

SHAPE FROM X

One image:

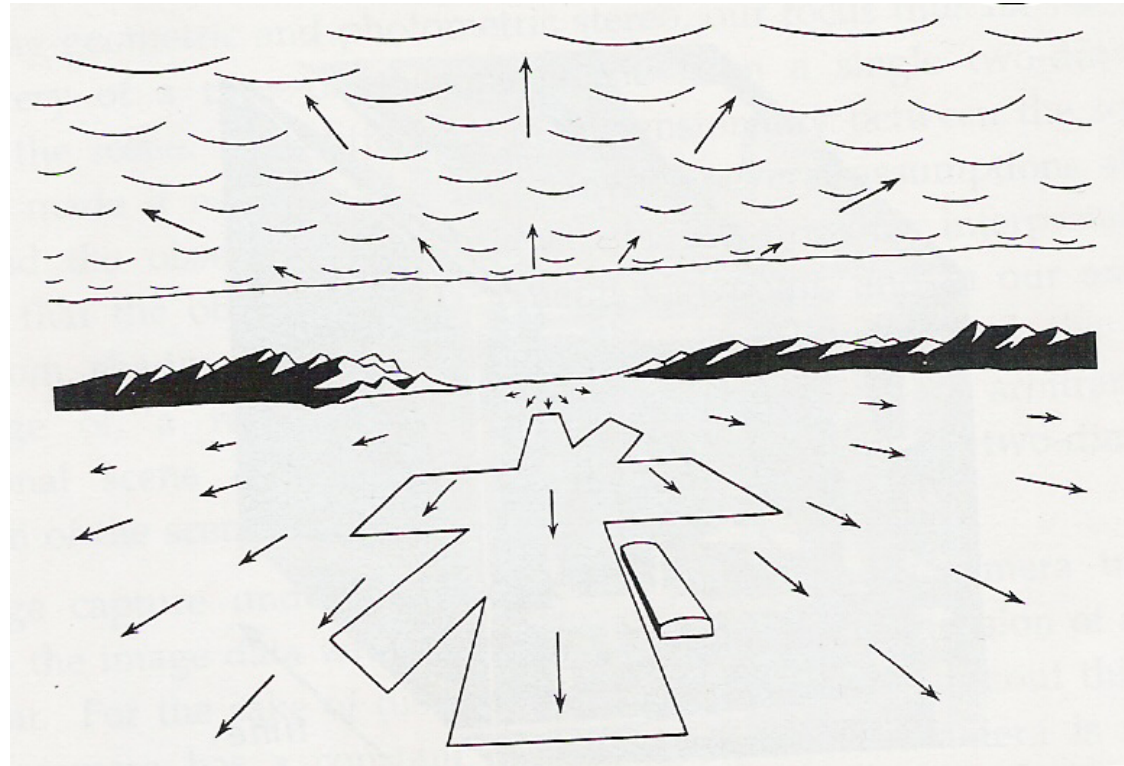
- Texture
- Shading

Two images or more:

- Stereo
- Contours
- **Motion**



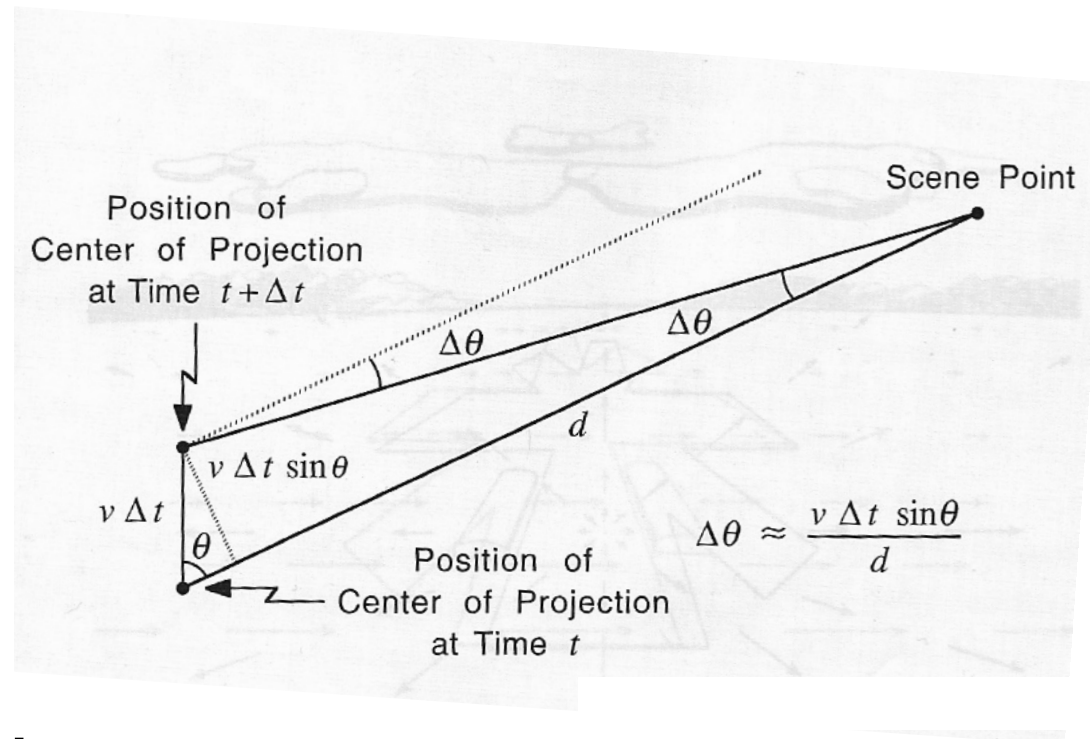
MOTION



When objects move at equal speed, those more remote seem to move more slowly.

Euclid, 300 BC

VELOCITY vs DISTANCE



Velocity is:

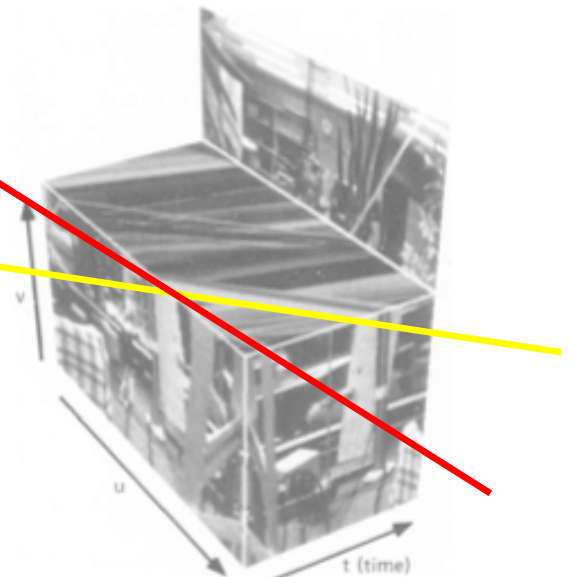
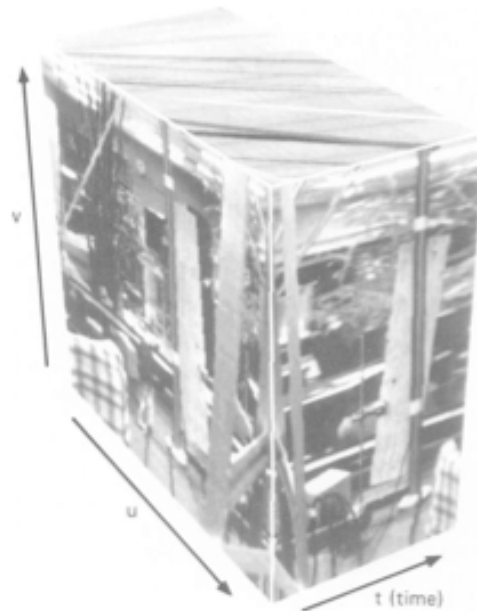
- Inversely proportional to the distance of the point to the observer.
- Proportional to the sine of the angle between the line of sight and the direction of translation.

EPIPOLAR PLANE ANALYSIS

Sequence:



Image cube:

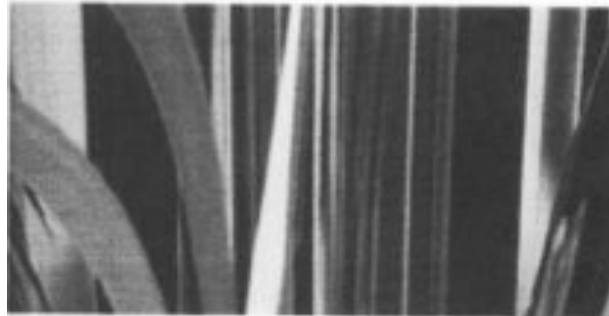


Bolles et al. , IJCV'87

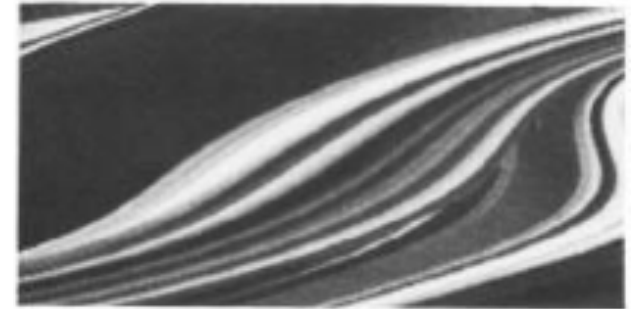
GENERALIZED MOTION



Orthogonal
viewing

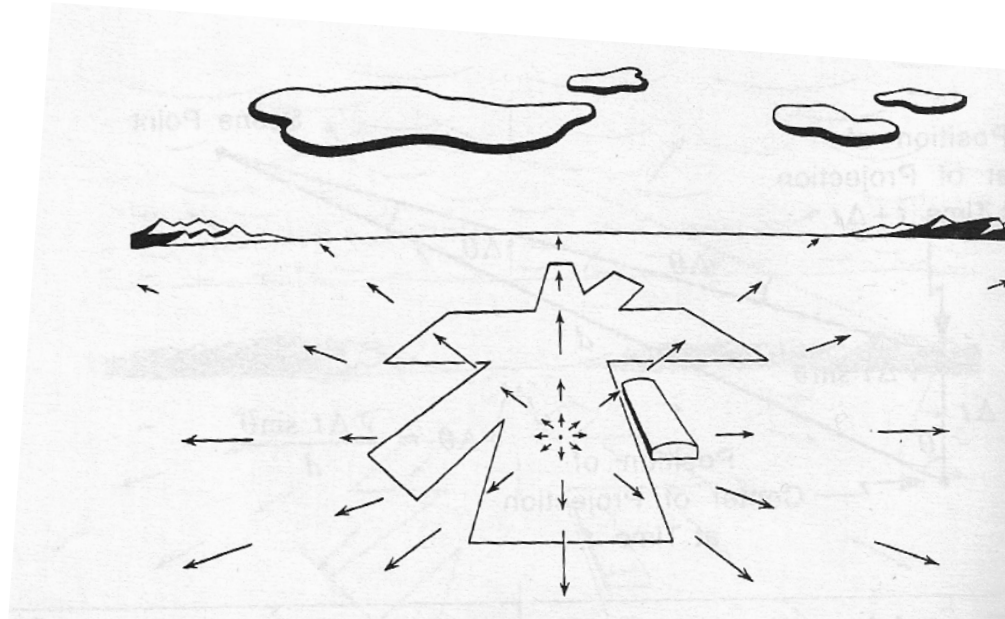


Non-orthogonal
viewing



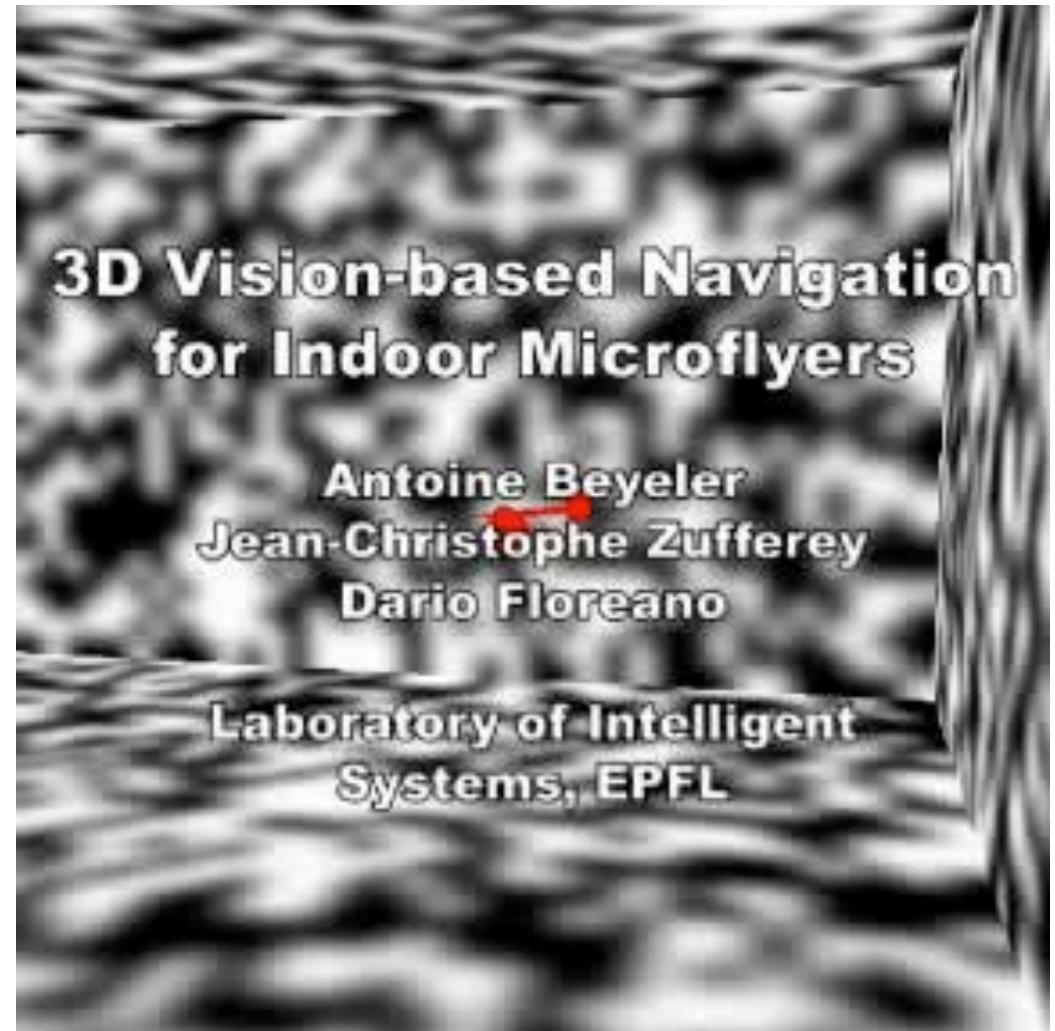
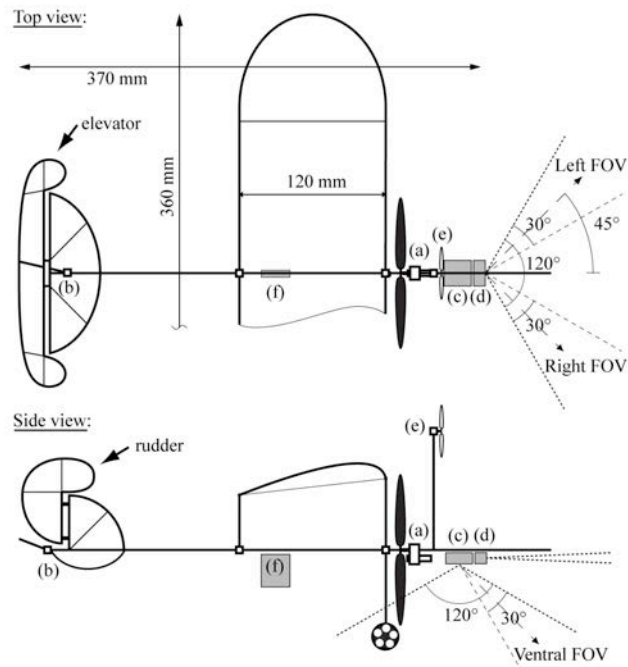
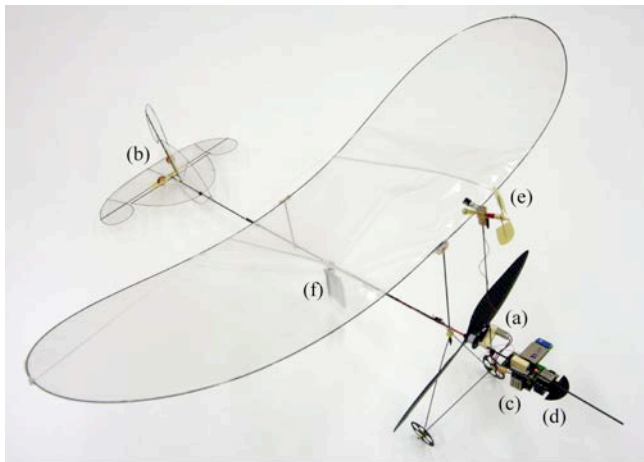
View direction
varying

FOCUS OF EXPANSION



For a translational motion of the camera, all the **motion-field** vectors converge or diverge from a single point: the focus of expansion (FOE) or contraction (FOC).

MICROFLYER



MOTION FIELD ESTIMATION



Approaches classified with respect to the assumptions they make about the scene:

- Images properties are preserved under relative motion between the camera and the scene.
- Feature points can be tracked across frames.

ASSUMPTION 1: BRIGHTNESS CONSTANCY

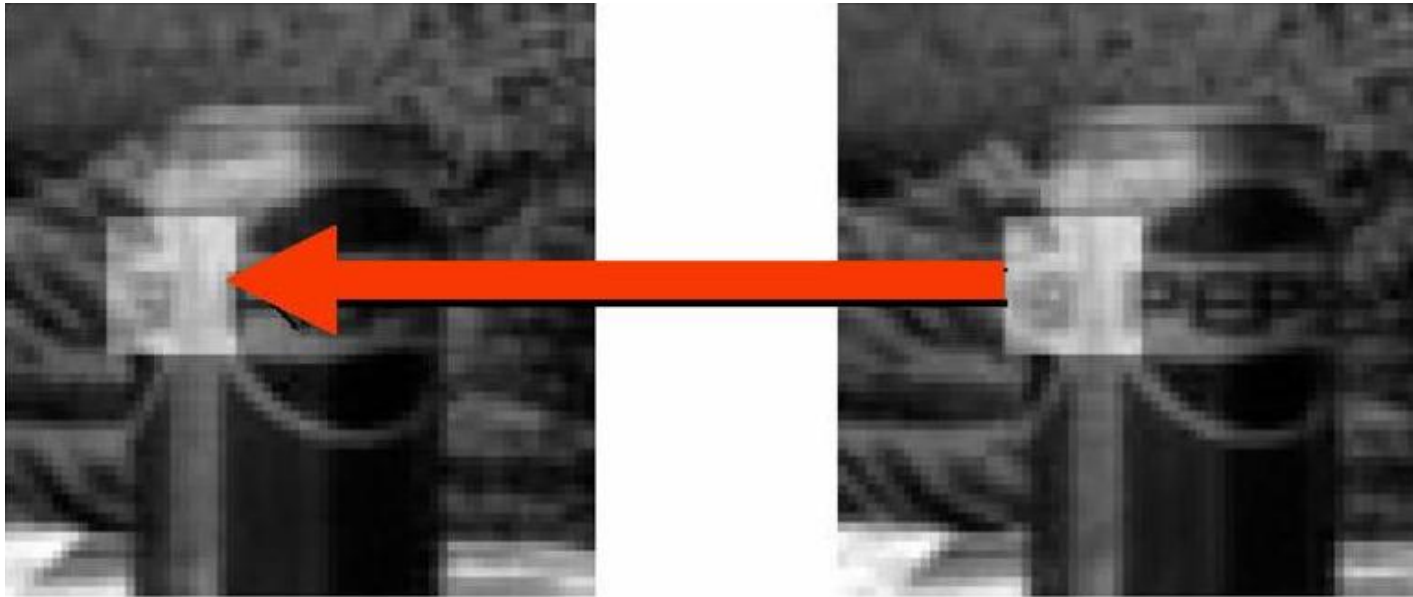
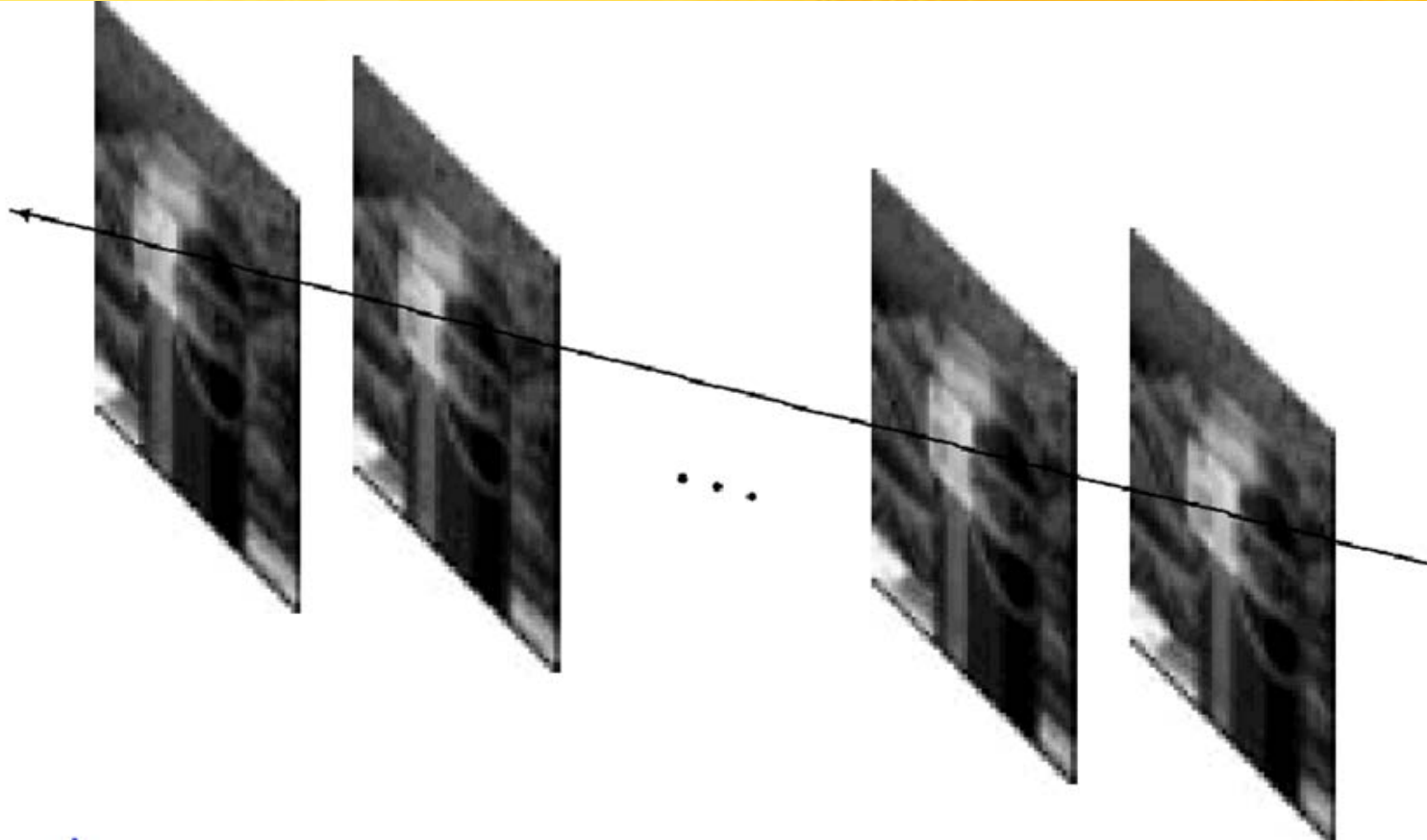


Image measurements (e.g. brightness) in a small region remain the same although its location may change.

$$I(x + u, y + v, t) = I(x, y, t)$$

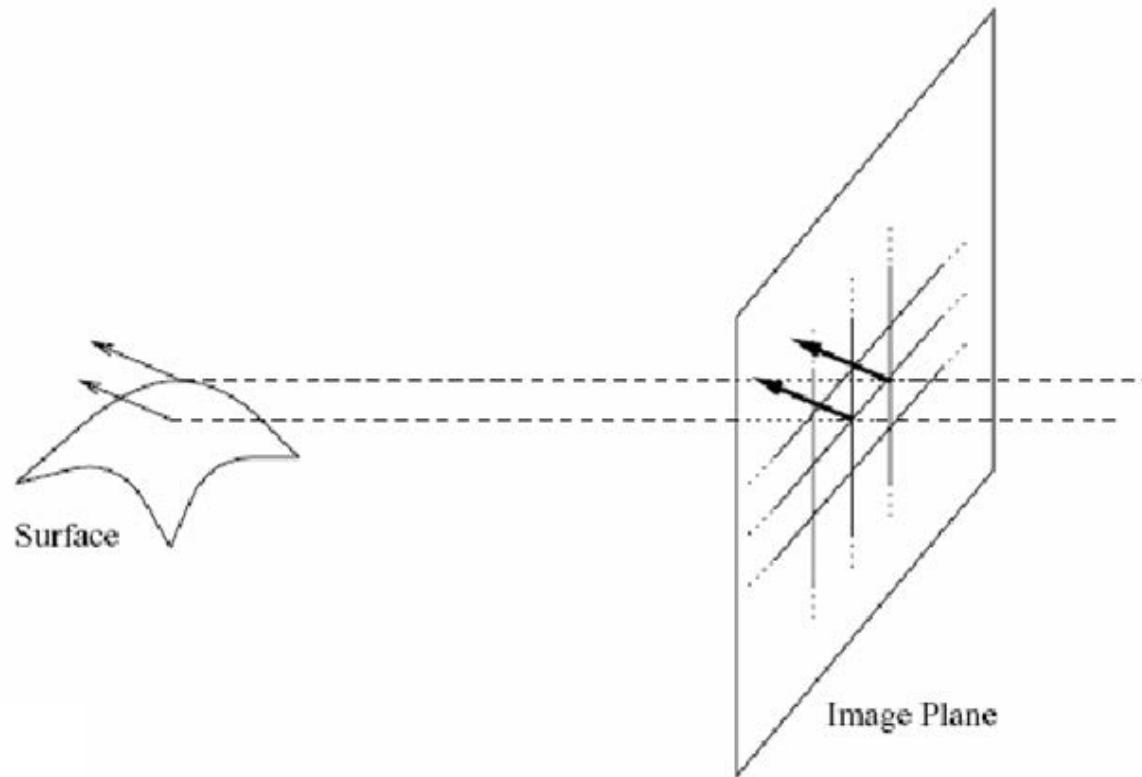
ASSUMPTION 2: TEMPORAL CONSISTENCY



The image motion of a surface patch changes gradually over time.

ASSUMPTION 3:

SPATIAL COHERENCE



- Neighboring points in the scene typically belong to the same surface and hence have similar motions.
- Since they also project to nearby image locations, we expect spatial coherence of the flow.

SPATIO TEMPORAL DERIVATIVES



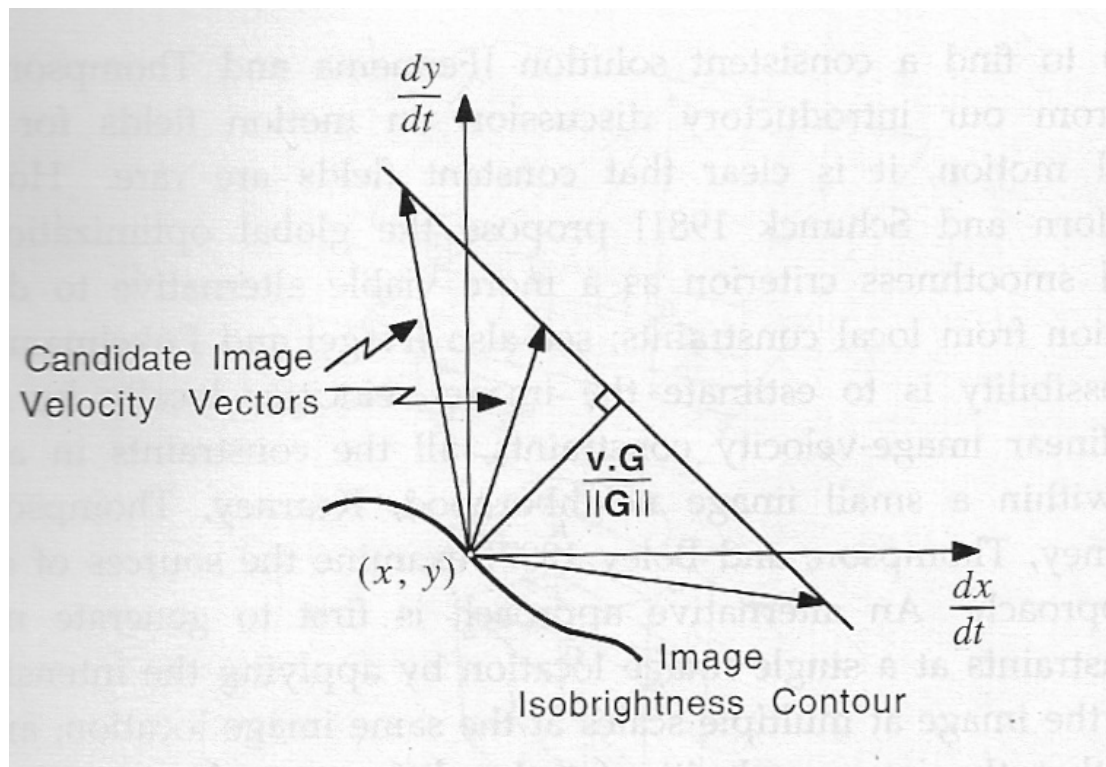
Under the assumptions of

- Brightness constancy,
- Temporal consistency,

we write:

$$\begin{aligned} \text{cst} &= I(x(t), y(t), t) \\ \Rightarrow 0 &= \frac{\delta I}{\delta x} \frac{dx}{dt} + \frac{\delta I}{\delta y} \frac{dy}{dt} + \frac{\delta I}{\delta t} \end{aligned}$$

NORMAL FLOW EQUATION



$$v \frac{G}{\|G\|} = - \frac{\frac{\partial I}{\partial t}}{\sqrt{\frac{\partial I^2}{\partial x} + \frac{\partial I^2}{\partial y}}}$$

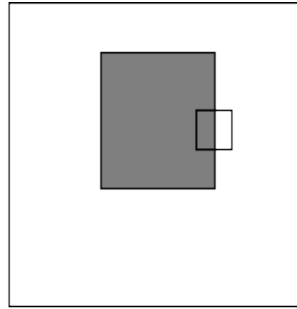
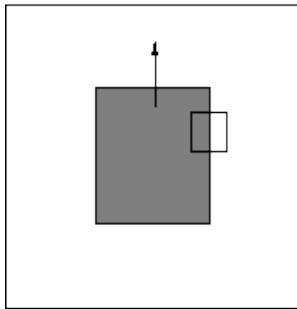
$$G = \left[\frac{\partial I}{\partial x}, \frac{\partial I}{\partial y} \right]$$

$$v = \left[\frac{dx}{dt}, \frac{dy}{dt} \right]$$

AMBIGUITIES

At each pixel, we have one equation and two unknowns.

--> Only the flow component in the gradient direction can be determined locally.



The motion is parallel to the edge, and it cannot be determined.

LOCAL CONSTANCY

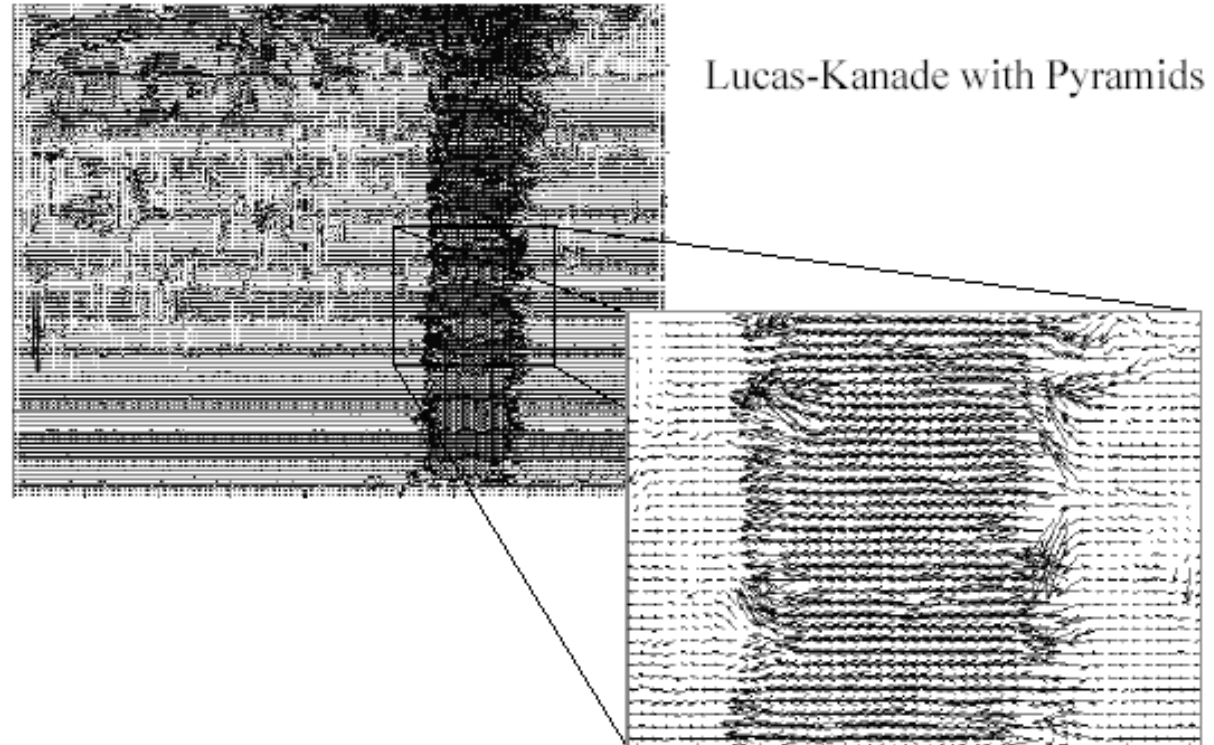


Assume the flow to constant is a 5x5 window:

$$\begin{bmatrix} I_x(p_1) & I_y(p_1) \\ I_x(p_2) & I_y(p_2) \\ \vdots & \vdots \\ I_x(p_{25}) & I_y(p_{25}) \end{bmatrix} \begin{bmatrix} u \\ v \end{bmatrix} = - \begin{bmatrix} I_t(p_1) \\ I_t(p_2) \\ \vdots \\ I_t(p_{25}) \end{bmatrix}$$

--> 25 equations for 2 unknown, which can be solved in the least squares sense.

ENFORCING CONSISTENCY



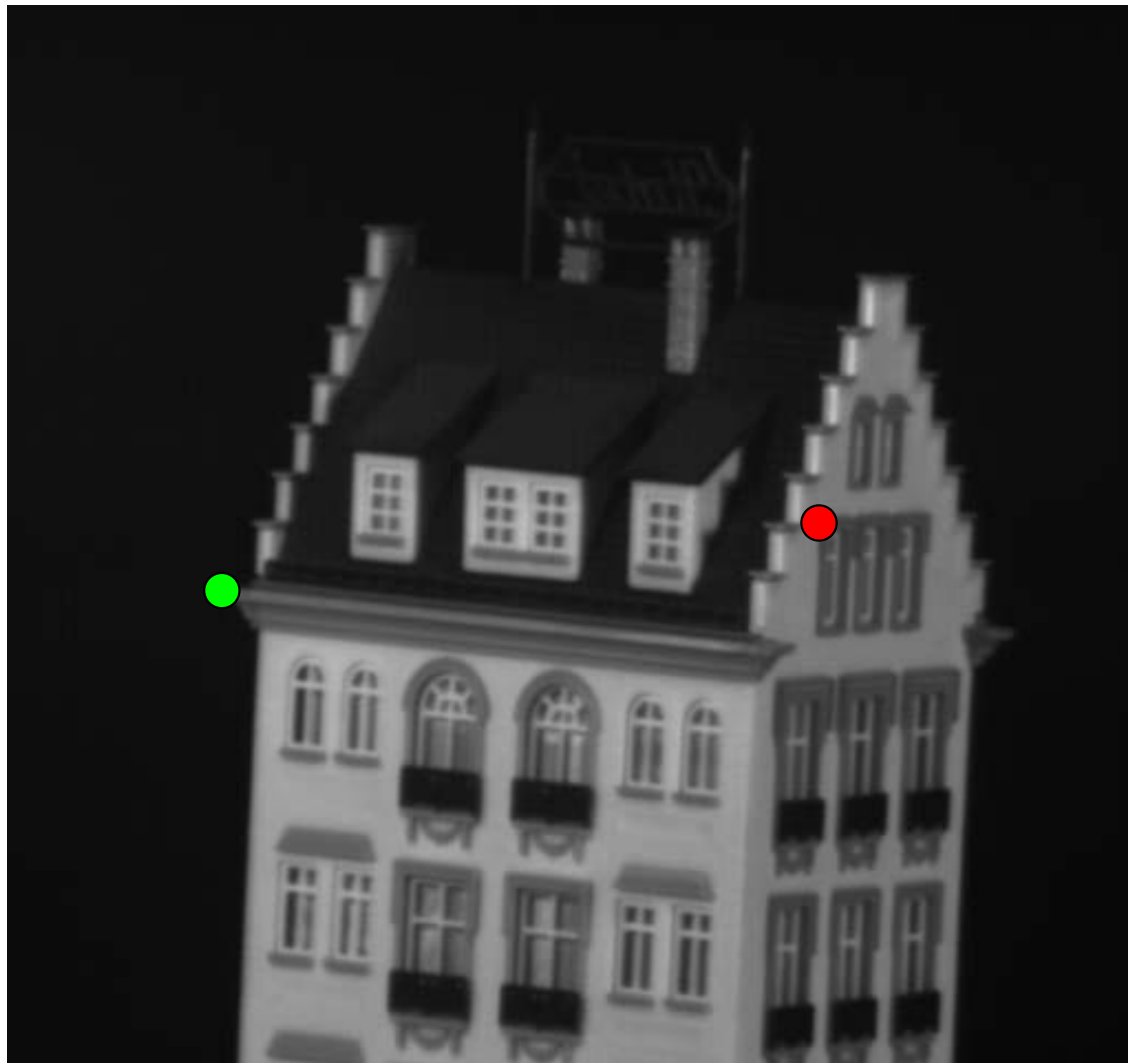
Under the assumption of spatial coherence:

- Hough Transform on the motion vectors.
- Regularization of the motion field.
- Multi scale approach.

But, the world is neither Lambertian nor smooth → Assumptions rarely valid.

TRACKING POINTS ACROSS IMAGES

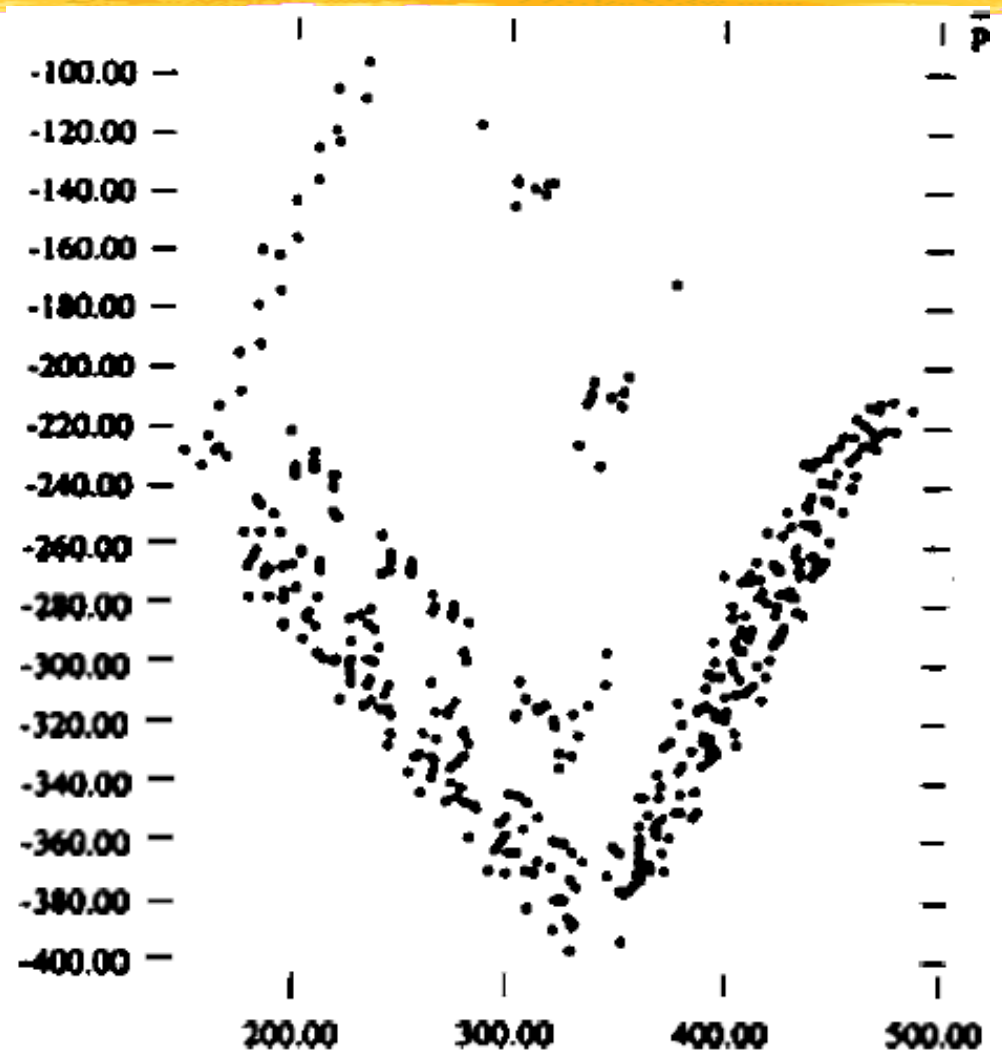
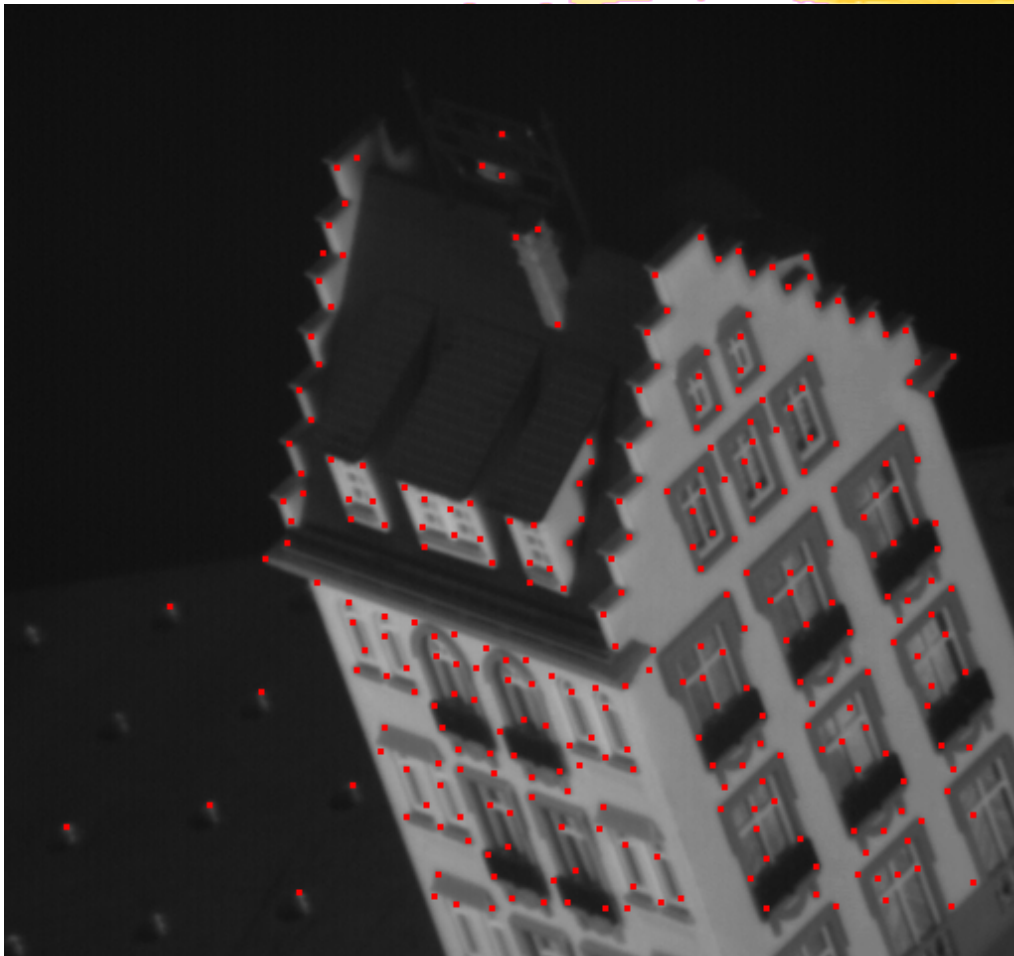
x_1^3



x_2^3



SHAPE RECONSTRUCTION



Factoring Image Sequences into Shape and Motion, C. Tomasi and T. Kanade, Proc. IEEE Workshop on Visual Motion (1991).

MULTI-VIEW PROJECTION

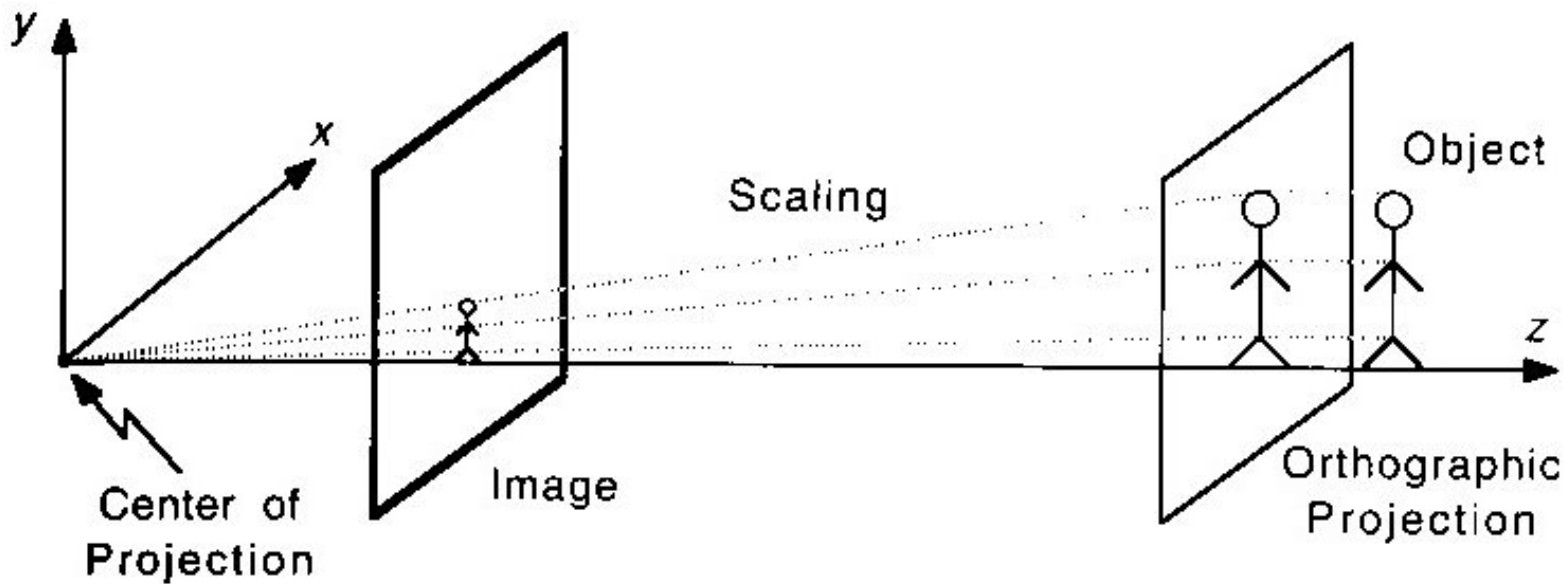


n image points are projected from 3-D scene points over m views via

$$\mathbf{x}_j^i = \mathbf{P}^i \mathbf{X}_j$$

where $i = 1, \dots, m$ and $j = 1, \dots, n$. Here each \mathbf{P}^i is a 3×4 matrix and each \mathbf{X}_j is a homogeneous 4-vector

ORTHOGRAPHIC PROJECTION



$$u = sx$$

$$v = sy$$

Special case of perspective projection:

- Large f
 - Objects close to the optical axis
- Parallel lines mapped into parallel lines.

MULTI-VIEW ORTHOGRAPHIC PROJECTION



The last row of each \mathbf{P}^i is $(0, 0, 0, 1)$ for affine cameras, so we can “ignore” it and write the orthographic projection as:

$$\mathbf{x}_j^i = \mathbf{M}^i \mathbf{X}_j + \mathbf{t}^i$$

where each \mathbf{X}_j is now an inhomogeneous 3-vector, each \mathbf{M}^i a 2×3 matrix, and each \mathbf{t}^i a 2-vector.

RECONSTRUCTION PROBLEM



Estimate affine cameras \mathbf{M}^i , translations \mathbf{t}^i , and 3-D points \mathbf{X}_j that minimize the geometric error in image coordinates:

$$\min_{\mathbf{M}^i, \mathbf{t}^i, \mathbf{X}_j} \sum_{i,j} \left(\mathbf{x}_j^i - (\mathbf{M}^i \mathbf{X}_j + \mathbf{t}^i) \right)^2$$

SIMPLIFYING THE PROBLEM



Normalization: We can eliminate the translation vectors \mathbf{t}^i by choosing the centroid of the image points in each image as the coordinate system origin

$$\mathbf{x}_j^i \leftarrow \mathbf{x}_j^i - \frac{1}{n} \sum_j \mathbf{x}_j^i$$

Working in “centered coordinates”, the minimization problem becomes:

$$\min_{\mathbf{M}^i, \mathbf{X}_j} \sum_{i,j} \left(\mathbf{x}_j^i - \mathbf{M}^i \mathbf{X}_j \right)^2$$

This works because the centroid of the 3-D points is preserved under affine transformations

MATRIX FORMULATION

Let the measurement matrix be:

$$W = \begin{pmatrix} x_1^1 & x_2^1 & \dots & x_n^1 \\ x_1^2 & x_2^2 & \dots & x_n^2 \\ \vdots & \vdots & \ddots & \vdots \\ x_1^m & x_2^m & \dots & x_n^m \end{pmatrix}$$

Since $x_j^i = M^i X_j$, minimizing means solving:

$$W = \begin{bmatrix} M^1 \\ \vdots \\ M^m \end{bmatrix} [X_1, \dots, X_n]$$

$2m \times 3$ $3 \times n$

SOLVING WITH SVD



There will be no exact solution with noisy points, so we want the nearest \mathbf{W}' to \mathbf{W} that is an exact solution

\mathbf{W}' is rank 3 since it's the product of a $2m \times 3$ motion matrix \mathbf{M}' and a $3 \times n$ structure matrix \mathbf{X}'

Use singular value decomposition to get rank 3 matrix \mathbf{W}' closest to \mathbf{W}

Let SVD of $\mathbf{W} = \mathbf{U}\mathbf{D}\mathbf{V}^T$

Then $\mathbf{W}' = \mathbf{U}_{2m \times 3} \mathbf{D}_{3 \times 3} \mathbf{V}_{n \times 3}^T$, where $\mathbf{U}_{2m \times 3}$ is the first 3 columns of \mathbf{U} , $\mathbf{D}_{3 \times 3}$ is an upper-left 3×3 submatrix of \mathbf{D} , and $\mathbf{V}_{n \times 3}^T$ is first three columns of \mathbf{V} .

STRUCTURE AND MOTION



Set stacked camera matrix as

$$\mathbf{M}' = \mathbf{U}_{2m \times 3} \text{sqrt}(\mathbf{D}_{3 \times 3})$$

and stacked 3-D structure matrix as

$$\mathbf{X}' = \text{sqrt}(\mathbf{D}_{3 \times 3}) \mathbf{V}_{n \times 3}^T$$

so that $\mathbf{W}' = \mathbf{M}' \mathbf{X}'$

METRIC UPGRADE



There is an affine ambiguity since an arbitrary 3 x 3 rank 3 matrix **A** can be inserted as:

$$\mathbf{W}' = (\mathbf{M}'\mathbf{A})(\mathbf{A}^{-1}\mathbf{X}')$$

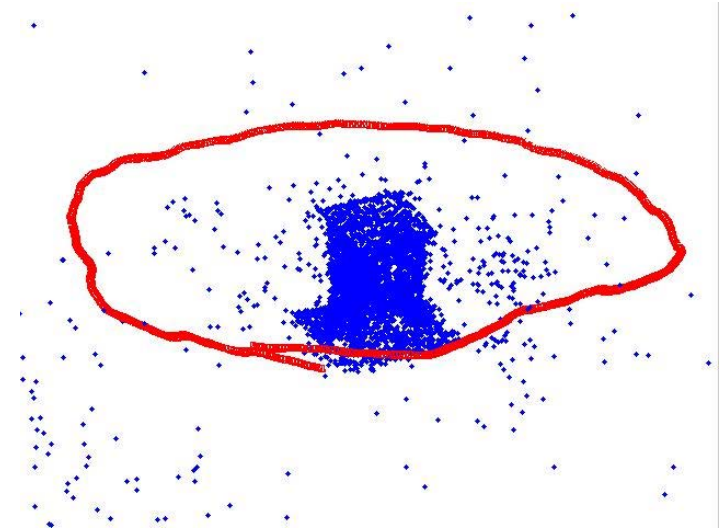
Get rid of ambiguity by finding **A** that performs “metric rectification”

Affine camera provides orthonormality constraints on **A**:

Rows of $\mathbf{M}=\mathbf{M}'\mathbf{A}$ are unit vectors: $\mathbf{m}_i \cdot \mathbf{m}_i = 1$.

Rows of $\mathbf{M}=\mathbf{M}'\mathbf{A}$ are orthogonal: $\mathbf{m}_i \cdot \mathbf{m}_j = 0$.

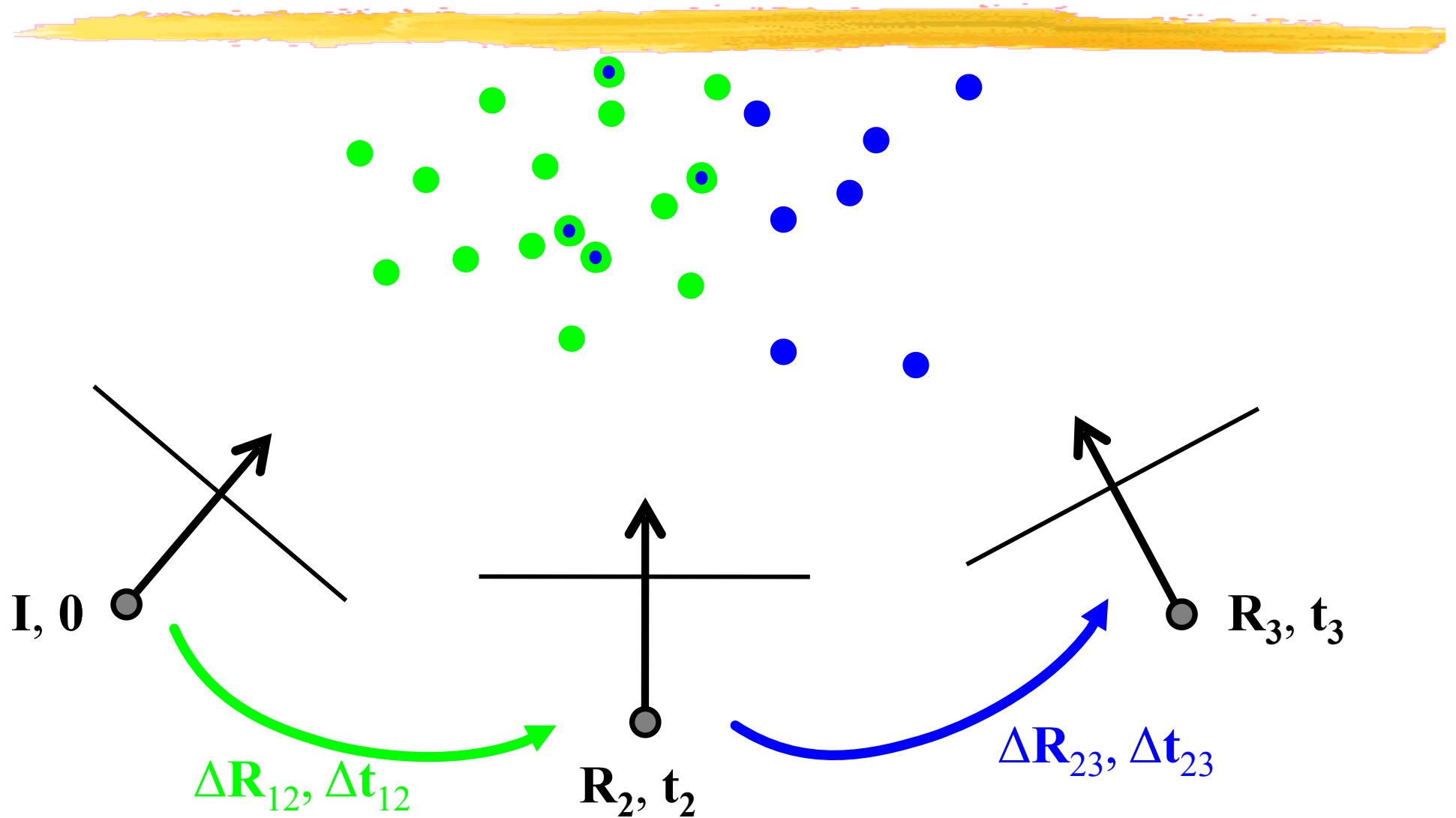
SIMULTANEOUS LOCALIZATION AND MAPPING



Steadly et al., ICCV'03

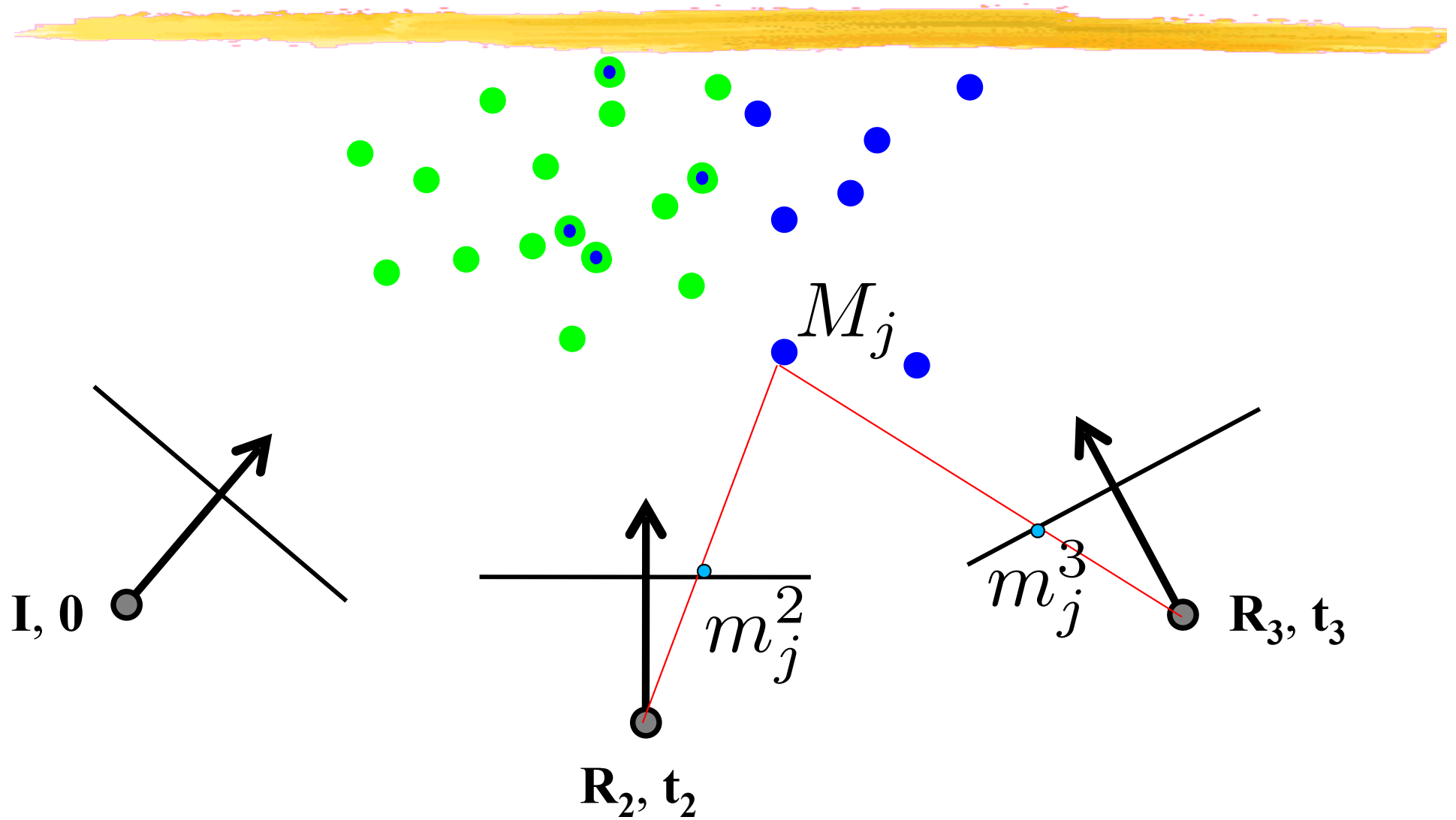
- Compute point tracks.
- Infer both camera motion and 3D structure.

SEQUENTIAL STRUCTURE FROM MOTION



-> Trajectory and 3D points defined up to a Euclidean motion and scale

BUNDLE ADJUSTMENT



$$\operatorname{argmin}_{R_i, t_i, M_j} \sum_i \sum_j \|\operatorname{proj}(R_i, t_i, M_j) - m_j^i\|^2$$

GLOBAL OPTIMIZATION



$$\operatorname{argmin}_{R_i, t_i, M_j} \sum_i \sum_j \|\operatorname{proj}(R_i, t_i, M_j) - m_j^i\|^2$$

- Often performed using the Levenberg-Marquardt algorithm.
- Many parameters to estimate, but sparse Jacobian matrix.
- Initial estimates computed using the eight point algorithm:
 - Given 8 point correspondences between a pair of images, ΔR and ΔT can be estimated in closed form by solving an SVD.

AUGMENTED REALITY

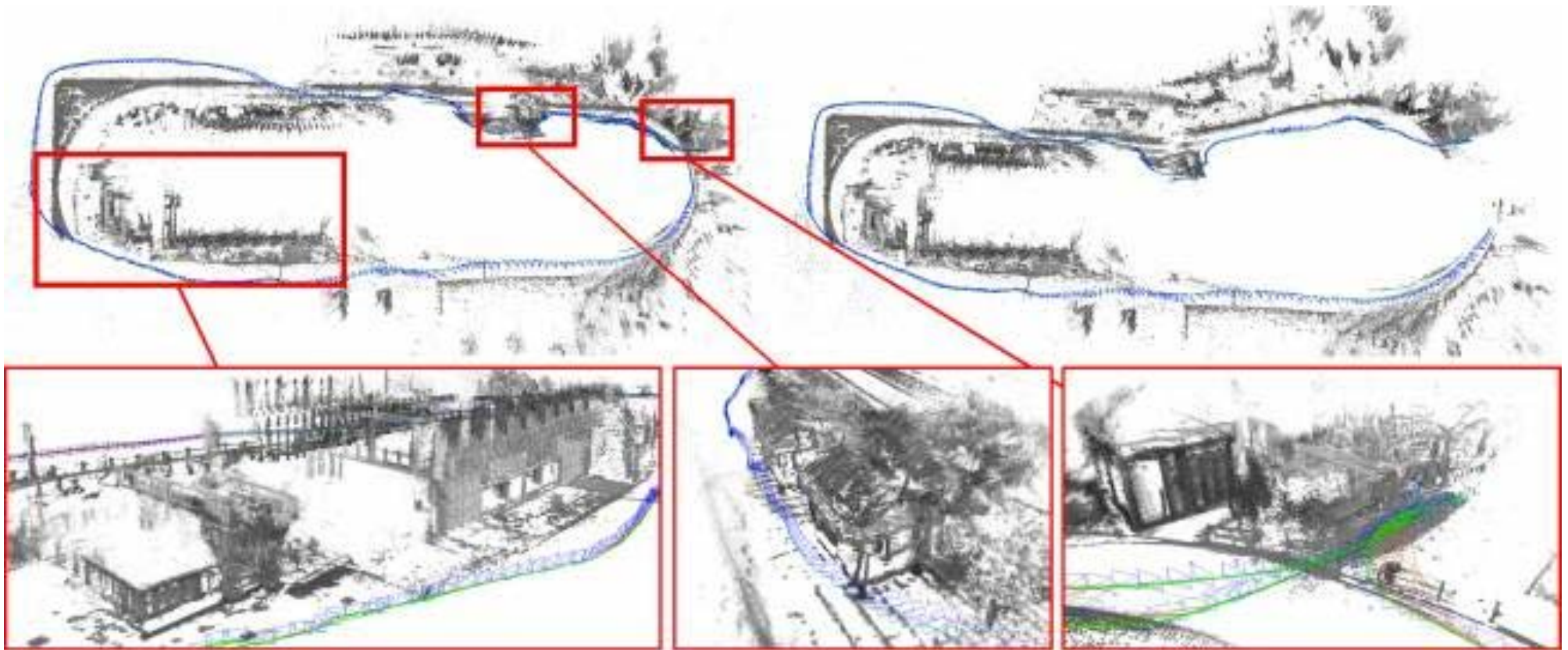


Parallel Tracking and Mapping
for Small AR Workspaces

Extra video results made for
ISMAR 2007 conference

Georg Klein and David Murray
Active Vision Laboratory
University of Oxford

SIMULTANEOUS LOCALIZATION AND MAPPING



A robot can reconstruct its environment and position itself at the same time.

FUSING DEPTH MAPS



- Both the depth camera and the person are moving.
- Use a deformable model to combine the data over time.
- Real-time implementation.

INTO THE COMMERCIAL WORLD

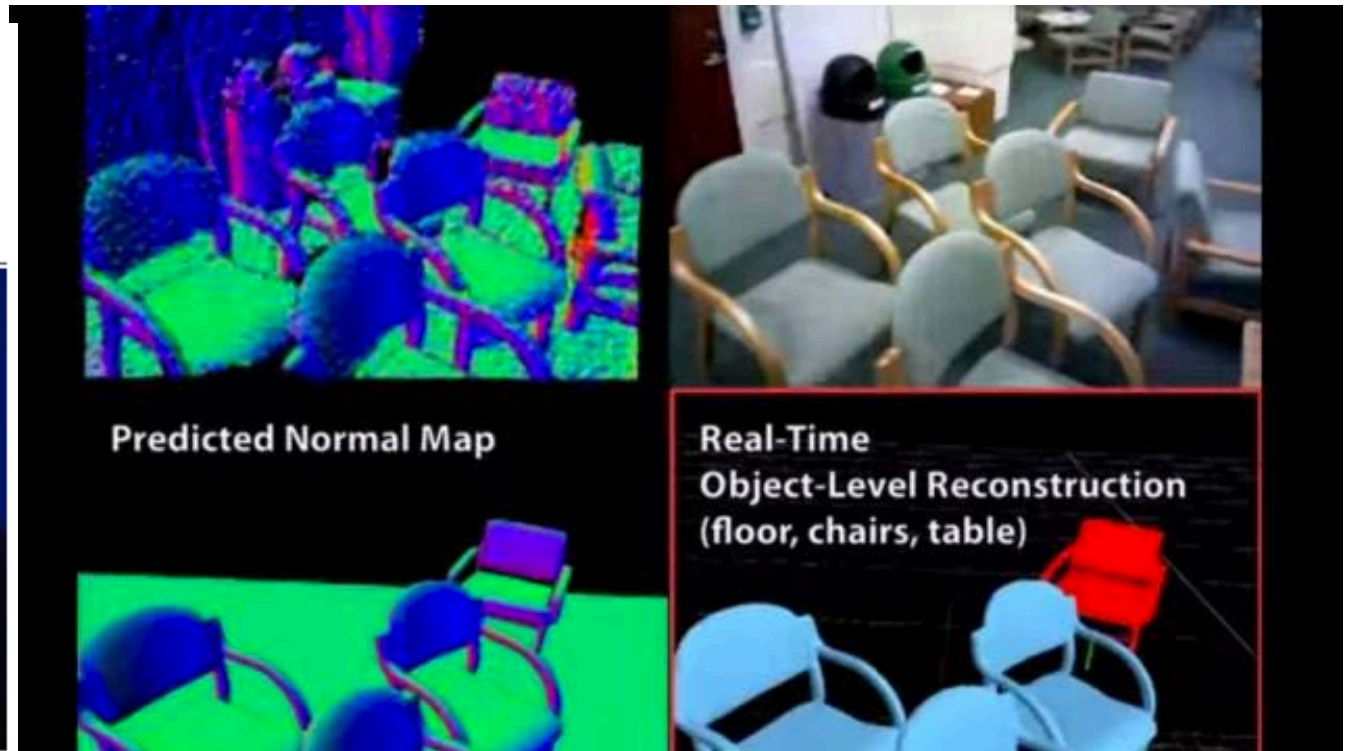
Facebook buys British virtual reality start-up Surreal Vision

Surreal Vision aims to make a computerised version of the world so real that users are unable to distinguish between the two

 614  156  0  27  797  Email



Oculus Rift is expected to be launched next year Photo: AFP



STRENGTHS AND LIMITATIONS



Strengths:

- Combine information from many images.

Limitations:

- Requires multiple views.
- Requires texture or depth camera.