

## IMAGE-BASED STEREO VISUAL SERVOING: 2D VS 3D FEATURES

E. Cervera \* F. Berry \*\* P. Martinet\*\*

\* *Robotic Intelligence Laboratory, Jaume-I University, 12071  
Castelló, Spain*

\*\* *LASMEA - GRAVIR, Blaise Pascal University of  
Clermont-Ferrand, 63177 Aubière - Cedex, France*

**Abstract:** This paper presents a visual servoing approach based on stereo vision. The pair of cameras is mounted on the end-effector of the manipulator arm. Theoretical developments are presented using either raw pixel coordinates, or 3D coordinates estimated from image features. No geometrical model is needed. The experimental setup is challenging: large rotations are involved, images are noisy, and cameras are coarsely calibrated. In this setup, the trajectory of the end-effector differs notably, sometimes leading the arm near its joint range limits. Experimental results demonstrate that using pixel coordinates is disadvantageous, compared with 3D coordinates estimated from the same pixel data.

**Keywords:** Robotic manipulators, Stereo vision

### 1. INTRODUCTION

Stereo visual information has been commonly considered as an alternative way to recover the depth, in the modeling phase of a vision system. The application of stereo vision in visual servoing was pioneered by Maru *et al.* (Maru *et al.*, 1993), and recent works (Hager *et al.*, 1995) (Lamiroy *et al.*, 2000) have awakened new interests, considering mainly the robustness and precision aspects.

Stereo visual servoing offers some advantages over the classical monocular 2D and 3D visual servoing approaches. Depth information can be recovered without need of any geometrical model of the observed object. It should be noted that even in 2D visual servoing, this information is needed for the computation of the image jacobian.

As pointed out in (Lamiroy *et al.*, 2000), a number of singularities exists in monocular visual servoing, making visual control impossible near those configurations. These singularities can be avoided by using a stereo rig, thus requiring less strict camera calibration.

The main goal of this paper is the study of image-based stereo visual servoing. We experimentally show that using 3D coordinates (estimated from the stereo images) in the feature vector performs better than using raw 2D image coordinates. In our experimental setup, the stereo rig is mounted on the end-effector of the arm. The programmed manipulation task is quite challenging: large rotations are involved, pixel noise is high, and camera calibration is coarse.

The rest of this paper is organized as follows: first, we consider the modeling of two stereo images of a set of points, in the simplified case where cameras are aligned.

Next, we develop visual control with three different features: in the first one, raw pixel coordinates are used. This is the so-called *image based approach* (Espiau *et al.*, 1992). Care must be taken with the definition of the coordinates frame of the cameras and the end-effector.

Image-based 3D features are then introduced: estimated coordinates, and a combination of pixel

data and stereo disparity. We show that this third approach exhibits the same nice properties as using coordinates, with regard to the end-effector trajectory.

Finally, we present experimental results of the presented approaches, with a comparison of image feature errors, the velocity screw, and the trajectory of the end-effector.

It should be noted that, in all of the approaches, the only source of information is the stereo rig. Thus, all the 3D information is estimated from these measurements, as well as from the intrinsic and extrinsic camera parameters (which are roughly known). Our interest is to compare the approaches to test whether there exists an advantage in using either the raw signals or the computed 3D features.

## 2. STEREO OBSERVATION OF A SET OF POINTS

Our setup consists of a stereo rig mounted on the end-effector of the manipulator. Let us define  $\mathcal{F}_e$  as the control frame attached to the end-effector,  $\mathcal{F}_l$  as the frame attached to the left camera, and  $\mathcal{F}_r$  as the frame attached to the right camera.

In this work, a segmented target defined by 5 points is considered. The corresponding raw feature vector is defined by

$$\underline{\mathbf{s}} = [u_1^l, v_1^l, u_1^r, v_1^r, \dots, u_5^l, v_5^l, u_5^r, v_5^r]^T \quad (1)$$

where  $\underline{\mathbf{U}}_i^l = (u_i^l, v_i^l)^T$  and  $\underline{\mathbf{U}}_i^r = (u_i^r, v_i^r)^T$  are the image coordinates of the  $i^{\text{th}}$  point, observed by the left and right cameras respectively and  $\underline{\mathbf{s}}_i$  is the  $i^{\text{th}}$  subvector of  $\underline{\mathbf{s}}$  such  $\underline{\mathbf{s}}_i = (\underline{\mathbf{U}}_i^l \ \underline{\mathbf{U}}_i^r)$ .

In our case, we consider a simplified configuration where both cameras are parallel with identical focal lengths ( $F_u, F_v$ ) and the control frame  $\mathcal{F}_e$  is located at the center of the both frames (Fig 1). Both cameras are aligned along the  $x$ -axis and the distance between them is  $b$ .

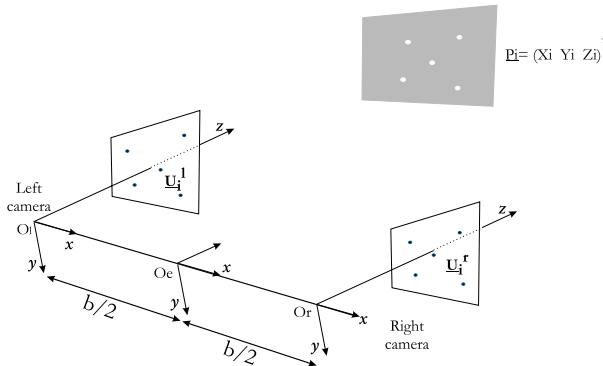


Fig. 1. The simplified configuration of our system.

Thus, the relationship between pixel and 3D coordinates is

$$\begin{cases} u_i^l = \frac{F_u X_i + F_u b/2}{Z_i} + u_0 \\ v_i^l = \frac{F_v Y_i}{Z_i} + v_0 \\ u_i^r = \frac{F_u X_i - F_u b/2}{Z_i} + u_0 \\ v_i^r = \frac{F_v Y_i}{Z_i} + v_0 \end{cases} \quad (2)$$

and the coordinates of the observed point can be easily deduced as

$$\hat{\underline{\mathbf{P}}}_i = \begin{pmatrix} \hat{X}_i \\ \hat{Y}_i \\ \hat{Z}_i \end{pmatrix} = \begin{pmatrix} \frac{b(u_i^l + u_i^r - 2u_0)}{2(u_i^l - u_i^r)} \\ \frac{bF_u(v_i^r + v_i^l - 2v_0)}{2(u_i^l - u_i^r)F_v} \\ \frac{bF_u}{u_i^l - u_i^r} \end{pmatrix} \quad (3)$$

These values are roughly estimated or are taken directly from their nominal values. No explicit calibration procedure has been undertaken.

## 3. VISUAL FEATURES

The essence of visual servoing is the computation of the matrix of derivatives (the jacobian) of the visual feature vector with respect to the velocity screw. Using the raw pixel data or the estimated 3D point coordinates is a matter of choice. Both approaches require an estimation of camera parameters. However, the resulting dynamic properties of the task may differ. In this section, we present the theoretical bases of both approaches, and a third feature vector which uses the stereo disparity, without fully estimating the real 3D coordinates.

### 3.1 Stereo 2D point

The feature vector is the raw image information (Eq. 1) and the jacobian matrix is

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_1^l \mathbf{M}_e^l \\ \mathbf{L}_1^r \mathbf{M}_e^r \\ \vdots \\ \mathbf{L}_5^l \mathbf{M}_e^l \\ \mathbf{L}_5^r \mathbf{M}_e^r \end{pmatrix} \quad (4)$$

where  $\mathbf{L}_i^l$  and  $\mathbf{L}_i^r$  are the interaction matrices for  $i^{\text{th}}$  point, relative to the left and right cameras respectively, as defined by Espiau *et al.* (Espiau

et al., 1992)

$$\mathbf{L}_i = \begin{pmatrix} -\frac{F_u}{z} & 0 & \frac{u_i}{z} & \frac{u_i v_i}{F_v} & -F_u - \frac{u_i^2}{F_u} & \frac{v_i F_u}{F_v} \\ 0 & -\frac{F_v}{z} & \frac{v_i}{z} & F v_+ \frac{v_i^2}{F_v} & -\frac{u_i v_i}{F_u} & -\frac{u_i F_v}{F_u} \end{pmatrix} \quad (5)$$

The dimension of the final image Jacobian  $\mathbf{L}$  is  $20 \times 6$ .  $\mathbf{M}_e^l$  and  $\mathbf{M}_e^r$  are the transformation matrices of the screw between the left and right camera frames and the end-effector frame. Given frames  $\mathcal{F}_e$  and  $\mathcal{F}_j$ , the relationship between the kinematic screws  $\mathbf{v}$  is

$$\underline{\mathbf{v}}^j = \mathbf{M}_e^j \underline{\mathbf{v}}^e \quad (6)$$

where the transformation matrix  $\mathbf{M}_e^j$  is

$$\mathbf{M}_e^j = \begin{pmatrix} \mathbf{R}_e^j & [\underline{\mathbf{t}}_e^j]_{\times} \mathbf{R}_e^j \\ \mathbf{O}_3 & \mathbf{R}_e^j \end{pmatrix} \quad (7)$$

It can be shown that the resulting interaction matrix (4) is the same as that obtained by Maru et al. (Maru et al., 1993).

### 3.2 Estimated 3D point

Since the 3D coordinates of the observed point can be computed from the image data (and an estimation of the extrinsic and intrinsic parameters of the cameras), they can also be used in the control law.

Thus, the feature vector consists of the estimated coordinates (eq. 3) and the jacobian matrix is

$$\mathbf{L} = \begin{pmatrix} -\mathbf{I}_3 & [\hat{\mathbf{P}}_1]_{\times} \\ \vdots & \\ -\mathbf{I}_3 & [\hat{\mathbf{P}}_5]_{\times} \end{pmatrix} \quad (8)$$

The main advantage of using 3D features is the linearity of the jacobian matrix. As a result, some theoretical properties of the trajectory of the end-effector can be obtained. Effectively, Cervera and Martinet (Cervera and Martinet, 1999) demonstrated, for a feature vector composed of a rather general set of 3D points, that the velocity screw of the camera is

$$\underline{\mathbf{v}} = -\lambda \begin{bmatrix} (\mathbf{P}_g^* - \mathbf{P}_g) + [\mathbf{P}_g]_{\times} \mathbf{R} \underline{\mathbf{u}} \sin \theta \\ \mathbf{R} \underline{\mathbf{u}} \sin \theta \end{bmatrix} \quad (9)$$

where  $\mathbf{P}_g$  is the center of gravity of the set of points,  $\mathbf{R}$  is the rotation between a Cartesian frame defined by the points and the end-effector frame, and  $\underline{\mathbf{u}}\theta$  are the axis and angle corresponding to the rotation matrix  $\mathbf{R}^T \mathbf{R}^*$ , that is, the rotation between the current and desired orientation of the set of points.

In addition, the center of gravity of the set of points translates along a *straight line trajectory* from its initial to its final position in the camera frame. As a consequence, the features are most likely to remain in the camera field of view during the whole task.

### 3.3 2D points and disparity

Instead of using the estimated 3D coordinates, we have experimented with the direct 2D image features, and the stereo disparity of the  $i^{th}$  point ( $u_i^l - u_i^r$ ). In the following control law, the feature vector is defined as

$$\underline{\mathbf{s}} = \left( \frac{u_1^l + u_1^r}{u_1^l - u_1^r}, \frac{v_1^l + v_1^r}{u_1^l - u_1^r}, \frac{1}{u_1^l - u_1^r} \dots \dots \frac{u_5^l + u_5^r}{u_5^l - u_5^r}, \frac{v_5^l + v_5^r}{u_5^l - u_5^r}, \frac{1}{u_5^l - u_5^r} \right)^T \quad (10)$$

It can be shown that this vector results from a linear combination of the 3D coordinates of the corresponding 3D point:

$$\underline{\mathbf{s}} = \left( (A \hat{\mathbf{P}}_1)^T \dots (A \hat{\mathbf{P}}_5)^T \right)^T \quad (11)$$

where

$$A = \begin{pmatrix} \frac{2}{b} & 0 & \frac{2u_0}{bF_u} \\ 0 & \frac{2F_v}{bF_u} & \frac{2v_0}{bF_u} \\ 0 & 0 & \frac{1}{bF_u} \end{pmatrix}$$

and the resulting Jacobian matrix for one point is as shown in Equation 12.

The interest in using this model is twofold: the 3D coordinates need not to be estimated, and the jacobian matrix is linear with respect to  $\underline{\mathbf{s}}$ . Effectively, as shown in (Cervera and Martinet, 1999), the jacobian matrix (12) can be expressed as

$$\mathbf{L} = \begin{pmatrix} -\mathbf{A} & \mathbf{A} [\mathbf{A}^{-1} \underline{\mathbf{s}}_1]_{\times} \\ \vdots & \\ -\mathbf{A} & \mathbf{A} [\mathbf{A}^{-1} \underline{\mathbf{s}}_5]_{\times} \end{pmatrix} \quad (13)$$

where  $\underline{\mathbf{s}}_i$  is the  $i^{th}$  element of  $\underline{\mathbf{s}}$  such  $\underline{\mathbf{s}}_i = (A \hat{\mathbf{P}}_i)^T$ . Additionally, some theoretical results from 3D points still hold for any linear combination: though the velocity screw (Eq. 9) is valid for small angles only, the trajectory of the center of gravity of the set of points *still translates along a straight path* during the task (see (Cervera and Martinet, 1999) for details).

$$\mathbf{L}_i = \begin{pmatrix} -\frac{2}{b} & 0 & -\frac{2u_0}{bF_u} & -\frac{u_0(v_i^l + v_i^r - 2v_0)}{F_v(u_i^l - u_i^r)} & \frac{u_0(u_i^l + u_i^r - 2u_0) - 2F_u^2}{F_u(u_i^l - u_i^r)} & \frac{F_u(v_i^l + v_i^r - 2v_0)}{F_v(u_i^l - u_i^r)} \\ 0 & -\frac{2F_v}{bF_u} & -\frac{2v_0}{bF_u} & -\frac{v_0(v_i^l + v_i^r - 2v_0) - 2F_v^2}{F_v(u_i^l - u_i^r)} & \frac{v_0(u_i^l + u_i^r - 2u_0)}{F_u(u_i^l - u_i^r)} & -\frac{F_v(u_i^l + u_i^r - 2u_0)}{F_u(u_i^l - u_i^r)} \\ 0 & 0 & -\frac{1}{bF_u} & -\frac{v_i^l + v_i^r - 2v_0}{2F_v(u_i^l - u_i^r)} & \frac{u_i^l + u_i^r - 2u_0}{2F_u(u_i^l - u_i^r)} & 0 \end{pmatrix} \quad (12)$$

#### 4. EXPERIMENTAL RESULTS

The mobile manipulator of the Robotic Intelligence Lab consists of a Nomad XR4000 platform and a Mitsubishi PA-10 arm. Attached to the end-effector of the arm is a stereo rig with two miniature CMOS color cameras, linked to two video boards which deliver the visual features at video rate (30 Hz).

The following table gives the estimation of the parameters (intrinsic and extrinsic) of both cameras, as used in the experiments.

$F_u$	$F_v$	$b$
300	450	118mm

These are nominal values, and no explicit calibration procedure has been carried out. Nevertheless, the system is robust with respect to this approximation.

The target object consists of four co-planar points located at the vertices of an 11cm square and the fifth point is located at the center of the square.

The velocity screw is computed from the pseudo-inverse of the jacobian matrix (Espiau *et al.*, 1992):

$$\underline{\mathbf{v}} = -\lambda \mathbf{L}^+(\underline{\mathbf{s}} - \underline{\mathbf{s}}^*) \quad (14)$$

with  $\lambda$  set to 0.5 in all the experiments. This value was chosen heuristically, being large enough for a relatively fast yet stable motion.

Image measurements are noisy, since the experiments are carried out in a standard office environment, without any special illumination. As a result, there is an almost-uniform noise whose magnitude is  $\pm 1$  for  $u_i^l$  and  $u_i^r$ , and  $\pm 2$  for  $v_i^l$  and  $v_i^r$ . Additionally, pixel coordinates are quantified to a resolution of  $200 \times 200$ .

Experimental results are depicted in Figures 2, 3, and 4. Each one consists of a set of plots (from top to bottom): the image trajectories of the points, the errors of the visual features, the velocity screw, and the 3D trajectory of the end-effector.

Convergence to the desired images is always achieved, but quality is worse with the stereo 2D features. As pointed out by Lamiroy *et al.* (Lamiroy *et al.*, 2000), the stereo jacobian is largely overconstrained, and the control data  $\underline{\mathbf{s}}$

and  $\underline{\mathbf{s}}^*$  are redundant. But this is not sufficient to explain the curvy trajectory of the end-effector (bottom of Fig. 2), which almost leads out of the range of robot joints.

Such trajectory is neither caused by a too high gain: with  $\lambda = 0.1$  a smoother but similar trajectory is obtained, as depicted in Fig. 5. This problem has not been addressed before since very few experiments with image-based stereo visual servoing have been carried out *with cameras mounted on the end-effector*. To our knowledge, only Maru *et al.* (Maru *et al.*, 1993) have worked with this setup, but their tasks involved rather small rotations  $(\phi, \theta, \psi) = (10, 10, 10)$  (degree). In our manipulation task, the rotation between the initial and destination poses is:  $(\phi, \theta, \psi) = (72, 57, 50)$ . Translational distance is 250 mm, as opposed to 173 mm in Maru *et al.* (Maru *et al.*, 1993).

Approaches based on 3D features work better due to the linearity of the jacobian matrix. As shown theoretically, not only the image points but the center of gravity of 3D points translates along a straight line. As a result, the trajectory of the end-effector frame is closer to a straight line too, even with large rotations between frames.

#### 5. SUMMARY AND CONCLUSIONS

This paper has presented several approaches to image-based stereo visual servoing. As a main result, it has been shown how the effectiveness of the servoing task can be improved if estimated 3D features are used instead of raw image data.

Theoretical developments show how 3D control features are extracted from stereo images, and the jacobian matrix is computed for raw pixels, estimated 3D coordinates, and a new feature vector which uses stereo disparity.

Real experiments with adverse conditions (large rotation, noisy images, coarse calibration) show that the trajectory of the end-effector strongly relies on the features chosen for the control loop.

Future work should state more precisely the robustness of the different approaches, with respect to camera parameters and signal loss. We are also interested in considering other visual features like lines, and studying the relationships between image data and estimated 3D features.

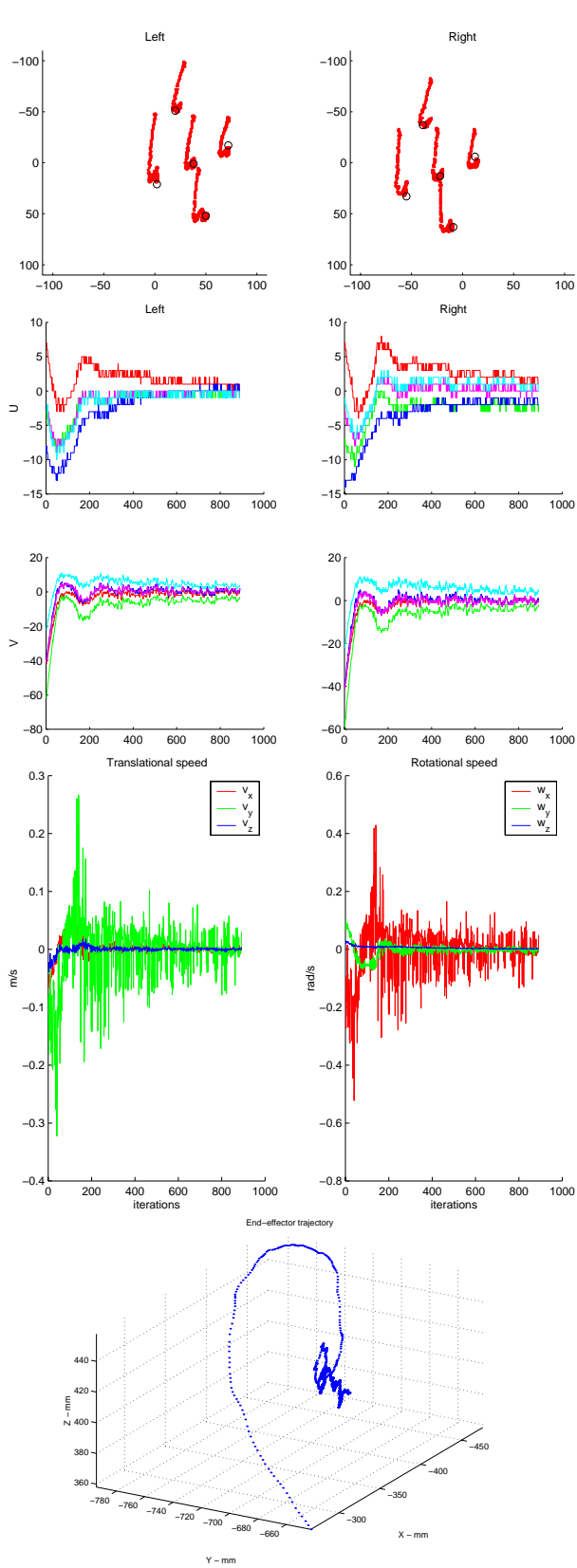


Fig. 2. Stereo 2D points: (from top to bottom) image trajectories, pixel errors, velocity screw, and trajectory of the end-effector.

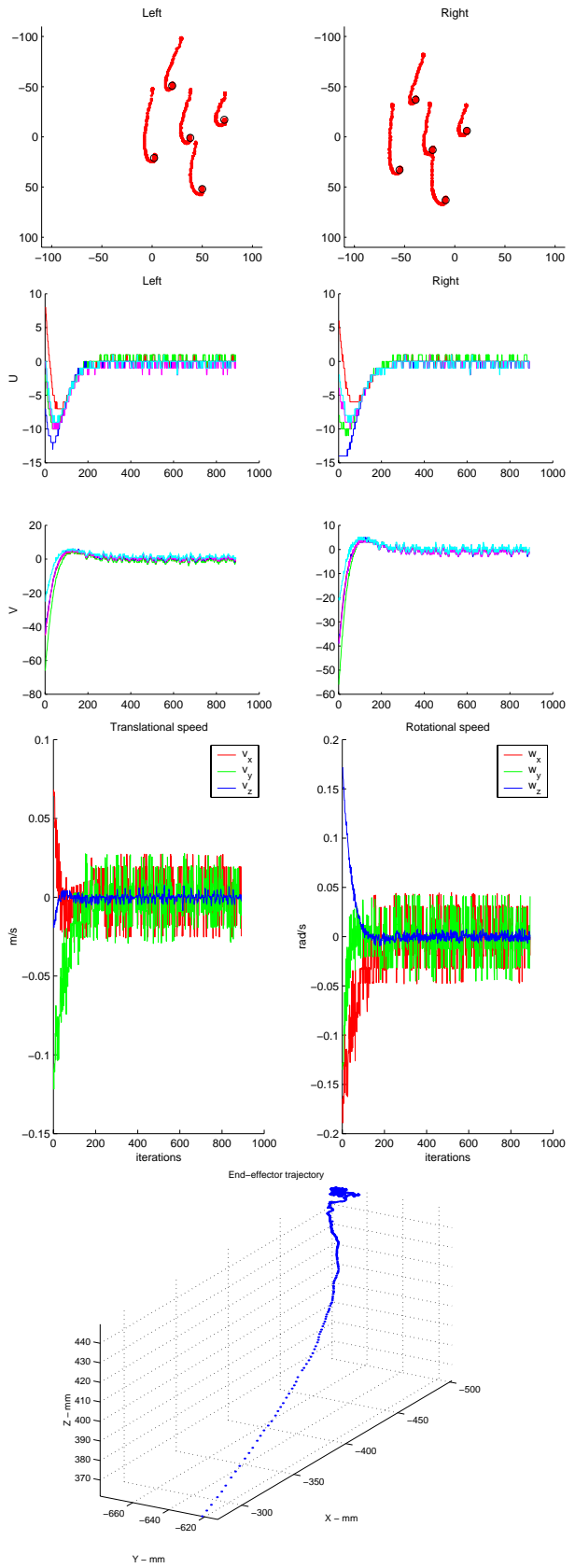


Fig. 3. Estimated 3D points: (from top to bottom) image trajectories, pixel errors, velocity screw, and trajectory of the end-effector

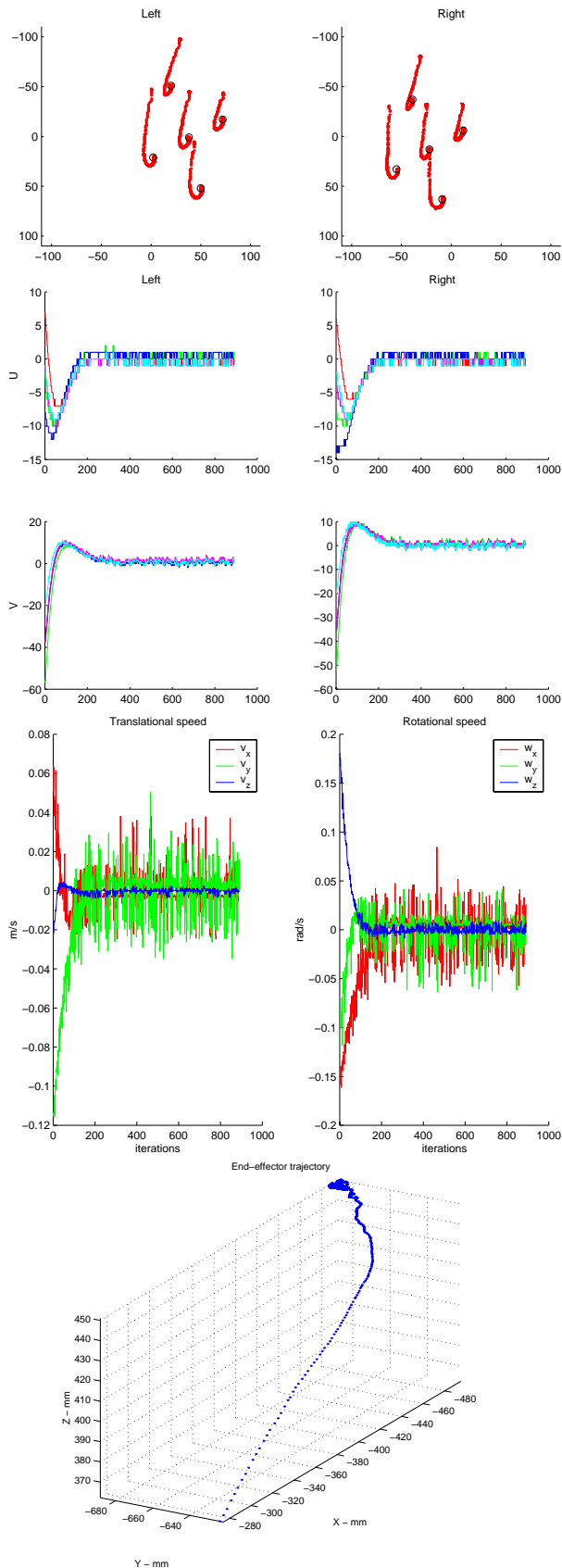


Fig. 4. Stereo 2D and disparity: (from top to bottom) image trajectories, pixel errors, velocity screw, and trajectory of the end-effector

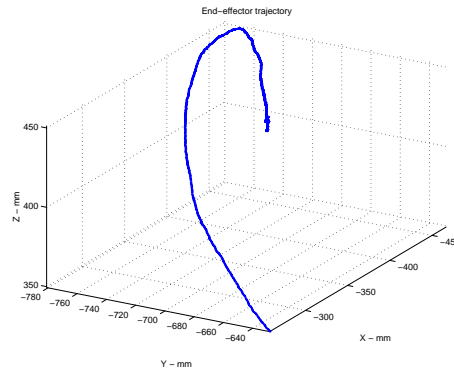


Fig. 5. Trajectory of the end-effector, with stereo 2D features, and  $\lambda = 0.1$ .

#### ACKNOWLEDGEMENT

This work is partially funded by the Valencian Government under grants GV99-67-1-14 and INV00-14-61, and by the CICYT under grant TAP98-0450.

#### 6. REFERENCES

Cervera, E. and P. Martinet (1999). Combining pixel and depth information in image-based visual servoing. In: *Proceedings of the International Conference on Advanced Robotics*. Vol. 1. ICAR'99. Tokyo, Japan.

Espiau, B., F. Chaumette and P. Rives (1992). A new approach to visual servoing in robotics. *IEEE Transactions on Robotics and Automation* 8(3), 313–326.

Hager, G., W. C. Chang and A. S. Morse (1995). Robot hand-eye coordination based on a stereo vision. *IEEE Control Systems Magazine* 15(1), 30–39.

Lamiroy, B., B. Espiau, N. Andreff and R. Horaud (2000). Controlling robots with two cameras: How to do it properly. In: *Proceedings of the IEEE International Conference on Robotics and Automation*. ICRA'2000. San Francisco, California, USA. pp. 2100–2105.

Maru, N., H. Kase, S. Yamada, A. Nishikawa and F. Miyazaki (1993). Manipulator control by visual servoing with stereo vision. In: *Proc. IROS'93*. Yokohama, Japan. pp. 1866–1870.