

VISUAL INERTIAL ODOMETRY (VO)

- WHY?

- Very accurate: relative position error is 0.1% - 2% travelled distance
- Not affected by wheel slippage
- Complementary to
 - wheel encoders (wheel odometry)
 - GPS (when degraded)
 - IMUs
 - Laser

- ASSUMPTIONS

- SUFFICIENT ILLUMINATION in the environment
- DOMINANCE OF STATIC SCENE over moving objects
- ENOUGH TEXTURE to allow apparent motion to be extracted
- Sufficient SCENE OVERLAP between consecutive frames

STRUCTURE FROM MOTION (SFM)

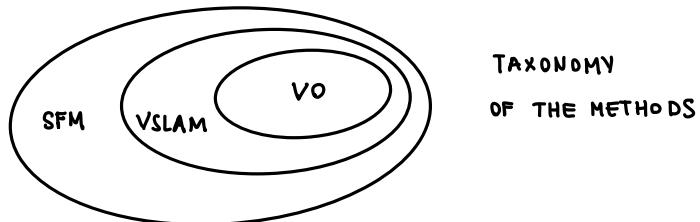
SFM is more general than VO and tackles the problem of 3D RECONSTRUCTION AND 6DOF POSE ESTIMATION from UNORDERED IMAGE SET

- VO is a particular case of SFM.
- VO focuses on estimating the 6DOF motion of the camera SEQUENTIALLY (as a new frame arrives) and in REAL TIME.
- Sometimes SFM is used as synonym of VO

VISUAL SIMULTANEOUS LOCALIZATION AND MAPPING (VSLAM)

SLAM = VO + LOOP DETECTION & CLOSURE

It guarantees global consistency: the estimated trajectory is globally correct from start to end



VO FLOW CHART

VO computes the camera path incrementally (pose after pose)

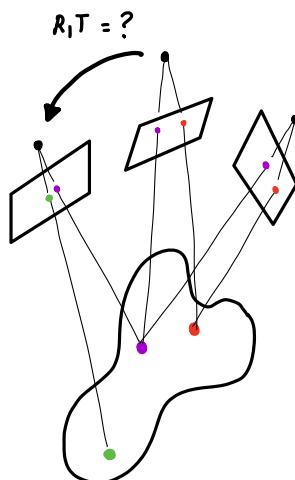
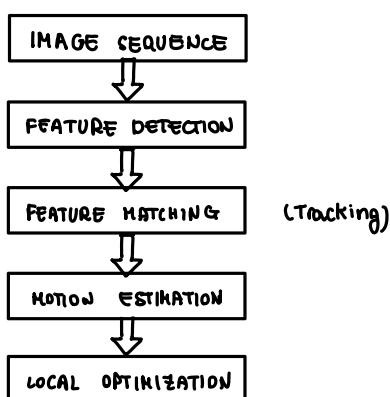
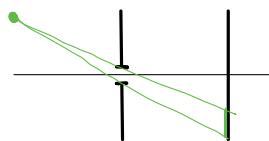


IMAGE FORMATION

THE PINHOLE CAMERA

In an **IDEAL PINHOLE**, only one ray of light reaches each point on the film



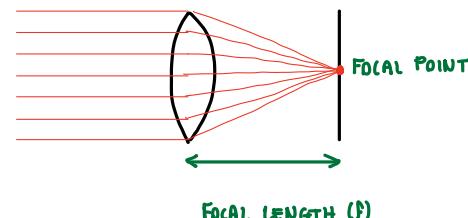
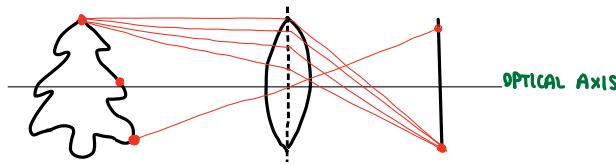
Shrinking the aperture helps but

- **LESS LIGHT** gets through (must increase the exposure)
- **DIFFRACTION EFFECTS** arise

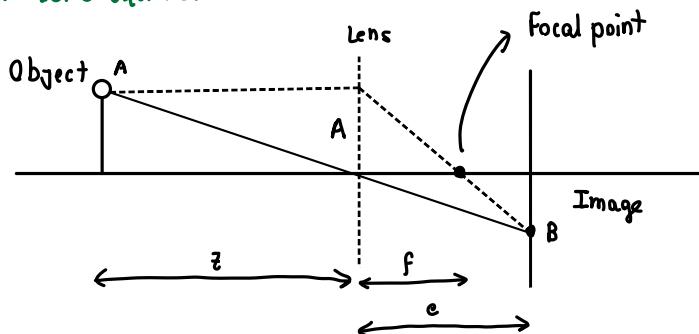
USING A CONVERGING LENS

We can use a lens that focuses light onto the film

- Ray passing through the **OPTICAL CENTER** are not deviated
- All rays parallel to the optical axis converge to the **FOCAL POINT**



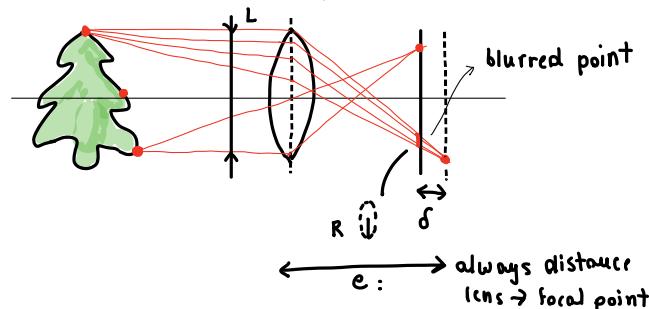
THIN LENS EQUATION



$$\left. \begin{aligned} \frac{B}{e} &= \frac{A}{z} \rightarrow \frac{B}{A} = \frac{e}{z} \\ \frac{B}{A} &= \frac{e-f}{f} = \frac{e}{f} - 1 \end{aligned} \right\} \frac{e}{f} - 1 = \frac{e}{z} \Rightarrow \frac{1}{f} - \frac{1}{e} = \frac{1}{z}$$

$$\boxed{\frac{1}{f} = \frac{1}{z} + \frac{1}{e}}$$

- Any object satisfying this equation is **IN FOCUS**
- Can we use this to measure distances?
- For a **FIXED FILM DISTANCE** e , there is a specific distance between the object and the lens (z) at which the object appears in "focus"
- Other points project to a "blur circle" in the image

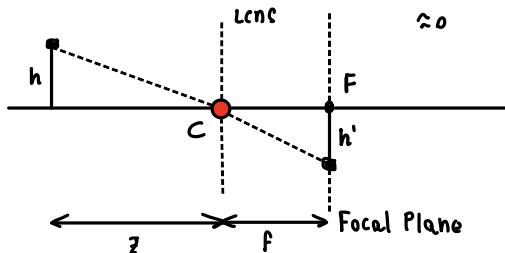


$$R = \frac{Ld}{2e}$$

- A small pinhole (L) gives a small blur (R)
- To capture a "good" image: adjust camera settings

THE PINHOLE APPROXIMATION

For distant objects $z \gg L, z \gg f$: $\frac{1}{f} = \frac{1}{z} + \frac{1}{e} \Rightarrow \frac{1}{f} \approx \frac{1}{e} \Rightarrow f \approx e$



C: Optical center or Center of Projection

$$\Rightarrow \frac{h'}{h} = \frac{f}{z} \Rightarrow h' = \frac{f}{z} h$$

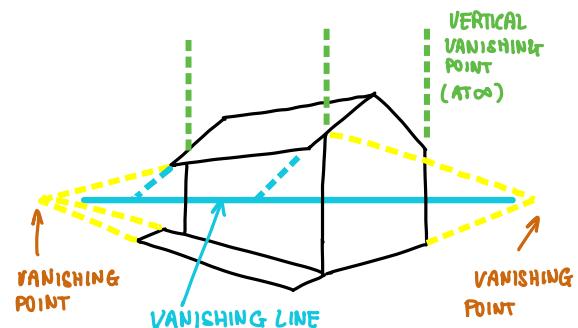
- The dependence of the image of an object on its depth (distance from camera) is known as **PERSPECTIVE**. (e.g. far objects appear small)

PERSPECTIVE

- straight lines are still straight
 - length is lost
 - angles are lost
- } parallelism & perpendicularity X

VANISHING POINTS AND LINES

- Parallel lines in the world intersect in the image at a **VANISHING POINT**
- Parallel planes in the world intersect in the image at a **VANISHING LINE**



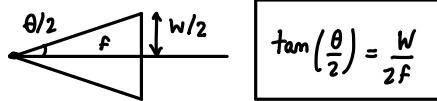
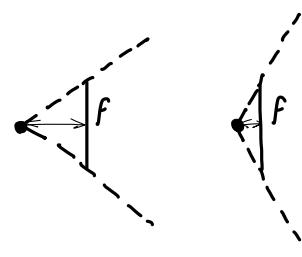
FOCUS AND DEPTH OF FIELD (DoF)

- The **DEPTH OF FIELD (DoF)** is the distance between the nearest and furthest object in the scene that appear acceptably sharp in an image
- Although sharpness is "perfect" at a single distance, in the DoF loss of sharpness is imperceptible
- A smaller aperture increases the DoF but reduces the amount of light into the camera

FIELD OF VIEW (FOV)

Angular measure of portion of 3D space seen by the camera

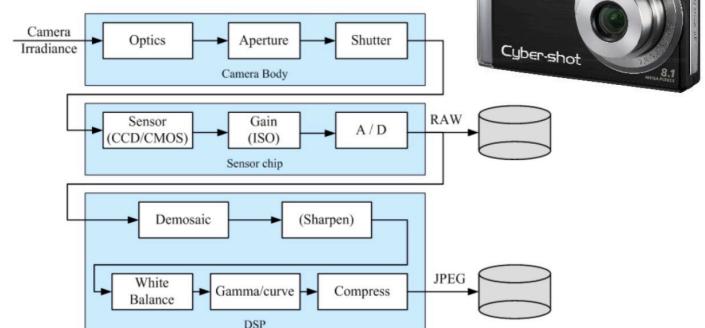
- As **f gets smaller**, image becomes more **wide** angle \rightarrow can see more points
- As **f gets larger**, image becomes more **narrow** angle \rightarrow can see less points



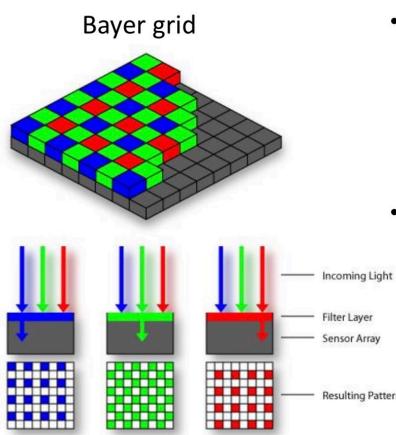
DIGITAL CAMERAS

The film is an array of CCD or CMOS light sensitive devices that convert photons into electrons

- Pixel intensity with 8 bits ranges between $[0, 255]$



COLOR SENSING IN DIGITAL CAMERAS



- The **BAYER PATTERN**:
 - 1/2 filter is **GREEN** → humans are much more sensitive to green
 - 1/4 filter is **RED**
 - 1/4 filter is **BLUE**
- For each pixel we estimate missing color component from neighboring values (**DEMOSAICING**)

TYPES OF SHUTTER

ROLLING SHUTTER:

- Rows of pixels are exposed at different times, one after the other
- Distort image for moving objects

GLOBAL SHUTTER:

- All pixels are exposed simultaneously
- Good for moving objects → no distortion

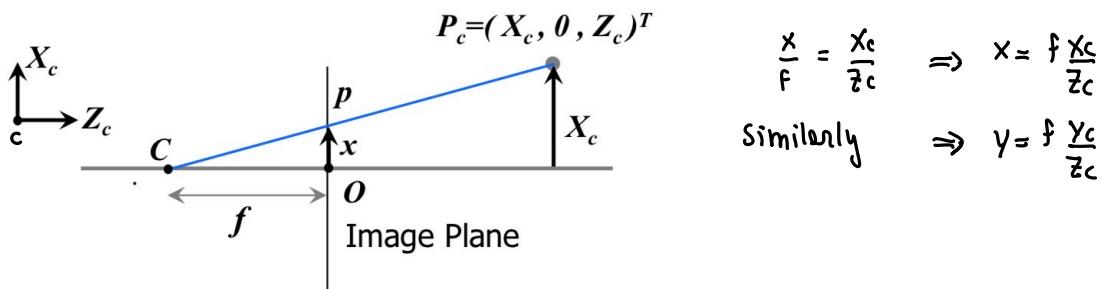
PERSPECTIVE CAMERA

Goal: Find pixel coordinates (u, v) of point P_w in the world frame:

- 1) Convert P_w to camera point P_c through extrinsics (transform $[R|t]$)
- 2) Convert P_c to image plane coordinates (x, y)
- 3) Convert (x, y) to discretized pixel coordinates (u, v)

$$1) \quad P_c = [R|t] P_w \quad P \text{ is in homogeneous coordinates and } [R|t] = \begin{bmatrix} R & t \end{bmatrix} \in \mathbb{R}^{3 \times 4}$$

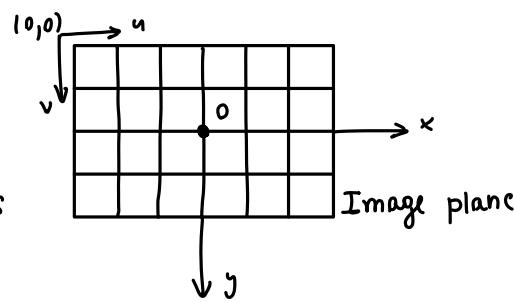
2)



3) Convert (x, y) to (u, v)

- pixel coords of O are (u_0, v_0)
- Scale factors K_u, K_v for the pixel size in both dimensions

$$\boxed{\begin{aligned} u &= u_0 + k_u x \\ v &= v_0 + k_v y \end{aligned}}$$



- Use HOMOGENEOUS COORDINATES for linear mapping from 3D to 2D, by introducing an extra element (scale)

$$p = \begin{pmatrix} u \\ v \end{pmatrix} \Rightarrow \tilde{p} = \begin{pmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{pmatrix} = \lambda \begin{pmatrix} u \\ v \\ 1 \end{pmatrix}$$

• It follows that

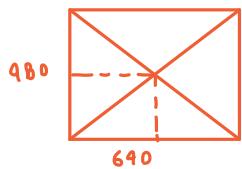
$$\begin{aligned} u &= u_0 + k_u f x_c / z_c \\ v &= v_0 + k_v f y_c / z_c \end{aligned} \Rightarrow \begin{bmatrix} \lambda u \\ \lambda v \\ \lambda \end{bmatrix} = \begin{bmatrix} k_u f & 0 & u_0 \\ 0 & k_v f & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix} = K \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

CALIBRATION MATRIX

EXERCISE 1

- Determin K of a camera 640×480 pixels with $\text{FoV} = 90^\circ$
- The principal point is in the center of the diagonals
- Also, pixels are square
- What is the vertical field of view?

$$\tan\left(\frac{\theta}{2}\right) = \frac{W}{2f} \Rightarrow \tan(45^\circ) = \frac{640}{2f} \Rightarrow f = 320 \text{ pixels}$$



$$u_0 = 320 \quad \text{Square pixels} \rightarrow k_u = k_v = 1 \Rightarrow K = \begin{bmatrix} 320 & 0 & 320 \\ 0 & 320 & 240 \\ 0 & 0 & 1 \end{bmatrix}$$

$$\theta_v = 2 \arctan\left(\frac{W}{2f}\right) = 2 \arctan\left(\frac{480}{2 \cdot 320}\right) = 2 \arctan\left(\frac{3}{4}\right) = 73.79^\circ$$

EXERCISE 2

- Prove that world parallel lines intersect at a vanishing point in the camera image

Equation of a 3D line: $\vec{p}_0 + \vec{d}t$, $t \in \mathbb{R}$ $\vec{d} = [l \ m \ n]$

$$\lim_{t \rightarrow \infty} Xv = \lim_{t \rightarrow \infty} \frac{f(x_0 + tl)}{(z_0 + tn)} = \frac{f l}{n}$$

$$\lim_{t \rightarrow \infty} Yv = \lim_{t \rightarrow \infty} \frac{f(y_0 + tm)}{(z_0 + tn)} = \frac{f m}{n}$$

PERSPECTIVE CAMERA (CONTINUED)

From the real world to the camera image

$$\begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix} = R \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix} + \begin{bmatrix} t_1 \\ t_2 \\ t_3 \end{bmatrix} = [R \mid t] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

PERSPECTIVE PROJECTION EQUATION

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R \mid t] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

We define the NORMALIZED IMAGE COORDINATES as $\begin{bmatrix} \bar{u} \\ \bar{v} \\ 1 \end{bmatrix} = K^{-1} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$ (as the focal length were 1)

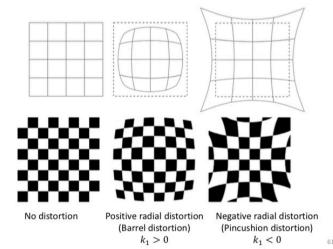
LENS DISTORTION

RADIAL DISTORTION

- For a given non distorted image point (u, v) the amount of distortion is non linear of its distance r from the principal point
- A simple quadratic model often works

$$\begin{bmatrix} u_d \\ v_d \end{bmatrix} = (1+k_1 r^2) \begin{bmatrix} u-u_0 \\ v-v_0 \end{bmatrix} + \begin{bmatrix} u_0 \\ v_0 \end{bmatrix}$$

$$r^2 = (u-u_0)^2 + (v-v_0)^2$$



- Depending on the distortion's amount (and thus the camera FOV), higher order terms can be introduced

TANGENTIAL DISTORTION

- If the lens is misaligned (not perfectly orthogonal to the image sensor) a non-radial distortion appears.

CAMERA CALIBRATION

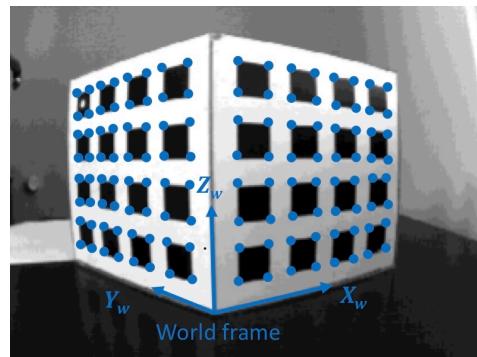
CALIBRATION: the process to determine the intrinsic K and extrinsic (R, t).
(optionally also the lens distortion)

There are 2 popular methods

- TSAI'S METHOD** : uses 3D objects
- ZHANG'S METHOD** : uses planar grids

TSAI'S METHOD

- measure the 3D position of $n \geq 6$ control points on a 3D calibration target and their 2D coordinates of their projection in the image.



- The method of Direct Linear Transform (DLT) rewrites the perspective projection as a HOMOGENEOUS LINEAR EQUATION and solves it by standard methods.

$$\lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = K[R|t] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} = \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & \ddots & & \\ \vdots & & \ddots & m_{34} \end{bmatrix}}_M \underbrace{\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}}_P = \begin{bmatrix} m_1^T \\ m_2^T \\ m_3^T \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

$$\rightarrow u = \frac{\lambda u}{\lambda} = \frac{m_1^T P}{m_3^T P} \quad \rightarrow (m_1^T - u m_3^T) P = 0$$

$$\rightarrow v = \frac{\lambda v}{\lambda} = \frac{m_2^T P}{m_3^T P} \quad \rightarrow (m_2^T - v m_3^T) P = 0$$

$$\rightarrow \begin{bmatrix} p_i^T & \theta^T - u_i P_i^T \\ \theta^T P_i^T - v_i P_i^T \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

- We can stack N points:
$$\begin{bmatrix} P_1^T & 0 & -u_1 P_1^T \\ 0 & P_1^T & -u_1 P_1^T \\ P_2^T & 0 & -u_2 P_2^T \\ 0 & P_2^T & -u_2 P_2^T \\ \vdots & & \end{bmatrix} \begin{bmatrix} m_1 \\ m_2 \\ m_3 \\ \vdots \end{bmatrix} = Q \cdot M = 0 \quad M \in \mathbb{R}^{12} \text{ (vectorized matrix)}$$

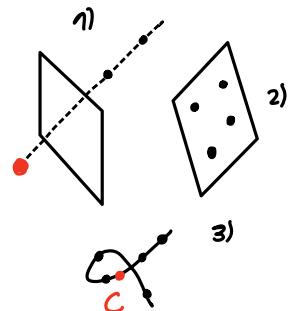
- $M \in \mathbb{R}^{4 \times 3}$: 12 unknowns
- $Q \in \mathbb{R}^{2n \times 12}$
- Each point provides 2 independent equations \rightarrow need ≥ 6 points
- If $n > 6$ find least square solution:

$$M = \underset{\|M\|^2=1}{\operatorname{argmin}} \|QM\|^2 = \underset{\|M\|^2=1}{\operatorname{argmin}} M^T Q^T Q M \quad \|M\|^2=1 \text{ prevents degenerate solution with } M=\emptyset$$

M^* is the eigenvector corresponding to the minimum eigenvalue

DEGENERATE CONFIGURATIONS

- Points lying on a **PLANE**
- Points along a single **LINE PASSING THROUGH THE CENTER OF PROJECTION**
- Camera and points on a **TWISTED CUBIC** (smooth curve in 3D of degree 3)
- Better use >20 points non coplanar



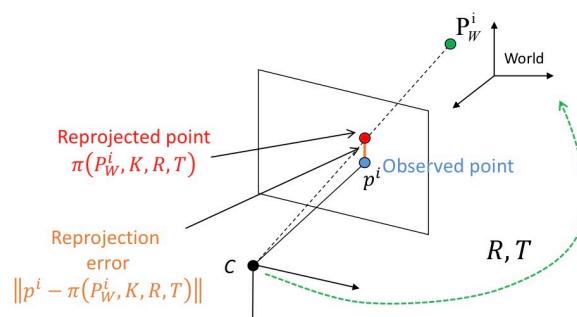
RECOVER K, R, T FROM M

$$\begin{bmatrix} m_{11} & m_{12} & m_{13} & m_{14} \\ m_{21} & m_{22} & m_{23} & m_{24} \\ m_{31} & m_{32} & m_{33} & m_{34} \end{bmatrix} = \begin{bmatrix} \alpha r_{11} + u_0 r_{31} & \alpha r_{12} + u_0 r_{32} & \alpha r_{13} + u_0 r_{33} & \alpha r_{14} + u_0 r_{34} \\ r_{31} & r_{32} & r_{33} & r_{34} \\ \vdots & & & \end{bmatrix} \quad \text{Not ensuring } R = R^T$$

- We can apply the QR factorization of M: $M = K \cdot R \cdot T$
upper orthogonal
triangular

RECONSTRUCTION ERROR

- The **REPROJECTION ERROR** is the euclidean distance between the observed point and the projected point (using estimated K, R, T).
- Refine K, R, T and introduce lens distortion solving non-linear optimization using DLT solution as guess



$$K, R, T, \text{dist} = \underset{K, R, T, \text{dist}}{\operatorname{argmin}} \sum_{i=1}^N \|p^i - \pi(P_W^i, K, R, T)\|^2$$

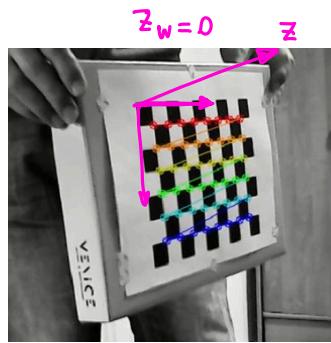
N.B. Levenberg - Marquardt method is more robust than Gauss-Newton to local minima.

ZHANG'S METHOD

- Tsai's method drawback: 3D points are not coplanar
- Today's calibration: multi view of a planar grid.

- In Zhang's method $z_w = 0$

$$\begin{aligned} \lambda \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} &= K[R|T] \begin{bmatrix} x_w \\ y_w \\ 0 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} du & 0 & u_0 \\ 0 & dv & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ 0 \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} du & 0 & u_0 \\ 0 & dv & v_0 \\ 0 & 0 & 1 \end{bmatrix} \underbrace{\begin{bmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{bmatrix}}_{H: \text{ homography}} \begin{bmatrix} x_w \\ y_w \\ 1 \end{bmatrix} \\ &= \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} h_1^T \\ h_2^T \\ h_3^T \end{bmatrix}}_{\text{Known}} \begin{bmatrix} x_w \\ y_w \\ 1 \end{bmatrix} \end{aligned}$$



→ Apply DLT for a single view $(h_1^T - u_i h_3^T) p_i = 0$ $\Rightarrow Q \cdot H = 0$
 $(h_2^T - v_i h_3^T) p_i = 0$
 \vdots
 $\vdots \forall i \in [1, n]$

- $Q \in \mathbb{R}^{2n \times 9}$, $H \in \mathbb{R}^9$
- each point provides 2 indep. equations
- At least 4 non collinear points required (more the better → use SVD as before)

RECOVER K, R, T

- multiple views needed (20-50 spanning the camera FoV)
- each view has the same K , but different R, T

$$H^j = K \cdot [R|T]^j$$

Algorithm

1. Estimate H_i ∀ i^{th} view using DLT algorithm

2. Determine K from a set of homographies

- ∀ homography

$$H_i = K[R|T] := [w_1 | w_2 | t]$$

- From orthogonality

$$r_1^T r_2 = \underbrace{r_1^T}_{w_1^T := B} \underbrace{K^T}_{\text{B}} \underbrace{K^{-1}}_{\text{B}} \underbrace{K}_{\text{B}} \underbrace{r_2}_{w_2^T} = w_1^T B w_2 = \emptyset$$

$$r_1^T r_1 = r_2^T r_2 = 1 = w_1^T B w_1 - w_2^T B w_2 = \emptyset$$

- Stack 2n equations for n views and estimate B

3. Determine K from B using Cholesky

4. Determine R, T ∀ homography

$$K r_1 = H_1 \text{ (first column)}$$

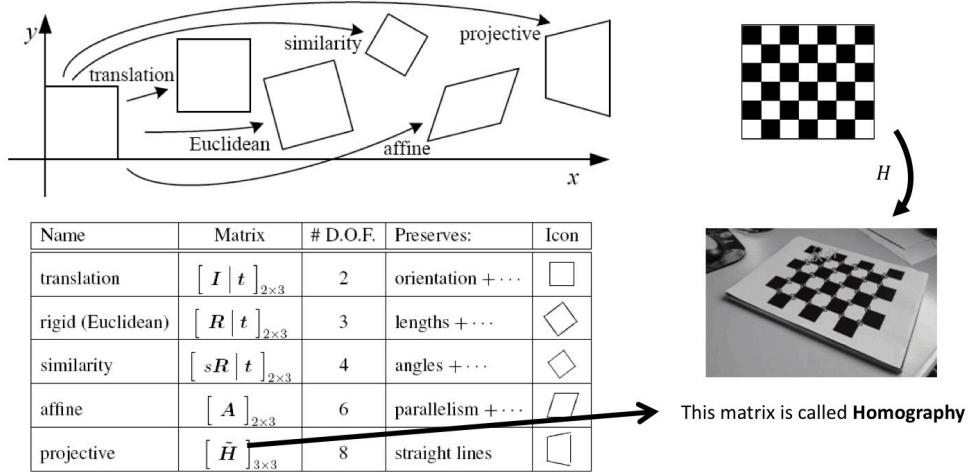
$$K T = H_3$$

$$K r_2 = H_2$$

$$\text{and } r_3 = r_1 \times r_2$$

} 2 equation ∀ homography

Types of 2D Transformations



Projective Transformation (Homography)

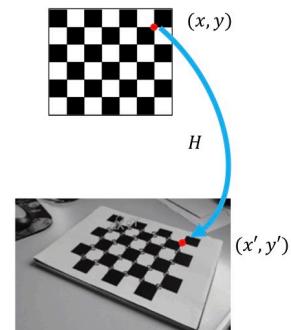
- A point (x, y) is transformed into (x', y') via:

$$x' = \frac{a_1x + a_2y + a_3}{a_7x + a_8y + 1}$$

$$y' = \frac{a_4x + a_5y + a_6}{a_7x + a_8y + 1}$$

- Homogeneous coordinates:

$$\lambda \begin{bmatrix} x' \\ y' \\ 1 \end{bmatrix} = \begin{bmatrix} a_1 & a_2 & a_3 \\ a_4 & a_5 & a_6 \\ a_7 & a_8 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

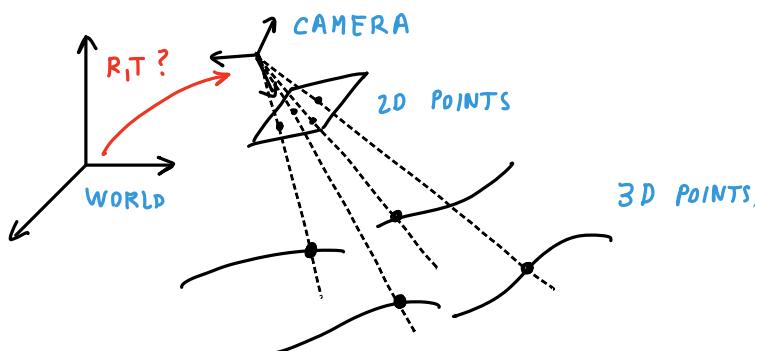


CAMERA LOCALIZATION (PERSPECTIVE FROM n POINTS = PnP)

PROBLEM: Determine 6 DoF pose of camera with respect to the world frame from a set of 3D-2D points correspondences

ASSUMPTIONS: Camera intrinsics are known

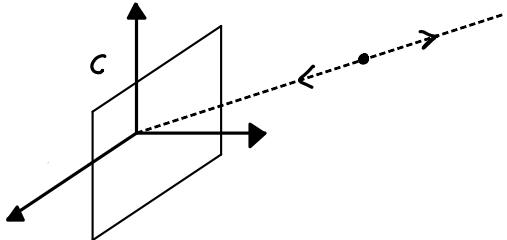
SOLUTIONS: DLT can be used but **suboptimal**. It exists an algebraic solution.



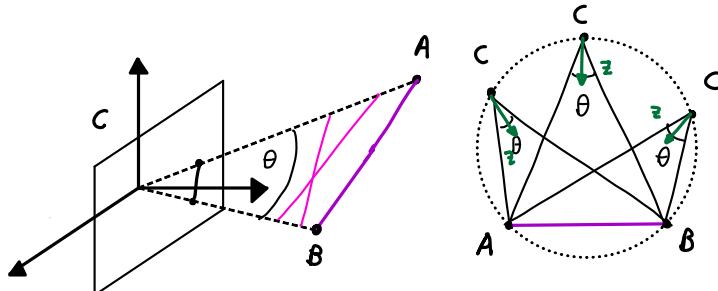
HOW MANY POINTS ARE ENOUGH?

- 1 POINT $\rightarrow \infty$ solutions
- 2 POINTS $\rightarrow \infty$ solutions but bounded
- 3 POINTS \rightarrow (if non collinear) up to 4 solutions
- 4 POINTS \rightarrow unique solution

1 POINT

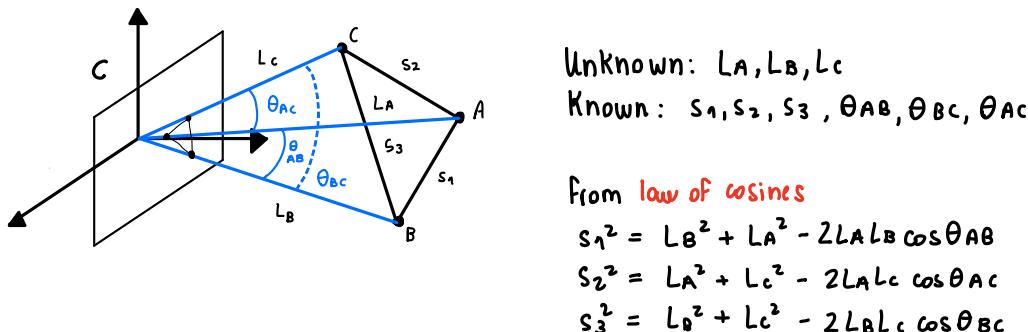


2 POINTS



bounded \rightarrow camera must lie on that circle but ∞ solutions

3 POINTS



FACTS

- n independent polynomials in n unknowns can have no more solutions than the product of their respective degrees
- If every term is constant or of degree 2, for each positive solution there is one negative

P3P Problem: 3 eq. of degree 2 $\rightarrow 2 \cdot 2 \cdot 2 = 8$ solutions

terms constant or degree 2 \rightarrow 4 solutions are positive $\Rightarrow 4$ solutions
4 solutions are negative (non valid)

- It can be shown that defining $x = L_B/L_A$ the system can be reduced to

$$G_0 + G_1 x + G_2 x^2 + G_3 x^3 + G_4 x^4 = 0$$

- A 4th point helps to disambiguate solutions. Their classification and recovering R and T with **Gao's algorithm** (OpenCV: `solvePnP P3P`)

MODERN SOLUTIONS TO PnP

- AP3P
- EPnP (efficient PnP): 10x more efficient and accurate than DLT (when camera is calibrated)

APPLICATIONS

- LOCALIZATION: given a PCL get camera pose
- MONOCULAR VISUAL ODOMETRY: given initial point cloud find correspondences and new camera pose

ROBUSTNESS TO OUTLIERS

- All PnP problems are error prone if outliers in 3D-2D correspondences
- The RANSAC algorithm can be used to remove outliers
- PnP + RANSAC in OpenCV: `solvePnP(Ransac)`
- The output of both EPnP and DLT can be refined via non linear optimization

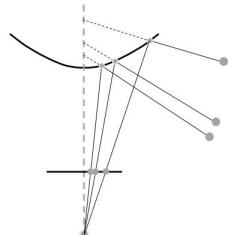
NON CONVENTIONAL CAMERA MODELS

FISHEYE : wide FOV ($> 130^\circ$)

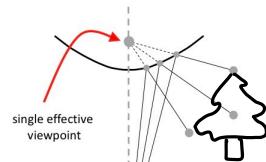
CATADIOPTRIC : 360° all around

Central vs Non-Central Omnidirectional Cameras

Non-Central projection system
Rays do not intersect in a single point

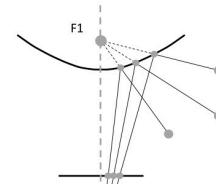


Central projection system
Rays intersect in a single point



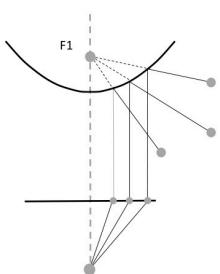
Central Omnidirectional Cameras

Hyperbola + Perspective camera



NB: one of the foci of the hyperbola must lie in the camera's center of projection

Parabola + Orthographic lens



CENTRAL CAMERAS are better

- We can apply standard algorithms valid for **perspective geometry**
- un warp part of the image in a perspective one
- transform image points into normalized vectors on the unit sphere

UNIFIED OMNIDIRECTIONAL CAMERA MODEL (for fisheye and catadioptric cameras)

- define a **image projection function** with polynomials whose coefficients need to be estimated
- found coeff + intrinsics + extrinsics via DLT.

$$\frac{\lambda}{\alpha} \begin{bmatrix} u - u_0 \\ v - v_0 \\ g(\rho) \end{bmatrix} = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix} \quad g(\rho) = \alpha + \sum_{i=1}^N \alpha_i \rho^i$$

$\alpha_i, i > 0 = 0 \rightarrow$ pinhole camera

- **OCamCalib**: toolbox for calibrating wide angle camera

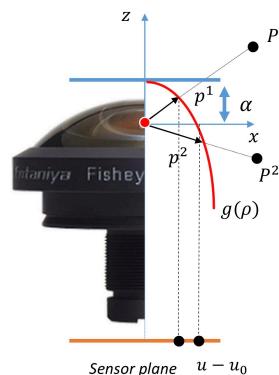


IMAGE FILTERING

- **FILTER**: process that remove certain "frequency components"
- **LOW-PASS FILTER**: smooth image (keeps low-frequency)
- **HIGH-PASS FILTER**: retains contours (keeps high-frequency = sharp intensity discontinuities)

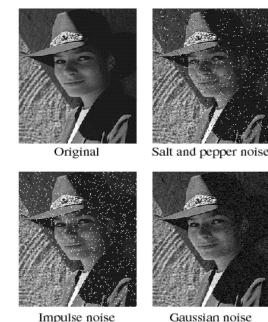
NOISE REDUCTION

- **SALT AND PEPPER NOISE**: random black and white pixels
- **IMPULSE NOISE**: random white pixels
- **GAUSSIAN NOISE**: intensity variation drawn from Gaussian

$$I(x,y) = I'(x,y) + \eta(x,y)$$

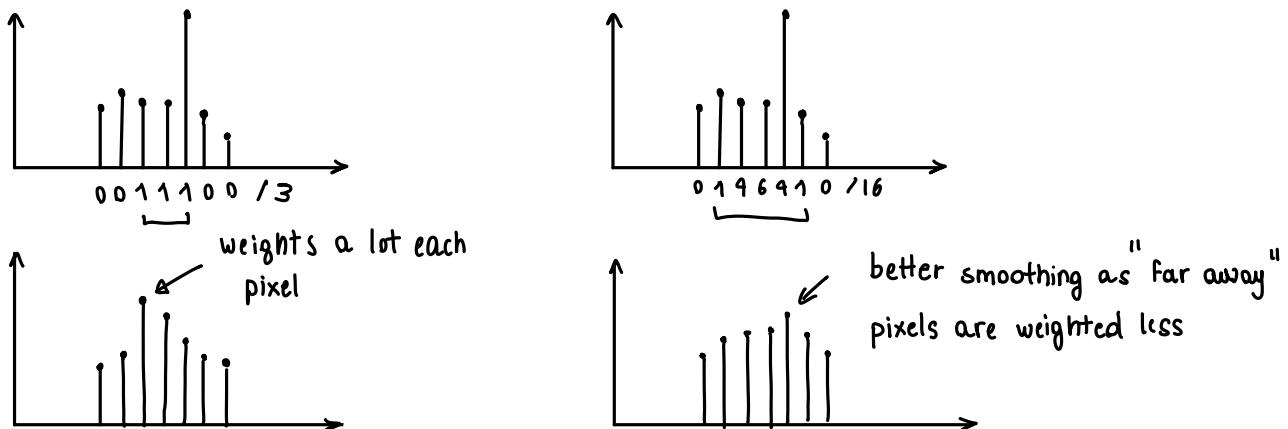
$$\eta(x,y) \sim \mathcal{N}(0,\sigma^2)$$

I' : ideal image



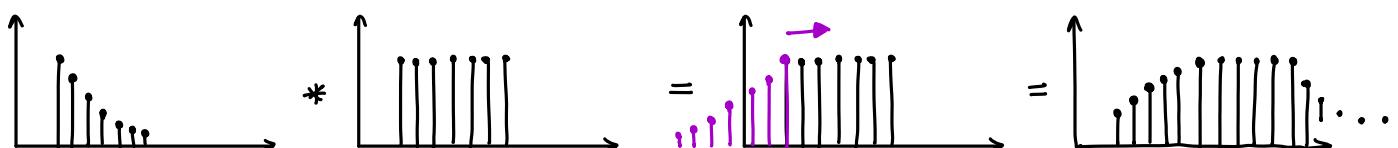
MOVING AVERAGE

- One possible method to reduce Gaussian noise
- **Assumptions**:
 - 1) pixels are like neighbours
 - 2) noise independent from pixel to pixel



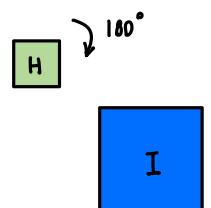
CONVOLUTION

- One of the sequences is flipped (right to left) before sliding over the other
- **Notation**: $a * b$
- **Properties**: linearity, associativity, commutativity



2D FILTERING VIA 2D CONVOLUTION

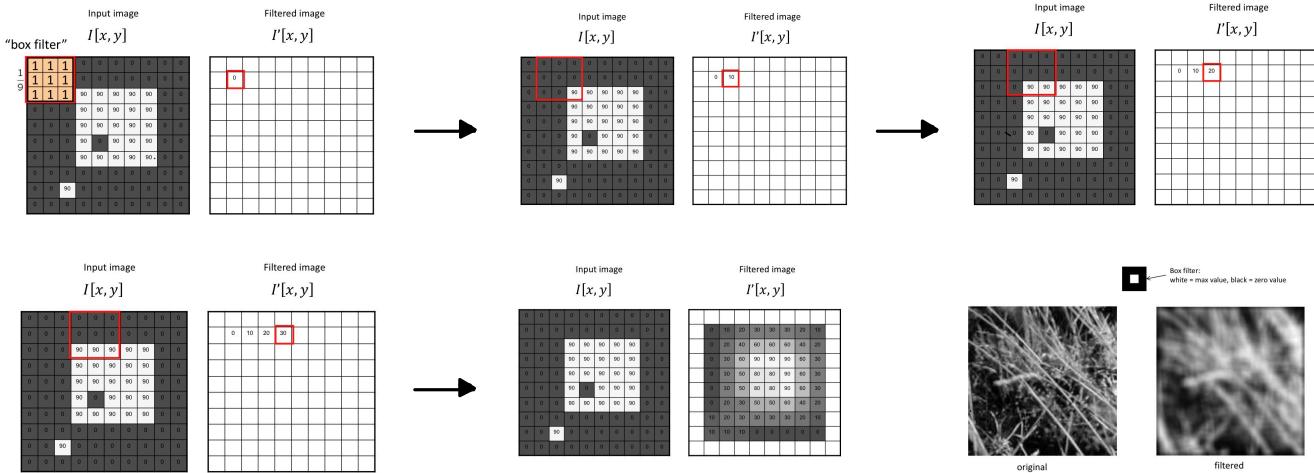
- Flip the filter in both dimensions (bottom to top, right to left = 180° turn)
- Then slide filter over the image and compute sum products
- Convolution replaces each pixel with a weighted sum of its neighbors
- The filter H is also called **Kernel** or **mask**



CONVOLUTION : $I' = I * H$ (linearity, associativity, commutativity) $I'[x,y] = \sum_{u=-k}^k \sum_{v=-k}^k I[x-u, y-v] H[u, v]$

CROSS-CORRELATION : $I' = I \otimes H$ (linearity but **not** associativity, commutat.) $I'[x,y] = \sum_{u=-k}^k \sum_{v=-k}^k I[x+u, y+v] H[u, v]$

→ The same if H is symmetric



NOTE : the box filter causes a "web-like" effect because of "aliasing": high frequency components seen as low frequency

GAUSSIAN FILTER

This filter is an approximation of a Gaussian 2D function $H[u,v] = \frac{1}{2\pi\sigma^2} e^{-\frac{u^2+v^2}{2\sigma^2}}$

- **Kernel size** : size of the filter → the larger, the better its discrete approximation
- **Variance** : controls the amount of smoothing
standard deviation = σ [pixels]
variance = σ^2 [pixels²]

0 0 0	0 0 0	0 0 0	0 0 0	0 0 0	0 0 0	0 0 0
0 0 0	0 0 0	0 0 0	0 0 0	0 0 0	0 0 0	0 0 0
0 0 0	90 90 90	90 90 90	90 90 90	0 0 0	0 0 0	0 0 0
0 0 0	90 90 90	90 90 90	90 90 90	0 0 0	0 0 0	0 0 0
0 0 0	90 90 90	90 90 90	90 90 90	0 0 0	0 0 0	0 0 0
0 0 0	90 90 90	90 90 90	90 90 90	0 0 0	0 0 0	0 0 0
0 0 0	90 90 90	90 90 90	90 90 90	0 0 0	0 0 0	0 0 0

1 1 2 1
16 1 2 1
 $H[u, v]$

SEPARABLE FILTERS

- A convolution with a 2D filter of $w \times w$ pixel size requires w^2 multiply-add operations per pixel
- 2D convolution can be sped up if the filter is **separable**:
- A Separable filter can be written as the product of 2 1-D filters : $H = v \cdot h^\top$

$I' = I * H = (I * h^\top) * v \rightarrow$ do 2 1D convolutions only = **2w multiply-add operations per pix.**

Examples:

• Box filter = $\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix} = \frac{1}{3} \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix} \cdot \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}$

• Gaussian filter = $\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \cdot \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \end{bmatrix}$

• Sobel filter = $\begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1 \\ 2 \\ 1 \end{bmatrix} \cdot \begin{bmatrix} -1 & 0 & 1 \end{bmatrix}$

BOUNDARY ISSUES

What do we do when the filter goes through the edge? We pad with different methods:

- zero padding (black contour)
- wrap around
- copy edge
- reflect across edge

NON LINEAR FILTERS

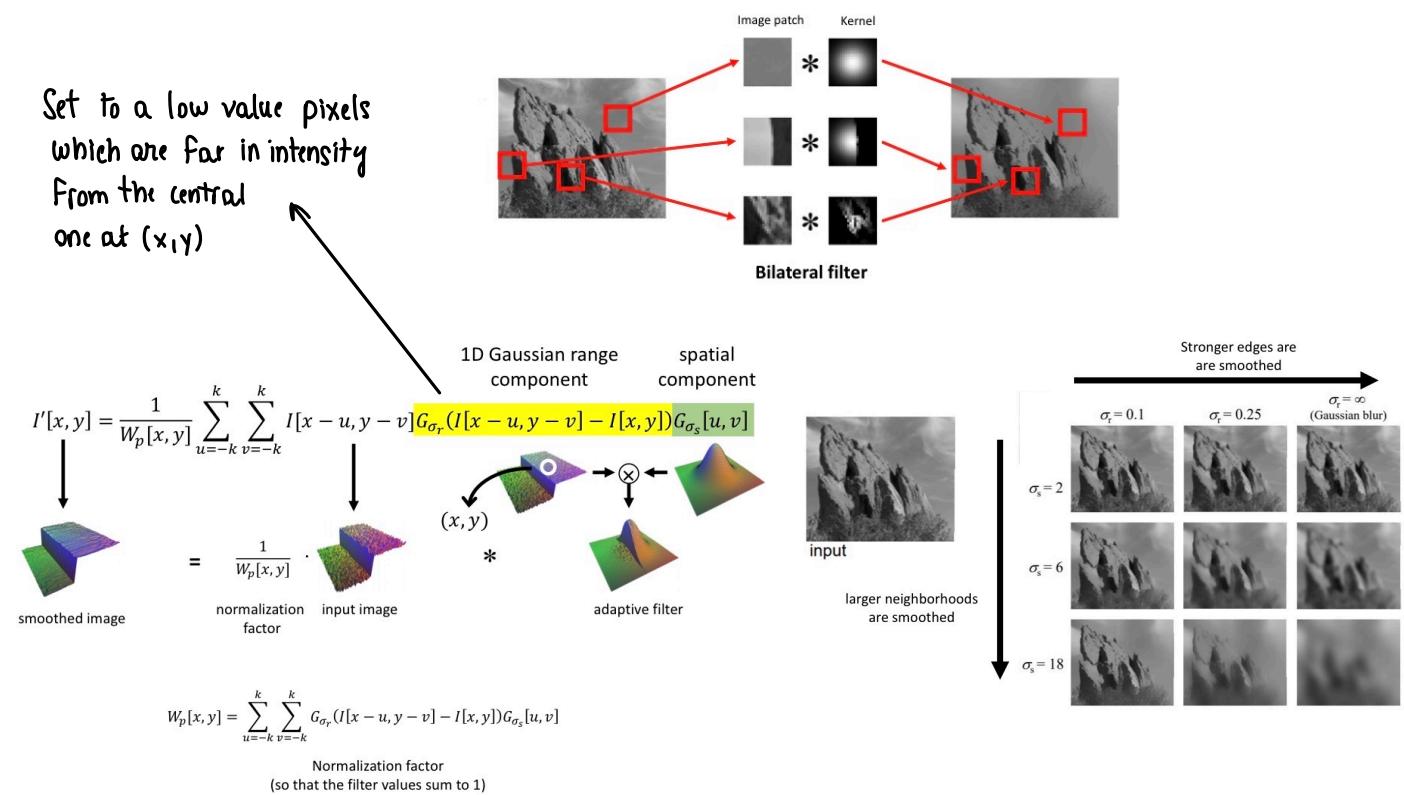
- **Linear filters** cannot remove effects like salt-and-pepper noise.
- Non-linear filters like the **median filter** can.

MEDIAN FILTER

- It's a non linear filter : sort the input patch and extract the median value
- Removes **spikes** : good for salt-and-pepper and impulse noise.
- Preserves strong edges and does not smooth (contrary to Gaussian filter)

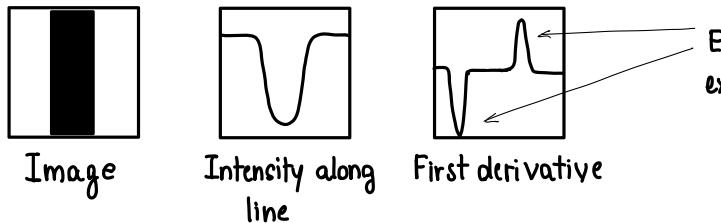
BILATERAL FILTER

- Solves the last problem by **adapting the Kernel locally to the intensity profile**, so they are **patch-content dependent**
- They smooth only pixels with the same brightness as the center pixel and ignore influence of pixels with different brightness across the discontinuity



EDGE DETECTION

- An **edge** is a place of fast change in the image intensity function



Edges correspond to local extrema of first derivative

$$\frac{\partial I(x,y)}{\partial x} = \lim_{\varepsilon \rightarrow 0} \frac{I(x+\varepsilon, y) - I(x, y)}{\varepsilon} \approx \frac{I(x+1, y) - I(x, y)}{1} \text{ adjacent derivative}$$

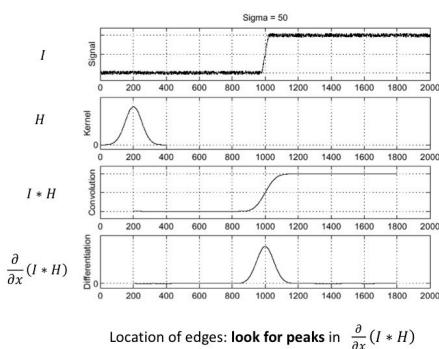
$$\approx \frac{I(x+1, y) - I(x-1, y)}{2} \text{ central derivative}$$

Central derivative filters:

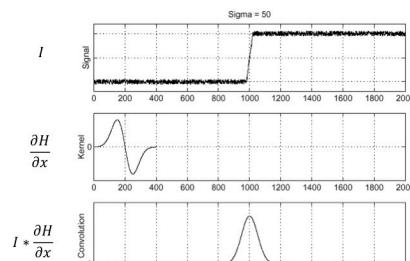
- Prewitt filter** $G_x = \begin{bmatrix} -1 & 0 & 1 \\ -1 & 0 & 1 \\ -1 & 0 & 1 \end{bmatrix}, G_y = \begin{bmatrix} -1 & -1 & -1 \\ 0 & 0 & 0 \\ 1 & 1 & 1 \end{bmatrix}$
- Sobel filter** $G_x = \begin{bmatrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{bmatrix}, G_y = \begin{bmatrix} -1 & -2 & -1 \\ 0 & 0 & 0 \\ 1 & 2 & 1 \end{bmatrix}$

IMAGE GRADIENT

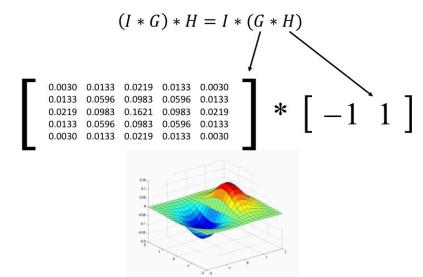
- The **image gradient** $\nabla I = [\partial I / \partial x, \partial I / \partial y]$
- The gradient direction is $\theta = \text{atan} 2(\partial I / \partial y, \partial I / \partial x)$
- The gradient points in the **direction of steepest ascent**
- The **edge strength** is given by the gradient magnitude $\|\nabla I\|$
- noise** has the effect of making gradient noisy \rightarrow edge difficult to find
- solution 1: smooth first and final peak of gradient



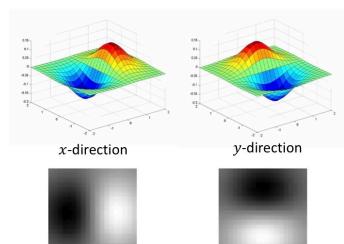
Differentiation property of convolution: $\frac{\partial}{\partial x}(I * H) = I * \frac{\partial H}{\partial x}$



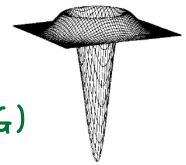
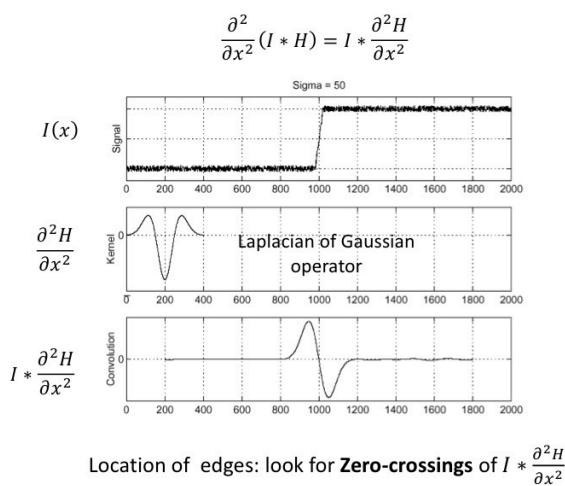
Derivative of Gaussian filter G along x



Derivative of Gaussian Filters



- solution 2: find zero-crossing of second order derivative (where there is a sign change = derivative change) of the intensity



LA PLACIAN OF GAUSSIAN (LoG)

- It is a circularly symmetric filter defined as:

$$\nabla^2 G_b = \frac{\partial^2 G_b}{\partial x^2} + \frac{\partial^2 G_b}{\partial y^2} \quad \nabla^2 \triangleq \frac{\partial^2}{\partial x^2} + \frac{\partial^2}{\partial y^2}$$

- two commonly used approximations

$$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix} \quad \begin{bmatrix} 1 & 1 & 1 \\ 1 & -8 & 1 \\ 1 & 1 & 1 \end{bmatrix}$$

RECAP FILTERS

Smoothing filters

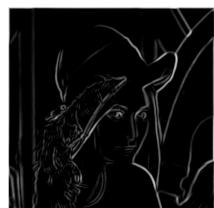
- has positive values
- sums to 1 → preserve brightness of constant regions
- removes "high frequency" components

Derivative filters

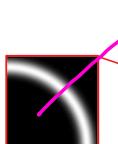
- has opposite signs used to get high response in regions of high contrast
- sums to 0 → no response in constant regions
- highlights "high frequency" components

THE CANNY EDGE DETECTION ALGORITHM

- Take a grayscale image. If RGB, take mean of its channels
- Convolve the image I with x and y derivative of Gaussians filter and compute the edge strength $\|\nabla I\|$
- Thresholding: set to 0 all pixels where $\|\nabla I\| < \text{threshold}$
- Thinning: look for local maxima in the edge strength in the direction of the gradient
 - take the directional derivative of the edge strength in the direction of the gradient and then looking for zero crossing (adjacent pixels where sign change = edge raising or falling)
 - Equivalent to $\nabla^2 G_b * I$



gradient
direction



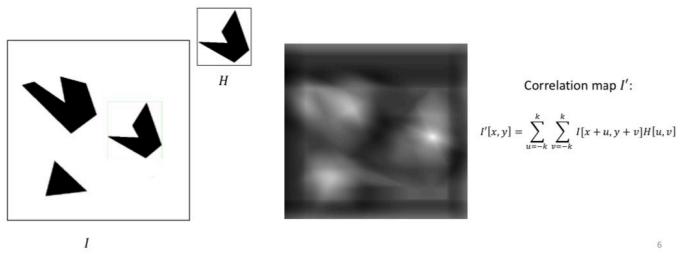
FILTERS FOR FEATURE DETECTION

Filters can be used to extract **features**

- keep what is important
- reduce later processing
- reduce memory requirements and compute

Examples are

- edge detection
- template matching
- Keypoint detection



6

TEMPLATE MATCHING

- Find locations in an image I that are similar to a **template** H
- We can use cross-correlation (180° rotated filter)
- Will work only if **scale**, **orientation**, **illumination** and in general the **appearance** of the template (including background) are **very similar**.

Correlation can be seen as **scalar product** between two vectorized image patches H and F : $\langle H, F \rangle = \|H\| \|F\| \cos \theta$

SIMILARITY MEASURES

$$\text{Normalized Cross Correlation (NCC)}: NCC = \cos \theta = \frac{\langle H, F \rangle}{\|H\| \|F\|} = \frac{\sum_u \sum_v H(u,v) F(u,v)}{\sqrt{\sum_u \sum_v H(u,v)^2} \sqrt{\sum_u \sum_v F(u,v)^2}}$$

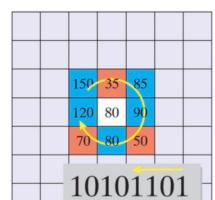
$$\text{Sum of Squared Differences (SSD)}: SSD = \sum_u \sum_v (H(u,v) - F(u,v))^2$$

$$\text{Sum of Absolute Differences (SAD)}: SAD = \sum_u \sum_v |H(u,v) - F(u,v)|$$

- To account for the difference in the average intensity of two images (typically caused by additive illumination changes) we subtract the mean value of each patch: $\mu_H = \frac{1}{n} \sum_u \sum_v H(u,v)$; $\mu_F = \frac{1}{m} \sum_u \sum_v F(u,v)$. We then replace H and F in the previous formula with $H' = H - \mu_H$ and $F' = F - \mu_F$ obtaining the **zero-mean** version of the metrics: **ZNCC**, **ZSSD** and **ZSAD**.
- Note that **ZNCC** is invariant to affine illumination changes: $I' = \alpha I + \beta$, the others not.

CENSUS TRANSFORM

- Maps an image patch to a bit string
 - If a pixel is greater than the center pixel its corresponding bit is set to 1, else 0
 - For a $w \times w$ patch the string will be $w^2 - 1$ long
- The two strings are **compared using the Hamming distance** (number of different bits = Σ XOR)
- **No square roots or divisions** \rightarrow efficient (especially on FPGA)
- Intensities relative to the center pixel \rightarrow **invariant to monotonic intensity changes**.



KEYPOINTS EXTRACTION AND MATCHING

Why? They are the feature we extract and match between consecutive frames to estimate motion. But Keypoints are also used for :

- Panorama stitching
- Object recognition
- 3D reconstruction
- Place recognition
- Indexing and database retrieval

Feature Matching Problem: find similar keypoints between two images of the same scene taken under **different conditions**.

PROBLEM 1

Can we **re-detect** features?

Repeatability: property of a "detector" to re-detect the same feature in different images of the same scene.

PROBLEM 2

Can we **match** to same feature?

Distinctiveness: a descriptor uniquely identifies a feature from other features without ambiguity. Must be **robust** to geometric and photometric changes

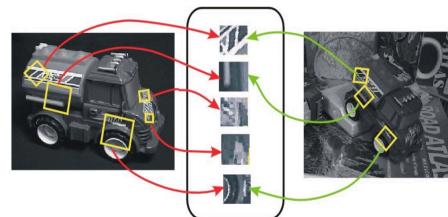
- Geometric changes : rotation, scale, view point (perspective change)
- Photometric changes : affine intensity transformations:

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

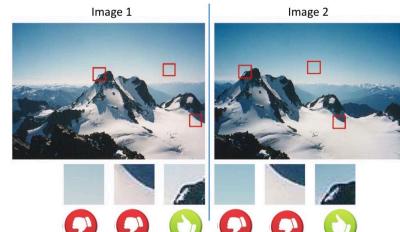
LOCAL INVARIANT FEATURES

Try to address these problems and be repeatable, distinctive and robust.

1. Detect repeatable and distinctive **interest points**,
2. Extract **invariant (robust) descriptors**.



REPEATABLE AND DISTINCTIVE FEATURE?



POINT FEATURES : CORNERS VS BLOB DETECTORS

Corner detector: intersection of two or more edges

(+) **high localization accuracy** : good for VO

(-) **less distinctive than blobs**

Ex.: Harris, Shi-Tomasi, SUSAN, FAST



Blob detector: any other image pattern that is **not a corner**

and differs significantly from its neighbors

(+) **blobs have less localization accuracy** than corners

(-) **blobs are more distinctive** : better for Place Recognition

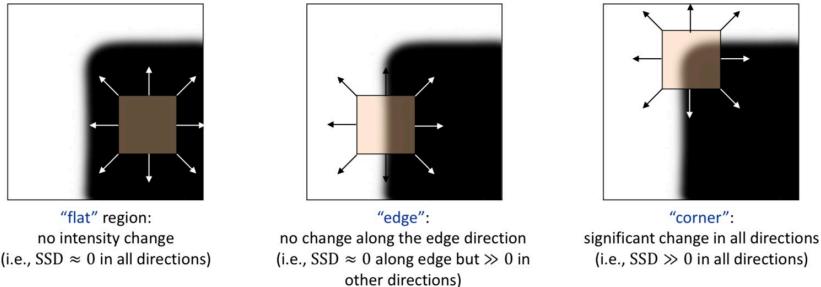
Ex.: MSER, LOG, DOG (SIFT), SURF, CenSurE, etc.



Connected region of pixels with similar color, a circle, etc...

CORNER DETECTION

- In the region around a corner, image gradients has two or more dominant directions.
- How?
 - We look at a region of pixels through a small window.
 - Shifting the window in any direction should cause large intensity changes (i.e. in SSD)



* Sum of Squared Distances

THE MORAVEC CORNER DETECTOR (1980)

- Consider the reference patch centered at (x, y) and the shifted window centered at $(x + \Delta x, y + \Delta y)$. The patch has size Ω . The Sum of Squared Distances is

$$SSD(\Delta x, \Delta y) = \sum_{x, y \in \Omega} (I(x, y) - I(x + \Delta x, y + \Delta y))^2$$

- Repeat this for the 4 directions (right, down, down right, down left = horizontal, vertical and two diagonals).
- The window's interest measure is the minimum of these 4.
- Features are chosen where interest measure has local maxima.
- Problem: 4 computations every time ...

THE HARRIS CORNER DETECTOR (1988)

- Implements Moravec without shifting but using derivative filters.
- The shifted image can be approximated with a 1st order Taylor expansion

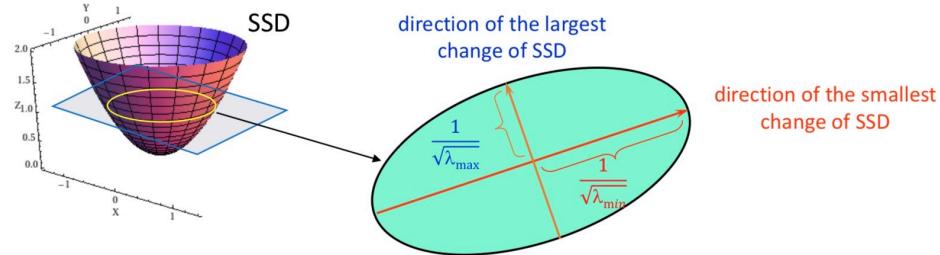
$$\begin{aligned} I(x + \Delta x, y + \Delta y) &\approx I(x, y) + I_x(x, y) \Delta x + I_y(x, y) \Delta y \Rightarrow SSD = \sum_{x, y \in \Omega} (I(x, y) - I(x + \Delta x, y + \Delta y))^2 \\ &= \sum_{x, y \in \Omega} (I_x(x, y) \Delta x + I_y(x, y) \Delta y)^2 \\ &= \sum_{x, y \in \Omega} [\Delta x \ \Delta y] \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} \end{aligned}$$

$$SSD(\Delta x, \Delta y) \approx [\Delta x \ \Delta y] \begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix}$$

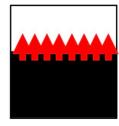
Quadratic form

M (2nd moment matrix): $I_x I_y$ are not matrix products but pixel-wise products.

- Since M is symmetric, it can always be decomposed into $M = R^{-1} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} R$
- We can visualize $\begin{bmatrix} \Delta x & \Delta y \end{bmatrix} \cdot \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} = \text{const}$ as an ellipse with axis length λ_1, λ_2 and orientation R .
- The two eigenvectors identify the directions of largest and smallest changes of SSD.



Example

 Edge

$$M = \begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & \lambda_2 \end{bmatrix}, \lambda_1 \approx 0$$

 Flat region

$$M = \begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \lambda_1 \approx 0 \wedge \lambda_2 \approx 0$$

 Corner

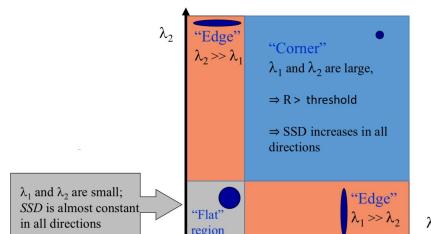
$$M = \begin{bmatrix} \sum I_x^2 & \sum I_x I_y \\ \sum I_x I_y & \sum I_y^2 \end{bmatrix} = \begin{bmatrix} \cos \frac{\pi}{4} & -\sin \frac{\pi}{4} \\ \sin \frac{\pi}{4} & \cos \frac{\pi}{4} \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} \begin{bmatrix} \cos \frac{\pi}{4} & \sin \frac{\pi}{4} \\ -\sin \frac{\pi}{4} & \cos \frac{\pi}{4} \end{bmatrix}, \lambda_1 \gg 0 \wedge \lambda_2 \gg 0$$

REVIEW: COMPUTE λ_1, λ_2, R FROM M

- $Ax = \lambda x$ for A square
- Non trivial eigenvalues ($x \neq 0$) Solving $\det(A - \lambda I) = 0$
- $A = M : \det \begin{bmatrix} m_{11} - \lambda & m_{12} \\ m_{21} & m_{22} - \lambda \end{bmatrix} = 0 \rightarrow \lambda_{1,2} = \frac{1}{2} [(m_{11} + m_{22}) \pm \sqrt{4m_{12}m_{21} + (m_{11} - m_{22})^2}]$
- Once λ is known, you find the two eigenvectors (two columns of R) solving $\begin{bmatrix} m_{11} - \lambda & m_{12} \\ m_{21} & m_{22} - \lambda \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = 0$

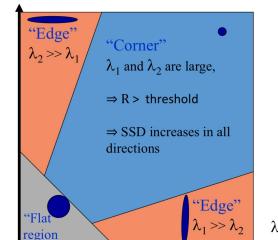
INTERPRETING EIGENVALUES

- A corner can be identified using a threshold : $R = \min(\lambda_1, \lambda_2) > \text{threshold}$
- R is called **cornerness function**
- The corner detector using this criterion is called "**Shi-Thomasi**" detector (1994).



- Computation of λ_1 and λ_2 is expensive \rightarrow Harris and Stephens suggested using a **different cornerness function**:

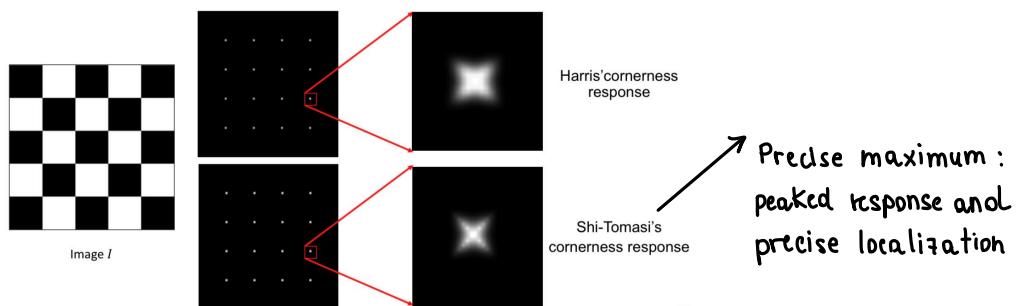
$$R = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2 = \det(M) - k \text{trace}^2(M)$$
 with k "magic number" $\in [0.04, 0.15]$
- The corner detector using this criterion is called "**Harris's detector**" (1988).



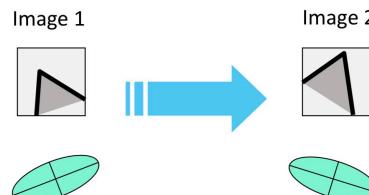
CORNER DETECTOR ALGORITHM (FULL)

- Compute derivatives in x and y (I_x, I_y) e.g. with Sobel filter.
- Compute I_x^2, I_y^2 and $I_x I_y$
- Convolve I_x^2, I_y^2 and $I_x I_y$ with a **box filter** to get $\sum I_x^2, \sum I_y^2$ and $\sum I_x I_y$ which are the entries of M (optionally use a **Gaussian filter** to avoid aliasing and give more weight to the central pixel).
- Compute cornerness R
- Threshold ($R > \text{threshold}$)
- Take the points of local maxima of R (**Non-Maxima Suppression**).

HARRIS VS SHI-TOMASI

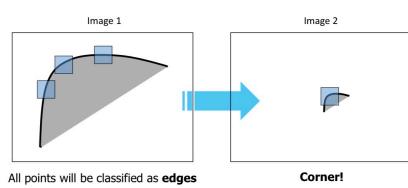


- Harris detector is **invariant to rotation**



- Harris detector is **NOT invariant to scaling**

Repeatability = $\frac{\# \text{correspondences detected}}{\# \text{correspondences present}}$ drop with scaling



In this case it is important to tune correctly the **Kernel size**

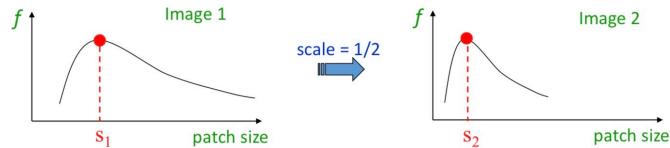
- Harris detector is **invariant to affine illumination changes** $\rightarrow I' = \alpha I + \beta \rightarrow I'_x = \alpha I_x, I'_y = \alpha I_y$
 α can affect the threshold
- Harris detector is **invariant to monotonic illumination changes** \rightarrow thanks to local non-maxima suppression
- Harris detector is **invariant to "small" viewpoint changes**

SCALE CHANGES

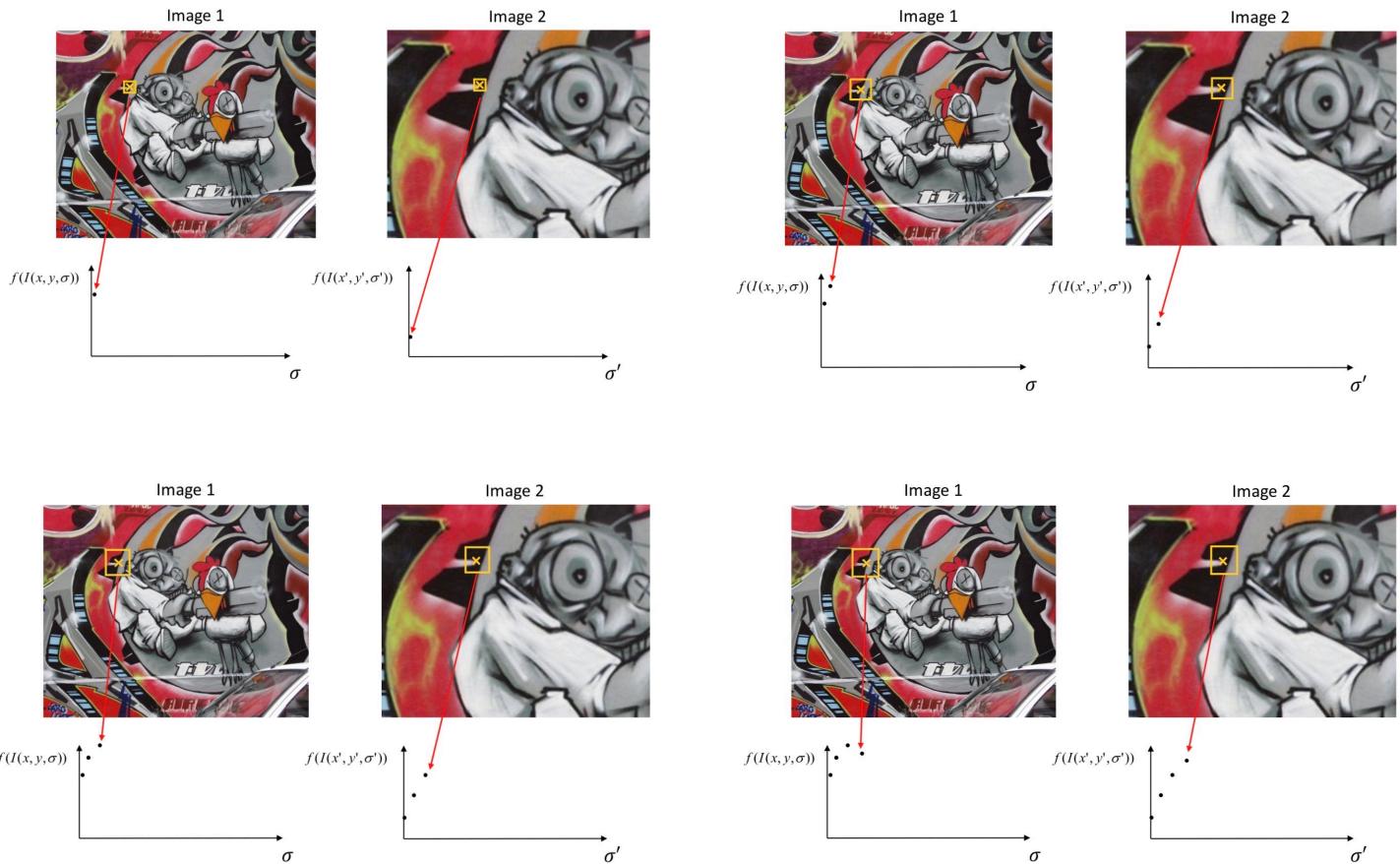
- **Question?** How can we scale image patches if they have different scale?
- **Solution 1:** Rescale the patch
 - Scale search is time consuming (needs to be done individually for all patches in one image)
 - Complexity is $(NS)^2$ assuming N features per image and S rescalings per feature
- **Solution 2:** automatic scale selection: automatically assign each feature its "own" scale

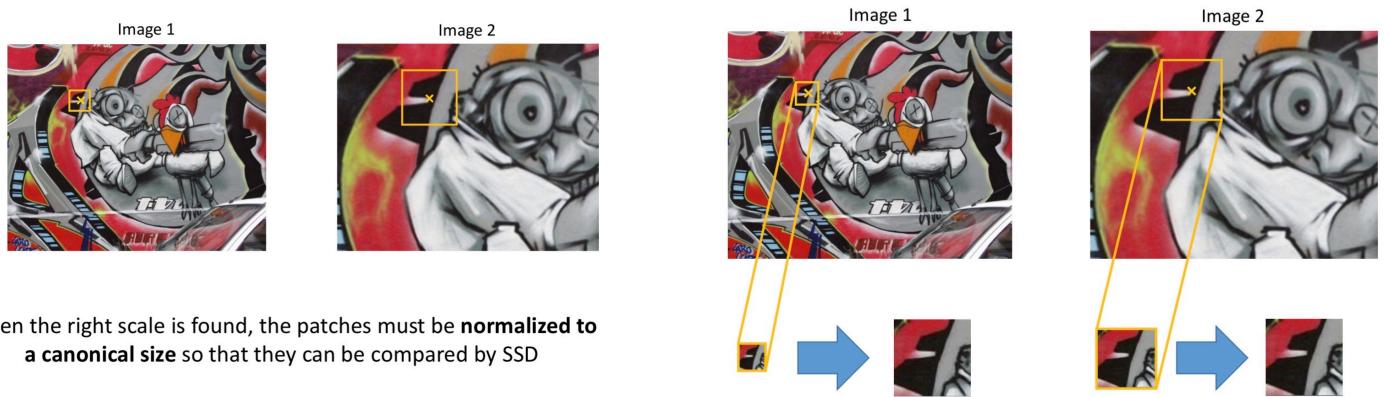
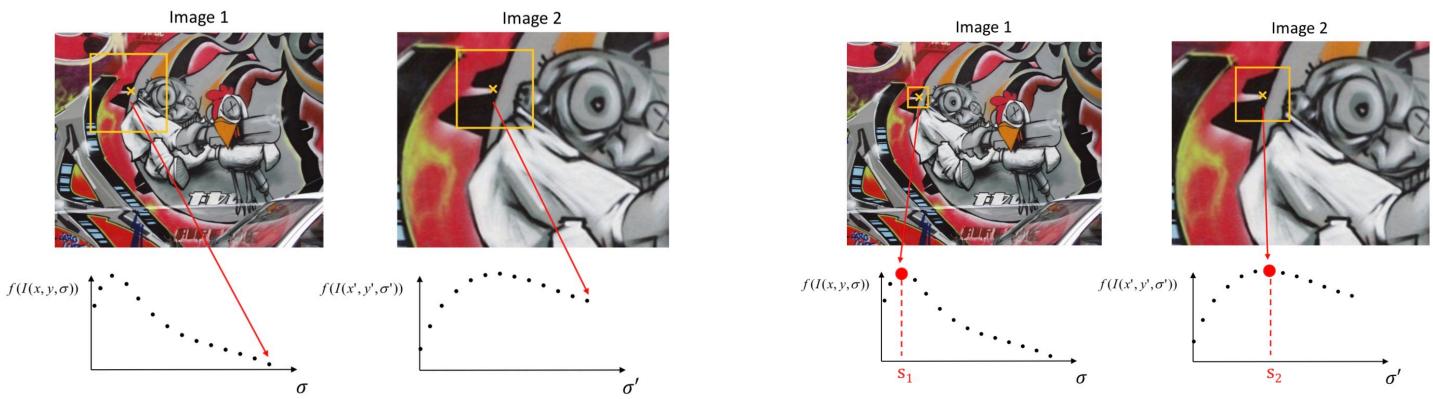
AUTOMATIC SCALE SELECTION

- Design a function on the image patch which is **scale invariant** (same value for corresponding patches even if they are at different scale)
- Find **local extrema** of this function
- The **patch size** at which the local extrema is reached should be **invariant** to image rescaling
- This scale invariant patch size is found in each image **independently**.



Example





- A "good function" for scale detection should have a single & sharp peak



- The ideal function for determining the scale is one that highlights **sharp discontinuities**.
- The solution is to convolve image with a **Kernel that highlights edges**.
- It has been shown that the **Laplacian of Gaussian** is optimal under certain assumptions:

$$\text{LoG}(x, y; b) = \nabla^2 G_b(x, y) = \frac{\partial^2 G_b(x, y)}{\partial x^2} + \frac{\partial^2 G_b(x, y)}{\partial y^2}$$

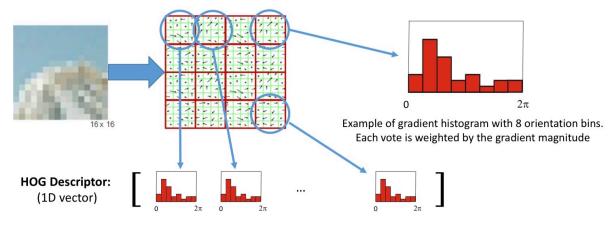
Works as a "blob" detector.

$$f = \text{Kernel} * \text{image}$$

- The correct scale is found as **local extrema** across consecutive smoothed patches (**increase b**).

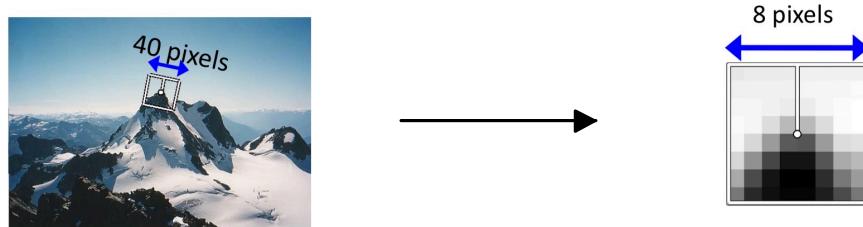
PATCH DESCRIPTORS

- "**Naive**" patch descriptor : patch of intensity (integer) values
- **Census descriptor** : see before (binary string)
- **HOG (Histogram of Oriented Gradients) descriptor** :
 - the image is divided into a grid of cells
 - for each cell a histogram of gradient directions is computed
 - the HOG descriptor is the concatenation of these histograms
 - differently from the previous two, it has float values



HOW TO ACHIEVE ROTATION AND SCALE INVARIANCE?

- Feature descriptors should be invariant to geometric and photometric changes
- **De-rotation**
 1. Determine the patch orientation (e.g. eigenvectors of M from Harris detector or dominant gradient direction (see next))
 2. De-rotate patch through "patch warping" (see next). This puts the patch in a canonical orientation.
- **Re-scaling**
 1. Detect scale using the LoG operator
 2. Rescale the patch to a canonical size (e.g. 8x8 pixels).



Example of patch after de-rotation and re-scaling

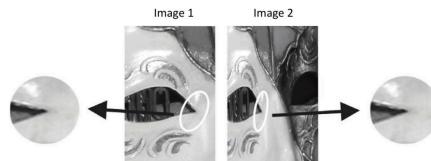
HOW TO DETERMINE PATCH ORIENTATION?

1. Multiply the patch with a Gaussian kernel to make the shape circular
2. Compute gradients vectors at each pixel
3. Build a histogram of gradient orientations **weighted by the gradient magnitude**. This histogram is the HOG descriptor.
4. Extract all local maxima in the histogram: each local maximum above a threshold is a candidate dominant direction
5. Construct a different Keypoint descriptor \forall dominant direction.

HOW TO ACHIEVE INVARIANCE TO SMALL VIEW-POINT CHANGES?

Affine warping provides invariance to small view-point changes

- The second moment matrix M of Harris can be used to identify the two directions of fastest and slowest SSD
- 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 canonical circular patch

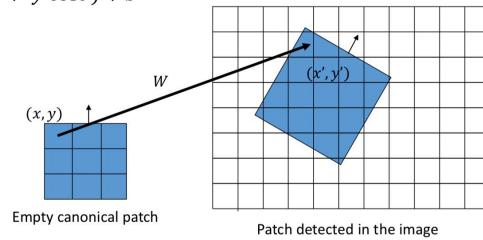


HOW TO WARP A PATCH?

- Start with an "empty" canonical patch (all pixels = 0)
- \forall pixel (x,y) apply the **warping function** $W(x,y)$ to compute the corresponding point in the source image. It will be in floating point and will fall between the image pixels.
- **Interpolate** the intensity values of the 4 closest pixels: **nearest neighbour, bilinear, bicubic interpolation**

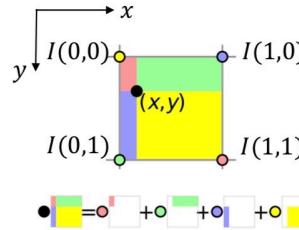
- Warping function W (counterclockwise rotation plus rescaling and translation):

$$\begin{aligned}x' &= s(x \cos\theta - y \sin\theta) + a \\y' &= s(x \sin\theta + y \cos\theta) + b\end{aligned}$$



BILINEAR INTERPOLATION

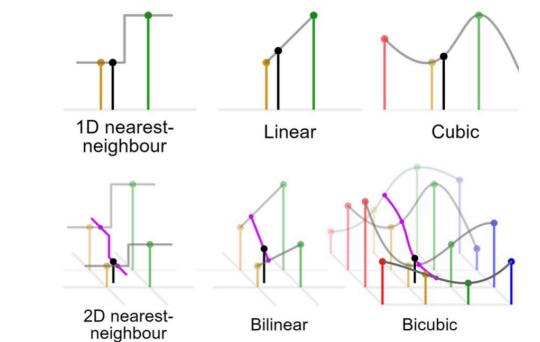
- Extension of linear interpolation in 2D
- Linear in each dimension but not a linear operator



Area of the opposite square

$$I(x, y) = I(0,0)(1-x)(1-y) + I(0,1)(1-x)y + I(1,0)x(1-y) + I(1,1)x(y)$$

$$\left. \begin{array}{l} \text{If } x = 0 \rightarrow I(0,y) = I(0,0)(1-y) + I(0,1)y \\ \text{If } y = 0 \rightarrow I(x,0) = I(0,0)(1-x) + I(1,0)x \end{array} \right\} \text{reduces to 1D linear (same for } x=1, y=1)$$



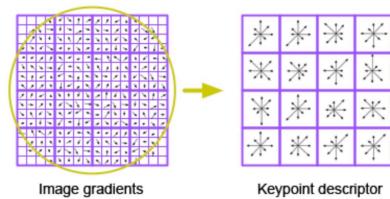
DISADVANTAGE OF PATCH DESCRIPTORS

- If the warp is not estimated accurately, very small errors in rotation, scale and view-point will affect matching score significantly.
- Computationally expensive: need to un warp every patch.

SCALE INVARIANT FEATURE TRANSFORM (SIFT, 2004)

Descriptor Computation:

- Multiply the patch by a Gaussian filter, compute dominant direction and de-rotate patch
- Consider a 16×16 pixel patch
- Compute HOG descriptor
 - Divide patch into 4×4 cells
 - Use 8 bins histograms
 - Concatenate all histograms into a single 1D vector
 - Resulting SIFT descriptor: $4 \times 4 \times 8 = 128$ values



Why 4×4 and 8 bins? Later the answer.

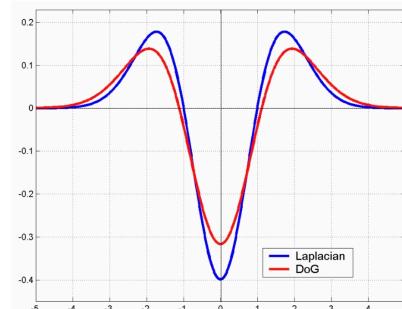
DESCRIPTOR NORMALIZATION

- The descriptor v is then normalized such that its ℓ_2 norm is 1: $\bar{v} = v / \|v\|$
- This guarantees that the descriptor is **invariant to linear illumination changes**.
- The descriptor is already **invariant to additive illumination** because it is based on gradients.
- ↳ It is **invariant to affine illumination changes**.

SIFT FACTS

- Can handle severe **viewpoint changes** (up to 50 deg out-of-plane rotations).
- Can handle even severe **non affine changes in illumination** (low to bright scenes).
- Computationally **expensive**: 10 fps on an i7 processor.
- SIFT uses the **Difference of Gaussians** Kernel instead of Laplacian of Gaussians (LoG) because computationally cheaper

$$\text{LoG}(x,y) \approx \text{DoG}(x,y) = G_{k^2}(x,y) - G_k(x,y)$$

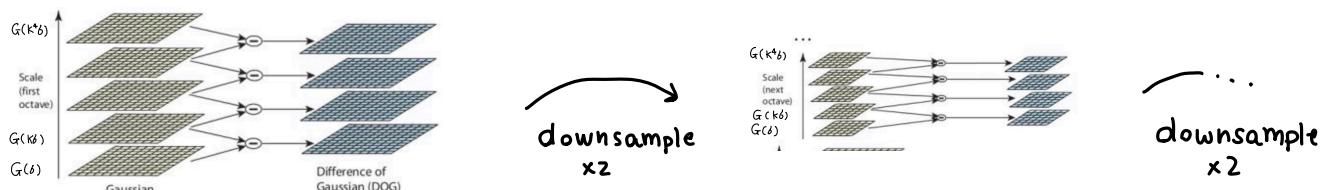


IMPLEMENTATION

1. Build a Space Scale Pyramid

The initial image is incrementally convolved with Gaussians $G(k^i b)$ to produce blurred images separated by a constant factor k in scale space:

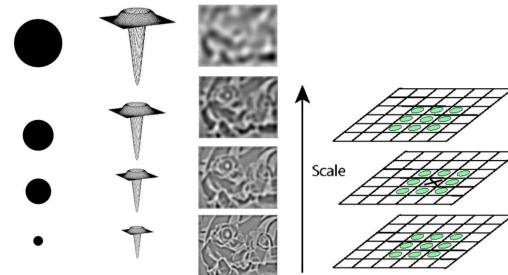
- The initial Gaussian $G(b)$ has $b = 1.6$.
- k is chosen $K = 2^{1/s}$ where s is the number of intervals into which each octave is divided.
- For efficiency, if $K^i = 2$, the image is downsampled by 2 and the procedure repeated 5 times.



2. Scale Space extrema detection

Detect local maxima and minima in scale spaces

- Each pixel is compared to 26 neighbors (in green below): its 8 neighbors in the current image + 9 neighbors in the adjacent upper and lower scale.
- If the pixel is a global maximum or minimum with respect to its 26 neighbors, then it is selected as SIFT features



SIFT RECAP

- A scale invariant feature transform
- An approach to detect and describe region of interest in an image.
- Invariant to 2D rotation and "reasonably" invariant to rescaling, viewpoint changes (up to 50°) and illumination
- Runs real-time but expensive (expensive steps are scale detection and descriptor extraction).
- SIFT outputs
 - the full histogram info: $4 \times 4 \times 8 = 128$ -element 1D vector,
 - location of the patch
 - scale (high scale = high blur = big blob)
 - orientation: angle of the patch.

FEATURE MATCHING

1. Define distance function ((z)SSD, (z)SAD,...) or Hamming distance for binary descriptors

2. Brute-force matching

- Compare each feature in I_1 with each feature in I_2 (N^2 comparisons)
- Take the closest descriptor \mapsto can return good score for false matches
- Better approach: compute ratio of distances to 1st to 2nd closest descriptor

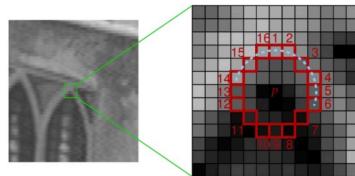
$$\text{accept if } \frac{d_1}{d_2} < \text{threshold } (\sim 0.8)$$

Why it works? If the closest is correct it is really correct while the second match is off. Instead when $d_1 \approx d_2$ then this is a noisy and ambiguous descriptor to match.

OTHER DESCRIPTORS

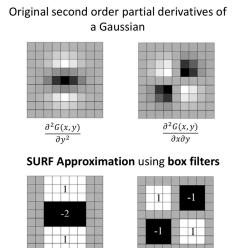
FEATURES FROM ACCELERATED SEGMENT TEST (FAST)

- Analyses intensities of 16 pixels around pixel p
- p is a FAST corner if a set of N contiguous pixels on the ring are brighter than the pixel intensity $I(p) + \text{threshold}$ or all darker than $I(p) - \text{threshold}$ ($N: 12$)
- A simple classifier is used to check the quality of the corners and reject the weak ones
- Fastest corner detector ever made : 100 million pixels/s
- Very sensitive to noise (why Harris is still common) → in fact thought as pre-processing for Harris.



SPEEDED UP ROBUST FEATURES (SURF) BLOB DETECTOR

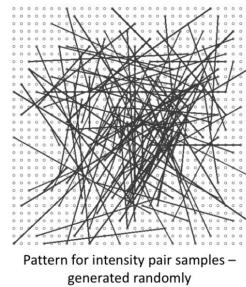
- Similar to SIFT but much faster
- Approximate DoG using box filters
- Results comparable with SIFT but
 - Faster computation
 - Generally shorter descriptors



BINARY ROBUST INDEPENDENT ELEMENTARY FEATURES (BRIEF)*

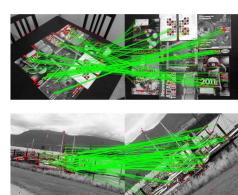
* can be applied to corner and blobs

- Goal: high speed computation and matching
- Binary description formation
 - smooth image
 - ↙ Keypoint (e.g. FAST)
 - sample 128 intensity pairs (p_1^i, p_2^i) $i \in [1, \dots, 128]$ within a squared patch around the keypoint
 - Create an empty 128-element descriptor
 - ↙ i^{th} pair:
 - if $I(p_1^i) < I(p_2^i)$, set i^{th} bit of descriptor to 1, else 0.
- Pattern is generated randomly (or learned) and reused for all patches
- PRO: binary descriptor → fast Hamming distance matching
- CON: not scale/rotation invariant.



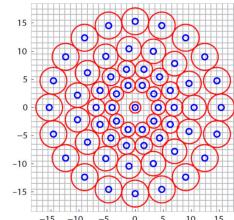
ORIENTED FAST AND ROTATED BRIEF (ORB)*

- Keypoint descriptor originally based on FAST.
- Binary descriptor based on BRIEF but adds an orientation component to make it rotation invariant



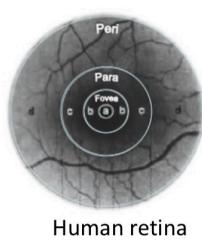
BINARY ROBUST INVARIANT SCALABLE KEYPOINTS (BRISK)*

- Keypoint detector based on FAST
- both rotation and scale invariant
- binary descriptor formed by pairwise intensity comparisons (like BRIEF) but on a radially symmetric sampling pattern
 - red circles: size of the smoothing Kernel applied
 - blue circles: smoothed pixel value used
- 10x faster than SURF, slower than BRIEF

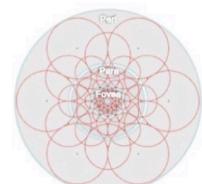


FAST RETINA KEYPOINT (FREAK)*

- rotation and scale invariant
- sampling pattern similar to BRISK but uses a more pronounced "retinal" (i.e. log-polar) sampling pattern inspired by the human retina: higher density at points near the center
- binary descriptor using pairwise comparison as in BRISK.
- pairs are learned (as in ORB).
- circles indicate size of smoothing Kernel
- coarse-to-fine matching (cascaded approach)
 - first compare first half of bits
 - if distance smaller than threshold proceed to compare next bits etc.
- Faster to compute, less memory than SIFT, SURF or BRISK.



Human retina



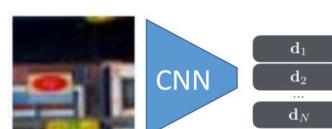
FREAK sampling pattern

LEARNED INVARIANT FEATURE TRANSFORM (LIFT)*

- learning-based descriptor
- rotation, scale, viewpoint and illumination invariant
- First a network predicts the patch orientation which is used to derotate the patch
- Then another NN is used to generate a patch descriptor ($\in \mathbb{R}^{128}$)
- Illumination invariance is achieved by randomizing illuminations during training
- LIFT > SIFT in repeatability.



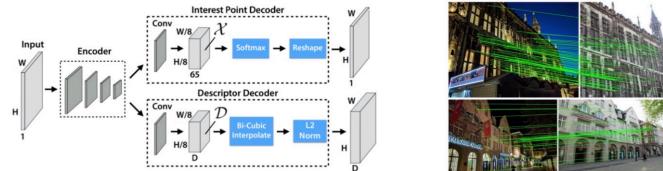
Keypoints with scales and orientations



neural network
predicts descriptor

SELF-SUPERVISED INTEREST POINT DETECTION AND DESCRIPTION (SUPERPOINT)

- Joint regression of keypoint location and descriptor.
- Self-supervised
- Trained on synthetic images and refined on homographies of real images
- Detector less accurate than SIFT and LIFT but descriptor outperforms SIFT and LIFT.
- But slower than SIFT and LIFT.



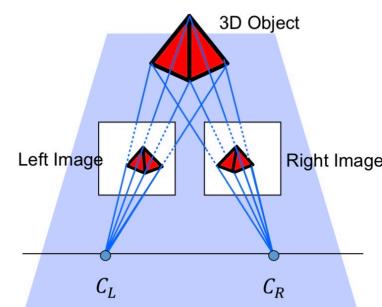
RECAP TABLE

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

STEREO VISION

Goal : recover the 3D structure by computing the intersection of corresponding rays.

- Tasks :
- Multiple View Geometry
 - 3D reconstruction (K, T, R known)
 - Structure from Motion (K, T, R unknown)
 - 2 View Geometry
 - Depth from stereo (K, T, R Known)
 - 2 view Structure from Motion (K, T, R unknown)

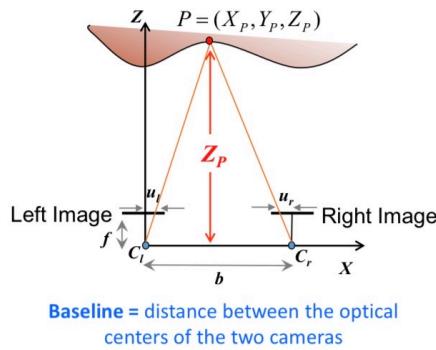


STEREO VISION

- **Goal:** find an expression of the 3D point coordinates as a function of the 2D image coords.
- **Assumptions:**
 - cameras are calibrated: both intrinsics and extrinsics are known
 - point correspondences are given

STEREO VISION -SIMPLIFIED CASE

- cameras are identical (**same intrinsics**);
- cameras are **aligned on the x-axis**.



From Similar Triangles:

$$\frac{f}{Z_p} = \frac{u_l}{X_p}$$

$$\frac{f}{Z_p} = \frac{-u_r}{b-X_p}$$

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

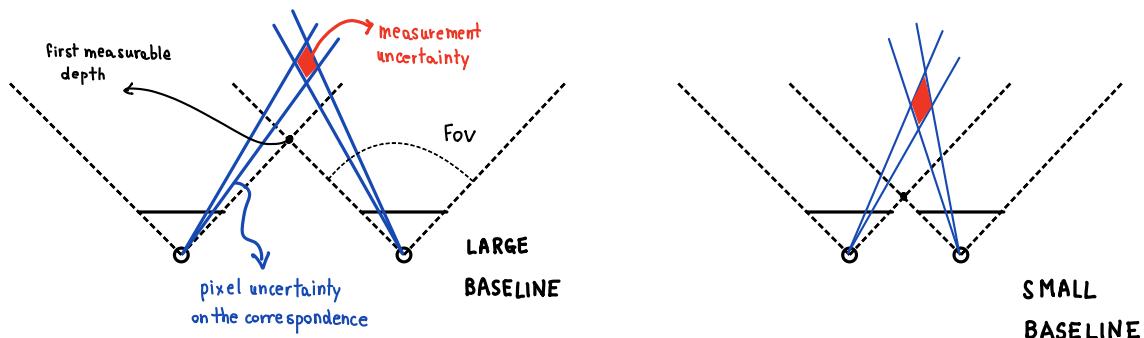
Disparity
difference in image location of the projection of a 3D point on two image planes

1. What's the max disparity of a stereo camera?
2. What's the disparity of a point at infinity?

LARGE VS SMALL BASELINE

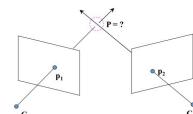
- **Large baseline**
 - small depth error
 - minimum measurable depth increases
↳ difficult for close objects

- **Small baseline**
 - Large depth error
 - minimum measurable depth increases
↳ easier for close objects.



STEREO VISION - GENERAL CASE

- Very hard to have perfectly aligned cameras
- Compute extrinsics and intrinsics with traditional methods (Tsai's or Zhang's method)
- **Triangulation**: problem of determining the 3D position of a point given a set of corresponding image locations and known camera points.
- In other words, we want to find the intersection of 2 rays but, because of noise and numerical errors, they can't meet exactly, so we need to compute an approximation.



$$- \alpha_z(\alpha_x \alpha_y) + \alpha_y(\alpha_z \alpha_x) \\ = 0$$

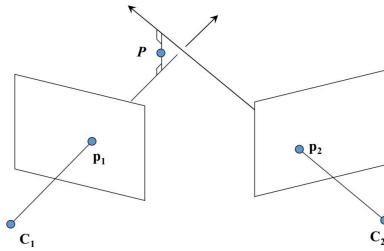
TRIANGULATION : LEAST SQUARE APPROXIMATION

- Left camera (here) assumed as world frame

$$\lambda_1 \begin{bmatrix} u_1 \\ v_1 \\ 1 \end{bmatrix} = \lambda_1 p_1 = k_1 [I|0] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} = M_1 \cdot P \Rightarrow p_1 \times \lambda_1 p_1 = p_1 \times M_1 \cdot P = [p_1]_x \cdot M_1 \cdot P = 0$$

$$\lambda_2 \begin{bmatrix} u_2 \\ v_2 \\ 1 \end{bmatrix} = \lambda_2 p_2 = k_2 [R|T] \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} = M_2 \cdot P \Rightarrow p_2 \times \lambda_2 p_2 = p_2 \times M_2 \cdot P = [p_2]_x \cdot M_2 \cdot P = 0$$

- geometric interpretation : P is computed as midpoint of the shortest segment connecting the two lines.



TRIANGULATION : NON LINEAR REFINEMENT

- Initialize P using least square approximation
- refine P by minimizing sum of left and right squared reprojection error

$$P = \arg \min_P \|p_1 - \pi(P, K_1, I, 0)\|^2 + \|p_2 - \pi(P, K_2, R, T)\|^2$$

- Can be minimized using Levenberg-Marquardt (more robust than Gauss-Newton to local minima)

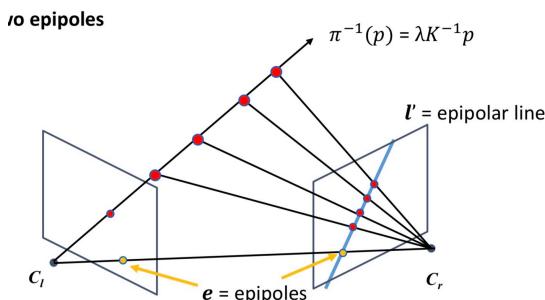
CORRESPONDENCE PROBLEM

Correspondence Problem : given a point p_L on the left image, find correspondence p_R on right image.

Search Problem : 2D exhaustive search is very expensive (N^2 comparisons for N descriptors).

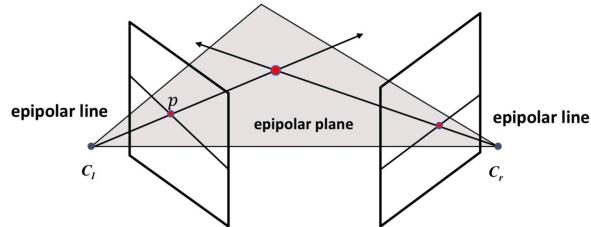
We can use geometrical consistency checks

- potential matches for p have to lie on the corresponding **epipolar line** l' .
- the **epipolar line** is the back-projection of the ray $\pi^{-1}(p)$ onto the camera image.
- the **epipole** is the projection of the optical center on the other camera image.
- a stereo camera has two epipoles.



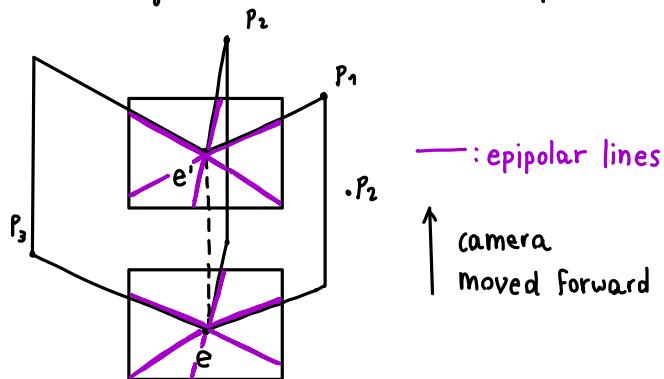
THE EPIPOLAR CONSTRAINT

- Camera centers C_l, C_r and the image point p define the **epipolar plane**.
- The intersection of the epipolar plane with the two image planes are called **epipolar lines**.
- **Epipolar constraint**: Corresponding points must lie along the epipolar line \rightarrow makes the correspondence problem a 1D search.



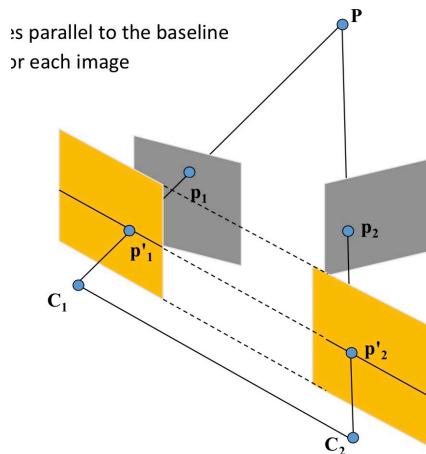
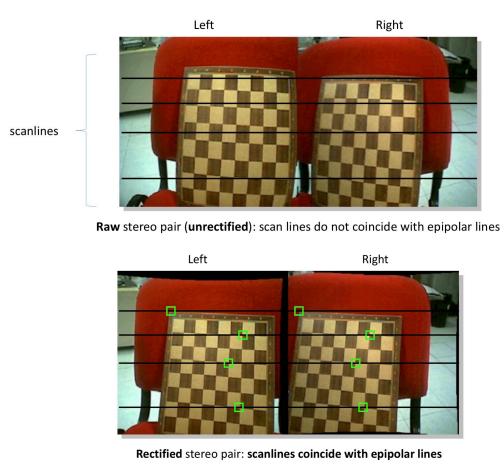
Example: Forward Motion

- two consecutive images from the same camera
- epipole has the same coordinate in both images
- epipolar lines "radiating" from the epipole (focus expansion)



STEREO RECTIFICATION

- no camera is perfectly calibrated
- more computationally efficient if scan-lines == epipolar lines
- **Stereo rectification**: wraps the left and right images into new "rectified" images such that the epipolar lines coincide with scanlines.
 - It works by computing 2 homographies which achieve the result



STEREO RECTIFICATION ALGORITHM

- The projection equation is $\lambda p = K[R|T]P$ where $R=R_{cw}$ and $T=T_{cw}$
- We use R_{cw} and T_{cw} as T_{cw} already gives us the camera position in world frame. Then the projection equation is $\lambda p = KR^{-1}(P-T) = KR^{-1}(P-C)$
- The goal of stereo rectification is to warp camera images such that they are **coplanar** (horizontal epipolar lines), so they have same \hat{R} and same intrinsics (vertically aligned) so they have same \hat{K}

$$\begin{aligned} \lambda_L \begin{bmatrix} u_L \\ v_L \\ 1 \end{bmatrix} &= K_L R_L^{-1} \left(\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} - C_L \right) & \text{OLD LEFT CAMERA} & \lambda_R \begin{bmatrix} u_R \\ v_R \\ 1 \end{bmatrix} &= K_R R_R^{-1} \left(\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix} - C_R \right) & \text{OLD RIGHT CAMERA} \\ \hat{\lambda}_L \begin{bmatrix} \hat{u}_L \\ \hat{v}_L \\ 1 \end{bmatrix} &= \hat{K} \hat{R}^{-1} \left(\underbrace{\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}}_{\text{New image plane}} - C_L \right) & \text{NEW LEFT CAMERA} & \hat{\lambda}_R \begin{bmatrix} \hat{u}_R \\ \hat{v}_R \\ 1 \end{bmatrix} &= \hat{K} \hat{R}^{-1} \left(\underbrace{\begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}}_{\text{New image plane}} - C_R \right) & \text{NEW RIGHT CAMERA} \\ \lambda_L R_L K_L^{-1} \begin{bmatrix} u_L \\ v_L \\ 1 \end{bmatrix} & & & \lambda_L R_R K_R^{-1} \begin{bmatrix} u_R \\ v_R \\ 1 \end{bmatrix} & & \\ \Rightarrow \hat{\lambda}_L \begin{bmatrix} \hat{u}_L \\ \hat{v}_L \\ 1 \end{bmatrix} &= \underbrace{\lambda_L \hat{K} \hat{R}^{-1} R_L K_L^{-1}}_{\text{homography of left camera}} \begin{bmatrix} u_L \\ v_L \\ 1 \end{bmatrix} & & \hat{\lambda}_R \begin{bmatrix} \hat{u}_R \\ \hat{v}_R \\ 1 \end{bmatrix} &= \underbrace{\lambda_L \hat{K} \hat{R}^{-1} R_R K_R^{-1}}_{\text{homography of right camera}} \begin{bmatrix} u_R \\ v_R \\ 1 \end{bmatrix} & \end{aligned}$$

- How do we choose \hat{K} and \hat{R} ?
- Good choices are

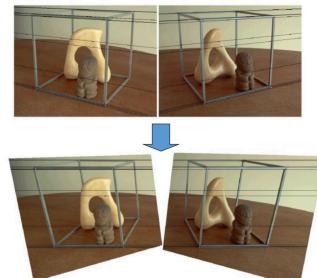
$$\hat{K} = \frac{K_L + K_R}{2} \quad \hat{R} = [\hat{r}_1, \hat{r}_2, \hat{r}_3] \quad \text{with} \quad \hat{r}_1 = \frac{C_R - C_L}{\|C_R - C_L\|} : \text{new image planes are parallel to baseline}$$

$$\hat{r}_2 = r_3 \times \hat{r}_1 : \text{where } r_3 \text{ is 3rd column of } R_L$$

$$\hat{r}_3 = \hat{r}_1 \times \hat{r}_2$$

Recap of the method

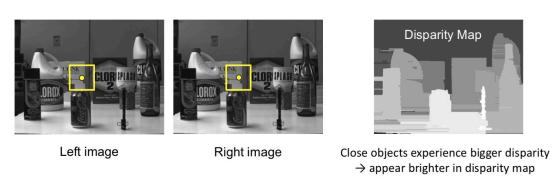
- First compute lens distortion
- then compute homographies and rectify
- use bilinear interpolation for warping.



DENSE STEREO CORRESPONDENCE : DISPARITY MAP

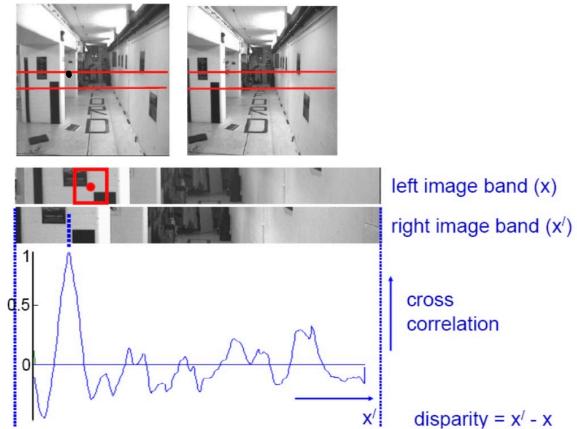
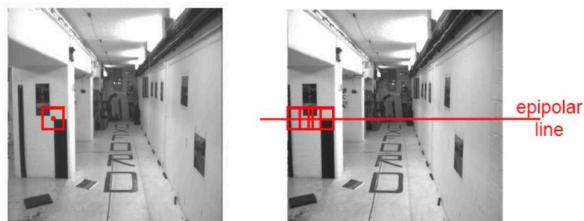
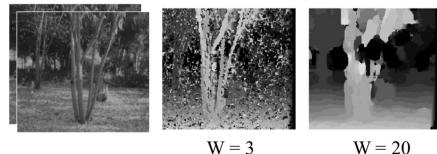
- Rectify stereo pair (if not already rectified)
- For each pixel in the left image, find its corresponding point in the right image.
- Compute the **disparity** for each found pair of correspondences
- Visualize it as a grayscale or color-coded image: **disparity map**

- Once the Stereo pair is rectified, the depth of each point can be computed recalling that $Z_p = b \cdot f / \text{disparity}$



CORRESPONDENCE PROBLEM - CONTINUED

- Once rectified correspondence search can be done along the same scanline
- To **average noise effects**, use a window around point of interest with the assumption that neighbouring pixels have similar intensity
- Find correspondence minimizing metric (\approx NCC, \approx SSD, \approx SAD, Census Transform ...)
- Texturless regions are not distinctive: high ambiguity for matches \mapsto increase window size
 - Smaller window
 - \times more noise
 - \checkmark more detail
 - Larger window
 - \times smoother disparity map
 - \checkmark less detail



CHALLENGES

- multiple possible matches for the same epipolar line
- occlusions and repetitive patterns
- non-lambertian surfaces (specularities) and texturless surfaces

SEMI-GLOBAL MATCHING (SGM)

- Popular open-source algo that computes dense disparity maps from a rectified stereo pair
- main idea**: perform coarse to fine block matching followed by regularization (smoothing)
- the estimated disparity map is a piece-wide smooth surface passing through the initial disparity map

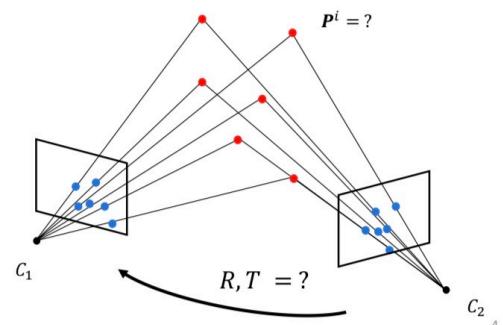


STRUCTURE FROM MOTION (SFM)

Problem formulation: given a set of n point correspondences between two images $\{p_1^i = (u_1^i, v_1^i), p_2^i = (u_2^i, v_2^i)\}$ where $i=1, \dots, n$ the goal is to simultaneously

- estimate the 3D points P_i
- the camera relative motion parameters (R, T)
- and the camera intrinsics that satisfy

$$\lambda_1 \begin{bmatrix} u_1^i \\ v_1^i \\ 1 \end{bmatrix} = K_1 [I|0] \begin{bmatrix} x_w^i \\ y_w^i \\ z_w^i \\ 1 \end{bmatrix} \quad \lambda_2 \begin{bmatrix} u_2^i \\ v_2^i \\ 1 \end{bmatrix} = K_2 [R|T] \begin{bmatrix} x_w^i \\ y_w^i \\ z_w^i \\ 1 \end{bmatrix}$$



CALIBRATED CASE

- K_1, K_2 known
- We use $\bar{P} = K^{-1}P$ for convenience

$$\lambda_1^i \begin{bmatrix} \bar{u}_1^i \\ \bar{v}_1^i \\ 1 \end{bmatrix} = [I|0]\bar{P}$$

$\underbrace{\bar{P}_1}_{\bar{P}}$

Vector in space
with same direction as ray with $f=1$

$$\lambda_2^i \begin{bmatrix} \bar{u}_2^i \\ \bar{v}_2^i \\ 1 \end{bmatrix} = [R|T]\bar{P}$$

$\underbrace{\bar{P}_2}_{\bar{P}}$

SCALE AMBIGUITY

- If we rescale the entire scene and camera views by a constant factor, the projections (in pixels) of the scene points in both images remain exactly the same.

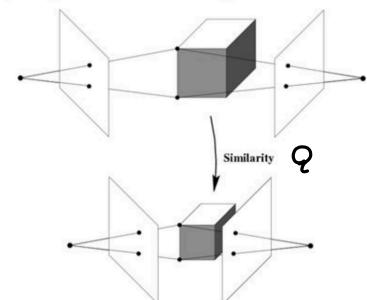
$$P_1 = K_1 T_1 P \quad T_1 = [R_1 | T_1]$$

$$P_2 = K_2 T_2 P$$

Given an affine scene transformation Q and its inverse Q^{-1}

$$P_1 = K_1 T_1 Q^{-1} Q P = K_1 \tilde{T}_1 \tilde{P}$$

$$P_2 = K_2 T_2 Q^{-1} Q P = K_2 \tilde{T}_2 \tilde{P}$$



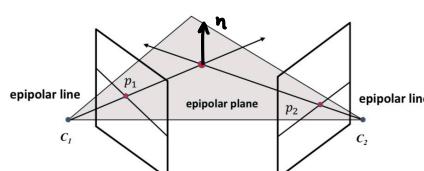
- In SFM it is therefore not possible to recover the absolute scale of the scene.
- Thus only **5 degrees of freedom** are measurable: Rotation (3 dofs) and translation direction up to scale (2 dofs)

→ given n correspondences

- $4n$ knowns (correspondence equations $\lambda P = \lambda' P'$)
- $5 + 3n$ unknowns

$$\rightarrow 4n \geq 5 + 3n \Rightarrow n \geq 5$$

- Can we estimate camera extrinsics without having to find the 3D position of points? Yes.
- How? We can use the **epipolar geometry**: the two rays and the baseline must be co-planar



$$\bar{p}_1 = \begin{bmatrix} \bar{u}_1 \\ \bar{v}_1 \\ 1 \end{bmatrix} \quad \bar{p}_2 = \begin{bmatrix} \bar{u}_2 \\ \bar{v}_2 \\ 1 \end{bmatrix} \quad T = C_2 - C_1 \quad \bar{p}_1, \bar{p}_2 \text{ and } T \text{ are coplanar}$$

\bar{p}_1 : normalized point in camera 1 reference frame
 \bar{p}_2 : normalized point in camera 2 reference frame

The extrinsics of \bar{p}_2 has $R_{C_2}w = R_{C_2}C_1$, $T_{C_2}w = T_{C_2}C_1$ since camera 1 is our world frame.
We can use these transformations to express \bar{p}_1 in C_2 frame and compute the normal to the plane:

$${}^C_2 \bar{p}_1 = R \bar{p}_1 + T$$

$$n = T \times (R \bar{p}_1 + T) = [T]_x R \bar{p}_1$$

And since \bar{p}_2 belongs to this plane:

$\bar{p}_2^T [T]_x R \bar{p}_1 = 0$	Epipolar Constraint
-------------------------------------	---------------------

$E = [T]_x R$	Essential Matrix
---------------	------------------

rotation doesn't change rank

Note that $\text{rank}(E) = \text{rank}(E^T) = \text{rank}([T]_x) = 2$

SOME HISTORY

- Kruppa showed in 1913 that 5 images correspondences is the minimal case and that there can be up to 11 solutions
- In 1988, Demazure shows that there are at most 10 distinct solutions.
- In 1996, Philipp proposed an iterative algo to find these solutions
- In 2004, Nister proposed the first efficient and non iterative solution.
- The first popular solution uses the 8-point algorithm or Longuet-Higgins algo.

THE 8 POINT ALGORITHM

Each pair of point correspondences \bar{p}_1 and \bar{p}_2 provide a linear equation $\bar{p}_2^T E \bar{p}_1$

$$\bar{u}_2 \bar{u}_1 e_{11} + \bar{u}_2 \bar{v}_1 e_{12} + \bar{u}_2 v_1 e_{21} + \bar{v}_2 \bar{u}_1 e_{22} + \bar{v}_2 \bar{v}_1 e_{32} + \bar{v}_2 v_1 e_{23} + \bar{u}_1 e_{31} + \bar{v}_1 e_{32} + e_{33} = \dots \quad q_i \quad \left[\begin{array}{c} e_{11} \\ e_{12} \\ e_{13} \\ \vdots \\ e_{33} \end{array} \right] \quad e_{ij} : \text{element of } E$$

For n correspondences q^i :

$$\begin{bmatrix} q^0 \\ q^1 \\ \vdots \\ q^n \end{bmatrix} \left[\begin{array}{c} e_{11} \\ e_{12} \\ e_{13} \\ \vdots \\ e_{33} \end{array} \right] = Q \cdot \bar{E} \quad \begin{matrix} \downarrow & \text{unknown} \\ \downarrow & \text{known} \end{matrix}$$

- Minimum solution: Q must have rank 8 for a unique (non trivial) solution \bar{E} . As each point (if not coplanar) provides an independent equation, we need more than $n \geq 8$ points.
 - An over-determined solution is to minimize $\|Q\bar{E}\|^2$ subject to $\|\bar{E}\|^2 \leq 1$ (avoid to select $\bar{E} = 0$)
 - The solution is the eigenvector corresponding to the smallest eigenvalue of the matrix $Q^T Q$ since it minimizes $\|Qx\| = x^T Q^T Q x = x^T \lambda x = \lambda_{\min} x^T x = \lambda_{\min} \|x\| = \lambda_{\min} = 1$
 - It can be solved via SVD
- $$Q = U S V^T \quad Q^T Q = V S^T U^T U S V = V S^2 V$$
- $$\rightarrow [U, S, V] = \text{svd}(Q) \quad x = V E[:, 9] \quad (\text{last eigenvector}) \rightarrow E = \text{reshape}(x, 3, 3)$$

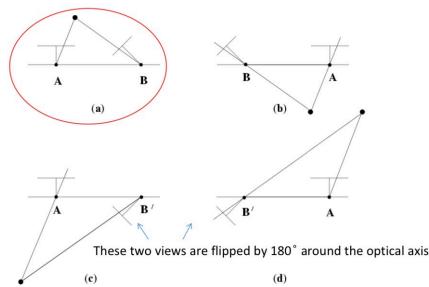
EXTRACT R AND T

- E has rank 2 : $\text{rank}(E) = \text{rank}([T]_x R) = \text{rank}([T]_x) = 2$
- Enforce rank=2 setting the minimum eigenvalue to 0 : $E = U \Sigma^{\frac{1}{2}} V^T = U \begin{bmatrix} b_1 & b_2 & \cancel{b_3} \end{bmatrix} V^T$

$$\hat{T} = U \begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \Sigma V^T \quad \hat{T} = \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & t_x \\ -t_y & t_x & 0 \end{bmatrix} \Rightarrow \hat{t} = \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$$

$$\hat{R} = U \begin{bmatrix} 0 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} V^T \quad T = K_2 \hat{t} \\ R = K_2 \hat{R} K_1^{-1}$$

- There are 4 possible solutions but only one solution where points are in front of both cameras.



STRUCTURE FROM MOTION (UNCALIBRATED CASE)

- So far we have used normalized pixel coordinates $\bar{p} = K^{-1} p$
- Then $\bar{p}_2^T E \bar{p}_1 = p_2^T K_2^{-1} E K_1^{-1} p_1 = 0$
- The same algorithm holds and we call $F = K_2^{-1} E K_1^{-1}$ the **Fundamental Matrix**.
- The advantage is that we work directly in pixel coordinates.

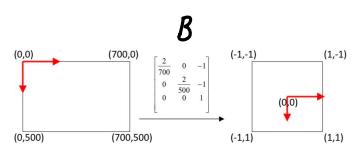
PROBLEMS WITH 8-POINT ALGORITHM

- Poor numerical conditioning which makes results very sensitive to noise.
- Can be fixed by rescaling the data : **normalized 8 point algorithm**.

NORMALIZED 8-POINT ALGORITHM

- Transform image coordinates so that they are in the range $[-1,1] \times [-1,1]$
 $\hat{p} = B p$
- Estimate normalized \hat{F} from \hat{p}_1 and \hat{p}_2
- Compute unnormalized F from \hat{F}

$$\hat{p}_2^T \hat{F} \hat{p}_1 = 0 \rightarrow p_2^T \underbrace{B_2^T \hat{F} B_1}_{F} p_1 \rightarrow F = B_2^T \hat{F} B_1$$



- In the original paper it was proposed to rescale each points sets such that each has centroid 0 and the standard deviation $\sqrt{2}$: $\hat{p}_i^i = \sqrt{2}/6 (p_i^i - \mu) ; \mu = \sum p_i^i / N ; \delta = 1/N \sum_{i=1}^N \|p_i^i - \mu\|^2$
- This transformation can be expressed in matrix form using homogeneous coordinates

$$\hat{p}_i^i = \begin{bmatrix} \sqrt{2}/6 & 0 & -\frac{\sqrt{2}}{6}\mu_x \\ 0 & \sqrt{2}/6 & -\frac{\sqrt{2}}{6}\mu_y \\ 0 & 0 & 1 \end{bmatrix} p_i^i$$

- Can we extract R, T, K_1, K_2 from F ?
In general no: infinite solutions exists.
- However if the coordinates of the principal points of each camera are known and the two camera have the same focal length f in pixels, then R, T and f can be determined uniquely

ERROR MEASURES

The quality of the estimated Essential or Fundamental matrix can be measured using different error metrics:

- Algebraic error
 - Directional error
 - Epipolar line distance
 - Reprojection error
- $\left. \begin{array}{l} \\ \\ \end{array} \right\}$ can only be ϕ for 8 points (or more and no noise)

Algebraic error: $\text{err} = \|Q\bar{E}\|^2 = \sum_{i=1}^N (\bar{p}_2^{i^T} E \bar{p}_1^i)$

Also

$$\|\bar{p}_2^{i^T} E \bar{p}_1^i\| = \|\bar{p}_2^T \cdot (E \bar{p}_1)\| = \|\bar{p}_2\| \|\bar{E} \bar{p}_1\| \cos(\theta) \quad \text{angle between } \bar{p}_2 \text{ and normal}$$

\Rightarrow if this is non-zero \bar{p}_1, \bar{p}_2 and n are not coplanar (but depends also on norms $\|\bar{p}_i\|$!)

Directional error: cosine of the angle from the epipolar plane $\text{err} = \sum_{i=1}^N \cos(\theta_i)^2$ (does not depend on norms)

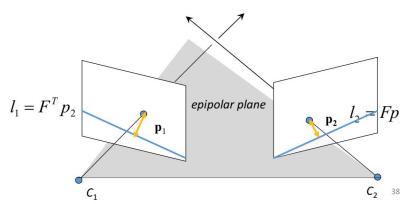
It is obtained normalizing the algebraic error: $\cos(\theta) = \bar{p}_2^T E \bar{p}_1 / \|\bar{p}_2^T E \bar{p}_1\|$

Epipolar line distance: sum of squared epipolar line to point distances $\rightarrow \text{err} = \sum_{i=1}^N d(p_1^i, l_1^i)^2 + d(p_2^i, l_2^i)^2$

Cheaper than reprojection error because it does not require point triangulation

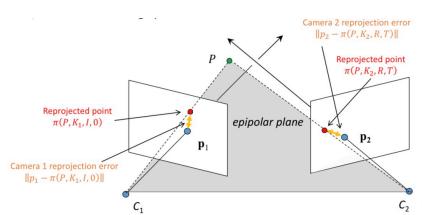
$\bar{p}_2^T F \bar{p}_1$ can be read as

- $\bar{p}_2^T \ell_2$ with ℓ the coefficients of epipolar line in the second image
- $\bar{p}_1^T F^T \bar{p}_2 = \bar{p}_1^T \ell_1$ with $\ell_1 = F^T \bar{p}_2$ the coefficients of the epipolar line in the first image



Reprojection error: sum of the squared Reprojection errors $\rightarrow \text{err} = \sum_{i=1}^N \|p_1^i - \pi(p_i^i, K_1, I, o)\|^2 + \|p_2^i - \pi(p_i^i, K_2, R, T)\|^2$

- More expensive to evaluate than previous because requires triangulation of points P^i in 3D!
- However most popular because it is the most accurate



ROBUST ESTIMATION

Often matches contain outliers :

- repetitive features
- changes in view point (scale also) and illumination
- image noise
- occlusions
- moving objects / blur

EXPECTATION MAXIMIZATION

Method for model fitting in presence of outliers.

Example: **line fitting**

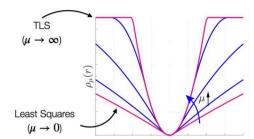
1. Estimate line parameters using least squares : $\min \sum r_i^2$ (r_i = point line distance)
2. Calculate new r_i V point and assign each a weight (e.g. $e^{-r_i^2}$) representing the likelihood that such an assignment is correct. (estimates the **expectation**)
3. Re-estimate line parameters (e.g. using weighted least squares : $\min \sum w_i r_i^2$). (**Maximization Step**)
4. Iterate 2-3 until convergence
5. Select the inliers as data points with weight higher than threshold.

Problems

- Very sensitive to the initial condition:
 - EM selects the initial condition minimizing sum of residuals
 - While this is a convex function, the result is influenced by a few large error values (distant outliers)
 - EM converges to the wrong solution if the initial condition is too far from the true one
- Alternatives
 - GNC algorithm
 - RANSAC algorithm

GRADUATED NON CONVEXITY (GNC) ALGORITHM

- Idea: optimize a surrogate function $\mathcal{L}_{\mu}(r_i)$ where μ controls the amount of non-convexity.
 - Start by solving the non robust convex least square ($\mu \rightarrow 0$)
 - At each iteration gradually increase non convexity ($\mu \rightarrow \infty$) and recompute w_i
 - It is shown to be robust up to 90% of outliers with five times fewer iterations than RANSAC.
 - However RANSAC can cope with more than 90% outliers



RANSAC (RANdom SAmple Consensus)

- RANSAC is the standard method for model fitting in the presence of outliers
- It is non deterministic: you get a different result every time you run it. (in contrast to EM and GNC)
- It is non sensitive to the initial condition and does not get stuck in local maxima

Example: line fitting

1. Select a sample of 2 points at random
2. Select the model parameters from the sample
3. Calculate residual \vee data point
4. Select data that support current hypothesis
5. Repeat 1-4 K times
6. Select the set with maximum number of inliers obtained with K iterations.

HOW MANY ITERATIONS?

- $w := \text{number of inliers}/N$ (probability of selecting an inlier)
- $N := \text{total number of data points}$
- Assuming that 2 points necessary to find a line:
 - $w^2 = P(\text{both selected points are inliers})$
 - $1-w^2 = P(\text{at least one of this point is an outlier})$
- Let K be the number of RANSAC iterations executed so far
 - $(1-w^2)^K = P(\text{RANSAC never selected two points that are both inliers})$
- $p := \text{probability of success}$
 $1-p = (1-w^2)^K$ and therefore

$$K = \frac{\log(1-p)}{\log(1-w^2)}$$

- In practice we only need a rough estimate of w
- More advanced techniques estimate w and adaptively update it at every iteration.

RANSAC FOR SFM

1. What is the model in SFM?
 - The Essential or Fundamental matrix
 - Alternatively R and T
2. What's the minimum number of points to estimate the model?
 - Theoretically 5 for calibrated cameras
 - 8 for the 8-point algo (both calibrated and not)
3. How do we compute the distance of a point from the model?
 - Algebraic error
 - Directional error
 - ...

RANSAC ITERATIONS K vs S.

Assume $p = 99\%$ and $\varepsilon = 1-w = 50\%$ (fraction of outliers)

- 8 point RANSAC $\rightarrow K = 1177$ Good method but too many iterations
- 5 point RANSAC $\rightarrow K = 145$ Efficient but can return up to 10 solutions of E
- 2 point RANSAC (e.g. line fitting) $\rightarrow K = 16$

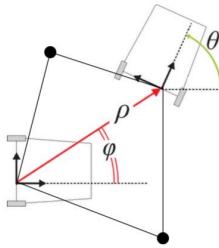
MOTION CONSTRAINTS

We can use motion constraints to reduce the solution space.

Planar Motion

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

$$R = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad T = \begin{bmatrix} \rho \cos \varphi \\ \rho \sin \varphi \\ 0 \end{bmatrix}$$



Observe that E has 2DoF (θ, φ , because ρ is the scale factor); thus, 2 correspondences are sufficient to estimate θ and φ [Ortin, 2001]

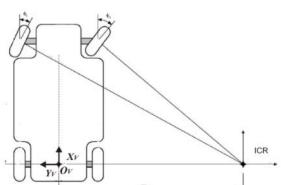
$$E = [T_x]R = \begin{bmatrix} 0 & 0 & \rho \sin(\varphi) \\ 0 & 0 & -\rho \cos(\varphi) \\ -\rho \sin(\varphi - \theta) & \rho \cos(\varphi - \theta) & 0 \end{bmatrix}$$

"2-Point RANSAC", Ortiz & Montiel, Indoor robot motion based on monocular images, Robotica, 2001. [PDF](#).

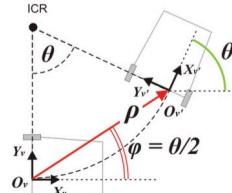
49

We can use less than 2 correspondences if we exploit wheeled vehicles with non-holonomic constraints.

- Wheeled vehicles like cars follow locally planar circular motion about the Instantaneous Center of Rotation (ICR)
- $\varphi = \theta/2 \rightarrow$ only 1 DoF, thus only 1 correspondence is needed
- This is the smallest parametrization possible



Example of Ackerman steering principle



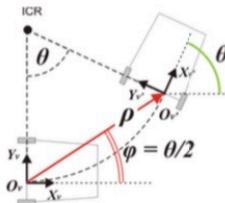
Locally-planar circular motion

Let's compute the Epipolar Geometry

$$E = [T_x]R \quad \text{Essential matrix}$$

$$R = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad T = \begin{bmatrix} \rho \cos \frac{\theta}{2} \\ \rho \sin \frac{\theta}{2} \\ 0 \end{bmatrix}$$

$$\bar{p}_2^T E \bar{p}_1 = 0 \quad \text{Epipolar constraint}$$



Locally-planar circular motion

Notice that ρ can be cancelled out

$$p_2^T E p_1 = 0 \Rightarrow \sin\left(\frac{\theta}{2}\right) \cdot (u_2 + u_1) + \cos\left(\frac{\theta}{2}\right) \cdot (v_2 - v_1) = 0$$

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

Scaramuzza, 1-Point-RANSAC Structure from Motion for Vehicle-Mounted Cameras by Exploiting Non-Holonomic Constraints, International Journal of Computer Vision, 2011. [PDF](#).

54

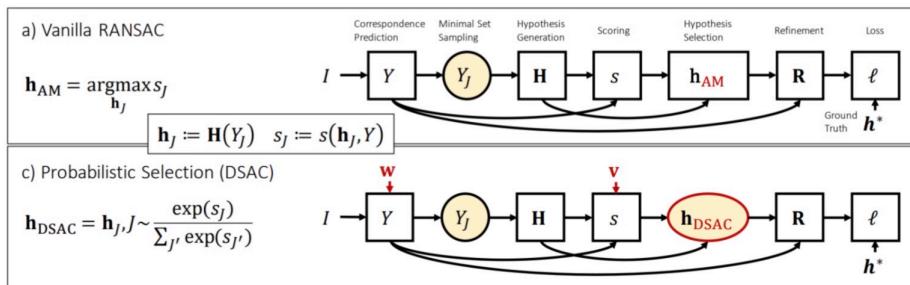
1-POINT RANSAC ALGORITHM

1. Compute $\theta \wedge$ point correspondence
2. Visualize histogram and take θ as the mean \rightarrow only 1 iteration for removing outliers ($\theta' \notin [\bar{\theta} \pm \theta]$) ↗ threshold
3. 1-Point RANSAC is only used to find inliers, motion is then estimated from them in 6 DoF.

DIFFERENTIABLE RANSAC

- RANSAC is not differentiable since it relies on selecting an hypothesis based on maximizing the number of inliers.
- DSAC introduces differentiable sample consensus

- Replaces **max** with **softmax** operator for model selection

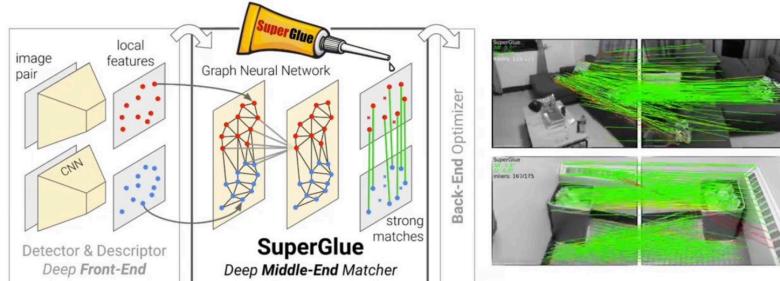


DEEP FUNDAMENTAL MATRIX ESTIMATION

- **Input**: two sets of noisy local features (coordinates + descriptors) contaminated by outliers
- **Output**: fundamental matrix
- **Idea**: solve a weighted homogeneous LS problem where robust weights are estimated using deep networks.
- **Robust**: handles extreme wide-baseline image pairs.

SUPER GLUE : LEARNING FEATURE MATCHING WITH GRAPH NEURAL NETWORKS

- **Input**: two sets of noisy local features (coordinates + descriptors) contaminated by outliers
- **Output**: strong & outlier free matching
- **Idea**: combines DL with classical optimization
- **Robust**: handles extreme wide baseline pairs.

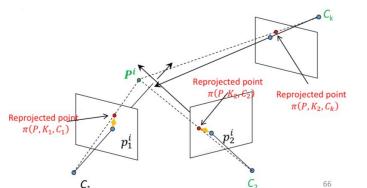


Sarlin, DeTone, Malisiewicz, Rabinovich, *SuperGlue: Learning Feature Matching with Graph Neural Networks*, International Conference on Computer Vision and Pattern Recognition (CVPR), 2020. [PDF](#) [Code](#).

2-VIEW BUNDLE ADJUSTMENT

- Non linear optimization of P^i and motion R, T
- Commonly used after LS (e.g 8-Point) estimation of R and T ($C_i = [R_i, T_i]$)
- Optimizes squared reprojection errors

$$P^i, C_2, \dots, C_n = \arg \min_{P^i, C_2, \dots, n} \sum_{k=1}^n \sum_{l=1}^N \| p_k^i - \pi(P^i, K_k, C_k) \|^2 \quad \begin{array}{l} \text{• when } n > \lambda \text{ is the } n\text{-view BA} \\ \text{• N.B : we assume } C_1 = [I, 0] \end{array}$$

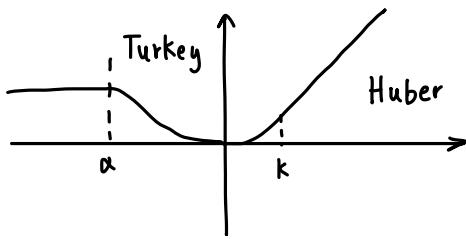


- First camera R, T can be arbitrary
- Occasionally residual terms are weighted
- Initial values P^i, R, T should be close to optimum
- Can be minimized using Levenberg - Marquardt (more robust than Gauss-Newton to local minima)
- Can be modified to also optimize intrinsics parameters.

HUBER AND TURKEY NORMS

To prevent that large reprojection errors can negatively impact the optimization, a more robust norm $\rho(\cdot)$ is used instead of L_2 .

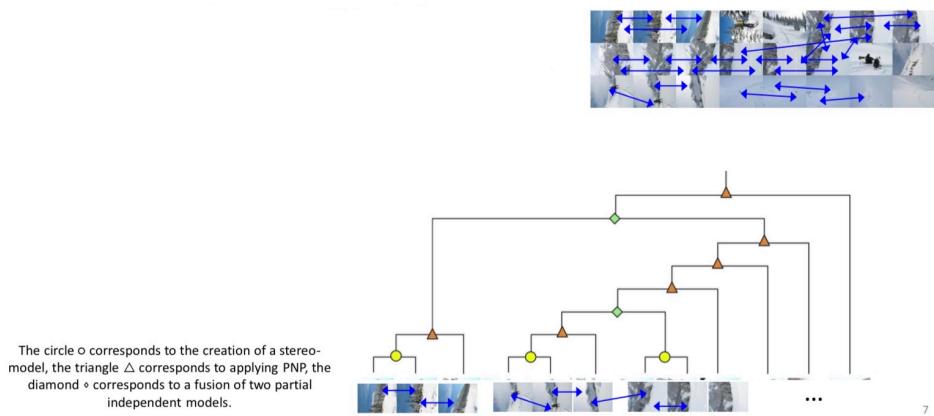
- Huber norm $\rho(x) = \begin{cases} x^2 & \text{if } |x| \leq k \\ k(2|x|-k) & \sim \end{cases}$
- Turkey norm $\rho(x) = \begin{cases} x^2 & \text{if } |x| \geq \alpha \\ \alpha^2(1 - (1 - (\frac{x}{\alpha})^2)^3) & \sim \end{cases}$



n-VIEW STRUCTURE FROM MOTION

HIERARCHICAL SFM

1. Extract and match features between nearby frames
2. Build cluster consisting of nearby frames
3. Extract topological tree (e.g. count number of SIFT features matches)
4. Start from the terminal nodes
 - compute 2-view SFM and build 3D model (=point cloud)
5. Iterate according to the tree structure:
 - 5.1 Merge new view by running 3-point RANSAC between 3D model and 3rd view
 - 5.2 Merge nearby models by running again 3 point RANSAC between one 3D model and one view of the other 3D model
 - 5.3 Bundle adjust

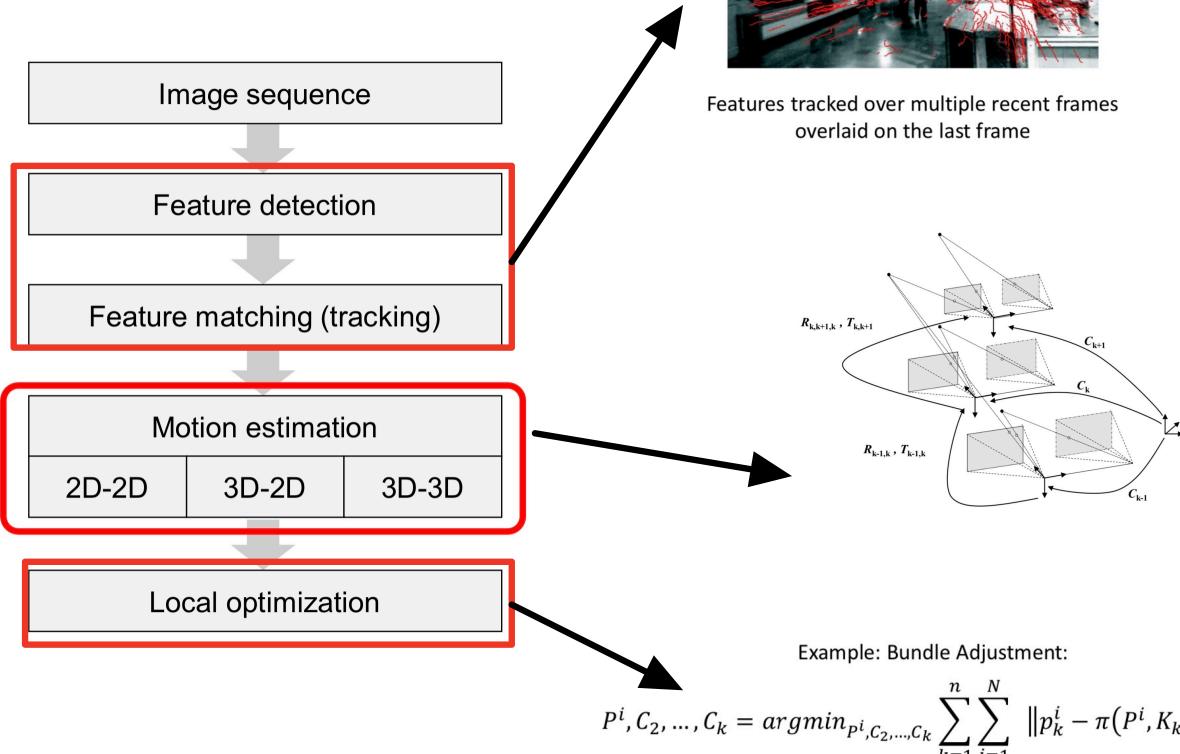


SEQUENTIAL STRUCTURE FROM MOTION (also called VISUAL ODOMETRY)

- Initialize structure from motion from 2 views (bootstrapping)
- For each additional view
 - Determine motion (pose) (localization)
 - Extend structure → can triangulate features during localization (mapping)
 - Refine structure and motion through bundle adjustment (optimization)

VO PIPELINE (RECAP)

VO computes the camera path incrementally (pose after pose)



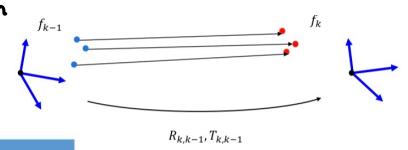
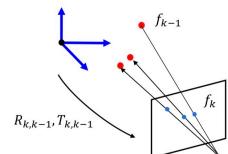
MOTION ESTIMATION

Or Pose-Graph Optimization (see later)

- Motion **From 2D to 2D** feature correspondences
 - both feature correspondences f_{k-1} and f_k are specified in image coordinates (2D)
 - the minimal case solution involves 5 features correspondences
 - popular algorithm 5 and 8-point algorithms
- Motion **From 3D to 2D** feature correspondences (i.e. Perspective from n points : PnP problem)
 - f_{k-1} is specified in 3D and f_k in 2D
 - DLT algorithm: minimal case is 6 points from 3D objects or 4 from planar objects
 - P3P algorithm: minimal case is 3 points (+1 for disambiguation)
 - EPNP algorithm: For more than 4 points
- Motion **From 3D to 3D** (also known as point cloud registration)
 - both f_{k-1} and f_k are specified in 3D. To do this it is necessary to triangulate the new features (need stereo camera)
 - the minimal case solution involves 3 non-collinear correspondences
 - popular algorithm consists of solving the system of equations with R and T unknown

$$\begin{bmatrix} X_{k-1}^i \\ Y_{k-1}^i \\ Z_{k-1}^i \end{bmatrix} = [R | T] \begin{bmatrix} X_k^i \\ Y_k^i \\ Z_k^i \end{bmatrix}$$

Type of correspondences	Monocular	Stereo
2D-2D	X	
3D-2D	X	X
3D-3D		X



CASE STUDY: MONOCULAR VO

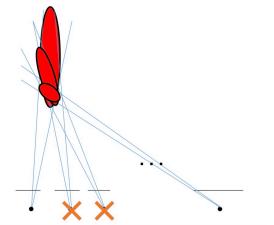
1. Bootstrapping (i.e. initialization)

- Initialize structure from motion from 2 views, e.g. 5-8 point RANSAC. The first extracted relative pose of the second camera allows to triangulate the detected features and obtain a first 3D point cloud (**structure**) of features.

- Refine structure and motion (Bundle Adjustment)

2. Keyframe selection (i.e skipping frames)

- When frames are taken at nearby positions compared to the scene distance, 3D points will exhibit large uncertainty
- One way to avoid this consists of skipping frames until the average uncertainty of the 3D points normalized by the average distance of the scene, falls below a certain threshold. The selected frames are called keyframes.
- Rule of thumb: add a keyframe when **keyframe distance / average-depth > threshold** ($\sim 10/20\%$)



3. Localization (i.e. pose estimation from a given point cloud)

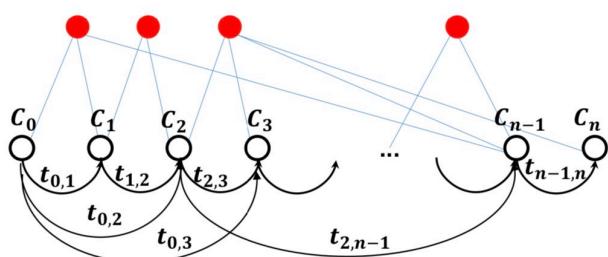
- Given a 3D point cloud (triangulated feature map) determine the pose of each additional view
- This is a case of 3D to 2D correspondences

4. Mapping (i.e. extend structure)

- Add the newly triangulate features to the point cloud
- It is not necessary to do this at every frame since neighboring frames have almost the same features
- More computation but can do loop-closure

LOCAL OPTIMIZATION

Refine structure and motion (often offline)



t_{ik} : transformation from camera pose C_k to camera pose C_i : $C_k \cdot t_{ki}$

- Pose graph optimization $\{C_1, \dots, C_n\} = \arg \min_{\{C_1, \dots, C_n\}} \sum_{i,j} \Sigma_i \Sigma_j \|C_i - C_j t_{ji}\|^2$

- Bundle adjustment : similar to pose-graph optimization but it also optimizes 3D points

$$P^i, C_1, \dots, C_n = \arg \min_{P^i, C_1, \dots, C_n} \sum_{k=1}^n \sum_{i=1}^m f(P^i - \pi(P^i, K_k, C_k)) \quad f(\cdot) : \text{Huber or Turkey norm}$$

BUNDLE ADJUSTEMENT VS POSE-GRAFH OPTIMIZATION

- BA is **more precise** than pose graph optimization it adds additional constraints (landmark constraints)
- BA is **more costly** $\Theta(qN + lm^2)$ with N being the number of points, m the number of camera poses and q and l the number of parameters for points and camera poses \rightarrow workarounds:
 - A **small window size** to make real time implementation possible
 - Not optimize over the 3D landmarks (motion only BA)

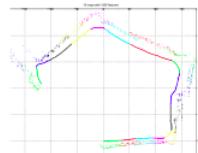
PLACE RECOGNITION

During VO two problems can occur:

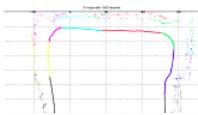
1. **Relocalization problem**: feature tracking can be lost (occlusion, low texture, quick illumination change)
or in case of monocular VO pure rotation followed by translation (multiple solutions exists?)
→ solution: relocalize camera pose and continue
2. **Loop closing problem**: when you go back to a previously mapped area
 - **Loop closure detection** : to avoid map duplication
 - **Loop correction** : to compensate the accumulated drift

VO vs. Visual SLAM (recap from Lecture 01)

- **Visual Odometry**
 - Focus on incremental estimation
 - **Guarantees local consistency** (i.e., estimated trajectory is locally correct, but not globally, i.e. from the start to the end)
- **Visual SLAM (Simultaneous Localization And Mapping)**
 - **SLAM = visual odometry + loop detection & closure**
 - **Guarantees global consistency** (the estimated trajectory is globally correct, i.e. from the start to the end)



Visual odometry



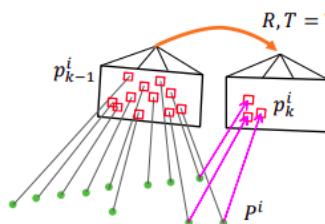
Visual SLAM

Indirect vs Direct Methods

Indirect methods

1. Extract & match features + 3-point RANSAC
2. Bundle Adjust by minimizing the **Reprojection Error**:

$$P^i, R, T = \arg \min_{P^i, R, T} \sum_{l=1}^N \rho(p_k^i - \pi(P^i, K, R, T))$$



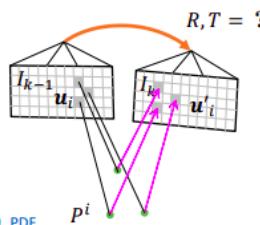
- ✓ Can cope with large frame-to-frame motions (large basin of convergence)

- ✗ Slow due to costly feature extraction, matching, and outlier removal (e.g., RANSAC)

Direct methods

1. No feature extraction & no RANSAC needed.
Instead, directly minimize **Photometric Error**:

$$P^i, R, T = \arg \min_{P^i, R, T} \sum_{l=1}^N \rho(I_{k-1}(p_{k-1}^i) - I_k(\pi(P^i, K, R, T)))$$



- ✓ All information in the image can be exploited (higher accuracy, higher robustness to motion blur and weak texture (i.e., weak gradients))

- ✓ Increasing the camera frame-rate reduces computational cost per frame (no RANSAC needed)

- ✗ Very sensitive to initial value → limited frame-to-frame motion (small basin of convergence)

Irani, Anandau, *All about direct methods*, Springer'99. [PDF](#)

Open Source Monocular VO and SLAM algorithms

- PTAM
- ORB-SLAM
- SVO
- LSD-SLAM
- DSO

Indirect methods: Minimize the feature reprojection error

Direct methods: Minimize the feature photometric error

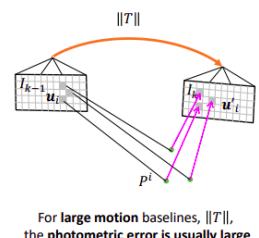
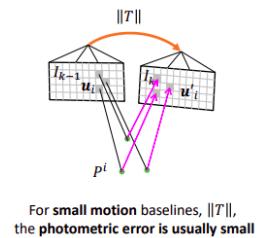
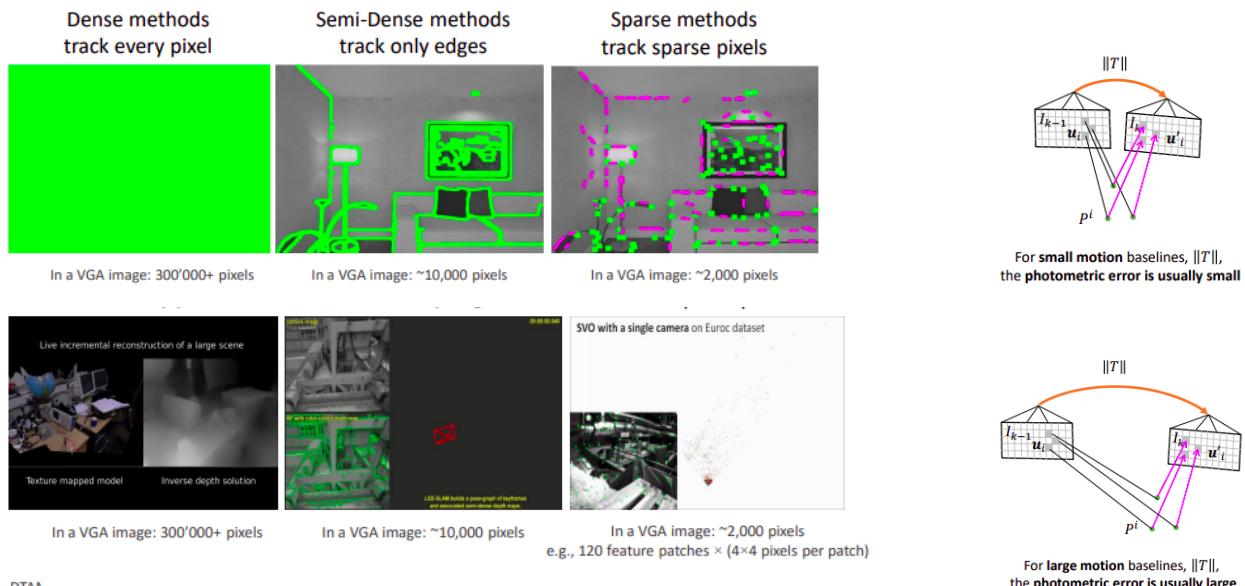
PTAM: Parallel Tracking and Mapping

- Monocular only
- **Feature based**
 - FAST corners + patch descriptors
 - Minimizes reprojection error
 - **Jointly optimizes poses & structure** (sliding window BA)
- First to propose **keyframe-based VO**
- First to propose **alternation of localization** (i.e., camera tracking) and **mapping**
- Tracking and mapping running in **two independent threads**: updated map is used by localization thread asynchronously, as soon it becomes available
- Includes:
 - Relocalization
 - No global optimization, only local
- **Real-time (30Hz)**, however global optimization is not done in real time but asynchronously every once in a while

ORB-SLAM

- Supports both **monocular and stereo cameras**
- **Feature based**
 - FAST corners + ORB descriptors
 - ORB: binary descriptor, very fast to compute and match (Hamming distance)
 - **Jointly optimizes poses & structure** (sliding window BA)
- **Same workflow as PTAM** (keyframe based, alternation of localization and mapping as independent threads)
- Includes:
 - Loop closing
 - Relocalization
 - Final optimization
- **Real-time (30Hz)**, however global optimization is not done in real time but asynchronously every once in a while

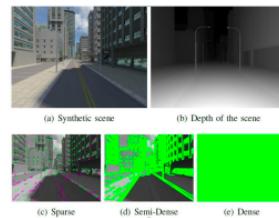
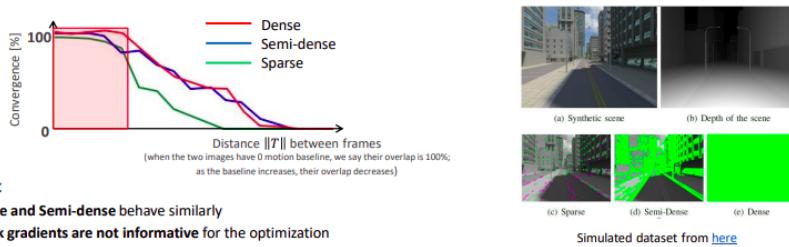
Direct Methods: Dense, Semi-dense, Sparse



Direct Methods: Dense, Semi-dense, Sparse

- What is the influence of the motion baseline on the convergence rate of direct methods?

We can use **photorealistic simulation** to answer this question by generating thousands of data



Simulated dataset from [here](#)

LSD-SLAM Large Scale Semi Dense - SLAM

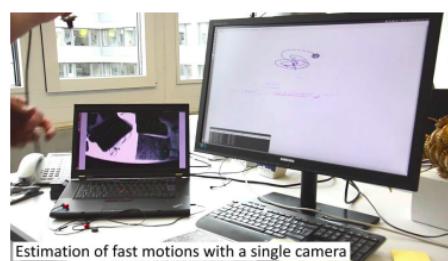
- Supports both **monocular** and **stereo** cameras
- Direct** (photometric error) + **Semi-Dense** formulation
 - 3D structure represented as **semi-dense** depth map
 - Minimizes **photometric error**
 - Separately optimizes poses & structure (sliding window)
- Same workflow as PTAM** (keyframe based, alternation of localization and mapping as independent threads)
- Includes:
 - Loop closing
 - Relocalization
 - Final optimization
- Real-time (30Hz)**, however global optimization is not done in real time but asynchronously every once in a while

DSO Direct Sparse Odometry

- Supports both **monocular** and **stereo** cameras
- Direct** (photometric error) + **Sparse** formulation
 - 3D structure represented as **sparse large gradients'** depth map
 - Minimizes **photometric error**
 - Jointly optimizes poses & structure** (sliding window)
 - Incorporates photometric correction to compensate exposure time change ($\Delta t_{k-1}, \Delta t_k$)
- $p^l, R, K = \arg \min_{p^l, R, K} \sum_{l=1}^N p \left(I_{k-1}(p^l_{k-1}) - \frac{\Delta t_{k-1}}{\Delta t_k} I_k \left(\pi(p^l, K, R, T) \right) \right)$
- Same workflow as PTAM** (keyframe based, alternation of localization and mapping as independent threads)
- Real-time (30Hz)**, however global optimization is not done in real time but asynchronously every once in a while

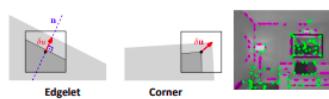
SVO Semi Direct Visual Odometry

- Supports both **monocular**, **stereo**, **multi-camera systems** as well as omnidirectional models (fisheye and catadioptric)
- Combines **indirect** + **direct** methods
 - Direct methods for frame-to-frame motion estimation
 - Indirect methods for frame-to-keyframe pose refinement
- Mapping**
 - Probabilistic depth estimation (heavy-tail Gaussian distribution)
- Includes:
 - Loop closing,
 - Relocalization,
 - Final optimization
- Same workflow as PTAM** (keyframe based, alternation of localization and mapping as independent threads)
- Faster than real-time: up to 400 fps** on i7 laptops and **100 fps** on smartphone PCs (Odroid (ARM)) or NVIDIA Jetson



	Mean	CPU@20 fps
SVO Mono	2.53	55 ± 10%
ORB Mono SLAM (No loop closure)	29.81	187 ± 32%
LSD Mono SLAM (No loop closure)	23.23	236 ± 37%
DSO	20.12	181 ± 27%

↑ Processing time in milliseconds ↑ CPU load (100% = 1 core)

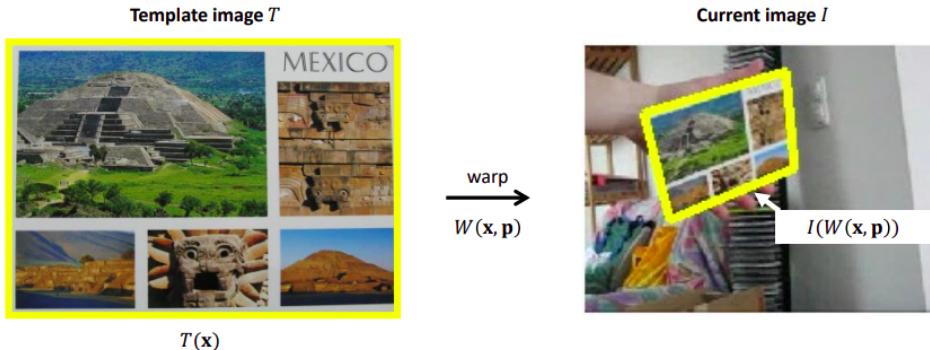


Template tracking

Goal: **follow** a template image in a video sequence

Problem Formulation

Goal: estimate the transformation W (warp) between a template T and the current image



Two possible approaches

- Indirect methods

- ✓ Can cope with large frame-to-frame motions (large basin of convergence)
- ✗ Slow due to costly feature extraction, matching, and outlier removal (e.g., RANSAC)

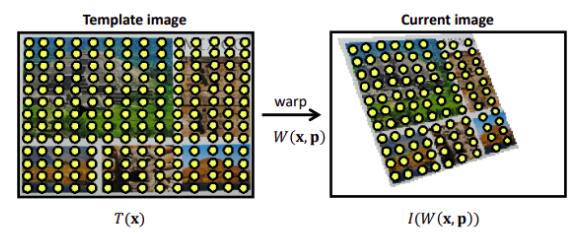
- Direct methods

- ✓ All information in the image can be exploited (higher accuracy, higher robustness to motion blur and weak texture (i.e., weak gradients))
- ✓ Increasing the camera frame-rate reduces computational cost per frame (no RANSAC needed)
- ✗ Very sensitive to initial value → limited frame-to-frame motion (small basin of convergence)

Direct methods work by directly processing pixel intensities

Goal: estimate the parameters \mathbf{p} of the transformation $W(\mathbf{x}, \mathbf{p})$ that minimize the Sum of Squared Differences:

$$\mathbf{p} = \operatorname{argmin}_{\mathbf{p}} \sum_{\mathbf{x} \in T} [I(W(\mathbf{x}, \mathbf{p})) - T(\mathbf{x})]^2$$



* Every yellow dot in this image denotes a pixel

Indirect methods work by detecting and matching features

1. Detect and match features that are invariant to scale, rotation, view point changes (e.g., SIFT)
2. Geometric verification (RANSAC) (e.g., 4-point RANSAC for planar objects, or 5 or 8-point RANSAC for 3D objects)
3. Refine estimate by minimizing the sum of squared reprojection errors between the observed feature f^l in the current image and the warped corresponding feature $W(x^l, \mathbf{p})$ from the template

$$\mathbf{p} = \operatorname{argmin}_{\mathbf{p}} \sum_{l=1}^N [W(x^l, \mathbf{p}) - f^l]^2$$

- **Pros:** can cope with large frame-to-frame motion and strong illumination changes
- **Cons:** computationally expensive



Assumptions

- **Brightness constancy**

The intensity of the pixels to track does not change much over consecutive frames → **It does not cope with strong illumination changes**



- **Temporal consistency**

Small frame-to-frame motion (1-2 pixels).

→ **It does not cope with large frame-to-frame motion.** However, this can be addressed using coarse-to-fine multi-scale implementations (see later)

- **Spatial coherency**

All pixels in the template undergo the same transformation (i.e., they all lay on the same 3D surface)

→ **No errors in the template image boundaries:** only the object to track appears in the template image

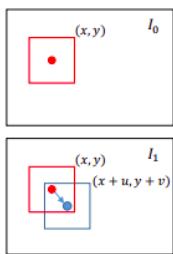
→ **No occlusion:** the entire template is visible in the input image



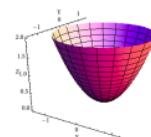
9

KLT tracking applied to pure translation

Consider the reference patch centered at (x, y) in image I_0 and the shifted patch centered at $(x + u, y + v)$ in image I_1 . The patch has size Ω . We want to find the motion vector (u, v) that minimizes the Sum of Squared Differences (SSD):



$$\begin{aligned} SSD(u, v) &= \sum_{x, y \in \Omega} (I_0(x, y) - I_1(x + u, y + v))^2 \\ &\quad I_1(x + u, y + v) \cong I_1(x, y) + I_x u + I_y v \\ \Rightarrow SSD(u, v) &\cong \sum_{x, y \in \Omega} (I_0(x, y) - I_1(x, y) - I_x u - I_y v)^2 \\ \Rightarrow SSD(u, v) &\cong \sum_{x, y \in \Omega} (\Delta I - I_x u - I_y v)^2 \end{aligned}$$



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

To minimize it, we differentiate it with respect to (u, v) and equate it to zero:

$$\frac{\partial SSD}{\partial u} = 0, \frac{\partial SSD}{\partial v} = 0$$

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

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

$$\sum I_x (\Delta I - I_x u - I_y v) = 0$$

$$\sum I_y (\Delta I - I_x u - I_y v) = 0$$

- Linear system of two equations in two unknowns

Haven't we seen this matrix already?

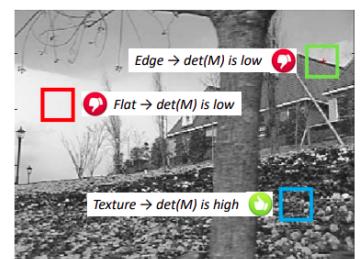
- We can write them in matrix form:

This is the M matrix of the Harris detector (Lecture 05)

$$\begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \sum I_x \Delta I \\ \sum I_y \Delta I \end{bmatrix} \Rightarrow \begin{bmatrix} u \\ v \end{bmatrix} = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix}^{-1} \begin{bmatrix} \sum I_x \Delta I \\ \sum I_y \Delta I \end{bmatrix}$$

$$M = \begin{bmatrix} \sum I_x I_x & \sum I_x I_y \\ \sum I_x I_y & \sum I_y I_y \end{bmatrix} = R^{-1} \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} R$$

For M to be invertible, $\det(M)$ should be non zero, which means that its eigenvalues should be large (i.e., not a flat region, not an edge) → in practice, it should be a corner or more generally contain texture!



15

Application to Optical Flow

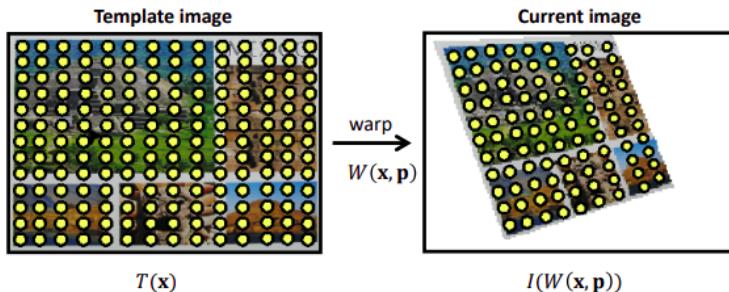
What if you track every single pixel in the image?



KLT applied to generic warps

Goal: estimate the parameters \mathbf{p} of the transformation $W(\mathbf{x}, \mathbf{p})$ that minimize the SSD:

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



* Every yellow dot in this image denotes a pixel

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

- Assume that an initial estimate of \mathbf{p} is known. Then, we want to find the increment $\Delta \mathbf{p}$ that minimizes

$$SSD = \sum_{\mathbf{x} \in T} [I(W(\mathbf{x}, \mathbf{p} + \Delta \mathbf{p})) - T(\mathbf{x})]^2$$

- First-order Taylor approximation of $I(W(\mathbf{x}, \mathbf{p} + \Delta \mathbf{p}))$ yields to:

$$I(W(\mathbf{x}, \mathbf{p} + \Delta \mathbf{p})) \cong I(W(\mathbf{x}, \mathbf{p})) + \nabla I \frac{\partial W}{\partial \mathbf{p}} \Delta \mathbf{p}$$

$\nabla I = [I_x, I_y]$ = Image gradient evaluated at $W(\mathbf{x}, \mathbf{p})$

Jacobian of the warp $W(\mathbf{x}, \mathbf{p})$

- By replacing $I(W(\mathbf{x}, \mathbf{p} + \Delta \mathbf{p}))$ with its 1st order approximation, we get

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

- How do we minimize it?
- We differentiate SSD with respect to $\Delta \mathbf{p}$ and we equate it to zero, i.e., $\frac{\partial SSD}{\partial \Delta \mathbf{p}} = 0$

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

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

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

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

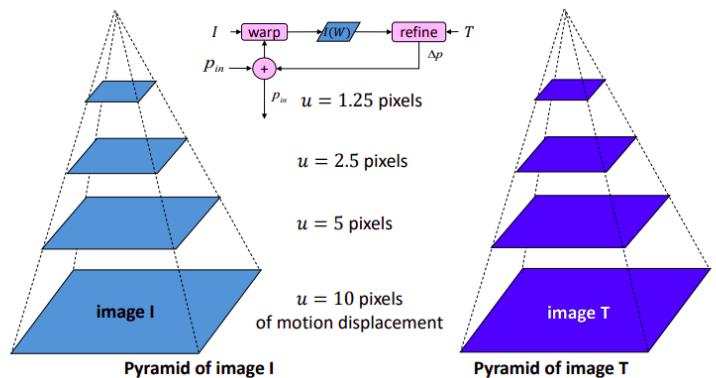
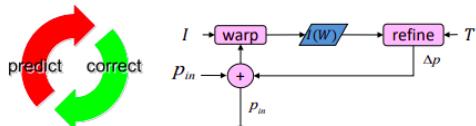
Second moment matrix (Hessian) of the warped image

What does H look like when the warp is a pure translation?

Coarse-to-fine estimation

KLT algorithm is iterated until convergence by following a **predict-correct cycle**

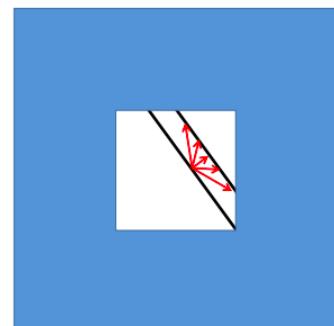
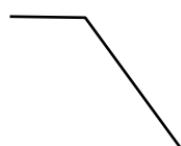
1. A prediction $I(W(x, p))$ of the warped image is computed from an initial estimate of p
2. The correction parameter Δp is then computed as a function of the error $T(x) - I(W(x, p))$ between the prediction and the template. The larger this error, the larger the correction applied
3. Steps 1 & 2 are iterated until the error is small and the output parameters are used as input for the next frame



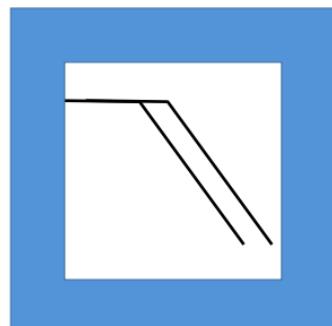
Aperture Problem

- Consider the motion of the following corner

- Now, look at the local brightness changes **through a small aperture**
- We **cannot always determine** the motion direction → **Infinite motion solutions** may exist!
- Solution?

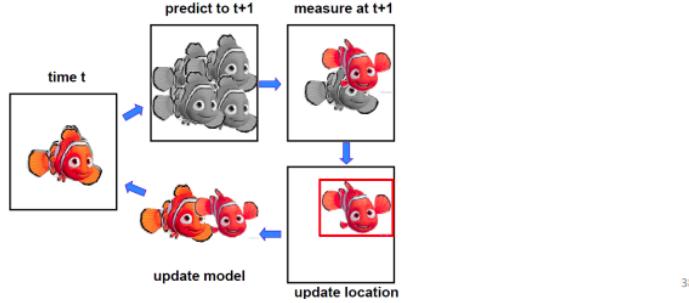


- Now, look at the local brightness changes **through a small aperture**
- We **cannot always determine** the motion direction → **Infinite motion solutions** may exist!
- Solution?
 - Increase aperture size!



Generalization of KLT

- The same concept (predict/correct) can be applied to tracking of 3D object (**in this case, what is the transformation to estimate? What is the template?**)
- In order to deal with wrong prediction, it can be implemented in a **Particle-Filter** fashion (using multiple hypotheses that need to be validated)



38

Displacement-model Jacobians ∇W_p

$$p = (a_1, a_2, \dots, a_n)$$

- Translation: $W(\mathbf{x}, \mathbf{p}) = \begin{bmatrix} x + a_1 \\ y + a_2 \end{bmatrix}$ $\frac{\partial W}{\partial \mathbf{p}} = \begin{bmatrix} \frac{\partial W_1}{\partial a_1} & \frac{\partial W_1}{\partial a_2} \\ \frac{\partial W_2}{\partial a_1} & \frac{\partial W_2}{\partial a_2} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$
- Euclidean: $W(\mathbf{x}, \mathbf{p}) = \begin{bmatrix} x\cos(a_3) - y\sin(a_3) + a_1 \\ x\sin(a_3) + y\cos(a_3) + a_2 \end{bmatrix}$ $\frac{\partial W}{\partial \mathbf{p}} = \begin{bmatrix} 1 & 0 & -x\sin(a_3) - y\cos(a_3) \\ 0 & 1 & x\cos(a_3) - y\sin(a_3) \end{bmatrix}$
- Affine: $W(\mathbf{x}, \mathbf{p}) = \begin{bmatrix} a_1x + a_3y + a_5 \\ a_2x + a_4y + a_6 \end{bmatrix}$ $\frac{\partial W}{\partial \mathbf{p}} = \begin{bmatrix} x & 0 & y & 0 & 1 & 0 \\ 0 & x & 0 & y & 0 & 1 \end{bmatrix}$

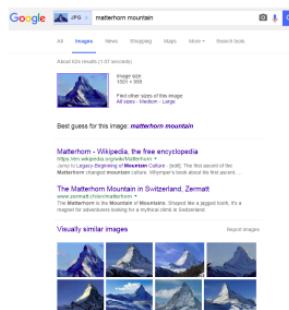
Place Recognition

- **Robotics:**

- Has the robot been to this place before?
- Which images were taken around the same location?

- **Image retrieval:**

- Have I seen this image before?
- Which images in my database look similar to it?
E.g., Google Reverse Image Search



Visual Place Recognition

- **Goal:** query an image in a database of N images
- **Complexity:** $O(NF^2)$ feature comparisons (assumes each image has F features)
 - Example:
 - assume 1,000 SIFT features per image $\rightarrow F = 1,000$
 - assume $N = 100,000,000$
 - $\rightarrow NF^2 = 100,000,000,000,000$ feature comparisons!
 - If we assume 10 microseconds per feature comparison \rightarrow 1 image query would take **32 years!**

Solution: Use an index file! Complexity reduces to $O(F)$

Text Retrieval

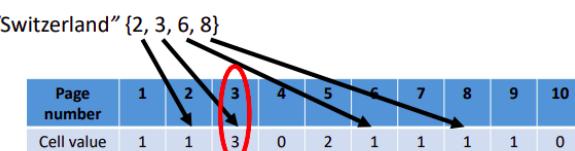
- Image retrieval takes inspiration from **text retrieval**
- For text documents, an efficient way to find all **pages** in which a **word** occurs is to use an **index file**
- To **retrieve a given text query**, it is then sufficient to use a **voting scheme**
- Suppose that we have a **document with 10 pages** and we want to determine on which page this sequence of words appears:

"Zurich is a city of Switzerland"

- "Zurich is a city of Switzerland" is the **query text**
- Suppose that this is our index file:

Word	Page numbers
Zurich	3, 5, 7
City	1, 3, 5, 9
Switzerland	2, 3, 6, 8

- The solution is to use a **voting array** that has as many cells as the number of pages in the document
- We first set all the cell values to 0
- Then, we add 1 to each cell corresponding to the page numbers where the words of the query text appear according to the index file



Bag of Words

Using the analogy from text retrieval, we need to:

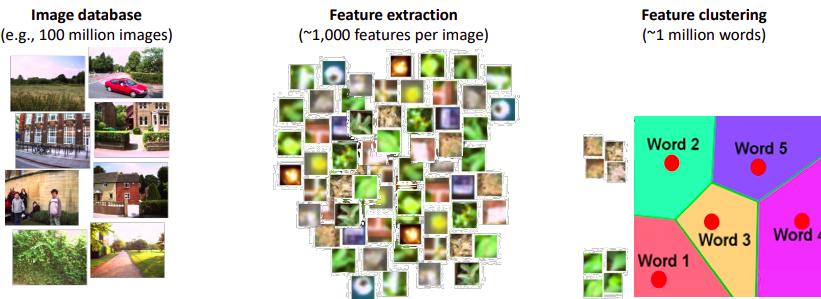
- define what a “*visual word*” is and
- define a “*vocabulary*” of visual words

This approach is known as “Bag of Words” (BOW)

- **Can a SIFT descriptor be used as a visual word? And a BRISK descriptor?**
 - SIFT $\rightarrow 128 \times 4 \text{ bytes float} = 512 \text{ bytes} = 4096 \text{ bits} = 2^{4096}$ possible SIFT descriptors!
 - BRISK-128 $\rightarrow 128 \text{ bits} = 2^{128}$ possible BRISK descriptors!
 - **Usually, 1 million visual words is enough**
 - Idea: cluster SIFT descriptors into visual words

How to extract Visual Words from descriptors

- **Collect a large enough dataset** that is representative of all possible images that are relevant to your application (e.g., for automotive place recognition, you may want to collect million of street images sampled around the world)
- Extract features and descriptors from each image and map them into the **descriptor space** (e.g., for SIFT, 128 dimensional descriptor space)
- **Cluster the descriptor space into K clusters**
- **The centroid of each cluster is a visual word.**
 - This is computed by taking the arithmetic average of all the descriptors within the same cluster:
 - e.g., for SIFT, each cluster contains SIFT features that are very similar to each other;
 - the visual word then is the average all the SIFT descriptors in that cluster



- Extracting SIFT features from a VGA images takes ~ 20 ms on an i7 CPU
- This means that the **extraction of features from 100 million images would take ~ 23 days** without accounting for the time needed for clustering all these features
- However, notice that **this is ok since this process only needs to be done once**, when the database has been created.
- If the database grows as new images are collected, **new features can be extracted and the visual words can be updated accordingly**. This update process **does not need to run in real time**

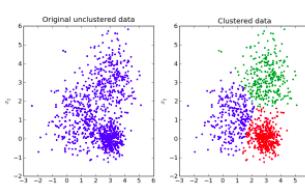
How do we cluster the descriptor space?

- **k-means clustering** is an algorithm to partition n data points into k clusters in which each data point x belongs to the cluster S_i with center m_i
- It minimizes the **squared Euclidean distance** between points x and their nearest cluster centers m_i

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

Algorithm:

- Randomly initialize k cluster centers
- Iterate until 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



Building the Image Vocabulary

- The **Image Vocabulary** is a data structure that lists all extracted visual words
- Each visual word is assigned a unique identifier (an **integer number**)
- Each visual word** in the image vocabulary **points to a list of images** (from the entire image database) in which that word appears
- If the database grows, the vocabulary is updated accordingly

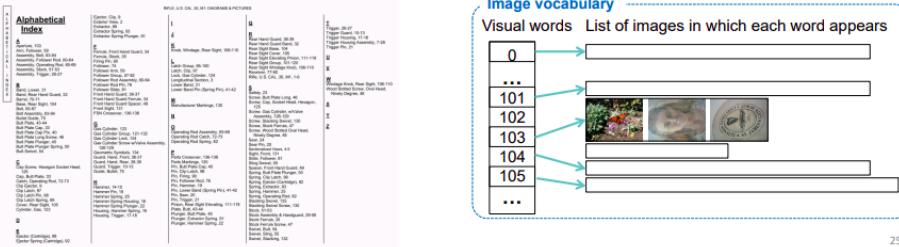
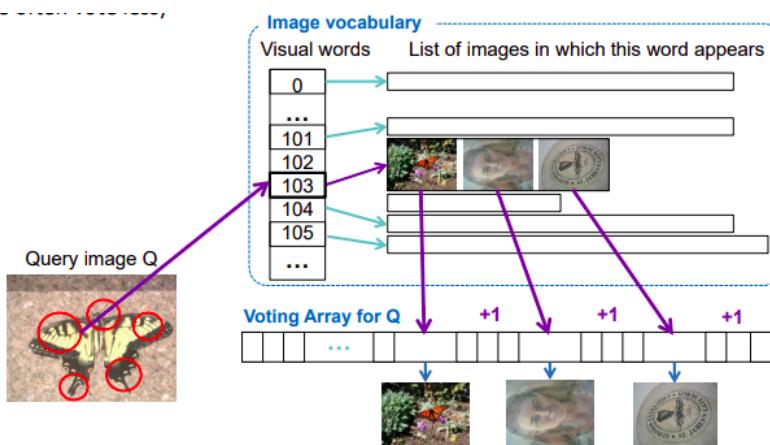


Image Retrieval

- To query an image in the database, we must first extract features from the query image (this takes about 20 ms for ~1,000 SIFT features)
- Then, we initialize the **voting array** to 0 (the voting array has as many cells as the number of images in the database).
- Then, we **look-up** each feature in the image vocabulary (basically, we look for the **closest visual word** in the vocabulary)
- Finally, **each visual word votes for multiple images** as we saw for the case of text retrieval; however, the **voting is not uniform but is weighted by the inverse of the frequency of the word** (i.e., words that repeat more often vote less)



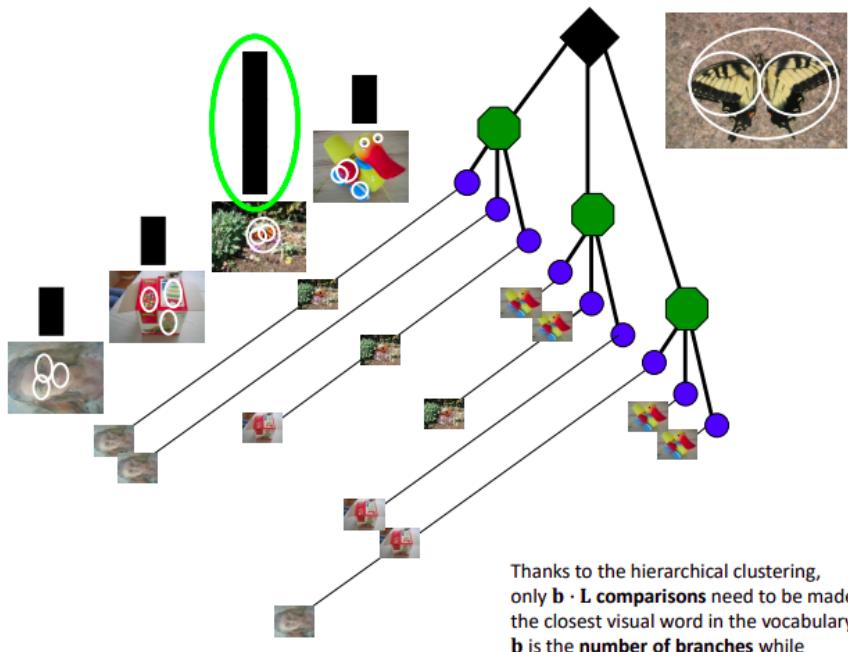
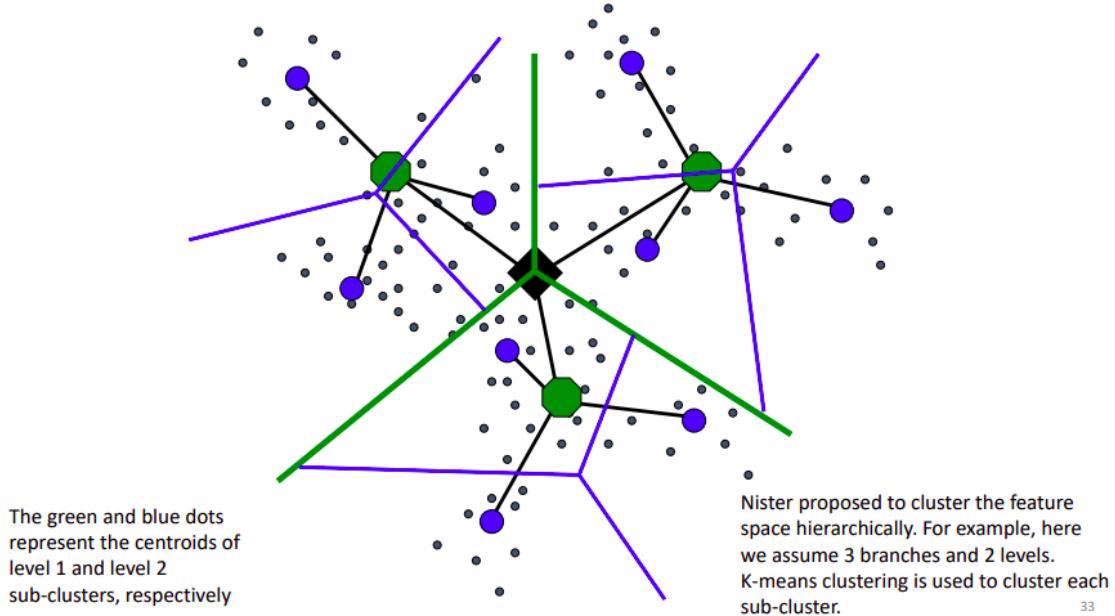
Issues

Every feature in the query image has to be compared against all visual words in the vocabulary:

- Example:
 - assume our query image has 1,000 features;
 - assume 1 million visual words → number of feature comparisons would be equal to **1 billion!**
 - If we assume 10 microseconds per feature comparison, then querying one image would take **~3 hours!**
- How can we make the comparison cheaper?
 - Idea: use **hierarchical clustering**

Hierarchical Clustering

- David Nister proposed to cluster the feature space in a **coarse-to-fine manner** so that visual words could be represented as the **terminal vertices of a search tree**.
- As we will see, this significantly reduces the feature-to-word association, bringing image retrieval within the reach of resource-constrained mobile devices!



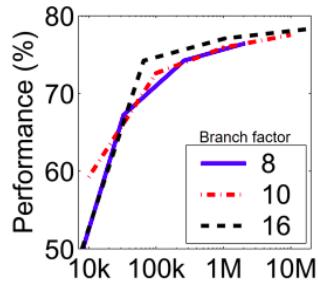
Example

Querying an image in a database of **100 million images**

- assume a query image with $F = 1,000$ features
- assume a tree structure with $b = 10$ branches per level and $L = 6$ levels (i.e., $b^L = 1,000,000$ visual words)
- Then, the **number of feature comparisons** = $F \cdot b \cdot L = 1,000 \cdot 10 \cdot 6 = 60,000$ instead of $F \cdot b^L = 1$ billion comparisons (which was the case without hierarchical clustering)
- If we assume 10 microseconds per feature comparison
 → 1 image query would take **0.6 seconds!**

How many visual words, branches, and levels?

- More words is better, but **1 million words** are used in practice
- Also, **6 branches** and **10 depth levels** are practically used (e.g., ORB-SLAM)



Geometric Verification

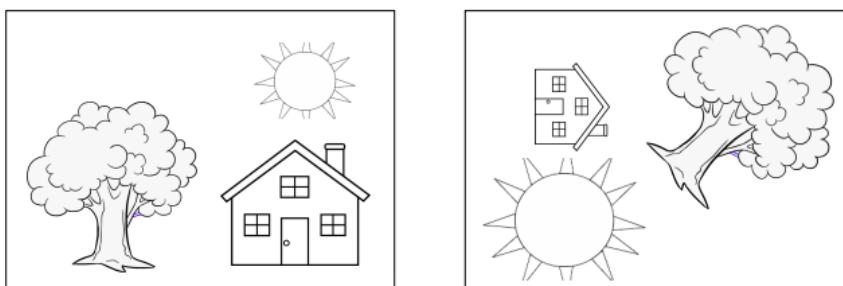
Bag of Words returns a list of images with score above a given threshold. **How do we pick the image that was captured closest to the query image?**



- Observe that the visual vocabulary discards the spatial relationships between features
- Solution: Test each returned image for geometric consistency against the query image (e.g. using 5- or 8-point RANSAC) and pick the image with **the smallest reprojection error and largest number of inliers**

Curiosity

Imagine to query two images with the same features shuffled around. Will the scores returned by Bag of Words be different?



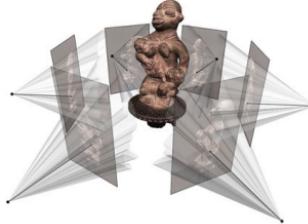
Open Challenges

When does place recognition fail and how would you address them?

Dense Reconstruction (or Multi-view stereo)

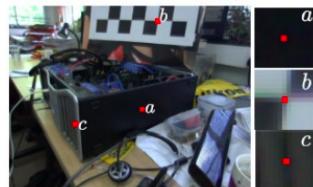
Problem definition:

- **Input:** calibrated images from several viewpoints (i.e., K, R, T are known for each camera, e.g., from SFM)
- **Output:** 3D object **dense** reconstruction (ideally of every pixel)



Challenges

- Dense reconstruction requires establishing dense correspondences
- But not all pixels can be matched reliably:
 - flat regions,
 - edges,
 - viewpoint and illumination changes,
 - occlusions



Idea: Take advantage of many small-baseline views where high-quality matching

Dense reconstruction workflow

Step 1: Local methods

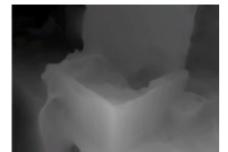
- Estimate depth independently for each pixel

How do we compute correspondences for every pixel?



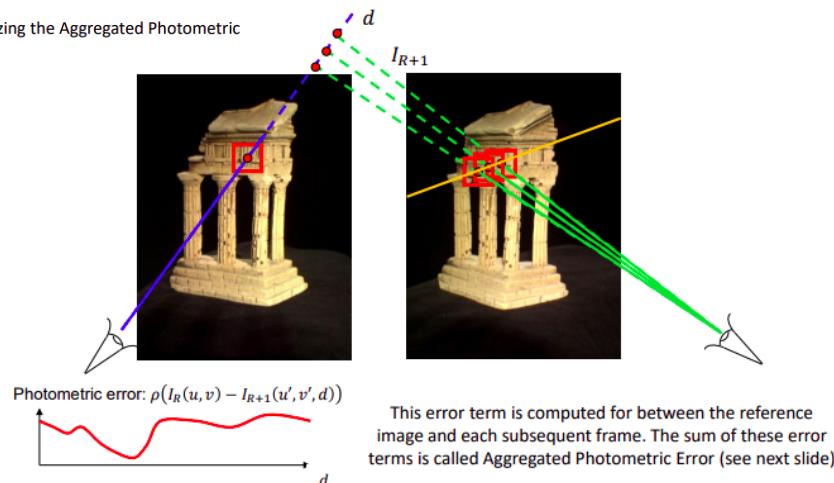
Step 2: Global methods

- Refine the *depth map as a whole* by enforcing smoothness. This process is called *regularization*



Solution: Aggregated Photometric Error

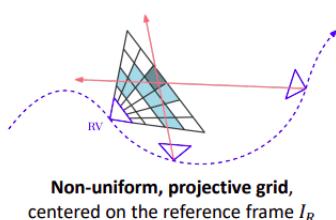
Set the first image as reference and estimate depth at each pixel by minimizing the Aggregated Photometric Error in all subsequent frames



Disparity Space Image (DSI)

- Image resolution: 240×180 pixels
- Number of disparity (depth) levels: 100
- DSI:
 - size: $240 \times 180 \times 100$ voxels; each voxel contains the Aggregated Photometric Error $C(u,v,d)$ (see next slide)
 - white = high Aggregated Photometric Error
 - blue = low Aggregated Photometric Error

DSI (dark means high)



Disparity Space Image (DSI)

- For a given image point (u, v) and for discrete depth hypotheses d , the **Aggregated 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 has value:

$$C(u, v, d) = \sum_{k=R+1}^{R+n-1} \rho(I_R(u, v) - I_k(u', v', d))$$

where n is the number of images considered and $I_k(u', v', d)$ 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 depth hypothesis d ; thus, formally:

$$I_k(u', v', d) = I_k\left(\pi(T_{k,R}(\pi^{-1}(u, v) \cdot d))\right)$$

where $T_{k,R}$ is the relative pose between frames R and K

- $\rho(\cdot)$ is the photometric error (SSD) (e.g. L_1 , L_2 , Tukey, or Huber norm)

Depth estimation

The solution to the depth estimation problem is to find a **function $d(u, v)$** (called **depth map**) in the DSI that minimizes the **aggregated photometric error**:

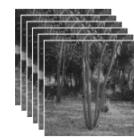
$$\text{depth map} = d(u, v) = \arg \min_d \sum_{(u,v)} C(u, v, d(u, v))$$

Influence of patch size

How does this relate to stereo vision?



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



$W = 3$

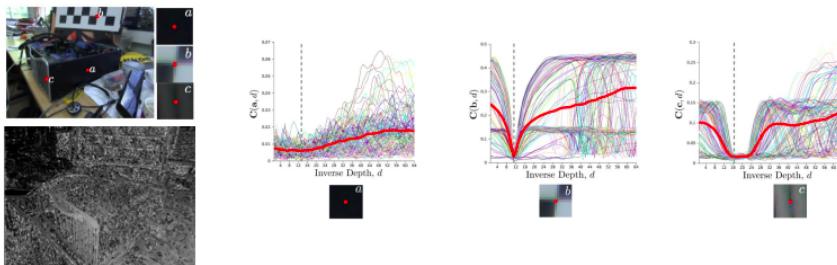
$W = 20$

Can we use a patch size of 1×1 pixels?

Depth map: each pixel intensity encodes the depth of that pixel. Dark = far, bright = close.

Influence of the patch appearance

- Aggregated photometric error for **flat regions** (a) and **edges parallel to epipolar lines** (c) **show flat valleys (plus noise)**
- For **textured areas** (e.g., corners (b) or blobs), the aggregated photometric error has one **distinctive minimum**
- Repetitive texture** shows **multiple minima**



Regularization

To penalize wrong reconstruction due to image noise and ambiguous texture, we add a smoothing term (called regularization term) to the optimization:

$$d(u, v) = \arg \min_d \sum_{(u,v)} C(u, v, d(u, v)) \quad (\text{local methods})$$

subject to

Piecewise smooth (global methods)

Regularization

- Formulated in terms of energy minimization
- The objective is to find a **surface** $d(u, v)$ that minimizes a global energy functional:

$$E(d) = \underbrace{E_d(d)}_{\text{Data term}} + \lambda \underbrace{E_s(d)}_{\text{Regularization term (i.e., smoothing)}}$$

where λ controls the tradeoff between data and regularization. **What happens as λ increases?**

Data term: $E_d(d) = \sum_{(u,v)} C(u, v, d(u, v))$

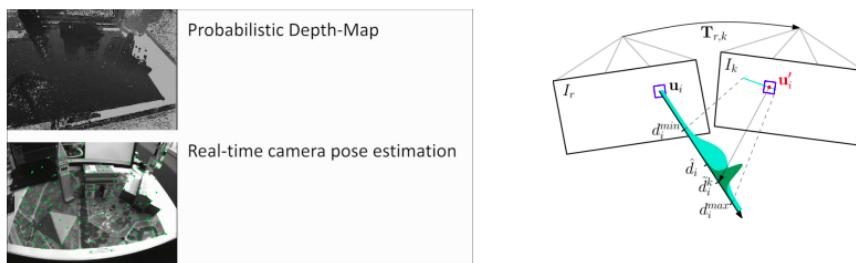
Regularization term: $E_s(d) = \sum_{(u,v)} \left(\frac{\partial d}{\partial u} \right)^2 + \left(\frac{\partial d}{\partial v} \right)^2$

How can we deal with depth discontinuities between separate objects?

The regularization term $E_s(d)$ basically fills the holes: it smooths the depth map by softening discontinuities

Probabilistic Monocular Dense Reconstruction

- Estimate depth uncertainty
- **Regularize only 3D points with low depth uncertainty** (does not fill holes if present, which is good for robotic applications)

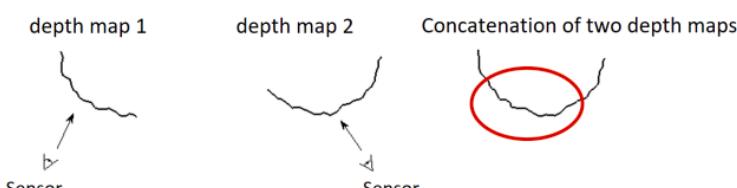


Pizzoli, Forster, Scaramuzza, REMODE: Probabilistic, Monocular Dense Reconstruction in Real Time, IEEE International Conference on Robotics and Automation (ICRA), 2014. [PDF](#) [Video](#) [Code](#).

23

What about large motions?

- When the **distance between the current frame and the reference frame** gets too large (e.g., $> 10\%$ of the average scene depth), then the Aggregated Photometric Error starts to diverge due to the **large view point changes**
- Solution: **create a new reference frame** (keyframe) and start a new depth map computation



GPU for 3D Dense Reconstruction

- **Image processing**
 - Filtering & Feature extraction (i.e., **convolutions**)
 - Warping (e.g., **epipolar rectification, homography**)
- **Multiple-view geometry**
 - Search for **dense correspondences**
 - Pixel-wise operations (SAD, SSD, NCC)
 - **Matrix and vector operations** (epipolar geometry)
 - Aggregated Photometric Error for multi-view stereo
- **Global optimization**
 - Variational methods (i.e., **regularization** (smoothing))
 - **Parallel, in-place** operations for gradient / divergence computation
- **Deep learning**

optimized for **serial processing** (i.e., only a few instructions can be executed during the same clock cycle)
forms calculations in **parallel** on thousands of cores (i.e., thousands of instructions can be executed in the same clock cycle)



ALU: Arithmetic Logic Unit

Supervised Learning

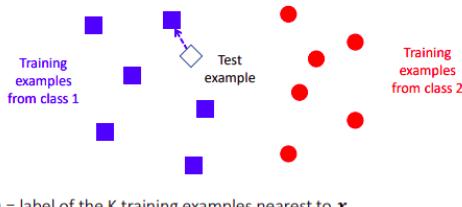
Find function $f(\mathbf{x}, \theta)$ that imitates a ground truth signal

Machine Learning Keywords

1. **Loss:** Quantify how good θ are
2. **Optimization:** The process of finding θ that minimize the loss
3. **Function:** Problem modelling → Deep networks are highly non-linear $f(\mathbf{x}, \theta)$

Classifiers: K-Nearest neighbor

Features are represented in the descriptor space

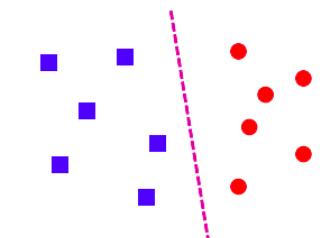


How fast is training? How fast is testing?

- $O(1), O(n)$

What is a good distance metric ? What K should be used? \circlearrowright

Find a *linear function* to separate the classes:



What is θ ? What is the dimensionality of images?

Multi layer perceptron (non linear classifier)

- Composition of linear layer with non linear activation functions
- **Back-propagation:** changes weights using gradients of loss with respect to weights (stochastic gradient descent)
- Too many parameters → possible **overfitting**

Convolutional Neural Networks

- Use masks (shared weights) for spatial consistency
- Going deep improve performances → extract progressively more sophisticated information

Supervised Learning

- We have access to both input data or images and ground truth labels
- Network trained with supervision usually perform the best
- Data generation and hand labeling is hard
- Usual tasks are image segmentation and image captioning

Unsupervised Learning

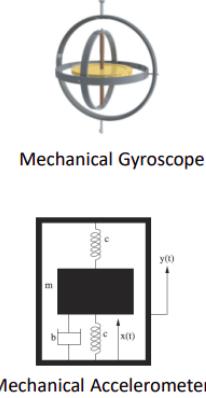
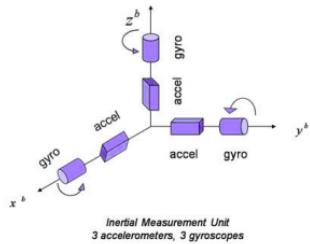
- only access to input data or images → can use larger datasets (no manual labeling needed)
- examples are monocular depth estimation, structure from motion

	Supervised	Unsupervised
Performance	Usually better for the same dataset size.	Usually worse, but can outperform supervised methods due to larger data availability.
Data availability	Low, due to manual labelling.	High, no labelling required.
Training	Simple, ground truth gives a strong supervision signal.	Sometimes difficult, loss functions have to be engineered to get good results.
Generalizability	Good, although sometimes the network learns to blindly copy the labels provided, leading to poor generalizability.	Better, since unsupervised losses often encode the task in a more fundamental way.

What is an IMU?

- **Inertial Measurement Unit**

- Gyroscope: Angular velocity
- Accelerometer: Linear Accelerations



- **Different categories**

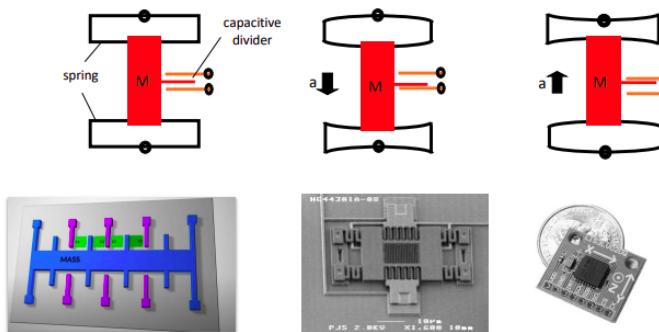
- Mechanical (\$100,000-1M)
- Optical (\$20,000-100k)
- MEMS (from 1\$ (phones) to 1,000\$ (higher cost because they have a microchip running a Kalman filter))

- **For small mobile robots & drones:**
MEMS IMU are mostly used

- Cheap
- Power efficient
- Light weight and solid state

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.



Why do we need an IMU?

- Monocular vision is **scale ambiguous**

- Pure vision is **not robust enough**
 - Low texture
 - High Dynamic Range (HDR)
 - High speed motion



Robustness is a critical issue: Tesla accident

"The autopilot sensors on the Model S failed to distinguish a white tractor-trailer crossing the highway against a bright sky." [The Guardian]



Why is an IMU alone not enough?

- Pure IMU integration will lead to **large drift** (especially in cheap IMUs)
- Example: **1D scenario**. Double integration of acceleration returns the position:

$$x(t) = x_0 + v_0(t - t_0) + \int \int_{t_0}^t a(\tau) d\tau^2$$

- If there is a **constant bias in the acceleration**, the error of position will be **proportional to t^2**
- Similarly for the orientation: if there is a **bias in angular velocity**, the error is **proportional to the time t**

GRADE/TIME	1 s	10 s	60 s	10 min	1 hr
Consumer	6 cm	6.5 m	400 m	200 km	39,000 km
Industrial	6 mm	0.7 m	40 m	20 km	3,900 km
Tactical	1 mm	8 cm	5 m	2 km	400 km
Navigation	<1 mm	1 mm	50 cm	100 m	10 km

Automotive, smartphones, and drones accelerometers

Table from Vectornav, one of the best IMU companies. Errors were computed assuming the device at rest:
<https://www.vectornav.com/resources/inertial-navigation-primer/specifications-and-error-budgets/specs-nserrorbudget>

Why visual inertial fusion?

IMU and vision are complementary

Cameras

- **Exteroceptive sensor:** measures light energy from the environment
- ✗ **Sensitive to motion blur, HDR, texture**
- ✓ **Drift is bounded when observing the same scene**
- ✓ **Precise in slow motion**
- ✗ **Limited output rate (~100 Hz)**
- ✗ **Scale ambiguity in monocular setup**

IMU

- **Proprioceptive sensor:** measures values internal to the system
- ✓ **Insensitive to motion blur, HDR, texture**
- ✗ **Drift grows unbounded regardless of the environment**
- ✗ **Less precise in slow motion (low signal-to-noise ratio)**
- ✓ **High output rate (1,000-10,000 Hz)**
- ✓ **Can be used as a prior to predict next feature positions**

What cameras and IMU have in common: both can be used to estimate the pose incrementally (known as dead-reckoning), which suffers from drift over time. **Solution: fuse them together to reduce drift (see later)**

IMU Measurement Model

The model measures angular velocity $\tilde{\omega}_B(t)$ and acceleration $\tilde{a}_B(t)$ in the body frame B :

$$\begin{aligned}\tilde{\omega}_B(t) &= \omega_B(t) + b^g(t) + n^g(t) && \text{IMU biases + noise in body frame} \\ \tilde{a}_B(t) &= R_{BW}(t)(a_W(t) - g) + b^a(t) + n^a(t)\end{aligned}$$

Raw IMU measurements true ω (in body frame) and true a (in world frame) to estimate

Notation:

- The superscript g stands for gyroscope and a for accelerometer
- R_{BW} is the rotation of the World frame W with respect to Body frame B
- The gravity g is expressed in the World frame
- Biases and noise are expressed in the body frame
- Additive, zero-mean Gaussian white noise: $n^g(t), n^a(t)$
- Biases: $b^g(t), b^a(t)$
 - The gyroscope and accelerometer biases are considered **slowly varying “constants”**. Their temporal fluctuation is modeled by saying that the derivative of the bias is a zero-mean Gaussian noise with standard deviation σ_b

12

$$\dot{b}(t) = \sigma_b w(t) \quad w(t) \sim N(0,1)$$

- Some facts about IMU biases:
 - They change with **temperature** and **mechanical and atmospheric pressure**
 - Thus they **may also be different every time the IMU is turned on**
 - Good news: **they can be estimated!** (see later)

IMU Integration Model

- The **IMU Integration Model** computes the position, orientation, and velocity of the IMU in the world frame. To do this, we must first compute the acceleration $a(t)$ in the world frame from the measured one $\tilde{a}(t)$ in the body frame (see Slide 13):

$$a(t) = R_{WB}(t)(\tilde{a}(t) - b(t)) + g$$

- The position p_k at time t_k can then be **predicted** from the position p_{k-1} at time t_{k-1} by integrating all the inertial measurements $\{\tilde{a}_j, \tilde{\omega}_j\}$ within that time interval:

$$p_k = p_{k-1} + v_{k-1}(t_k - t_{k-1}) + \int_{t_{k-1}}^{t_k} (R_{WB}(t)(\tilde{a}(t) - b(t)) + g) dt^2$$

- A similar expression can be obtained to predict the velocity v_k and orientation R_{WB} of the IMU in the world frame as functions of both \tilde{a}_j and $\tilde{\omega}_j$

For convenience, the **IMU Integration Model** is normally written as

$$\begin{pmatrix} p_k \\ q_k \\ v_k \end{pmatrix} = f \begin{pmatrix} p_{k-1} \\ q_{k-1} \\ v_{k-1} \end{pmatrix} \quad \text{or, more compactly: } x_k = f(x_{k-1}, u)$$

where:

- $x = \begin{bmatrix} p \\ q \\ v \end{bmatrix}$ represents the **IMU state**, i.e., **position, orientation, and velocity**
- q is the IMU **orientation** R_{WB} (usually represented using quaternions)
- $u = \{\tilde{a}_j, \tilde{\omega}_j\}$ are the accelerometer and gyroscope measurements in the time interval $[t_{k-1}, t_k]$

Visual Inertial Odometry

Different paradigms exist:

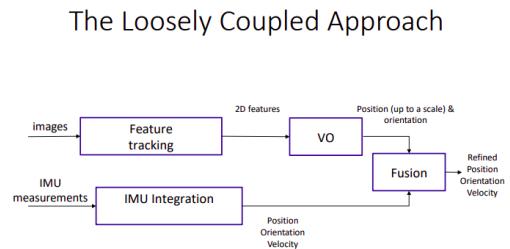
L are the landmarks (3D points)

- **Loosely coupled:**
 - Treats VO and IMU as **two separate black boxes** (not coupled)
 - Each black box estimates pose and velocity from visual (up to a scale) and inertial data (absolute scale)
 - **Easy to implement**
 - **Inaccurate. Should not be used**

- **Tightly coupled:**

- Makes use of the **raw sensors' measurements** (2D features and IMU readings)
 - More accurate
 - More implementation effort

In this lecture, we will only see **tightly coupled approaches**

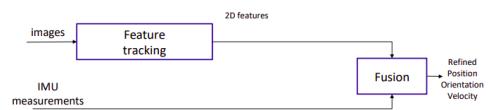


- System states:

The Tightly Coupled Approach

Tightly coupled: $X = [p(t); q(t); v(t); b^a(t); b^g(t); L_1; L_2; \dots; L_k]$

Loosely coupled: $X = [p(t); q(t); v(t); b^a(t); b^g(t)]$



Closed-form Solution (1D case)

- From a **single camera** we only get the absolute position x up to a **unknown scale factor s** , thus

$$x = s\tilde{x}$$

where \tilde{x} is the relative motion estimated by the camera.

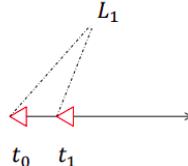
- From the IMU

$$x = x_0 + v_0(t_1 - t_0) + \iint_{t_0}^{t_1} a(t) dt^2$$

- By equating them

$$s\tilde{x} = x_0 + v_0(t_1 - t_0) + \iint_{t_0}^{t_1} a(t) dt^2$$

- As shown in [Martinelli'14], if we assume to know x_0 (usually we set it to 0), then, for 6DOF, both s and v_0 can be determined in closed form from a single feature observation and 3 views



Non-linear Optimization Methods

VIO is solved as a non-linear Least Square optimization problem over:

where

- $X = \{x_1, \dots x_N\}$: set of state **estimates** x_k (**position, velocity, orientation**) at frame times k
 - $L = \{L_1, \dots, L_M\}$: **3D landmarks**
 - $f(x_{k-1}, u)$: **state prediction** obtained by integrating IMU measurements $u = \{\tilde{a}_j, \tilde{\omega}_j\}$
 - $\pi(x_k, l_i)$: **expected measurement at state estimates** x_k from projection of landmark L_i onto camera frame I_k
 - z_{i_k} : **observed features**
 - Λ_k : **state covariance** from the IMU integration $f(x_{k-1}, u)$
 - Σ_{i_k} : is the covariance from the noisy 2D feature measurements

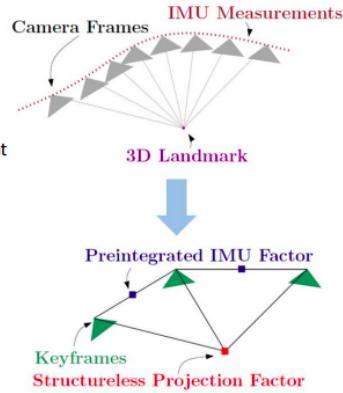
Case Study 1: OKVIS

Because the complexity of the optimization is cubic with respect to the number of cameras poses and features, real-time operation becomes infeasible as the trajectory and the map grow over time, OKVIS proposed to only **optimize the current pose and a window of past keyframes**

Case Study 2: SVO+GTSAM

Solves the same optimization problem as OKVIS but:

- Keeps all the frames (from the start to the end of the trajectory)
- To make the optimization efficient
 - Marginalizes 3D landmarks (rather than triangulating landmarks, only their visual bearing measurements are kept and the visual constraint are enforced by the Epipolar Line Distance (Lecture 08))
 - pre-integrates the IMU data between keyframes (see later)
- Optimization solved using Factor Graphs via [GTSAM](#)
 - Very fast because it only optimizes the poses that are affected by a new observation



Problem with IMU integration

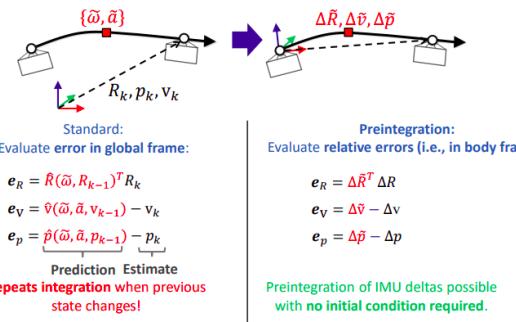
IMU Pre-Integration

- The integration of IMU measurements, $f(x_{k-1}, u)$, from $k - 1$ to k is related to the state estimate at time $k - 1$
- During optimization, every time the linearization point at $k - 1$ changes, the integration between $k - 1$ and k must be re-evaluated, thus slowing down the optimization

$$\{X, L, b^a, b^g\} = \operatorname{argmin}_{\{X, L, b^a, b^g\}} \left\{ \sum_{k=1}^N \|f(x_{k-1}, u) - x_k\|_{A_k}^2 + \sum_{k=1}^N \sum_{l=1}^M \|\pi(x_k, L_l) - z_{i_k}\| \right\}$$

IMU residuals Reprojection residuals
(Bundle Adjustment term)

- **Idea: Preintegration**
 - defines relative motion increments, expressed in body frame, which are independent on position, orientation, and velocity at k [1]
 - [2] uses this theory by leveraging the manifold structure of the rotation group $\text{SO}(3)$



Non-linear Optimization methods	Filtering methods
Optimize a window of multiple states (or all the states) using non-linear Least-Squares optimization	Solve the same problem by running only one iteration of the optimization function (e.g., using Extended Kalman Filter (EKF))
<ul style="list-style-type: none"> ✓ Multiple iterations (it re-linearizes at each iteration) ✓ Achieves the highest accuracy ✗ Slower 	<ul style="list-style-type: none"> ✗ One iteration only ✗ Sensitive to linearization point ✓ Fastest

Case Study 1: ROVIO

- EKF state: $\mathbf{X} = [p(t); q(t); v(t); b^a(t); b^g(t); L_1; L_2; \dots; L_k]$
- Basic idea:
 1. **Prediction step:** predicts next position, velocity, orientation, and features using IMU integration model
 2. **Measurement update:** refines state by leveraging visual constraint (ROVIO minimizes the photometric error between corresponding points (alternative would be the reprojection error))
- Complexity of the EKF grows **quadratically** in the number of estimated landmarks
 - Thus, **max 20 landmarks** are tracked to allow real-time operation
- Only updates the most recent state

Case Study 2: MSCKF

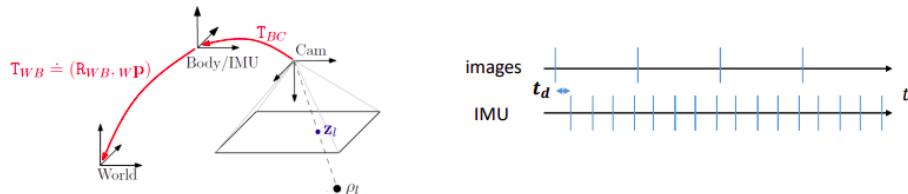
- **MSCKF (Multi-State Constraint Kalman Filter)** updates **multiple past poses** $\{p_{C_1}, q_{C_1}, \dots, p_{C_N}, q_{C_N}\}$ in addition to the current state $\{p(t), q(t), v(t)\}$. State vector:

$$\mathbf{X} = [p(t); q(t); v(t); b^a(t); b^g(t); p_{C_1}; q_{C_1}; \dots; p_{C_N}; q_{C_N}]$$

- **Prediction step:** same as ROVIO
- **Measurement update:**
 - Differently from ROVIO,
 - landmark positions are **not added to the state vector**, thus can run very fast independently of the number of features
 - Visual constraint is obtained from the **Epipolar Line Distance** (Lecture 08)
- Used in spacecraft landing (**NASA/JPL Moon and Mars landing**), **DJI drones**, **Google ARCore**
- Released open source within the OpenVins project: https://github.com/rpng/open_vins

Camera-IMU Calibration

- **Goal:** estimate the **rigid-body transformation** T_{BC} and time **offset** t_d between the **camera** and the **IMU** caused by communication delays and the internal sensor delays (introduced by filters and logic).
- **Assumptions:** Camera and IMU rigidly attached. Camera intrinsically calibrated.
- **Data:**
 - Image points from calibration pattern (checkerboard or QR board)
 - IMU measurements: accelerometer $\{a_k\}$ and gyroscope $\{\omega_k\}$



Kalibr Toolbox

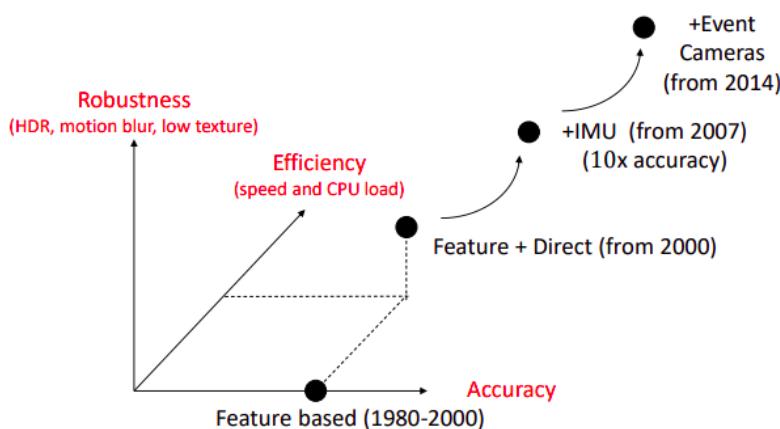
- Solves a non-linear Least Square optimization problem similar to that seen before but also optimizes over T_{BC}, t_d :

$$\{X, L, T_{BC}, t_d, b^a, b^g\} = \underset{\{X, L, T_{BC}, t_d, b^a, b^g\}}{\operatorname{argmin}} \left\{ \sum_{k=1}^N \|f(x_{k-1}, u) - x_k\|_{A_k}^2 + \sum_{k=1}^N \sum_{i=1}^M \|\pi(x_k, L_i) - z_{i_k}\|_{\Sigma_{i_k}}^2 \right\}$$

IMU residuals Reprojection residuals
 (Bundle Adjustment term)

- Continuous-time modelling using splines for X
- Numerical solver: Levenberg-Marquardt

A Short Recap of the last 30 years of VIO



Challenges are:

- latency and motion blur
- dynamic range

Event cameras do not suffer from these problems.

They are inspired by human eyes that have many photoreceptors but relatively less axioms.

They offer higher sensitivity (for example in low light conditions)

What is an Event Camera

First commercialized by Prof. T. Delbrück in 2008 at the Institute of Neuroinformatics of UZH & ETH under the name of Dynamic Vision Sensor (DVS)

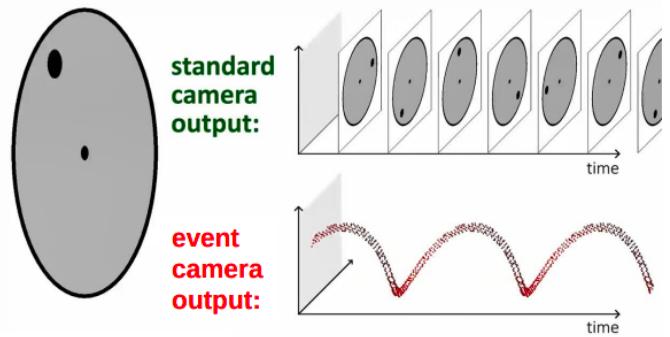
Advantages

- Low-latency (~1 micro-seconds)
- High-dynamic range (HDR) (140 dB instead 60 dB)
- High updated rate (1 MHz)
- Low power (10mW instead 1W)

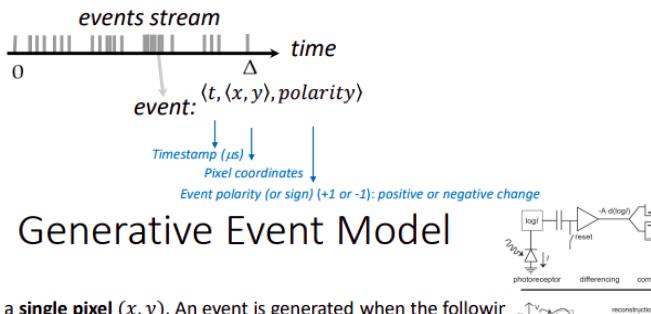
Challenges

- Paradigm shift: Requires totally new vision algorithms because:
 - Asynchronous pixels
 - No intensity information (only binary intensity changes)

- A traditional camera outputs frames at fixed time intervals:



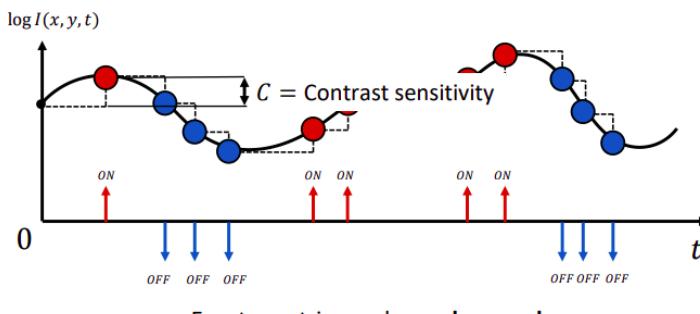
- By contrast, an event camera outputs asynchronous events at microsecond resolution. An event is generated each time a single pixel detects a change of intensity



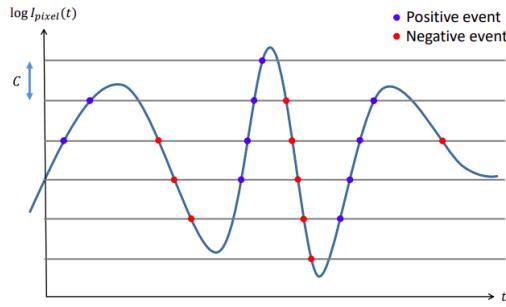
Generative Event Model

- Consider the intensity at a single pixel (x, y) . An event is generated when the following is satisfied:

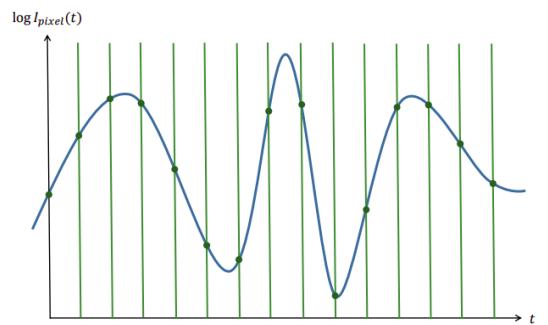
$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$



Event cameras sample intensity when this crosses a threshold
(Level-crossing sampling)



Standard cameras sample intensity at uniform time intervals
(uniform time sampling)



Current commercial applications

- **Monitoring and surveillance**
 - Action and gesture recognition in HDR scenes
- **Industrial automation**
 - Fast object counting
- **Computational photography**
 - Deblurring, super resolution, HDR, slow-motion video
- **Automotive:**
 - low-latency detection, object classification, low-power and low-memory storage

Calibration of an Event Camera

- Standard **pinhole camera model** still valid (same optics)
- Standard passive calibration patterns **cannot be used**
 - need to move the camera → inaccurate corner detection
- **Blinking patterns** (computer screen, LEDs)
- ROS DVS driver + intrinsic and extrinsic mono & stereo calibration: https://github.com/uzh-rpg/rpg_dvs_ros

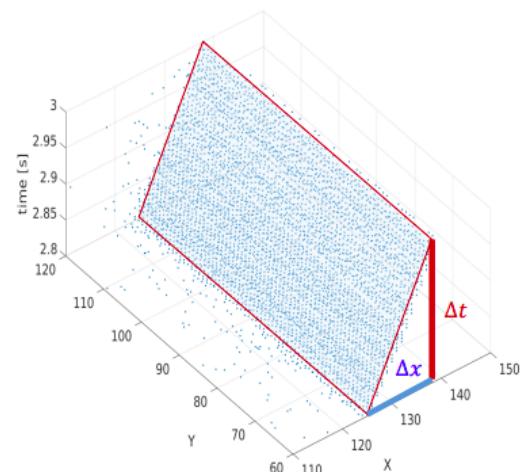
A Simple Optical Flow Algorithm

- Let's assume pure horizontal motion
- White pixels become black → brightness decrease → negative events (in black color)



- The same edge, visualized in space-time
- Events are represented by dots

The edge is moving at
a speed of:
 $v = \frac{\Delta x}{\Delta t}$

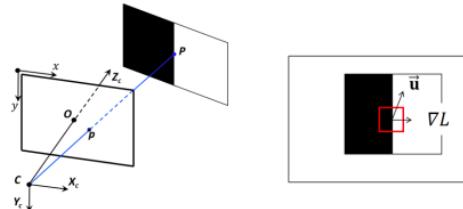


1st order approximation of the Generative Event Model

- An event is generated when the following condition is satisfied:

$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$

- For many applications, it is convenient to derive a 1st order approximation
- Let us define $L(x, y, t) = \log(I(x, y, t))$
- Consider a given pixel $p(x, y)$ with gradient $\nabla L(x, y)$ undergoing the motion $\mathbf{u} = (u, v)$ in pixels, induced by a moving 3D point P



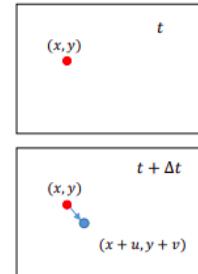
25

- Let's apply the **brightness constancy assumption**, which says that the intensity value of p before and after the motion must remain unchanged:

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

- By replacing the right-hand term with its 1st order approximation at $t + \Delta t$, we get:

$$\begin{aligned} L(x, y, t) &= L(x, y, t + \Delta t) + \frac{\partial L}{\partial x} u + \frac{\partial L}{\partial y} v \\ \Rightarrow L(x, y, t + \Delta t) - L(x, y, t) &= -\frac{\partial L}{\partial x} u - \frac{\partial L}{\partial y} v \\ \Rightarrow \pm C &= -\nabla L \cdot \mathbf{u} \end{aligned}$$



- This formula shows that **maximum generation of events** (i.e., higher event rate) occurs when the **relative motion of the camera is perpendicular to the edge** and is **minimum when parallel to the edge**.

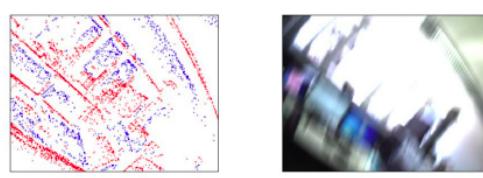
Application 1: Image Reconstruction from events

- Probabilistic simultaneous gradient reconstruction and rotation estimation from $\pm C = -\nabla L \cdot \mathbf{u}$
- Obtain **image intensity from gradient** via Poisson reconstruction
- The reconstructed image has **super-resolution and High Dynamic Range (HDR)**
- Can run in **real time on a GPU**

Application 2: 6DoF Tracking from Photometric Map

- Probabilistic **6DoF motion estimation** from $\pm C = -\nabla L \cdot \mathbf{u}$
- Assumes **photometric map** (x, y, z , grayscale Intensity) is **given**
- Useful for **VR/AR applications** (low-latency, HDR, no motion blur)
- Can run in **real time on a GPU**

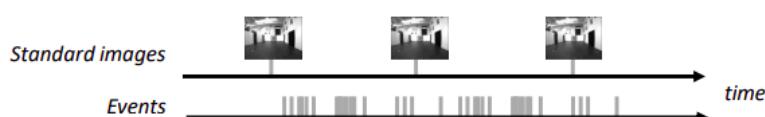
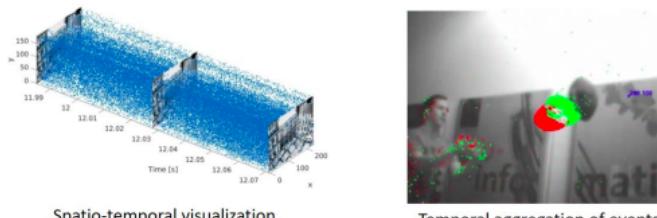
Combining Standard Cameras with Event Cameras



	Event Camera	Standard Camera
Update rate	High (asynchronous): 1 MHz	Low (synchronous)
Dynamic Range	High (140 dB)	Low (60 dB)
Motion Blur	No	Yes
Static motion	No (event camera is a high pass filter)	Yes
Absolute intensity	No (but reconstructable up to a constant)	Yes
Maturity	< 10 years of research	> 60 years of research!

DAVIS sensor: Events + Images + IMU

- Combines an **event** and a **standard camera** in the **same pixel array** (\rightarrow the same pixel can both trigger events and integrate light intensity).
- It also has an **IMU**



Application 1: Deblurring a blurry video

- Idea:** A **blurry image** can be regarded as the **integral of a sequence of latent images** during the exposure time, while the **events** indicate the **changes between the latent images**
- Solution:** sharp image obtained by subtracting the double integral of event from input image

$$\log \text{Input blur image} - \iint \text{Input events} = \log \text{Output sharp image}$$

Application 3: Event-based KLT Tracking

- Goal:** Extract **features from standard frames** and track them using only **events** in the **blind time between two frames**
- Uses the event generation model via joint estimation of patch warping and optic flow

Recap

- All the approaches seen so far use the **generative event model**

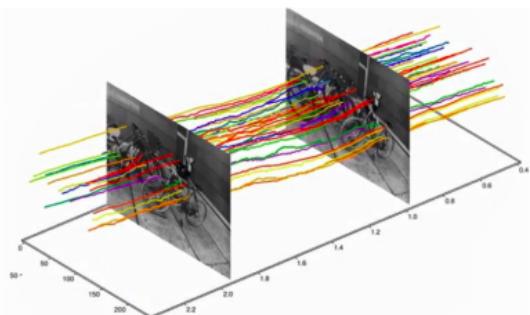
$$\log I(x, y, t + \Delta t) - \log I(x, y, t) = \pm C$$

- or its 1st order approximation

$$\pm C = -\nabla L \cdot \mathbf{u}$$

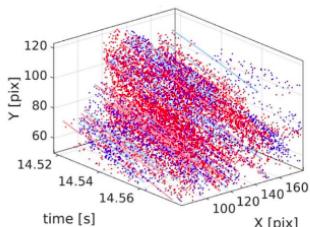
which **requires knowledge of the contrast sensitivity C**

- Unfortunately, C is **scene dependent** and might differ from pixel to pixel
- Alternative approach:** Contrast maximization framework

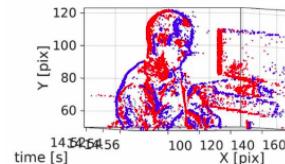


Contrast Maximization Framework

Idea: Warp spatio-temporal volume of events to **maximize contrast** (e.g., sharpness) of the resulting image

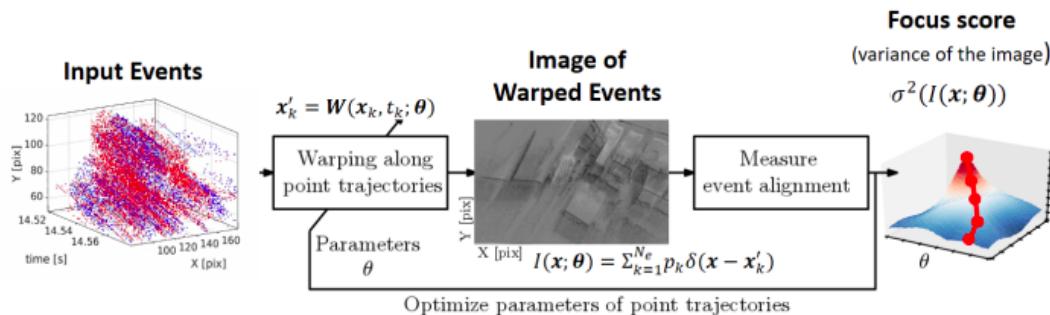


Aggregated image
without motion correction



Aggregated image
with motion correction

Contrast Maximization Framework



- $x'_k = W(x_k, t_k; \theta)$: This warps the (x, y) pixels coordinates of each event, not their time. Possible warps: roto-translation, affine, homography.
- $I(x; \theta) = \sum_{k=1}^{N_e} p_k \delta(x - x'_k)$: This builds a grayscale image, where the intensity of each pixel at the warped location (x', y') is equal to the summation of the polarity p (i.e., positive and negative events (+1, -1))
- $\sigma^2(I(x; \theta))$: The assumption here is that if an image contains **high variance** then there is a wide **spread of responses, both edge-like and non-edge like**, representative of a normal, in-focus image. But if there is **very low variance**, then there is a tiny spread of responses, indicating there are very little edges in the image. As we know, the more an image is blurred, the less edges there are.

Application 1: Image Stabilization



- Goal: **Estimate rotational motion (3DoF)** of an event camera
- Can process millions of events per second in real time on a smartphone PC (e.g., OdroidXU4)
- Works up to over $\sim 1,000$ deg/s

Application 2: Motion Segmentation

Application 3: Dynamic Obstacle Avoidance

- Works with relative speeds of up to **10 m/s**
- Perception latency: **3.5 ms**

Application 4: “Ultimate SLAM”

Goal: combining **events, images**, and **IMU** for robust visual SLAM in HDR and high speed scenarios

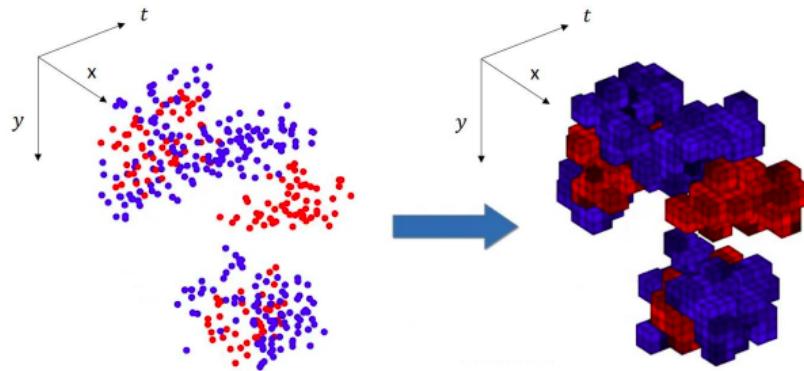
85% accuracy gain over standard VIO in HDR and high speed scenarios → enable autonomous navigation in low light



Input representation

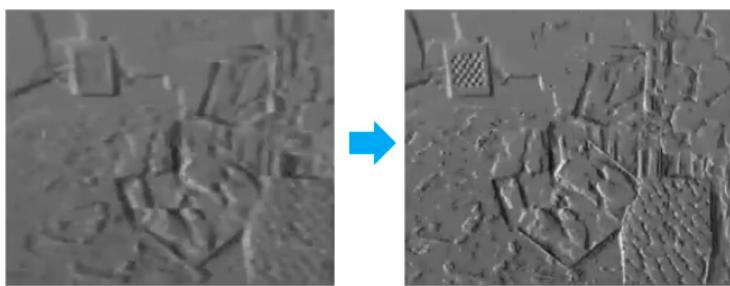
How do we pass sparse events into a convolutional neural network designed for standard images?

Represent events in space-time into a 3D voxel grid (x, y, t): each voxel contains sum of positive and negative events falling within the voxel



Contrast as Loss Function for Unsupervised Learning

- Idea: maximize sharpness of the aggregated event image



Applications:

1. Unsupervised learning of optical flow → ideas is to maximize the sharpness of the aggregated event image
2. **Image reconstruction** from events → see next slide
3. **Slow motion video**: we can combine an event camera with an HD RGB camera. We use events to upsample low-frame rate with low memory footprint

Overview

- Recurrent neural network (main module: Unet)
- Input: sequences of *event tensors* (3D spatio-temporal volumes of events^[3])
- Trained in simulation only, without seeing a single real image
- To improve robustness we randomize the contrast sensitivity during simulation.
- Event camera simulator (ESIM): <http://rpg.ifi.uzh.ch/esim.html>

