

ICRA'2018 Tutorial on Vision-based Robot Control



Geometric and Photometric Vision-based Robot Control: Modeling Approach

François Chaumette

Rainbow group, Inria at Irisa
Rennes, France

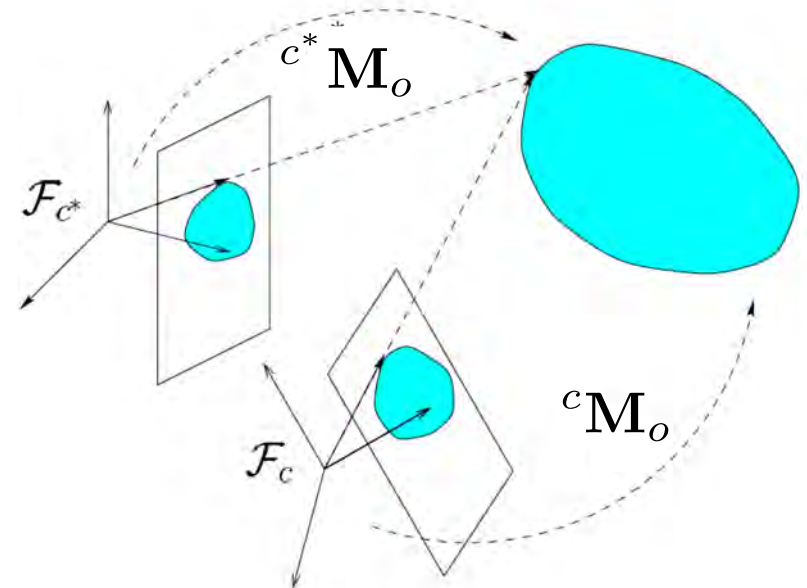
<http://team.inria.fr/rainbow>



How to control robot motion from vision?

1st basic idea: determine only once the displacement to be done
(open loop/saccade)

$${}^c M_c = {}^c M_o {}^c M_o^{-1}$$



Advantages:

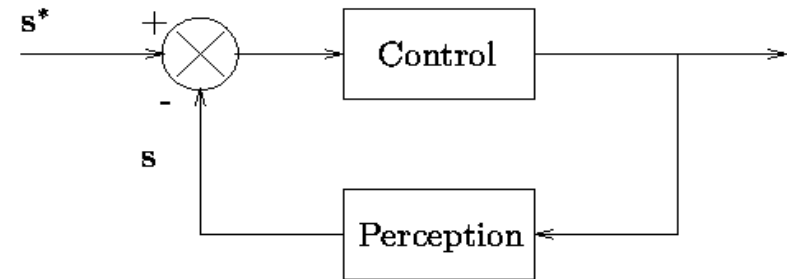
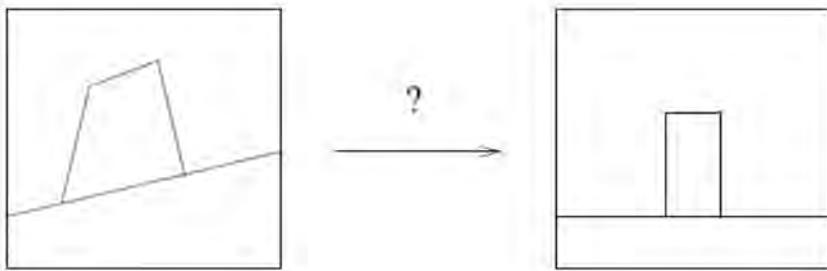
- Only one image to be processed and one very fast displacement to be achieved if the full system is perfectly calibrated

Drawbacks:

- Not robust to modeling and calibration errors
- Iterating may help, or not... Object detection for each new image

Vision-based robot control/visual servoing

Closed loop control of a dynamic system
by iterative minimization of a visual error (Lyapunov function)



Advantages:

- Positioning accuracy
- Robustness with respect to modeling and calibration errors
- Reactive to changes (target tracking)

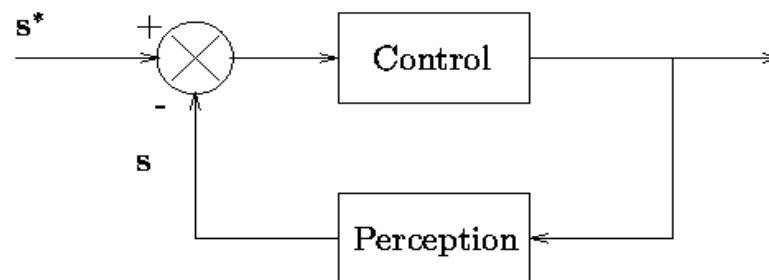
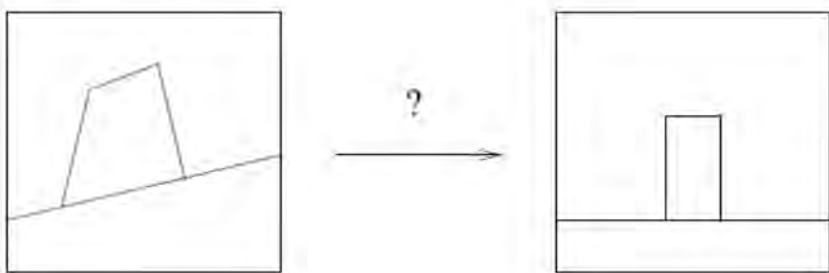
Drawbacks:

- Need many images to be processed

How to proceed?

Usual steps:

- extract and track visual measurements near video rate
- design visual features from the available measurements
- design a control scheme taking into account the system and environment constraints for an adequate system behavior (stability, robustness, ...)



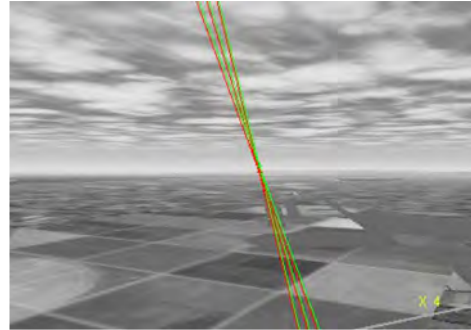
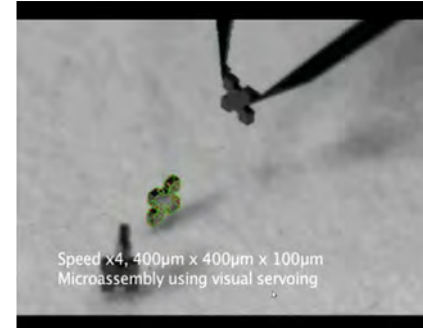
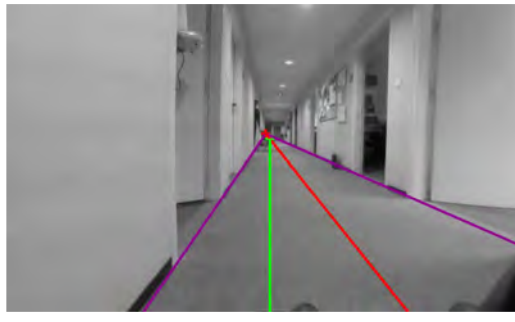
Alternative to SLAM:

- Achieve a task with the minimal information required

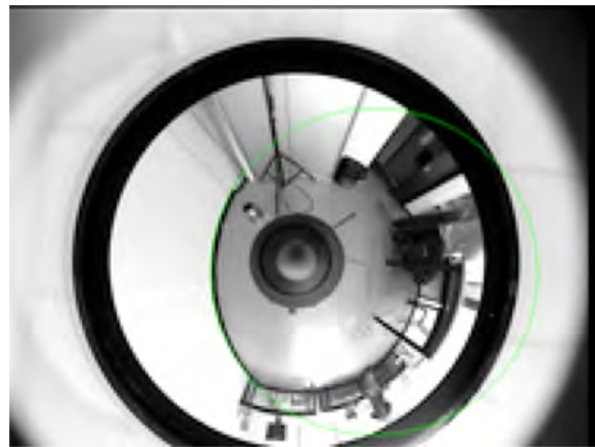
Action part for the perception-action cycle (active vision)

A wide spectrum of applications

Just need a camera and a robot



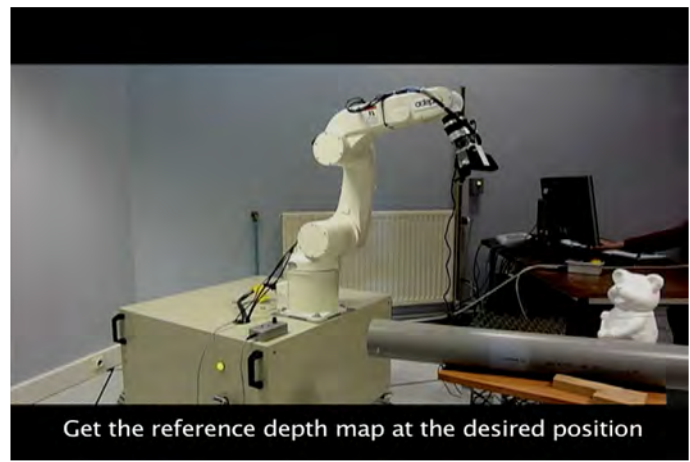
Whatever sort of vision sensor



Omnidirectional camera



2D US probe



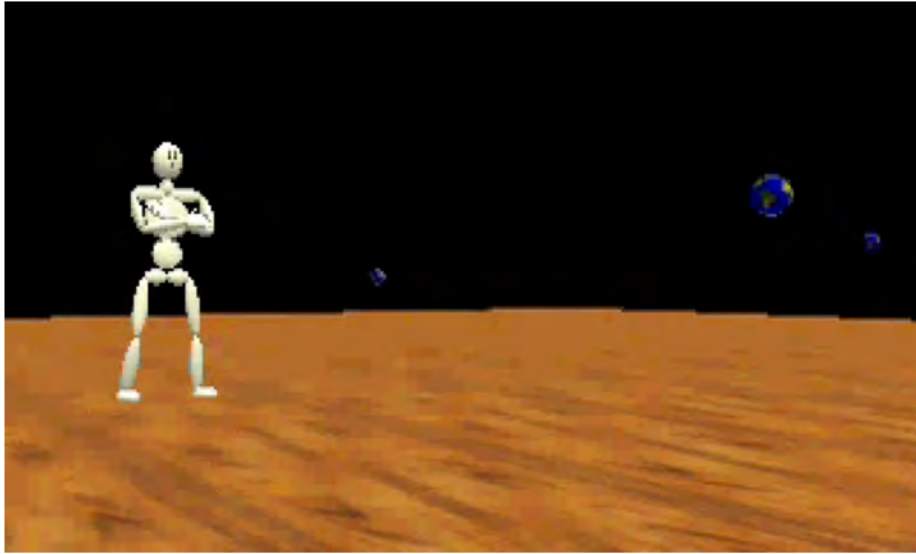
RGB-D sensor

Just need a camera

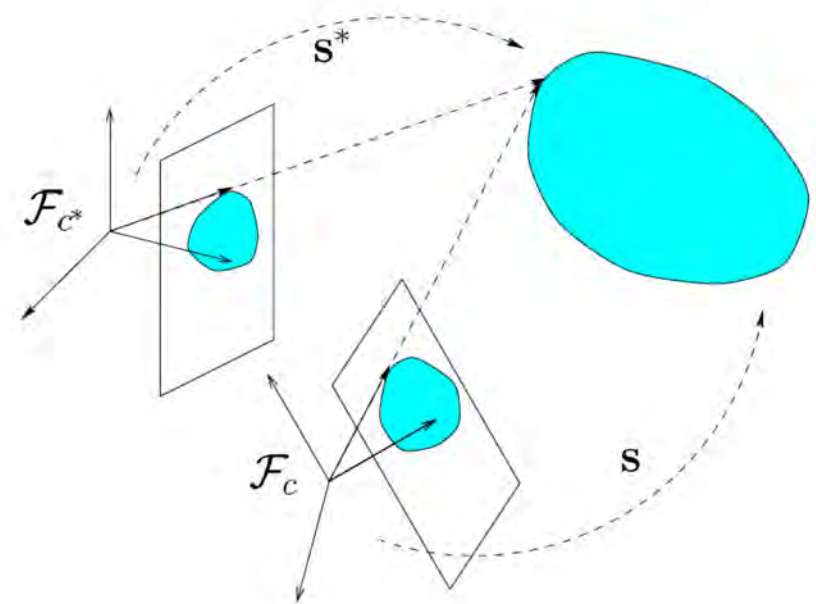
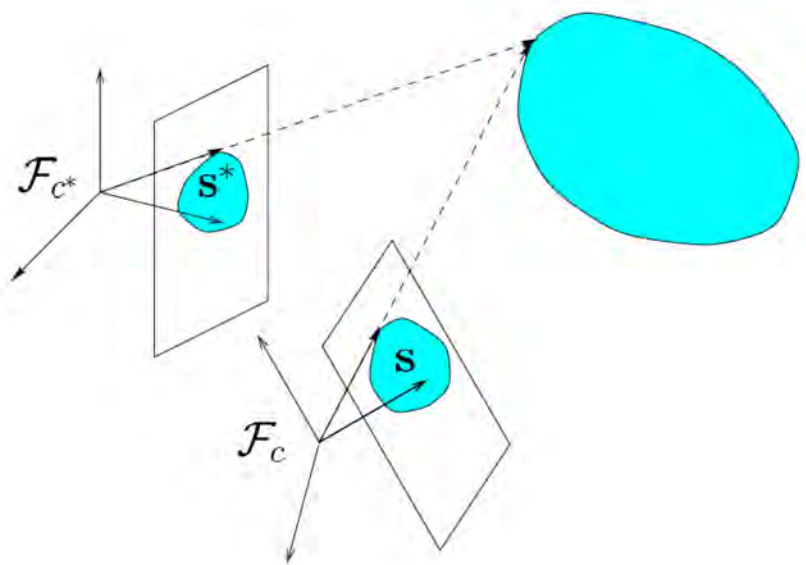


Pose estimation / 3D tracking can be formulated as Virtual Visual Servoing

Just need a computer



The basic tools: Modeling



2D visual features (IBVS) /

3D visual features (PBVS)

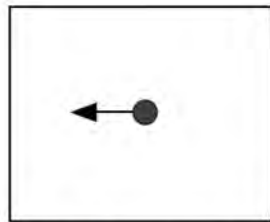
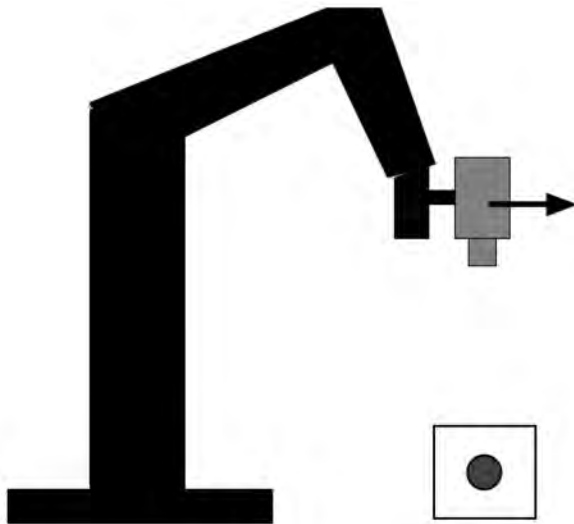
Same principle in both cases (but not same properties)

$$\text{Visual features: } s = s(\mathbf{p}(t)) \quad \Rightarrow \quad \dot{s} = \mathbf{L}_s \mathbf{v} = \mathbf{J}_s \dot{\mathbf{q}}$$

- \mathbf{L}_s : interaction matrix, \mathbf{J}_s : feature Jacobian
- $\mathbf{v} = (\mathbf{v}, \boldsymbol{\omega}) \in se_3$: instantaneous camera velocity in camera frame

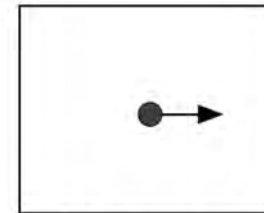
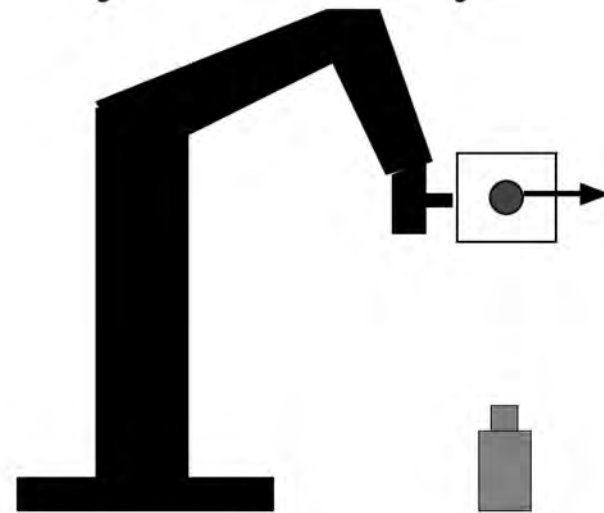
The basic tools: the feature Jacobian

Eye-in-hand system



$$\begin{aligned} \dot{\mathbf{s}} &= \mathbf{L}_s^c \mathbf{V}_n^n \mathbf{J}_n(\mathbf{q}) \dot{\mathbf{q}} + \frac{\partial \mathbf{s}}{\partial t} \\ &= \mathbf{J}_s \dot{\mathbf{q}} + \frac{\partial \mathbf{s}}{\partial t} \end{aligned}$$

Eye-to-hand system



$$\begin{aligned} \dot{\mathbf{s}} &= -\mathbf{L}_s^c \mathbf{V}_n^n \mathbf{J}_n(\mathbf{q}) \dot{\mathbf{q}} + \frac{\partial \mathbf{s}}{\partial t} \\ &= -\mathbf{L}_s^c \mathbf{V}_\emptyset^\emptyset \mathbf{V}_n^n \mathbf{J}_n(\mathbf{q}) \dot{\mathbf{q}} + \frac{\partial \mathbf{s}}{\partial t} \end{aligned}$$

The basic tools

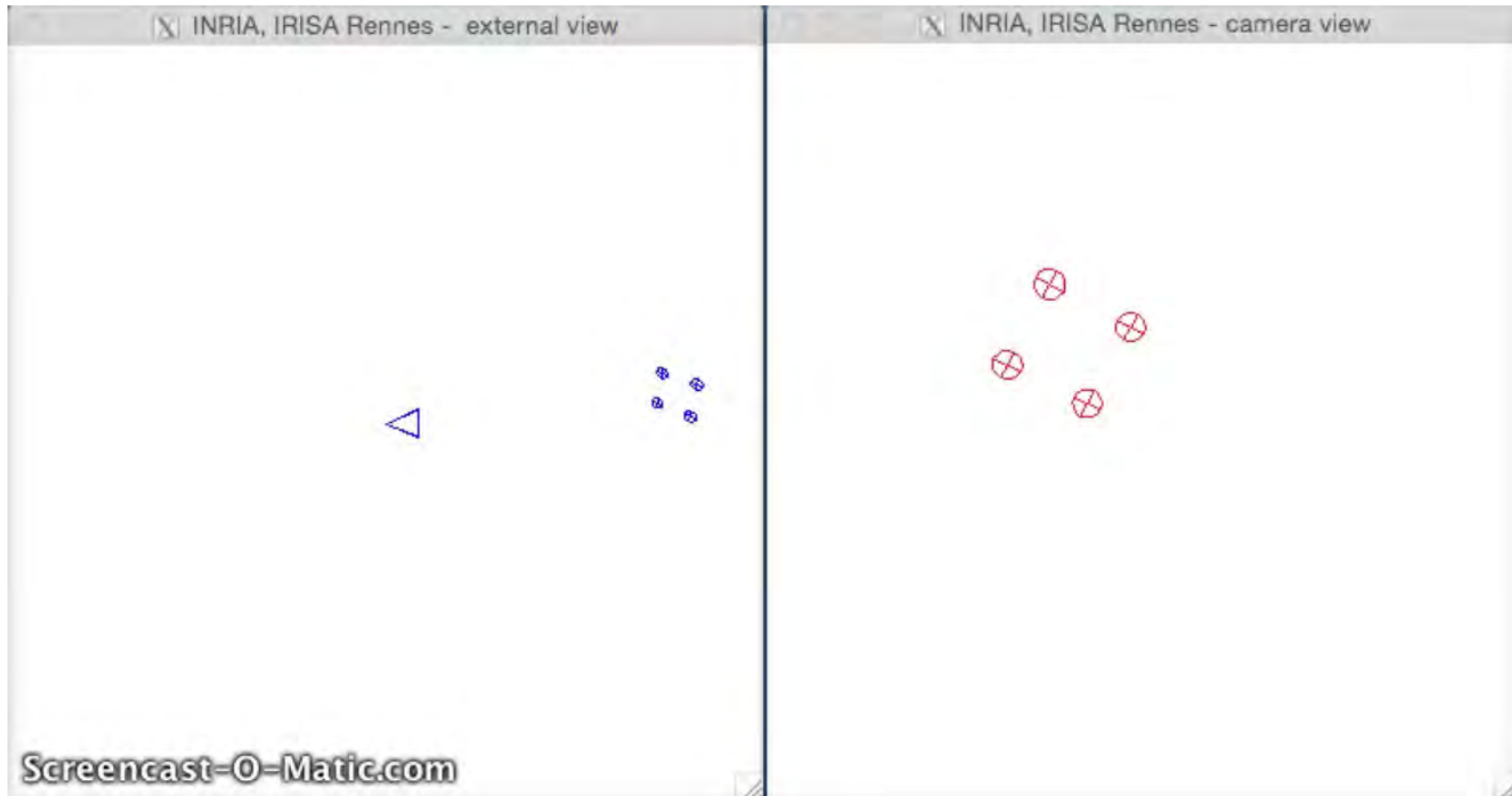
Modeling: $\dot{\mathbf{s}} = \mathbf{L}_s \mathbf{v}$

Control: $\mathbf{v} = -\lambda \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$ to try to ensure $\dot{\mathbf{s}} = -\lambda(\mathbf{s} - \mathbf{s}^*)$
(exponential decoupled decrease)

Stability analysis: $\mathcal{L} = \frac{1}{2} \|\mathbf{s} - \mathbf{s}^*\|^2$

$\dot{\mathcal{L}} = -\lambda (\mathbf{s} - \mathbf{s}^*)^T \mathbf{L}_s \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$ Usually, LAS only

Usually LAS only: potential local minimum (for 6 dof)



$$\mathbf{s} = (x_1, y_1, \dots, x_4, y_4) \quad \mathbf{v} = -\lambda \mathbf{L}_s^+ (\mathbf{s} - \mathbf{s}^*)$$

This local minimum can be avoided with another choice of \mathbf{s} or $\widehat{\mathbf{L}}_s^+$

The basic tools

Modeling: $\dot{\mathbf{s}} = \mathbf{L}_s \mathbf{v}$

Control: $\mathbf{v} = -\lambda \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$

For an image point: $\mathbf{s} = \mathbf{x} = (x, y)$

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

The depth Z_i of each point appears for the 3 translational dof (true $\forall \mathbf{s} \in 2D$)

- Can be approximated: $Z_i(t) = Z_i^*$
- Can be estimated: $Z_i(t) = \widehat{Z}_i(t)$
 - by triangulation with stereovision
 - from pose if 3D object model available
 - up to a scale factor from epipolar geometry/homography with current & desired images
 - from structure from known motion

The basic tools

Modeling: $\dot{\mathbf{s}} = \mathbf{L}_s \mathbf{v}$

Control: $\mathbf{v} = -\lambda \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$

For an image point: $\mathbf{s} = \mathbf{x} = (x, y)$

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

For same image point but with cylindrical coordinates: $\mathbf{s} = (\rho, \theta)$

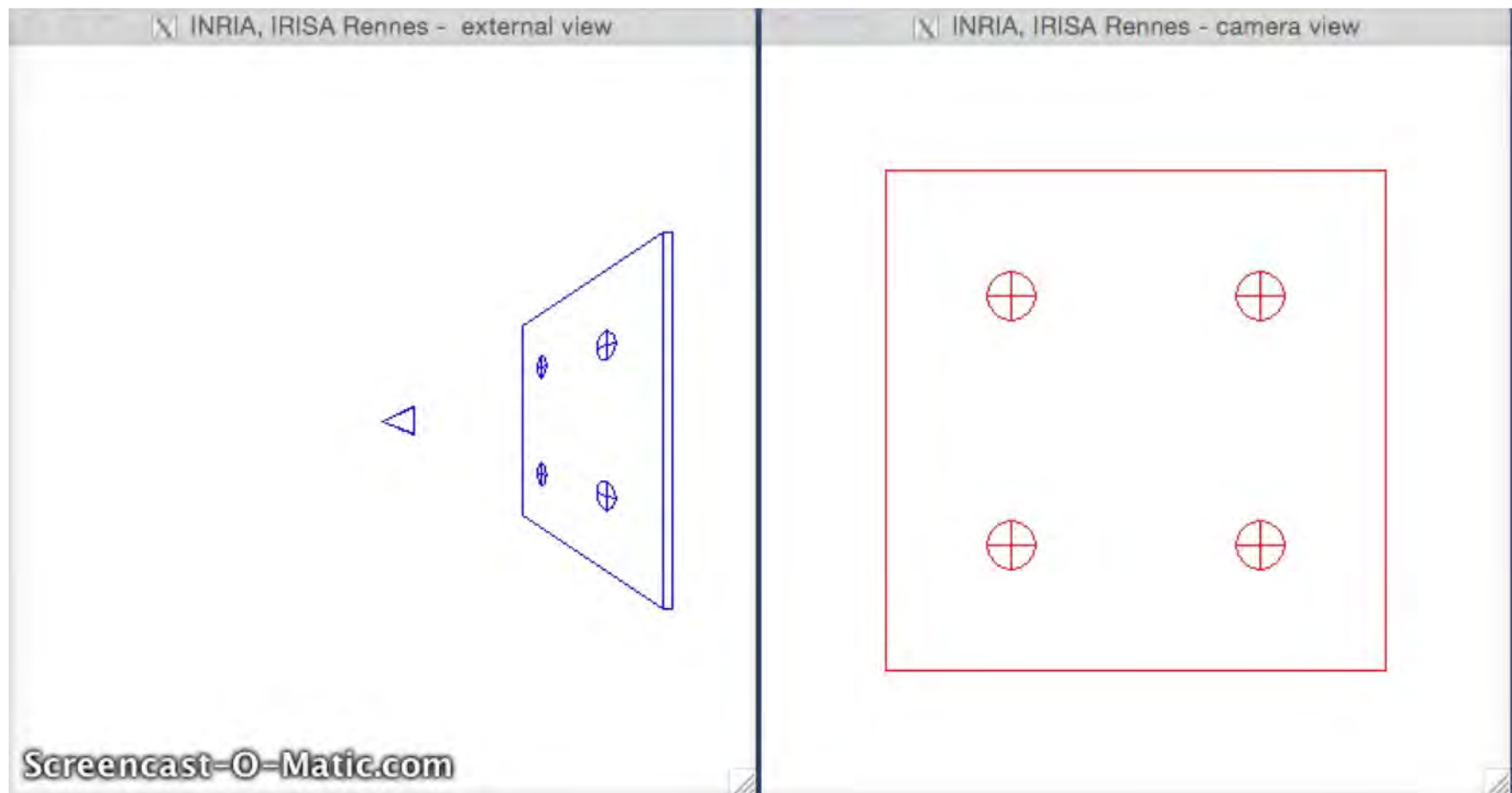
$$\mathbf{L}_{(\rho, \theta)} = \begin{bmatrix} \frac{-\cos \theta}{Z} & \frac{-\sin \theta}{Z} & \frac{\rho}{Z} & (1 + \rho^2) \sin \theta & -(1 + \rho^2) \cos \theta & 0 \\ \frac{\sin \theta}{\rho Z} & \frac{-\cos \theta}{\rho Z} & 0 & \frac{\cos \theta}{\rho} & \frac{\sin \theta}{\rho} & -1 \end{bmatrix}$$

Different choices of \mathbf{s} will induce different image & robot behaviors

Open problem for 6 dof (solved for 4 dof)

What are the visual features for an optimal behavior?

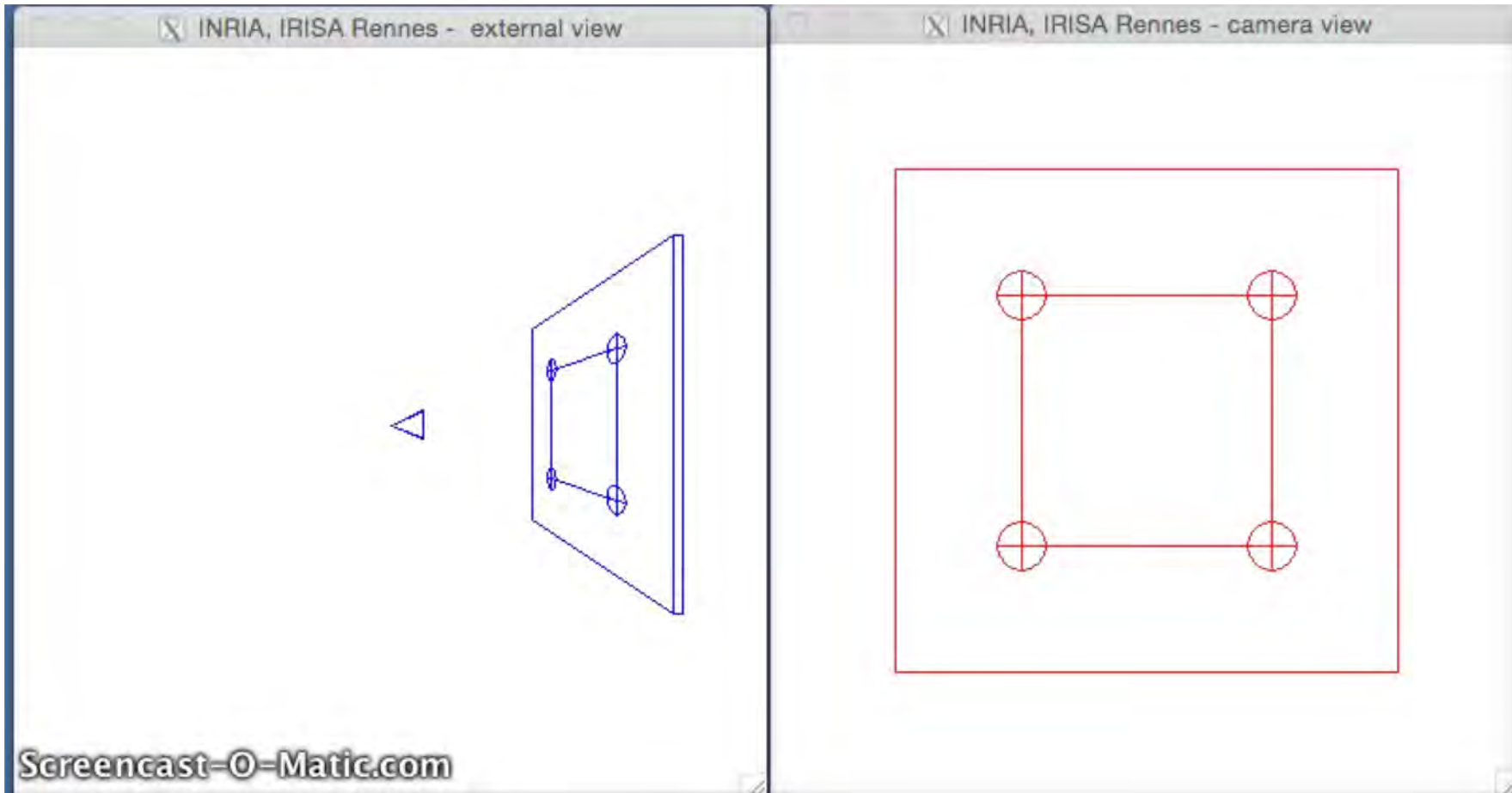
A very bad choice



$$s = (x_1, y_1, \dots, x_4, y_4)$$

What are the good visual features?

A perfect choice for this particular configuration



$$s = (\rho_1, \theta_1, \dots, \rho_4, \theta_4)$$

The basic tools

Modeling: $\dot{\mathbf{s}} = \mathbf{L}_s \mathbf{v}$

Control: $\mathbf{v} = -\lambda \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$

\mathbf{L}_s known for many visual features:

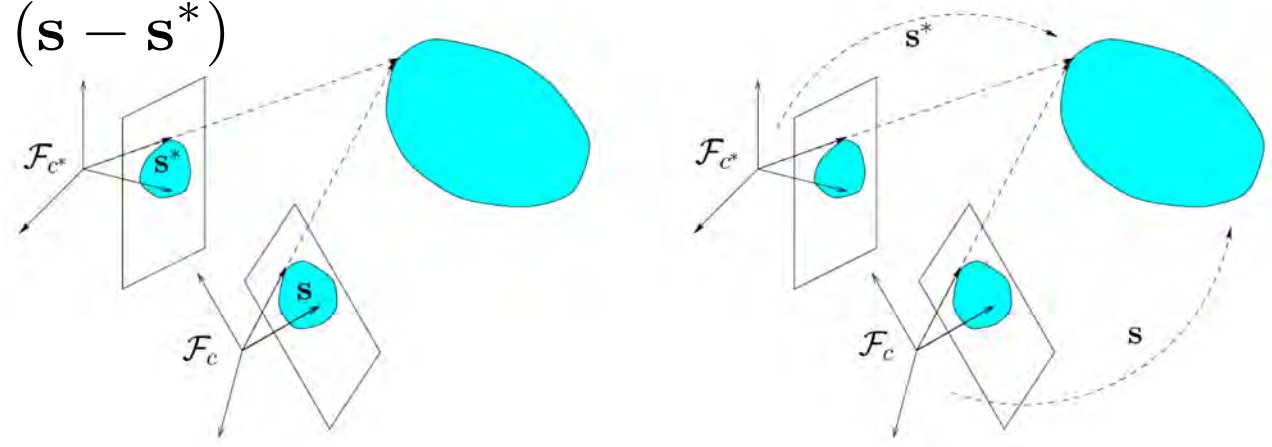
- In 3D, directly from kinematics: ${}^c \mathbf{t}_o$, ${}^{c^*} \mathbf{t}_c$, $\theta \mathbf{u}$: GAS if 3D is perfect
- In 2D:
 - Point, segment, straight line, circle, cylinder, sphere, ...
 - Moments for planar or almost planar shapes

If \mathbf{L}_s unknown, it can be estimated (off-line, on-line, by learning)
but be careful to non-linearity and stability

From your application (robot dof, object, task), search for the best choice

My 2 cents on the endless debate: IBVS vs PBVS

$$v = -\lambda \widehat{\mathbf{L}}_s^+ (s - s^*)$$



2D visual features (IBVS) / 3D visual features (PBVS)

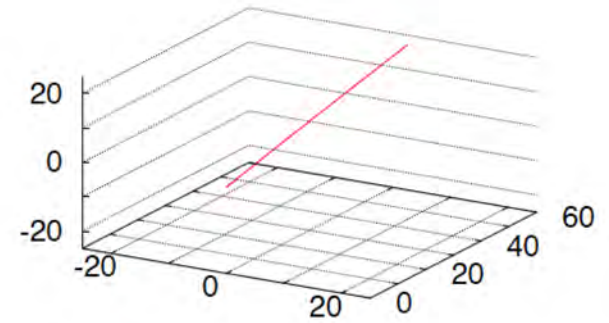
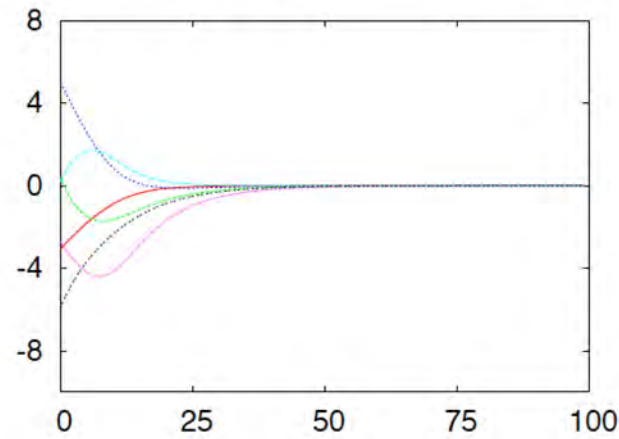
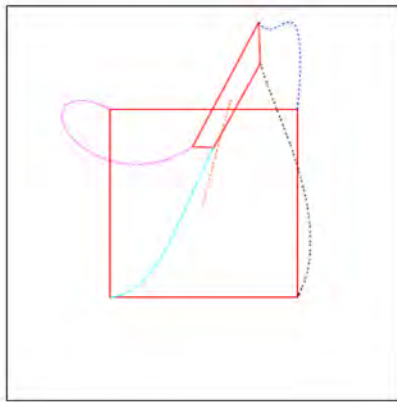
For IBVS, 3D appears in $\widehat{\mathbf{L}}_s$ but not in s

So 3D noise will affect the transient, but not the accuracy at the goal

This is not the case for PBVS

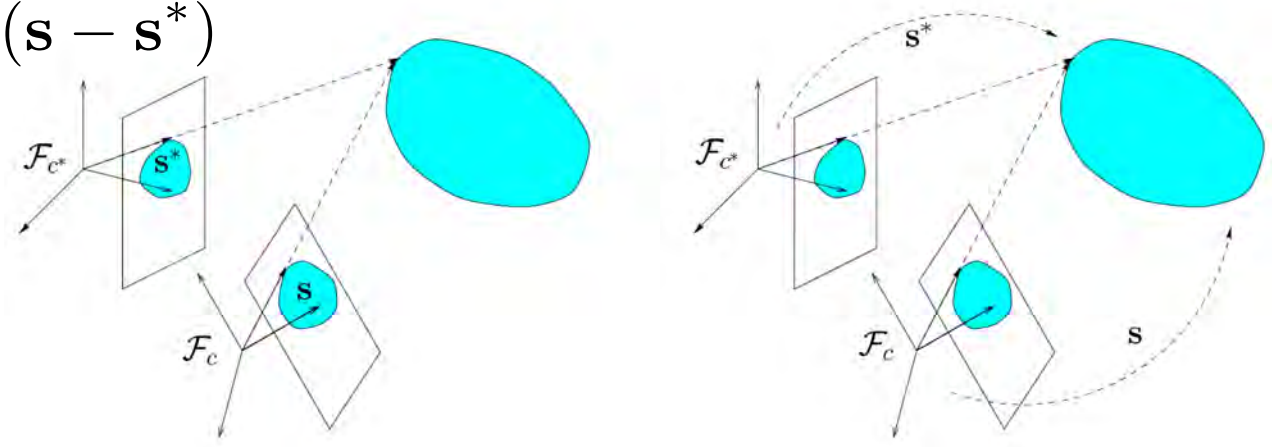
And the winner was to combine 2D and 3D visual features (2 1/2 D VS)...

Results using $s = (c^* t_c, \mathbf{x}_g, \theta u_z)$



My 2 cents on the endless debate: IBVS vs PBVS

$$\mathbf{v} = -\lambda \widehat{\mathbf{L}}_s^+ (\mathbf{s} - \mathbf{s}^*)$$



2D visual features (IBVS) / 3D visual features (PBVS)

For IBVS, 3D appears in $\widehat{\mathbf{L}}_s$ but not in \mathbf{s}

So 3D noise will affect the transient, but not the accuracy at the goal

This is not the case for PBVS

And the winner was to combine 2D and 3D visual features (2 ½ D VS)...

Now, try to design IBVS with PBVS behavior (search for \mathbf{s} such that $\mathbf{L}_s \approx \mathbf{I}$)

A new family of visual servoing: photometric VS

Remove the image processing part:

- No more extraction nor tracking visual measurements near video rate

Advantages:

- Robustness to image processing errors and noise!
- End-to-end control (here without deep learning)

Photometric/direct/dense visual servoing

Visual features: intensity of each pixel $s = \mathbf{I}(\mathbf{x}(t))$

\mathbf{I}^*

\mathbf{I}

$\mathbf{I} - \mathbf{I}^*$



Modeling: $\mathbf{L}_I = -\nabla \mathbf{I}_x \mathbf{L}_x$ (function of the image content)

$$\mathcal{L} = \frac{1}{2} \|\mathbf{I} - \mathbf{I}^*\| \text{ highly non linear}$$

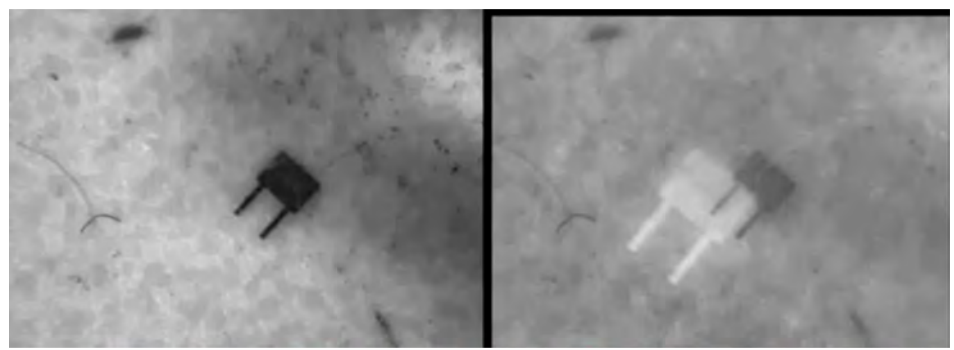
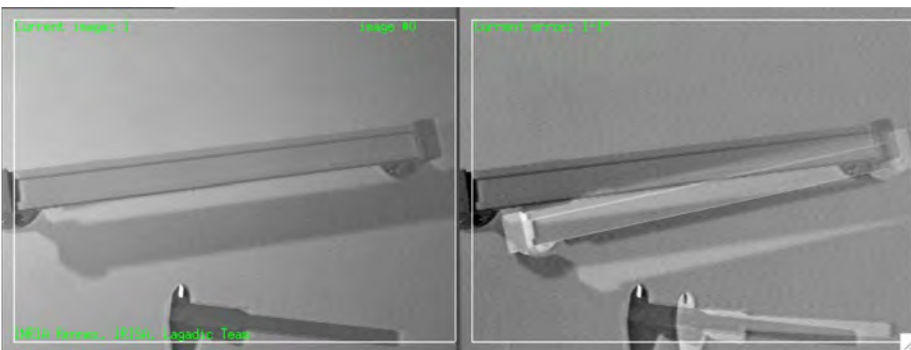
Drawbacks: small convergence domain, strange robot trajectory

But no feature extraction, tracking nor matching
+ excellent positioning accuracy

Photometric visual servoing

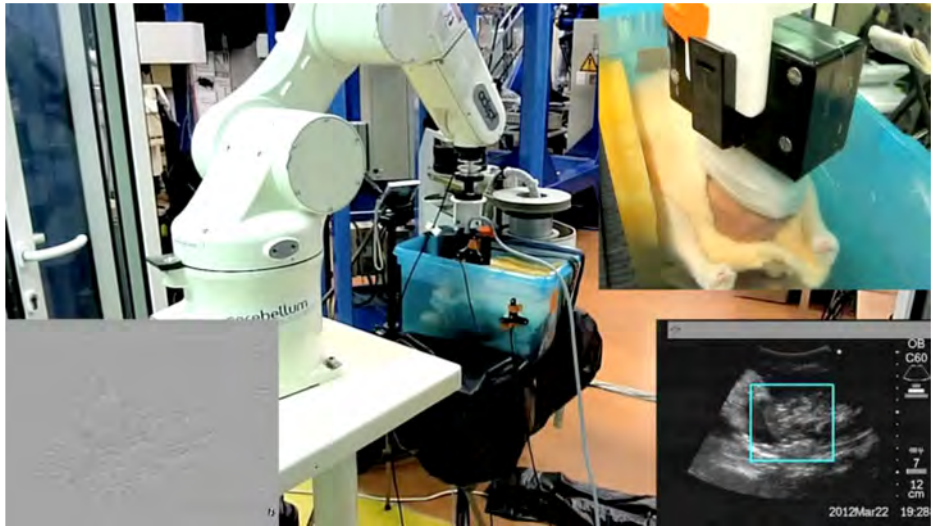
Robustness to global illumination changes by using $s = (\mathbf{I} - \bar{\mathbf{I}}) / \sigma_{\mathbf{I}}^2$

Robustness to outliers (occlusion) by using $s = \rho_{\mathbf{I}} \mathbf{I}$

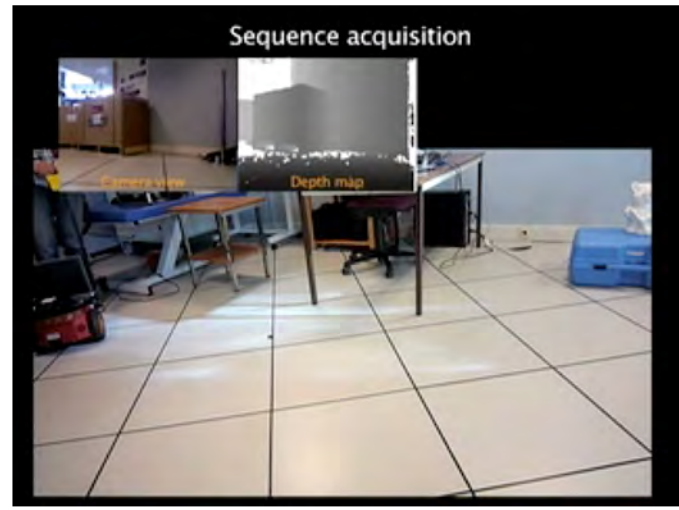
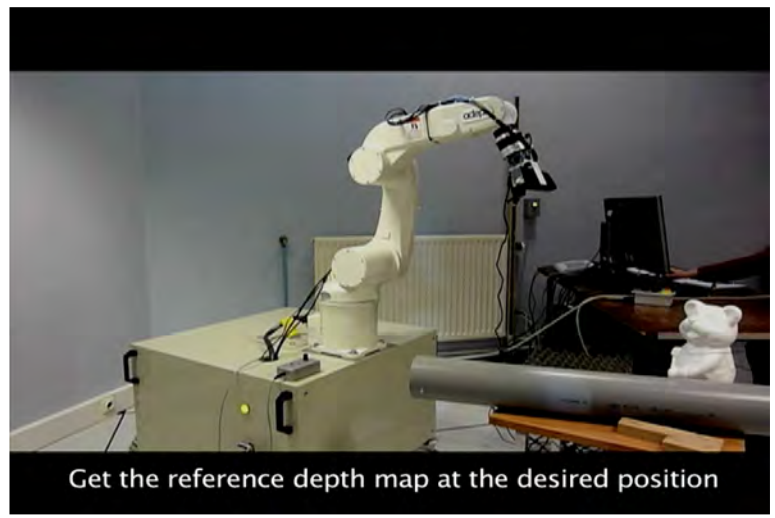


Accuracy < 0.1 μm

Similar on ultrasound images



Similar on depth map from RGB-D sensor: $s = \rho_Z Z$



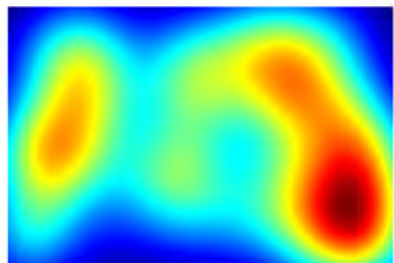
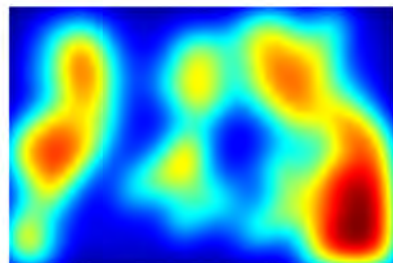
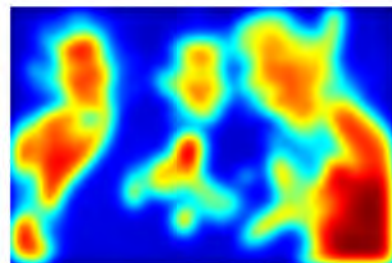
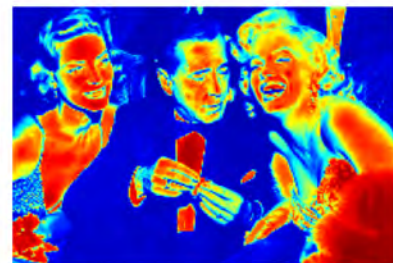
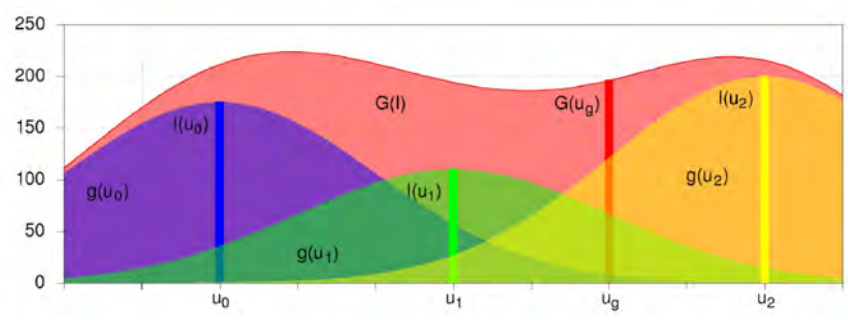
In the same spirit:

- RGB components, spatial gradient image
- Sum of conditional variance: $\mathcal{L} = \|\mathbf{I}(\mathbf{x}) - \hat{\mathbf{I}}(\mathbf{x})\|$ with $\hat{\mathbf{I}}(\mathbf{x}) = \epsilon(\mathbf{I}(\mathbf{x}), \mathbf{I}^*(\mathbf{x}))$
- Maximize mutual information between current and desired image
- Histogram-based visual servoing
- Wavelet
- ...

Mixture of Gaussians

Enlarge the convergence domain

$$G(\mathbf{u}_g, \lambda_g) = \sum_{\mathbf{u}_i \in \mathbf{I}} I(\mathbf{u}_i) \exp \left(-\frac{(u_g - u_i)^2 + (v_g - v_i)^2}{2\lambda_g^2} \right)$$



$\lambda_g = 0.1$

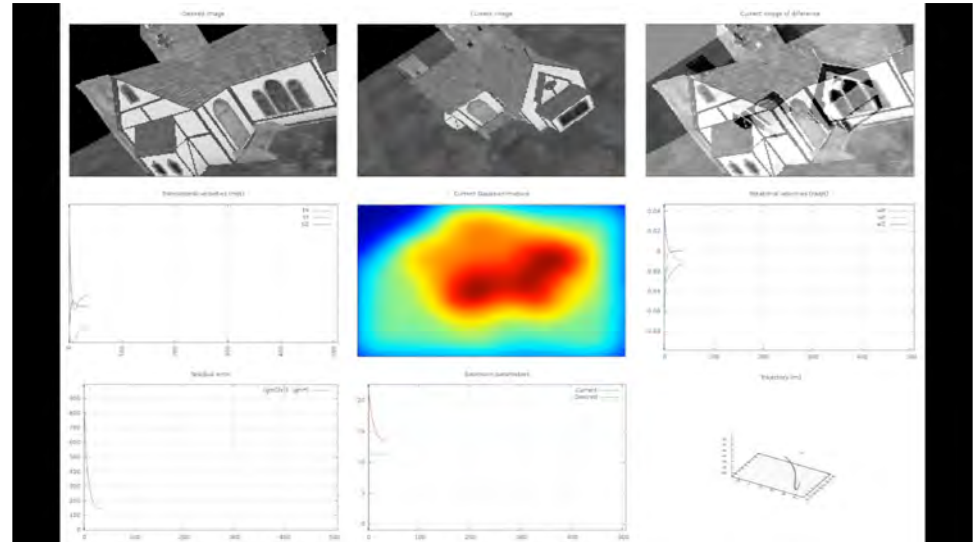
$\lambda_g = 5$

$\lambda_g = 10$

$\lambda_g = 20$

Control simultaneously

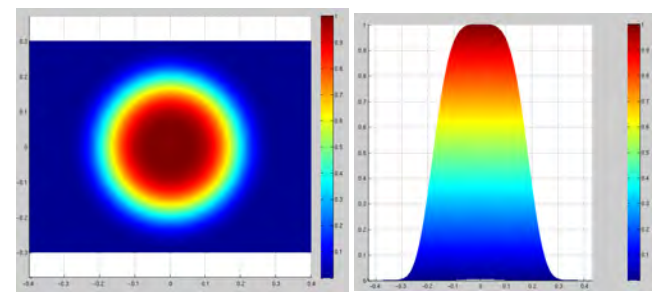
- the camera motion
- the expansion parameter λ_g (large to small)



Photometric moments

Going back to geometric features for enlarging the convergence domain and improving the robot trajectory

$$m_{pq} = \iint_{\pi} x^p y^q w(\mathbf{x}) I(\mathbf{x}, t) dx dy$$



Then select adequate moments (area, cog, main orientation, ...)



I*



I



I - I*

To go further

- Target tracking
 - PI controller
 - Estimate, predict and compensate the target motion (feed forward)
- Consider constraints:
 - visibility, occlusion, obstacles
 - joint limits, singularities
 - dynamics: non holonomy, under-actuation
 - Path planning in the image, optimal control, MPC
 - Redundancy, task sequencing, stack of tasks
- Multi sensor-based control
 - Modeling, fusion

Thanks for your attention

Acknowledgments: Lagadic colleagues and the VS worldwide community

