# Estimating camera pose from a single image and a known map
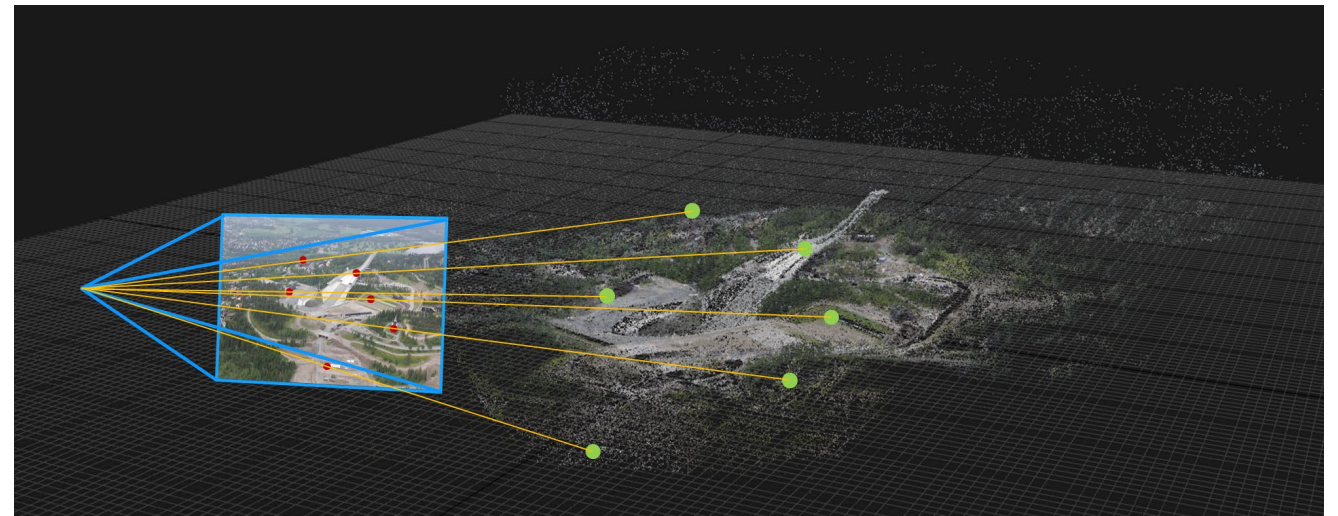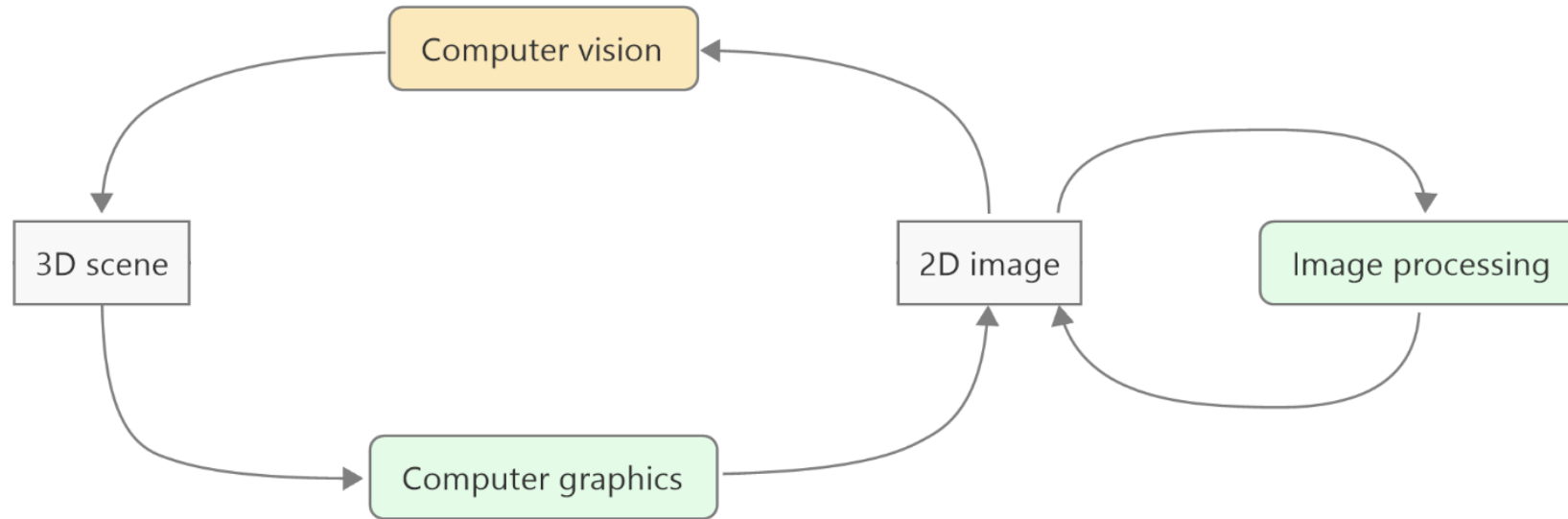
Trym Vegard Haavardsholm
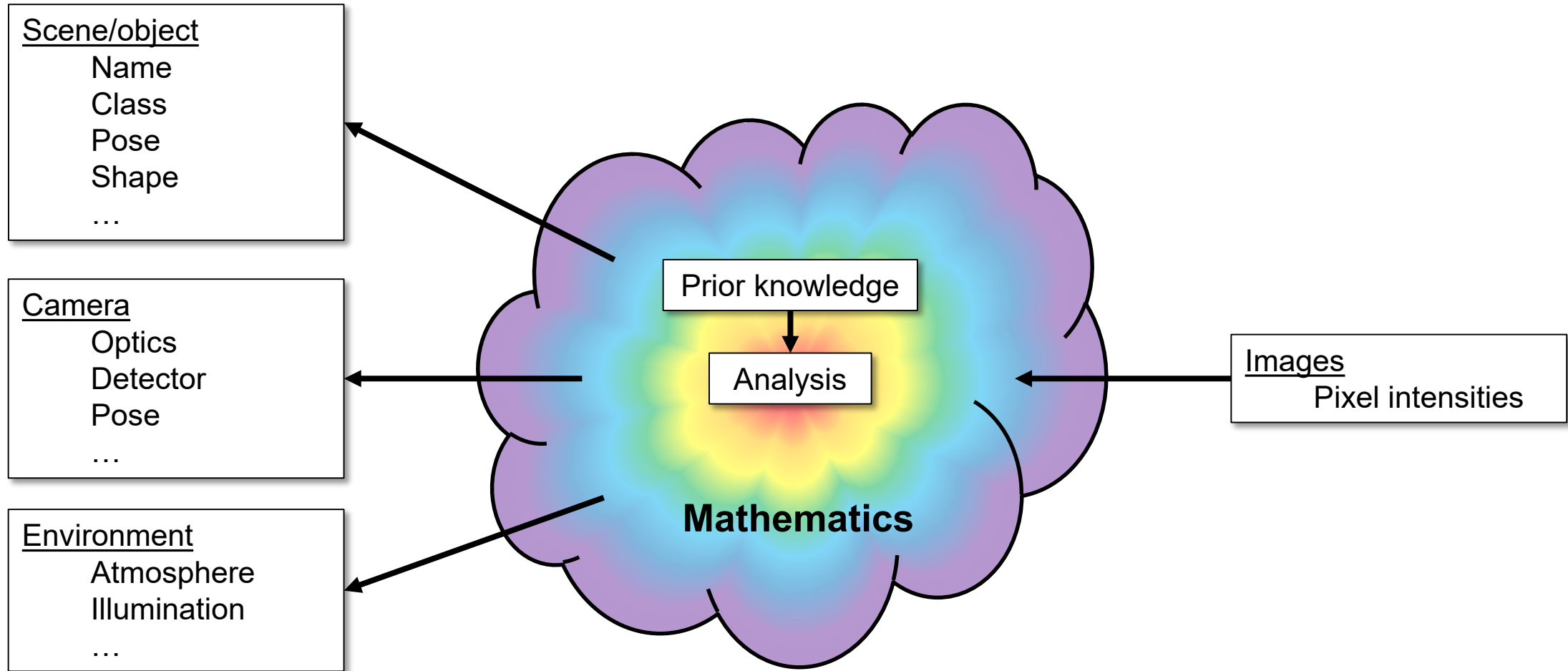
2023



**TEK**5030

# Computer vision is an *inverse problem*!

**TEK**5030

# The *inverse* analysis process



Scene/object
Name
Class
Pose
Shape
…

Camera
Optics
Detector
Pose
…

Environment
Atmosphere
Illumination
…

Prior knowledge

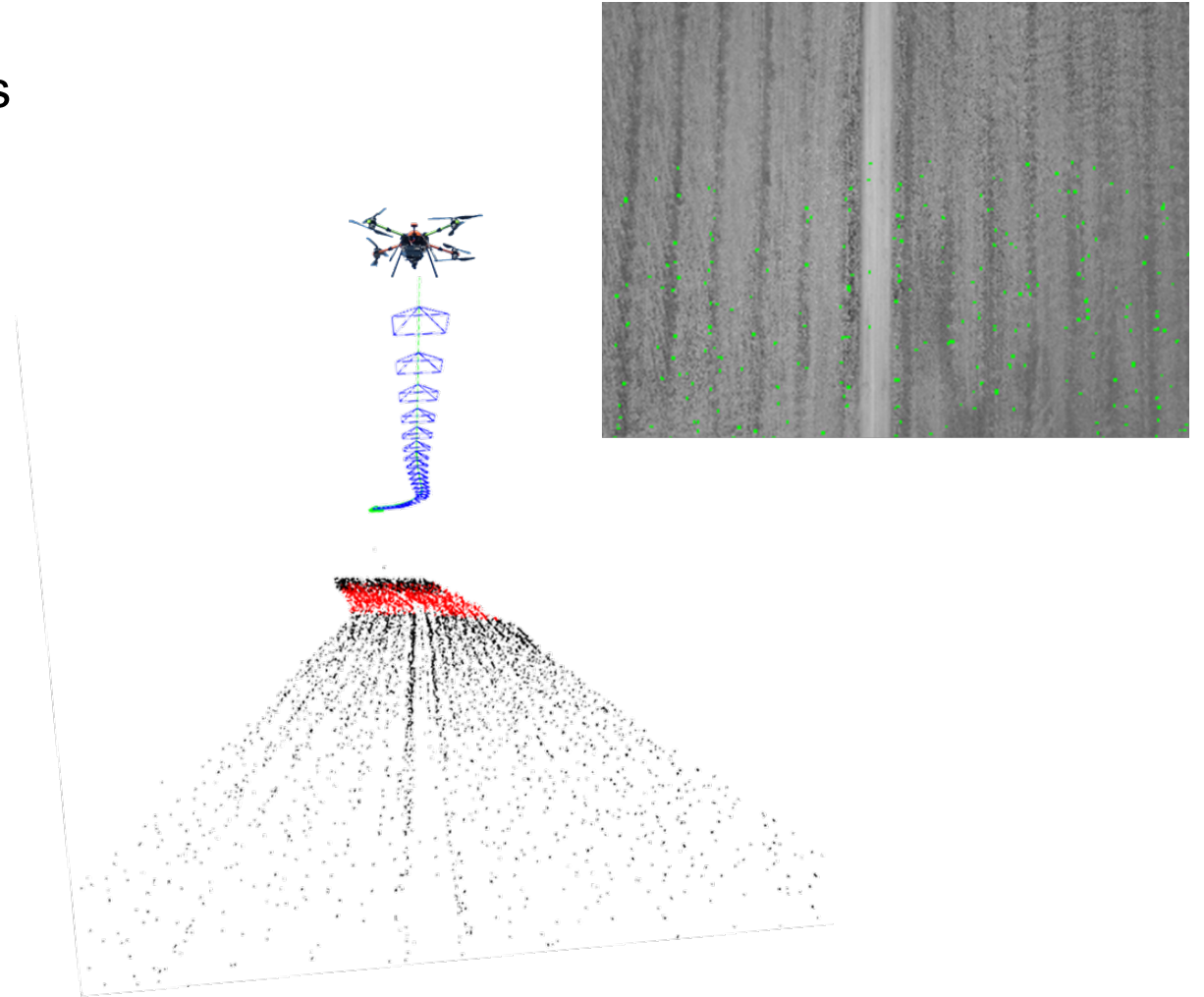Analysis

**Mathematics**

Images
Pixel intensities

# Localisation
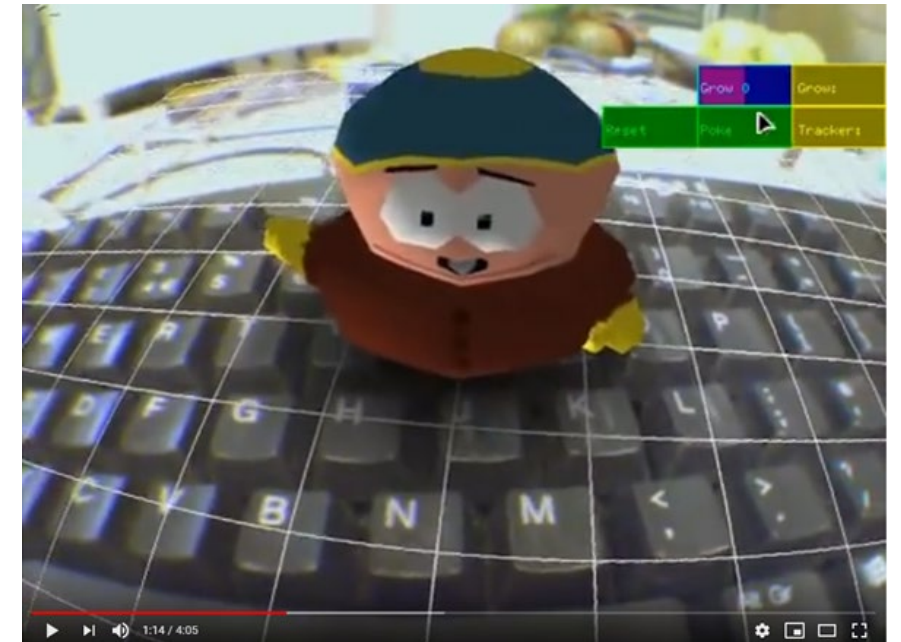
Pose estimation based on correspondences with a known map is called **localisation**

In **visual localisation**,
this is also sometimes called **tracking**
- Tracking the map in the image frames

**TEK5030**

# Why learn about localisation?





From PTAM by Georg Klein and David Murray (2007)
https://www.youtube.com/watch?v=F3s3M0mokNc

**TEK5030**

# How can we track a map with a camera?

# Pose from 2D correspondences with known 3D points

# Pose from 2D correspondences with known 3D points

**TEK**5030

# Pose from 2D correspondences with known 3D points

Minimise **geometric error**

$$\mathbf{T}^*_{wc} = \underset{\mathbf{T}_{wc}}{\operatorname{argmin}} \sum_i \left\| \pi(\mathbf{T}^{-1}_{wc} \cdot \mathbf{x}^w_i) - \mathbf{u}_i \right\|^2$$

# Pose from 2D correspondences with known 3D points

Minimise **geometric error**

$$\mathbf{T}_{wc}^{*} = \operatorname*{argmin}_{\mathbf{T}_{wc}} \sum_{i} \left\| \pi(\mathbf{T}_{wc}^{-1} \cdot \mathbf{x}_i^w) - \mathbf{u}_i \right\|^2$$

also called **reprojection error**

# Pose estimation

We will solve the indirect tracking problem

$$\mathbf{T}_{wc}^* = \underset{\mathbf{T}_{wc}}{\operatorname{argmin}} \sum_i \left\| \pi(\mathbf{T}_{wc}^{-1} \cdot \mathbf{x}_i^w) - \mathbf{u}_i \right\|^2$$

in the next few videos.

But lets first solve a simpler problem,
when we can assume that the world is planar!

# Pose estimation relative to a world plane

Choose the world coordinate system
so that the $xy$-plane corresponds to
a plane $\Pi$ in the scene

$$\mathbf{x}_\Pi^w = \begin{bmatrix} x \\ y \\ 0 \end{bmatrix} \qquad \mathbf{x}^\Pi = \begin{bmatrix} x \\ y \end{bmatrix}$$

**TEK5030**

# Pose estimation relative to a world plane

We can map points on the world plane into image coordinates by using the perspective camera model

$$\tilde{\mathbf{u}} = \mathbf{K}\begin{bmatrix}\mathbf{R}\ \mathbf{t}\end{bmatrix}\tilde{\mathbf{x}}_{\Pi}^{w}$$

$$\mathbf{T}_{cw} = \begin{bmatrix}\mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1\end{bmatrix}$$

**TEK**5030

# Pose estimation relative to a world plane

We can map points on the world plane
into image coordinates by using
the perspective camera model

$$\tilde{\mathbf{u}} = \mathbf{K}\begin{bmatrix} \mathbf{R} & \mathbf{t} \end{bmatrix}\tilde{\mathbf{x}}_{\Pi}^{w}$$

$$\mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

$$= \mathbf{K}\begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3, \mathbf{t} \end{bmatrix}\begin{bmatrix} x \\ y \\ 0 \\ 1 \end{bmatrix}$$

**TEK5030**

# Pose estimation relative to a world plane

We can map points on the world plane into image coordinates by using the perspective camera model

$$\tilde{\mathbf{u}} = \mathbf{K} \begin{bmatrix} \mathbf{R} \ \mathbf{t} \end{bmatrix} \tilde{\mathbf{x}}_\Pi^w$$

$$= \mathbf{K} \begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3, \mathbf{t} \end{bmatrix} \begin{bmatrix} x \\ y \\ 0 \\ 1 \end{bmatrix}$$

$$= \mathbf{K} \begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{t} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$\mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$



$$\tilde{\mathbf{u}} = \mathbf{K}[\mathbf{R} \ \ \mathbf{t}]\tilde{\mathbf{x}}_\Pi^w$$

$\mathcal{F}_i$

$\mathbf{u}$

$\mathbf{K}$

$\mathcal{F}_c$

$\mathbf{x}_n$

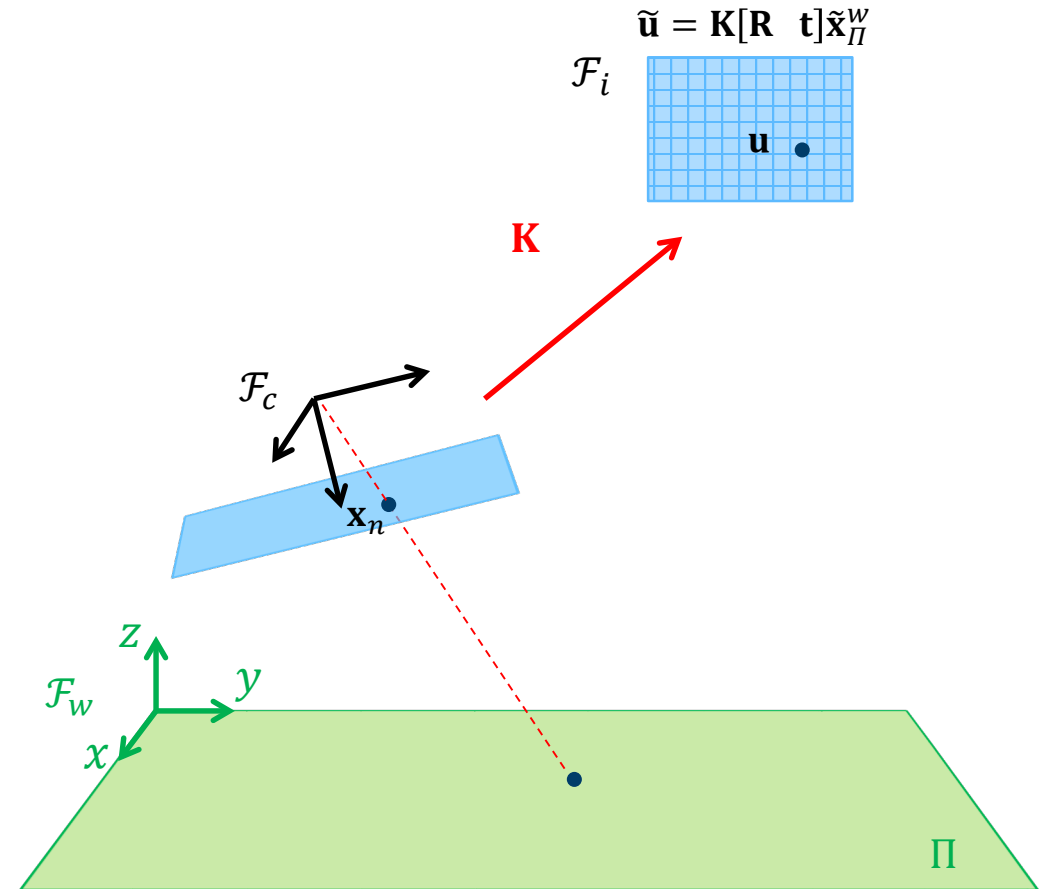$z$

$y$

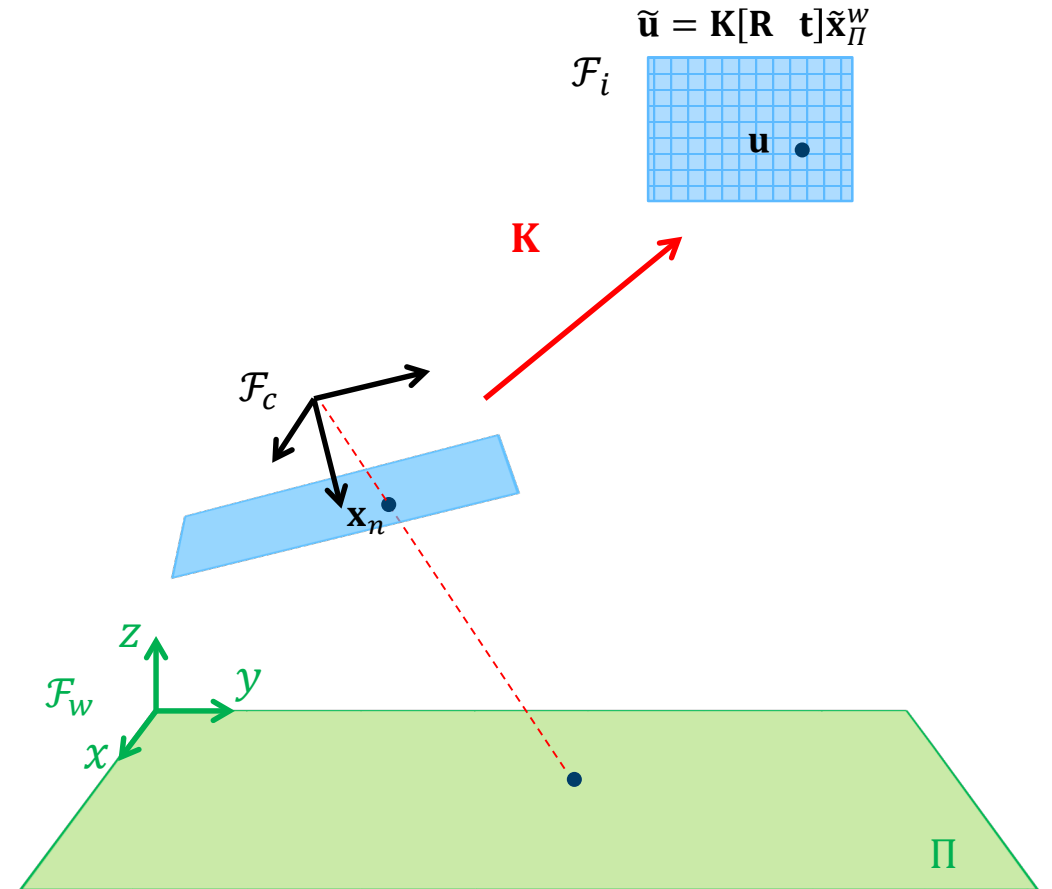$\mathcal{F}_w$
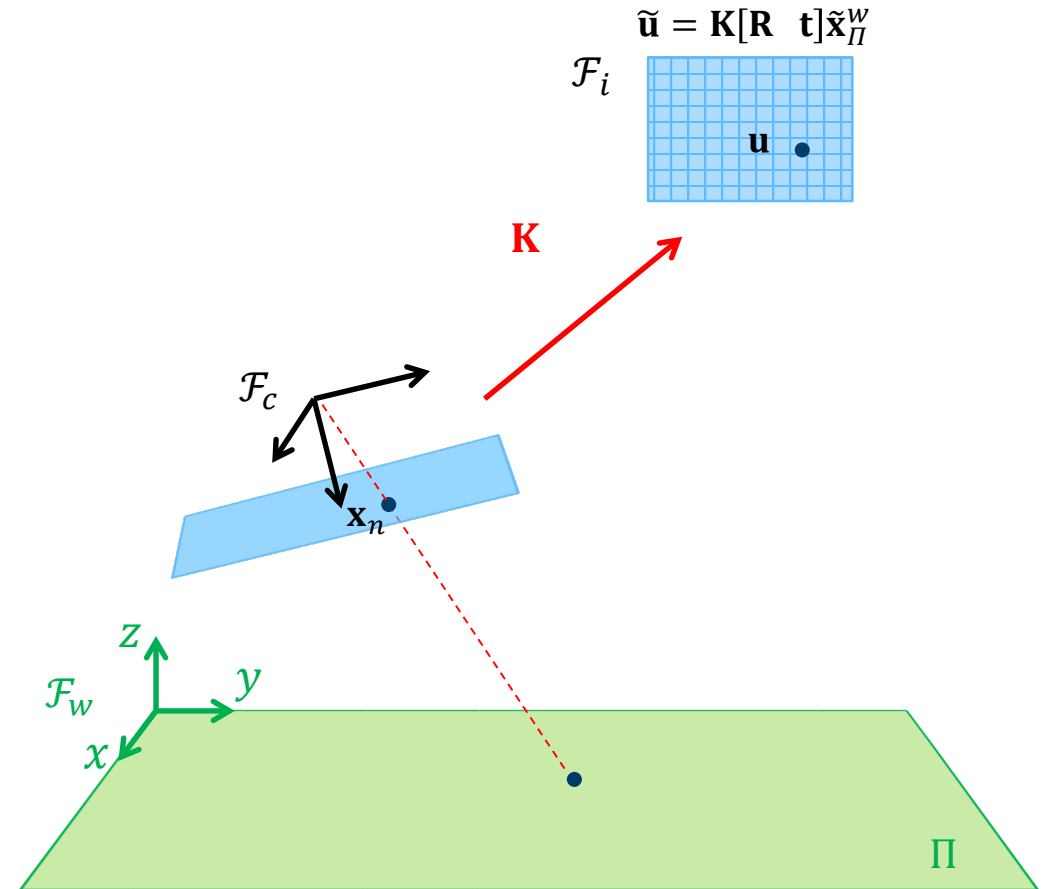
$x$

$\Pi$

**TEK5030**

# Pose estimation relative to a world plane

We can map points on the world plane into image coordinates by using the perspective camera model

$$\tilde{\mathbf{u}} = \mathbf{K} \begin{bmatrix} \mathbf{R}\ \mathbf{t} \end{bmatrix} \tilde{\mathbf{x}}_{\Pi}^{w}$$

$$\mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

$$= \mathbf{K} \begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3, \mathbf{t} \end{bmatrix} \begin{bmatrix} x \\ y \\ 0 \\ 1 \end{bmatrix}$$

$$= \mathbf{K} \begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{t} \end{bmatrix} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix}$$

$$= \mathbf{H}_{i\Pi} \tilde{\mathbf{x}}^{\Pi}$$



$$\tilde{\mathbf{u}} = \mathbf{K}[\mathbf{R}\ \ \mathbf{t}]\tilde{\mathbf{x}}_{\Pi}^{w}$$

$\mathcal{F}_i$

$\mathbf{u}$

$\mathbf{K}$

$\mathbf{H}_{i\Pi}$

$\mathcal{F}_c$

$\mathbf{x}_n$

$z$

$y$

$\mathcal{F}_w$

$x$

$\Pi$

# Pose estimation relative to a world plane

⇒ For a calibrated camera,
we have a relationship between the camera pose
and the homography between the world plane and the image!

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right] \qquad \mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

# Pose estimation relative to a world plane

$\Rightarrow$ For a calibrated camera,
we have a relationship between the camera pose
and the homography between the world plane and the image!

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right] \qquad \mathbf{T}_{cw} = \begin{bmatrix} \mathbf{R} & \mathbf{t} \\ \mathbf{0} & 1 \end{bmatrix}$$

How can we use this to
estimate camera pose
given a homography?



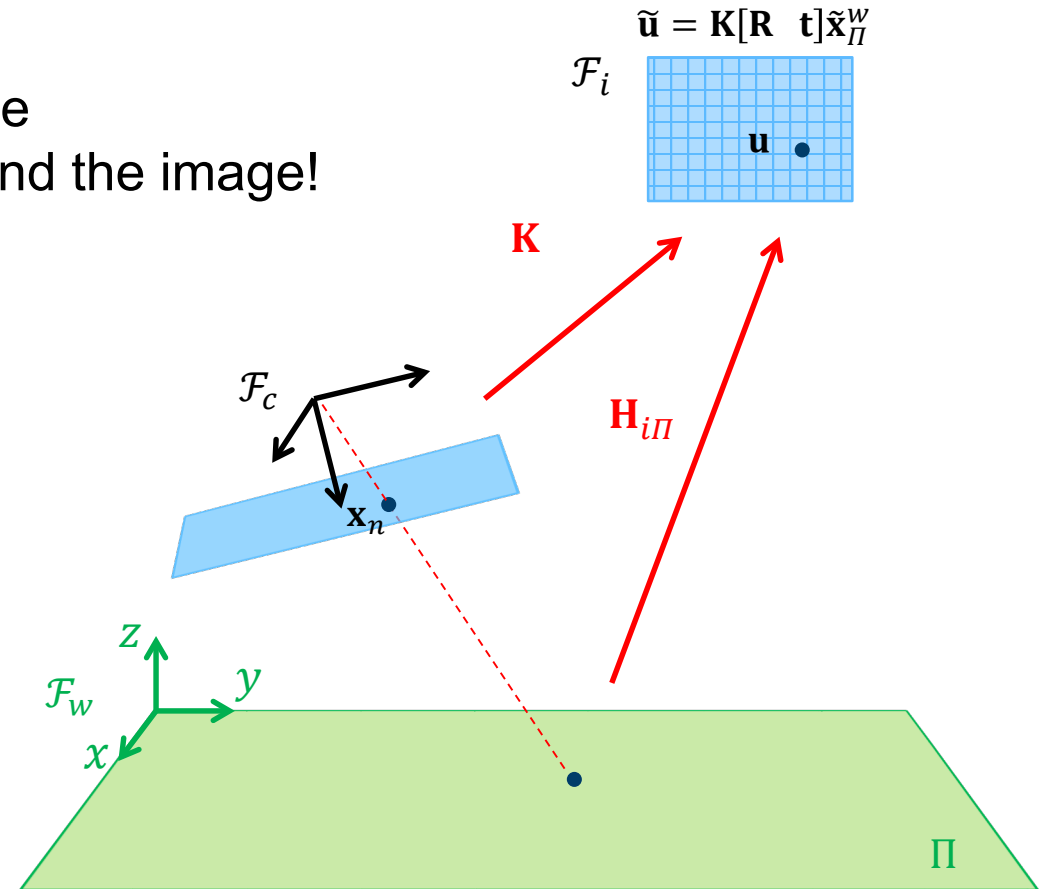$$\tilde{\mathbf{u}} = \mathbf{K}[\mathbf{R} \quad \mathbf{t}]\tilde{\mathbf{x}}_{\Pi}^{W}$$

$\mathcal{F}_i$

$\mathbf{u}$

$\mathbf{K}$

$\mathbf{H}_{i\Pi}$

$\mathcal{F}_c$

$\mathbf{x}_n$

$z$

$y$

$\mathcal{F}_w$

$x$

$\Pi$

# Pose estimation relative to a world plane

Assume a perfect, noise-free homography between the world plane and the image:

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right]$$

# Pose estimation relative to a world plane

Assume a perfect, noise-free homography between the world plane and the image:

$$\mathbf{H}_{i\Pi} = \mathbf{K}\begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{t} \end{bmatrix}$$

Then, because of scale ambiguity:

$$\begin{bmatrix} \mathbf{r}_1, \mathbf{r}_2, \mathbf{t} \end{bmatrix} \sim \mathbf{K}^{-1}\mathbf{H}_{i\Pi} = \mathbf{M}$$

**TEK5030**

# Pose estimation relative to a world plane

Assume a perfect, noise-free homography between the world plane and the image:

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right]$$

Then, because of scale ambiguity:

$$\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right] \sim \mathbf{K}^{-1}\mathbf{H}_{i\Pi} = \mathbf{M}$$

Since the columns of rotation matrices have unit norm,
we find a scale factor $\lambda$ so that the first two columns of $\mathbf{M}$ also get unit norm.
We then have the two possible solutions:

$$\left[\hat{\mathbf{r}}_1, \hat{\mathbf{r}}_2, \hat{\mathbf{t}}\right] = \pm\lambda\mathbf{M}$$

**TEK**5030

# Pose estimation relative to a world plane

Assume a perfect, noise-free homography between the world plane and the image:

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right]$$

Then, because of scale ambiguity:

$$\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right] \sim \mathbf{K}^{-1}\mathbf{H}_{i\Pi} = \mathbf{M}$$

Since the columns of rotation matrices have unit norm,
we find a scale factor $\lambda$ so that the first two columns of $\mathbf{M}$ also get unit norm.
We then have the two possible solutions:

$$\left[\hat{\mathbf{r}}_1, \hat{\mathbf{r}}_2, \hat{\mathbf{t}}\right] = \pm\lambda\mathbf{M}$$

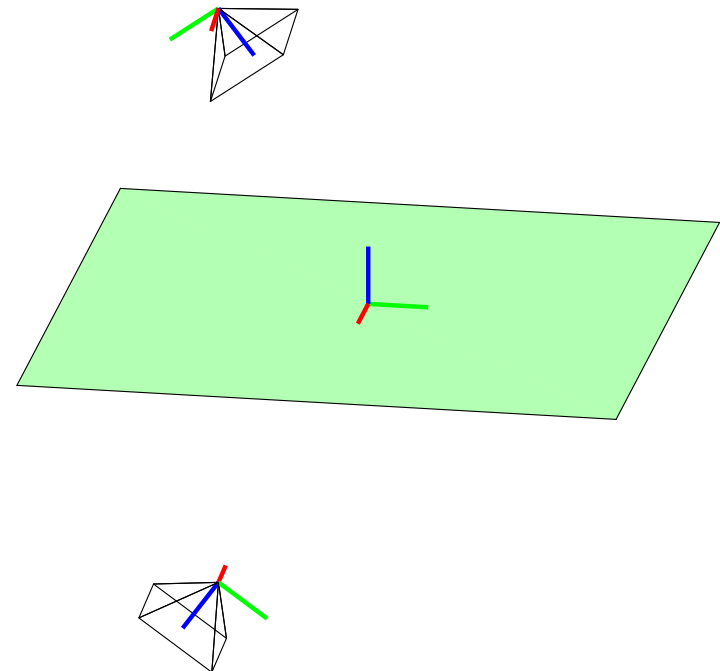The last column in $\widehat{\mathbf{R}}$ is given by the cross product of the two first columns:

$$\hat{\mathbf{r}}_3 = \pm\left(\hat{\mathbf{r}}_1 \times \hat{\mathbf{r}}_2\right), \text{ where the sign is chosen so that } \det\left(\hat{\mathbf{R}}\right) = 1$$

**TEK5030**

# Pose estimation relative to a world plane

We are now able to reconstruct the camera pose in the world coordinate system for each of the two solutions:

$$\hat{\mathbf{T}}_{wc} = \hat{\mathbf{T}}_{cw}^{-1} = \begin{bmatrix} \hat{\mathbf{R}}^T & -\hat{\mathbf{R}}^T\hat{\mathbf{t}} \\ \mathbf{0} & 1 \end{bmatrix} \quad \text{where} \quad \hat{\mathbf{R}} = \begin{bmatrix} \hat{\mathbf{r}}_1, \hat{\mathbf{r}}_2, \hat{\mathbf{r}}_3 \end{bmatrix}$$

It is easy to find the correct solution in practice because only one side of the plane is typically visible

**TEK5030**

# Pose estimation with planar correspondences

With a homography estimated from point correspondences,
this approach will typically not give proper rotation matrices because of noise

$$\hat{\mathbf{R}} \notin SO(3)$$

# Pose estimation with planar correspondences

With a homography estimated from point correspondences,
this approach will typically not give proper rotation matrices because of noise

$$\hat{\mathbf{R}} \notin SO(3)$$

But it is possible to find the closest rotation matrix with SVD!

$$\hat{\mathbf{R}} \rightarrow \hat{\mathbf{R}}^* \in SO(3)$$

**TEK5030**

# Pose estimation with planar correspondences

Let $\overline{\mathbf{M}}$ be the matrix with the two first columns of $\mathbf{M}$:

$$\overline{\mathbf{M}} = \begin{bmatrix} \mathbf{m}_1, \mathbf{m}_2 \end{bmatrix}$$

# Pose estimation with planar correspondences

Let $\bar{\mathbf{M}}$ be the matrix with the two first columns of $\mathbf{M}$:

$$\bar{\mathbf{M}} = \begin{bmatrix} \mathbf{m}_1, \mathbf{m}_2 \end{bmatrix}$$

With SVD we can get the decomposition $\bar{\mathbf{M}} = \mathbf{U}_{3 \times 2} \mathbf{\Sigma}_{2 \times 2} \mathbf{V}_{2 \times 2}^T$.

**TEK5030**

# Pose estimation with planar correspondences

Let $\overline{\mathbf{M}}$ be the matrix with the two first columns of $\mathbf{M}$:

$$\overline{\mathbf{M}} = \left[\mathbf{m}_1, \mathbf{m}_2\right]$$

With SVD we can get the decomposition $\overline{\mathbf{M}} = \mathbf{U}_{3\times 2}\boldsymbol{\Sigma}_{2\times 2}\mathbf{V}^T_{2\times 2}$.
The first two columns $\overline{\mathbf{R}}^*$ of a proper $\widehat{\mathbf{R}}^*$ is then

$$\overline{\mathbf{R}}^* = \mathbf{U}\mathbf{V}^T$$

**TEK5030**

# Pose estimation with planar correspondences

Let $\bar{\mathbf{M}}$ be the matrix with the two first columns of $\mathbf{M}$:

$$\bar{\mathbf{M}} = \left[ \mathbf{m}_1, \mathbf{m}_2 \right]$$

With SVD we can get the decomposition $\bar{\mathbf{M}} = \mathbf{U}_{3 \times 2} \mathbf{\Sigma}_{2 \times 2} \mathbf{V}_{2 \times 2}^T$.
The first two columns $\bar{\mathbf{R}}^*$ of a proper $\widehat{\mathbf{R}}^*$ is then

$$\bar{\mathbf{R}}^* = \mathbf{U}\mathbf{V}^T$$

The corresponding scale $\lambda$ can be computed as:

$$\lambda = \frac{\operatorname{trace}\left( \bar{\mathbf{R}}^{*T} \bar{\mathbf{M}} \right)}{\operatorname{trace}\left( \bar{\mathbf{M}}^T \bar{\mathbf{M}} \right)} = \frac{\sum_{i=1}^{3} \sum_{j=1}^{2} r_{ij}^* m_{ij}}{\sum_{i=1}^{3} \sum_{j=1}^{2} m_{ij}^2}$$

**TEK5030**

# Pose estimation with planar correspondences

Let $\bar{\mathbf{M}}$ be the matrix with the two first columns of $\mathbf{M}$:

$$\bar{\mathbf{M}} = \begin{bmatrix} \mathbf{m}_1, \mathbf{m}_2 \end{bmatrix}$$

With SVD we can get the decomposition $\bar{\mathbf{M}} = \mathbf{U}_{3\times2}\mathbf{\Sigma}_{2\times2}\mathbf{V}^T_{2\times2}$.
The first two columns $\bar{\mathbf{R}}^*$ of a proper $\widehat{\mathbf{R}}^*$ is then

$$\bar{\mathbf{R}}^* = \mathbf{U}\mathbf{V}^T$$

The corresponding scale $\lambda$ can be computed as:

$$\lambda = \frac{\text{trace}\left(\bar{\mathbf{R}}^{*T}\bar{\mathbf{M}}\right)}{\text{trace}\left(\bar{\mathbf{M}}^{T}\bar{\mathbf{M}}\right)} = \frac{\sum_{i=1}^{3}\sum_{j=1}^{2} r_{ij}^* m_{ij}}{\sum_{i=1}^{3}\sum_{j=1}^{2} m_{ij}^2}$$

With $\bar{\mathbf{R}}^*$ and $\lambda$, we can now compute the pose with ambiguity as we did in the error-free case

**TEK5030**

# Summary

2D-3D pose estimation:

– Homography-based method

$$\mathbf{H}_{i\Pi} = \mathbf{K}\left[\mathbf{r}_1, \mathbf{r}_2, \mathbf{t}\right]$$

– Minimising geometric/reprojection error

$$\mathbf{T}_{wc}^* = \underset{\mathbf{T}_{wc}}{\operatorname{argmin}} \sum_i \left\| \pi(\mathbf{T}_{wc}^{-1} \cdot \mathbf{x}_i^w) - \mathbf{u}_i \right\|^2$$