

## Computer Vision



## The goal of computer vision



La Gare Montparnasse, 1895

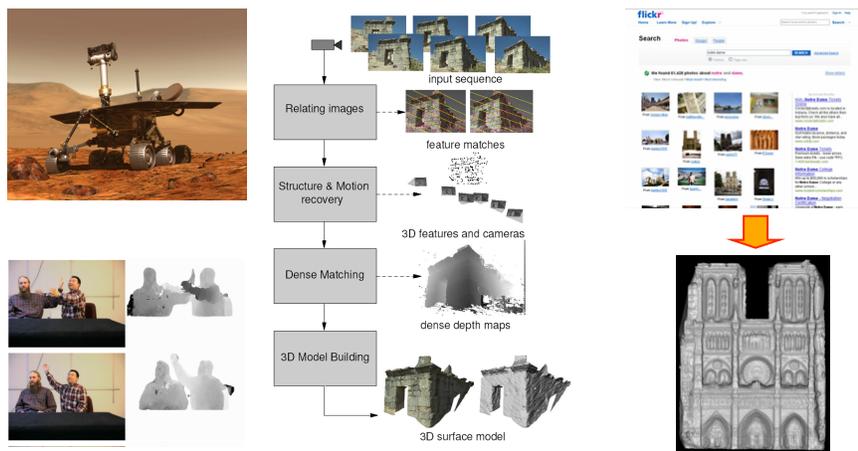
betw

0	3	2	5	4	7	6	9	8
3	0	1	2	3	4	5	6	7
2	1	0	3	2	5	4	7	6
5	2	3	0	1	2	3	4	5
4	3	2	1	0	3	2	5	4
7	4	5	2	3	0	1	2	3
6	5	4	3	2	1	0	3	2
9	6	7	4	5	2	3	0	1
8	7	6	5	4	3	2	1	0

## What kind of information can we extract from an image?

- Metric 3D information
- Semantic information

## Vision as measurement device



## Vision as a source of semantic information



## Object categorization



## Scene and context categorization

- outdoor
- city
- traffic
- ...

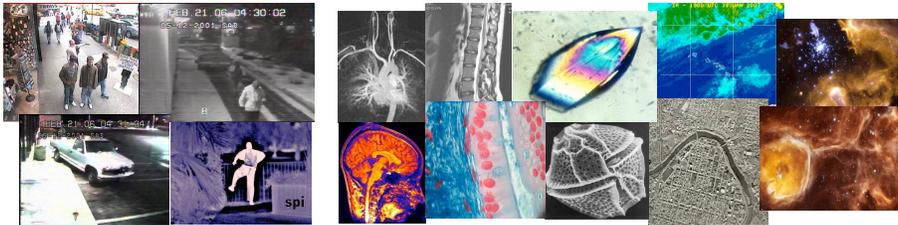


## Qualitative spatial information



## Why study computer vision?

- Vision is useful: Images and video are everywhere!



## Why study computer vision?

- Vision is useful
- Vision is interesting
- Vision is difficult
  - Half of primate cerebral cortex is devoted to visual processing

### Why is computer vision difficult?

### Challenges of recognizing objects in images



## Challenges: viewpoint variation



Michelangelo 1475-1564

slide credit: Fei-Fei, Feras

## Challenges: illumination

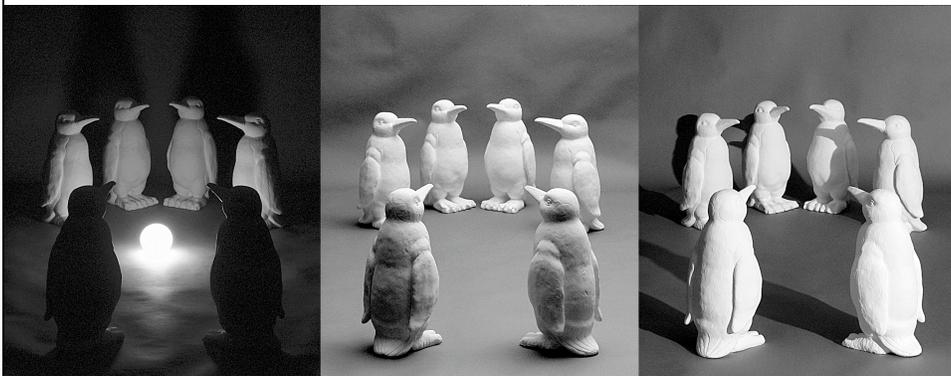


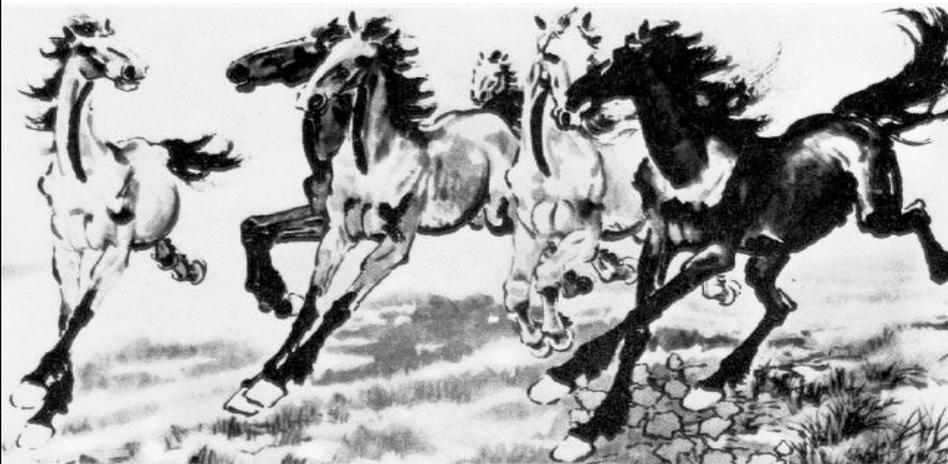
image credit: J. Koender

Challenges: scale



slide credit: Fei-Fei, Feraus

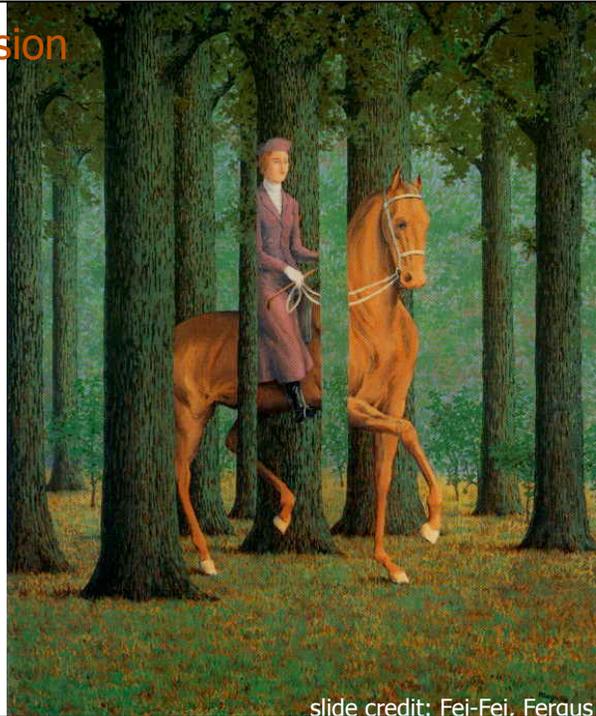
Challenges: deformation



Xu, Beihong 1943

slide credit: Fei-Fei, Feraus

Challenges: occlusion



Magritte, 1957

slide credit: Fei-Fei, Feras

Challenges: background clutter



Emperor shrimp and commensal crab on a sea cucumber in Fiji  
Photograph by Tom Larkin

NATIONAL GEOGRAPHIC

© 2007 National Geographic Society. All rights reserved.

## Challenges: Motion

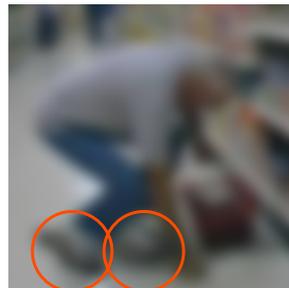
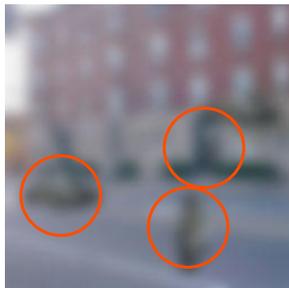
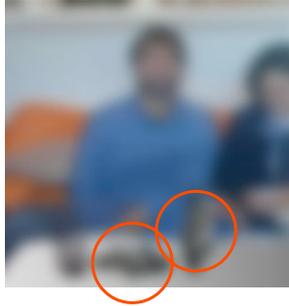


## Challenges: object intra-class variation



slide credit: Fei-Fei, Feras

## Challenges: local ambiguity



slide credit: Fei-Fei, Fergus

## Challenges or opportunities?

- Images are confusing, but they also reveal the structure of the world through numerous cues
- Our job is to interpret the cues!



### Depth cues: Linear perspective



NATIONALGEOGRAPHIC.COM

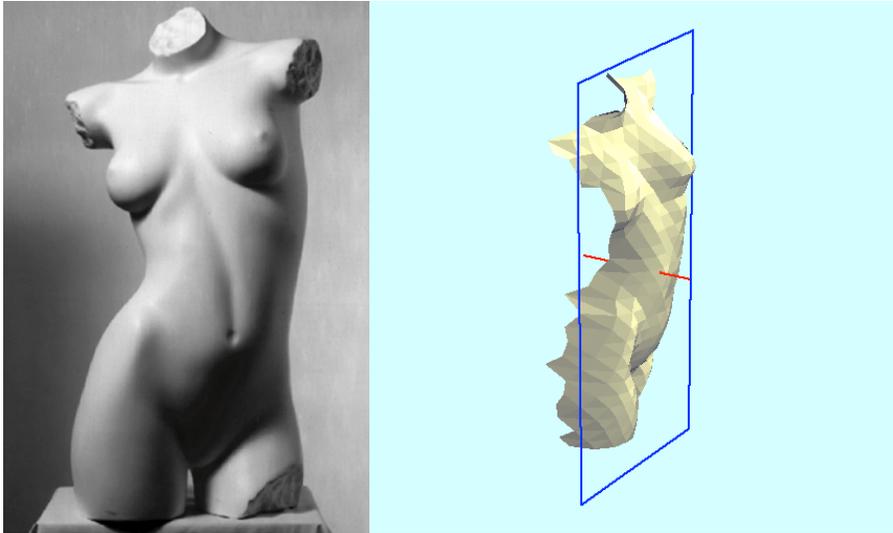
© 2003 National Geographic Society. All rights reserved.

### Depth cues: Aerial perspective



© 2002 National Geographic Society. All rights reserved. NATIONALGEOGRAPHIC.COM

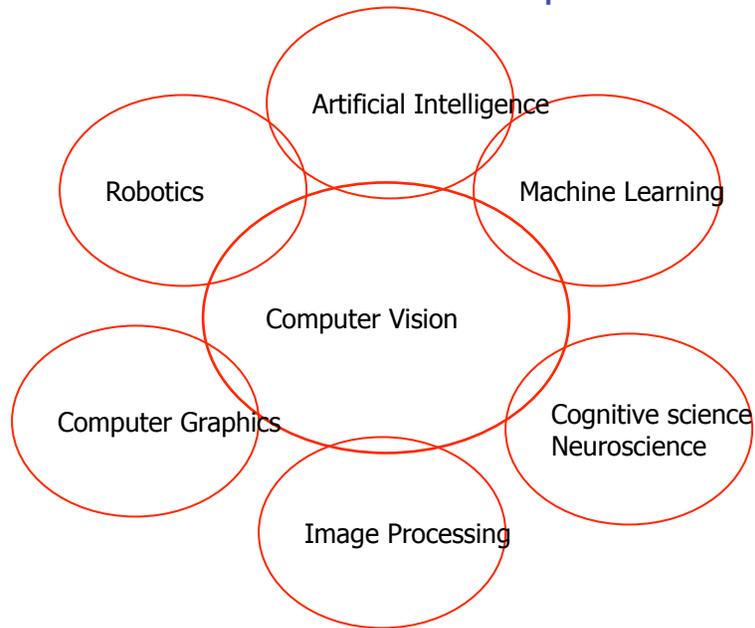
### Shape and lighting cues: Shading



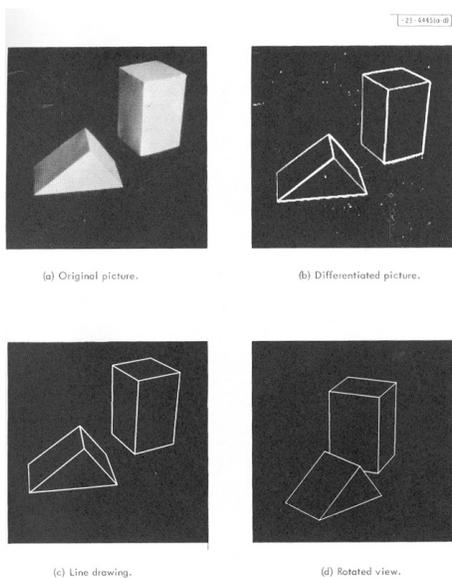
### Grouping cues: "Common fate"



## Connections to other disciplines

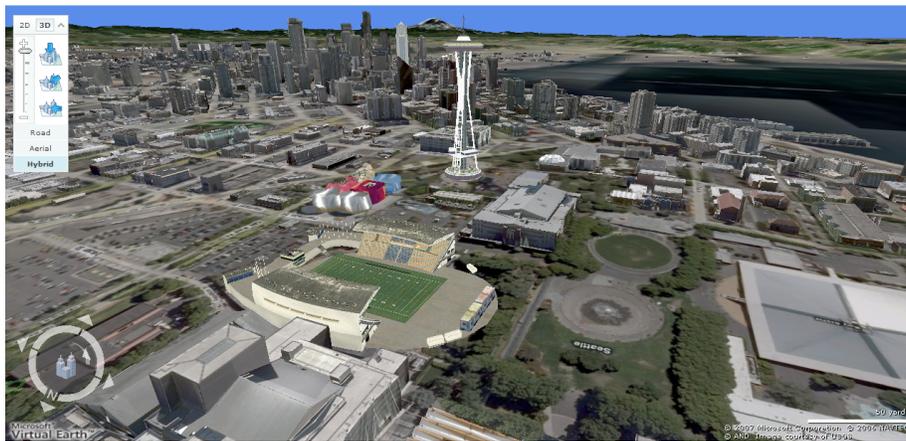


## Origins of computer vision



## Computer Vision in the Real World

### 3D urban modeling



[Bing maps](#), Google Streetview

Source: S. Se

## 3D urban modeling: Microsoft Photosynth



<http://labs.live.com/photosynth/>

Source: S. Se

## Face detection



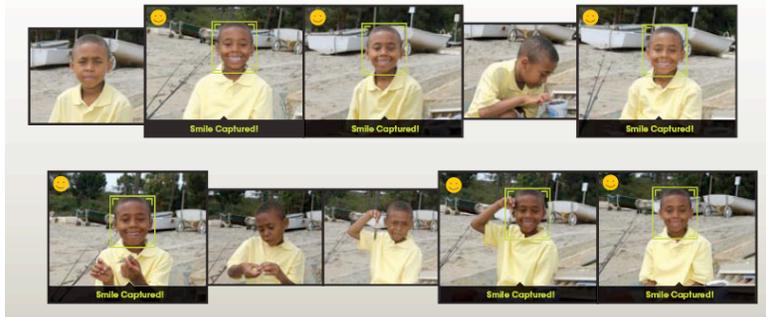
- Many new digital cameras now detect faces
  - Canon, Sony, Fuji, ...

Source: S. Se

## Smile detection

### The Smile Shutter flow

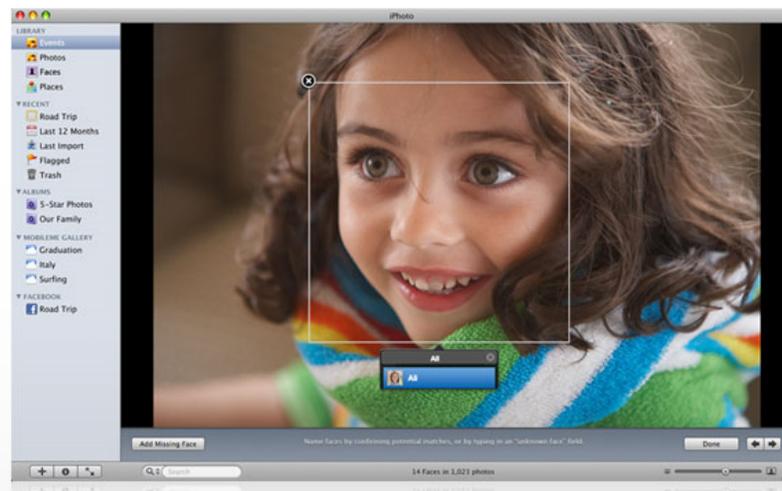
Imagine a camera smart enough to catch every smile! In Smile Shutter Mode, your Cyber-shot® camera can automatically trip the shutter at just the right instant to catch the perfect expression.



[Sony Cyber-shot® T70 Digital Still Camera](#)

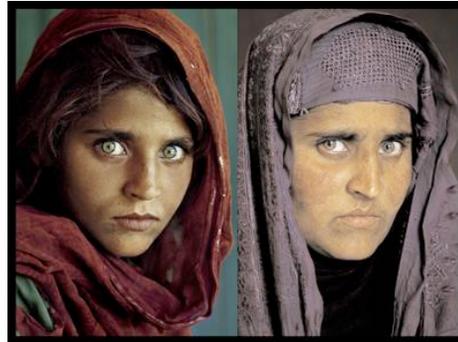
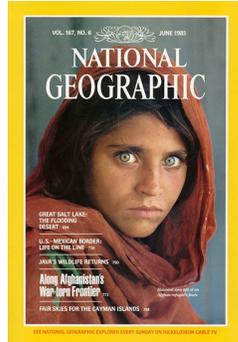
Source: S. Se

## Face recognition: Apple iPhoto software

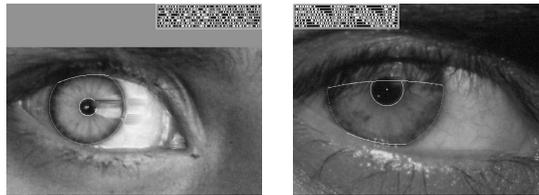


<http://www.apple.com/ilife/iphoto/>

## Biometrics



### How the Afghan Girl was Identified by Her Iris Patterns

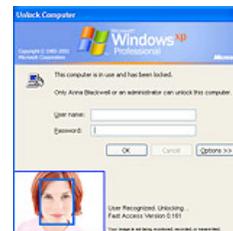
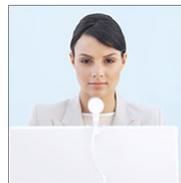


Source: S. Seitz

## Biometrics



Fingerprint scanners on many new laptops, other devices



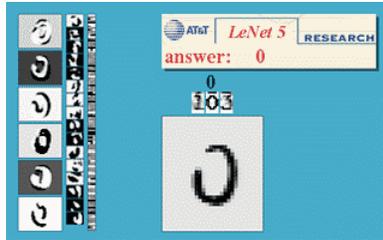
Face recognition systems now beginning to appear more widely <http://www.sensiblevision.com/>

Source: S. Seitz

## Optical character recognition (OCR)

Technology to convert scanned docs to text

- If you have a scanner, it probably came with OCR software



Digit recognition, AT&T labs

License plate readers

[http://en.wikipedia.org/wiki/Automatic\\_number\\_plate\\_recognition](http://en.wikipedia.org/wiki/Automatic_number_plate_recognition)

Source: S. Seitz

## Mobile visual search: Google Goggles

Google Goggles in Action

Click the icons below to see the different ways Google Goggles can be used.



## Mobile visual search: iPhone Apps





Query Images



Perspective



Zoom



Rotation



Occlusion



Lighting



Blur



Zoom

Matched Image



## Automotive safety

▶ manufacturer products
consumer products ◀

### Our Vision. Your Safety.



▶ **EyeQ** Vision on a chip



[read more](#)

▶ **Vision Applications**

Road, Vehicle, Pedestrian Protection and more



[read more](#)

▶ **AWS** Advance Warning System



[read more](#)

**News**

- ▶ Mobileye Advanced Technologies Power Volvo Cars World First Collision Warning With Auto Brake System
- ▶ Volvo: New Collision Warning with Auto Brake Helps Prevent Rear-end

[all news](#)

**Events**

- ▶ Mobileye at Equip Auto, Paris, France
- ▶ Mobileye at SEMA, Las Vegas, NV

[read more](#)

- **Mobileye:** Vision systems in high-end BMW, GM, Volvo models
  - "In mid 2010 Mobileye will launch a world's first application of full emergency braking for collision mitigation for pedestrians where vision is the key technology for detecting pedestrians."

Source: A. Shashua, S. Seitz

## Vision in supermarkets



### LaneHawk by EvolutionRobotics

“A smart camera is flush-mounted in the checkout lane, continuously watching for items. When an item is detected and recognized, the cashier verifies the quantity of items that were found under the basket, and continues to close the transaction. The item can remain under the basket, and with LaneHawk, you are assured to get paid for it...”

Source: S. Seitz

## Vision-based interaction (and games)



Nintendo Wii has camera-based IR tracking built in. See [Lee's work at CMU](#) on clever tricks on using it to create a [multi-touch display](#)!



Sony EyeToy



Assistive technologies

Source: S. Seitz

## Vision for robotics, space exploration



[NASA'S Mars Exploration Rover Spirit](#) captured this westward view from atop a low plateau where Spirit spent the closing months of 2007.

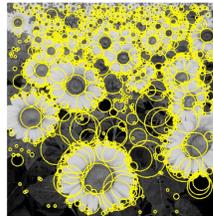
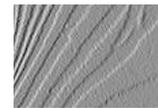
### Vision systems (JPL) used for several tasks

- Panorama stitching
- 3D terrain modeling
- Obstacle detection, position tracking
- For more, read "[Computer Vision on Mars](#)" by Matthies et al.

Source: S. Seitz

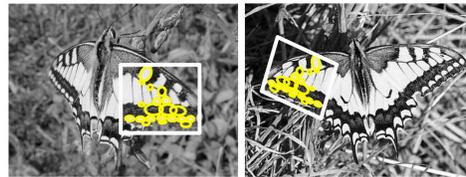
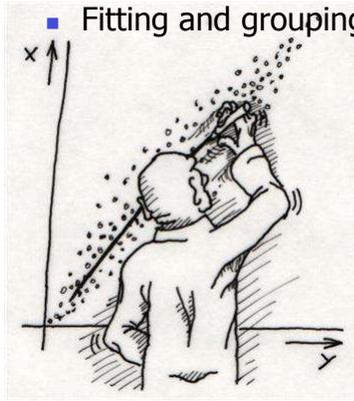
## I. Early vision

- Basic image formation and processing

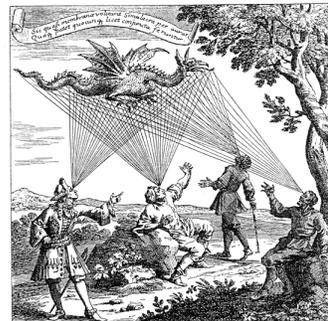
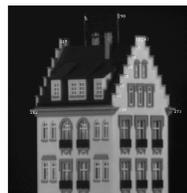
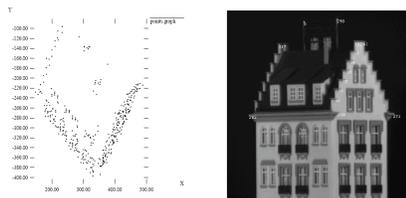
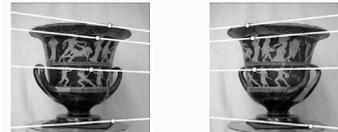


## II. "Mid-level vision"

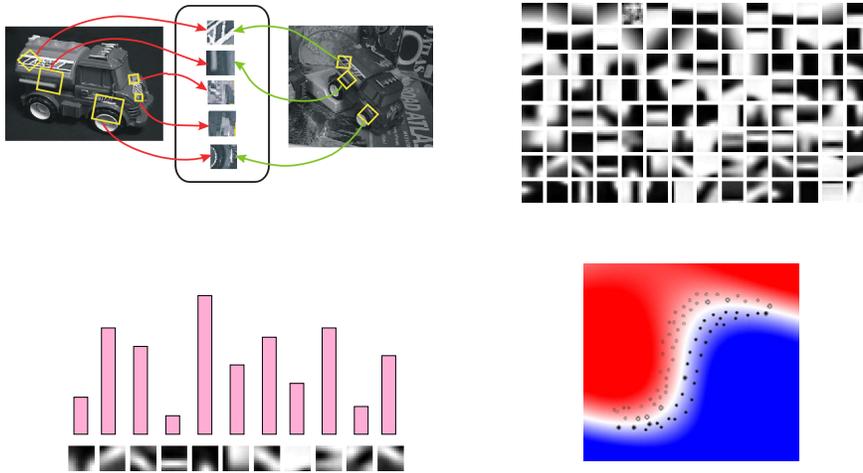
### ■ Fitting and grouping



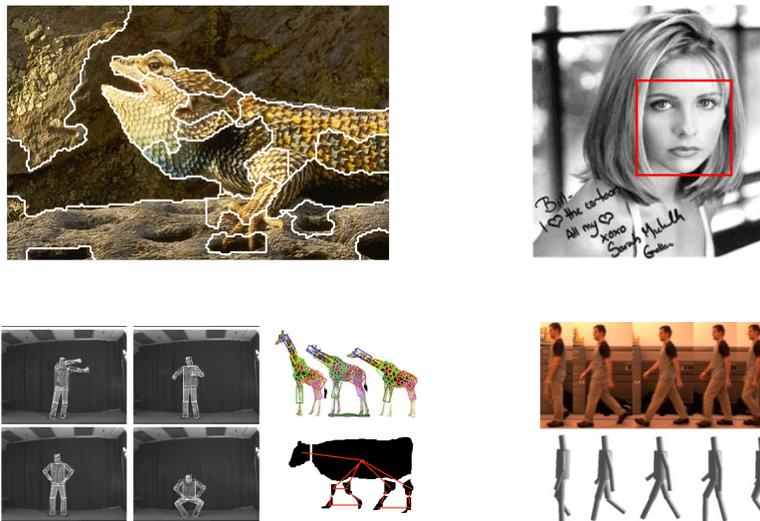
## III. Multi-view geometry



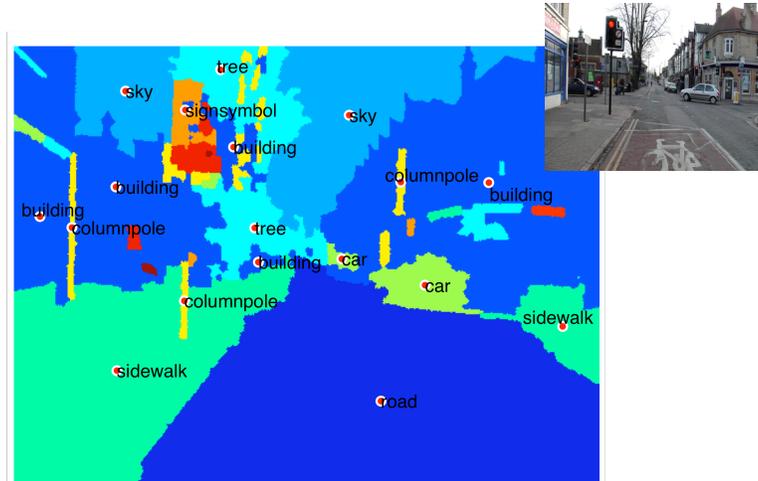
### IV. Recognition



### V. Advanced Topics



## Semantic Labeling of Street Scenes



- Using image segmentation and scene and object recognition for automated image labeling

## Role of Perception in Robotics

- Where am I relative to the world?
  - sensors: vision, stereo, range sensors, acoustics
  - problems: scene modeling/classification/recognition
  - integration: localization/mapping algorithms (e.g. SLAM)
- What is around me?
  - sensors: vision, stereo, range sensors, acoustics, sounds, smell
  - problems: object recognition, structure from x, qualitative modeling
  - integration: collision avoidance/navigation, learning

Jana Kosecka



## Visual Perception Topics

### Techniques

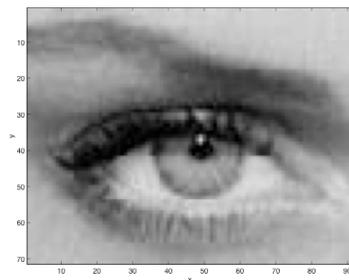
- Computational Stereo
- Feature detection and matching
- Motion tracking and visual feedback

### Applications in Robotics:

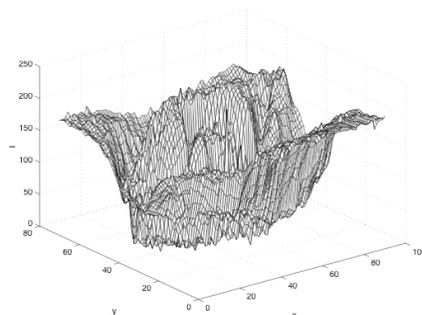
- range sensing, Obstacle detection, environment interaction
- Mapping, registration, localization, recognition
- Manipulation

## Image - Apperance

Image



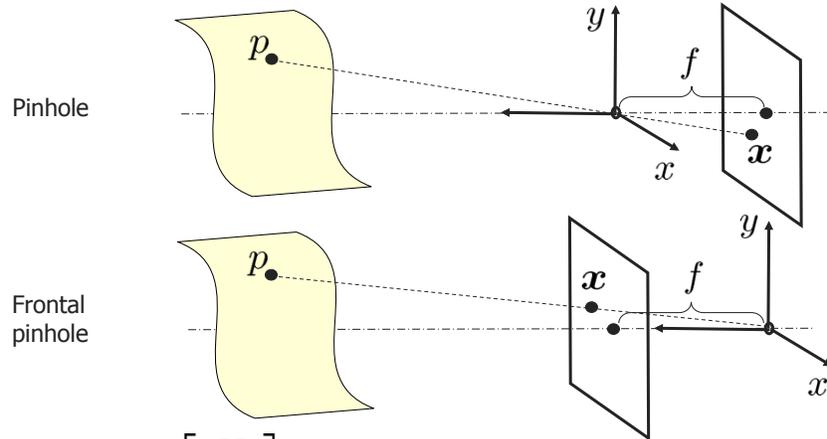
Brightness values



$I(x,y)$

Jana Kosecka

## Image Formation



$$\mathbf{X} = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} \rightarrow \mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} = \frac{f}{Z} \begin{bmatrix} X \\ Y \end{bmatrix}$$

J. Kosecka, GMU

## Pinhole Camera Model

- Image coordinates are nonlinear function of world coordinates
- Relationship between coordinates in the camera frame and sensor plane

2-D coordinates  $\mathbf{x} = \begin{bmatrix} x \\ y \end{bmatrix} = \frac{f}{Z} \begin{bmatrix} X \\ Y \end{bmatrix}$

Homogeneous coordinates

$$\mathbf{x} \rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \frac{1}{Z} \begin{bmatrix} fX \\ fY \\ Z \end{bmatrix}, \quad \mathbf{X} \rightarrow \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix},$$

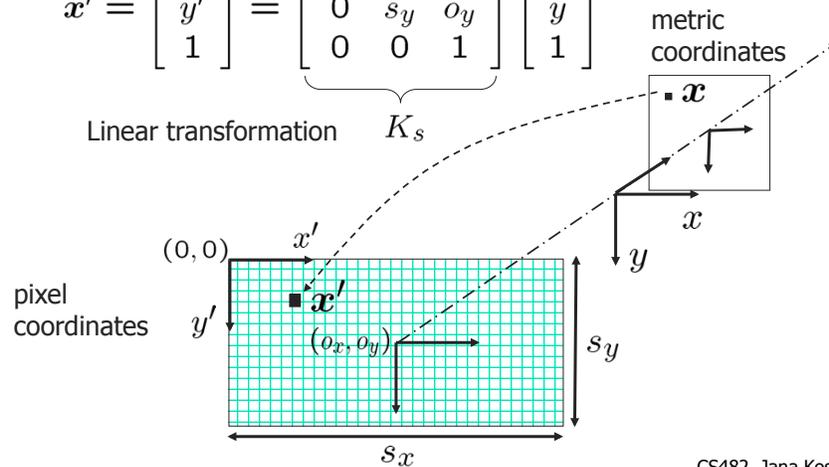
$$Z \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} f & 0 & 0 \\ 0 & f & 0 \\ 0 & 0 & 1 \end{bmatrix}}_{K_f} \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{\Pi_0} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

## Image Coordinates

- Relationship between coordinates in the sensor plane and image

$$\mathbf{x}' = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} s_x & s_\theta & o_x \\ 0 & s_y & o_y \\ 0 & 0 & 1 \end{bmatrix}}_{K_s} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Linear transformation



CS482, Jana Kosecka

## Calibration Matrix and Camera Model

- Relationship between coordinates in the world frame and image
- Intrinsic parameters

Pinhole camera

Pixel coordinates

$$\lambda \mathbf{x} = K_f \Pi_0 \mathbf{X} \quad \mathbf{x}' = K_s \mathbf{x}$$

- Adding transformation between camera coordinate systems and world coordinate system
- Extrinsic Parameters

$$\lambda \mathbf{x}' = \begin{bmatrix} f s_x & f s_\theta & o_x \\ 0 & f s_y & o_y \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} R & T \\ 0 & 1 \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

$$\lambda \mathbf{x} = K_f \Pi_0 g \mathbf{X} = \Pi \mathbf{X}$$

## Image of a Point

Homogeneous coordinates of a 3-D point  $p$

$$\mathbf{X} = [X, Y, Z, W]^T \in \mathbb{R}^4, \quad (W = 1)$$

Homogeneous coordinates of its 2-D image

$$\mathbf{x} = [x, y, z]^T \in \mathbb{R}^3, \quad (z = 1)$$

Projection of a 3-D point to an image plane

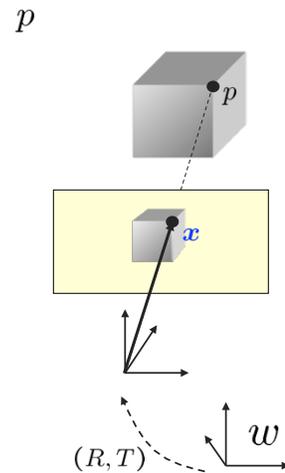
$$\lambda \mathbf{x} = \Pi \mathbf{X}$$

$$\lambda \in \mathbb{R}, \quad \Pi = [R, T] \in \mathbb{R}^{3 \times 4}$$

$$\lambda \mathbf{x}' = \Pi \mathbf{X}$$

$$\lambda \in \mathbb{R}, \quad \Pi = [KR, KT] \in \mathbb{R}^{3 \times 4}$$

Jana Kosecka, CS 685



58

## Image of a Line

Homogeneous representation of a 3-D line  $L$

$$\mathbf{X} = \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix} = \begin{bmatrix} X_o \\ Y_o \\ Z_o \\ 1 \end{bmatrix} + \mu \begin{bmatrix} V_1 \\ V_2 \\ V_3 \\ 0 \end{bmatrix}, \quad \mu \in \mathbb{R}$$

Homogeneous representation of its 2-D image

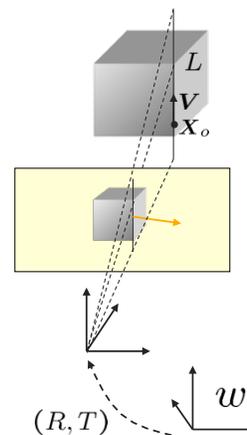
$$\mathbf{l} = [a, b, c]^T \in \mathbb{R}^3$$

Projection of a 3-D line to an image plane

$$\mathbf{l}^T \mathbf{x} = \mathbf{l}^T \Pi \mathbf{X} = 0$$

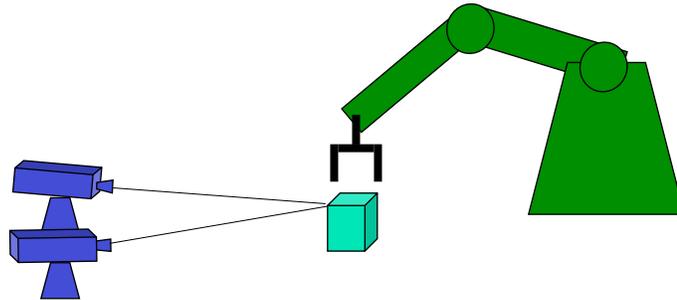
$$\Pi = [KR, KT] \in \mathbb{R}^{3 \times 4}$$

Jana Kosecka, CS 685



59

## What is Computational Stereo?



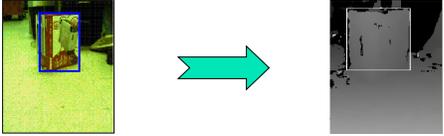
Viewing the same physical point from two different viewpoints allows depth from triangulation

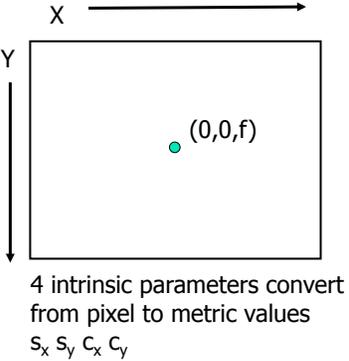
## Computational Stereo

- Much of geometric vision is based on information from 2 (or more) camera locations
- Hard to recover 3D information from a single 2D image without extra knowledge
- Motion and stereo (multiple cameras) are both common in the world
- Stereo vision is ubiquitous in nature (oddly, nearly 10% of people are stereo blind)
- Stereo involves the following *three problems*:
  1. calibration
  2. matching (*correspondence problem*)
  3. reconstruction (*reconstruction problem*)

## Binocular Stereo System: Geometry

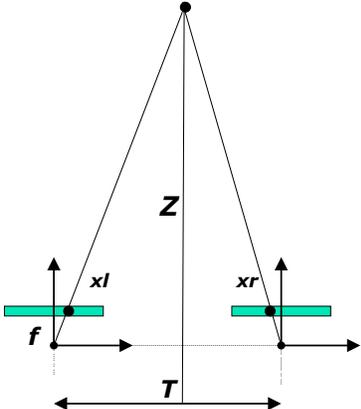
- **GOAL:** Passive 2-camera system using triangulation to generate a depth map of a world scene.
- **Depth map:**  $z=f(x,y)$  where  $x,y$  are coordinates one of the image planes and  $z$  is the height above the respective image plane.
  - Note that for stereo systems which differ only by an offset in  $x$ , the  $v$  coordinates (projection of  $y$ ) is the same in both images!
  - Note we must convert from image (pixel) coordinates to external coordinates -- **requires calibration**





## Stereo Configuration

- Images are scan-aligned
- Disparity between two images – inversely proportional to depth
- Disparity – difference between  $x$ -coordinates of a feature
- Triangle similarity



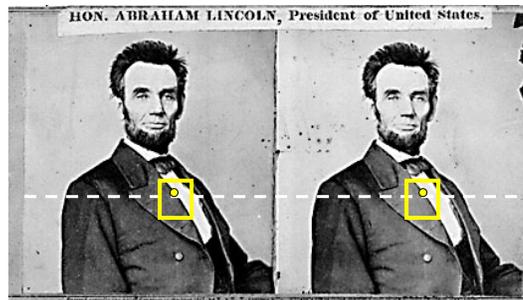
$$\frac{Z}{T} = \frac{Z-f}{T-x_l-x_r}$$

$$Z = \frac{fT}{\text{disparity}}$$

## Stereo Vision

- Distance is inversely proportional to disparity
  - closer objects can be measured more accurately
- Disparity is proportional to baseline
  - For a given disparity error, the accuracy of the depth estimate increases with increasing baseline
  - However, as baseline is increased, some objects may appear in one camera, but not in the other.
  - Image resolution is also a factor

## Stereo Matching – Stereo Correspondence



For each epipolar line (scanline)

For each pixel in the left image

- compare with every pixel on same epipolar line in right image
- pick pixel with minimum match cost
- This will never work, so:
- Match Windows

## Region based Similarity Metric

- Sum of squared differences

$$SSD(h) = \sum_{\tilde{x} \in W(x)} \|I_1(\tilde{x}) - I_2(h(\tilde{x}))\|^2$$

- Normalize cross-correlation

$$NCC(h) = \frac{\sum_{W(x)} (I_1(\tilde{x}) - \bar{I}_1)(I_2(h(\tilde{x})) - \bar{I}_2)}{\sqrt{\sum_{W(x)} (I_1(\tilde{x}) - \bar{I}_1)^2 \sum_{W(x)} (I_2(h(\tilde{x})) - \bar{I}_2)^2}}$$

- Sum of absolute differences

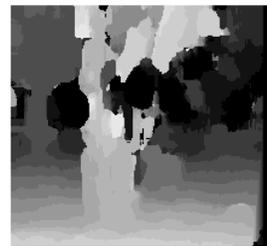
$$SAD(h) = \sum_{\tilde{x} \in W(x)} |I_1(\tilde{x}) - I_2(h(\tilde{x}))|$$

66

## Window size



W = 3



W = 20

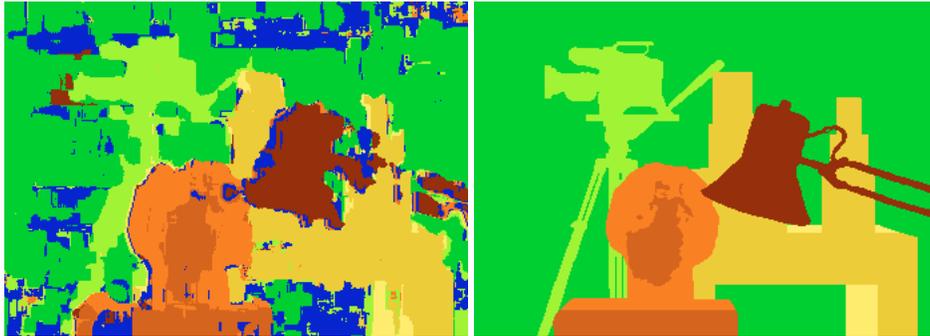
- Effect of window size

### With adaptive window

- T. Kanade and M. Okutomi, [A Stereo Matching Algorithm with an Adaptive Window: Theory and Experiment](#), Proc. International Conference on Robotics and Automation, 1991.
- D. Scharstein and R. Szeliski, [Stereo matching with nonlinear diffusion](#), International Journal of Computer Vision, 28(2): 155-174, July 1998

(S. Seitz) Jana Kosecka

## Results with window correlation



Window-based matching  
(best window size)

Ground truth

(slide courtesy S. Seitz)

Jana Kosecka

## Results with better method



State of the art method

Ground truth

Boykov et al., [Fast Approximate Energy Minimization via Graph Cuts](#),  
International Conference on Computer Vision, September 1999.

(slide courtesy S. Seitz)

Jana Kosecka

## Applications of Real-Time Stereo

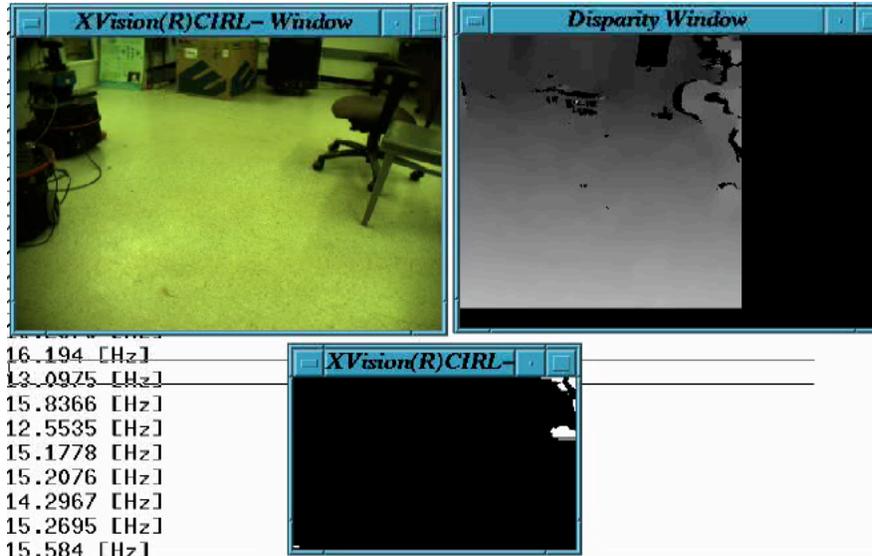
- Mobile robotics
  - Detect the structure of ground; detect obstacles; conveying
- Graphics/video
  - Detect foreground objects and matte in other objects (super-matrix effect)
- Surveillance
  - Detect and classify vehicles on a street or in a parking garage
- Medical
  - Measurement (e.g. sizing tumors)
  - Visualization (e.g. register with pre-operative CT)

## Obstacle Detection (cont'd)

Observation: Removing the ground plane immediately exposes obstacles



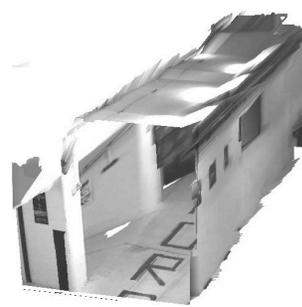
## Applications of Real-Time Stereo



## Oxford corridor



using 6 images



3D model

## Feature based stereo

- Instead of matching each pixel
- Match features in the image
- What are good features ? – next lecture
- Examples of features – line matching, point matching, region matching

## Uncalibrated Camera

$$\mathbf{x}' = \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = K \mathbf{x} = \underbrace{\begin{bmatrix} fs_x & fs_\theta & o_x \\ 0 & fs_y & o_y \\ 0 & 0 & 1 \end{bmatrix}}_K \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Linear transformation  $K$

pixel coordinates  $(0,0)$   $x'$   $y'$   $(o_x, o_y)$   $s_x$   $s_y$

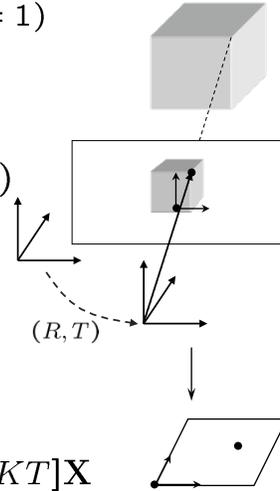
calibrated coordinates  $x$   $y$

## Uncalibrated Camera

$$\mathbf{X} = [X, Y, Z, W]^T \in \mathbb{R}^4, \quad (W = 1)$$

### Calibrated camera

- Image plane coordinates  $\mathbf{x} = [x, y, 1]^T$
- Camera extrinsic parameters  $g = (R, T)$
- Perspective projection  $\lambda \mathbf{x} = [R, T] \mathbf{X}$



### Uncalibrated camera

- Pixel coordinates  $\mathbf{x}' = K \mathbf{x}$
- Projection matrix  $\lambda \mathbf{x}' = \Pi \mathbf{X} = [KR, KT] \mathbf{X}$

Jana Kosecka, CS 685

76

## Calibration with a Rig

Use the fact that both 3-D and 2-D coordinates of feature points on a pre-fabricated object (e.g., a cube) are known.



Jana Kosecka, CS 685

77

## Calibration with a Rig

- Given 3-D coordinates on known object

$$\lambda \mathbf{x}' = [KR, KT]\mathbf{X} \longrightarrow \lambda \mathbf{x}' = \Pi \mathbf{X}$$

$$\lambda \begin{bmatrix} x^i \\ y^i \\ 1 \end{bmatrix} = \begin{bmatrix} \pi_1^T \\ \pi_2^T \\ \pi_3^T \end{bmatrix} \begin{bmatrix} X^i \\ Y^i \\ Z^i \\ 1 \end{bmatrix}$$

- Eliminate unknown scales

$$x^i (\pi_3^T \mathbf{X}) = \pi_1^T \mathbf{X},$$

$$y^i (\pi_3^T \mathbf{X}) = \pi_2^T \mathbf{X}$$

- Recover projection matrix  $\Pi = [KR, KT] = [R', T']$

$$\min \|\Pi^s\|^2 \quad \text{subject to} \quad \|\Pi^s\|^2 = 1$$

$$\Pi^s = [\pi_{11}, \pi_{21}, \pi_{31}, \pi_{12}, \pi_{22}, \pi_{32}, \pi_{13}, \pi_{23}, \pi_{33}, \pi_{14}, \pi_{24}, \pi_{34}]^T$$

- Factor the into  $R \in SO(3)$  and using QR decomposition

- Solve for translation  $T = K^{-1}T'$

78

## Alternative camera models/projections

Orthographic projection

$$\mathbf{x}' = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

Scaled orthographic projection

$$\mathbf{x}' = s \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

Affine camera model

$$\mathbf{x}' = \begin{bmatrix} a_{11} & a_{12} & a_{13} & a_{14} \\ a_{21} & a_{22} & a_{23} & a_{24} \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

83

## Stereo

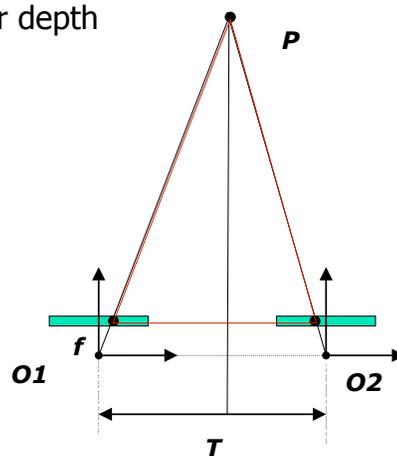
- What if the motion between cameras is not known ?

Jana Kosecka, CS 685

85

## Canonical Stereo Configuration

- Assumes (two) cameras
- Known positions and focal lengths
- Recover depth



$$\frac{Z}{T} = \frac{Z-f}{T-x_l-x_r}$$

$$Z = \frac{fT}{\text{disparity}}$$

86

## Rigid Body Motion – Two Views

$\mathbf{X} = [X, Y, Z, 1]^T$

$\mathbf{x} = [x, y, 1]^T$

$\lambda_1 \mathbf{x}_1 = \mathbf{X}$

$\lambda_2 \mathbf{x}_2 = R \lambda_1 \mathbf{x}_1 + T$

$\lambda_2 \mathbf{x}_2 = R \lambda_1 \mathbf{x}_1 + T$

Jana Kosecka, CS 685 87

## 3D Structure and Motion Recovery

Euclidean transformation

$$\lambda_2 \mathbf{x}_2 = R \lambda_1 \mathbf{x}_1 + T$$

measurements
unknowns

$$\sum_{j=1}^n \|\mathbf{x}_1^j - \pi(R_1, T_1, \mathbf{X})\|^2 + \|\mathbf{x}_2^j - \pi(R_2, T_2, \mathbf{X})\|^2$$

Find such **Rotation** and **Translation** and **Depth** that the reprojection error is minimized

Two views  $\sim$  200 points

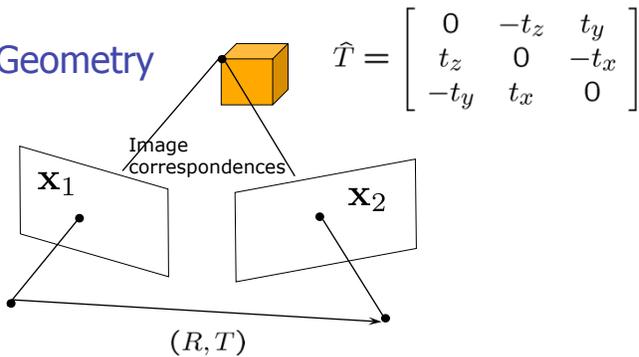
6 unknowns – **Motion** 3 Rotation, 3 Translation

- **Structure** 200x3 coordinates
- (-) universal scale

**Difficult optimization problem**

Jana Kosecka, CS 685 88

## Epipolar Geometry



$$\lambda_2 \mathbf{x}_2 = R \lambda_1 \mathbf{x}_1 + T \quad / \hat{\mathbf{x}}_2^T T$$

- Algebraic Elimination of Depth [Longuet-Higgins '81]:

$$\mathbf{x}_2^T \underbrace{\hat{T} R}_{E} \mathbf{x}_1 = 0$$

- Essential matrix  $E = \hat{T} R$

## Characterization of Essential Matrix

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

Essential matrix  $E = \hat{T} R$  special 3x3 matrix

$$\mathbf{x}_2^T \begin{bmatrix} e_1 & e_2 & e_3 \\ e_4 & e_5 & e_6 \\ e_7 & e_8 & e_9 \end{bmatrix} \mathbf{x}_1 = 0$$

### (Essential Matrix Characterization)

A non-zero matrix  $E$  is an essential matrix iff its SVD:  $E = U \Sigma V^T$  satisfies:  $\Sigma = \text{diag}([\sigma_1, \sigma_2, \sigma_3])$  with  $\sigma_1 = \sigma_2 \neq 0$  and  $\sigma_3 = 0$  and  $U, V \in SO(3)$



## Pose Recovery

- There are **two** pairs  $(R, T)$  corresponding to essential matrix  $E$ .
- There are **two** pairs  $(R, T)$  corresponding to essential matrix  $-E$ .
- Positive depth constraint disambiguates the impossible solutions
- Translation has to be non-zero, can be recovered up to scale
- Points have to be in general position
  - degenerate configurations – planar points
  - quadratic surface
- Linear 8-point algorithm
- Nonlinear 5-point algorithms yields up to 10 solutions

## 3D Structure Recovery

$$\lambda_2 \mathbf{x}_2 = R \lambda_1 \mathbf{x}_1 + \gamma T \quad \text{unknowns}$$

- Eliminate one of the scale's

$$\lambda_1^j \widehat{\mathbf{x}}_2^j R \mathbf{x}_1^j + \gamma \widehat{\mathbf{x}}_2^j T = 0, \quad j = 1, 2, \dots, n$$

- Solve LLSE problem

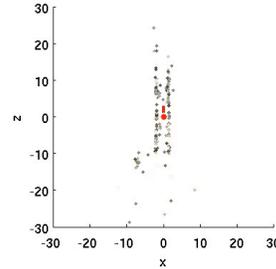
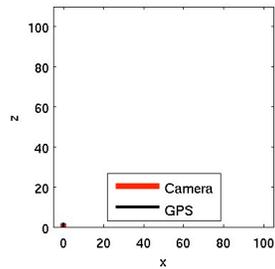
$$M^j \bar{\lambda}^j \doteq \begin{bmatrix} \widehat{\mathbf{x}}_2^j R \mathbf{x}_1^j, & \widehat{\mathbf{x}}_2^j T \end{bmatrix} \begin{bmatrix} \lambda_1^j \\ \gamma \end{bmatrix} = 0$$

If the configuration is non-critical, the Euclidean structure of the points and motion of the camera can be reconstructed up to a universal scale.

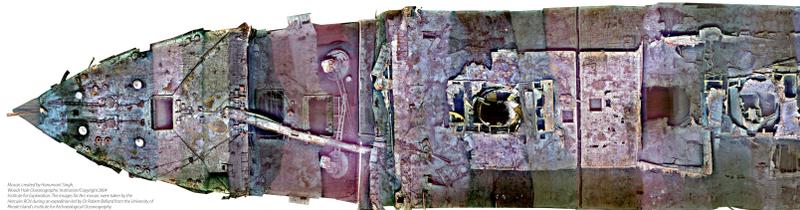
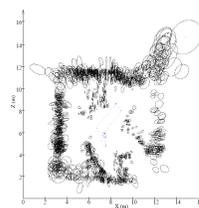
- Alternatively recover each point depth separately

# Visual Odometry

estimate motion from image correspondences



# Mapping, Localization, Recognition



Model created by Neil Alford. Digitally  
 rendered by the author. All rights reserved.  
 This is a Creative Commons Attribution-NonCommercial-ShareAlike  
 license. All images are property of the author. All rights reserved.  
 Please do not use this for any other purpose.

## Dealing with correspondences

- Previous methods assumed that we have exact correspondences
- Followed by linear least squares estimation
- Correspondences established either by tracking (using affine or translational flow models)
- Or wide-baseline matching (using scale/rotation invariant features and their descriptors)
- In many cases we get incorrect matches/tracks

102

## The RANSAC algorithm

- Generate  $M$  (a predetermined number) model hypotheses, each of them is computed using a minimal subset of points
- Evaluate each hypothesis
- Compute its residuals with respect to all data points.
- Points with residuals less than some threshold are classified as its inliers
- The hypothesis with the maximal number of inliers is chosen. Then re-estimate the model parameter using its identified inliers.

103



## RANSAC – Practice

- The theoretical number of samples needed to ensure 95% confidence that at least one outlier free sample could be obtained.

$$\rho = 1 - (1 - (1 - \epsilon)^k)^s$$

- Probability that a point is an outlier  $1 - \epsilon$
- Number of points per sample  $k$
- Probability of at least one outlier free sample  $\rho$
- Then number of samples needed to get an outlier free sample with probability  $\rho$

$$s = \frac{\log(1 - \rho)}{\log(1 - (1 - \epsilon)^k)}$$

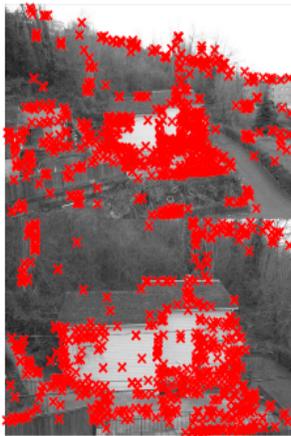
104

## Adaptive RANSAC

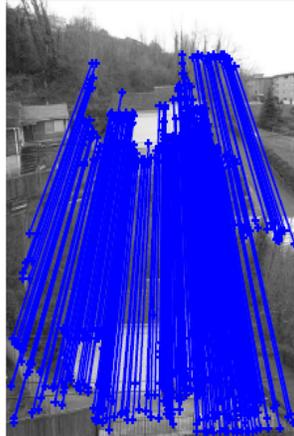
- $s = \text{infinity}$ ,  $\text{sample\_count} = 0$ ;
- While  $s > \text{sample\_count}$  repeat
  - choose a sample and count the number of inliers
  - set  $\epsilon = 1 - (\text{number\_of\_inliers}/\text{total\_number\_of\_points})$
  - set  $s$  from  $\epsilon$  and  $\rho = 0.99$
  - increment  $\text{sample\_count}$  by 1
- terminate

107

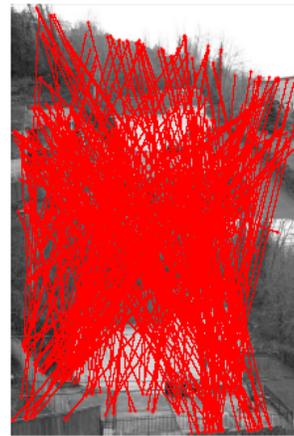
## Robust technique



(a) correspondences.



(b) identified inliers.



(c) identified outliers.

108

## Robust matching

- Select set of putative correspondences  $x_1^j, x_2^j$   

$$x_2^T F x_1 = 0$$
- Repeat
  1. Select at random a set of 8 successful matches
  2. Compute fundamental matrix
  3. Determine the subset of inliers, compute distance to epipolar line

$$d_j^2 \doteq \frac{(x_2^{jT} F_k x_1^j)^2}{\|\hat{e}_3 F x_1^j\|^2 + \|x_2^{jT} F \hat{e}_3\|^2} \quad d_j \leq \tau_d$$

4. Count the number of points in the consensus set

109