

## Structure from motion

Compute the 3D scene structure and camera motion from a sequence of frames.

INPUT: video of a scene, sequence of frames.

OUTPUT: 3D scene structure

① Detect feature points: corners, SIFT points, ...

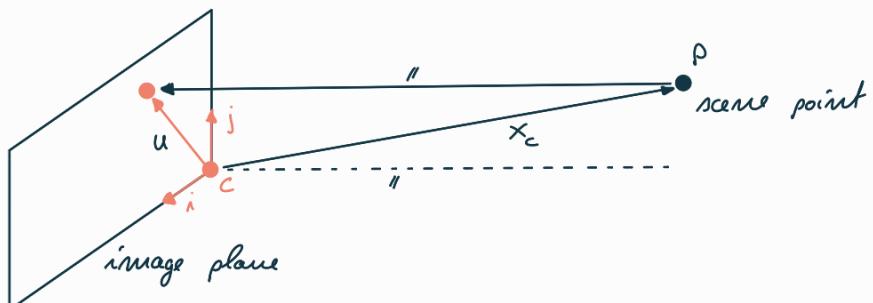
② Track the features through the entire video: template matching, optical flow, comparing SIFT descriptors, ...  
→ set of tracked features  
N points in each frame 1, ..., F.

INPUT: set of corresponding image points (2D):  $(u_{f,p}, v_{f,p}) \quad \forall p \in 1, \dots, N$

OUTPUT: find scene points (3D)

A set is a sequence of points & feature (in time)

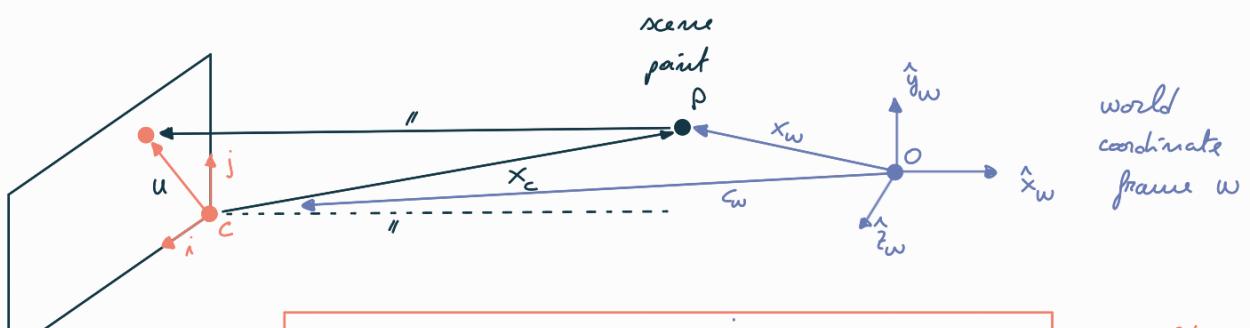
Assume we are using an **orthographic camera**: distance between points in the scene is small compared to the distance of the scene from the camera: the magnification of the camera is constant for all points in the scene and between scenes.



$$u = (u, v)$$

$$\begin{aligned} u &= i \cdot x_c = i^T x_c \\ v &= j \cdot x_c = j^T x_c \end{aligned}$$

camera coordinate frame



$$u = (u, v)$$

$$\begin{aligned} u &= i \cdot x_c = i^T x_c = i^T (x_w - c_w) = i^T (p - c) \\ v &= j \cdot x_c = j^T x_c = j^T (x_w - c_w) = j^T (p - c) \end{aligned}$$

world coordinate frame

Given corresponding image points (2D)  $(u_{f,p}, v_{f,p})$   
 Find scene points  $\{P_p\}$ .

Camera positions  $\{c_f\}$  and camera orientations  $\{\hat{i}_f, \hat{j}_f\}$  are unknown.

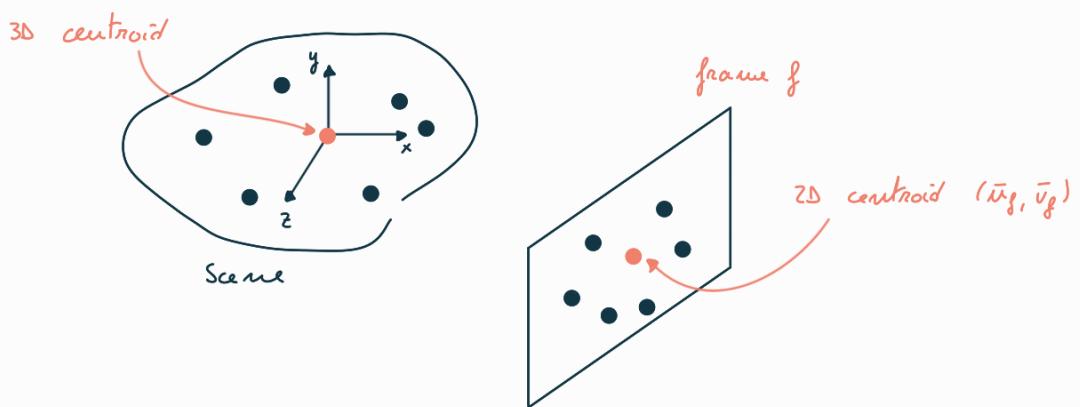
Image of point  $P_p$  in camera frame  $f$ : 
$$\begin{cases} u_{f,p} = \hat{i}_f^T (P_p - c_f) \\ v_{f,p} = \hat{j}_f^T (P_p - c_f) \end{cases}$$

Known      Unknown

We remove  $c_f$  from the equations to simplify the problem using the **centering trick**.

**Centering trick**: assuming origin of the world at centroid of scene points :

$$\frac{1}{N} \sum_{p=1}^N P_p = \bar{P} = 0$$



Centroid  $(\bar{u}_f, \bar{v}_f)$  of the image points in frame  $f$ :

$$\bar{u}_f = \frac{1}{N} \sum_{p=1}^N u_{f,p} = \frac{1}{N} \sum_{p=1}^N \hat{i}_f^T (P_p - c_f)$$

$$= \frac{1}{N} \hat{i}_f^T \sum_{p=1}^N P_p - \frac{1}{N} \sum_{p=1}^N \hat{i}_f^T c_f$$

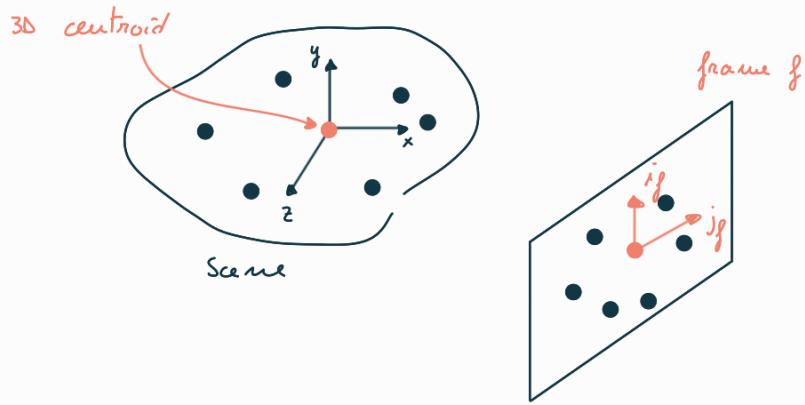
we set the origin  
on the centroid      constant

$$\bar{u}_f = -\hat{i}_f^T c_f$$

the  $\bar{v}_f$  component can be computed  
the same way

$$\bar{v}_f = -\hat{j}_f^T c_f$$

The centroid of the feature points  
is not a function of the locations  
of the scene points



Shift camera origin to the centroid  $(\bar{u}_f, \bar{v}_f)$

New image points :

Now we have the image coordinates of same points but we no longer have the camera center in the expression

$$\tilde{u}_{f,p} = u_{f,p} - \bar{u}_f = i_f^T (p_p - c_f) - i_f^T c_f = i_f^T p_p$$

$$\tilde{v}_{f,p} = v_{f,p} - \bar{v}_f = j_f^T (p_p - c_f) - j_f^T c_f = j_f^T p_p$$

Now for each point in a particular frame, we get this expression:

$$\begin{aligned} \tilde{u}_{f,p} &= i_f^T p_p \\ \tilde{v}_{f,p} &= j_f^T p_p \end{aligned} \quad \longrightarrow \quad \begin{vmatrix} \tilde{u}_{f,p} \\ \tilde{v}_{f,p} \end{vmatrix} = \begin{vmatrix} i_f^T \\ j_f^T \end{vmatrix} p_p$$

We have  $N$  points and  $F$  frames. We can organize them into a single matrix:

$$\begin{array}{c}
 \text{point 1} \quad \text{point 2} \quad \dots \quad \text{point } N \\
 \hline
 \text{image 1} \quad \tilde{u}_{1,1} \quad \tilde{u}_{1,2} \quad \dots \quad \tilde{u}_{1,N} \quad | \quad i_1^T \quad | \quad \text{point 1} \\
 \text{image 2} \quad \tilde{u}_{2,1} \quad \tilde{u}_{2,2} \quad \dots \quad \tilde{u}_{2,N} \quad | \quad i_2^T \quad | \quad \text{point 2} \\
 \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad | \quad \vdots \quad | \quad \vdots \\
 \text{image } F \quad \tilde{u}_{F,1} \quad \tilde{u}_{F,2} \quad \dots \quad \tilde{u}_{F,N} \quad | \quad i_F^T \quad | \quad \text{point } N \\
 \hline
 \text{image 1} \quad \tilde{v}_{1,1} \quad \tilde{v}_{1,2} \quad \dots \quad \tilde{v}_{1,N} \quad | \quad j_1^T \quad | \quad p_1 \\
 \text{image 2} \quad \tilde{v}_{2,1} \quad \tilde{v}_{2,2} \quad \dots \quad \tilde{v}_{2,N} \quad | \quad j_2^T \quad | \quad p_2 \\
 \vdots \quad \vdots \quad \vdots \quad \vdots \quad \vdots \quad | \quad \vdots \quad | \quad \vdots \\
 \text{image } F \quad \tilde{v}_{F,1} \quad \tilde{v}_{F,2} \quad \dots \quad \tilde{v}_{F,N} \quad | \quad j_F^T \quad | \quad p_N
 \end{array}$$

$W_{2F \times N}$

$M_{2F \times 1}$

$S_{3 \times N}$

Centroid-Subtracted  
Feature points  
(known)

Camera  
motion  
(unknown)

Scene  
structure  
(unknown)

Observation matrix

Can we find  $M$  and  $S$  from  $W$ ?

Yes because the observation matrix  $W$  has a special property: it has a very low rank.

linear independence of vectors

A set of vectors  $\{v_1, v_2, \dots, v_m\}$  is said to be linearly independent if no vector can be represented as a weighted linear sum of the others.

Column rank

The column rank of a matrix is the number of linearly independent columns of the matrix

Row rank

The row rank of a matrix is the number of linearly independent rows of the matrix

Rank

For any matrix:

$$\text{columnRank}(A) = \text{rowRank}(A) = \text{rank}(A)$$
$$\text{rank}(A) \leq \min(m, n)$$

$$m \left| \begin{array}{c} \\ A \\ m \end{array} \right| \begin{array}{l} \text{columnRank}(A) \leq m \\ \text{rowRank}(A) \leq m \end{array}$$

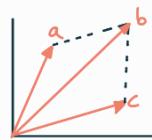
Rank is the dimensionality of the space spanned by column or row vectors of the matrix.

$$A = \begin{vmatrix} a_{11} & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \end{vmatrix} = \begin{vmatrix} a & b & c \end{vmatrix}$$

$$\text{rank}(A) = 1$$



$$\text{rank}(A) = 2$$



$$\text{rank}(A) = 3$$

in a 3-dimensional space  
we need three vectors to  
express any arbitrary vector

The rank has some properties :

- $\text{rank}(A^T) = \text{rank}(A)$
- $\text{rank}(A_{m \times m} B_{m \times p}) = \min(\text{rank}(A_{m \times m}), \text{rank}(B_{m \times p})) \leq \min(m, m, p)$
- $\text{rank}(AA^T) = \text{rank}(A^TA) = \text{rank}(A^T) = \text{rank}(A)$
- $A_{m \times m}$  is invertible iff  $\text{rank}(A_{m \times m}) = m$

Back to the observation matrix

$$\begin{array}{c}
 \begin{array}{ccc}
 \text{point 1} & \text{point 2} & \text{point } N
 \end{array} \\
 \begin{array}{c|ccc|c}
 \text{image 1} & \tilde{u}_{1,1} & \tilde{u}_{1,2} & \cdots & \tilde{u}_{1,N} \\
 \text{image 2} & \tilde{u}_{2,1} & \tilde{u}_{2,2} & \cdots & \tilde{u}_{2,N} \\
 \vdots & \vdots & \vdots & & \vdots \\
 \text{image } F & \tilde{u}_{F,1} & \tilde{u}_{F,2} & \cdots & \tilde{u}_{F,N}
 \end{array} = \begin{array}{c|ccc|c}
 \begin{array}{c}
 i_1^T \\
 i_2^T \\
 \vdots \\
 i_F^T
 \end{array} & \text{point 1} & \text{point 2} & \cdots & \text{point } N
 \end{array} \\
 \begin{array}{c|ccc|c}
 \text{image 1} & \tilde{v}_{1,1} & \tilde{v}_{1,2} & \cdots & \tilde{v}_{1,N} \\
 \text{image 2} & \tilde{v}_{2,1} & \tilde{v}_{2,2} & \cdots & \tilde{v}_{2,N} \\
 \vdots & \vdots & \vdots & & \vdots \\
 \text{image } F & \tilde{v}_{F,1} & \tilde{v}_{F,2} & \cdots & \tilde{v}_{F,N}
 \end{array} = \begin{array}{c|ccc|c}
 \begin{array}{c}
 j_1^T \\
 j_2^T \\
 \vdots \\
 j_F^T
 \end{array} & p_1 & p_2 & \cdots & p_N
 \end{array}
 \end{array}$$

$W_{2F \times N}$        $M_{2F \times 2}$        $S_{3 \times N}$

Centroid-Subtracted  
 Feature points  
 (known)

Camera  
 motion  
 (unknown)

Scene  
 structure  
 (unknown)

Observation matrix

$$W = M \times S$$

$$2F \times N \quad 2F \times 3 \quad 3 \times N$$

$$\text{rank}(MS) \leq \text{rank}(M) \longrightarrow \text{rank}(RS) \leq \min(3, 2F)$$

$$\text{rank}(MS) \leq \text{rank}(S) \longrightarrow \text{rank}(RS) \leq \min(3, N)$$

$$\text{rank}(W) = \text{rank}(MS) \leq \min(3, 2F, N) = 3$$

## Singular Value Decomposition (SVD)

For any matrix  $A$  there exists a factorization:

$$A_{m \times m} = U_{m \times m} \cdot \Sigma_{m \times m} \cdot V^T_{m \times m}$$

where  $U$  and  $V$  are orthonormal and  $\Sigma$  is diagonal.

$$\Sigma_{m \times m} = \begin{vmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{vmatrix} \quad \sigma_1 \leq \sigma_2 \leq \cdots \leq \sigma_n \text{ are the singular values}$$

If  $\text{rank}(A) = r$  then  $A$  has  $r$  non-zero singular values.

## Tomari - Kanade Factorization Method

We can "factorize"  $W$  into  $U$  and  $S$ .

We will use a technique called Singular Value Decomposition (SVD)

Using SVD:  $W = U \Sigma V^T$

$$= \begin{vmatrix} U & \begin{vmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \sigma_n \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 0 \end{vmatrix} & V^T \end{vmatrix}_{2F \times 2F \quad 2F \times N \quad N \times N}$$

since  $\text{rank}(W) \leq 3$ ,  $\text{rank}(\Sigma) \leq 3$  ( $\Sigma$  has at most 3 non-zero singular values).

$$= \begin{vmatrix} \begin{matrix} 3 \\ 2F-3 \end{matrix} & U_1 & U_2 & \begin{matrix} \sigma_1 & 0 & 0 \\ 0 & \sigma_2 & 0 \\ 0 & 0 & \sigma_3 \end{matrix} & \begin{matrix} \cdots & 0 \\ 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \cdots & \vdots \\ 0 & 0 & 0 & \cdots & 0 \end{matrix} & \begin{matrix} V_1^T \\ V_2^T \end{matrix} & 3 \\ N-3 & & & & & & \end{vmatrix}_{2F \times 2F \quad 2F \times N \quad N \times N}$$

$$W = U_1 \begin{matrix} \sum \\ (2F \times 3) \end{matrix} \begin{matrix} V_1^T \\ (3 \times N) \end{matrix}$$

## Factorization (Finding $M, S$ )

$$W = U_1 \Sigma_1 V_1^T$$

$$W = U_1 (\Sigma_1)^{\frac{1}{2}} (\Sigma_1)^{\frac{1}{2}} V_1^T \quad \text{Two possible factors for } W$$

$$W = U_1 (\Sigma_1)^{\frac{1}{2}} Q Q^{-1} (\Sigma_1)^{\frac{1}{2}} V_1^T \quad \text{Also two possible valid factors for } W$$

now we can use the orthonormality of  $M$

$$M = \begin{vmatrix} i_1^T \\ i_2^T \\ \vdots \\ i_F^T \\ j_1^T \\ j_2^T \\ \vdots \\ j_F^T \end{vmatrix} = U_1 (\Sigma_1)^{\frac{1}{2}} Q = \begin{vmatrix} \hat{i}_1^T \\ \vdots \\ \hat{i}_F^T \\ \hat{j}_1^T \\ \vdots \\ \hat{j}_F^T \end{vmatrix} Q = \begin{vmatrix} \hat{i}_1^T Q \\ \vdots \\ \hat{i}_F^T Q \\ \hat{j}_1^T Q \\ \vdots \\ \hat{j}_F^T Q \end{vmatrix}$$

computed    unknown

orthonormality constraints :

$$\left. \begin{array}{l} i_f \cdot i_f = \hat{i}_f^T i_f = 1 \\ j_f \cdot j_f = \hat{j}_f^T j_f = 1 \\ i_f \cdot j_f = \hat{i}_f^T j_f = 0 \end{array} \right\} \longrightarrow$$

for each frame  
we get these  
3 equations

$$\left\{ \begin{array}{l} \hat{i}_f^T Q \bar{Q}^{-1} \hat{i}_f = 1 \\ \hat{j}_f^T Q \bar{Q}^{-1} \hat{j}_f = 1 \\ \hat{i}_f^T Q \bar{Q}^{-1} \hat{j}_f = 0 \end{array} \right.$$

$Q$  is  $3 \times 3$ , thus it has 3 variables  
we have  $3F$  of these constraints  
 $Q$  can be solved with 3 or more images ( $F \geq 3$ ) using Newton's method.

After finding  $Q$ :

$$M = U_1 (\Sigma_1)^{\frac{1}{2}} Q \quad \text{camera motion}$$

$$S = Q^{-1} (\Sigma_1)^{\frac{1}{2}} V_1^T \quad \text{scene structure}$$

Summary : orthographic structure from motion

- ① Detect and track feature points
- ② Create the centroid-subtracted matrix  $W$  of corresponding feature points
- ③ Compute SVD of  $W$  and enforce rank constraint

$$W = U \Sigma V^T = \begin{matrix} U_1 \\ (2F \times 3) \end{matrix} \quad \begin{matrix} \Sigma_1 \\ (3 \times 3) \end{matrix} \quad \begin{matrix} V_1^T \\ (3 \times N) \end{matrix}$$

- ④ Set  $M = U_1 (\Sigma_1)^{\frac{1}{2}} Q$  and  $S = Q^{-1} (\Sigma_1)^{\frac{1}{2}} V_1^T$
- ⑤ Find  $\Omega$  by enforcing the orthonormality constraint