

# Visual Servoing

---

- Visual measures give geometrical and qualitative information on the environment without physical interaction
  - Employed for high-level control
    - e.g. task planning
  - Employed for low-level control
    - e.g. feedback control
- Look-and-move approach (e.g. visual grasping)
  - Trajectory planned on the visual information
  - Simple motion controller
  - Open loop control → very sensitive to the uncertainties
- Vision-based control approach (also termed visual servoing approach)
  - Visual measurements used in the control loop → robust!
  - Error between the current estimated pose and the end-effector pose
- Camera provides one or more 2D matrices of values of light intensity
  - No direct measures!
  - → Extraction of image feature parameters
  - → Pose estimation (geometric relationships)
  - → Calibration (intrinsic and extrinsic parameters)

- Visual-based control schemes can be divided in the following categories:
  - Position-based visual servoing
  - Image-based visual servoing
  - Hybrid visual servoing

## *Configuration of the visual system*

### *Number of cameras*

- 3D vision or stereo vision
  - More cameras are used to observe the same object
  - Depth can be directly measured
    - Triangulation (e.g. human sight)
- Monocular vision
  - Only one camera is used to observe the object
  - Depth can be estimated
    - The same object is observed by the same camera from different angles
    - Estimation on the basis of certain geometrical characteristics of the object known in advance

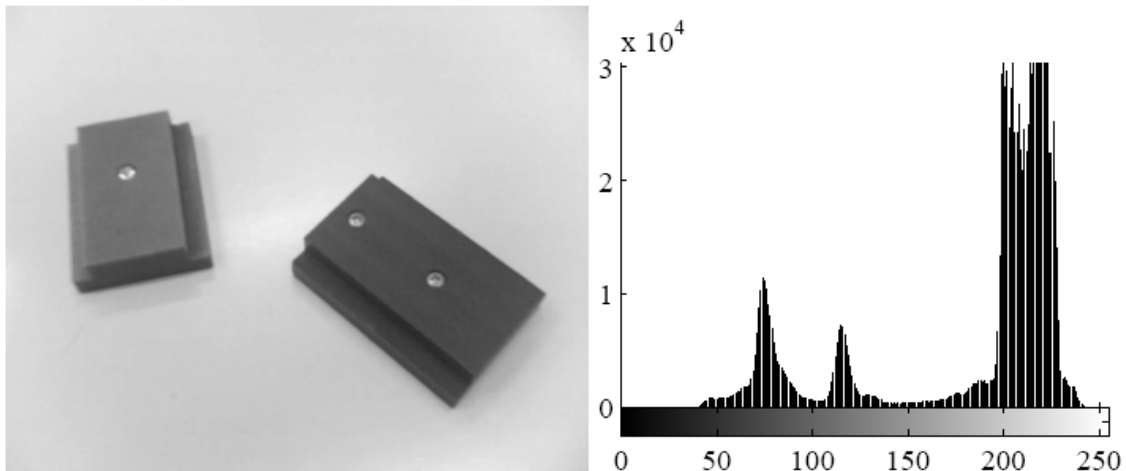
### *Placement of cameras*

- Fixed configuration (eye-to-hand)
  - Field of view does not change during the execution of the task
  - Problems about vision occlusions
- Mobile configuration (eye-in-hand)
  - Field of view changes during task execution
  - Occlusions are virtually absent
- Hybrid configuration
- Robotics heads
  - Typically, two cameras mounted on motorized mechanism (pan-tilt cameras)

## Image Processing

- Visual information is very rich and varied
  - Complex and computational expensive elaboration
  - Extraction of numerical synthetic information
    - Image feature parameters
    - Image segmentation and interpretation
  
- Image function
  - Vector function:  $I(X, Y)$
  - Defined on the set of pixels of the stored image
  - It represents physical quantities related to the pixel in a sampled and quantized form
    - Light intensity (red, green and blue or shades of gray)
  
- Color images
  - The image function has three components corresponding to the light intensity in the wavelengths of red, green and blue
  
- Monochrome black-and-white image
  - The image function is scalar and coincides with the light intensity in shades of gray (gray-level)
  - The number of gray levels depends on the adopted grey-scale resolution
    - e.g. 256 levels → 1 byte of memory

- Gray-level histogram
  - Synthetic representation of the monochrome image
  - It provides the frequency of occurrence of each gray level in the image
  - Normalized histogram (frequencies divided by the total number of pixels)

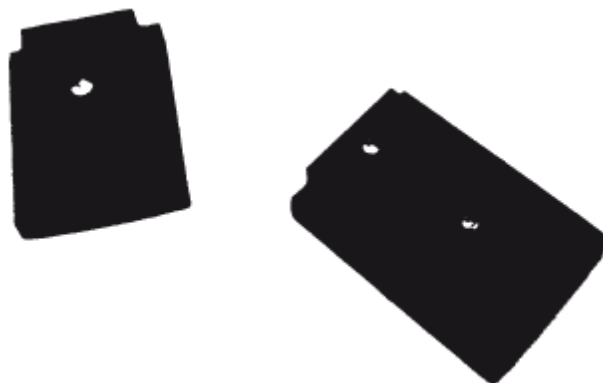


## *Image segmentation*

- Goal: division of the image in groups (referred to as “segments”)
    - Components of each group are similar with respect to one or more characteristics
    - Distant segments correspond to distinct objects or homogenous object parts
  
  - Segmentation approaches
    - Region-based segmentation
      - Grouping set of pixels sharing common features into 2D connected areas
      - High memory usage
      - Low computational load
    - Boundary-based segmentation
      - Identifying the pixels corresponding to object contours and isolating them
      - Low memory usage
      - High computational load
- Complementary approaches

### *Region-based segmentation*

- Connected regions obtained by continuous merging of initially small groups of adjacent pixels into larger ones
  - Check of the “uniformity predicate” for adjacent regions
    - e.g. Check on the gray level
- Binarization (or binary segmentation)
  - Binary light intensity scale
    - Binary image. Only two possible values: 0 black, 1 white
    - Obtained by comparing the gray level of each pixel with a threshold
    - Separation of objects from the background → labeling
  - Selection of the threshold
    - Use of the gray-level histogram
    - Threshold corresponds to the minimum among two peaks (named *modes*), one of which corresponding to the background
    - Appropriate filtering of noisy images



*Boundary-based segmentation:*

- Boundary obtained by grouping many single “local” edges, corresponding to local discontinuities of image gray level
  - Derivation of an intermediate image based on local edges
  - Construction of short-curve segments by edge linking
  - Achievement of the boundaries by joining short-curve segments through geometric primitive known in advance
  - Effectiveness depends on reliability and authenticity of local edges
  
- Local edges detection
  - Filtering process (e.g. convolution)
    - Implemented via hardware
  - Computation of the spatial gradient of the image function
    - Information about rate of change of the gray level and about the direction of the maximum variation of the gray level
    - Different solutions for:
      - Direction for derivatives computation
      - Operators for the derivatives approximation
        - First differences, Roberts, Sobel, ...
    - Threshold to detect maximum (threshold sensitivity!)
  - Computation of the image function Laplacian
    - Second derivatives of the image function along two orthogonal directions (discrete operators)
    - Threshold to detect minimum
    - Sensitive to noise
    - Not provide directional information



## Image interpretation

- Goal: computation of the image feature parameters
  - It works on segmented image (in terms of boundaries or in terms of regions)
- Feature parameters often used in visual servoing: *Moments*
  - Defined on a region of the image
  - Can be used to characterize easily position, orientation and shape of the 2D object correspondent to the region of definition
  - General definition of *moment* of  $i, j$  order

$$m_{i,j} = \sum_{X_I, Y_I \in \mathcal{R}} I(X_I, Y_I) X_I^i Y_I^j$$

- In the case of binary image:  $I(X_I, Y_I) = 1$  for  $X_I, Y_I \in \mathcal{R}$
- $m_{0,0}$  coincides with the area (“mass”) of the region
- *Centroid (position)* coordinates of the region

$$\bar{x} = \frac{m_{1,0}}{m_{0,0}}, \quad \bar{y} = \frac{m_{0,1}}{m_{0,0}}$$

correspondent to the center of  $\mathcal{R}$

- General definition of the *central moments* of  $i, j$  order

$$\mu_{i,j} = \sum_{X_I, Y_I \in \mathcal{R}} (X_I - \bar{x})^i (Y_I - \bar{y})^j$$

that are invariant with respect to translation

- According to a mechanical analogy,  $\mu_{2,0}$  e  $\mu_{0,2}$  represent the inertia moments with respect to  $X$  and  $Y$  axes, while  $\mu_{1,1}$  is an inertia product, then the matrix

$$I = \begin{bmatrix} \mu_{2,0} & \mu_{1,1} \\ \mu_{1,1} & \mu_{0,2} \end{bmatrix}$$

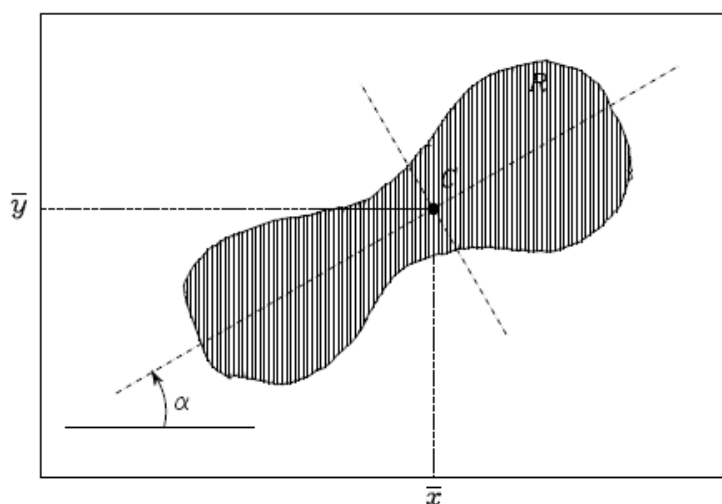
has the meaning of inertia tensor relative to the center of mass

- The eigenvalues of the matrix  $I$  define the *principal moments of inertia*, while the corresponding eigenvectors define the *principal axes of inertia* of the region
- If the region is asymmetric, the angle between the principal axis corresponding to the maximum moment and the axis  $X$

$$\alpha = \frac{1}{2} \tan^{-1} \left( \frac{2\mu_{1,1}}{\mu_{2,0} - \mu_{0,2}} \right)$$

can be used to characterize the *orientation* of the region

- The moments and the corresponding parameters can be computed also from the boundaries



- Feature points
  - Relevant points of the object (e.g. corners)
- Geometric primitives
  - Lines and lines segment, which are projections of linear edges or solids of revolution
  - Ellipses, obtained as projections of circles of spheres
  - These primitives can be characterized on the image plane in terms of a minimum set of parameters

## Pose Estimation

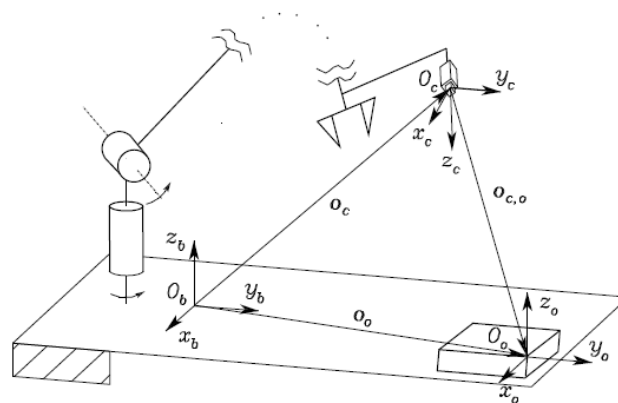
- For visual servoing purposes, a mapping between the feature parameters of an object measured in the image plane and the operational space variables is needed, defining the relative pose of the object with respect to the camera (or vice versa)
- Differential mapping in terms of velocity
  - Linear mapping!
  - Numerical integration algorithms can be used
- Feature parameters vector
  - Set of feature parameters of an image
  - Vector  $s$  ( $k \times 1$ ), termed *feature vector*
  - Normalized coordinates will be used
    - Intrinsic parameters of the camera are needed!
  - The feature vector of a point is defined as:

$$s = \begin{bmatrix} X \\ Y \end{bmatrix} \quad \tilde{s} = \begin{bmatrix} X \\ Y \\ 1 \end{bmatrix}$$

- Object pose with respect to the camera is defined as

$$T_o^c = \begin{bmatrix} R_o^c & o_{c,o}^c \\ \mathbf{0}^T & 1 \end{bmatrix}$$

with  $o_{c,o}^c = o_o^c - o_c^c$



## Analytic solution

- Problem: computation of elements of matrix  $\mathbf{T}_o^c$  from the feature vector  $\mathbf{s}$  containing the projections of  $n$  points of the object (e.g. eye-in-hand camera)

$$\mathbf{s}_i = \begin{bmatrix} X_i \\ Y_i \end{bmatrix} \quad \mathbf{s} = \begin{bmatrix} s_1 \\ \vdots \\ s_n \end{bmatrix}$$

Knowing the object model, the homogenous coordinates of the points of the object with respect to the camera frame can be expressed as

$$\mathbf{r}_{o,i}^o = \mathbf{p}_i^o - \mathbf{o}_o^o \quad \tilde{\mathbf{r}}_{o,i}^c = \mathbf{T}_o^c \tilde{\mathbf{r}}_{o,i}^o$$

and using the *perspective transformation* yields

$$\lambda_i \tilde{\mathbf{s}}_i = \mathbf{\Pi} \mathbf{T}_o^c \tilde{\mathbf{r}}_{o,i}^o$$

with  $\lambda_i > 0$  and

$$\mathbf{\Pi} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

- $n$  correspondences  $\rightarrow$  PnP Problem (Perspective- $n$ -Points)
  - System of  $n$  equations for the unknown elements of  $\mathbf{T}_o^c$
  - Multiple solutions may exist!
- P3P problem has four solutions, in the case of three non-collinear points
- P4P and P5P problems each have:
  - at least two solutions, in the case of non-coplanar points
  - unique solution, in the case of at least four coplanar points and no triplets of collinear points

- PnP problem, with  $n \geq 6$  non-coplanar points has only one solution
- In the particular case of coplanar points, the derivation of the solution to the PnP problem becomes simpler
  - Assuming that the points lie in the plane of equation  $z_o = 0$
  - Multiplying both sides of equations system by the skew-symmetric matrix  $\mathbf{S}(\tilde{\mathbf{s}}_i)$ , the following homogeneous equation is obtained

$$\mathbf{S}(\tilde{\mathbf{s}}_i)\mathbf{H} [r_{x,i} \quad r_{y,i} \quad 1]^T = \mathbf{0}$$

where  $r_{x,i}$  and  $r_{y,i}$  are the two non-null components of vector  $\mathbf{r}_{o,i}^o$ , and where  $\mathbf{H}$  is known as **planar homography** that is a (3x3) matrix

$$\mathbf{H} = [\mathbf{r}_1 \quad \mathbf{r}_2 \quad \mathbf{o}_{c,o}^c]$$

with  $\mathbf{r}_1$  and  $\mathbf{r}_2$  the first and second column of rotation matrix  $\mathbf{R}_o^c$  respectively

- Notice that the planar homography is linear with respect to  $\mathbf{H}$ , and then it can be rewritten in the form

$$\mathbf{A}_i(\mathbf{s}_i)\mathbf{h} = \mathbf{0}$$

where  $\mathbf{h}$  is the (9x1) column vector obtained by stacking the columns of the matrix  $\mathbf{H}$ , while  $\mathbf{A}_i$  is the (3x9) matrix

$$\mathbf{A}_i(\mathbf{s}_i) = [r_{x,i}\mathbf{S}(\tilde{\mathbf{s}}_i) \quad r_{y,i}\mathbf{S}(\tilde{\mathbf{s}}_i) \quad \mathbf{S}(\tilde{\mathbf{s}}_i)]$$

- Since the rank of  $\mathbf{S}(\tilde{\mathbf{s}}_i)$  is at most 2, then the rank of matrix  $\mathbf{A}_i$  is also at most 2; therefore it is necessary to consider at least 4 equations written for 4 points of the plane, leading to the following system of 12 equations with 9 unknowns

$$\begin{bmatrix} A_1(s_1) \\ A_2(s_2) \\ A_3(s_3) \\ A_4(s_4) \end{bmatrix} \mathbf{h} = \mathbf{A}(s)\mathbf{h} = \mathbf{0}$$

- It can be shown that, considering four points with no triplets of collinear points, matrix  $\mathbf{A}$  has rank 8 and the system admits a non-null solution  $\zeta\mathbf{h}$ , defined up to a scaling factor  $\zeta$ , from which it can be obtained

$$\begin{aligned} \mathbf{r}_1 &= \zeta\mathbf{h}_1 \\ \mathbf{r}_2 &= \zeta\mathbf{h}_2 \\ \mathbf{o}_{c,o}^c &= \zeta\mathbf{h}_3 \end{aligned}$$

where  $\mathbf{h}_i$  denotes the  $i$ -th column of matrix  $\mathbf{H}$

- The value of the constant  $\zeta$  can be computed by imposing the unit norm constraint to vectors  $\mathbf{r}_1$  and  $\mathbf{r}_2$

$$|\zeta| = \frac{1}{\|\mathbf{h}_1\|} = \frac{1}{\|\mathbf{h}_2\|}$$

while the sign can be determined by choosing the solution corresponding to the object in front of the camera. Finally, the vector  $\mathbf{r}_3$  can be computed as  $\mathbf{r}_3 = \mathbf{r}_1 \times \mathbf{r}_2$

- In a real situation, to reduce *noise* effects, a number  $n > 4$  of correspondences can be considered and the solution can be computed using a *least-squares* technique
  - In general, the resulting matrix  $\mathbf{Q} = [\mathbf{r}_1 \quad \mathbf{r}_2 \quad \mathbf{r}_3]$  does not guarantee the property of being a rotation matrix
  - A possible solution consists of computing the rotation matrix “closest” to  $\mathbf{Q}$  with respect to a given norm (e.g. Frobenius norm)

$$\|R_o^c - Q\|_F = \text{Tr} \left( (R_o^c - Q)^T (R_o^c - Q) \right)$$

with the constraint that  $R_o^c$  is a rotation matrix

- The problem of minimizing the Frobenius norm is equivalent to that of maximizing the trace of matrix  $R_o^c{}^T Q$ , and the problem has the following solution

$$R_o^c = U \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & \sigma \end{bmatrix} V^T$$

where  $U$  and  $V^T$  are, respectively, the left and the right orthogonal matrices of the singular value decomposition of  $Q = U\Sigma V^T$  and  $\sigma = \det(UV^T)$  ensures that the determinant of  $R_o^c$  is equal to 1

- The above derivation is a particular case of a method termed *direct linear transformation*
  - It allows the computation of the elements of  $T_o^c$  by solving a system of linear equations obtained using  $n$  correspondences relative to points in generic configuration
  - The system is obtained from the equalities

$$S(\tilde{s}_i) [R_o^c \quad o_{c,o}^c] \tilde{r}_{o,i}^o = \mathbf{0}$$

from which two independent linear equations with 12 unknowns are obtained ( $n$  correspondences  $\rightarrow 2n$  equations)

- It can be shown that, by considering a set of 6 points not all coplanar, the corresponding system of 12 equations with 12 unknowns has rank 11 and it admits a solution defined up to a scaling factor



- In practical applications, due to the presence of noise, the system of equations has rank 12 and admits only the null solution
  - It is necessary to consider  $n > 6$  of correspondences and a least-square technique
- In conclusion, it is necessary to know:
  - Object geometry (positions of the points with respect to the object frame)
  - Camera intrinsic parameters, to compute normalized coordinates from pixel ones
  - Camera extrinsic parameters, to compute the object pose with respect to the base frame, both with fixed and mobile cameras:
    - Camera eye-to-hand:  $\mathbf{T}_o^b = \mathbf{T}_c^b \mathbf{T}_o^c$
    - Camera eye-in-hand:  $\mathbf{T}_o^e = \mathbf{T}_c^e \mathbf{T}_o^c$

## Interaction matrix

- If the object is in motion with respect to the camera, the feature vector  $\mathbf{s}$  is time-varying; hence it is possible to define a ( $k \times 1$ ) velocity vector in the image plane  $\dot{\mathbf{s}}$
- The motion of the object with respect to the camera is characterized by the relative velocity

$$\mathbf{v}_{c,o}^c = \begin{bmatrix} \dot{\mathbf{o}}_{c,o}^c \\ \mathbf{R}_c^T(\boldsymbol{\omega}_o - \boldsymbol{\omega}_c) \end{bmatrix}$$

where  $\dot{\mathbf{o}}_{c,o}^c$  is the time derivative of vector  $\mathbf{o}_{c,o}^c = \mathbf{R}_c^T(\mathbf{o}_o - \mathbf{o}_c)$

- The equation relating these two velocities is

$$\dot{\mathbf{s}} = \mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) \mathbf{v}_{c,o}^c$$

where  $\mathbf{J}_s$  is a ( $k \times 6$ ) matrix, termed *image Jacobian*

- It is useful to consider the mapping between the image plane velocity  $\dot{\mathbf{s}}$ , the absolute velocity of the camera frame

$$\mathbf{v}_c^c = \begin{bmatrix} \mathbf{R}_c^T \dot{\mathbf{o}}_c \\ \mathbf{R}_c^T \boldsymbol{\omega}_c \end{bmatrix}$$

and the absolute velocity of the object frame

$$\mathbf{v}_o^c = \begin{bmatrix} \mathbf{R}_c^T \dot{\mathbf{o}}_o \\ \mathbf{R}_c^T \boldsymbol{\omega}_o \end{bmatrix}$$

- The vector  $\dot{\mathbf{o}}_{c,o}^c$  can be expressed as

$$\dot{\mathbf{o}}_{c,o}^c = \mathbf{R}_c^T(\dot{\mathbf{o}}_o - \dot{\mathbf{o}}_c) + \mathbf{S}(\mathbf{o}_{c,o}^c) \mathbf{R}_c^T \boldsymbol{\omega}_c$$

from which

$$\mathbf{v}_{c,o}^c = \mathbf{v}_o^c + \boldsymbol{\Gamma}(\mathbf{o}_{c,o}^c) \mathbf{v}_c^c \quad \boldsymbol{\Gamma}(\cdot) = \begin{bmatrix} -\mathbf{I} & \mathbf{S}(\cdot) \\ \mathbf{O} & -\mathbf{I} \end{bmatrix}$$

and the differential mapping can be rewritten as

$$\dot{\mathbf{s}} = \mathbf{J}_s \mathbf{v}_o^c + \mathbf{L}_s \mathbf{v}_c^c$$

where the ( $k \times 6$ ) matrix

$$\mathbf{L}_s = \mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) \boldsymbol{\Gamma}(\mathbf{o}_{c,o}^c)$$

is termed *interaction matrix*

- The interaction matrix defines the linear mapping between the absolute velocity of the camera and the corresponding image plane velocity, in the case the object is fixed with respect to the base frame ( $\mathbf{v}_o^c = \mathbf{0}$ )
- The analytic expression of the interaction matrix is, in general, easier than that of image Jacobian, so that the latter can be computed from the interaction matrix using the following equation

$$\mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) = \mathbf{L}_s \Gamma(-\mathbf{o}_{c,o}^c)$$

obtained from the inversion of the previous equation defining the interaction matrix

## Interaction matrix of a point

- Consider a point  $P$  of the object characterized by the vector

$$\mathbf{r}_c^c = \mathbf{R}_c^T(\mathbf{p} - \mathbf{o}_c)$$

- As the feature vector  $\mathbf{s}$  of the point, choose the vector of normalized coordinates

$$\mathbf{s} = \mathbf{s}(\mathbf{r}_c^c) = \frac{1}{z_c} \begin{bmatrix} x_c \\ y_c \end{bmatrix} \quad \mathbf{r}_c^c = [x_c \quad y_c \quad z_c]^T$$

- Computing the time derivative yields

$$\dot{\mathbf{s}} = \frac{\partial \mathbf{s}(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} \dot{\mathbf{r}}_c^c$$

with

$$\frac{\partial \mathbf{s}(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -x_c/z_c \\ 0 & 1 & -y_c/z_c \end{bmatrix} = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -X \\ 0 & 1 & -Y \end{bmatrix}$$

- Computing the time derivative of the vector  $\mathbf{r}_c^c$ , under the assumption that  $\mathbf{p}$  is constant ( $\dot{\mathbf{v}}_o^c = \mathbf{0}$ )

$$\dot{\mathbf{r}}_c^c = -\mathbf{R}_c^T \dot{\mathbf{o}}_c + \mathbf{S}(\mathbf{r}_c^c) \mathbf{R}_c^T \boldsymbol{\omega}_c = [-\mathbf{I} \quad \mathbf{S}(\mathbf{r}_c^c)] \mathbf{v}_c^c$$

- Combining the previous two equations, the following expression of interaction matrix of a point can be found

$$\mathbf{L}_s(\mathbf{s}, z_c) = \begin{bmatrix} -\frac{1}{z_c} & 0 & \frac{X}{z_c} & XY & -(1+X^2) & Y \\ 0 & -\frac{1}{z_c} & \frac{Y}{z_c} & 1+Y^2 & -XY & -X \end{bmatrix}$$

- The image Jacobian can be computed from the interaction matrix

$$\mathbf{J}_s(\mathbf{s}, \mathbf{T}_o^c) = \frac{1}{z_c} \begin{bmatrix} 1 & 0 & -X & -r_{o,y}^c X & r_{o,z}^c + r_{o,x}^c X & -r_{o,y}^c \\ 0 & 1 & -Y & -(r_{o,z}^c + r_{o,y}^c Y) & r_{o,x}^c Y & r_{o,x}^c \end{bmatrix}$$

where  $r_{o,x}^c$ ,  $r_{o,y}^c$  e  $r_{o,z}^c$  are the components of the vector

$$\mathbf{r}_o^c = \mathbf{r}_c^c - \mathbf{o}_{c,o}^c = \mathbf{R}_o^c \mathbf{r}_o^o$$

## Interaction matrix of a set of points

- The interaction matrix of a set of  $n$  points of the object  $P_1, P_2, \dots, P_n$  can be built by considering the  $(2n \times 1)$  feature vector
- By denoting with  $L_{s_i}(s_i, z_{c,i})$  the interaction matrix corresponding to point  $P_i$ , the interaction matrix of the set of points will be the  $(2n \times 6)$  matrix

$$L_s(s, z_c) = \begin{bmatrix} L_{s_1}(s_1, z_{c,1}) \\ \vdots \\ L_{s_n}(s_n, z_{c,n}) \end{bmatrix}$$

with  $z_c = [z_{c,1}, z_{c,2}, \dots, z_{c,n}]^T$

## Interaction matrix of a line segment

- A line segment is the part of the line connecting two points  $P_1, P_2$
- The projection on the image plane is still a line segment that can be characterized by the middle point coordinates  $\bar{x}, \bar{y}$ , the length  $L$  and the angle  $\alpha$  formed by the line with respect to axis  $X$ :

$$\mathbf{s} = \begin{bmatrix} \bar{x} \\ \bar{y} \\ L \\ \alpha \end{bmatrix} = \begin{bmatrix} (X_1 + X_2)/2 \\ (Y_1 + Y_2)/2 \\ \sqrt{\Delta X^2 + \Delta Y^2} \\ \tan^{-1}(\Delta Y/\Delta X) \end{bmatrix} = \mathbf{s}(s_1, s_2)$$

where  $\Delta X = X_2 - X_1$ ,  $\Delta Y = Y_2 - Y_1$  and  $\mathbf{s}_i = [X_i \ Y_i]^T$ , with  $i = 1, 2$

- Computing the time derivative of the feature vector yields

$$\begin{aligned} \dot{\mathbf{s}} &= \frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \dot{\mathbf{s}}_1 + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \dot{\mathbf{s}}_2 \\ &= \left( \frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \mathbf{L}_{s_1}(s_1, z_{c,1}) + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \mathbf{L}_{s_2}(s_2, z_{c,2}) \right) \mathbf{v}_c^c \end{aligned}$$

where  $\mathbf{L}_{s_i}$  is the interaction matrix of point  $P_i$  (under the assumption that the line segment is fixed with respect to the base frame)

- Then, the interaction matrix of a line segment is

$$\mathbf{L}_s(\mathbf{s}, z_c) = \frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} \mathbf{L}_{s_1}(s_1, z_{c,1}) + \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} \mathbf{L}_{s_2}(s_2, z_{c,2})$$

with

$$\frac{\partial \mathbf{s}}{\partial \mathbf{s}_1} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \\ -\Delta X/L & -\Delta Y/L \\ \Delta Y/L^2 & -\Delta X/L^2 \end{bmatrix}, \quad \frac{\partial \mathbf{s}}{\partial \mathbf{s}_2} = \begin{bmatrix} 1/2 & 0 \\ 0 & 1/2 \\ \Delta X/L & \Delta Y/L \\ -\Delta Y/L^2 & \Delta X/L^2 \end{bmatrix}$$

- The vectors  $\mathbf{s}_1$  and  $\mathbf{s}_2$  can be computed as functions of the parameters  $\bar{x}, \bar{y}, L$  and  $\alpha$ ; therefore the interaction matrix can be expressed as a function of the feature vector  $\mathbf{s}$ , besides the components  $z_{c,1}$  and  $z_{c,2}$  of the endpoints  $P_1$  and  $P_2$  of the line segment

## Algorithmic solution

- In general, the interaction matrix  $L_s$  is a matrix of dimension  $(k \times m)$ 
  - $k$  is the number of feature parameters of the image
  - $m$  is dimension of the velocity vector  $v_c^c$ 
    - $m < 6$  when the relative motion of the object with respect to the camera is constrained
    - $m = 6$  in the other cases
- The image Jacobian  $J_s$  is also of dimension  $(k \times m)$  and its rank coincides with that of the interaction matrix
- In the case that  $L_s$  is full rank, it is possible to compute  $v_{c,o}^c$  from  $\dot{s}$
- If  $k = m$

$$v_{c,o}^c = \Gamma(o_{c,o}^c) L_s^{-1} \dot{s}$$

- If  $k > m$ , the interaction matrix has more rows than columns  $\rightarrow$  least-squares solution

$$v_{c,o}^c = \Gamma(o_{c,o}^c) (L_s^T L_s)^{-1} L_s^T \dot{s}$$

where the left pseudo-inverse of  $L_s$  has been used

- If  $k < m$ , the interaction matrix has more columns than rows  $\rightarrow$  infinite solutions
  - Number of parameters is not sufficient to determine uniquely the relative motion of the object with respect to the camera
  - Relative motions of the object with respect to the camera (or vice versa) do not produce variations of the image feature parameters  $\rightarrow$  velocities belong to the null subspace of  $J_s$
  - This case has no interest if the problem is that of computing uniquely the relative pose of the object with respect to the camera from feature parameters in the image plane

- The pose estimation problem may be cast in a form analogous to that of inverse kinematics algorithms for robot manipulators
- To represent the relative pose of the object with respect to the camera using a minimum number of coordinates, consider the  $(m \times 1)$  vector

$$\mathbf{x}_{c,o} = \begin{bmatrix} o_{c,o}^c \\ \phi_{c,o} \end{bmatrix}$$

- If Euler angles are used to represent orientation, the transformation matrix between  $\mathbf{v}_{c,o}^c$  and  $\dot{\mathbf{x}}_{c,o}$  is

$$\mathbf{v}_{c,o}^c = \begin{bmatrix} \mathbf{I} & \mathbf{O} \\ \mathbf{O} & \mathbf{T}(\phi_{c,o}) \end{bmatrix} \dot{\mathbf{x}}_{c,o} = \mathbf{T}_A(\phi_{c,o}) \dot{\mathbf{x}}_{c,o}$$

- From which the mapping between  $\dot{\mathbf{s}}$  and  $\dot{\mathbf{x}}_{c,o}$  is given by

$$\dot{\mathbf{s}} = \mathbf{J}_{A_s}(\mathbf{s}, \mathbf{x}_{c,o}) \dot{\mathbf{x}}_{c,o}$$

where the matrix

$$\mathbf{J}_{A_s}(\mathbf{s}, \mathbf{x}_{c,o}) = \mathbf{L}_s \mathbf{\Gamma}(-o_{c,o}^c) \mathbf{T}_A(\phi_{c,o})$$

has a meaning analogous to that of the analytic Jacobian of a manipulator

- By denoting with  $\hat{\mathbf{s}} = \mathbf{s}(\hat{\mathbf{x}}_{c,o})$  the current estimate of feature vector  $\mathbf{s}$ , computed from the pose specified by  $\hat{\mathbf{x}}_{c,o}$  (current estimate of vector  $\mathbf{x}_{c,o}$ ), it can be defined an algorithm to minimize the error

$$\mathbf{e}_s = \mathbf{s} - \hat{\mathbf{s}}$$

- Computing the time derivative of the previous equation yields

$$\dot{\mathbf{e}}_s = -\dot{\hat{\mathbf{s}}} = -\mathbf{J}_{A_s}(\hat{\mathbf{s}}, \hat{\mathbf{x}}_{c,o}) \dot{\hat{\mathbf{x}}}_{c,o}$$

- Assumed that matrix  $\mathbf{J}_{A_s}$  is square and nonsingular, the choice

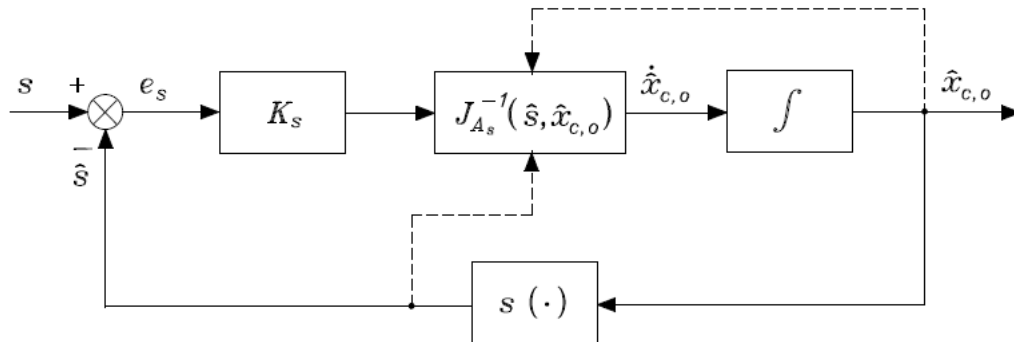
$$\dot{\hat{\mathbf{x}}}_{c,o} = \mathbf{J}_{A_s}^{-1}(\hat{\mathbf{s}}, \hat{\mathbf{x}}_{c,o}) \mathbf{K}_s \mathbf{e}_s$$

leads to the equivalent linear system

$$\dot{\mathbf{e}}_s + \mathbf{K}_s \mathbf{e}_s = \mathbf{0}$$



If  $K_s$  is a positive definite matrix (usually diagonal), the previous system is asymptotically stable and the error tends to zero with a convergence speed that depends on the eigenvalues of the matrix  $K_s$



- The function  $s(\cdot)$  denotes the function computing the feature vector of the “virtual” image corresponding to the current estimate  $\hat{x}_{c,o}$  of the object pose with respect to the camera
- The convergence properties depend on the choice of the image feature parameters, on the initial condition  $\hat{x}_{c,o}(0)$  and on instability problems related to singularities of matrix  $J_{A_s}$
- Singularities of the matrix  $J_{A_s}$ 
  - Representation singularities of the orientation
  - Singularities of the interaction matrix (!)
- To separate the effects of the two types of singularities, it is convenient to realize the computation in two steps
 
$$\hat{v}_{c,o}^c = \Gamma(o_{c,o}^c) L_s^{-1} K_s e_s \qquad \hat{x}_{c,o} = T_A^{-1}(\phi_{c,o}) \hat{v}_{c,o}^c$$
- Under the assumption of working far from representation singularities, the problem of singularities of the interaction matrix can be overcome by using a number  $k$  of feature parameters greater than the minimum required  $m$  ( $k > m$ )
  - Reduction of the effects of measurement noises
  - The resulting estimation algorithm requires the use of the left pseudo-inverse of interaction matrix

$$\hat{v}_{c,o}^c = \Gamma(o_{c,o}^c)(L_s^T L_s)^{-1} L_s^T K_s e_s$$

- The convergence can be shown using the direct Lyapunov method based on the following positive definite function

$$V(e_s) = \frac{1}{2} e_s^T K_s e_s$$

Computing the time derivative follows

$$\dot{V} = -e_s^T K_s L_s (L_s^T L_s)^{-1} L_s^T K_s e_s$$

which is negative semi-definite, because  $\mathcal{N}(L_s^T) \neq \emptyset$  being  $L_s^T$  a matrix with more columns than row  $\rightarrow$  the system is stable, but not asymptotically stable!

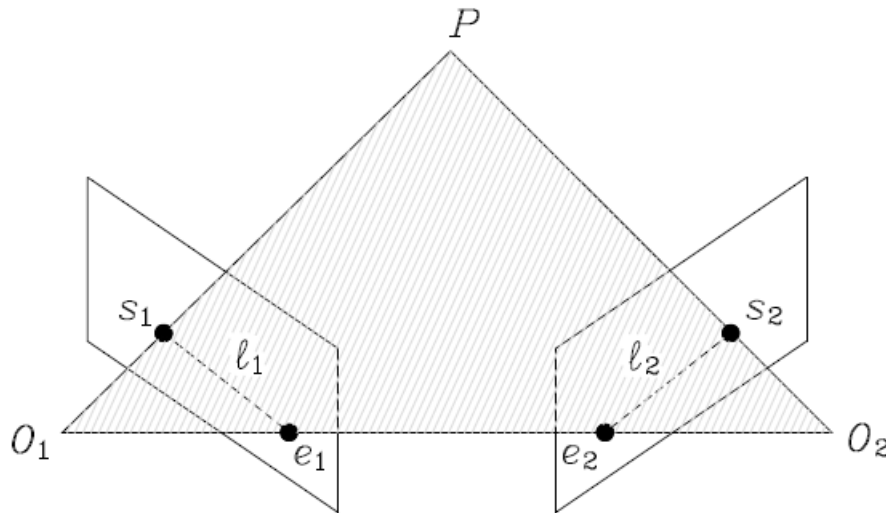
- The algorithm can get stuck with  $e_s \neq 0$  and  $K_s e_s \in \mathcal{N}(L_s^T)$
- These methods for their performance are mainly adopted for real-time “visual-tracking” applications
  - The estimate has as initial value the estimate computed at the previous step time

## Stereo Vision

- 2D Image → Does not give any explicit information on depth
- This information can be recovered:
  - Indirect way from the geometric model of the object, assumed to be known
  - Direct way from two or more images from different points of view (Stereo Vision)
    - 2 or more cameras
    - One moving camera (fixed object!)
- Fundamental problems in the framework of stereo vision
  - Correspondence problem
    - *Definition*: identification of the points of the two images that are projections of the same point of the scene (conjugate or corresponding points)
    - Existence of geometric constraints between two images of the same point
    - Existence of details of the scene which appear to be similar in the different images
  - 3D reconstruction
    - *Definition*: computation of the relative pose of the cameras (calibrated and not) and thus, starting from this pose, the position in the 3D space of the points of the observed object

## Epipolar geometry

- It has the basis on the *epipolar constraint*
  - Consider two cameras, noted as 1 and 2



- The coordinates of a point  $P$  expressed in the two frames are related by the following equation

$$p^1 = o_{1,2}^1 + R_2^1 p^2$$

- Let  $s_1$  and  $s_2$  be the coordinates of the projections of  $P$  on the image planes

$$\lambda_i \tilde{s}_i = \Pi \tilde{p}^i = p^i, \quad i = 1, 2$$

from which, substituting the previous equation in the first one

$$\lambda_1 \tilde{s}_1 = o_{1,2}^1 + \lambda_2 R_2^1 \tilde{s}_2$$

Premultiplying both sides by  $S(o_{1,2}^1)$  gives

$$\lambda_1 S(o_{1,2}^1) \tilde{s}_1 = \lambda_2 S(o_{1,2}^1) R_2^1 \tilde{s}_2$$

Premultiplying both sides by  $\tilde{s}_1^T$

$$\lambda_2 \tilde{s}_1^T S(o_{1,2}^1) R_2^1 \tilde{s}_2 = 0$$

which has to be satisfied for any value of the scalar  $\lambda_2$ ; therefore this equality is equivalent to the so-called *epipolar constraint*

$$\tilde{\mathbf{s}}_1^T \mathbf{E} \tilde{\mathbf{s}}_2 = 0$$

with  $\mathbf{E} = \mathbf{S}(\mathbf{o}_{1,2}^1) \mathbf{R}_2^1$  known as a (3x3) matrix termed *essential matrix*

- The epipolar constraint expresses in analytic form the geometric constraint existing between the projections of the same point on the image planes of two different cameras
  - The segments  $O_1P$ ,  $O_2P$  and  $O_1O_2$  (base line) must lie on the same plane (epipolar plane)
- The projections of the optical centers  $O_1$ ,  $O_2$  of each camera on the image plane of the other camera  $\mathbf{e}_1$ ,  $\mathbf{e}_2$  are termed *epipoles*
  - The epipoles and the base line do not change if the point P varies
- The segments  $\ell_1$  and  $\ell_2$  are termed *epipolar lines* (intersection of the epipolar plane with the image planes of the two cameras)
- For the purpose of *computing the correspondences*
  - Finding correspondences reduces to searching for a point along the epipolar line and not on the whole image plane

## Triangulation

- Purely geometric method of 3D reconstruction which can be used in the case that both intrinsic and extrinsic parameters of the camera are known (calibrated cameras)
- This method allows the computation of the coordinates  $\mathbf{p} = [p_x \ p_y \ p_z]^T$  of a point  $P$  with respect to the base frame, starting from the normalized coordinates  $\mathbf{s}_i = [X_i \ Y_i]^T$  of the projections of  $P$  on the image planes

- If the base frame coincides with Frame 1, it can be easily shown

$$\begin{aligned} \mathbf{p} &= \lambda_1 \tilde{\mathbf{s}}_1 \\ \mathbf{p} &= \mathbf{o} + \lambda_2 \mathbf{R} \tilde{\mathbf{s}}_2 \end{aligned}$$

which can be resolved with respect to  $\mathbf{p}$

- Computing  $\lambda_1$  from the third scalar equation of the first vectorial equation and replacing its value into the other two yields

$$\begin{bmatrix} 1 & 0 & -X_1 \\ 0 & 1 & -Y_1 \end{bmatrix} \mathbf{p} = \mathbf{0}$$

- Using a similar derivation for the second vectorial equation obtained by premultiplying both sides by  $\mathbf{R}^T$ , the following system is obtained

$$\begin{bmatrix} \mathbf{r}_1^T - X_2 \mathbf{r}_3^T \\ \mathbf{r}_2^T - Y_2 \mathbf{r}_3^T \end{bmatrix} \mathbf{p} = \begin{bmatrix} o_x - o_z X_2 \\ o_y - o_z Y_2 \end{bmatrix}$$

with  $\mathbf{R} = [\mathbf{r}_1 \ \mathbf{r}_2 \ \mathbf{r}_3]$  and  $\mathbf{R}^T \mathbf{o} = [o_x \ o_y \ o_z]^T$

- A system of 4 linear equations with respect to  $\mathbf{p}$  with 3 unknowns has been obtained, of which only 3 equations are linear independent
- In practical applications, because of noise, these 4 equations are independent and the system has only one approximated solution: suitable algorithms based on least-squares techniques have to be adopted to compute it

- In the case of cameras with parallel and aligned image planes, the solution can be greatly simplified
  - $\mathbf{R} = \mathbf{I}$  and  $\mathbf{R}^T \mathbf{o} = [b \ 0 \ 0]^T$ , with  $b > 0$  being the distance between the origins of the two cameras frame
  - The system solution is

$$p_x = \frac{X_1 b}{X_1 - X_2}$$

$$p_y = \frac{Y_1 b}{X_1 - X_2} = \frac{Y_2 b}{X_1 - X_2}$$

$$p_z = \frac{b}{X_1 - X_2}.$$

## Absolute orientation

- *Definition:* Computation of the relative motion between a system of two calibrated cameras and a known rigid object
  - Under the assumption of rigid motion, let us consider the base frame coincident with Frame 1
  - If the stereo camera system is moving and the object is fixed, let  $\mathbf{p}_1, \dots, \mathbf{p}_n$  be the position vectors of  $n$  points measured at time  $t$ , and let  $\mathbf{p}'_1, \dots, \mathbf{p}'_n$  be the position vectors of the same points measured at time  $t'$ , using triangulation
  - With the assumption of rigid motion

$$\mathbf{p}_i = \mathbf{o} + \mathbf{R}\mathbf{p}'_i \quad i = 1, \dots, n$$

where  $\mathbf{o}$  and  $\mathbf{R}$  are the unknowns

- From rigid body mechanics it is known that this problem has a unique solution in the case of three non-collinear points
  - Since points are obtained using triangulation, the measurements are affected by error and the system may have no solution
  - It is convenient to consider a number  $n > 3$ , and the computed solution is that which minimizes the function

$$\sum_{i=1}^n \|\mathbf{p}_i - \mathbf{o} - \mathbf{R}\mathbf{p}'_i\|^2$$

with the constraint that  $\mathbf{R}$  is a rotation matrix



- The problem can be separated in two parts
  - Observe that the value  $\mathbf{o} = \bar{\mathbf{p}} - \mathbf{R}\bar{\mathbf{p}}'$  minimizes the function if  $\bar{\mathbf{p}}$  and  $\bar{\mathbf{p}}'$  are chosen as the centroids of the set of points  $\{\mathbf{p}_i\}$  and  $\{\mathbf{p}'_i\}$ , and then defined as

$$\bar{\mathbf{p}} = \frac{1}{n} \sum_{i=1}^n \mathbf{p}_i, \quad \bar{\mathbf{p}}' = \frac{1}{n} \sum_{i=1}^n \mathbf{p}'_i$$

- The problem becomes that of computing  $\mathbf{R}$  which minimizes the linear quadratic function

$$\sum_{i=1}^n \|\bar{\mathbf{p}}_i - \mathbf{R}\bar{\mathbf{p}}'_i\|^2 \quad \begin{array}{l} \bar{\mathbf{p}}_i = \mathbf{p}_i - \bar{\mathbf{p}} \\ \bar{\mathbf{p}}'_i = \mathbf{p}'_i - \bar{\mathbf{p}}' \end{array}$$

- The solution is the matrix  $\mathbf{R}$  which minimizes the trace of the matrix  $\mathbf{R}^T \mathbf{K}$ , with

$$\mathbf{K} = \sum_{i=1}^n \bar{\mathbf{p}}_i \bar{\mathbf{p}}'_i{}^T$$

already presented in the analytic solution for the pose estimation

### 3D reconstruction from planar homography

- *Goal*: evaluate the motion of a cameras system which observes an unknown planar object
  - Rotation and displacement defined up to a scaling factor
- *Assumption*: feature points of the observed object lie on the same plane
  - Additional constraint to the epipolar constraint → planar homography
- By denoting with  $d_2 > 0$  the distance of the plane from the origin of Frame 2, with  $\mathbf{n}^2$  the unit vector orthogonal to the plane containing the feature points and with  $\mathbf{p}^2$  the position vector of a point  $P$  expressed with respect to Frame 2, it is

$$\frac{1}{d_2} \mathbf{n}^{2T} \mathbf{p}^2 = 1$$

From this relation and from the relation of the representations between the two cameras, it follows

$$\mathbf{p}^1 = \mathbf{H} \mathbf{p}^2 \quad \mathbf{H} = \mathbf{R}_2^1 + \frac{1}{d_2} \mathbf{o}_{1,2}^1 \mathbf{n}^{2T}$$

from which

$$\tilde{\mathbf{s}}_1 = \lambda \mathbf{H} \tilde{\mathbf{s}}_2 \quad \lambda = \lambda_2 / \lambda_1 > 0$$

By premultiplying both sides by  $\mathbf{S}(\tilde{\mathbf{s}}_1)$ , the following equality yields

$$\mathbf{S}(\tilde{\mathbf{s}}_1) \mathbf{H} \tilde{\mathbf{s}}_2 = \mathbf{0}$$

representing a planar homography defined by matrix  $\mathbf{H}$

- The system admits solution up to a scaling factor, if coordinates of 4 points in the plane are known
  - The constant value can be computed using a numerical algorithm based on the expression of the matrix  $\mathbf{H}$
  - Once  $\mathbf{H}$  is known, it is possible to compute  $\mathbf{n}^2$ ,  $\mathbf{R}_2^1$  and  $\mathbf{o}_{1,2}^1/d_2$ 
    - The solution is defined up to a scalar factor  $d_2$
    - Two admissible solutions exist!

## Camera Calibration

- Camera calibration
  - Estimation of the intrinsic parameters (matrix  $\Omega$ )
  - Estimation of the extrinsic parameters (pose of the camera frame with respect to the base frame or the end-effector frame)
- Various calibration techniques exist which are based on algorithms similar to those used for the pose estimation of an object with respect to the camera
  - If the intrinsic parameters are known, the solution to the  $PnP$  problem with  $n$  coplanar points can be directly used to compute the extrinsic parameters
    - With a fixed camera:  $T_c^b = T_o^b(T_o^c)^{-1}$
    - With a non-fixed camera:  $T_c^e = T_o^e(T_o^c)^{-1}$
- In a general case where the intrinsic parameters are not known, a suitable method in three phases can be adopted
  - **Phase 1: planar homography computation**
    - A planar homography can be computed starting from the pixel coordinates

$$c_i = \begin{bmatrix} X_{Ii} \\ Y_{Ii} \end{bmatrix}$$

from which the following equations system is obtained

$$S(c_i)H' [r_{x,i} \quad r_{y,i} \quad 1]^T = 0$$

$$H' = \Omega H$$

where  $\Omega$  is the intrinsic parameters matrix and  $H$  is the matrix

$$H = [r_1 \quad r_2 \quad o_{c,o}^e]$$

- With the method exposed in the case of analytic solution,  $\zeta H'$  can be computed in the case of  $n \geq 4$

○ **Phase 2: intrinsic parameters  $\Omega$  computation**

- Matrix  $\Omega$  can be computed from the elements of matrix  $\zeta \mathbf{H}'$

$$\zeta [\mathbf{h}'_1 \quad \mathbf{h}'_2 \quad \mathbf{h}'_3] = \Omega [r_1 \quad r_2 \quad o_{c,o}^c]$$

where  $\mathbf{h}'_i$  denotes the  $i$ -th column of  $\mathbf{H}'$

- Computing  $r_1$  and  $r_2$  and imposing the orthogonality and unit norm constraints on these vectors yields

$$\mathbf{h}'_1{}^T \Omega^{-T} \Omega^{-1} \mathbf{h}'_2 = 0$$

$$\mathbf{h}'_1{}^T \Omega^{-T} \Omega^{-1} \mathbf{h}'_1 = \mathbf{h}'_2{}^T \Omega^{-T} \Omega^{-1} \mathbf{h}'_2$$

which, being linear, can be rewritten in the form

$$\mathbf{A}' \mathbf{b} = \mathbf{0}$$

where  $\mathbf{A}'$  is a (2x6) matrix depending on  $\mathbf{h}'_1$  and  $\mathbf{h}'_2$ , while

$$\mathbf{b} = [b_{11} \quad b_{12} \quad b_{22} \quad b_{13} \quad b_{23} \quad b_{33}]^T$$

with  $b_{i,j}$  the generic element of the symmetric matrix

$$\mathbf{B} = \Omega^{-T} \Omega^{-1} = \begin{bmatrix} 1/\alpha_x^2 & 0 & -X_0/\alpha_x^2 \\ 0 & 1/\alpha_y^2 & -Y_0/\alpha_y^2 \\ -X_0/\alpha_x^2 & -Y_0/\alpha_y^2 & 1 + X_0^2/\alpha_x^2 + Y_0^2/\alpha_y^2 \end{bmatrix}$$

- By repeating the Phase 1  $k$  times, with the same plane placed in different pose, the system  $\mathbf{A}' \mathbf{b} = \mathbf{0}$  is obtained, where  $\mathbf{A}'$  is a  $2k \times 6$  matrix
- In the case  $k \geq 3$  a unique solution  $\gamma \mathbf{b}$  exists up to a scaling factor  $\gamma$
- From matrix  $\gamma \mathbf{B}$  the following expression for the intrinsic parameters can be found

$$X_0 = -b'_{13}/b'_{11}$$

$$Y_0 = -b'_{23}/b'_{22}$$

$$\alpha_x = \sqrt{\gamma/b'_{11}}$$

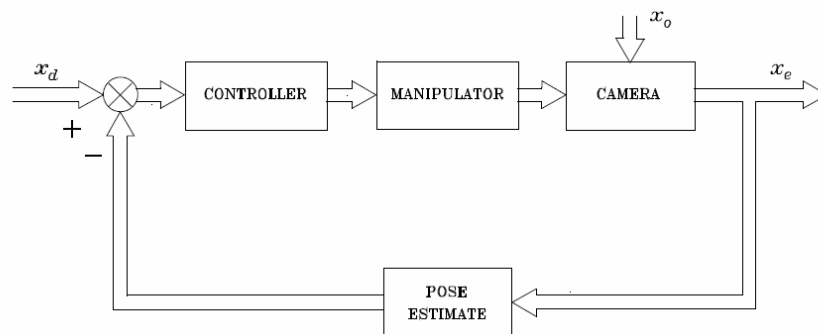
$$\alpha_y = \sqrt{\gamma/b'_{22}}$$

with  $b'_{i,j} = \gamma b_{i,j}$  and  $\gamma = b_{13} + b_{23} + b_{33}$

- **Phase 3: extrinsic parameters computation**
  - Once matrix  $\mathbf{\Omega}$  of intrinsic parameters has been computed, it is possible to evaluate the matrix  $\mathbf{H}$  (up to a constant  $\zeta$ ) from matrix  $\mathbf{H}'$  for one of the  $k$  poses of the plane
  - Having evaluated  $\mathbf{H}$ , the computation of matrix  $\mathbf{T}_o^c$  can be performed as for the case of the solution of a PnP problem
  - It is then possible to compute
    - In eye-to-hand case:  $\mathbf{T}_c^b = \mathbf{T}_o^b (\mathbf{T}_o^c)^{-1}$
    - In eye-in-hand case:  $\mathbf{T}_c^e = \mathbf{T}_o^e (\mathbf{T}_o^c)^{-1}$
- The method illustrated above is merely conceptual
  - Does not provide satisfactory solutions in the presence of measurement noise
  - Lens distortion is not considered
  - The accuracy of the results can be improved using models that take in account the distortion phenomena of the lenses, together with nonlinear optimization techniques
    - A technique similar to that proposed with the assumption of absence of lens distortion and of a linear model can be used
    - Then, nonlinear techniques for optimization can be used for the estimation of the parameters of the complete nonlinear model
- Calibration plane with a “chessboard” pattern
- A calibration method can also be set up starting from the solution of a PnP nonplanar problem using the direct linear transformation method

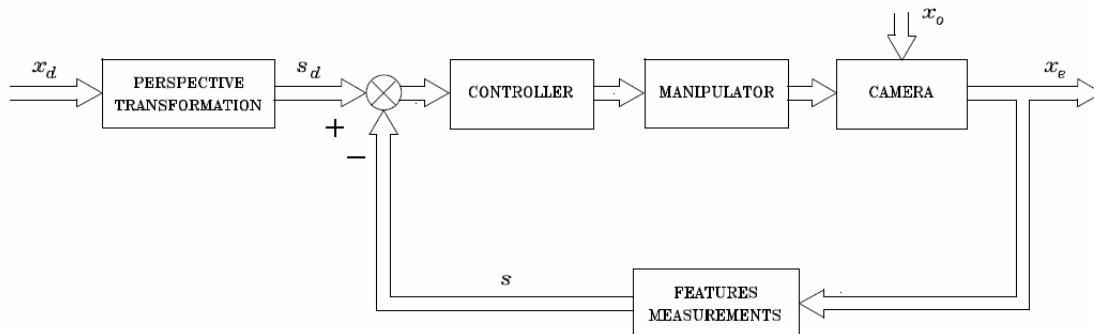
## The Visual Servoing Problem

- *Goal:* The objective of visual servoing is to ensure that the end-effector, on the basis of the visual measures elaborated in real-time, reaches and keeps a (constant or time-varying) desired pose with respect to the observed object
  - Only regulation problem will be considered
- The visual system provides feature parameters in the image plane as direct measures, while the robotic task is defined in the operational space, in terms of the relative pose of the end-effector with respect to the object
- Two kinds of control approaches can be defined
  - Visual Servoing in the Operational Space (Position Based Visual Control)
  - Visual Servoing in the Image Space (Image Based Visual Control)
- Visual servoing in the operational space



- Real-time estimation of the relative pose between the object and the camera
  - Estimation can be performed analytically or using iterative numerical algorithms
- Direct control of the operational space variables
- The absence of a direct control of the image features can lead the object to exit from the camera field of view → open loop!
- Very sensitive with respect to calibration errors

- If one camera is used, the object model is needed
- Visual servoing in the image space



- The control action is computed on the basis of the error defined as the difference between the values of image features in the desired configuration and in the current one
- The object pose estimation with respect to the camera is not required
- It is possible to keep the object within the camera field of view, but singular configurations may occur
- The end-effector trajectories are not easily predictable → collisions, joints limits violations
- Robust with respect to calibration errors
- The observed object model is not needed

## Position-based Visual Servoing

- Only regulation problem will be considered (constant reference)
- Assumptions:
  - Calibrated camera
  - End-effector frame  $\equiv$  Camera frame
  - Fixed object with respect to the base frame
- Features measurements for on-line estimation of matrix  $\mathbf{T}_o^c \rightarrow \mathbf{x}_{c,o}$
- The problem can be formulated by imposing a desired value to the relative pose of the object frame with respect to the camera  $\mathbf{T}_o^d \rightarrow \mathbf{x}_{d,o}$
- The displacement matrix of the camera frame in the current pose with respect to the desired pose can be evaluated

$$\mathbf{T}_c^d = \mathbf{T}_o^d (\mathbf{T}_o^c)^{-1} = \begin{bmatrix} \mathbf{R}_c^d & \mathbf{o}_{d,c}^d \\ \mathbf{0}^T & 1 \end{bmatrix}$$

and a suitable error vector in operational space can be computed

$$\tilde{\mathbf{x}} = - \begin{bmatrix} \mathbf{o}_{d,c}^d \\ \phi_{d,c} \end{bmatrix}$$

which does not depend on the object pose and does not coincide with the difference between  $\mathbf{x}_{d,o} - \mathbf{x}_{c,o}$

- The control has to be designed so that the operational space error  $\tilde{\mathbf{x}}$  tends to zero asymptotically
- The control objective can be satisfied provided that  $\mathbf{x}_{d,o}$ , corresponding to the homogeneous transformation matrix

$$\mathbf{T}_d = \mathbf{T}_c (\mathbf{T}_c^d)^{-1} = \begin{bmatrix} \mathbf{R}_d & \mathbf{o}_d \\ \mathbf{0}^T & 1 \end{bmatrix}$$

belongs to the dexterous workspace of the manipulator



## PD control with gravity compensation

- Computation of the time derivative of  $\tilde{x}$  yields

$$\dot{o}_{d,c}^d = \dot{o}_c^d - \dot{o}_d^d = \mathbf{R}_d^T \dot{o}_c$$

$$\dot{\phi}_{d,c} = \mathbf{T}(\phi_{d,c})^{-1} \omega_{d,c}^d = \mathbf{T}(\phi_{d,c})^{-1} \mathbf{R}_d^T \omega_c$$

then taking in account that  $\dot{o}_d^d = \mathbf{0}$  and  $\omega_d^d = \mathbf{0}$ , being  $\mathbf{o}_d$  and  $\mathbf{R}_d$  constant yields

$$\dot{\tilde{x}} = -\mathbf{T}_A(\phi_{d,c})^{-1} \begin{bmatrix} \mathbf{R}_d^T & \mathbf{O} \\ \mathbf{O} & \mathbf{R}_d^T \end{bmatrix} v_c$$

- Since the end-effector frame and the camera frame coincide

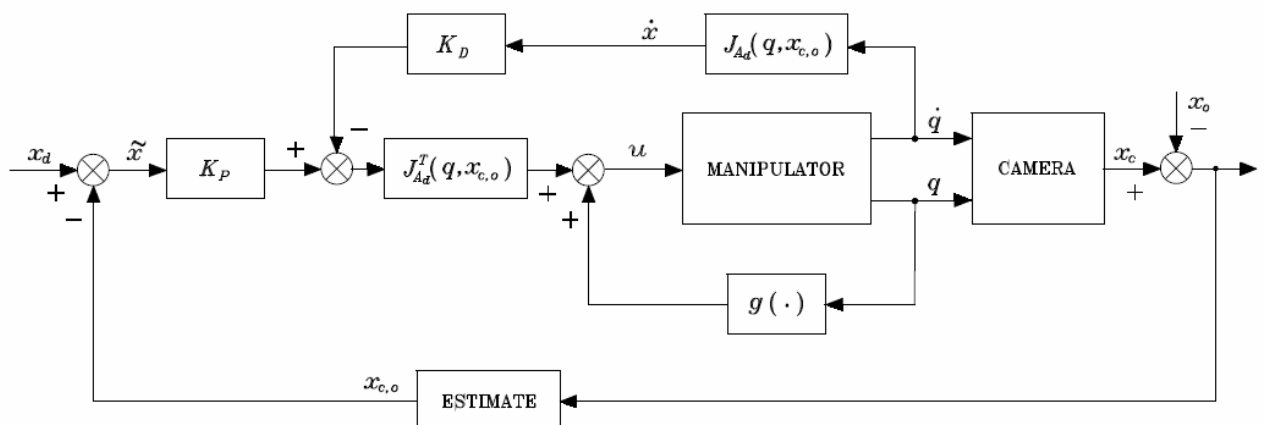
$$\dot{\tilde{x}} = -\mathbf{J}_{A_q}(q, \tilde{x}) \dot{q}$$

$$\mathbf{J}_{A_q}(q, \tilde{x}) = \mathbf{T}_A(\phi_{d,c})^{-1} \begin{bmatrix} \mathbf{R}_d^T & \mathbf{O} \\ \mathbf{O} & \mathbf{R}_d^T \end{bmatrix} \mathbf{J}(q)$$

- Position-based visual servoing of PD type with gravity compensation has the expression

$$u = g(q) + \mathbf{J}_{A_q}^T(q, \tilde{x}) \mathbf{K}_P \tilde{x} - \mathbf{J}_{A_q}^T \mathbf{K}_D \mathbf{J}_{A_q}(q, \tilde{x}) \dot{q}$$

which has the following block scheme (the sum blocks have a purely conceptual meaning!)



- The asymptotic stability of the equilibrium pose corresponding to  $\tilde{x} = \mathbf{0}$ , under the assumption of symmetric and positive definite matrices  $\mathbf{K}_P$  and  $\mathbf{K}_D$ , can be proven using the following Lyapunov function

$$V = \frac{1}{2} \dot{q}^T B(q) \dot{q} + \frac{1}{2} \tilde{x}^T K_P \tilde{x}$$

- The derivative term can also be chosen as  $-\mathbf{K}_D \dot{q}$

## Resolved-velocity control

- Performance limits due to
  - Information deriving from visual measurements is computed at frequency lower or equal to the camera frame rate
  - Frequencies are at least of one order of magnitude lower than the typical frequencies used for motion control of robot manipulators
  - The control gains must be set to values much lower than those typically used for motion control to preserve stability of the closed-loop system
    - Low speed of convergence
    - Poor disturbance rejection capability
- This problem can be avoided assuming that the manipulator is equipped with a high-gain motion controller
  - Effects on the tracking errors deriving from manipulator dynamics and disturbances are neglected (manipulator as an ideal positioning device)

$$\mathbf{q}(t) \simeq \mathbf{q}_r(t).$$

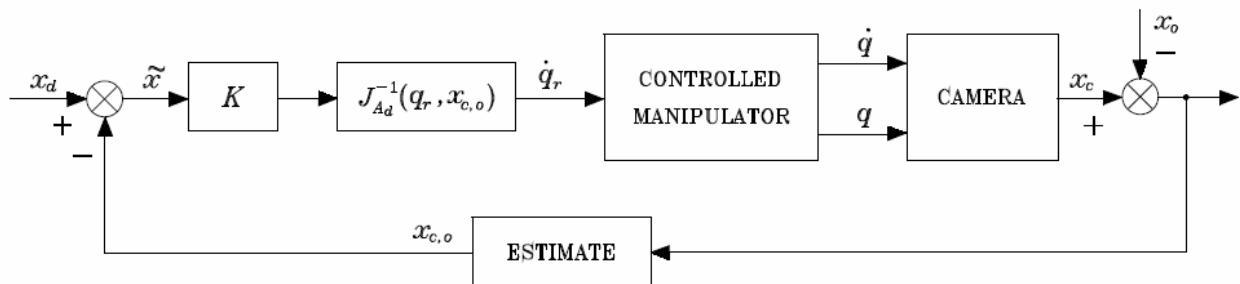
- $\mathbf{q}_r$  is computed from visual measurements so that the operational space tracking error goes asymptotically to zero
- From the equation about the error time derivative, it can be chosen

$$\dot{\mathbf{q}}_r = \mathbf{J}_{A_q}^{-1}(\mathbf{q}_r, \tilde{\mathbf{x}}) \mathbf{K} \tilde{\mathbf{x}}$$

which yields the linear equation

$$\dot{\tilde{\mathbf{x}}} + \mathbf{K} \tilde{\mathbf{x}} = \mathbf{0}$$

which is asymptotically stable under the assumption that  $\mathbf{K}$  is a positive definite matrix



- It is a *resolved-velocity control*, because it is based on the computation of velocity from the operational space error
- The choice of  $\mathbf{K}$  influences the transient behavior of the trajectory of the camera frame
  - If  $\mathbf{K}$  is a diagonal matrix with the same gains for the positional part, the origin of the camera frame follows the line segment connecting the initial position to the desired position
  - The orientation trajectory depends on the choice of  $\mathbf{K}$  and on the particular choice of Euler angles
  - The possibility of knowing in advance the trajectory of the camera is important to avoid that object may exit from the camera field of view, making visual measurements unavailable

## Image-based Visual Servoing

- A desired constant value  $\mathbf{s}_d$  of the feature parameters is imposed, under the assumption that the object is fixed with respect to the base frame
  - It is assumed that the desired pose  $\mathbf{s}_d = \mathbf{s}(\mathbf{x}_{d,o})$  belongs to the dexterous workspace of the manipulator and univocally determined
    - $n \geq 4$  points for coplanar points (no triplet of collinear points)
    - $n \geq 6$  points for non-coplanar points
    - In the case of  $m < 6$ , a reduced number of points can be used
  - $\mathbf{s}_d$  can be directly computed by measuring the feature parameters when the object is in the desired pose with respect to the camera
    - $\mathbf{x}_{d,o}$  does not need to be known
  - It is worth recalling that  $\mathbf{L}_s(\mathbf{s}, \mathbf{z}_c)$  depends on  $\mathbf{z}_c = [z_{c,1}, \dots, z_{c,n}]$
- The control law must be designed so as to guarantee that the image space error

$$\mathbf{e}_s = \mathbf{s}_d - \mathbf{s}$$

tends asymptotically to zero

*PD control with gravity compensation*

- Consider the following Lyapunov function

$$V(\dot{q}, e_s) = \frac{1}{2} \dot{q}^T B(q) \dot{q} + \frac{1}{2} e_s^T K_{Ps} e_s > 0 \quad \forall \dot{q}, e_s \neq 0$$

with  $K_{Ps}$  symmetric and positive definite ( $k \times k$ ) matrix

Computing the time derivative and taking into account the expression of the joint space dynamic model of the manipulator gives

$$\dot{V} = -\dot{q}^T F \dot{q} + \dot{q}^T (u - g(q)) - \dot{e}_s^T K_{Ps} e_s$$

Since  $\dot{s}_d = 0$  and the object is fixed with respect to the base frame, under the assumption that the camera frame coincides with the end-effector frame, it follows

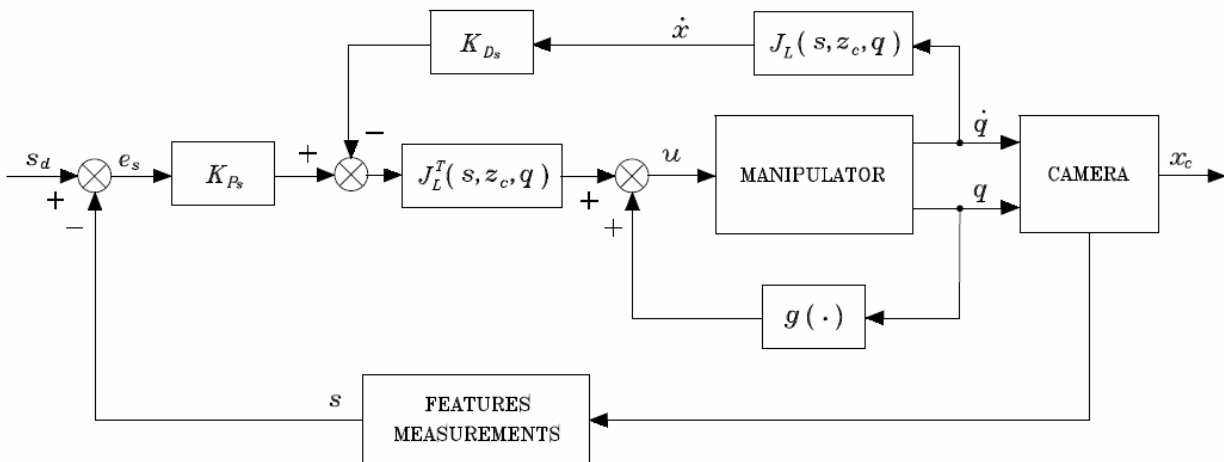
$$\begin{aligned} \dot{e}_s &= -\dot{s} = -J_L(s, z_c, q) \dot{q} \\ J_L(s, z_c, q) &= L_s(s, z_c) \begin{bmatrix} R_c^T & O \\ O & R_c^T \end{bmatrix} J(q) \end{aligned}$$

With the choice of the following control law

$$u = g(q) + J_L^T(s, z_c, q) K_{Ps} e_s - J_L^T K_{Ds} J_L(s, z_c, q) \dot{q}$$

with  $K_{Ds}$  symmetric and positive definite ( $k \times k$ ) matrix, it is

$$\dot{V} = -\dot{q}^T F \dot{q} - \dot{q}^T J_L^T K_{Ds} J_L \dot{q}$$



- The last term corresponds to a derivative action in the image space and has to be added to increase damping (derivative action in the image space)
  - It can be replaced with  $-\mathbf{K}_{D_s}\dot{\mathbf{s}}$ , but being  $\dot{\mathbf{s}}$  not available, the term  $-\mathbf{K}_D\dot{\mathbf{q}}$  can be used, with  $\mathbf{K}_D$  symmetric and positive definite ( $n \times n$ ) matrix

- The system Lyapunov function decreases until  $\dot{\mathbf{q}} \neq \mathbf{0} \rightarrow$  the system reaches the equilibrium state characterized by

$$\mathbf{J}_L^T(\mathbf{s}, \mathbf{z}_c, \mathbf{q})\mathbf{K}_{P_s}\mathbf{e}_s = \mathbf{0}$$

and if the interaction matrix and the geometric Jacobian are full rank, it will result  $\mathbf{e}_s = \mathbf{0}$

- The control law requires the computation of  $\mathbf{z}_c$ , which in the image-based visual servoing philosophy should be avoided
  - In some applications  $\mathbf{z}_c$  is known with good approximation (e.g. object belonging to a plane)
  - It is assumed  $\mathbf{z}_c$  constant, as the value in the initial configuration or that in the desired configuration  $\rightarrow$  it is equivalent to using an estimate  $\hat{\mathbf{L}}_s$  of the interaction matrix
    - The stability proof becomes much more complex!

## Resolved-velocity control

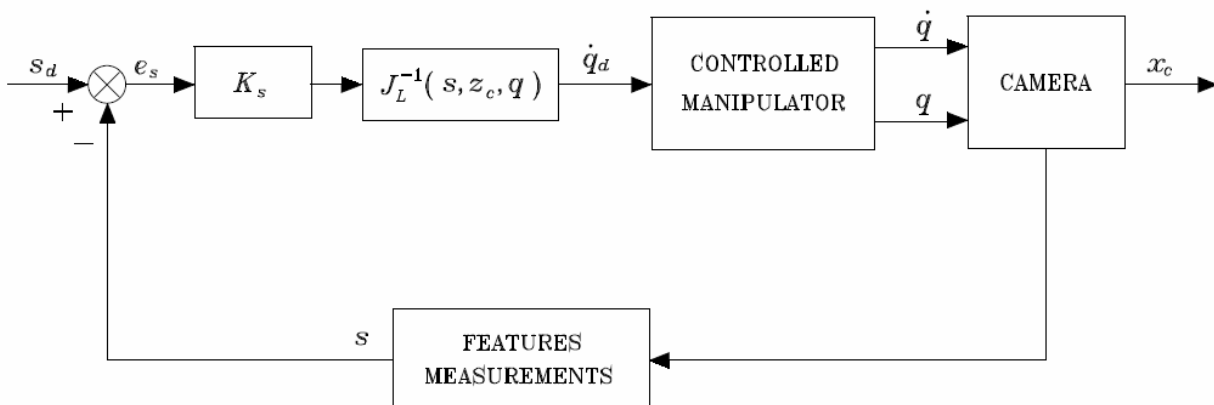
- If the following choice of the reference velocity in the joint space is assumed

$$\dot{q}_r = J_L^{-1}(s, z_c, q_r) K_s e_s$$

under the assumption of invertible matrix  $J_L$ , the following linear equation of the error holds

$$\dot{e}_s + K_s e_s = 0$$

If  $K_s$  is a positive definite matrix, the error tends asymptotically to zero



- The system is affected by problems related to the singularities of  $J_L$ 
  - The singularities are both those of the geometric Jacobian and those of the interaction matrix (!)
  - The control law is computed in two steps

$$v_r^c = L_s^{-1}(s, z_c) K_s e_s$$

$$\dot{q}_r = J^{-1}(q) \begin{bmatrix} R_c & O \\ O & R_c \end{bmatrix} v_r^c$$

- Using a number  $k$  of feature parameters greater than the minimum required  $m$ , the left pseudo-inverse can be used

$$v_r^c = (L_s^T L_s)^{-1} L_s^T K_s e_s$$



- Stability of the closed-loop system can be shown using the Lyapunov function

$$V(e_s) = \frac{1}{2} e_s^T K_s e_s$$

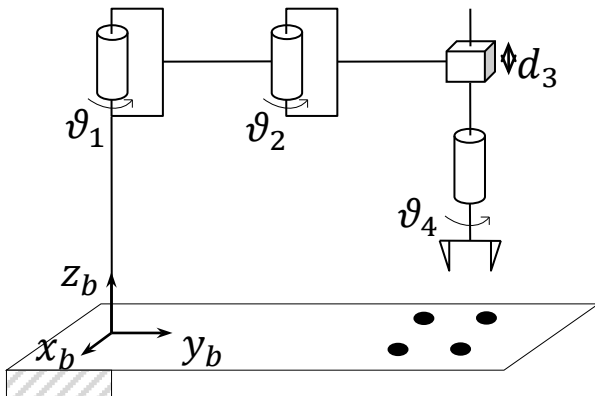
$$\dot{V} = -e_s^T K_s L_s (L_s^T L_s)^{-1} L_s^T K_s e_s$$

which is negative semi-definite because  $\mathcal{N}(L_s^T) \neq \emptyset$ , being  $L_s^T$  a matrix with more columns than rows  $\rightarrow$  the system is stable but not asymptotically stable

- The error is bounded, but in some cases the system may reach an equilibrium with  $e_s \neq 0$  and  $K_s e_s \in \mathcal{N}(L_s^T)$
- The computation of  $L_s$  requires the knowledge of  $z_c$ 
  - An estimate of matrix  $\hat{L}_s^{-1}$  (or of the pseudo-inverse) can be used
  - The Lyapunov method can be used to prove that the control scheme remains stable provided that matrix  $L_s \hat{L}_s^{-1}$  is positive definite
  - Notice that  $z_c$  is the only information depending on the object geometry
- The choice of  $K_s$  influences the trajectories of the feature parameters
  - In the case of features points, if a scalar matrix  $K_s$  is set, the projections of these points on the image plane will follow line segments
    - The corresponding camera motion cannot be easily predicted!

## Comparison Among Various Control Schemes

- SCARA manipulator ( $m = 4$ ) with planar object (4 points)



- A dynamic model is considered for the manipulator
- The control is implemented with a sample time of 0.04s (25Hz), which coincides with the minimum frame rate of analog cameras
- Visual servoing in the operational space:  $\mathbf{x}_{c,o}$  estimated with a numerical algorithm using only  $P_1$  and  $P_2$  ( $\mathbf{J}_{A_s}$  is a (4x4) square matrix), integration step of 1ms and  $\mathbf{K}_s = 160\mathbf{I}_4$  (converges in 0.03s)
- Visual servoing in the image space: only points  $P_1$  and  $P_2$  ( $\mathbf{J}_L$  (4x4) square matrix)
- Initial and desired values:
 
$$\mathbf{x}_c(0) = [1 \quad 1 \quad 1 \quad \pi/4]^T$$

$$\mathbf{x}_{c,o}(0) = [0 \quad 0 \quad 0.5 \quad 0]^T$$

$$\mathbf{x}_{d,o} = [-0.1 \quad 0.1 \quad 0.6 \quad -\pi/3]^T$$

$$\mathbf{s}_d = [-0.1667 \quad 0.1667 \quad -0.0833 \quad 0.0223]^T$$
- The dynamics of the velocity-controlled manipulator has been neglected
- $\hat{\mathbf{L}}_s = \mathbf{L}_s(\mathbf{s}, z_d)$ , with  $z_d$  the third component of  $\mathbf{x}_{d,o}$

- **Scheme A:** Position-based visual servoing of PD type with gravity compensation

$$\mathbf{K}_P = \text{diag}\{500, 500, 10, 10\}$$

$$\mathbf{K}_D = \text{diag}\{500, 500, 10, 10\}$$

- **Scheme B:** Resolved-velocity position-based visual servoing

$$\mathbf{K} = \text{diag}\{1, 1, 1, 2\}$$

corresponding to a time constant of 1s for the three position variables and of 0.5s for the orientation variable

- **Scheme C:** Image-based visual servoing of PD type with gravity compensation

$$\mathbf{K}_{P_s} = 300\mathbf{I}_4 \quad \mathbf{K}_{D_s} = 330\mathbf{I}_4$$

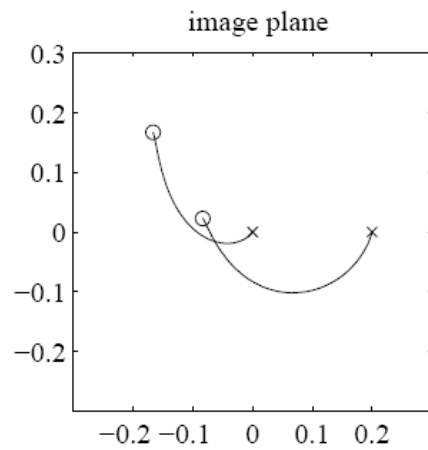
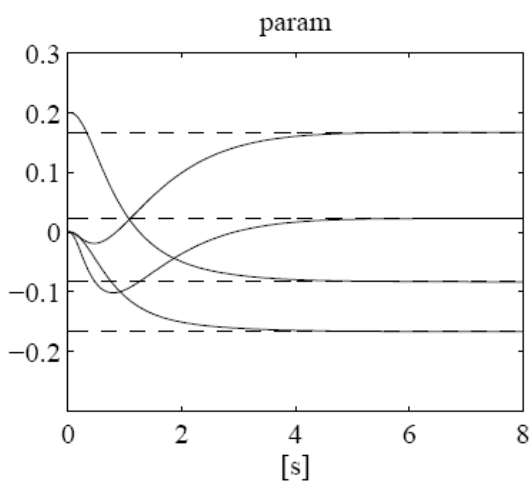
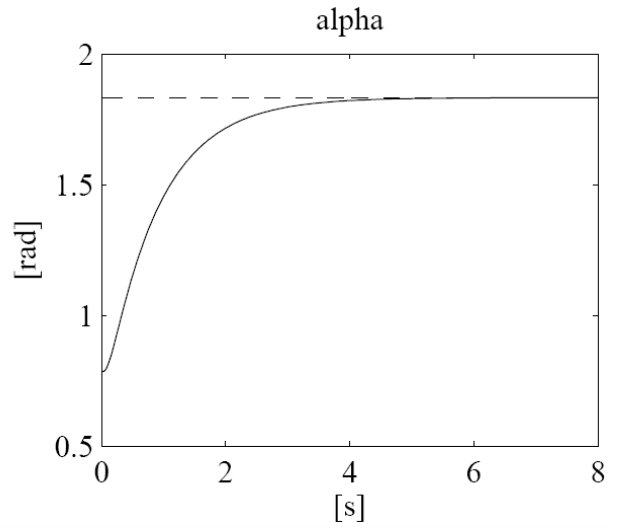
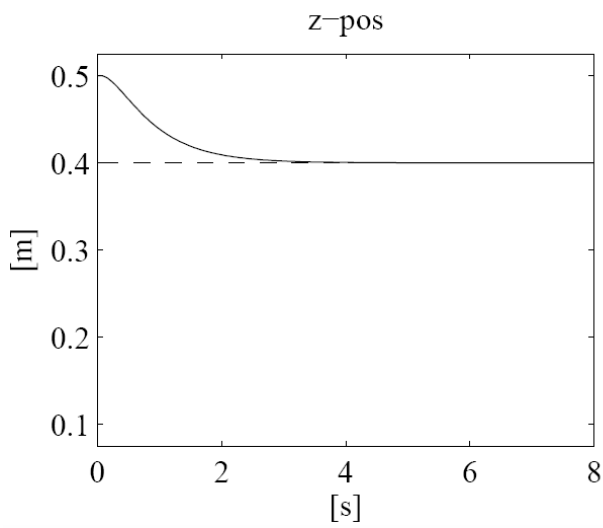
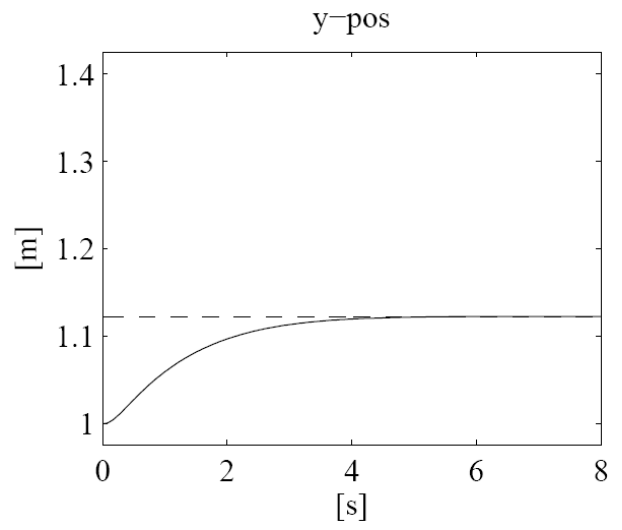
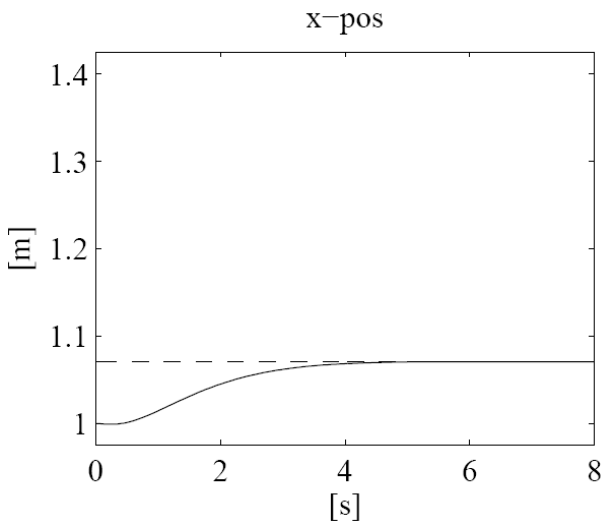
- **Scheme D:** Resolved-velocity image-based visual servoing

$$\mathbf{K}_s = \mathbf{I}_4$$

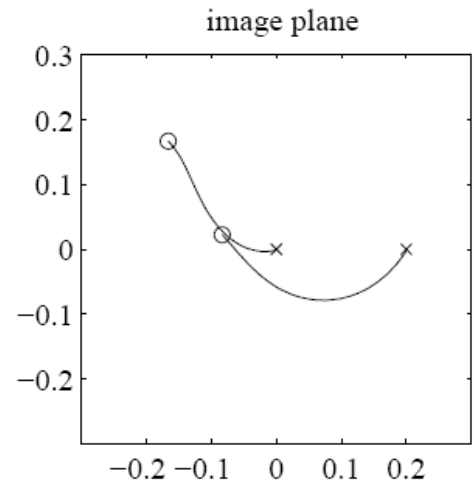
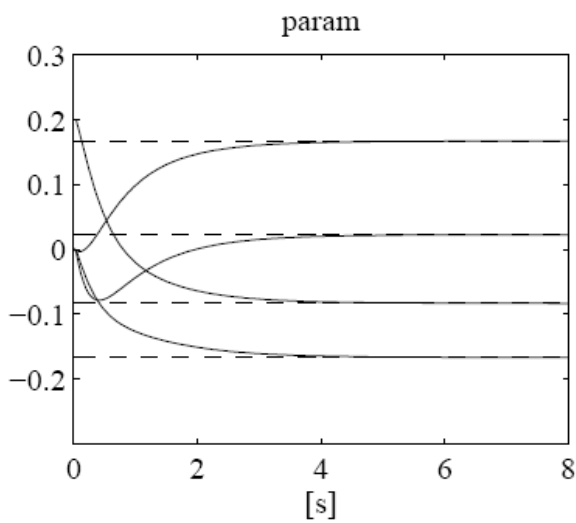
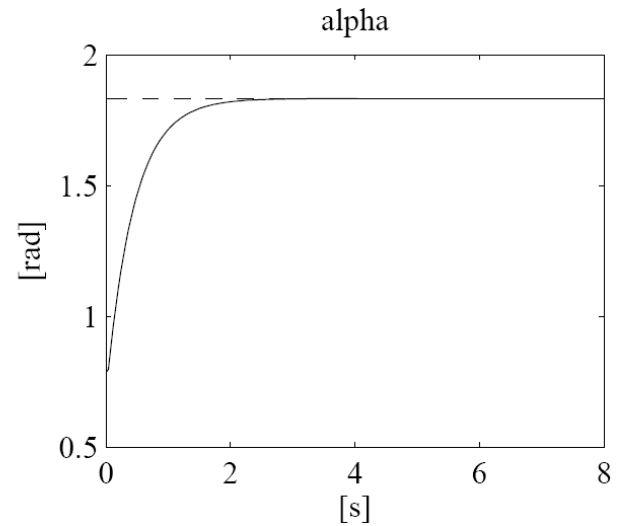
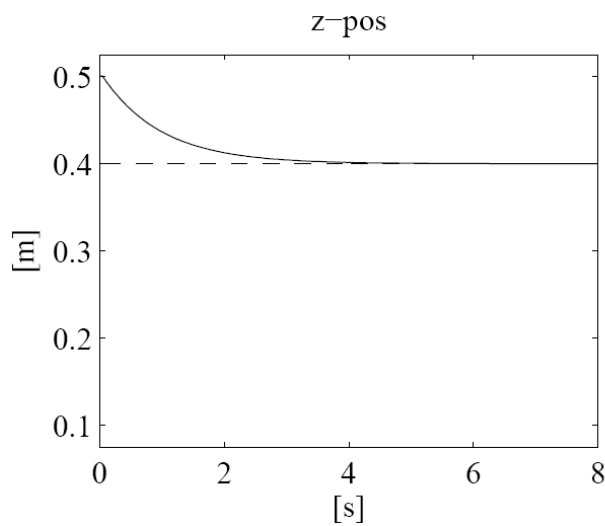
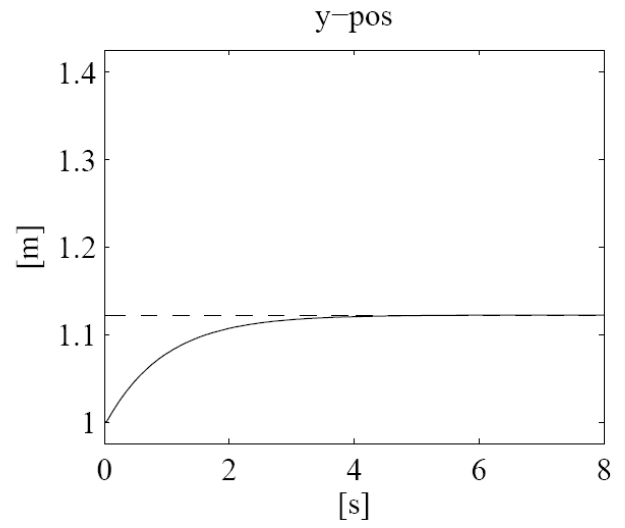
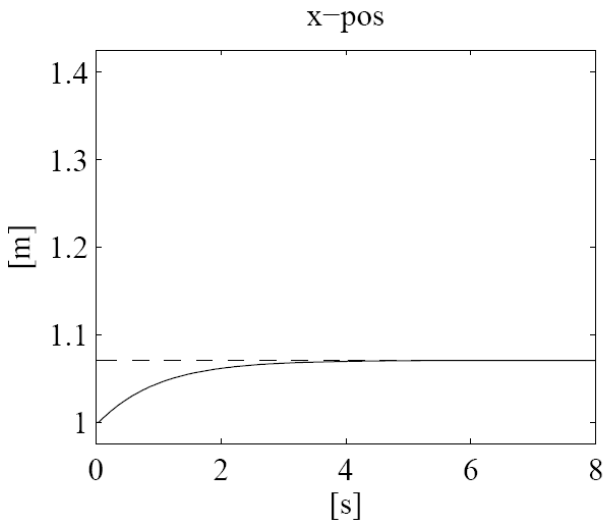
corresponding to a time constant of 1s for the feature parameters

The parameters of the various controllers have been chosen in such a way as to show the particular features of the different control laws and, at the same time, to allow a significant comparison of the performance of each scheme

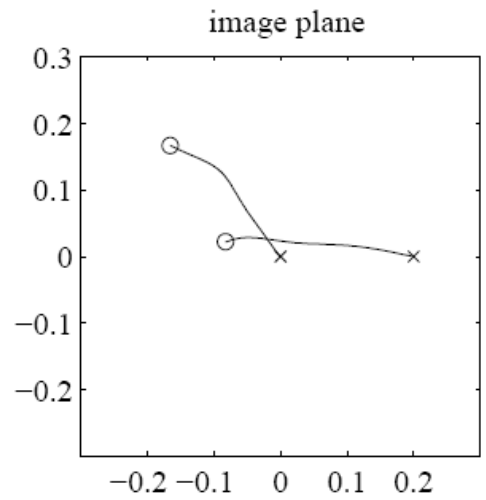
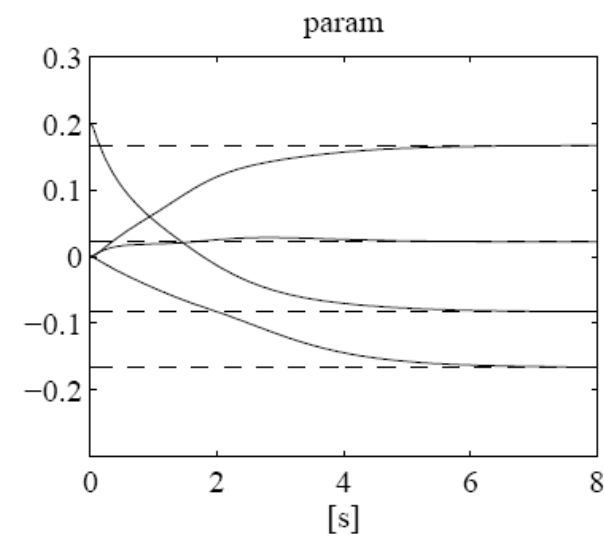
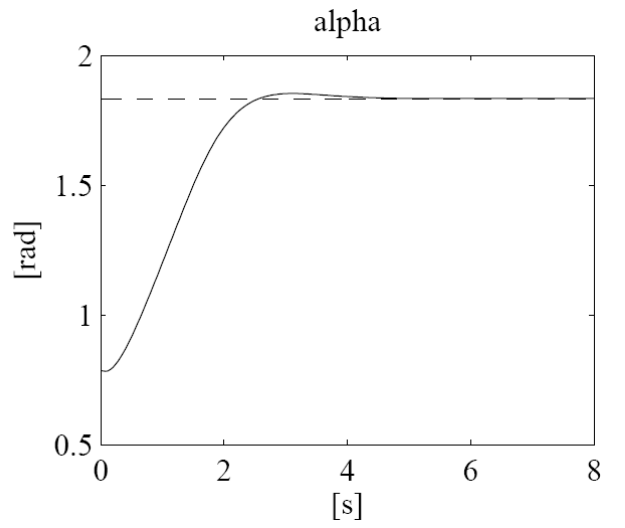
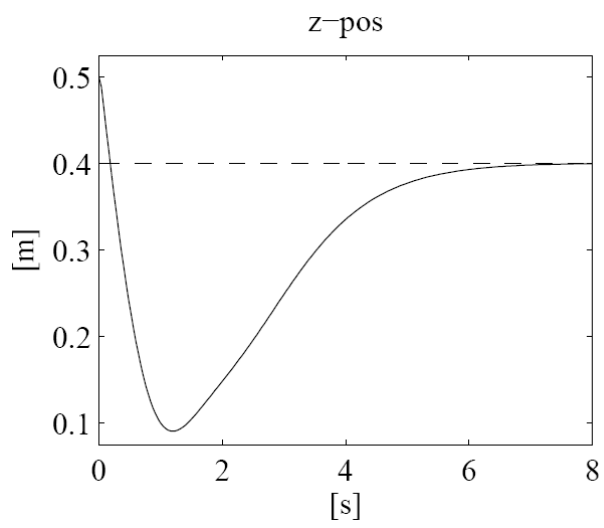
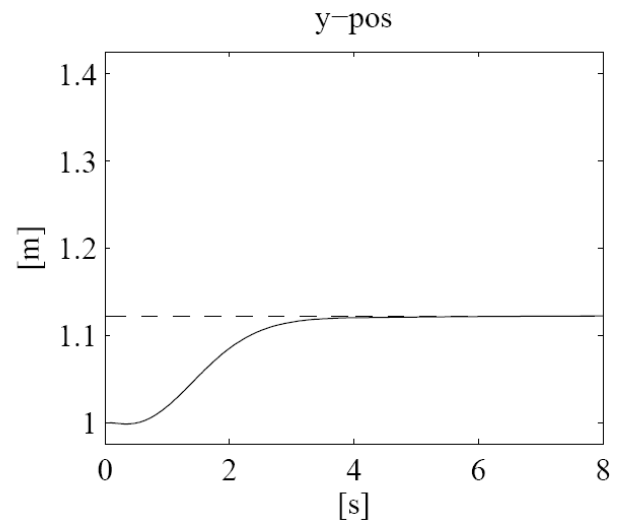
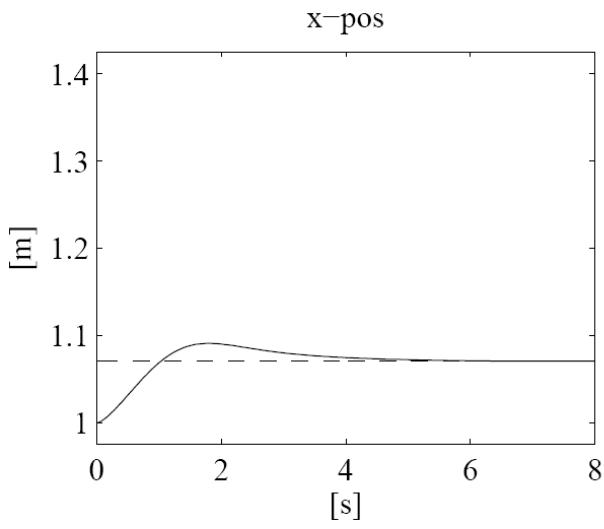
Scheme A



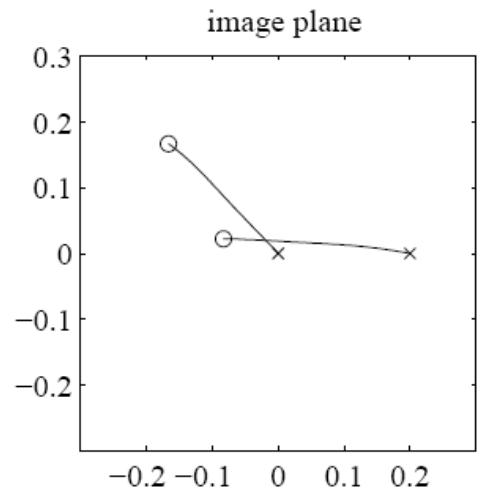
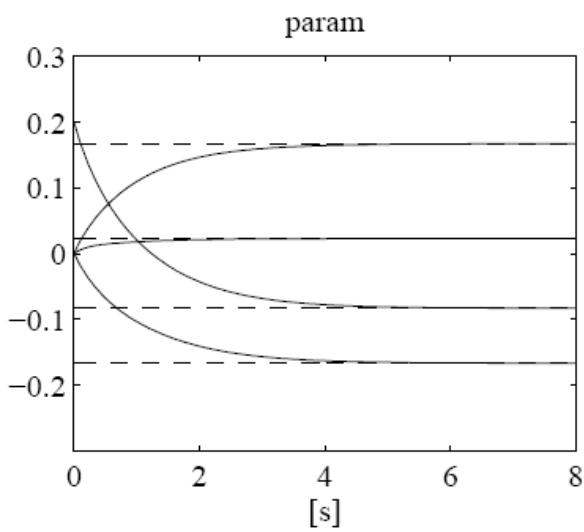
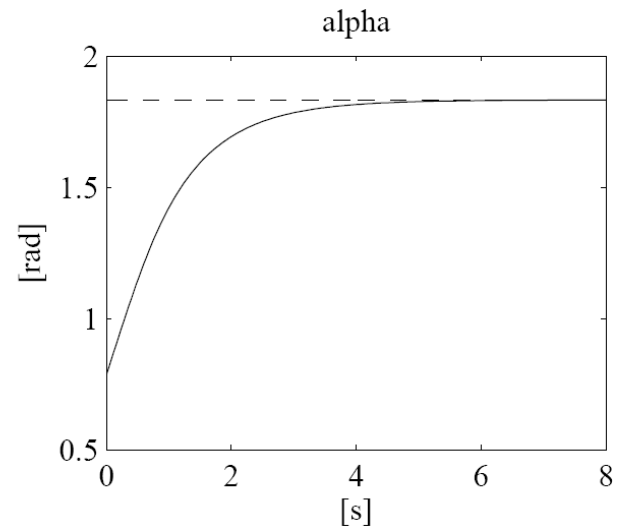
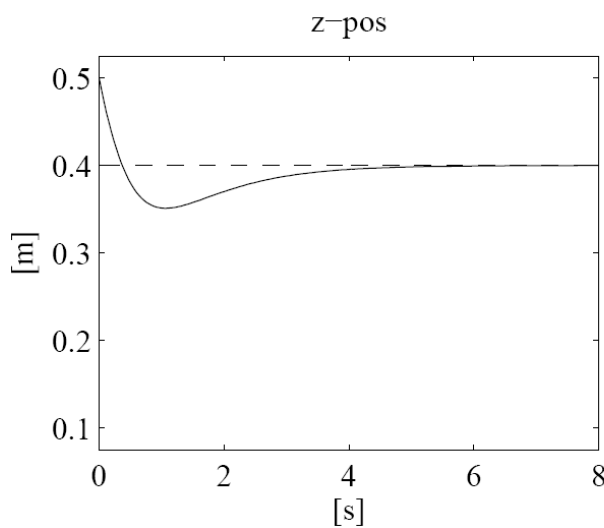
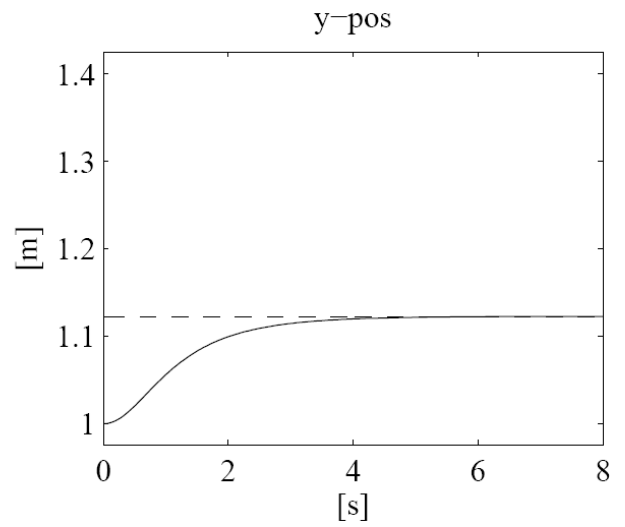
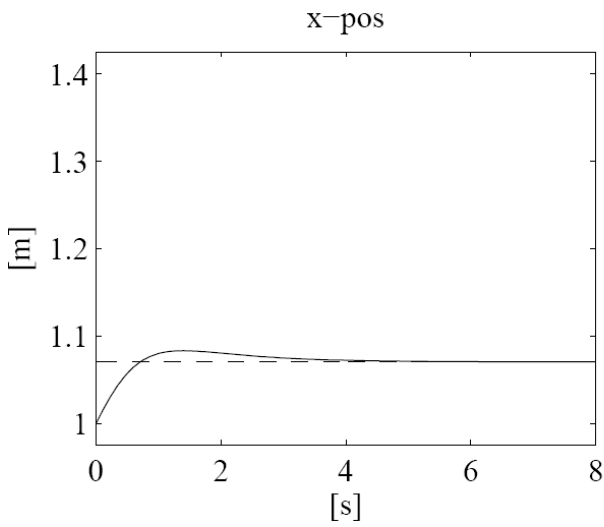
Scheme B



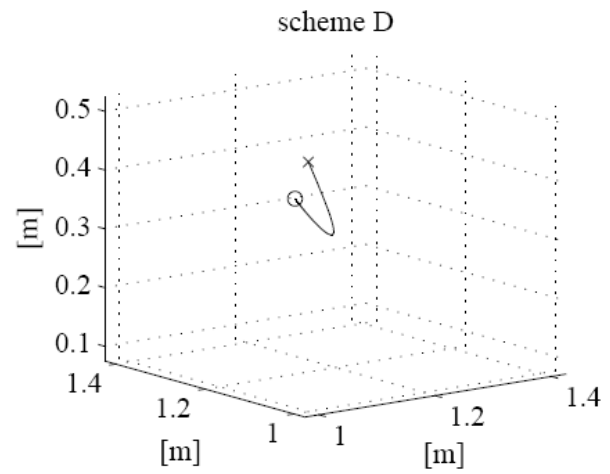
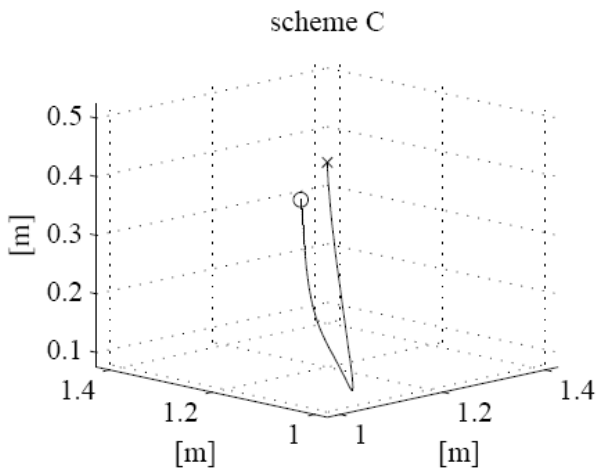
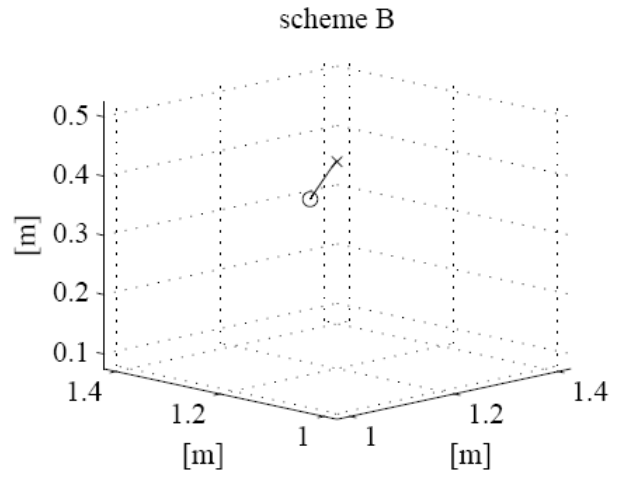
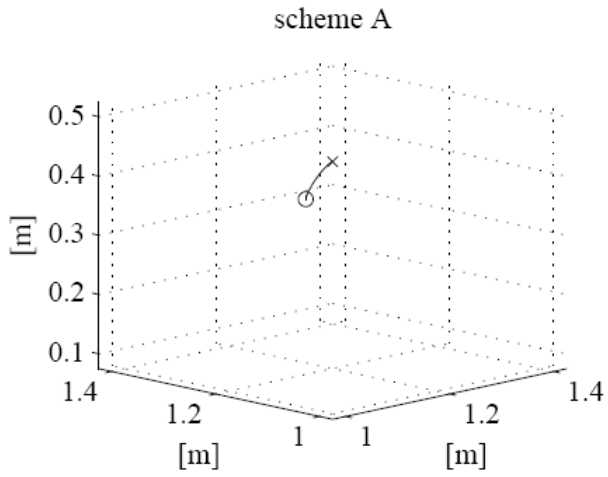
Scheme C



Scheme D



### Paths of camera frame origin





## Hybrid Visual Servoing

- It combines the benefits of position-based and image-based visual servoing
  - The control error is defined in the operational space for some components and in the image space for the others
    - Control on the operational space trajectories
    - Control on keeping image features in the camera field of view
- Estimation of some operational space variables is required
  - e.g. using planar homography in the case of planar object with at least 4 points (only  $\mathbf{s}$  and  $\mathbf{s}_d$  are necessary)
- Imposing the planar homography between the desired camera frame and the current one yields

$$\mathbf{H} = \mathbf{R}_d^c + \frac{1}{d_d} \mathbf{o}_{c,d}^c \mathbf{n}^{dT}$$

From  $\mathbf{H}$  it is possible to compute  $\mathbf{n}^{dT}$ ,  $\mathbf{R}_d^c$  e  $\mathbf{o}_{c,d}^c/d_c$  at each sampling time

- By adopting a resolved-velocity approach, the control objective consists of computing the reference absolute velocity of the camera frame

$$\mathbf{v}_r^c = \begin{bmatrix} \nu_r^c \\ \boldsymbol{\omega}_r^c \end{bmatrix}$$

- From matrix  $\mathbf{R}_d^c$  the orientation error can be computed
  - Having defined with  $\boldsymbol{\phi}_{c,d}$  the vector of the Euler angles extracted from  $\mathbf{R}_d^c$ , the control vector can be chosen as

$$\boldsymbol{\omega}_r^c = -\mathbf{T}(\boldsymbol{\phi}_{c,d}) \mathbf{K}_o \boldsymbol{\phi}_{c,d}$$

where  $\mathbf{K}_o$  is a symmetric positive definite (3x3) matrix, from which

$$\dot{\boldsymbol{\phi}}_{c,d} + \mathbf{K}_o \boldsymbol{\phi}_{c,d} = \mathbf{0}$$

implying that the orientation error tends to zero asymptotically with convergence of exponential type

- The control vector  $\mathbf{v}_r^c$  should be selected so that the positional part of the error between the desired and the current camera pose converges to zero
  - It is not possible to define the error as the difference of the coordinates of a point of the object (these cannot be directly measured)
    - The corresponding coordinates in the image plane are available

$$\begin{aligned} \mathbf{s}_{p,d} &= [X_d \ Y_d]^T = [x_d/z_d \ y_d/z_d]^T \\ \mathbf{s}_p &= [X \ Y]^T = [x_c/z_c \ y_c/z_c]^T \end{aligned}$$

- From the computation of the homography it can be possible to compute the values

$$\rho_z = z_c/z_d = \frac{d_c}{d_d} \frac{\mathbf{n}^{dT} \tilde{\mathbf{s}}_{p,d}}{\mathbf{n}^{cT} \tilde{\mathbf{s}}_p}$$

$$\frac{d_c}{d_d} = 1 + \mathbf{n}^{cT} \frac{\mathbf{o}_{c,d}^c}{d_d} = \det(\mathbf{H})$$

with  $\mathbf{n}^c = \mathbf{R}_d^c \mathbf{n}^d$

- The position error can be defined as

$$\mathbf{e}_p(\mathbf{r}_d^c, \mathbf{r}_c^c) = \begin{bmatrix} X_d - X \\ Y_d - Y \\ \ln \rho_z \end{bmatrix}$$

noticing that convergence to zero of  $\mathbf{e}_p$  implies convergence to zero of  $\mathbf{r}_d^c - \mathbf{r}_c^c$  and vice versa

- Computing the time derivative of  $\mathbf{e}_p$ , being  $\mathbf{r}_d^c$  constant, yields

$$\dot{\mathbf{e}}_p = \frac{\partial \mathbf{e}_p(\mathbf{r}_c^c)}{\partial \mathbf{r}_c^c} \dot{\mathbf{r}}_c^c = -\mathbf{J}_p \mathbf{v}_c^c - \mathbf{J}_o \boldsymbol{\omega}_c^c$$

where

$$\mathbf{J}_p = \frac{1}{z_d \rho_z} \begin{bmatrix} -1 & 0 & X \\ 0 & -1 & Y \\ 0 & 0 & -1 \end{bmatrix} \quad \mathbf{J}_o = \begin{bmatrix} XY & -1 - X^2 & Y \\ 1 + Y^2 & -XY & -X \\ -Y & X & 0 \end{bmatrix}$$

- The previous equation suggests the following choice (being  $\mathbf{J}_p(z_d)$  invertible)

$$\mathbf{v}_r^c = \mathbf{J}_p^{-1}(\mathbf{K}_p \mathbf{e}_p - \mathbf{J}_o \boldsymbol{\omega}_r^c)$$

- Under the assumption that  $z_d$  is known, and  $\mathbf{v}_c^c \approx \mathbf{v}_r^c$  and  $\boldsymbol{\omega}_c^c \approx \boldsymbol{\omega}_r^c$ , it yields

$$\dot{\mathbf{e}}_p + \mathbf{K}_p \mathbf{e}_p = \mathbf{0}$$

- If  $z_d$  is not known, the estimate can be adopted

$$\hat{\mathbf{J}}_p^{-1} = \frac{\hat{z}_d}{z_d} \mathbf{J}_p^{-1}$$

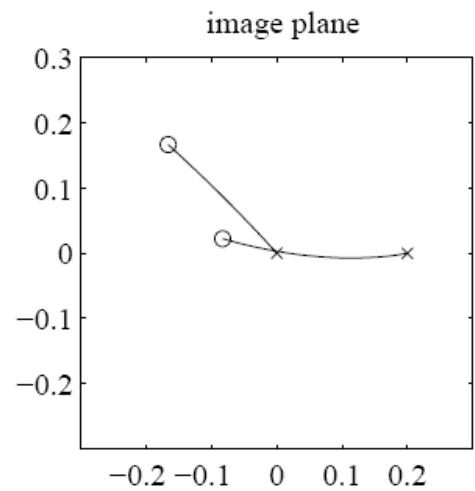
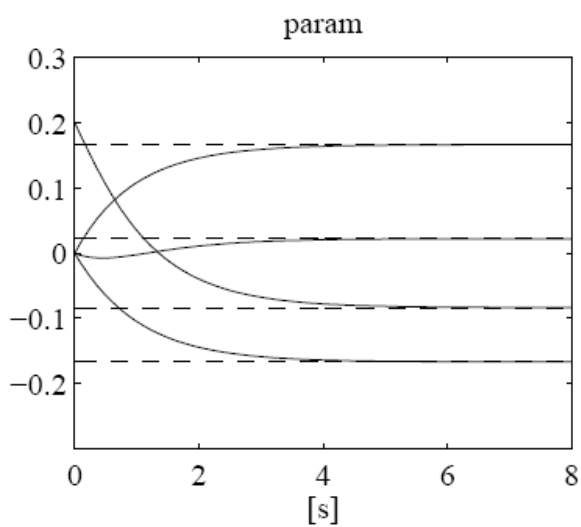
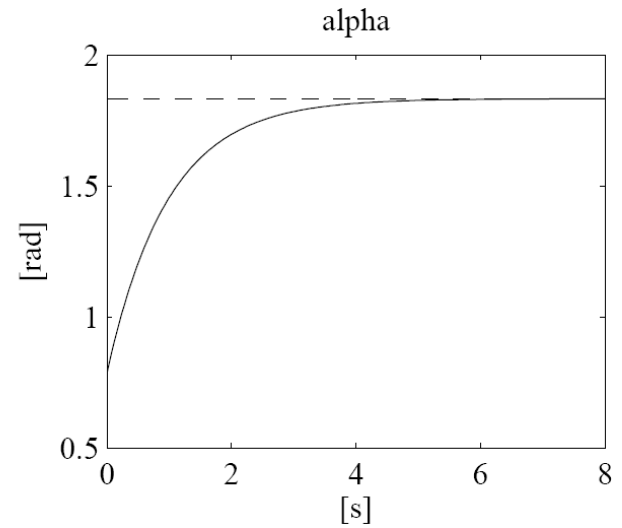
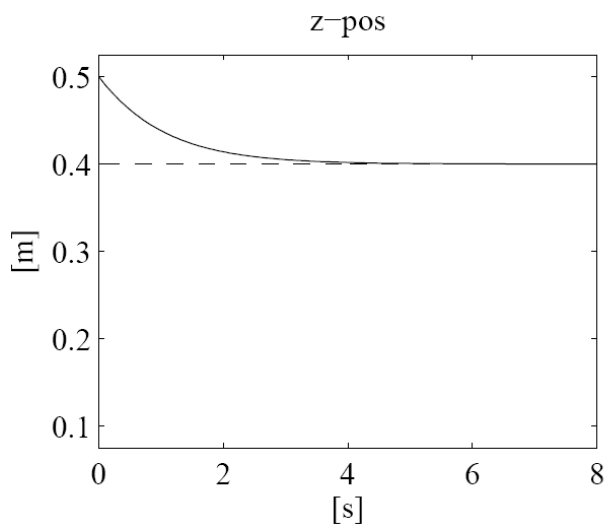
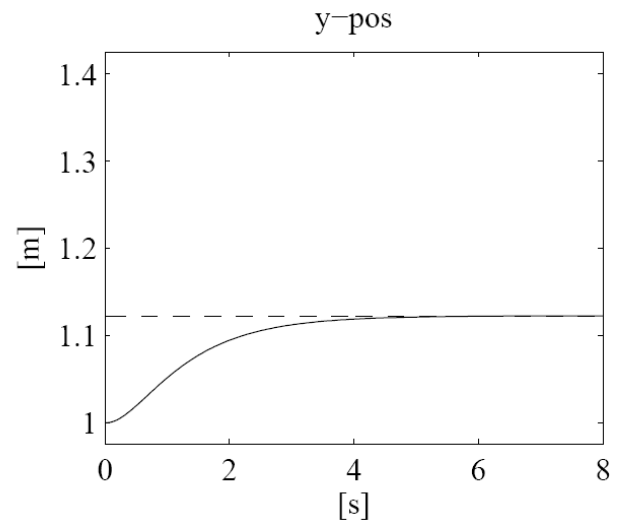
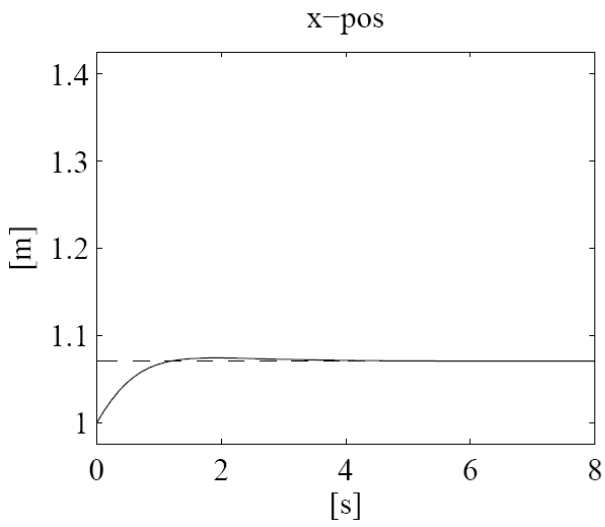
from which the following error equation is obtained

$$\dot{\mathbf{e}}_p + \frac{\hat{z}_d}{z_d} \mathbf{K}_p \mathbf{e}_p = \left(1 - \frac{\hat{z}_d}{z_d}\right) \mathbf{J}_o \boldsymbol{\omega}_r^c$$

where the gain is scaled and the time history of  $\mathbf{e}_p$  depends on  $\boldsymbol{\omega}_r^c$

- This method is only one of the possible visual servoing approaches based on the computation of the planar homography and on its decomposition

Hybrid control scheme



## Path of camera frame origin

