

# Mutual-Information Based Visual Servoing

*Master of Science in Engineering in AI and Robotics*

*Master of Science in Control Engineering*

*Course of Medical Robotics*

*Giuseppe andrea Cuzzo, Francesca Palermo*



**SAPIENZA**  
UNIVERSITÀ DI ROMA

# Introduction (1)

What is Visual Servoing and what is its goal?

Visual Servoing uses the information acquired by a vision sensors (such as cameras) for feedback control of the pose/motion of a robot (or of parts of it). In the past, this approach required the extraction of visual information from the image in order to design the control law.

Nowadays we can find different kind of methods but the goal remains the same for all of them:

From its current pose  $\mathbf{r}$ , the robot has to reach the desired pose  $\mathbf{r}^*$

# Introduction (2)

## The approach of the paper

The authors of the paper (A. Dame, E. Marchand) use a method which involves information contained in the images, since the frames' intensities are quite sensitive to modification of the environments.

They use the Mutual Information (MI) defined by Shannon, which consists in measuring the mutual dependence between two random variable and compares the distribution of the information in the images.

The higher the MI, the better the alignment between two images is.

In this way we have a control law which doesn't require any tracking or matching step.

# Direct Approaches

## Kernel Visual Servoing

- Large convergence domain
- No precise alignment information
- Task limited to 4 DOF
- Very sensitive to light's variations

## • Photometric Visual Servoing

- Doesn't rely on any tracking or matching process
- The result doesn't suffer from measurement errors
- Performs a very accurate positioning task

$$\hat{r} = \arg \min_r \sum_x (I(r, x) - I^*(x))^2$$

# Proposed Approach

Define an alignment function  $f$  as the MI between the two images.

$$\hat{\mathbf{r}} = \arg \max_{\mathbf{r}} MI(\mathbf{I}(\mathbf{r}), \mathbf{I}^*)$$

- Very accurate positioning task
- Large convergence area
- Robust to occlusions
- Robust to illumination variations
- Robust to the alignment between images acquired using different modalities

# Shannon Mutual Information – Entropy

The entropy  $H(\mathbf{I})$  is a measure of variability of a random variable given by:

$$H(\mathbf{I}) = \sum_{i=0}^{N_{c_I}} p_I(i) \log(p_I(i))$$

where:

- $i$  is a positive value of  $\mathbf{I}(x)$
- $p_I(i) = \Pr(\mathbf{I}(x) = i)$  is the probability distribution of  $i$

Since  $-\log(p_I(i))$  is a measure of uncertainty of the event  $i$ , then  $H(\mathbf{I})$  is a weighted mean of the uncertainties.

So  $H(\mathbf{I})$  is the variability of  $\mathbf{I}$

# Shannon Mutual Information – Joint Entropy

The joint entropy  $H(\mathbf{I}, \mathbf{I}^*)$  can be defined as the variability of the couple of the variables  $(\mathbf{I}, \mathbf{I}^*)$  and is given by:

$$H(\mathbf{I}, \mathbf{I}^*) = - \sum_{i=0}^{N_{c_I}} \sum_{j=0}^{N_{c_{I^*}}} p_{II^*}(i, j) \log(p_{II^*}(i, j))$$

where:

- $i$  and  $j$  are the possible values of the  $\mathbf{x}$  pixel the variables  $\mathbf{I}(\mathbf{x})$  and  $\mathbf{I}^*(\mathbf{x})$
- $p_{II^*}(i, j) = \Pr(\mathbf{I}(\mathbf{x}) = i \cap \mathbf{I}^*(\mathbf{x}) = j)$  is the joint probability distribution function

# Shannon Mutual Information – Original Mutual Information

The MI of two random variables  $I$  and  $I^*$  is given by:

$$MI(I, I^*) = H(I) + H(I^*) - H(I, I^*)$$

Where MI measures the quantity of information shared by two random variables



# Shannon Mutual Information – Link with Visual Servoing

The MI can be written with respect to the pose  $\mathbf{r}$ :

$$MI(\mathbf{r}) = MI(\mathbf{I}(\mathbf{r}), \mathbf{I}^*) = H(\mathbf{I}(\mathbf{r})) + H(\mathbf{I}^*) - H(\mathbf{I}(\mathbf{r}), \mathbf{I}^*)$$

Now, substituting these entropies, we have:

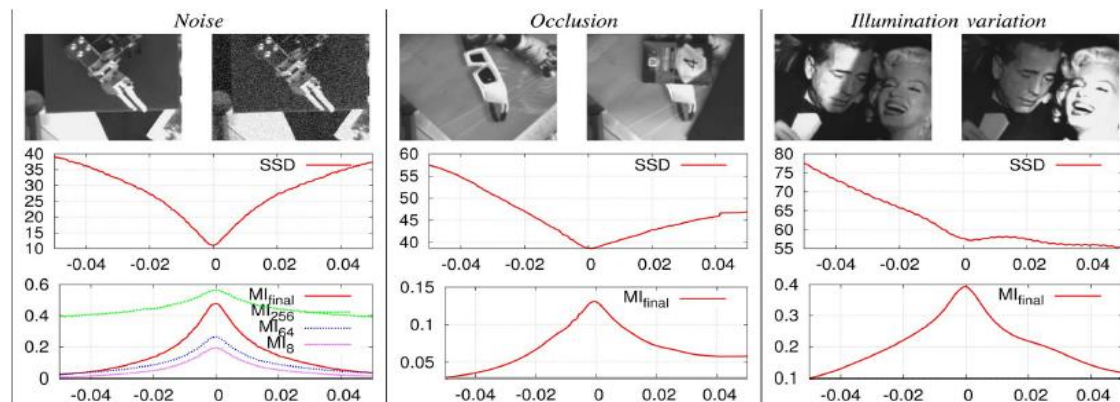
$$MI(\mathbf{r}) = \sum_{i,j} p_{II^*}(i,j,\mathbf{r}) \log \frac{p_{II^*}(i,j,\mathbf{r})}{p_I(i,\mathbf{r})p_{I^*}(j)}$$

# Adapting the Mutual Information Formulation - Histograms Binning (1)

To reduce the number of bins, since it corresponds to the maximum gray-level intensity of the image  $N_{c_I} = 255$ , the image's intensities are simply scaled as follows, where  $N_c$  is equal to 8:

$$\bar{I}(\mathbf{r}, x) = I(\mathbf{r}, x) \frac{N_c}{N_{c_I}}$$

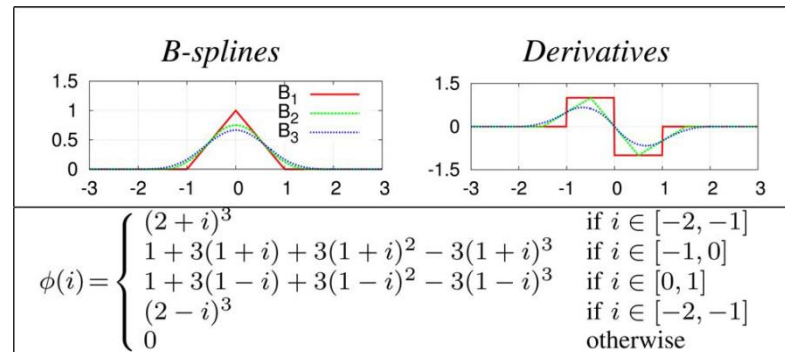
The choice of having 8 colors is given by the fact that the original definition of MI requires the computation of large histograms that are highly time consuming, not differentiable, and yields local maxima.



# Adapting the Mutual Information Formulation - Histograms Binning (2)

To calculate each probability, the following function B-spline has been chosen because of its advantage of being well adapted to histogram computation, since it doesn't require renormalization, and for optimization problems, since its derivatives is easy and inexpensive to compute:

$$\phi(i) = \begin{cases} (2+i)^3 & \text{if } i \in [-2, -1] \\ 1 + 3(1+i) + 3(1+i)^3 - 3(2+i)^3 & \text{if } i \in [-1, 0] \\ 1 + 3(1-i) + 3(1-i)^3 - 3(2-i)^3 & \text{if } i \in [0, 1] \\ (2-i)^3 & \text{if } i \in [1, 2] \\ 0 & \text{otherwise} \end{cases}$$



# Adapting the Mutual Information Formulation - Histograms Binning (3)

Using the scaled images and the B-splines function  $\phi$ , we obtain:

$$p_I(i, \mathbf{r}) = \frac{1}{N_x} \sum_x \phi(i - \bar{I}(\mathbf{r}, x))$$

$$p_{I^*}(j) = \frac{1}{N_x} \sum_x \phi(j - \bar{I}^*(x))$$

$$p_{II^*}(i, j, \mathbf{r}) = \frac{1}{N_x} \sum_x \phi(i - \bar{I}(\mathbf{r}, x)) \phi(j - \bar{I}^*(x))$$

# Mutual Information-Based Control Law (1)

The regulation of a task function  $\mathbf{e}$  to zero is done using this control law:

$$\mathbf{v} = -\lambda \widehat{\mathbf{L}}_e^+ \mathbf{e}^T$$

where:

- $\lambda$  is a positive scalar factor used to tune the convergence rate
- $\widehat{\mathbf{L}}_e^+$  is an estimation of the pseudoinverse of the interaction matrix associated with the task

The pseudoinverse can be replaced by the inverse leading to

$$\mathbf{v} = -\lambda \mathbf{H}_{MI}^{-1} \mathbf{L}_{MI}^T$$

where  $\mathbf{H}_{MI}$  is the interaction matrix of  $\mathbf{L}_{MI}$

# Mutual Information-Based Control Law (2)

The expressions of the Gradient and Hessian are

$$\mathbf{L}_{MI} = -\frac{\partial MI(I(\mathbf{r}), I^*)}{\partial \mathbf{r}} = -\sum_{i,j} \frac{\partial p_{II^*}(i, \mathbf{r})}{\partial \mathbf{r}} \left( 1 + \log \left( \frac{p_{II^*}}{p_I} \right) \right)$$

$$\mathbf{H}_{MI} = -\frac{\partial \mathbf{L}_{MI}}{\partial \mathbf{r}} = -\sum_{i,j} \frac{\partial p_{II^*}^T}{\partial \mathbf{r}} \frac{\partial p_{II^*}}{\partial \mathbf{r}} \left( \frac{1}{p_{II^*}} - \frac{1}{p_I} \right) + \frac{\partial^2 p_{II^*}}{\partial \mathbf{r}^2} \left( \frac{p_{II^*}}{p_I} \right) \quad (1)$$

In the approach proposed in the paper, the full Hessian matrix is computed using the second-order derivatives.

# Mutual Information-Based Control Law (3)

Using the previous expression regarding the joint probability:

$$\frac{\partial p_{II^*}(i, j, \mathbf{r})}{\partial \mathbf{r}} = \frac{1}{N_x} \sum_x \frac{\partial \phi}{\partial \mathbf{r}} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) \phi(j - \bar{\mathbf{I}}^*(\mathbf{x}))$$

$$\frac{\partial^2 p_{II^*}(i, j, \mathbf{r})}{\partial \mathbf{r}^2} = \frac{1}{N_x} \sum_x \frac{\partial^2 \phi}{\partial \mathbf{r}^2} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) \phi(j - \bar{\mathbf{I}}^*(\mathbf{x}))$$

$\phi$  is chosen as a third order B-spline.

# Mutual Information-Based Control Law (4)

The derivative of the function  $\phi$  respect to the pose  $\mathbf{r}$  can be written as

$$\frac{\partial \phi}{\partial \mathbf{r}} = (i - \bar{I}(\mathbf{r}, \mathbf{x})) = -\frac{\partial \phi}{\partial i} (i - \bar{I}(\mathbf{r}, \mathbf{x})) \nabla \bar{I} \mathbf{L}_x$$

With  $\nabla \bar{I} = (p_x \nabla \bar{I}_x, p_y \nabla \bar{I}_y)$  that is the image gradient expressed in the metric space that are obtained using the classical image gradients and the camera parameters  $(p_x, p_y)$ .

$\mathbf{L}_x$  is the interaction matrix that links the displacement of a point in the image plan to the camera velocity and is given by

$$\mathbf{L}_x = \begin{bmatrix} -1/Z & 0 & x/Z & xy & -(1+x^2) & y \\ 0 & -1/Z & y/Z & 1+y^2 & -xy & -x \end{bmatrix}$$



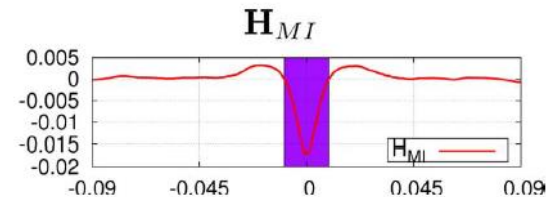
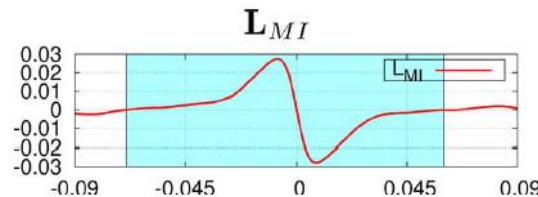
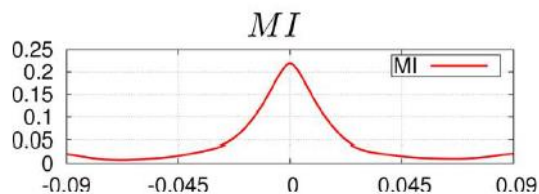
# Mutual Information-Based Control Law (5)

The second-order derivative of the  $\phi$  function is given by

$$\begin{aligned} \frac{\partial^2 \phi}{\partial \mathbf{r}^2} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) &= \frac{\partial^2 \phi}{\partial i^2} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) (\nabla \bar{\mathbf{I}} \mathbf{L}_x)^T (\nabla \bar{\mathbf{I}} \mathbf{L}_x) \\ &\quad - \frac{\partial \phi}{\partial i} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) (\nabla \bar{\mathbf{I}}_x \mathbf{H}_x + \nabla \bar{\mathbf{I}}_y \mathbf{H}_y) \\ &\quad - \frac{\partial \phi}{\partial i} (i - \bar{\mathbf{I}}(\mathbf{r}, \mathbf{x})) \mathbf{L}_x^T \nabla^2 \bar{\mathbf{I}} \mathbf{L}_x \end{aligned}$$

Where:

- $\nabla^2 \bar{\mathbf{I}} \in \mathbf{R}^{2 \times 2}$  is the gradient of  $\nabla \bar{\mathbf{I}}$  in the metric space
- $\mathbf{H}_x$  and  $\mathbf{H}_y$  are, respectively, the derivatives of the first and second line of the interaction matrix  $\mathbf{L}_x$



# Optimization approaches (1)

The optimization approach proposed in the paper is commonly known as a preconditioning approach and it consists in studying the valley shape of the cost function at convergence. Once you know the shape, it is possible to modify the steepest gradient-descent direction to make it follow the valley.

To characterize the valley shape at convergence, it's possible to simply estimate the Hessian matrix of the cost function computed at convergence.

Good assumption:

Consider that the current image at convergence is similar to the desired image, and the Hessian matrix at convergence  $\mathbf{H}_{MI}^*$  is given by (1) [slide 15] using  $\mathbf{I} = \mathbf{I}^*$ .

## Optimization approaches (2)

Since it is computed at the max of the MI function, the resulting Hessian matrix is a negative matrix. It is ideal to adapt the direction of the gradient to make it follow the valley using

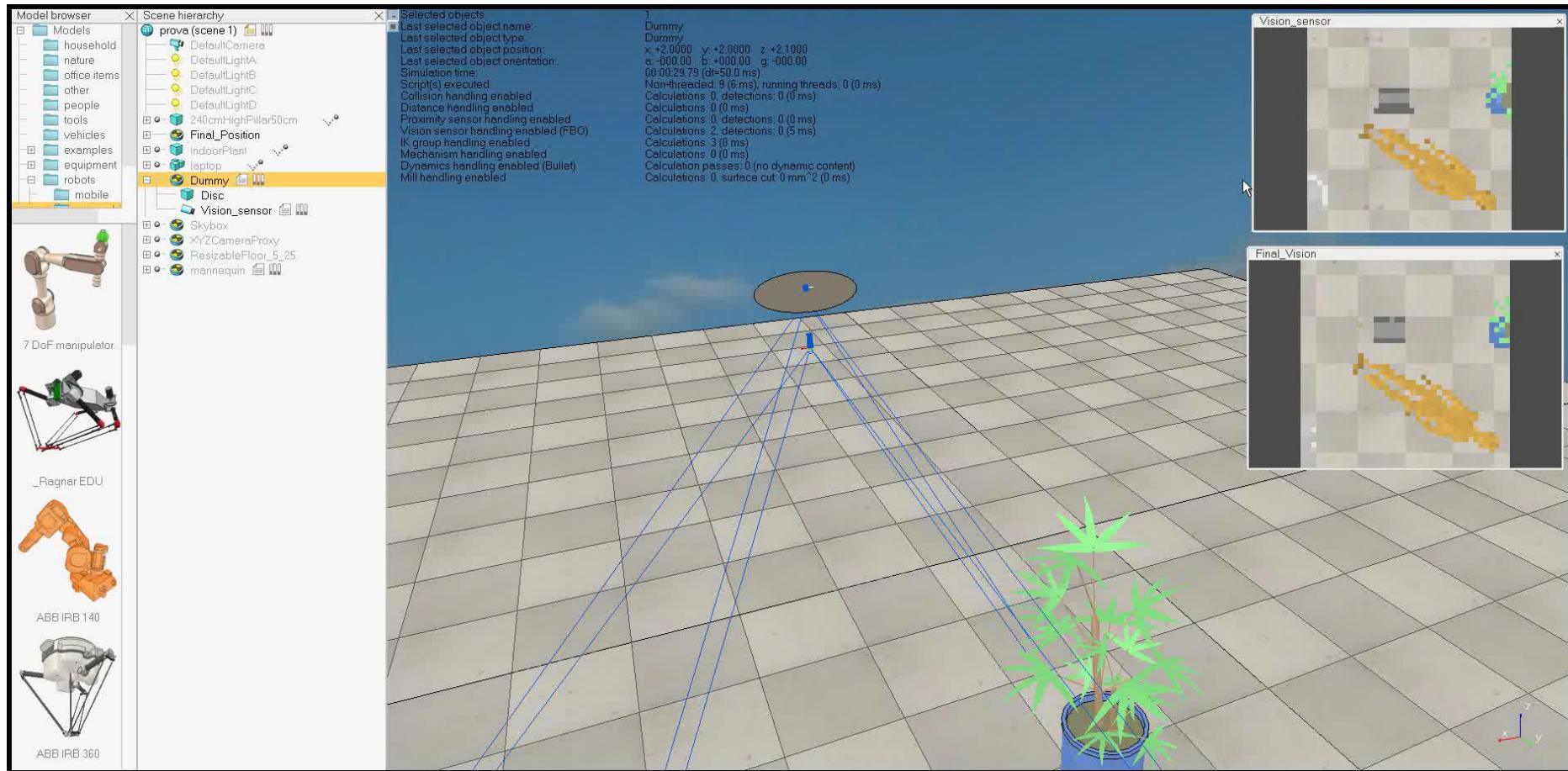
$$\mathbf{v} = -\lambda \mathbf{H}_{MI}^{*-1} \mathbf{L}_{MI}^T$$

The matrix  $\mathbf{H}_{MI}^*$  gives an ideal norm to the velocity that brings it to a null value at convergence.

With this approach, at each iteration, the computation of the interaction matrix  $\mathbf{L}_{MI}$  is the only one required.

Thanks to that, we can say that the control law computation is very fast.

# Application of the Control Law



# Real use of the control law

In the medical application this paper has the property to be robust to multimodal alignment.

Like compare two C.A.T. for more purpose like seeing the growth of a cancer mass or comparing the fracture after some time.