



Computer Vision 13 – Multiple view geometry, 3D and stereo

doc. dr. Janez Perš
(with contributions by prof. Stanislav Kovačič)

Laboratory for Machine Intelligence
Faculty of Electrical Engineering
University of Ljubljana

Regarding the exam

- Mandatory:
 - Check that you can log into <https://studij.fe.uni-lj.si>
 - Apply for an exam **at least five days** before exam date!
- 1st written exam date: Monday, 21.1. 18.00
- 2nd written exam date: Friday, 1.2. (time TBA)
- Results announced 2-3 days later
- Oral exam 1-2 days after results
- You **must** have your labs done and graded before sitting written exam!

Q+A regarding the exam

Q: Will there be questions regarding Matlab/code?

A: No. Matlab is only a tool we use for labs.

Q: Will there be theoretical questions in written exam?

A: Not in the sense “explain this in X sentences”. But not all questions will be purely computational. I may provide you with 2 graphs and ask you to draw the 3rd one, based on your knowledge of theory.

Q: What tools will be allowed during a written exam?

A: **Calculators** will be allowed, but NO PHONES.

Q: Which lectures should we study for the exam?

A: All of them. And the labs to refresh your knowledge.

Quick recap of the previous lectures

- Image formation
- Image analysis
- Image alignment and matching

Outline

- 3D sensing & perception
- Stereo vision
- Lateral stereo vision model
- Epipolar geometry
- Stereo matching

Motivation

- With projection from 3D world to 2D image, one dimension (depth, range, object distance) is lost.
- Fundamental questions of CV:
 - How to recover 3D representation of scene, based on one, two, or more images
- Note:
 - Many CV problems can be solved without explicit 3D reconstruction!
 - (remember – until now, we did not talk about it!)

Motivation

- Recovering 3D information is the domain of 3D computer vision, such as
 - Stereo Vision,
 - Photometric stereo,
 - Shape from X,
 - Structured Lighting, sometimes extended to “Active Stereo”
 - and “Structure from Motion” – SfM.

3D structure from 2D data

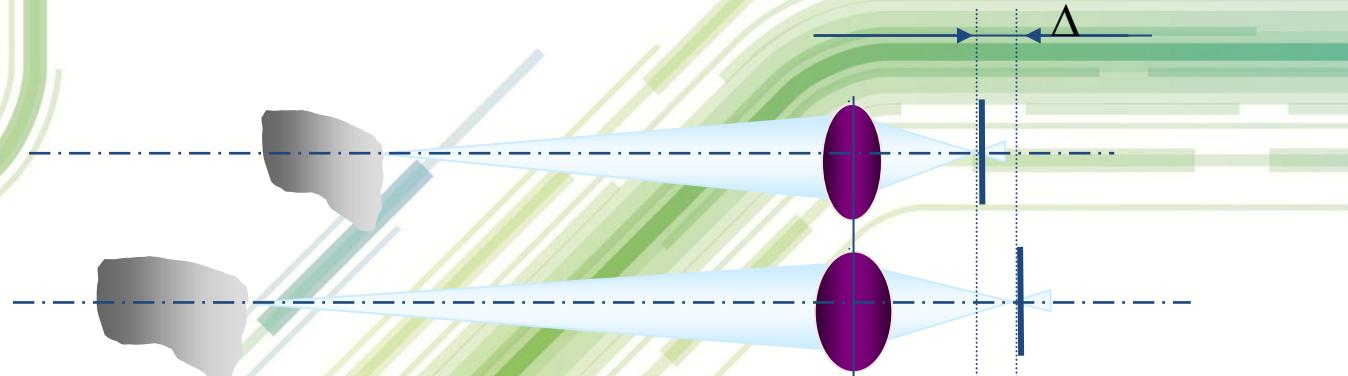
- 3D based on one image:
 - Focusing (Depth from Focus),
 - Zooming (Depth from Zooming),
 - Texture analysis (Depth from texture),
 - Shading (Depth from Shading), “Shape from X”.
 - Time of Flight cameras

3D structure from 2D data

- 3D based on one image plus “controlled structured light”
 - Photometric stereo,
 - Range finders, or scanners.
- 3D from two or more images (“multiple view”):
 - Stereo vision, “stereo”
 - Active stereo (stereo + structured lighting)
 - Depth from Motion, and Structure from Motion - SfM.

3D from focusing, zoom and texture

3D from one image - focusing

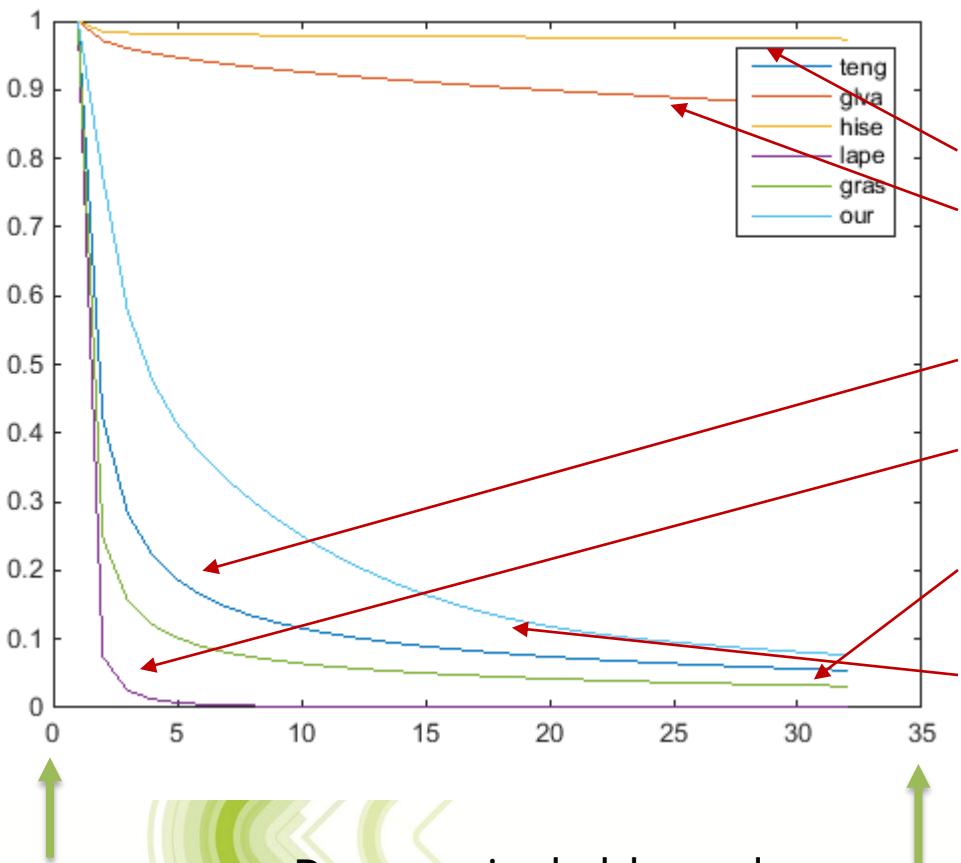


- When image is not in focus, it looks blurred, of low contrast, i.e., it *lacks high-frequency components*.
- That means that sensor is not in the imaging plane. If we could evaluate “out of focus”, be able to bring the image ‘in focus’, then, having calibrated lens, we should be able to recover object distance from Δ
- For depth recovery we need
 - Autofocusing, which needs
 - contrast measure
 - controlled lens/sensor positioning
 - calibrated lens

3D from one image - focusing

- Measures of focus:
 - grey value variance
 - if variance is high, diversity of grey values should be large, thus, "image is in focus"
 - entropy-based measures
 - compute grey value histogram, compute entropy
 - if entropy is high, contrast should be high, thus, "image is in focus"
 - gradient-based measures
 - compute gradient magnitude, sum over region
 - if gradient is high, image must be sharp, therefore, "image is in focus"
- Our work:
 - Prof. Matej Kristan (now at FRI) has devised well accepted focusing method while part of our group.
 - Check: <http://vision.fe.uni-lj.si/RESEARCH/focus/index.html>

Focusing - an example



Only a few, among many focus measures are shown!

- **hise:** graylevel entropy
- **glva:** graylevel variance
- **teng:** (mean) sum of grad magnitude squared (sobel)
- **lape:** laplacian energy, mean of laplacian squared
- **gras:** (mean) sum of squared x derivatives
- **our:** Bayes entropy and DCT
- Focus measures based on derivatives are simple and work well.

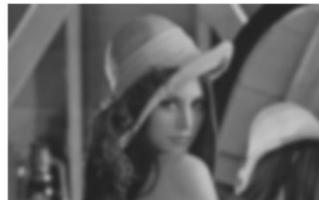


Progressively blurred

images

- • • •

Computer

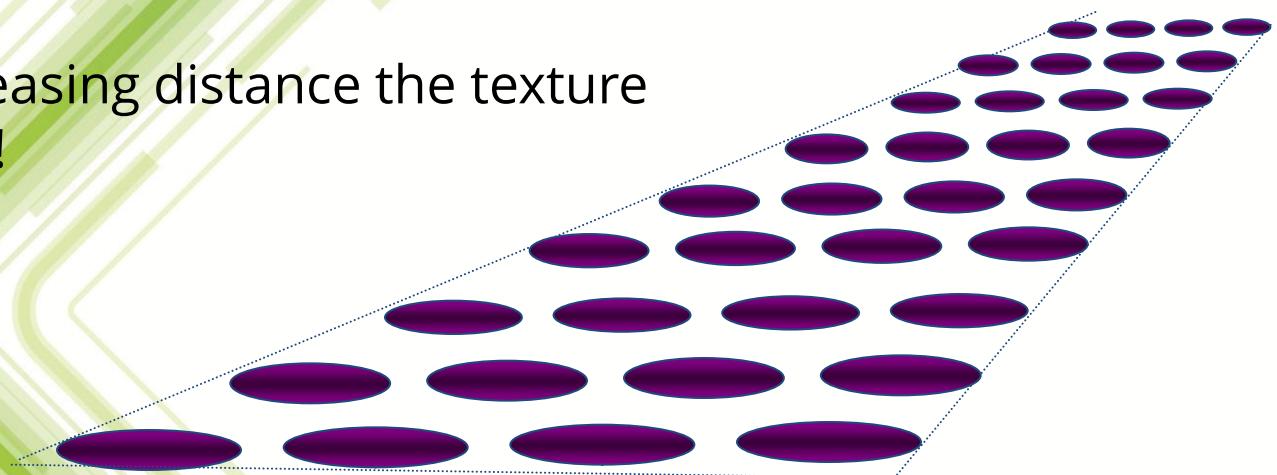


3D from one image - zooming

- Depth from zooming:
 - due to zooming pixels in image ,move', but pixels from more distant objects move less.
$$Z = (r_1 - r_2) \frac{f_1 f_2}{r_1 f_2 - r_2 f_1}$$
 - Therefore, if we could estimate these displacements in terms of depth, we should be able to recover depth.
- What is needed:
 - to evaluate relative movement of points, correspondence problem (e.g. optical flow, matching), should be solved, motorized and calibrated
 - zoom lenses are needed.

3D from focus, zoom and texture

- Note:
 - depth from focus and depth from zoom are more of qualitative value, rather than quantitative, e.g., this object is closer, or more distant,
- Another qualitative technique is “depth from texture”
 - With increasing distance the texture gets finer!



3D from structured lighting

3D from structured lighting

- Terminology:
 - (1) Structured lighting, Laser range finding, LIDAR (Light Detection and Ranging)
 - (2) Active stereo - just like two-camera stereo, but the stereo correspondence is simplified by the *controlled light source*.
 - All in all, it largely simplifies the *correspondence problem*

Polar coordinates: (ρ, β)

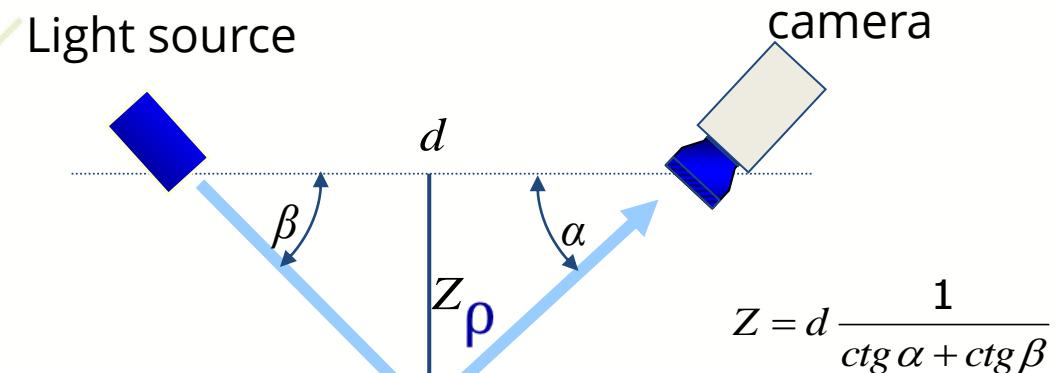
$$Z = \rho \sin(\alpha)$$

$$X = \rho \cos(\alpha)$$

$$\rho = d \sin(\beta) / \sin(\alpha + \beta)$$

$$Z = d \frac{\sin \alpha \sin \beta}{\sin(\alpha + \beta)}$$

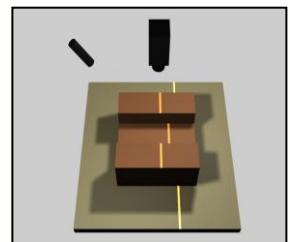
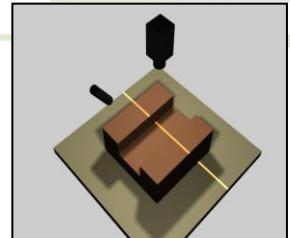
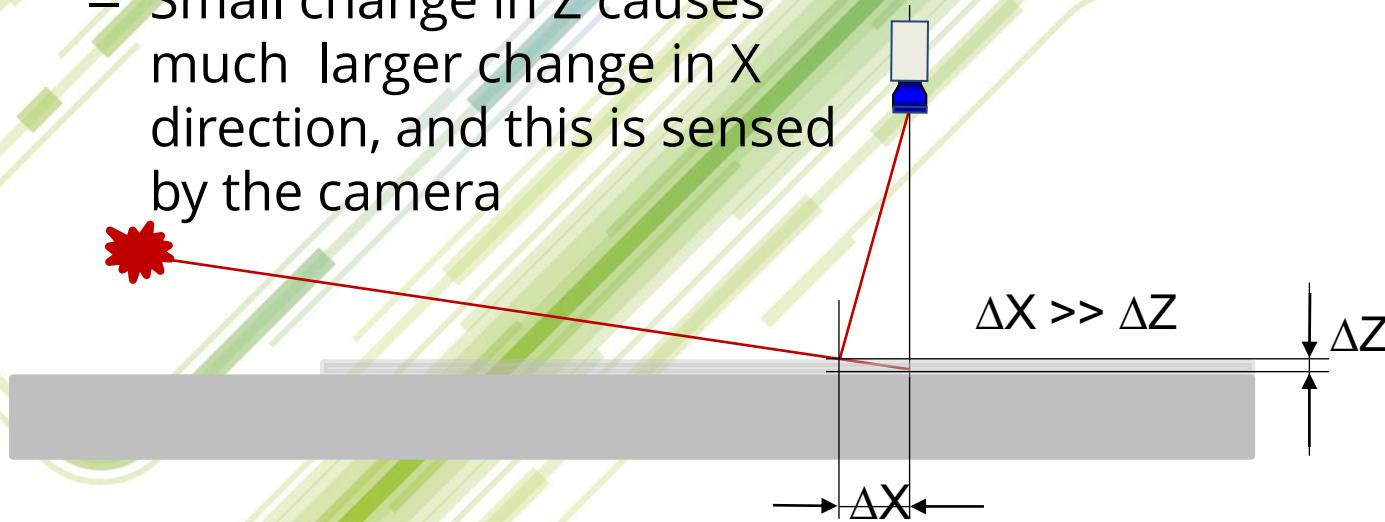
8.1.2019



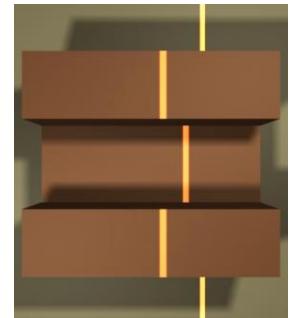
3D from structured lighting

- “Optical amplification”

- Small change in Z causes much larger change in X direction, and this is sensed by the camera

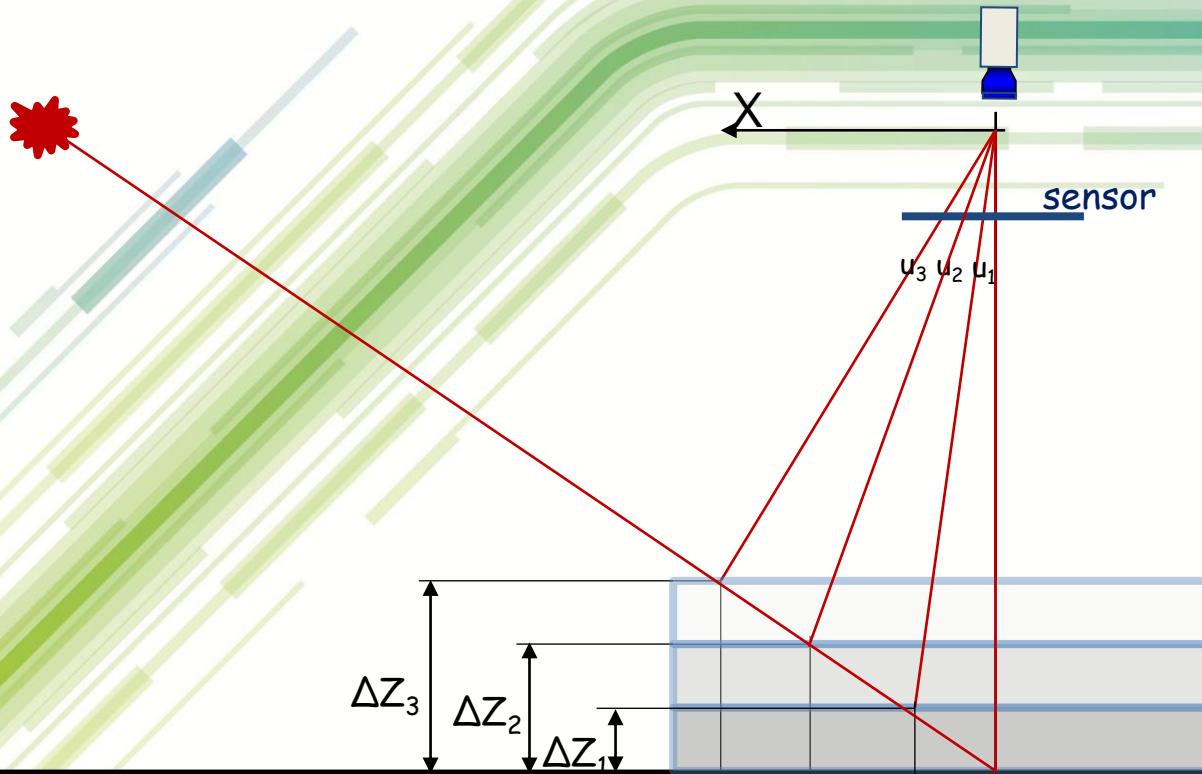


- Extensively used in industry due to high measuring speed and accuracy, $\sim 10 \mu\text{m}$ and better
 - In many industrial applications we only need to estimate the relative distances in depth
 - Therefore, the system can be calibrated in a simple way



3D from structured lighting

ΔZ_1	u_1
ΔZ_2	u_2
ΔZ_3	u_3
...	...



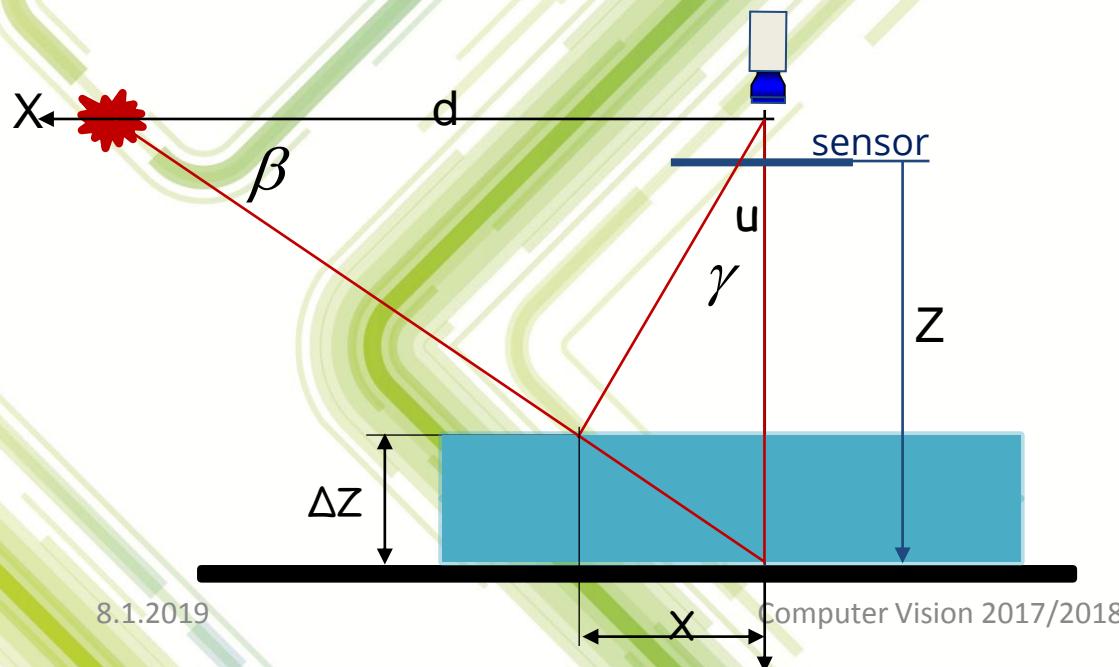
- Calibration: Stack objects of known thickness (or design calibration object) and measure line positions in the image.
- Tabulate (and interpolate) between measured positions

3D from structured lighting

$$\frac{X}{Z - \Delta Z} = \frac{u}{f} = \tan \gamma \rightarrow X = (Z - \Delta Z) \tan \gamma$$

$$Z = d \tan \beta$$

$$\frac{\Delta Z}{X} = \tan \beta \rightarrow \Delta Z = X \tan \beta = (Z - \Delta Z) \tan \gamma \tan \beta \rightarrow \Delta Z = \frac{Z \tan \gamma \tan \beta}{1 + \tan \gamma \tan \beta} = \frac{d \tan \gamma}{1 + \tan \gamma \tan \beta}$$



To measure depth at the exact position, exchange the position of camera and light source (equivalent setup!)

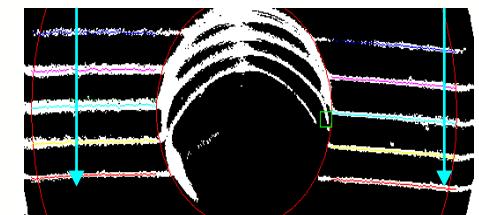
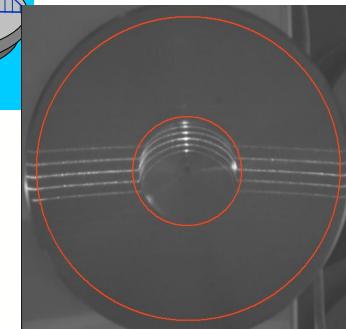
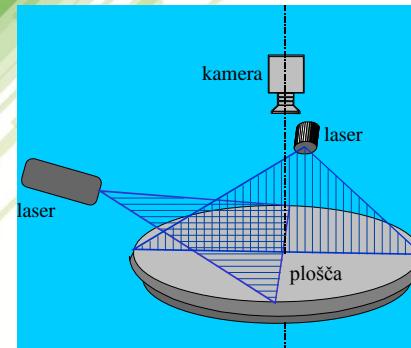
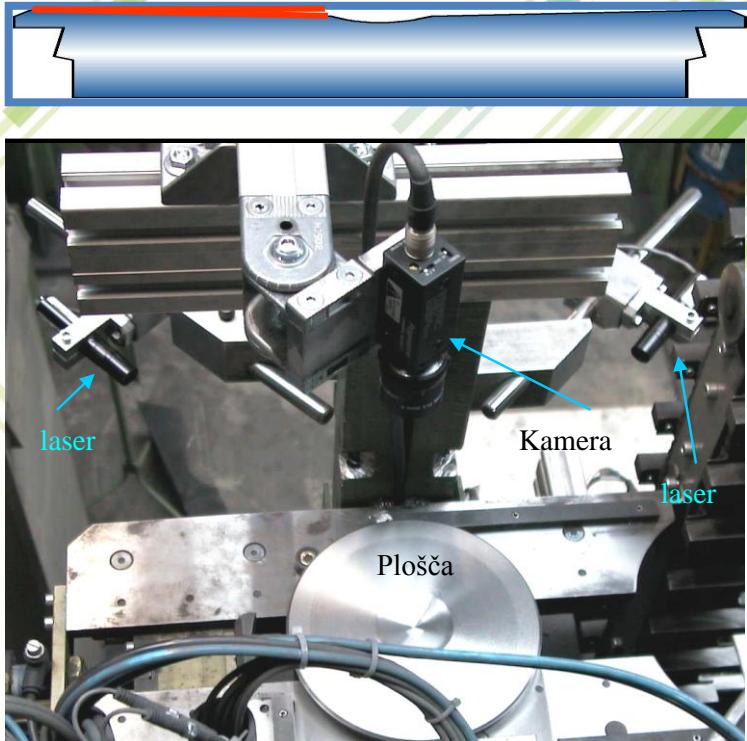
Structured lighting



Various solutions on the market!

Examples of custom design

- Concavity grading $0^\circ - 0,5^\circ$ in steps of $0,05^\circ$



Images by Franci Lahajnar

Examples of custom design

- Measuring profiles



EUREKA
Doing business
through technology

Raising the productivity and competitiveness of European technology. Boosting national economies on the international market and strengthening the basis for sustainable prosperity and well-being.

ABOUT EUREKA NEWS & MEDIA ACTIVITIES

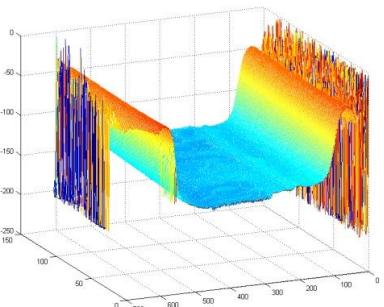
Projects Success Stories

EUREKA SUCCESS STORY > EI 3450 QSPAI

'SMART' QUALITY CONTROL SYSTEM CUTS RISK OF HUMAN ERROR ON ASSEMBLY LINES

A EUREKA-backed project has demonstrated that assembly line errors can be avoided by enabling the equipment to learn from previous actions.

A photograph of a person wearing a white protective suit and mask, working at a workstation. The workstation has various electronic displays and controls. The background shows industrial shelving and equipment.



Space-coded structured light

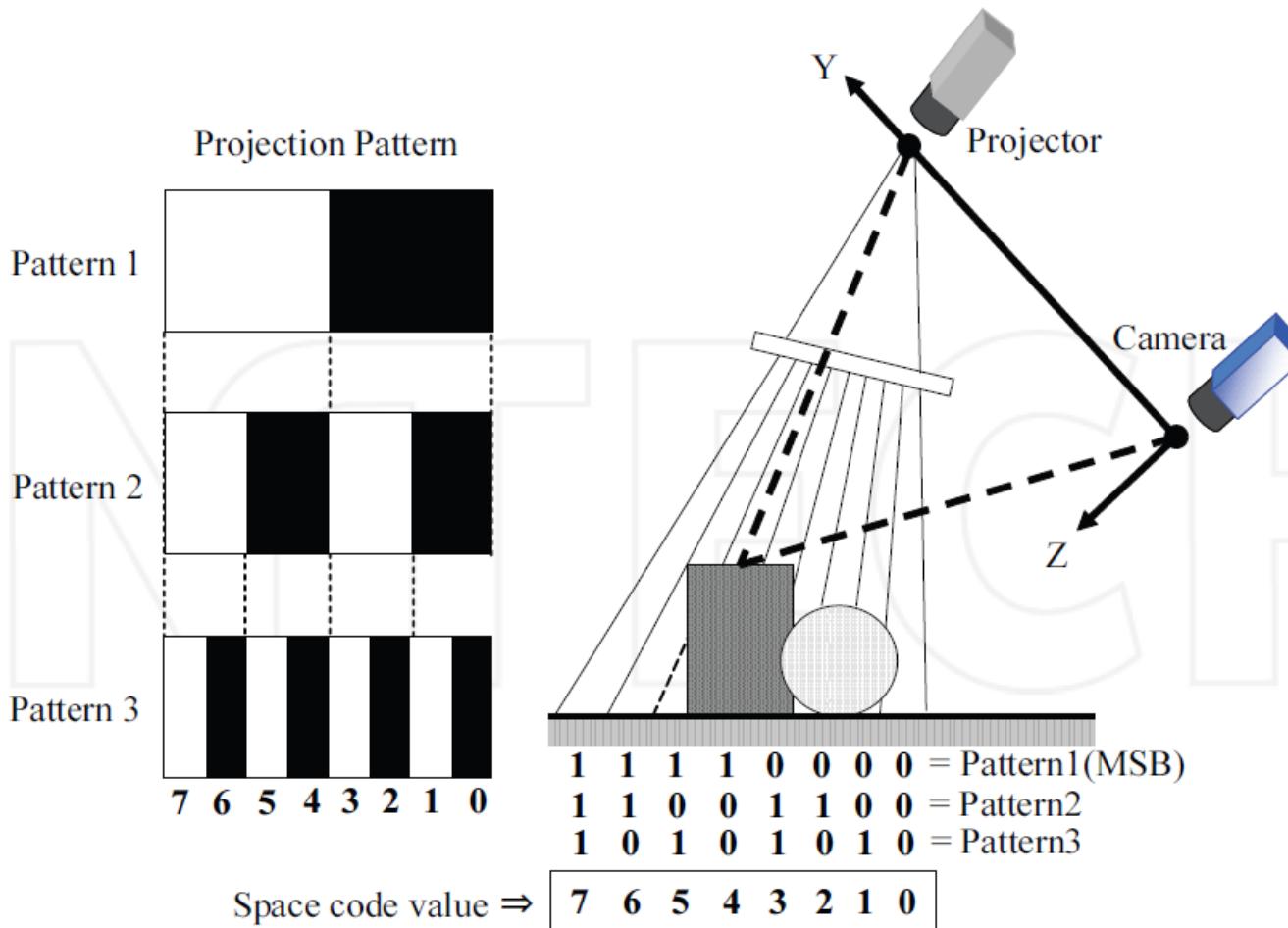
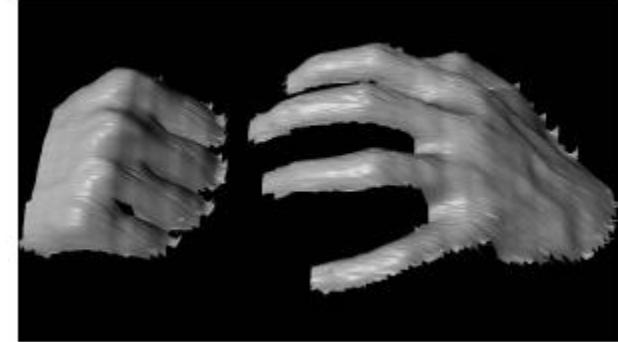
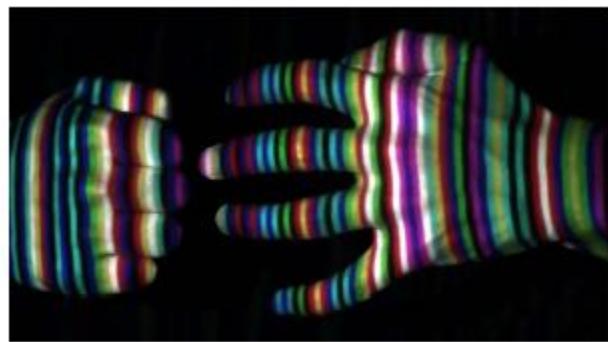


Fig. 1. Coded structured light projection method.

Source: <http://cdn.intechopen.com/pdfs-wm/32185.pdf>

Color-coded structured light



Color coding and space coding simplify/alleviate the correspondence problem.

L. Zhang, B. Curless, and S. M. Seitz. [Rapid Shape Acquisition Using Color Structured Light and Multi-pass Dynamic Programming](#). 3DPVT 2002

Kinect 360: Structured infrared light



Source: <http://bbzippo.wordpress.com/2010/11/28/kinect-in-infrared>

- Note that this is the “old” Kinect
- “new” Kinect (Kinect v2, Kinect for Windows) uses *entirely different* technology from the different company

Kinect V2: Time of flight camera

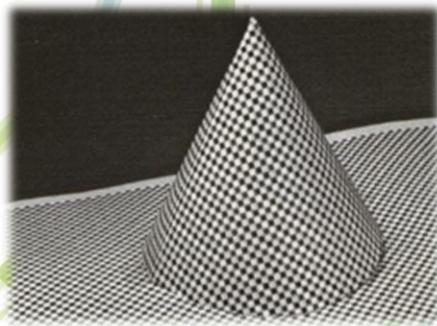
- Also known as “Kinect for windows”, Kinect for Xbox One.



- It emits modulated infrared laser light
- It is able to measure *how long does it take for light to travel to the object and back.*
- It is able to measure this for *each of the pixels!*
- The time delay (“time of flight”) directly translates into the distance of object from the sensor for each pixel!

Stereo vision

Task of stereo vision

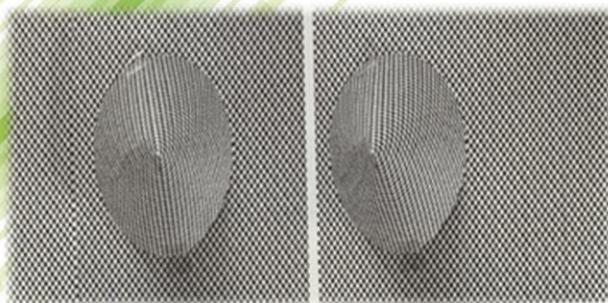


I.

This is the scene, with
3D object of a particular interest

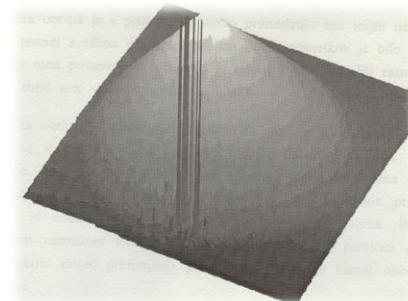
II.

Take two (or more)
2D images of
that scene

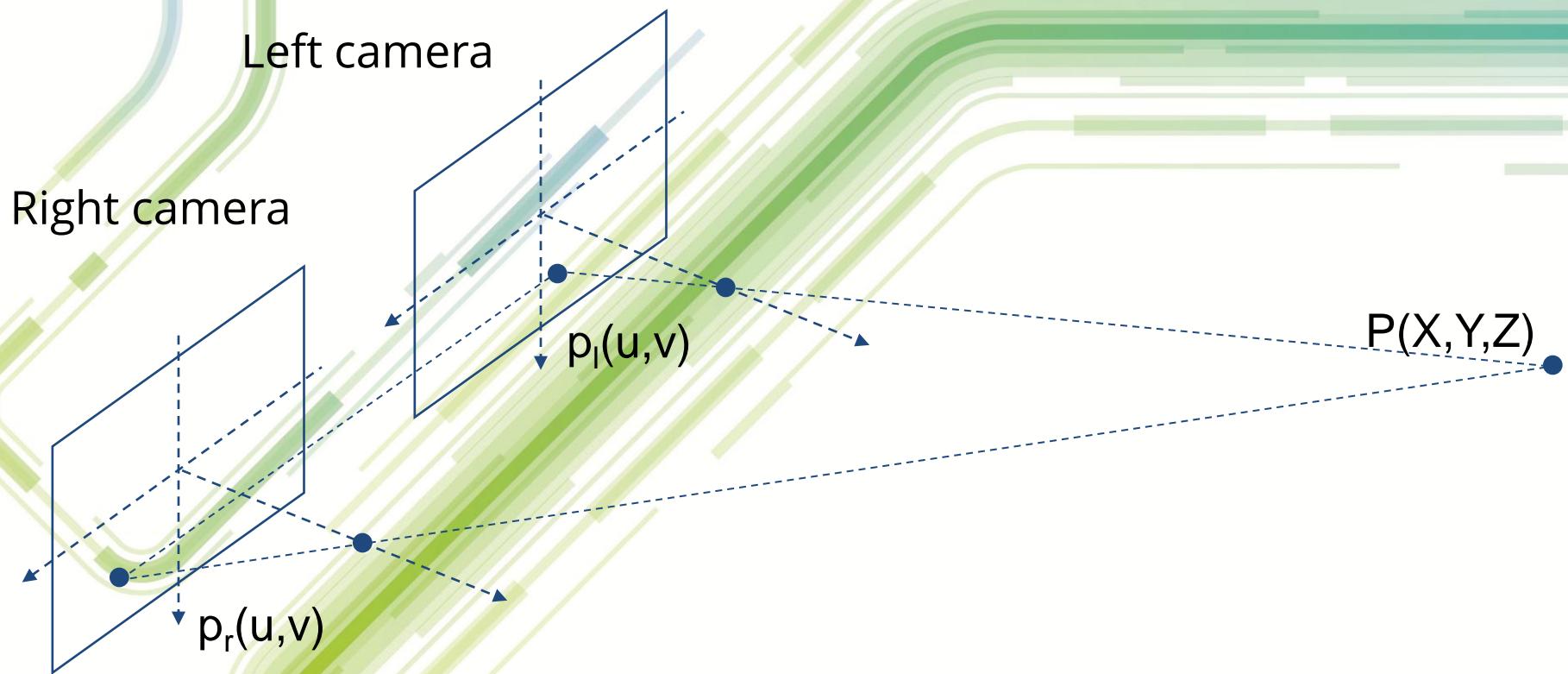


III.

Based on 2D images recover
3D structure of the scene



Basic concept



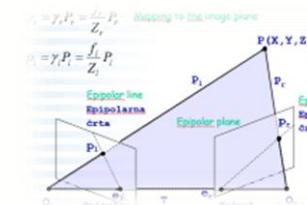
- It is possible to infer the depth from the stereo disparity (stereo parallax) $p_l(u,v) \neq p_r(u,v)$

Why stereo?

- It is biologically motivated
- Extensively studied and fairly well understood
- Solutions and technology exist
- Many needs and uses in
 - manufacturing/quality control,
 - entertainment, ...

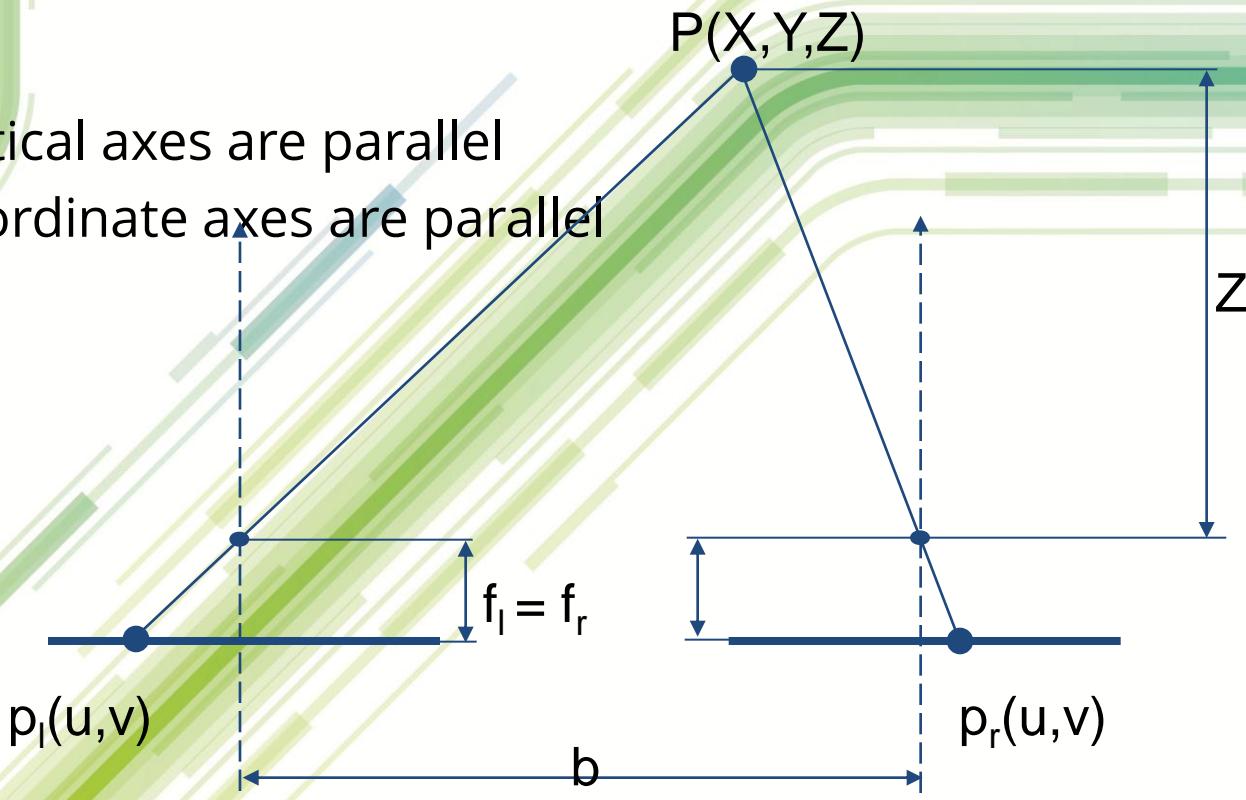
Stereo lecture outline

- Basic stereo vision model (lateral stereo model)
- Epipolar geometry, epipolar lines
- Essential matrix and Fundamental matrix
- Stereo matching



Lateral stereo model

- Optical axes are parallel
- Coordinate axes are parallel

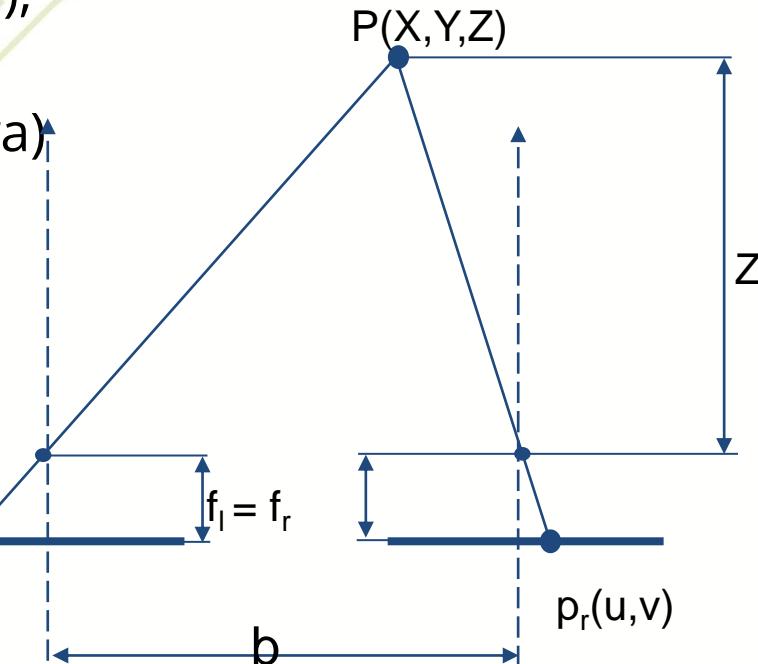


$$Z = f \frac{b}{d}, \quad d = u_r - u_l, \quad v_r = v_l$$

d = disparity!

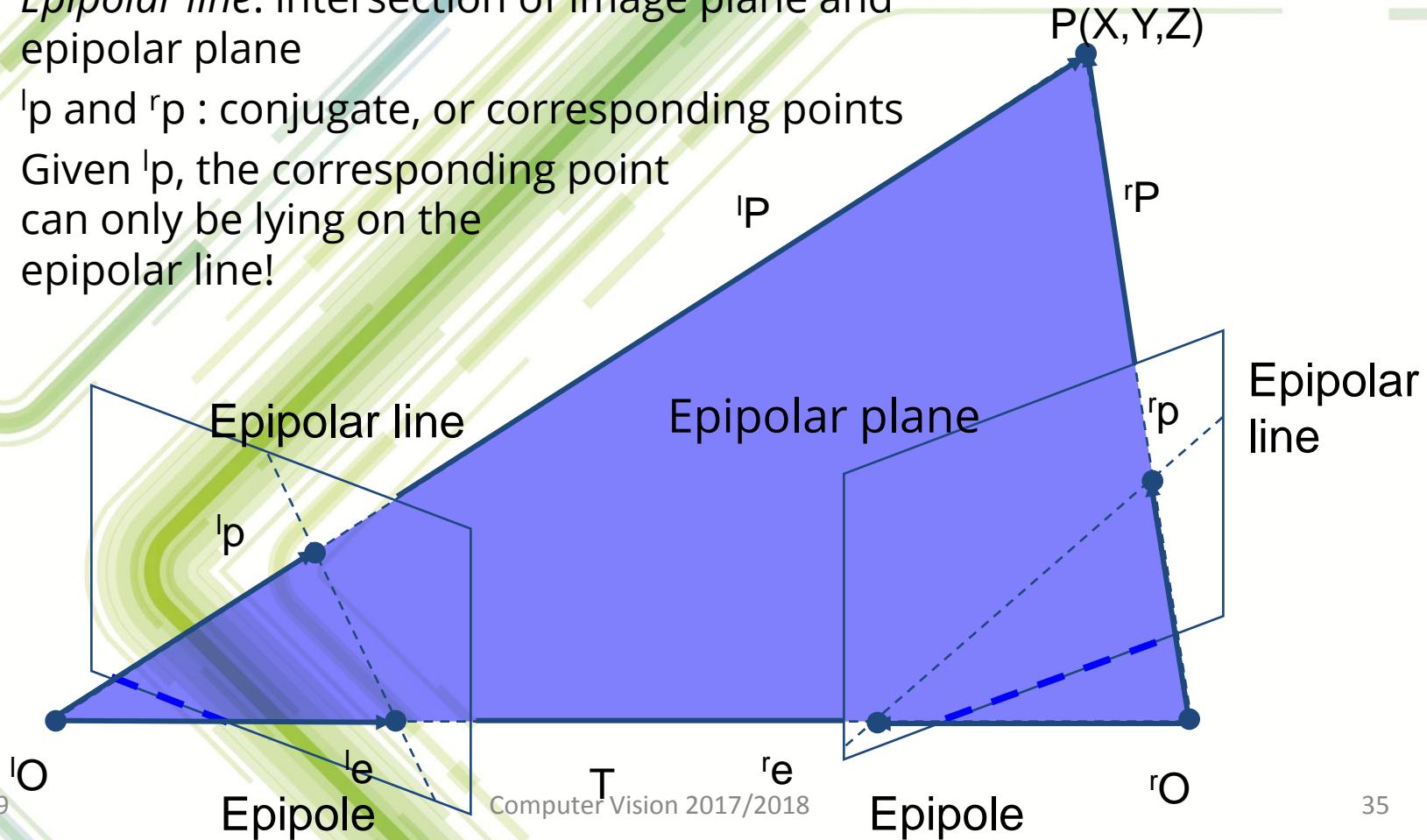
Basic parameters of the stereo setup

- Basic parameters
 - Focal lengths (left and right camera), preferably equal
 - Image centers (left and right camera)
 - Pixel size
 - Baseline distance
- In general
 - Intrinsic (internal) parameters:
 - Mapping from camera coordinates to image coordinates: focal lengths, image centers, pixel size
 - Extrinsic (external) parameters:
 - mapping between camera coordinate systems (relative position and orientation)



General case - nonparallel cameras

- *Epipolar plane*: defined by points P , ${}^l O$ and ${}^r O$
- *Epipolar line*: intersection of image plane and epipolar plane
- ${}^l p$ and ${}^r p$: conjugate, or corresponding points
- Given ${}^l p$, the corresponding point can only be lying on the epipolar line!



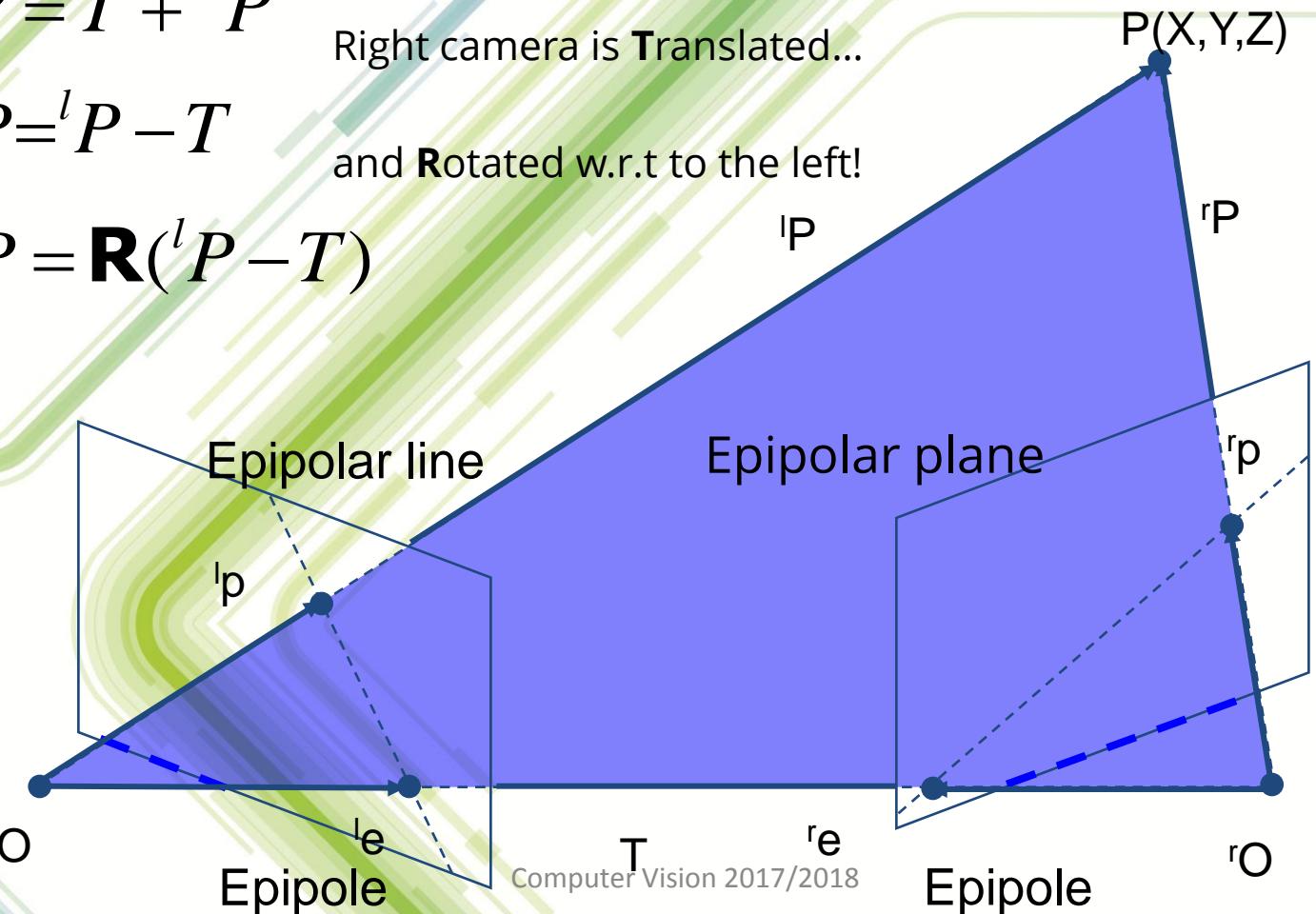
General case - nonparallel cameras

$$^l P = T + ^r P$$

$$^r P = ^l P - T$$

$$^r P = \mathbf{R} (^l P - T)$$

Right camera is **Translated...**
and **Rotated w.r.t to the left!**



General case - nonparallel cameras

$$({}^l P - T)^T (T \times {}^l P) = 0$$

Epipolar plane equation!

$$(\mathbf{R}^T {}^r P)^T (T \times {}^l P) = 0$$

Note that $({}^l P - T) = \mathbf{R}^T {}^r P$

$${}^r P^T \mathbf{R} (T \times {}^l P) = 0$$

Normal vector

$$T \times P_l$$

Epipolar line

$${}^l p$$

${}^l O$

Epipole

Epipolar plane

$${}^r P$$

Epipolar line

${}^r O$

Epipole

General case - nonparallel cameras

Op.: $T \times {}^l P = \mathbf{S} {}^l P$

$${}^r P^T \mathbf{R} \mathbf{S} {}^l P = 0$$

$${}^r P^T \mathbf{E} {}^l P = 0$$

epipolar
plane
equation

Normal
vector

$$T \times P_l$$

Epipolar line

Epipole

Epipolar plane

Epipolar
line

$$T$$

Epipole

$${}^r O$$

Essential matrix

$${}^r P^T \mathbf{E} {}^l P = 0$$

$$\mathbf{E} = \mathbf{R} \mathbf{S}$$

${}^r P$ and ${}^l P$ represent the same scene point P w.r.t left and right camera coordinate systems

- External (extrinsic) parameters
- \mathbf{E} : Essential matrix, (3×3) of rank 2
 - \mathbf{S} : translation matrix (rank 2)
 - \mathbf{R} : rotation matrix (rank 3)
- \mathbf{E} is singular, has two equal non-zero singular values and one zero singular value

Essential matrix

$${}^r P^T \mathbf{E} {}^l P = 0 \quad \mathbf{E} = \mathbf{R} \mathbf{S}$$

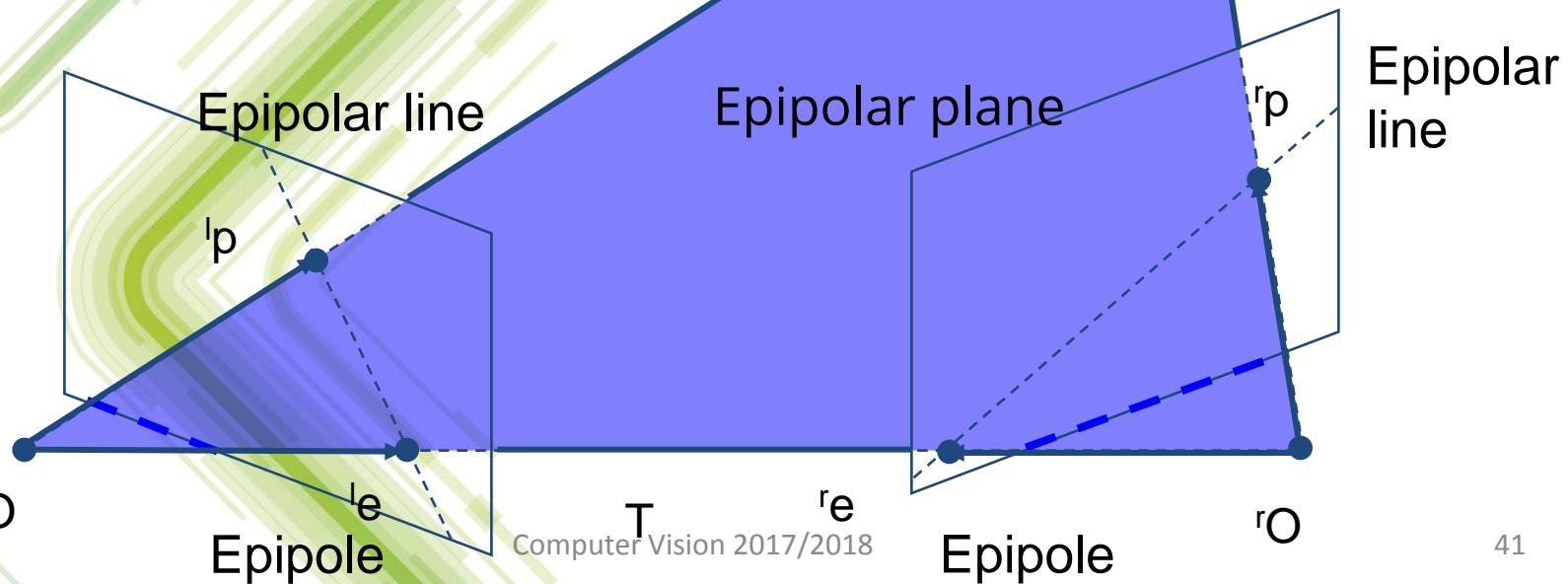
- Note,
 - given essential matrix \mathbf{E} and point position ${}^l P$ in the left camera, point position ${}^r P$ in the right camera (or the opposite), can be determined up to a scale factor!
 - \mathbf{R} has three degrees of freedom (3 angles of rotation)
 - \mathbf{S} depends on three translation components, T_x, T_y, T_z , but determined only up to a scale, thus $\mathbf{T} = [T_x, T_y, 1]$.
 - That leaves 2 DOF
 - Therefore, \mathbf{E} has $3+2 = 5$ DOF

Epipolar geometry

- Mapping to the image plane

$${}^r p = \gamma_r {}^r P = \frac{f_r}{{}^r Z} {}^r P$$

$${}^l p = \gamma_l {}^l P = \frac{f_l}{{}^l Z} {}^l P$$



Fundamental matrix

Mapping to the image plane

$${}^r p = \gamma_r {}^r P = \frac{f_r}{{}^r Z} {}^r P$$

$${}^l p = \gamma_l {}^l P = \frac{f_l}{{}^l Z} {}^l P$$

${}^l p$ and ${}^r p$ are in real coordinates

K: internal (intrinsic) parameters

$${}^r p = \mathbf{K}_r {}^r P \rightarrow {}^r P = \mathbf{K}_r^{-1} {}^r p$$

$${}^l p = \mathbf{K}_l {}^l P \rightarrow {}^l P = \mathbf{K}_l^{-1} {}^l p$$

${}^l p$ and ${}^r p$
in homogeneous
coordinates!

$${}^r P^T \mathbf{E} {}^l P = 0$$

$${}^r p^T \mathbf{K}_r^{-T} \mathbf{E} \mathbf{K}_l^{-1} {}^l p = 0$$

F: fundamental matrix

$${}^r p^T \mathbf{F} {}^l p = 0$$

Fundamental matrix

$${}^r p^T \mathbf{F} {}^l p = 0$$

$$\mathbf{F} = \mathbf{K}_r^{-T} \mathbf{E} \mathbf{K}_l^{-1}$$

$${}^r l = \mathbf{F} {}^l p$$

$${}^r p^T {}^r l = 0$$

- F: **fundamental matrix**, 3x3 of rank 2
- F depends on internal and external parameters, has 7 DOF
- l is a vector of line coefficients
- Line equation, i.e., **epipolar line equation**

$$p_r^T \mathbf{F} p_l = ({}^r p^T \mathbf{F} {}^l p)^T = {}^l p^T ({}^r p^T \mathbf{F})^T = {}^l p^T \mathbf{F}^T {}^r p = {}^l p^T {}^l l = 0$$

- Given point in one image, by knowing F, we can determine the corresponding epipolar line in the second image!
 - Finding stereo *correspondences* is much easier this way!

E, notations and conventions

Derivation of E in previous slides directly followed T.V. book.

$${}^r P^T \mathbf{E} {}^l P = 0 \quad \text{with} \quad \mathbf{E} = \mathbf{R} \mathbf{S}$$

where S is defined with respect to the left camera, and R is defined with respect to the right camera.

More often, however, you encounter different, though equivalent definition of E,

$${}^r P^T \mathbf{E} {}^l P = 0 \quad \text{with} \quad \mathbf{E} = \mathbf{S} \mathbf{R}$$

wheras R and S now define a “pose” of the left camera with respect to the right camera.

Also note: $({}^r P^T \mathbf{E} {}^l P)^T = (\mathbf{E} {}^l P)^T {}^r P = {}^l P^T \mathbf{E}^T {}^r P = 0$

E, notations and conventions

Frequently, epipolar constraint equation is written like

$${}^r \tilde{p}^T \mathbf{E} {}^l \tilde{p} = 0$$

where lower p stands for left and right image coordinates, respectively (in homogeneous form)

In this case we assume 'canonical' camera with $f = 1$ and image origin aligned with the principal point. Thus, internal parameter matrix is identity matrix.

$$\mathbf{K} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

In other words, we assume that internal parameters are known.

Derivation of E, algebraic approach

Projection matrices for the two cameras, note that $K = I$, and
 $[{}^W \mathbf{R}_l \ {}^W \mathbf{t}_l]$ and $[{}^W \mathbf{R}_r \ {}^W \mathbf{t}_r]$ define poses of cameras wrt world (W)

$${}^l \gamma \ {}^l \tilde{p} = {}^l \mathbf{K} [{}^W \mathbf{R}_l^T / -{}^W \mathbf{R}_l^{TW} \mathbf{t}_l] {}^w P = \mathbf{I} [\mathbf{I} | \mathbf{0}] {}^w \tilde{P}$$

$${}^r \gamma \ {}^r \tilde{p} = {}^r \mathbf{K} [{}^W \mathbf{R}_r^T / -{}^W \mathbf{R}_r^{TW} \mathbf{t}_r] {}^w P = \mathbf{I} [{}^W \mathbf{R}_r^T / -{}^W \mathbf{R}_r^{TW} \mathbf{t}_r] {}^w \tilde{P}$$

Because $W \equiv l$

$$[{}^W \mathbf{R}_r^T / -{}^W \mathbf{R}_r^{TW} \mathbf{t}_r] = [{}^l \mathbf{R}_r^T / -{}^l \mathbf{R}_r^{Tl} \mathbf{t}_r] = [{}^r \mathbf{R}_l / {}^r \mathbf{t}_l]$$

This is the pose of left camera wrt right camera

Derivation of E, algebraic approach

Therefore

$${}^r \gamma {}^r \tilde{p} = [{}^r \mathbf{R}_l / {}^r \mathbf{t}_l] {}^w \tilde{P} = {}^r \mathbf{R}_l {}^w P + {}^r \mathbf{t}_l = {}^r \mathbf{R}_l {}^l \gamma {}^l \tilde{p} + {}^r \mathbf{t}_l \quad / \times {}^r \mathbf{t}_l$$

Note that
 ${}^r \mathbf{t}_l \times {}^r \mathbf{t}_l = 0$

$${}^r \gamma {}^r \mathbf{t}_l \times {}^r \tilde{p} = {}^r \mathbf{t}_l \times {}^r \mathbf{R}_l {}^l \gamma {}^l \tilde{p} + 0 / \cdot {}^r \tilde{p}$$

$${}^r \gamma {}^r \tilde{p}^T \cdot {}^r \mathbf{t}_l \times {}^r \tilde{p} = {}^r \tilde{p}^T \cdot {}^r \mathbf{t}_l \times {}^r \mathbf{R}_l {}^l \gamma {}^l \tilde{p}$$

$${}^r \tilde{p}^T \cdot {}^r \mathbf{t}_l \times {}^r \mathbf{R}_l {}^l \gamma {}^l \tilde{p} = {}^r \tilde{p}^T \cdot {}^r \mathbf{S}_l \cdot {}^r \mathbf{R}_l {}^l \gamma {}^l \tilde{p}$$

Note that $\mathbf{t} \times \mathbf{v} = \mathbf{S}\mathbf{v}$

$${}^r \tilde{p}^T \mathbf{E} {}^l \tilde{p} = 0$$

Note that $\mathbf{E} = \mathbf{S}\mathbf{R}$

Simple (synthetic) example

Camera 1 parameters (,left' camera):

$f = 8 \text{ mm}$, pixel size = $10 \mu\text{m}$, image size: 1280×960 (pixels)

image center: (640,480), no lens distortion

translation: (-0.1m, 0.0, 0.0), rotation about y (vertical) axis: 8 deg

Camera 2 parameters (,right' camera), same intrinsic parameters:

$f = 8 \text{ mm}$, pixel size = $10 \mu\text{m}$, image size: 1280×960 (pixels)

image center: (640,480), no lens distortion

translation: (0.1m, 0.0, 0.0), rotation about y (vertical) axis: -8 deg

In other words: the same camera

Once rotated ,right' for 8 deg. And translated ,left' for 10cm,

Once rotated ,left' for 8 deg. And translated ,right' for 10cm.

Simple (synthetic) example

Left Camera (intrinsic) K_l :

800	0	640
0	800	480
0	0	1

Right Camera (intrinsic) K_r :

800	0	640
0	800	480
0	0	1

Left Camera pose (extrinsic) C_l :

0.9903	0	0.1392	-0.1000
0	1.0000	0	0
-0.1392	0	0.9903	0
0	0	0	1.0000

Right Camera pose (extrinsic) C_r :

0.9903	0	-0.1392	0.1000
0	1.0000	0	0
0.1392	0	0.9903	0
0	0	0	1.0000

Left Camera projection matrix M_l :

881.2852	0	522.4331	88.1285
66.8031	800.0	475.3287	6.6803
0.1392	0	0.9903	0.0139

Right Camera projection matrix M_r :

703.1437	0	745.1100	-70.3144
-66.8031	800.0	475.3287	6.6803
-0.1392	0	0.9903	0.0139

Simple (synthetic) example

Take a single scene point wrt world coordinate system

$$P = \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \begin{bmatrix} 0.08 \\ 0.15 \\ 0.40 \end{bmatrix}$$

Transform it to the left camera coordinate system,

$$\tilde{P}_l = C_l^{-1} \begin{bmatrix} P \\ 1 \end{bmatrix} = \begin{bmatrix} 0.9903 & 0 & -0.1392 & 0.0990 \\ 0 & 1 & 0 & 0 \\ 0.1392 & 0 & 0.9903 & 0.0139 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.08 \\ 0.15 \\ 0.40 \\ 1 \end{bmatrix} = \begin{bmatrix} 0.1226 \\ 0.1500 \\ 0.4212 \\ 1 \end{bmatrix}$$

Project it into the left image coordinate system

$$\tilde{p}_l = KP_l = \begin{bmatrix} 800 & 0 & 640 \\ 0 & 800 & 480 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 0.1226 \\ 0.1500 \\ 0.4212 \end{bmatrix} = \begin{bmatrix} 367.6046 \\ 322.1560 \\ 0.4212 \end{bmatrix} = \begin{bmatrix} 872.8416 \\ 764.9284 \\ 1 \end{bmatrix}$$

Point position in the left image $p_l = [873 \ 765]$

Simple (synthetic) example

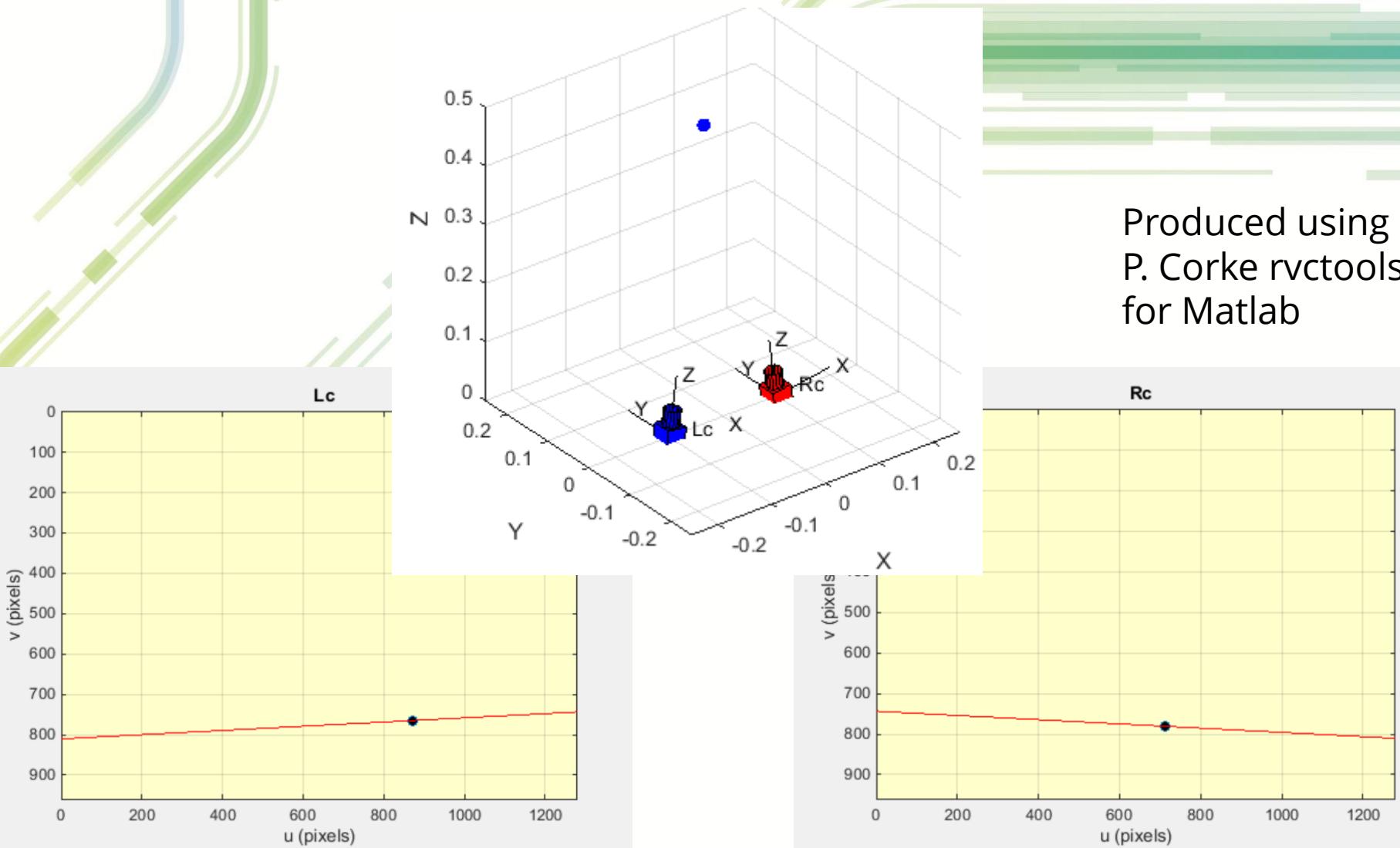
Of course, we could simply apply the camera projection matrix \mathbf{M} ,

$${}^l\tilde{\mathbf{p}} = \mathbf{M}_l \tilde{\mathbf{P}} = \begin{bmatrix} 881.2852 & 0 & 522.4331 & 88.1285 \\ 66.8031 & 800.0 & 475.3287 & 6.6803 \\ 0.1392 & 0 & 0.9903 & 0.0139 \end{bmatrix} \begin{bmatrix} 0.08 \\ 0.15 \\ 0.40 \\ 1 \end{bmatrix} = \begin{bmatrix} 367.6046 \\ 322.1560 \\ 0.4212 \end{bmatrix} = 0.4212 \begin{bmatrix} 872.84 \\ 764.93 \\ 1 \end{bmatrix}$$

And the same for the right camera

$${}^r\tilde{\mathbf{p}} = \mathbf{M}_r \tilde{\mathbf{P}} = \begin{bmatrix} 703.1437 & 0 & 745.1100 & -70.3144 \\ 66.8031 & 800.0 & 475.3287 & 6.6803 \\ 0.1392 & 0 & 0.9903 & 0.0139 \end{bmatrix} \begin{bmatrix} 0.08 \\ 0.15 \\ 0.40 \\ 1 \end{bmatrix} = \begin{bmatrix} 283.9811 \\ 311.4675 \\ 0.3989 \end{bmatrix} = 0.3989 \begin{bmatrix} 711.93 \\ 780.83 \\ 1 \end{bmatrix}$$

Simple (synthetic) example



Simple (synthetic) example

Let's take left camera as a reference. Relative orientation of the right camera wrt the left camera is:

$${}^l C_r = C_l^{-1} C_r = \begin{bmatrix} 0.9613 & 0 & -0.2756 & 0.1981 \\ 0 & 1 & 0 & 0 \\ 0.2756 & 0 & 0.09613 & 0.0278 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

with $R = \begin{bmatrix} 0.9613 & 0 & -0.2756 \\ 0 & 1 & 0 \\ 0.2756 & 0 & 0.09613 \end{bmatrix}$, $t = \begin{bmatrix} 0.1981 \\ 0 \\ 0.0278 \end{bmatrix}$, $S = \begin{bmatrix} 0 & -0.0278 & 0 \\ 0.0278 & 0 & 0.1981 \\ 0 & 0.1981 & 0 \end{bmatrix}$

The essential matrix E is,

$$E = R^T \cdot S = \begin{bmatrix} 0 & -0.0278 & 0 \\ -0.0278 & 0 & 0.1981 \\ 0 & -0.1981 & 0 \end{bmatrix}$$

Simple (synthetic) example

But, with the definition of E as $E = S R$,
we compute the pose of the first (left) camera wrt second (right) camera.

$$\begin{bmatrix} {}^r \mathbf{R}_l & {}^r \mathbf{t}_l \\ 0 & 1 \end{bmatrix} = \mathbf{C}_r^{-1} \mathbf{C}_l = \begin{bmatrix} 0.9903 & 0 & 0.1392 & -0.1 \\ 0 & 1.0000 & 0 & 0 \\ -0.1392 & 0 & 0.9903 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}^{-1} \begin{bmatrix} 0.9903 & 0 & -0.1392 & 0.1 \\ 0 & 1.0000 & 0 & 0 \\ 0.1392 & 0 & 0.9903 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

$$\begin{bmatrix} {}^r \mathbf{R}_l & {}^r \mathbf{t}_l \\ 0 & 1 \end{bmatrix} = \begin{bmatrix} 0.9613 & 0 & 0.2756 & -0.1981 \\ 0 & 1.0000 & 0 & 0 \\ -0.2756 & 0 & 0.9613 & 0.0278 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \mathbf{S} = \begin{bmatrix} 0 & -0.0278 & 0 \\ 0.0278 & 0 & 0.1981 \\ 0 & -0.1981 & 0 \end{bmatrix}$$

$$\mathbf{E} = \mathbf{S} \mathbf{R} = \begin{bmatrix} 0 & -0.0278 & 0 \\ 0.0278 & 0 & 0.1981 \\ 0 & -0.1981 & 0 \end{bmatrix} \begin{bmatrix} 0.9613 & 0 & 0.2756 \\ 0 & 1.0000 & 0 \\ -0.2756 & 0 & 0.9613 \end{bmatrix} = \begin{bmatrix} 0 & -0.0278 & 0 \\ -0.0278 & 0 & 0.1981 \\ 0 & -0.1981 & 0 \end{bmatrix}$$

This gives us the same result as before, as it should.

Making the sense of it

- Why does epipolar geometry matter?
- It constrains the search for correspondences,
 - 1D search along a line instead of 2D search within image!
- Improves robustness,
 - false matches are less likely to occur, matches that do not satisfy epipolar constraint, are rejected.

Stereo calibration

- Solving for F is non-trivial computational problem.
$$\mathbf{p}_r^T \mathbf{F} \mathbf{p}_l = 0$$
 - In principle, we need 8 scene points, i.e., 8 pairs of corresponding image points, but more than that are recommended.
 - Each scene point contributes a single (homogeneous) equation.
 - This is an optimization problem that can be solved using:
 - LS
 - SVD
 - RANSAC

Stereo calibration

- 8 point algorithm (conceptually)
 - Construct system of N equations formed by N correspondences ($N \geq 8$)
$$\mathbf{A}\mathbf{x} = \mathbf{0}$$
 - \mathbf{A} is $N \times 9$ matrix, \mathbf{x} is vector of unknowns (i.e., \mathbf{F} is $3 \times 3 = 9$ unknowns)
 - Compute SVD of \mathbf{A} ,
$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$$
 - The solution \mathbf{x} is column of \mathbf{V} for the smallest singular value in Σ

Stereo calibration

- 8 point algorithm (conceptually)
 - Construct system of N equations formed by N correspondences ($N \geq 8$)

$$\mathbf{A}\mathbf{x} = \mathbf{0}$$

- A is $N \times 9$ matrix, x is vector of unknowns (i.e., F is $3 \times 3 = 9$ unknowns)
- Compute SVD of A,

$$\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$$

- The solution x is column of V for the smallest singular value in Σ
- (but this is not the end)

Stereo calibration

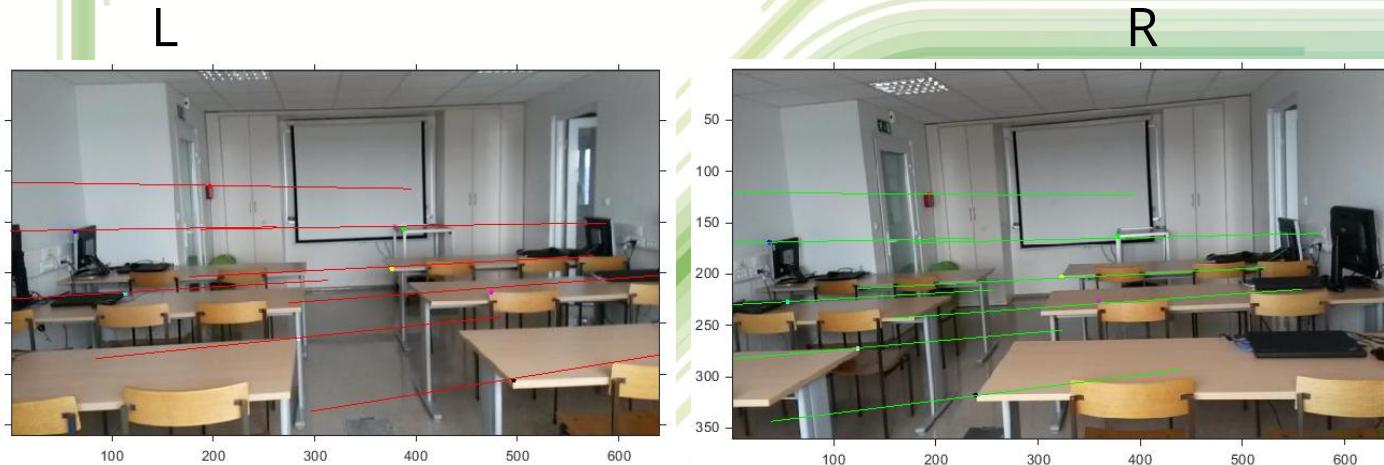
- To enforce singularity of F , compute SVD of F ,

$$\mathbf{F} = \mathbf{U}_F \mathbf{D}_F \mathbf{V}_F^T$$

- Set the smallest singular value in D to zero, $D \rightarrow D_c$, compute corrected F

$$\mathbf{F} = \mathbf{U}_F \mathbf{D}_C \mathbf{V}_F^T$$

Simple example, estimate F

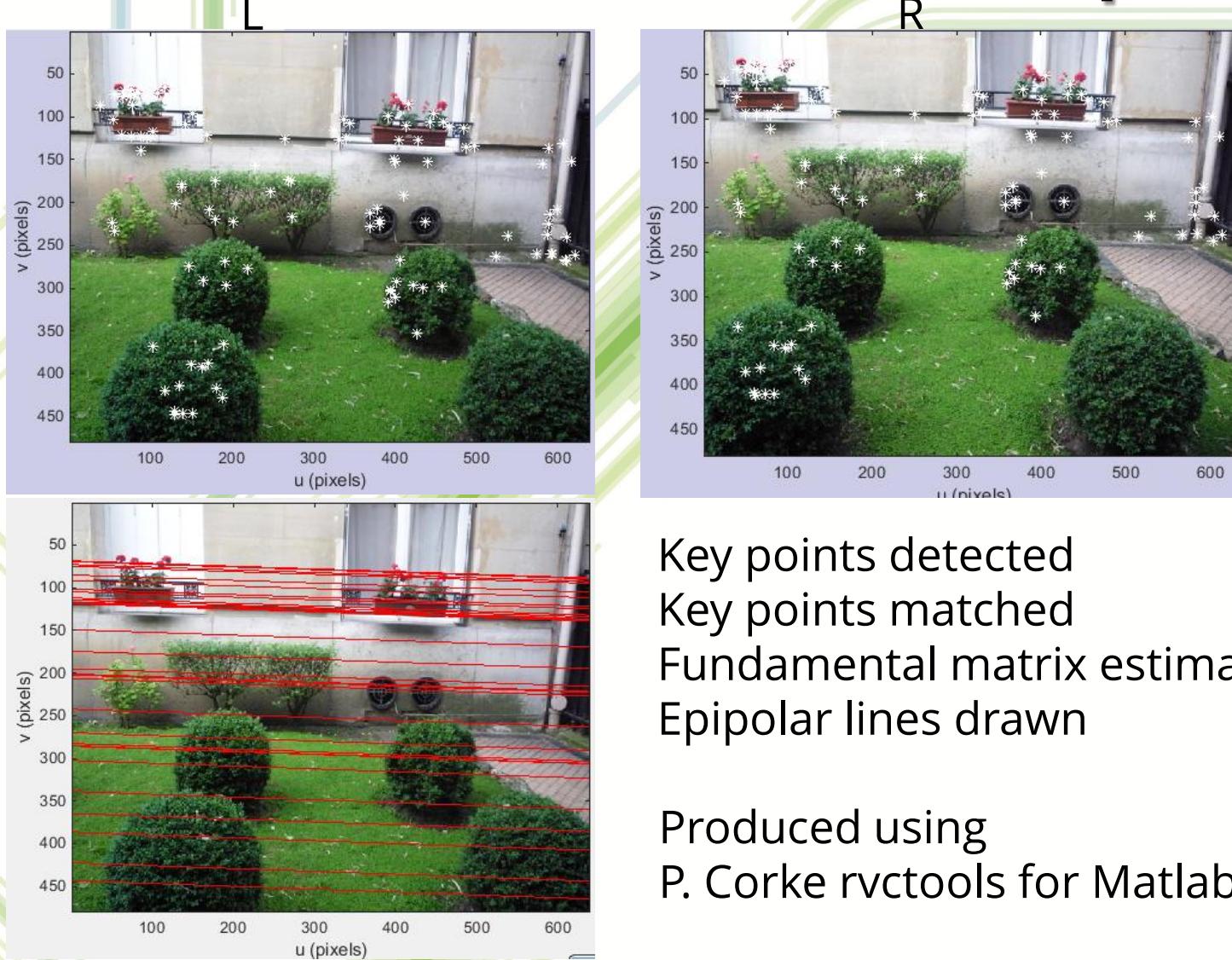


F was estimated from 8 correspondences (marked points), then epipolar lines were drawn for each point pair.



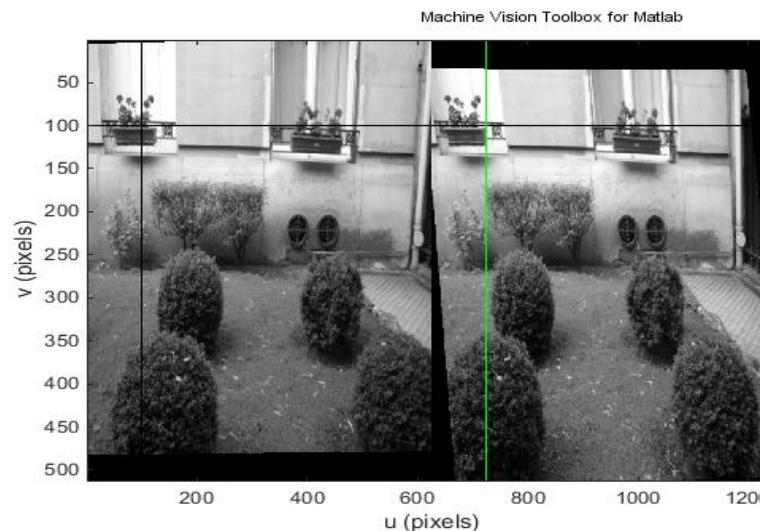
L and R image shown in red and green color plane, along with epipolar lines.

Stereo calibration example



Rectification

- The task:
 - Vergent cameras (non-parallel) -> transform images to ,ideal stereo pair'
 - Result: Horizontal lines are also epipolar lines. The images look like the images would have been taken using parallel cameras



Produced using P. Corke rvctools for Matlab

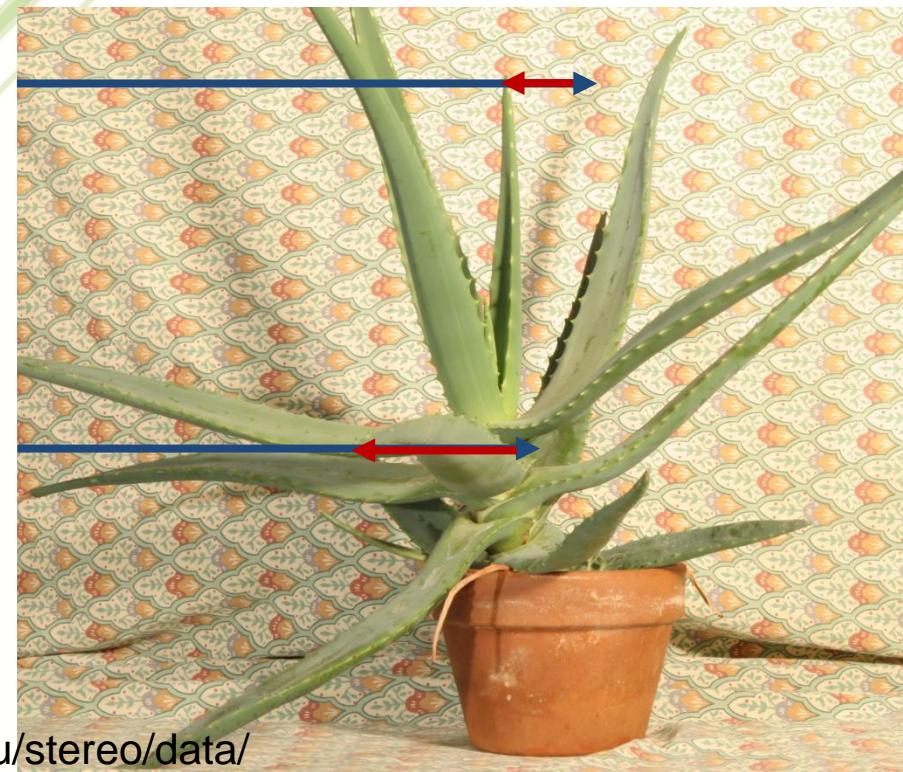
Stereo matching

The goal: reconstruction

- Based on two slightly different images of the same scene, caused by camera translation, recover the ‚lost‘ (3rd) dimension (distance, range, depth).



<http://vision.middlebury.edu/stereo/data/>



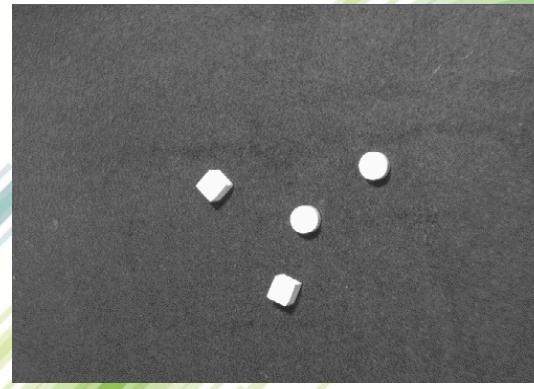
The stereo problem

- Basic problems (“issues”) of the stereo vision:
 - Find correspondence, and therefore disparity, i.e., find points in the left and the right image that are projections of the same scene point.
 - The solution to this problem is called *,stereo matching’*
 - Once the correspondence is solved, either for small or large number of points, compute 3D representation (‘structure’) of the scene. That is called *reconstruction*.
 - Nevertheless, to do that, stereo *calibration* is needed.

Is stereo matching hard?

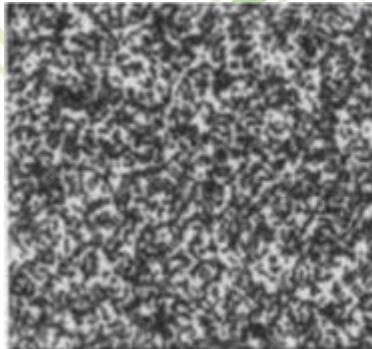


L

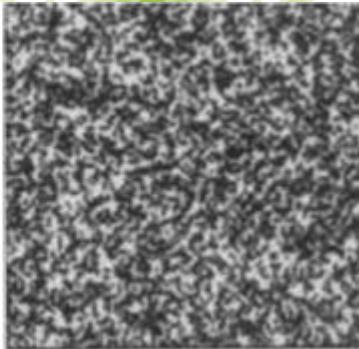


R

Trivial

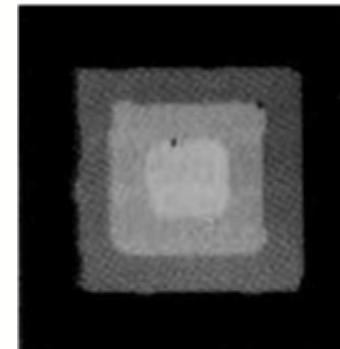


L



R

Non-trivial



Disparity map / field

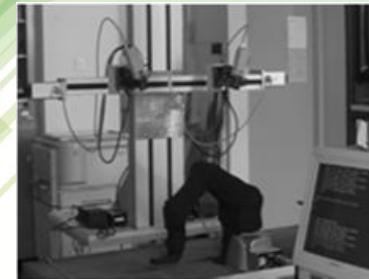
D. Marr 82: Vision

Stereo cameras (rigs, heads)

TRC ~ 95
(active stereo head)



FE ~ 95



Point Grey
Bumblebee



IDS
(stereo+structured light)

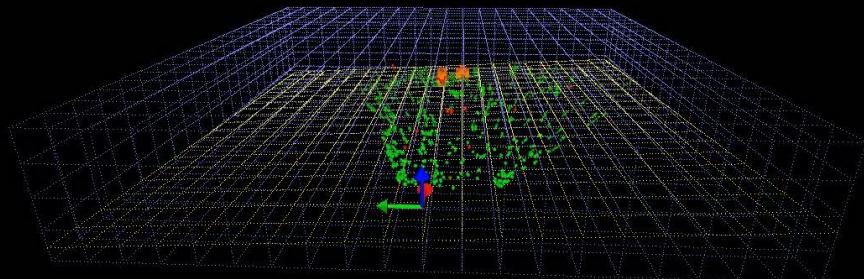
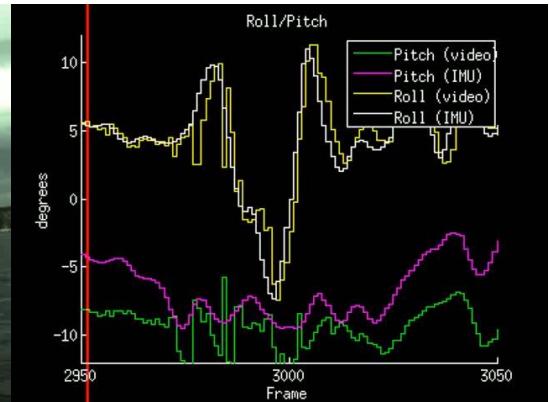
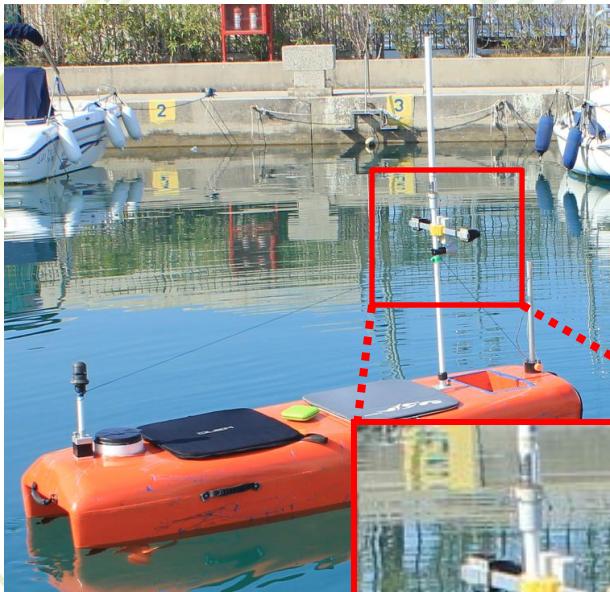


Nao

Computer Vision 2017/2018

Stereo cameras (rigs, heads)

- FE LSV & Harpha Sea (2012)



Stereo matching

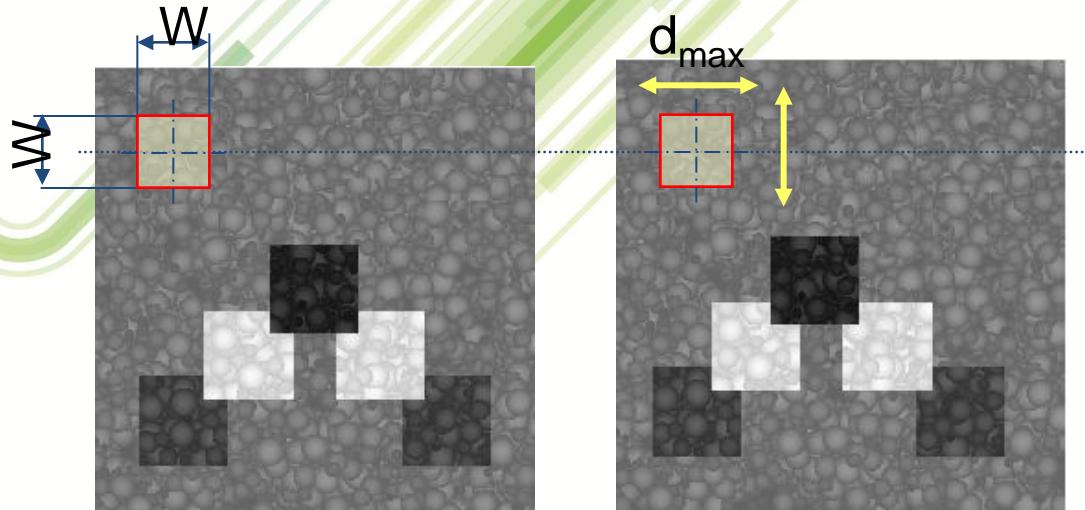
- For a point in one (left) image find the corresponding point in the other (right) image
- Once the correspondence is found, compute disparity
- Questions:
 - which points to match?
 - how?
- What is needed:
 - matching strategy
 - similarity measure

Stereo matching strategies

- Stereo matching strategies:
 - Correlation based (region based, ,dense' stereo).
 - Take a (square) patch in one image and search for a match in the other image.
 - Choose similarity measure (SM): SAD, SSD, NCC, ZNCC, ...
 - search for match, find maximum, or minimum of SM
 - Feature based (,local image features based')
 - detect key points (distinct image points), edges, corners, ...
 - do that independently in both images
 - Perform key-point matching, search for matching point pairs,
 - using predefined distance/similarity measure.

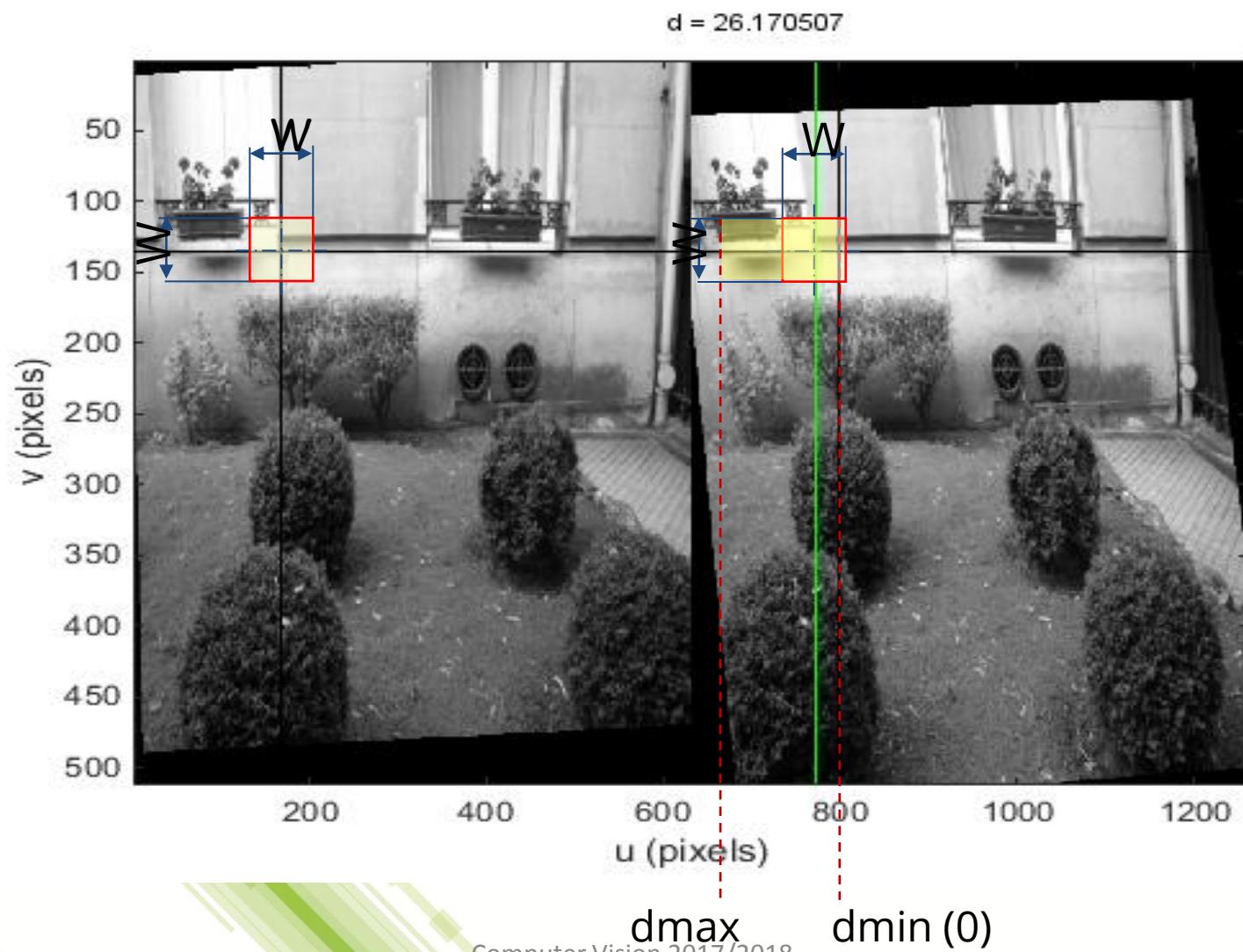
Stereo matching

- Correlation based – what do we need?
 - select window size, the size of regions involved in matching
 - select search area, i.e., how far to search for correspondences
 - search for best match using predefined similarity measure



- Note: synthetic images are quite convenient for the evaluation of matching algorithms.

Stereo matching



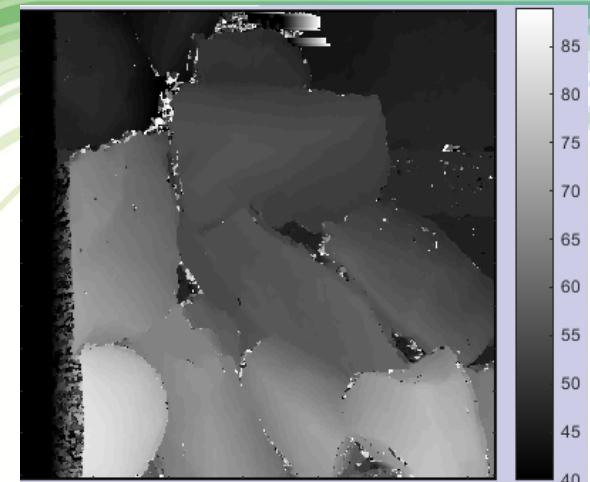
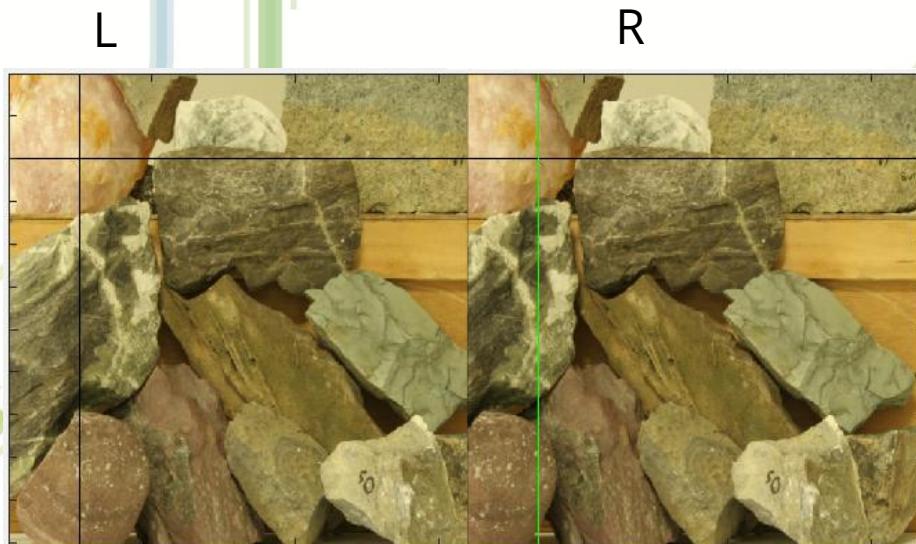
Stereo matching

- Correlation based (we talked about that before)
 - Select similarity measure, $S(d)$, $d = (u, v)$:
 - Sum of absolute differences (within window), SAD,
 - Sum of square differences, SSD,
 - Correlation coefficient, CC,
 - Normalized correlation, NCC, ZNCC.
 - compute similarity measure as a function of displacement d , $S(d)$,
 - search for the maximum

$$S(d(i, j)) = \sum_{k=-w}^{k=+w} \sum_{l=-w}^{l=+w} I_l(i + k, j + l) \bullet I_r(i + k + u, j + l + v)$$

$$d(i, j) = \arg \max_d \{S(d)\}$$

Stereo matching example



Similarity measure: ZNCC

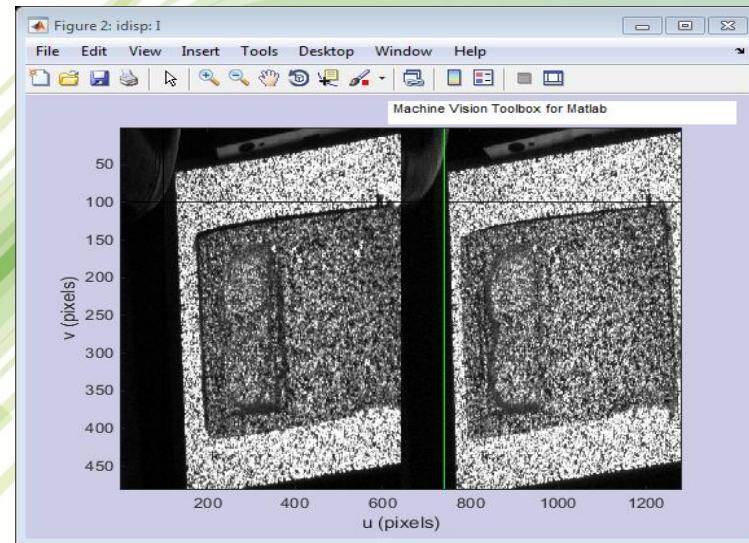
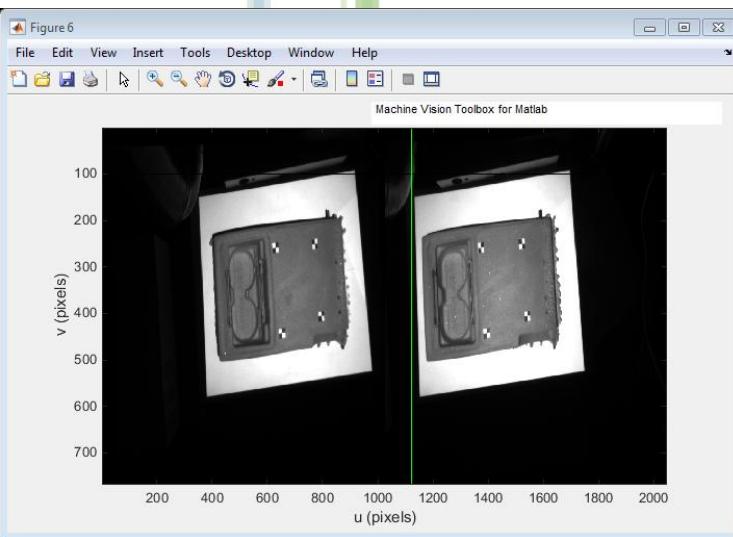
Window size: 7 x 7

Produced using P. Corke RVCTools for Matlab

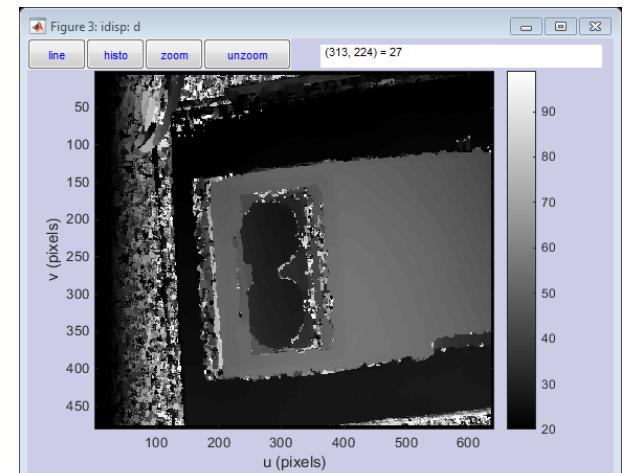
Stereo matching

- Correlation based approaches:
 - produce dense disparity field (map), and therefore dense depth (+)
 - require textured surfaces
 - produce many false matches (-)
- Note:
 - Regions not covered by both cameras cannot be matched
 - Regions hidden/self occluded in one image cannot be matched.

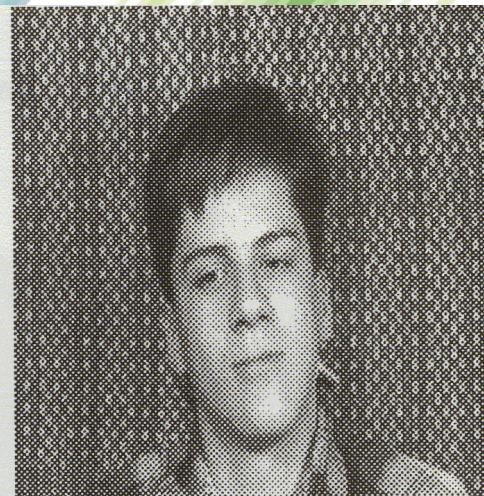
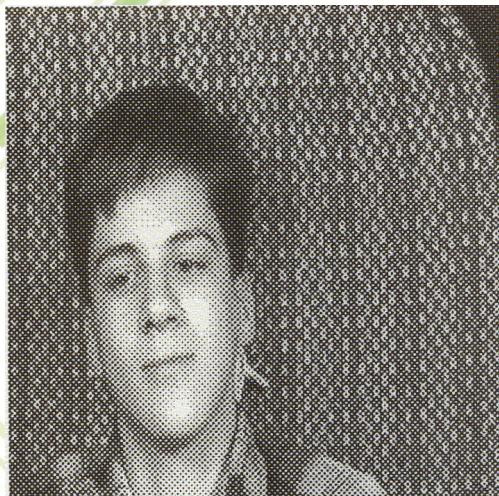
Stereo matching example



- Regions without texture (homogeneous regions) cannot be matched.
- Solution: *Stereo + Structured lighting!*
- But, occlusions are always problematic
- Furthermore, false matches due to
- ambiguity are likely to occur.



Stereo example - a face

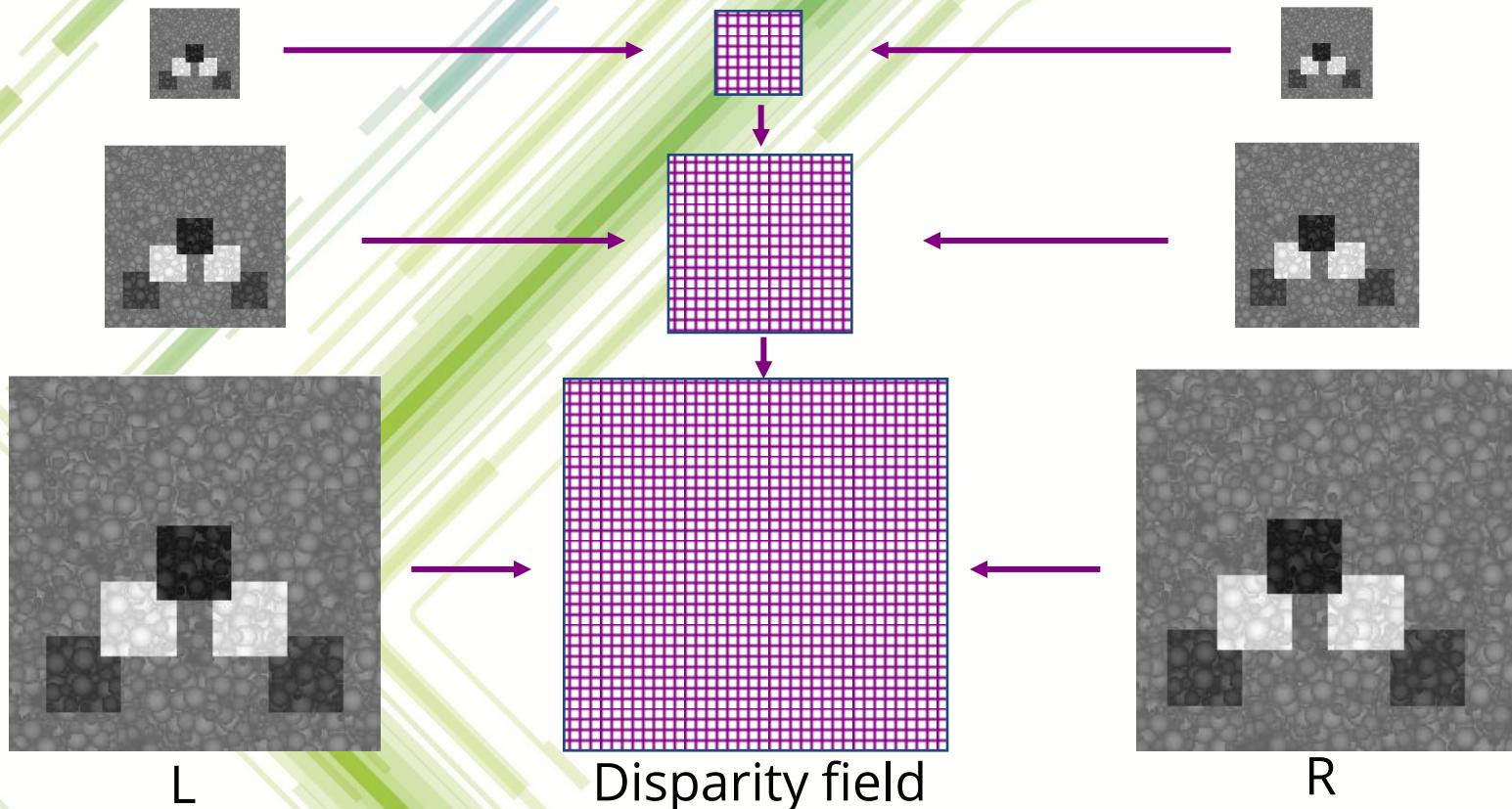


Correlation based approaches

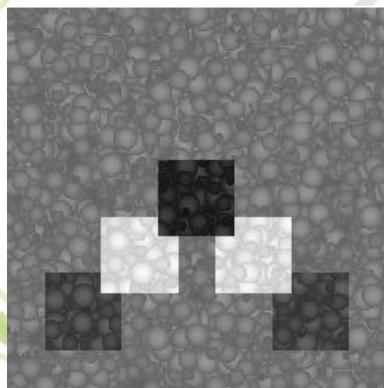
- Questions:
 - Window size ($3 \times 3, 5 \times 5, 7 \times 7, \dots$) ?
 - Search area ?
 - In practice derived experimentally, and additional knowledge.
 - Alternatives:
 - Multiresolution approaches
(start with large window, then reduce the window size)
 - Pyramidal approaches
(start with small resolution images, then proceed to higher resolutions, window size is kept constant, however.)

Stereo matching

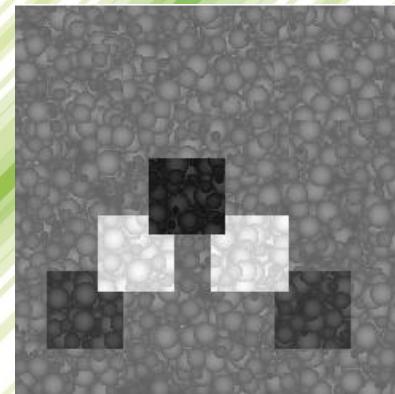
Image pyramid (used in coarse-to-fine approach)



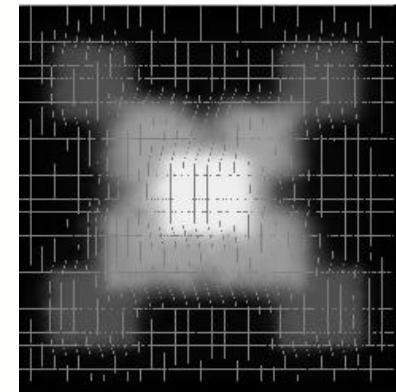
Stereo matching - an example



Left image



Right image

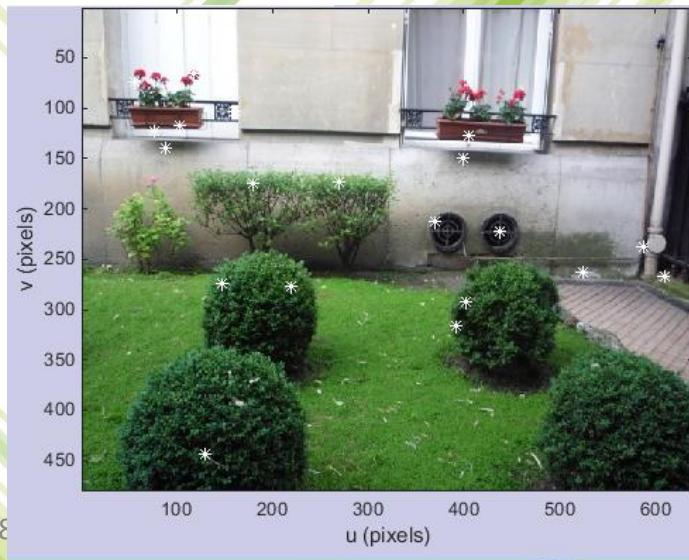
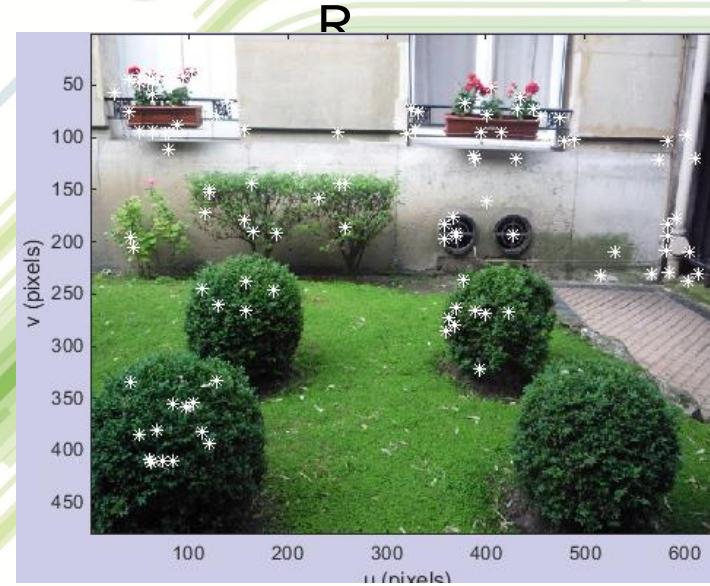
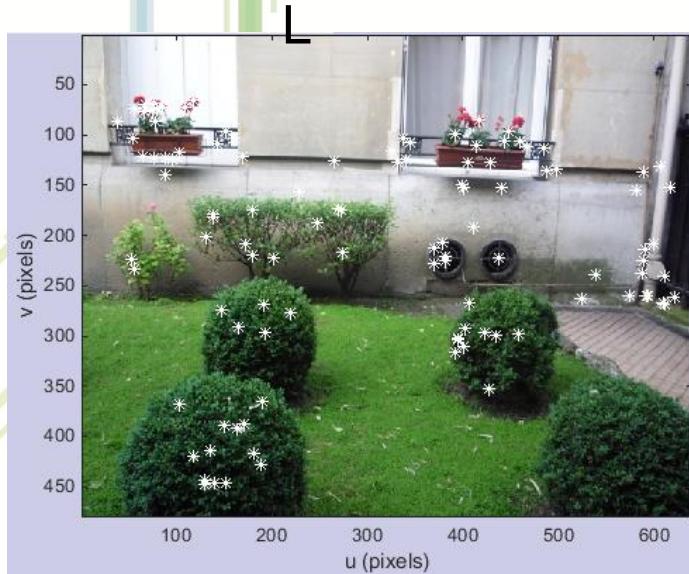


Depth image

Stereo matching

- Feature based approaches:
 - Detect distinctive image points in both images
 - various ,detectors', and ,descriptors' exist
 - edges, lines
 - Corners
 - match these features for each point in one image search for the corresponding point in the other image.
 - Produce sparse depth map,
 - interpolation is needed
 - less false matches
 - less sensitive to photometric variability

Stereo matching - an example



Key points detected
Epipolar geometry computed
Stereo matching performed
Matched points shown

Produced using P. Corke rvctools for Matlab

Stereo vision

- Observations:
 - Disparity decreases as distance increases (in case of parallel cameras)
 - Error grows with square of distance (Z)
 - (stereo works well for small distances)
 - As baseline (b) increases, error decreases, but due to perspective, the images are more and more different and consequently, false matches are more likely to occur.
- Some solutions employ 'multi-baseline' stereo approach.

Stereo vision - analysis of errors

$$Z = \frac{fb}{d}$$

$$\Delta Z = \frac{\partial Z}{\partial b} \Delta b + \frac{\partial Z}{\partial f} \Delta f + \frac{\partial Z}{\partial d} \Delta d$$

$$\Delta Z = \left| \frac{f}{d} \right| \Delta b + \left| \frac{b}{d} \right| \Delta f + \left| \frac{fb}{d^2} \right| \Delta d$$

$$\Delta Z = \left| \frac{z}{b} \right| \Delta b + \left| \frac{z}{f} \right| \Delta f + \left| \frac{z}{d} \right| \Delta d$$

$\Delta b, \Delta f \ll \Delta d$ ($\Delta d = 1 \text{ pixel}$) discretization error

$$\Delta Z \approx \left| \frac{z}{d} \right| \Delta d = \frac{z}{bf} \Delta d = \frac{z}{bf_u} \lambda_u \Delta d$$

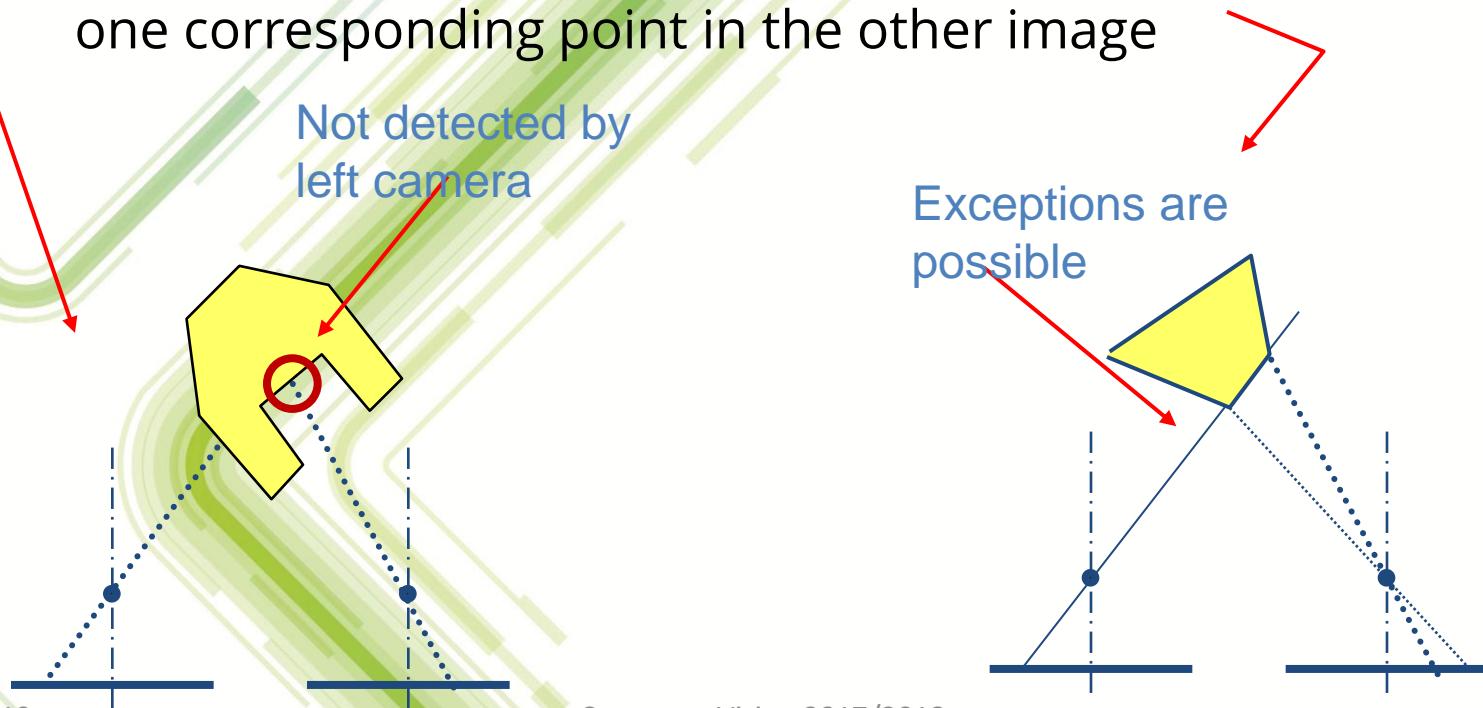
pixel size

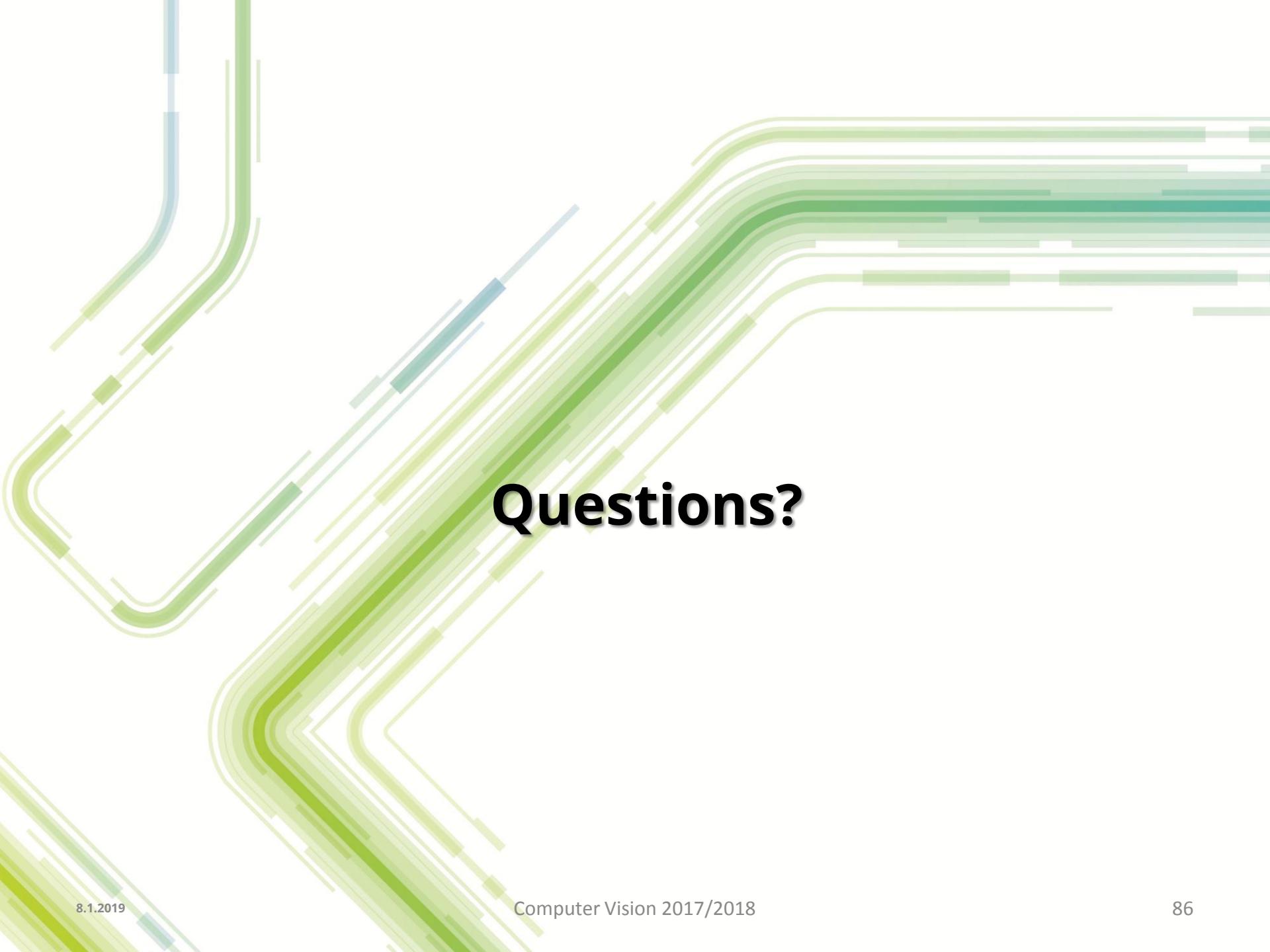
focal length in pixels

(discretization error = 1 pixel)

Key stereo matching problems

- perfect match does not exist, images are different
- Therefore, additional constraints are helpful
- Epipolar geometry - limits search and false matches
- Uniqueness constraint - point in one image can have at most one corresponding point in the other image





Questions?