

Lecture:

Non-Rigid Structure from Motion

**3D Vision
Universitat Pompeu Fabra**

Discussion

Non-Rigid Shapes?

- Can we obtain non-rigid 3D information from images?

Structure from Motion

(Rigid) Structure from Motion

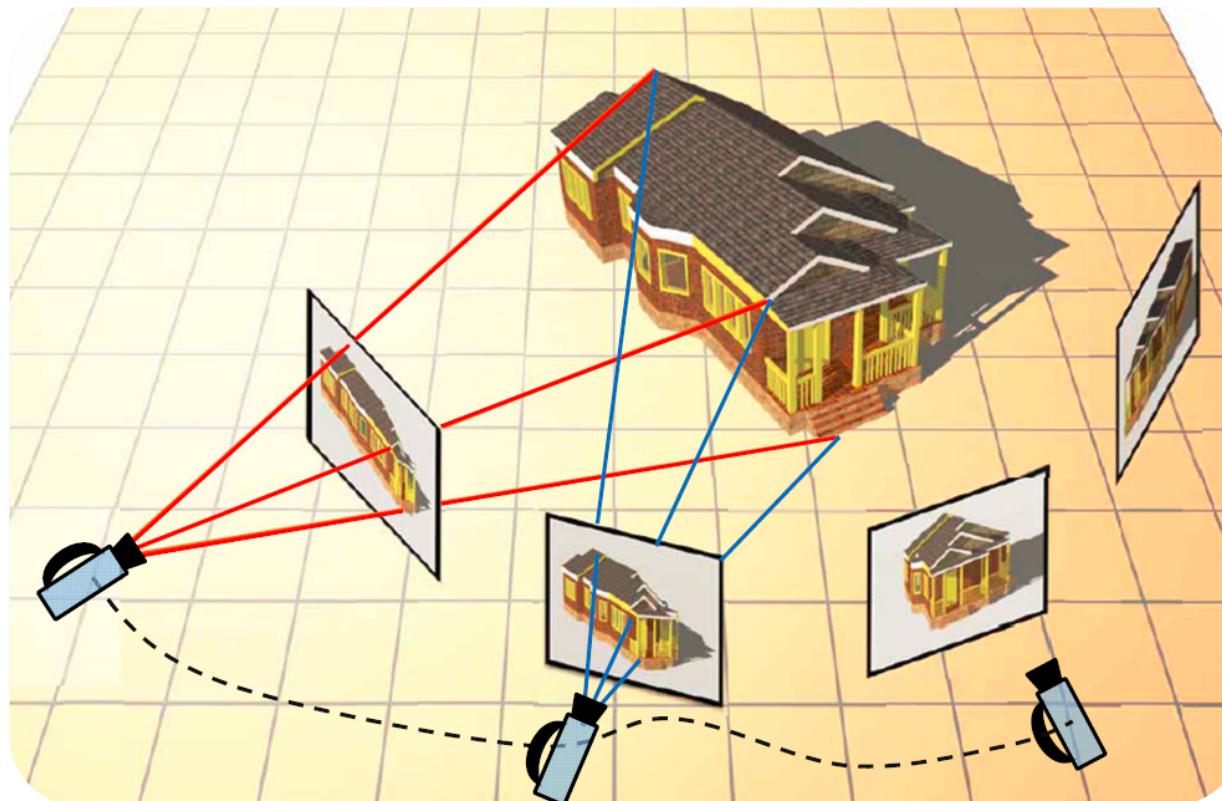
Given: a monocular video (or a collection of pictures)

We want: simultaneously recovering the 3D shape and the camera motion



Epipolar geometry can be used

The assumption of rigidity is enough to make the problem well-posed

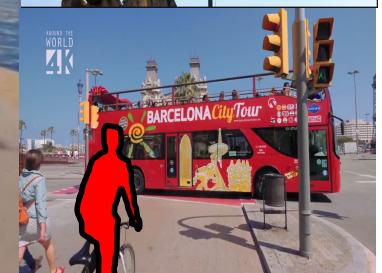
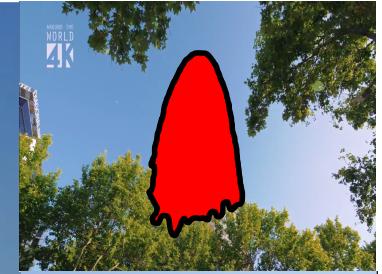


What about non-rigid motion?

Input 2D tracks form an image sequence:

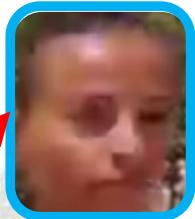
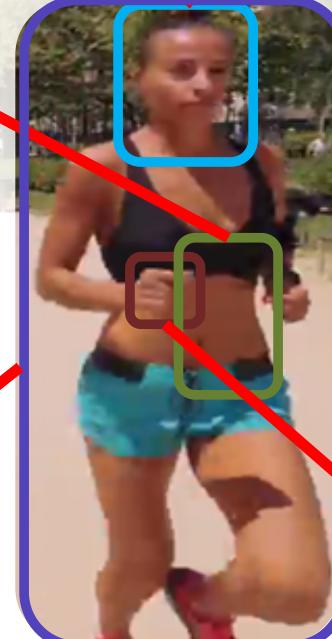
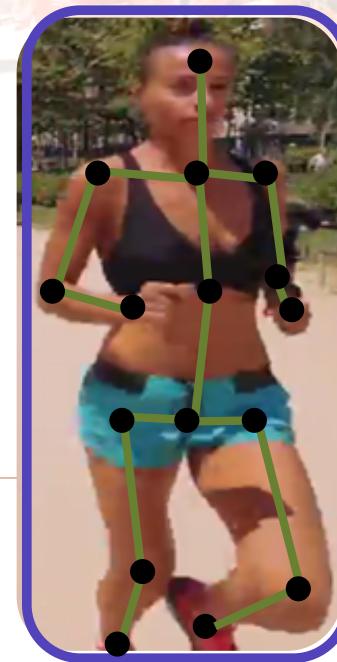


Our world is Non Rigid!



No external markers!

One or Many

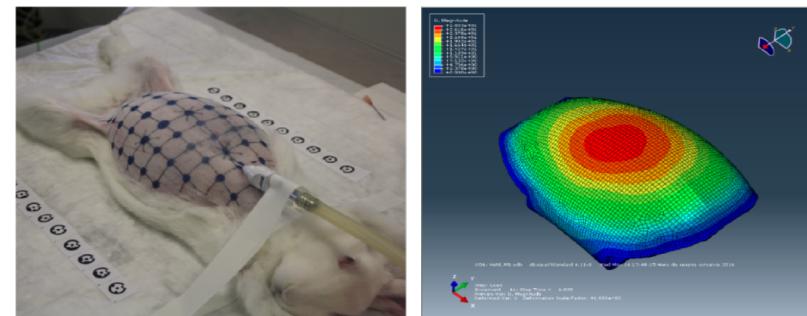


Why is this important?

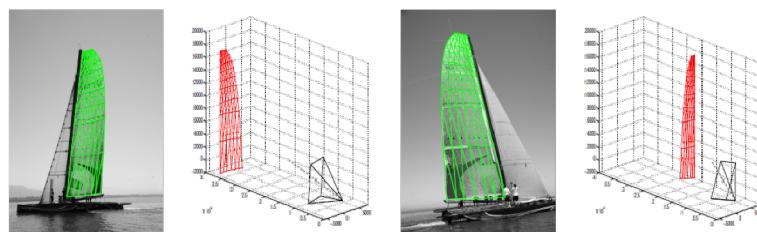
The world is non-rigid! Too many everyday applications in many different domains



Movie industry, augmented reality



Experimental industry



Sport industry: sailing

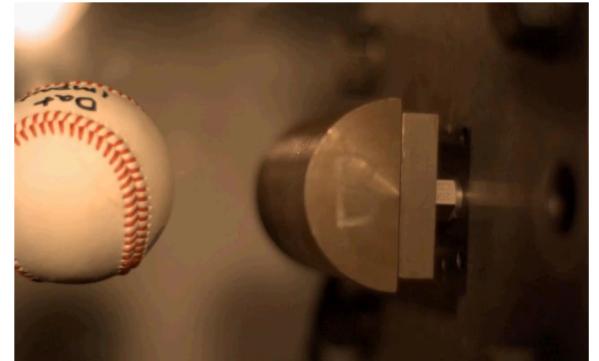


Endoscopy

The movie industry



Even more details



Animal Reconstruction



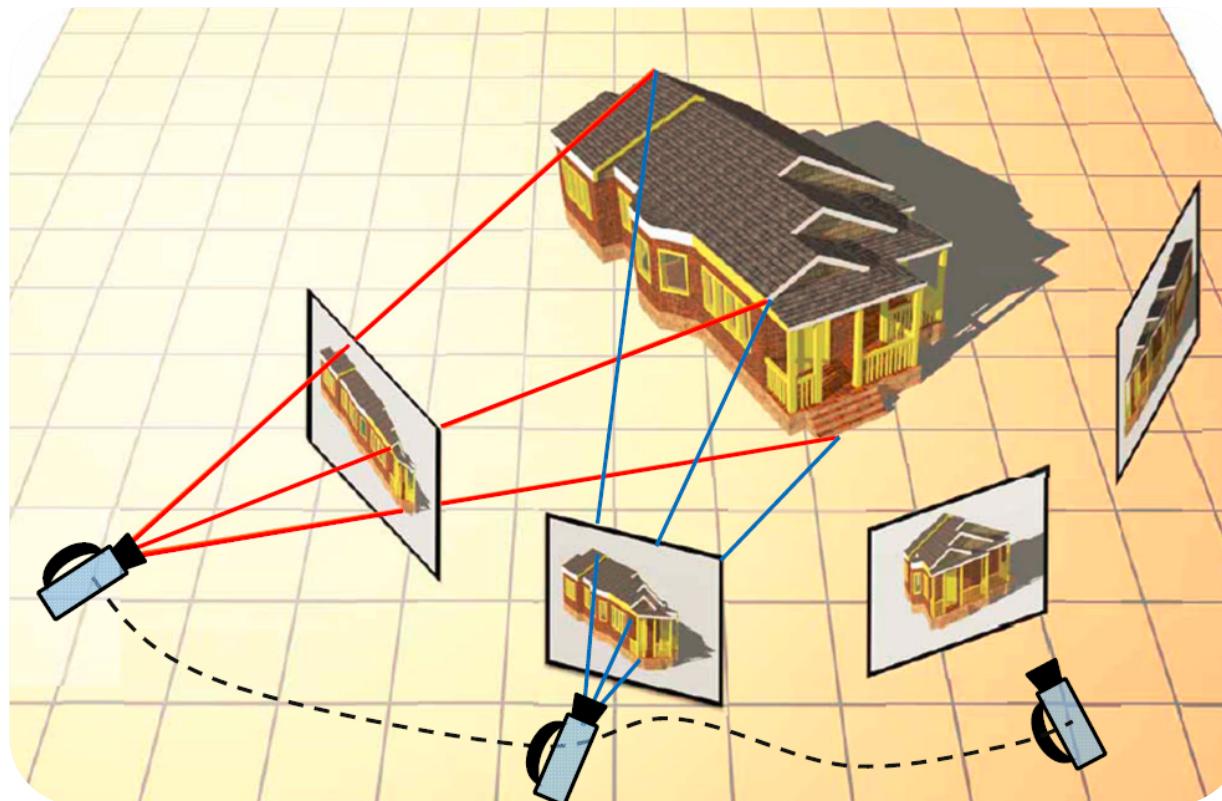
...produces better robots?

Courtesy of BBC



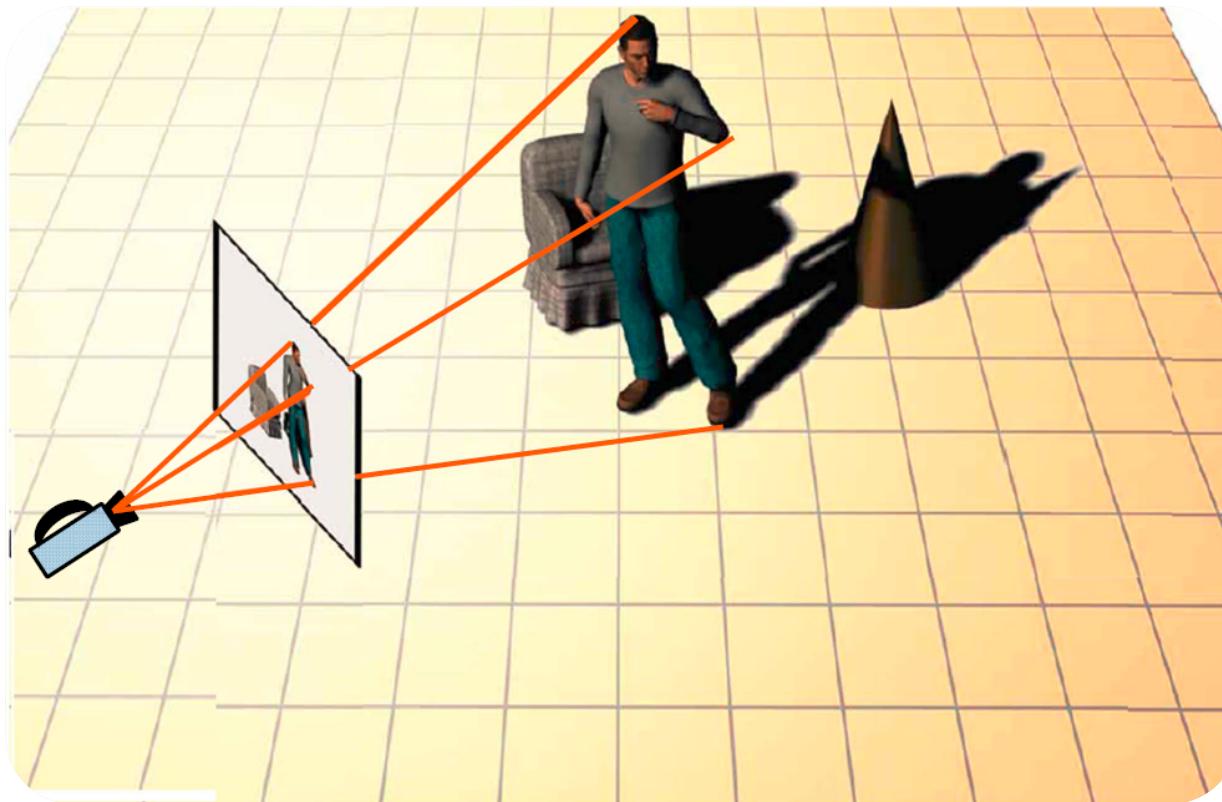
Epipolar geometry can be used

The assumption of rigidity is enough to make the problem well-posed



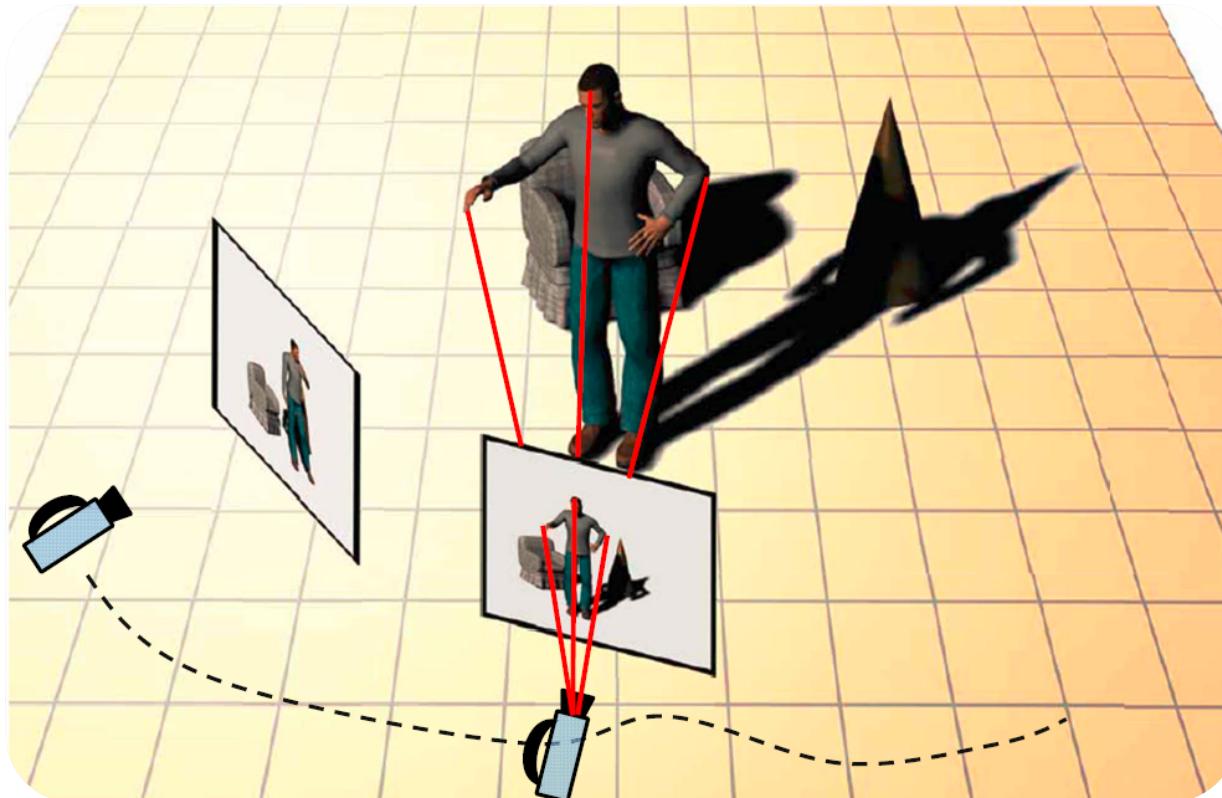
Can Epipolar geometry be used?

Considering only one image, we obtain the same 3D constraint



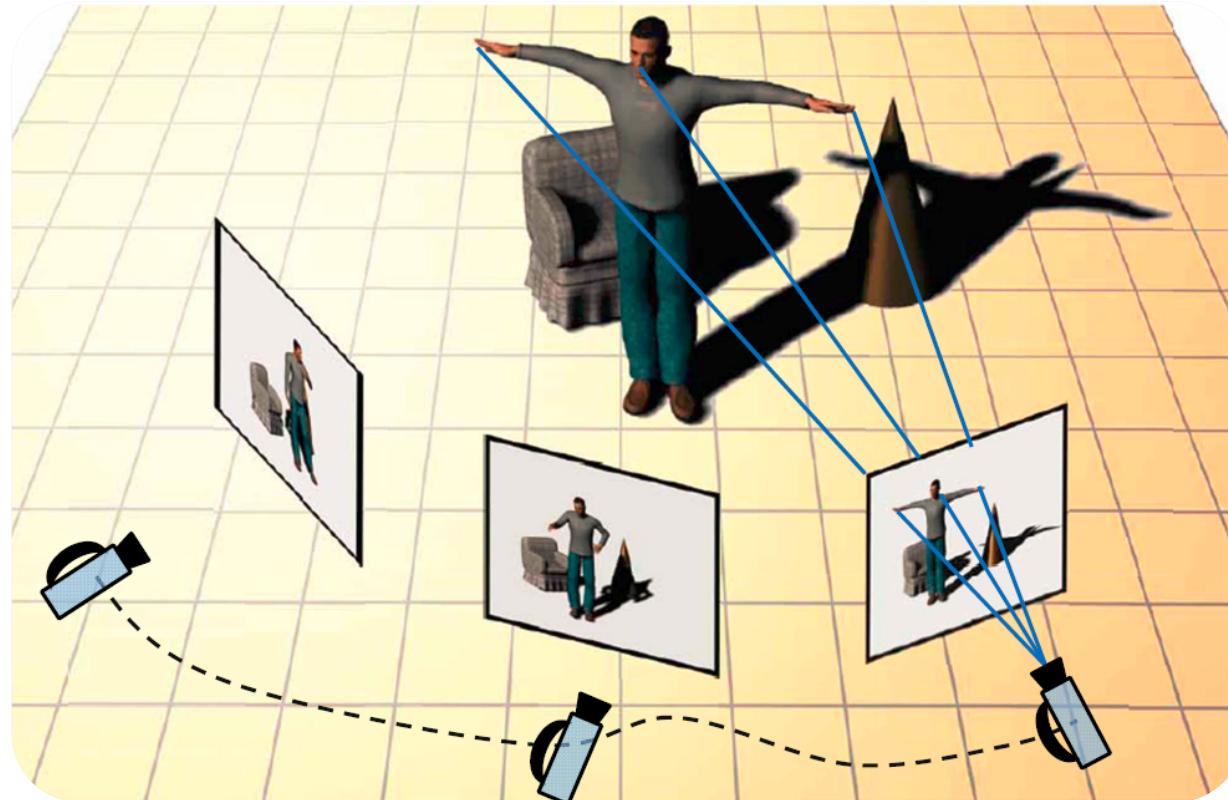
Epipolar geometry ~~can be used~~

After acquiring a new image, we obtain a similar constraint but now triangulation is not available since the shape is non rigid



Epipolar geometry ~~can be used~~

After acquiring a new image, we obtain a similar constraint but now triangulation is not available since the shape is non rigid

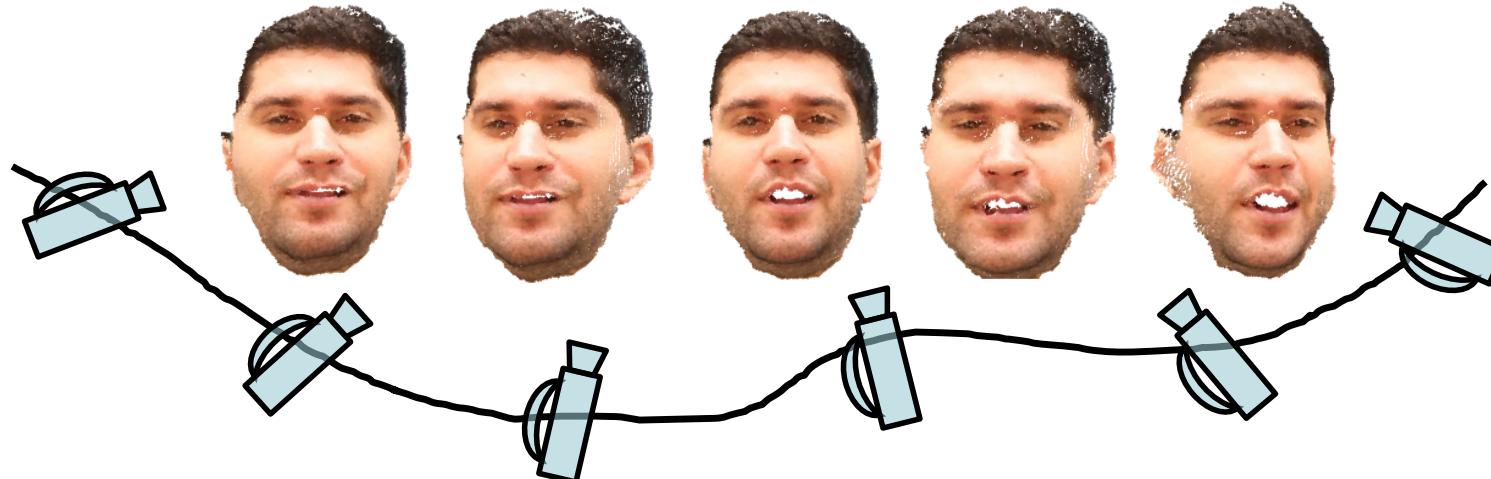


Non-Rigid Structure from Motion

Non-Rigid Structure from Motion

Given: a monocular video (or a collection of pictures)

We want: simultaneously recovering the 3D shape of a time-varying object (4D estimation) and the camera motion



Some Results



Neural Dense Non-Rigid Structure from Motion with Latent Space Constraints

[with voice-over]

Vikramjit Sidhu^{1,2} Edgar Tretschk¹

Vladislav Golyanik¹ Antonio Agudo³ Christian Theobalt¹

¹ Max Planck Institute for Informatics, SIC ² Saarland University

³ Institut de Robòtica i Informàtica Industrial, CSIC-UPC

Solving the problem

The problem can be solved by:

- **Factorization:** a closed-form solution can be achieved by using SVD factorization, enforcing a specific rank (this can change as a function of the type of camera model, or the type of scene). In theory, it is hard to accurately enforce constraints
- **Non-linear Optimization:** the solution is achieved iteratively, the computational cost can be bigger but additional priors can be enforced accurately

In terms of processing, the problem can be solved:

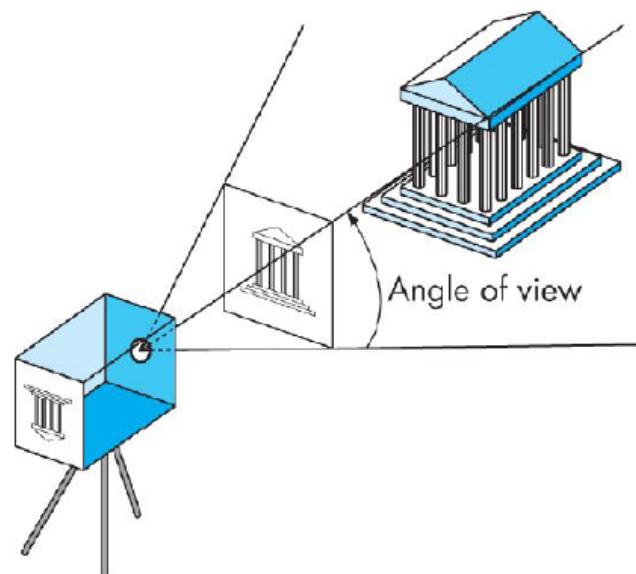
- **Offline:** all the frames are processed at once, after video capture
- **Online:** the frames are processed as the data arrive, frame by frame. More real applications, but can become less accurate

Non-Rigid Structure from Motion

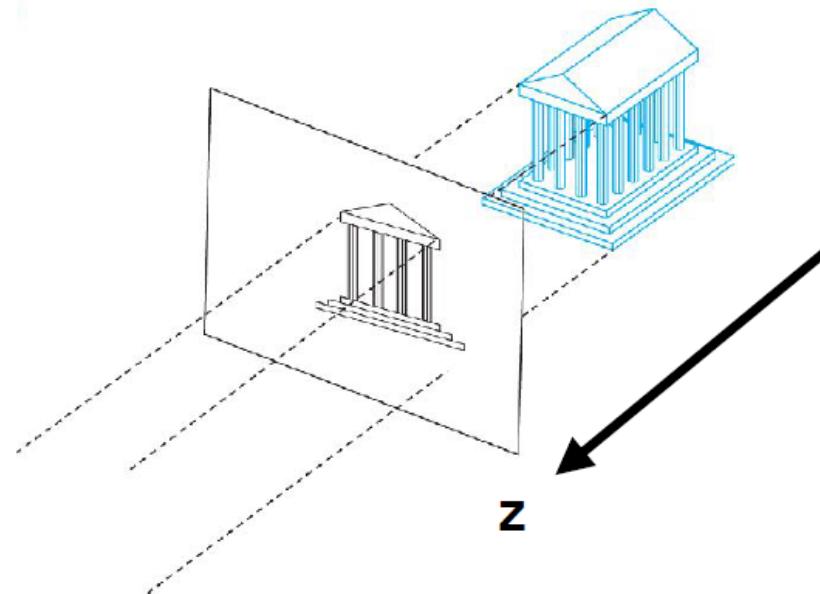
Problem Statement

A Reminder of Camera Models

- Perspective camera: All rays converge to the optical center
- Orthographic camera: All rays are parallel. Z-coordinate is irrelevant in the projection



Perspective camera



Orthographic camera

3D-to-2D: Perspective Model

A p -th 3D point $\mathbf{X}_p = [X_p, Y_p, Z_p]^T$ in homogeneous coordinates can be related with its 2D projection $\mathbf{x}_p = [x_p, y_p]^T$ by means of a matrix \mathbf{P}^i for the i -th image, such as:

$$\begin{bmatrix} \mathbf{x}_p^i \\ 1 \end{bmatrix} = [\mathbf{P}^i] \begin{bmatrix} \mathbf{X}_p \\ 1 \end{bmatrix}$$

where \mathbf{P}^i is 3x4 matrix as:

$$\mathbf{P}^i = \begin{bmatrix} f_x & s & u_0 \\ 0 & f_y & v_0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \end{bmatrix}$$

3D-to-2D: Orthographic Model

A p -th 3D point $\mathbf{X}_p = [X_p, Y_p, Z_p]^T$ can be related with its 2D projection $\mathbf{x}_p = [x_p, y_p]^T$ by means of a matrix \mathbf{R}^i for the i -th image, such as:

$$[\mathbf{x}_p^i] = [\mathbf{R}^i] [\mathbf{X}_p] + \mathbf{t}_p^i$$

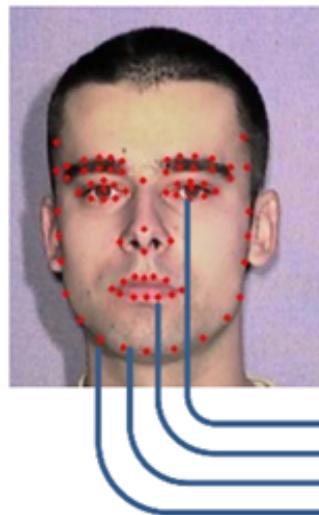
where \mathbf{R}^i is 2x3 matrix and \mathbf{t}^i is a 2x1 translation vector as:

$$\mathbf{R}^i = \begin{bmatrix} r_{11} & r_{12} & r_{13} \\ r_{21} & r_{22} & r_{23} \end{bmatrix} \quad \mathbf{t}^i = \begin{bmatrix} t_x \\ t_y \end{bmatrix}$$

In practice, we subtract the translations by assuming centered observations (i.e., they are equivalent to the mean values of \mathbf{x}_p). For later computations, we will approximate $\mathbf{x}_p = \mathbf{x}_p - \mathbf{t}^i$

Problem Statement

Orthographic camera



A 2D image point is represented as a 2-vector containing the image coordinates at the given frame

$$\begin{pmatrix} \text{[Red Box]} \\ \text{[Red Box]} \\ \text{[Red Box]} \\ \text{[Red Box]} \end{pmatrix} = \begin{pmatrix} \text{[Grey Box]} \end{pmatrix}_{2 \times P} -$$

$$\underbrace{\begin{pmatrix} \text{[Grey Box]} \end{pmatrix}_{2 \times P}}_{\text{Measurement Matrix}} = \underbrace{\begin{pmatrix} \text{[Green Box]} \end{pmatrix}_{2 \times 3}}_{\text{Camera Matrix}} \underbrace{\begin{pmatrix} \text{[Yellow Box]} \end{pmatrix}_{3 \times P}}_{\text{3D shape Matrix}} + \underbrace{\begin{pmatrix} \text{[Purple Box]} \end{pmatrix}_{2 \times P}}_{\text{2D translation}}$$

In the rigid case, it is the same

Full Linear Relation

Orthographic camera



$$\begin{bmatrix} \text{Gray Box} \\ \vdots \\ \text{Gray Box} \end{bmatrix} = \begin{bmatrix} \text{Green Box} \\ \ddots \\ \text{Green Box} \end{bmatrix} \quad \boxed{\begin{bmatrix} \text{Yellow Box} \\ \text{Yellow Box} \end{bmatrix}}$$

2IxP 2Ix3I 3IxP

A different 3D configuration per image

Each image of a non-rigid body increases the dimensionality of the problem and the number of unknowns.

Measurement Matrix

Considering P non-rigid 3D points observed in I RGB images, we can collect all observations to obtain a linear system such as:

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} \mathbf{P}^1 & & \\ & \ddots & \\ & & \mathbf{P}^I \end{bmatrix}}_{\mathbf{R}} \underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{X}}$$
$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} \mathbf{R}^1 & & \\ & \ddots & \\ & & \mathbf{R}^I \end{bmatrix}}_{\mathbf{R}} \underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \end{bmatrix}}_{\mathbf{X}}$$

$3I \times P$

$3I \times 4I$

$4I \times P$

Perspective camera

$2I \times P$

Orthographic camera

where \mathbf{W} is a $3I \times P$ matrix, \mathbf{P} is $3I \times 4I$, and \mathbf{X} is $4I \times P$ for the perspective case (relation on the left); and \mathbf{W} is a $2I \times P$ matrix, \mathbf{R} is $2I \times 3I$, and \mathbf{X} is $3I \times P$ for the perspective case (relation on the right)

What about the rank of \mathbf{W} ?

Considering P non-rigid 3D points observed in I RGB images, we can collect all observations to obtain a linear system such as:

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} \mathbf{P}^1 & & \\ & \ddots & \\ & & \mathbf{P}^I \end{bmatrix}}_{\mathbf{R}} \underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{X}}$$

$3I \times P$

$3I \times 4I$

$4I \times P$

Perspective camera

$$\text{rank}(\mathbf{W}) \leq \min(3I, P)$$

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} \mathbf{R}^1 & & \\ & \ddots & \\ & & \mathbf{R}^I \end{bmatrix}}_{\mathbf{R}} \underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \end{bmatrix}}_{\mathbf{X}}$$

$2I \times P$

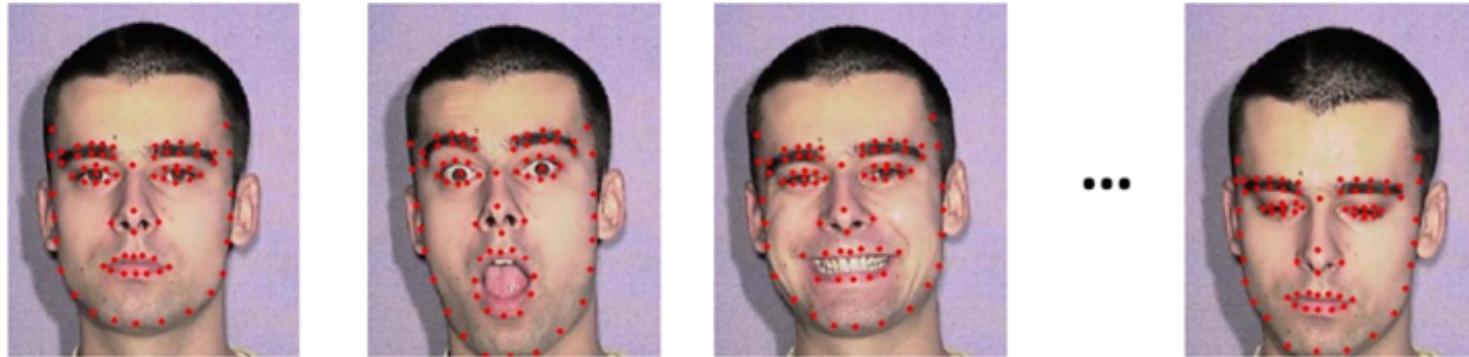
$2I \times 3I$

Orthographic camera

$$\text{rank}(\mathbf{W}) \leq \min(2I, P)$$

A severely ill-posed problem

Orthographic camera



$$\begin{bmatrix} \text{Gray block} \\ \vdots \\ \text{Gray block} \end{bmatrix}_{2I \times P} = \begin{bmatrix} \text{Green block} & & \\ & \ddots & \\ & & \text{Green block} \end{bmatrix}_{2I \times 3I} \begin{bmatrix} \text{Yellow block} \\ \vdots \\ \text{Yellow block} \end{bmatrix}_{3I \times P}$$

Measurement matrix is known Camera matrix is unknown 3D shape matrix is unknown

2ip entries << 6i variables + 3ip variables

This is an explosion of variables

A Toy Comparison

Let us assume a 1 minute video with just 100 tracked points, and considering only the estimation of the 3D shape

Rigid Case

Input data:

100 points x 60 sec x 30 Hz x 2

= 360,000 measurements

Unknowns:

100 points x 3

= 300 unknowns

well-posed problem

Non-Rigid Case

Input data:

100 points x 60 sec x 30 Hz x 2

= 360,000 measurements

ill-posed problem

How can I solve the problem?

The art of priors

Including deformation priors is substantially more difficult than using simple rigidity

Many possibilities were presented

$$F_s(L) = k_s |L - L_0|$$



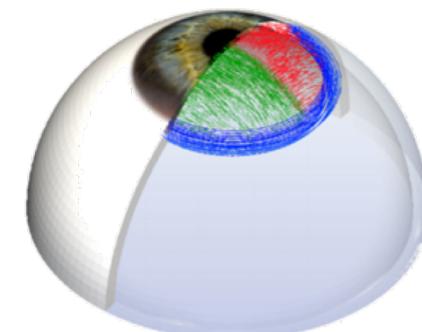
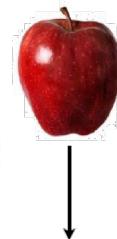
A wide variety of priors in literature:

- **Physical priors.** Particle dynamics, elasticity, finite elements, and many others
- **Probabilistic priors.** Low-rank models on shape, trajectory, shape-trajectory or force domains. Union of subspaces, Gaussian priors
- **Geometric priors:** isometric, as rigid as possible, bone lengths, quadratic models
- **Temporal priors:** temporal-coherent deformations
- **Piecewise priors**
- **Many others**

$$\frac{\nu E}{1 - \nu^2} \nabla(\nabla^\top \mathbf{u}) + \frac{E}{2(1 + \nu)} \nabla^2 \mathbf{u} = -\mathbf{f}_{xy}$$

$$F_g = -mg$$

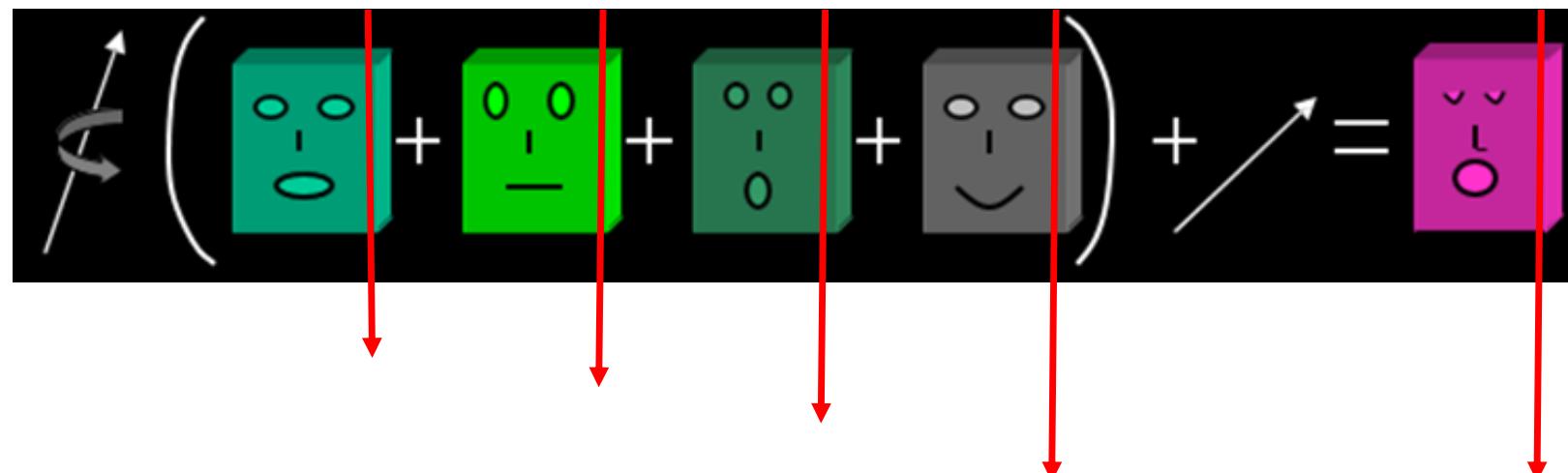
$$g = (0, 0, -9.8) \text{ m/s}^2$$



Shape Linear Subspace (a probabilistic prior)

A Low-Rank Shape Model

Basically, the non-rigid 3D shape can be obtained as a *linear combination of fixed shape vectors*. For every combination of weight coefficients, a different solution can be achieved:



Rotation

Linear combination of
some shapes

Translation

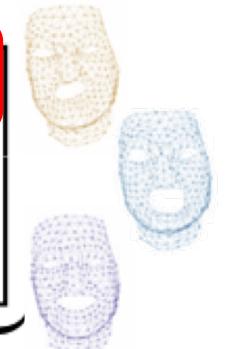
Your estimation

Including the low-rank shape model

We approximate the 3D shape by a linear combination of K shape vectors \mathbf{b} (normally, $K \ll P$ or I). For every k -th component, a weight coefficient l_k is needed. As the shape is non-rigid, by modifying the coefficients for every i -th image, we will change the 3D shape as:

$$\underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \end{bmatrix}}_{\mathbf{X}} = \underbrace{\begin{bmatrix} l_1^1 \mathbf{I}_3 & \dots & l_K^1 \mathbf{I}_3 \\ \vdots & \ddots & \vdots \\ l_1^I \mathbf{I}_3 & \dots & l_K^I \mathbf{I}_3 \end{bmatrix}}_L \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{bmatrix}}_B$$

\mathbf{B}_k



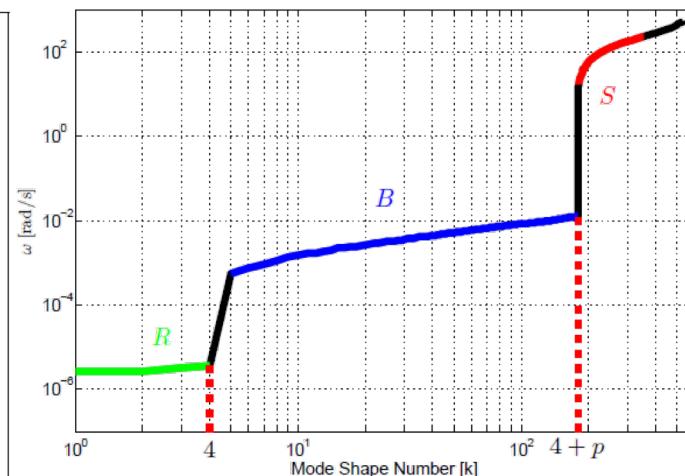
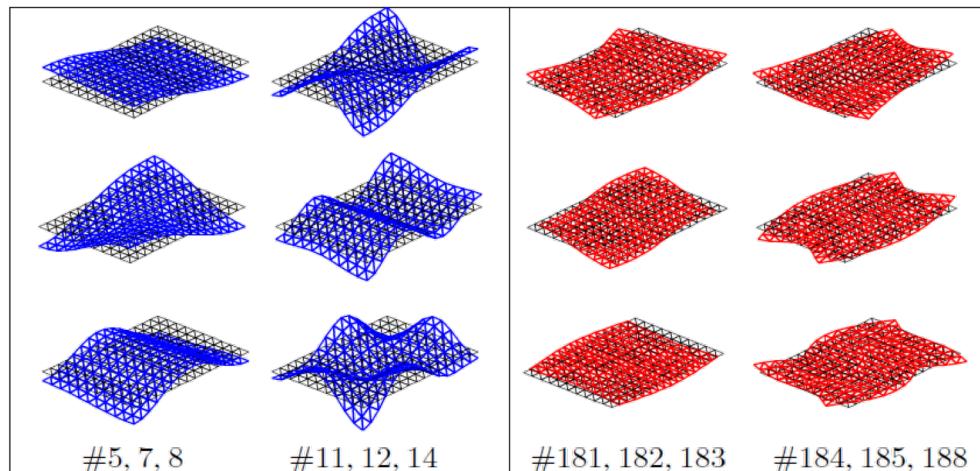
3IxP 3Ix3K 3KxP

Another type of expression for the i -th image: $\mathbf{X}^i = \sum_{k=1}^K l_k^i \mathbf{B}_k$

Shape Basis Estimation

In non-rigid structure from motion, we have some alternatives to estimate the shape basis:

- The most natural is to learn it on the fly, using only the input data
- The input data can also be used to estimate a shape basis from a shape at rest (like a mean shape) by applying:
 - Modal analysis based on physical models
 - Spectral analysis based on a distance matrix
- If training data are assumed, we learn it by means of a learning approach (PCA, deep based, etc.). This approach is supervised



Non-Rigid Structure from motion by factorization

Including the low-rank shape model

Thanks to the relation between the 3D shape and the shape basis:

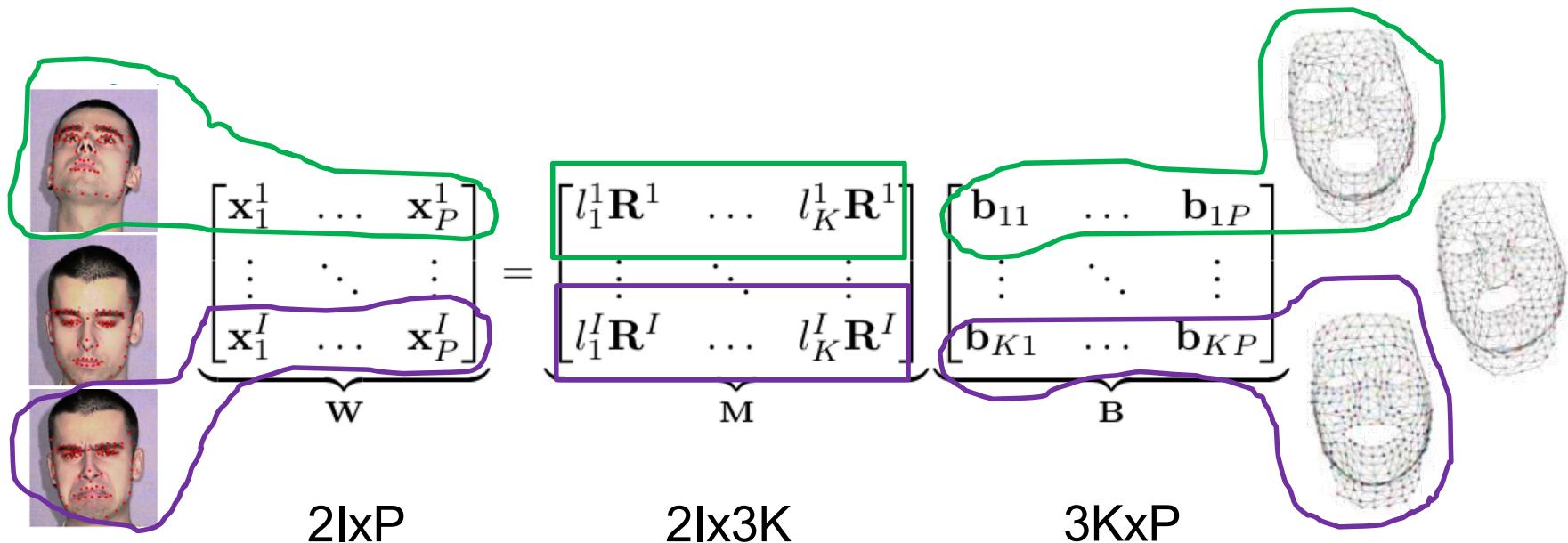
$$\underbrace{\begin{bmatrix} \mathbf{X}_1^1 & \dots & \mathbf{X}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{X}_1^I & \dots & \mathbf{X}_P^I \end{bmatrix}}_{\mathbf{X} \quad 3I \times P} = \underbrace{\begin{bmatrix} l_1^1 \mathbf{I}_3 & \dots & l_K^1 \mathbf{I}_3 \\ \vdots & \ddots & \vdots \\ l_1^I \mathbf{I}_3 & \dots & l_K^I \mathbf{I}_3 \end{bmatrix}}_{\mathbf{L} \quad 3I \times 3K} \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{bmatrix}}_{\mathbf{B} \quad 3K \times P}$$

we obtain the projection equation by using the low-rank shape model as:

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \end{bmatrix}}_{\mathbf{W} \quad 2I \times P} = \underbrace{\begin{bmatrix} l_1^1 \mathbf{R}^1 & \dots & l_K^1 \mathbf{R}^1 \\ \vdots & \ddots & \vdots \\ l_1^I \mathbf{R}^I & \dots & l_K^I \mathbf{R}^I \end{bmatrix}}_{\mathbf{M} \quad 2I \times 3K} \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{bmatrix}}_{\mathbf{B} \quad 3K \times P}$$

Including the low-rank shape model

Orthographic camera



What about the perspective case?

A similar analysis can be followed, but now, considering homogeneous coordinates. We can obtain:

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_W_{3I \times P} = \underbrace{\begin{bmatrix} l_1^1 \hat{\mathbf{P}}^1 & \dots & l_K^1 \hat{\mathbf{P}}^1 & \bar{\mathbf{P}}^1 \\ \vdots & \ddots & \dots & \vdots \\ l_1^I \hat{\mathbf{P}}^I & \dots & l_K^I \hat{\mathbf{P}}^I & \bar{\mathbf{P}}^I \end{bmatrix}}_M_{3I \times 3K+1} \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \\ 1 & \dots & 1 \end{bmatrix}}_B_{3K+1 \times P}$$

$$\mathbf{P}^I = [\hat{\mathbf{P}}^I | \bar{\mathbf{P}}^I]$$

$$3 \times 3 \quad 3 \times 1$$

Factorization

In both cases, the goal is to infer the motion factor (\mathbf{P} or \mathbf{R}) and the 3D coordinates \mathbf{X} of the observed non-rigid object from 2D point tracks in a monocular video \mathbf{W} :

Orthographic camera

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} l_1^1 \mathbf{R}^1 & \dots & l_K^1 \mathbf{R}^1 \\ \vdots & \ddots & \vdots \\ l_1^I \mathbf{R}^I & \dots & l_K^I \mathbf{R}^I \end{bmatrix}}_{\mathbf{M}} \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{bmatrix}}_{\mathbf{B}}$$

$2I \times P$ $2I \times 3K$ $3K \times P$

Perspective camera

$$\underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ 1 & \dots & 1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{W}} = \underbrace{\begin{bmatrix} l_1^1 \hat{\mathbf{P}}^1 & \dots & l_K^1 \hat{\mathbf{P}}^1 & \bar{\mathbf{P}}^1 \\ \vdots & \ddots & \dots & \vdots \\ l_1^I \hat{\mathbf{P}}^I & \dots & l_K^I \hat{\mathbf{P}}^I & \bar{\mathbf{P}}^I \end{bmatrix}}_{\mathbf{M}} \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \\ 1 & \dots & 1 \end{bmatrix}}_{\mathbf{B}}$$

$3I \times P$ $3I \times 3K+1$ $3K+1 \times P$

The full linear system

$$\mathbf{W} = \mathbf{M}\mathbf{B}$$

Two factors: motion factor \mathbf{M} (camera rotation and weight coefficients) and shape one as a product of \mathbf{B} and the coefficients

More on factorization

Orthographic camera

Because \mathbf{M} is a $2I \times 3K$ matrix and \mathbf{B} is a $3K \times P$ matrix, the rank of \mathbf{W} is $3K$. If we apply SVD to \mathbf{W} , we will have only **$3K$ non-zero singular values**

However, measurements are normally noisy, and in practice the rank will not be $3K$. We have to impose it

Applying SVD factorization, we have:

$$\mathbf{W} = \mathbf{U}\mathbf{A}\mathbf{V}^T = [\mathbf{U}\sqrt{\mathbf{A}}][\sqrt{\mathbf{A}}\mathbf{V}^T] = [\mathbf{U}\sqrt{\mathbf{A}}\mathbf{Q}][\mathbf{Q}^{-1}\sqrt{\mathbf{A}}\mathbf{V}^T]$$

i.e., $\mathbf{M} = \mathbf{U}\sqrt{\mathbf{A}}\mathbf{Q}$ and $\mathbf{B} = \mathbf{Q}^{-1}\sqrt{\mathbf{A}}\mathbf{V}^T$ (the two factors we look for)

Many solutions can be achieved by modifying \mathbf{Q} . Of course, for all invertible $3K \times 3K$ \mathbf{Q} matrices

We need to tune the rank K a priori

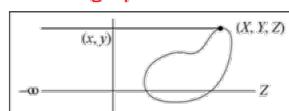
Metric Upgrade

How is \mathbf{Q} computed?

Enforcing orthogonality constraints on the camera rotation. A rotation matrix always has some properties (it is not a random matrix), since lies in the $\text{SO}(3)$ manifold

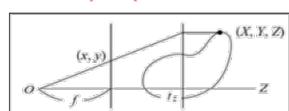
Be careful. Now, matrix \mathbf{M} also includes the weight coefficients in addition to the camera rotations!

1.Orthographic:



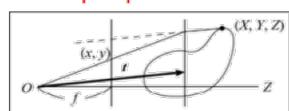
$$\begin{matrix} \mathbf{R}_i \\ 2 \times 3 \end{matrix} \quad \begin{matrix} \mathbf{R}_i^T \\ 3 \times 2 \end{matrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad \text{2x2 identity matrix}$$

2.Weak perspective:



$$\begin{matrix} \mathbf{R}_i \\ 2 \times 3 \end{matrix} \quad \begin{matrix} \mathbf{R}_i^T \\ 3 \times 2 \end{matrix} = \begin{bmatrix} a & 0 \\ 0 & a \end{bmatrix} \quad \text{2x2 diagonal matrix}$$

3.Para perspective:



$$\begin{matrix} \mathbf{R}_i \\ 2 \times 3 \end{matrix} \quad \begin{matrix} \mathbf{R}_i^T \\ 3 \times 2 \end{matrix} = \begin{bmatrix} a & c \\ c & b \end{bmatrix} \quad \text{2x2 full matrix}$$

**But in many cases, we
cannot observe all the
points in all the images**

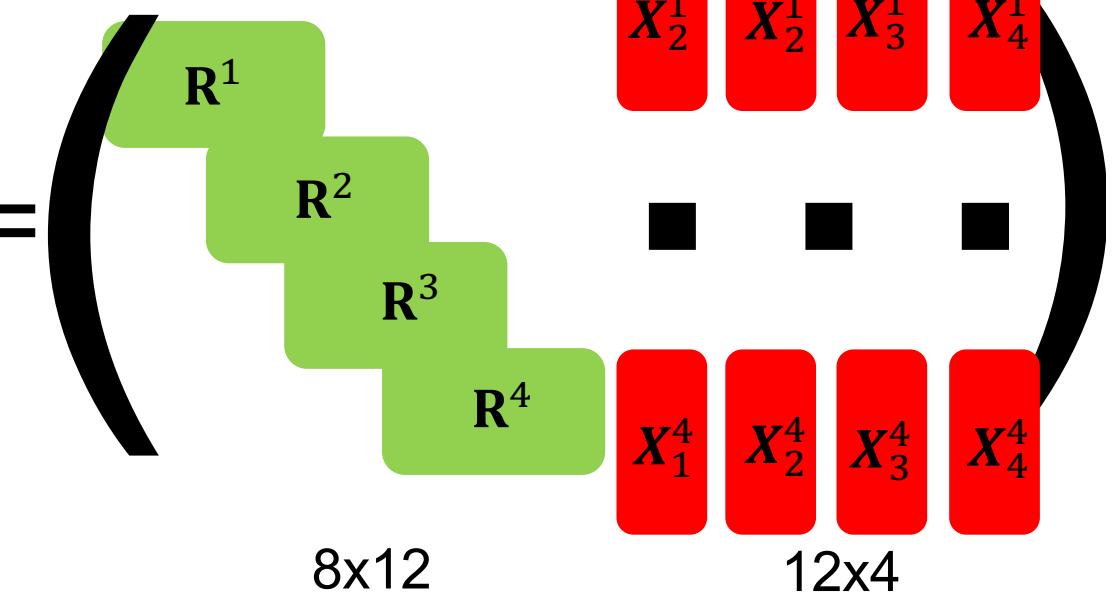
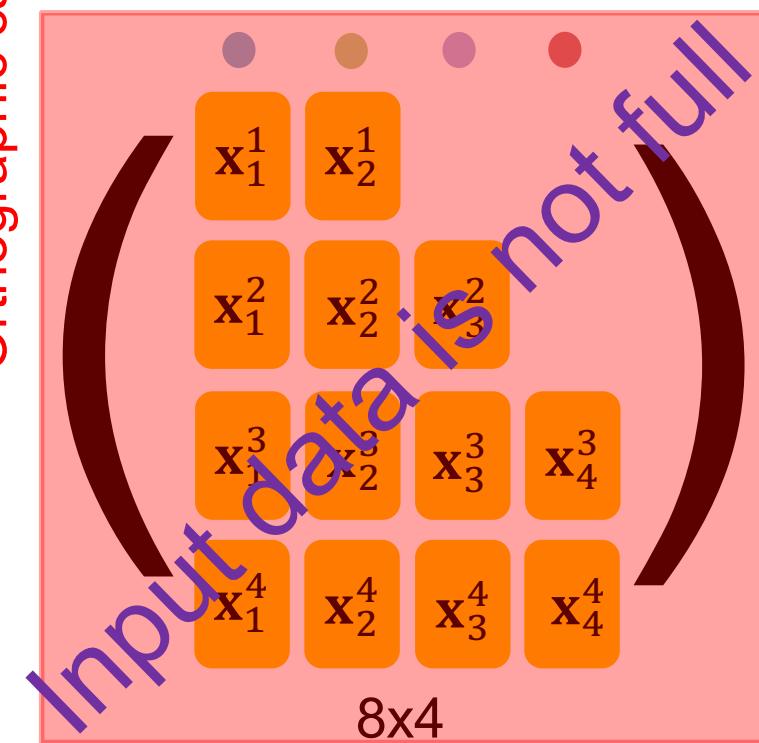
==

Missing tracks

A toy example with missing tracks



Orthographic camera



Handling missing tracks

Two alternatives are possible:

- Applying a *matrix completion* algorithm to infer the missing entries, and then run factorization over the full measurement matrix
- No consider missing entries in the formulation by applying non-linear optimization. Once the 3D model and camera pose are computed, the 2D missing tracks can be inferred too



Non-Rigid Structure from Motion by Non-Linear Optimization

Problem Statement

For an orthographic camera, we have:

$$\hat{\mathbf{W}}^i = \underbrace{\begin{bmatrix} \mathbf{x}_1^1 & \dots & \mathbf{x}_P^1 \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^I & \dots & \mathbf{x}_P^I \end{bmatrix}}_W = \underbrace{\begin{bmatrix} l_1^1 \mathbf{R}^1 & \dots & l_K^1 \mathbf{R}^1 \\ \vdots & \ddots & \vdots \\ l_1^I \mathbf{R}^I & \dots & l_K^I \mathbf{R}^I \end{bmatrix}}_M \underbrace{\begin{bmatrix} \mathbf{b}_{11} & \dots & \mathbf{b}_{1P} \\ \vdots & \ddots & \vdots \\ \mathbf{b}_{K1} & \dots & \mathbf{b}_{KP} \end{bmatrix}}_B \mathbf{B}_k$$

The problem (compacting over the points) can be formulated as:

$$\arg \min_{\mathbf{R}^i, \mathbf{B}_k, l_k^i, \mathbf{t}^i} \sum_{i=1}^I \left\| \hat{\mathbf{W}}^i - \mathbf{R}^i \sum_{k=1}^K l_k^i \mathbf{B}_k - \mathbf{t}^i \right\|^2$$

and we perform non-linear optimization by minimizing a geometric error cost function. Translation \mathbf{t}^i is optional

Bundle Adjustment

Normally, the Levenberg-Marquardt method is used to minimize the problem. We need a Jacobian matrix \mathbf{J} as the derivative of the function with respect to the unknowns (\mathbf{R} , \mathbf{B} and the set of weight \mathbf{l}_k)

Again, there are many variants on how to proceed to reduce the computational complexity of the problem:

- Alternate minimization of motion and shape parameters
- *Sparse methods.* The computation of \mathbf{J} is complex, but it can be approximated by considering a *binary pattern*

Initialization: The optimization can be initialized assuming a rigid shape, i.e., using rigid factorization or non-linear optimization for a rigid shape

Bundle Adjustment

The bundle adjustment method:

- Minimize the cost function with *Levenberg-Marquadt*
- Exploit the sparseness of the Jacobian function matrix to *decrease computation and memory requirements*

The Levenberg-Marquadt algorithm does:

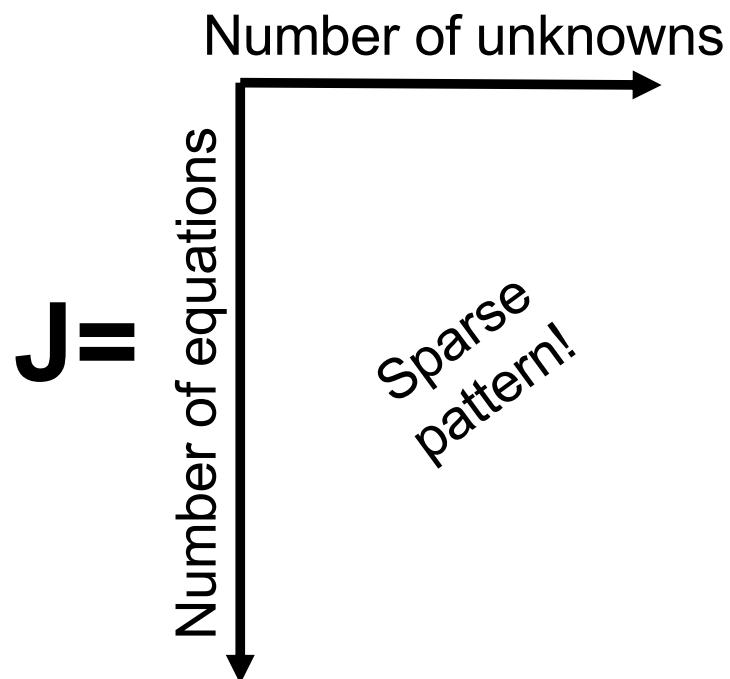
- Mixture of Gauss-Newton and Gradient descent
- Behaves like Gauss-Newton when close to the minimum (quadratic region)
- Gradient descent when the prediction is poor
- Depends on a parameter θ that controls the *mixture* of Gauss-Newton and Gradient descent as:

$$(JJ^T + \theta I) \delta p = -g$$

Parameters we
want to estimate

Exercise

Let us assume a monocular video of 3 images, where 6 points are observed. Considering the map is non-rigid and the visibility is full, define the corresponding Jacobian matrix. A low-rank shape model of rank 2 can be considered



Including priors

As in the rigid case, we can apply temporal smoothness priors, but now, in both camera motion and shape deformation (be careful when input data are a collection of pictures). To this end, we may consider the expression:

$$\arg \min_{\mathbf{R}^i, \mathbf{B}_k, l_k^i, \mathbf{t}^i} \sum_{i=1}^I \|\hat{\mathbf{W}}^i - \mathbf{R}^i \sum_{k=1}^K l_k^i \mathbf{B}_k - \mathbf{t}^i\|^2 + \gamma \sum_{i=1}^{I-1} \|\mathbf{L}^i - \mathbf{L}^{i+1}\|_{\mathcal{F}}^2 + \phi \sum_{i=1}^{I-1} \|\mathbf{R}^i - \mathbf{R}^{i+1}\|_{\mathcal{F}}^2$$

where \mathbf{L}^i includes all K weight coefficients in the i-th image

How can we obtain a sequential solution?



We solve the optimization in a sequential manner, considering the information as *the data arrive*. Future frames are not available. Two options:

- Pure sequential (frame by frame)
- Sliding window (from 3 to 5 consecutive frames)

Initialization is performed by rigid estimation (assuming just the initial frames). The problem is actually challenging

**Then, can we infer 3D information
from a single RGB image?**

Yes... but with some constraints

Maybe, the most used approach to handle this scenario is assuming training data. For instance, a shape basis can be computed from training data, and on testing, we only have to adjust the weight coefficients. Supervised approach.

Obtaining training data is not easy in some scenarios

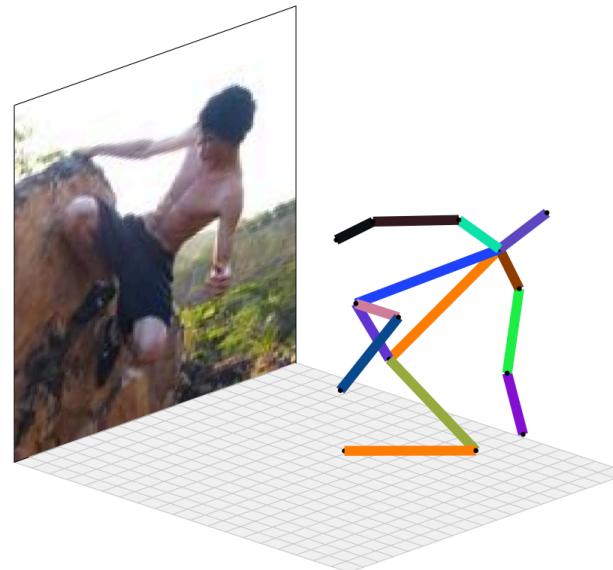
Training data for a type of object or deformation are not available to handle other cases

Alternative approaches can be also found in literature, by exploiting silhouettes or a mean shape in combination with geometrical priors to establish correspondences

Exploiting Training Data

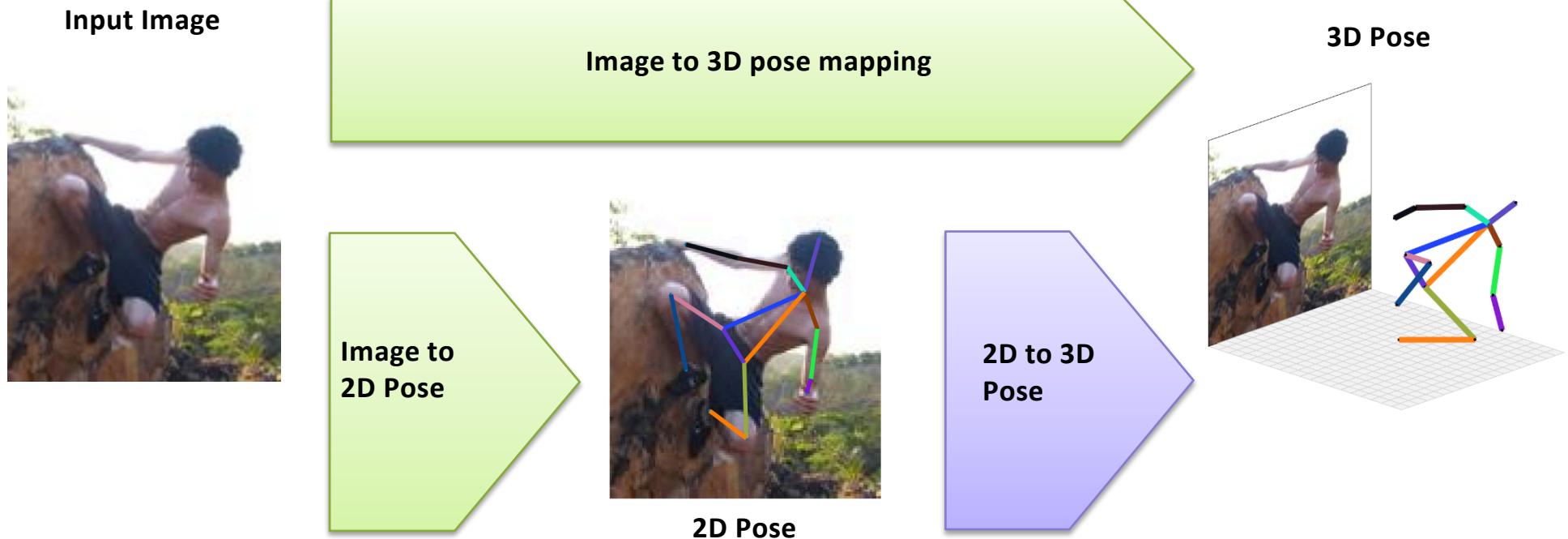
Problem Formulation

Objective: Given a **single image** of an object (a person in the example), estimate the 3D position of its body represented by a skeleton with $P=14$ joints:

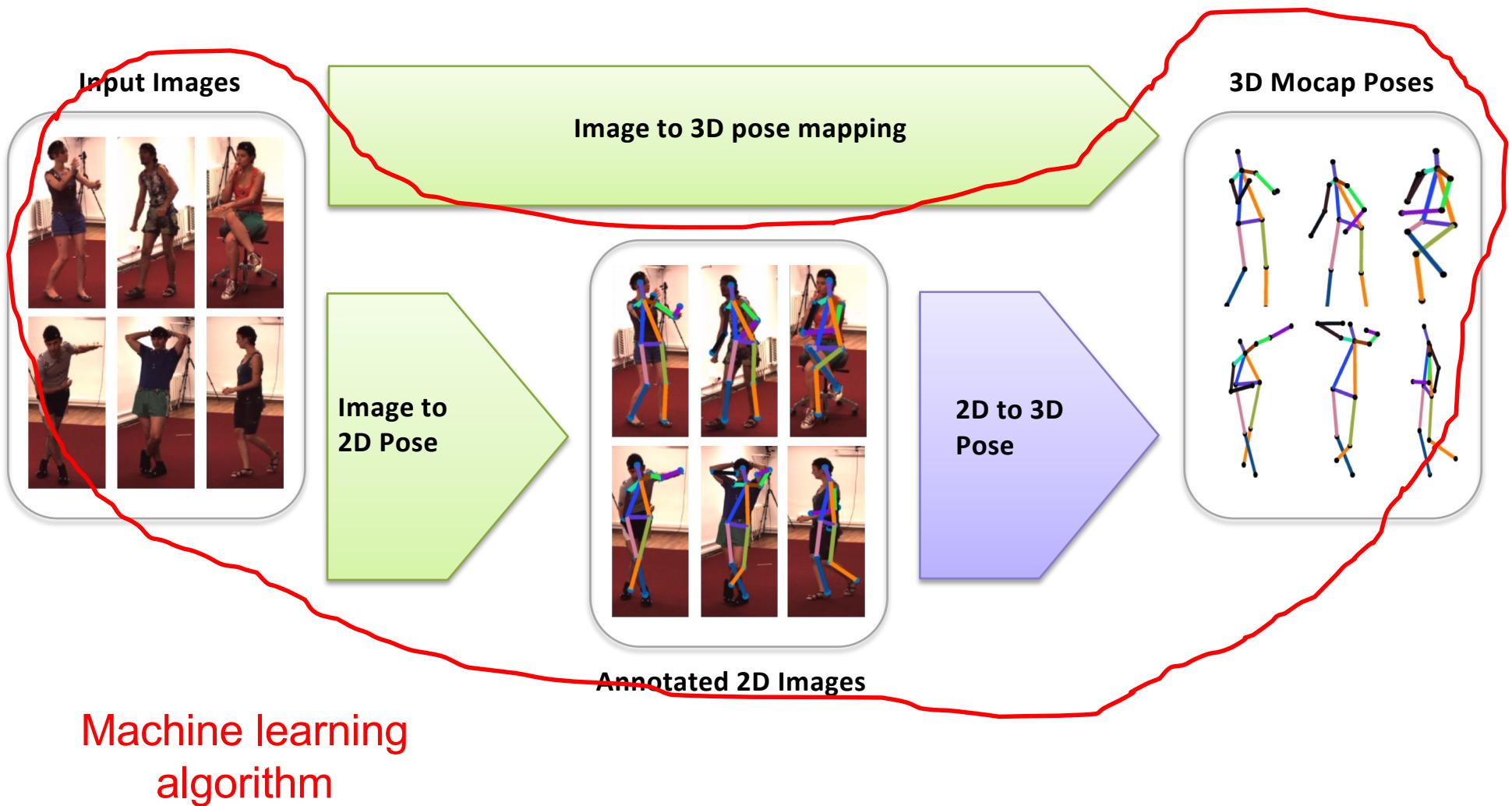


Challenge: This is an ill-posed problem; many different 3D poses have very similar 2D projections

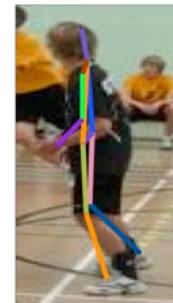
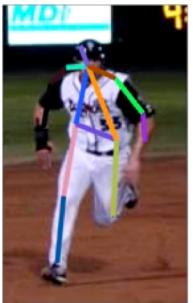
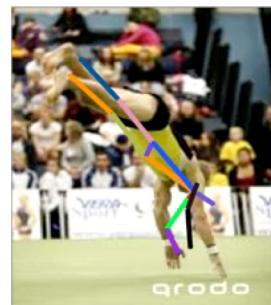
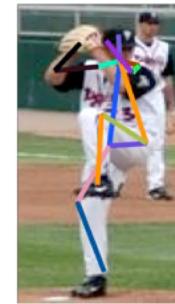
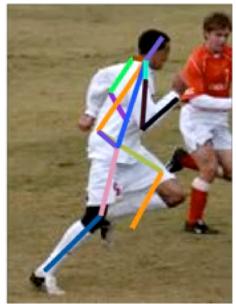
From image to 3D



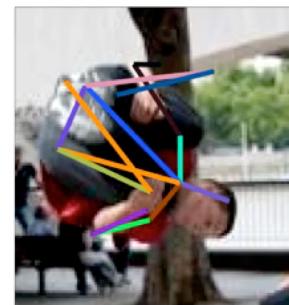
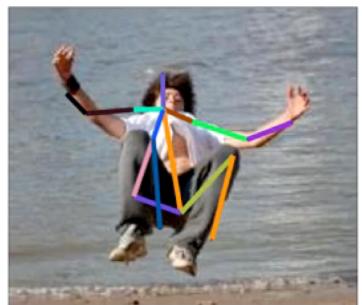
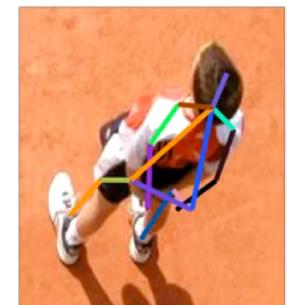
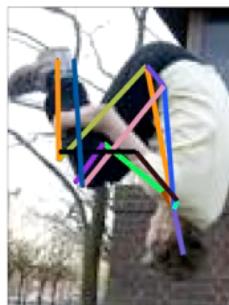
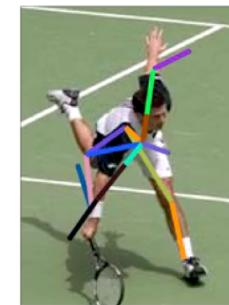
The Learning Process (Training)



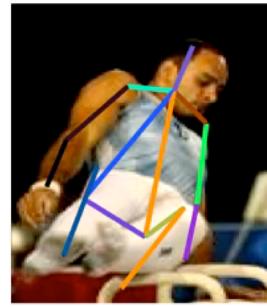
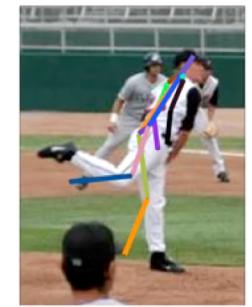
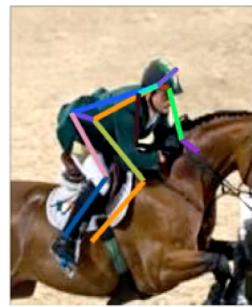
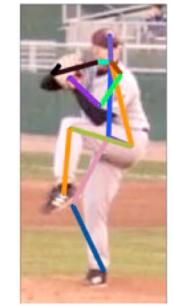
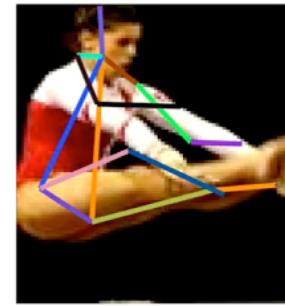
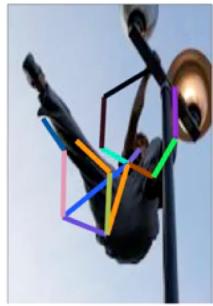
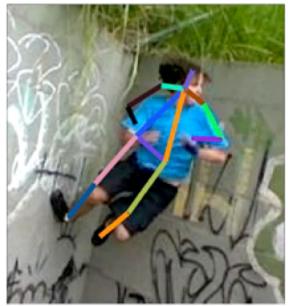
Some Results



Some Results

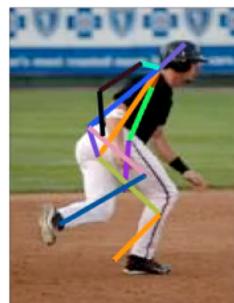
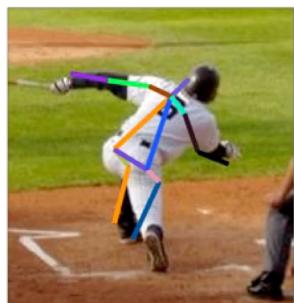


Some Results



Some Results: Failure Cases

Many 3D configurations can produce the same 2D projection



Things to remember

3D and 4D information can be obtained from a sequence of images

For rigid objects, the problem is well-posed. For non-rigid ones, it is inherently ill-posed (additional priors are necessary)

Model-based approaches can handle a wide variety of deformations. They are normally universal and generic. No supervision is needed

Data-based approaches require a lot of data to constrain the solution space. Obtaining *good* data can become hard. Only for a particular object or deformation (depending on the training data)

Future must be unsupervised, and probably, combining both model- and data-based approaches. With a hand-held camera, performing the estimation of multiple scenarios

Acknowledgments

Thanks to Kris Kitani, Yaser Sheikh, Alessio del Bue, Lourdes Agapito