

Perception II: Pinhole camera and Stereo Vision Autonomous Mobile Robots

Davide Scaramuzza

Margarita Chli, Paul Furgale, Marco Hutter, Roland Siegwart

ETH zürich University of Zurich

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Mobile Robot Control Scheme





Computer vision | definition

University of Zurich[™]

ETHzürich

Automatic extraction of "meaningful" information from images and videos



Semantic information



Geometric information

ETH zürich () University of Zurich

Computer vision | applications

- 3D reconstruction and modeling
- Recognition
- Motion capture
- Augmented reality:
- Video games and tele-operation
- Robot navigation and automotive
- Medical imaging



ROBOTICS &



Google Earth, Microsoft's Bing Maps



Mars rover Spirit used cameras for visual odometry







Sony Cybershot WX1

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart ROBOTICS & ASL PERCEPTION Autonomous Systems Lab GROUP



• If we place a piece of film in front of an object, do we get a reasonable image?



ETH zürich University of ZurichTM ASL GROUP Autonomous Systems Lab

The camera | image formation

- If we place a piece of film in front of an object, do we get a reasonable image?
- Add a barrier to block off most of the rays
 - This reduces blurring
 - The opening is known as the aperture



The camera | camera obscura (pinhole camera)

- Pinhole model:
 - Captures beam of rays all rays through a single point
 - The point is called Center of Projection or Optical Center
 - An "inverted" image is formed on the Image Plane
- We will use the pinhole camera model to describe how the image is formed





ROBOTICS &

PERCEPTION Autonomous Systems Lab

Gemma-Frisius (1508–1555)



University of Zurich^{™™}



Home-made pinhole camera



What can we do to reduce the blur?

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart Based on slide by Steve Seitz



0.35 mm

0.6mm

Why not make the aperture



Lecture 5 Thukiversity of a Model

Shrinking the aperture





Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart Why not make the aperture as small as possible?

- Less light gets through (must increase the exposure)
- Diffraction effects...



ETH zürich University of Zurich^{uth} ASL Autonomous Systems Lab

The camera | why use a lens?

- The ideal pinhole: only one ray of light reaches each point on the film
 - ⇒ image can be very dim; gives rise to diffraction effects
- Making the pinhole bigger (i.e. aperture) makes the image blurry





- A lens focuses light onto the film
- Rays passing through the optical center are not deviated





The camera | why use a lens?

- A lens focuses light onto the film
- Rays passing through the optical center are not deviated
- All rays parallel to the optical axis converge at the focal point



ETHzürich **University** of **Zurich**^{was}

The camera | pinhole approximation

• What happens if $z \gg f$?



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ASL

Autonomous Systems Lab

ROBOTICS &

PERCEPTION

GROUP

ETH zürich University of ZurichTM

Perspective effects

Far away objects appear smaller



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart





T

Perspective effects





Projective Geometry

What is lost?

- Length
- Angles



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter



Projective Geometry

What is preserved?

ETH zürich

• Straight lines are still straight

University of Zurich^{™™}



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter



Vanishing points and lines

University of Zurich[™]

ETH zürich

Parallel lines in the world intersect in the image at a "vanishing point"







Vanishing points and lines



ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Perspective and art

ETH zürich

University of Zurich[™]

- Use of correct perspective projection indicated in 1st century B.C. frescoes
- Skill resurfaces in Renaissance: artists develop systematic methods to determine perspective projection (around 1480-1515)



Raphael



Durer, 1525

Playing with Perspective

University of Zurich[™]

ETHzürich

- Perspective gives us very strong depth cues
 ⇒ hence we can perceive a 3D scene by viewing its 2D representation (i.e. image)
- An example where perception of 3D scenes is misleading:



"Ames room"

A clip from "The computer that ate Hollywood" documentary. Dr. Vilayanur S. Ramachandran. **ETH**zürich **University of** Zurich^{™™}



Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm

CS4733 Class Notes, Stereo Imaging



Perspective Projection

Figure 1: Perspective imaging geometry showing relationship between 3D points and image plane points.

1 Stereo Imaging: Camera Model and Perspective Transform

We typically use a pinhole camera model that maps points in a 3-D camera frame to a 2-D projected image frame. In figure 1, we have a 3D camera coordinate frame X_c, Y_c, Z_c with origin O_c , and an image coordinate frame X_i, Y_i, Z_i with origin O_i . The focal length is f. Using similar triangles, we can relate image plane and world space coordinates. We have a 3D point P = (X, Y, Z) which projects onto the image plane at P' = (x, y, f). O_c is the origin of the camera coordinate system, known as the *center of projection* (COP) of the camera.

Using similar triangles, we can write down the following relationships:

$$\frac{X}{x} = \frac{Z}{f} \quad ; \quad \frac{Y}{y} = \frac{Z}{f} \quad ; \quad x = f \cdot \frac{X}{Z} \quad ; \quad y = f \cdot \frac{Y}{Z}$$

If f = 1, note that perspective projection is just scaling a world coordinate by its Z value. Also note that all 3D points along a line from the COP through a designated position (x, y) on the image plane will have the same image plane coordinates.

We can also describe perspective projection by the matrix equation:

$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \stackrel{\triangle}{=} \begin{bmatrix} s \cdot x \\ s \cdot y \\ s \end{bmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

where s is a scaling factor and $[x, y, 1]^T$ are the projected coordinates in the image plane.

We can generate *image space* coordinates from the projected camera space coordinates. These are the actual pixels values that you use in image processing. Pixels values (u, v) are derived by scaling the camera image plane coordinates in the x and y directions (for example, converting mm to pixels), and adding a translation to the origin of the image space plane. We can call these scale factors D_x and D_y , and the translation to the origin of the image plane as (u_0, v_0) .

If the pixel coordinates of a projected point (x,y) are (u,v) then we can write:

$$\frac{x}{D_x} = u - u_0; \quad \frac{y}{D_y} = v - v_0;$$
$$u = u_0 + \frac{x}{D_x}; \quad v = v_0 + \frac{y}{D_y}$$

where D_x, D_y are the physical dimensions of a pixel and (u_0, v_0) is the origin of the pixel coordinate system. $\frac{x}{D_x}$ and $\frac{y}{D_y}$ are simply the number of pixels, and we center them at the pixel coordinate origin. We can also put this into matrix form as:

$$\begin{bmatrix} s \cdot u \\ s \cdot v \\ s \end{bmatrix} = \begin{bmatrix} \frac{1}{D_x} & 0 & u_0 \\ 0 & \frac{1}{D_y} & v_0 \\ 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} s \cdot x \\ s \cdot y \\ s \end{bmatrix}$$
$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} \stackrel{\triangle}{=} \begin{bmatrix} s \cdot u \\ s \cdot v \\ s \end{bmatrix} = \begin{bmatrix} \frac{1}{D_x} & 0 & u_0 \\ 0 & \frac{1}{D_y} & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \cdot \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

 $P^{image} = T^{image}_{persp}T^{persp}_{camera}P^{camera}$

In the above, we assumed that the point to be imaged was in the camera coordinate system. If the point is in a previously defined world coordinate system, then we also have to add in a standard 4x4 transform to express the world coordinate point in camera coordinates:

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} s \cdot u \\ s \cdot v \\ s \end{bmatrix} = \begin{bmatrix} \frac{1}{D_x} & 0 & u_0 \\ 0 & \frac{1}{D_y} & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} wX \\ wY \\ wZ \\ 1 \end{bmatrix}$$

$$P^{image} = T^{image}_{persp} T^{persp}_{camera} T^{camera}_{world} P^{world}$$

Summing all this up, we can see that we need to find the following information to transform an arbitrary 3D world point to a designated pixel in a computer image:

- 6 parameters that relate the 3D world point to the 3D camera coordinate system (standard 3 translation and 3 rotation): (R, T)
- Focal Length of the camera: f
- Scaling factors in the x and y directions on the image plane: (D_x, D_y)
- Translation to the origin of the image plane: (u_0, v_0) .

This is 11 parameters in all. We can break these parameters down into *Extrinsic* parameters which are the 6-DOF transform between the camera coordinate system and the world coordinate system, and the *Intrinsic* parameters which are unique to the actual camera being used, and include the focal length, scaling factors, and location of the origin of the pixel coordinate system.

2 Camera Calibration

Camera calibration is used to find the mapping from 3D to 2D image space coordinates. There are 2 approaches:

- Method I: Find both extrinsic and intrinsic parameters of the camera system. However, this can be difficult to do. The instinsic parameters of the camera may be unknown (i.e. focal length, pixel dimension) and the 6-DOF transform also may be difficult to calculate directly.
- Method 2: An easier method is the "Lumped" transform. Rather than finding individual parameters, we find a composite matrix that relates 3D to 2D. Given the equation below:

$$P^{image} \ = \ T^{image}_{persp} T^{persp}_{camera} T^{camera}_{world} P^{world}$$

we can lump the 3 T matrices into a 3x4 calibration matrix C:

$$P^{image} = C P^{world}$$
$$C = T^{image}_{persp} T^{camera}_{camera} T^{camera}_{world}$$

• C is a single 3×4 transform that we can calculate empirically.



• Multiplying out the equations, we get:

$$c_{11}x + c_{12}y + c_{13}z + c_{14} = u$$

$$c_{21}x + c_{22}y + c_{23}z + c_{24} = v$$

$$c_{31}x + c_{32}y + c_{33}z + c_{34} = w$$

- Substituting u = u'w and v = v'w, we get:
 - 1. $c_{11}x + c_{12}y + c_{13}z + c_{14} = u'(c_{31}x + c_{32}y + c_{33}z + c_{34})$
 - 2. $c_{21}x + c_{22}y + c_{23}z + c_{24} = v'(c_{31}x + c_{32}y + c_{33}z + c_{34})$
- How to interpret <u>1</u> and <u>2</u>:
 - 1. If we know all the c_{ij} and x, y, z, we can find u', v'. This means that if we know calibration matrix C and a 3-D point, we can predict its image space coordinates.
 - 2. If we know x, y, z, u', v', we can find c_{ij} . Each 5-tuple gives 2 equations in c_{ij} . This is the basis for empirically finding the calibration matrix C (more on this later).
 - 3. If we know c_{ij} , u', v', we have 2 equations in x, y, z. They are the equations of 2 planes in 3-D. 2 planes form an intersection which is a line. These are the equations of the line emanating from the center of projection of the camera, through the image pixel location u', v' and which contains point x, y, z.

• We can set up a linear system to solve for c_{ij} : AC = B

x_1	y_1	z_1	1	0	0	0	0	$-u_1'x$	$-u_1'y$	$-u_1'z$		c_{11}		[u'_1	
0	0	0	0	x_1	y_1	z_1	1	$-v_1'x$	$-v_1'y$	$-v_1'z$		c_{12}			v_1'	
x_2	y_2	z_2	1	0	0	0	0	$-u_2'x$	$-u_2'y$	$-u_2'z$		C_{13}			u_2'	
0	0	0	0	x_2	y_2	z_2	1	$-v_2'x$	$-v_2'y$	$-v_2'z$		c_{14}			v_2'	
•												c_{21}			u'_3	
•												c_{22}		=	v'_3	
•												c_{23}			•	
•												C_{24}			•	
•												c_{31}			•	
•												c_{32}			u'_N	
•										-] [c_{33}	ļ		v'_N	
											We can	assur	ne $c_{34}=1$	L		

- Each set of points x, y, z, u', v' yields 2 equations in <u>11</u> unknowns (the c_{ij} 's).
- To solve for C, A needs to be invertible (square). We can <u>overdetermine</u> A and find a Least-Squares fit for C by using a pseudo-inverse solution.

If A is $N \times 11$, where N > 11,

$$AC = B$$

$$A^{T}AC = A^{T}B$$

$$C = \underbrace{(A^{T}A)^{-1}}_{\text{pseudo inverse}} A^{T}B$$

3 COMPUTATIONAL STEREO

Stereopsis is an identified human vision process. It is a passive, simple procedure that is robust to changes in lighting, scale, etc. Humans can fuse random dot stereograms that contain no high-level information about the objects in the fused images, yet they can infer depth from these stereograms. The procedure is:

- Camera-Modeling/Image-acquisition
- Feature extraction identify edges, corners, regions etc.
- Matching/Correspondence find same feature in both images
- Compute depth from matches use calibration information to back project rays from each camera and intersect them (triangulation)

• Interpolate surfaces - Matches are sparse, and constraints such as smoothness of surfaces are needed to "fill in" the depth between match points.

Camera Modeling: An important consideration in computational stereo is the setup of the cameras. The **baseline** between the camera centers determines the accuracy of the triangulation. Large baseline means more accuracy; however as the baseline gets larger, the same physical event in each image may not be found.

The cameras also have to be calibrated and registered. Calibration is relatively straightforward, and a variety of methods exist. Some methods extend the simple least squares model we discussed to include non-linear effects of lens distortion (particularly true with short a focal length lens).

Registration is needed to make use of the epipolar constraint. This constraint consists of a plane that includes both camera's optical centers and a point in 3-D space. This **epilolar plane** intersects both image planes in a straight line.

Feature Extraction: Identifying features in each image that can be matched is an important part of the stereo process. It serves 2 purposes: 1) data reduction so we are not forced to deal with every single pixel as a potential match, and 2) stability - features are seen to be more stable than a single gray level pixel.

There are 2 approaches: feature-based methods which find primitives such as edges, corners, lines, arcs in each image and match them; and area-based methods that identify regions or areas of pixels that can be matched using correlation based methods. Sometimes both methods are used, with feature-based methods proposing a match and area-based methods centered on the feature used to verify it.

Correspondence: The heart of the stereo problem is a search procedure. Given a pixel in image 1, it can potentially match each of N^2 pixels in the other image. To cut down this search space, cameras are often registered along scan lines. This means that he epipolar plane intersects each image plane along the same scan line. A pixel in image 1 can now potentially match only a pixel along the corresponding scan line in image 2, reducing the search from $O(N^2)$ to O(N). The match criteria can include not only the location of a feature like an edge, but also the edge direction and polarity.

Problems in Matching: A number of problems occur during matching to create false matches: These are occlusions, periodic features such as texture, homogeneous regions without features, baseline separation errors, and misregistered images. Stereo can usually only provide sparse 3-D data at easily identified feature points.

The camera | perspective camera

 For convenience, the image plane is usually represented in front of C such that the image preserves the same orientation (i.e. not flipped)



ETH zürich University of Zurich

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Perspective projection | from scene points to pixels

- The Camera point $\mathbf{P}_{\mathbf{C}} = (X_C, 0, Z_C)^{\mathrm{T}}$ projects to $\mathbf{p} = (x, y)$ onto the image plane
- From similar triangles:

$$\frac{x}{f} = \frac{X_c}{Z_c} \Longrightarrow x = \frac{fX_c}{Z_c}$$

• Similarly, in the general case:

$$\frac{y}{f} = \frac{Y_c}{Z_c} \Longrightarrow y = \frac{fY_c}{Z_c}$$



Perspective projection from scene points to pixels

- To convert p, from the local image plane coordinates (x, y) to the pixel coordinates (u, v), we need to account for:
 - The pixel coordinates of the camera optical center $0 = (u_0, v_0)$
 - Scale factor k for the pixel-size

University of Zurich[™]

ETH zürich

$$u = u_0 + kx \Rightarrow u_0 + k \frac{f X_C}{Z_C}$$
$$v = v_0 + ky \Rightarrow v_0 + k \frac{f Y_C}{Z_C}$$

 Use Homogeneous Coordinates for linear mapping from 3D to 2D, by introducing an extra element (scale):

p =

$$\begin{pmatrix} u \\ v \end{pmatrix} \qquad \qquad \widetilde{p} = \begin{bmatrix} u \\ \widetilde{v} \\ \widetilde{w} \end{bmatrix} = \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$



| 31

ETH zürich () University of Zurich



Perspective projection | from scene points to pixels

$$u = u_0 + kx \Rightarrow u_0 + k \frac{f_{C}}{z_C}$$
$$v = v_0 + ky \Rightarrow v_0 + k \frac{f_{C}}{z_C}$$

• Expressed in matrix form and homogeneous coordinates:

$$\begin{bmatrix} \lambda u \\ \lambda v \\ \lambda \end{bmatrix} = \begin{bmatrix} kf & 0 & u_0 \\ 0 & kf & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} X_c \\ Y_c \\ Z_c \end{bmatrix}$$

• Or alternatively





university of Zurich[™]

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Perspective projection | from scene points to pixels

$$\begin{bmatrix} X_{c} \\ Y_{c} \\ Z_{c} \end{bmatrix} = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \begin{bmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{bmatrix} + \begin{bmatrix} t_{1} \\ t_{2} \\ t_{3} \end{bmatrix} = \begin{bmatrix} R & | T \end{bmatrix} \cdot \begin{bmatrix} X_{w} \\ Y_{w} \\ Z_{w} \\ 1 \end{bmatrix}$$

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} X_{c} \\ Y_{c} \\ Z_{c} \end{bmatrix}$$

$$P_{c} = P_{v}$$

Perspective Projection_Matrix





ETHzürich **University of** ZurichTH



Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision




Perspective projection | radial distortion



No distortion

Barrel distortion

Pincushion

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ETH zürich () University of Zurich



Perspective projection | radial distortion

- The standard model of radial distortion is a transformation from the ideal coordinates (u, v) (i.e., undistorted) to the real observable coordinates (distorted) (u_d, v_d)
- The amount of distortion of the coordinates of the observed image is a nonlinear function of their radial distance. For most lenses, a simple quadratic model of distortion produces good results

where

$$\begin{bmatrix} u_d \\ v_d \end{bmatrix} = (1+k_1r^2) \begin{bmatrix} u-u_0 \\ v-v_0 \end{bmatrix} + \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}$$

$$r^{2} = (u - u_{0})^{2} + (v - v_{0})^{2}$$

ETH zürich University of Zurich

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Summary: Perspective projection equations

• To recap, a 3D world point $P = (X_w, Y_w, Z_w)$ projects into the image point p = (u, v)

$$\lambda p = \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} R \mid T \end{bmatrix} \cdot \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} \quad \text{where} \quad K = \begin{bmatrix} \alpha & 0 & u_0 \\ 0 & \alpha & v_0 \\ 0 & 0 & 1 \end{bmatrix}$$

and λ is the depth ($\lambda = Z_w$) of the scene point

• If we want to take into account for the radial distortion, then the distorted coordinates (u_d, v_d) (in pixels) can be obtained as

$$\begin{bmatrix} u_{d} \\ v_{d} \end{bmatrix} = (1 + k_{1}r^{2}) \begin{bmatrix} u - u_{0} \\ v - v_{0} \end{bmatrix} + \begin{bmatrix} u_{0} \\ v_{0} \end{bmatrix}$$

where $r^{2} = (u - u_{0})^{2} + (v - v_{0})^{2}$

Autonomous Mobile Robots

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ETHzürich **University of** Zurich^{™™}



Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision



Camera Calibration

TH zürich

University of Zurich[™]

Procedure to determine the *intrinsic parameters* of a camera



Camera Calibration

ETHzürich

University of Zurich[™]

- Use camera model to interpret the projection from world to image plane
- Using known correspondences of $p \Leftrightarrow P$, we can compute the unknown parameters K, R, T by applying the perspective projection equation Projection Matrix X_w
- ... so associate known, physical distances in the world to pixel-distances in image





ROBOTICS &

GROUE

PERCEPTION Autonomous Systems Lab



ETH zürich () University of Zurich



Camera Calibration (Direct Linear Transform (DLT) algorithm)

• We know that :
$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K \begin{bmatrix} R | T \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

• So there are 11 values to estimate: $\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$
e.g. m_{34} could be set to 1)

Each observed point gives us a pair of equations:

$$u_{i} = \frac{\lambda u_{i}}{\lambda} = \frac{m_{11}X_{i} + m_{12}Y_{i} + m_{13}Z_{i} + m_{14}}{m_{31} + m_{32} + m_{33} + m_{34}}$$
$$v_{i} = \frac{\lambda v_{i}}{\lambda} = \frac{m_{21}X_{i} + m_{22}Y_{i} + m_{23}Z_{i} + m_{24}}{m_{31} + m_{32} + m_{33} + m_{34}}$$

To estimate 11 unknowns, we need at least [?]; points to calibrate the camera ⇒ solved using linear least squares

ETH zürich () University of Zurich



Camera Calibration (Direct Linear Transform (DLT) algorithm)

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix} = K[R \mid T] \begin{bmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{bmatrix}$$

what we obtained: the 3x4 projection matrix,
 what we need: its decomposition into the camera calibration matrix *K*, and the rotation *R* and position *T* of the camera.

- Use QR factorization to decompose the 3x3 submatrix ($m_{11:33}$) into the product of an upper triangular matrix **K** and a rotation matrix **R** (orthogonal matrix)
- The translation *T* can subsequently be obtained by:

$$T = K^{-1} \begin{bmatrix} m_{14} \\ m_{24} \\ m_{34} \end{bmatrix}$$





DLT algorithm applied to multi-robot mutual localization



In this case, the camera has been pre-calibrated (i.e., K is known). Can you think of how the DLT algorithm could be modified so that only R and T need to determined and not K?

Autonomous Mobile Robots

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ETHzürich **University of** Zurich^{™™}



Outline of this lecture

- Perspective camera model
- Lens distortion
- Camera calibration
 - DLT algorithm
- Stereo vision

ETH zürich () University of Zurich



Stereo Vision versus Structure from Motion

• Stereo vision:

is the process of obtaining **depth information** from a pair of images coming from two cameras that look at the same scene from different but **known** positions

• Structure from Motion:

is the process of obtaining **depth and motion information** from a pair of images coming from the same camera that looks at the same scene from different positions

ROBOTICS & University of Zurich[™] **ETH** zürich PERCEPTION Autonomous Systems Lab GROUP

Depth from Stereo

Autonomous Mobile Robots

- From a single camera, we can only deduct the ray on which each image point lies
- With a stereo camera (binocular), we can solve for the intersection of the rays and recover the 3D structure



ETH zürich () University of Zurich

The "human" binocular system

- Stereopsys: the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial disotion is also removed. This process is called «rectification»



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ROBOTICS &

PERCEPTION Autonomous Systems Lab

ETH zürich

University of Zurich[™]

ROBOTICS & PERCEPTION Autonomous Systems Lab

The "human" binocular system

- **Stereopsys:** the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial disotion is also removed. This process is called «rectification»





Make a simple test:

- 1. Fix an object
- 2. Open and close alternatively the left and right eyes.
- The horizontal displacement is called **disparity** •
- The smaller the disparity, the farther the object •

ETHzürich

University of Zurich[™]

ROBOTICS & PERCEPTION Autonomous Systems Lab

The "human" binocular system

- **Stereopsys:** the brain allows us to see the left and right retinal images as a single 3D image
- The images project on our retina up-side-down but our brains lets us perceive them as «straight». Radial disotion is also removed. This process is called «rectification»





Make a simple test:

- 1. Fix an object
- 2. Open and close alternatively the left and right eyes.
- The horizontal displacement is called **disparity** •
- The smaller the disparity, the farther the object •

University of Zurich[™] **ETH** zürich

Stereo Vision | simplified case

An ideal, simplified case assumes that both cameras are **identical** and aligned with the x-axis

 $Z_P = \frac{DJ}{u_l - u_r}$

- Can we find an expression for the depth Z_P of point P_W ?
- From similar triangles:
 - $\frac{f}{Z_{P}} = \frac{u_{l}}{X_{P}}$ $\frac{f}{Z_{P}} = \frac{-u_{r}}{b X_{P}}$ **Disparity Disparity** is the difference in image location of the projection of a 3D point in two image planes
- **Baseline** is the distance between the two cameras





⁴⁰ Stereo Vision - The simplified case

 The simplified case is an ideal case. It assumes that both cameras are identical and are aligned on a horizontal axis



4c

From Similar Triangles:



Disparity difference in image location of the projection of a 3D point in two image planes

Baseline distance between the optical centers of the two cameras

© R. Siegwart , D. Scaramuzza and M.Chli, ETH Zurich - ASL

^{4c} ⁴¹ Stereo Vision facts

$$Z_P = \frac{bf}{u_l - u_r}$$

- 1. Depth is inversely proportional to disparity $(u_l u_r)$
 - Foreground objects have bigger disparity than background objects
- 2. Disparity is proportional to stereo-baseline b
 - The smaller the baseline b the more uncertain our estimate of depth
 - However, as b is increased, some objects may appear in one camera, but not in the other (remember both cameras have parallel optical axes)
- **3.** The projections of a single 3D point onto the left and the right stereo images are called **'correspondence pair'**

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Stereo Vision | general case

University of Zurich[™]

ETH zürich

- Two identical cameras do not exist in nature!
- Aligning both cameras on a horizontal axis is very difficult
- In order to use a stereo camera, we need to know the intrinsic extrinsic parameters of each camera, that is, the relative pose between the cameras (rotation, translation) ⇒ We can solve for this through camera calibration





Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

Auto



Stereo Vision | general case

University of Zurich[™]

ETH zürich

- To estimate the 3D position of P_W we can construct the system of equations of the left and right camera
- Triangulation is the problem of determining the 3D position of a point given a set of corresponding image locations and known camera poses.





Correspondence Search | the problem

University of Zurich^{uz#}

ETH zürich

- Goal: identify corresponding points in the left and right images, which are the reprojection of the same 3D scene point
 - Typical similarity measures: Normalized Cross-Correlation (NCC), Sum of Squared Differences (SSD), Sum of Absolute Differences (SAD), Census Transform
 - Exhaustive image search can be computationally very expensive! Can we make the correspondence search in 1D?





ETH zürich

University of Zurich^{™™}

ROBOTICS & PERCEPTION Autonomous Systems Lab

Correspondence Problem

- Exhaustive image search can be computationally very expensive!
- Can we make the correspondence search in 1D?
- Potential matches for p have to lie on the corresponding epipolar line l'
 - The **epipolar line** is the projection of the infinite ray $\pi^{-1}(p)$ corresponding to p in the other camera • image
 - The **epipole** e' is the projection of the optical center in the other camera image ٠





Correspondence Search | the epipolar constraint

- The epipolar plane is defined by the image point **p** and the optical centers
- Impose the epipolar constraint to aid matching: search for a correspondence along the epipolar line



ROBOTICS &

PERCEPTION Autonomous Systems Lab



Correspondence Search | the epipolar constraint

Thanks to the epipolar constraint, corresponding points can be searched for, along epipolar lines ⇒ computational cost reduced to 1 dimension!



ROBOTICS &

PERCEPTION Autonomous Systems Lab



Example: converging cameras

- Remember: all the epipolar lines intersect at the epipole
- As the position of the 3D point varies, the epipolar lines "rotate" about the baseline



Left image

Right image

Autonomous Mobile Robots

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

University of Zurich



Example: horizontally aligned cameras





Left image



Right image

ETH zürich University of Zurich

ROBOTICS & ASL PERCEPTION GROUP

Example: forward motion (parallel to the optical axis)

- Epipole has the **same coordinates** in both images
- Points move along lines radiating from e: "Focus of expansion"





Left image



Right image

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Stereo Rectification

ETH zürich

University of Zurich[™]

- Even in commercial stereo cameras the left and right image are never perfectly aligned
- In practice, it is convenient if image scanlines are the epipolar lines
- Stereo rectification warps the left and right images into new "rectified" images, whose epipolar lines are aligned to the baseline

ROBOTICS & University of Zurich[™] ETHzürich PERCEPTION Autonomous Systems Lab **Stereo Rectification** Reprojects image planes onto a common plane parallel to the baseline It works by computing two homographies (image warping), one for each input image reprojection As a result, the new epipolar lines are horizontal and the scanlines of the left and right image are aligned



Epipolar Rectification - Example

• First, remove radial distortion

University of Zurich^{™™}

ETH zürich





Epipolar Rectification - Example

First, remove radial distortion

University of Zurich[™]

ETH zürich

Then, compute homographies (warping) and rectify







Stereo Rectification: example



Autonomous Mobile Margarita Chli, Paul F

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Stereo Vision | disparity map

University of Zurich[™]

ETH zürich

- The disparity map holds the disparity value at every pixel:
 - Identify correspondent points of all image pixels in the original images
 - Compute the disparity $(u_l u_r)$ for each pair of correspondences
- Usually visualized in gray-scale images
- Close objects experience bigger disparity; thus, they appear brighter in disparity map



Left image

Right image





ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Stereo Vision | disparity map

University of Zurich^{uz#}

ETH zürich

- The disparity map holds the disparity value at every pixel:
 - Identify correspondent points of all image pixels in the original images
 - Compute the disparity $(u_l u_r)$ for each pair of correspondences
- Usually visualized in gray-scale images
- Close objects experience bigger disparity; thus, they appear brighter in disparity map
- From the disparity, we can compute the depth *Z* as:

$$Z = \frac{bf}{u_1 - u_r}$$

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart





66



- 1. Stereo camera calibration ⇒ compute camera relative pose
- 2. Epipolar rectification ⇒ align images & epipolar lines
- 3. Search for correspondences
- 4. Output: compute stereo triangulation or disparity map



Correspondence problem

University of Zurich^{uz∺}

ETH zürich

 Now that the left and right images are rectified, the correspondence search can be done along the same scanlines




Correspondence problem

University of Zurich^{uz∺}

If we look at the intensity profiles of two corresponding scanlines, there is a clear correspondence between intensities but also noise and ambiguities









250

200

1.00

Autonomous Mobile Robots Margarita Chli, Paul Furgale, I

ETH zürich

ETH zürich **University of** ZurichTH ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Correspondence problem

- To average noise effects, use a window around the point of interest
- Neighborhood of corresponding points are similar in intensity patterns
- Similarity measures:
 - Zero-Normalized Cross-Correlation (ZNCC)
 - Sum of Squared Differences (SSD),
 - Sum of Squared Differences (SAD)
 - **Census Transform** (Census descriptor plus Hamming distance)





epipolar line

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter **ETH**zürich **University of** ZurichTH ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Correlation-based window matching



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart



Correspondence Problems: Textureless regions (the aperture problem)



ROBOTICS &

PERCEPTION Autonomous Systems Lab

Autono Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart





Solution: increase window size



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart



Effects of window size W

University of Zurich^{™™}







W = 3

W = 20

- Smaller window
 - + More detail
 - More noise
- Larger window
 - + Smoother disparity maps
 - Less detail

Autonomous Mobile Robots

THzürich

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart



Textureless surfaces

University of Zurich^{™™}

THzürich

Occlusions, repetition

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart ROBOTICS & ASL PERCEPTION Autonomous Systems Lab GROUP University of Zurich[™]

How can we improve window-based matching?

- Beyond the epipolar constraint, there are "soft" constraints to help identify corresponding points
 - Uniqueness
 - Only one match in right image for every point in left image
 - Ordering

ETH zürich

- Points on same surface will be in same order in both views
- Disparity gradient
 - Disparity changes smoothly between points on the same surface

ROBOTICS &

PERCEPTION Autonomous Systems Lab



ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Results with window search

Data



Window-based matching



Ground truth





Better methods exist...

University of Zurich^{uz∺}



Graph cuts

Ground truth

Y. Boykov, O. Veksler, and R. Zabih, Fast Approximate Energy Minimization via Graph Cuts, PAMI 2001

For code, datasets, and comparisons all the algorithms: <u>http://vision.middlebury.edu/stereo/</u>

EHzürich

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

Sparse correspondence search

- Restrict search to sparse set of detected features
- Rather than pixel values (or lists of pixel values) use *feature descriptor* and an associated similarity metrics
- Still use epipolar geometry to narrow the search further



Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roiand Siegward

ETH zürich University of ZurichTH

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Template matching

- Find locations in an image that are similar to a *template*
- If we look at filters as **templates**, we can use correlation to detect these locations







Template



Template matching

ETH zürich

University of Zurich^{uz∺}

- Find locations in an image that are similar to a *template*
- If we look at filters as **templates**, we can use correlation to detect these locations



Detected template



Correlation map



Similarity measures

ETH zürich

Sum of Squared Differences (SSD)

University of Zurich^{uz∺}

$$SSD = \sum_{u=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - F(u,v))^{2}$$

Sum of Absolute Differences (SAD) (used in optical mice)

$$SAD = \sum_{u=-k}^{k} \sum_{v=-k}^{k} |H(u,v) - F(u,v)|$$

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

Similarity measures

ETHzürich

University of Zurich^{uz∺}

 For *slight* invariance to intensity changes, the Zero-mean Normalized Cross Correlation (ZNCC) is widely used

$$ZNCC = \frac{\sum_{u=-k\nu=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_{H}) (F(u,v) - \mu_{F})}{\sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_{H})^{2}} \sqrt{\sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} (F(u,v) - \mu_{F})^{2}}} \begin{cases} \mu_{H} = \frac{\sum_{u=-k\nu=-k}^{k} \sum_{v=-k}^{k} H(u,v)}{(2N+1)^{2}} \\ \mu_{F} = \frac{\sum_{u=-k\nu=-k}^{k} \sum_{v=-k}^{k} F(u,v)}{(2N+1)^{2}} \end{cases}$$

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

University of Zurich^{uz#} EHzürich

Correlation as an inner product

Considering the filter H and the portion of the image F_x as vectors \Rightarrow their correlation is: H.

$$\langle H, F_x \rangle = \|H\| \|F_x\| \cos \theta$$

In **ZNCC** we consider the unit vectors of H and F_x , hence we measure their similarity based on the angle θ . Alternatively, ZNCC maximizes $cos\theta$

$$\cos\theta = \frac{\langle H, F_x \rangle}{\|H\|\|F_x\|} = \frac{\sum_{u=-kv=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_H)(F(u,v) - \mu_F)}{\sqrt{\sum_{u=-kv=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_H)^2} \sqrt{\sum_{u=-kv=-k}^{k} \sum_{v=-k}^{k} (F(u,v) - \mu_F)^2}}$$

Autonomous Mobile Robots

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ROBOTICS &

PERCEPTION Autonomous Systems Lab



Choosing the Baseline

University of Zurich^{uz∺}

- What's the optimal baseline?
 - Too small:

ETHzürich

- Large depth error
- Can you quantify the error as a function of the disparity?
- Too large:
 - Minimum measurable distance increases
 - Difficult search problem for close objects



Autonomous Mobile Robots Large DaseIIIIe Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart





- 1. Stereo camera calibration ⇒ compute camera relative pose
- 2. Epipolar rectification ⇒ align images & epipolar lines
- 3. Search for correspondences
- 4. Output: compute stereo triangulation or disparity map
- 5. Consider how baseline & image resolution affect accuracy of depth estimates

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

ETH zürich University of Zurich

ROBOTICS & ASL PERCEPTION Autonomous Systems Lab

SFM: Structure From Motion (watch video segment)

- Given image point correspondences, $\mathbf{x}_i \leftrightarrow \mathbf{x}'_i$, determine R and T
- Keep track of point trajectories over multiple views to reconstruct scene structure and motion







Multiple-view structure from motion



Image courtesy of Nader Salman

Autonomous Mobile Robots Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart





Multiple-view structure from motion

Results of Structure from motion from user images from flickr.com

[Seitz, Szeliski ICCV 2009]



Colosseum, Rome 2,106 images, 819,242 points

Autonomous Mobile Robots

Margarita Chli, Paul Furgale, Marco Hutter, Martin Rufli, Davide Scaramuzza, Roland Siegwart

San Marco square, Venice 14,079 images, 4,515,157 points