

Camera calibration is an essential process in computer vision that aims to determine the intrinsic and extrinsic parameters of a camera.

These parameters are crucial for accurately mapping 3D points in the scene to 2D points in the image

# Lecture 2

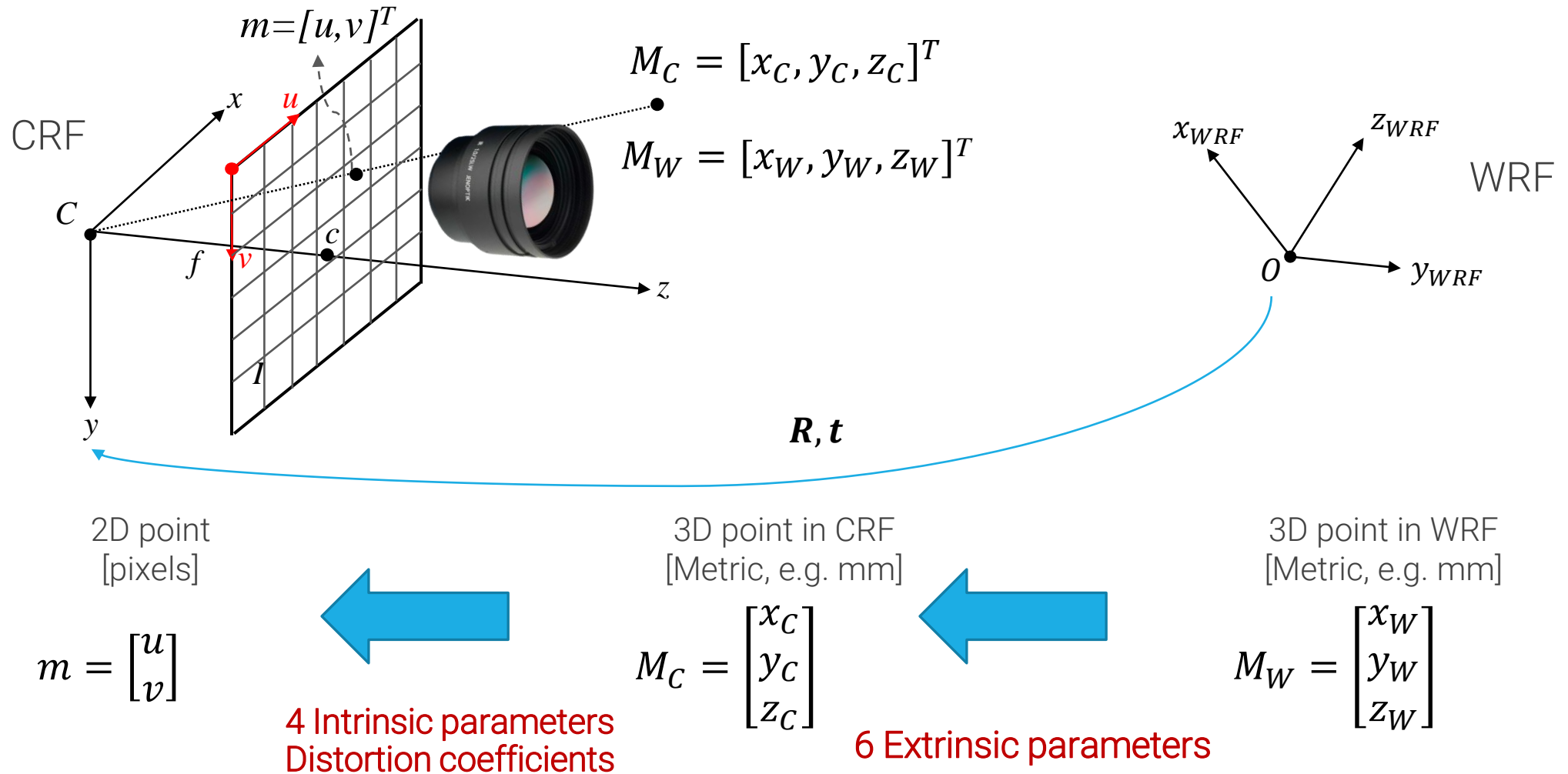
## Camera Calibration

---

IMAGE PROCESSING AND COMPUTER VISION – PART 2

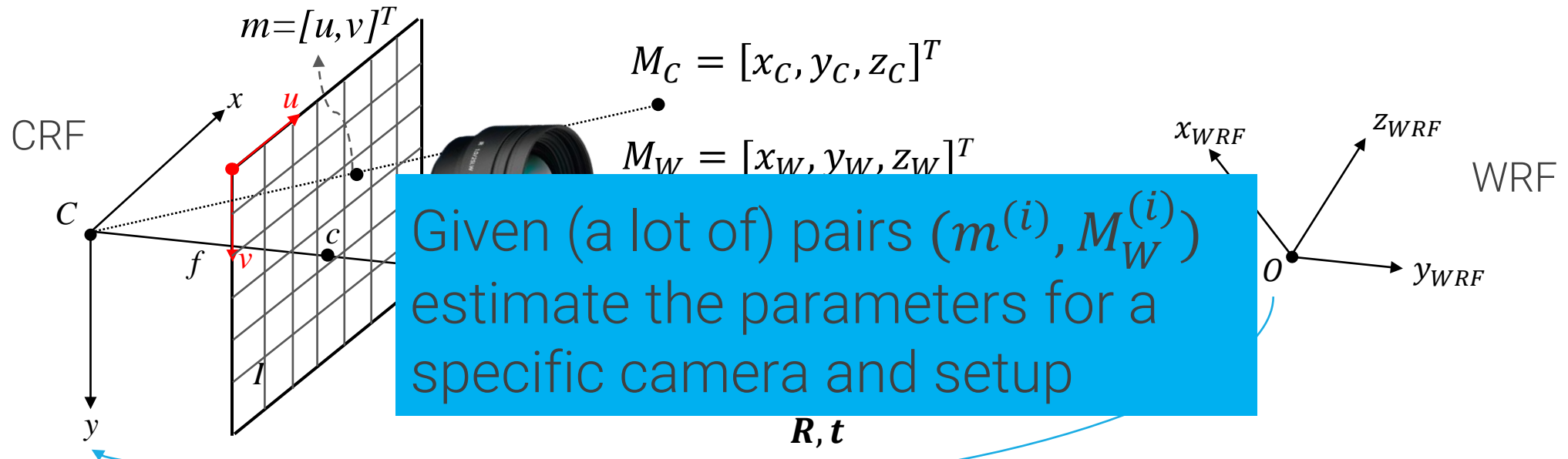
SAMUELE SALTI

# Complete camera model



# Camera calibration

given projected points and WRF points (their coordinates) estimate the parameters involved in the image formation process



Known

2D point  
[pixels]

$$m = \begin{bmatrix} u \\ v \end{bmatrix}$$

4 Intrinsic parameters  
Distortion coefficients

3D point in CRF  
[Metric, e.g. mm]

$$M_C = \begin{bmatrix} x_C \\ y_C \\ z_C \end{bmatrix}$$

Unknown

6 Extrinsic parameters

3D point in WRF  
[Metric, e.g. mm]

$$M_W = \begin{bmatrix} x_W \\ y_W \\ z_W \end{bmatrix}$$

Known

how retrieve correspondances between  $m$  and  $M$ ?  
- using calibration patterns

# Calibration patterns

the calibration depends on the type of calibration objects and the calibration patterns used

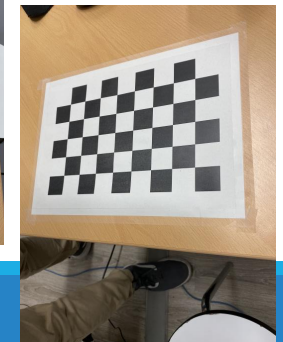
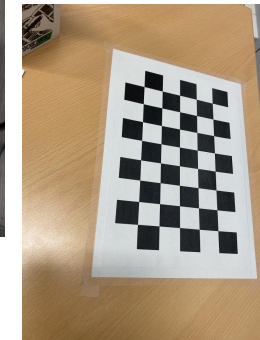
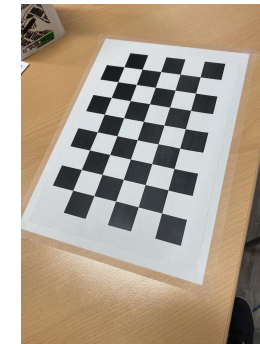
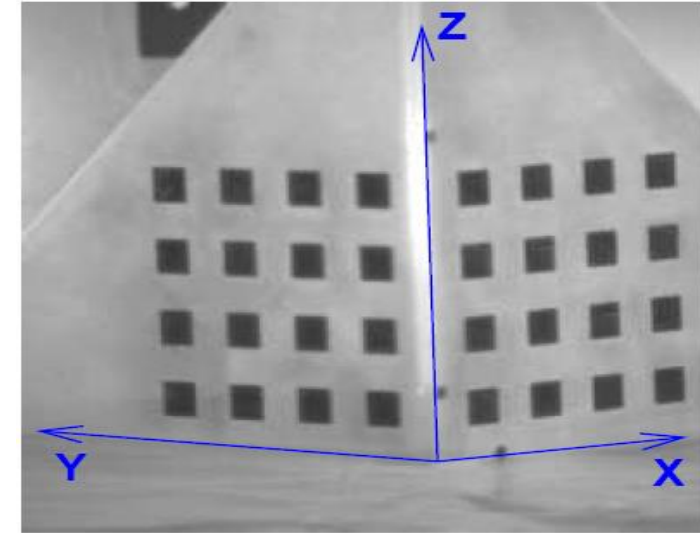
Camera calibration approaches can be split into two main categories:

- Those relying on **a single image of a 3D calibration object**
  - featuring several (at least 2) planes containing a known pattern like a checkerboard which has a recognizable pattern
- Those relying on **several (at least 3) different images of one given planar pattern** from different positions and angles

In practice, 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

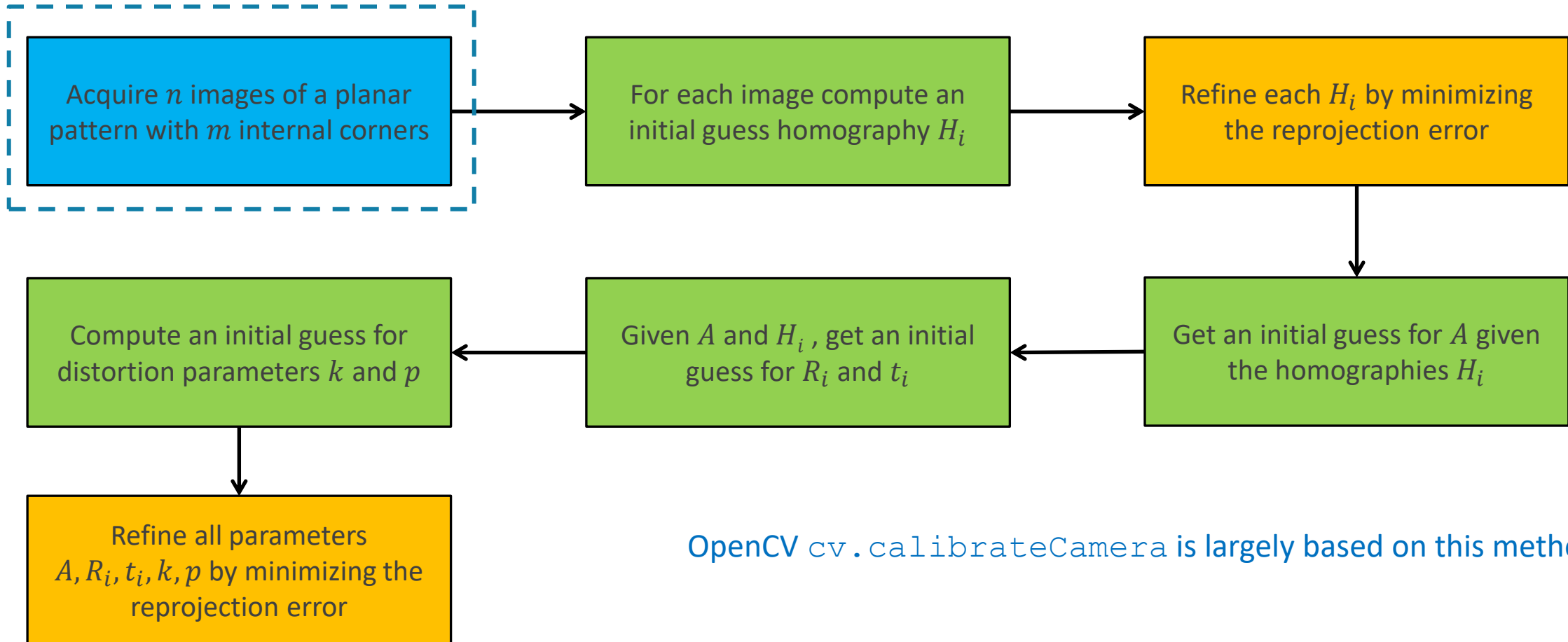
- the main Computer Vision toolboxes include specific functions (OpenCV, Matlab CC Toolbox)



# Zhang's Method

green: minimize algebraic error  
solve a linear system with least square

yellow: minimize geometric error  
solve for minimization of L2-norms



OpenCV `cv.calibrateCamera` is largely based on this method

Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.  
<https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/tr98-71.pdf>

1° step: acquire input images

# Calibration pattern

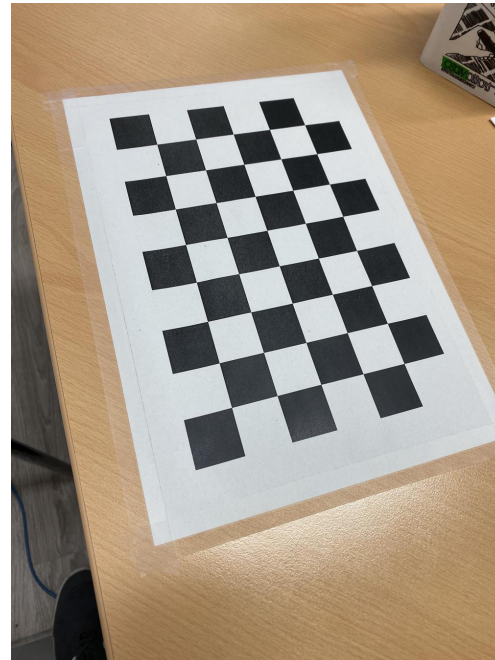
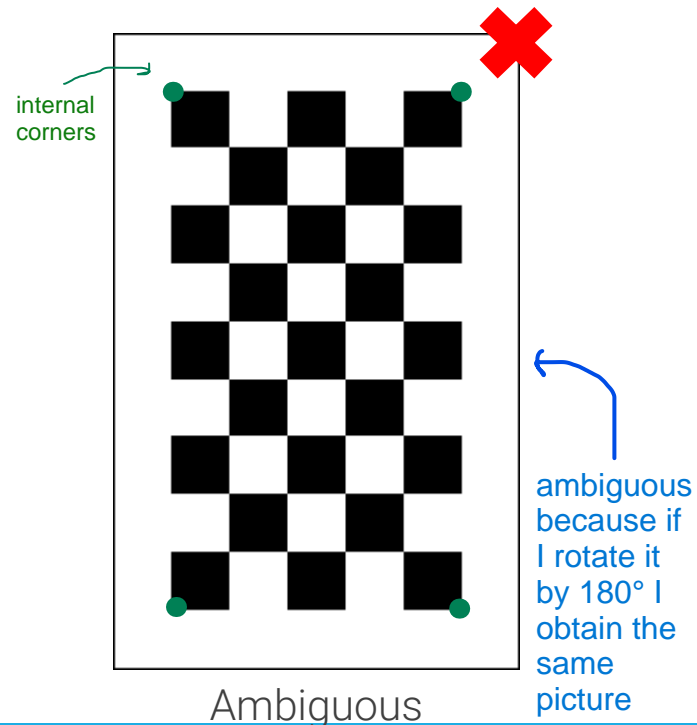
ambiguous: if transformation applied, no clear distinction  
e.g. chessboard if rotated we cannot distinguish differences wrt the two views

Given a chessboard **pattern**, we know:

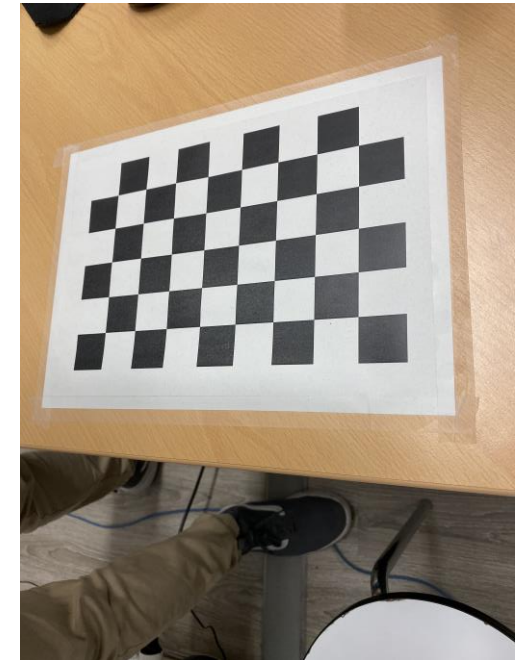
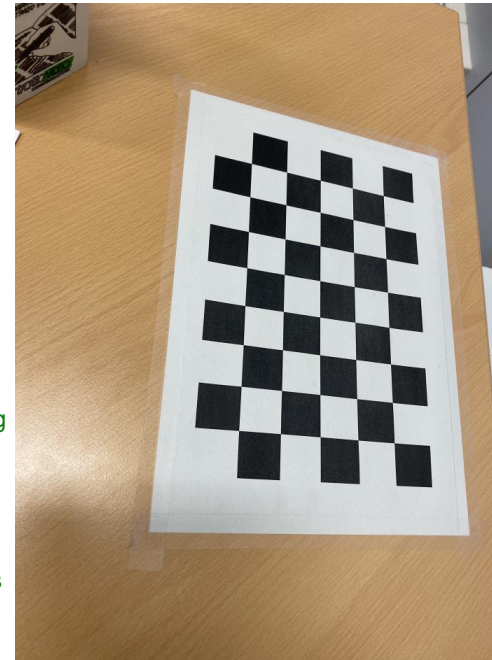
- The **number of internal corners** of the pattern, usually odd along one dimension and even along the other to remove rotation ambiguities.
- The **size of the squares** that form the pattern (in mm, cm...)

doing this we retrieve the 2d coordinates we need (m)

Internal corners can be detected easily by standard algorithms (e.g. the Harris corner detector)



these images are easier to disambiguate after transformations like rotations





in the previous slide we have found 2D coordinates (m) now we want  $M_W$ , the 3D coords from WRF

# 3D coordinates

attach the WRF to the pattern itself in order to conveniently define the 3d points

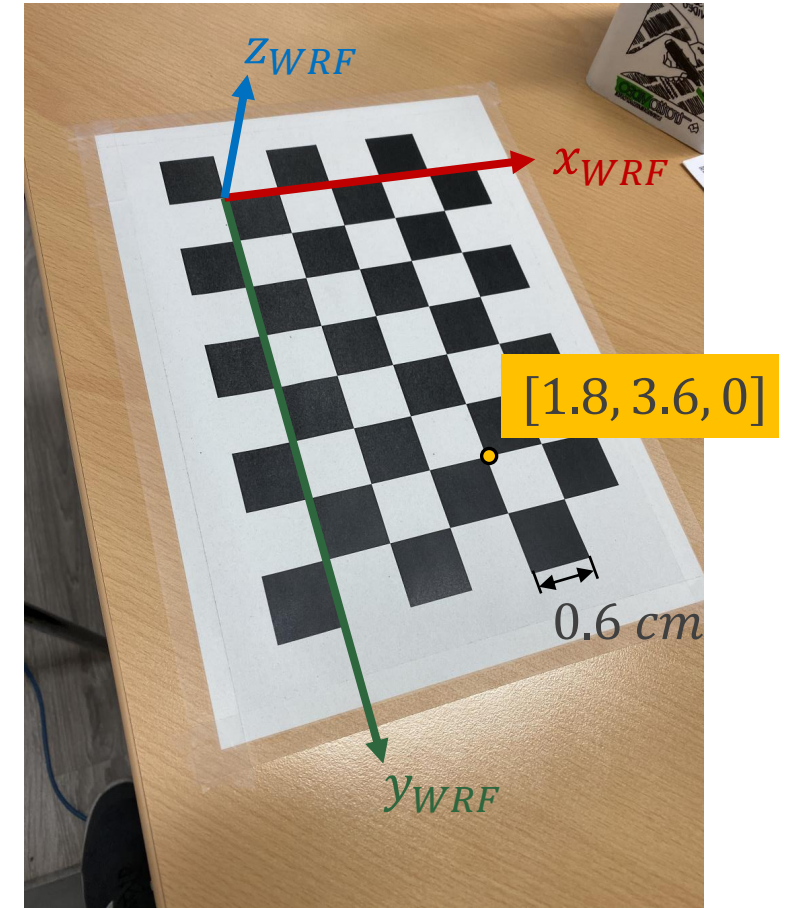
The World Reference Frame can be conveniently defined.

In an unambiguous pattern, the WRF can be defined so that

- it has its origin always in the same corner (e.g., the one next to the dark square on the right of the chessboard if both dark squares are on top);
- its plane  $z = 0$  is the pattern itself => the third coordinate is always 0;
- the  $x, y$  axes are aligned to the chessboard (e.g.,  $x$  along the short side and  $y$  along the long one).

Given such rules and the known square side, it is possible to define 3D coordinates for all corners in an image of the pattern.

This setup simplifies the problem because all the points on the pattern lie on a single plane, and we can set the  $z$ -coordinate of these points to 0. This assumption allows us to work with 2D coordinates on the plane, which can be extended to 3D homogeneous coordinates.



in order to estimate properly the coefficients of radial distortion, the pattern should appear in the corners, because I can get correspondences only where the pattern is and I can estimate params

# Extrinsic parameters

since we collect several images

The World Reference Frame is different for each calibration image.

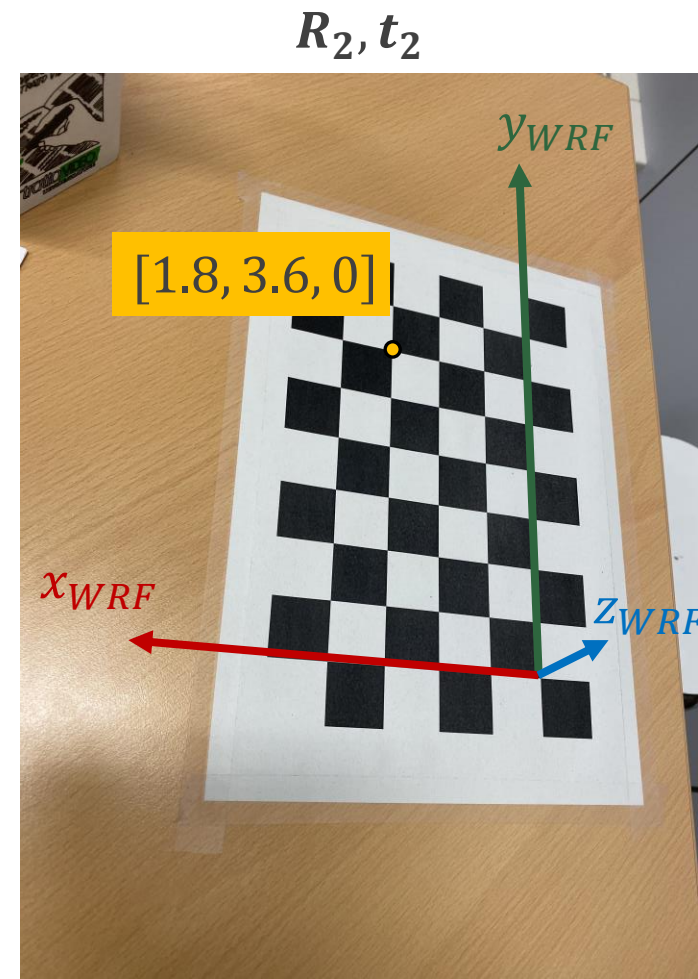
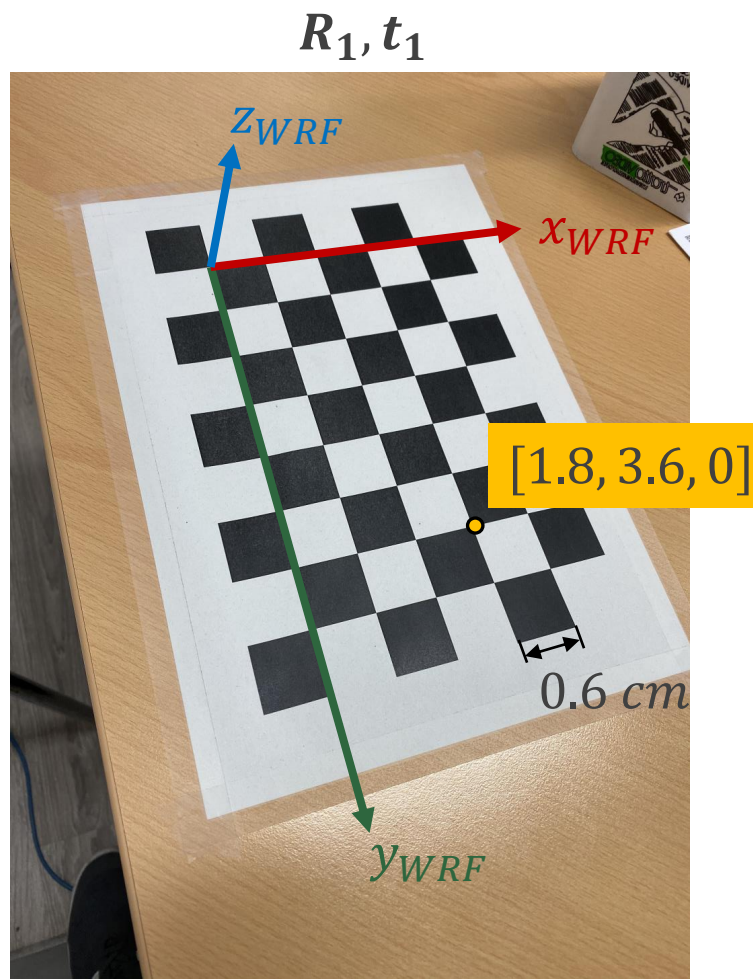
The  $[R \ t]$  are estimated wrt the World Reference Frame attached to the target, which moves with the pattern. what it change is  $[R \ t]$  !

Therefore, we estimate as many extrinsic matrices as the number of images used for calibration (usually 10 to 20 images of the pattern)

we have to estimate, if we collect 'n' frames, '1' intrinsic matrix and 'n' extrinsic matrices

the WRF is always the same

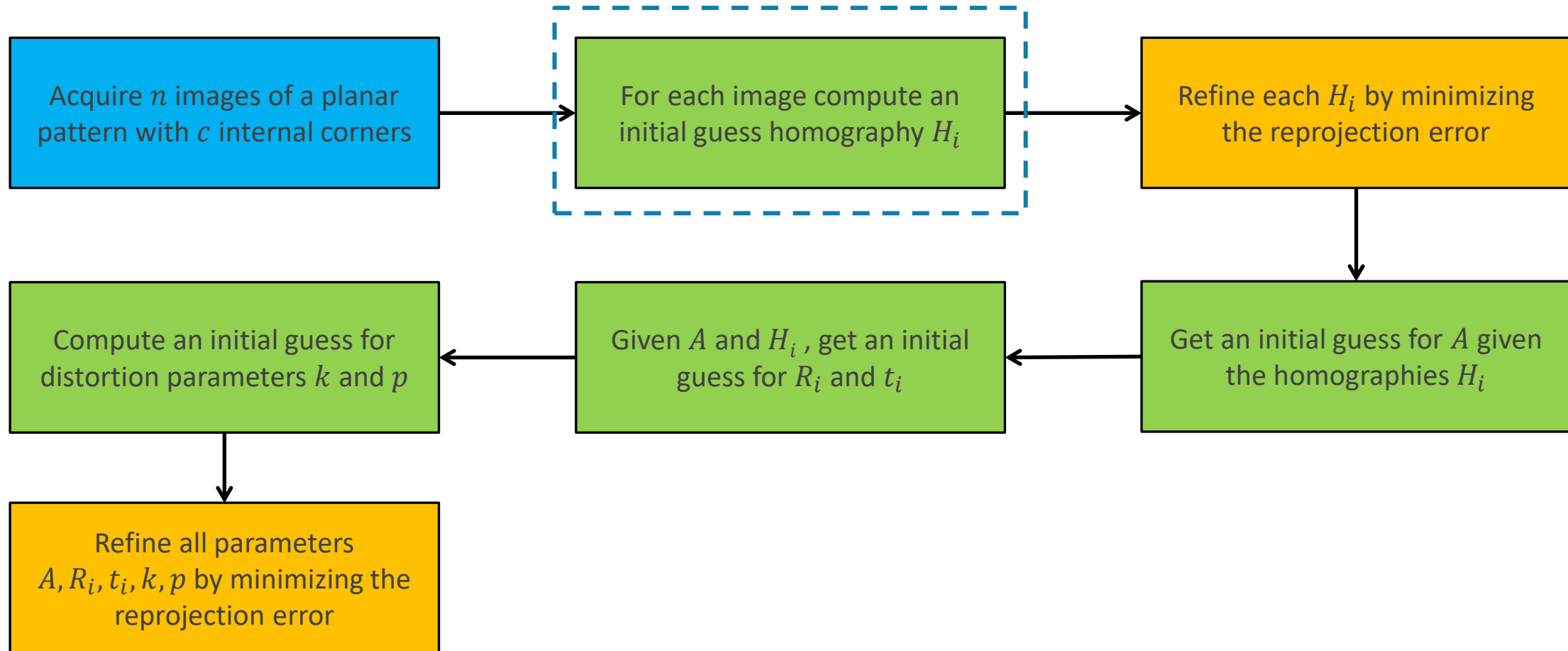
- what is different is the position of the WRF wrt the Camera, according to different translations and rotations in each image, because we have to take different pictures where the pattern is located at different positions wrt the camera





remember: PPMs are always the same in the image plane

# Zhang's Method



Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.

convert both the 2d and 3d coords to homogeneous coordinates

# $P$ as a Homography

For each image, compute the homography matrix  $H$  that relates the points on the planar pattern in the WRF to the points in the image plane.

Due to the choice of the WRF associated with calibration images, in each of them we consider only 3D points with  $z = 0$

Accordingly, the PPM for points on the pattern can be simplified to a 3x3 matrix:

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

*Annotations:*

- A blue 'X' is drawn over the third column of the matrix  $\mathbf{P}$ .
- A blue box highlights the '0' in the third row of the vector  $\begin{bmatrix} x \\ y \\ 0 \\ 1 \end{bmatrix}$ , with the label "translation" below it.
- A green arrow points from the text "the 3d coords of the point in the planar pattern" to the vector  $\begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$ .
- A blue box contains the text "We look at a planar scene".

Such a transformation, denoted here as  $\mathbf{H}$ , is known as **homography** and represents a general transformation between **projective planes** the world plane and the image plane

$\mathbf{H}$  can be thought of as a simplification of  $\mathbf{P}$  in case the imaged object is **planar**

The homography matrix  $H$  transforms points on the world plane (pattern plane) to the image plane. It can be viewed as a plane-to-plane projective transformation, similar to how the projection matrix (PPM) relates 3D points to 2D image points in general.

# Estimating $H_i$ (DLT algorithm)

Given the  $i$ -th image of a pattern with  $c$  corners, we can write **3 linear equations for each corner  $j$**  where:

- 3D coordinates are known due to the WRF definition
- 2D coordinates are known due to corners having been detected in the  $i$ -th image
- the unknowns are the 9 elements in  $H_i$

$$\begin{array}{c}
 \text{u in the frame i} \\
 \text{corner j} \\
 \tilde{m}_{ij} \equiv \begin{bmatrix} u_{ij} \\ v_{ij} \\ 1 \end{bmatrix} \equiv \begin{bmatrix} p_{i,1,1} & p_{i,1,2} & p_{i,1,4} \\ p_{i,2,1} & p_{i,2,2} & p_{i,2,4} \\ p_{i,3,1} & p_{i,3,2} & p_{i,3,4} \end{bmatrix} \begin{bmatrix} x_j \\ y_j \\ 1 \end{bmatrix} \equiv H_i \tilde{w}_j \equiv \begin{bmatrix} -h_{i1}^T & - \\ -h_{i2}^T & - \\ -h_{i3}^T & - \end{bmatrix} \tilde{w}_j \equiv \begin{bmatrix} h_{i1}^T \tilde{w}_j \\ h_{i2}^T \tilde{w}_j \\ h_{i3}^T \tilde{w}_j \end{bmatrix}
 \end{array}$$

Unknowns
break down H into three rows
dot product between  $h_i$  and  $w_j$

Coordinates of a corner in the image

Coordinates of the same corner in the WRF

what changes is [R t]

Stacking  $3c$  such equations for the  $c$  corners we get a system of equations, but...

each corner will give me 3 equations

How do we solve a system of equations in a projective space?

then up to an unknown scale factor?

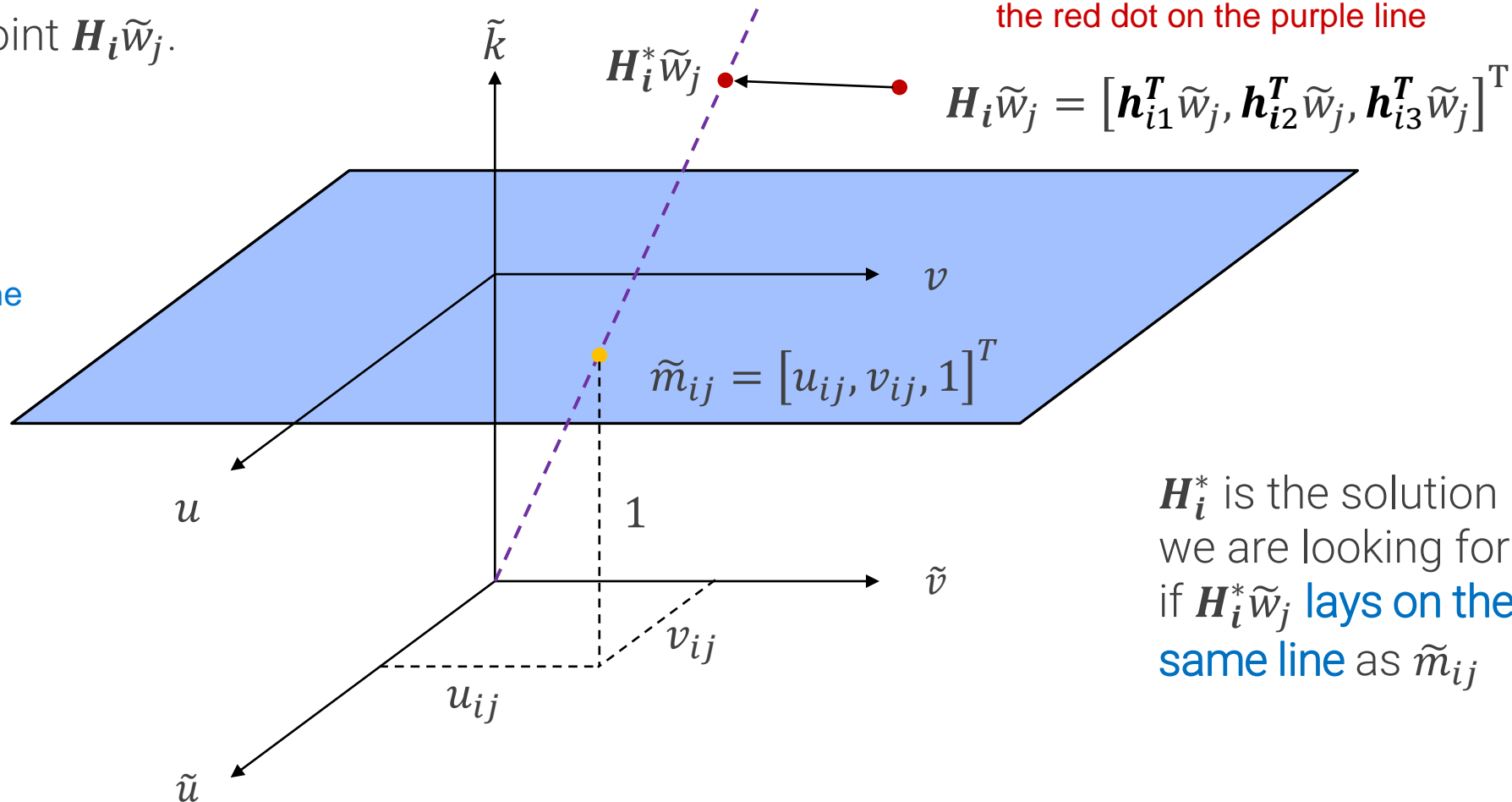
we basically have to solve a linear system

In projective geometry, two points are considered equivalent if they lie on the same line, meaning they are collinear

# When are two 3D points equivalent in $\mathbb{P}^2$ ?

By deciding a value for the 9 entries of  $\mathbf{H}_i$  we move the point  $\mathbf{H}_i \tilde{\mathbf{w}}_j$ .

we want that red dot and the yellow one to be the same



so that the transformation that maps the two dots up to an unknown scalar factor

we have to find the transformation  $\mathbf{H}$  that put the red dot on the purple line

$\mathbf{H}_i^*$  is the solution we are looking for if  $\mathbf{H}_i^* \tilde{\mathbf{w}}_j$  lays on the same line as  $\tilde{\mathbf{m}}_{ij}$

# Estimating $H_i$ (DLT algorithm)

the cross product measures the area of the parallelogram between two vectors

Two points lay on the same line if their cross product is the zero vector.

$$\tilde{m}_{ij} \equiv \mathbf{H}_i \tilde{\mathbf{w}}_j \Rightarrow \tilde{m}_{ij} \times \mathbf{H}_i \tilde{\mathbf{w}}_j = \mathbf{0} \Rightarrow \tilde{m}_{ij} \times \mathbf{H}_i \tilde{\mathbf{w}}_j = \begin{bmatrix} u_{ij} \\ v_{ij} \\ 1 \end{bmatrix} \times \begin{bmatrix} \mathbf{h}_{i1}^T \tilde{\mathbf{w}}_j \\ \mathbf{h}_{i2}^T \tilde{\mathbf{w}}_j \\ \mathbf{h}_{i3}^T \tilde{\mathbf{w}}_j \end{bmatrix} = \begin{bmatrix} v_{ij} \mathbf{h}_{i3}^T \tilde{\mathbf{w}}_j - \mathbf{h}_{i2}^T \tilde{\mathbf{w}}_j \\ \mathbf{h}_{i1}^T \tilde{\mathbf{w}}_j - u_{ij} \mathbf{h}_{i3}^T \tilde{\mathbf{w}}_j \\ u_{ij} \mathbf{h}_{i2}^T \tilde{\mathbf{w}}_j - v_{ij} \mathbf{h}_{i1}^T \tilde{\mathbf{w}}_j \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

cross product

3x9 matrix

9x1 column vector

Only 2 equations are linearly independent (multiply first one by  $-u$  and second one by  $-v$ , then sum them to get the third one)

$$\begin{bmatrix} \mathbf{0}_{3 \times 1}^T & -\tilde{\mathbf{w}}_j^T & v_{ij} \tilde{\mathbf{w}}_j^T \\ \tilde{\mathbf{w}}_j^T & \mathbf{0}_{3 \times 1}^T & -u_{ij} \tilde{\mathbf{w}}_j^T \\ -v_{ij} \tilde{\mathbf{w}}_j^T & u_{ij} \tilde{\mathbf{w}}_j^T & \mathbf{0}_{3 \times 1}^T \end{bmatrix} \begin{bmatrix} \mathbf{h}_{i1} \\ \mathbf{h}_{i2} \\ \mathbf{h}_{i3} \end{bmatrix} = \mathbf{0}_{3 \times 1}$$

because each  $\mathbf{h}_i$  is a 3x1 vector

$$\mathbf{h}_k^T \tilde{\mathbf{w}}_j = \tilde{\mathbf{w}}_j^T \mathbf{h}_k$$

2x9 matrix

$$\begin{bmatrix} \mathbf{0}_{3 \times 1}^T & -\tilde{\mathbf{w}}_j^T & v_{ij} \tilde{\mathbf{w}}_j^T \\ \tilde{\mathbf{w}}_j^T & \mathbf{0}_{3 \times 1}^T & -u_{ij} \tilde{\mathbf{w}}_j^T \end{bmatrix} \begin{bmatrix} \mathbf{h}_{i1} \\ \mathbf{h}_{i2} \\ \mathbf{h}_{i3} \end{bmatrix} = \mathbf{0}_{2 \times 1}$$



# Estimating $H_i$ (DLT algorithm)

Given  $c$  corners, we can create a  $Ax = 0$  form **homogeneous, overdetermined** linear system of equations more equations than unknowns

$$\begin{bmatrix} \mathbf{0}_{3 \times 1}^T & -\tilde{\mathbf{w}}_1^T & v_{i1} \tilde{\mathbf{w}}_1^T \\ \tilde{\mathbf{w}}_1^T & \mathbf{0}_{3 \times 1}^T & -u_{i1} \tilde{\mathbf{w}}_1^T \\ \vdots & \vdots & \vdots \\ \mathbf{0}_{3 \times 1}^T & -\tilde{\mathbf{w}}_c^T & v_{ic} \tilde{\mathbf{w}}_c^T \\ \tilde{\mathbf{w}}_c^T & \mathbf{0}_{3 \times 1}^T & -u_{ic} \tilde{\mathbf{w}}_c^T \end{bmatrix} \begin{bmatrix} h_{i1} \\ h_{i2} \\ h_{i3} \end{bmatrix} = \mathbf{0}_{2c \times 1} \Rightarrow L_i \mathbf{h}_i = \mathbf{0}$$

overdetermined systems have always one exact solution, that is 0

9x1 column vector

2cx9 matrix

To avoid the trivial solution  $\mathbf{h} = \mathbf{0}$  we look for solutions with an additional constraint, e.g.,  $\|\mathbf{h}\| = 1$ .

generally, in strict math terms, this is an impossible system to solve  
then we relax our constraints saying that H doesn't have to map perfectly the points onto the two different projective planes, but as close as possible, geometrically speaking.

the norm of h equal to one to avoid 0 solution trivial

algebraically, we solve this problem using Least Square, as we see in the following slide

# Singular Value Decomposition

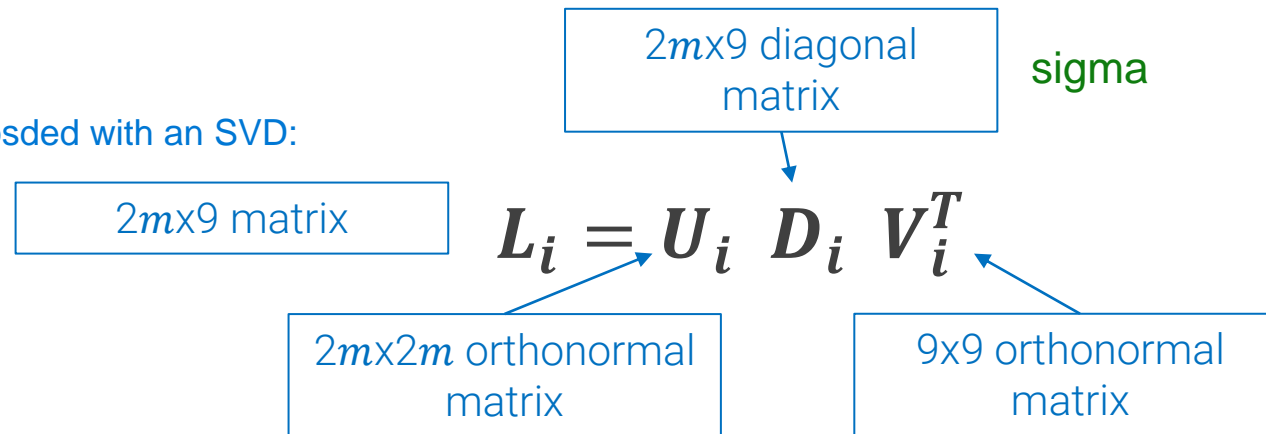
The solution  $\mathbf{h}^*$  is found by minimizing the norm of the vector  $\mathbf{L}_i \mathbf{h}_i$  it becomes a Least Square problem

$$\mathbf{h}_i^* = \underset{\mathbf{h}_i \in \mathbb{R}^9}{\operatorname{argmin}} \|\mathbf{L}_i \mathbf{h}_i\| \text{ s.t. } \|\mathbf{h}_i\| = 1$$

It is known from linear algebra that the solution to such problem can be found via Singular Value Decomposition of  $\mathbf{L}_i$ . In particular, the solution is  $\mathbf{h}_i^* = \mathbf{v}_9$ , i.e., the last column of  $\mathbf{V}_i$

every matrix, which transforms between spaces, can be decomposed with an SVD:

- rotate to a convenient ref. sys
- scale axis independently
- rotate to another position

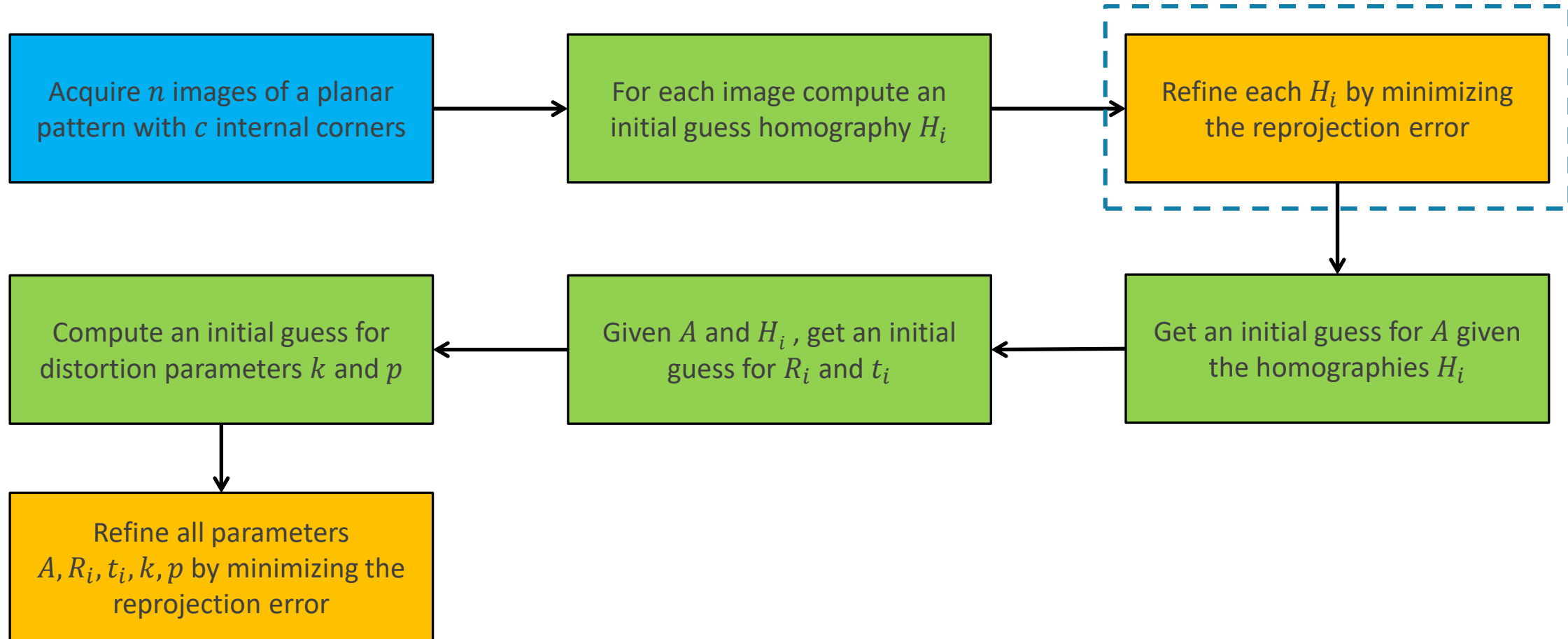


every image will have a different homography

# Zhang's Method

so DLT gave an initial guess for  $H_i$

the reprojection error is what we have interest into. It quantifies how accurately the estimated camera parameters and 3D world coordinates can reproduce the observed image points.

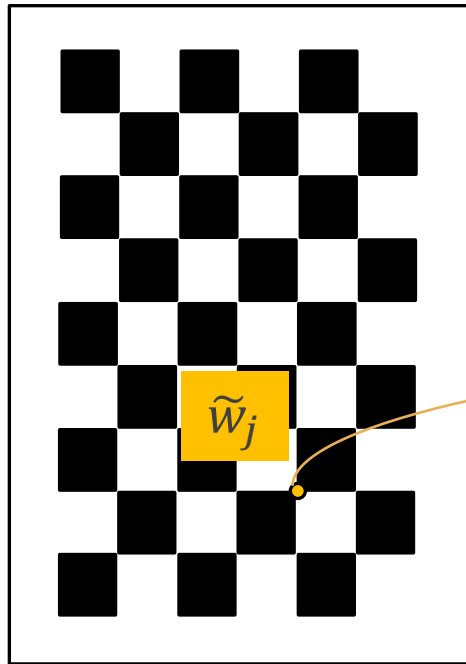


Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.

the mapping from  $w_j$  and  $m_j$  cannot be linear, because we have consider also noise and lens distortions -> then the red dot cannot coincide with the blue dot -> then we approximate this mapping ignoring them just to have a starting point

# What error should we minimize?

I am approximating the mapping as linear because atm I am not considering lenses, which inrtoduce non-lienar distortions



the blue dot has two coordinates computed through harris corner detector

$H_i$

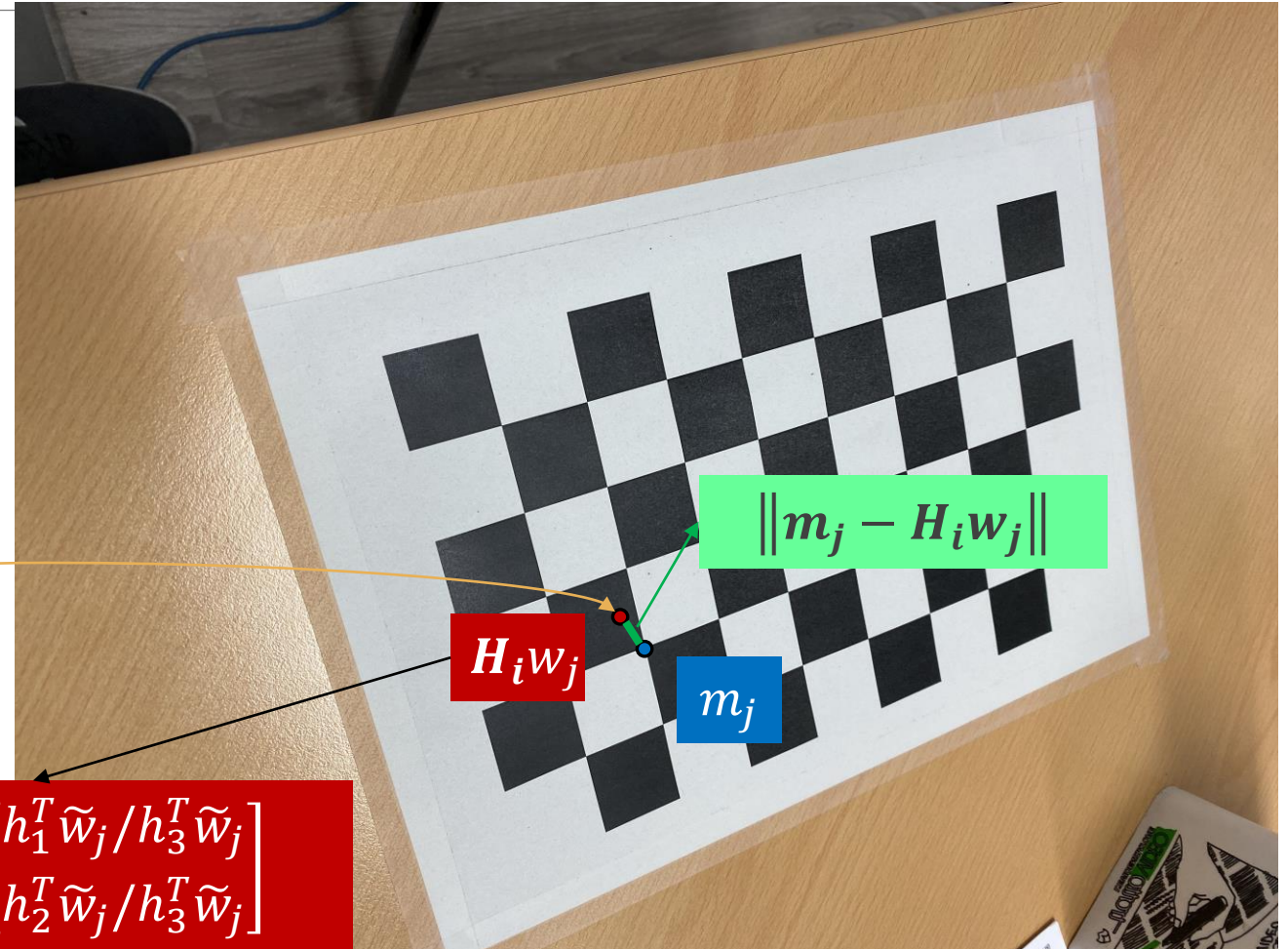
divide by third-coord  
to get out the  
projective space and  
returning bak to the  
euclidean

$$H_i w_j = \begin{bmatrix} h_1^T \tilde{w}_j / h_3^T \tilde{w}_j \\ h_2^T \tilde{w}_j / h_3^T \tilde{w}_j \end{bmatrix}$$

is a 2x1 Euclidean vector

$H$  is a mapping between a pattern (and a corner) and the corner in another image

$H_i w_j$  is the projected corner and  $m_j$  is the real corner



the ideal mapping maps as close as possible the red dot and the blue dot

minimize a geometric error

# Non-linear refinement of $H_i$

but it is non-linear, how to solve it?  
- iterative algo, like GD, SGD etc.

Given the initial guess for  $H_i$ , we can refine it by a non-linear minimization problem:

the main drawback that regards iterative algorithms is that they depend too much on the starting point, namely the initial guess

$$\mathbf{H}_i^* = \underset{H_i}{\operatorname{argmin}} \sum_{j=1}^c \|\mathbf{m}_{ij} - \mathbf{H}_i \mathbf{w}_j\|^2 \quad i = 1, \dots, n$$

which can be solved for by using an iterative algorithm, like the Levenberg-Marquardt algorithm.

This additional optimization step corresponds to the minimization of the reprojection error (typically referred to as geometric error) measured for each of the 3D corners 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.

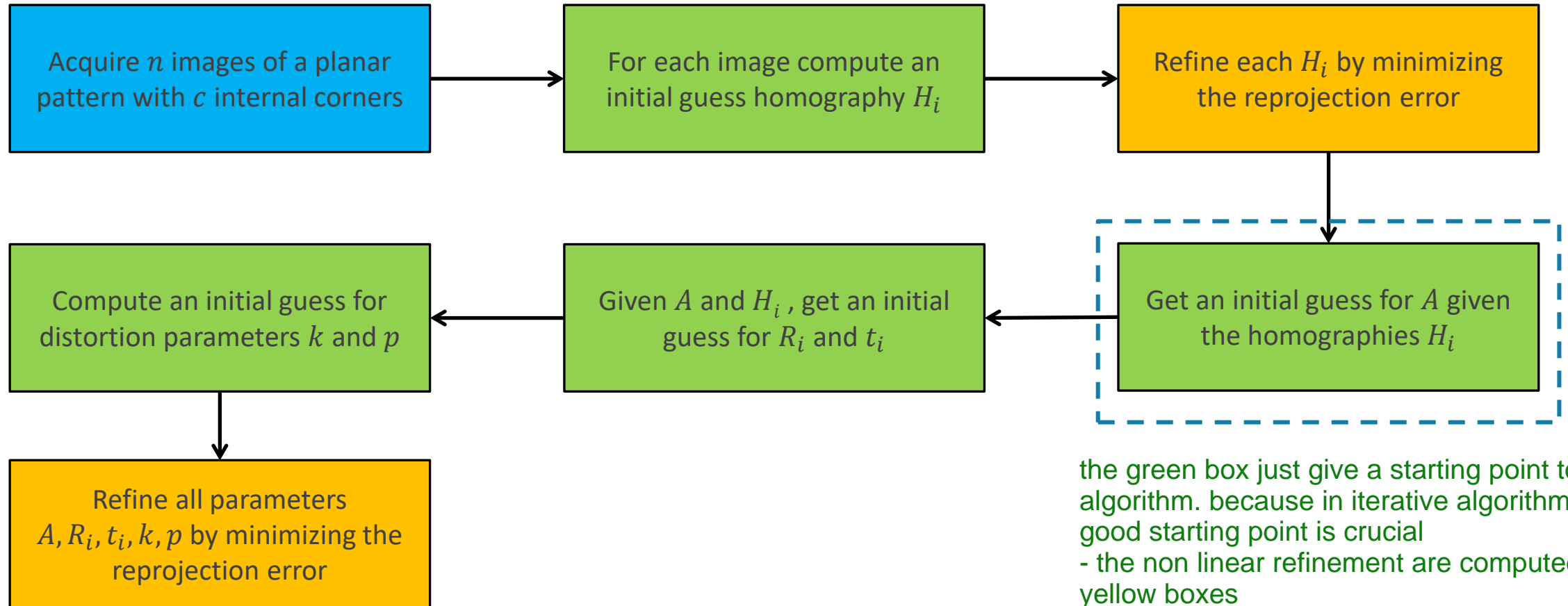
The error minimized to estimate the initial guess when solving the linear system is instead referred to as algebraic error or distance. Solutions based on minimization of the algebraic error may not be aligned with our intuition, yet there exist a unique solution, which is cheap to compute. Hence, they are a good starting point for a geometric, non-linear minimization, which effectively minimize the distance we care about.



# Zhang's Method

until now we have estimated 3 of the 4 columns of the PPM

in order to estimate the last one, we have to estimate the intrinsic and extrinsic parameters independently



the green box just give a starting point to the algorithm. because in iterative algorithms finding a good starting point is crucial  
- the non linear refinement are computed into the yellow boxes

Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.

camera calibration wants to estimate all the parameters. we already estimated the homography  $H$ , then let's estimate the last column of  $P$ , that contains the intrinsic parameters of the lens, which we have neglected before.

# Estimation of the intrinsic parameters

All the images acquired for calibration share the same **intrinsic parameters** parameters of the lenses

We can establish the following relations between them <sup>the images</sup> and the extrinsic and intrinsic parameters <sup>because we said that  $H$  is  $P$  when we have planar objects</sup>

the PPM of the  $i$ -th calibration frame

$$P_i \equiv A[R_i | t_i] = A[r_{i1} \quad r_{i2} \quad \cancel{r_{i3}} \quad t_i] \Rightarrow H_i = [h_{i1} \quad h_{i2} \quad h_{i3}] = [kAr_{i1} \quad kAr_{i2} \quad kAt_i]$$

remember: we have only one set of intrinsic params ( $A$ ) and different extrinsic parameters for each frame

$$\Rightarrow k r_{i1} = A^{-1} h_{i1} \quad k r_{i2} = A^{-1} h_{i2}$$

<sup>at this step we don't know these</sup> <sup>we have only estimated  $H$  up to now</sup>

<sup>why do I need to put a scale factor  $k$ ? because the equivalence between projective spaces holds by a scale factor  $k$</sup>

<sup>$k r_{i1}$  and  $k r_{i2}$  has to be orthogonal to each other because they are columns of an orthogonal rotation matrix  $R_i$</sup>

Since the column vectors of each  $R_i$  are **orthonormal**, we get the following constraints

<sup>so if they are orthogonal also the RHS must be orthogonal</sup>

$$1) \langle r_{i1}, r_{i2} \rangle = 0 \Rightarrow \langle A^{-1} h_{i1}, A^{-1} h_{i2} \rangle = 0 \Rightarrow h_{i1}^T A^{-T} A^{-1} h_{i2} = 0$$

<sup>i'll have  $h_i$  for each image</sup>

$$2) \langle r_{i1}, r_{i1} \rangle = \langle r_{i2}, r_{i2} \rangle \Rightarrow \langle A^{-1} h_{i1}, A^{-1} h_{i1} \rangle = \langle A^{-1} h_{i2}, A^{-1} h_{i2} \rangle \Rightarrow h_{i1}^T A^{-T} A^{-1} h_{i1} = h_{i2}^T A^{-T} A^{-1} h_{i2}$$

<sup>they also have unit length</sup> <sup>the norm of these two vectors must be the same</sup> <sup>therefore in order to calibrate I need at least three images</sup>

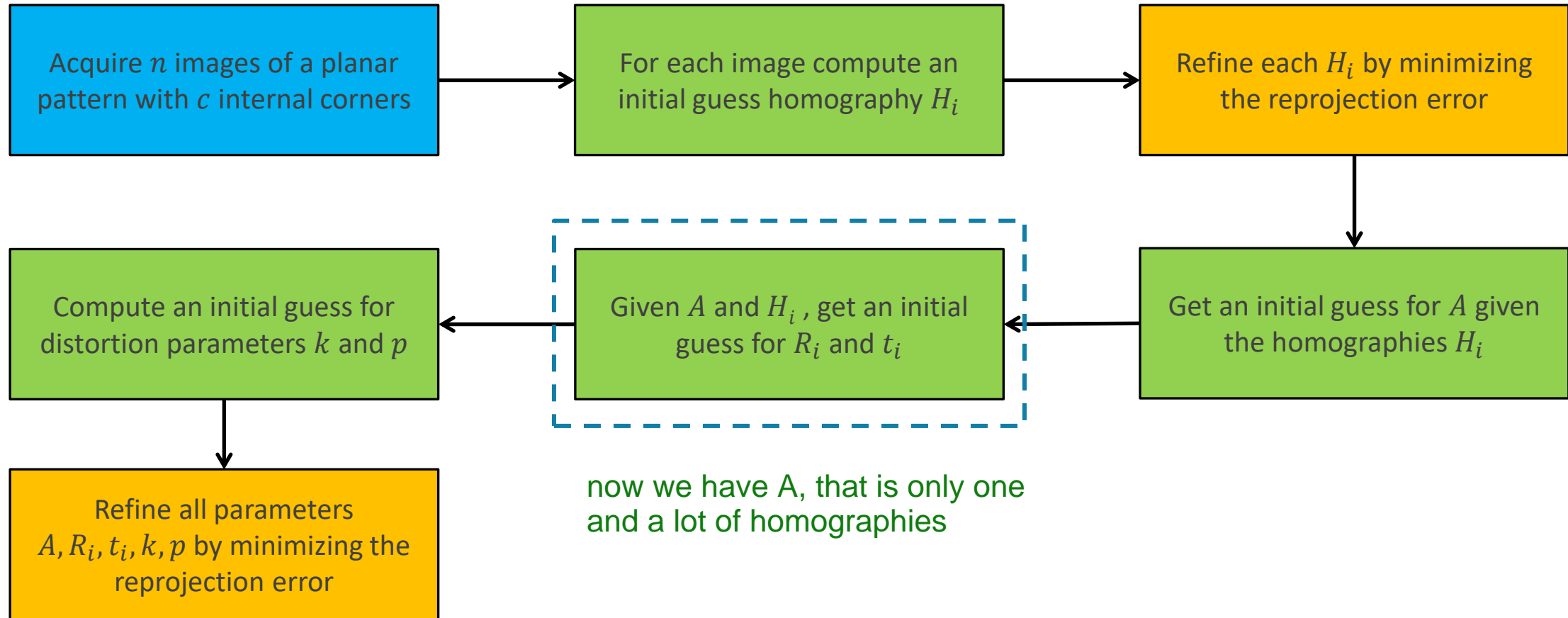
When you have multiple images of the same scene or object, each image provides additional constraints on the parameters to be estimated.

By stacking these **two constraints for each image**, we get a homogeneous system of equations which can be solved again by SVD if  $n \geq 3$  images are collected (6 unknowns since  $A^{-T} A^{-1}$  is symmetric).

I cannot say that the norm of  $A^{-1} h_i$  is 1 because I would neglect  $k \rightarrow$  but I can say that the norms of  $k r_{i1}$  and  $k r_{i2}$  must be the same (obv not 1) because  $k$  is the same

I need at least 3 images because the system of equations formed by stacking the constraints may not be overdetermined, meaning there may not be enough constraints to uniquely determine the unknown parameters. Adding more images allows for redundancy in the equations, making the system overdetermined and improving the stability and accuracy of the solution.

# Zhang's Method



Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.

# Estimation of the extrinsic parameters

let's retrieve an initial guess for rotation and translation parameters

Once  $\mathbf{A}$  has been estimated, it is possible to compute  $\mathbf{R}_i$  and  $\mathbf{t}_i$  (for each image) given  $\mathbf{A}$  and the previously computed homography  $\mathbf{H}_i$ :

$$\mathbf{H}_i = [\mathbf{h}_{i1} \quad \mathbf{h}_{i2} \quad \mathbf{h}_{i3}] = [k\mathbf{A}\mathbf{r}_{i1} \quad k\mathbf{A}\mathbf{r}_{i2} \quad k\mathbf{A}\mathbf{t}_i] \Rightarrow \mathbf{r}_{i1} = \frac{1}{k}\mathbf{A}^{-1}\mathbf{h}_{i1}$$

the only thing i miss to compute R is k

As  $\mathbf{r}_{i1}$  is a unit vector, the normalization constant can be computed as  $k = \|\mathbf{A}^{-1}\mathbf{h}_{i1}\|$

Then, the same constant can be used to compute

not perfectly unit vector because of noise etc.

$$\mathbf{r}_{i2} = \frac{1}{k}\mathbf{A}^{-1}\mathbf{h}_{i2} \quad \text{and} \quad \mathbf{t}_i = \frac{1}{k}\mathbf{A}^{-1}\mathbf{h}_{i3}$$

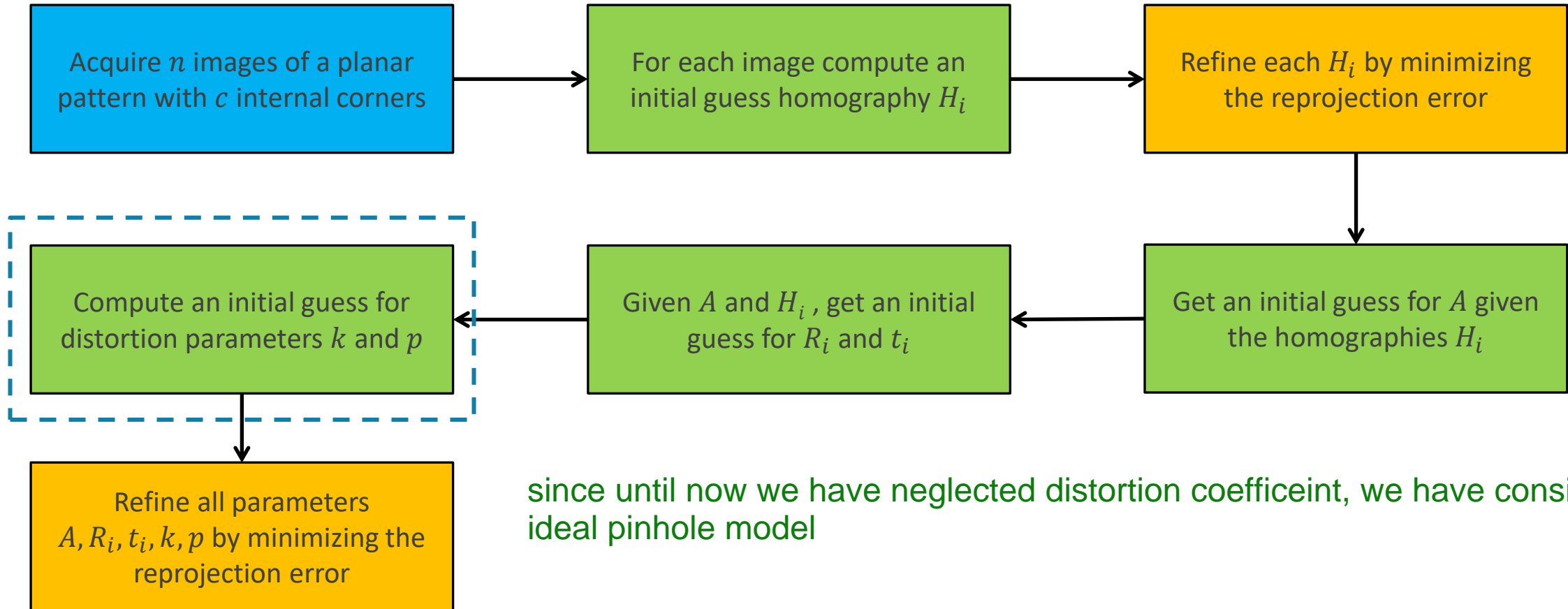
Finally, enforcing again orthonormality of  $\mathbf{R}_i$ ,  $\mathbf{r}_{i3} = \mathbf{r}_{i1} \times \mathbf{r}_{i2}$  the cross product results in a vector orthogonal to the plane where the multiplying vectors lie

Yet, the resulting matrix  $\mathbf{R}_i$  will not be exactly orthonormal since  $\mathbf{r}_{i1}$  and  $\mathbf{r}_{i2}$  are not necessarily orthogonal and  $\mathbf{r}_{i2}$  does not necessarily have unit length since  $k$  was computed for  $\mathbf{r}_{i1}$ .

However, SVD of  $\mathbf{R}_i$  allows to find the closest orthonormal matrix to it by substituting  $\mathbf{D}$  with  $\mathbf{I}$ .

since we computed the norm of  $\mathbf{r}_{i1}$  and applied to  $\mathbf{r}_{i2}$ , we are not sure it is exactly a unit vector  
therefore we computed an ALMOST orthonormal matrix -> we can however compute an orthonormal matrix through SVD, throwing away the sigma and multiplying U and V, namely substituting it with the identity I

# Zhang's Method



since until now we have neglected distortion coefficient, we have considered the ideal pinhole model



# Lens distortion coefficients

strong assumption

So far, we have neglected lens distortion and calibrated a pure pinhole model. The coordinates predicted by the homographies starting from points in the WRFs correspond to the ideal (undistorted) pixel coordinates of the chessboard corners  $m_{undist}$ . The measured coordinates of the corners in the images are the real (distorted) coordinates  $m$ .

Original Zhang's method deploys such information to estimate coefficients  $k_1, k_2$  of the radial distortion function:

$$\begin{array}{l} \text{the real coord we observe,} \\ \text{computed with harris,} \\ \text{which are "warped"} \end{array} \begin{array}{c} \left[ \begin{array}{c} x \\ y \end{array} \right] \\ \text{radial} \\ \text{function} \end{array} = L(r) \begin{array}{c} \left[ \begin{array}{c} x_{undist} \\ y_{undist} \end{array} \right] \\ \text{we take this from the homographies} \\ \text{computed} \end{array} = (1 + k_1 r^2 + k_2 r^4) \begin{array}{c} \left[ \begin{array}{c} x_{undist} \\ y_{undist} \end{array} \right] \\ \text{we take this from the homographies} \\ \text{computed} \end{array}$$

LINEAR SYSTEM

OpenCV uses a different method for estimating the distortion parameters:

- 3 coefficients for radial distortion ( $k_1, k_2, k_3$ )
- 2 coefficients for tangential distortion ( $p_1, p_2$ )

we have to approximate  $x_{undist}$ ,  $y_{undist}$ , doing a strong assumption: the homography  $H$  we have computed is a valid transformation, a good model

# Metric image coordinates

we have to move into pixel space

**Recall:** lens distortion takes place **before** we change metric image coordinates to pixel coordinates. But we measure and predict pixel coordinates.

we want metric image coordinates, not the pixel coordinates  $u, v$

We can transform back pixel coordinates  $\begin{bmatrix} u \\ v \end{bmatrix}$  to metric image coordinates  $\begin{bmatrix} x \\ y \end{bmatrix}$  thanks to the estimated intrinsic matrix  $\mathbf{A}$

$$\begin{matrix} \text{pixels} \\ \begin{bmatrix} ku \\ kv \\ k \end{bmatrix} \end{matrix} \equiv \mathbf{A} \begin{matrix} \text{metric img coords} \\ \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} \end{matrix} \Rightarrow \begin{bmatrix} ku \\ kv \\ k \end{bmatrix} \equiv \begin{bmatrix} f_u x + u_0 \\ f_v y + v_0 \\ 1 \end{bmatrix} \Rightarrow \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} \frac{u - u_0}{f_u} \\ \frac{v - v_0}{f_v} \end{bmatrix}$$

The same transformation holds between  $u_{undist}, v_{undist}$  and  $x_{undist}, y_{undist}$

Then, the distortion equation in pixel coordinates become


$$\begin{bmatrix} \frac{u - u_0}{f_u} \\ \frac{v - v_0}{f_v} \end{bmatrix} = (1 + k_1 r^2 + k_2 r^4) \begin{bmatrix} \frac{u_{undist} - u_0}{f_u} \\ \frac{v_{undist} - v_0}{f_v} \end{bmatrix} \Rightarrow \begin{bmatrix} u - u_0 \\ v - v_0 \end{bmatrix} = (1 + k_1 r^2 + k_2 r^4) \begin{bmatrix} u_{undist} - u_0 \\ v_{undist} - v_0 \end{bmatrix}$$

$u, v$  are real coords of our pixels warped by lens distortion, computed by harris

$u/v_{undist}$  are the coords computed by projecting the 3d points according to our homographies

# Lens distortion coefficients

we know everything but  $k_1$  and  $k_2$ ,  
the distortion coefficients

$$\begin{bmatrix} u - u_0 \\ v - v_0 \end{bmatrix} = (1 + k_1 r^2 + k_2 r^4) \begin{bmatrix} u_{undist} - u_0 \\ v_{undist} - v_0 \end{bmatrix} \Rightarrow \begin{bmatrix} u - u_0 \\ v - v_0 \end{bmatrix} - \begin{bmatrix} u_{undist} - u_0 \\ v_{undist} - v_0 \end{bmatrix} = (\underline{k_1} r^2 + \underline{k_2} r^4) \begin{bmatrix} u_{undist} - u_0 \\ v_{undist} - v_0 \end{bmatrix}$$


all of this is  
known,  
except for  
 $k_1, k_2$

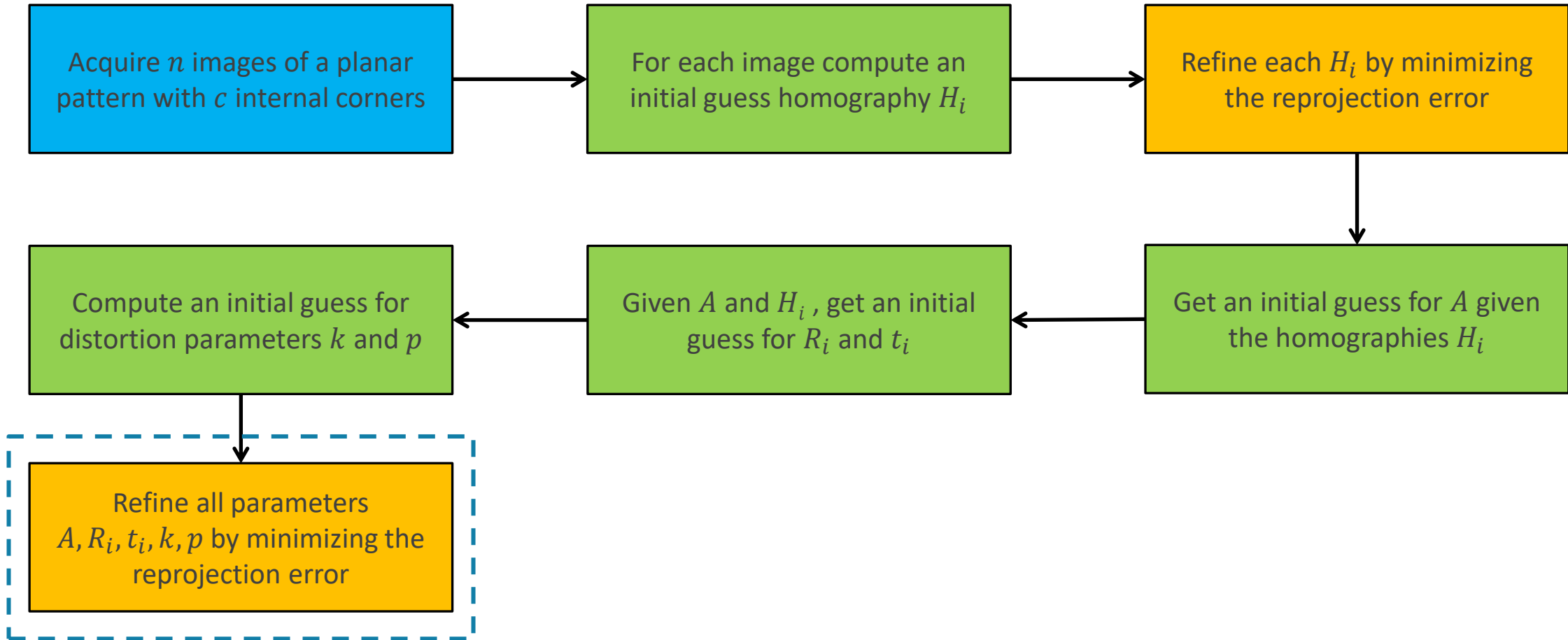
$$\begin{bmatrix} u - u_{undist} \\ v - v_{undist} \end{bmatrix} = \begin{bmatrix} (u_{undist} - u_0)r^2 & (u_{undist} - u_0)r^4 \\ (v_{undist} - v_0)r^2 & (v_{undist} - v_0)r^4 \end{bmatrix} \begin{bmatrix} k_1 \\ k_2 \end{bmatrix}$$

We get a **linear, non-homogeneous** system of linear equations  $\mathbf{D}\mathbf{k} = \mathbf{d}$  in the unknowns  $\mathbf{k} = [k_1 \ k_2]^T$ .  
With  $c$  corners in  $n$  images we get  $2nc$  equations in 2 unknowns, which can be solved in a least square  
sense, i.e., minimizing  $\|\mathbf{D}\mathbf{k} - \mathbf{d}\|_2$ , by computing the pseudo-inverse matrix  $\mathbf{D}^\dagger$  as

"D dagger"

$$\mathbf{k}^* = \min_{\mathbf{k}} \|\mathbf{D}\mathbf{k} - \mathbf{d}\|_2 = \mathbf{D}^\dagger \mathbf{d} = (\mathbf{D}^T \mathbf{D})^{-1} \mathbf{D}^T \mathbf{d}$$

# Zhang's Method



Zhengyou Zhang. A flexible new technique for camera calibration. IEEE Trans. on PAMI, 2000.

# Refinement by non-linear optimization

Final non-linear refinement of the estimated parameters. As for homographies, the procedure highlighted so far seeks to minimize an **algebraic error**, without any real physical meaning.

A more accurate solution can instead be found by a so called Maximum Likelihood Estimate (**MLE**) aimed at minimization of the geometric (i.e. reprojection) error.

We use all the values estimated so far as **initial guesses**.

Under the hypothesis of i.i.d. (independent identically distributed) noise, the MLE for our models is obtained by minimization of the error

$$A^*, k^*, R_i^*, t_i^* = \underset{A, k, R_i, t_i}{\operatorname{argmin}} \sum_{i=1}^n \sum_{j=1}^c \|\tilde{m}_{ij} - \hat{m}(A, k, R_i, t_i, \tilde{w}_j)\|^2$$

real ones  
computed  
by Harris

this time we will optimize not for H, but for those parameters

with respect to all the unknown camera parameters, which can be solved again by using an **iterative algorithm**, like the Levenberg-Marquardt algorithm.



get rid of

# Compensate lens distortion

After calibration, we have a precise mathematical model that maps points in the 3D world to points in the image plane. Yet, lens distortion makes the system non-linear and therefore cumbersome to use.

Hence, it is common, once a camera has been calibrated, to **warp** the images it takes so to simulate a camera without lens distortion, whose camera formation model is linear, i.e. it is the estimated PPM.



undistort



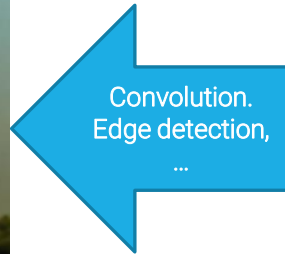
<https://docs.nvidia.com/isaac/archive/2020.2/packages/perception/doc/warp.html>

# Image **Warping** $\neq$ Image **filtering**

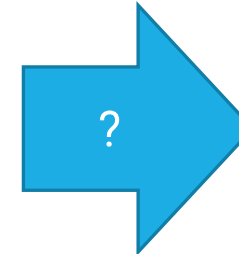
**Warping** refers to transformations of the **spatial domain** of images, while **filtering** refers to changes in the **RGB values** of images **more related to intensity, the content**



**Filtered** (Gaussian blur)



Original

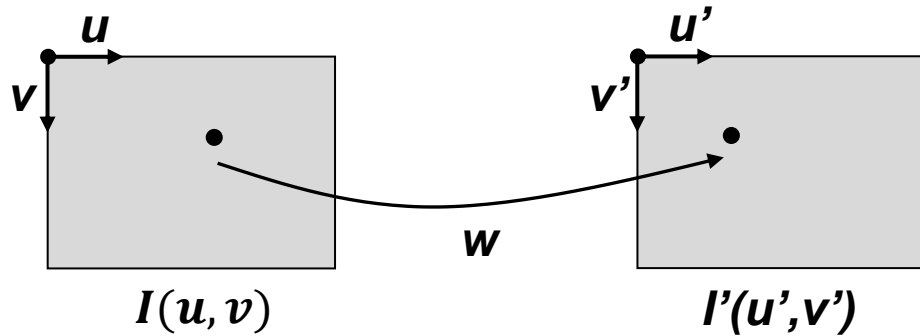


**Warped** (Rotation & scale)

# Image Warping

it is a relationship, a mapping, between coords  
create a new image  $I'$

If we have a function that computes point in image  $I'$  starting from point in image  $I$ , we can copy the value



warping function

$$\begin{cases} u' = w_u(u, v) \\ v' = w_v(u, v) \end{cases}$$

$$I'(w_u(u, v), w_v(u, v)) = I(u, v)$$

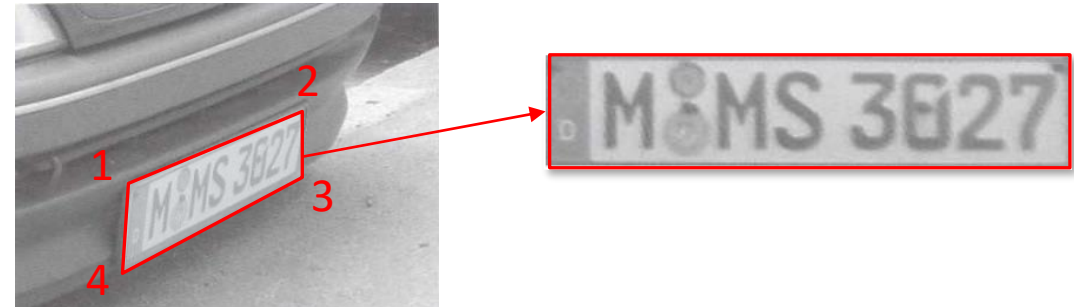
Can be just a rotation

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

transf. coord

Or it can be a full homography

$$k \begin{bmatrix} u' \\ v' \\ 1 \end{bmatrix} = \begin{bmatrix} h_{11} & h_{12} & h_{13} \\ h_{21} & h_{22} & h_{23} \\ h_{31} & h_{32} & h_{32} \end{bmatrix} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix}$$



# Forward mapping

going from  $I$  to  $I'$  is called Forward mapping

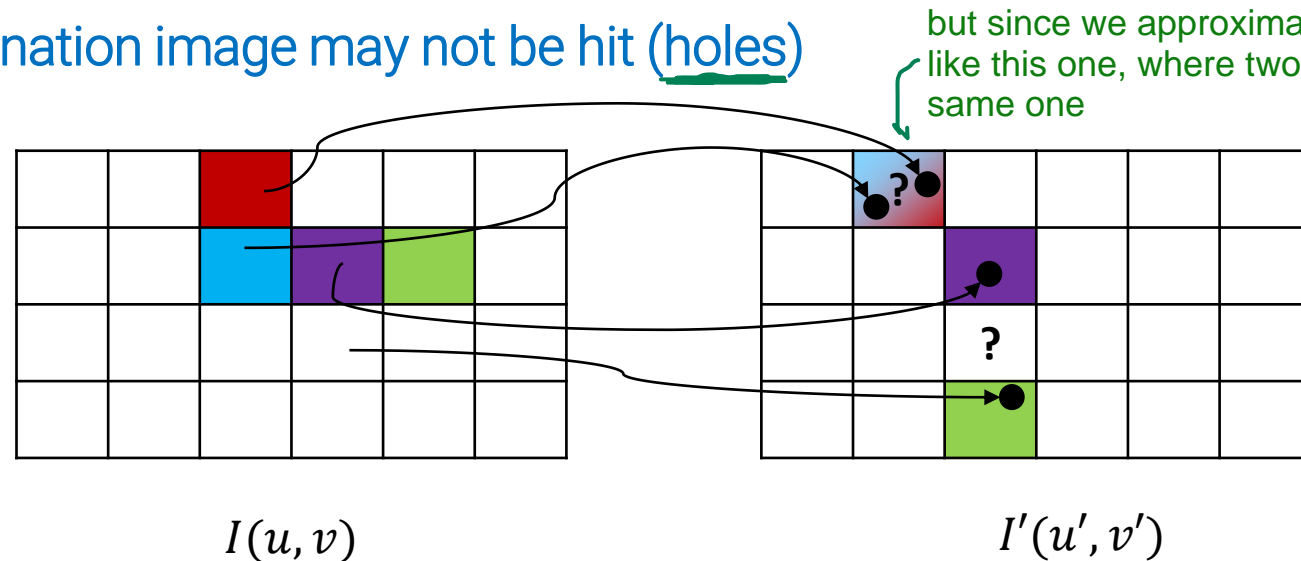
If we start from the input image coordinates, after applying the warping function in general we get **continuous coordinates in the output image**, not discrete ones. This is called a **forward mapping**. we want discrete, not real continuous numbers -> therefore we do roundings

Possible choices to make coordinates integer numbers: truncate, **nearest neighbor** (i.e. rounding), ... Regardless of the discretization function, due to rounding

- more than one pixel can go to one position (folds)
- some pixels of the destination image may not be hit (holes)

two main problems  
of forw. mapping

$$\begin{cases} u' = w_u(u, v) \\ v' = w_v(u, v) \end{cases}$$



we want an injective  
warping function



# Backward mapping

we start from  $I'$  the target image using inverse mapping  $w^{-1}$

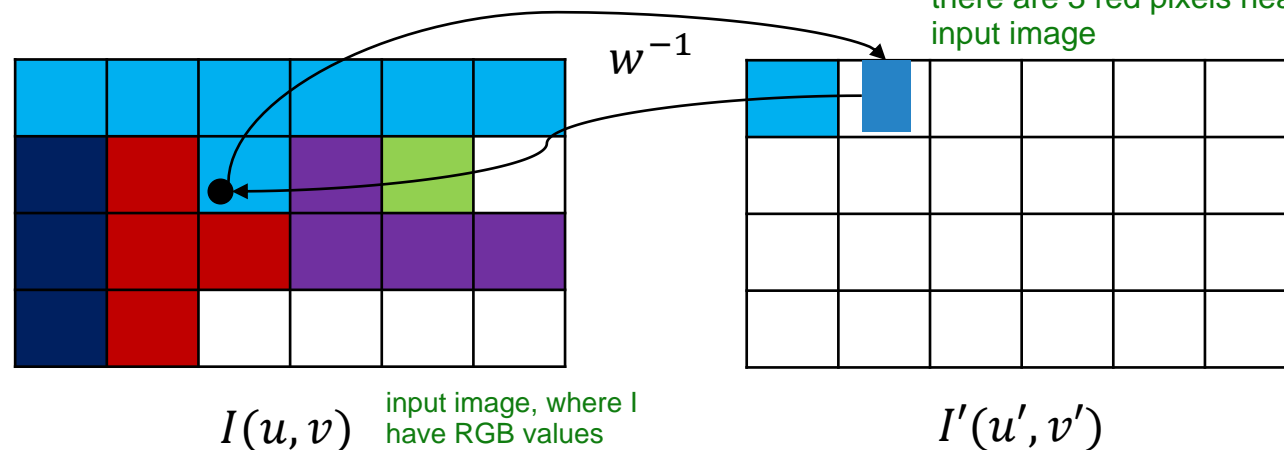
doing this, we start for sure from an integer number, cause we start from a certain pixel and go back to the source image  $I$ , and then we'll have there a non-integer value, but it is not a problem because even if we discretize it with different techniques we'll always end up with a value for the pixel in the target image

We can avoid these problems if we compute input coordinates corresponding to **each pair of integer coordinates in the output image**, by using the **inverse mapping**  $w^{-1}$ . This strategy is referred to as **backward mapping**.

Yet, the input coordinates are still continuous values. Which discretization strategy are used?

- Truncate
- Nearest Neighbour just choose the closest one
- Interpolate between the 4 closest point (bilinear, bicubic, etc...)

$$\begin{cases} u = w_u^{-1}(u', v') \\ v = w_v^{-1}(u', v') \end{cases}$$



we compute a linear interpolation of intensities  $I_a$  and  $I_b$ , but I have to compute them doing another linear interpolation between  $I_1$  and  $I_2$

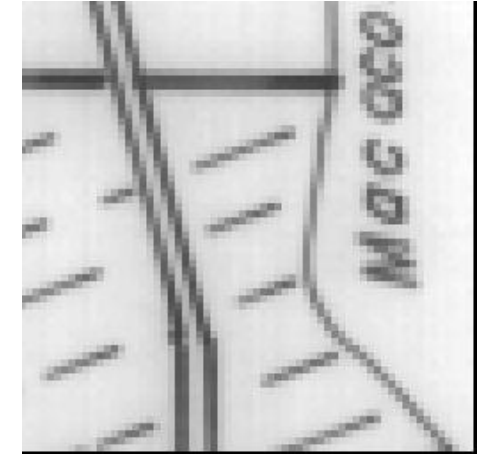
in order to get lin. inter. of  $I_a$  and  $I_b$  we compute first the horizontal interp. between  $I_1$  and  $I_2$  and  $I_3$  and  $I_4$

# Bilinear Interpolation



Input

Warp (Zoom) by nearest neighbor



sharper edges

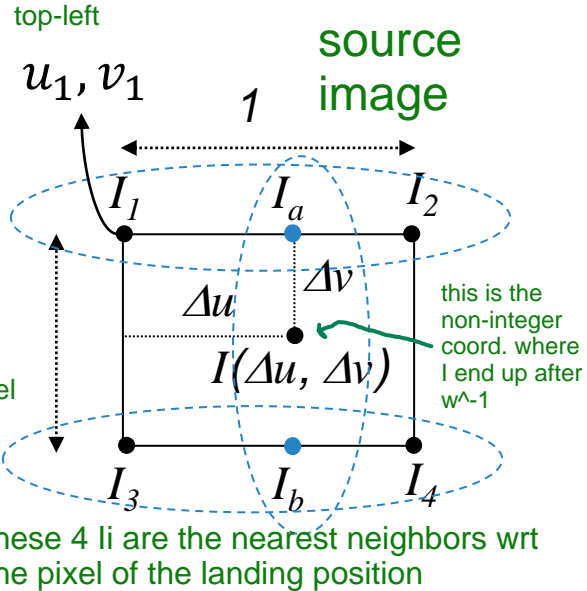
use this in order to preserve e.g. binary masks/classes/categories/orical variables, etc.

Warp (Zoom) by Bilinear Interpolation

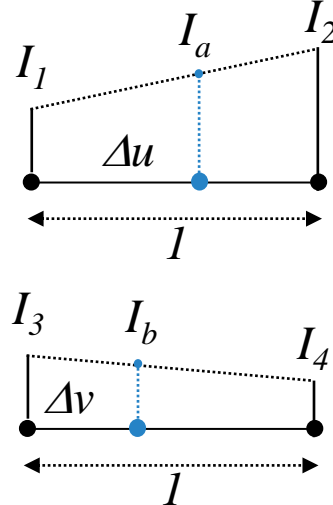


smoother image

use e.g. for natural images



the height is the intensity value: in this case  $I_2$  is brighter than  $I_1$



$$\Delta u = u - u_1$$

$$\Delta v = v - v_1$$

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

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

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

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

The closer a point to a pixel the smaller the other weights become

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

$$I'(u', v') = (1 - \Delta u)(1 - \Delta v)I_1 + \Delta u(1 - \Delta v)I_2 + (1 - \Delta u)\Delta vI_3 + \Delta u\Delta vI_4$$

it is a linear combination of the "corner" points: depending on the position  $I$ , e.g. if the point is close to  $I_4$ , the resulting point will be really similar to  $I_4$  and the other points  $I_1, I_2, I_3$ , will be negligible



# «Undistort» warping

$u, v$  is the point where the point should be landed without the distortion caused by lens, which bend the image

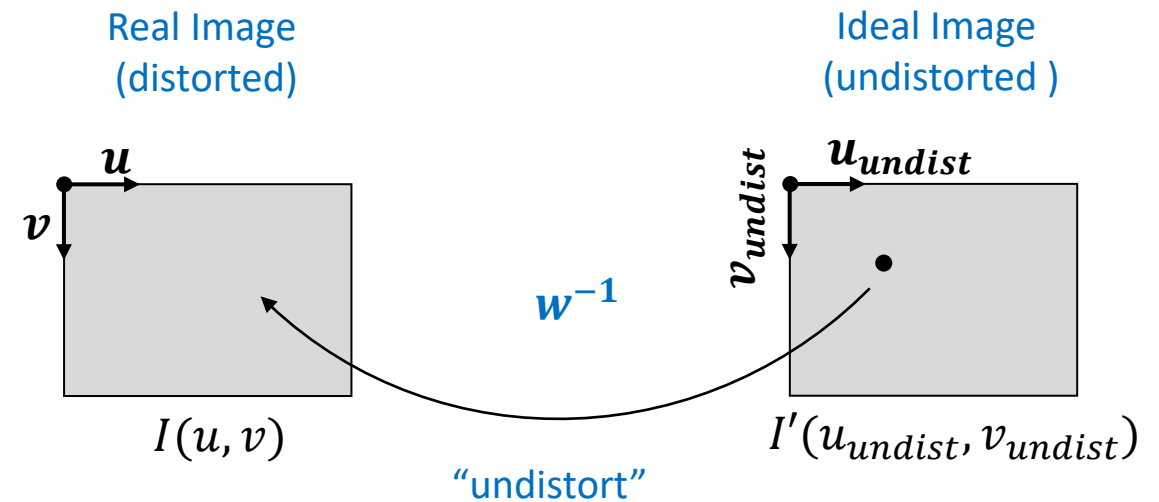
Once the lens distortion parameters have been computed by camera calibration, the image can be corrected by a backward warp from the undistorted to the distorted image based on the adopted lens distortion model. **For this images, the image formation model is linear, i.e. the PPM.**

$$\forall(u_{undist}, v_{undist}): I'(u_{undist}, v_{undist}) = I(w_u^{-1}(u_{undist}, v_{undist}), w_v^{-1}(u_{undist}, v_{undist}))$$

$w^{-1}$  is given by the zhang's radial distortion function that we use in camera calibration to find the real coordinates of pixels exposed to distortion effects

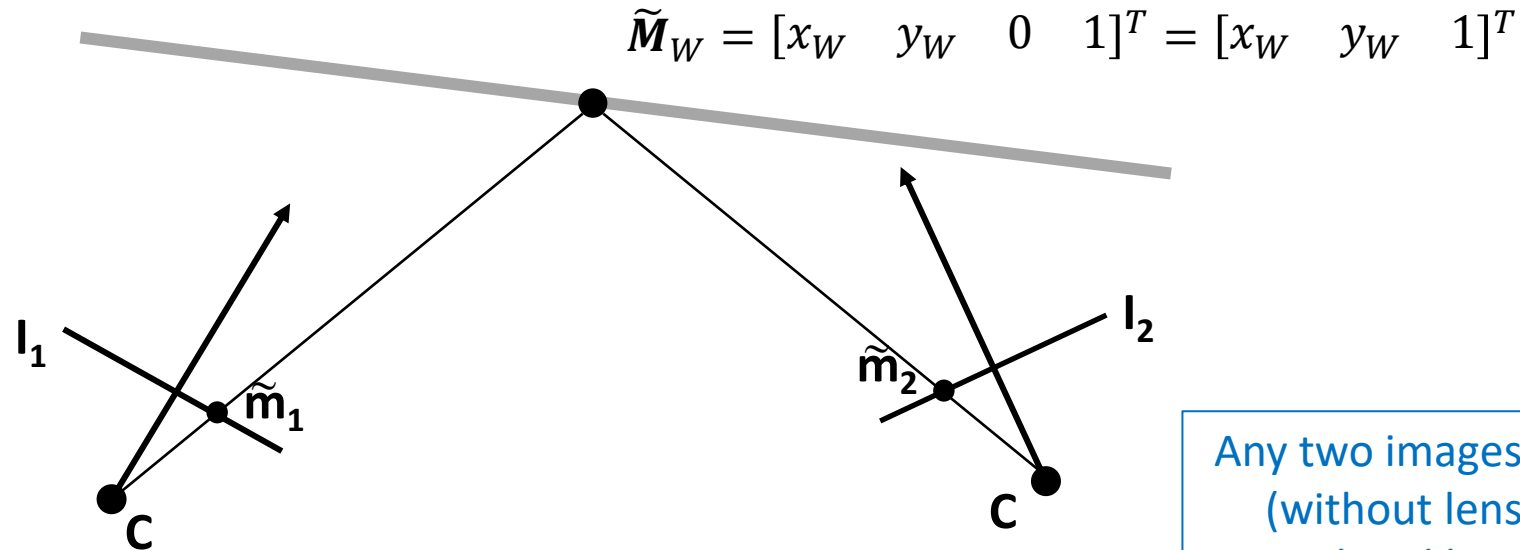
$$w = \begin{cases} u = u_{undist} + (k_1 r^2 + k_2 r^4)(u_{undist} - u_0) \\ v = v_{undist} + (k_1 r^2 + k_2 r^4)(v_{undist} - v_0) \end{cases}$$

Zhang's  
Radial distortion



# Warping with homographies

when we project  $M$  in WRF to camera 1 we obtain  $m_1$



Any two images of a planar scene (without lens distortion) are related by a homography

the homography  $H$  maps  $M_W$  to  $m$

$$\tilde{m}_1 = H_1 \tilde{M}_W$$

$$\tilde{m}_2 = H_2 \tilde{M}_W$$

we are talking about linear camera model now, without distortion

which are the  $m_1$  coord in the other pov of  $m_2$ ?

$$\tilde{m}_1 = H_1 H_2^{-1} \tilde{m}_2$$

i can get a relationship between pixel coords, without considering WRF

$$\tilde{m}_2 = H_2 H_1^{-1} \tilde{m}_1$$

$$H_{21}$$

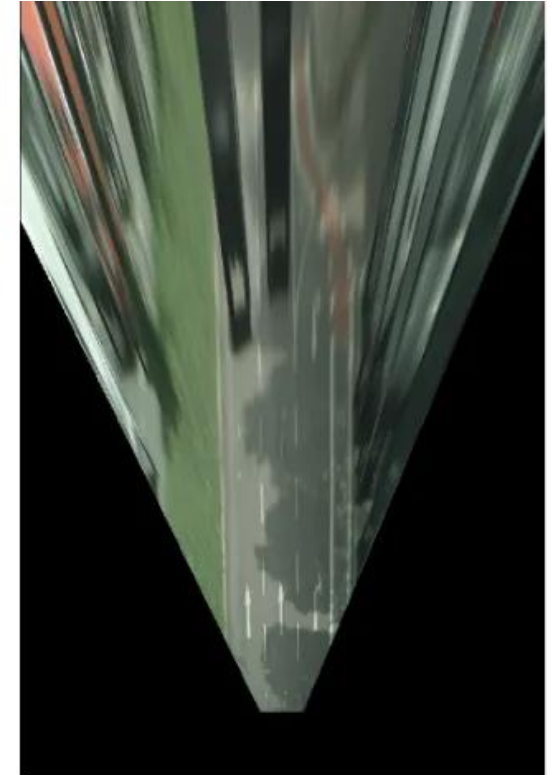
Homography

$$H_{12} = H_{21}^{-1}$$

this is the homography which goes from image plane 1 to image plane 2

# Example: Inverse Perspective Mapping

in this case the second image plane is a virtual sensor, it's like we simulate to have another image e.g. taking interesting points on the street



<https://towardsdatascience.com/a-hands-on-application-of-homography-ipm-18d9e47c152f>

another useful use case: rotate camera around its optical center

I can compute images as we have one in a different orientation -> really powerful operation

# Warping with homographies

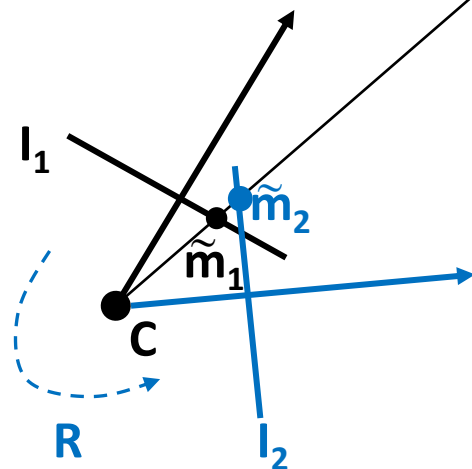
in this case the relationship between images are homographies and these homographies are valid for each  $M_w$ , for each 3D point in the WRF

I can compute images as if I had a camera mounted in different orientation

same camera, just rotated around the center

generic world point

$$\tilde{M}_W = \begin{bmatrix} x_W \\ y_W \\ z_W \\ 1 \end{bmatrix}$$



the relation between  $m_1$  and  $m_2$  is a rotation

I can pass from image plane 1 to image plane 2 computing a rotation matrix

Any two images taken by a camera rotating about its optical center are related by a homography (if lens distortion has been removed)

we get rid of the homog. part  $[I | 0]$

what if I want to sim. a virtual sensor and project  $M_w$  onto  $I_2$ ? I have to apply a rotation to match the CRF from WRF

$$\tilde{m}_1 = A [I | 0] \tilde{M}_W = A M_W \quad \tilde{m}_2 = A [R | 0] \tilde{M}_W = A R M_W$$

in this case there is no rotation or transl, no extrinsic params, the WRF and CRF coincides

$$\tilde{m}_1 = A R^{-1} A^{-1} \tilde{m}_2$$

$$\tilde{m}_2 = A R A^{-1} \tilde{m}_1$$

I can express pixel points as other pixel points, as before

$$H_{21} \xleftarrow{\text{Homography}} H_{12} = H_{21}^{-1}$$

pitch is the rotation around the horizontal x axis

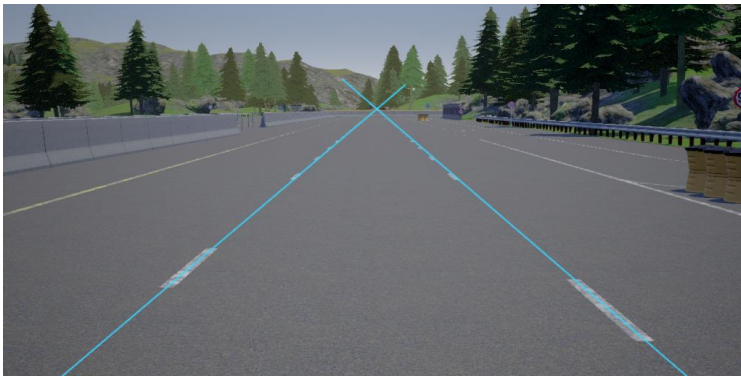
# Example: compensate pitch or yaw

Autonomous driving cameras should be ideally mounted with the optical axis parallel to the road plane and aligned with the direction of motion. It is however difficult to obtain perfect alignment by mechanical mounting only.

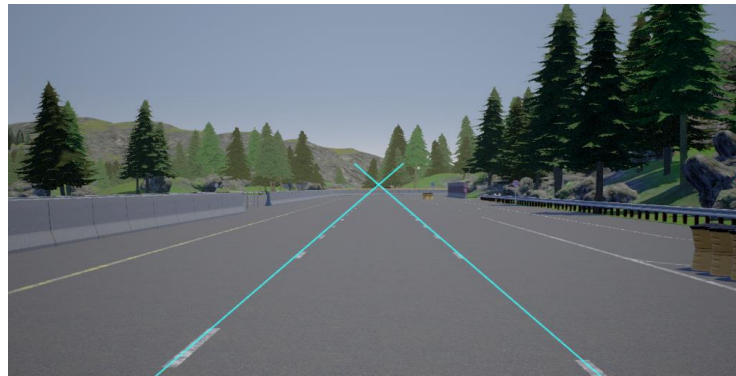
With a calibrated camera, however, it is possible to **compensate pitch** (rotation around the x axis) and **yaw** (rotation around the y axis) **by estimating the vanishing point of the lane lines when the vehicle is travelling straight in a lane.**

we can compute the straight lines with the Hough Transform and their intersection, namely the vanishing point, obtaining  $m_{inf}$  and its coordinates  $u,v$

*Pitch =  $-5^\circ$*



*Pitch =  $0^\circ$*



the vanishing point is perfectly at the center of the image

*Pitch =  $+5^\circ$*



<https://thomasfermi.github.io/Algorithms-for-Automated-Driving/CameraCalibration/VanishingPointCameraCalibration.html>

# Example: compensate pitch or yaw

When the vehicle is driving straight with respect to the lines, their orientation in a world reference frame attached to the vehicle is parallel to the z axis, i.e. it is  $[0 \ 0 \ 1]^T$  and their point at infinity in  $\mathbb{P}^3$  is  $[0 \ 0 \ 1 \ 0]^T$ .

The vanishing point is then  $m_\infty \equiv A[R|0] \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix} \equiv Ar_3 \equiv A \begin{bmatrix} 0 \\ \sin \beta \\ \cos \beta \end{bmatrix}$  since a rotation around the  $x$  axis of an angle  $\beta$  has expression  $R_{pitch} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \cos \beta & \sin \beta \\ 0 & -\sin \beta & \cos \beta \end{bmatrix}$  we want to estimate the pitch wrt the ideal image with 0 pitch

remember: is the projection of the point at infinity  
our unknown

Then, if we estimate the coordinates of the vanishing point in the image, we can compute  $r_3 = \frac{A^{-1}m_\infty}{\|A^{-1}m_\infty\|_2}$  and from it  $R_{pitch}$  and finally the warping homography from 0 pitch to input image as  $AR_{pitch}A^{-1}$

normalization

With the same procedure it is possible to correct simultaneously yaw and pitch.