

## Forward Imaging Model : 3D to 2D.

world coordinates.

$$\underline{x}_w = \begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix}$$

camera coordinates.

$$\underline{x}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \end{bmatrix}$$

image coordinates.

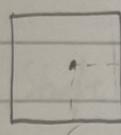
$$\underline{x}_i = \begin{bmatrix} x_i \\ y_i \end{bmatrix}$$

coordinates transformation

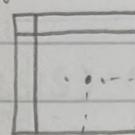
perspective projection

image plane  $\rightarrow$  image sensor mapping

image plane has  $(\hat{x}_i, \hat{y}_i)$   $\rightarrow$  image sensor has  $(u, v)$



unit: mm  $\rightarrow$

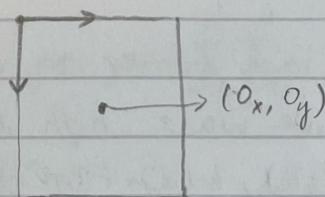


(pixel can be unit pixel.  
rectangle).

$m_x, m_y$ : pixel densities (pixels/mm) in x & y directions.

$$u = m_x x_i = m_x f \frac{x_c}{z_c}, \quad v = m_y y_i = m_y f \frac{y_c}{z_c}$$

we assume the location of principle point (centre of the image) as  $(0_x, 0_y)$ . This is unknown, and usually origin is default to top left corner.



a scale factor which converts metres in 3D world coord into distances in pixels.

$$u = f_x \frac{x_c}{z_c} + 0_x, \quad v = f_y \frac{y_c}{z_c} + 0_y$$

where  $(f_x, f_y) = (m_x f, m_y f)$  are focal lengths in pixels in x & y directions.

params which align lens projection centre (the point of intersection from pinhole/optical axis ray to the sensor plane) wrt centre of the sensor plane.

$f_x, f_y, 0_x, 0_y$ : intrinsic params of the camera. They represent camera's internal geometry.

intrinsic matrix:

perspective  
projection.

$$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} f_x & 0 & o_x & 0 \\ 0 & f_y & o_y & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$$

camera intrinsic  
calibration.  
use  $[u]$ ,  $[v]$  to

estimate  
focal lengths in pixels

$f_x$ ,  $f_y$

& principal point in  
pixels (usually close to  
centre of the image)

$o_x$ ,  $o_y$ .

where  $M_{int}$  is the intrinsic matrix,

$$M_{int} = [K|0]$$

&  $K$  is the calibration matrix,

$$K = \begin{bmatrix} f_x & 0 & o_x \\ 0 & f_y & o_y \\ 0 & 0 & 1 \end{bmatrix}$$
 which is an upper right  
triangular matrix.

$\begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$  is the pixel coord.  $\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$  is the 'camera coord.' use homogeneous

3D coord.  $\begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$  & homogeneous 4D coord.  $\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$  to construct a matrix for  
easy computation.

in world coord.

$$\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$$

► in pixel coordinates. This is an approx. w/ few assumptions:

1. we use a pinhole camera model  $\Rightarrow$  i.e. all light rays go through a single point (i.e.  $(o_x, o_y)$ ) instead of a lens.
2.  $\Rightarrow$  the camera is reduced to a 3D point  $\overset{\text{(location at)}}{\underset{\text{the optical centre}}{\text{}}}.$  (as the optical centre is where all rays 'originate' from the pinhole camera).

optical axis is the vector from optical centre  $\perp$  to the image plane.

$\begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix}$  describe the 3D location of the obj. relative to the camera.

extrinsic parameters.

rotation matrix  $R$ :

$$\begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{bmatrix} \rightarrow \text{direction of } \hat{x}_c \text{ in world coordinate frame.}$$

$R$  is orthonormal:

NB:  $u$  and  $v$  are orthonormal iff:

$$1. \underline{u}^T \underline{v} = 0$$

$$2. \underline{u}^T \underline{u} = \underline{v}^T \underline{v} = 1$$

orthonormal matrix

$$R^{-1} = R^T \quad R^T R = R R^T = I$$

extrinsic matrix:  $M_{ext}$

$$\tilde{\underline{x}}_c = \begin{bmatrix} x_c \\ y_c \\ z_c \\ 1 \end{bmatrix} = \underbrace{\begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{M_{ext}} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

$$M_{ext} = \begin{bmatrix} R_{3 \times 3} & \underline{t} \\ 0_{1 \times 3} & 1 \end{bmatrix}$$

$$\tilde{\underline{x}}_c = M_{ext} \tilde{\underline{x}}_w \leftarrow \text{coordinate transformation}$$

$$\begin{bmatrix} x_w \\ y_w \\ z_w \end{bmatrix} = \text{world coord. position of obj. relative to world origin.}$$

The transformation,  $M_{ext}$ , is a rigid body transformation (i.e. no shearing / deformation) that first rotates all points in world coord (using matrix  $R$ ) then translates all points (using translation vector  $\underline{t}$ ). s.t. From the camera's perspective, the world is now transformed as if the camera is fixed at the origin looking straight ahead.

world to camera , camera to pixel.

$$M_{ext} \tilde{\underline{x}}_w = \tilde{\underline{x}}_c$$

$$M_{int} \tilde{\underline{x}}_c = \tilde{\underline{u}}$$

projection matrix  $P$ :

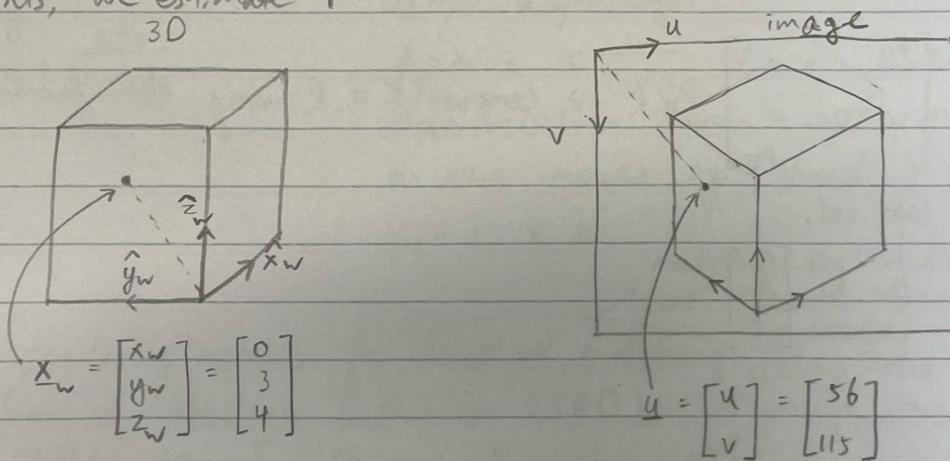
$$\tilde{\mathbf{x}} = M_{int} M_{ext} \tilde{\mathbf{x}}_w = P \tilde{\mathbf{x}}_w$$

$$\begin{bmatrix} \tilde{u} \\ \tilde{v} \\ \tilde{w} \end{bmatrix} = \begin{bmatrix} P_{11} & P_{12} & P_{13} & P_{14} \\ P_{21} & P_{22} & P_{23} & P_{24} \\ P_{31} & P_{32} & P_{33} & P_{34} \end{bmatrix} \begin{bmatrix} x_w \\ y_w \\ z_w \\ 1 \end{bmatrix}$$

we only need to know  $P$  to map all points in world coord. into the pixels in an image.

Camera Calibration example

to identify the correspondence between 3D scene objects & image points, we estimate  $P$ :



for  $i^{th}$  point, unknown:  $\begin{bmatrix} u^i \\ v^i \\ 1 \end{bmatrix} = P \begin{bmatrix} x_w^i \\ y_w^i \\ z_w^i \\ 1 \end{bmatrix}$

$$\begin{bmatrix} u^i \\ v^i \\ 1 \end{bmatrix} = P \begin{bmatrix} x_w^i \\ y_w^i \\ z_w^i \\ 1 \end{bmatrix}$$

known

We rearrange s.t. we put all  $\begin{bmatrix} u^i \\ v^i \end{bmatrix}$  &  $\begin{bmatrix} x_w^i \\ y_w^i \\ z_w^i \end{bmatrix}$  into a single matrix  $A$ ,  
 $A_i \in \{1, \dots, n\}$

and we rearrange matrix  $P$  into a vector  $p$ :

$$\begin{bmatrix} P_{11} & P_{12} \\ & \ddots \\ & & P_{34} \end{bmatrix} \rightarrow p = \begin{bmatrix} P_{11} \\ \vdots \\ P_{14} \\ P_{21} \\ \vdots \\ P_{24} \\ P_{31} \\ \vdots \\ P_{34} \end{bmatrix} \quad Ap = 0$$

to estimate  $p$ ,  $\min_p \|Ap\|^2$  s.t.  $\|p\|^2 = 1$

$A^T A p = \lambda p$   $p$  is the eigenvector corresponding to the smallest eigenvalue  $\lambda$  of matrix  $A^T A$ .

?  $\begin{bmatrix} P_{11} & & \\ & \ddots & \\ & & P_{33} \end{bmatrix} = KR \rightarrow$  compute  $K$  &  $R$  using QR factorisation

(last col,

$$t = K^{-1} \begin{bmatrix} P_{14} \\ P_{24} \\ P_{34} \end{bmatrix}$$

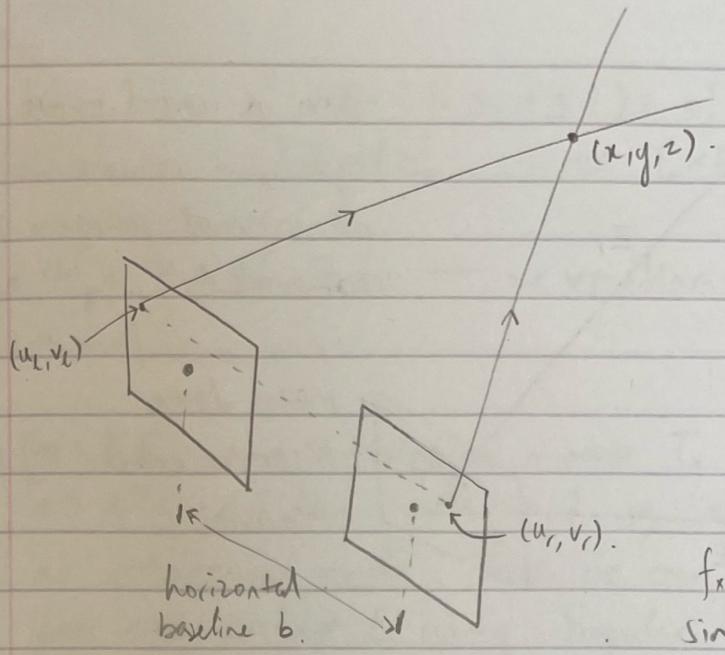
computing 2D  $\rightarrow$  3D: triangulation (need 2 imgs)

binocular vision: 2 cameras used to perceive depth by comparing 2 slightly different views of the same scene.

(Stereo system: 2 camera system).

(Simplified stereo calibration): 2 cameras place at the same height,  $b$  apart  $\Rightarrow$  horizontal baseline.

$$0 \leftarrow b \rightarrow 0 \quad vs.$$



perspective projection eqn:

$$u_L = f_x \frac{x}{z} + o_x$$

$$v_L = f_y \frac{y}{z} + o_y$$

$$u_r = f_x \frac{x-b}{z} + o_x$$

$$v_r = f_y \frac{y}{z} + o_y$$

$f_x, f_y, o_x, o_y, b$  are known  
since we've already calculated  
camera calibration at this  
point.

solve for  $(x, y, z)$ :

$$\begin{cases} x = \frac{b(u_L - o_x)}{u_L - u_r} \\ y = \dots \\ z = \frac{bf_x}{u_L - u_r} \end{cases}$$

$u_L - u_r$ : disparity.

depth  $z \propto \frac{1}{u_L - u_r}$  (i.e. further away,  
then location doesn't appear that different  
on the images. e.g. bud.).

$u_L - u_r \propto b$ .

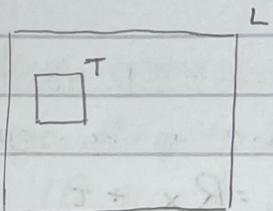
matching the smaller 'patch'  
i.e. the template image with

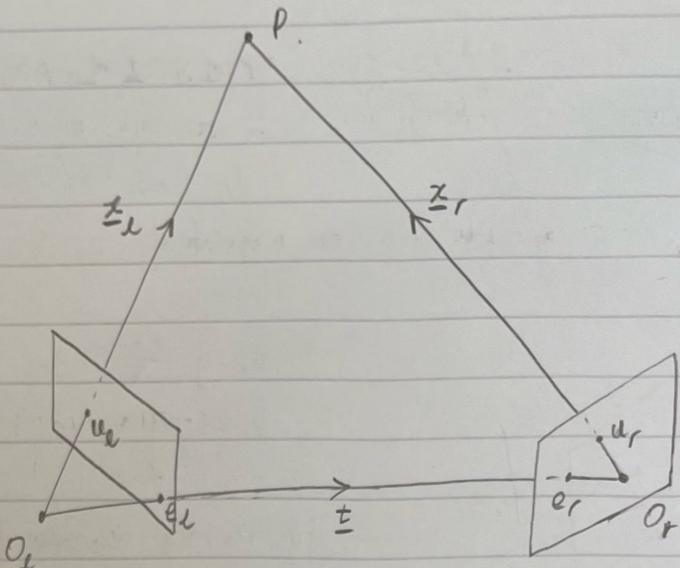
use similarity metrics for template matching: sliding the template  
eg. Sum of Absolute Differences (SAD) through the img.

$$SAD(k, l) = \sum_{(i,j) \in T} |E_e(i, j) - E_e(i+k, j+l)|$$

for pixel  $(k, l) \in L$  (large img);  $(i, j) \in T$  (template).

SSD, NCC.





even in uncalibrated stereo system, we assume we know all internal params of 2 cameras,  $f_x^{(l)}, f_y^{(l)}, o_x^{(l)}, o_y^{(l)}$ ,  $\dots, \dots, \dots, \dots$ , from the image's meta data  
 $\Rightarrow$  camera matrix  $K_l, K_r$  is known.

$O_l, O_r$  are camera centres (pinhole camera).

point  $P$  of the scene,  $O_l, O_r$  construct a unique epipolar plane. Plane  $O_l P O_r$  intersect the image planes at  $u_l, u_r$  are  $u_l, u_r$ , these 2 lines are known as epipolar lines.

camera's projection onto the other plane (i.e. the intersection w/ the image plane), at points  $u_l$  and  $u_r$ , are epipoles.

$\overrightarrow{O_l O_r}$  is the translation vector,  $t$ .

we have vector  $n$  normal to the epipolar plane,  $n = t \times z_l$ . cross prod

$$\text{epipolar constraint: } \underline{x}_l \cdot (t \times \underline{z}_l) = 0 = [x_l \ y_l \ z_l] \begin{bmatrix} t_y z_l - t_z y_l \\ t_z x_l - t_x z_l \\ t_x y_l - t_y x_l \end{bmatrix}$$

$$= [x_l \ y_l \ z_l] \begin{bmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ -t_y & t_x & 0 \end{bmatrix} \begin{bmatrix} x_l \\ y_l \\ z_l \end{bmatrix} = 0$$

$\downarrow$   
translation matrix,  $T_x$ .

$\hookrightarrow$  the coord of  $P$   
in left camera frame.

3D-to-3D translation of a point from left camera frame to its representation in the right camera frame.

$$\underline{x}_l = R \underline{x}_r + t$$

$$\begin{bmatrix} x_l \\ y_l \\ z_l \end{bmatrix} = \begin{bmatrix} r_{11} & & \\ & \ddots & \\ & & r_{33} \end{bmatrix} \begin{bmatrix} x_r \\ y_r \\ z_r \end{bmatrix} + \begin{bmatrix} t_x \\ t_y \\ t_z \end{bmatrix}$$

$$\begin{aligned} \underline{x}_e^T \cdot (T_x \times (R \underline{x}_r + \underline{t})) &= \underline{x}_e^T \cdot (T_x \times R \underline{x}_r + \underline{t} \times \underline{t}) \\ &= \underline{x}_e^T \cdot T_x \times R \underline{x}_r \quad \text{let } T_x \times R = E \text{ be the essential matrix.} \\ &= \underline{x}_e^T \cdot E \underline{x}_r \leftarrow \text{our epipolar constraint.} \end{aligned}$$

$E = T_x R$ , where translation matrix  $T_x$  is skew-symmetric matrix.  
 $\Leftrightarrow A^T = -A$  (aka,  $a_{ij} = -a_{ji}$ ) eg.  $a_{12} = -a_{21} = -a_{21}$ .  $R$  is  
orthonormal matrix. Thus, we can decompose  $T_x$  and  $R$   
from their product  $E$  using Singular Value Decomposition.

$\underline{x}_e, \underline{x}_r$  are unknown. we use image coord  $u_e, u_r$ .

$$[u_e \ v_e \ 1] \underbrace{K_e^{-1} E K_r^{-1}}_{\text{a } 3 \times 3 \text{ matrix.}} \begin{bmatrix} u_r \\ v_r \\ 1 \end{bmatrix} = 0.$$

$$\text{fundamental matrix } F := K_e^{-1} E K_r^{-1} \Leftrightarrow E = K_e^T F K_r \Leftrightarrow T_x R.$$

(first calc.  $F$  using  $[u_e \ v_e \ 1] F \begin{bmatrix} u_r \\ v_r \\ 1 \end{bmatrix} = 0$ . then calc.  $E = K_e^T F K_r$ ,

(lastly calc.  $T_x$  &  $R$  using SVD, once we work out  $t$  &  $R$ ,  
we have calibrated the uncalibrated camera system).

(also, we find corresponding point pairs in l d r images using  
SIFT:  $(u_e^{(1)}, v_e^{(1)})$ , ...,  $(u_e^{(m)}, v_e^{(m)})$ .

Finding epipolar lines (aka. finding dense correspondence)

any point on left image,  $(u_e^{(i)}, v_e^{(i)})$ , corresponding to 1 point on  
the right epipolar line. So, we perform 1D search on the  
right epipolar line, to find the exact matching point.

given we calculated  $F$ , and  $(u_e^{(i)}, v_e^{(i)})$ , right epipolar line

is given by:  $a_e u_r + b_e v_r + c_e = 0$

where  $a_e = f_{11} u_e + f_{21} v_e + f_{31}$ ,  $b_e = f_{12} u_e + f_{22} v_e + f_{32}$ ,

$c_e = f_{13} u_e + f_{23} v_e + f_{33}$ .

computing depth

Given the system of cameras are calibrated (relative position / translation from one frame to the other known, by  $\underline{x}_e^T E \underline{x}_r$ ), we can work out the depth information from camera B given image produced by A

→ distortion coefficients.

$M_{int} \rightarrow M_{ext}$ .

rectification image.

→ stereo

→ camera calibration: openCV documentation.

Detected 7.2.19

## Keypoint detectors

Keypoint: a distinct location in an image.

Descriptor: a vector which summarises the local structure around the keypoint.

Keypoint detectors: → tells me positions of local distinct points.

1. Harris.

2. Shi-Tomasi (SOTA, openCV)

3. Förstner

} Structure matrix approaches.

4. Difference of Gaussians ← compare between blurred images.

for structure matrix based approaches: for each pixel  $(x, y)$ , compute the SSD in the local area window centred around  $(x, y)$

$W_{xy}$  as:

$$f(x, y) = \sum_{(u, v) \in W_{xy}} (I(u, v) - I(u + \delta u, v + \delta v))^2$$

where  $I(u + \delta u, v + \delta v) \approx I(u, v) + [J_x \ J_y] \begin{bmatrix} \delta u \\ \delta v \end{bmatrix}$

$$f(x, y) \approx \sum \left( [J_x \ J_y] \begin{bmatrix} \delta u \\ \delta v \end{bmatrix} \right)^2 = \sum \begin{bmatrix} \delta u \\ \delta v \end{bmatrix}^T \underbrace{\begin{bmatrix} J_x^2 & J_x J_y \\ J_y J_x & J_y^2 \end{bmatrix}}_{\text{structure matrix, } M} \begin{bmatrix} \delta u \\ \delta v \end{bmatrix}$$

structure matrix,  $M$ .

To detect the corner, criterion is that eigenvalues  $\lambda_1, \lambda_2$  of the structure matrix are both big.

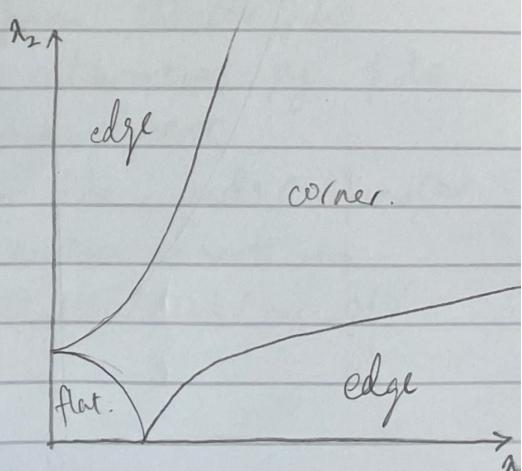
1. Harris corner detection:

$$\begin{aligned} R &= \det(M) - k(\text{trace}(M))^2 \\ &= \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2, \quad k = 0.04 \sim 0.06 \end{aligned}$$

$R \gg 0$ : corner ( $\lambda_1 \approx \lambda_2 \gg 0$ )

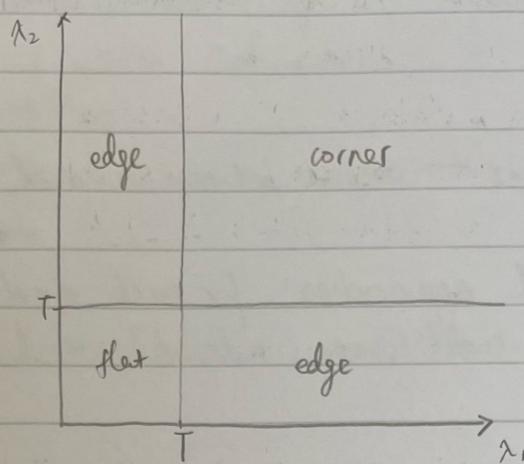
$R \ll 0$ : edge ( $\lambda_1 \gg \lambda_2$  or  $\lambda_2 \gg \lambda_1$ )

$|R| \approx 0$ : flat ( $\lambda_1 \approx \lambda_2 \approx 0$ )



2. suppose we set threshold  $T$   
if smallest eigenvalue  $\lambda_{\min}(M) \geq T \cdot \text{corner}$ .

$$\lambda_{\min}(M) = \frac{\text{trace}(M)}{2} - \frac{1}{2} \sqrt{(\text{trace}(M))^2 - 4 \det(M)}$$



usual practice: 1. RGB  $\rightarrow$  greyscale 2. apply smoothing before doing corner detection

4. DoG1:
1. apply Gaussian smoothing of different kernel size, get imgs with different "degree of blurriness"
  2. compare 2 blurry imgs and get their difference  
(not only compare locally on one img; but also w/ 2 blurry imgs, to find the locally distinct key points).
  3. compute maxima suppression

## Feature descriptors

Feature descriptors: a vector which describes the patch around a keypoint.

1. SIFT: 16 patch  $\times$  8-bin histograms for gradient directions = 128-long vector.

We need to filter out incorrect match: use Lowe's Ratio Test  
given we want to find best match for feature descriptor  $q$ ,  
with either  $p_1$  (best match) &  $p_2$  ( $2^{\text{nd}}$  best match),

1.  $d(q, p_1) < T$ : Euclidean difference between  $q$  &  $p_1$  must be less than threshold  $T$ .
2.  $\frac{d(q, p_1)}{d(q, p_2)} < \frac{1}{2}$ : only accept if  $p_1$  is substantially better than  $p_2$

2. Binary Descriptor: fast, good for online resources / videos.  
short binary str. 01011100

- select a patch around a keypoint.
- select 2 pixels in that patch.
- compare the intensity of those 2 pixels: either 0 or 1.

$$b = \begin{cases} 1 & \text{if } I(s_1) < I(s_2) \\ 0 & \text{o/w} \end{cases}$$

Fast to compare: use Hamming distance

$$d_H(B_1, B_2) = \sum (\text{xor}(B_1, B_2))$$

bitstring 1

example: BRIEF (256-bit binary feature descriptor)

cons: can't identify rotated same features.

solution: ORB (extension of BRIEF)

compared to SIFT, SIFT respond better to scale. ORB (100x faster to compute than SIFT)  
4096-bit

## RANSAC

RANSAC - an algo which separate in-lies (correctly identified pair of keypoints) w/ outliers (incorrectly identified pair of keypoints)  
number of trials to succeed w/ prob.  $p$  is given by:

$$T = \frac{\log(1-p)}{\log(1-(1-e)^s)}$$

most all sampled points  $s$   
are in-lies.

$s$ : # sampled data point (this depend on my image-matching algo. e.g. if I use 3-point algo. to determine the relationship between 2 images,  $s = 5$ . (larger the  $s$ , more difficult to draw because all datapoints  $\leftarrow$  these in-lies using trial-n-error approach). must be in-lies for RANSAC to draw correct results.)  
 $e$ : outlier ratio:  $\frac{\# \text{ outliers}}{\# \text{ total datapoints}}$ . (we don't always know this)

→ in terms of RANSAC,  $T \uparrow$  exponentially with  $e$  &  $s$ ! so, we would prefer an algo. with small  $s$ . e.g. prefer 5-point algo. compared to 8-point algo.

→ RANSAC can't do multiple hypothesis - it's for finding the single best match (of sample points  $s$ , therefore single best correspondence between 2 images).

→ RANSAC used for computing e.g. fundamental matrix, essential matrix. for computing visual odometry information.

## DLT - extrinsic & intrinsic matrix calculation

DLT (direct linear transformation) map 3D world coord onto 2D image plane, assuming perfect lens. (since not counting in distortion, DLT is just an approx.).

camera calibration: compute intrinsic matrix / information

camera localisation: compute extrinsic information

we have 4 word systems:

$$1. \text{ world, } S_o = \begin{bmatrix} x \\ y \\ z \end{bmatrix}$$

$$x_o = \begin{bmatrix} x_o \\ y_o \\ z_o \end{bmatrix}$$

$R$ : rotation matrix

$$2. \text{ camera: origin @ projection centre of the camera } S_k = \begin{bmatrix} kx \\ ky \\ kz \end{bmatrix}$$

$$3. \text{ image plane (2D plane which 3D world project onto. analogue' } S_c = \begin{bmatrix} c_x \\ c_y \\ c_z \end{bmatrix}$$

$$4. \text{ sensor (pixel location: 'digital'.) } S_s = \begin{bmatrix} s_x \\ s_y \end{bmatrix}$$

camera location:  $(x_o, R)$

$x = PX$ .  $x$ : 2D pixel coord  $X$ : 3D world

$$P = K \begin{bmatrix} R & -RX_o \\ 0 & 1 \end{bmatrix} \quad \text{aka. Mat: } \begin{bmatrix} 1 & 0 & t_x \\ 0 & 1 & t_y \\ 0 & 0 & 1 \end{bmatrix}$$

↓  
intrinsic  
 $3 \times 3$

↓  
projection  
 $4 \times 4$

↓  
extrinsics.

we chain the transformations:

$$\begin{bmatrix} s_x \\ s_y \\ 1 \end{bmatrix} = {}^s H_c {}^c P_k \begin{bmatrix} x \\ y \\ z \\ 1 \end{bmatrix} \quad {}^c H_o = \begin{bmatrix} R & -RX_o \\ 0^T & 1 \end{bmatrix}$$

image word  
to  
sensor(pixel)  
word

transform  
object coord  
to  
camera coord.  
(simplified)

camera coord  
to  
image coord.

calibration:  
matrix  
(aka. intrinsic  
matrix)

$$P = {}^s H_c {}^c K [R| -R X_o] = {}^s H_c {}^c K R [I_3 | -X_o] = K R [I_3 | -X_o]$$

Camera matrix:  
map world coord  
→ image plane.

(entire)  
calibration  
matrix

where  $K^c$  is the calibration matrix for ideal camera,  
aka. 'origin of the camera',  
\*in world coord\*.

$-X_o = \begin{bmatrix} -x_o \\ -y_o \\ -z_o \end{bmatrix}$ : projection centre aka. pinhole location  
aka. 'origin of the camera',  
\*in world coord\*.

3D world coord → 3D camera coord is a rigid body transformation,  
done by translation & rotation only.

(distortion)

3. lens are imperfect : (eg. barrel effect: straight line mapped to curves on image plane. to address this issue, we apply non-linear transformation after all linear transformations. This non-linear transformation depend on location of the pixel as well as parameters.

eg. for barrel effect, the transformation required :

$$\begin{cases} {}^a x = x(1 + q_1 r^2 + q_2 r^4) \\ {}^a y = y(1 + q_1 r^2 + q_2 r^4) \end{cases}$$

$r$ : radius from principle point in the image to the pixel.

$q_1, q_2$ : params we need to calculate.

$x, y$ : location of the pixel.

2. In the lecture, calibration matrix  $K$  is breakdown into :  $K = {}^s H_c {}^e K$

$$= \begin{bmatrix} 1 & s & x_H \\ 0 & 1+m & y_H \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} c & 0 & 0 \\ 0 & c & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} c & cs & x_H \\ 0 & c(1+m) & y_H \\ 0 & 0 & 1 \end{bmatrix},$$

$$P = {}^s R [I_3 | -X_o]$$

$c$ : camera constant

$s$ : sheer

$m$ : scale difference

$x_H, y_H$ : principle point.

intrinsics

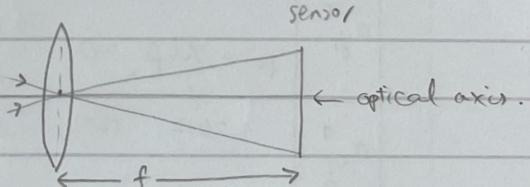
$$\text{where } X_o = \begin{bmatrix} x_o \\ y_o \\ z_o \end{bmatrix}$$

extrinsics

calibration matrix:  $K = \begin{bmatrix} f_x & s & c_x \\ 0 & f_y & c_y \\ 0 & 0 & 1 \end{bmatrix}$   $f_x, f_y$  are focal length in pixels  
 $f$  is the focal length, in mm. - distance between the pinhole (optical centre) to the sensor plane.

$f_x, f_y$  are focal length in pixel units.

$s_x, s_y$ : sensor's pixel pitch: distance between the centre of adjacent pixels, measured in mm.



$$f_x = \frac{t}{s_x}, \quad f_y = \frac{t}{s_y}$$

e.g. my sensor is 4000x3000 pixels, I may have  $c_x = 1980 \text{ px}$ ,  $c_y = 1510 \text{ px}$ .

$c_x, c_y$ : principal-point coordinates, the point which is projected from pinhole along optical axis onto the sensor plane should be v. close to centre of the sensor/image, but usually off set by few mm. in pixel units.

## DLT for calibration

(11 degrees of freedom 6 extrinsic, 5 intrinsic)  
 when we estimate extrinsic & intrinsic matrix, we're assuming affine camera:  
 i.e. no non-linear distortion due to lens. Note when calibrating,  
 should select 3D keypoints in the world that doesn't lie on the  
 same plane. (otherwise estimation becomes unstable)

$$x = Px \Rightarrow \begin{bmatrix} u \\ v \\ w \end{bmatrix} = P \begin{bmatrix} u \\ v \\ w \\ 1 \end{bmatrix}, \text{ 2D 3D coord in homogeneous coord.}$$

(since  $P$  is  $3 \times 4$  matrix. 3D 4D vector makes it convenient that single matrix  $P$  exist).

normalize it

$$\begin{bmatrix} \frac{u}{w} \\ \frac{v}{w} \\ \frac{x}{w} \end{bmatrix} = P \begin{bmatrix} \frac{u}{1} \\ \frac{v}{1} \\ \frac{x}{1} \end{bmatrix} \Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = P \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

$$\text{let } x_i = Px_i = \begin{bmatrix} p_{11} \\ p_{21} \\ p_{31} \end{bmatrix}$$

$X_i$  → vector A  
 $B_i$  → vector B  
 $C_i$  → vector C

$$\Rightarrow x_i = \begin{bmatrix} A^T \\ B^T \\ C^T \end{bmatrix} X_i = \begin{bmatrix} A^T X_i \\ B^T X_i \\ C^T X_i \end{bmatrix}$$

$$\Rightarrow x_i = \frac{u_i}{w_i} = \frac{A^T X_i}{C^T X_i}, \quad y_i = \frac{v_i}{w_i} = \frac{B^T X_i}{C^T X_i}$$

$$\Rightarrow \begin{cases} -X_i^T A + x_i X_i^T C = 0 \\ -X_i^T B + y_i X_i^T C = 0 \end{cases} \quad \text{let } P = \begin{bmatrix} A & B & C \end{bmatrix}$$

rewrite as

$$\begin{cases} a_{x_i}^T P = 0 \\ a_{y_i}^T P = 0 \end{cases} \quad \text{where } a_{x_i}^T = (-X_i^T, 0^T, x_i X_i^T) \\ a_{y_i}^T = (0^T, -X_i^T, y_i X_i^T)$$

stack it into one...

$$\begin{bmatrix} a_{x_1}^T \\ a_{y_1}^T \\ \vdots \\ a_{x_I}^T \\ a_{y_I}^T \end{bmatrix} P = M P = 0$$

2Ix12.      12x1

note  $M$  is known,  $P$  is what we want to solve.

realistically...  $M_P = w$  for small  $w$  which we want min.

$$\Rightarrow \hat{P} = \arg \min_P w^T w = \arg \min_P f^T M^T M f$$

using SVD,  $M = USV^T$  where  $S$  is a diagonal matrix w/ largest first, smallest last (0 would be ideal).  
 $\downarrow$  (now  $P$  is found!)  
 $V^T$  store the singular vectors (eigenvectors)  $\Rightarrow$  choose smallest eigenvector for  $f$ .

# DLT for calibration

now, we obtained  $3 \times 4$  matrix  $P$ . how do we calc.  $K$ ,  $R$ ,  $X_0$ ?

$$P = [KR| - KRX_0] = [H|h]$$

$$\Rightarrow X_0 \text{ (projection centre)} = -H^{-1}h$$

(aka.  $M_{int}$ ).  $\downarrow$  camera location  
in world coord.

$H = KR$ , decompose using QR decomposition.

$$\text{QR decomposition of } H^{-1} = (KR)^{-1} = R^{-1}K^{-1} = R^T K^{-1}$$

then, rotation matrix  $R$ :  $R(z, \alpha)R$  where  $\alpha \rightarrow R(z, \alpha) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & -1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$   
finally, normalise  $K$ :  $\frac{1}{K_{33}} KR(z, \alpha)$ .

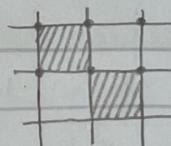
Zhang's method (find intrinsics) - calibration

$$\text{for } x = PX = KR[I_3| - X_0]X$$

↑ image point (sensor)    ↑ world location    ↓ intrinsics, which we want to find out using Zhang's method. (5 params)

using checkerboard for camera calibration (assuming we know size & structure - aka. flat surface)

For every image, we define a coord s.t. the checkerboard forms the XY plane  $\Rightarrow$  All points on the checkerboard, has  $Z=0$ .



$$\begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} c & cs & x_H \\ 0 & cl + tm & y_H \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_1 \\ r_{21} & r_{22} & r_{23} & t_2 \\ r_{31} & r_{32} & r_{33} & t_3 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \\ 1 \end{bmatrix}$$

now, for checkerboard points,  $Z=0$ .

$$\Rightarrow \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \begin{bmatrix} -K \\ I \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & t_1 \\ r_{21} & r_{22} & t_2 \\ r_{31} & r_{32} & t_3 \end{bmatrix} \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}, \text{ let } \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix} \text{ unknown which we need to calc.}$$

$$\Rightarrow \text{let } H = \begin{bmatrix} | & | & | \\ h_1 & h_2 & h_3 \\ | & | & | \end{bmatrix} = K \begin{bmatrix} | & | & | \\ r_1 & r_2 & t \\ | & | & | \end{bmatrix} \Rightarrow \begin{bmatrix} x_i \\ y_i \\ 1 \end{bmatrix} = H \begin{bmatrix} X_i \\ Y_i \\ 1 \end{bmatrix}$$

# Zhang's method

using DLT,  $\begin{cases} \underline{a_x}^T h = 0 \\ \underline{a_y}^T h = 0 \end{cases}$

with  $\underline{a_x}^T = (-X_i, -Y_i, H, 0, 0, 0, x_i X_i, y_i Y_i, x_i)$

$\underline{a_y}^T = (0, 0, 0, -X_i, -Y_i, -1, y_i X_i, y_i Y_i, y_i)$

we need at least 4 points to determine estimate for H.

$H = K [I_3, r_2, I_3]$ , let  $B = K^T K^{-1}$ .  $K$  is upper triangular  $\Rightarrow K$  is invertible.

$[I_3, r_2, I_3] = K^{-1} [h_1, h_2, h_3]$

$$\begin{cases} r_1 = K^{-1} h_1 \\ r_2 = K^{-1} h_2 \end{cases}$$

$r_1, r_2, r_3$  are orthonormal  $\Rightarrow r_1^T r_2 = 0, \|r_1\| = \|r_2\| = 1$

B

$\Rightarrow$  constraints:  $\begin{cases} \underline{h}_1^T (K^T K^{-1}) \underline{h}_2 = 0 \\ \underline{h}_1^T K^T K^{-1} \underline{h}_1 - \underline{h}_2^T K^T K^{-1} \underline{h}_2 = 0 \end{cases}$

(Cholesky decomposition:  $\text{chol}(B) = AA^T$  where  $A = K^{-T}$ )

define  $b = (b_{11}, b_{12}, b_{13}, b_{22}, b_{23}, b_{33})$

$$B = \begin{bmatrix} b_{11} & & \\ & \ddots & \\ & & b_{33} \end{bmatrix} \quad \text{coefficients}$$

2x6 matrix

$$\Rightarrow \text{constraints: } \begin{cases} \underline{v}_{12}^T b = 0 \\ \underline{v}_{11}^T b - \underline{v}_{22}^T b = 0 \end{cases} \Rightarrow \begin{pmatrix} \underline{v}_{12}^T \\ \underline{v}_{11}^T - \underline{v}_{22}^T \end{pmatrix} b = 0$$

↑  
1x6

Given n images, define matrix V s.t.  $V^T b = 0$ .

$$\text{img 1} \rightarrow \begin{pmatrix} \underline{v}_{11}^T \\ \underline{v}_{11}^T - \underline{v}_{22}^T \\ \vdots \end{pmatrix} b = 0. \quad \text{impose } \|b\| = 1,$$

$$\text{img } n \rightarrow \begin{pmatrix} \underline{v}_{12}^T \\ \vdots \\ \underline{v}_{11}^T - \underline{v}_{22}^T \end{pmatrix}$$

using SVD, smallest eigenvalue =  $b^* = \underset{b}{\operatorname{argmin}} \|Vb\|$