

## Computer Vision 2 WS 2018/19

### Part 12 – Visual Odometry 04.12.2018

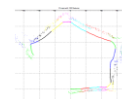
Prof. Dr. Bastian Leibe

RWTH Aachen University, Computer Vision Group  
<http://www.vision.rwth-aachen.de>



## Course Outline

- Single-Object Tracking
- Bayesian Filtering
- Multi-Object Tracking
  - Introduction
  - MHT, (JPDAF)
  - Network Flow Optimization
- Visual Odometry
  - Sparse interest-point based methods
  - Dense direct methods
- Visual SLAM & 3D Reconstruction
- Deep Learning for Video Analysis



2

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II



image source: [Zhang, Li, Nevatia, CVPR08]

## Topics of This Lecture

- Visual Odometry
  - Definition, Motivation
- Geometry Background
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- Point-based Visual Odometry
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

3

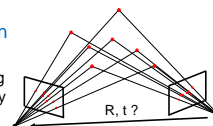
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II



## Recap: What is Visual Odometry ?

### Visual odometry (VO)...

- ... is a variant of **tracking**
  - Track motion (position and orientation) of the camera from its images
  - Only considers a limited set of recent images for real-time constraints
- ... also involves a **data association** problem
  - Motion is estimated from corresponding interest points or pixels in images, or by correspondences towards a local 3D reconstruction



4

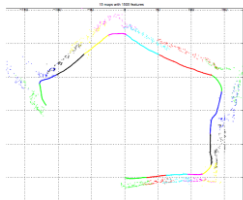
Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stuckler



## Recap: What is Visual Odometry ?

### Visual odometry (VO)...

- ... is prone to **drift** due to its local view
- ... is primarily concerned with estimating camera motion
  - Not all approaches estimate a 3D reconstruction of the associated interest points/ pixels explicitly.
  - If so it is **only locally consistent**



5

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stuckler



image source: [Clemente et al., RSS 2007]

## Visual Odometry Example

### SVO: Fast Semi-Direct Monocular Visual Odometry

Christian Forster, Matia Pizzoli, Davide Scaramuzza



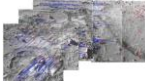
6

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stuckler



## Visual Odometry Term

- **Odometry**
  - Greek: „hodos“ – path, „metron“ – measurement
  - Motion or position estimation from measurements or controls
  - Typical example: wheel encoders
- **Visual Odometry**
  - 1980-2004: Prominent research by NASA JPL for Mars exploration rovers (Spirit and Opportunity in 2004)
  - David Nistér's „Visual Odometry“ paper from 2004 about keypoint-based methods for monocular and stereo cameras



7

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Jörg Stückler



RWTH AACHEN UNIVERSITY  
Image source: NASA (Chen et al., RAM, 2006)

## Why Visual Odometry?

- VO is often used to complement other motion sensors
  - GPS
  - Inertial Measurement Units (IMUs)
  - Wheel odometry
  - etc.
- VO is much more accurate than wheel odometry and not prone to wheel slippage.
- VO is important in GPS-denied environments (indoors, close to buildings, etc.)

8

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Jörg Stückler



RWTH AACHEN UNIVERSITY

## Sensor Types for Visual Odometry

- **Monocular cameras**
  - Pros: Low-power, light-weight, low-cost, simple to calibrate and use
  - Cons: requires motion parallax and texture, scale not observable
- **Stereo cameras**
  - Pros: depth without motion, less power than active structured light
  - Cons: requires texture, accuracy depends on baseline, synchronization and extrinsic calibration of the cameras
- **Active RGB-D sensors**
  - Pros: no texture needed, similar to stereo processing
  - Cons: active sensing consumes power, blackbox depth estimation



9

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Jörg Stückler



RWTH AACHEN UNIVERSITY  
Image source: IDS, PointGrey, ASUS

## Topics of This Lecture

- **Visual Odometry**
  - Definition, Motivation
- **Geometry Background**
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- **Point-based Visual Odometry**
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

10

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II



RWTH AACHEN UNIVERSITY

## A Note about Notation

- This course material originated from the 2016 CV 2 lecture held together with Jörg Stückler (now Prof. @ MPI Tübingen)
  - The notation follows the **MASKS** textbook and is slightly different from the notation used in the CV 1 lecture.
  - We'll stick with this notation in order to be consistent with the later lectures
  - *In case you get confused by notation, please interrupt me and ask...*



An Invitation to  
3D Vision,  
Y. Ma, S. Soatto,  
J. Kosecka, and  
S. S. Sastry,  
Springer, 2004

11

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II



RWTH AACHEN UNIVERSITY

## Geometric Point Primitives

- |                           | 2D  | 3D   |
|---------------------------|---|--|
| • Point                   | $\mathbf{x} = \begin{pmatrix} x \\ y \end{pmatrix} \in \mathbb{R}^2$                                      | $\mathbf{x} = \begin{pmatrix} x \\ y \\ z \end{pmatrix} \in \mathbb{R}^3$  |
| • Augmented vector        | $\bar{\mathbf{x}} = \begin{pmatrix} x \\ y \\ 1 \end{pmatrix} \in \mathbb{R}^3$                           | $\bar{\mathbf{x}} = \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \in \mathbb{R}^4$                                   |
| • Homogeneous coordinates | $\tilde{\mathbf{x}} = \begin{pmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{w} \end{pmatrix} \in \mathbb{P}^2$ | $\tilde{\mathbf{x}} = \begin{pmatrix} \tilde{x} \\ \tilde{y} \\ \tilde{z} \\ \tilde{w} \end{pmatrix} \in \mathbb{P}^3$ |
- $\tilde{\mathbf{x}} - \tilde{w}\bar{\mathbf{x}}$

12

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Jörg Stückler



RWTH AACHEN UNIVERSITY

### Euclidean Transformations

- Euclidean transformations apply rotation and translation
 
$$\mathbf{x}' = \mathbf{R}\mathbf{x} + \mathbf{t} \quad \bar{\mathbf{x}}' = \begin{pmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{pmatrix} \bar{\mathbf{x}}$$
- Rigid-body motion: preserves distances and angles when applied to points on a body

13 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### Special Orthogonal Group SO(n)

- Rotation matrices have a special structure
 
$$\mathbf{R} \in \text{SO}(n) \subset \mathbb{R}^{n \times n}, \det(\mathbf{R}) = 1, \mathbf{R}\mathbf{R}^T = \mathbf{I}$$
 i.e. orthonormal matrices that preserve distance and orientation
- They form a group denoted as Special Orthogonal Group  $\text{SO}(n)$ 
  - The group operator is matrix multiplication – associative, but not commutative!
  - Inverse and neutral element exist
- 2D rotations only have 1 degree of freedom (DoF), i.e. angle of rotation
- 3D rotations have 3 DoFs, several parametrizations exist such as Euler angles and quaternions

14 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### 3D Rotation Representations – Matrix

- Straight-forward: **Orthonormal matrix**

$$\mathbf{R} = \begin{pmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{pmatrix} \in \mathbb{R}^{3 \times 3}$$
- Pro:
  - Easy to concatenate and invert
 
$$\mathbf{R}_C^A = \mathbf{R}_B^A \mathbf{R}_C^B \quad \mathbf{R}_A^A = (\mathbf{R}_B^B)^{-1}$$
- Con:
  - Overparametrized (9 parameters for 3 DoF) – problematic for optimization

15 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### 3D Rotation Representations – Euler Angles

- Euler Angles:** 3 consecutive rotations around coordinate axes  
Example: roll-pitch-yaw angles  $\alpha, \beta, \gamma$  (X-Y-Z):
 
$$\mathbf{R}_{XYZ}(\alpha, \beta, \gamma) = \mathbf{R}_Z(\gamma) \mathbf{R}_Y(\beta) \mathbf{R}_X(\alpha)$$
 with
 
$$\mathbf{R}_X(\alpha) = \begin{pmatrix} 1 & 0 & 0 \\ 0 & \cos(\alpha) & -\sin(\alpha) \\ 0 & \sin(\alpha) & \cos(\alpha) \end{pmatrix}$$

$$\mathbf{R}_Y(\beta) = \begin{pmatrix} \cos(\beta) & 0 & \sin(\beta) \\ 0 & 1 & 0 \\ -\sin(\beta) & 0 & \cos(\beta) \end{pmatrix}$$

$$\mathbf{R}_Z(\gamma) = \begin{pmatrix} \cos(\gamma) & -\sin(\gamma) & 0 \\ \sin(\gamma) & \cos(\gamma) & 0 \\ 0 & 0 & 1 \end{pmatrix}$$
- 12 possible orderings of rotation axes (f.e. Z-X-Z)
- Pro: Minimal with 3 parameters
- Con: Singularities (gimbal lock), concatenation/inversion via conversion from/to matrix

16 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### 3D Rotation Representations – Axis-Angle

- Axis-Angle:** Rotate along axis  $\mathbf{n} \in \mathbb{R}^3$  by angle  $\theta \in \mathbb{R}$ :
 
$$\mathbf{R}(\mathbf{n}, \theta) = \mathbf{I} + \sin(\theta)\hat{\mathbf{n}} + (1 - \cos(\theta))\hat{\mathbf{n}}^2 \quad \|\mathbf{n}\|_2 = 1$$
 where  $\hat{\mathbf{x}} := \begin{pmatrix} 0 & -z & y \\ z & 0 & -x \\ -y & x & 0 \end{pmatrix} \quad \hat{\mathbf{x}}\mathbf{y} = \mathbf{x} \times \mathbf{y}$
- Reverse:  $\theta = \cos^{-1}\left(\frac{\text{tr}(\mathbf{R}) - 1}{2}\right) \quad \mathbf{n} = \frac{1}{2\sin(\theta)} \begin{pmatrix} r_{32} - r_{23} \\ r_{13} - r_{31} \\ r_{21} - r_{12} \end{pmatrix}$
- 4 parameters:  $\{\mathbf{n}, \theta\}$
- 3 parameters:  $\omega - \theta \mathbf{n}$
- Pro: minimal representation for 3 parameters
- Con:  $\{\mathbf{n}, \theta\}$  has unit norm constraint on  $\mathbf{n}$  - problematic for optimization; both parametrizations not unique; concatenation/inversion via  $\text{SO}(3)$

17 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### 3D Rotation Representations – Quaternions

- Unit Quaternions:**  $\mathbf{q} = (q_x, q_y, q_z, q_w)^T \in \mathbb{R}^4, \|\mathbf{q}\|_2 = 1$
- Relation to axis-angle representation:
  - Axis-angle to quaternion:
 
$$\mathbf{q}(\mathbf{n}, \theta) = \left( \mathbf{n}^T \sin\left(\frac{\theta}{2}\right), \cos\left(\frac{\theta}{2}\right) \right)$$
  - Quaternion to axis-angle:
 
$$\mathbf{n}(\mathbf{q}) = \begin{cases} (q_x, q_y, q_z)^T / \sin(\theta/2), & \theta \neq 0 \\ \mathbf{0}, & \theta = 0 \end{cases}$$

$$\theta = 2 \arccos(q_w)$$

18 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticicler

### 3D Rotation Representations – Quaternions cont.

- Pros:
  - Unique up to opposing sign  $q \quad -q$
  - Direct rotation of a point:
 
$$p' = Rp = q(R)pq(R)^{-1}$$
  - Direct concatenation of rotations:
 
$$q(R_2R_1) = q(R_2)q(R_1)$$
  - Direct inversion of a rotation:
 
$$q(R^{-1}) = q(R)^{-1}$$
- with  $q^{-1} = (-q_{xyz}^T, q_w)^T$ ,
- $q_1q_2 = (q_{1,w}q_{2,xyz} + q_{2,w}q_{1,xyz} + q_{1,xyz} \times q_{2,xyz}, q_{1,w}q_{2,w} - q_{1,xyz}q_{2,xyz})$
- Con: Normalization constraint is problematic for optimization

19 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

### Special Euclidean Group SE(3)

- Euclidean transformation matrices have a special structure as well:
 
$$T = \begin{pmatrix} R & t \\ 0 & 1 \end{pmatrix} \in SE(3) \subset \mathbb{R}^{4 \times 4}$$
- Translation  $t$  has 3 degrees of freedom
- Rotation  $R \in SO(3)$  has 3 degrees of freedom
- They also form a group which we call  $SE(3)$ . The group operator is matrix multiplication:
 
$$\cdot : SE(3) \times SE(3) \rightarrow SE(3)$$

$$T_B^A \cdot T_C^B \mapsto T_C^A$$

20 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

### Definition of Visual Odometry

- Visual odometry is the process of estimating the egomotion of an object using only inputs from visual sensors on the object
- Inputs: images at discrete time steps  $t$ ,
  - Monocular case: Set of images  $I_{0:t} = \{I_0, \dots, I_t\}$
  - Stereo case: Left/right images  $I_{0:t}^l = \{I_0^l, \dots, I_t^l\}$ ,  $I_{0:t}^r = \{I_0^r, \dots, I_t^r\}$
  - RGB-D case: Color/depth images  $I_{0:t} = \{I_0, \dots, I_t\}$ ,  $Z_{0:t} = \{Z_0, \dots, Z_t\}$
- Output: relative transformation estimates  $T_t^{t-1} \in SE(3)$  between frames
- Conventions:
  - Let  $T_t \in SE(3)$  be the camera pose at time  $t$  in the world frame
  - $T_t^{-1}$  transforms points from camera frame at time  $t$  to  $t-1$ , i.e.
 
$$T_t = T_0T_1^{-1} \dots T_t^{-1}$$

21 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

### Direct vs. Indirect Methods

- Direct methods
  - formulate alignment objective in terms of photometric error (e.g. intensities)
 
$$p(T_2 | T_1, \xi) \rightarrow E(\xi) = \int_{u \in \Omega} |I_1(u) - I_2(u, \xi)| du$$
- Indirect methods
  - formulate alignment objective in terms of reprojection error of geometric primitives (e.g. points, lines)
 
$$p(Y_2 | Y_1, \xi) \rightarrow E(\xi) = \sum_i |y_{1,i} - \omega(y_{2,i}, \xi)|$$


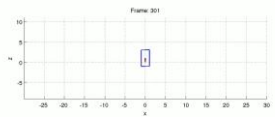
22 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

### Topics of This Lecture

- Visual Odometry
  - Definition, Motivation
- Geometry Background
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- Point-based Visual Odometry
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

23 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

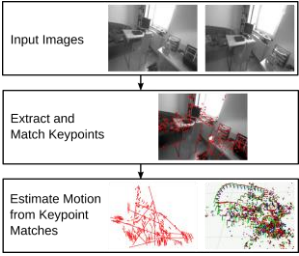
### Point-based (Indirect) Visual Odometry Example

24 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ivo Stukler

### Point-based Visual Odometry Pipeline

- Keypoint detection and local description (CV I)
- Robust keypoint matching (CV I)
- Motion estimation
  - 2D-to-2D: motion from 2D point correspondences
  - 2D-to-3D: motion from 2D points to local 3D map
  - 3D-to-3D: motion from 3D point correspondences (e.g., stereo, RGB-D)



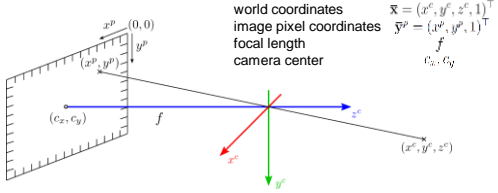
25 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Topics of This Lecture

- Visual Odometry
  - Definition, Motivation
- Geometry Background
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- Point-based Visual Odometry
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

26 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II

### Recap: Pinhole Projection Camera Model



world coordinates  $\bar{x} = (x^w, y^w, z^w, 1)^T$   
 image pixel coordinates  $\bar{y}^p = (x^p, y^p, 1)^T$   
 focal length  $f$   
 camera center  $(c_x, c_y)$

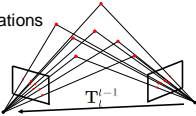
$$\begin{pmatrix} x^p \\ y^p \\ 1 \end{pmatrix} = \begin{pmatrix} f & 0 & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{pmatrix} \begin{pmatrix} x^w/z^w \\ y^w/z^w \\ 1 \end{pmatrix}$$

camera matrix  $C'$   $\Rightarrow \pi(\bar{x}) = \bar{y}$  (normalized image coordinates)

27 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

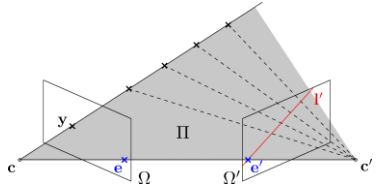
### 2D-to-2D Motion Estimation

- Given corresponding image point observations
  - $Y_t = \{y_{t,1}, \dots, y_{t,N}\}$
  - $Y_{t-1} = \{y_{t-1,1}, \dots, y_{t-1,N}\}$
 of unknown 3D points  $X = \{x_1, \dots, x_N\}$
- determine relative motion  $T_t^{t-1}$  between frames
- Obvious try: minimize reprojection error using least squares
 
$$E(T_t^{t-1}, X) = \sum_{i=1}^N \|\bar{y}_{t,i} - \pi(\bar{x}_i)\|_2^2 + \|\bar{y}_{t-1,i} - \pi(T_t^{t-1}\bar{x}_i)\|_2^2$$
- Convexity? Uniqueness (scale-ambiguity)?
- Alternative algebraic approaches: 8-point / 5-point algorithm



28 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

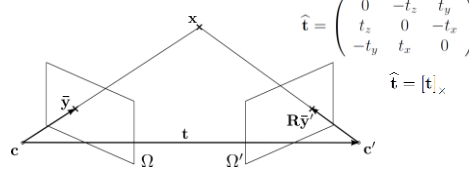
### Recap: Epipolar Geometry



- Camera centers  $c, c'$  and image point  $y \in \Omega$  span the epipolar plane  $\Pi$
- The ray from camera center  $c$  through point  $y$  projects as the epipolar line  $l'$  in image plane  $\Omega'$
- The intersections of the line through the camera centers with the image planes are called epipoles  $c, c'$

29 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Essential Matrix



$$\hat{t} = \begin{pmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{pmatrix}$$

$$\hat{t} = [t]_{\times}$$

- The rays to the 3D point and the baseline  $t$  are coplanar
 
$$\bar{y}^T (t \times R\bar{y}') = 0 \Leftrightarrow \bar{y}^T \hat{t} R\bar{y}' = 0$$
- The Essential matrix  $E := \hat{t}R$  captures the relative camera pose
- Each point correspondence provides an „epipolar constraint“
- 5 correspondences suffice to determine  $E$
- (Simpler: 8-point algorithm)

30 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

## Eight-Point Algorithm for Essential Matrix Estimation

- First proposed by Longuet and Higgins, 1981
- Algorithm:
  1. Rewrite epipolar constraints as a linear system of equations  
 $\tilde{y}_i E \tilde{y}_i' = a_i E e_i = 0 \rightarrow A E e = 0 \quad A = (a_1^T, \dots, a_N^T)^T$   
 using Kronecker product  $a_i = \tilde{y}_i \otimes \tilde{y}_i'$  and  $E e = (e_{11}, e_{12}, e_{13}, \dots, e_{33})^T$
  2. Apply singular value decomposition (SVD) on  $A = U_A S_A V_A^T$  and unstack the 9th column of  $V_A$  into  $\tilde{E}$
  3. Project the approximate  $\tilde{E}$  into the (normalized) essential space:  
 Determine the SVD of  $\tilde{E} = U \text{diag}(\sigma_1, \sigma_2, \sigma_3) V^T$  with  $U, V \in \text{SO}(3)$   
 and replace the singular values  $\sigma_1 \geq \sigma_2 \geq \sigma_3$  with  $1, 1, 0$  to find

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

31

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticchiè



RWTH AACHEN  
UNIVERSITY

## Eight-Point Algorithm cont.

- Algorithm (cont.):
  - Determine one of the following 2 possible solutions that intersects the points in front of both cameras:

$$R = UR_Z^T \left( \pm \frac{\pi}{2} \right) V^T \quad \hat{t} = UR_Z \left( \pm \frac{\pi}{2} \right) \text{diag}(1, 1, 0) U^T$$

- A derivation can be found in the MASKS textbook, Ch. 5

### Remarks

- Algebraic solution does not minimize geometric error
- Refine using non-linear least-squares of reprojection error
- Alternative: formulate epipolar constraints as „distance from epipolar line“ and minimize this non-linear least-squares problem

32

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticchiè



RWTH AACHEN  
UNIVERSITY

## Triangulation

- Goal: Reconstruct 3D point  $\tilde{x} = (x, y, z, w)^T \in \mathbb{P}^3$  from 2D image observations  $\{y_1, \dots, y_N\}$  for known camera poses  $\{T_1, \dots, T_N\}$

- Linear solution: Find 3D point such that reprojections equal its projections

$$y_i' = \pi(T_i, \tilde{x}) = \begin{pmatrix} r_{11}x + r_{12}y + r_{13}z + r_{14}w \\ r_{21}x + r_{22}y + r_{23}z + r_{24}w \\ r_{31}x + r_{32}y + r_{33}z + r_{34}w \end{pmatrix}$$

- Each image provides one constraint  $y_i - y_i' = 0$
- Leads to system of linear equations  $A \tilde{x} = 0$ , two approaches:
  - Set  $w = 1$  and solve nonhomogeneous system
  - Find nullspace of  $A$  using SVD (this is what we did in CV I)

- Non-linear solution: Minimize least squares reprojection error (more accurate)

$$\min_{\tilde{x}} \left\{ \sum_{i=1}^N \|y_i - y_i'\|_2^2 \right\}$$

33

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticchiè



RWTH AACHEN  
UNIVERSITY

## Relative Scale Recovery

### Problem:

- Each subsequent frame-pair gives another solution for the reconstruction scale

### Solution:

- Triangulate overlapping points  $Y_{t-2}, Y_{t-1}, Y_t$  for current and last frame pair

$$\Rightarrow X_{t-2,t-1}, X_{t-1,t}$$

- Rescale translation of current relative pose estimate to match the reconstruction scale with the distance ratio between corresponding point pairs

$$r_{i,j} = \frac{\|X_{t-2,t-1,i} - X_{t-2,t-1,j}\|_2}{\|X_{t-1,t,i} - X_{t-1,t,j}\|_2}$$

- Use mean or robust median over available pair ratios

34

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticchiè



RWTH AACHEN  
UNIVERSITY

## Algorithm: 2D-to-2D Visual Odometry

**Input:** image sequence  $I_{0:t}$

**Output:** aggregated camera poses  $T_{0,t}$

### Algorithm:

For each current image  $I_k$ :

1. Extract and match keypoints between  $I_{k-1}$  and  $I_k$
2. Compute relative pose  $T_{k-1}^{k-1}$  from essential matrix between  $I_{k-1}, I_k$
3. Compute relative scale and rescale translation of  $T_{k-1}^{k-1}$  accordingly
4. Aggregate camera pose by  $T_k = T_{k-1} T_{k-1}^{k-1}$

35

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II  
Slide credit: Ugo Sticchiè



RWTH AACHEN  
UNIVERSITY

## Topics of This Lecture

### Visual Odometry

- Definition, Motivation

### Geometry Background

- Euclidean Transformations
- 3D Rotation representations
- Definition of Visual Odometry
- Direct vs. Indirect methods

### Point-based Visual Odometry

- 2D-to-2D Motion Estimation
- 2D-to-3D Motion Estimation
- 3D-to-3D Motion Estimation
- Further Considerations

36

Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 – Multi-Object Tracking II

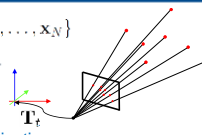


RWTH AACHEN  
UNIVERSITY

### 2D-to-3D Motion Estimation

- Given a local set of 3D points  $X = \{x_1, \dots, x_N\}$  and corresponding image observations  $Y_t = \{y_{t,1}, \dots, y_{t,N}\}$  determine camera pose  $T_t$  within the local map
- Minimize least squares **geometric reprojection error**

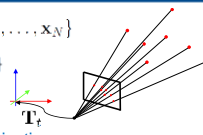
$$E(T_t) = \sum_{i=1}^N \|y_{t,i} - \pi(T_t x_i)\|_2^2$$
- Perspective-n-Points (PnP)** problem, many approaches exist, e.g.,
  - Direct linear transform (DLT)
  - EPnP [Lepetit et al., An accurate O(n) Solution to the PnP problem, IJCV 2009]
  - OPnP [Zheng et al., Revisiting the PnP Problem: A Fast, General and Optimal Solution, ICCV 2013]



37 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticchiè

### Direct Linear Transform for PnP

- Goal: determine projection matrix  $P = (R \ t) \in \mathbb{R}^{3 \times 4} = \begin{pmatrix} P_1 \\ P_2 \\ P_3 \end{pmatrix}$
- Each 2D-to-3D point correspondence 3D:  $\tilde{x}_i = (x_i, y_i, z_i, w_i)^T \in \mathbb{P}^3$  2D:  $\tilde{y}_i = (x'_i, y'_i, w'_i)^T \in \mathbb{P}^2$  gives two constraints
 
$$\begin{pmatrix} 0 & -w'_i \tilde{x}_i^T & y'_i \tilde{x}_i^T \\ w'_i \tilde{x}_i^T & 0 & -x'_i \tilde{x}_i^T \end{pmatrix} \begin{pmatrix} P_1^T \\ P_2^T \\ P_3^T \end{pmatrix} = 0$$
 through  $\tilde{y}_i \times (P \tilde{x}_i) = 0$
- Form linear system of equation  $A p = 0$  with  $p := \begin{pmatrix} P_1^T \\ P_2^T \\ P_3^T \end{pmatrix} \in \mathbb{R}^9$  from  $N \geq 6$  correspondences
- Solve for  $p$ : determine unit singular vector of  $A$  corresponding to its smallest eigenvalue

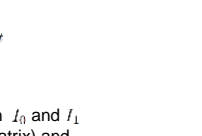


38 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticchiè

### Algorithm: 2D-to-3D Visual Odometry

**Input:** image sequence  $I_{0:t}$   
**Output:** aggregated camera poses  $T_{0:t}$

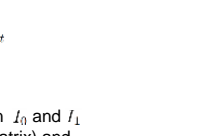
**Algorithm:**  
Initialize:  
1. Extract and match keypoints between  $I_0$  and  $I_1$   
2. Determine camera pose (essential matrix) and triangulate 3D keypoints  $X_1$   
For each current image  $I_k$ :  
1. Extract and match keypoints between  $I_{k-1}$  and  $I_k$   
2. Compute camera pose  $T_k$  using PnP from 2D-to-3D matches  
3. Triangulate all new keypoint matches between  $I_{k-1}$  and  $I_k$  and add them to the local map  $X_k$



39 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticchiè

### Topics of This Lecture

- Visual Odometry
  - Definition, Motivation
- Geometry Background
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- Point-based Visual Odometry
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

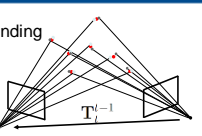


40 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II

### 3D-to-3D Motion Estimation

- Given 3D point coordinates of corresponding points in two camera frames
 
$$X_{t-1} = \{x_{t-1,1}, \dots, x_{t-1,N}\}$$

$$X_t = \{x_{t,1}, \dots, x_{t,N}\}$$
 determine relative camera pose  $T_t^{t-1}$
- Idea: determine rigid transformation that aligns the 3D points
- Geometric least squares error:  $E(T_t^{t-1}) = \sum_{i=1}^N \|\tilde{x}_{t-1,i} - T_t^{t-1} \tilde{x}_{t,i}\|_2^2$
- Closed-form solutions available, e.g., [Arun et al., 1987]
- Applicable, e.g., for calibrated stereo cameras (triangulation of 3D points) or RGB-D cameras (measured depth)

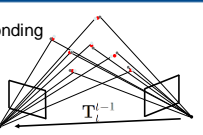


41 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticchiè

### 3D Rigid-Body Motion from 3D-to-3D Matches

- [Arun et al., Least-squares fitting of two 3-d point sets, IEEE PAMI, 1987]
- Corresponding 3D points,  $N \geq 3$ 

$$X_{t-1} = \{x_{t-1,1}, \dots, x_{t-1,N}\} \quad X_t = \{x_{t,1}, \dots, x_{t,N}\}$$
- Determine means of 3D point sets
 
$$\mu_{t-1} = \frac{1}{N} \sum_{i=1}^N x_{t-1,i} \quad \mu_t = \frac{1}{N} \sum_{i=1}^N x_{t,i}$$
- Determine rotation from
 
$$A = \sum_{i=1}^N (x_{t-1,i} - \mu_{t-1})(x_{t,i} - \mu_t)^T \quad A = USV^T \quad R_{t-1}^t = VU^T$$
- Determine translation as  $t_{t-1}^t = \mu_t - R_{t-1}^t \mu_{t-1}$



42 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticchiè

## Algorithm: 3D-to-3D Stereo Visual Odometry

**Input:** stereo image sequence  $I_{0:t}^l, I_{0:t}^r$   
**Output:** aggregated camera poses  $T_{0:t}$

### Algorithm:

For each current stereo image  $I_{k-1}^l, I_{k-1}^r$ :

1. Extract and match keypoints between  $I_{k-1}^l$  and  $I_{k-1}^r$
2. Triangulate 3D points  $X_k$  between  $I_{k-1}^l$  and  $I_{k-1}^r$
3. Compute camera pose  $T_k^{k-1}$  from 3D-to-3D point matches  $X_k$  to  $X_{k-1}$
4. Aggregate camera poses by  $T_k = T_{k-1} T_k^{k-1}$

43

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II  
 Slide credit: Ugo Sticker



RWTH AACHEN  
 UNIVERSITY

## Topics of This Lecture

- Visual Odometry
  - Definition, Motivation
- Geometry Background
  - Euclidean Transformations
  - 3D Rotation representations
  - Definition of Visual Odometry
  - Direct vs. Indirect methods
- Point-based Visual Odometry
  - 2D-to-2D Motion Estimation
  - 2D-to-3D Motion Estimation
  - 3D-to-3D Motion Estimation
  - Further Considerations

44

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II



RWTH AACHEN  
 UNIVERSITY

## Further Considerations

- How to detect keypoints?
- How to match keypoints?
- How to cope with outliers among keypoint matches?
- How to cope with noisy observations?
- When to create new 3D keypoints? Which keypoints to use?
- 2D-to-2D, 2D-to-3D or 3D-to-3D?
- Optimize over more than two frames?
- ...

45

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II  
 Slide credit: Ugo Sticker



RWTH AACHEN  
 UNIVERSITY

## Recap: Keypoint Detectors

- Corners
  - Image locations with locally prominent intensity variation
  - Intersections of edges
- Blobs
  - Image regions that stick out from their surrounding in intensity/texture
  - Circular high-contrast regions
- Examples: Harris, FAST
- Scale-selection: Harris-Laplace
- E.g.: LoG, DoG (SIFT), SURF
- Scale-space extrema in LoG/DoG



Harris Corners



DoG (SIFT) Blobs

46

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II  
 Slide credit: Ugo Sticker



RWTH AACHEN  
 UNIVERSITY

Image source: Svetlana Lazebnik

## Recap: Keypoint Detectors

- Desirable properties of keypoint detectors for VO:
  - High repeatability,
  - Localization accuracy,
  - Robustness,
  - Invariance,
  - Computational efficiency



Harris Corners



DoG (SIFT) Blobs

47

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II  
 Slide credit: Ugo Sticker



RWTH AACHEN  
 UNIVERSITY

Image source: Svetlana Lazebnik

## Recap: Keypoint Detectors

- Corners vs. blobs for visual odometry:
  - Typically corners provide higher spatial localization accuracy, but are less well localized in scale
  - Corners are typically detected in less distinctive local image regions
  - Highly run-time efficient corner detectors exist (e.g., FAST)



Harris Corners



DoG (SIFT) Blobs

48

Visual Computing Institute | Prof. Dr. Bastian Leibe  
 Computer Vision 2  
 Part 11 – Multi-Object Tracking II  
 Slide credit: Ugo Sticker

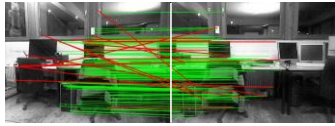


RWTH AACHEN  
 UNIVERSITY

Image source: Svetlana Lazebnik



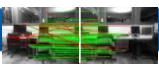
### Recap: Keypoint Matching



- Desirable properties for VO:
  - High recall,
  - Precision,
  - Robustness,
  - Computational efficiency

49 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

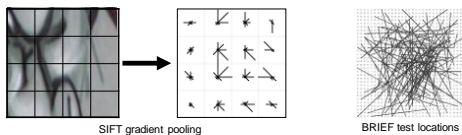
### Recap: Keypoint Matching



- Several data association principles:
  - Matching by **reprojection error / distance to epipolar line**
    - Assumes an initial guess for camera motion
    - (e.g., Kalman filter prediction, IMU, or wheel odometry)
  - Detect-then-track** (e.g., KLT-tracker):
    - Correspondence search by local image alignment
    - Assumes incremental small (but unknown) motion between images
  - Matching by **descriptor**:
    - Scale/viewpoint-invariant local descriptors allow for wider image baselines
- Robustness through **RANSAC** for motion estimation

50 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Recap: Local Feature Descriptors

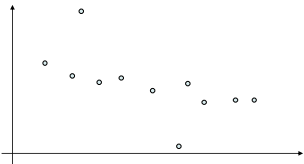


- Extract signatures that describe local image regions:
  - Histograms over image gradients (SIFT)
  - Histograms over Haar-wavelet responses (SURF)
  - Binary patterns (BRIEF, BRISK, FREAK, etc.)
  - Learning-based descriptors (e.g., Calonder et al., ECCV 2008)
- Rotation-invariance: Align with dominant orientation
- Scale-invariance: Adapt local region extent to keypoint scale

51 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Recap: RANSAC

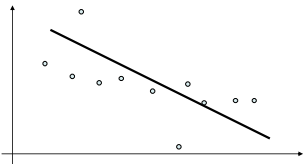
- Model fitting in presence of noise and **outliers**
- Example: fitting a line through 2D points



52 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Recap: RANSAC

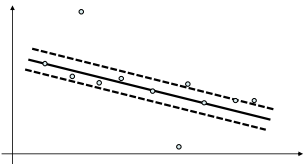
- Least-squares solution, assuming constant noise for all points



53 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Recap: RANSAC

- We only need 2 points to fit a line. Let's try 2 random points



54 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Ugo Sticker

### Recap: RANSAC

- Let's try 2 other random points

Quite bad..  
3 inliers  
8 outliers

55 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Igor Stukler

### Recap: RANSAC

- Let's try yet another 2 random points

Quite good!  
9 inliers  
2 outliers

56 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Igor Stukler

### Recap: RANSAC

- Let's use the inliers of the best trial to perform least squares fitting

Even better!

57 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Igor Stukler

### Recap: RANSAC

- R**ANdom **S**Ample **C**onsensus algorithm formalizes this idea
- Algorithm:**  
Input: data  $D$ ,  $s$  required data points for fitting, success probability  $p$ , outlier ratio  $\epsilon$   
Output: inlier set

  - Compute required number of iterations  $N = \frac{\log(1-p)}{\log(1-(1-\epsilon)^s)}$
  - For  $N$  iterations do:
    - Randomly select a subset of  $s$  data points
    - Fit model on the subset
    - Count inliers and keep model/subset with largest number of inliers
  - Refit model using found inlier set

58 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Igor Stukler

### Recap: RANSAC

- Required number of iterations  
-  $N$  for  $p = 0.99$

	Req. #points $s$	Outlier ratio $\epsilon$						
		10%	20%	30%	40%	50%	60%	70%
Line	2	3	5	7	11	17	27	49
Plane	3	4	7	11	19	35	70	169
Essential matrix	8	9	26	78	272	1177	7025	70188

59 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II  
Slide credit: Igor Stukler

### Textbooks

- More background on Algebraic Geometry and Visual Odometry can be found in

An Invitation to 3D Vision,  
Y. Ma, S. Soatto,  
J. Kosecka, and  
S. S. Sastry,  
Springer, 2004

60 Visual Computing Institute | Prof. Dr. Bastian Leibe  
Computer Vision 2  
Part 11 - Multi-Object Tracking II