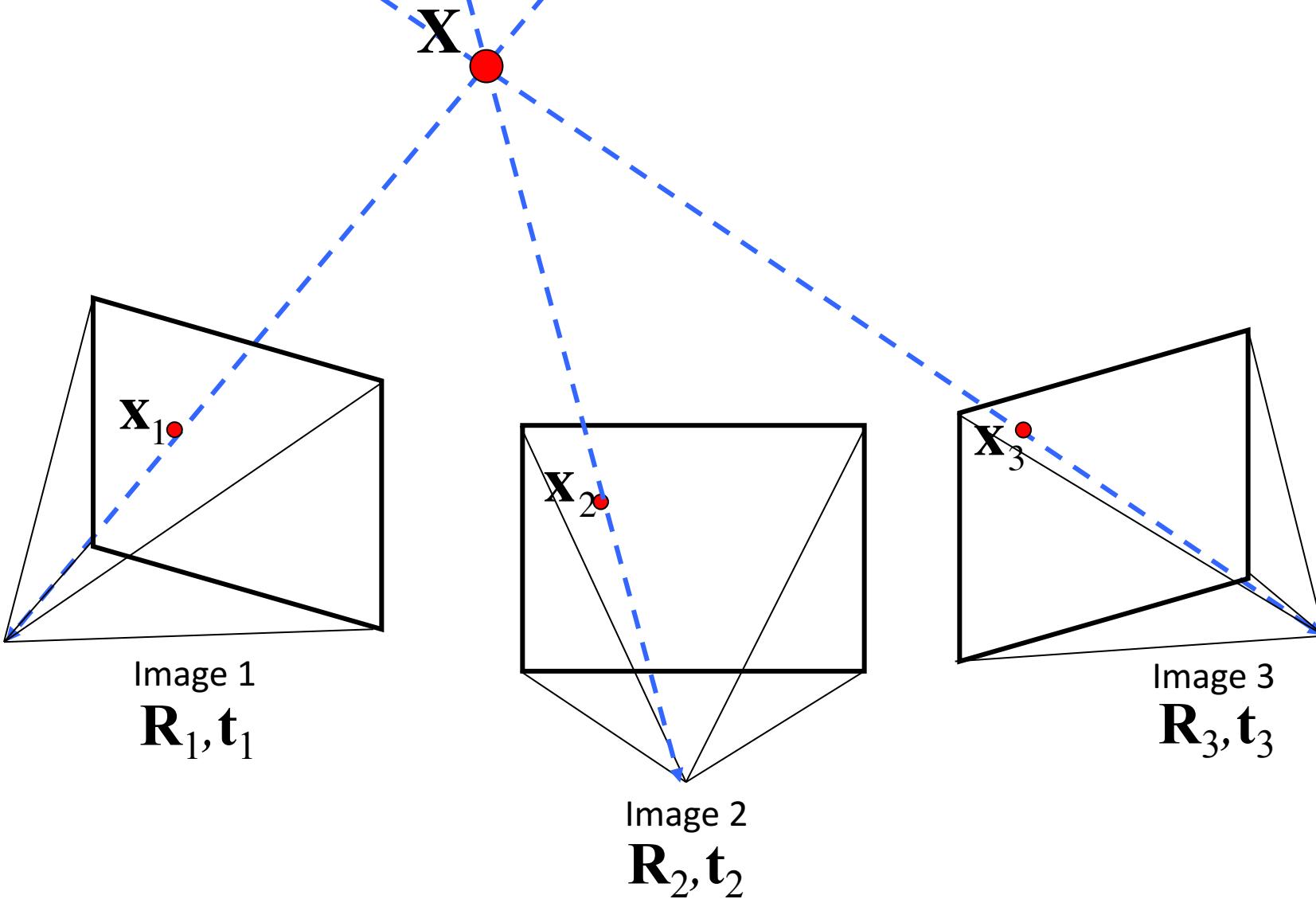


- [1] Learning Two-view Stereo Matching, J Xiao, J Chen, DY Yeung, and L Quan, 2008.
- [2] Accurate, Dense, and Robust Multi-View Stereopsis, Y Furukawa and J Ponce, 2007.
- [3] Multi-view stereo revisited. M Goesele, B Curless, SM Seitz. 2006.
- [4] Robust Dense Matching Using Local and Global Geometric Constraints, M Lhuillier & L Quan, 2000.

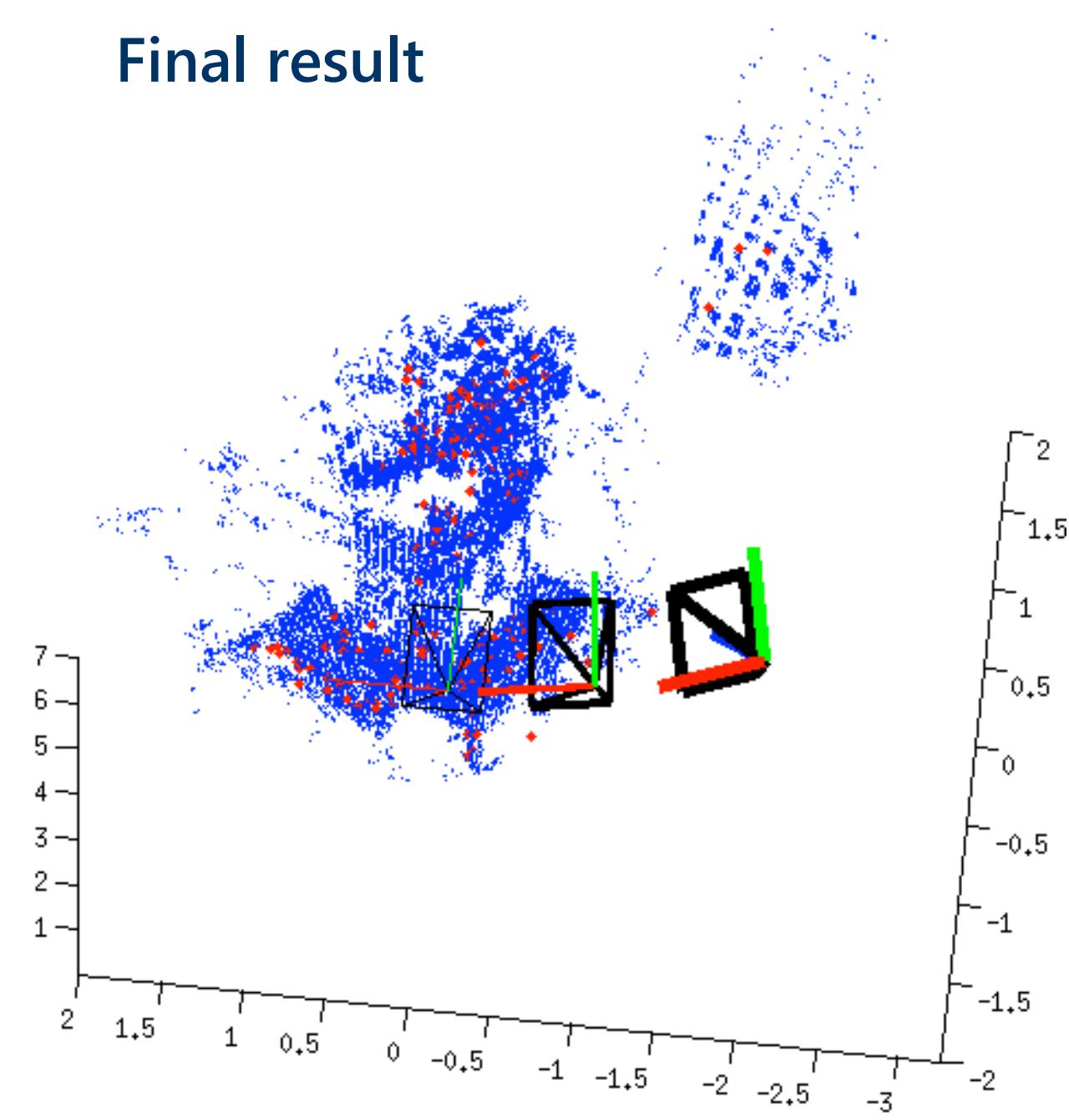
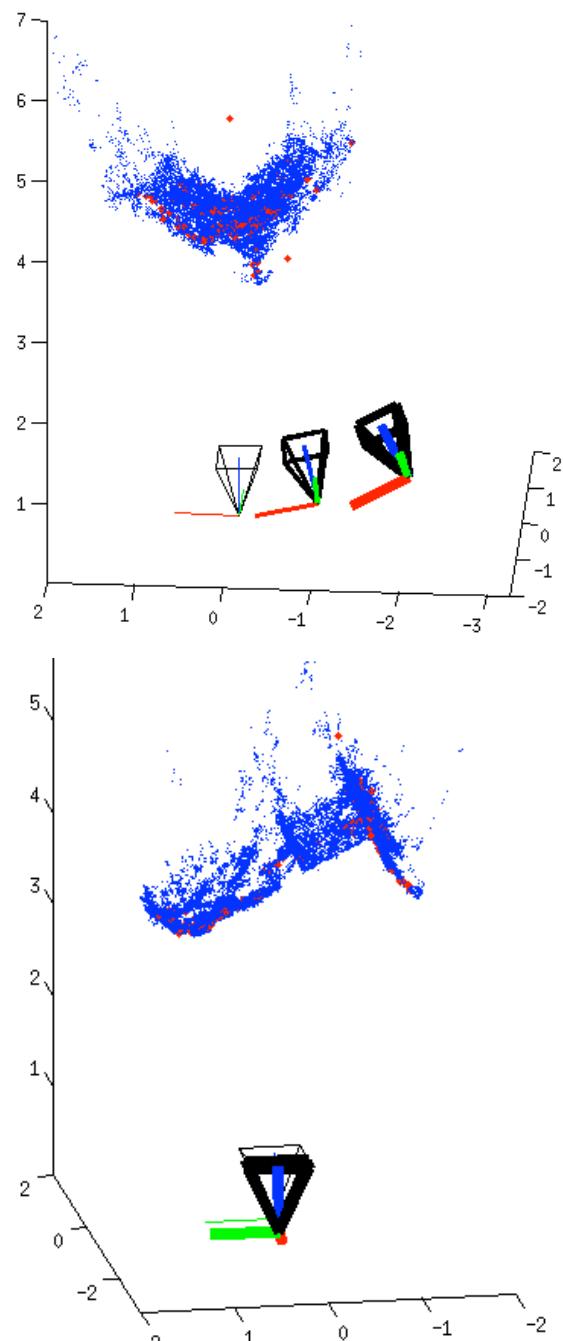
In another context:

- [5] PatchMatch: A Randomized Correspondence Algorithm for Structural Image Editing, C Barnes, E Shechtman, A Finkelstein, and DB Goldman, 2009

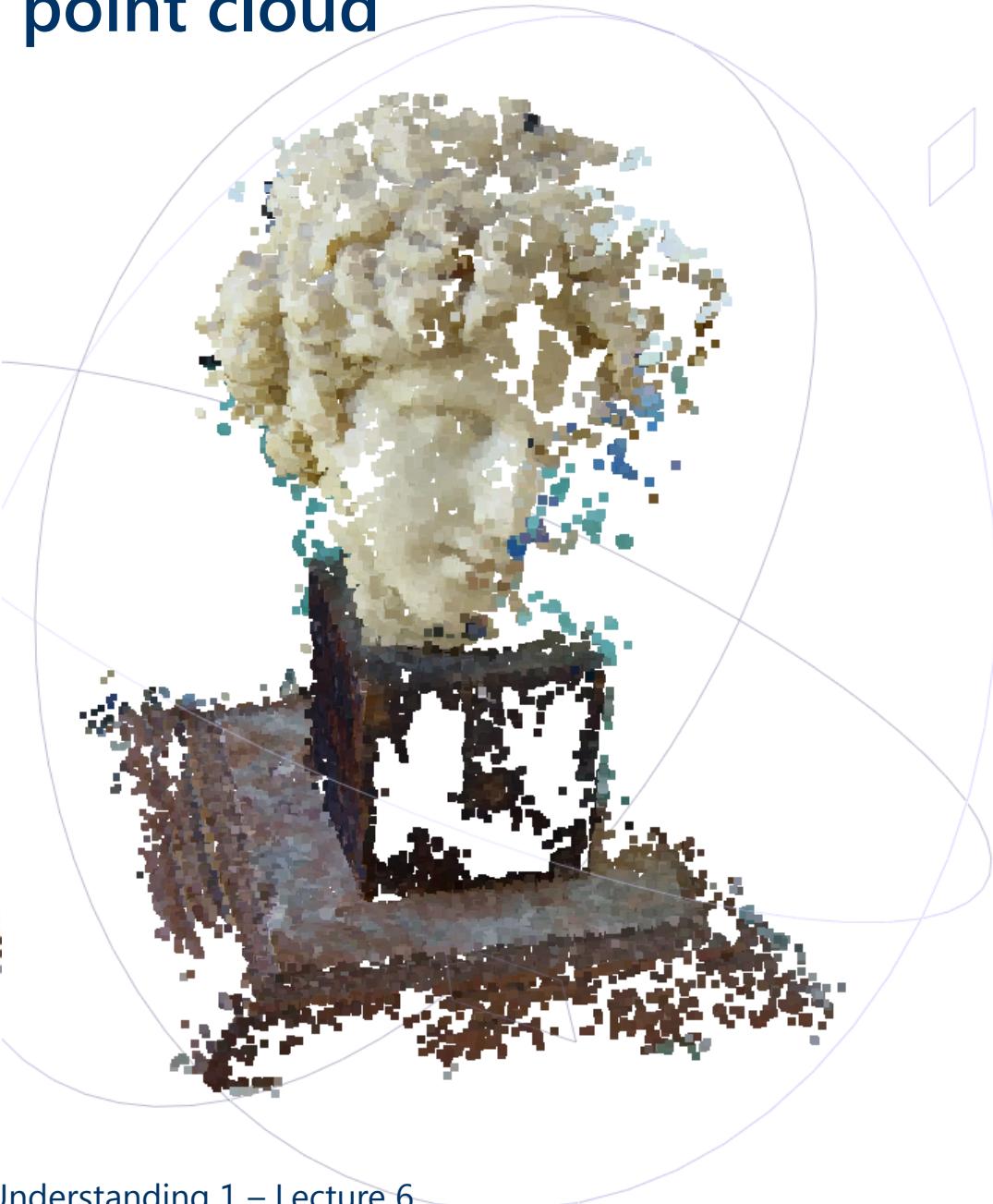
Triangulation



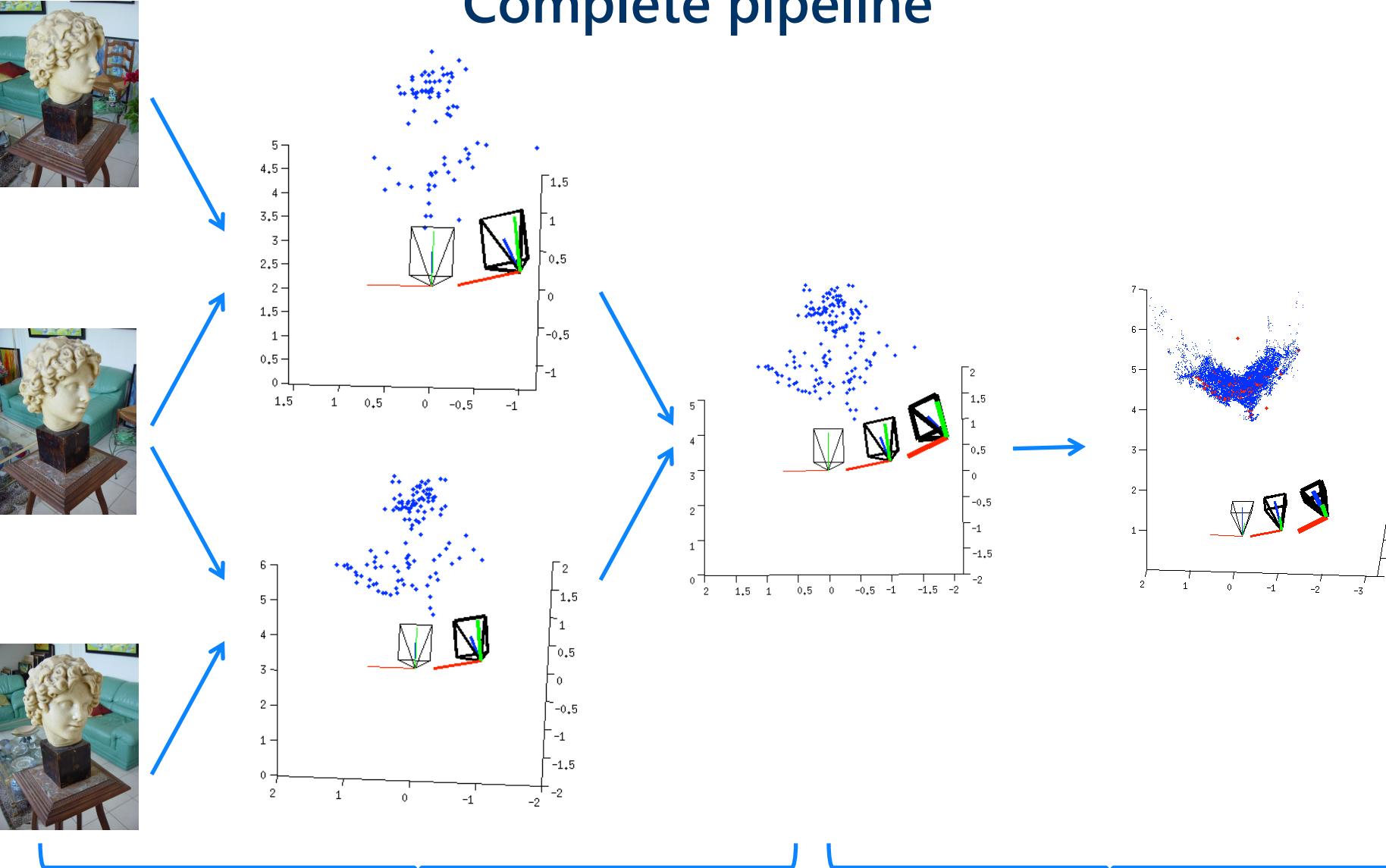
Final result



Colorize the point cloud



Complete pipeline



Structure from Motion (SFM)

Multi-view Stereo (MVS)

Wait — How to get the focal length?

- Auto-calibration
 - Self-calibration and metric reconstruction inspite of varying and unknown intrinsic camera parameters, M Pollefeys, R Koch and L Van Gool, IJCV 1999
- Grid search to look for the solution with minimal reprojection error
 - for $f = \text{min}_f : \text{max}_f$
 - do everything, then obtain reprojection error after bundle adjustment
- Optimize for this value in bundle adjustment
- Camera calibration (with checkerboard)
- EXIF of JPEG file recorded from digital camera
 - Read Bundler code to understand how to convert EXIF into focal length value
<http://phototour.cs.washington.edu/bundler/>

Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

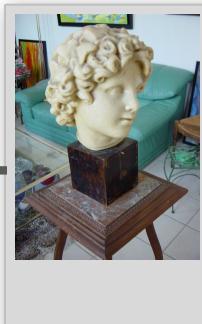
Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

= Images → Models: Image-based Modeling



+



+



=



Steps

Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

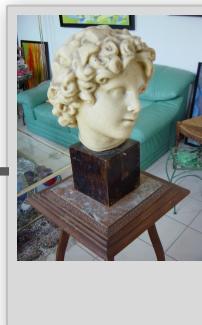
Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

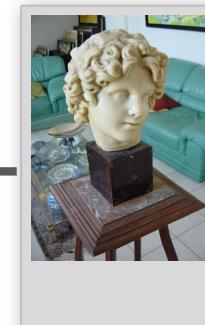
= Images → Models: Image-based Modeling



+



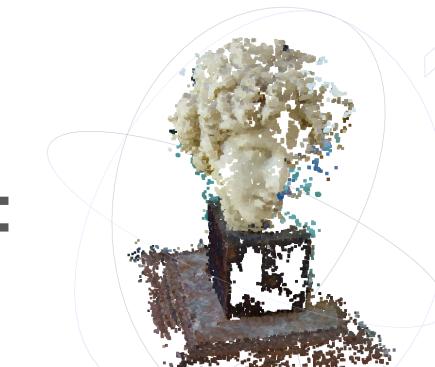
+



+



=

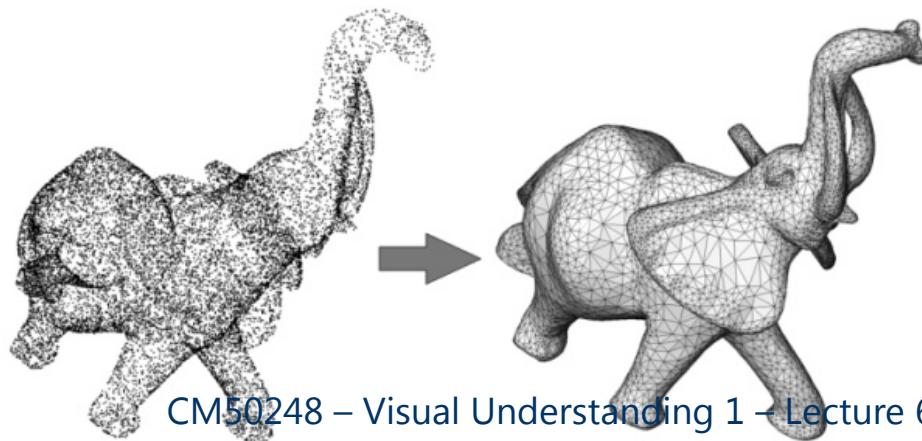


Real-world applications

- Streetview Reconstruction and Recognition
 - <http://vision.princeton.edu/projects/2009/ICCV/>
 - <http://vision.princeton.edu/projects/2009/TOG/>
- Photo Tourism
 - <http://phototour.cs.washington.edu/>
- Microsoft Photosynth
 - <http://photosynth.net/>
- 2d3, boujou (Matchmovers) and movies
 - <http://www.2d3.com/>
 - <http://www.vicon.com/boujou/>
- Robotics: SLAM (Simultaneous Localization And Mapping)
 - <http://openslam.org/>

Steps

- Images → Points: Structure from Motion
- Points → More points: Multiple View Stereo
- Points → Meshes: Model Fitting
- + Meshes → Models: Texture Mapping
-
- = Images → Models: Image-based Modeling



Point cloud → 3D mesh model

- Surface reconstruction:

- Marching cube

- http://en.wikipedia.org/wiki/Marching_cubes

- Poisson surface reconstruction

- <http://research.microsoft.com/en-us/um/people/hoppe/proj/poissonrecon/>

- Model fitting:

- RANSAC & J-linkage

- <http://www.diegm.uniud.it/fusiello/demo/jlk/>

- InverseCSG

- <http://research.google.com/pubs/pub39963.html>

- GlobFit

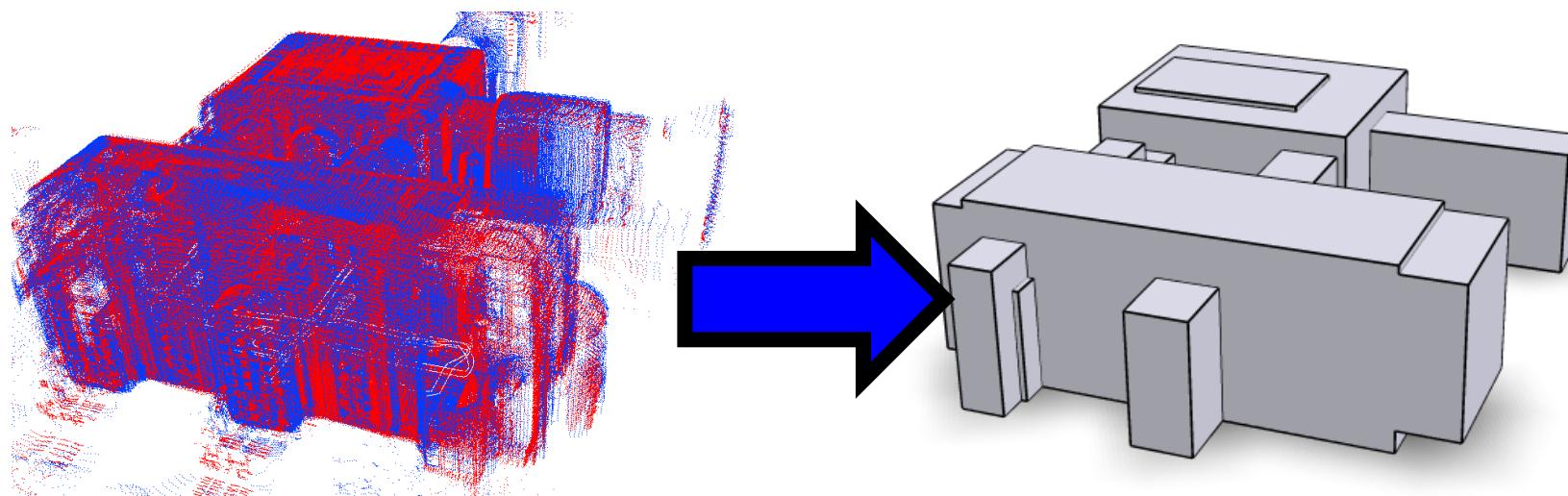
- <http://code.google.com/p/globfit/>

InverseCSG Algorithm

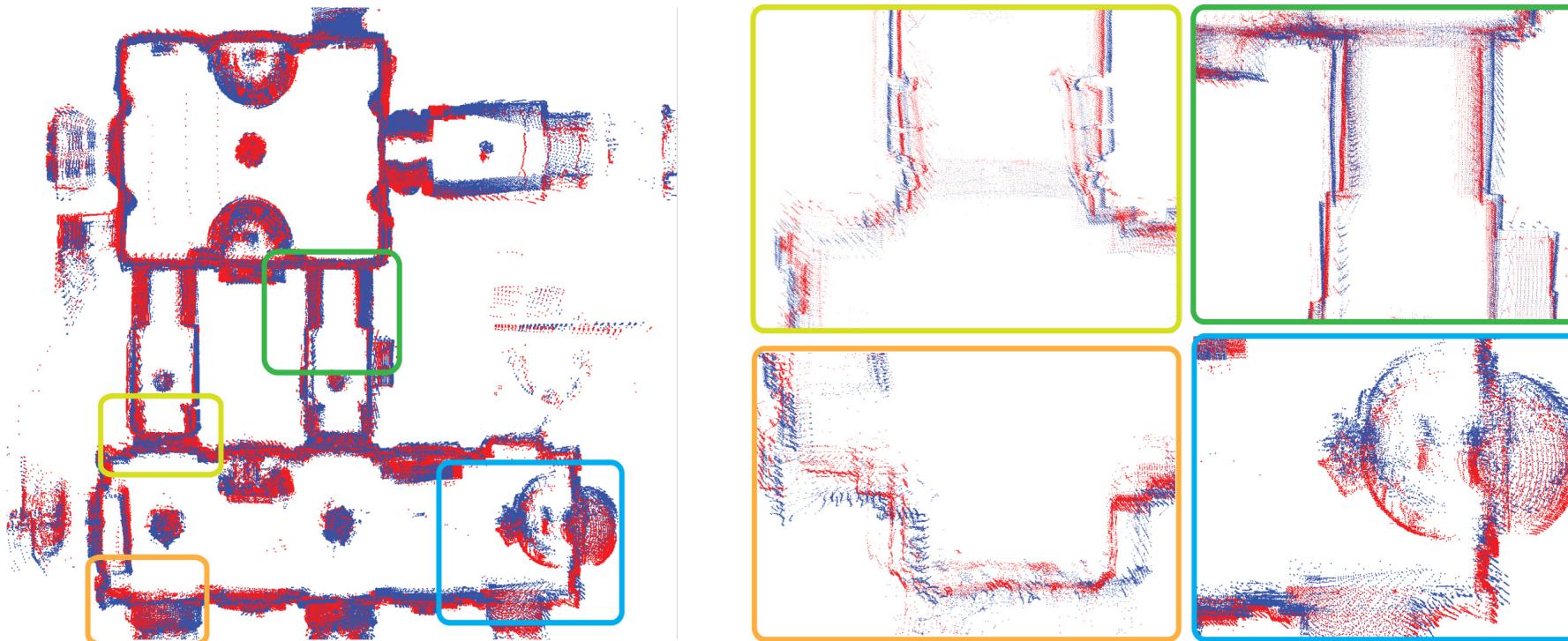


Reconstructing the World's Museums, J. Xiao and Y. Furukawa, 2012.

InverseCSG Algorithm



Noisy Points



top-down view of input points

Steps

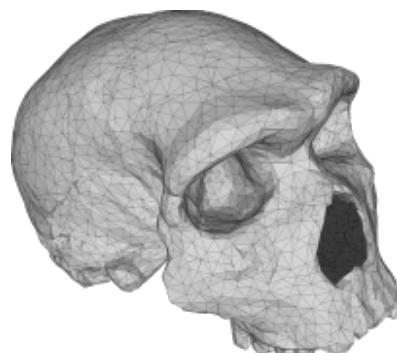
Images → Points: Structure from Motion

Points → More points: Multiple View Stereo

Points → Meshes: Model Fitting

+ Meshes → Models: Texture Mapping

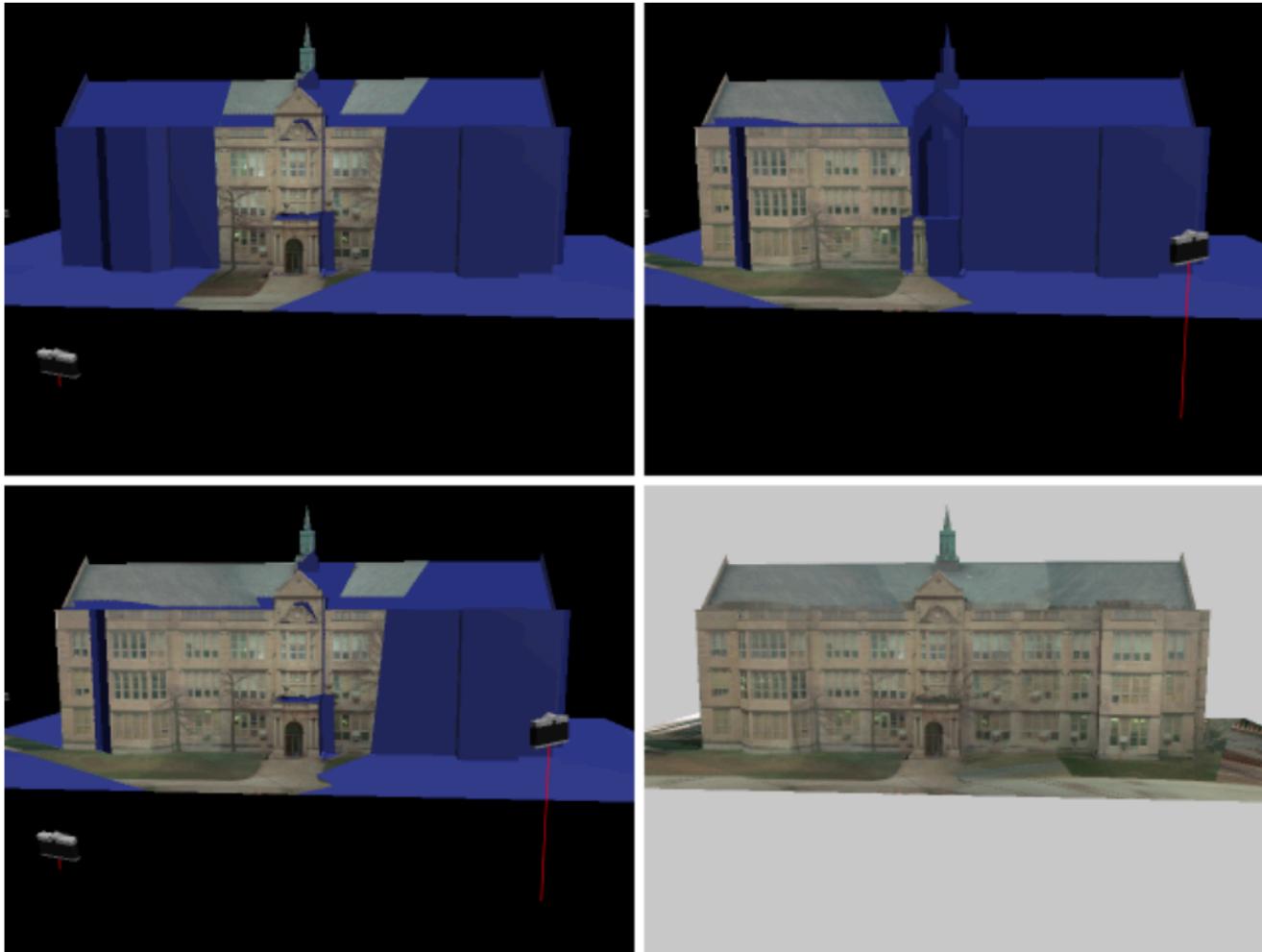
= Images → Models: Image-based Modeling



Texture mapping

- Texture stitching
- View-based texture mapping
- View interpolation and warping
- Interactive visualization

Texture stitching



The process of assembling projected images to form a composite rendering
Image from Paul Debevec's Thesis

View-Dependent Texture Mapping



(a)



(b)



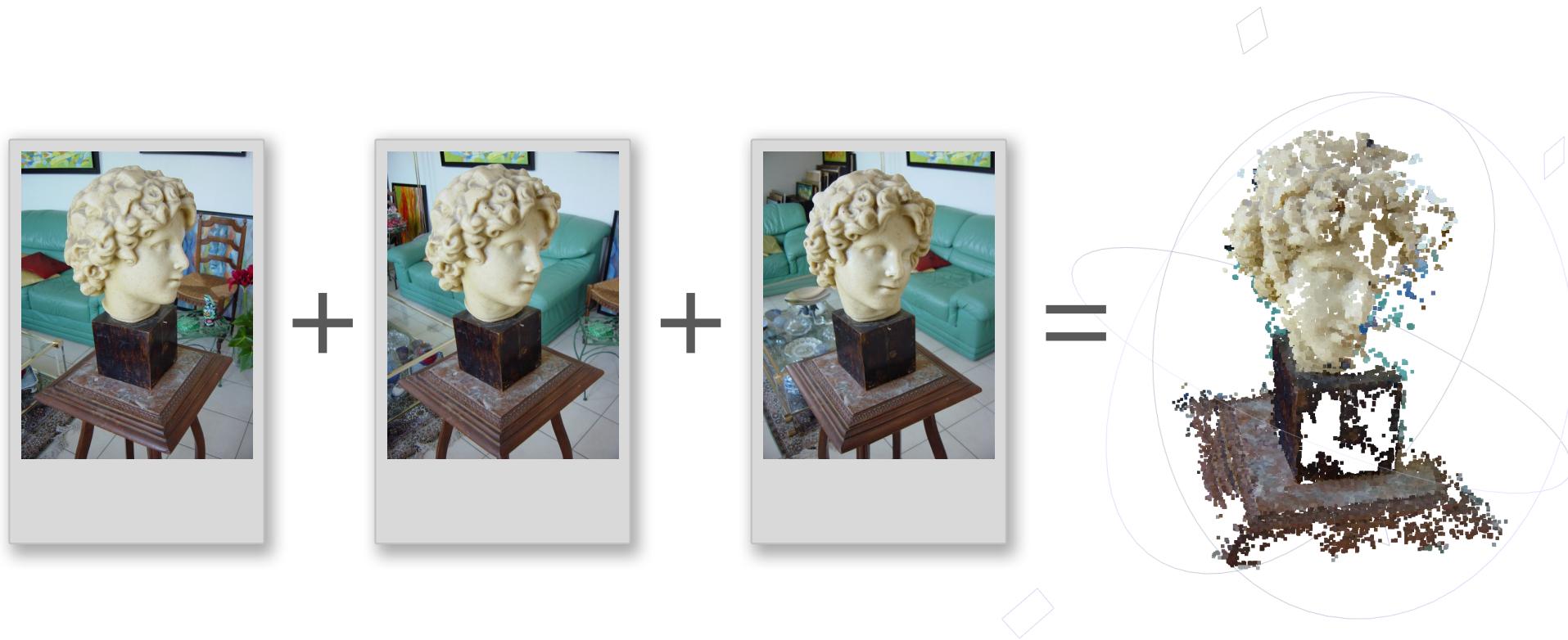
(c)



(d)

Modeling and Rendering Architecture from Photographs,
Paul Debevec, PhD Thesis, UC Berkeley, 1996

Q & A



Recognition

- How about something completely different?



Bag of words

■ Which is the vision text book?

Document	13747	the
A	12479	and
	7765	a
	7423	of
	...	
	2456	image
	...	
	1711	vision
	...	
	919	algorithms
	...	
	773	ieee
	...	
	699	recognition
	...	

Document	18062	the
B	8969	and
	7144	a
	6695	to
	...	
	1532	lord
	...	
	1311	ser
	...	
	800	jon
	...	
	745	ned
	...	
	590	men
	579	like
	572	tyrion

Bag of words

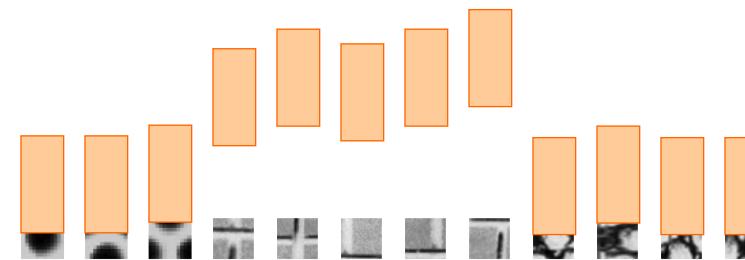
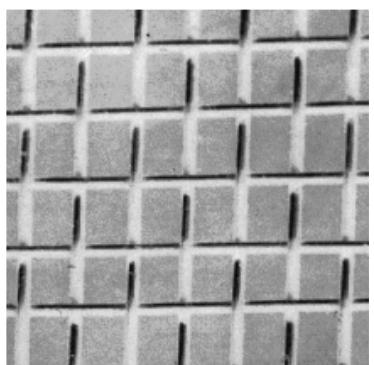
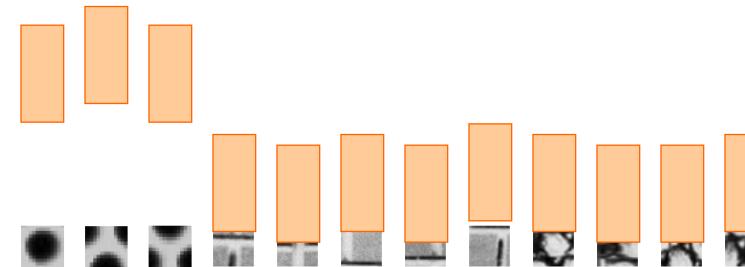
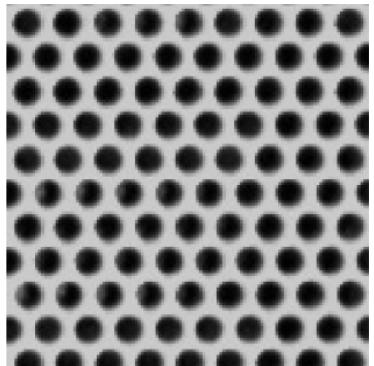
- Two-stage technique:
 - Learning/Training
 - Testing/Application
- Databases of classification test data (with ground truth)



Bag of words

- Create a shared dictionary (English language)
- Find histogram of word counts
 - Including words that aren't present
- Classify based on histogram
- "ser" → Game of Thrones
- "recognition" → Vision text-book

Bag of features

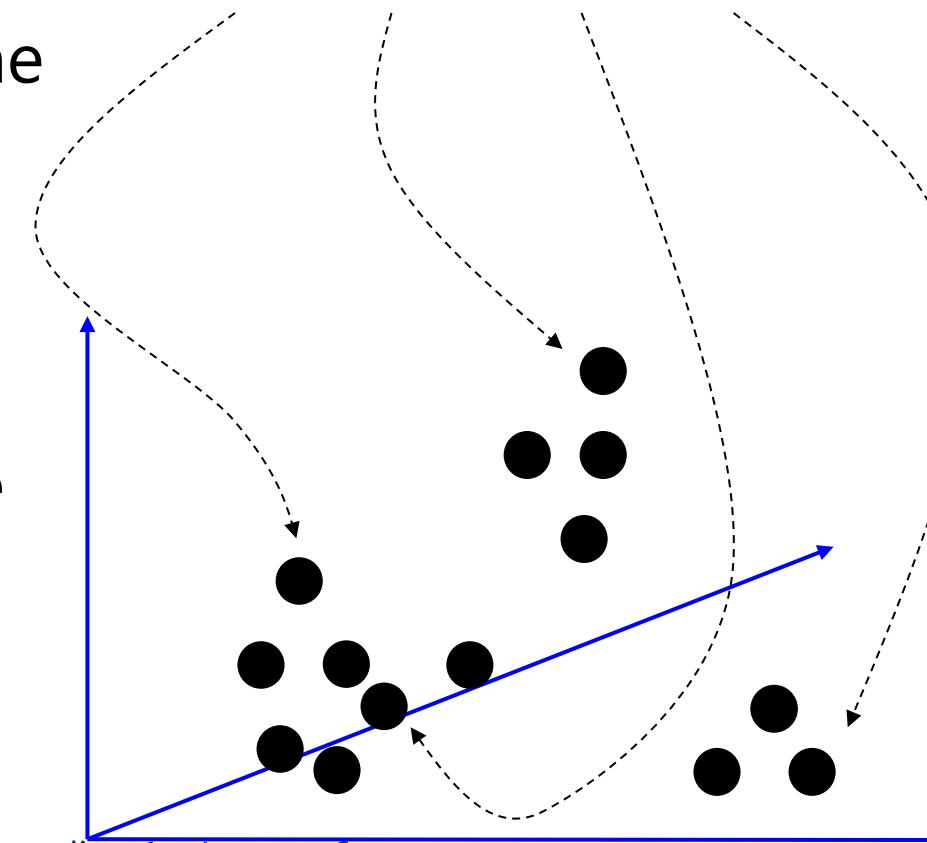
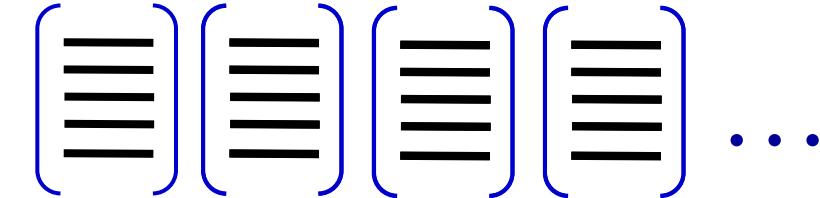


Bag of features

Imagine features in n-dimensional space, one dimension per element in the vector.

SIFT features have 128 elements.

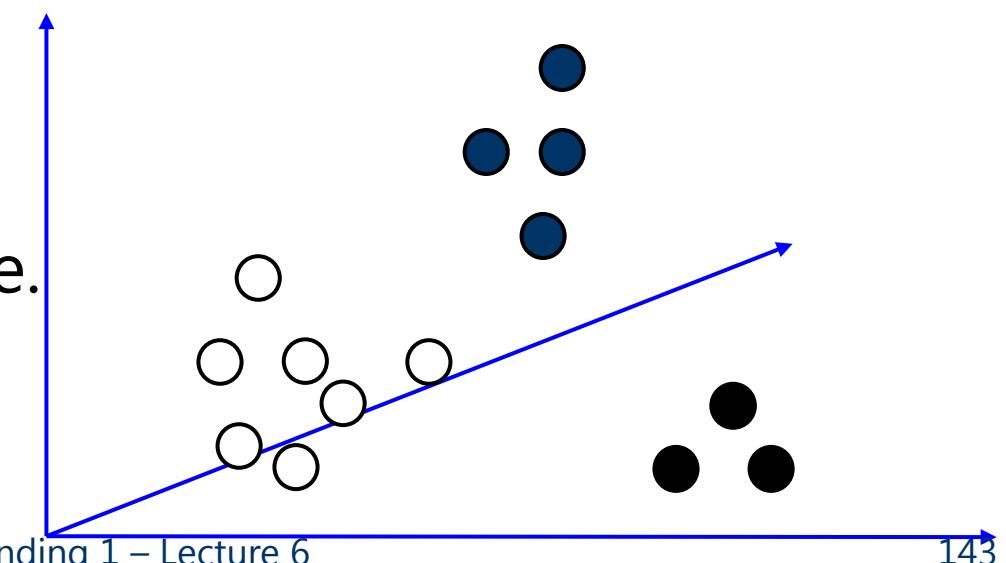
Higher dimensional representations are very useful in vision.



Bag of features

Cluster the points in n-space. This is the dictionary equivalent. Each cluster represents a word.

For each image, work out the number of instances of each word (features that belong to the given cluster). Use this to build the histogram describing the image.



Bag of features

- Now your image is represented by a histogram of words, another vector, probably 1000+ elements.
- Train a Support Vector Machine to classify.

Questions?

Well done!

You've made it to the end!