



ICVSS 2013

*Calabria ~ 14-20 July*

International Computer Vision Summer School



The devil is in the details:  
anatomy of a structure-from-motion  
pipeline

Andrea Fusiello (andrea.fusiello@uniud.it)  
University of Udine

# Not all the details are equally important

- Science, after Leibniz, is about sorting out what is **necessary** from what is **contingent**
  - Contingent: it is A, but it could have been B
  - Necessary: there is a reason why it had to be A
- We should be aware of this distinction in our work, and be consistent with it when we write scientific papers (the speaker is no exception)
- Therefore, I will try to point out what is necessary in the flood of information that will follow.

# 3D modeling from images

- ▣ The holy grail of computer vision for ~20 years
- ▣ Input: pictures (nothing else)
- ▣ Output: a 3D model (points, lines, surfaces)



# Why image-based modeling?

- ❑ The process of modeling is cumbersome
- ❑ Requires well trained personnel
- ❑ Image-based modeling is simple (to use)
- ❑ Consider the the “Web 2.0” revolution: everyone can be a (3D) contents creator.

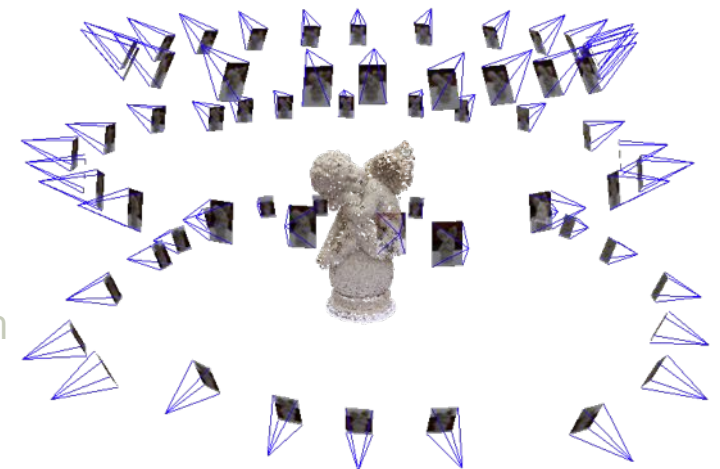
# A short introduction

- A camera is modeled by a  $3 \times 4$  matrix  $P$
- Given  $P$  and corresponding points, the 3D structure can be reconstructed by triangulation



# A short introduction

- ▣ Problems:
  - ▣ **Matching:** compute corresponding points
  - ▣ **Structure and Motion:** recover camera matrices
  - ▣ **Bridging the semantic gap:** upgrade from cloud of points to a high-level model



# Early works



Debevec, *et al* - SIGGRAPH 1996



Schaffalitzky and Zisserman - ECCV 2002

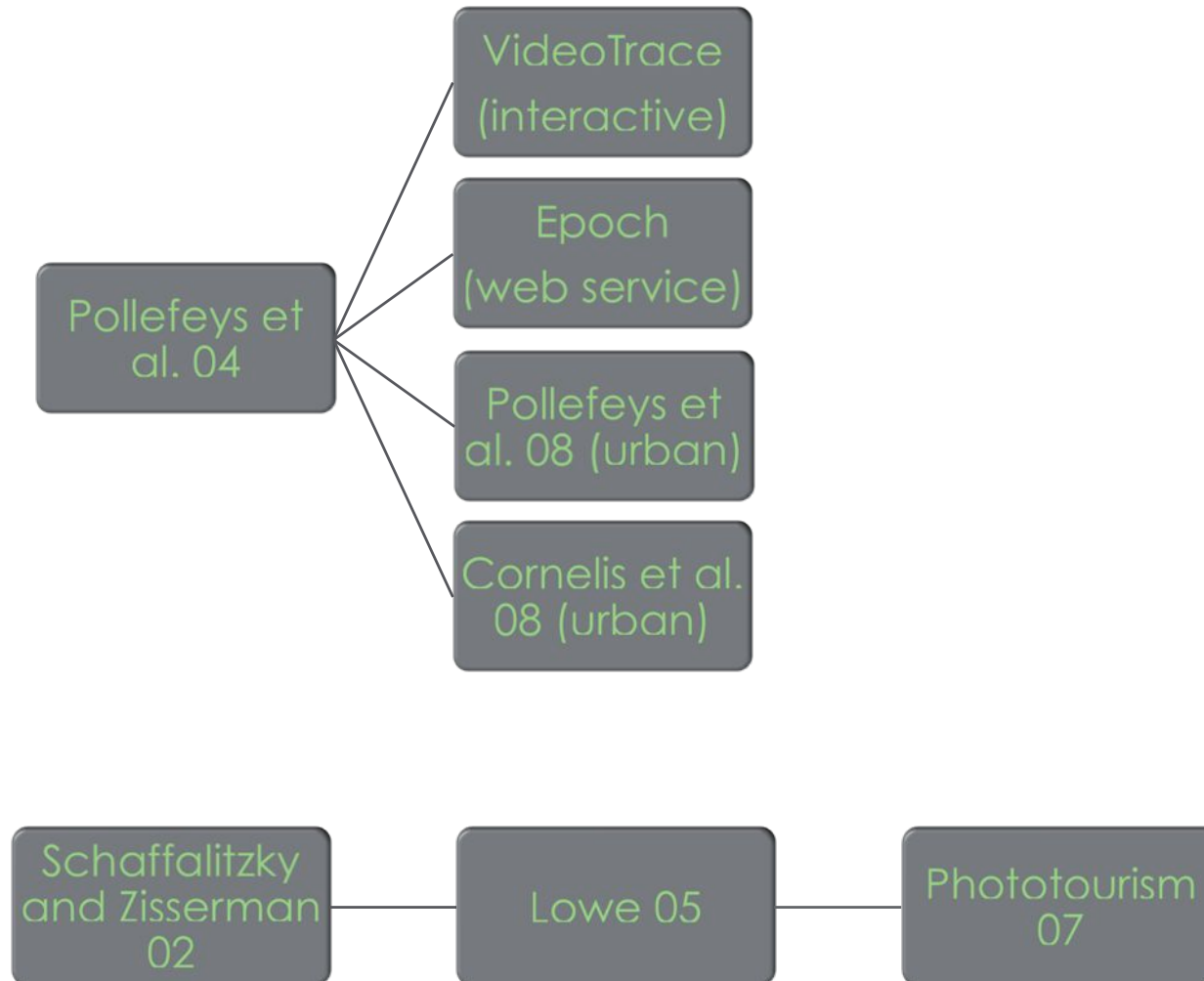


Pollefeys *et al.* - IJCV 2004



Brown and Lowe - 3DIM 2005

# A short genealogy





## Some references

- Debevec, P. E., Taylor, C. J., and Malik, J. 1996. Modeling and rendering architecture from photographs: a hybrid geometry- and image-based approach. In Proceedings of the 23rd Annual Conference on Computer Graphics and interactive Techniques SIGGRAPH '96. ACM, New York, NY, 11-20.
- Matthew Brown and David G. Lowe, "Unsupervised 3D object recognition and reconstruction in unordered datasets," International Conference on 3-D Digital Imaging and Modeling (3DIM 2005), Ottawa, Canada (June 2005).
- Frederik Schaffalitzky, Andrew Zisserman: Multi-view Matching for Unordered Image Sets, or "How Do I Organize My Holiday Snaps?". ECCV (1) 2002: 414-43.
- M. Pollefeys, L. Van Gool, M. Vergauwen, F. Verbiest, K. Cornelis, J. Tops, R. Koch, Visual modeling with a hand-held camera, International Journal of Computer Vision 59(3), 207-232, 2004
- Maarten Vergauwen and Luc Van Gool, "Web-Based 3D Reconstruction Service", Machine Vision Applications, 17, pp. 411-426, 2006.
- A. van den Hengel, A. Dick, T. Thormählen, B. Ward, and P. H. S. Torr . VideoTrace: Rapid interactive scene modelling from video. ACM Transactions on Graphics, 26(3), Article No. 86, July 2007.
- Noah Snavely, Steven M. Seitz, Richard Szeliski, "Photo tourism: Exploring photo collections in 3D," ACM Transactions on Graphics (SIGGRAPH Proceedings), 25(3), 2006, 835-846.
- Sameer Agarwal, Yasutaka Furukawa, Noah Snavely, Ian Simon, Brian Curless, Steven M. Seitz and Richard Szeliski Building Rome in a Day. Communications of the ACM, Vol. 54, No. 10, Pages 105-112, October 2011.



## Background theory

Where we fill our box with the (geometric) tools

## 2 Projective Geometry

The physical space is the Euclidean 3-D space  $\mathbb{E}^3$ , a real 3-dimensional affine space endowed with the inner product.

Our ambient space is the projective 3-D space  $\mathbb{P}^3$ , obtained by completing  $\mathbb{E}^3$  with a projective plane, known as plane at infinity  $\Pi_\infty$ . In this ideal plane lie the intersections of the planes parallel in  $\mathbb{E}^3$ .

The projective (or homogeneous) coordinates of a point in  $\mathbb{P}^3$  are 4-tuples defined up to a scale factor. We write

$$\mathbf{M} \simeq (x, y, z, t) \quad (1)$$

where  $\simeq$  indicates equality to within a multiplicative factor.

The affine points are those of  $\mathbb{P}^3$  which do not belong to  $\Pi_\infty$ . Their projective coordinates are of the form  $(x, y, z, 1)$ , where  $(x, y, z)$  are the usual Cartesian coordinates.

$\Pi_\infty$  is defined by its equation  $t = 0$ .

The linear transformations of a projective space into itself are called collineations or homographies. Any collineation of  $\mathbb{P}^3$  is represented by a generic  $4 \times 4$  invertible matrix.

Affine transformations are the subgroup of collineations of  $\mathbb{P}^3$  that preserves the plane at infinity (i.e., parallelism).

Similarity transformations are the subgroup of affine transformations that leave invariant a very special curve, the *absolute conic*, which is in the plane at infinity and whose equation is:

$$x^2 + y^2 + z^2 = 0 = t \quad (2)$$

Similarity transformations preserves the angles.

### 3 Pin-hole Camera Geometry

The pin-hole camera is described by its *optical centre*  $C$  (also known as *camera projection centre*) and the *image plane*.

The distance of the image plane from  $C$  is the *focal length*  $f$ .

The line from the camera centre perpendicular to the image plane is called the *principal axis* or *optical axis* of the camera.

The plane parallel to the image plane containing the optical centre is called the *principal plane* or *focal plane* of the camera.

The relationship between the 3-D coordinates of a scene point and the coordinates of its projection onto the image plane is described by the *central* or *perspective projection*.

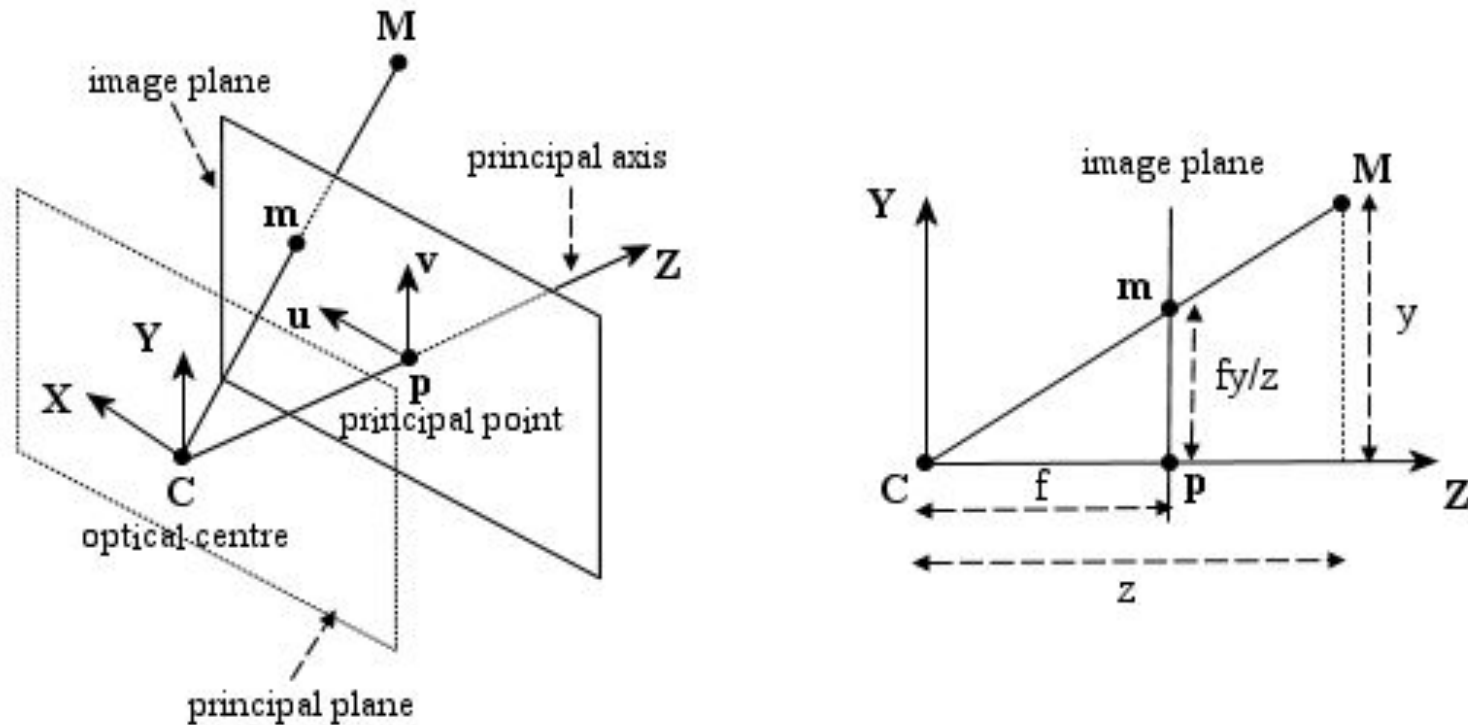


Fig. 2. Pin-hole camera geometry. The left figure illustrates the projection of the point **M** on the image plane by drawing the line through the camera centre **C** and the point to be projected. The right figure illustrates the same situation in the **YZ** plane, showing the similar triangles used to compute the position of the projected point **m** in the image plane.

A 3-D point is projected onto the image plane with the line containing the point and the optical centre (see Figure 2).

Let the centre of projection be the origin of a Cartesian coordinate system wherein the  $z$ -axis is the principal axis.

By similar triangles it is readily seen that the 3-D point  $(x, y, z)^T$  is mapped to the point  $(fx/z, fy/z)^T$  on the image plane.

### 3.1 The camera projection matrix

If the world and image points are represented by homogeneous vectors, then perspective projection can be expressed in terms of matrix multiplication as

$$\begin{pmatrix} fx \\ fy \\ z \end{pmatrix} = \begin{bmatrix} f & 0 & 0 & 0 \\ 0 & f & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix} \begin{pmatrix} x \\ y \\ z \\ 1 \end{pmatrix} \quad (3)$$

The matrix describing the mapping is called the *camera projection matrix*  $P$ .

Equation (3) can be written simply as:

$$z\mathbf{m} = P\mathbf{M} \quad (4)$$

where  $\mathbf{M} = (x, y, z, 1)^T$  are the homogeneous coordinates of the 3-D point and  $\mathbf{m} = (fx/z, fy/z, 1)^T$  are the homogeneous coordinates of the image point.

The projection matrix  $P$  in Eq. (3) represents the simplest possible case, as it only contains information about the focal distance  $f$ .



## General camera

The above formulation assumes a special choice of world coordinate system and image coordinate system. It can be generalized by introducing suitable changes of the coordinates systems.

Changing coordinates in space is equivalent to multiplying the matrix  $P$  to the right by a  $4 \times 4$  matrix:

$$G = \begin{bmatrix} R & \mathbf{t} \\ 0 & 1 \end{bmatrix} \quad (5)$$

$G$  is composed by a rotation matrix  $R$  and a translation vector  $\mathbf{t}$ . It describes the position and orientation of the camera with respect to an external (world) coordinate system. It depends on six parameters, called *extrinsic* parameters.

The rows of  $R$  are unit vectors that, together with the optical centre, define the *camera reference frame*, expressed in world coordinates.

Changing coordinates in the image plane is equivalent to multiplying the matrix  $P$  to the left by a  $3 \times 3$  matrix:

$$K = \begin{bmatrix} f/s_x & f/s_x \cot \theta & o_x \\ 0 & f/s_y & o_y \\ 0 & 0 & 1 \end{bmatrix} \quad (6)$$

$K$  is the *camera calibration matrix*; it encodes the transformation in the image plane from the so-called *normalized camera coordinates* to *pixel coordinates*.

It depends on the so-called *intrinsic* parameters:

- focal distance  $f$  (in mm),
- principal point (or image centre) coordinates  $o_x, o_y$  (in pixel),
- width ( $s_x$ ) and height ( $s_y$ ) of the pixel footprint on the camera photosensor (in mm),
- angle  $\theta$  between the axes (usually  $\pi/2$ ).

The ratio  $s_y/s_x$  is the aspect ratio (usually close to 1).

Thus the camera matrix, in general, is the product of three matrices:

$$P = K[I|\mathbf{0}]G = K[R|\mathbf{t}] \quad (7)$$

In general, the projection equation writes:

$$\zeta \mathbf{m} = PM \quad (8)$$

where  $\zeta$  is the distance of  $\mathbf{M}$  from the focal plane of the camera (this will be shown after), and  $\mathbf{m} = (u, v, 1)^T$ .

Note that, except for a very special choice of the world reference frame, *this “depth” does not coincide with the third coordinate of  $\mathbf{M}$ .*

If  $P$  describes a camera, also  $\lambda P$  for any  $0 \neq \lambda \in \mathbb{R}$  describes the same camera, since these give the same image point for each scene point.

In this case we can also write:

$$\mathbf{m} \simeq P\mathbf{M} \quad (9)$$

where  $\simeq$  means “equal up to a scale factor.”

In general, the camera projection matrix is a  $3 \times 4$  full-rank matrix and, being homogeneous, it has 11 degrees of freedom.

Using QR factorization, it can be shown that any  $3 \times 4$  full rank matrix  $P$  can be factorised as:

$$P = \lambda K[R|\mathbf{t}], \quad (10)$$

( $\lambda$  is recovered from  $K(3, 3) = 1$ ).

## 3.2 Camera anatomy

### Projection centre

The camera projection centre  $\mathbf{C}$  is the only point for which the projection is not defined, i.e.:

$$P\mathbf{C} = P \begin{pmatrix} \tilde{\mathbf{C}} \\ 1 \end{pmatrix} = \mathbf{0} \quad (11)$$

where  $\tilde{\mathbf{C}}$  is a 3-D vector containing the Cartesian (non-homogeneous) coordinates of the optical centre.

After solving for  $\tilde{\mathbf{C}}$  we obtain:

$$\tilde{\mathbf{C}} = -P_{1:3}^{-1}P_4 \quad (12)$$

where the matrix  $P$  is represented by the block form:  $P = [P_{1:3}|P_4]$  (the subscript denotes a range of columns).

## Optical ray

The projection can be geometrically modelled by a ray through the optical centre and the point in space that is being projected onto the image plane (see Fig. 2).

The *optical ray* of an image point  $\mathbf{m}$  is the locus of points in space that projects onto  $\mathbf{m}$ .

It can be described as a parametric line passing through the camera projection centre  $\mathbf{C}$  and a special point (at infinity) that projects onto  $\mathbf{m}$ :

$$\mathbf{M} = \begin{pmatrix} -P_{1:3}^{-1}P_4 \\ 1 \end{pmatrix} + \zeta \begin{pmatrix} P_{1:3}^{-1}\mathbf{m} \\ 0 \end{pmatrix}, \quad \zeta \in \mathbb{R}. \quad (14)$$

If  $\lambda = 1$  the parameter  $\zeta$  in Eq. (14) represent the the depth of the point  $\mathbf{M}$ .

Knowing the intrinsic parameters is equivalent to being able to trace the optical ray of any image point (with  $P = [K|\mathbf{0}]$ ).

### 3.3 Camera calibration (or resection)

A number of point correspondences  $\mathbf{m}_i \leftrightarrow \mathbf{M}_i$  is given, and we are required to find a camera matrix  $P$  such that

$$\mathbf{m}_i \simeq P\mathbf{M}_i \quad \text{for all } i. \quad (15)$$

The equation can be rewritten in terms of the cross product as

$$\mathbf{m}_i \times P\mathbf{M}_i = \mathbf{0}. \quad (16)$$

This form will enable a simple a simple linear solution for  $P$  to be derived. Using the properties of the Kronecker product ( $\otimes$ ) and the `vec` operator [18], we derive:

$$\begin{aligned} \mathbf{m}_i \times P\mathbf{M}_i = \mathbf{0} &\iff [\mathbf{m}_i]_{\times} P\mathbf{M}_i = \mathbf{0} \iff \text{vec}([\mathbf{m}_i]_{\times} P\mathbf{M}_i) = \mathbf{0} \iff \\ &\iff (\mathbf{M}_i^T \otimes [\mathbf{m}_i]_{\times}) \text{vec } P = \mathbf{0} \end{aligned}$$

These are three equations in 12 unknown.

Although there are three equations, only two of them are linearly independent: Indeed, the rank of  $(\mathbf{M}_i^T \otimes [\mathbf{m}_i]_{\times})$  is two because it is the Kronecker product of a rank-1 matrix by a rank-2 matrix.

From a set of  $n$  point correspondences, we obtain a  $2n \times 12$  coefficient matrix  $A$  by stacking up two equations for each correspondence.

In general  $A$  will have rank 11 (provided that the points are not all coplanar) and the solution is the 1-dimensional right null-space of  $A$ .

The projection matrix  $P$  is computed by solving the resulting linear system of equations, for  $n \geq 6$ .

If the data are not exact (noise is generally present) the rank of  $A$  will be 12 and a least-squares solution is sought.

The least-squares solution for  $\text{vec}(P)$  is the singular vector corresponding to the smallest singular value of  $A$ .

This is called the Direct Linear Transform (DLT) algorithm [10].



## 4 Two-View Geometry

The two-view geometry is the intrinsic geometry of two different perspective views of the same 3-D scene (see Figure 3). It is usually referred to as *epipolar geometry*.

The two perspective views may be acquired simultaneously, for example in a stereo rig, or sequentially, for example by a moving camera. From the geometric view-point, the two situations are equivalent, provided that the scene do not change between successive snapshots.

Most 3-D scene points must be visible in both views simultaneously. This is not true in case of occlusions, i.e., points visible only in one camera. Any unoccluded 3-D scene point  $\mathbf{M} = (x, y, z, 1)^T$  is projected to the left and right view as  $\mathbf{m}_\ell = (u_\ell, v_\ell, 1)^T$  and  $\mathbf{m}_r = (u_r, v_r, 1)^T$ , respectively (see Figure 3).

Image points  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$  are called *corresponding points* (or conjugate points) as they represent projections of the same 3-D scene point  $\mathbf{M}$ .

The knowledge of image correspondences enables scene reconstruction from images.

The concept of correspondence is a cornerstone of multiple-view vision. In this notes we assume *known correspondences*, and explore their use in geometric algorithms. Techniques for computing dense correspondences are surveyed in [23, 2].

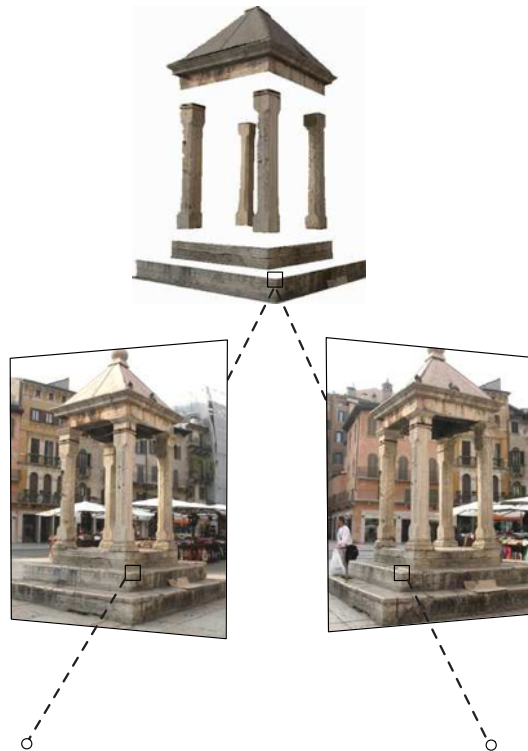


Fig. 3. Two perspective views of the same 3-D scene with conjugate points highlighted

We will refer to the camera projection matrix of the left view as  $P_\ell$  and of the right view as  $P_r$ . The 3-D point  $\mathbf{M}$  is then imaged as (17) in the left view, and (18) in the right view:

$$\zeta_\ell \mathbf{m}_\ell = P_\ell \mathbf{M} \quad (17)$$

$$\zeta_r \mathbf{m}_r = P_r \mathbf{M}. \quad (18)$$

Geometrically, the position of the image point  $\mathbf{m}_\ell$  in the left image plane  $I_\ell$  can be found by drawing the optical ray through the left camera projection centre  $\mathbf{C}_\ell$  and the scene point  $\mathbf{M}$ . The ray intersects the left image plane  $I_\ell$  at  $\mathbf{m}_\ell$ .

Similarly, the optical ray connecting  $\mathbf{C}_r$  and  $\mathbf{M}$  intersects the right image plane  $I_r$  at  $\mathbf{m}_r$ .

The relationship between image points  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$  is given by the epipolar geometry, described in Section 4.1.

## 4.1 Epipolar Geometry

The epipolar geometry describes the geometric relationship between two perspective views of the same 3-D scene.

The key finding, discussed below, is that *corresponding image points must lie on particular image lines*, which can be computed without information on the calibration of the cameras.

This implies that, given a point in one image, one can search the corresponding point in the other along a line and not in a 2-D region, a significant reduction in complexity.

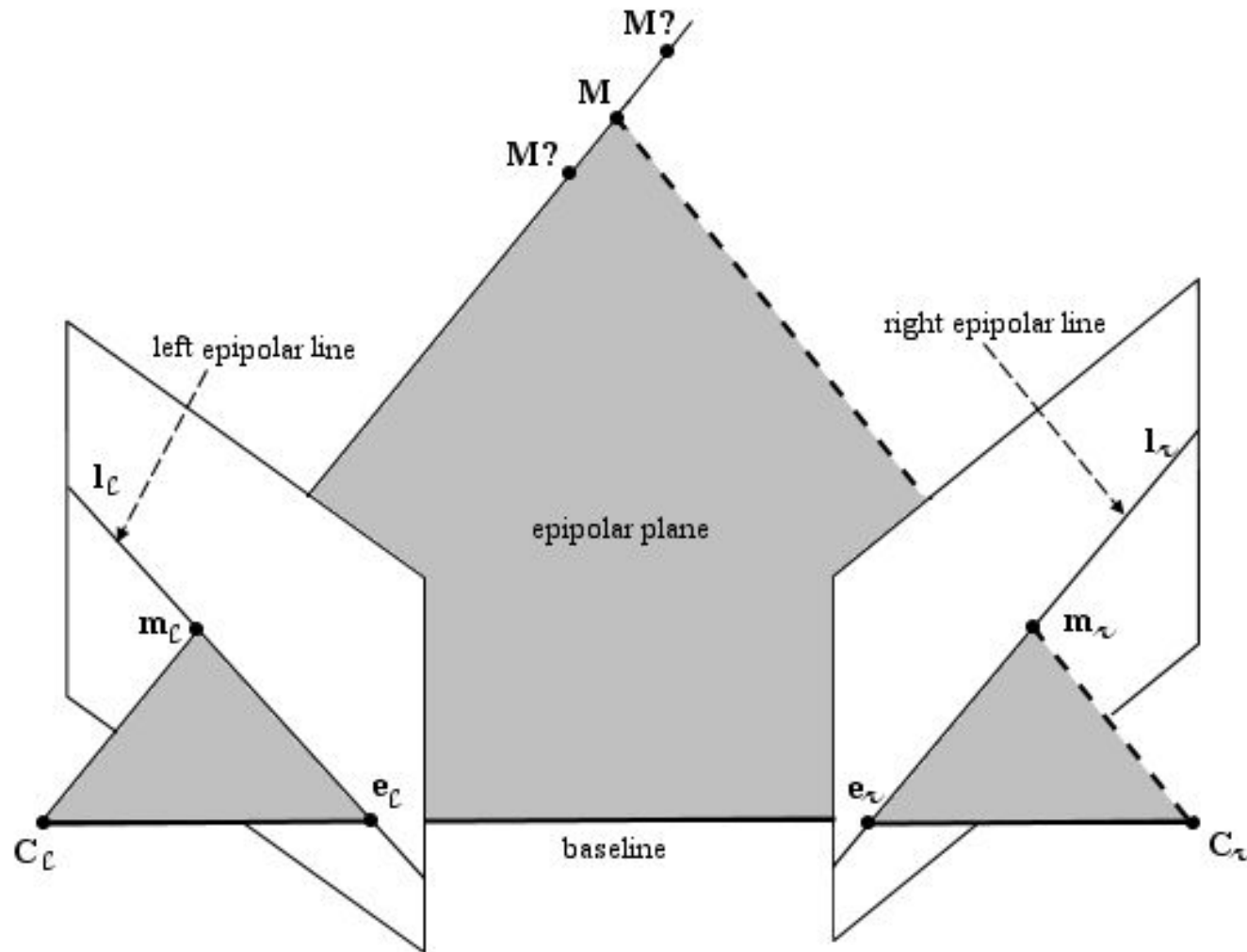


Fig. 4. The epipolar geometry and epipolar constraint.

Any 3-D point  $\mathbf{M}$  and the camera projection centres  $\mathbf{C}_\ell$  and  $\mathbf{C}_r$  define a plane that is called *epipolar plane*.

The projections of the point  $\mathbf{M}$ , image points  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$ , also lie in the epipolar plane since they lie on the rays connecting the corresponding camera projection centre and point  $\mathbf{M}$ .

The conjugate epipolar lines,  $\mathbf{l}_\ell$  and  $\mathbf{l}_r$ , are the intersections of the epipolar plane with the image planes. The line connecting the camera projection centres  $(\mathbf{C}_\ell, \mathbf{C}_r)$  is called the *baseline*.

The baseline intersects each image plane in a point called *epipole*.

By construction, the left epipole  $\mathbf{e}_\ell$  is the image of the right camera projection centre  $\mathbf{C}_r$  in the left image plane. Similarly, the right epipole  $\mathbf{e}_r$  is the image of the left camera projection centre  $\mathbf{C}_\ell$  in the right image plane.

All epipolar lines in the left image go through  $\mathbf{e}_\ell$  and all epipolar lines in the right image go through  $\mathbf{e}_r$ .

## The epipolar constraint.

An epipolar plane is completely defined by the camera projection centres and one image point.

Therefore, given a point  $\mathbf{m}_\ell$ , one can determine the epipolar line in the right image on which the corresponding point,  $\mathbf{m}_r$ , must lie.

The equation of the epipolar line can be derived from the equation describing the optical ray. As we mentioned before, the right epipolar line corresponding to  $\mathbf{m}_\ell$  geometrically represents the projection (Eq. (8)) of the optical ray through  $\mathbf{m}_\ell$  (Eq. (14)) onto the right image plane:

$$\zeta_r \mathbf{m}_r = P_r \mathbf{M} = P_r \underbrace{\begin{pmatrix} -P_{\ell 1:3}^{-1} P_{\ell 4} \\ 1 \end{pmatrix}}_{\mathbf{e}_r} + \zeta_\ell P_r \begin{pmatrix} P_{\ell 1:3}^{-1} \mathbf{m}_\ell \\ 0 \end{pmatrix} \quad (19)$$

---


$$\zeta \mathbf{m} = P \mathbf{M} \quad (8)$$

$$\mathbf{M} = \begin{pmatrix} -P_{1:3}^{-1} P_4 \\ 1 \end{pmatrix} + \zeta \begin{pmatrix} P_{1:3}^{-1} \mathbf{m} \\ 0 \end{pmatrix}, \quad \zeta \in \mathbb{R}. \quad (14)$$

If we now simplify the above equation we obtain the description of the right epipolar line:

$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \zeta_\ell \underbrace{P_{r1:3} P_{\ell1:3}^{-1} \mathbf{m}_\ell}_{\mathbf{m}'_\ell} \quad (20)$$

This is the equation of a line through the right epipole  $\mathbf{e}_r$  and the image point  $\mathbf{m}'_\ell$  which represents the projection onto the right image plane of the point at infinity of the optical ray of  $\mathbf{m}_\ell$ .

The equation for the left epipolar line is obtained in a similar way.





Fig. 5. Left and right images with epipolar lines.

## 4.2 Triangulation

Given the camera matrices  $P_\ell$  and  $P_r$ , let  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$  be two corresponding points satisfying the epipolar constraint. It follows that  $\mathbf{m}_r$  lies on the epipolar line  $F\mathbf{m}_\ell$  and so the two rays back-projected from image points  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$  lie in a common epipolar plane. Since they lie in the same plane, they will intersect at some point. This point is the reconstructed 3-D scene point  $\mathbf{M}$ .

Analytically, the reconstructed 3-D point  $\mathbf{M}$  can be found by solving for parameter  $\zeta_\ell$  or  $\zeta_r$  in Eq. (20). Let us rewrite it as:

$$\mathbf{e}_r = \zeta_r \mathbf{m}_r - \zeta_\ell \mathbf{m}'_\ell \quad (21)$$

The depth  $\zeta_r$  and  $\zeta_\ell$  are unknown. Both encode the position of  $\mathbf{M}$  in space, as  $\zeta_r$  is the depth of  $\mathbf{M}$  wrt the right camera and  $\zeta_\ell$  is the depth of  $\mathbf{M}$  wrt the left camera.

---


$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \underbrace{\zeta_\ell P_{r1:3} P_{\ell1:3}^{-1} \mathbf{m}_\ell}_{\mathbf{m}'_\ell} \quad (20)$$

However, triangulation can also be cast as a null-space problem.

Let us consider  $\mathbf{m} = [u, v, 1]^T$ , the projection of the 3D point  $M$  according to the perspective projection matrix  $P$ . From (8) one obtains:

$$\begin{cases} (\mathbf{p}_1 - u\mathbf{p}_3)^T \mathbf{M} = 0 \\ (\mathbf{p}_2 - v\mathbf{p}_3)^T \mathbf{M} = 0 \end{cases} \quad (22)$$

and then, in matrix form:

$$\begin{bmatrix} (\mathbf{p}_1 - u\mathbf{p}_3)^T \\ (\mathbf{p}_2 - v\mathbf{p}_3)^T \end{bmatrix} \mathbf{M} = \mathbf{0}_{2 \times 1} \quad (23)$$

Hence, one point gives two homogeneous equations.

Let us consider now  $\mathbf{m}' = [u', v', 1]^T$ , the corresponding point of  $\mathbf{m}$  in the second image, and let  $P'$  be the second perspective projection matrix.

---


$$\zeta \mathbf{m} = P\mathbf{M} \quad (8)$$

Being both projection of the same 3D point  $\mathbf{M}$ , the equations provided by  $\mathbf{m}$  and  $\mathbf{m}'$  can be stacked:

$$\begin{bmatrix} (\mathbf{p}_1 - u\mathbf{p}_3)^T \\ (\mathbf{p}_2 - v\mathbf{p}_3)^T \\ (\mathbf{p}'_1 - u'\mathbf{p}'_3)^T \\ (\mathbf{p}'_2 - v'\mathbf{p}'_3)^T \end{bmatrix} \mathbf{M} = \mathbf{0}_{4 \times 1} \quad (24)$$

The solution is the null-space of the  $4 \times 4$  coefficient matrix, which must then have rank three, otherwise only the trivial solution  $\mathbf{M} = \mathbf{0}$  would be possible. In the presence of noise this rank condition cannot be fulfilled exactly, so a least squares solution is sought, typically via SVD, as in calibration with DLT. In [12] this method is called “*linear-eigen*”.

This method generalizes to the case of  $N > 2$  cameras: each one gives two equations and one ends up with  $2N$  equations in four unknowns.

Triangulation is addressed in more details in [1, 12, 10].

## Levels of description...

The epipolar geometry can be described analytically in several ways, depending on the amount of the *a priori* knowledge about the stereo system. We can identify three general cases.

- (i) If both *intrinsic* and *extrinsic* camera parameters are known, we can describe the epipolar geometry in terms of the projection matrices (Equation (20)).
- (ii) If only the *intrinsic* parameters are known, we work in normalized camera coordinates and the epipolar geometry is described by the *essential matrix*.
- (iii) If neither intrinsic nor extrinsic parameters are known the epipolar geometry is described by the *fundamental matrix*.

---


$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \zeta_l \underbrace{P_{r1:3} P_{l1:3}^{-1}}_{\mathbf{m}'_l} \mathbf{m}_l \quad (20)$$

### **...and ambiguity in reconstruction.**

Likewise, what can be reconstructed (by triangulation) depends on what is known about the scene and the stereo system. We can identify three cases.

- (i) *If both the intrinsic and extrinsic camera parameters are known*, we can solve the reconstruction problem unambiguously.
- (ii) *If only the intrinsic parameters are known*, we can estimate the extrinsic parameters and solve the reconstruction problem up to an unknown scale factor (+ a rigid transformation that correspond to the arbitrariness in fixing the world reference frame). In other words,  $R$  can be estimated completely, and  $t$  up to a scale factor.
- (iii) *If neither intrinsic nor extrinsic parameters are known*, i.e., the only information available are pixel correspondences, we can still solve the reconstruction problem but only up to an unknown, global projective transformation of the world. This ambiguity w may be reduced if additional information is supplied on the cameras or the scene (see Sec .5).

## 4.3 The calibrated case

Suppose that a set of image correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$  are given. It is assumed that these correspondences come from a set of 3-D points  $\mathbf{M}_i$ , which are unknown.

The intrinsic parameters are known, i.e. the cameras are *calibrated*, but the position and attitude of the cameras are unknown.

The situation – discussed previously – when the intrinsic *and* extrinsic parameters are known will be referred to as *full calibrated* for the sake of clarity.

We will see that the epipolar geometry is described by the essential matrix and that, starting from the essential matrix, only a reconstruction up to a similarity transformation (rigid+uniform scale) can be achieved. Such a reconstruction is referred to as “Euclidean”.

$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \zeta_\ell \underbrace{P_{r1:3} P_{\ell1:3}^{-1} \mathbf{m}_\ell}_{\mathbf{m}'_\ell} \quad (20)$$

### 4.3.1 The Essential Matrix E

As the intrinsic parameters are known, we can switch to *normalized camera coordinates*:  $\mathbf{m} \leftarrow K^{-1} \mathbf{m}$  (please note that this change of notation will hold throughout this section).

Consider a pair of cameras  $P_\ell$  and  $P_r$ . Without loss of generality, we can fix the world reference frame onto the first camera, hence:

$$P_\ell = [I|0] \quad \text{and} \quad P_r = [R|\mathbf{t}]. \quad (25)$$

With this choice, the unknown extrinsic parameters have been made explicit.

If we substitute these two particular instances of the camera projection matrices in Equation (20), we get

$$\zeta_r \mathbf{m}_r = \mathbf{t} + \zeta_\ell R \mathbf{m}_\ell; \quad (26)$$

in other words, the point  $\mathbf{m}_r$  lies on the line through the points  $\mathbf{t}$  and  $R\mathbf{m}_\ell$ . In homogeneous coordinates, this can be written as follows:

$$\mathbf{m}_r^T (\mathbf{t} \times R \mathbf{m}_\ell) = 0, \quad (27)$$



as the homogeneous line through two points is expressed as their cross product, and a dot product of a point and a line is zero if the point lies on the line.

The cross product of two vectors can be written as a product of a skew-symmetric matrix and a vector. Equation (27) can therefore be equivalently written as

$$\mathbf{m}_r^T [\mathbf{t}]_{\times} R \mathbf{m}_\ell = 0, \quad (28)$$

where  $[\mathbf{t}]_{\times}$  is the skew-symmetric matrix of the vector  $\mathbf{t}$ . Let us define the *essential matrix*  $E$ :

$$E \triangleq [\mathbf{t}]_{\times} R. \quad (29)$$

In summary, the relationship between the corresponding image points  $\mathbf{m}_\ell$  and  $\mathbf{m}_r$  in normalized camera coordinates is the bilinear form:

$$\mathbf{m}_r^T E \mathbf{m}_\ell = 0. \quad (30)$$

$E$  encodes only information on the rigid displacement between cameras. It has five degrees of freedom: a 3-D rotation and a 3-D translation direction.

$E$  is singular, since  $\det[\mathbf{t}]_{\times} = 0$ , and it is a homogeneous quantity.

### 4.3.2 Reconstruction up to a Similarity

If a sufficient number of point correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$  is given, we can use Equation (30) to compute the unknown matrix  $E$  (see Sec. 4.4.2).

The reconstruction is achieved starting from the essential matrix, which contains – entangled – the unknown extrinsic parameters.

Unlike the fundamental matrix, the only property of which is to have rank two, the essential matrix is characterised by the following theorem [15].

**Theorem 4.1** *A real  $3 \times 3$  matrix  $E$  can be factorised as product of a nonzero skew-symmetric matrix and a rotation matrix if and only if  $E$  has two identical singular values and a zero singular value.*

The rotation  $R$  and translation  $\mathbf{t}$  are then used to instantiate a camera pair as in Equation (25), and this camera pair is subsequently used to reconstruct the structure of the scene by triangulation.

The rigid displacement ambiguity arises from the arbitrary choice of the world reference frame, whereas the scale ambiguity derives from the fact that  $\mathbf{t}$  can be scaled arbitrarily in Equation (29) and one would get the same essential matrix ( $E$  is defined up to a scale factor).

Therefore translation can be recovered from  $E$  only up to an unknown scale factor which is inherited by the reconstruction.

This is also known as *depth-speed ambiguity* (in a context where points are moving and camera is stationary): a large motion of a distant point and a small motion of a nearby point produces the same motion in the image.

---

$$P_\ell = [I|0] \quad \text{and} \quad P_r = [R|\mathbf{t}]. \quad (25)$$

$$E \triangleq [\mathbf{t}]_\times R. \quad (29)$$

## 4.4 The weakly calibrated case

Suppose that a set of image correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$  are given. It is assumed that these correspondences come from a set of 3-D points  $\mathbf{M}_i$ , which are unknown.

Similarly, the position, attitude and calibration of the cameras are not known.

This situation is usually referred to as *weak calibration*, and we will see that the epipolar geometry is described by the *fundamental matrix* and the scene may be reconstructed up to a projective ambiguity.

$$P_\ell = [I|0] \quad \text{and} \quad P_r = [R|\mathbf{t}]. \quad (25)$$

#### 4.4.1 The Fundamental Matrix F

The fundamental matrix can be derived in a similar way to the essential matrix. All camera parameters are assumed unknown; we write therefore a more general version of Equation (25):

$$P_\ell = K_\ell[I|0] \quad \text{and} \quad P_r = K_r[R|\mathbf{t}]. \quad (33)$$

Inserting these two projection matrices into Equation (20), we get

$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \zeta_\ell K_r R K_\ell^{-1} \mathbf{m}_\ell \quad \text{with} \quad \mathbf{e}_r = K_r \mathbf{t}, \quad (34)$$

which states that point  $\mathbf{m}_r$  lies on the line through  $\mathbf{e}_r$  and  $K_r R K_\ell^{-1} \mathbf{m}_\ell$ . As in the case of the essential matrix, this can be written in homogeneous coordinates as:

$$\mathbf{m}_r^T [\mathbf{e}_r]_\times K_r R K_\ell^{-1} \mathbf{m}_\ell = 0. \quad (35)$$

The matrix

$$F = [\mathbf{e}_r]_\times K_r R K_\ell^{-1} \quad (36)$$

is the *fundamental matrix*  $F$ , giving the relationship between the corresponding image points in pixel coordinates.

Therefore, the bilinear form that links corresponding points writes:

$$\mathbf{m}_r^T F \mathbf{m}_\ell = 0. \quad (37)$$

$F$  is the algebraic representation of the epipolar geometry in the least information case. It is a  $3 \times 3$ , rank-two homogeneous matrix. It has only seven degrees of freedom since it is defined up to a scale and its determinant is zero. Notice that  $F$  is completely defined by pixel correspondences only (the intrinsic parameters are not needed).

For any point  $\mathbf{m}_\ell$  in the left image, the corresponding epipolar line  $\mathbf{l}_r$  in the right image can be expressed as

$$\mathbf{l}_r = F \mathbf{m}_\ell. \quad (38)$$

Similarly, the epipolar line  $\mathbf{l}_\ell$  in the left image for the point  $\mathbf{m}_r$  in the right image can be expressed as

$$\mathbf{l}_\ell = F^T \mathbf{m}_r. \quad (39)$$

The left epipole  $\mathbf{e}_\ell$  is the right null-vector of the fundamental matrix and the right epipole is the left null-vector of the fundamental matrix:

$$F \mathbf{e}_\ell = 0 \quad (40)$$

$$\mathbf{e}_r^T F = 0 \quad (41)$$

One can see from the derivation that the essential and fundamental matrices are related through the camera calibration matrices  $K_\ell$  and  $K_r$ :

$$F = K_r^{-T} E K_\ell^{-1}. \quad (42)$$

Consider a camera pair. Using the fact that if  $F$  maps points in the left image to epipolar lines in the right image, then  $F^T$  maps points in the right image to epipolar lines in the left image, Equation (34) gives:

$$\zeta_r F^T \mathbf{m}_r = \zeta_\ell (\mathbf{e}_\ell \times \mathbf{m}_\ell). \quad (43)$$

This is another way of writing the epipolar constraint: the epipolar line of  $\mathbf{m}_r$  ( $F^T \mathbf{m}_r$ ) is the line containing its corresponding point ( $\mathbf{m}_\ell$ ) and the epipole in the left image ( $\mathbf{e}_\ell$ ).

---


$$\zeta_r \mathbf{m}_r = \mathbf{e}_r + \zeta_\ell K_r R K_\ell^{-1} \mathbf{m}_\ell \quad \text{with} \quad \mathbf{e}_r = K_r \mathbf{t}, \quad (34)$$

## 4.4.2 Estimating $F$ : the eight-point algorithm

If a number of point correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$  is given, we can use Equation (37) to compute the unknown matrix  $F$ .

We need to convert Equation (37) from its bilinear form to a form that matches the null-space problem. To this end we use again the `vec` operator, as in the DLT algorithm:

$$\mathbf{m}_r^T F \mathbf{m}_\ell = 0 \iff \text{vec}(\mathbf{m}_r^T F \mathbf{m}_\ell) = 0 \iff (\mathbf{m}_r^T \otimes \mathbf{m}_\ell^T) \text{vec}(F) = 0.$$

Each point correspondence gives rise to one linear equation in the unknown entries of  $F$ . From a set of  $n$  point correspondences, we obtain a  $n \times 9$  coefficient matrix  $A$  by stacking up one equation for each correspondence.

In general  $A$  will have rank 8 and the solution is the 1-dimensional right null-space of  $A$ .

---


$$\mathbf{m}_r^T F \mathbf{m}_\ell = 0. \tag{37}$$



The fundamental matrix  $F$  is computed by solving the resulting linear system of equations, for  $n \geq 8$ .

If the data are not exact and more than 8 points are used, the rank of  $A$  will be 9 and a least-squares solution is sought.

The least-squares solution for  $\text{vec}(F)$  is the singular vector corresponding to the smallest singular value of  $A$ .

This method does not explicitly enforce  $F$  to be singular, so it must be done *a posteriori*.

Replace  $F$  by  $F'$  such that  $\det F' = 0$ , by forcing to zero the least singular value.

It can be shown that  $F'$  is the closest singular matrix to  $F$  in Frobenius norm.

Geometrically, the singularity constraint ensures that the epipolar lines meet in a common epipole.

### 4.4.3 Reconstruction up to a Projective Transformation

The reconstruction task is to find the camera matrices  $P_\ell$  and  $P_r$ , as well as the 3-D points  $\mathbf{M}_i$  such that

$$\mathbf{m}_\ell^i = P_\ell \mathbf{M}^i \quad \text{and} \quad \mathbf{m}_r^i = P_r \mathbf{M}^i, \quad \forall i \quad (44)$$

If  $T$  is any  $4 \times 4$  invertible matrix, representing a collineation of  $\mathbb{P}_3$ , then replacing points  $\mathbf{M}^i$  by  $T\mathbf{M}^i$  and matrices  $P_\ell$  and  $P_r$  by  $P_\ell T^{-1}$  and  $P_r T^{-1}$  does not change the image points  $\mathbf{m}_\ell^i$ . This shows that, if nothing is known but the image points, the structure  $\mathbf{M}^i$  and the cameras can be determined only up to a projective transformation.

The procedure for reconstruction follows the previous one. Given the weak calibration assumption, the fundamental matrix can be computed (using the algorithm described in Section 4.4.1), and from a (non-unique) factorization of  $F$  of the form

$$F = [\mathbf{e}_r]_\times A \quad (45)$$

two camera matrices  $P_\ell$  and  $P_r$ :

$$P_\ell = [I | \mathbf{0}] \quad \text{and} \quad P_r = [A | \mathbf{e}_r], \quad (46)$$

can be created in such a way that they yield the fundamental matrix  $F$ , as can be easily verified. The position in space of the points  $\mathbf{M}^i$  is then obtained by triangulation.

The matrix  $A$  in the factorization of  $F$  can be set to  $A = -[\mathbf{e}_r]_{\times} F$  (this is called the *epipolar projection matrix* [17]).

Unlike the essential matrix,  $F$  does not admit a unique factorization, whence the projective ambiguity follows.

Indeed, for any  $A$  satisfying Equation (45), also  $A + \mathbf{e}_r \mathbf{x}^T$  for any vector  $\mathbf{x}$ , satisfies Equation (45).

More in general, any homography induced by a plane can be taken as the  $A$  matrix.

---

$$F = [\mathbf{e}_r]_{\times} A \quad (45)$$

## 4.5 More on calibration

Calibration problems can be cast as determining the transformation between two reference frames. Depending on the nature of the available measures and where the reference frames are attached we have four calibration or *orientation* problems<sup>1</sup>:

**Relative orientation** is the problem of determining the position and attitude of one perspective camera with respect to another camera from correspondences between points in 2-D images (See. Sec. 4.3.2).

**Absolute orientation** is the problem of aligning two sets of points, whose 3-D locations have been measured (or reconstructed) in two different reference frames.

**Exterior orientation** is the problem of determining the position and attitude of a perspective camera from correspondences between 3-D points and their 2-D images.

**Interior orientation** is the problem of determining the (affine) transformation from normalized camera coordinates to pixel coordinates, i.e., the intrinsic parameters of the camera.

---

<sup>1</sup>this terminology comes from Photogrammetry.

### 4.5.1 Absolute orientation (with scaling)

Given two sets of 3-D points  $\mathbf{X}^i$  and  $\mathbf{Y}^i$ , related by<sup>2</sup>

$$\mathbf{X}^i = s(R\mathbf{Y}^i + t) \quad \text{for all } i = 1 \dots N \quad (47)$$

we are required to find the rotation matrix  $R$ , the vector  $t$  and the scalar  $s$ .

Summing these equations for all  $i$  and dividing by  $N$  shows that the translation is found with:

$$t = \frac{1}{s} \left( \frac{1}{N} \sum_{i=1}^N \mathbf{X}^i \right) - R \left( \frac{1}{N} \sum_{i=1}^N \mathbf{Y}^i \right)$$

Combining this with Eq. (47) gives

$$\bar{\mathbf{X}}^i = sR\bar{\mathbf{Y}}^i$$

where  $\bar{\mathbf{X}}^i = \mathbf{X}^i - \frac{1}{N} \sum_{i=1}^N \mathbf{X}^i$  and  $\bar{\mathbf{Y}}^i = \mathbf{Y}^i - \frac{1}{N} \sum_{i=1}^N \mathbf{Y}^i$ .

Because the rotation matrix does not change the length of the vectors, we can immediately solve for the scale from  $\|\bar{\mathbf{X}}^i\| = s\|\bar{\mathbf{Y}}^i\|$ .

We are left with the problem of estimating the unknown rotation between two sets of points.

Let  $\bar{X}$  be the  $3 \times N$  matrix formed by stacking the points  $\bar{X}^i$  side by side and  $\bar{Y}$  be the matrix formed likewise by stacking the scaled points  $s\bar{Y}^i$ . In presence of noise, we would like to minimize the sum of the square of the errors, or

$$\sum_{i=1}^N \|\bar{X}^i - sR\bar{Y}^i\|^2 = \|\bar{X} - R\bar{Y}\|_F^2$$

where  $\|\cdot\|_F$  is the Frobenius norm.

This problem is known as the Orthogonal Procrustes Problem and the solution is given by [16]

$$R = V \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \det(VU^T) \end{bmatrix} U^T$$

where  $UDV^T = \bar{Y}\bar{X}^T$  is the SVD of the  $3 \times 3$  matrix  $\bar{Y}\bar{X}^T$ .

## 4.5.2 Exterior orientation

*Exterior orientation* (also called *Perspective  $n$ -Points* camera pose) is a problem that appears repeatedly in computer vision, but in context different from 3-D reconstruction, such as visual servoing and augmented reality.

Given a number of point correspondences  $\mathbf{m}^i \leftrightarrow \mathbf{M}^i$  and the intrinsic camera parameters  $K$ , we are required to find a rotation matrix  $R$  and a translation vector  $\mathbf{t}$  (which specify attitude and position of the camera) such that:

$$\zeta^i K^{-1} \mathbf{m}^i = [R|\mathbf{t}]\mathbf{M}^i = (R\tilde{\mathbf{M}}^i + \mathbf{t}) \quad \text{for all } i. \quad (48)$$

One could immediately solve this problem by doing camera resection with DLT in normalized camera coordinates. The algorithm is linear, but it does not enforce the orthonormality constraints on the rotation matrix.

Instead, we present here the linear method proposed by Fiore [6]. He first recovers the unknown depths  $\zeta^i$ , and then observes that what is left is an absolute orientation problem, whose solution yields a rotation matrix which is inherently orthonormal.

In order to recover the depths, let's write Eq. (48) in matrix form:

$$K^{-1} \underbrace{[\zeta^1 \mathbf{m}^1, \zeta^2 \mathbf{m}^2, \dots, \zeta^n \mathbf{m}^n]}_W = [R|\mathbf{t}] \underbrace{[\mathbf{M}^1, \mathbf{M}^2, \dots, \mathbf{M}^n]}_M. \quad (49)$$

Let  $r = \text{rank } M$ . Take its SVD:  $M = UDV^T$  and let  $V_2$  be a matrix composed by the last  $n - r$  columns of  $V$ , which spans the null-space of  $M$ . Then,  $MV_2 = 0_{3 \times (n-r)}$ , and also

$$K^{-1} W V_2 = 0_{3 \times (n-r)} \quad (50)$$

By taking  $\text{vec}$  on both sides we get:

$$(V_2^T \otimes K^{-1}) \text{vec}(W) = \mathbf{0}. \quad (51)$$

Let us observe that:

$$\text{vec}(W) = \begin{bmatrix} \zeta^1 \mathbf{m}^1 \\ \zeta^1 \mathbf{m}^2 \\ \vdots \\ \zeta^n \mathbf{m}^n \end{bmatrix} = \underbrace{\begin{bmatrix} \mathbf{m}^1 & 0 & \dots & 0 \\ & & \ddots & \\ 0 & 0 & \dots & \mathbf{m}^n \end{bmatrix}}_D \underbrace{\begin{bmatrix} \zeta^1 \\ \vdots \\ \zeta^n \end{bmatrix}}_{\zeta} \quad (52)$$



Hence

$$((V_2^T \otimes K^{-1})D) \zeta = \mathbf{0}. \quad (53)$$

From the last equation the depths  $\zeta$  can be recovered (up to a scale factor) by solving a null-space problem.

The size of the coefficients matrix is  $3(n - r) \times n$ , and in order to determine a one-parameter family of solutions, it must have rank  $n - 1$ , hence  $3(n - r) \geq n - 1$ .

Therefore, at least  $n \geq (3r - 1)/2$  points are needed. If points are in general position, 6 are sufficient, but if they are on a plane, only 4 suffices.

Now that the left side of Eq. (48) is known, up to a scale factor, we are left with an absolute orientation (with scale) problem:

$$\zeta^i K^{-1} \mathbf{m}^i = s(R\tilde{\mathbf{M}}^i + \mathbf{t}) \quad \text{for all } i. \quad (54)$$

which we solve using the algorithm of Sec. 4.5.1. As a result, the rotation matrix estimate is orthonormal by construction.

---


$$\zeta^i K^{-1} \mathbf{m}^i = [R|\mathbf{t}]\mathbf{M}^i = (R\tilde{\mathbf{M}}^i + \mathbf{t}) \quad \text{for all } i. \quad (48)$$

## 5 Autocalibration

The aim of *autocalibration* is to compute the intrinsic parameters, starting from weakly calibrated cameras.

More in general, the task is to recover metric properties of camera and/or scene, i.e., to compute a Euclidean reconstruction.

There are two classes of methods:

1. Direct: solve directly for the intrinsic parameters.
2. Stratified: first obtain a projective reconstruction and then transform it to a Euclidean reconstruction (in some cases an affine reconstruction is obtained in between).

The reader is referred to [7] for a review of autocalibration, and to [19, 26, 14, 21, 20, 9] for classical and recent work on the subject.

## 5.1 Counting argument

Consider  $m$  cameras. The difference between the d.o.f. of the multifocal geometry (e.g. 7 for two views) and the d.o.f. of the rigid displacements (e.g. 5 for two views) is the number of independent constraints available for the computation of the intrinsic parameters (e.g. 2 for two views).

The multifocal geometry of  $m$  cameras (represented by the m-focal tensor) has  $11m - 15$  d.o.f. Proof: a set of  $m$  cameras have  $11m$  d.o.f., but they determine the m-focal geometry up to a collineation of  $\mathbb{P}_3$ , which has 15 d.o.f. The net sum is  $11m - 15$  d.o.f.

On the other hand, the rigid displacements in  $m$  views are described by  $6m - 7$  parameters:  $3(m - 1)$  for rotations,  $2(m - 1)$  for translations, and  $m - 2$  ratios of translation norms.

Thus,  *$m$  weakly calibrated views give  $5m - 8$  constraints available for computing the intrinsic parameters.*

Let us suppose that  $m_k$  parameters are known and  $m_c$  parameters are constant.

The first view introduces  $5 - m_k$  unknowns. Every view but the first introduces  $5 - m_k - m_c$  unknowns.

Therefore, the unknown intrinsic parameters can be computed provided that

$$5m - 8 \geq (m - 1)(5 - m_k - m_c) + 5 - m_k. \quad (55)$$

For example, if the intrinsic parameters are constant, three views are sufficient to recover them.

If one parameter (usually the skew) is known and the other parameters are varying, at least eight views are needed.

## 5.2 A simple direct method

If we consider two views, two independent constraints are available for the computation of the intrinsic parameters from the fundamental matrix.

Indeed,  $F$  has 7 d.o.f, whereas  $E$ , which encode the rigid displacement, has only 5 d.o.f. There must be two additional constraint that  $E$  must satisfy, with respect to  $F$ .

In particular, these constraints stem from the equality of two singular values of the essential matrix (Theorem 4.1) which can be decomposed in two independent polynomial equations.

Let  $F_{ij}$  be the (known) fundamental matrix relating views  $i$  and  $j$ , and let  $K_i$  and  $K_j$  be the respective (unknown) intrinsic parameter matrices.

The idea of [20] is that the matrix

$$E_{ij} = K_i^T F_{ij} K_j, \quad (56)$$

satisfies the constraints of Theorem 4.1 only if the intrinsic parameters are correct.

## 5.3 Stratification

We have seen that a projective reconstruction can be computed starting from points correspondences only (weak calibration), without any knowledge of the camera matrices.

Projective reconstruction differs from Euclidean by an unknown projective transformation in the 3-D projective space, which can be seen as a suitable change of basis.

Starting from a projective reconstruction the problem is computing the transformation that “straighten” it, i.e., that upgrades it to an Euclidean reconstruction.

To this purpose the problem is *stratified* [17, 4] into different representations: depending on the amount of information and the constraints available, it can be analyzed at a projective, affine, or Euclidean level.

Let us assume that a projective reconstruction is available, that is a sequence  $P_i$  of  $m + 1$  camera matrices and a set  $\mathbf{M}^j$  of  $n + 1$  3-D points such that:

$$\mathbf{m}_i^j \simeq P_i \mathbf{M}^j \quad i = 0 \dots m, \quad j = 0 \dots n. \quad (58)$$

Without loss of generality, we can assume that camera matrices writes:

$$P_0 = [I \mid \mathbf{0}]; \quad P_i = [A_i \mid \mathbf{e}_i] \quad \text{for } i = 1 \dots m \quad (59)$$

We are looking for the a  $4 \times 4$  non-singular matrix  $T$  that upgrades the projective reconstruction to Euclidean:

$$\mathbf{m}_i^j \simeq \underbrace{P_i T}_{P_i^E} \underbrace{T^{-1} \mathbf{M}^j}_{\text{structure}}, \quad (60)$$

$P_i^E = P_i T$  is the Euclidean camera,

We can choose the first Euclidean-calibrated camera to be  $P_0^E = K_0 [I \mid \mathbf{0}]$ , thereby fixing arbitrarily the world reference frame:

$$P_0^E = K_0 [I \mid \mathbf{0}] \quad P_i^E = K_i [R_i \mid \mathbf{t}_i] \quad \text{for } i = 1 \dots m. \quad (61)$$

With this choice, it is easy to see that  $P_0^E = P_0T$  implies

$$T = \begin{bmatrix} K_0 & \mathbf{0} \\ \mathbf{r}^T & s \end{bmatrix} \quad (62)$$

where  $\mathbf{r}^T$  is a 3-D vector and  $s$  is a scale factor, which we will arbitrarily set to 1 (the Euclidean reconstruction is up to a scale factor).

Under this parametrization  $T$  is clearly non singular, and it depends on eight parameters.

Substituting (62) in  $P_i^E \simeq P_iT$  gives

$$P_i^E = [K_i R_i \mid K_i \mathbf{t}_i] \simeq P_i T = [A_i K_0 + \mathbf{e}_i \mathbf{r}^T \mid \mathbf{e}_i] \quad \text{for } i > 0 \quad (63)$$

and, considering only the leftmost  $3 \times 3$  submatrix, gives

$$K_i R_i \simeq A_i K_0 + \mathbf{e}_i \mathbf{r}^T = P_i \begin{bmatrix} K_0 \\ \mathbf{r}^T \end{bmatrix} \quad (64)$$



Rotation can be eliminated using  $RR^T = I$ , leaving:

$$K_i K_i^T \simeq P_i \begin{bmatrix} K_0 K_0^T & K_0 \mathbf{r} \\ \mathbf{r}^T K_0^T & \mathbf{r}^T \mathbf{r} \end{bmatrix} P_i^T \quad (65)$$

This is the basic equation for autocalibration (called *absolute quadric constraint*), relating the unknowns  $K_i$  ( $i = 0 \dots m$ ) and  $\mathbf{r}$  to the available data  $P_i$  (obtained from weakly calibrated images).

Note that (65) contains five equations, because the matrices of both members are symmetric, and the homogeneity reduces the number of equations with one.

# Geometric toolbox review

- Resection:
  - given image to 3D correspondences, compute  $P$
  
- Exterior orientation:
  - given image to 3D correspondences, compute  $(R,t)$
  
- Absolute orientation:
  - Given 3D-3D correspondences, compute similarity or projectivity linking the two sets
  
- Relative orientation
  - Essential matrix  $E$  can be computed from image correspondences + camera parameters  $K$
  - Rigid motion  $(R,t)$  can be extracted from  $E$

# Geometric toolbox review

- Triangulation (or intersection)
  - Given  $P$  ( $\geq 2$ ) and image correspondences, compute 3D points.
- Autocalibration.
  - Recover  $K$  from image correspondences.
- Bundle adjustment.
  - Optimize for  $P$  and 3D points simultaneously.

## 6 Dealing with errors

In this section we will approach estimation problems from a more “practical” point of view.

First, we will discuss how the presence of errors in the data affects our estimates and describe the countermeasures that must be taken to obtain a good estimate.

Errors can be small (the match is correct but the position of the point has a limited accuracy), and these are usually modeled as Gaussian noise, or large (the match is wrong) and these are called *outliers*.

Since they have a different nature, they will be treated differently.

## 6.1 Pre-conditioning

In presence of noise (or errors) on input data, the accuracy of the solution of a linear system depends crucially on the *condition number* of the system. The lower the condition number, the less the input error gets amplified (the system is more stable).

As [11] pointed out, it is crucial for linear algorithms (as the DLT algorithm) that input data is properly pre-conditioned, by a suitable coordinate change (origin and scale): points are translated so that their centroid is at the origin and are scaled so that their average distance from the origin is  $\sqrt{2}$ .

This improves the condition number of the linear system that is being solved.

Apart from improved accuracy, this procedure also provides invariance under similarity transformations in the image plane.

## 6.2 Algebraic vs geometric error

Measured data (i.e., image or world point positions) is noisy.

Usually, to counteract the effect of noise, we use more equations than necessary and solve with least-squares.

What is actually being minimized by least squares?

In a typical null-space problem formulation  $Ax = 0$  (like the DLT algorithm) the quantity that is being minimized is the square of the residual  $\|Ax\|$ .

In general, if  $\|Ax\|$  can be regarded as a distance between the geometrical entities involved (points, lines, planes, etc..), then what is being minimized is a geometric error, otherwise (when the error lacks a good geometrical interpretation) it is called an algebraic error.

All the linear algorithm (DLT and others) we have seen so far minimize an algebraic error. Actually, there is no justification in minimizing an algebraic error apart from the ease of implementation, as it results in a linear problem.

Usually, the minimization of a geometric error is a non-linear problem, that admit only iterative solutions and requires a starting point.

So, why should we prefer to minimize a geometric error? Because:

- The quantity being minimized has a meaning
- The solution is more stable
- The solution is invariant under Euclidean transforms

Often linear solution based on algebraic residuals are used as a starting point for a non-linear minimization of a geometric cost function, which “gives the solution a final polish” [10].

## Geometric error for resection

The goal is to estimate the camera matrix, given a number of correspondences  $(\mathbf{m}^j, \mathbf{M}^j)$   $j = 1 \dots n$

The geometric error associated to a camera estimate  $\hat{P}$  is the distance between the measured image point  $\mathbf{m}^j$  and the re-projected point  $\hat{P}_i \mathbf{M}^j$ :

$$\min_{\hat{P}} \sum_j d(\hat{P} \mathbf{M}^j, \mathbf{m}^j)^2 \quad (75)$$

where  $d()$  is the Euclidean distance between the homogeneous points.

The DLT solution is used as a starting point for the iterative minimization (e.g. Gauss-Newton)



## Geometric error for triangulation

The goal is to estimate the 3-D coordinates of a point  $\mathbf{M}$ , given its projection  $\mathbf{m}_i$  and the camera matrix  $\mathbf{P}_i$  for every view  $i = 1 \dots m$ .

The geometric error associated to a point estimate  $\hat{\mathbf{M}}$  in the  $i$ -th view is the distance between the measured image point  $\mathbf{m}_i$  and the re-projected point  $P_i\hat{\mathbf{M}}$ :

$$\min_{\hat{\mathbf{M}}} \sum_i d(P_i\hat{\mathbf{M}}, \mathbf{m}_i)^2 \quad (76)$$

where  $d()$  is the Euclidean distance between the homogeneous points.

The linear solution is used as a starting point for the iterative minimization (e.g. Gauss-Newton).

## Geometric error for $F$

The goal is to estimate  $F$  given a number of point correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$ .

The geometric error associated to an estimate  $\hat{F}$  is given by the distance of conjugate points from conjugate lines (note the symmetry):

$$\min_{\hat{F}} \sum_j d(\hat{F} \mathbf{m}_\ell^j, \mathbf{m}_r^j)^2 + d(\hat{F}^T \mathbf{m}_r^j, \mathbf{m}_\ell^j)^2 \quad (77)$$

where  $d()$  here is the Euclidean distance between a line and a point (in homogeneous coordinates).

The eight-point solution is used as a starting point for the iterative minimization (e.g. Gauss-Newton).

Note that  $F$  must be suitably parameterized, as it has only seven d.o.f. 11

## Geometric error for H

The goal is to estimate  $H$  given a number of point correspondences  $\mathbf{m}_\ell^i \leftrightarrow \mathbf{m}_r^i$ .

The geometric error associated to an estimate  $\hat{H}$  is given by the symmetric distance between a point and its transformed conjugate:

$$\min_{\hat{H}} \sum_j d(\hat{H} \mathbf{m}_\ell^j, \mathbf{m}_r^j)^2 + d(\hat{H}^{-1} \mathbf{m}_r^j, \mathbf{m}_\ell^j)^2 \quad (78)$$

where  $d()$  is the Euclidean distance between the homogeneous points. This is also called the *symmetric transfer error*.

The linear solution is used as a starting point for the iterative minimization (e.g. Gauss-Newton).

## Bundle adjustment (reconstruction)

If measurements are noisy, the projection equation will not be satisfied exactly by the reconstructed camera matrices and structure.

We wish to minimize the image distance between the re-projected point  $\hat{P}_i \hat{\mathbf{M}}^j$  and measured image points  $\mathbf{m}_i^j$  for every view in which the 3-D point appears:

$$\min_{\hat{P}_i, \hat{\mathbf{M}}^j} \sum_{i,j} d(\hat{P}_i \hat{\mathbf{M}}^j, \mathbf{m}_i^j)^2 \quad (125)$$

where  $d()$  is the Euclidean distance between the homogeneous points.

If the reconstruction is projective  $\hat{P}_i$  is parameterized with its 11 d.o.f. whereas if the reconstruction is Euclidean, one should use  $\hat{P}_i = \hat{K}_i [\hat{R}_i | \hat{\mathbf{t}}_i]$  where the rotation has to be suitably parameterized with 3 d.o.f.

As  $m$  and  $n$  increase, this becomes a very large minimization problem.

However the Jacobian of the residual has a specific structure that can be exploited to gain efficiency.

Primary structure: on the row corresponding to  $\mathbf{m}_i^j$ , only the two elements corresponding to camera  $\hat{P}_i$  and to point  $\hat{M}^j$  are nonzero.

Secondary structure: not all points are seen in all views (data-dependent).

See [53] for a review and a more detailed discussion on bundle adjustment.

## 6.3 Robust estimation

Up to this point, we have assumed that the only source of error affecting correspondences is in the measurements of point's position. This is a small-scale noise that gets averaged out with least-squares.

In practice, we can be presented with *mismatched* points, which are *outliers* to the noise distribution (i.e., rogue measurements following a different, unmodelled, distribution).

These outliers can severely disturb least-squares estimation (even a single outlier can totally offset the least-squares estimation, as illustrated in Fig. 6.)

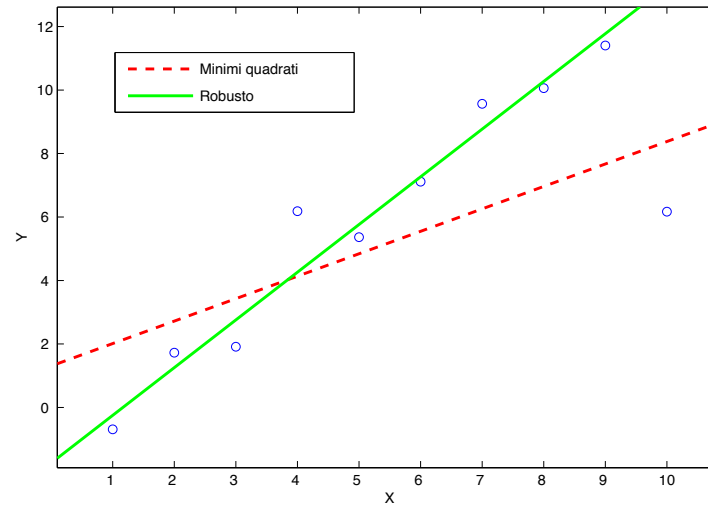


Fig. 6. A single outlier can severely offset the least-squares estimate (red line), whereas the robust estimate (green line) is unaffected.

The goal of robust estimation is to be insensitive to outliers (or at least to reduce sensitivity).

## M-estimators

Least squares:

$$\min_{\theta} \sum_i (r_i/\sigma_i)^2 \quad (80)$$

where  $\theta$  are the regression coefficient (what is being estimated) and  $r_i$  is the residual. M-estimators are based on the idea of replacing the squared residuals by another function of the residual, yielding

$$\min_{\theta} \sum_i \rho(r_i/\sigma_i) \quad (81)$$

$\rho$  is a symmetric function with a unique minimum at zero that grows sub-quadratically, called *loss function*.

Differentiating with respect to  $\theta$  yields:

$$\sum_i \frac{1}{\sigma_i} \rho'(r_i/\sigma_i) \frac{dr_i}{d\theta} = 0 \quad (82)$$

The M-estimate is obtained by solving this system of non-linear equations.



## RANSAC

Given a model that requires a minimum of  $p$  data points to instantiate its free parameters  $\theta$ , and a set of data points  $S$  containing outliers:

1. Randomly select a subset of  $p$  points of  $S$  and instantiate the model from this subset
2. Determine the set  $S_i$  of data points that are within an error tolerance  $t$  of the model.  $S_i$  is the consensus set of the sample.
3. If the size of  $S_i$  is greater than a threshold  $T$ , re-estimate the model (possibly using least-squares) using  $S_i$  (the set of inliers) and terminate.
4. If the size of  $S_i$  is less than  $T$ , repeat from step 1.
5. Terminate after  $N$  trials and choose the largest consensus set found so far.

Three parameters need to be specified:  $t$ ,  $T$  and  $N$ .

Both  $T$  and  $N$  are linked to the (unknown) fraction of outliers  $\epsilon$ .

$N$  should be large enough to have a high probability of selecting at least one sample containing all inliers. The probability to randomly select  $p$  inliers in  $N$  trials is:

$$P = 1 - (1 - (1 - \epsilon)^p)^N \quad (83)$$

By requiring that  $P$  must be near 1,  $N$  can be solved for given values of  $p$  and  $\epsilon$ .

$T$  should be equal to the expected number of inliers, which is given (in fraction) by  $(1 - \epsilon)$ .

At each iteration, the largest consensus set found so far gives a lower bound on the fraction of inliers, or, equivalently, an upper bound on the number of outliers. This can be used to adaptively adjust the number of trials  $N$ .

$t$  is determined empirically, but in some cases it can be related to the probability that a point under the threshold is actually an inlier [10].

As pointed out in [25], RANSAC can be viewed as a particular M-estimator.

The objective function that RANSAC maximizes is the number of data points having absolute residuals smaller than a predefined value  $t$ . This may be seen as minimizing a binary loss function that is zero for small (absolute) residuals, and 1 for large absolute residuals, with a discontinuity at  $t$ .

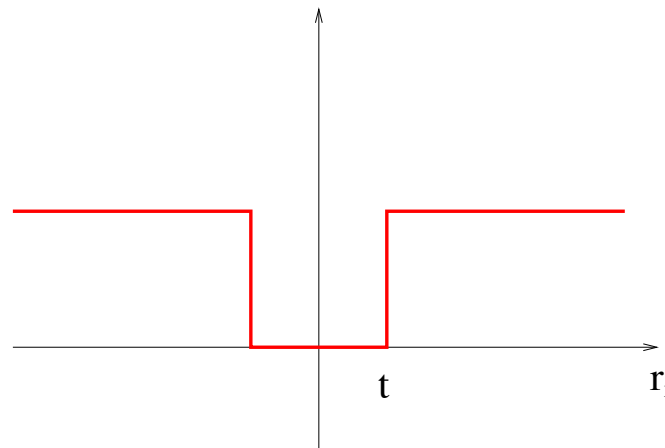


Fig. 7. RANSAC loss function

By virtue of the prespecified inlier band, RANSAC can fit a model to data corrupted by substantially more than half outliers.

However, if more than 50% of the data are outliers they may conspire to produce a smaller total residual than the true inliers.

## LMedS

Another popular robust estimator is the Least Median of Squares. It is defined by:

$$\min_{\theta} \text{med}_i r_i \quad (84)$$

It can tolerate up to 50% of outliers, as up to half of the data point can be arbitrarily far from the “true” estimate without changing the objective function value.

Since the median is not differentiable, a random sampling strategy similar to RANSAC is adopted. Instead of using the consensus, each sample of size  $p$  is scored by the median of the residuals of all the data points. The model with the least median (lowest score) is chosen.

A final weighted least-squares fitting is used.

With respect to RANSAC, LMedS can tolerate “only” 50% of outliers, but requires no setting of thresholds.

## References

- [1] S. Avidan and A. Shashua. Novel view synthesis in tensor space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1034–1040, 1997.
- [2] P. Beardsley, A. Zisserman, and D. Murray. Sequential update of projective and affine structure from motion. *International Journal of Computer Vision*, 23(3):235–259, 1997.
- [3] B. S. Boufama. The use of homographies for view synthesis. In *Proceedings of the International Conference on Pattern Recognition*, pages 563–566, 2000.
- [4] Myron Z. Brown, Darius Burschka, and Gregory D. Hager. Advances in computational stereo. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(8):933–1008, August 2003.
- [5] O. Faugeras. *Three-Dimensional Computer Vision: A Geometric Viewpoint*. The MIT Press, Cambridge, MA, 1993.
- [6] O. Faugeras. Stratification of 3-D vision: projective, affine, and metric representations. *Journal of the Optical Society of America A*, 12(3):465–484, 1994.
- [7] O. Faugeras and Q-T Luong. *The geometry of multiple images*. MIT Press, 2001.
- [8] O. D. Faugeras and L. Robert. What can two images tell us about a third one? In *Proceedings of the European Conference on Computer Vision*, pages 485–492, Stockholm, 1994.
- [9] Paul D. Fiore. Efficient linear solution of exterior orientation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 23(2):140–148, 2001.

- [10] A. Fusiello. Uncalibrated Euclidean reconstruction: A review. *Image and Vision Computing*, 18(6-7):555–563, May 2000.
- [11] A. Fusiello, A. Benedetti, M. Farenzena, and A. Busti. Globally convergent autocalibration using interval analysis. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(12):1633–1638, December 2004.
- [12] A. Fusiello, E. Trucco, and A. Verri. A compact algorithm for rectification of stereo pairs. *Machine Vision and Applications*, 12(1):16–22, 2000.
- [13] Andrea Fusiello. A matter of notation: several uses of the kronecker product in computer vision. Submitted to Pattern Recognition Letters.
- [14] R. Hartley, E. Hayman, L. de Agapito, and I. Reid. Camera calibration and the search for infinity. In *Proceedings of the International Conference on Computer Vision*, 1999.
- [15] R. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2nd edition, 2003.
- [16] R. I. Hartley. In defence of the 8-point algorithm. In *Proceedings of the International Conference on Computer Vision*, pages 1064–1071, Washington, DC, USA, 1995. IEEE Computer Society.
- [17] R. I. Hartley and P. Sturm. Triangulation. *Computer Vision and Image Understanding*, 68(2):146–157, November 1997.
- [18] R.I. Hartley. Theory and practice of projective rectification. *International Journal of Computer Vision*, 35(2):1–16, November 1999.
- [19] A. Heyden. Projective structure and motion from image sequences using subspace methods. In *Scandinavian Conference on Image Analysis*, pages 963–968, 1997.

- [20] A. Heyden. A common framework for multiple-view tensors. In *Proceedings of the European Conference on Computer Vision*, Freiburg, Germany,, 1998.
- [21] A. Heyden. Tutorial on multiple view geometry. In conjunction with ICPR00, September 2000.
- [22] A. Heyden and K. Åström. Euclidean reconstruction from constant intrinsic parameters. In *Proceedings of the International Conference on Pattern Recognition*, pages 339–343, Vienna, 1996.
- [23] A. Heyden and K. Åström. Minimal conditions on intrinsic parameters for Euclidean reconstruction. In *Proceedings of the Asian Conference on Computer Vision*, page XXX, Hong Kong, 1998.
- [24] T.S. Huang and O.D. Faugeras. Some properties of the E matrix in two-view motion estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 11(12):1310–1312, December 1989.
- [25] F. Isgrò and E. Trucco. Projective rectification without epipolar geometry. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1:94–99, Fort Collins, CO, June 23-25 1999.
- [26] F. Isgrò, E. Trucco, P. Kauff, and O. Schreer. 3-D image processing in the future of immersive media. *IEEE Transactions on Circuits and Systems for Video Technology*, 14(3):288–303, 2004.
- [27] S. Ivekovic, A. Fusiello, and E. Trucco. Fundamentals of multiple view geometry. In O. Schreer, P. Kauff, and T. Sikora, editors, *3D Videocommunication. Algorithms, concepts and real-time systems in human centered communication*, chapter 6. John Wiley & Sons, 2005. ISBN: 0-470-02271-X.
- [28] K. Kanatani. *Geometric Computation for Machine Vision*. Oxford University Press, 1993.
- [29] S. Laveau and O. Faugeras. 3-D scene representation as a collection of images and fundamental matrices. Technical Report 2205, INRIA, Institut National de Recherche en Informatique et an Automatique, February 1994.

- [30] Jed Lengyel. The convergence of graphics and vision. *IEEE Computer*, 31(7):46–53, July 1998.
- [31] D. Liebowitz and A. Zisserman. Metric rectification for perspective images of planes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 482–488, 1998.
- [32] C. Loop and Z. Zhang. Computing rectifying homographies for stereo vision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages I:125–131, Fort Collins, CO, June 23-25 1999.
- [33] Q.-T. Luong and T. Viéville. Canonical representations for the geometries of multiple projective views. *Computer Vision and Image Understanding*, 64(2):193–229, 1996.
- [34] Yi Ma, Stefano Soatto, Jana Kosecka, and Shankar S. Sastry. *An Invitation to 3-D Vision*. Springer, November 2003.
- [35] J. R. Magnus and H. Neudecker. *“Matrix Differential Calculus with Applications in Statistics and Econometrics”*. John Wiley & Sons, revised edition, 1999.
- [36] S. Mahamud, M. Hebert, Y. Omori, and J. Ponce. Provably-convergent iterative methods for projective structure from motion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages I:1018–1025, 2001.
- [37] S. J. Maybank and O. Faugeras. A theory of self-calibration of a moving camera. *International Journal of Computer Vision*, 8(2):123–151, 1992.
- [38] P.R.S. Mendonça and R. Cipolla. A simple technique for self-calibration. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages I:500–505, 1999.
- [39] J. Oliensis. Fast and accurate self-calibration. In *Proceedings of the International Conference on Computer Vision*, 1999.



- [40] M. Pollefeys, R. Koch, and L. Van Gool. Self-calibration and metric reconstruction in spite of varying and unknown internal camera parameters. In *Proceedings of the International Conference on Computer Vision*, pages 90–95, Bombay, 1998.
- [41] M. Pollefeys, F. Verbiest, and L. Van Gool. Surviving dominant planes in uncalibrated structure and motion recovery. In *Proceedings of the European Conference on Computer Vision*, pages 837–851, 2002.
- [42] L. Robert, C. Zeller, O. Faugeras, and M. Hébert. Applications of non-metric vision to some visually-guided robotics tasks. In Y. Aloimonos, editor, *Visual Navigation: From Biological Systems to Unmanned Ground Vehicles*, chapter 5, pages 89–134. Lawrence Erlbaum Associates, 1997.
- [43] D. Scharstein and R. Szeliski. A taxonomy and evaluation of dense two-frame stereo correspondence algorithms. *International Journal of Computer Vision*, 47(1):7–42, May 2002.
- [44] A. Shashua and N. Navab. Relative affine structure: Canonical model for 3D from 2D geometry and applications. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 18(9):873–883, September 1996.
- [45] Noah Snavely, Steven M. Seitz, and Richard Szeliski. Modeling the world from internet photo collections. *International Journal of Computer Vision*, 80(2):189–210, 2008.
- [46] C. V. Stewart. Robust parameter estimation in computer vision. *SIAM Review*, 41(3):513–537, 1999.
- [47] P. Sturm and B. Triggs. A factorization based algorithm for multi-image projective structure and motion. In *Proceedings of the European Conference on Computer Vision*, pages 709–720, Cambridge, UK, 1996.
- [48] R. Szeliski. Video mosaics for virtual environments. *IEEE Computer Graphics and Applications*, 16(2):22–30, March 1996.
- [49] C. Tomasi and T. Kanade. Shape and motion from image streams under orthography – a factorization method. *International Journal of Computer Vision*, 9(2):137–154, November 1992.

- [50] Philip H. S. Torr and Andrew Zisserman. Robust computation and parametrization of multiple view relations. In *ICCV*, pages 727–732, 1998.
- [51] B. Triggs. Factorization methods for projective structure from motion. In *CVPR*, pages 845–851, 1996.
- [52] B. Triggs. Autocalibration and the absolute quadric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 609–614, Puerto Rico, 1997.
- [53] Bill Triggs, Philip F. McLauchlan, Richard I. Hartley, and Andrew W. Fitzgibbon. Bundle adjustment - a modern synthesis. In *Proceedings of the International Workshop on Vision Algorithms*, pages 298–372. Springer-Verlag, 2000.
- [54] E. Trucco and A. Verri. *Introductory Techniques for 3-D Computer Vision*. Prentice-Hall, 1998.
- [55] Cha Zhang and Tsuhan Chen. A survey on image-based rendering - representation, sampling and compression. Technical Report AMP 03-03, Electrical and Computer Engineering - Carnegie Mellon University, Pittsburgh, PA 15213, June 2003.
- [56] Z. Zhang. Determining the epipolar geometry and its uncertainty: A review. *International Journal of Computer Vision*, 27(2):161–195, March/April 1998.
- [57] A. Zisserman. Single view and two-view geometry. Handout, EPSRC Summer School on Computer Vision, 1998. available from [http://www.dai.ed.ac.uk/CVonline/LOCAL\\_COPIES/EPSRC\\_SSAZ/epsrc\\_ssz.html](http://www.dai.ed.ac.uk/CVonline/LOCAL_COPIES/EPSRC_SSAZ/epsrc_ssz.html).

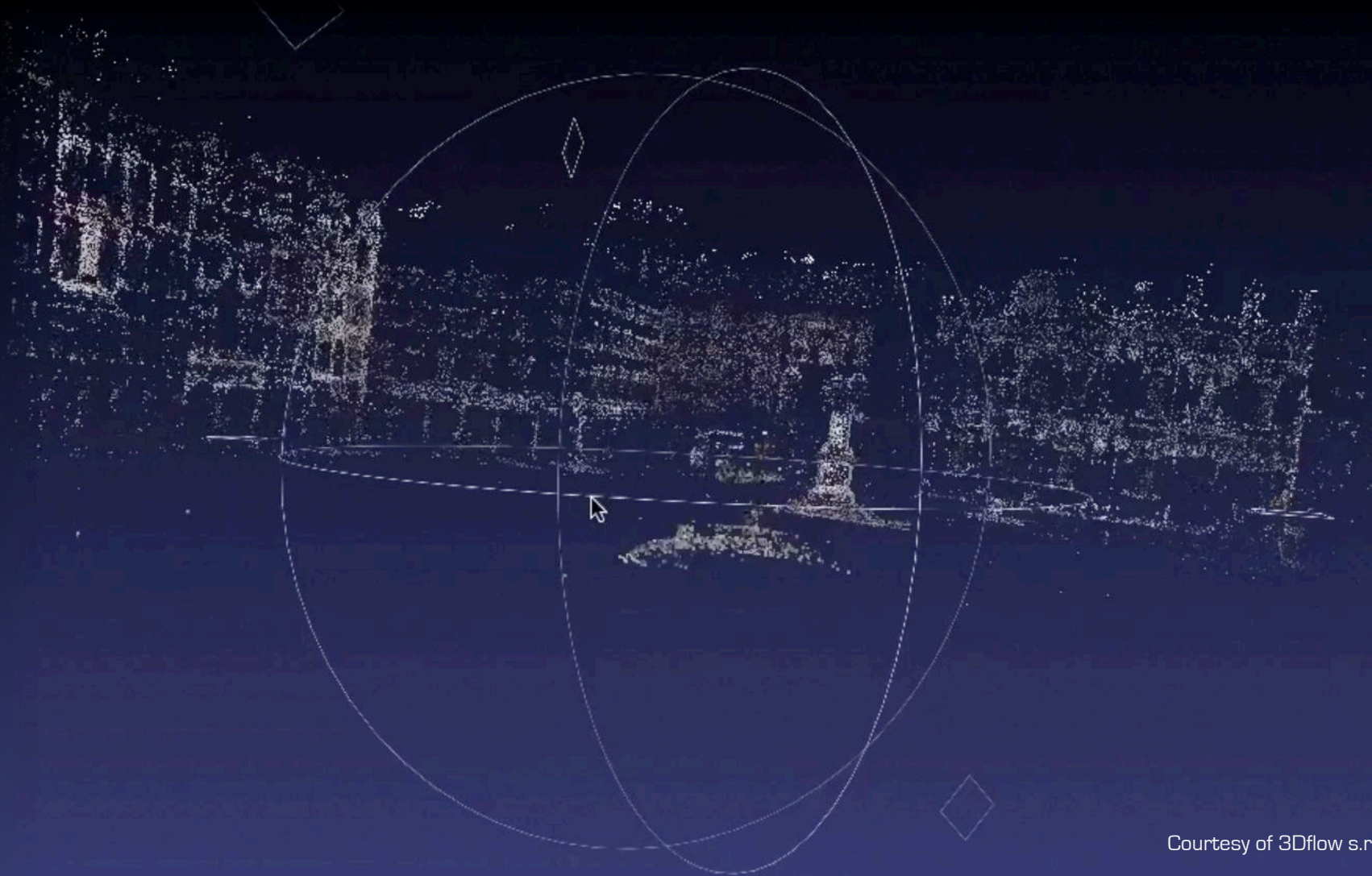


## Structure from motion (in practice)

Where we describe a working pipeline

Credits to R. Toldo (3Dflow s.r.l.) , R. Gherardi and M.Farenzena (Univ. Verona)

# Piazza delle erbe, Udine ~200 views

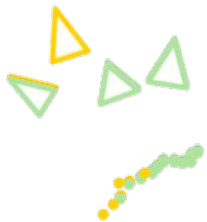


# Approaches to structure from motion

- ❑ Resection-intersection cycle (use 2-views tools, incremental)
- ❑ Global rotation first (2-views tools, starts from a network of E matrices, solve for R of each camera then compute t) [Govindu 11]
- ❑ Multiple-view geometry: generalizes 2-views but needs all points seen in all views, in principle) [Sturm, Triggs 96]

# Sequential structure from motion

- 1. Match or track points over the whole image sequence
- 2. Initialize the structure and motion recovery
  - (a) Select two views that are suited for initialization;
  - (b) Solve **relative orientation** and set up the initial frame;
  - (c) Reconstruct the initial structure (**intersection**)
- 3. For every additional view,
  - (a) Infer matches to the existing 3D structure,
  - (b) Compute the camera orientation (**exterior orientation**) using a robust algorithm
  - (c) Refine the existing structure (**intersection**);
  - (d) Initialize new structure points (**intersection**);
  - (e) Refine the structure and motion through BA.



Resection-intersection cycle

# Sequence... who said sequence?

- ❑ We assumed images come in a sequence, but usually they are unordered.
- ❑ We need to:
  - ❑ Discover putative matching images without doing the actual match
    - ❑ Based on image content
  - ❑ Sequence the images (after matching):
    - ❑ Define the seed pair (critical!)
    - ❑ Define next view to process based on match

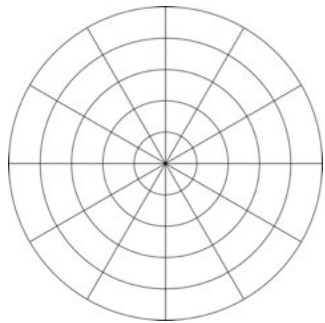
# Pre-processing summary

- Keypoint extraction
- Matching - broad phase: select  $O(n)$  views to be matched
- Matching - narrow phase: match keypoints between pair
- Sequencing: determine processing order (can be on-line)



# Keypoint extraction

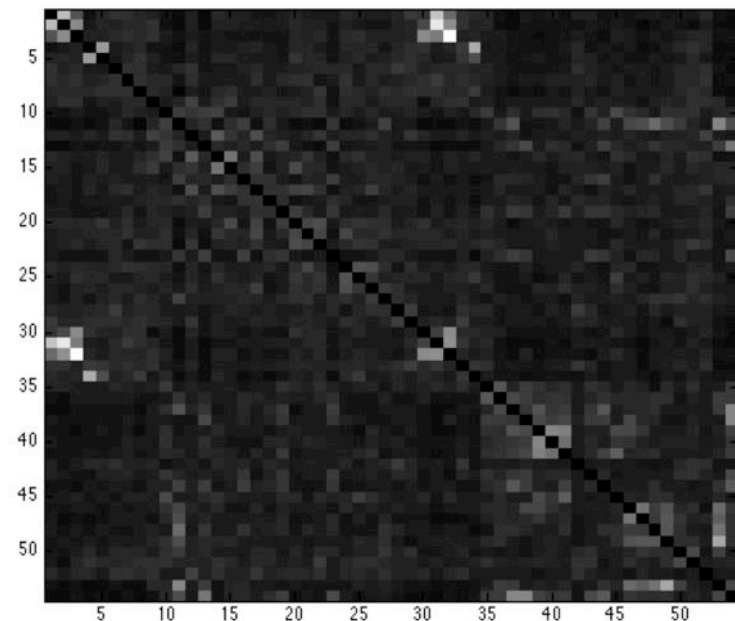
- ❑ Detector: scale-space extrema of the scale-normalized Laplacian [T. Lindeberg, 1994]. SIFT is an variation on this theme.
- ❑ Eg. , we used a 8-level scale-space and in each level the Laplacian is computed by convolution (in CUDA) with a  $3 \times 3$  kernel.
- ❑ Key: multiresolution pyramid based on derivative operator (LoG, DoG, DoH).



- ❑ Descriptor: 128- dimensional radial descriptor based on the accumulated response of steerable derivative filters (similar to GLOH and SIFT).
- ❑ Key: derivatives, directions histogram,
- ❑ Important details: subpixel estimation in scale-space, votes with radial weights and interpolation.

# Which images are to be matched?

- ❑ Recover the **image graph**, i.e., the graph that tells which image overlaps (or can be matched) with which other.
- ❑ For each key-point descriptor its approximate  $k$  nearest neighbours in feature space are computed (via ANN)
- ❑ A 2D histogram is then built: increment  $\text{bin}(i,j)$  whenever a keypoint of image  $i$  has a keypoint of image  $j$  in its  $k$ -neighbourhood;



# Which images are to be matched?

- Consider the complete weighted graph  $G = (V, E)$  where  $V$  are views and the weighted adjacency matrix is the 2D histogram.
- This graph has  $|V| = O(n^2)$ . The objective is to extract a subgraph  $G'$  with a number of edges that is linear in  $n$ .
- Lowe's approach: every image is connected (node) to the  $m$  ( $=8$ ) images that have the greatest number of keypoints matches in common. This creates a graph with  $mn = O(n)$  edges, being  $m$  constant.
- When the number of images is large, however, this tends to create cliques of very similar images with weak inter-cliques connections

# Which images are to be matched?

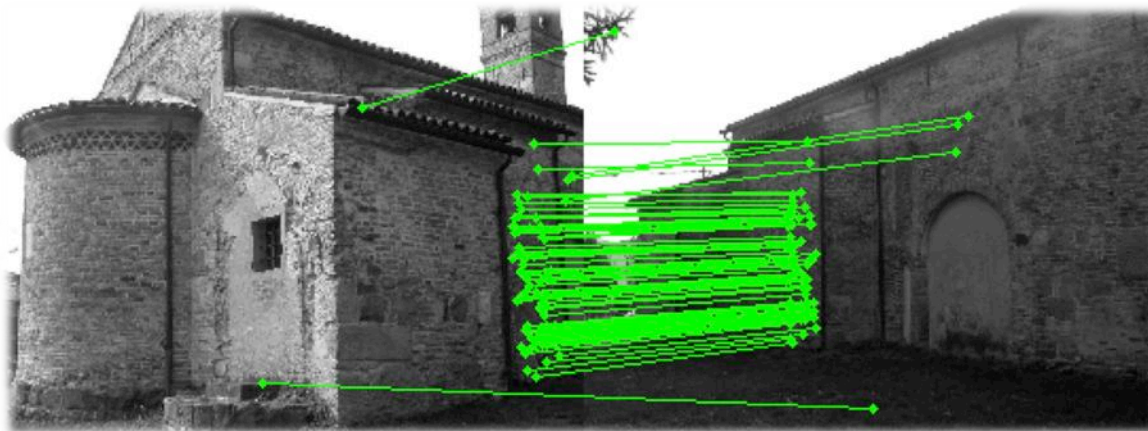
- ❑ In graph theory, a graph is  $k$ -edge-connected if it remains connected whenever fewer than  $k$  edges are removed.
- ❑ The graph produced by the original approach has a low  $k$ , while one would like to have  $k$  as high as possible (ideally  $k = m$ ).
- ❑ We devised a strategy that builds a subgraph  $G'$  of  $G$  which is  $m$ -edge-connected by construction.
  - ❑ Build the maximum spanning tree of  $G$ : (has  $n - 1$  edges);
  - ❑ remove them from  $G$  and add them to  $G'$
  - ❑ repeat  $m$  times.
- ❑ The resulting graph has  $(n-1)m = O(n)$  edges and is  $m$ -edge-connected

# Matching and model selection

- Match SIFT keypoints (images connected in the graph)
- Compute F and H with RANSAC
- Non linear refinement
- Compute GRIC-F and GRIC-H

$$\text{GRIC} = \sum \rho(e_i^2) + nd \log(r) + k \log(rn)$$

$$\rho(e) = \min \left( \frac{e^2}{\sigma^2}, 2(r - d) \right)$$



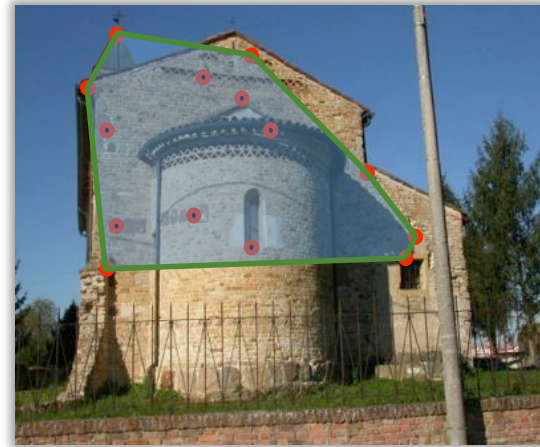
# Tracks

- From pairwise matches to tracks
  - a lot of bookkeeping...



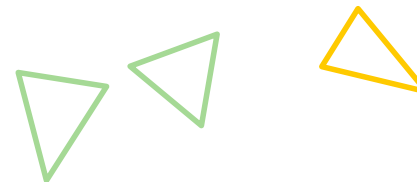
# Sequencing

- Initial pair selection.
- Score based on:
  - # features in common
  - god image coverage
  - good geometry (F vs H)



- Next view:
  - # visible points of the partial reconstruction

$$S_{i,j} = \frac{CH_i}{A_i} + \frac{CH_j}{A_j} + \frac{\text{gric}(F_{i,j})}{\text{gric}(H_{i,j})}$$



Many equivalent heuristics can be devised!

# Sequential Structure from Motion

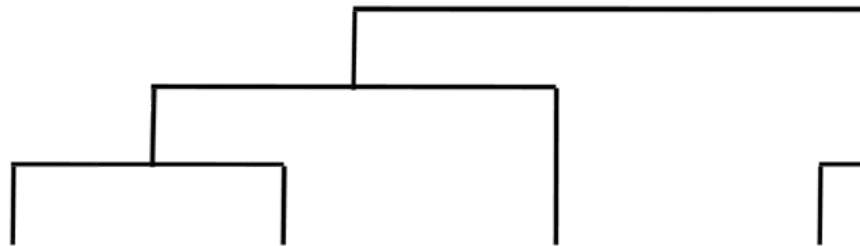


Sequential  
paradigm

- ⊙ Incremental addition, one view at a time
- ⊙ Always grows a single reconstruction



# Hierarchical Structure from Motion



- ✓ Easily parallelizable
- ✓ No initial pair dependency
- ✓ Less drift, error containment
- ✓ Provably lower complexity
- ✓ Graceful failure



Hierarchical  
paradigm

- ⊙ Leaves correspond to images
- ⊙ Internal nodes are partial reconstructions
- ⊙ Root node is the final reconstruction

# Tree construction

- Agglomerative, hierarchical clustering problem

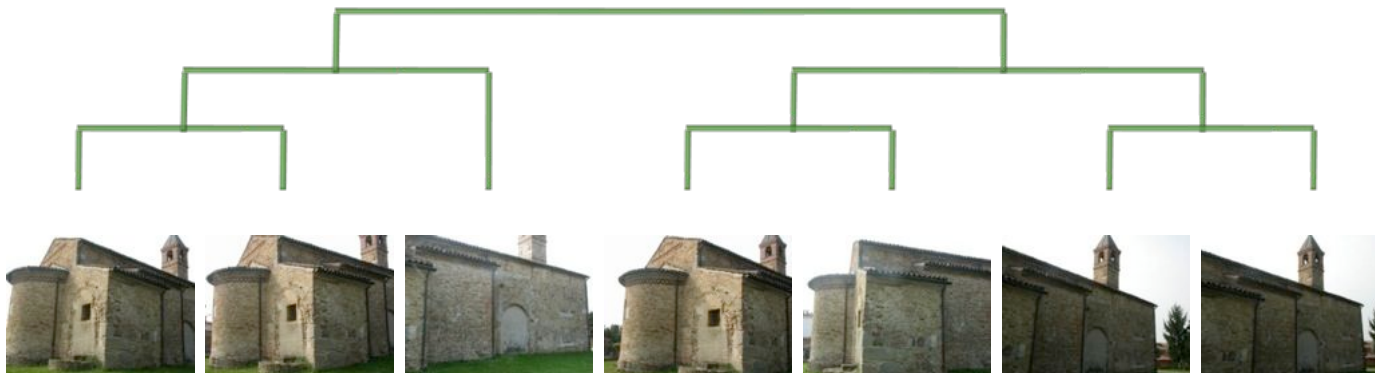
- Simple (complete, average, ward's) linkage

- Distance needed:

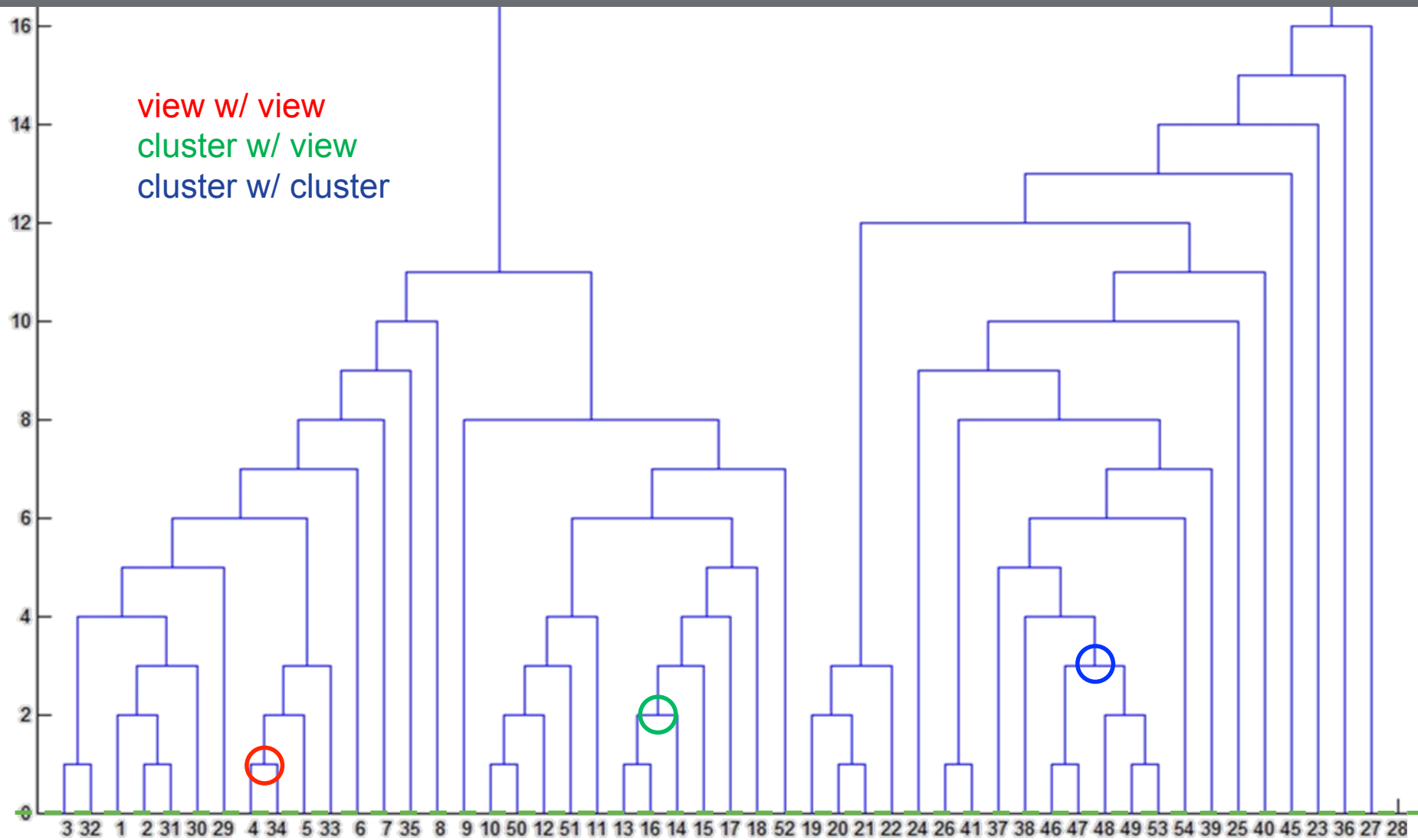
- Common matches  $a_{i,j} = \frac{1}{2} \frac{|S_i \cap S_j|}{|S_i \cup S_j|} + \frac{1}{2} \frac{CH(S_i) + CH(S_j)}{A_i + A_j}$

- Good coverage

- Pairs only:  $\text{gric}(F_{i,j}) < \alpha \text{gric}(H_{i,j})$  with  $\alpha \geq 1$

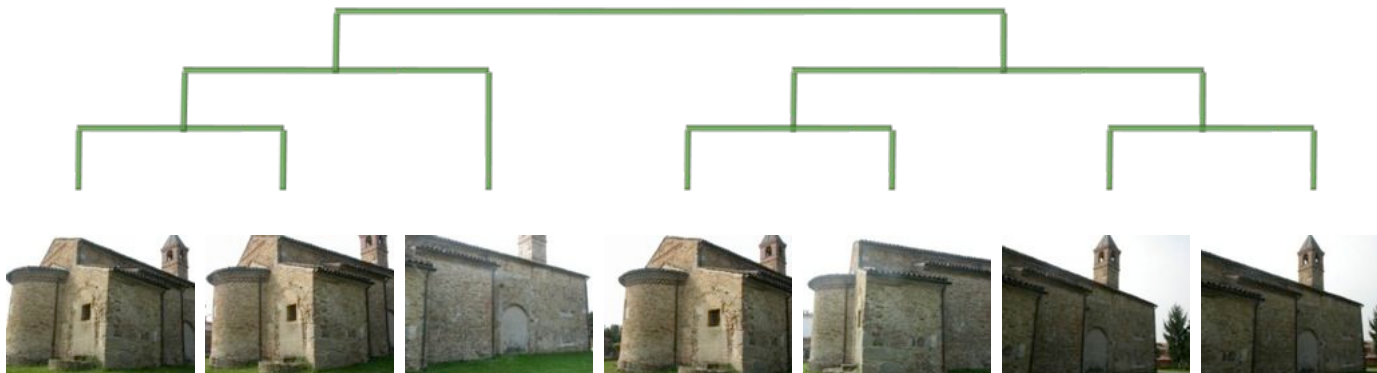


# Tree traversal



# Improving on hierarchical SfM

- ❑ Local bundle adjustment
- ❑ Balanced tree
- ❑ Autocalibration
- ❑ Intrinsic locking

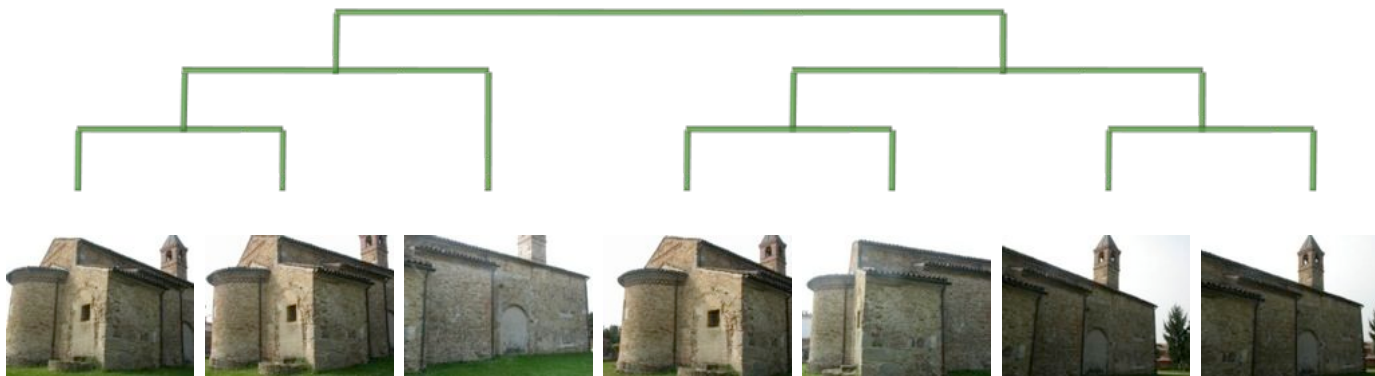


## Local BA

- Consider two clusters  $A$  and  $B$ , where  $|A| < |B|$ .
- Merge by transforming  $A$  onto  $B$ .
- The bundle adjustment involves all the views of  $A$  and the subset  $B'$  of views of  $B$  that share some tracks with  $A$
- As a consequence, the points not visible by views in  $A \cup B'$  are excluded from bundle adjustment.

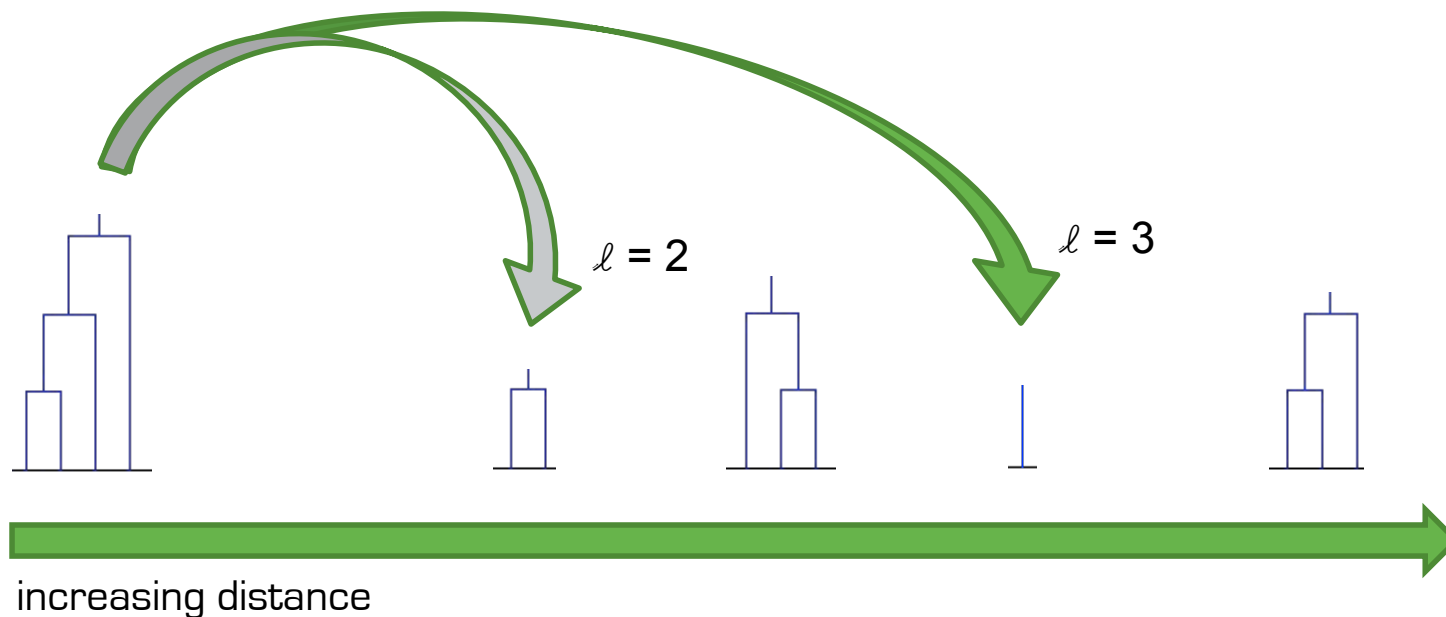
# Balancing the tree

- ❑ Hierarchical SfM is less computationally demanding than the sequential one by one order of magnitude (best case)
- ❑ Complexity analysis assumed the tree perfectly balanced
- ❑ Worst case is equivalent to the sequential paradigm
- ❑ Can we obtain a balanced tree, approximating best complexity?

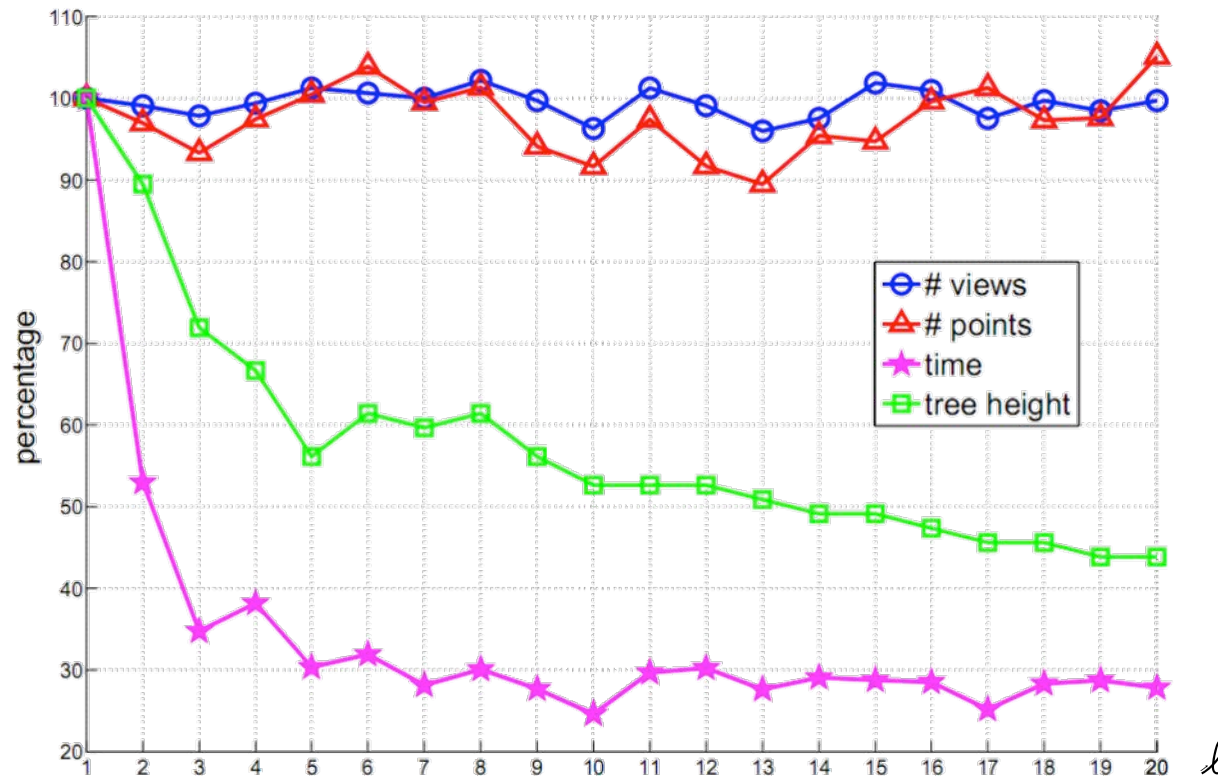


# Balancing the tree

- Agglomerative clustering: **closest first**
- Balanced tree: **smallest first**
- Idea: prefer smaller clusters in a distance neighbourhood of size  $\ell$



# Balancing the tree



- Points and cameras unchanged
- Running time more than halved



# Practical autocalibration

- Given a guess on the intrinsic parameters of two cameras estimate (in closed form) a consistent upgrading collineation. This yields an estimate of all cameras but the first.
- Score the intrinsic parameters of these  $n-1$  cameras based on the likelihood of skew, aspect ratio and principal point.
- Iterate through all possible (finite) combinations of intrinsic parameters looking for the best upgrading homography

## Estimation of the plane at infinity

- Given a projective reconstruction in canonical form:

$$P_1 = [I \mid \mathbf{0}] \quad P_2 = [Q_2 \mid \mathbf{q}_2] \quad H = \begin{bmatrix} K_1 & \mathbf{0} \\ \mathbf{v}^\top & \lambda \end{bmatrix}$$

- The Euclidean perspective projection matrices are:

$$P_1^E = [K_1 \mid \mathbf{0}] \simeq P_1 H$$

$$P_2^E = K_2 [R_2 \mid \mathbf{t}_2] \simeq P_2 H = [Q_2 K_1 + \mathbf{q}_2 \mathbf{v}^\top \mid \lambda \mathbf{q}_2]$$

- Hence  $R_2$  is the sum of a 3 by 3 matrix and a rank 1 term:

$$R_2 \simeq K_2^{-1} (Q_2 K_1 + \mathbf{q}_2 \mathbf{v}^\top) = K_2^{-1} Q_2 K_1 + \mathbf{t}_2 \mathbf{v}^\top$$

## Estimation of the plane at infinity

- There always exist  $R^*$  such that:  $R^* \mathbf{t}_2 = [\|\mathbf{t}_2\| \ 0 \ 0]^T$

$$R^* R_2 \simeq \overbrace{R^* K_2^{-1} Q_2 K_1}^W + [\|\mathbf{t}_2\| \ 0 \ 0]^T \mathbf{v}^T$$

$$R^* R_2 = \begin{bmatrix} \mathbf{w}_1^T + \|\mathbf{t}_2\| \mathbf{v}^T \\ \mathbf{w}_2^T \\ \mathbf{w}_3^T \end{bmatrix} / \|\mathbf{w}_3\|$$

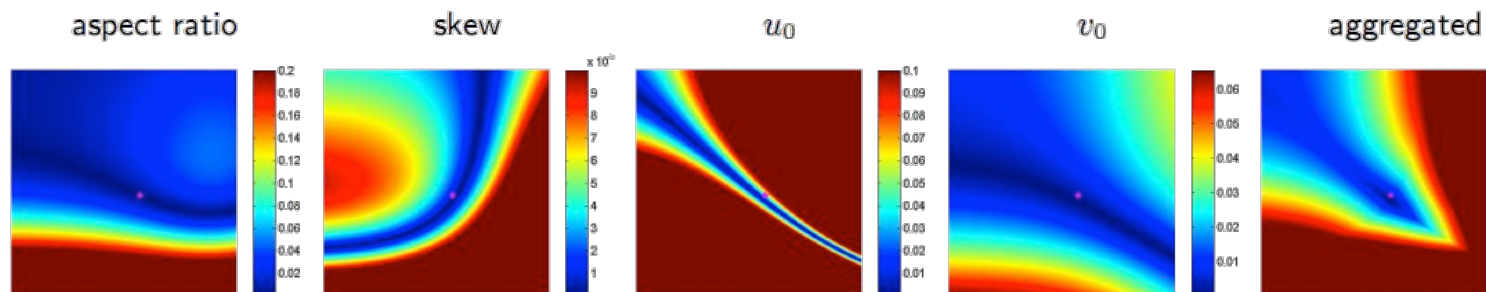
$$\mathbf{v} = (\mathbf{w}_2 \times \mathbf{w}_3 / \|\mathbf{w}_3\| - \mathbf{w}_1) / \|\mathbf{t}_2\|$$

- From  $\mathbf{v}$  the upgrading collineatin H is obtained

# Scoring rectifying homographies

- We score each sample  $(f_1, f_2)$  with a cost function based on the expectation of intrinsic parameters of the transformed cameras

$$\mathcal{C}(K) = \overbrace{w_{sk} |k_{1,2}|}^{\text{skew}} + \overbrace{w_{ar} |k_{1,1} - k_{2,2}|}^{\text{aspect ratio}} + \overbrace{w_{u_o} |k_{1,3}| + w_{v_o} |k_{2,3}|}^{\text{principal point}}$$



- Plots shows very clear valleys, aggregated cost has a unambiguous minima (even from 2 images!)

# Exhaustive search

- Space of internal parameters is inherently bounded by the finiteness of the acquisition devices

$$0.3 (W + H) < f < 3 (W + H)$$

$$sk \approx 0$$

$$ar \approx 1$$

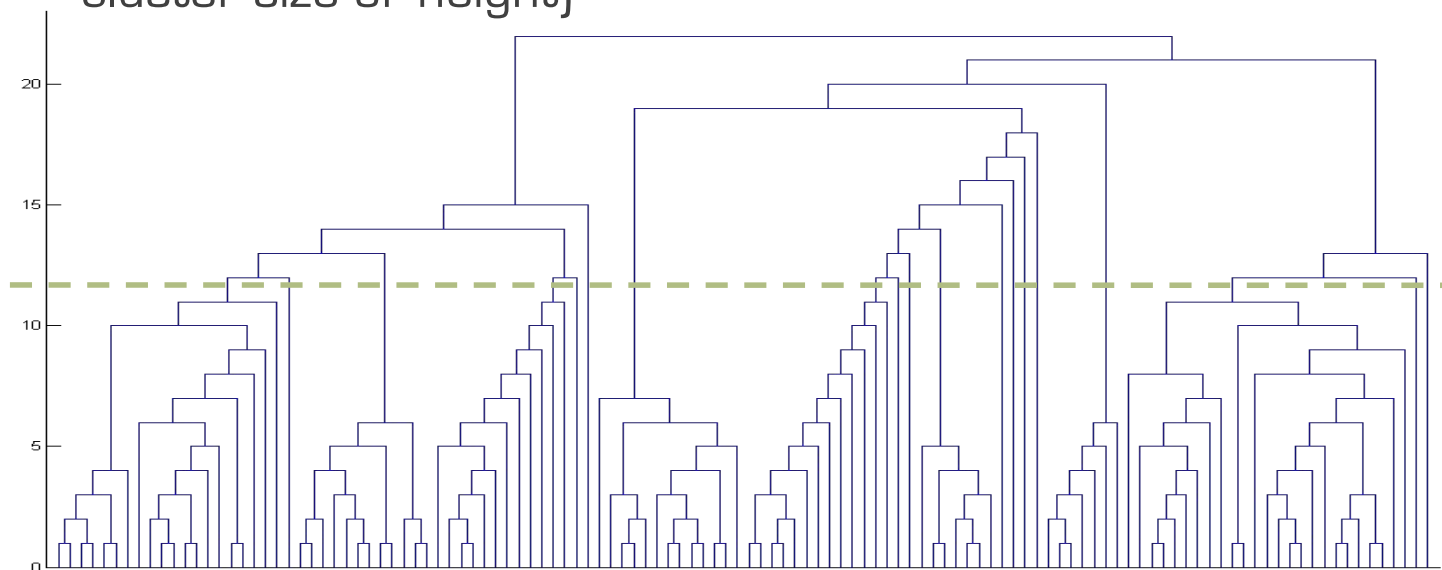
$$W / 2 - 0.1 W < u_0 < W / 2 + 0.1 W$$

$$H / 2 - 0.1 H < v_0 < H / 2 + 0.1 H$$

- Approximating  $(u_0, v_0)$  with the image centre, the search space is reduced to a bounded region of  $R^2$

# Intrinsics locking

- ▣ Bundle adjustment optimizes internal camera parameters
- ▣ Idea: lock the intrinsic parameters after stabilization (threshold on cluster size or height)



## Things we learnt

- Use good features: SIFT are the best ones. Scale-normalized Laplacian works as good as; Descriptors based on i) directional histograms of ii) derivatives work as good as BUT careful tuning is crucial.
- Be selective on accepting matches, reject when uncertain (you will be able to recover uncertain point later). The use of tracks consistency is crucial (not just pairwise matching)
- Error containment: BA is great but has local convergence. Since you never know when you are in the convergence basin you should always keep drift under control: refine geometric objects with proper geometric error minimization and do BA as often as possible.
- Autocalibration: it works, however better results when subsets of images with the same internal parameters are singled out.

## References (on hierarchical SfM)

- R. Gherardi, M. Farenzena, A. Fusiello. Improving the Efficiency of Hierarchical Structure-and-Motion. CVPR 2010, San Francisco, USA, June 13-18, 2010.
- M. Farenzena, A. Fusiello, R. Gherardi. Structure-and-motion pipeline on a hierarchical cluster tree. In 3DIM'09, ICCV Workshops, pages 1489–1496, Kyoto, Japan, October 3-4 2009.
- M. Farenzena, A. Fusiello, R. Gherardi, and R. Toldo. Towards unsupervised reconstruction of architectural models. In VMV 2008, pages 41–50, Konstanz, DE, October 8-10 2008. IOS Press
- Code (author: R. Toldo) is available at: [samantha.3dflow.net](http://samantha.3dflow.net)