

University of Bologna

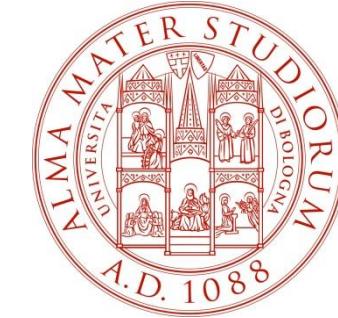


Image Formation and Acquisition (Part 3)

Luigi Di Stefano (luigi.distefano@unibo.it)

Camera Calibration (1)

- We have now at disposal a camera model thoroughly explaining the image formation and digitization process.
- Such a model includes the PPM, which, in turn, can be decomposed into 3 independent items: intrinsic parameter matrix (A), rotation matrix (R), translation vector (T), as well as the lens distortion parameters.
- **Camera calibration** is the process whereby all parameters defining the camera model are - as accurately as possible- estimated for a specific camera device. Depending on the application, either the PMM only or also its independent components (A , R , T) need to be estimated.

$$\tilde{P} = [P_1 \ P_2 \ P_3 \ P_4] \quad P = A[1 \ 1 \ 0]G \rightarrow \text{estimate}$$

- Many camera calibration algorithms do exist. The basic process, though, relies always on setting up a linear system of equations given a set of known 3D-2D correspondences, so as to then solve for the unknown camera parameters.
- To obtain the required correspondences specific physical objects (referred to as calibration targets) providing easily detectable image features (such as, e.g., chessboard or dot patterns) are typically deployed.

$$\tilde{m}_j^{2D} = \underbrace{A[R \mid T]}_{\begin{array}{l} 6 \text{ DOF} \\ + \\ 4 \text{ cam.} \\ \text{param}(A) \end{array}} \tilde{m}_j^{3D}$$

↓
pixels

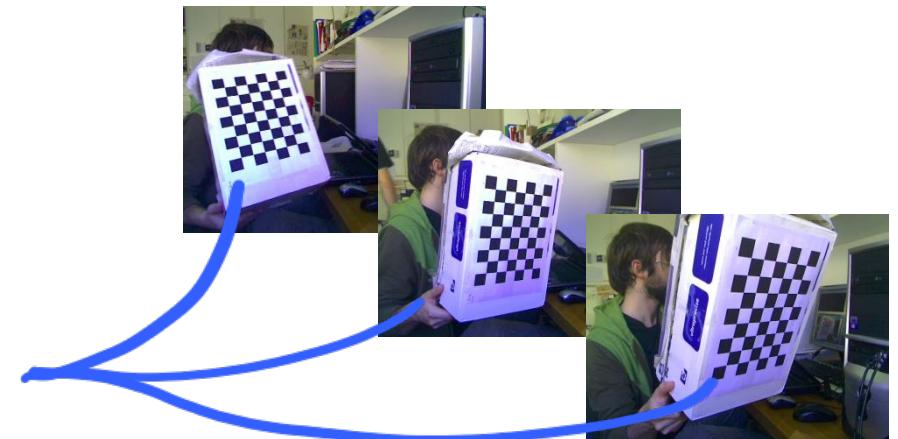
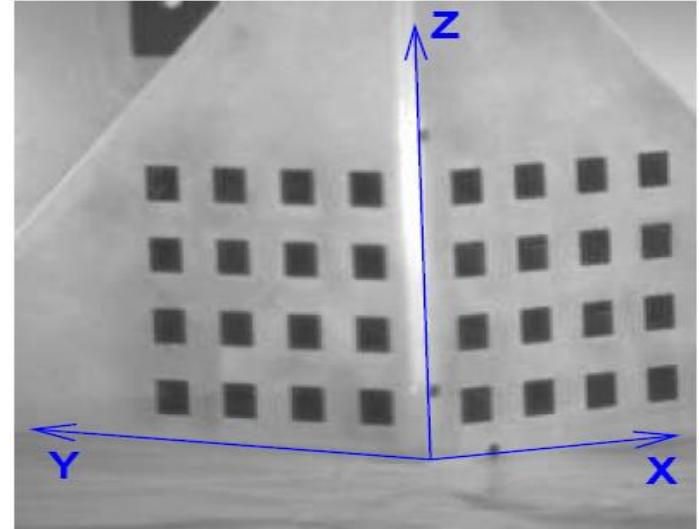
↓
World
ref. frame

to estimate the parameter, we
 need to collect much \tilde{m}_j
 of an object that we're already
 measured

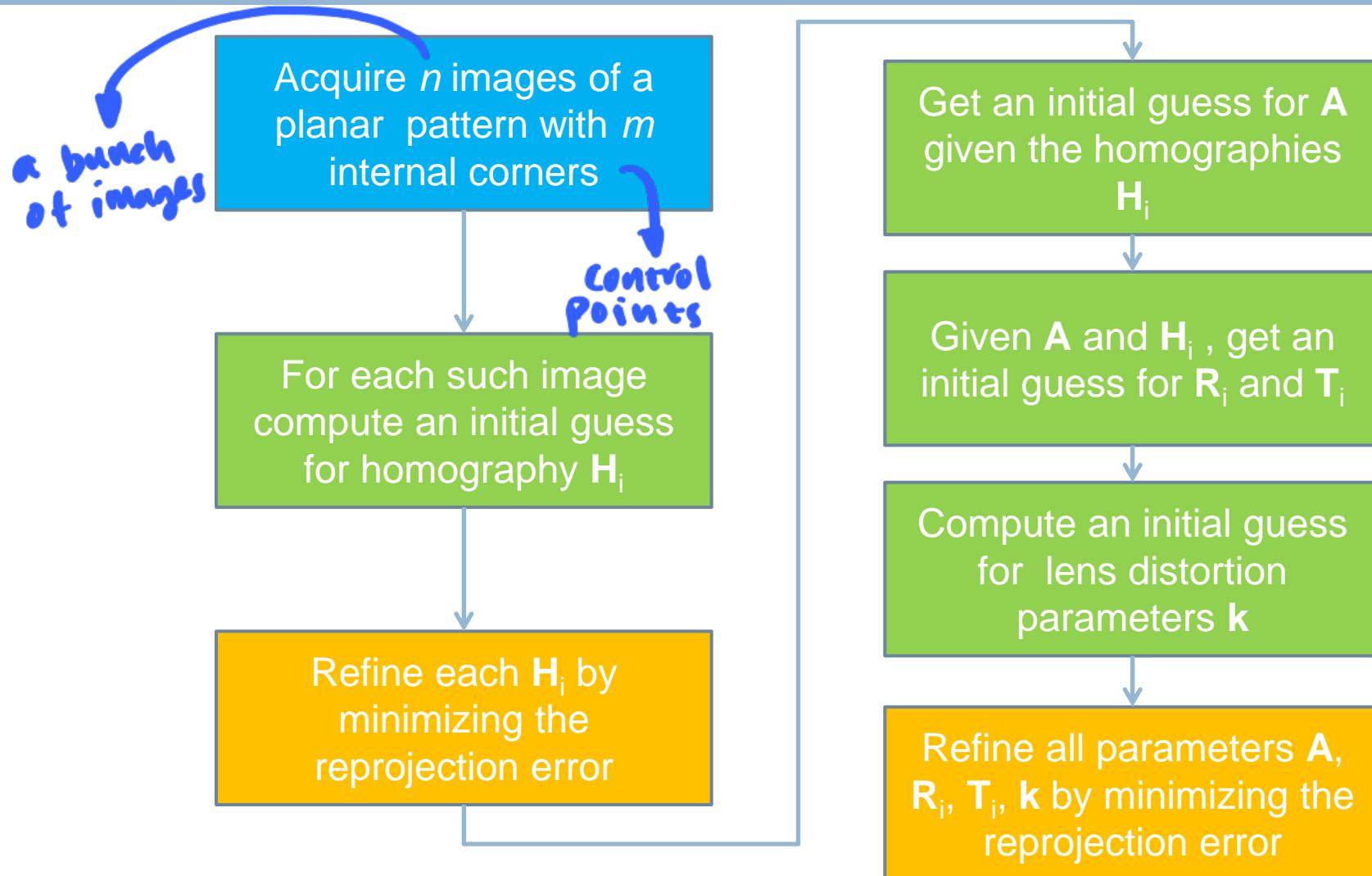
Camera Calibration (2)

- Camera calibration approaches can be split into two main categories:
 - Those relying on a **single image** featuring **several** (at least 2) **planes** containing a known patterns.
 - Those relying on **several** (at least 3) different **images** of **one given planar pattern**.
- In practise, it is difficult to build accurate targets containing multiple planes, while an accurate planar target can be attained rather easily.
- Implementing a camera calibration software requires a significant effort. However, the main Computer Vision toolboxes include specific functions (OpenCV, Matlab CC Toolbox, Halcon,...)

internal corners
as control points



Zhang's Method (1)



Zhang's Method (2)



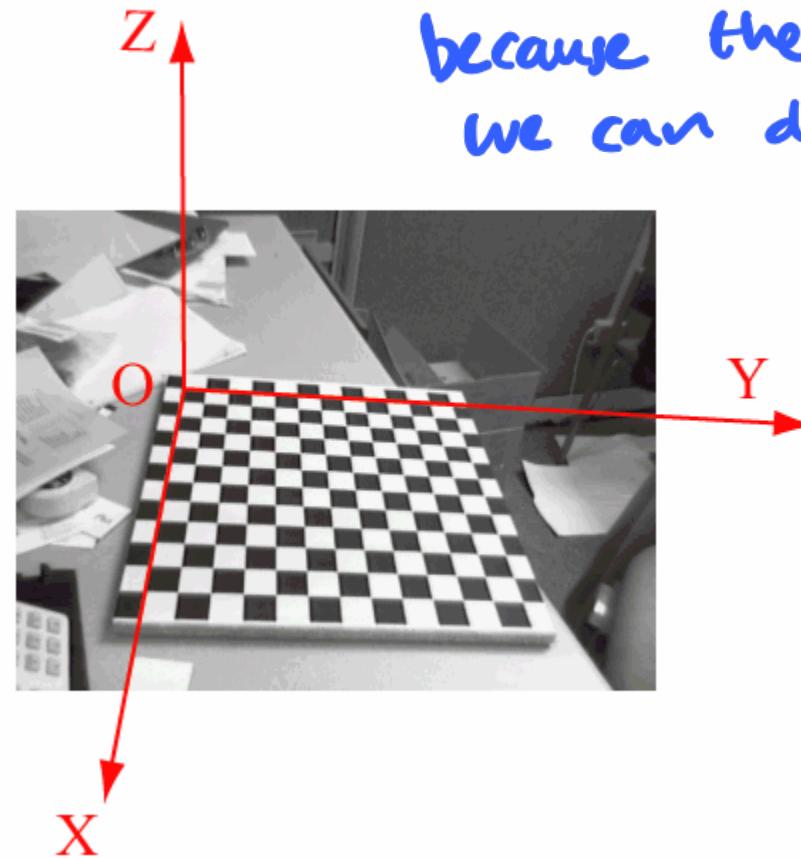
- Given a planar chessboard pattern, known are:
 - The number of internal corners of the pattern, different along the two orthogonal directions for the sake of disambiguation (i.e. rows, columns).
 - The size of the squares which form the pattern.
- Internal corners can be detected easily by standard algorithms (e.g. the Harris corner detector, possibly with sub-pixel refinement for improved accuracy).
- In each image, the 3D WRF is taken at the top-left corner of the pattern, with plane $z=0$ given by the pattern itself and the x,y axis aligned to the two orthogonal directions, in particular so as to keep always the same association between axis and directions (e.g. $x=\text{rows}$, $y = \text{columns}$). Thus, as for 3D points:
 - The third coordinate is always 0.
 - x and y are determined by the known size of the squares forming the chessboard.

Estrinsic Parameters

- It is worth pointing out that each image requires its own estimate of the extrinsic parameters, as they are different from one to the other.
 - A global WRF can be taken to coincide with that associated with one of the calibration images, e.g. the first.

Needed :

- 3D coordinates of control points represented
- pixel coordinates in the image



because the target is planar,
we can do the homography

Planar Calibration Target → Homographies

- Due to the calibration pattern being planar and the choice of the WRF associated with the calibration images, for each of them we can consider only 3D control points (i.e. corners) with $z=0$. Therefore, the mapping between the 3D coordinates of the control points and their projections in the calibration images is a **homography**:

$$k\tilde{\mathbf{m}} = \tilde{\mathbf{P}} \tilde{\mathbf{W}} \begin{bmatrix} p_{1,1} & p_{1,2} & p_{1,3} & p_{1,4} \\ p_{2,1} & p_{2,2} & p_{2,3} & p_{2,4} \\ p_{3,1} & p_{3,2} & p_{3,3} & p_{3,4} \end{bmatrix} \begin{bmatrix} x \\ y \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} p_{1,1} & p_{1,2} & p_{1,4} \\ p_{2,1} & p_{2,2} & p_{2,4} \\ p_{3,1} & p_{3,2} & p_{3,4} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \mathbf{H}\tilde{\mathbf{w}}'$$

$k\tilde{\mathbf{m}} = \mathbf{H}\tilde{\mathbf{w}}'$

pixel coord. of corner *corner coord. of corner* *vector of coordinates of the corner* *homography (estimate)*

- Given a pattern with m corner, we can write m systems of 3 linear equations as above, wherein both 3D (with $z=0$) as well as 2D coordinates are known due to corners having been detected in the calibration image and the unknowns being, therefore, the elements in \mathbf{H} .

Estimating H (DLT Algorithm)

- Starting from the previous homography equation we can write:

$k\tilde{m} = H\tilde{w}' \Rightarrow \tilde{m} \times H\tilde{w}' = 0 \Rightarrow$
 $\begin{bmatrix} v\mathbf{h}_3^T \tilde{\mathbf{w}}' - \mathbf{h}_2^T \tilde{\mathbf{w}}' \\ \mathbf{h}_1^T \tilde{\mathbf{w}}' - u\mathbf{h}_3^T \tilde{\mathbf{w}}' \\ u\mathbf{h}_2^T \tilde{\mathbf{w}}' - v\mathbf{h}_1^T \tilde{\mathbf{w}}' \end{bmatrix} = 0, H = \begin{bmatrix} \mathbf{h}_1^T \\ \mathbf{h}_2^T \\ \mathbf{h}_3^T \end{bmatrix}$

$\begin{bmatrix} \mathbf{0}^T & -\tilde{\mathbf{w}}'^T & v\tilde{\mathbf{w}}'^T \\ \tilde{\mathbf{w}}'^T & \mathbf{0}^T & -u\tilde{\mathbf{w}}'^T \\ -v\tilde{\mathbf{w}}'^T & u\tilde{\mathbf{w}}'^T & \mathbf{0}^T \end{bmatrix} \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{bmatrix} = Ah = 0$

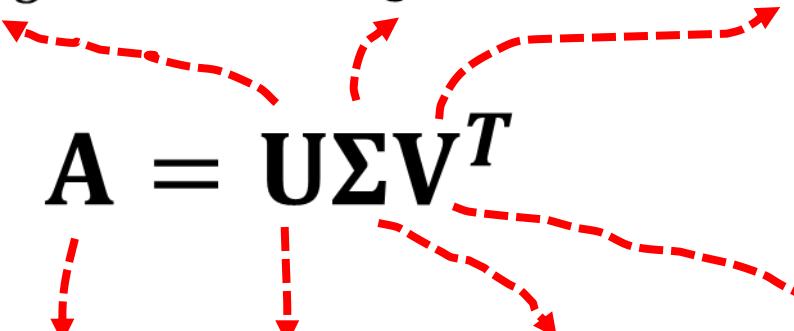
$\left[\begin{array}{ccccccccc} h_{11} & h_{12} & \dots & & & & & & \\ h_{21} & \dots & \dots & & & & & & \\ h_{31} & \dots & \dots & & & & & & \\ i & j & k & & & & & & \\ u & v & 1 & & & & & & \\ h_1^{T\sim} & h_2^{T\sim} & h_3^{T\sim} & & & & & & \\ h_1^{T\sim} w' & h_2^{T\sim} w' & h_3^{T\sim} w' & & & & & & \end{array} \right]$

Estimating homography in one calibr. image

- Between the previous 3 equations in 9 unknowns, only 2 are linearly independent, so that typically only the first 2 are kept. By deploying all correspondences, for each image we build a linear system of $2m$ equations in 9 unknowns. This allows for estimating H by minimizing the “algebraic error” represented by the norm of vector Ah : $\mathbf{h}^* = \operatorname{argmin}_{\mathbf{h} \in R^9} \|Ah\|$, subject to $\|\mathbf{h}\| = 1$ to avoid the trivial solution ($\mathbf{h} = 0$).
- It is known from *linear algebra* that the solution of the above over-constrained system of equations can be obtained in closed form by the Singular Value decomposition (SVD) of matrix A .

Estimating H by the SVD

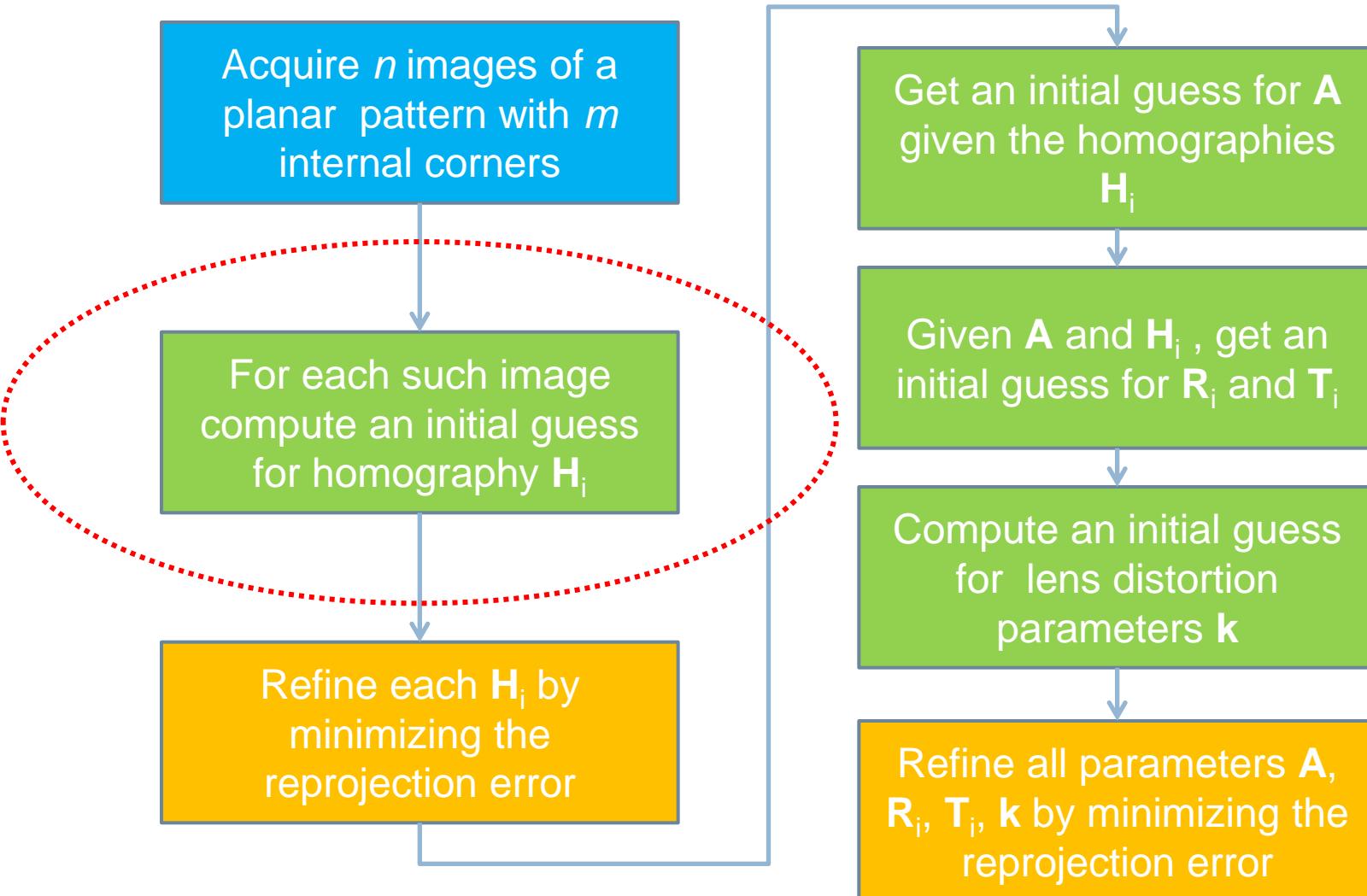
$$A = U \Sigma V^T$$

Orthogonal Diagonal Orthogonal

 2m × 9 2m × 2m 2m × 9 9 × 9

$$V = [v_1 \dots v_9] \rightarrow h^* = v_9$$

See Appendix 1 for the solution of the homography estimation problems in case only 4 correspondences are available (i.e. the minimum number of required correspondences).

Zhang's Method



Non-linear Refinement



- Given the previous initial estimation, \mathbf{H} is later refined by the least-squares solution of the non-linear minimization problem:

$$\mathbf{H} = \operatorname{argmin}_{\mathbf{H} \in R^9} \sum_{j=1}^m \left\| \tilde{\mathbf{m}}_j - \mathbf{H}\tilde{\mathbf{w}}'_j \right\|^2$$

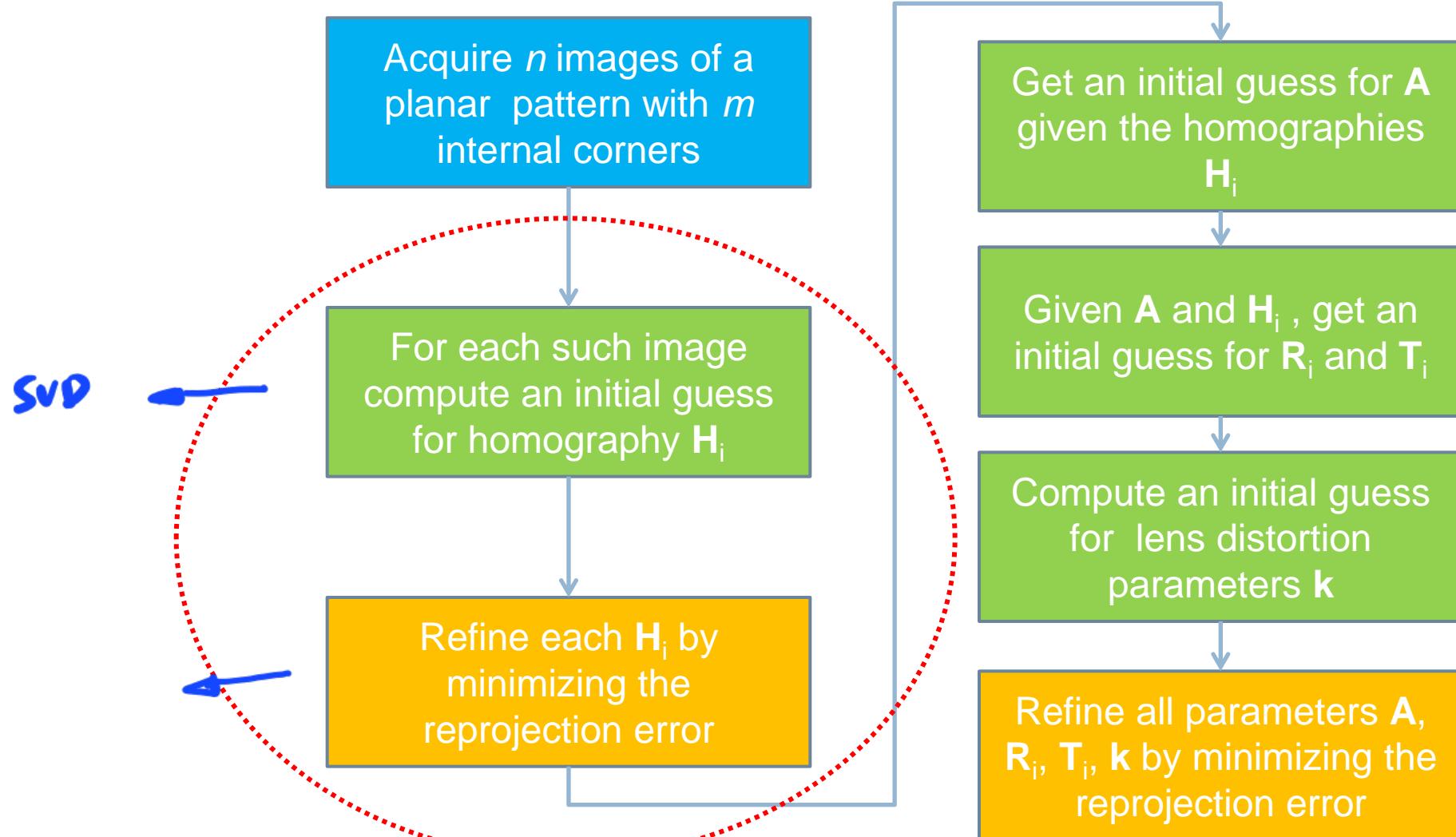
minimizes the distance
world corner
observed corner
estimated / measured corner

The diagram shows the non-linear refinement equation $\mathbf{H} = \operatorname{argmin}_{\mathbf{H} \in R^9} \sum_{j=1}^m \left\| \tilde{\mathbf{m}}_j - \mathbf{H}\tilde{\mathbf{w}}'_j \right\|^2$. A blue arrow points from the term $\left\| \cdot \right\|^2$ to the text "minimizes the distance". A yellow arrow points from the term $\tilde{\mathbf{m}}_j$ to the text "observed corner". Another yellow arrow points from the term $\mathbf{H}\tilde{\mathbf{w}}'_j$ to the text "estimated / measured corner".

- which can be obtained in practice using by the Levenberg-Marquardt (LM) algorithm.
- This additional optimization steps corresponds to the minimization of the reprojection error measured for each of the 3D corners (with $z=0$) of the pattern by comparing the pixel coordinates predicted by the estimated homography to the pixel coordinates of the corresponding corner extracted in the image. The rationale is that the “best” homography would predict with the best accuracy the positions of the corner features actually found in the image. Such a reprojection error is typically referred to as “geometric error”.



Zhang's Method



Estimation of the intrinsic parameters (1)



- As \mathbf{H} is known up to a scale factor, we can establish the following relation between \mathbf{H} and the PPM:

$$\mathbf{H} = [\mathbf{h}_1 \quad \mathbf{h}_2 \quad \mathbf{h}_3] = [\overset{S+1}{\mathbf{p}_1} \quad \overset{S+1}{\mathbf{p}_2} \quad \overset{S+1}{\mathbf{p}_4}] = \lambda \mathbf{A} [\mathbf{r}_1 \quad \mathbf{r}_2 \quad \mathbf{T}] \rightarrow \lambda \mathbf{A} [\mathbf{R} \mid \mathbf{T}]$$

- \mathbf{R} is an orthogonal matrix, its vectors being orthonormal. This constrains the intrinsic parameters to obey to the following relations:

$$\mathbf{r}_1^T \mathbf{r}_2 = 0$$

\Rightarrow

$$\mathbf{h}_1^T \mathbf{A}^{-T} \mathbf{A}^{-1} \mathbf{h}_2 = 0$$

$$\|\mathbf{r}_1\| = 1$$

$$\|\mathbf{r}_2\| = 1$$

$$\mathbf{r}_1^T \mathbf{r}_1 = \mathbf{r}_2^T \mathbf{r}_2$$

\Rightarrow

$$\mathbf{h}_1^T \mathbf{A}^{-T} \mathbf{A}^{-1} \mathbf{h}_1 = \mathbf{h}_2^T \mathbf{A}^{-T} \mathbf{A}^{-1} \mathbf{h}_2$$

$$\|\mathbf{r}_1\| = \|\mathbf{r}_2\|$$

$$\mathbf{r}_1^T \cdot \mathbf{r}_2 = 0 \rightarrow \frac{1}{\lambda^2} (\mathbf{A}^{-T} \mathbf{h}_1)^T \mathbf{A}^{-T} \mathbf{A}^{-1} \mathbf{h}_2 = 0$$

$$\mathbf{h}_1 = \lambda \mathbf{A} \mathbf{r}_1$$

$$\mathbf{r}_1 = \frac{1}{\lambda} \mathbf{A}^{-1} \mathbf{h}_1$$

$$\mathbf{r}_2 = \frac{1}{\lambda} \mathbf{A}^{-1} \mathbf{h}_2$$

where the unknowns are the entries of $\mathbf{B} = \mathbf{A}^{-T} \mathbf{A}^{-1}$. As \mathbf{A} is upper triangular, \mathbf{B} turns out to be symmetric, so that the unknowns are just 6.

- By stacking the above two equations provided by each calibration image, we obtain a $2n \times 6$ linear system, which can be solved in case at least 3 images are available (more details on next slide)

Estimation of the extrinsic parameters (2)



- By posing

$$\mathbf{B} = \mathbf{A}^{-T} \mathbf{A}^{-1} = \begin{pmatrix} B_{11} & B_{12} & B_{13} \\ B_{12} & B_{22} & B_{23} \\ B_{13} & B_{23} & B_{33} \end{pmatrix}, \quad b = [B_{11}, B_{12}, B_{22}, B_{13}, B_{23}, B_{33}]^T$$

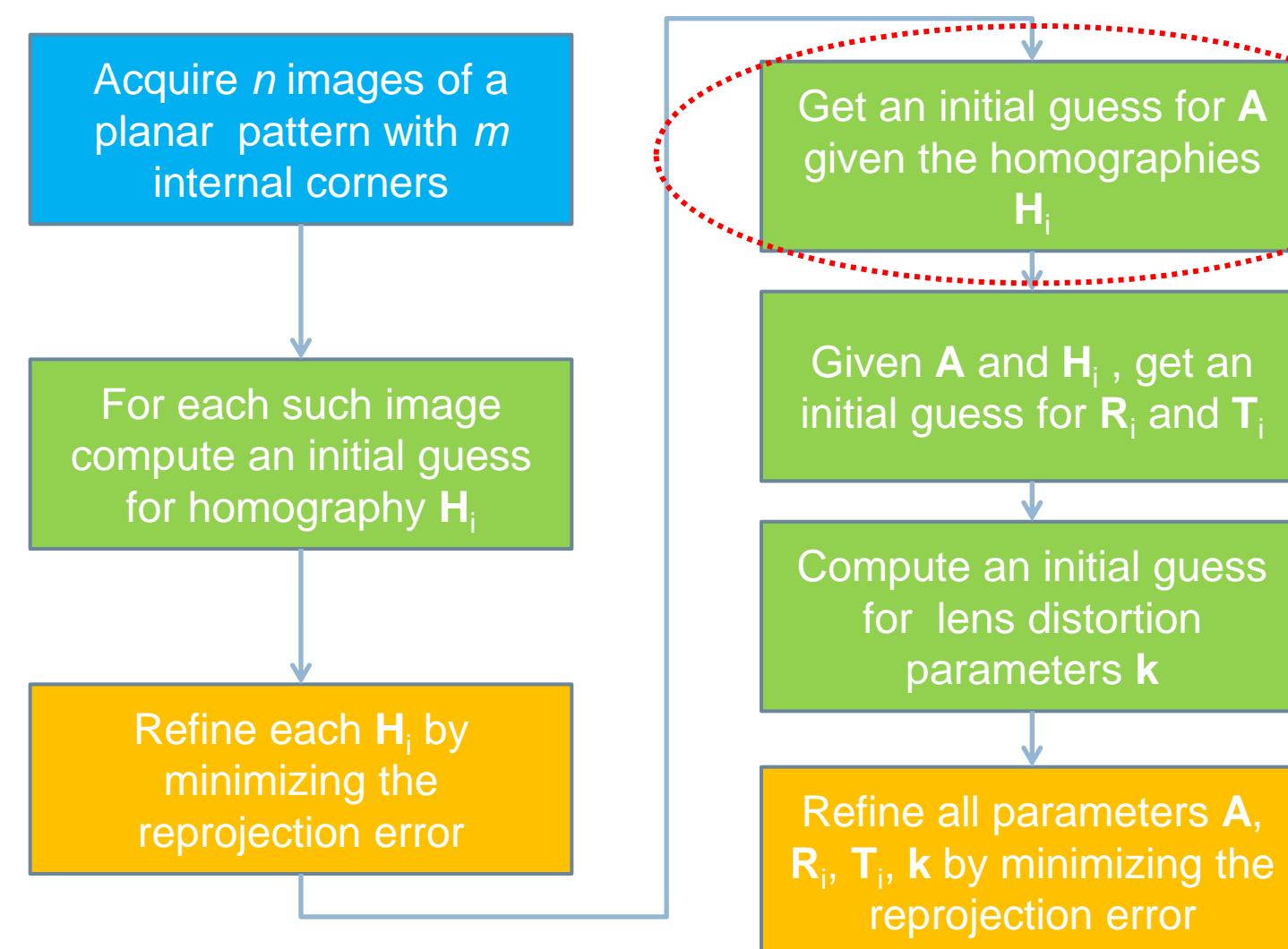
$$\mathbf{h}_i^T = [h_{i1}, h_{i2}, h_{i3}], \quad \mathbf{v}_{ij}^T = [h_{i1}h_{j1}, h_{i1}h_{j2} + h_{i2}h_{j1}, h_{i2}h_{j2}, h_{i3}h_{j1} + h_{i1}h_{j3}, h_{i3}h_{j2} + h_{i2}h_{j3}, h_{i3}h_{j3}]$$

it can be noticed that

$$\mathbf{h}_i^T \mathbf{B} \mathbf{h}_j = \mathbf{v}_{ij}^T \mathbf{b} \quad \xrightarrow{\hspace{1cm}} \quad \begin{aligned} \mathbf{h}_1^T \mathbf{B} \mathbf{h}_2 &= 0 \Rightarrow \mathbf{v}_{12}^T \mathbf{b} = 0 \\ \mathbf{h}_1^T \mathbf{B} \mathbf{h}_1 &= \mathbf{h}_2^T \mathbf{B} \mathbf{h}_2 \Rightarrow \mathbf{v}_{11}^T \mathbf{b} = \mathbf{v}_{22}^T \mathbf{b} \Rightarrow (\mathbf{v}_{11} - \mathbf{v}_{22})^T \mathbf{b} = 0 \end{aligned}$$

- Therefore, each image provides 2 equations in the 6 independent unknowns in \mathbf{B} , so that with n calibration images we get a homogeneous linear system of equations in the form $\mathbf{Vb} = 0$, which can be solved by the SVD.
- Once \mathbf{b} has been calculated, the intrinsic parameters (i.e. \mathbf{A}) can be obtained in closed form.

Zhang's Method



$$\begin{bmatrix} h_1 & h_2 & h_3 \end{bmatrix} = \lambda A \begin{bmatrix} r_1 & r_2 & T \end{bmatrix}$$

Estimation of the extrinsic parameters

- Once \mathbf{A} has been estimated, for each image it is possible to compute \mathbf{R} and then \mathbf{T} given the previously computed homography \mathbf{H} :

$$\mathbf{H} = \begin{bmatrix} \mathbf{h}_1 & \mathbf{h}_2 & \mathbf{h}_3 \end{bmatrix} = \lambda \mathbf{A} \begin{bmatrix} \mathbf{r}_1 & \mathbf{r}_2 & \mathbf{T} \end{bmatrix}$$

$\mathbf{h}_1 = \lambda \mathbf{A} \mathbf{r}_1 \rightarrow \mathbf{r}_1 = \frac{1}{\lambda} \mathbf{A}^{-1} \mathbf{h}_1$

- As \mathbf{r}_1 is a unit vector:

$$\lambda = \|\mathbf{A}^{-1} \mathbf{h}_1\|, \quad \mathbf{r}_2 = \frac{1}{\lambda} \mathbf{A}^{-1} \mathbf{h}_2$$

- \mathbf{r}_3 can be derived from \mathbf{r}_1 and \mathbf{r}_2 by exploiting orthonormality:

$\text{It has to be perpendicular to } \mathbf{r}_1 \text{ and } \mathbf{r}_2$

- Finally, we can get \mathbf{T}



$$\mathbf{T} = \frac{1}{\lambda} \mathbf{A}^{-1} \mathbf{h}_3$$

how to estimate this?

λ needs to be decided first. So, λ must be the norm of $\|\mathbf{A}^{-1} \mathbf{h}_1\|$

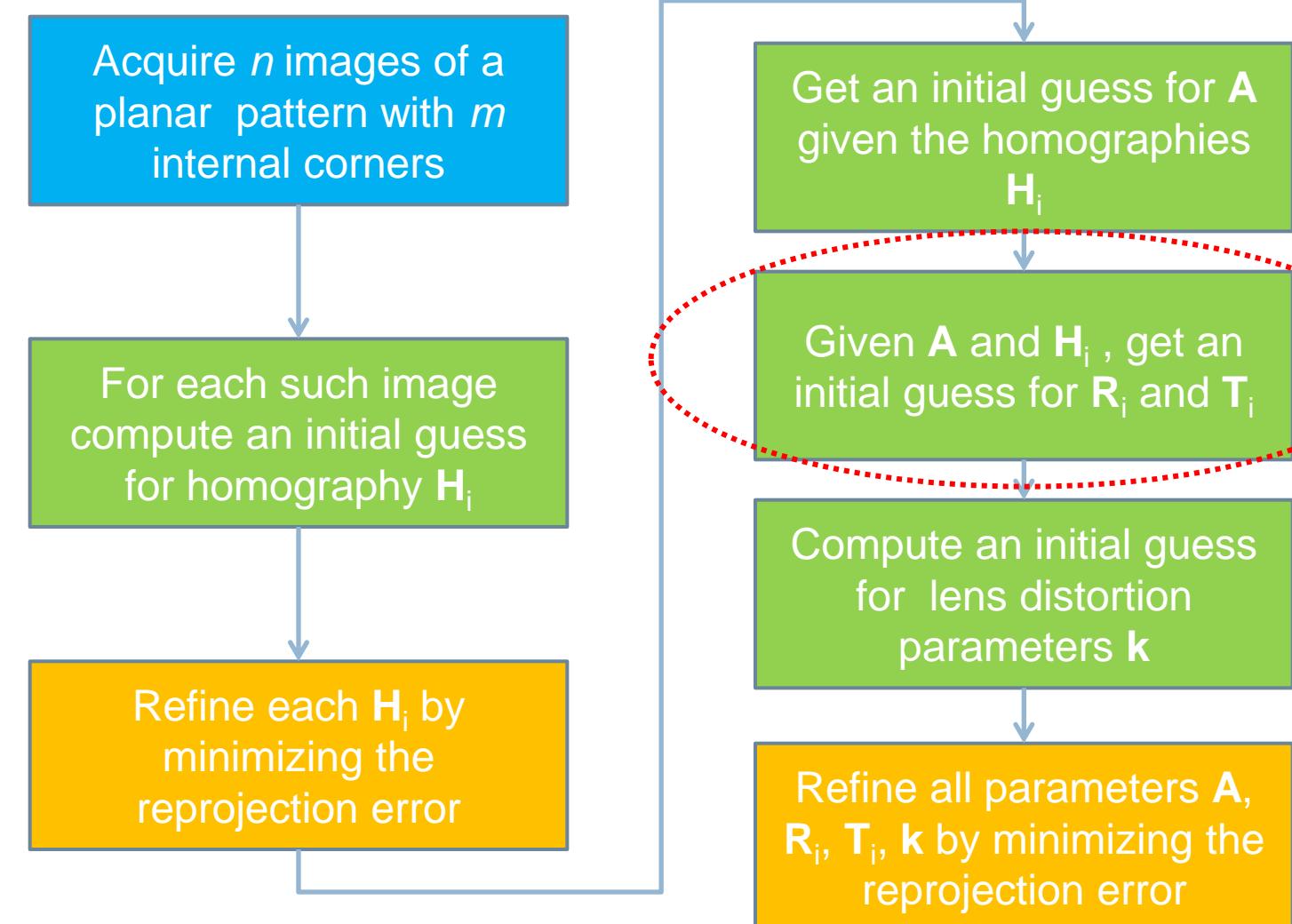
$$\mathbf{r}_3 = \mathbf{r}_1 \times \mathbf{r}_2$$

$$\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$$

$R^* = V I V^T$ using SVD, where :

$$\Sigma = \begin{bmatrix} r_1 & & 0 \\ 0 & r_2 & \\ & & r_3 \end{bmatrix}$$

Zhang's Method



Lens distortion parameters (1)



- Given the homographies, we have both the real (distorted) coordinates of the corner features found in the images as well as the ideal (undistorted) coordinates predicted by the homographies. Zhang's method deploys such information to estimate parameters k_1, k_2 of the radial distortion function.
- Given the already known intrinsic parameter matrix, A , and the lens distortion model reported below:

$$\begin{pmatrix} x' \\ y' \end{pmatrix} = L(r) \begin{pmatrix} \tilde{x} \\ \tilde{y} \end{pmatrix} = (1 + k_1 r^2 + k_2 r^4) \begin{pmatrix} \tilde{x} \\ \tilde{y} \end{pmatrix}$$

Annotations in blue:

- "distorted coordinate" points to (x', y')
- "undistorted" points to $\begin{pmatrix} \tilde{x} \\ \tilde{y} \end{pmatrix}$
- "however, we don't know undistorted c." points to the coefficient $(1 + k_1 r^2 + k_2 r^4)$
- "the unknown" points to the right side of the equation

we need first to establish the relationship between distorted (u', v') and ideal (\tilde{u}, \tilde{v}) pixel coordinates:

there's no way to obtain the real undistorted coordinate

$$\begin{cases} u = \alpha_u \frac{x}{z} + u_0 \\ v = \alpha_v \frac{y}{z} + v_0 \end{cases}$$

Annotations in blue:

- "focal length in horizontal" points to α_u
- "focal length in vertical" points to α_v

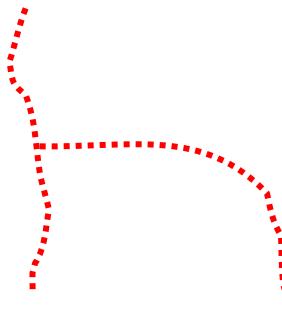


$$\begin{cases} \tilde{x} = \frac{\tilde{u} - u_0}{\alpha_u} \\ \tilde{y} = \frac{\tilde{v} - v_0}{\alpha_v} \end{cases} \quad \begin{cases} x' = \frac{u' - u_0}{\alpha_u} \\ y' = \frac{v' - v_0}{\alpha_v} \end{cases}$$

Lens distortion parameters (2)



$$\begin{bmatrix} x' \\ y' \end{bmatrix} = (1 + k_1 r^2 + k_2 r^4) \begin{bmatrix} \tilde{x} \\ \tilde{y} \end{bmatrix} \rightarrow \begin{cases} \frac{u' - u_0}{\alpha_u} = (1 + k_1 r^2 + k_2 r^4) \left(\frac{\tilde{u} - u_0}{\alpha_u} \right) \\ \frac{v' - v_0}{\alpha_v} = (1 + k_1 r^2 + k_2 r^4) \left(\frac{\tilde{v} - v_0}{\alpha_v} \right) \end{cases}$$



- It is therefore possible to set up a linear system where the unknowns are the distortion coefficients:

$$\begin{cases} u' = \tilde{u} + (k_1 r^2 + k_2 r^4)(\tilde{u} - u_0) \\ v' = \tilde{v} + (k_1 r^2 + k_2 r^4)(\tilde{v} - v_0) \end{cases}$$

$$\begin{bmatrix} (\tilde{u} - u_0)r^2 & (\tilde{u} - u_0)r^4 \\ (\tilde{v} - v_0)r^2 & (\tilde{v} - v_0)r^4 \end{bmatrix} \begin{bmatrix} k_1 \\ k_2 \end{bmatrix} = \begin{bmatrix} u' - \tilde{u} \\ v' - \tilde{v} \end{bmatrix}$$

handy

$$r^2 = \left(\frac{\tilde{u} - u_0}{\alpha_u} \right)^2 + \left(\frac{\tilde{v} - v_0}{\alpha_v} \right)^2$$

$$x \quad \quad \quad y^2$$

$$z = \sqrt{(\tilde{x} - x_c)^2 + (\tilde{y} - y_c)^2}$$

center

$$(x_c, y_c) = (0, 0)$$

$$z = \tilde{x}^2 + \tilde{y}^2$$

- With m corner features in n images, we get a linear system with $2mn$ equations in 2 unknowns, which admits a least-squares solution:

$$2nm \times 2$$

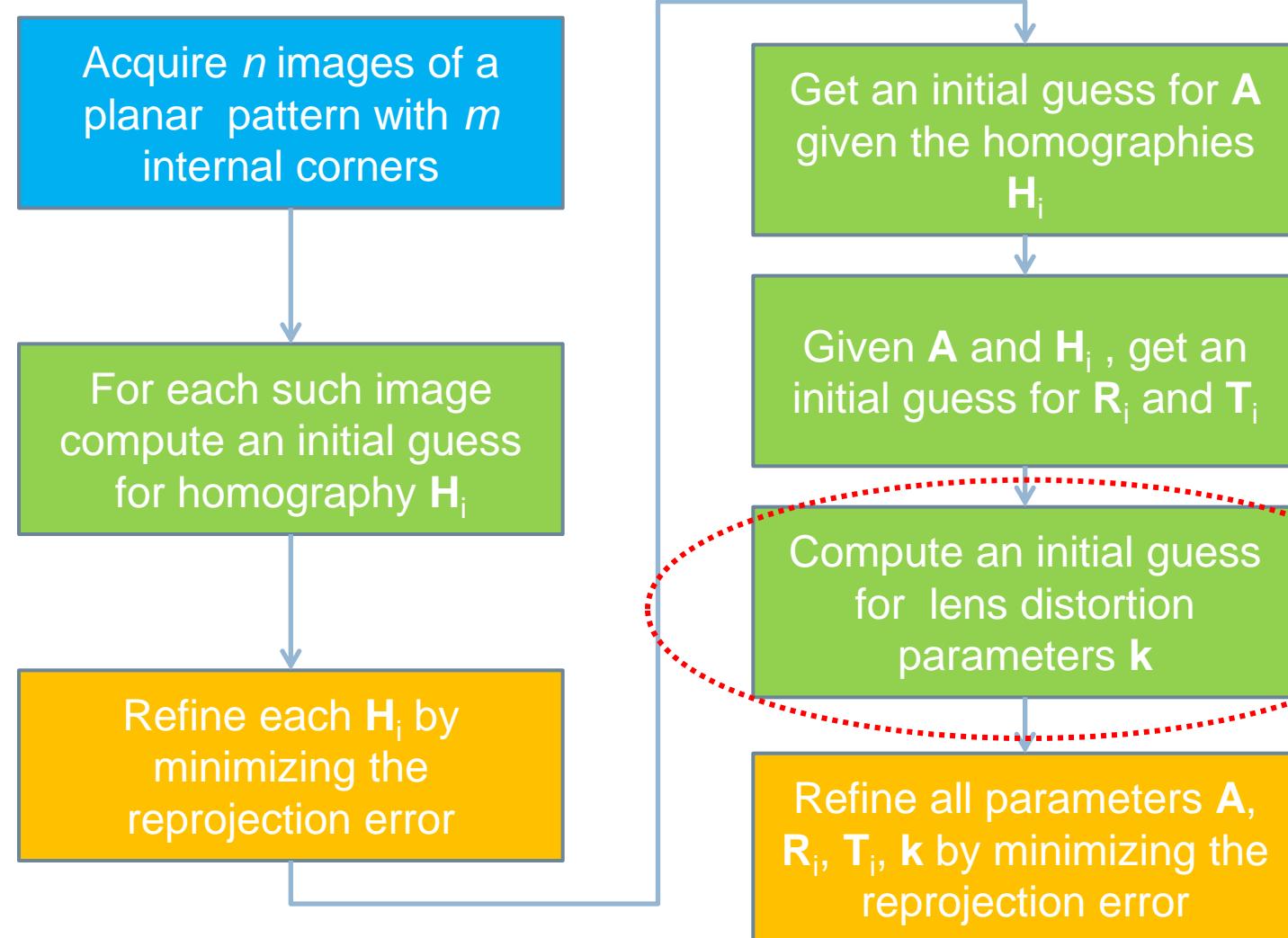
$$\mathbf{D}\mathbf{k} = \mathbf{d}$$

$$2 \times 1$$

$$\mathbf{D}^+ \mathbf{d} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{d}$$

$$\mathbf{D}^+ = \mathbf{V} \boldsymbol{\Sigma}^+ \mathbf{U}^T$$

Zhang's Method



Refinement by non-linear optimization



- Again, the procedure highlighted so far seeks to minimize an algebraic error rather than the actual reprojection error.
- A more accurate solution can instead be found by a so called Maximum Likelihood Estimate (MLE) aimed at minimization of the reprojection (i.e. geometric) error.
- Under the hypothesis of *i.i.d.* (independent identically distributed) noise, the MLE for our models is obtained by minimization of the error

$$\sum_{i=1}^n \sum_{j=1}^m \left\| \mathbf{m}_{ij} - \hat{\mathbf{m}}(\mathbf{A}, \mathbf{k}, \mathbf{R}_i, \mathbf{T}_i, \mathbf{w}_j) \right\|^2$$

all corners of an image

all images *it allows slight error*

The equation shows a double summation over all images and corners. A blue circle highlights the inner summation index j , with a blue arrow pointing from it to the handwritten note "all corners of an image". Another blue circle highlights the outer summation index i , with a blue arrow pointing from it to the handwritten note "all images". To the right of the equation, another blue arrow points to the handwritten note "it allows slight error".

with respect to all the unknown parameters.

- The solution of the above non-linear optimization problem is again provided by the Levenberg-Marquardt algorithm.

Some important additional considerations regarding PPMs are reported in Appendix 2.

Calibration of a stereo camera (1)

to determine the distance of
the obj.
from the
camera



- Given a stereo system (aka *rig* or *head*), we are able to calibrate separately the two cameras. This would provide us with:
 - Intrinsic parameters and lens distortion coefficients
 - Extrinsic parameters, i.e. the rigid motion between the CRF and a given WRF (e.g. associated with the calibration target).
- However, calibrating a stereo system requires determining also the rigid motion, \mathbf{R} and \mathbf{T} , between the two cameras (i.e. between CRF_L and CRF_R) *we need to know the correlation between two cams*
- Knowing the extrinsic parameters wrt the same WRF would allow to compute the rigid motion between the two cameras by simply chaining the transformations: $\mathbf{G}_i(\mathbf{G}_j)^{-1}$, with either ($i=L, j=R$) or ($i=R, j=L$). With Zhang's method this can be achieved, e.g., by taking the same static pattern position as the first image in both calibration sequences.
- However, this is not a robust solution, especially wrt noise and can easily yield to quite inaccurate results.

Calibration of a stereo camera (2)



- A more accurate solution is achieved by starting from an initial guess for \mathbf{R} and \mathbf{T} and then refining the estimation by non-linear minimization of the reprojection error.
- To obtain a robust initial guess, the same position of the planar calibration target is seen in each image of both sequences. This allows to estimate the rigid motion wrt to each such pattern positions as $\mathbf{G}_i(\mathbf{G}_j)^{-1}$. [$\mathbf{R}_i | \mathbf{T}_i$] ↪
- The initial guess is then taken as the median over the values \mathbf{R}_i and \mathbf{T}_i computed as $\mathbf{G}_i(\mathbf{G}_j)^{-1}$ from each pair of corresponding calibration images. Purposely, each \mathbf{R}_i is converted into a 3 elements vector (rotation angles). median rotational matrix
- Finally, the estimation is refined by non-linear minimization of the reprojection error in both images of each pair by Levenberg-Marquardt's:

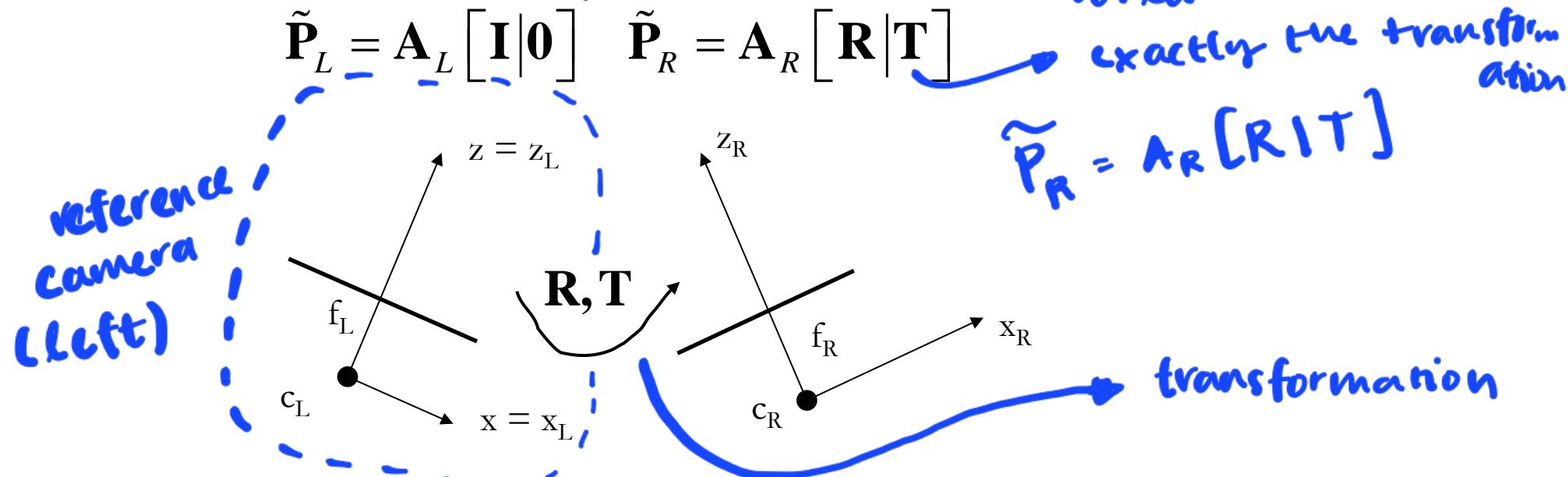
$$\sum_{k=L,R} \sum_{i=1}^n \sum_{j=1}^m \left\| \mathbf{m}_{i,j}^k - \tilde{\mathbf{m}}(\mathbf{A}_L, \mathbf{A}_R, \mathbf{K}_L, \mathbf{K}_R, \mathbf{R}, \mathbf{T}, \mathbf{w}_j) \right\|^2$$

Stereo Reference Frame (SRF)



*Sensing the depth (or distance in the world)
becomes possible*

- Once the stereo rig has been calibrated, it is often useful to define a RF attached to the system, which can be thought of as the “WRF” of the Stereo Camera and will be denoted hereinafter as **SRF** (Stereo Reference Frame). *the master ↪*
- The standard choice is just to pick-up the CRF of one of the two cameras, e.g. the **left camera**. Accordingly, the extrinsic parameters of the right camera in the SRF are given by the previously computed **R** and **T**, with the corresponding PPMs given by:



Rectification

- It is well known that a standard (i.e. lateral) stereo geometry is the most convenient choice to search for corresponding pixels:

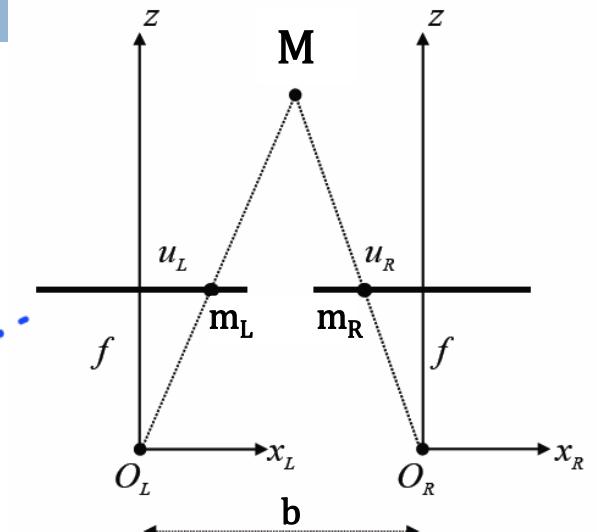
$$(1): M_L = M_R + \begin{bmatrix} b \\ 0 \\ 0 \end{bmatrix}, (2): f_L = f_R = f \rightarrow A_L = A_R = A$$

must be identical

because there's pixelization

- Unfortunately, it is impossible to achieve such geometry by mechanical alignment.

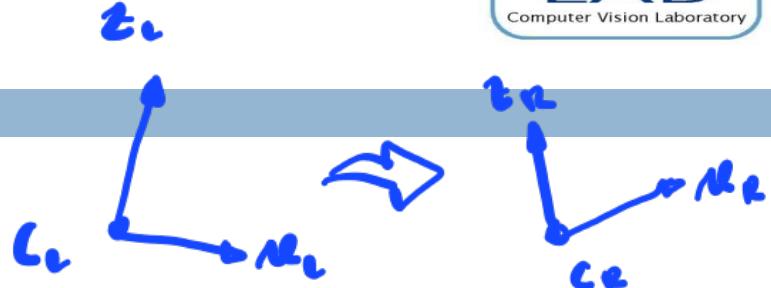
- However, once calibration has been performed, one can define two new cameras (i.e. two new PPMs) in standard geometry by virtually rotating the original ones about their optical centres (1) and constraining the intrinsic parameters of the new cameras to be the same (2).
- Then, the images acquired by each of the two original cameras can be rectified, i.e. transformed into those that would have been acquired by the new ones, by a **homography**.
- Rectification needs to happen after compensation of lens distortion (i.e. on undistorted images). An initial transformation to remove lens distortion is thus always applied before transforming again the images according to the **rectification homographies**.



The new PPMs



- Given $\mathbf{A}_L, \mathbf{A}_R, \mathbf{R}, \mathbf{T}$ the new PPMs can be achieved as follows:
- A new intrinsic parameter matrix, \mathbf{A}_{new} , is arbitrarily chosen (e.g. the mean between $\mathbf{A}_L, \mathbf{A}_R$).
- The new rotation matrix for both cameras, \mathbf{R}_{new} , is then determined. As the row vectors in \mathbf{R}_{new} represent the X, Y, e Z axis of the new CRFs into the WRF, taking as WRF the previously introduced SRF, a lateral stereo system is achieved as follows.
- The new X axis is taken parallel to the baseline vector $\mathbf{B} = \mathbf{C}_R - \mathbf{C}_L$. As we wish to express \mathbf{B} into the SRF: $\mathbf{B} = -\mathbf{R}^T \mathbf{T} - \mathbf{0} = -\mathbf{R}^T \mathbf{T} = [B_x \ B_y \ B_z]^T$. Then, the first row of \mathbf{R}_{new} , is given by the unit vector parallel to the baseline vector: $\mathbf{r}_1 = \frac{\mathbf{B}}{\|\mathbf{B}\|}$
- The new Y axis is orthogonal to X and to arbitrary unit vector \mathbf{k} , e.g. chosen to be parallel to the old Z axis of the left camera ($\mathbf{k} = [0 \ 0 \ 1]^T$):
- The new Z axis is, of course, orthogonal to X and Y: $\mathbf{r}_3 = \mathbf{r}_1 \wedge \mathbf{r}_2$
- Thus, the new PPMs are: $\tilde{\mathbf{P}}'_L = \mathbf{A}_{\text{new}}[\mathbf{R}_{\text{new}} | 0]$, $\tilde{\mathbf{P}}'_R = \mathbf{A}_{\text{new}}[\mathbf{R}_{\text{new}} | -\mathbf{R}_{\text{new}} \mathbf{C}_R]$



$$\mathbf{r}_2 = \mathbf{k} \wedge \mathbf{r}_1 = \frac{[-B_y \ B_x \ 0]^T}{\sqrt{B_x^2 + B_y^2}}$$

Rectification Homographies



- Both cameras undergo a rotation about the optical center and a change of intrinsics. We have already studied that, in such a case, the images (i.e. the original and rectified ones) are related by a homography and thus we know how to compute these transformations.
- Considering the left camera first:

$$\tilde{\mathbf{m}}_L = \mathbf{A}_L [\mathbf{I} \ 0] \tilde{\mathbf{M}}$$

$$\tilde{\mathbf{m}}'_L = \mathbf{A}_{new} [\mathbf{R}_{new} \ 0] \tilde{\mathbf{M}}$$

$$\tilde{\mathbf{m}}_L = \mathbf{A}_L \mathbf{R}_{new}^{-1} \mathbf{A}_{new}^{-1} \tilde{\mathbf{m}}'_L$$

 \mathbf{H}_L

**Rectification
Homography for the left
image**

- To compute the other rectification homography, we may conveniently think of moving the origin of the WRF into the optical centre of the right camera:

$$\tilde{\mathbf{m}}_R = \mathbf{A}_R [\mathbf{R} \ 0] \tilde{\mathbf{M}}_R$$

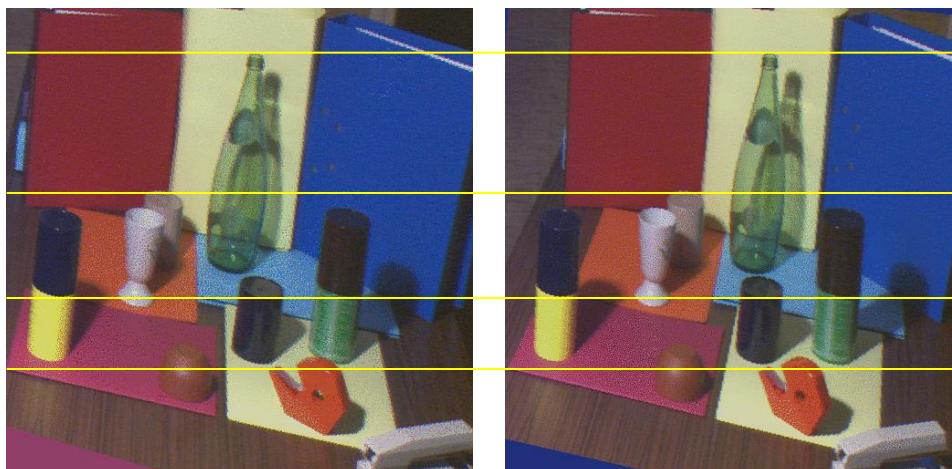
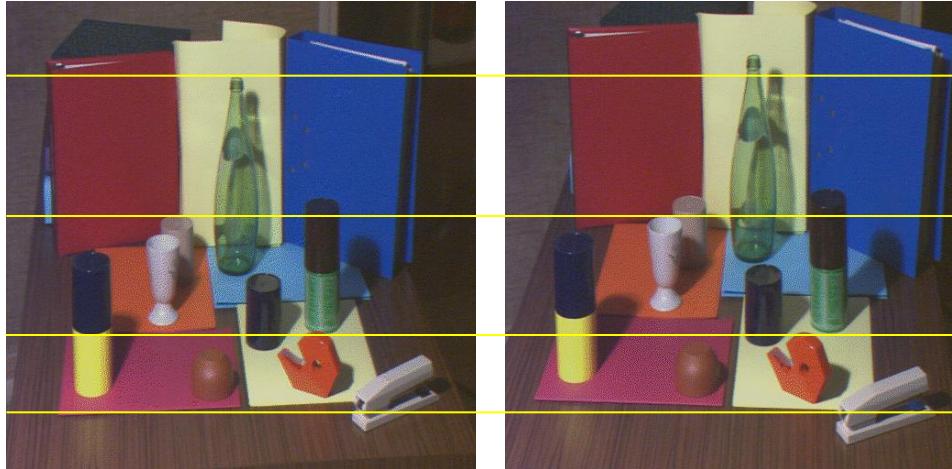
$$\tilde{\mathbf{m}}'_R = \mathbf{A}_{new} [\mathbf{R}_{new} \ 0] \tilde{\mathbf{M}}_R$$

$$\tilde{\mathbf{m}}_R = \mathbf{A}_R \mathbf{R} \mathbf{R}_{new}^{-1} \mathbf{A}_{new}^{-1} \tilde{\mathbf{m}}'_R$$

 \mathbf{H}_R

**Rectification
Homography for the
right image**

An example



From pixels to 3D coordinates



*emits light
so the reflection
is measured in time
of flight*

*obtained from
depth images*

If the camera is calibrated, i.e. matrix \mathbf{A} is known, and depth is also known, as it is the case, e.g., for the reference image of a stereo camera or a TOF or Structured Light (e.g. Kinect) sensor, we can estimate the **3D coordinates of pixels in the CRF** (Camera Reference Frame) in order to obtain a so called **point cloud**, a representation widely used in **3D Computer Vision**.

$$\begin{array}{ccccc}
 p^* & & \begin{bmatrix} x \\ y \\ z \end{bmatrix} & \xrightarrow{\quad \mathbf{A} \quad} & p^* = \mathbf{A} P \\
 & & \mathbf{P} & & \\
 & & & & \frac{p^*}{z} = \begin{bmatrix} \alpha_u \frac{x}{z} + u_0 \\ \alpha_v \frac{y}{z} + v_0 \\ 1 \end{bmatrix} = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = p
 \end{array}$$

$\mathbf{P} = \mathbf{A}^{-1} p^*$
 $\mathbf{P} = z \mathbf{A}^{-1} \frac{p^*}{z}$

Thus:

3D coordinates in the CRF depth inverse of A pixel coordinates

with:

$$A^{-1} = \begin{bmatrix} \frac{1}{\alpha_u} & 0 & -\frac{u_0}{\alpha_u} \\ 0 & \frac{1}{\alpha_v} & -\frac{v_0}{\alpha_v} \\ 0 & 0 & 1 \end{bmatrix}$$

From pixels to 3D coordinates



If the camera is calibrated, i.e. matrix \mathbf{A} is known, and depth is also known, as it is the case, e.g., for the reference image of a stereo camera or a TOF or Structured Light (e.g. Kinect) sensor, we can estimate the **3D coordinates of pixels in the CRF** (Camera Reference Frame) in order to obtain a so called **point cloud**, a representation widely used in **3D Computer Vision**.

$$\begin{aligned}
 \begin{bmatrix} \alpha_u x + u_0 z \\ \alpha_v y + v_0 z \\ z \end{bmatrix} &= \begin{bmatrix} \alpha_u & 0 & u_0 \\ 0 & \alpha_v & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} \quad \rightarrow \quad \mathbf{p}^* = \mathbf{A} \mathbf{P} \\
 \mathbf{p}^* &= \mathbf{A} \mathbf{P} \quad \frac{\mathbf{p}^*}{z} = \begin{bmatrix} \alpha_u \frac{x}{z} + u_0 \\ \alpha_v \frac{y}{z} + v_0 \\ 1 \end{bmatrix} = \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{p} \\
 \mathbf{P} &= \mathbf{A}^{-1} \mathbf{p}^* \\
 \mathbf{P} &= z \mathbf{A}^{-1} \frac{\mathbf{p}^*}{z}
 \end{aligned}$$

3D coordinates in the CRF

depth

inverse of A

pixel coordinates

Thus: $\mathbf{P} = z \mathbf{A}^{-1} \frac{\mathbf{p}^*}{z}$

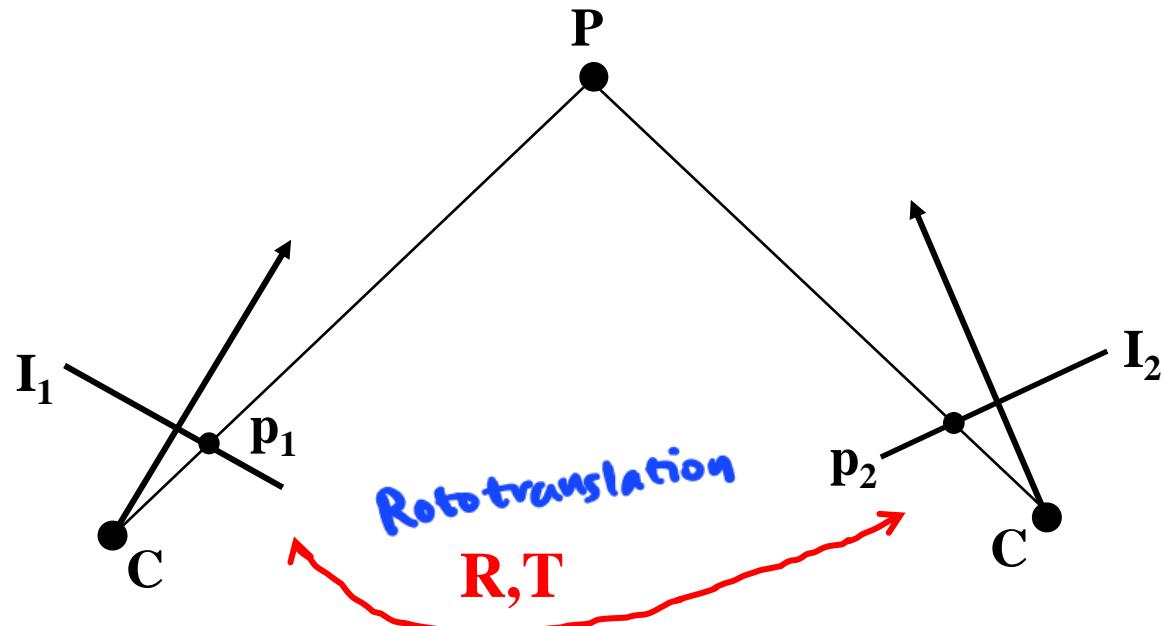
with:

$$A^{-1} = \begin{bmatrix} \frac{1}{\alpha_u} & 0 & -\frac{u_0}{\alpha_u} \\ 0 & \frac{1}{\alpha_v} & -\frac{v_0}{\alpha_v} \\ 0 & 0 & 1 \end{bmatrix}$$

..and back to pixels



We may wish to find where a certain pixel does project into another image taken by the same camera, typically a moving camera used to scan a static scene. Purposefully, alongside depth and camera intrinsics, we need to know the rigid motion (roto-translation) between the two views (i.e. CRFs).



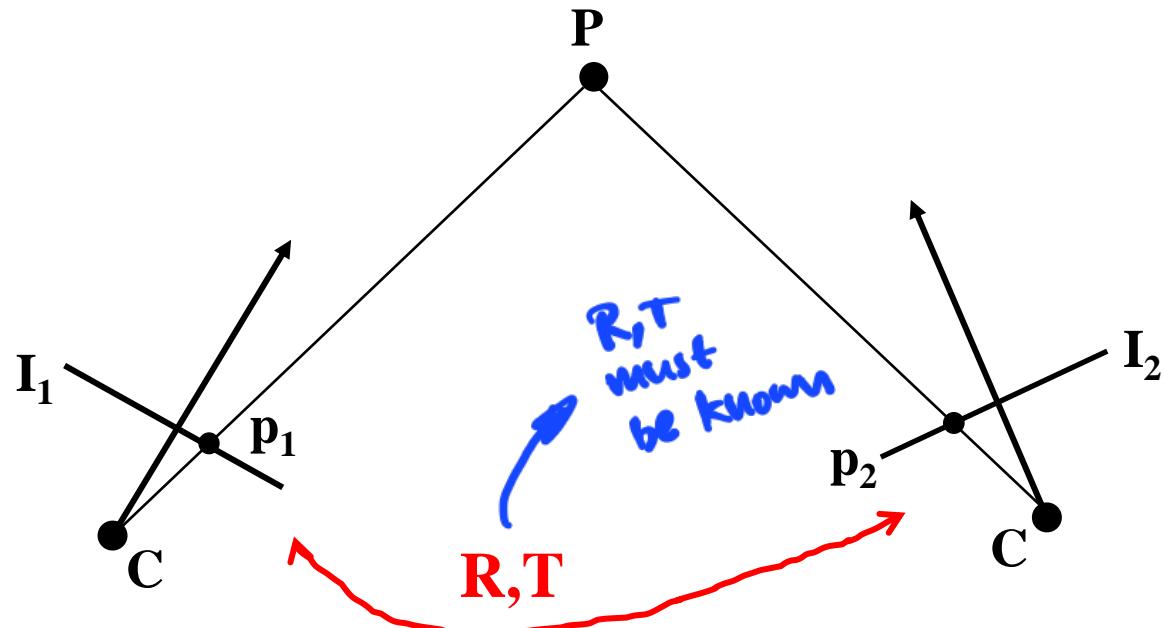
$$zA^{-1}p_1$$

$$\text{with: } T_{1 \rightarrow 2}(P_1) = RP_1 + T$$

..and back to pixels



We may wish to find where a certain pixel does project into another image taken by the same camera, typically a moving camera used to scan a static scene. Purposefully, alongside depth and camera intrinsics, we need to know the rigid motion (roto-translation) between the two views (i.e. CRFs).

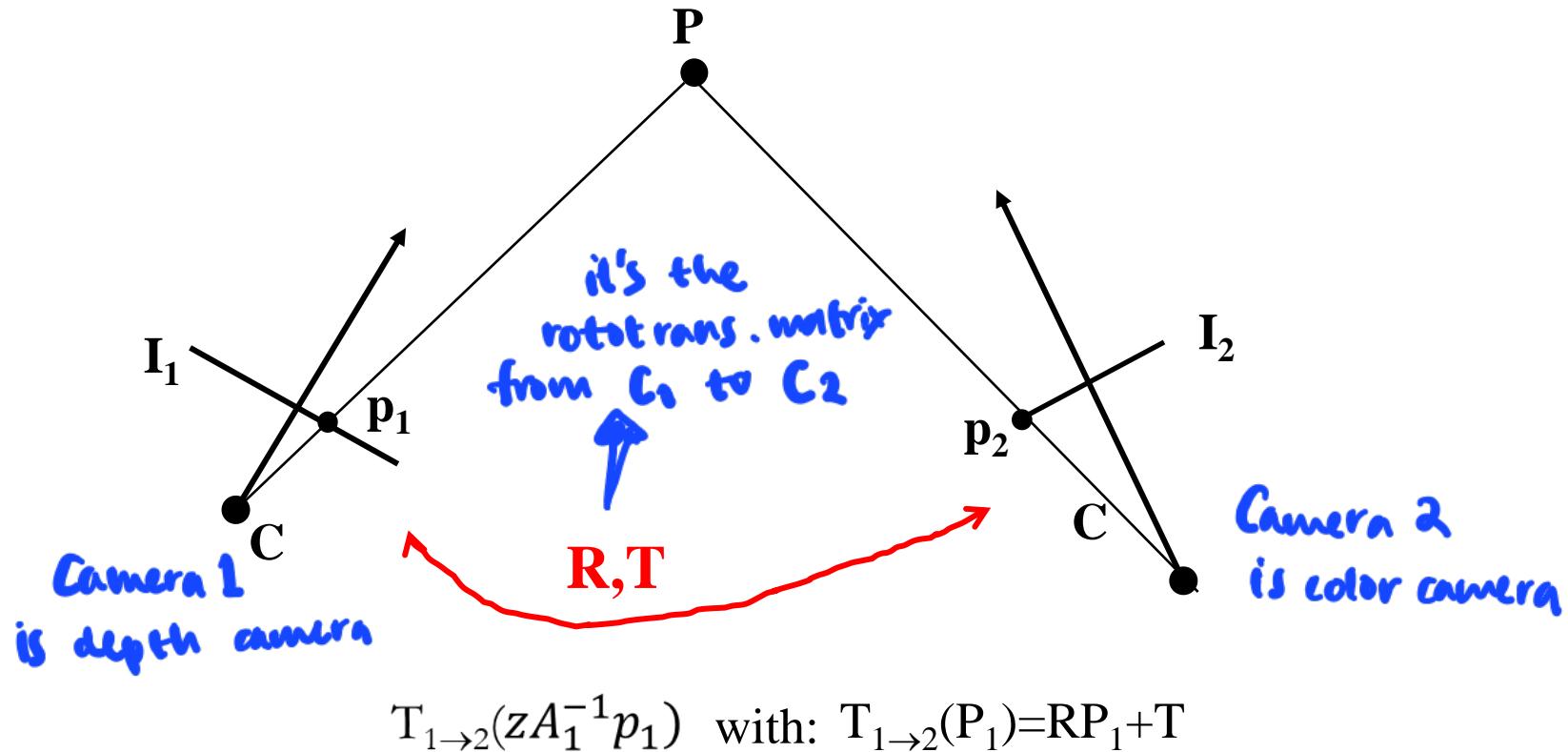


$$p_2 = A T_{1 \rightarrow 2}(z A^{-1} p_1) \quad \text{with: } T_{1 \rightarrow 2}(P_1) = RP_1 + T$$

Cameras may be different



Similarly, the other image may be taken by a different camera. We would then need to know the intrinsic parameters of the other camera too.



Cameras may be different



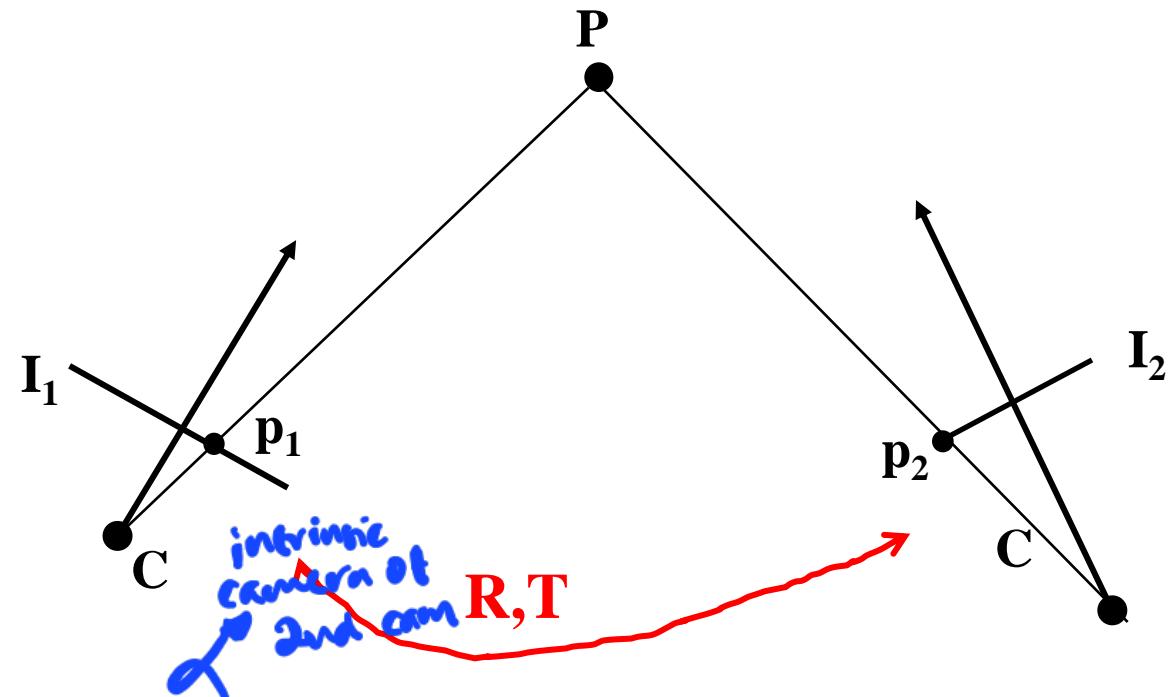
Similarly, the other image may be taken by a different camera. We would then need to know the intrinsic parameters of the other camera too.



Basler – Blaze ToF + RGB



Microsoft Azure Kinect



$$p_2 = A_2 T_{1 \rightarrow 2}(z A_1^{-1} p_1) \quad \text{with: } T_{1 \rightarrow 2}(P_1) = RP_1 + T$$



Luxonis – OAK-D-PoE



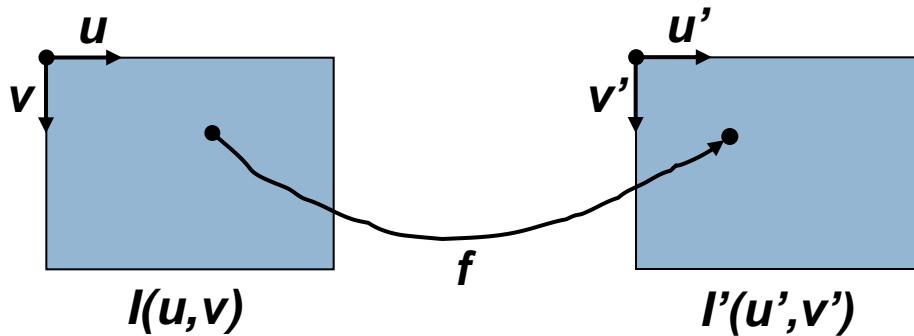
Lucid – Helios2 & Triton

This is how an **RGB-D** image is obtained in many *depth cameras* (e.g. MS Kinect).

Image Warping

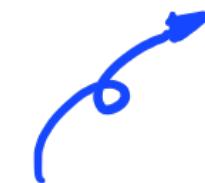


transformation



$$\begin{cases} u' = f_u(u, v) \\ v' = f_v(u, v) \end{cases}$$

$$I'(f_u(u, v), f_v(u, v)) = I(u, v)$$

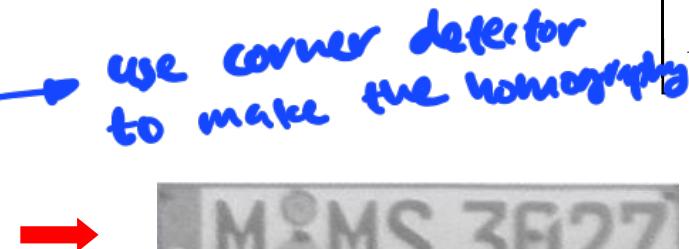


copy the pixel coordinate and put it on a new coord. system

Examples of image warping:

Rotation $\rightarrow \begin{bmatrix} u' \\ v' \end{bmatrix} = \begin{pmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{pmatrix} \begin{bmatrix} u \\ v \end{bmatrix}$

Removal of perspective deformation



use corner detector to make the homography

$$s \begin{bmatrix} x' \\ y' \\ w' \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{33} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

Other examples:

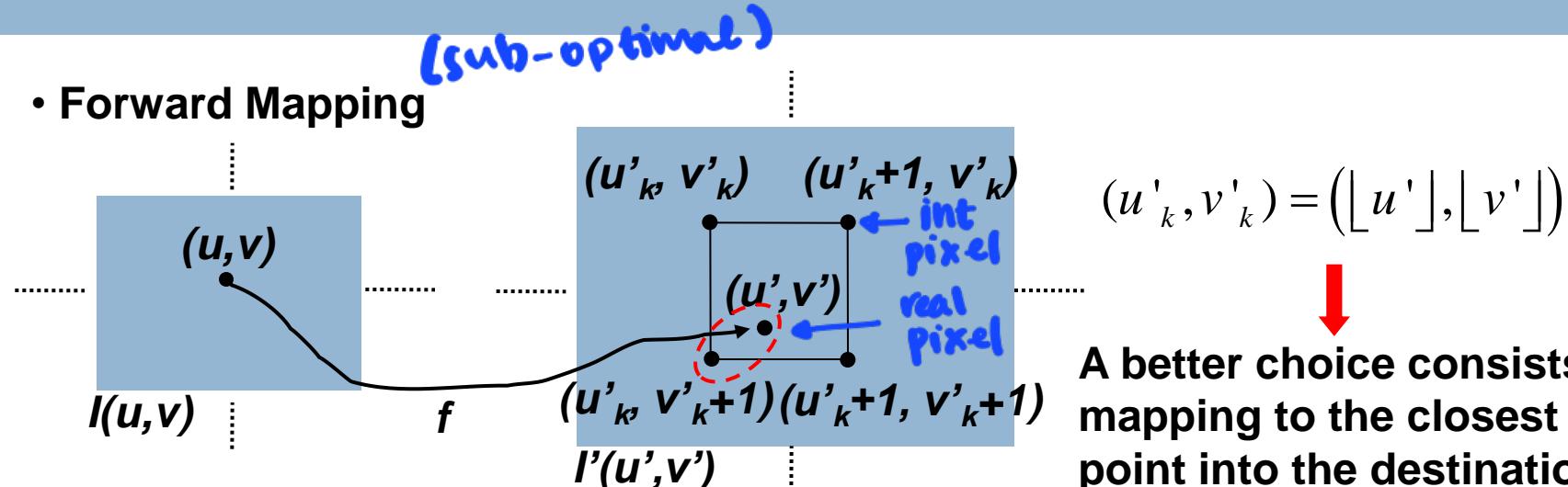
- Remove lens distortion
- Stereo rectification
-

need correspondence to map pixels between 2 images

Forward/Backward Mapping



- Forward Mapping



$$(u'_k, v'_k) = (\lfloor u' \rfloor, \lfloor v' \rfloor)$$

A better choice consists in mapping to the closest point into the destination image

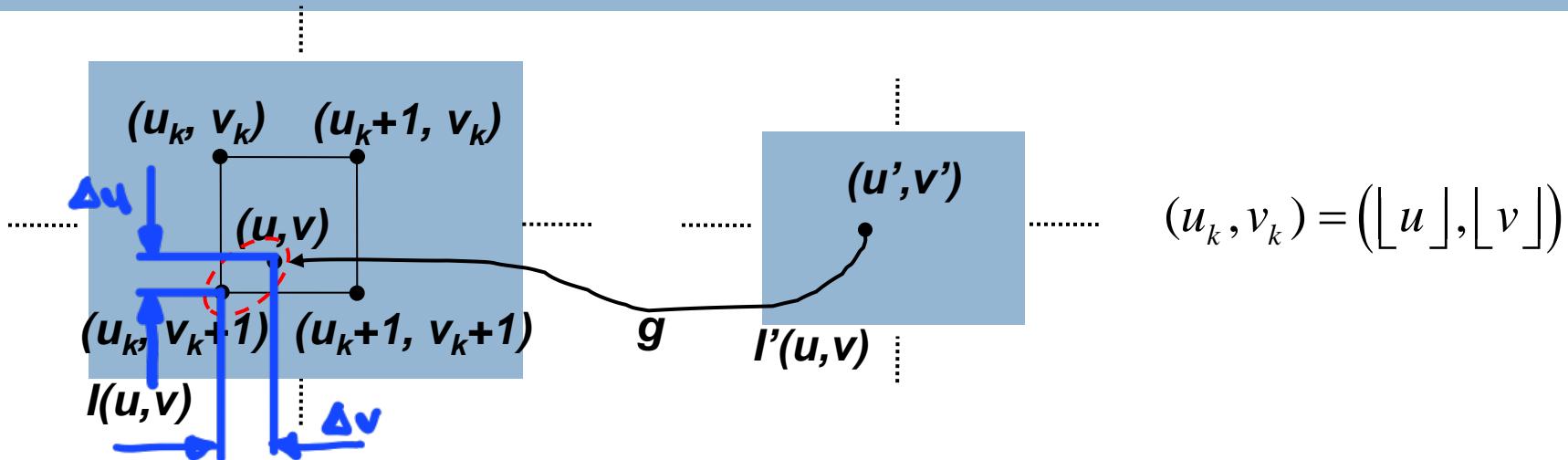
because assigned pixels are conflicting.

- Backward Mapping

$$\begin{cases} u = g_u(u', v') \\ v = g_v(u', v') \end{cases} \rightarrow \forall (u', v'): I'(u', v') = I(g_u(u', v'), g_v(u', v'))$$

No holes and folds in $I'(u', v')$, which mapping strategy ?

Mapping Strategies



- Mapping from the closest point (Nearest Neighbour Mapping)
- Interpolation between the 4 closest points (*bilinear, bicubic,...*)

$$\Delta u = u - u_k$$

$$I_1 = I(u_k, v_k)$$

$$I_2 = I(u_k + 1, v_k)$$

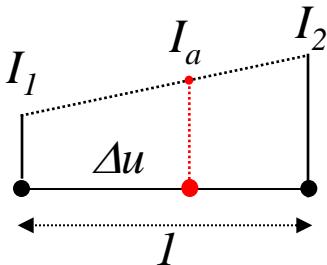
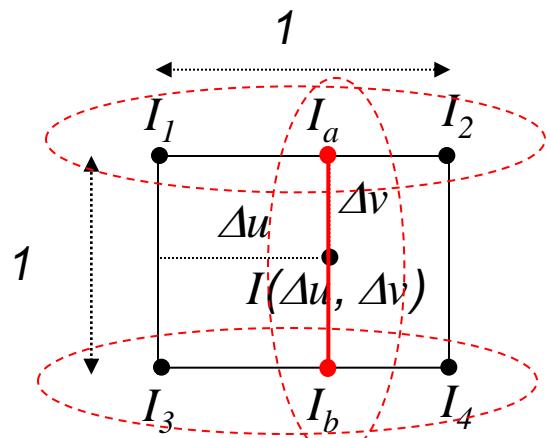
$$\Delta v = v - v_k$$

$$I_3 = I(u_k, v_k + 1)$$

$$I_4 = I(u_k + 1, v_k + 1)$$

Bilinear Interpolation (1)

*we interpolate
two times :
horizontally & vertically*



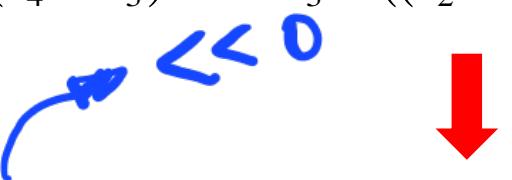
$$\frac{I_a - I_1}{\Delta u} = I_2 - I_1$$

$$I_a = \frac{(I_2 - I_1)\Delta u + I_1}{\Delta u}$$

$$I_b = (I_4 - I_3)\Delta u + I_3$$

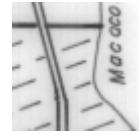
$$I(\Delta u, \Delta v) = (I_b - I_a)\Delta v + I_a$$

$$I(\Delta u, \Delta v) = ((I_4 - I_3)\Delta u + I_3 - ((I_2 - I_1)\Delta u + I_1))\Delta v + (I_2 - I_1)\Delta u + I_1$$

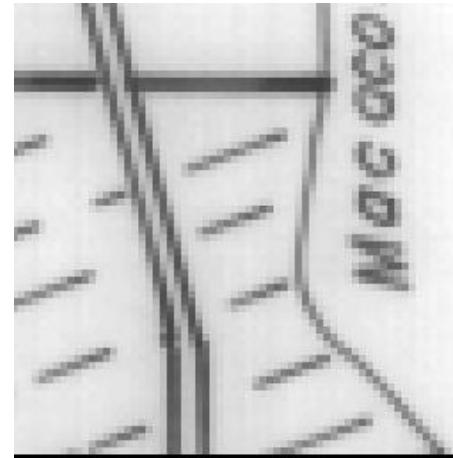


$$I(u', v') = (1 - \Delta u)(1 - \Delta v)I_1 + \Delta u(1 - \Delta v)I_2 + (1 - \Delta u)\Delta v I_3 + \Delta u\Delta v I_4$$

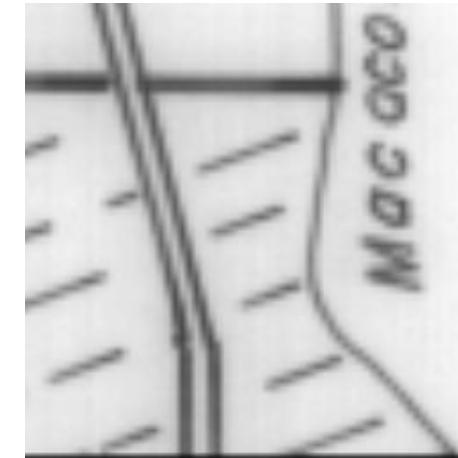
Bilinear Interpolation (2)



Input

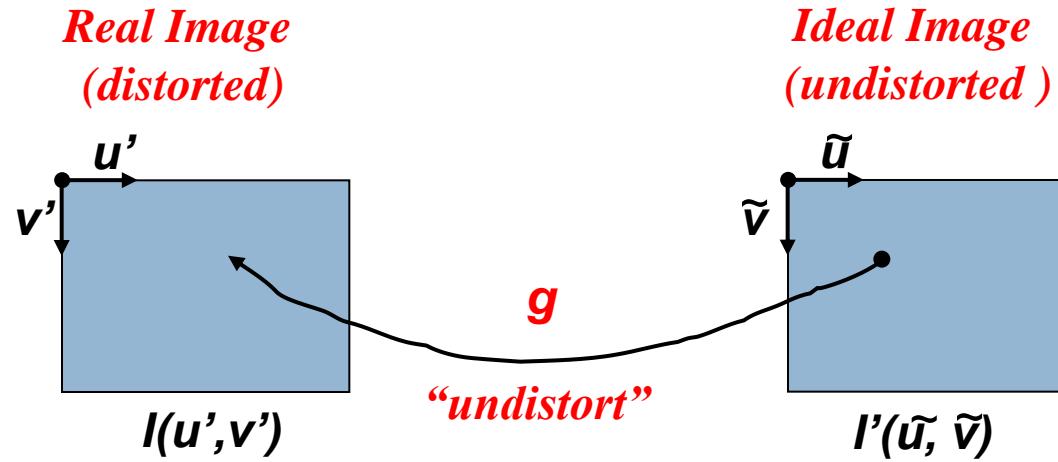


*Warp (Zoom)
by NNM*

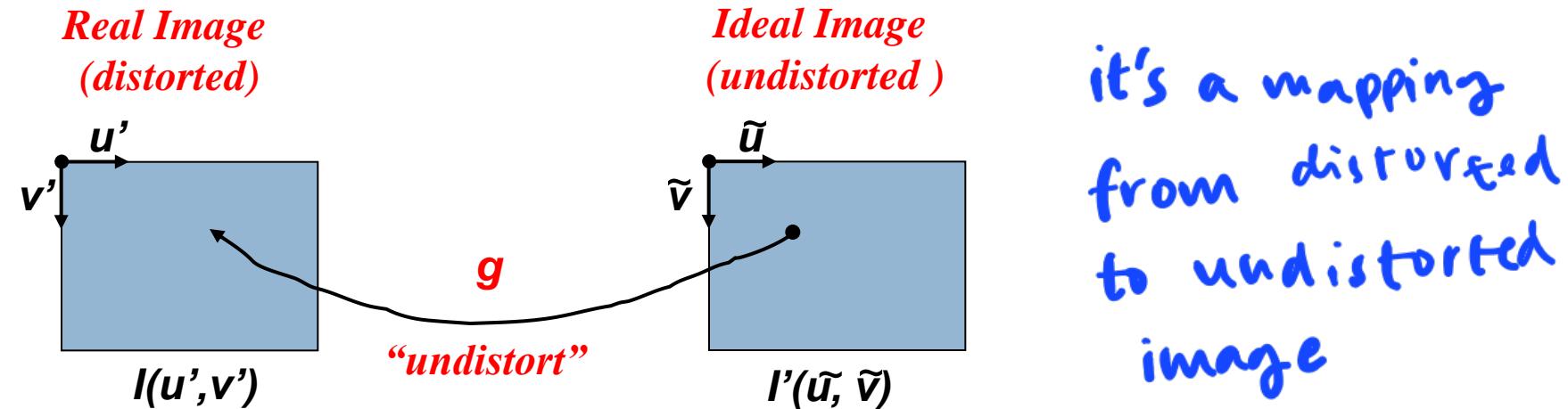


*Warp (Zoom) by
Bilinear Interpolation*

Warping to compensate lens distortion (1)



Warping to compensate lens distortion (2)



Once the lens distortion parameters have been computed by camera calibration, the image can be corrected by a backward warp from the undistorted to the distorted image based on the adopted lens distortion model:

$$\rightarrow \forall (\tilde{u}, \tilde{v}): I'(\tilde{u}, \tilde{v}) = I(g_u(\tilde{u}, \tilde{v}), g_v(\tilde{u}, \tilde{v}))$$

For example, using Zhang's calibration method

$$\begin{cases} u' = \tilde{u} + (k_1 r^2 + k_2 r^4)(\tilde{u} - u_0) \\ v' = \tilde{v} + (k_1 r^2 + k_2 r^4)(\tilde{v} - v_0) \end{cases}$$

Appendix 1 - DLT Algorithm: 4 points case (1)



- By considering 4 point pairs, we obtain a system of 8 equations in 9 unknowns :

$$\mathbf{A}\mathbf{h} = 0, \quad \mathbf{A} = [\mathbf{A}_1 \ \mathbf{A}_2 \dots \mathbf{A}_9], \quad \mathbf{h} = \begin{bmatrix} \mathbf{h}_1 \\ \mathbf{h}_2 \\ \mathbf{h}_3 \end{bmatrix} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_9 \end{bmatrix}$$

where \mathbf{A} is a 8×9 matrix, $\mathbf{A}_1 \dots \mathbf{A}_9$ are 8×1 vectors, \mathbf{h} is a 9×1 vector, $\mathbf{h}_1, \mathbf{h}_2, \mathbf{h}_3$ are 3×1 vectors representing the rows of the 3×3 matrix \mathbf{H} defining the homography, $h_1 \dots h_9$ are the 9 elements of such matrix.

As \mathbf{H} is defined up to a scale factor, we can set $h_9=1$, so as to obtain a non homogeneous linear system with 8 equations and 8 unknowns

$$\tilde{\mathbf{A}}\tilde{\mathbf{h}} = \mathbf{b}, \quad \tilde{\mathbf{A}} = [\mathbf{A}_1 \ \mathbf{A}_2 \dots \mathbf{A}_8], \quad \tilde{\mathbf{h}} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_8 \end{bmatrix}, \quad \mathbf{b} = -\mathbf{A}_9$$

which can be solved easily by standard methods (Cramer's rule, Inversion of the coefficient matrix, Gaussian Elimination)

Appendix 1 - DLT Algorithm: 4 points case (2)



- Alternatively, it is possible to constrain the norm of \mathbf{h} , e.g. so as to render it equal to 1. Purposely, we can assume again h_9 as fixed, but no longer equal to one :

$$\tilde{\mathbf{A}}\tilde{\mathbf{h}} = \mathbf{b}, \quad \tilde{\mathbf{A}} = [\mathbf{A}_1 \ \mathbf{A}_2 \dots \mathbf{A}_8], \quad \tilde{\mathbf{h}} = \begin{bmatrix} h_1 \\ h_2 \\ \vdots \\ h_8 \end{bmatrix}, \quad \mathbf{b} = -h_9 \mathbf{A}_9$$

$$\tilde{\mathbf{h}} = -h_9 \tilde{\mathbf{A}}^{-1} \mathbf{A}_9 \rightarrow \mathbf{h} = \begin{bmatrix} -h_9 \tilde{\mathbf{A}}^{-1} \mathbf{A}_9 \\ h_9 \end{bmatrix} = h_9 \begin{bmatrix} -\tilde{\mathbf{A}}^{-1} \mathbf{A}_9 \\ 1 \end{bmatrix}$$

$$\|\mathbf{h}\| = 1 \rightarrow h_9 \sqrt{\|\tilde{\mathbf{A}}^{-1} \mathbf{A}_9\|^2 + 1} = 1 \rightarrow h_9 = \frac{1}{\sqrt{\|\tilde{\mathbf{A}}^{-1} \mathbf{A}_9\|^2 + 1}}$$

- According to this method, the coefficient matrix of the linear system needs to be inverted and then, using the formulas shown above it is possible to compute h_9 first and then \mathbf{h}

Appendix 2 - Decomposition of the PPM

- Zhang's method provides separately \mathbf{A} , \mathbf{R} e \mathbf{T} . Other methods, though, compute directly the PPM (see, e.g. Hartley&Zisserman, Chapter 7).
- It is however always possible to decompose any given PPM into its elementary components. Purposely, we can deploy the so called **QR factorization**, so as to factorize a real square matrix into the product between an ***orthogonal*** and ***upper triangular*** matrices.
- Let's then $\mathbf{P}^{-1} = \mathbf{UL}$ be the factorization of the inverse of the square matrix \mathbf{P} , with \mathbf{U} orthogonal and \mathbf{L} upper triangular.
- Hence, to obtain the rotation matrix and the intrinsic parameter matrix:

$$\tilde{\mathbf{P}} = [\mathbf{P} | \mathbf{p}_4] = [\mathbf{AR} | \mathbf{AT}] \Rightarrow \mathbf{P}^{-1} = \mathbf{R}^{-1} \mathbf{A}^{-1} \Rightarrow$$

$$\mathbf{R} = \mathbf{U}^{-1}, \mathbf{A} = \mathbf{L}^{-1}$$

- Finally, the translation vector \mathbf{T} is given by:

$$\mathbf{T} = \mathbf{A}^{-1} \mathbf{p}_4 = \mathbf{L} \mathbf{p}_4$$

Appendix 2 - Optical centre from the PPM



- It can be shown that a PPM is always a full rank (i.e. rank 3) matrix, and, alike, that every full-rank 3x4 matrix defines a perspective projection.
 - Intuitively: full rank is mandatory, as otherwise the projection of a 3D point will not be an image point but a higher dimensional subspace (such as a line or plane)
 - Recalling that (0,0,0) does not represent any valid point in homogeneous coordinates, we can observe that a point in space belonging to the matrix null space (or kernel, i.e. the solutions of the associated homogenous system) must be a point whose projection is undefined.
- Such a point is the optical centre, which is indeed the only point in space whose perspective projection onto the image plane is undefined (one cannot define a projection ray from the optical centre and itself). We can therefore write :

$$\tilde{\mathbf{P}} \begin{bmatrix} \mathbf{c} \\ 1 \end{bmatrix} = \mathbf{0} \text{ and thus } \mathbf{c} = -\mathbf{P}^{-1}\mathbf{p}_4 \text{ with } \tilde{\mathbf{P}} = [\mathbf{P} | \mathbf{p}_4]$$

3D coordinates of the optical center in the WRF

Appendix 2 - Optical ray from the PPM



- The optical ray of an image point, $\tilde{\mathbf{m}}$, is the 3D line between $\tilde{\mathbf{m}}$ and the optical centre, , $\tilde{\mathbf{C}} = \begin{bmatrix} \mathbf{C} \\ 1 \end{bmatrix}$.
Indeed, the optical ray is the locus of points, $\tilde{\mathbf{M}}$, satisfying the equation $k\tilde{\mathbf{m}} = \tilde{\mathbf{P}}\tilde{\mathbf{M}}$ for a given $\tilde{\mathbf{m}}$.
int is the optical centre, which is indeed the only point in space whose perspective projection onto
the image plane is undefined (one cannot define a projection ray from the optical centre and itself)
- The equation of the optical ray in Euclidean coordinates can be determined as follows:

$$\tilde{\mathbf{M}}_\infty = \begin{bmatrix} a \\ b \\ c \\ 0 \end{bmatrix} = \begin{bmatrix} \mathbf{V} \\ 0 \end{bmatrix}, \quad \tilde{\mathbf{P}} \begin{bmatrix} \mathbf{P}^{-1}\tilde{\mathbf{m}} \\ 0 \end{bmatrix} = [\mathbf{P} \ p_4] \begin{bmatrix} \mathbf{P}^{-1}\tilde{\mathbf{m}} \\ 0 \end{bmatrix} = \tilde{\mathbf{m}}$$

$\mathbf{M} = \mathbf{C} + \lambda \mathbf{V}$

optical centre → any vector parallel to the optical ray
 $\mathbf{M} = \mathbf{C} + \lambda \mathbf{P}^{-1}\tilde{\mathbf{m}}$
 $\tilde{\mathbf{M}} = \tilde{\mathbf{C}} + \lambda \begin{bmatrix} \mathbf{P}^{-1}\tilde{\mathbf{m}} \\ 0 \end{bmatrix}$ (in projective coordinates)

is a point at infinity AND belongs to the optical ray
 $\mathbf{V} = \mathbf{P}^{-1}\tilde{\mathbf{m}}$

Appendix 3 - Rectification Homographies by optical rays (Appendix 2)

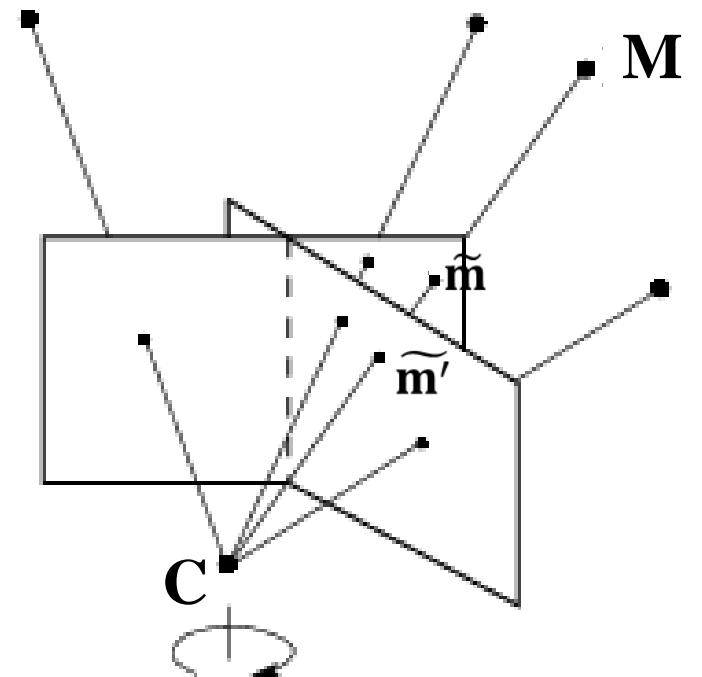
- We must compute the transformations that map the original images, i.e. those acquired by the cameras with PPMs, $\tilde{\mathbf{P}}_L$ and $\tilde{\mathbf{P}}_R$, into a rectified pair, i.e. the two images that would have been acquired by the cameras with PPMs $\tilde{\mathbf{P}}'_L$ and $\tilde{\mathbf{P}}'_R$.
- Given a 3D point in the WRF (i.e. the SRF), for each of the two cameras we can write the equation of the optical ray through the point based on both $\tilde{\mathbf{P}}$ and $\tilde{\mathbf{P}}'$ (subscripts L,R are omitted for the sake of simplicity).

$$\begin{aligned}\tilde{\mathbf{P}} = [\mathbf{P} \ p_4] &\rightarrow \mathbf{M} = \mathbf{C} + \lambda \mathbf{P}^{-1} \tilde{\mathbf{m}} \\ \tilde{\mathbf{P}}' = [\mathbf{P}' \ p'_4] &\rightarrow \mathbf{M} = \mathbf{C} + \lambda \mathbf{P}'^{-1} \tilde{\mathbf{m}}'\end{aligned}$$

$$\mathbf{P}^{-1} \tilde{\mathbf{m}} = \lambda \mathbf{P}'^{-1} \tilde{\mathbf{m}}'$$

$$\tilde{\mathbf{m}} = \lambda \mathbf{P} \mathbf{P}'^{-1} \tilde{\mathbf{m}}'$$

$$H_L = \mathbf{P}_L \mathbf{P}'_L^{-1}, \quad H_R = \mathbf{P}_R \mathbf{P}'_R^{-1}$$



Homography

Main References



- 1) R. Hartley and A. Zisserman, “Multiple View Geometry in Computer Vision”, 2° Edition, Cambridge University Press, 2003.
- 2) Z. Zhang, “A flexible new technique for camera calibration”, IEEE Trans. On Pattern Analysis and Machine Intelligence, 22(11):1330-1334, November 2000.
- 3) A. Fusiello, “Visione Computazionale – Tecniche di ricostruzione tridimensionale”, Collana Informatica Franco Angeli, 2013.
- 4) Bradski, Adrian Kaehler, “Learning OpenCV - Computer Vision with the OpenCV Library”, O'Reilly Media, 2008.