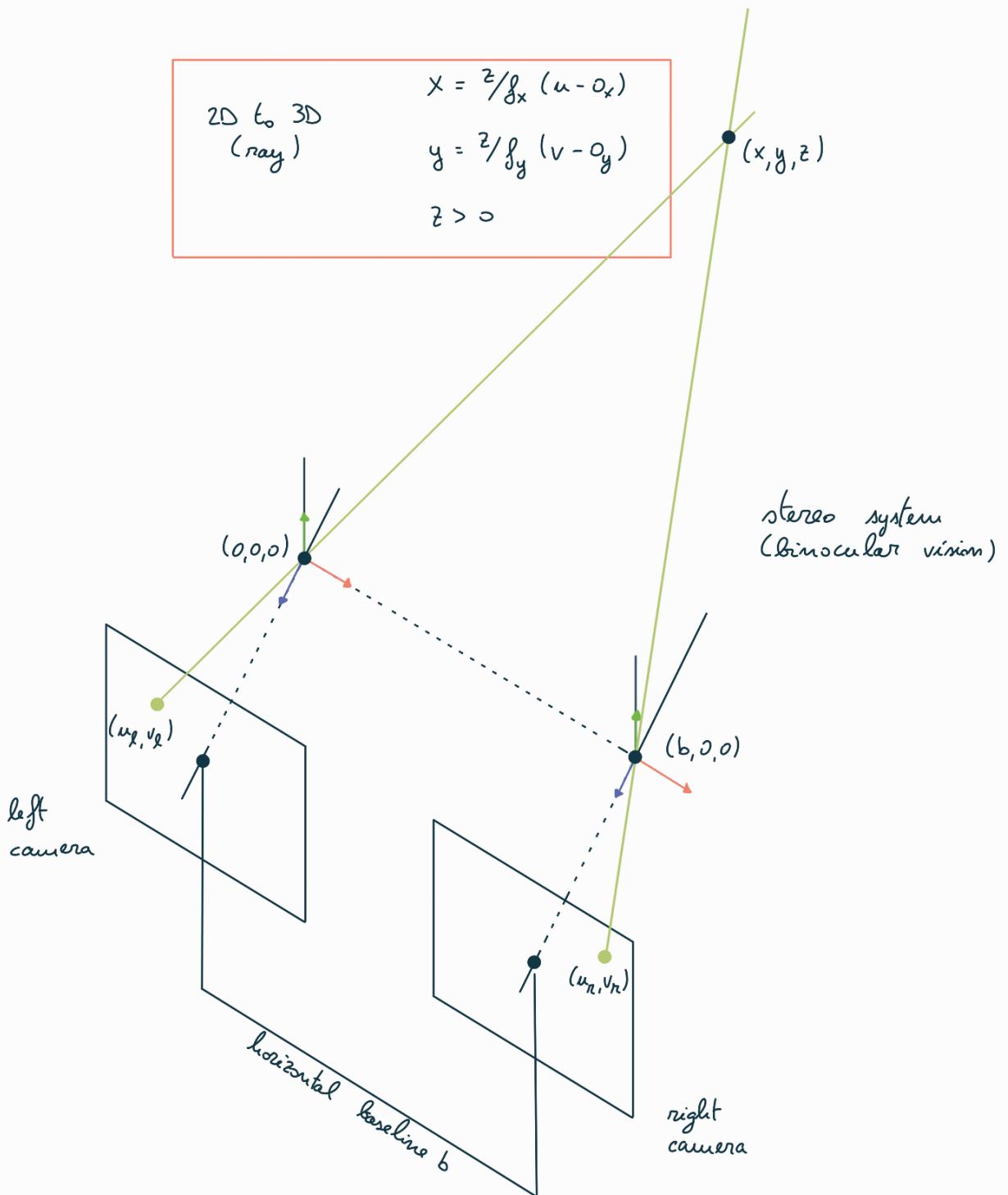


Simple stereo

We now have a calibrated camera. let's look to a simple method for recovering the 3-dimensional structure of a scene from two images.
 If we have a fully calibrated camera, we can say that each point in the image lies on a outgoing ray. We know the equation of these rays.

$$\boxed{\begin{array}{l} \text{3D to 2D} \\ \text{(point)} \end{array} \quad \begin{array}{l} u = f_x \frac{x_c}{z_c} + o_x \\ v = f_y \frac{y_c}{z_c} + o_y \end{array}}$$

$$\boxed{\begin{array}{l} \text{2D to 3D} \\ \text{(ray)} \end{array} \quad \begin{array}{l} x = z/f_x (u - o_x) \\ y = z/f_y (v - o_y) \\ z > 0 \end{array}}$$



If we have the two corresponding points we are left with 4 equations:

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

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

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

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

left point

right point

We have everything (f_x, f_y, o_x, o_y from the calibration procedure, u_l, v_l, u_r, v_r are supposed as given) except x, y, z .

Solving for x, y, z :

$$x = \frac{b(u_l - o_x)}{(u_l - u_r)}$$

$$y = \frac{f_x(v_l - o_y)}{f_y(u_l - u_r)}$$

if the point in the second image is not given we need to find it. This is called **stereo matching**.

$$\text{depth of the point } (z) = \frac{bf_x}{(u_l - u_r)}$$

where $u_l - u_r$ is called **disparity**. The disparity is inversely proportional to the depth and proportional to the baseline.

Uncalibrated stereo

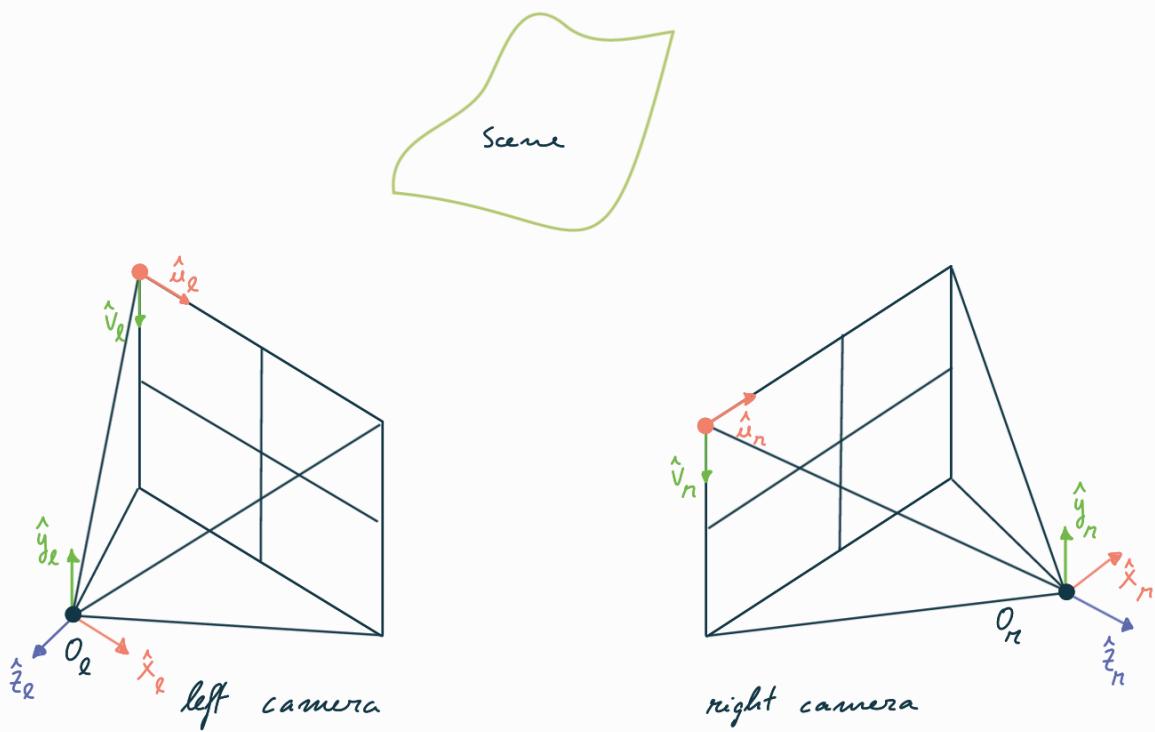
If you know the internal parameters of the two cameras it is possible to compute the translation and rotation of one camera with respect to the other one from two arbitrary views: we can then estimate the 3D structure of a static scene from two arbitrary views.

Often these days the internal parameters of an image are stored in its metadata; what we don't know is the translation and rotation among the points where two images of the same 3D scene are taken (two different cameras or same one).

To reconstruct the scene we need to formulate the relationship between left and right cameras. This is very concisely described in what's called **epipolar geometry**. Epipolar geometry tells us that points in the left and right image are related to each other through a sing 3x3 matrix called a **fundamental matrix**. If we can find the fundamental matrix we can find rotation and translation of one camera with respect to the other one.

We can develop a method to **estimate** the fundamental matrix F by using a small number of corresponding **points** in these two arbitrary views.

Once we have F we can find rotation and translation and the uncalibrated stereo is now fully calibrated. Now we can recover the structure of the scene.



- ① We assume that the camera matrix K (intrinsic parameters) is known for each camera
- ② Find a small number of reliable Corresponding Points between those two arbitrary views (e.g. SIFT). We will get a set of points:

$$\begin{matrix} (u_L^{(1)}, v_L^{(1)}) \\ \vdots \\ (u_L^{(m)}, v_L^{(m)}) \end{matrix}$$

$$\begin{matrix} (u_R^{(1)}, v_R^{(1)}) \\ \vdots \\ (u_R^{(m)}, v_R^{(m)}) \end{matrix}$$

Turns out we don't need much points, we need a minimum of 8.

③ Find Relative Camera Position t and Orientation R

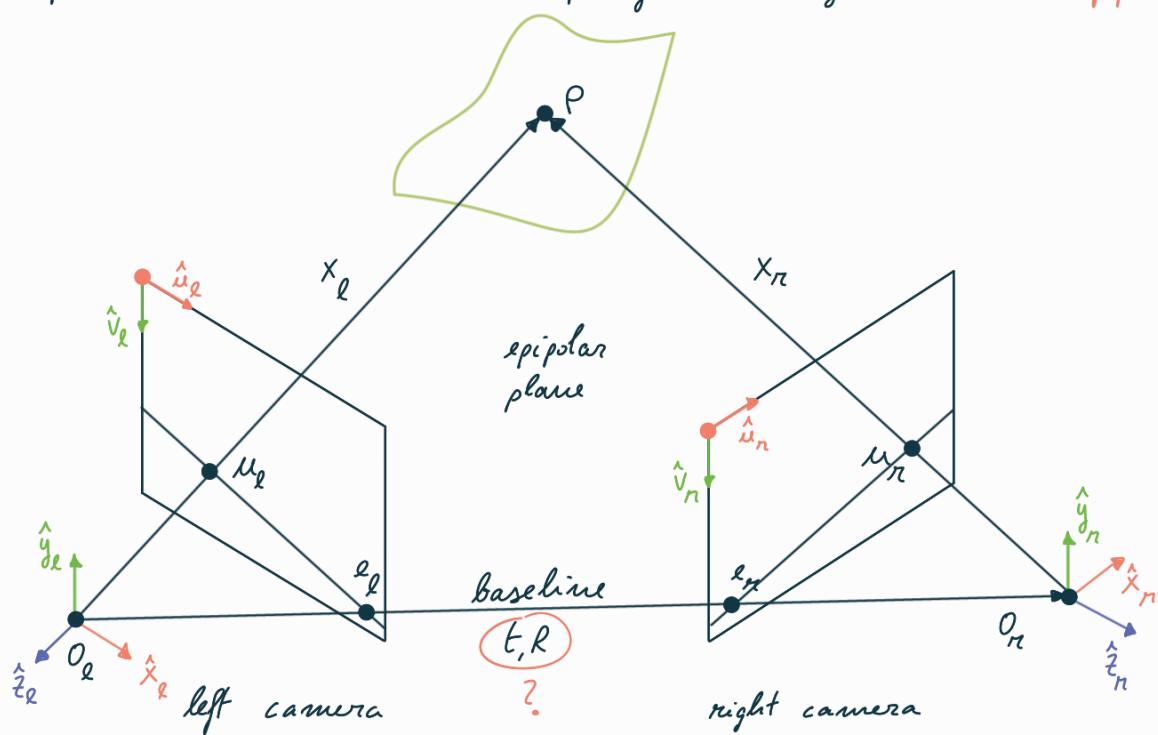
now the uncalibrated
stereo system is calibrated

④ Find Dense Correspondance (1D search)

⑤ Compute Depth using Triangulation

Epipolar geometry

Goal: find the relative position and orientation of one camera with respect to the other; this is the process of **calibrating an uncalibrated stereo system**. Relative position and orientation are completely described by what's called **epipolar geometry**.

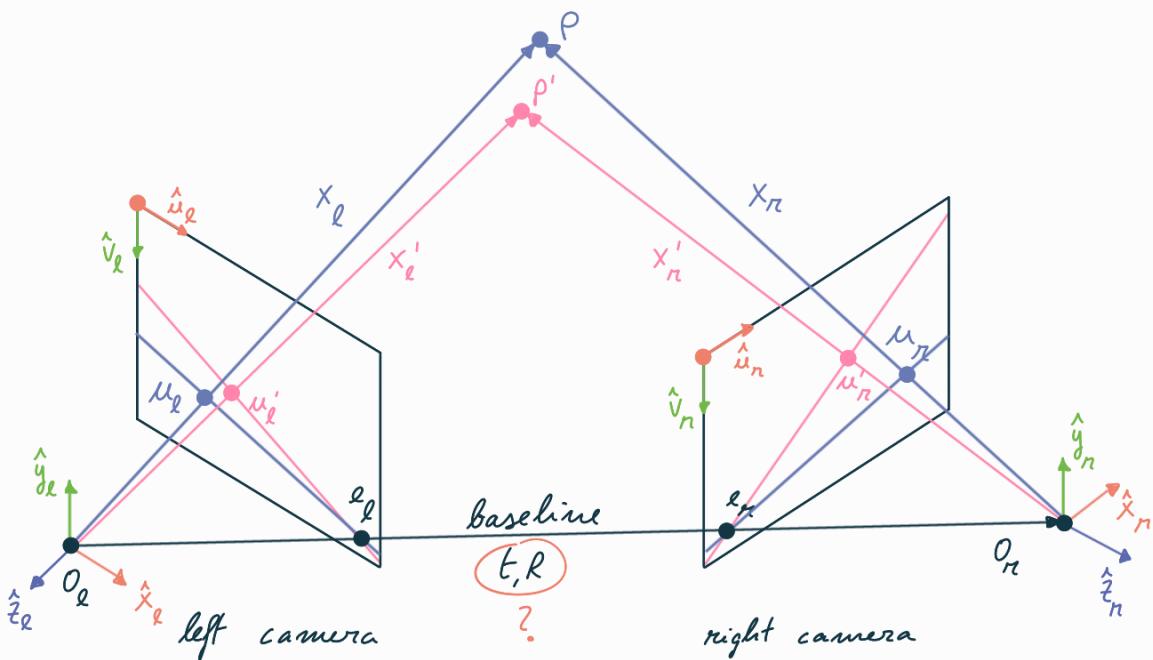


Baseline : line connecting the two camera centers

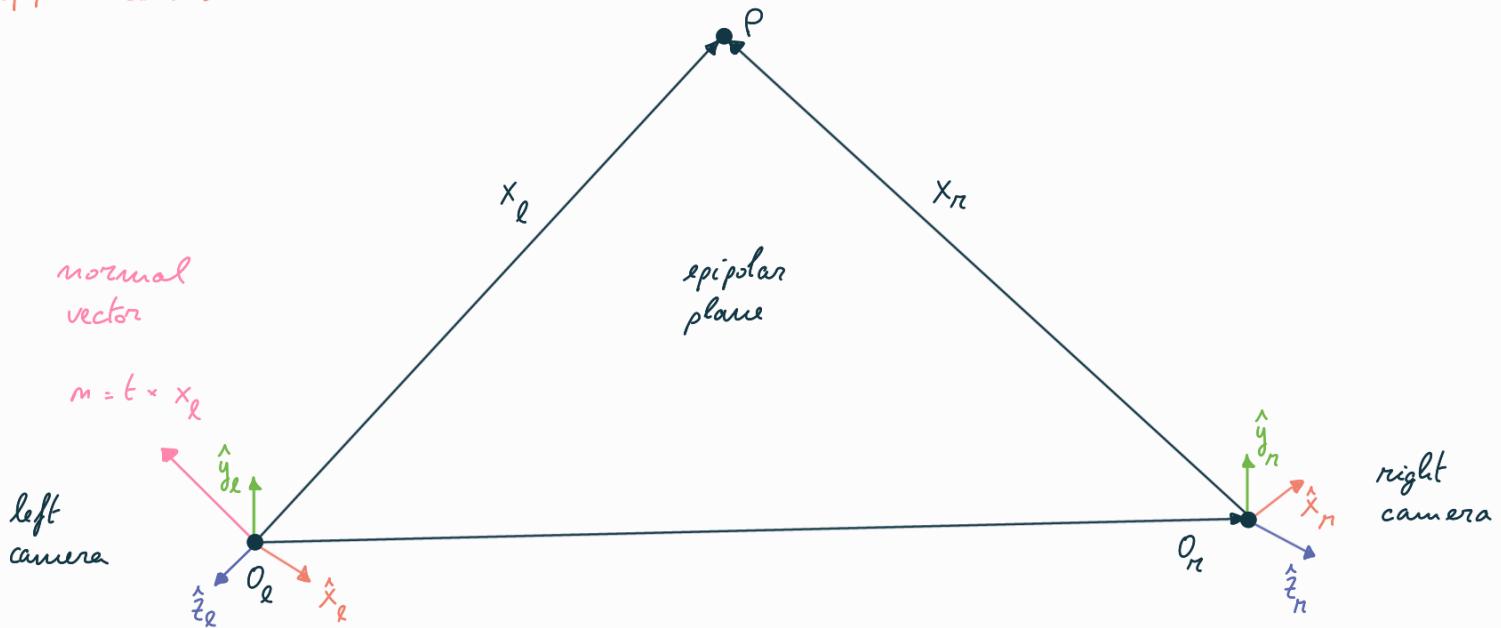
Epipole : image point of origin/pinhole of one camera as viewed by the other camera; every stereo system has a unique pair of epipoles. Equivalently, each epipole is the intersection point between the baseline and the image frame.

Epipolar plane : plane formed by camera origins (O_l and O_n), epipoles (e_l and e_n) and scene point P . Each point in the scene lies on a unique epipolar plane. There is a one parameter family, or a pencil, of epipolar planes.

Epipolar line : intersection of an epipolar plane with the image plane. All epipolar lines intersect at the epipole



Epipolar constraint



Vector normal to the epipolar plane : $m = t \times x_L$

The normal vector m should be perpendicular to x_L : $m \cdot x_L = 0$

Epipolar constraint : dot product of m and x_L (perpendicular vectors) is zero :

$$x_L \cdot (t \times x_L) = 0$$

writing it in vector form :

$$\begin{vmatrix} x_L & y_L & z_L \end{vmatrix} \begin{vmatrix} t_y z_L - t_z y_L \\ t_z x_L - t_x z_L \\ t_x y_L - t_y x_L \end{vmatrix} = 0 \quad (\text{cross product definition})$$

writing it in matrix form :

$$\boxed{\begin{vmatrix} x_L & y_L & z_L \end{vmatrix} \begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} x_L \\ y_L \\ z_L \end{vmatrix} = 0} \quad (\text{matrix vector form})$$

epipolar constraint

We have another constraint :

$t_{3 \times 1}$: Position of Right Camera in Left Camera's Frame

$R_{3 \times 3}$: Orientation of Left Camera in Right Camera's Frame

$$x_L = R x_R + t \rightarrow \boxed{\begin{vmatrix} x_L \\ y_L \\ z_L \end{vmatrix} = \begin{vmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{vmatrix} \begin{vmatrix} x_R \\ y_R \\ z_R \end{vmatrix} + \begin{vmatrix} t_x \\ t_y \\ t_z \end{vmatrix}}$$

relationship among x_L and x_R

Putting these two constraints together we have:

$$\begin{vmatrix} x_e & y_e & z_e \end{vmatrix} \begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} x_e \\ y_e \\ z_e \end{vmatrix} = 0$$

$$\begin{vmatrix} x_e \\ y_e \\ z_e \end{vmatrix} = \boxed{\begin{vmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{vmatrix} \begin{vmatrix} x_r \\ y_r \\ z_r \end{vmatrix} + \begin{vmatrix} t_x \\ t_y \\ t_z \end{vmatrix}}$$

$$\begin{vmatrix} x_e & y_e & z_e \end{vmatrix} \left(\begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{vmatrix} \begin{vmatrix} x_r \\ y_r \\ z_r \end{vmatrix} + \begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} t_x \\ t_y \\ t_z \end{vmatrix} \right) = 0$$

$$\begin{vmatrix} x_e & y_e & z_e \end{vmatrix} \boxed{\begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{vmatrix}} \begin{vmatrix} x_r \\ y_r \\ z_r \end{vmatrix} = 0$$

$$\begin{vmatrix} x_e & y_e & z_e \end{vmatrix} \boxed{\begin{vmatrix} l_{11} & l_{12} & l_{13} \\ l_{21} & l_{22} & l_{23} \\ l_{31} & l_{32} & l_{33} \end{vmatrix}} \begin{vmatrix} x_r \\ y_r \\ z_r \end{vmatrix} = 0$$

Essential Matrix E

The essential matrix expresses the relationship between x_e and x_r

$$E = T_x R$$

It is possible to decompose E back into T_x and R

$$\begin{vmatrix} l_{11} & l_{12} & l_{13} \\ l_{21} & l_{22} & l_{23} \\ l_{31} & l_{32} & l_{33} \end{vmatrix} = \begin{vmatrix} 0 & -t_z & t_y \\ t_z & 0 & -t_x \\ t_y & t_x & 0 \end{vmatrix} \begin{vmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \\ r_{31} & r_{32} & r_{33} \end{vmatrix}$$

Given that T_x is a Skew-Symmetric matrix ($a_{ij} = -a_{ji}$) and R is an Orthonormal matrix, it is possible to decompose E back into T_x and R from their product using Single Value Decomposition (SVD).

Once you have E you can compute T and R and calibrate the uncalibrated stereo.

$$\overbrace{x_e^T E x_r = 0}^{\text{essential matrix}} \quad \begin{array}{l} \text{3D position in left camera coordinates} \\ \text{3D position in right camera coordinates} \end{array}$$

We don't have x_e and x_r but we have their projections

We will try to use the projections of x_l and x_n to estimate the essential matrix.

Back to perspective projection equations

$$u_l = f_x \frac{x_l}{z_l} + o_x^{(l)} \rightarrow z_l u_l = f_x x_l + z_l o_x^{(l)}$$

$$v_l = f_y \frac{x_l}{z_l} + o_y^{(l)} \rightarrow z_l v_l = f_y x_l + z_l o_y^{(l)}$$

Representing in matrix form:

$$\begin{vmatrix} u_l \\ v_l \\ 1 \end{vmatrix} = \begin{vmatrix} z_l u_l \\ z_l v_l \\ z_l \end{vmatrix} = \begin{vmatrix} f_x^{(l)} x_l + z_l o_x^{(l)} \\ f_y^{(l)} x_l + z_l o_y^{(l)} \\ z_l \end{vmatrix} = \begin{vmatrix} f_x & 0 & o_x^{(l)} \\ 0 & f_y & o_y^{(l)} \\ 0 & 0 & 1 \end{vmatrix} \begin{vmatrix} x_l \\ y_l \\ z_l \end{vmatrix}$$

we have the internal parameters
Known camera matrix K_l

left camera

$$\begin{vmatrix} u_l \\ v_l \\ 1 \end{vmatrix} = \begin{vmatrix} f_x^{(l)} & 0 & o_x^{(l)} \\ 0 & f_y^{(l)} & o_y^{(l)} \\ 0 & 0 & 1 \end{vmatrix} \begin{vmatrix} x_l \\ y_l \\ z_l \end{vmatrix}$$

right camera

$$\begin{vmatrix} u_n \\ v_n \\ 1 \end{vmatrix} = \begin{vmatrix} f_x^{(n)} & 0 & o_x^{(n)} \\ 0 & f_y^{(n)} & o_y^{(n)} \\ 0 & 0 & 1 \end{vmatrix} \begin{vmatrix} x_n \\ y_n \\ z_n \end{vmatrix}$$

$$x_l^T = \begin{vmatrix} u_l & v_l & 1 \end{vmatrix} z_l K_l^{-1 T}$$

$$x_n = k_n^{-1} z_n \begin{vmatrix} u_n \\ v_n \\ 1 \end{vmatrix}$$

$$x_l^T E x_n = 0$$

$$\begin{vmatrix} u_l & v_l & 1 \end{vmatrix} z_l K_l^{-1 T} E k_n^{-1} z_n \begin{vmatrix} u_n \\ v_n \\ 1 \end{vmatrix} = 0$$

We don't have z_l . We know z_l (depth of the 3D point in the scene) $\neq 0$: if it is 0, it is on the same plane of the camera center

$$\begin{vmatrix} u_l & v_l & 1 \end{vmatrix} \underset{3 \times 3}{\cancel{z_l}} K_l^{-1 T} E \underset{3 \times 3}{\cancel{k_n^{-1}}} z_n \begin{vmatrix} u_n \\ v_n \\ 1 \end{vmatrix} = 0$$

$K_l^{-1 T} E k_n^{-1}$ is a 3×3 matrix and it's called fundamental matrix F .

The fundamental matrix is an expression of the epipolar constraints.

$$\begin{vmatrix} x_l & y_l & 1 \end{vmatrix} \begin{vmatrix} f_{14} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{vmatrix} \begin{vmatrix} u_n \\ v_n \\ 1 \end{vmatrix} = 0$$

Once we have F :

$$\bar{E} = K_L^T F K_R$$

Once we have E :

$$E = T_x R \text{ (with Single Value Decomposition)}$$

We only need to find a single 3x3 matrix, the fundamental matrix, to calibrate the uncalibrated stereo system

Stereo Calibration Procedure

The epipolar constraint tells us that there is a single 3×3 matrix, the fundamental matrix, that we need to find to calibrate our uncalibrated stereo system. Our goal is to estimate the fundamental matrix.

- Find a set of corresponding features in left and right image (e.g. using SIFT)

$$\begin{array}{l} (u_l^{(1)}, v_l^{(1)}) \\ \vdots \\ (u_l^{(m)}, v_l^{(m)}) \end{array} \quad \begin{array}{l} (u_r^{(1)}, v_r^{(1)}) \\ \vdots \\ (u_r^{(m)}, v_r^{(m)}) \end{array}$$

- For each correspondence i , write the epipolar constraint

$$\left| \begin{array}{ccc} x_l^{(i)} & y_l^{(i)} & 1 \end{array} \right| \left| \begin{array}{ccc} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{array} \right| \left| \begin{array}{c} u_r^{(i)} \\ v_r^{(i)} \\ 1 \end{array} \right| = 0$$

Known	Unknown	Known
-------	---------	-------

Expanding the expression we get a linear equation :

$$(f_{11} u_r^{(i)} + f_{12} v_r^{(i)} + f_{13}) u_l^{(i)} + (f_{21} u_r^{(i)} + f_{22} v_r^{(i)} + f_{23}) v_l^{(i)} + (f_{31} u_r^{(i)} + f_{32} v_r^{(i)} + f_{33}) = 0$$

we get one such linear equation for each same point.

- Stack them to form a linear system

$$\left| \begin{array}{ccccccccc} u_r^{(1)} u_l^{(1)} & v_r^{(1)} u_l^{(1)} & u_l^{(1)} & u_r^{(1)} v_l^{(1)} & v_r^{(1)} v_l^{(1)} & v_l^{(1)} & u_r^{(1)} & v_r^{(1)} & 1 \\ \vdots & \vdots \\ u_r^{(i)} u_l^{(i)} & v_r^{(i)} u_l^{(i)} & u_l^{(i)} & u_r^{(i)} v_l^{(i)} & v_r^{(i)} v_l^{(i)} & v_l^{(i)} & u_r^{(i)} & v_r^{(i)} & 1 \\ \vdots & \vdots \\ u_r^{(m)} u_l^{(m)} & v_r^{(m)} u_l^{(m)} & u_l^{(m)} & u_r^{(m)} v_l^{(m)} & v_r^{(m)} v_l^{(m)} & v_l^{(m)} & u_r^{(m)} & v_r^{(m)} & 1 \end{array} \right| = \left| \begin{array}{c} f_{11} \\ f_{21} \\ f_{31} \\ f_{12} \\ f_{22} \\ f_{32} \\ f_{13} \\ f_{23} \\ f_{33} \end{array} \right| = 0$$

Known	Unknown
-------	---------

$$A f = 0$$

Fundamental matrix acts on homogeneous coordinates

$$\left| \begin{array}{ccc} x_l & y_l & 1 \end{array} \right| \left| \begin{array}{ccc} f_{11} & f_{12} & f_{13} \\ f_{21} & f_{22} & f_{23} \\ f_{31} & f_{32} & f_{33} \end{array} \right| \left| \begin{array}{c} u_r \\ v_r \\ 1 \end{array} \right| = 0 = \left| \begin{array}{ccc} x_l & y_l & 1 \end{array} \right| \left| \begin{array}{ccc} K f_{11} & K f_{12} & K f_{13} \\ K f_{21} & K f_{22} & K f_{23} \\ K f_{31} & K f_{32} & K f_{33} \end{array} \right| \left| \begin{array}{c} u_r \\ v_r \\ 1 \end{array} \right|$$

Fundamental matrix is only defined up to a scale ($F = kF$).
We can set it to some arbitrary scale

$$\|f\|^2 = 1 \quad (\text{for convenience})$$

- ④ Find least squares solution for fundamental matrix F .

The 8 point algorithm

We want Af as close as possible to 0 and $\|f\|^2 = 1$:

$$\min_f \|Af\|^2 \text{ such that } \|f\|^2 = 1$$

constrained linear least squares problem

- ⑤ Rearrange f to get matrix F .

- ⑥ Compute essential matrix from fundamental matrix

$$E = k_l^T F k_r$$

- ⑦ Extract R and t from E using Singular Value Decomposition

$$E = T_x R$$

Fully calibrated stereo system