Lecture Summary

Lecture 01: Introduction

Definition 1. Computer Vision: Automatic extraction of meaningful information from images and videos.

- Semantic information: meaning of objects.
- Geometric information: shapes.

Vision in Humans

Vision is the most powerful sense:

- Retina $\approx 1000 mm^2$.
- Contains 130 million **photoreceptors** (120 mio rods and 10 mio cones for color sampling)
- 3GBytes/s information flow, would need 500 Megapixel camera (8 Megapixel and range 15 degrees)

Why hard?

• only numbers, viewpoint variations, illumination challenges, motion, intra-class variations, scale and shape ambiguities

Origin: L.G. Roberts, MIT, 1963, Solids Perception.

Visual Odometry (VO)

Definition 2. VO is the process of incrementally estimating the pose of the vehicle by examining the changes that motion induces on the images of its onboard cameras.

Why VO?

- ≠ Wheel Odometry, VO **not affected** by wheel slippage and adverse conditions in general.
- More accurate. Relative position error 0.1-2 %.
- Can be used as a complement to wheel encorders, GPS, IMUs, etc.
- Crucial for flying, walking and underwater robots.

Assumptions

- Sufficient illumination,
- dominance of static scene over moving objects,
- enough texture to allow apparent motion,
- sufficient scene overlap between consecutive frames.

History

- 1980: NASA, Moraveck, Mars Rovers, sliding camera.
 - Mars rover with single camera sliding on a rail horizontally, taking 9 pictures at equidistant intervals.
 - Corners were detected with his algorithm and normalized using correlation.
 - Outliers were removed by checking depth inconsistencies.
 - Although one camera, stereo becuase used triangulation between different positions of the robot.
 - Single camera without stereo vision has disadvantage that motion cab be recovered up to a scale factor.
- 1980-2000: NASA: mission to mars.
- 2004: David Nister: Visual Odometry paper.

VO vs VSLAM vs SFM

Structure From Motion (SFM): more general than VO, tackles the problem of 3D reconstruction and 6DOF pose estimation from unordered image sets (e.g. reconstruction through flickr).

 \rightarrow VO focuses on estimating 3D motion **sequentially** and in **real time**.

Visual Simultaneous Localization And Mapping (VSLAM):

- VO focuses on incremental estimation and **local consistency**.
- focus on **globally consistent** estimation.
- VSLAM = VO + loop detection + graph optimization.
- Tradeoff between performance and consistency, simplificity of implementation.
- VO doesn't need to keep track of all previous history of the camera. Good for real-time.

Working Principle of VO:

1. Compute the relative motion T_k from images I_{k-1} to image I_k

$$T_k = \begin{pmatrix} R_{k,k-1} & t_{k,k-1} \\ 0 & 1 \end{pmatrix}.$$
 (1)

2. Concatenate them to recover full trajectory

$$C_n = C_{n-1} \cdot T_n \tag{2}$$

3. An optimization over the last m poses can be done to refine locally the trajectory (Pose-Graph or Bundle Adjustment).

Direct Image Alignment

First of all, there is a difference between direct and indirect methods. Indirect methods first extract features and then use them for localization and building of the map. Direct methods instead, try to recover the environment depth and structure and the camera pose through an optimisation on the map and camera parameters together (affected from light changes).

 \rightarrow The whole direct problem is a minimization of the per-pixel intensity difference, i.e.

$$T_{k,k-1} = \operatorname{argmin}_T \sum_{i} ||I_k(u'_i - I_{k-1}(u_i))||_{\sigma}^2$$
(3)

- Dense: \approx 300'000 pixels. Require powerful hardware (GPU)
- Semi-Dense: $\approx 10'000$ pixels.
- Sparse: $\approx 2'000$ pixels. (100-200 features x 4x4 patch). Basically pointclouds, used mainly to do localization (camera pose).

VO Flowchart

VO computes the camera path incremental, pose per pose

- 1. Image sequence,
- 2. Feature detection (front end),
- 3. Feature matching (tracking) (front end).
- 4. Motion estimation (2D-2D,3D-3D,3D,2D) (front end),
- 5. Local optimization (back end).

Lecture 02: Image Formation

How to Form an Image

Placing a film in front of an object and illuminating the object, the light is then reflected on the film. Not **reasonable** image! The rays don't converge in the same point, unsharp, blur.

Circle Blur or Circle of Confusion

In optics, a circle of confusion is an optical spot caused by a cone of light rays from a lens, not coming to a perfect focus when imaging a point source.

Pinhole Camera

Adding a barrier with a pinhole \rightarrow camera obscura:

- Opening is the **aperture**.
- Reduces blurring.
- **ideal pinhole**: only one ray of light reaches each point on the film. Bigger aperture, blurry image.

Why not as small as possible? Diffraction effects cause interference in waves as we near the wavelenth and less light gets through.

Converging Lens

- Rays passing through the **Optical Center** are not deviated.
- All rays parallel to the **Optical Axis** converge at the **Focal point**

Using similar triangles and Figure 1 one gets

$$\frac{B}{A} = \frac{e}{z} \text{ and}$$

$$\frac{B}{A} = \frac{e-f}{f} = \frac{e}{f} - 1.$$
(4)

Toghether we get the **thin lens equation**.

$$\frac{e}{f} - 1 = \frac{e}{z} \tag{5}$$

Remark.

• Any object point satisfying this equation is in **focus**. This is used for retrieving the distance to an object, depth from focus.



Figure 1: Thin lens

• z >> f:

$$\frac{1}{f} = \frac{1}{e} \Rightarrow f \approx e. \tag{6}$$

This is why we can think of a lens of focal length f as being equivalent to a pinhole distance f from the focal plane. \rightarrow we need to adjust the image plane such that objects at infinity are in focus, namely

$$\frac{h'}{h} = \frac{f}{z} \Rightarrow h' = \frac{f}{z}h.$$
(7)

The dependence of the apparent size of the object on its depth is known as **per-spective**.

In Focus and Blur Circle

- There is a specific distance from the lens at which world points are in focus in the image.
- Other points project to a **blur circle** in the image with radius

$$R = \frac{L\delta}{2e} \tag{8}$$

 \rightarrow a minimal pinhole gives minimal R and

 $\rightarrow R$ should remain smaller than image resolution.

Projective Geometry

- Straight lines are still straight.
- Length and angles are lost.
- Parallel lines in the world intersect in the image as a **vanishing point**.
- Parallel planes in the world intersect in the image at a **vanishing line**.

Other Parameters

Focus and Depth of Field

- DOF is the distance between the nearest and farthest objects in a scene that appear acceptably sharp.
- Decrease in sharpness is gradual on each side of the focused distance.
- Smaller aperture increases the range in which the object appears in focus but reduces the light.
- As f gets smaller: wide angle image.
- As f gets bigger: **narrow angle** image.

With Figure 2, one can compute the **field of view**



Figure 2: Scheme for angles.

$$\tan\left(\frac{\theta}{2}\right) = \frac{W}{2f} \Rightarrow f = \frac{W}{2} \left[\tan\left(\frac{\theta}{2}\right) \right]^{-1} \tag{9}$$

 \rightarrow smaller FOV = larger focal length.

Digital Camera

Instead of using a film we use a sensor array and we convert informations in numbers (e.g. [0, 255] for 8 bytes), as in Figure 3.

Color Sensing

- Bayer pattern (1976) places green filters over half of the sensors and red and blue over remaining ones. Humans are more sensitive to high frequency detail in luminance than chrominance. This method has the disadvantage that the number of detected pixels is cut by 4.
- The three-chip color camera splits in three color filters the light and each chip measures light intensity for one color. Here, resolution is preserved.
- Estimate missing components from neighboring values: **demosaicing**



Figure 3: Procedure digital cameras.

Perspective Camera Model

The procedure reads

- 1. Change of coordinates from 3D world to adapted frame.
- 2. Projection from the camera frame to the image plane.
- 3. Change in pixel.

For convenience, the image plane is usially represented in front of C such that the image



preserves the same orientation.

Figure 4: Frames.

Perspective Projection

We use the notation from euclidean to homogeneous:

$$\underbrace{\begin{pmatrix} u \\ v \\ w \\ w \end{pmatrix}}_{Homog.} \to \underbrace{\begin{pmatrix} u/w \\ v/w \\ 1 \\ Eucl. \end{pmatrix}}_{Eucl.}$$
(10)

From similar triangles and Figure 5 one gets

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

$$\frac{y}{f} = \frac{Y_c}{Z_c} \Rightarrow y = \frac{fY_c}{Z_c}$$
(11)



Figure 5: Figure for perspective projection.

Pixel Coordinates

From local image plane coords (x, y) to the pixel coords (u, v), with scale factors $k_{u,v}$ (inverse of the effective pixel site along the u(v) direction, measured in pixel m^{-1}):

$$u = u_0 + k_u x \Rightarrow u = u_0 + \frac{k_u f X_c}{Z_c}$$

$$v = v_0 + k_v y \Rightarrow v = v_0 + \frac{k_v f Y_c}{Z_c},$$
(12)

which expressed in matrix form reads

$$\begin{pmatrix} \lambda u \\ \lambda v \\ \lambda \end{pmatrix} = \begin{pmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} X_c \\ Y_c \\ Z_c \end{pmatrix} = K \begin{pmatrix} X_c \\ Y_c \\ Z_c \end{pmatrix},$$
(13)

with $\alpha_{u,v} = k_{u,v} f$ focal lengths and K calibration matrix (intrinsic parameter matrix). At the end, it holds

$$\lambda \cdot \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = \underbrace{(K)}_{3 \times 3, \ intrinsic} \cdot \underbrace{(\mathbb{I}_{3 \times 3}|0)}_{3 \times 4} \cdot \underbrace{\begin{pmatrix} R & \vec{t} \\ 0 & 1 \end{pmatrix}}_{4 \times 4, \ extrinsic} \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix}.$$
(14)

Lens Distortion

Radial Distortion

The amount of distortion of the coordinates of the observed image is a nonlinear function of their radial distance. This is a transformation from ideal to distorted coordinates. For most lenses, one writes a simple quadratic model

$$\begin{pmatrix} u_d \\ v_d \end{pmatrix} = (1+k_1r^2) \cdot \begin{pmatrix} u-u_0 \\ v-v_0 \end{pmatrix} + \begin{pmatrix} u_0 \\ v_0 \end{pmatrix},$$
(15)

with

$$r^{2} = (u - u_{0})^{2} + (v - v_{0})^{2}.$$
(16)

Depending on the amount of distortion, one can introduce higher order terms:

$$\begin{pmatrix} u_d \\ v_d \end{pmatrix} = (1 + k_1 r^2 + k_2 r^4 + k_3 r^6) \cdot \begin{pmatrix} u - u_0 \\ v - v_0 \end{pmatrix} + \begin{pmatrix} u_0 \\ v_0 \end{pmatrix}$$
(17)

How to build an AMES room



Figure 6: Thin lens

Lecture 03: Image Formation 2

Pose determination from n Points (PnP) Problem

Given the realitve spatial locations of n control points and given the angle to every pair of control points from an additional point called the Center of Perspective C_P , find the lengths of the line segments joining C_P to each of the control points.

We assume we know the camera intrinsic parameters. Given known 3D landmarks in the world and their image correspondence in the camera frame, determine the 6DOF pose of the camera in the world frame. Where is the camera?

- Given 1 point: ∞ solutions.
- Given 2 points: ∞ bounded solutions.
- Given 3 non collinear points: Finitely many (up to 4) solutions.
- Given 4 points: unique solution.

With 3 points one can use the fact that the angles inscribed in the triangle are the same: the Carnot's theorem reads

$$s_{1,2,3}^2 = L_{B,A,A}^2 + L_{C,C,B}^2 - 2L_{B,A,A}L_{C,C,B}\cos(\theta_{BC,AC,AB})$$
(18)

In general, n independent polynomials with n unknowns, can have no more solution than the **product** of their degrees: here 8.

 \rightarrow fourth point to **disambiguate** the solutions! By defining $x = \frac{L_B}{L_A}$ we can reduce the system to a 4th order equation

$$G_0 + G_1 x + G_2 x^2 + G_3 x^3 + G_4 x^4 = 0.$$
⁽¹⁹⁾

This applies to camera pose estimation from known 3D - 2D correspondences (e.g. hololens).

Camera Calibration

Determine **intrinsic** and **extrinsic** parameters of the camera model.

Tsai, 1987: Measure the 3D position of more than 6 points on a 3D calibration target and the 2D coordinates of their projection. We can do that, by recalling the perspective projection equation, by neglecting the radial distortion.

Direct Linear Transform

image point =
$$\tilde{p} = \begin{pmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{pmatrix} = \lambda \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K[R|T] \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{pmatrix} \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$
indep. elements
$$= \underbrace{\begin{pmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{pmatrix}}_{M} \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} m_1^T \\ m_2^T \\ m_3^T \end{pmatrix} \cdot \underbrace{\begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}}_{P}.$$
(20)

It follows

assuming

$$u = \frac{\tilde{u}}{\tilde{w}} = \frac{m_1^T \cdot P}{m_3^T \cdot P}$$

$$v = \frac{\tilde{v}}{\tilde{w}} = \frac{m_2^T \cdot P}{m_3^T \cdot P}$$
(21)

and hence

$$(m_1^T - u_i m_3^T) \cdot P_i = 0 (m_2^T - v_i m_3^T) \cdot P_i = 0.$$
 (22)

Rearranging the terms you have

$$\begin{pmatrix} P_1^T & 0^T & -u_1 P_1^T \\ 0^T & P_1^T & -v_1 P_1^T \end{pmatrix} \cdot \begin{pmatrix} m_1 \\ m_2 \\ m_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}.$$
 (23)

For *n* points we have a big $2n \times 12$ matrix *Q*. The problem hence reads

$$Q \cdot M = 0, \tag{24}$$

where Q is known and M is unknown.

Minimal Solution:

- Rank 11 to have unique non-trivial (up to scale) solution M (Q known!).
- $\bullet\,$ Each 3D/2D correspondence provides 2 independent equations.
- $5 + \frac{1}{2}$ correspondences are needed (in fact 6).

Overdetermined Solution:

- More than 6 points.
- Minimize $||QM||^2$ with the constraint $||M|| = 1 \rightarrow \text{SVD}$. The solution is the eigenvector corresponding to the smallest eigenvalue of Q^TQ . That's because this is the unit vector x that minimizes $||Qx||^2 = x^TQ^TQx$. This can be done in matlab with

[U,S,V] = svd(Q); M = V(:,12);

Degenerated Configurations:

- Points lying on a plane and or along a line passing through the projection center.
- Camera and points on a twisted cubic (degree 3).

Once we have M, we know from its definition

$$M = K(R|T). (25)$$

Remark.

- We are not enforcing orthogonality of R.
- QR factorization of M, whith R (orth.) and T (upper triangular matrix).

Tsai Method 1987

- 1. Edge detection.
- 2. Straight line fitting to the detected edges.
- 3. Intersecting the lines to obtain the image corners (<0.1 pixels accuracy).
- 4. Use more than 6 points (more than 20) and **not** all on same plane.

Originally pixels were not squared (parallelogramms, skew, no rectangle). Most cameras today are well manufactured: $\frac{\alpha_u}{\alpha_v} = 1$ and $K_{12} = 0$. **Residual**: Average reprojection error, computed as the distance (in pixels) between the

Residual: Average reprojection error, computed as the distance (in pixels) between the observed point and the camera-reprojected 3D point. Accuracy of calibration. What if K is known? Nothing changes!

Calibration from Planar Grids (Homographies)

Zhang, Microsoft Use a planar grid (chessboard) and a few image of it at different orientations. Setting $Z_w = 0$ we get

$$\begin{pmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{pmatrix} = \begin{pmatrix} \alpha_u & 0 & v_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{pmatrix} \cdot \begin{pmatrix} X_w \\ 0 \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} \alpha_u & 0 & 0 \\ 0 & \alpha_v & \nu_0 \\ 0 & 0 & 1 \end{pmatrix} \cdot \begin{pmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{pmatrix} \cdot \begin{pmatrix} X_w \\ Y_w \\ 1 \end{pmatrix}$$

$$= H \cdot \begin{pmatrix} X_w \\ Y_w \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} h_1^T \\ h_2^T \\ h_3^T \end{pmatrix} \cdot \begin{pmatrix} X_w \\ Y_w \\ 1 \end{pmatrix} .$$

$$(26)$$

Hence, one more time

$$u = \frac{\tilde{u}}{\tilde{w}} = \frac{h_1^T \cdot P}{h_3^T \cdot P}$$

$$v = \frac{\tilde{v}}{\tilde{w}} = \frac{h_2^T \cdot P}{h_3^T \cdot P}$$
(27)

and hence

Rearranging the terms you have

$$\begin{pmatrix} P_1^T & 0^T & -u_1 P_1^T \\ 0^T & P_1^T & -v_1 P_1^T \end{pmatrix} \cdot \begin{pmatrix} h_1 \\ h_2 \\ h_3 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix},$$
(29)

where Q is known and H is unknown. For n points we have a big $2n \times 9$ matrix Q. The problem hence reads

$$Q \cdot H = 0. \tag{30}$$

Minimal Solution:

- Rank 8 to have unique non-trivial solution H (Q known!).
- Each point correspondence provides 2 independent equations.
- 4 non-collinear points are needed.

Overdetermined Solution:

- More than 4 points.
- Minimize $||QM||^2 \rightarrow$ SVD. The solution is the eigenvector corresponding to the smallest eigenvalue of $Q^T Q$. We can decompose as before.

This **projective** transformation is called **Homography**. Applications are

- Augmented reality
- Beacon-based localization.

DLT vs. PnP

• If the camera is calibrated, only R and T need to be determined. Pnp leads to smaller error.

Non Conventional Camera Models

Omnidirectional Cameras

- Wide FOV dioptric cameras (e.g. fisheye (180)).
- \bullet Catadioptric cameras (e.g. mirrors (>180)) combine a standard camera with a shaped mirror.
 - Mirror: central, mirror (surface of revolution of a conic), single effective view point.
 - Perspetive: hyperbola+perspective / parabola+orthographic lens.
- Polydioptric cameras (e.g. multiple overlapping cameras) $\approx 360.$

Central: A vision system is said to be central when the optical rays to the viewd objects intersect in a single point in 3D called projection center or *single effective viewpoint*. For hyperbolic and elliptical mirrors, the single viewpoint property is achieved by ensuring that the camera center coincides with one of the foci of the hyperbola (ellipse). For this, see Figure 7

Class of **rotated (swept) conic shapes** (hyperbolical, parabolical, elliptical mirrors) **Why is it important that the camera is central?** If the camera is central, we can unwarp parts of omnidirectional image into perspective. We can transform image points in the unit sphere. We can apply algorithms for perspective geometry. Perspective and omnidirectional model are equal!



Fig. 3. Central catadioptric cameras can be built by using hyperbolic and parabolic mirrors. The parabolic mirror requires the use of an orthographic lens.

Figure 7: How should the mirrors be?

Lecture 04: Image Filtering

Difference between convolution and correlation: correlation is a metric for similarity between two different signals. Convolution applies one signal to the other. *Filtering*: Accepting or rejecting certain frequency components.

- low-pass filter smooths an image.
- high-pass filter retains the edges of an image.

Low-pass Filtering

We want to reduce noise! There different types of noise:

- Salt and pepper noise: randum occurences of black and white pixels.
- Impulse noise: random occurences of white pixels.
- Gaussian noise: variations in intensity drawn from a Gaussian normal distribution.

How can we reduce the noise to try to recover the ideal image?

Moving Average

Replaces each pixel with an average of all the values in its neighborhood.

- Pixels like neighbors.
- Assumption: noise process independent from pixel to pixel.

Weighted Moving Average

Can add weights to moving average

- Uniform weights.
- Non-uniform weights.

Nothing else than convolution!

Convolution

One of the sequences is flipped before sliding over the other. Linearity, associativity, commutativity. Notation f * g. In 1D:

$$f * g = \int_{-\infty}^{\infty} f(\tau)g(t-\tau)\mathrm{d}\tau.$$
 (31)

In 2D:

$$G[i, j] = H * F$$

= $\sum_{u=-k}^{k} \sum_{v=-k}^{k} H[u, v] F[i - u, j - v].$ (32)

In other words: replacing each pixel with a linear combination of its neighbors. The filter H is also called kernel or mask. One can change the weights

Gaussian Filter

What if we want the **closest** pixels to have a higher influence on the output?

$$h(u,v) = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2 + v^2}{2\sigma^2}}.$$
(33)

What parameter matter?

- Size of the kernel. The Gaussian has generally infinite support but discrete filters use finite kernels.
- The **variance** of gaussian: determines extent of smoothing (larger variance, larger smoothing).

Boundary Issues

The filter window falls off the edge of the image. We need to pad the image borders with

- Zero padding (black)
- Wrap around
- Copy edge
- Reflect across edge

Median Filter

Linear smoothing filters do not alleviate salt and pepper noise! Non-linear filter. Removes spikes: good for impulse and salt and pepper noise.

Computes the median value and replaces the high value with that.

- +: Preserves sharp transitions.
- -: Removes small brightness variations.

High-pass Filtering (edge detection)

We want an idealized line drawing. The edge contours in the image correspond to important scene contours. Edges are nothing else than sharp intensity changes. Images can be expressed as functions f(x, y). Edges correspond to **extrema of derivative**.

Differentiation and Convolution

For discrete data, it holds

$$\frac{\mathrm{d}f(x,y)}{\mathrm{d}x} \approx \frac{f(x+1,y) - f(x,y)}{1}.$$
(34)

Partial derivatives of an image are in x (-1,1), in $y \begin{pmatrix} -1 \\ 1 \end{pmatrix}$. Other finite differences methods are the

• Prewitt Filter

$$G_x = \begin{pmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{pmatrix}, \qquad G_y = \begin{pmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{pmatrix}$$
(35)

• Sobel Filter

$$G_x = \begin{pmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{pmatrix}, \qquad G_y = \begin{pmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{pmatrix}$$
(36)

The **Gradient** of an image if given by

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial x} & \frac{\partial f}{\partial y} \end{pmatrix}$$
(37)

The gradient direction is given as

$$\Theta = \tan^{-1} \left(\frac{\frac{\partial f}{\partial y}}{\frac{\partial f}{\partial x}} \right).$$
(38)

The **edge strength** is given by

$$||\nabla f|| = \sqrt{\left(\frac{\partial f}{\partial y}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}.$$
(39)

Handling Noise

If we differentiate a noisy signal, we get infinite many peaks. Solutions are

- first **smooth** the signal (with a convolution with h). Then differentiate.
- Combining the two, convolute with $\left(\frac{d}{dx}h\right) * f$.

Laplacian of a Gaussian

Consider

$$\frac{\partial^2}{\partial x^2}(h*f)\tag{40}$$

Where is the edge? Zero-crossin of bottom graph.

Summary

- Smoothing filters:
 - Has positive values.
 - Sums to $1 \rightarrow$ preserves brightness of constant regions.
 - Removes high frequency components.
- Derivative Filters
 - Has opposite signs, used to get high response in regions of high constrast.
 - Sums to $0 \rightarrow$ no response in constant regions.
 - highlights high frequency components.

The Canny Edge-Detection Algorithm

- We compute the gradient of smoothed image in both directions. (convolve the image wieh x and y derivatives of Gaussian filters)
- We discard pixels whose gradient magnitude is below a certain threshold.
- Non-maximal suppression: local maxima along gradient diretion. High intensity means high probability of the presence of an edge: this is not enough. Only local maxima can be considered as part of an edge. A local maxima can be found where the gradient derivative is 0.
 - Compare the edge strength of the current pixel with the edge strength of the pixel in the positive and negative gradient directions.
 - If the edge strength of the current pixel is the largest compared to the other pixels in the mask with the same direction (i.e., the pixel that is pointing in the y-direction, it will be compared to the pixel above and below it in the vertical axis), the value will be preserved. Otherwise, the value will be suppressed.

Lecture 05: Feature Detection 1

Goal: reduce amount of data to process in later stages, discard redoundancy to preserve only what is useful (lower bandwidth and memory storage). In general

- Edge detection.
- Template matcing.
- Keypoint detection.

Filters for Template Matching

We want to 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. What if the template is not identical to the object we want to detect? This works only if

- Scale,
- orientation,
- illumination,
- appearance of the template and the object are similar.

What about the objects in the background?

Correlation as Scalar Product

We consider images H and F as vectors and express the correlation between them as

$$\langle H, F \rangle = ||H|| \cdot ||F|| \cdot \cos(\theta). \tag{41}$$

If we use **Normalized Cross Correlation** (NCC) (highest complexity), we consider the unit vectors of H and F, hence we measure their similarity based on the angle θ . For identical vectors one gets NCC = 1. NCC is invariant to linear intensity changes! It holds

$$\cos(\theta) = \frac{\langle H, F \rangle}{||H|| \cdot ||F||} = \frac{\sum_{u=-k}^{k} \sum_{v=-k}^{k} H(u, v) F(u, v)}{\sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} H(u, v)^{2}} + \sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} F(u, v)^{2}}} > 0.8.$$
(42)

Other methods are the Sum of Absolute Differences (SAD) (simplest)

$$SAD = \sum_{u=-k}^{k} \sum_{v=-k}^{k} |H(u,v) - F(u,v)|,$$
(43)

the Sum of Squared Differences (SSD) (high computational complexity)

$$\sum_{u=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - F(u,v))^2.$$
(44)

The normalized cross correlation (NCC) takes values between -1 and 1, 1 equals identical.

To account for the difference in mean of the two images (caused principally by illumination changes), we substract the mean value of each image:

• Zero-mean Sum of Absolute Differences (ZSAD)

$$ZSAD = \sum_{u=-k}^{k} \sum_{v=-k}^{k} |(H(u,v) - \mu_H) - (F(u,v) - \mu_F)|.$$
(45)

• Zero-mean Sum of Squared Differences (ZSSD)

$$\sum_{u=-k}^{k} \sum_{v=-k}^{k} ((H(u,v) - \mu_H) - (F(u,v) - \mu_F))^2.$$
(46)

• Zero-mean ormalized Cross Correlation (ZNCC)

$$ZNCC = \frac{\sum_{u=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_H) \cdot (F(u,v) - \mu_F)}{\sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} (H(u,v) - \mu_H)^2} + \sqrt{\sum_{u=-k}^{k} \sum_{v=-k}^{k} (F(u,v) - \mu_F)^2}}.$$
(47)
with $\mu_H = \frac{\sum_{u=-k}^{k} \sum_{v=-k}^{k} H(u,v)}{(2N+1)^2}$

Remark. ZNCC is **invariant** to affine intensity changes.

Census Transform

It maps an image patch to a bit string. The general rule is that **if a pixel is greater than the center pixel** its corresping bit is set to 1, else to 0. For a $n \times n$ window the string will be $n^2 - 1$ bits long. The 2 bit strings are compared using the **Hamming distance** (if bigger than previous 1, else 0, starting from right). The **Advantages** are

- More robust to object background problem
- No square roots or divisions are required. Efficient!
- Intensities are considered relative to the center pixel of the patch making it **invariant to monotonic intensity** changes.

Point-feature Extraction and Matching

Keypoint extraction is the key ingredient of motion estimation! Furthermore, used for panorama stitching, object recognition, 3D reconstruction, place recognition, google images.

Why problematic? We need to align images! How? Detect point features in both images and find corresponding pairs to align them. Two big problem arise

• Problem 1: Detect the **same** points **independently** in both images. No chance to match, need **repeatable** feature detector.

• Problem 2: for each point, identify its correct correspondence. Need **reliable** and **distinctive** feature descriptor. Robust to **geometric** and **illumination** changes.

Geometric changes: rotation, scale and viewpoint (i.e. perspective changes). **Illumination changes:** Affine illumination changes

$$I'(x,y) = \alpha I(x,y) + \beta.$$
(48)

Invariant local features: Subset of local feature types designed to be invariant to common geometric and photometric transformations. In general

- 1. Detect distinctive interest points,
- 2. extract invariant descriptors.

Distinctive features s.t. repeatable? Some features are better than others (angles, not uniform color,...).

Corners: a corner is defined as the intersection of one or more edges. It has **high localization accuracy**. Less distinctive than a **blob**. Examples of corner detectors are Harris, Shi-Tomasi, SUSAN, FAST.

Blob: is any other image pattern which is not a corner, that differs significantly from its neighbors in intensity and texture. This has less localization accuracy, but better for place recognition because more distinctive than a corner E.g. MSER, LOG, DOG, SIFT, SURF, CenSurE.

Corner Detection

In the region around a corner, the image gradient has two or more dominant directions. Corners are repeatable and distinctive.

The Moravec Corner Detector (1980)

We can easily recognize the point by looking through a small window: by shifting the window, one can give large change in intensity. Moravec used SSD, whith

- Flat region: no intensity change! (SSD ≈ 0 in all directions).
- Edge: no change along the edge direction (SSD ≈0 along adge but >> 0 in other directions).
- Corner: significant change in at least two directions (SSD >> 0 in at least 2 directions.

Sums of squares of differences of pixels adjacent in each of four directions (horizontal, vertical and two diagonals) over each window are calculated, and the window's interest measure is the minimum of these four sums.[Moravec,80]

The Harris Corner Detector (1988)

Implements Moravec corner detector without physically shifting the window and hence just by looking at the patch itself: using **differential calculus**.

Implementation:LEt I be a grayscale image. We consider the reference patch centered at (x, y) and the shifted window centered at $(x + \Delta x, y + \Delta y)$. The patch has size P. We compute

$$SSD(\Delta x, \Delta y) = \sum_{x,y \in P} \left(I(x,y) - I(x + \Delta x, y + \Delta y) \right)^2.$$
(49)

We define

$$I_x = \frac{\partial I(x, y)}{\partial x}, \quad I_y = \frac{\partial I(x, y)}{\partial y}, \tag{50}$$

and approximate with first order Taylor expansion:

$$I(x + \Delta x, y + \Delta y) \approx I(x, y) + I_x(x, y)\Delta x + I_y(x, y)\Delta y$$

$$\Rightarrow SSD(\Delta x, \Delta y) \approx \sum_{x,y \in P} (I_x(x, y)\Delta x + I_y(x, y)\Delta y)^2$$
(51)

Simple quadratic function in the deltas! We can write this in matrix form:

$$SSD(\Delta x, \Delta y) \approx \begin{pmatrix} \Delta x & \Delta y \end{pmatrix} \cdot \underbrace{\sum_{x,y \in P} \begin{pmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{pmatrix}}_{M} \cdot \begin{pmatrix} \Delta x \\ \Delta y \end{pmatrix},$$
(52)

where M is the second moment matrix. Let's analyze some special cases:

- Edge along x: $M = \begin{pmatrix} 0 & 0 \\ 0 & \lambda_2 \end{pmatrix}$.
- Flat region: $M = \begin{pmatrix} 0 & 0 \\ 0 & 0 \end{pmatrix}$
- Aligned corner: $M = \begin{pmatrix} \cos(45) & -\sin(45) \\ \sin(45) & \cos(45) \end{pmatrix} \cdot \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \cdot \begin{pmatrix} \cos(45) & \sin(45) \\ -\sin(45) & \cos(45) \end{pmatrix}$

What if the corner is not aligned with the image axis? The general case has M symmetric, which can always be decomposed into

$$M = R^{-1} \cdot \begin{pmatrix} \lambda_1 & 0\\ 0 & \lambda_2 \end{pmatrix} \cdot R.$$
(53)

Claim 1. One can visualize this as an ellipse with axis lengths determined by eigenvalues $(1/\sqrt{\lambda_{\max,\min}})$ and two axes determined by the eigenvectors of M (columns of R).

Proof. Let's consider

$$M = \begin{pmatrix} v_1 & v_2 \end{pmatrix} \cdot \begin{pmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{pmatrix} \cdot \begin{pmatrix} v_1^T \\ v_2^T \end{pmatrix} = 1$$
(54)

Then, using the quadratic form one gets

$$x^{T} \cdot (v_{1} \quad v_{2}) \cdot \begin{pmatrix} \lambda_{1} & 0\\ 0 & \lambda_{2} \end{pmatrix} \cdot \begin{pmatrix} v_{1}^{T}\\ v_{2}^{T} \end{pmatrix} \cdot x = 1$$

$$\lambda_{1}x^{T}v_{1}v_{1}^{T}x + \lambda_{2}x^{T}v_{2}v_{2}^{T}x = 1$$

$$\lambda_{1}(v_{1}^{T}x)^{T}(v_{1}^{T}x) + \lambda_{2}(v_{2}^{T}x)^{T}(v_{2}^{T}x) = 1$$

$$\frac{(v_{1}^{T}x)^{2}}{\left(\frac{1}{\sqrt{\lambda_{1}}}\right)^{2}} + \frac{(v_{2}^{T}x)^{2}}{\left(\frac{1}{\sqrt{\lambda_{2}}}\right)^{2}} = 1,$$
(55)

from which is clear that the eigenvectors v_1, v_2 represent the axis directions of the ellipse and $\frac{1}{\sqrt{\lambda_1}}, \frac{1}{\sqrt{\lambda_2}}$ their length.

Remark. Large ellipses denote flat region, small ones a corner!

Interpreting the eigenvalues A corner can then identified by checking whether the minimum of the two eigenvalues of M is larger than a certain user-defined threshold. Mathematically, this is the **Shi-Tomasi detector**

$$R = \min(\lambda_1, \lambda_2) > \text{threshold.}$$
(56)

- Corner: $\lambda_{1,2}$ are large, R > threshold, SSD increases in each direction.
- Edges: $\lambda_1 >> \lambda_2$ or vice-versa.
- Both small: flat region.

Problem: The eigenvalues are expensive to compute: *Harris and Stephens* suggested to use the fact

$$R = \lambda_1 \cdot \lambda_2 - k(\lambda_1 + \lambda_2)^2 = det(M) - k \cdot trace^2(M), \quad k \in (0.04, 0.15)$$
(57)

Algorithm:

- (I) Compute derivatives in x and y directions e.g. with Sobel filter.
- (II) Compute $I_x^2, I_y^2, I_x I_y$.
- (III) Convolve $I_x^2, I_y^2, I_x I_y$, with a box filter to get the sums of each element, which are the entries of the matrix M. (optionally use Gaussian filter instead of box filter to give more imoportance to central pixels).
- (IV) Compute Harris Corner Measure R (with Shi-tomasi or Harris).
- (V) Find points with large corner response (R > threshold).
- (VI) Take the points of local maxima of R.

Harris vs. Shi-Tomasi???

Repeatability Can it re-detect the same image patches (Harris corners) when the image exhibits changes?

- Corner response *R* is **invariant to image rotation**. Shape (eigenvalues) remains the same. Isotropic!
- Invariant to **affine intensity changes**: eigenvalues are scaled by a constant factor but position of the maxima remains the same.
- Not invariant to **image scale**. Scaling the image by ×2 results in 18 % of correspondences get matched.

Lecture 06: Point Feature Detection 2

Scale Changes

A possible solution is to rescale the patch, i.e. bring it to the canonical scale. The problem by scale search is that it is scale consuming: we need to do it **individually** for all patches in one image. In fact, the complexity would be $(NM)^2$ (assuming N features per image and M scale levels for each image). \Rightarrow a possible solution is to assign each feature its own scale.

Automatic Scale Selection

- Design a function on the image patch, which is scale invariant, i.e., which has the same value for corresponding regions, even if they are at different scales.
- For a point in one image, we can consider it is a function of region size.

Approach: We take a local maximum of the function: the region size for which the maximum is achieved, should be invariant to image scale. **This scale invariant region size is found in each image independently**. When the right scale is found, the patch must be **normalized**.

- Good function: single and sharp peaks!
- If multiple peaks: assign **more** region sizes to have a unique feature. Blobs and corners are the ideal locations!

Function:

Convolve image with kernel to identify sharp discontinuities:

$$f = \text{Kernel} * \text{Image} \tag{58}$$

It has been shown that the Laplacian of Gaussian kernel is optimal under certain assumptions

$$LoF = \nabla^2 G(x, y) = \frac{\partial G(x, y)}{\partial x^2} + \frac{\partial G(x, y)}{\partial y^2},$$
(59)

Then, the correct scale is found as local maxima across consecutive smoothed images. This should be done for **severals** region sizes (15).

Note that an efficient implementation of multisale detection uses the so called scale-space pyramid: instead of varying the window size of the feature detector, the idea is to generate upsampled (enlarge the image, interpolating) or downsampled versions of the same image.

Feature Descriptors

We already know how to detect points, but how can we describe them for matching?

- **Simplest Descriptor**: **Intensity** values within a squared patch or gradient histogram.
- Census transform or Histograms of Oriented Gradients.

Then, descriptor mathicing can be done using Hamming Distance (Census) or (Z)SSD,(Z)SAD, (Z)NCC. We would like to fame the same features regardless of the transformation that is applied to them: most feature methods are designed to be invariant of

- 2D translation,
- 2D rotation,
- scale.

Some of them can also handle

- Small view point invariance (e.g. SIFT works up to about 60 degrees).
- Linear illumination changes.

How to achieve Invariance?

Step 1: Re-scaling and De-rotation:

- Find the **correct scale** using LoG operator.
- **Rescale** the patch to a default size (e.g. 8×8 pixels).
- Find the **local orientation** (e.g. with dominant direction, Harris eigenvectors)
- **De-rotate** the patch.

In order to de-rotate the patch, one uses **patch-warping**:

Patch Warping

- 1. Start with **empty** canonical patch (all pixels set to 0).
- 2. For each (x, y) in the empty patch, apply the **warping function** W(x, y) to compute the corresponding position in the detected image. It will be in floating point and will fall between the image pixels.
- 3. Interpolate the intensity values of the 4 closest pixels in the detected image with
 - Nearest neighbor,
 - Bilinear interpolation

Example 1: Rotational Warping

Counterclockwise rotation:

$$\begin{aligned} x' &= x\cos(\theta) - y\sin(\theta) \\ y' &= x\sin(\theta) + y\cos(\theta) \end{aligned} \tag{60}$$

Bilinear Interpolation

It is an extension of the linear interpolation, for interpolating functions of two variables on a rectilinear 2D grid. The key idea is to perform linear interpolation in one direction and then, again, in the other direction. We have that

$$I(x,y) = I(0,0) \cdot (1-x) \cdot (1-y) + I(0,1) \cdot (1-x) \cdot y + I(1,0) \cdot x \cdot (1-y) + I(1,1) \cdot xy$$
(61)

Example 2: Affine Warping

To achieve **slight view-point invariance**:

- The second moment matrix M can be used to identify the two directions of fastest and slowest change of intensity around the feature.
- Out of these two directions, an elliptic patch is extracted at the scale computed with the LoG operator.
- The region inside the ellipse is normalized to a circular one.

There are however disadvantages:

- If not warped patches, very small errors in rotation, scale and view-point will affect matching score significantly.
- Computationally expensive (need to unwarp every patch)

A better solution nowadays is HOGs

Histogram of Oriented Gradients

- Compute a histogram of orientations of intensity gradients.
- Peaks in histogram are dominant orientations.
- **Keypoint orientation**= **histogram peak**. If there are multiple candidate peaks, construct a different keypoint for each such orientation.
- Rotate patch according to this angle: this puts the patches into a canonical form.

Scale Invariant Feature Transform (SIFT) Descriptor

The uniqueness of SIFT is that these features are extremely distinctive and can be successfully mathed between images with very different illumination, rotation, viewpoint, and scale-changes. The SIFT algorithm does:

- Identification of Keypoint location and scale
- Orientation assignment
- Generation of keypoint descriptor

Descriptor computation:

- 1. Divide the patch into 4×4 sub-patches=16 cells.
- 2. Compute HOG (8 bins, i.e. 8 directions) for all pixels inside each sub-patch.
- 3. Concatenate all HOGs into a single 1D vector. This is the resulting SIFT descriptor: $4 \times 4 \times 8 = 128$ values.
- 4. Descriptor matching: SSD (euclidean-distance).

Intensity Normalization

The descriptor vector v is then normalized such that its l_2 norm is 1:

$$\bar{v} = \frac{v}{\sqrt{\sum_{i}^{n} v_i^2}}.$$
(62)

Remark. This guarantees that the descriptor is invariant to linear illumination changes. This was already invariant to additive illumination because it is based on gradients.

SIFT matching robustness

- Can handle changes in viewpoint (up to 60 degree out-of-plane rotation).
- Can handle significant changes in illumination (low to bright scenes).
- Expensive: 10fps.

In order to reduce the computational cost, one can use difference of Gaussian instead of Lapiacian:

$$LOG \approx DOG = G_{k\sigma}(x, y) - G_{\sigma}(x, y)$$
(63)

SIFT Detector

SIFT keypoints are local extrema (maxima and minima) in both **space and scale** of the DoG images:

- Detect maxima and minima of difference-of-Gaussian in scale space.
- Each point is compared to its 8 neighbors in the current image and 9 neighbors each in the scales above and below.
- For each max and min found, the output is the location and the scale: this is a candidate keypoint.

Similar to Harris? While in Harris (keypoint location) the keypoint is identified in the image plane as local maximum of the corner function, in SIFT the keypoint is a local minimum or maximum of the DoG image in both position and scale. Implementation:

- 1. The initial image is **incrementally convolved with Gaussians** $G(k\sigma)$ to produce images separated by a constant factor k in scale space.
 - (a) The initial Gaussian $G(\sigma)$ has $\sigma = 1.6$.
 - (b) k is chosen such that $k = 2^{\frac{1}{s}}$, where s in an integer (typically s = 3).
 - (c) For efficiency reasons, when k reaches 2, the image is downsampled by a factor of 2 and then the procedure is repeated up to 4 or 6 octaves (pyramid levels).
- 2. Adjacent image scales are then subtracted to produce the difference-of-Gaussian (DoG) images.



Figure 8: Difference of Gaussian.

Summary

- An approach to detect and describe regions of interest in an image.
- SIFT detector = DoG detector.
- SIFT features are reasonably invariant to changes in rotation, scaling, and changes in viewpoint (up to 60deg) and illumination.
- Real time but still **slow** (10Hz on an i7 laptop).

The **repeatability** can be expressed as

$$\frac{\text{number of correspondences detected}}{\text{number correspondences present}}$$
(64)

The highest repeatability is obtained when sampling 3 scales per octave.

Influence of Number of Orientation and NUmber of Sub-Patches: Single orientation histogram is poor at discriminating, but the results continue to improve up to a 4x4 array of histograms with 8 orientations.

- **Descriptor**: 4x4x8 = 128-element 1D vector.
- Location: 2D vector.
- Scale of the patch. 1 scalar value.
- Orientation (angle of the patch). 1 scalar value.

SIFT for object recognition

Can be simply implemented by returning as best object match the one with the largest number of correspondences with the template (object to detect). 4 or 5 point RANSAC can be used to remove outliers.

Feature Matching

Given a feature I_1 , how to find the best match in I_2 ?

- 1. Define distance function that compares two descriptors ((Z)SSD,SAD,NCC, or Hamming distance for binary descriptors (e.g. Census, BRIEF, BRISK).
- 2. Brute-force matching:
 - Test all the features in I_2 .
 - Take the one at min distance.

Issues with closest descriptor: Can give good scores to very ambiguous (bad) matches (curse of dimensionality). A better approach would be to compute the ratio of distances between the *first* and *second* match:

$$\frac{d(f_1)}{d(f_2)} < \text{Threshold (usually 0.8)},\tag{65}$$

where

$$d(f_1)$$
 is the distance of the closest neighbor
 $l(f_1)$ is the distance of the closest neighbor (66)

 $d(f_2)$ is the distance of the second closest neighbor

Explanation for distance ratio: In SIFT, the nerest neighbor is defined as the keypoint with minimum euclidean distance. However, many features from an image 1 may not have any correct match in image 2 because they arise from background cluttere or were *not detected at all* in the image 1. An effective measure is otained by comparing the distance of the closest **neighbor** to that of the second closest neighbor. Why? Correct matches need to have the closest neighbor significantly closer than the closest incorrect match, to achieve reliable matching. Moreover, for **false matches**, there will likely be a number of other false matches within similar distances due to the high dimensionality of the feature space. (aka **curse of dimensionality**). We can think of the second closest match as providing an estimate of the density of false matches within this portion of the feature space, and at the same time identifying specific instances of feature ambiguity. **Why 0.8**?

- Eliminates 90% of the false matches,
- $\bullet\,$ Discards less than 5% of the correct matches.

SURF (Speeded Up Robust Features)

- Based on ideas similar to SIFT.
- Approximated computation for detection and descriptor using box filters.

- Results are comparable with SIFT but
 - Faster computation and
 - Shorter descriptors.

FAST detector (Features from Accelerated Segment Test)

- Studies intensity of pixels around candidate pixel C.
- C is FAST corner **if** a set on N contiguous pixels on circle are
 - all brighter than intensity(C) + threshold or
 - all darker than intensity(C) + threshold.
- Typically tests for 9 contiguous pixels in a 16 pixel circumference.
- Very fast detector (100 Mega pixel/second).

BRIEF descriptor (Binary Robust Independent Elementary Features)

- Goal: high speed.
- **Binary** descriptor formation
 - Smooth image,
 - For each detected keypoint (e.g. with FAST) sample 256 intensity pairs $(p_1^i, p_2^i), i = 1 256$ within a squared patch around keypoint.
 - Create an empty 256-element descriptor.
 - For each i-th pair:
 - * if $I_{p_1^i} < I_{p_2^i}$, then set i-th bit of descriptor to 1.
 - $\ast\,$ else to 0.
- The **pattern is generated randomly** (or by ML) only once: then, same pattern is used for all patches.
- Pros: **Binary Descriptor** allows very fast Hamming distance matching: count the number of bits that are different in the descriptors matched.
- Cons: Not scale/rotation invariant.

ORB descriptor (Oriented FAST and Rotated BRIEF)

- Keypoint detector based on FAST.
- BRIEF descriptors are steered(?) according to keypoint orientation (to provide rotation invariance).
- Good binary features are learned by minimizing the correlation on a set of training patches.

BRISK descriptor (Binary Robust Invariant Scalable Keypoints)

- Binary: formed by pairwise intensity comparisons.
- Pattern defines intensity comparisons in the keypoint neighborhood.
- **Red circles**: size of the smoothing kernel applied.
- Blue circles: smoothed pixel values used.
- Compare short- and long-distance pairs for orientation assignment and descriptor formation.
- Detector and descriptor speed: circa 10 times faster than SURF.
- Slower than BRIEF, but scale- and rotation- invariant.

Recap Table				
Detector	Descriptor that can be used	Localization Accuracy of the detector	Relocalization & Loop closing	Efficiency
Harris	Patch SIFT BRIEF ORB BRISK	++++	+ +++++ +++ ++++ +++	+++ + ++++ ++++ ++++
Shi-Tomasi	Patch SIFT BRIEF ORB BRISK	++++	+ ++++ +++ ++++ ++++	++ + ++++ ++++ ++++
FAST	Patch SIFT BRIEF ORB BRISK	++++	+++ +++++ +++ +++	+++++ + +++++ +++++
SIFT	SIFT	+++	++++	+
SURF	SURF	+++	++++	++

Figure 9: Recap for detectors and descriptors.

Things to remember

- Similarity metrics: NCC (ZNCC), SSD (ZSSD), SAD (ZSAD), Census Transform.
- Point feature detection:
 - Properties and invariance to transformations: Challenges are rotation, scale, view-point and illumination changes.
 - Extraction:

- * Moravec
- * Harris and Shi-Tomasi: rotation invariance.
- Automatic Scale Detection
- Descriptor
 - * Intensity patches:
 - \cdot Canonical representation: how to make them invariant to transformations: rotation, scale, illumination and view point (affine).
 - $\ast\,$ Better solution: Histogram of oriented gradients: SIFT descriptor.
- Matching:
 - * (Z)SSD, SAD, NCC, Hamming distance (last one only for binary descriptors), ration first/second closest descriptor.
- Depending on the task, you may want to trade off repeatability and robustness for speed: approximated solutions, combinations of efficient detectors and descriptors.
 - $\ast\,$ Fast corner detector,
 - * Keypoint descriptors faster than SIFT: SURF, BRIEF, ORB, BRISK.

Lecture 07: Multiple View Geometry 1

We have different problem statements:

• 3D reconstruction from multiple views

- Assumption: K, T, R are known.
- Goal: Recover the 3D structure from images.
- Structure from motion
 - Assumption: K, T, R are unknown.
 - **Goal:** Recover simultaneously 3D scene structure and camera poses (up to scale) from multiple images.

For a 2-view geometry, we can define the same problems as

- Depth from stereo (stereo vision)
 - Assumption: K, T, R are known.
 - Goal: Recover the 3D structure from images.
- 2-view structure from motion
 - Assumption: K, T, R are unknown.
 - Goal: Recover simultaneously 3D scene structure and camera poses (up to scale) and intrinsic parameters from two different views of the scene.

Depth from stereo

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

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 let us perceive them as straight. Radial distortion is removed. This process is also known as **rectification**. The distance of two seen images is called **disparity** (allows us to perceive the depth, i.e. the **smaller** the disparity, the **farther** the object). An application of this concept are **stereograms**.

Another applications is stereo photography, stereo viewers: Take two pictures of the same subject from two different viewpoints and display them so that each eye sees only one of the images.

Stereo Vision: basics

The basic principle behind stereo vision is **triangulation**.

- Gives reconstruction as intersection of two rays.
- Requires camera pose (calibration) and point correspondence (matching pairing points of the two images which are the projection of the same point in the scene).

There are basically two cases

- 1. Simplified case: identical cameras are **aligned**.
- 2. General case: different cameras are **not aligned**.

Simplified Case



Figure 10: Simplified case for stereo vision.

The two cameras are identical, meaning that they have the same **focal length**, and are **aligned** with the x-axis. If we have a world point P_w , a distance from the axis to the point Z_P , a distance between the cameras b and focal length f, we can use similar triangles from Figure 10 and get

$$\frac{f}{Z_P} = \frac{u_l}{X_P}$$

$$\frac{f}{Z_P} = \frac{-u_r}{b - X_P}$$

$$\Rightarrow X_P = \frac{u_l \cdot b}{u_l - u_r}$$

$$\Rightarrow Z_P = \frac{b \cdot f}{u_l - u_r}.$$
(67)

The difference $u_l - u_r$ is called **disparity**: difference in image location of the projection of a 3D point on two image planes. (QUESTIONS) Observation from this equation are

- Distance is inversely proportional to disparity: the distance to near objects can be measured more accurately than that to distant objects.
- Disparity is proportional to b. For a given disparity error, the accuracy of depth estimate increases with increasing baseline b.
- As b is increased, since the distance of the two cameras is increased, some objects may appear in one camera but not in the other (field of view of cameras). These objects won't have disparity.
- If the baseline b is unknown it is possible to reconstruct the scene up to a scale (structure from motion!)

What is the optimal baseline?

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

General Case: triangulation

In reality no cameras are identical and aligning both cameras on a horizontal axis is impossible. Why?

- There will always be differences in focal lengths due to manufacturing.
- The internal orientation of the CCD (elements which accumulate charge) in the camera package is unknown, theoretically orientated, but not in reality.
- . In order to be able to use a stereo camera, we need to compute
 - The extrinsic parameters (relative rotation and translation)
 - the intrinsic parameters (focal length, optical center, radial distortion of each camera).

In order to do this we use calibration methods (Tsai or homographies). How do we get the **relative pose?** Lines do not always interset in the 3D space: we want to minimize the error. For the two cameras we have

$$\tilde{p}_l = \lambda_l \cdot \begin{pmatrix} u_l \\ v_l \\ 1 \end{pmatrix} = K_l \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix}, \quad \tilde{p}_r = \lambda_r \cdot \begin{pmatrix} u_r \\ v_r \\ 1 \end{pmatrix} = K_r \cdot R \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix} + T.$$
(68)

Triangulation: The problem of determining the 3D position of a point given a set of corresponding image locations and known camera poses.

In order to triangulate, we use **least-squares** approximation:

$$\lambda_{1} \cdot \begin{pmatrix} u_{1} \\ v_{1} \\ 1 \end{pmatrix} = K \cdot [I|0] \cdot \begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \\ 1 \end{pmatrix} \Rightarrow \lambda_{1}p_{1} = M_{1} \cdot P \quad \text{left camera.}$$

$$\lambda_{2} \cdot \begin{pmatrix} u_{2} \\ v_{2} \\ 1 \end{pmatrix} = K \cdot [R|T] \cdot \begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \\ 1 \end{pmatrix} \Rightarrow \lambda_{2}p_{2} = M_{2} \cdot P \quad \text{right camera.}$$
(69)

We solve for P and get a system $A \cdot \begin{pmatrix} \lambda_1 \\ \lambda_2 \end{pmatrix} = b$, which cannot be inverted (A is 3×2). We use pseudoinverse approximation (least squares) and get

$$A^{T} \cdot A \cdot \begin{pmatrix} \lambda_{1} \\ \lambda_{2} \end{pmatrix} = A^{T} \cdot b \Rightarrow \begin{pmatrix} \lambda_{1} \\ \lambda_{2} \end{pmatrix} = (A^{T} \cdot A)^{-1} \cdot A^{T} \cdot b.$$
(70)

Remark. This is a problem with 6 equations and 5 unknowns: the 3 element of the coordinates of the world point and the two depth factors λ_1, λ_2 .

Interpretation: Given the projections $p_{1,2}$ of a 3D point P in two or more images, we want to find the coordinates of the 3D point by intersecting the two rays corresponding to the projections. We want to find the shortest segment connecting the two viewing rays and let P be the **midpoint** of the segment. The two rays won't meet exactly because of **noise** and **numerical errors**.

Triangulation: nonlinear approach

We want to find P that minimizes the sum of squared reprojection error

$$SSRE = d^{2}(p_{1}, \pi_{1}(P)) + d^{2}(p_{2}, \pi_{2}(P)),$$
(71)

where

$$d(p_1, \pi_1(P)) = ||p_1 - \pi_1(P)||$$
(72)

is called **reprojection error**. The observed point is p_1, p_2 and the reprojected one is M_1P, M_2P . In practice, this is done by initializing P using linear approach and then minimize SSRE using Gauss-Newton of Levenberg-Marquardt.

Correspondence Problem

Given a point p in a first image, where is its corresponding point p' in the right image? **Correspondence Search**: Identify image patches in the left and in the right images, corresponding to the same scene structure. Similarity measures:

- (Z)ZNCC
- (Z)SSD
- (Z)SAD
• Census Transform

Problem: Exhaustive image search can be computationally very expensive! Can we do that in 1D?

 \rightarrow 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. Since $\pi(p) = \lambda K$, $\pi(p)^{-1} = \lambda K^{-1}p$. This makes sense: if we observe a projection p_1 in the left camera, this can correspond to each world point lying on the infinite ray (every one of these points would project into p_1). All these points have a different projection on the right camera, which in 2D forms the epipolar line.
- The **epipole** is the projection of the optical center on the other camera image.
- A stereo camera has two epipoles!



Figure 11: Epipolar lines and epipoles.

The Epipolar Constraint

- The epipolar plane is uniquely defined by the two optical centers C_l, C_r and one image point p.
- The **epipolar constraint** constraints the location, in the second view, of the corresponding point to a given point in the first view.
- \Rightarrow This reduces the search to 1D problem along conjugate epipolar lines.

Example: converging cameras

Important: All epipolar lines intersect at the epipole! As the position of the 3D point varies, the epipolar line *rotates* about the baseline.

Example: identical and horizontally-aligned cameras

e and e^\prime at infinity

Example: forward motion

Epipoles have the same coordinates in both images: points move along lines radiating from e: Focus of expansion

Stereo Rectification

- Even in commercial stereo cameras, left and right images are not aligned.
- It is convenient if image scanlines are the epipolar lines
- Stereo rectification warps left and right images into new *rectified* images, whose epipolar lines are aligned to the baseline.



Figure 12: Stereo Rectification.

- Reprojects image planes onto a common plane parallel to the baseline.
- It works by computing two homographies, one for each input image reprojection.
- \Rightarrow Then, scanlines are **aligned** and epipolar lines are **horizontal**.

Idea: we define two new Perspective Projection Matrices obtained by rotating the old ones around their optical centers, until focal planes become coplanar (containing the baseline). This ensures **epipoles at infinity**. In order to have horizontal epipolar lines, the baseline should be *parallel* to the new X axis of both cameras. Moreover, corresponding points should have **the same vertical coordinate**. This is obtained by having same intrinsic parameters for the new cameras. Since the the focal length is the same, the new image planes are coplanar. PPMs are the same as the old cameras, whereas the new orientation (the same for both cameras) differs from the old ones by suitable rotations. For both cameras, **intrinsic parameters are the same**.

Implementation

1. The perspective equation for a point in the world is

image point =
$$\tilde{p} = \begin{pmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{pmatrix} = \lambda \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K[R|T] \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}.$$
 (73)

This can be rewritten in a more convenient way by considering [R|T] as the transformation from the world to the Camera frame (T expressed as C):

$$\lambda \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K \cdot R^{-1} \cdot \left(\begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix} - C \right)$$
(74)

2. We can then write the Perspective Equation for the Left and Right cameras: we assume for generality that they have the same intrinsic parameters:

$$\lambda_{L} \cdot \begin{pmatrix} u_{L} \\ v_{L} \\ 1 \end{pmatrix} = K_{L} \cdot R_{L}^{-1} \cdot \begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{L}$$
 left camera
$$\lambda_{R} \cdot \begin{pmatrix} u_{R} \\ v_{R} \\ 1 \end{pmatrix} = K_{R} \cdot R_{R}^{-1} \cdot \begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{R}$$
right camera (75)

3. The goal of stereo rectification is to warp left and right camera images such that their focal planes are **coplanar** and the intrinsic parameters are identical. It follows

$$\lambda_{L} \cdot \begin{pmatrix} u_{L} \\ v_{L} \\ 1 \end{pmatrix} = K_{L} \cdot R_{L}^{-1} \cdot \left(\begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{L} \right) \rightarrow \overline{\lambda}_{L} \cdot \begin{pmatrix} \overline{u}_{L} \\ \overline{v}_{L} \\ 1 \end{pmatrix} = \overline{K} \cdot \overline{R}^{-1} \cdot \left(\begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{L} \right)$$
$$\lambda_{R} \cdot \begin{pmatrix} u_{R} \\ v_{R} \\ 1 \end{pmatrix} = K_{R} \cdot R_{R}^{-1} \cdot \left(\begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{R} \right) \rightarrow \overline{\lambda}_{R} \cdot \begin{pmatrix} \overline{u}_{R} \\ \overline{v}_{R} \\ 1 \end{pmatrix} = \overline{K} \cdot \overline{R}^{-1} \cdot \left(\begin{pmatrix} X_{w} \\ Y_{w} \\ Z_{w} \end{pmatrix} - C_{R} \right)$$
(76)

4. By solving for $\begin{pmatrix} X_w \\ Y_w \\ Z_w \end{pmatrix}$ for each camera, we can compute the Homography (or warping) that needs to be applied to rectify each camera image:

$$\overline{\lambda}_{L} \cdot \begin{pmatrix} \overline{u}_{L} \\ \overline{v}_{L} \\ 1 \end{pmatrix} = \lambda_{L} \cdot \underbrace{\overline{K} \cdot \overline{R}^{-1} \cdot R_{L} \cdot K_{L}^{-1}}_{\text{homography left camera}} \cdot \begin{pmatrix} u_{L} \\ v_{L} \\ 1 \end{pmatrix}$$

$$\overline{\lambda}_{R} \cdot \begin{pmatrix} \overline{u}_{R} \\ \overline{v}_{R} \\ 1 \end{pmatrix} = \lambda_{R} \cdot \underbrace{\overline{K} \cdot \overline{R}^{-1} \cdot R_{R} \cdot K_{R}^{-1}}_{\text{homography right camera}} \cdot \begin{pmatrix} u_{R} \\ v_{R} \\ 1 \end{pmatrix}$$
(77)

5. The new $\overline{K}, \overline{R}$ can be chosen as

$$\overline{K} = \frac{K_L + K_R}{2}$$

$$\overline{R} = [\overline{r}_1, \overline{r}_2, \overline{r}_3],$$
(78)

where $\overline{r}_1, \overline{r}_2, \overline{r}_3$ are the column vectors of \overline{R} . These can be computed as

$$\overline{r}_1 = \frac{C_2 - C_1}{||C_2 - C_1||}$$

$$\overline{r}_2 = r_3 \times \overline{r_1} \text{ where } r_3 \text{ is the third column of } R_L$$

$$\overline{r}_3 = \overline{r}_1 \times \overline{r}_2.$$
(79)

For more details have a look at A compact alg. for rectification of stereo pairs.

Example: First, you remove the radial distortion (use e.g. bilinear interpolation). Then, compute homographies and rectify (bilinear interpolation).

If the images are rectified, we can perform the correspondence search along the same scanlines!

In order to solve the correspondence problem, we want to average noise effects. One can use a window around the point of interest and use similarity measures to find neighborhoods (these should be similar in intensity patterns).

Correlation-based window matching

- Problem: textureless regions (**aperture problem**), high ambiguity! **Solution:** increase window size to distinguish!
 - Smaller window: more detail but more noise!
 - Larger window: smoother disparity maps but less detail

Disparity Map

A disparity map apear as a grayscale image where the intensity of every pixel point is proportional to the disparity of that pixel in the left and right image: objects that are closer to the camera appear lighter, whiler farther objects appear darker. Input to dense 3D reconstruction

- 1. For each pixel in the left image, find its corresponding point in the right image.
- 2. Compute the disparity for each pair of correspondences.
- 3. Visualized in gray-scale or color coded image.

Close objects experience bigger disparity. They appear brighter in disparity map. The depth Z can be computed from the disparity by recalling that

$$Z_P = \frac{bf}{u_l - u_r} \tag{80}$$

This is really useful for obstacle avoidance. Challenges include: occlusion, repetition, nonlambertian surfaces (specularities), textureless surfaces. This is important for obstacle avoidance.

Improvements:

Multiple matches could satisfy the epipolar constraint. A good way to address the problem would be to have

- Uniqueness: only one match in right image for every point in left image.
- Ordering: points on same surface will be in same order in both views
- Disparity gradient: disparity changes smoothly between points on the same surface.

Sparse Stereo Correspondence restrict search to sparse set.

Lecture 08 - Multiple View Geometry 2

Two-view Structure from Motion

The camera relative pose is unknown: this is e.g. the case when the two images are taken from the same camera but at different times and positions.

Problem formulation: Given *n* point correspondences between two images, $\{p_1^i = (u_1^i, v_1^i), p_2^i = (u_2^i, v_2^i)\}$, simultaneously estimate the 3D points P^i , the camera relativemotion parameters (R, T), and the camera intrinsics K_1, K_2 that satisfy:

$$\lambda_{1} \cdot \begin{pmatrix} u_{1}^{i} \\ v_{1}^{i} \\ 1 \end{pmatrix} = K_{1} \cdot [I|0] \cdot \begin{pmatrix} X_{w}^{i} \\ Y_{w}^{i} \\ Z_{w}^{i} \\ 1 \end{pmatrix},$$

$$\lambda_{2} \cdot \begin{pmatrix} u_{2}^{i} \\ v_{2}^{i} \\ 1 \end{pmatrix} = K_{2} \cdot [R|T] \cdot \begin{pmatrix} X_{w}^{i} \\ Y_{w}^{i} \\ Z_{w}^{i} \\ 1 \end{pmatrix}$$
(81)

We have two cases then:

Calibrated Cameras $(K_1, K_2 \text{ known})$

For convenience, we use normalized image coordinates

$$\begin{pmatrix} \bar{u} \\ \bar{v} \\ 1 \end{pmatrix} = K^{-1} \cdot \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}.$$
(82)

We want to find R, T, P^i which satisfy

$$\lambda_{1} \cdot \begin{pmatrix} \bar{u}_{1}^{i} \\ \bar{v}_{1}^{i} \\ 1 \end{pmatrix} = K_{1} \cdot [I|0] \cdot \begin{pmatrix} X_{w}^{i} \\ Y_{w}^{i} \\ Z_{w}^{i} \\ 1 \end{pmatrix},$$

$$\lambda_{2} \cdot \begin{pmatrix} \bar{u}_{2}^{i} \\ \bar{v}_{2}^{i} \\ 1 \end{pmatrix} = K_{2} \cdot [R|T] \cdot \begin{pmatrix} X_{w}^{i} \\ Y_{w}^{i} \\ Z_{w}^{i} \\ 1 \end{pmatrix}.$$
(83)

Scale Ambiguity: If we rescale the entire scene by a constant factor (i.e. similarity transformation), the projections (in pixels) of the scene points in both images remain the same (because the angles remain the same).

- In *monocular* vision it is **not possible** to recover the absolute scale of the scene.
- In stereo vision, only 5 degrees of freedom are measurable:
 - 3 parameters to describe the **rotation**.
 - 2 parameters for the translation up to a scale (we can only compute the direction of translation but not its length (magnitude)).

How many knowns and unknowns?

- 4n knowns: n correspondences, each one (u_1^i, v_1^i) and (u_2^i, v_2^i) , i = 1, ..., n.
- 5 + 3n unknowns: 5 for the motion up to a scale (3 rotation and 2 translation) and 3n which is the number of coordinates of the n 3D points.

It should hold

$$4n \ge 5 + 3n$$

$$\Rightarrow n \ge 5.$$
 (84)

The first analytical solution for 5 points was given by Kruppa in 1913 (10 degree order polynomial, up to 10 solution with complex ones).

Let's define the **cross product** as a matrix multiplication

$$a \times b = \begin{pmatrix} 0 & -a_z & a_y \\ a_z & 0 & -a_x \\ -a_y & a_x & 0 \end{pmatrix} \cdot \begin{pmatrix} b_x \\ b_y \\ b_z \end{pmatrix} = [a]_x \cdot b$$
(85)

Epipolar Geometry

$$\bar{p}_1 = \begin{pmatrix} \bar{u}_1 \\ \bar{v}_1 \\ 1 \end{pmatrix}, \quad \bar{p}_2 = \begin{pmatrix} \bar{u}_2 \\ \bar{v}_2 \\ 1 \end{pmatrix}$$
(86)

We can observe that p_1, p_2, T are coplanar:

$$p_2^T \cdot n = 0$$

$$p_2^T \cdot (T \times p_1') = 0$$

$$p_2^T \cdot (T \times (Rp_1)) = 0$$

$$p_2^T \cdot [T]_x \cdot Rp_1 = 0$$

$$p_2^T \cdot E \cdot p_1 = 0 \text{ is the epipolar constraint, where } E = [T]_x \cdot R \text{ is the essential matrix.}$$
(87)

This is also called the Longuet-Higgins equation. Applying the constraints results in *four* different solutions for R and T.



Figure 13: Epipolar constraint

How to compute the Essential Matrix

Kruppa's solution (with at least 5 correspondences) is not efficient. In 1996, Philipp proposed an iterative solution. In 2004, the first efficient non iterative solution was proposed. This uses **Groebner Decomposition**. The first popular solution uses 8 points and is called the **8 point algorithm or Longuet-Higgins algorithm** (still used in NASA rovers).

The 8-point Algorithm

The Essential matrix is defined by

$$\bar{p}_2^T \cdot E \cdot \bar{p}_1 = 0. \tag{88}$$

Each pair of point correspondences provides a linear equation. For n points we can write

This problem can be written as

$$Q \cdot \bar{E} = 0. \tag{90}$$

Two types of solution

- Minimal Solution
 - $-Q_{n\times 9}$ should have rank 8 to have unique (up to scale) non trivial solution \overline{E} .
 - Each point correspondence provides 1 independent equation.
 - Thus, 8 point correspondences are needed.

• Over-determined Solution

- -n > 8 points.
- A solution is to minimize $||Q \cdot \overline{E}||^2$ subject to the constraint $||\overline{E}||^2 = 1$. The solution is the eigenvector corresponding to the smallest eigenvalue of matrix $Q^T \cdot Q$.
- This can be solved with Singular Value Decomposition.
- **Degenerate Solution** if 3D points are coplanar. There is the 5 point algorithm which holds also for coplanar points.

Interpretation

With the algorithm we try to minimize the algebraic error

$$\sum_{i=1}^{N} (\bar{p}_2^{i^T} \cdot E \cdot \bar{p}_1^i)^2, \tag{91}$$

where

$$\bar{p}_2^T \cdot E \cdot \bar{p}_1 = ||\bar{p}_2|| \cdot ||E \cdot \bar{p}_1|| \cdot \cos(\theta)$$
(92)

which is not zero is p_1, p_2, T are not coplanar. When extracting solutions, four results are available. We look only at results with points in front of both cameras (Cheirality Constraint).

Uncalibrated Cameras $(K_1, K_2 \text{ unknown})$

It holds

$$\bar{p}_2^T \cdot E \cdot \bar{p}_1 = 0, \tag{93}$$

where

$$\begin{pmatrix} \bar{u}_{1}^{i} \\ \bar{v}_{1}^{i} \\ 1 \end{pmatrix} = K_{1}^{-1} \cdot \begin{pmatrix} u_{1}^{i} \\ v_{1}^{i} \\ 1 \end{pmatrix}, \quad \begin{pmatrix} \bar{u}_{2}^{i} \\ \bar{v}_{2}^{i} \\ 1 \end{pmatrix} = K_{2}^{-1} \cdot \begin{pmatrix} u_{2}^{i} \\ v_{2}^{i} \\ 1 \end{pmatrix}.$$
(94)

By rewriting the constraint, one obtains

$$\begin{pmatrix} u_2^i \\ v_2^i \\ 1 \end{pmatrix}^T \cdot K_2^{-T} \cdot E \cdot K_1^{-1} \cdot \begin{pmatrix} u_1^i \\ v_1^i \\ 1 \end{pmatrix} = 0$$

$$\begin{pmatrix} u_2^i \\ v_2^i \\ 1 \end{pmatrix}^T \cdot F \cdot \begin{pmatrix} u_1^i \\ v_1^i \\ 1 \end{pmatrix} = 0,$$

$$(95)$$

where F is the **fundamental matrix**, which can be computed as

$$F = K_2^{-T} \cdot E \cdot K_1^{-1} = K_2^{-T} \cdot [T]_x \cdot R \cdot K_1^{-1}.$$
(96)

The same 8-point algorithm can be used to compute the fundamental matrix:

There are **orders of magnitude** of difference, which leads to poor results with least-squares. How to solve?

Normalized 8-point algorithm

This estimates the Fundamental matrix on a set of **Normalized correspondences** (with better numerical properties) and then **unnormalizes** the result to obtain the fundamental matrix for the original given correspondences.

Idea: Transform image coordinates so that they are in the range $[-1, 1] \times [-1, 1]$. One way is to apply the following rescaling and shift A more popular is to rescale the two



Figure 14: Shift for normalized algorithm.

point sets such that the centroid of each is 0 and the mean standard deviation $\sqrt{2}$. This can be done for every point as follows

$$\hat{p}^i = \frac{\sqrt{2}}{\sigma} \cdot (p^i - \mu), \tag{98}$$

where

$$\mu = \frac{1}{N} \sum_{i=1}^{n} p^i \tag{99}$$

is the centroid of the set and $\sigma = \frac{1}{N} \sum_{i=1}^{n} ||p^{i} - \mu||^{2}$ is the mean standard deviation. This transformation can be expressed in matrix form

$$\hat{p}^{i} = \begin{pmatrix} \frac{\sqrt{2}}{\sigma} & 0 & -\frac{\sqrt{2}}{\sigma} \mu^{x} \\ 0 & \frac{\sqrt{2}}{\sigma} & -\frac{\sqrt{2}}{\sigma} \mu^{y} \\ 0 & 0 & 1 \end{pmatrix} \cdot p^{i}.$$
(100)

The algoritm at the end reads

- 1. Normalize point correspondences: $\hat{p}_1 = B_1 \cdot p_1$, $\hat{p}_2 = B_2 \cdot p_2$.
- 2. Estimate \hat{F} using normalized coordinates \hat{p}_1, \hat{p}_2 .
- 3. Compute F from \hat{F} :

$$\hat{p}_{2}^{T} \cdot F \cdot \hat{p}_{1} = 0$$

$$p_{2}^{T} \cdot B_{2}^{T} \cdot \hat{F} \cdot B_{1} \cdot p_{1} = 0$$

$$\Rightarrow F = B_{2}^{T} \cdot \hat{F} \cdot B_{1}$$
(101)

Error Measures

The quality of the estimated fundamental matrix can be measured by looking at cost functions. The first is defined using the Epipolar Constraint

$$err = \sum_{i=1}^{N} (\bar{p}_2^{i^T} \cdot E \cdot p_1^i)^2$$
(102)

This error will exactly be 0 if computed from 8 points. For more points not 0 because of image noise/outliers. Better methods are

Directional Error

Sum of the angular distances to the Epipolar plane: $err = \sum_{i} (\cos(\theta_i))^2$, where

$$\cos(\theta) = \left(\frac{p_2^T \cdot E \cdot p_1}{||p_2^T|| \cdot ||E \cdot p_1||}\right)$$
(103)

Epipolar Line Distance

Squared Epipolar-Line-to-point Distances

$$err = \sum_{i=1}^{N} d^2(p_1^i, l_1^i) + d^2(p_2^i, l_2^i).$$
(104)

Cheaper than reprojection error: does not require point triangulation!



Figure 15: Epipolar Line Distance.

Reprojection Error

Sum of the Squared Reprojection Errors

$$err = \sum_{i=1}^{N} ||p_1^i - \pi_1(P^i)||^2 + ||p_2^i - \pi_2(P^i, R, T)||^2$$
(105)

Computation is expensive because of point triangulation, but is the most accurate!



Figure 16: Reprojection Error.

Robust Structure From Motion

Matched points are usually contaminated by **outliers**. Causes for this are

- Change in view point and illumination
- Image noise
- Occlusions
- Blur

The task of removing them is for **Robust Estimation**. Since error is integrating over time, this increases and is really bad.

RANSAC (Random Sample Consensus)

Ransac is the standard method for **model fitting in the presence of outliers** (noise points or wrong data). It can be applied to all problems where the goal is to estimate parameters of a model from the data. An easy example is RANSAC for **line fitting**:

- 1. Select sample of 2 points at random.
- 2. Calculate model parameters that fit the data in the sample.
- 3. Calculate error function for each data point.
- 4. Select data that supports current hypothesis.
- 5. Repeat.
- 6. Select the set with the maximum number of inliers obtained within k iterations.

How many iterations? All pairwise combinations : $\frac{N \cdot (N-1)}{2}$. This is computationally unfeasible if N is too large. With a probabilistic approach, one can reduce this: Let w be the number of inliers/N, N be the total number of data points. We can think of w as

w = P(selecting an inlier-point out of the dataset). (106)

We assume that the 2 points necessary to estimate a line are selected independently, i.e.

$$w^{2} = P(\text{both selected points are inliers})$$

$$1 - w^{2} = P(\text{at least one of these two points is outlier})$$
(107)

Let k indicate the number of RANSAC iterations so far, then

$$(1 - w^2)^k = P(\text{RANSAC never selected two points both inliers})$$
 (108)

Let p be the probability of success:

$$1 - p = (1 - w^{2})^{k}$$

$$\Rightarrow k = \frac{\log(1 - p)}{\log(1 - w^{2})}.$$
(109)

Remark. Think of having p = 0.99 and w = 0.5, then k = 16, which is damatically fewer than all combinations. The number of points does not influence that!.

RANSAC applied to general model fitting is

- 1. Initial: let A be a set of N points.
- 2. Repeat.
- 3. Randomly select a sample of s points from A.
- 4. Fit a model from the s points.
- 5. Compute the distances of all other points from this model.
- 6. Construct the inlier set (i.e. count the number of points whose distance is < d).
- 7. Store these inliers.
- 8. Until maximum number of iterations k is reached.
- 9. The set with the maximum number of inliers is chosen as solution to the problem.

$$k = \frac{\log(1-p)}{\log(1-w^s)}.$$
(110)

In order to implement RANSAC for Structure From Motion (SFM), we need three key ingredients

- a) What's the **model** in SFM? \rightarrow the Essential Matrix (for calibrated cameras) or the Fundamental Matrix (for uncalibrated cameras). Alternatively, R and T.
- b) What's the **minimum number of points** to estimate the model? \rightarrow We know that 5 points is the theoretical minimum number of points. However, 8-point algorithm, then 8 is the minimum.
- c) How do we compute the **distance** of a point from the model. \rightarrow We can use the epipolar constraint to measure how well a point correspondence verifies the model E or F, respectively. However, the **Directional error**, the **Epipolar line distance**, or the **Reprojection error (even better)** are used.

- 1. Randomly select 8 point correspondences.
- 2. Fit the model to all other points and count the inliers.
- 3. Repeat from 1 for k times.

With s points

$$k = \frac{\log(1-p)}{\log(1-(1-\varepsilon)^{s})}$$
(111)

Remark. No 6 DOF estimation for the 2-point RANSAC. k increases exponentially with the fraction of outliers ε .

- As observed, k is exponential in the number of points s necessary to estimate the model. We can see that k increases exponentially with the fraction of outliers ε .
- The 8-point algorithm is extremely simple and was very successful; however it requires more than 1177 iterations.
- The 5-point algorithm only requires 145 iterations, but can return up to 10 solutions of E.

Can we use less than 5 points? With planar motion

Planar Motion

Planar motion is described by three parameters ϑ, φ, ρ

)

$$R = \begin{pmatrix} \cos(\theta) & -\sin(\theta) & 0\\ \sin(\theta) & \cos(\theta) & 0\\ 0 & 0 & 1 \end{pmatrix}, \qquad T = \begin{pmatrix} \rho \cos(\varphi)\\ \rho \sin(\varphi)\\ 0 \end{pmatrix}$$
(112)

Let's compute the Epipolar Geometry



Figure 17: Planar motion.

$$E = [T]_x \cdot R$$

$$= \begin{pmatrix} 0 & 0 & \rho \sin(\varphi) \\ 0 & 0 & -\rho \cos(\varphi) \\ -\rho \sin(\varphi) & \rho \cos(\varphi) & 0 \end{pmatrix} \cdot \begin{pmatrix} \cos(\theta) & -\sin(\theta) & 0 \\ \sin(\theta) & \cos(\theta) & 0 \\ 0 & 0 & 1 \end{pmatrix}$$
(113)
$$= \begin{pmatrix} 0 & 0 & \rho \sin(\varphi) \\ 0 & 0 & -\rho \cos(\varphi) \\ -\rho \sin(\varphi - \theta) & \rho \cos(\varphi - \theta) & 0 \end{pmatrix}.$$

E has 2 DoF (θ , φ), because ρ is the scale factor. Thus, 2 correspondences are sufficient to estimate them.

But: can we use less than 2 point correspondences? Yes, if we exploid wheeled vehicles with **non-holonomic** constraints. Wheeled vehicles like cars, follow locally-planar circular motion about the instantaneous Center of Rotation (ICR). Since $\varphi = \theta/2$, meaning that we have only 1 DoF. Only 1 point correspondence is needed. This is the smallest parametrization possible and results in the most efficient algorithm for removing outliers (Scaramuzza). This updates the problem to be

$$R = \begin{pmatrix} \cos(\theta) & -\sin(\theta) & 0\\ \sin(\theta) & \cos(\theta) & 0\\ 0 & 0 & 1 \end{pmatrix}, \qquad T = \begin{pmatrix} \rho \cos(\frac{\theta}{2})\\ \rho \sin(\frac{\theta}{2})\\ 0 \end{pmatrix}$$
(114)

and

$$E = [T]_x \cdot R$$

$$= \begin{pmatrix} 0 & 0 & \rho \sin(\frac{\theta}{2}) \\ 0 & 0 & -\rho \cos(\frac{\theta}{2}) \\ \rho \sin(\frac{\theta}{2}) & -\rho \cos(\frac{\theta}{2}) & 0 \end{pmatrix}.$$
(115)

With the Epipolar Geometry constraint leads to

$$\theta = -2\tan^{-1}\left(\frac{v_2 - v_1}{u_2 + u_1}\right).$$
(116)

Only one iteration: compute θ for every point correspondence. Up to 1000 Hz, 1-point RANSAC in only used to find the inliers. Motion is then estimated from them in 6DOF.

ICR



Example of Ackerman steering principle



Locally-planar circular motion

Figure 18: Non-holonomic.

Lecture 09 - Multiple View Geometry 3

Bundle Adjustement (BA)

Nonlinear, simultaneous refinement of structure and motion (i.e. R, T, P^i). It is used after linear estimation of R and T. This, computes R, T, P^i by minimizing the Sum of Squared Reprojection Errors:

$$(R, T, P^{i}) = \operatorname{argmin}_{R, T, P^{i}} \sum_{i=1}^{N} ||p_{1}^{i} - \pi_{1}(P^{i}, C_{1})||^{2} + ||p_{2}^{i} - \pi_{2}(P^{i}, C_{2})||, \qquad (117)$$

where C_1, C_2 are the **pose** of the camera in the **world frame**. This can be minimized using *Lavenberg-Marquardt* (more robust than Gauss-Newton to local minima). It is better to initialize it close to the minimum. Same for multiple views!

Hierarchical SFM

- 1. Extract and match features between nearby frames.
- 2. Identify clusters consisting of 3 nearby frames:
- 3. Compute SFM for the 3 frames:
 - Compute SFM between 1 and 2 and build pointcloud.
 - Merge 3rd view running 3-point RANSAC between point cloud and 3rd view.
- 4. Merge clusters pairwise and refine (BA) both structure and motion.

Example is building Rome in one day.

Sequential SFM

With n views. Also called Visual Odometry (VO).

- 1. Initialize structure and motion from 2 views (bootstrapping).
- 2. For each additional view:
 - Determine pose (localization).
 - Extend structure (i.e. extract and triangulate new features).
 - Refine both pose and structure (BA).

 $2\mathbf{D}$ to $2\mathbf{D}$ Motion from Image feature correspondences

- Both feature points f_{k-1} and f_k are specified in 2D.
- The minimal-case solution involves 5-point correspondences
- The solution is found by minimizing the reprojection error:

$$T_k = \begin{pmatrix} R_{k,k-1} & t_{k,k_1} \\ 0 & 1 \end{pmatrix} = \operatorname{argmin}_{T_k} \sum_i ||p_k^i - \hat{p}_{k-1}^i||^2.$$
(118)

Popular algorithms: 8-/5-point.

3D to 2D Motion from 3D structure and Image correspondences

- f_{k-1} is given in 3D, f_k in 2D.
- This problem is known as *camera resection* or PnP (perspective from n points).
- The minimal-case solution involves 3 correspondences (+1 for disambiguating the four solutions).
- The solution is found by minimizing the reprojection error:

$$T_k = \begin{pmatrix} R_{k,k-1} & t_{k,k_1} \\ 0 & 1 \end{pmatrix} = \operatorname{argmin}_{X^i,C_k} \sum_{i,k} ||p_k^i - g(X^i,C_k)||^2.$$
(119)

Popular algorithms: P3P.

3D to **3D** Motion from 3D-3D Point correspondences (point cloud registration).

- Both f_{k-1} and f_k are specified in 3D. To do this, it is necessary to triangulate 3D points (e.g. use a stereo camera).
- The minimal case-solution involves 3 non collinear correspondences
- The solution is found by minimizing the 3D-3D euclidean distance

$$T_{k} = \begin{pmatrix} R_{k,k-1} & t_{k,k_{1}} \\ 0 & 1 \end{pmatrix} = \operatorname{argmin}_{T_{k}} \sum_{i} ||\tilde{X}_{k}^{i} - T_{k} \cdot \tilde{X}_{k-1}^{i}||.$$
(120)

Popular algorithms: Arun 87, ICP, BA.

Case Study: Monocular Visual Odometry (one camera!)

Bootstrapping:

- Initialize structure and motion from 2 views: e.g. 8-point algorithm + RANSAC.
- Refine structure and motion (BA)
- How far should the frames be? If too small baseline, large depth uncertainty. If too large baseline, small depth uncertainty.
- → One way to avoid this consists of skipping frames until average uncertainty of the 3D points decreases below a certain threshold. The selected frames are called keyframes. In general

$$\frac{\text{keyframe distance}}{\text{average-depth}} > \text{threshold } (10 - 20\%).$$
(121)

Remember the picture with ellipses which describe the depth uncertainty. If baseline increases, the ellipses decrease in size.

Localization

• Compute camera pose from known 3D-to-2D feature correspondence.

- Extract correspondences by solving for R and t (K is known).

$$\lambda \cdot \begin{pmatrix} u \\ v \\ 1 \end{pmatrix} = K \cdot [R|T] \cdot \begin{pmatrix} X_w \\ Y_w \\ Z_w \\ 1 \end{pmatrix}$$
(122)

- What is the minimal number of required point correspondences
 - 6 for linear solution (DLT algorithm).
 - -3 for a non linear solution (P3P algorithm).
 - 3 point RANSAC.

Extend Structure

• Extract and triangulate new features.

By denoting the relative motion between adjacent keyframes as

$$T_k = \begin{pmatrix} R_{k,k-1} & t_{k,k_1} \\ 0 & 1 \end{pmatrix}, \tag{123}$$

we can concatenate transformations to find the full trajectory of the camera as

$$C_k = T_{k,k-1} \cdot C_{k-1} \tag{124}$$

A non-linear refinement (BA) over the last m poses (+visible structure) can be performed to get a more accurate estimate of the local trajectory.

Loop Closure Detection (i.e. Place Recognition)

- Relocalization problem: during VO, tracking can be lost (due to occlusions, low tecture, quick motion, illumination change).
- Solution is to re-localize camera pose and continue.
- Loop closing problem: when go back where you already have been:
 - Loop detection: to avoid map duplication (e.g. same crossing rotated)
 - Loop correction: to compensate the accumulated drift!
- In both cases places recognition is needed (lecture 12)

VO vs. Visual SLAM

- VO: Focus on incremental estimation/local consistency. VO sacrifies consistency for real-time performance, without need to take into account all previous history of the camera (as SLAM does).
- Visual VSLAM: Simultaneous Localizazion and mapping. Focus on globally consistent estimation. Practically VO + loop detection + graph optimization.

Feature-based Methods

- 1. Extract and match features (+RANSAC)
- 2. Minimize Reprojection Error:

$$T_{k,k-1} = \operatorname{argmin}_{T} \sum_{i} ||u_{i}' - \pi(p_{i})||_{\Sigma}^{2}$$
(125)

Good: Large frame-to-frame motions, accuracy and efficient optimization of SFM (BA).

Bad: Slow due to costly feature extraction and matching, matching outliers (RANSAC).

Direct Methods (all pixels)

1. Minimize photometric error:

$$T_{k,k-1} = \operatorname{argmin}_T \sum_i ||I_k(u_i') - I_{k-1}(u_i)||_{\sigma}^2,$$
(126)

where

$$u'_{i} = \pi(T \cdot (\pi^{-1}(u_{i}) \cdot d))$$
(127)

Good: All information in the image can be exploited. Increasing camera frame-rate reduces computational cost per frame.

Bad: Limited frame to frame motion. Joint optimization of dense structures and motion too expensive.

ORB-SLAM

- Feature based:
 - Fast corner + Oriented Rotated Brief descriptor.
 - Binary descriptor.
 - Very fast to compute and compare.
 - Minimizes reprojection error.
- Includes:
 - Loop closing.
 - Relocalization.
 - Final optimization.
- Real time: 30Hz

LSD-SLAM

- Direct based + Semi-dense formulation:
 - 3D geometry represented as semi dense depth maps.
 - minimizes photometric error.

- **Separately** optimizes poses and structures.
- Includes:
 - Loop closing.
 - Relocalization.
 - Final optimization.
- Real time: 30Hz

DSO

- Direct based + sparse formulation:
 - 3D geometry represented as sparse large gradients.
 - Minimizes photometric error.
 - Jointly optimizes poses and structures (sliding window).
 - Incorporate photometric correction to compensate exposure time change
- Real time: 30Hz

SVO

- Direct based :
 - Corners and edgelets.
 - Frame to frame motion estimation.
- Feature based :
 - Frame to Keyframe pose refinement.
- Mapping:
 - Probabilistic depth estimation.
 - Multi camera system
- $\bullet~400~{\rm fps}$ on i7 laptops, 100 fps on smartphone PC

Direct Methods: Dense, Semidense, Sparse: Dense and Semidense behave similarly. Dense is only useful if one has motion blur and defocus.

Lecture 10 - Dense 3D Reconstruction

For the 3D reconstruction from multiple views we assume that camera are calibrated

- intrinsically (K is known for each camera), and
- extrinsically (T and R between cameras are known, for instance, from SFM).

For the multi-view stereo, we have as:

Input: calibrated images from several viewpoints.

Output: 3D object dense reconstruction.

Recall: The two camera centers and the image point p determine the epipolar plane, which intersects each camera image plane in the epipolar lines. Since we use the epipolar constraint, corresponding points only need to be searched along epipolar lines.

Dense Reconstruction

We want to estimate the structure from a *dense* region of pixels (hence not only from corners). The workflow is:

- 1. Local methods: estimate depth for every pixel independently.
- 2. Global methods: refine the depth surface as a whole by enforcing smoothness constraint.

We use the **photometric error** (SSD): this is derived for every combination of the reference image and any further image. **IDEA**: optimal depth minizes the photometric error in all images as a function of the depth in the first image.

Aggregated Photometric Error

The Dense reconstruction requires establishing dense correspondences. These are computed basing on the photometric error (SSD between corresponding patches of intensity values (min patch size: 1×1 pixels). Pros and cons of large and small patches?

- Small window: Pro: more detail, Cons: more noise.
- Large window: Pro: smoother disparity maps, Cons: less detail.

Not all the pixels can be matched reliably, due to viewpoint changes, occlusions. We take advantage of *many* small baseline views, where *high quality* matching is possible. Important facts:

- Repetitive texture shows multiple minima.
- The aggregated photometric error for *flat regions* and *edges* parallel to the epipolar line show **flat valleys** (noise!).
- For distinctive features, the aggregated photometric error has one clear minimum.

Disparity Space Image (DSI)

For a given image point (u, v) and for discrete depth hypotheses d, the aggregate photometric error C(u, v, d) with respect to the reference image I_r can be stored in a volumetric 3D grid called the *Disparity Space Image (DSI)*, where each voxel (group of u, v, d) has value

$$C(u, v, d) = \sum_{k} \rho\left(\tilde{I}_{k}(u', v', d) - I_{r}(u, v)\right),$$
(128)

where $\tilde{I}_k(u', v', d)$ is the patch of intensity values in the k-th image centered on the pixel (u', v') corresponding to the patch $I_r(u, v)$ in the reference image I_r and epth hypothesis d. Furthermore ρ is the photometric error (SSD).

Solution to depth estimation problem

Is a function d(u, v) in the DSI that satisfies:

Minimum aggregated photometric error (i.e.
$$argmin_dC$$
)
AND (129)
Piecewise smooth (global methods)

Interpolating while not overfitting!

Global Methods: We formulate them in terms of energy minimization. The objective is to find a surface d(u, v) that minimizes a global energy

$$E(d) = \underbrace{E_d(d)}_{\text{data term}} + \underbrace{\lambda \cdot E_s(d)}_{\text{regularization term}}, \qquad (130)$$

where

$$E_d(d) = \sum_{(u,v)} C(u, v, d(u, v))$$
(131)

and

$$E_s(d) = \sum_{(u,v)} \rho_d(d(u,v) - d(u+1,v)) + \rho_d(d(u,v) - d(u,v+1)).$$
(132)

 ρ_d is a norm (e.g. the $L_{1,2}$ or Huber norm). λ controls the tradeoff data (regularization). What happens as λ increases? Higher **smoothing**!

Regularized depth maps

- The regularization term $E_s(d)$
 - Smooths non smooth surfaces (result of noisy measurements) as well as discontinuities.
 - Fills the holes.
- Popular assumption: discontinuities in intensity **coincide** with discontinuities in depth.
- Control **smoothness penalties** according to image gradient (discrete)

$$\rho_d(d(u,v) - d(u+1,v)) \cdot \rho_I(||I(u,v) - I(u+1,v)||)$$
(133)

• ρ_I is some monotically *decreasing* function of intensity differences: **lower** smoothness cost for **high intensity gradients** (if there are high intensity gradients, you don't want to smooth them as they are a crucial information in your image.

Choosing the stereo baseline

What is the optimal baseline?

- Too small: large depth error.
- Too large: difficult search problem.

A possible approach is **depth map fusion** (different depth maps with different perspectives gives a complete image).

GPGPU for Dense Reconstruction

General Purpose Computing on Graphics Processing Unit. Perform demanding calculations on the GPU instead of the CPU. We can run processes in **parallel** on thousands of cores (CPU is optimized for serial processing). More transistors for data processing.

- Fast pixel processing (ray tracing, draw textures, shaded triangles,..)
- Fast matrix/vector operations (transform vertices)
- Programmable (shading, bump mapping)
- Floating-point support (accurate computations)
- Deep learning.

And

- Image processing
 - Filtering and feature extractions (e.g. convolutions)
 - Warping (e.g. epipolar rectification, homography).

• Multiple-view geometry

- Search for dense correspondences (pixel wise operations, matrix and vector operations (epipolar geometry).
- Aggregated photometric error.
- Global Optimation
- Variational methods (i.e. regularization (smoothing)) (divergence computation)

Typically on consumer hardware: 1024 threads per multiprocessor, 30 multiprocessors: 30000 threads. CPU with 4 cores which supports 32 threads. High arithmetic intensity. Have a look at Scaramuzza work!

Lecture 11 - Tracking

Point Tracking

Problem: Given two images, estimate the motion of a pixel from image I_0 to image I_1 . Two approaches exist, depending on the amount of motion between the frames:

- Block-based methods
- Differential methods

Template Tracking Given two images, estimate the warping that defines the motion and/or the distortion of a template from image I_0 to image I_1 .

Block-based methods

- Search for the corresponding patch in a *neighborhood* around the point.
- Use SSD, SAD, NCC to search for corresponding patches in a local neighborhood of the point. The search region usially is a $D \times D$ squared patches. We have to perform $D \times D$ comparisons, computationally demanding.

Differential Methods

• Look at the local brightness changes at the **same** location. **NO patch shift** is performed. (centered in the same point!)

Spatial Coherency

We assume that all the pixels in the patch undergo the same motion (same u and v). Also, assume that the time interval between the two images I_0 and I_1 is small. We want to find the motion vector (u, v) that minimizes the Sum of Squared Differences (SSD). As we did for Harris, we look at the first order Taylor approximation of the sum:

$$SSD = \sum (I_0(x, y) - I_1(x + u, y + v))^2 \approx \sum (I_0(x, y) - I_1(x, y) - I_x \cdot u - I_y \cdot v)^2 = \sum (\Delta I - I_x \cdot u - I_y \cdot v)^2,$$
(134)

which is a simple quadratic function in two variables (u, v).

Motion Vector

To minimize the E, we differentiate with respect to (u, v) and equate to 0.

$$\frac{\partial E}{\partial u} = 0 \Rightarrow -2I_x \sum (\Delta I - I_x \cdot u - I_y \cdot v) = 0$$

$$\frac{\partial E}{\partial v} = 0 \Rightarrow -2I_y \sum (\Delta I - I_x \cdot u - I_y \cdot v) = 0$$

(135)

Linear system of two equations in two unknowns. We can write this in matrix form

$$\begin{pmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{pmatrix} \cdot \begin{pmatrix} u \\ v \end{pmatrix} = \begin{pmatrix} \sum I_x \cdot \Delta I \\ \sum I_y \cdot \Delta I \end{pmatrix}$$
$$\begin{pmatrix} u \\ v \end{pmatrix} = \underbrace{\left(\sum I_x I_x \sum_{i=1}^{-1} I_x I_y \right)}_{M} \cdot \begin{pmatrix} \sum I_x \cdot \Delta I \\ \sum I_y \cdot \Delta I \end{pmatrix}$$
(136)

These are not matrix products, but pixel-wise products! For M to be invertible, its determinant should be non 0. From the decomposition

$$M = R^{-1} \cdot \begin{pmatrix} \lambda_1 & 0\\ 0 & \lambda_2 \end{pmatrix} \cdot R, \tag{137}$$

we know that det(M) is non zero when its eigenvalues are large (i.e. not a flat region and not an edge. In practice, it should be a corner or in general contain **texture**.

Aperture Problem

If we look at local brightness changes through a small aperture, we cannot always determine the motion direction, because **infinite** motion directions (solutions) may exist. The solution is to **increase** the aperture size.

Application of Differential Methods: Optical Flow

Optical flow is the pattern of appearent motion of objects in a visual scene, caused by the relative motion between the observer (eye or camera) and the scene. It tracks the motion of every pixel between two consecutive frames. For each pixel, we compute

- The vector direction and,
- The vector length (amount of movement).

An issue could be the choice of the right patch size.

Block-based vs. Differential Methods

Block-based Methods: search for the corresponding patch in a neighborhood of the point to be tracked. The search region is usually a square of $n \times n$ pixels.

- +: robust to large motions
- -: Can be computationally expensive $(n \times n \text{ comparisons for a single point track})$

Differential Methods:

- +: Much more efficient than block-based methods. Thus, can be used to track the motion of every pixel in the image. It avoids searching in the neighborhood of the point by analyzing the **local intensity changes** of an image patch at a specific location (no search is performed).
- -: Works only for **small** motions (high frame rate). For larger motion, multiscale implementations are used, but are then expensive.

Transformations

• Translation: 2 DOF

$$W(x,p) = \begin{pmatrix} x+a_1\\ y+a_2 \end{pmatrix} = \begin{pmatrix} 1 & 0 & a_1\\ 0 & 1 & a_2 \end{pmatrix} \cdot \begin{pmatrix} x\\ y\\ 1 \end{pmatrix}.$$
 (138)

• Euclidean: 3 DOF

$$W(x,p) = \begin{pmatrix} x\cos(\alpha) - y\sin(\alpha) + a_1 \\ x\sin(\alpha) + y\cos(\alpha) + a_2 \end{pmatrix} = \begin{pmatrix} \cos(\alpha) & \sin(\alpha) & a_1 \\ \sin(\alpha) & \cos(\alpha) & a_2 \end{pmatrix} \cdot \begin{pmatrix} x \\ y \\ 1 \end{pmatrix}.$$
 (139)

• Affine: 6 DOF

$$W(x,p) = \begin{pmatrix} a_1x + a_3y + a_5\\ a_2x + a_4y + a_6 \end{pmatrix} = \begin{pmatrix} a_1 & a_3 & a_5\\ a_2 & a_4 & a_6 \end{pmatrix} \cdot \begin{pmatrix} x\\ y\\ 1 \end{pmatrix}.$$
 (140)

• Projective (homography): 8 DOF

$$x' = \frac{a_1 x + a_2 y + a_3}{a_7 x + a_8 y + 1},$$

$$y' = \frac{a_4 x + a_5 y + a_6}{a_7 x + a_8 y + 1}.$$
(141)

Recalling that the Jacobian of a function

$$F(x_1, x_2, \dots, x_n) = \begin{pmatrix} f_1(x_1, x_2, \dots, x_n) \\ \vdots \\ f_m(x_1, x_2, \dots, x_n) \end{pmatrix}$$
(142)

is

$$J(F) = \nabla F = \begin{pmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \vdots \\ \frac{\partial f_n}{\partial x_1} & \cdots & \frac{\partial f_n}{\partial x_n} \end{pmatrix}$$
(143)

Template Tracking

Follow a template image in a video sequence by estimating the warp. Basically, we track the **distortion** function.

Problem Formulation

Template Warping Given the template image T(x), take all the pixels from the template image and warp them using the function W(x, p) parameterized in terms of parameters p. The goal of template-based tracking is to find the set of warp parameters p such that

$$I(W(x,p)) = T(x).$$
 (144)

This is solved by determining p that minimizes the Sum of Squared Differences

$$E = SSD = \sum_{x \in T} \left[I(W(x, p)) - T(x) \right]^2.$$
(145)

Assumptions are:

- No errors in the template image boundaries: only the appearance of the object to be tracket appears in the template image.
- No occlusion: the entire template is visible in input image.
- Brightness consistency assumption: the intensity of the object appearance is always the same across different views.

Lucas-Kanade Tracker

Uses the Gauss-Newton method for minimization, i.e.

- Applies a first order approximation of the warp,
- Attempts to minimize the SSD iteratively.

Derivation: Starting from

$$E = SSD = \sum_{x \in T} \left[I(W(x, p)) - T(x) \right]^2,$$
(146)

we assume that an initial estimate of p is **known**. Then, we want to find the increment Δp that minimizes

$$\sum_{x \in T} \left[I(W(x, p + \Delta p)) - T(x) \right]^2.$$
 (147)

The first order Taylor approximation of the term in brackets reads

$$I(W(x, p + \Delta p)) \approx I(W(x, p)) + \nabla I \frac{\partial W}{\partial p} \Delta p.$$
(148)

By replacing that in the equation we get

$$E = \sum_{x \in T} \left[I(W(x, p)) + \nabla I \frac{\partial W}{\partial p} \Delta p - T(x) \right]^2.$$
(149)

In order to minimize it, we differentiate and equate to 0, i.e.

$$\frac{\partial E}{\partial \Delta p} = 0. \tag{150}$$

0.0

It holds

$$\frac{\partial E}{\partial \Delta p} = 0$$

$$2\sum_{x \in T} \left(\nabla I \frac{\partial W}{\partial p}\right)^T \left[I(W(x,p)) + \nabla I \frac{\partial W}{\partial p} \Delta p - T(x) \right] = 0$$
(151)
$$\Delta p = H^{-1} \sum_{x \in T} \left(\nabla I \frac{\delta W}{\delta p}\right)^T (T(x) - I(W(x,p))),$$

where $H = \sum_{x \in T} \left(\nabla I \frac{\partial W}{\partial p} \right)^T \left(\nabla I \frac{\partial W}{\partial p} \right)$ is the second moment matrix of the warped image (Hessian). So, the algorithm reads

- 1. Warp I(x) with W(x, p).
- 2. Compute the error.
- 3. Compute warped gradients, ∇I evaluated at W(x, p).
- 4. Evaluate the Jacobian of the warping $\frac{\partial W}{\partial p}$.
- 5. Compute the inverse Hessian H^{-1} .
- 6. Multiply steepest descend with error.
- 7. Comptue Δp .
- 8. Update parameters $p \leftarrow p + \Delta p$.
- 9. Repeat until $\Delta p < \varepsilon$.

The algorithm follows a predict-correct cycle. A prediction I(W(x, p)) of the warped image is computed from an initial estimate. The correction parameter Δp is computed. The larger the error, the larger the correction. Challenges are

- How to get initial p?
- If the initial estimate is too far, the linear approximation does not longer hold. Solution are **pyramidal implementations**
- Illumination changes, object deformations,
- Occlusions,
- A solution can be to update the template with newest image.

Coarse to Fine estimation - Pyramidal implementations: Because the small motion assumption, regular optical flow methods work bad if the object we are tracking moves a long distance. Building image pyramids for each image and doing optical flow on each layer of the pyramid (to get rid of small motion constraints). Decrease resolution to get smaller image.

Generalization of Lucas-Kanade The same concept (predict/correct) can be applied to tracking of 3D objects. In order to deal with wrong prediction, it can be implemented in a Particle-Filter fashion.

Tracking by detection of local image features

- 1. Keypoint detection and matching. Invariant to scale, rotation or perspective.
- 2. Geometric verification (RANSAC)

Issues are:

- How to segment the object to track from background?
- How to initialize the warping?
- How to handle occlusions?
- How to handle illumination changes and no modeled effects?

Lecture 12 - Recognition

Recognition applications are: large scale image retrieval, recognition for smartphones, face recognition (four basic types of feature detectors, white areas are subtracted from black ones, deep learning), technology to convert scanned docs to text, pedestrian recognition. Challenges are: intra class variations (how to detect ANY car), context and human experience.

There are essentially two schools of approaches:

- Model based: tries to fit a model (2D or 3D) using a set of corresponding features (SIFT + RANSAC). **Example:** Is this book present in the scene? Look for corresponding matches: if most of the book's keypoints are present in the scene, the book is present in the scene.
- Appearance based: the model is defined by a set of images representing the object (template matching for example). **Example:** The model of the object is simply an image (template). Shift the template over the image and compare (NCC, SSD). This works only if the object and the template are identical.

The main goal of object recognition is to **classify**. Say yes or no to presence of an oject or categorize it. Either with bounding box or full segmentation.

Detection via Classification: main idea

We need

- Obtain training data.
- Define features.
- Define classifier.

Consider all subwindows in an image and sample a multiple scales and positions. Make a decision per window.

Generalization: Machine Learning approach

Apply a prediction function to a feature representation of the image to get the desired output.

$$\underbrace{y}_{\text{output}} = \underbrace{f}_{\text{prediction function input: image features}} (\underbrace{x}_{\text{input: image features}}).$$
(152)

Training: Given a training set of labeled (with correct answers!) examples, estimate the prediction function by minimizing the prediction error on the set.

Testing: Apply f to a never-before-seen test example x and output the predicted value y = f(x).

Features could be blob features, image histograms or histograms of oriented gradients.

Classifiers: Nearest neighbor

Features are represented in the descriptor space. Brute force matching, compute distances. Important facts are

- No training required.
- All we need is a distance function for our inputs
- Problem: need to compute distances to all training examples.

Features could be blob features, image histograms or histograms of oriented gradients.

Classifiers: Linear

Find a linear function to separate the classes

$$f(x) = \operatorname{sgn}(wx + b). \tag{153}$$

Classifiers: Nonlinear

Find a nonlinear function to separate the classes.

How can one define a classifier? We need to cluster the training data, then we need a distance function to determine to which cluster the query image belongs to.

K-Means Clustering

This is an algorithm to partition n observations into k clusters in which each observation x belongs to the cluster S_i with centroid m_i . It minimizes the sum of squared Euclidean distances between points x and their nearest cluster centers m_i

$$D(X,M) = \sum_{i=1}^{k} \sum_{x \in S_i} (x - m_i)^2.$$
 (154)

The algorithm reads

- 1. Randomly initialize k cluster centers
- 2. Iterate untl convergence:
 - Assign each data point x_i to the nearest center m_i .
 - Recompute each cluster center as the mean of all points assigned to it.

Bag of Words

This is used e.g. for large-scale image retrieval.

Visual Place Recognition We want to find the most similar images of a query image in a database of N images. The complexity is $\frac{N^2M^2}{2}$ feature comparisons (assuming each image has M features. Each image must be compared with all other images. Only with M = N = 100, you get 50 million images. The solution is to use an **inverted file index**, which reduces the complexity to $N \times M$.

Inverted File Text

For text documents, an index is important. We want to find every image in which a feature occurs. How many SIFT or BRISK features exist ?

- SIFT: infinite.
- BRISK-128: $2^{128} = 2.4 \cdot 10^{38}$.

For this reason, we need to create visual words, then put them into a vocabulary. Basically, we collect images and we extract features. A visual word is the centroid of a cluster. We then cluster the descriptors with the different words.

An inverted file index lists all visual words in the vocabulary (extracted at training time). Each word points to a list of images from the all image Data Base, in which which that word appears. The DB grows as the robot navigates and collects new images. Voting Array: has as many cells as the images in the DB. Each word in the query image votes for an image.

Robust object/scene recognition

Visual Vocabulary discards the spatial relationships between features: two images with the same features shuffled around will return a 100% match when using only appearance information. This can be overcome with **geometric verification**: test the h most similar images to the query image for geometric consistency (5,8 point RANSAC) and retain the image with the smallest reprojection error and largest number of inliers. More words is better!

Lecture 13 - Visual Inertial Fusion

Pose Graph Optimization

So far we assumed that the transformations are between consecutive frames, but transformation can be computed also between non adjacent frames T_{ij} (e.g. when features from previous keyframes are still observed). They can be used as additional constraints to improve cameras poses by minimizing the following:

$$C_k = \operatorname{argmin}_{c_k} \sum_{i} \sum_{j} ||C_i - C_j \cdot T_{ij}||^2$$
(155)

- For efficiency, only the last m keyframes are used.
- Gauss-Newton or Levenber-Marquadt are typically used to minimize it. For large graphs, there are open source things.



Figure 19: Pose graph optimization.

Bundle Adjustment (BA)

This incorporates the knowledge of landmarks (3D points).

$$X^{i}, C_{k} = \operatorname{argmin}_{X^{i}, C_{k}} \sum_{i} \sum_{k} \rho\left(p_{k}^{i} - \pi(X^{i}, C_{k})\right).$$
(156)

Outliers are a problem, how can we penalize them? In order to penalize wrong matches, we can juse the Huber or Turkey cost.

$$\begin{aligned} \mathbf{Huber} \qquad \rho(x) &= \begin{cases} x^2, & \text{if } |x| \le k \\ k \cdot (2|x| - k) & \text{if } |x| \ge k \text{ linear} \end{cases} \end{aligned}$$

$$\begin{aligned} \mathbf{Tukey} \qquad \rho(x) &= \begin{cases} \alpha^2 & \text{if } |x| \ge \alpha \\ \alpha^2 \cdot \left(1 - (1 - (\frac{x}{\alpha})^2)^3\right) & \text{if } |x| \le \alpha. \end{cases} \end{aligned}$$

$$\end{aligned}$$

$$\begin{aligned} (157)$$

Bundle Adjustment vs Pose-graph Optimization

- BA is more precise than pose-graph optimization because it adds additional constraints (landmark constraints).
- But more costly: $O((qM + lN)^3)$ with M and N being the number of points and camera poses and q and l the number of parameters for points and camera poses. The Jacobian is cubic in q and l. Workarounds are

- A small window size limits the number of parameters for the optimization and thus makes real-time bundle adjustment possible.
- It is possible to reduce the computational complexity by just optimizing the camera parameters and keeping the 3D landmarks fixed, e.g. freeze the 3D points and adjust the poses



Figure 20: Tukey vs. Huber norm.



Figure 21: Bundle Adjustment.

IMUs

Inertial Mesaurement Unit. Measures angular velocity and linear accelerations.

- Mechanical: spring/damper system.
- Optical: Phase shift projected laser beams is proportional to angular velocity.
- MEMS (accelerometer): a spring-like structure connects the device to a seismic mass vibrating in a capacitive divider. A capacitive divider converts the displacement of the seismic mass into an electric signal. Damping is created by the gas sealed in the device.
- MEMS (gyroscopes): measure the Coriolis forces acting on MEMS vibrating structures. Their working principle is similar to the haltere of a fly. Have a look!

Why IMU?

- Monocular vision is scale ambiguous.
- Pure vision is not robust enough (Tesla accident):
 - Low texture.
 - High dynamic range.
 - High speed motion.

Why not just IMU? Why Vision?

Pure IMU integration will lead to large drift (especially cheap IMUs). Integration of angular velocity to get orientation: error **proportional to** t. Double integration to get position: if there is a bias in acceleration, the error of position is **proportional to** t^2 . The actually position error also depends on the error of orientation.

Why visual inertial fusion?

- Cameras
 - + Precise in slow motion.
 - + Rich information for other purposes
 - Limited output rate ($\sim 100Hz$)
 - Scale ambiguity in monocular setup.
 - Lack of robustness
- IMU
 - + Robust.
 - + High output rate ($\sim 1000Hz$).
 - + Accurate at high acceleration.
 - Large relative uncertainty when at low acceleration/angular velocity.
 - Ambiguity in gravity / acceleration.

Together, they can work for state estimation: loop detection and loop closure.

IMU: Measurement Model

$$\tilde{\omega}_{WB}^{B}(t) = \omega_{WB}^{B}(t) + b^{g}(t) + n^{g}(t)$$

$$\tilde{a}_{WB}^{B}(t) = R_{BW}(t) \cdot \left(a_{WB}^{W}(t) - g^{W}\right) + b^{a}(t) + n^{a}(t)$$
(158)

where g stands for Gyroscope and a for accelerometer. The noise is additive Gaussian white noise. The bias has own dynamics

$$\dot{b}(t) = \sigma_b \cdot w(t), \tag{159}$$

i.e. the derivative of the bias is white Gaussian noise (random walk). In discrete time, one writes

$$b[k] = b[k-1] + \sigma_{bd} \cdot w[k], \quad w[k] \sim \mathcal{N}(0,1), \quad \sigma_{bd} = \sigma_b \cdot \sqrt{t}$$
(160)

IMU biases

- Can be estimated,
- Can change due to temperature change, mechanical pressure,...
- Can change everytime the IMU is started.

Integration leads to

$$p_{Wt_2} = P_{Wt_1} + (t_2 - t_1)v_{Wt_1} + \int \int_{t_1}^{t_2} R_{Wt}(t) \left(\tilde{a}(t) - b^a(t) + g^w\right) dt^2, \quad (161)$$

which depends on initial position and velocity. The rotation R(t) can be computed with a giroscope.

Different Paradigms

Loosely Coupled Approach

Treats VO and IMU as two separate (not coupled black boxes). Each block estimates **pose and velocity** from visual and inertial data (pose and velocity up to a scale and inertial data in absolute scale).



Figure 22: Loosely Coupled Approach.

Tightly Coupled Approach

Makes use of the raw sensors' measurements: 2D features, IMU readings, more accurate, more implementation effort.

Filtering: Visual Inertial Formulation

System states are:

- Tightly Coupled: $X = (p_W(t); q_{WB}(t); v_W(t); b^a(t); b^g(t); L_{w,1}; \ldots; L_{w,K})$, with L Landmarks.
- Loosely Coupled $X = (p_W(t); q_{WB}(t); v_W(t); b^a(t); b^g(t))$



Figure 23: Tightly Coupled Approach.

Closed-form Solution (1D case)

The absolute pose x is known up to a scale s, thus

$$x = s\tilde{x}.\tag{162}$$

From the IMU we get

$$x = x_0 + v_0 \cdot (t_1 - t_0) + \int \int_{t_0}^{t_1} a(t) dt$$
(163)

By equating them we get

$$s\tilde{x} = x_0 + v_0 \cdot (t_1 - t_0) + \int \int_{t_0}^{t_1} a(t) dt$$
(164)

As shown, for 6DOF both s and v_0 can be determined from a single feature observation and 3 views. x_0 can be set to 0. It holds

$$s\tilde{x}_{1} = v_{0} \cdot (t_{1} - t_{0}) + \int \int_{t_{0}}^{t_{1}} a(t) dt$$

$$s\tilde{x}_{2} = v_{0} \cdot (t_{2} - t_{0}) + \int \int_{t_{0}}^{t_{2}} a(t) dt \qquad (165)$$

$$\Rightarrow \begin{pmatrix} \tilde{x}_{1} & (t_{0} - t_{1}) \\ \tilde{x}_{2} & (t_{0} - t_{2}) \end{pmatrix} \cdot \begin{pmatrix} s \\ v_{0} \end{pmatrix} = \begin{pmatrix} \int \int_{t_{0}}^{t_{1}} a(t) dt \\ \int \int_{t_{0}}^{t_{2}} a(t) dt \end{pmatrix}.$$

Closed-form Solution (general case)

Consider N feature observations and 6DOF case. Can be used to initialize filter and smoothers. One can show hat a linear system of equations can be achieved and solved using the pseudoinverse:

$$AX = S, (166)$$

where X is the vector of unknowns (3D point distances, absolute scle, initial velocity, gravity vector, biases). A and S contain 2D feature coordinates, acceleration, and angular velocity measurements.
Filtering	Fixed-lag Smoothing	Full smoothing	
Only updates the most recent states (e.g., extended Kalman filter)	Optimizes window of statesMarginalizationNonlinear least squares optimization	Optimize all statesNonlinear Least squares optimization	
×1 Linearization	✓ Re-Linearize	✓ Re-Linearize	
×Accumulation of linearization errors	*Accumulation of linearization errors	✓ Sparse Matrices✓ Highest Accuracy	
×Gaussian approximation of marginalized states	*Gaussian approximation of marginalized states		
✓Fastest	✓Fast	×Slow (but fast with GTSAM)	

Figure 24: Different Paradigms.

Different Paradigms

E.g. ROVIO, minimizes the photometric error instead of the reprojection error.

Filtering: Problems

- Wrong linearization point: linearization depends on the current estimates of states, which can be wrong.
- Complexity of the EKF grows quadratically in the number of landmarks. Few Landmarks are usually tracked to allow real time operation.
- Alternative: MSCKF: keeps a window of recent states and updates them using EKF. Incorporate visual observation without including point positions into the states.

Maximum A Posteriori (MAP) Estimation

Fusion solved as a non-linear optimization problem. Increased accuracy over filtering methods. We have

$$x_k = f(x_{k-1}), \quad z_k = h(x_{i_k}, l_{i_j}),$$
(167)

where X are the robot states, L the 3D points and Z the features and IMU measurements. It holds

$$\{X^*, L^*\} = \operatorname{argmax}_{X,L} P(X, L|Z)$$

= $\operatorname{argmin}_{X,L} \{\underbrace{\sum_{k=1}^{N} ||f(x_{k-1}) - x_k||_{\Lambda_k}^2}_{\text{IMU residuals}} + \underbrace{\sum_{i=1}^{M} ||h(x_{i_k}) - z_i||_{\Sigma_i}^2}_{\text{Reprojection residuals}}$ (168)

An open problem is consistency:

- Filters: Linearization around different values of the same variable may lead to error.
- Smoothing methods: may get stuck in local minima.

Camera-IMU calibration

Goal: Estimate the rigid body transformation T_{BC} and delay t_d between a camera and an IMU rigidly attached. Assume that the camera has already been intrinsically calibrated. **Data:** Image points of detected calibration pattern and IMU measurements (accelerometer and gyroscope).

Approach: Minimize a cost function

$$J(\theta) = J_{\text{feat}} + J_{\text{acc}} + J_{\text{gyro}} + J_{\text{bias}_{acc}} + J_{\text{bias}_{gyro}}$$
(169)

Event Based Vision

Feature based vs. Photometric (direct) methods

Feature based methods

- 1. Extract and match features (+RANSAC)
- 2. Minimize reprojection error.

They are

- +: Large frame-to-frame motions.
- -: slow due to costly feature extraction and matching.
- -: matching outliers.

Direct (photometric) errors are

- 1. The pixel is the feature to track.
- 2. Minimize photometric error.

They are

- +: All information in the image can be exploited (precision, robustness).
- +: Increasing camera frame rate reduces computational cost.
- -: Limited frame-to-frame motion.

Pure vision is not robust enough to low texture, HDR, high speed motion.

Accuracy, Efficiency, Robustness



Figure 25: SLAM research.

Event-based Cameras

Motivation: current flight maneuvers achieved with onboard cameras are still too slow compared with those attainable by birds. We need low latency sensors and algorithms. The average robot-vision algorithms have latencies of 50-200 ms, which puts a hard bound on the agility of the platform. Event cameras enable **low-latency sensory motor control** < 1ms.

Human Vision System

130 million photoreceptors (similar to pixels) but only 2 millions axons (wires that connect).

DVS

Advantages:

- Low latency (1 micro second)
- High dynamic range (140 dB instead of 60 dB)
- Low power: 10mW instead of 1W

Disadvantages:

- Paradigm shift: requires totally new vision algorithms
 - Asynchronus pixels,
 - No intensity information (only binary intensity changes).

A traditional camera outputs frames at fixed time intervals. By contrast, a DVS outputs asynchronous events at microsecond resolution. An event is generated each time a single pixel detects an intensity changes value:

event:
$$\langle t, \langle x, y \rangle, \operatorname{sign}\left(\frac{\mathrm{d}I(x, y)}{\mathrm{d}t}\right)\rangle$$
 (170)

All pixels are independent from another. Implements **level-crossing** sampling. Reacts to **logarithmic** brightness changes.

DVS Operating Principle

Each pixel is independent of all the other pixels. Events are generated everytime a single pixel sees a change of the logarithm of the brightness that is equal to C, i.e.

$$|\log(I)| = |\log(I(t + \Delta t) - \log(I(t)))| = C,$$
(171)

where $C \in [0.15, 0.20]$ is called **contrast sensitivity** and can be tuned by the user. Since brightness can be either positive or negative, we have ON event if = C and OFF event if = -C. Traditional sampling is performed with the discriminant (time) on x-axis. Level-crossing sampling works with the change in intensity, in the y-axis.



Figure 26: DVS circuit



Figure 27: DVS.

Applications

Low power monitoring, fast closed-loop contro, high dynamic range imaging, low power gesture recognition, high speed flow speed estimation.

DVS vs High speed cameras

Calibration of a DVS

The standard pinhole camera model is still valid (same optics). Standard passive calibration cannot be used: we would need to move the camera. **Blinking patterns** (computer screen, LEDs).

A simple optical flow algorithm: a moving edge

White pixels become black, i.e. the brightness decrease, i.e. negative events (black color). Events are represented by dots. At what speed is the edge moving? $v = \frac{\Delta x}{\Delta t}$.

How many events should be used?

Two different approaches

• Event-by-event processing (i.e. estimate the state event by event): Pros: low latency, Cons: with high speed motion, there are dozens of millions of events per seconds (GPU)

		<u>.</u>	
	Photron Fastcam SA5	Matrix Vision Bluefox	DVS
Max fps or measurement rate	1MHz	90 Hz	1MHz
Resolution at max fps	64x16 pixels	752x480 pixels	346x260 pixels
Bits per pixels	12 bits	8-10	1 bits
Weight	6.2 Kg	30 g	30 g
Active cooling	yes	No cooling	No cooling
Data rate	1.5 GB/s	32MB/s	~1MB/s on average
Power consumption	150 W + Ilighting	1.4 W	20 mW
Dynamic range	n.a.	60 dB	140 dB

Figure 28: DVS vs High Speed Cameras.

• Event-packet processing (i.e. process the last N events): Pros: N can be tuned to allow real-time performance on a CPU. Cons: no longer microsecond resolution (when is this really necessary=)

Event-by-event based Processing

Let's start with an approximation:

$$\Delta \log(I) = \frac{\partial \log(I)}{\delta t} \Delta t$$

= $\frac{1}{I} \frac{\delta I}{\delta t} \Delta t$
= $\frac{\partial I}{I}$. (172)

Claim 2. To simplify the notation, let's assume that $I(x, y, t) = \log(I(x, y, t))$. Consider a given pixel p(x, y) moving with apparent motion $\vec{u} = (u, v)$ (i.e. induced by a moving 3D patch). It can be shown, that an event is generated if the scalar product between the gradient and the appearent motion vector u is equal to C.

$$-\nabla I \cdot u = C \tag{173}$$

Proof. The proof comes from the brightness constancy assumption, which says that the intensity value of p, before and after the motion, must remained unchanged

$$I(x, y, t) = I(x + u, y + v, t + \Delta t)$$
(174)

By replacing the right-hand term by its first order approximation at $t + \Delta t$, we get

$$I(x, y, t) = I(x, y, t + \Delta t) + \frac{\partial I}{\partial x}u + \frac{\partial I}{\partial y}v$$

$$I(x, y, t + \Delta t) - I(x, y, t) = -\frac{\partial I}{\partial x}u - \frac{\partial I}{\partial y}v$$

$$\Rightarrow \Delta I = C = -\nabla \cdot u.$$
(175)

This equation described the linearized event generation equation for an event generated by a gradient ∇I that moved by a motion vector u (optical flow) during a time interval Δt . 1 Equation, 2 Unknowns, solution is to add events.

Case Study 1: Image Intensity Reconstruction

The intensity signal at the event time can be reconstructed by integration of $\pm C$. Given the events and the camera motion (rotation), recover the absolute brightness. **Explanation:** An event camera naturally responds to edges, hence, if we know the motion, we can relate the events to world coordinates to get an edge/gradient map. Then, just integrate the gradient map to get absolute intensity.

1. Recover the gradient map of the scene. Let $L = \log(I)$. Then

$$\Delta L(t) = L(t) - L(t - \Delta t) = C.$$
(176)

In terms of the brightness map M(x, y):

$$M(p_m(t)) - M(p_m(t - \Delta t)) \approx g \cdot v \cdot \Delta t, \qquad (177)$$

with $g = \nabla M(p_m(t))$.

2. Integrate the gradient to obtain brightness. Poisson reconstruction: integrate the gradient map g to get absolute brightness M.

Case Study 2: Event-based Corner Detection

FAST-like event-based corner detection: operates on surface of active events. The event is considered a corner if

- 3-6 contiguous pixels on red ring are newer than all other pixels on the same ring and,
- 4-6 contiguous pixels on blue ring are newer than all other pixels on the same ring

1.1 Event-packet based processing

EVO: parallel tracking and mapping in real-time. Tracking, 6DOF pose, Mapping, 3D Map. How does a 3D mapping works? An event camera reacts to strong gradients in the scene. Areas of high ray-density likely indicate the presence of 3D structures. The ray-density can be seen as tue Disparity Space Image (DSI). This is a projective sampling grid (with adaptive thresholding).



Figure 29: DVS vs High Speed Cameras.

DAVIS: Dynamic and Active-pixel Vision Sensor

Combines an event sensor (DVS) with a standard camera in the same pixel array. Output are frames (at 30 Hz) and events (asynchronous). One can them perform SLAM with an IMU, which increases robustness and accuracy.

Open problems for DVS are: noise modeling, asynchronous feature and object detection and tracking, sensor fusion, asynchronous learning and recognition, estimation and control, low power computation.