

# Models and Control Strategies for Visual Servoing

Nils T Siebel, Dennis Peters and Gerald Sommer  
*Christian-Albrechts-University of Kiel*  
*Germany*

## 1. Introduction

Visual servoing is the process of steering a robot towards a goal using visual feedback in a closed control loop as shown in Figure 1. The *output*  $u_n$  of the controller is a *robot movement* which steers the robot towards the goal. The state  $x_n$  of the system cannot be directly observed. Instead a visual measurement process provides feedback data, the vector of *current image features*  $y_n$ . The *input* to the controller is usually the difference between desired ( $y^*$ ) and actual values of this vector—the *image error vector*  $\Delta y_n$ .

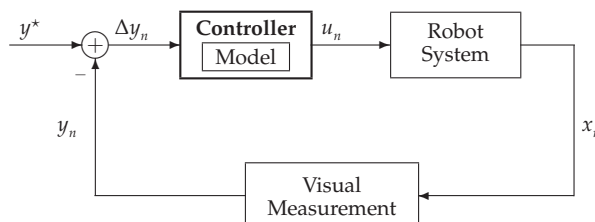


Fig. 1. Closed-loop image-based visual servoing control

In order for the controller to calculate the necessary robot movement it needs two main components:

1. a *model* of the environment—that is, a model of how the robot/scene will change after issuing a certain control command; and
2. a *control law* that governs how the next robot command is determined given current image measurements and model.

In this chapter we will look in detail on the effects different models and control laws have on the properties of a visual servoing controller. Theoretical considerations are combined with experiments to demonstrate the effects of popular models and control strategies on the behaviour of the controller, including convergence speed and robustness to measurement errors.

## 2. Building Models for Visual Servoing

### 2.1 Task Description

The aim of a visual servoing controller is to move the end-effector of one or more robot arms such that their configuration in relation to each other and/or to an object fulfils certain task-specific conditions. The feedback used in the controller stems from visual data, usually taken



Fig. 2. Robot Arm with Camera and Object

from one or more cameras mounted to the robot arm and/or placed in the environment. A typical configuration is shown in Figure 2. Here a camera is mounted to the robot's gripper ("eye-in-hand" setup), looking towards a glass jar. The controller's task in this case is to move the robot arm such that the jar can be picked up using the gripper. This is the case whenever the visual appearance of the object in the image has certain properties. In order to detect whether these properties are currently fulfilled a camera image can be taken and image processing techniques applied to extract the image positions of object markings. These image positions make up the *image feature vector*.

Since the control loop uses visual data the goal configuration can also be defined in the image. This can be achieved by moving the robot and/or the object in a suitable position and then acquiring a camera image. The image features measured in this image can act as *desired image features*, and a comparison of actual values at a later time to these desired values ("image error") can be used to determine the degree of agreement with the desired configuration. This way of acquiring desired image features is sometimes called "*teaching by showing*".

From a mathematical point of view, a successful visual servoing control process is equivalent to solving an optimisation problem. In this case a measure of the image error is minimised by moving the robot arm in the space of possible configurations. Visual servoing can also be regarded as practical feedback stabilisation of a dynamical system.

## 2.2 Modelling the Camera-Robot System

### 2.2.1 Preliminaries

The *pose* of an object is defined as its position and orientation. The *position* in 3D Euclidean space is given by the 3 Cartesian coordinates. The *orientation* is usually expressed by 3 angles, i.e. the rotation around the 3 coordinate axes. Figure 3 shows the notation used in this chapter, where *yaw*, *pitch* and *roll* angles are defined as the mathematically positive rotation around the  $x$ ,  $y$  and  $z$  axis. In this chapter we will use the  $\{\cdot\}$ -notation for a *coordinate system*, for example  $\{W\}$  will stand for the world coordinate system. A variable coordinate system—one which changes its pose to over time—will sometimes be indexed by the *time index*  $n \in \mathbb{N} =$

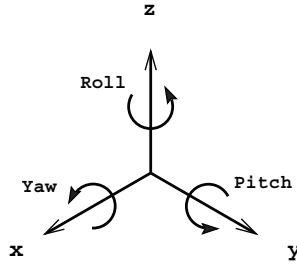


Fig. 3. Yaw, pitch and roll

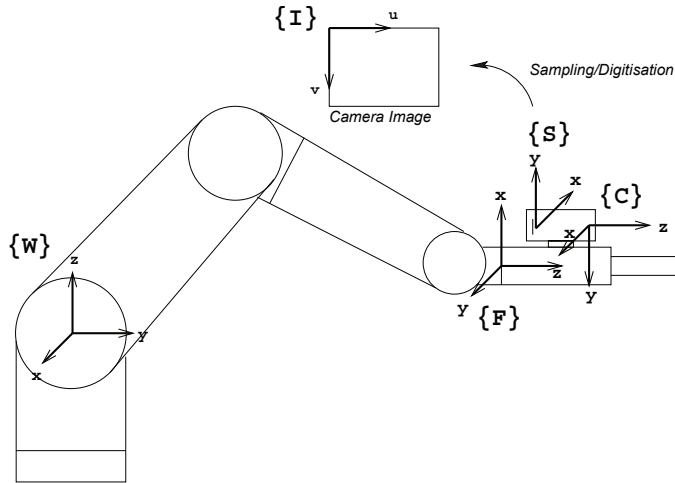


Fig. 4. World, Flange, Camera, Sensor and Image coordinate systems

0, 1, 2, . . . . An example is the camera coordinate system  $\{C_n\}$ , which moves relative to  $\{W\}$  as the robot moves since the camera is mounted to its hand.

Figure 4 lists the coordinate systems used for modelling the camera-robot system. The *world coordinate system*  $\{W\}$  is fixed at the robot base, the *flange coordinate system*  $\{F\}$  (sometimes called “*tool coordinate system*”, but this can be ambiguous) at the flange where the hand is mounted. The *camera coordinate system*  $\{C\}$  (or  $\{C_n\}$  at a specific time  $n$ ) is located at the optical centre of the camera, the *sensor coordinate system*  $\{S\}$  in the corner of its CCD/CMOS chip (sensor); their orientation and placement is shown in the figure. The *image coordinate system* which is used to describe positions in the digital image is called  $\{I\}$ . It is the only system to use pixel as its unit; all other systems use the same length unit, e.g. mm.

Variables that contain coordinates in a particular coordinate system will be marked by a superscript left of the variable, e.g.  ${}^A x$  for a vector  $x \in \mathbb{R}^n$  in  $\{A\}$ -coordinates. The coordinate transform which transforms a variable from a coordinate system  $\{A\}$  to another one,  $\{B\}$ , will be written  ${}^B_A T$ . If  ${}^A x$  and  ${}^B x$  express the pose of the same object then

$${}^A x = {}^A_B T {}^B x, \quad \text{and always} \quad {}^A_B T = ({}^B_A T)^{-1}. \tag{1}$$

The *robot’s pose* is defined as the pose of  $\{F\}$  in  $\{W\}$ .

### 2.2.2 Cylindrical Coordinates

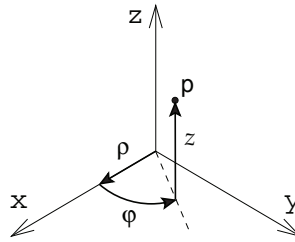


Fig. 5. A point  $p = (\rho, \varphi, z)$  in cylindrical coordinates.

An alternative way to describe point positions is by using a cylindrical coordinate system as the one in Figure 5. Here the position of the point  $p$  is defined by the distance  $\rho$  from a fixed axis (here aligned with the Cartesian  $z$  axis), an angle  $\varphi$  around the axis (here  $\varphi = 0$  is aligned with the Cartesian  $x$  axis) and a height  $z$  from a plane normal to the  $z$  axis (here the plane spanned by  $x$  and  $y$ ). Using the commonly used alignment with the Cartesian axes as in Figure 5 converting to and from cylindrical coordinates is easy. Given a point  $p = (x, y, z)$  in Cartesian coordinates, its cylindrical coordinates  $p = (\rho, \varphi, z) \in \mathbb{R} \times ]-\pi, \pi] \times \mathbb{R}$  are as follows:

$$\begin{aligned} \rho &= \sqrt{x^2 + y^2} \\ \varphi &= \text{atan2}(y, x) \\ &\stackrel{*}{=} \begin{cases} 0 & \text{if } x = 0 \text{ and } y = 0 \\ \arcsin(\frac{y}{\rho}) & \text{if } x \geq 0 \\ \arcsin(\frac{y}{\rho}) + \pi & \text{if } x < 0 \end{cases} \quad (2) \\ z &= z, \end{aligned}$$

(\* up to multiples of  $2\pi$ ), and, given a point  $p = (\rho, \varphi, z)$  in cylindrical coordinates:

$$\begin{aligned} x &= \rho \cos \varphi \\ y &= \rho \sin \varphi \\ z &= z. \end{aligned} \quad (3)$$

### 2.2.3 Modelling the Camera

A simple and popular approximation to the way images are taken with a camera is the *pinhole camera model* (from the pinhole camera/camera obscura models by Ibn al-Haytham “Alhacen”, 965–1039 and later by Gérard Desargues, 1591–1662), shown in Figure 6. A light ray from an object point passes an aperture plate through a very small hole (“pinhole”) and arrives at the sensor plane, where the camera’s CCD/CMOS chip (or a photo-sensitive film in the 17th century) is placed. In the digital camera case the sensor elements correspond to picture elements (“pixels”), and are mapped to the image plane. Since pixel positions are stored in the computer as unsigned integers the centre of the  $\{I\}$  coordinate system in the image plane is shifted to the upper left corner (looking towards the object/monitor). Therefore the centre  $c \neq (0, 0)^T$ .

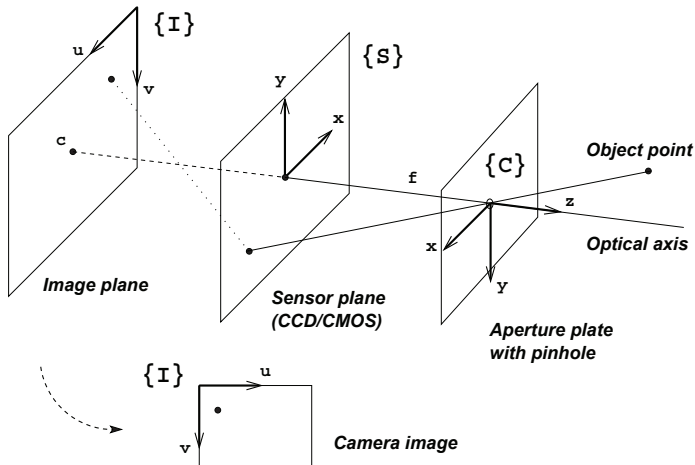


Fig. 6. Pinhole camera model

Sometimes the sensor plane is positioned in front of the aperture plate in the literature (e.g. in Hutchinson et al., 1996). This has the advantage that the  $x$ - and  $y$ -axis of  $\{S\}$  can be (directionally) aligned with the ones in  $\{C\}$  and  $\{I\}$  while giving identical coordinates. However, since this alternative notation has also the disadvantage of being less intuitive, we use the one defined above.

Due to the simple model of the way the light travels through the camera the object point's position in  $\{C\}$  and the coordinates of its projection in  $\{S\}$  and  $\{I\}$  are proportional, with a shift towards the new centre in  $\{I\}$ . In particular, the sensor coordinates  ${}^s p = ({}^s x, {}^s y)^T$  of the image of an object point  ${}^c p = ({}^c x, {}^c y, {}^c z)^T$  are given as

$${}^s x = \frac{{}^c x \cdot f}{{}^c z} \quad \text{and} \quad {}^s y = \frac{{}^c y \cdot f}{{}^c z}, \quad (4)$$

where  $f$  is the distance the aperture plate and the sensor plane, also called the "focal length" of the camera/lens.

The pinhole camera model's so-called "perspective projection" is not an exact model of the projection taking place in a modern camera. In particular, lens distortion and irregularities in the manufacturing (e.g. slightly tilted CCD chip or positioning of the lenses) introduce deviations. These modelling errors may need to be considered (or, corrected by a lens distortion model) by the visual servoing algorithm.

### 2.3 Defining the Camera-Robot System as a Dynamical System

As mentioned before, the camera-robot system can be regarded as a dynamical system. We define the state  $x_n$  of the robot system at a time step  $n \in \mathbb{N}$  as the current robot pose, i.e. the pose of the flange coordinate system  $\{F\}$  in world coordinates  $\{W\}$ .  $x_n \in \mathbb{R}^6$  will contain the position and orientation in the  $x, y, z, \text{yaw, pitch, roll}$  notation defined above. The set of possible robot poses is  $\mathcal{X} \subset \mathbb{R}^6$ . The output of the system is the image feature vector  $y_n$ . It contains pairs of image coordinates of object markings viewed by the camera, i.e.  $({}^s x_1, {}^s y_1, \dots, {}^s x_M, {}^s y_M)^T$  for  $M = \frac{m}{2}$  object markings (in our case  $M = 4$ , so  $y_n \in \mathbb{R}^8$ ).

Let  $\mathcal{Y} \subset \mathbb{R}^m$  be the set of possible output values. The *output (measurement) function* is  $\eta : \mathcal{X} \rightarrow \mathcal{Y}$ ,  $x_n \mapsto y_n$ . It contains the whole measurement process, including projection onto the sensor, digitisation and image processing steps.

The *input (control) variable*  $u_n \in \mathcal{U} \subset \mathbb{R}^6$  shall contain the desired pose change of the camera coordinate system. This robot movement can be easily transformed to a new robot pose  $\tilde{u}_n$  in  $\{W\}$ , which is given to the robot in a move command. Using this definition of  $u_n$  an input of  $(0, 0, 0, 0, 0, 0)^T$  corresponds to no robot movement, which has advantages, as we shall see later. Let  $\varphi : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$ ,  $(x_n, u_n) \mapsto x_{n+1}$  be the corresponding *state transition (next-state) function*.

With these definitions the camera-robot system can be defined as a time invariant, time discrete input-output system:

$$\begin{aligned} x_{n+1} &= \varphi(x_n, u_n) \\ y_n &= \eta(x_n). \end{aligned} \quad (5)$$

When making some mild assumptions, e.g. that the camera does not move relative to  $\{F\}$  during the whole time, the state transition function  $\varphi$  can be calculated as follows:

$$\begin{aligned} \varphi(x_n, u_n) &= x_{n+1} = {}^W x_{n+1} = {}^W \tilde{u}_n \hat{=} {}_{F_{n+1}}^W \mathbf{T} \\ &= \underbrace{{}^W \mathbf{T}}_{\hat{=} x_n} \circ \underbrace{{}_{C_n}^{F_n} \mathbf{T}}_{\star} \circ \underbrace{{}_{C_{n+1}}^{C_n} \mathbf{T}}_{\hat{=} u_n} \circ \underbrace{{}_{F_{n+1}}^{C_{n+1}} \mathbf{T}}_{\star}, \end{aligned} \quad (6)$$

where  $\{F_n\}$  is the flange coordinate system at time step  $n$ , etc., and the  $\hat{=}$  operator expresses the equivalence of a pose with its corresponding coordinate transform.

$\star = \text{external ("extrinsic") camera parameters}$ ;  ${}_{C_n}^{F_n} \mathbf{T} = {}_{C_{n+1}}^{T_{n+1}} \mathbf{T} = ({}_{T_{n+1}}^{C_{n+1}} \mathbf{T})^{-1} \quad \forall n \in \mathbf{N}$ .

For  $m = 2$  image features corresponding to coordinates  $({}^s x, {}^s y)$  of a projected object point  ${}^W p$  the equation for  $\eta$  follows analogously:

$$\begin{aligned} \eta(x) = y &= {}^s y = {}_c^s \mathbf{T} {}_c^c p \\ &= {}_c^s \mathbf{T} \circ {}_T^c \mathbf{T} \circ {}_W^T \mathbf{T} {}^W p, \end{aligned} \quad (7)$$

where  ${}_c^s \mathbf{T}$  is the mapping of the object point  ${}^c p$  depending on the focal length  $f$  according to the pinhole camera model / perspective projection defined in (4).

## 2.4 The Forward Model—Mapping Robot Movements to Image Changes

In order to calculate necessary movements for a given desired change in visual appearance the relation between a robot movement and the resulting change in the image needs to be modelled. In this section we will analytically derive a *forward model*, i.e. one that expresses image changes as a function of robot movements, for the eye-in-hand setup described above. This forward model can then be used to predict changes effected by controller outputs, or (as it is usually done) simplified and then inverted. An *inverse model* can be directly used to determine the controller output given actual image measurements.

Let  $\Phi : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{Y}$  the function that expresses the system output  $y$  depending on the state  $x$  and the input  $u$ :

$$\Phi(x, u) := \eta \circ \varphi(x, u) = \eta(\varphi(x, u)). \quad (8)$$

For simplicity we also define the function which expresses the behaviour of  $\Phi(x_n, \cdot)$  at a time index  $n$ , i.e. the dependence of image features on the camera movement  $u$ :

$$\Phi_n(u) := \Phi(x_n, u) = \eta(\varphi(x_n, u)). \quad (9)$$

This is the forward model we wish to derive.

$\Phi_n$  depends on the camera movement  $u$  and the current system state, the robot pose  $x_n$ . In particular it depends on the position of all object markings in the current camera coordinate system. In the following we need assume the knowledge of the camera's focal length  $f$  and the  $^c_z$  component of the positions of image markings in  $\{C\}$ , which cannot be derived from their image position  $(^s_x, ^s_y)$ . Then with the help of  $f$  and the image coordinates  $(^s_x, ^s_y)$  the complete position of the object markings in  $\{C\}$  can be derived with the pinhole camera model (4).

We will first construct the model  $\Phi_n$  for the case of a single object marking,  $M = \frac{m}{2} = 1$ . According to equations (6) and (7) we have for an object point  $^w_p$ :

$$\begin{aligned} \Phi_n(u) &= \eta \circ \varphi(x_n, u) \\ &= \begin{matrix} s \\ c_{n+1} \end{matrix} \mathbf{T} \circ \begin{matrix} c_{n+1} \\ c_n \end{matrix} \mathbf{T} \circ \begin{matrix} c_n \\ T \end{matrix} \mathbf{T} \circ \begin{matrix} T \\ w \end{matrix} \mathbf{T} \begin{matrix} w \\ p \end{matrix} \\ &= \begin{matrix} s \\ c_{n+1} \end{matrix} \mathbf{T} \circ \begin{matrix} c_{n+1} \\ c_n \end{matrix} \mathbf{T} \begin{matrix} c_n \\ x \end{matrix}, \end{aligned} \quad (10)$$

where  $^{c_n}x$  are the coordinates of the object point in  $\{C_n\}$ .

In the system state  $x_n$  the position of an object point  $^{c_n}x =: p = (p_1, p_2, p_3)^T$  can be derived with  $(^s_x, ^s_y)^T$ , assuming the knowledge of  $f$  and  $^c_z$ , via (4). Then the camera changes its pose by  $^c_u =: u = (u_1, u_2, u_3, u_4, u_5, u_6)^T$ ; we wish to know the new coordinates  $(^s_{\tilde{x}}, ^s_{\tilde{y}})^T$  of  $p$  in the image. The new position  $\tilde{p}$  of the point in new camera coordinates is given by a translation by  $u_1$  through  $u_3$  and a rotation of the camera by  $u_4$  through  $u_6$ . We have

$$\begin{aligned} \tilde{p} &= \text{rot}_x(-u_4) \text{rot}_y(-u_5) \text{rot}_z(-u_6) \begin{pmatrix} p_1 - u_1 \\ p_2 - u_2 \\ p_3 - u_3 \end{pmatrix} \\ &= \begin{pmatrix} c_5 c_6 & c_5 s_6 & -s_5 \\ s_4 s_5 c_6 - c_4 s_6 & s_4 s_5 s_6 + c_4 c_6 & s_4 c_5 \\ c_4 s_5 c_6 + s_4 s_6 & c_4 s_5 s_6 - s_4 c_6 & c_4 c_5 \end{pmatrix} \begin{pmatrix} p_1 - u_1 \\ p_2 - u_2 \\ p_3 - u_3 \end{pmatrix} \end{aligned} \quad (11)$$

using the short notation

$$s_i := \sin u_i, \quad c_i := \cos u_i \quad \text{for } i = 4, 5, 6. \quad (12)$$

Again with the help of the pinhole camera model (4) we can calculate the  $\{S\}$  coordinates of the projection of the new point, which finally yields the model  $\Phi_n$ :

$$\begin{aligned} \begin{bmatrix} ^s_{\tilde{x}} \\ ^s_{\tilde{y}} \end{bmatrix} &= \Phi(x_n, u) \\ &= \Phi_n(u) \\ &= f \cdot \begin{bmatrix} \frac{c_5 c_6 (p_1 - u_1) + c_5 s_6 (p_2 - u_2) - s_5 (p_3 - u_3)}{(c_4 s_5 c_6 + s_4 s_6) (p_1 - u_1) + (c_4 s_5 s_6 - s_4 c_6) (p_2 - u_2) + c_4 c_5 (p_3 - u_3)} \\ \frac{(s_4 s_5 c_6 - c_4 s_6) (p_1 - u_1) + (s_4 s_5 s_6 + c_4 c_6) (p_2 - u_2) + s_4 c_5 (p_3 - u_3)}{(c_4 s_5 c_6 + s_4 s_6) (p_1 - u_1) + (c_4 s_5 s_6 - s_4 c_6) (p_2 - u_2) + c_4 c_5 (p_3 - u_3)} \end{bmatrix}. \end{aligned} \quad (13)$$

## 2.5 Simplified and Inverse Models

As mentioned before, the controller needs to derive necessary movements from given desired image changes, for which an inverse model is beneficial. However,  $\Phi_n(u)$  is too complicated to invert. Therefore in practice usually a linear approximation  $\hat{\Phi}_n(u)$  of  $\Phi_n(u)$  is calculated and then inverted. This can be done in a number of ways.

### 2.5.1 The Standard Image Jacobian

The simplest and most common linear model is the *Image Jacobian*. It is obtained by Taylor expansion of (13) around  $u = 0$ :

$$\begin{aligned}
 y_{n+1} &= \eta(\varphi(x_n, u)) \\
 &= \Phi(x_n, u) \\
 &= \Phi_n(u) \\
 &= \Phi_n(0 + u) \\
 &= \Phi_n(0) + J_{\Phi_n}(0) u + \mathcal{O}(\|u\|^2).
 \end{aligned} \tag{14}$$

With  $\Phi_n(0) = y_n$  and the definition  $J_n := J_{\Phi_n}(0)$  the image change can be approximated

$$y_{n+1} - y_n \approx J_n u \tag{15}$$

for sufficiently small  $\|u\|_2$ .

The Taylor expansion of the two components of (13) around  $u = 0$  yields the Image Jacobian  $J_n$  for one object marking ( $m = 2$ ):

$$J_n = \begin{pmatrix} -\frac{f}{c_z} & 0 & \frac{s_x}{c_z} & \frac{s_x s_y}{f} & -f - \frac{s_x^2}{f} & s_y \\ 0 & -\frac{f}{c_z} & \frac{s_y}{c_z} & f + \frac{s_y^2}{f} & -\frac{s_x s_y}{f} & -s_x \end{pmatrix} \tag{16}$$

where again image positions were converted back to sensor coordinates.

The Image Jacobian for  $M$  object markings,  $M \in \mathbb{N}_{>1}$ , can be derived analogously; the change of the  $m = 2M$  image features can be approximated by



$$y_{n+1} - y_n \approx J_n u$$

$$= \begin{pmatrix} -\frac{f}{c_{z_1}} & 0 & \frac{s_{x_1}}{c_{z_1}} & \frac{s_{x_1} s_{y_1}}{f} & -f - \frac{s_{x_1}^2}{f} & s_{y_1} \\ 0 & -\frac{f}{c_{z_1}} & \frac{s_{y_1}}{c_{z_1}} & f + \frac{s_{y_1}^2}{f} & -\frac{s_{x_1} s_{y_1}}{f} & -s_{x_1} \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ -\frac{f}{c_{z_M}} & 0 & \frac{s_{x_M}}{c_{z_M}} & \frac{s_{x_M} s_{y_M}}{f} & -f - \frac{s_{x_M}^2}{f} & s_{y_M} \\ 0 & -\frac{f}{c_{z_M}} & \frac{s_{y_M}}{c_{z_M}} & f + \frac{s_{y_M}^2}{f} & -\frac{s_{x_M} s_{y_M}}{f} & -s_{x_M} \end{pmatrix} \begin{pmatrix} u_1 \\ \vdots \\ u_6 \end{pmatrix}, \quad (17)$$

for small  $\|u\|_2$ , where  $(s_{x_i}, s_{y_i})$  are the sensor coordinates of the  $i$ th projected object marking and  $c_{z_i}$  their distances from the camera,  $i = 1, \dots, M$ .

### 2.5.2 A Linear Model in the Cylindrical Coordinate System

Iwatsuki and Okiyama (2005) suggest a formulation of the problem in cylindrical coordinates. This means that positions of markings on the sensor are given in polar coordinates,  $(\rho, \varphi)^T$  where  $\rho$  and  $\varphi$  are defined as in Figure 5 ( $z = 0$ ). The Image Jacobian  $J_n$  for one image point is given in this case by

$$J_n = \begin{pmatrix} -\frac{f c_\varphi}{c_z} & -\frac{f s_\varphi}{c_z} & \frac{c_y s_\varphi + c_x c_\varphi}{c_z} & \left(f + \frac{c_y^2}{f}\right) s_\varphi + \frac{c_x c_y c_\varphi}{f} & \left(-f - \frac{c_x^2}{f}\right) c_\varphi - \frac{c_x c_y s_\varphi}{f} & c_y c_\varphi - c_x s_\varphi \\ \frac{f s_\varphi}{c_z} & -\frac{f c_\varphi}{c_z} & \frac{c_y c_\varphi + c_x s_\varphi}{c_z} & \left(f + \frac{c_y^2}{f}\right) c_\varphi - \frac{c_x c_y s_\varphi}{f} & \left(f + \frac{c_x^2}{f}\right) s_\varphi - \frac{c_x c_y c_\varphi}{f} & -c_y s_\varphi - c_x c_\varphi \end{pmatrix} \quad (18)$$

with the short notation

$$s_\varphi := \sin \varphi \quad \text{and} \quad c_\varphi := \cos \varphi. \quad (19)$$

and analogously for  $M > 1$  object markings.

### 2.5.3 Quadratic Models

A quadratic model, e.g. a quadratic approximation of the system model (13), can be obtained by a Taylor expansion; a resulting approximation for  $M = 1$  marking is

$$y_{n+1} = \begin{bmatrix} s_x \\ s_y \end{bmatrix} = \Phi_n(0) + J_{\Phi_n}(0) u + \frac{1}{2} \begin{bmatrix} u^T H_{S_x} u \\ u^T H_{S_y} u \end{bmatrix} + \mathcal{O}(\|u\|^3). \quad (20)$$

where again  $\Phi_n(0) = y_n$  and  $J_{\Phi_n}(0) = J_n$  from (16), and the Hessian matrices are

$$H_{S_x} = \begin{pmatrix} 0 & 0 & \frac{-f}{c_z^2} & \frac{-s_y}{c_z} & \frac{2^s x}{c_z} & 0 \\ 0 & 0 & 0 & \frac{-s_x}{c_z} & 0 & \frac{-f}{c_z} \\ \frac{-f}{c_z^2} & 0 & \frac{2^s x}{c_z^2} & \frac{2^s x^s y}{f c_z} & \frac{-2^s x^2}{f c_z} & \frac{s_y}{c_z} \\ \frac{-s_y}{c_z} & \frac{-s_x}{c_z} & \frac{2^s x^s y}{f c_z} & s_x \left(1 + 2 \left(\frac{s_y}{f}\right)^2\right) & -s_y \left(1 + 2 \left(\frac{s_x}{f}\right)^2\right) & \frac{s_y^2 - s_x^2}{f} \\ \frac{2^s x}{c_z} & 0 & \frac{-2^s x^2}{f c_z} & -s_y \left(1 + 2 \left(\frac{s_x}{f}\right)^2\right) & 2^s x \left(1 + \frac{s_x^2}{f}\right)^2 & -2^s x \frac{s_y}{f} \\ 0 & \frac{-f}{c_z} & \frac{s_y}{c_z} & \frac{s_y^2 - s_x^2}{f} & \frac{-2^s x^s y}{f} & -s_x \end{pmatrix} \quad (21)$$

as well as

$$H_{S_y} = \begin{pmatrix} 0 & 0 & 0 & 0 & \frac{s_y}{c_z} & \frac{f}{c_z} \\ 0 & 0 & \frac{-f}{c_z^2} & \frac{-2^s y}{c_z} & \frac{s_x}{c_z} & 0 \\ 0 & \frac{-f}{c_z^2} & \frac{2^s y}{c_z^2} & \frac{2^s y^2}{f c_z} & \frac{-2^s x^s y}{f c_z} & \frac{-s_x}{c_z} \\ 0 & \frac{-2^s y}{c_z} & \frac{2^s y^2}{f c_z} & 2^s y \left(1 + \frac{s_y^2}{f}\right)^2 & \left(\frac{s_y}{f}\right) \left(-2^s x \frac{s_y}{f}\right) & \frac{-2^s x^s y}{f} \\ \frac{s_y}{c_z} & \frac{s_x}{c_z} & \frac{-2^s x^s y}{f c_z} & \left(\frac{s_y}{f}\right) \left(-2^s x \frac{s_y}{f}\right) & s_y \left(1 + 2 \left(\frac{s_x}{f}\right)^2\right) & \frac{s_x^2 - s_y^2}{f} \\ \frac{f}{c_z} & 0 & \frac{-s_x}{c_z} & \frac{-2^s x^s y}{f} & \frac{s_x^2 - s_y^2}{f} & -s_y \end{pmatrix}. \quad (22)$$

#### 2.5.4 A Mixed Model

Malis (2004) proposes a way of constructing a mixed model which consists of different linear approximations of the target function  $\Phi$ . Let  $x_n$  again be the current robot pose and  $x^*$  the teach pose. For a given robot command  $u$  we set again  $\Phi_n(u) := \Phi(x_n, u)$  and now also  $\Phi^*(u) := \Phi(x^*, u)$  such that  $\Phi_n(0) = y_n$  and  $\Phi^*(0) = y^*$ . Then Taylor expansions of  $\Phi_n$  and  $\Phi^*$  at  $u = 0$  yield

$$y_{n+1} = y_n + J_{\Phi_n}(0)u + \mathcal{O}(\|u\|^2) \quad (23)$$

and

$$y_{n+1} = y_n + J_{\Phi^*}(0)u + \mathcal{O}(\|u\|^2). \quad (24)$$

In other words, both Image Jacobians,  $J_n := J_{\Phi_n}(0)$  and  $J^* := J_{\Phi^*}(0)$  can be used as linear approximations of the behaviour of the robot system. One of these models has its best validity

at the current pose, the other at the teach pose. Since we are moving the robot from one towards the other it may be useful to consider both models. Malis proposes to use a mixture of these two models, i.e.

$$y_{n+1} - y_n \approx \frac{1}{2}(J_n + J^*) u. \quad (25)$$

In his control law (see Section 3 below) he calculates the pseudoinverse of the Jacobians, and therefore calls this approach “Pseudo-inverse of the Mean of the Jacobians”, or short “PMJ”. In a variation of this approach the computation of mean and pseudo-inverse is exchanged, which results in the “MPJ” method. See Section 3 for details.

### 2.5.5 Estimating Models

Considering the fact that models can only ever approximate the real system behaviour it may be beneficial to use measurements obtained during the visual servoing process to update the model “online”. While even the standard models proposed above use current measurements to estimate the distance  $\hat{z}$  from the object to use this estimate in the Image Jacobian, there are also approaches that estimate more variables, or construct a complete model from scratch. This is most useful when no certain data about the system state or setup are available. The following aspects need to be considered when estimating the Image Jacobian—or other models:

- How precise are the measurements used for model estimation, and how large is the sensitivity of the model to measurement errors?
- How many measurements are needed to construct the model? For example, some methods use 6 robot movements to measure the 6-dimensional data within the Image Jacobian. In a static look-and-move visual servoing setup which may reach its goal in 10-20 movements with a given Jacobian the resulting increase in necessary movements, as well as possible mis-directed movements until the estimation process converges, need to be weighed against the flexibility achieved by the automatic model tuning.

The most prominent approach to estimation methods of the whole Jacobian is the *Broyden approach* which has been used by Jägersand (1996). The Jacobian estimation uses the following update formula for the current estimate  $\hat{J}_n$ :

$$\hat{J}_n := \frac{c_n}{c_{n-1}} \mathbf{T} \left( \hat{J}_{n-1} + \frac{(y_n - y_{n-1} - \hat{J}_{n-1} u_n) u_n^T}{u_n^T u_n} \right), \quad (26)$$

with an additional weighting of the correction term

$$J_n := \gamma \hat{J}_{n-1} + (1 - \gamma) \hat{J}_n, \quad 0 \leq \gamma < 1 \quad (27)$$

to reduce the sensitivity of the estimate to measurement noise.

In the case of Jägersand’s system using an estimation like this makes sense since he worked with a dynamic visual servoing setup where many more measurements are made over time compared to our setup (“static look-and-move”, see below).

In combination with a model-based measurement a non-linear model could also make sense. A number of methods for the estimation of quadratic models are available in the optimisation literature. More on this subject can be found e.g. in Fletcher (1987, chapter 3) and Sage and White (1977, chapter 9).

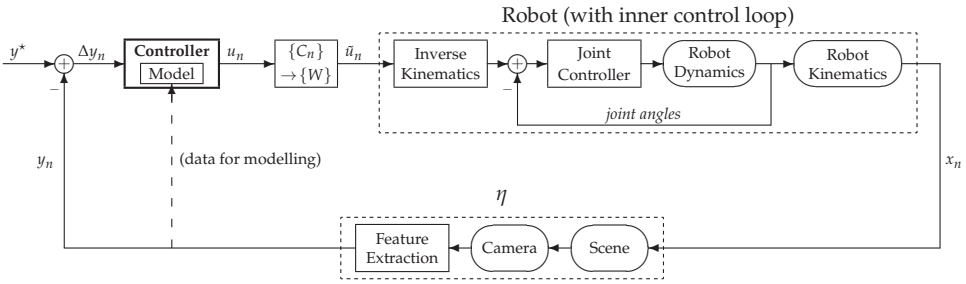


Fig. 7. Typical closed-loop image-based visual servoing controller

### 3. Designing a Visual Servoing Controller

Using one of the models defined above we wish to design a controller which steers the robot arm towards an object of unknown pose. This is to be realised in the visual feedback loop depicted in Figure 7. Using the terminology defined by Weiss et al. (1987) the visual servoing controller is of the type “*Static Image-based Look-and-Move*”. “Image-based” means that goal and error are defined in image coordinates instead of using positions in normal space (that would be “position-based”). “Static Look-and-Move” means that the controller is a sampled data feedback controller and the robot does not move while a measurement is taken. This traditionally implies that the robot is controlled by giving world coordinates to the controller instead of directly manipulating robot joint angles (Chaumette and Hutchinson, 2008; Hutchinson et al., 1996).

The object has 4 circular, identifiable markings. Its appearance in the image is described by the *image feature vector*  $y_n \in \mathbb{R}^8$  that contains the 4 pairs of image coordinates of these markings in a fixed order. The desired pose relative to the object is defined by the object’s appearance in that pose by measuring the corresponding *desired image features*  $y^* \in \mathbb{R}^8$  (“*teaching by showing*”). Object and robot are then moved so that no Euclidean position of the object or robot is known to the controller. The *input* to the controller is the *image error*  $\Delta y_n := y^* - y_n$ . The current image measurements  $y_n$  are also given to the controller for adapting its internal model to the current situation. The *output* of the controller is a relative movement of the robot in the camera coordinate system, a 6-dimensional vector  $(x, y, z, yaw, pitch, roll)$  for a 6 DOF movement.

Controllers can be classified into approaches where the control law (or its parameters) are adapted over time, and approaches where they are fixed. Since these types of controllers can exhibit very different controlling behaviour we will split our considerations of controllers into these two parts, after some general considerations.

#### 3.1 General Approach

Generally, in order to calculate the necessary camera movement  $u_n$  for a given desired image change  $\Delta \tilde{y}_n := \tilde{y}_{n+1} - y_n$  we again use an approximation  $\hat{\Phi}_n$  of  $\Phi_n$ , for example the image Jacobian  $J_n$ . Then we select

$$u_n \in \underset{u \in \mathcal{U}(x_n)}{\operatorname{argmin}} \|\Delta \tilde{y}_n - \hat{\Phi}_n(u)\|_2^2. \quad (28)$$

where a given algorithm may or may not enforce a restriction  $u \in \mathcal{U}(x_n)$  on the admissible movements when determining  $u$ . If this restriction is inactive and we are using a Jacobian,  $\hat{\Phi}_n = J_n$ , then the solution to (28) with minimum norm  $\|u_n\|_2$  is given by

$$u_n = J_n^+ \Delta \tilde{y}_n \quad (29)$$

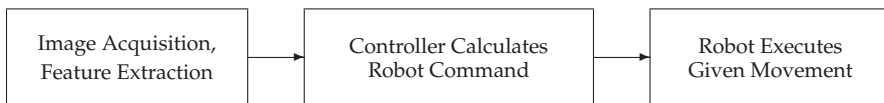
where  $J_n^+$  is the *pseudo-inverse* of  $J_n$ .

With 4 coplanar object markings  $m = 8$  and thereby  $J_n \in \mathbb{R}^{8 \times 6}$ . One can show that  $J_n$  has maximum rank<sup>1</sup>, so  $\text{rk } J_n = 6$ . Then the pseudo-inverse  $J_n^+ \in \mathbb{R}^{6 \times 8}$  of  $J_n$  is given by:

$$J_n^+ = (J_n^T J_n)^{-1} J_n^T \quad (30)$$

(see e.g. Deuffhard and Hohmann, 2003, chapter 3).

When realising a control loop given such a controller one usually sets a fixed error threshold  $\varepsilon > 0$  and repeats the steps



until

$$\|\Delta y_n\|_2 = \|y^* - y_n\|_2 < \varepsilon, \quad (31)$$

or until

$$\|\Delta y_n\|_\infty = \|y^* - y_n\|_\infty < \varepsilon \quad (32)$$

if one wants to stop only when the maximum deviation in *any* component of the image feature vector is below  $\varepsilon$ . Setting  $\varepsilon := 0$  is not useful in practice since measurements even in the same pose tend to vary a little due to small movements of the robot arm or object as well as measurement errors and fluctuations.

## 3.2 Non-Adaptive Controllers

### 3.2.1 The Traditional Controller

The most simple controller, which we will call the “*Traditional Controller*” due to its heritage, is a straightforward proportional controller as known in engineering, or a dampened Gauss-Newton algorithm as it is known in mathematics.

Given an Image Jacobian  $J_n$  we first calculate the full Gauss-Newton step  $\Delta u_n$  for a complete movement to the goal in one step (desired image change  $\Delta \tilde{y}_n := \Delta y_n$ ):

$$\Delta u_n := J_n^+ \Delta y_n \quad (33)$$

without enforcing a restriction  $u \in \mathcal{U}(x_n)$  for the admissibility of a control command.

In order to ensure a convergence of the controller the resulting vector is then scaled with a dampening factor  $0 < \lambda_n \leq 1$  to get the controller output  $u_n$ . In the traditional controller the factor  $\lambda_n$  is constant over time and the most important parameter of this algorithm. A typical value is  $\lambda_n = \lambda = 0.1$ ; higher values may hinder convergence, while lower values also significantly slow down convergence. The resulting controller output  $u_n$  is given by

<sup>1</sup> One uses the fact that no 3 object markings are on a straight line,  $\zeta_i > 0$  for  $i = 1, \dots, 4$  and all markings are visible (in particular, neither all four  $x_i$  nor all four  $y_i$  are 0).

$$u_n := \lambda \cdot J_n^+ \Delta y_n. \quad (34)$$

### 3.2.2 Dynamical and Constant Image Jacobians

As mentioned in the previous section there are different ways of defining the Image Jacobian. It can be defined in the current pose, and is then calculated using the current distances to the object,  ${}^c z_i$  for marking  $i$ , and the current image features. This is the *Dynamical Image Jacobian*  $J_n$ . An alternative is to define the Jacobian in the teach (goal) pose  $x^*$ , with the image data  $y^*$  and distances at that pose. We call this the *Constant Image Jacobian*  $J^*$ . Unlike  $J_n$ ,  $J^*$  is constant over time and does not require image measurements for its adaptation to the current pose.

From a mathematical point of view the model  $J_n$  has a better validity in the current system state and should therefore yield better results. We shall later see whether this is the case in practice.

### 3.2.3 The Retreat-Advance Problem

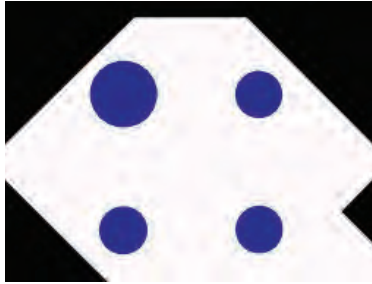


Fig. 8. Camera view in the start pose with a pure rotation around the  ${}^c z$  axis

When the robot's necessary movement to the goal pose is a pure rotation around the optical axis ( ${}^c z$ , approach direction) there can be difficulties when using the standard Image Jacobian approach (Chaumette, 1998). The reason is that the linear approximation  $J_n$  models the relevant properties of  $\Phi_n$  badly in these cases. This is also the case with  $J^*$  if this Jacobian is used. The former will cause an unnecessary movement away from the object, the latter a movement towards the goal. The larger the roll angle, the more pronounced is this phenomenon, an extreme case being a roll error of  $\pm\pi$  (all other pose elements already equal to the teach pose) where the Jacobians suggest a pure movement along the  ${}^c z$  axis. Corke and Hutchinson (2001) call this the "*Retreat-Advance Problem*" or the "*Chaumette Conundrum*".

### 3.2.4 Controllers using the PMJ and MPJ Models

In order to overcome the Retreat-Advance Problem the so-called "*PMJ Controller*" (Malis, 2004) uses the pseudo-inverse of the mean of the two Jacobians  $J_n$  and  $J^*$ . Using again a dampening factor  $0 < \lambda \leq 1$  the controller output is given by

$$u_n = \lambda \cdot \left( \frac{1}{2} (J_n + J^*) \right)^+ \Delta y_n. \quad (35)$$

Analogously, the “MPJ Controller” works with the mean of the pseudo-inverse of the Jacobians:

$$u_n = \lambda \cdot \left( \frac{1}{2} (J_n^+ + J^{*+}) \right) \Delta y_n. \quad (36)$$

Otherwise, these controllers work like the traditional approach, with a constant dampening  $\lambda$ .

### 3.2.5 Defining the Controller in the Cylindrical Coordinate System

Using the linear model by Iwatsuki and Okiyama (2005) in the cylindrical coordinate system as discussed in Section 2.5.2 a special controller can also be defined. The authors define the image error for the  $i$ th object marking as follows:

$$e_i := \begin{pmatrix} \rho^* - \rho \\ \rho(\varphi^* - \varphi) \end{pmatrix} \quad (37)$$

where  $(\rho, \varphi)^T$  is the current position and  $(\rho^*, \varphi^*)$  the teach position. The control command  $u$  is then given by

$$u = \lambda \tilde{J}^+ e, \quad (38)$$

$\tilde{J}^+$  being the pseudo-inverse of the Image Jacobian in cylindrical coordinates from equation (18).  $e$  is the vector of pairs of image errors in the markings, i.e. a concatenation of the  $e_i$  vectors.

It should be noted that even if  $e$  is given in cylindrical coordinates, the output  $u$  of the controller is in Cartesian coordinates.

Due to the special properties of cylindrical coordinates, the calculation of the error and control command is very much dependent on the definition of the origin of the coordinate system. Iwatsuki and Okiyama (2005) therefore present a way to shift the origin of the coordinate system such that numerical difficulties are avoided.

One approach to select the origin of the cylindrical coordinate system is such that the current pose can be transformed to the desired (teach) pose with a pure rotation around the axis normal to the sensor plane, through the origin. For example, the general method given by Kanatani (1996) can be applied to this problem.

Let  $l = (l_x, l_y, l_z)^T$  be the unit vector which defines this rotation axis, and  $o = (o_x, o_y)^T$  the new origin, obtained by shifting the original origin  $(0,0)^T$  in  $\{S\}$  by  $(\eta, \xi)^T$ .

If  $|l_z|$  is very small then the rotation axis  $l$  is almost parallel to the sensor. Then  $\eta$  and  $\xi$  are very large, which can create numerical difficulties. Since the resulting cylindrical coordinate system approximates a Cartesian coordinate system as  $\eta, \xi \rightarrow \infty$ , the standard Cartesian Image Jacobian  $J_n$  from (17) can therefore be used if  $|l_z| < \delta$  for a given lower limit  $\delta$ .

### 3.3 Adaptive Controllers

Using adaptive controllers is a way to deal with errors in the model, or with problems resulting from the simplification of the model (e.g. linearisation, or the assumption that the camera works like a pinhole camera). The goal is to ensure a fast convergence of the controller in spite of these errors.

### 3.3.1 Trust Region-based Controllers

*Trust Region methods* are known from mathematics as globally convergent optimisation methods (Fletcher, 1987). In order to optimise “difficult” functions one uses a model of its properties, like we do here with the Image Jacobian. This model is adapted to the current state/position in the solution space, and therefore only valid within some region around the current state. The main idea in trust region methods is to keep track of the validity of the current system model, and adapt a so-called “*Trust Region*”, or “*Model Trust Region*” around the current state within which the model does not exhibit more than a certain pre-defined “acceptable error”.

To our knowledge the first person to use trust region methods for a visual servoing controller was Jägersand (1996). Since the method was adapted to a particular setup and cannot be used here we have developed a different trust region-based controller for our visual servoing scenario (Siebel et al., 1999). The main idea is to replace the constant dampening  $\lambda$  for  $\Delta u_n$  with a variable dampening  $\lambda_n$ :

$$u_n := \lambda_n \cdot \Delta u_n = \lambda_n \cdot J_n^+ \Delta y_n. \quad (39)$$

The goal is to adapt  $\lambda_n$  before each step to balance the avoidance of model errors (by making small steps) and the fast movement to the goal (by making large steps).

In order to achieve this balance we define an *actual model error*  $e_n$  which is set in relation to a *desired (maximum) model error*  $e_{des}^2$  to adapt a bound  $\alpha_n$  for the movement of projected object points on the sensor. Using this purely image-based formulation has advantages, e.g. having a measure to avoid movements that lead to losing object markings from the camera’s field of view.

Our algorithm is explained in Figure 9 for one object marking. We wish to calculate a robot command to move such that the current point position on the sensor moves to its desired position. In step ①, we calculate an undamped robot movement  $\Delta u_n$  to move as close to this goal as possible ( $\Delta \tilde{y}_n := \Delta y_n$ ) according to an Image Jacobian  $J_n$ :

$$\Delta u_n := J_n^+ \Delta y_n. \quad (40)$$

This gives us a predicted movement  $\ell_n$  on the sensor, which we define as the maximum movement on the sensor for all  $M$  markings:

$$\ell_n := \max_{i=1, \dots, M} \left\| \begin{bmatrix} (J_n \Delta u_n)_{2i-1} \\ (J_n \Delta u_n)_{2i} \end{bmatrix} \right\|_2, \quad (41)$$

where the subscripts to the vector  $J_n \Delta u_n$  signify a selection of its components.

Before executing the movement we restrict it in step ② such that the distance on the sensor is less or equal to a current limit  $\alpha_n$ :

$$\begin{aligned} u_n &:= \lambda_n \cdot \Delta u_n \\ &= \min \left\{ 1, \frac{\alpha_n}{\ell_n} \right\} \cdot J_n^+ \Delta y_n. \end{aligned} \quad (42)$$

<sup>2</sup> While the name “desired error” may seem unintuitive the name is chosen intentionally since the  $\alpha$  adaptation process (see below) can be regarded as a control process to have the robot system reach exactly this amount of error, by controlling the value of  $\alpha_n$ .



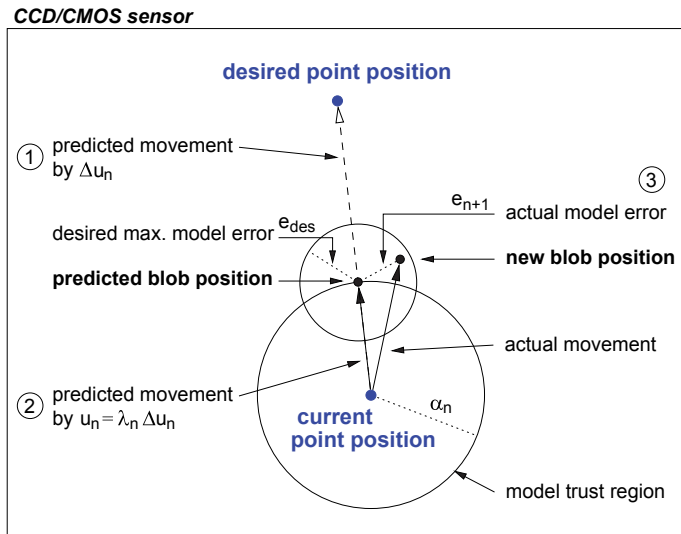


Fig. 9. Generation of a robot command by the trust region controller: view of the image sensor with a projected object marking

After this restricted movement is executed by the robot we obtain new measurements  $y_{n+1}$  and thereby the actual movement and model (prediction) error  $e_{n+1}$  (3), which we again define as the maximum deviation on the sensor for  $M > 1$  markings:

$$e_{n+1} := \max_{i=1, \dots, M} \left\| \begin{bmatrix} (\hat{y}_{n+1})_{2i-1} \\ (\hat{y}_{n+1})_{2i} \end{bmatrix} - \begin{bmatrix} (y_{n+1})_{2i-1} \\ (y_{n+1})_{2i} \end{bmatrix} \right\|_2. \quad (43)$$

where  $\hat{y}_{n+1}$  is the vector of predicted positions on the sensor,

$$\hat{y}_{n+1} := y_n + J_n u_n. \quad (44)$$

The next step is the adaptation of our restriction parameter  $\alpha_n$ . This is done by comparing the model error  $e_{n+1}$  with a given desired (maximum admissible) error  $e_{des}$ :

$$r_{n+1} := \frac{e_{n+1}}{e_{des}} \quad (45)$$

where  $r_n$  is called the *relative model error*. A small value signifies a good agreement of model and reality. In order to balance model agreement and a speedy control we adjust  $\alpha_n$  so as to achieve  $r_n = 1$ . Since we have a linear system model we can set

$$\alpha_{n+1} := \alpha_n \cdot \frac{e_{des}}{e_{n+1}} = \frac{\alpha_n}{r_{n+1}} \quad (46)$$

with an additional restriction on the change rate,  $\frac{\alpha_{n+1}}{\alpha_n} \leq 2$ . In practice, it may make sense to define minimum and maximum values  $\alpha_{\min}$  and  $\alpha_{\max}$  and set  $\alpha_0 := \alpha_{\min}$ .

In the example shown in Figure 9 the actual model error is smaller than  $e_{des}$ , so  $\alpha_{n+1}$  can be larger than  $\alpha_n$ .

```

Let  $n := 0$ ;  $\alpha_0 := \alpha_{\text{start}}$ ;  $y^*$  given
Measure current image features  $y_n$  and calculate  $\Delta y_n := y^* - y_n$ 
WHILE  $\|\Delta y_n\|_\infty \geq \varepsilon$ 
  Calculate  $J_n$ 
  IF  $n > 0$ 
    Calculate relative model error  $r_n$  via (43)
    Adapt  $\alpha_n$  by (46)
  END IF

  Calculate  $u_{\text{sd}_n} := J_n^T \Delta y_n$ ,  $\lambda_n := \frac{\|u_{\text{sd}_n}\|}{\ell_{\text{sd}_n}}$  and  $u_{\text{gn}_n} := J_n^+ \Delta y_n$ 

  Calculate  $u_{\text{dl}_n}$  via (52)
  Send control command  $u_{\text{dl}_n}$  to the robot
  Measure  $y_{n+1}$  and calculate  $\Delta y_{n+1}$ ; let  $n := n + 1$ 
END WHILE

```

Fig. 10. Algorithm: Image-based Visual Servoing with the Dogleg Algorithm

### 3.3.1.1 Remark:

By restricting the movement on the sensor we have implicitly defined the set  $\mathcal{U}(x_n)$  of admissible control commands in the state  $x_n$  as in equation (33). This  $\mathcal{U}(x_n)$  is the trust region of the model  $J_n$ .

### 3.3.2 A Dogleg Trust Region Controller

Powell (1970) describes the so-called *Dogleg Method* (a term known from golf) which can be regarded as a variant of the standard trust region method (Fletcher, 1987; Madsen et al., 1999). Just like in the trust region method above, a current model error is defined and used to adapt a trust region. Depending on the model error, the controller varies between a Gauss-Newton and a gradient (steepest descent) type controller.

The undamped Gauss-Newton step  $u_{\text{gn}_n}$  is calculated as before:

$$u_{\text{gn}_n} = J_n^+ \Delta y_n, \quad (47)$$

and the steepest descent step  $u_{\text{sd}_n}$  is given by

$$u_{\text{sd}_n} = J_n^T \Delta y_n. \quad (48)$$

The dampening factor  $\lambda_n$  is set to

$$\lambda_n := \frac{\|u_{\text{sd}_n}\|_2^2}{\ell_{\text{sd}_n}} \quad (49)$$

where again

$$\ell_{\text{sd}_n} := \max_{i=0, \dots, M} \left\| \begin{pmatrix} (\Delta \hat{y}_{\text{sd}_n})_{2i-1} \\ (\Delta \hat{y}_{\text{sd}_n})_{2i} \end{pmatrix} \right\|_2^2 \quad (50)$$



Fig. 11. Experimental setup with Thermo CRS F3 robot, camera and marked object

is the maximum predicted movement on the sensor, here the one caused by the steepest descent step  $u_{sd_n}$ . Analogously, let

$$\ell_{gn_n} := \max_{i=0, \dots, M} \left\| \begin{pmatrix} (\Delta \hat{y}_{gn_n})_{2i-1} \\ (\Delta \hat{y}_{gn_n})_{2i} \end{pmatrix} \right\|_2^2 \quad (51)$$

be the maximum predicted movement by the Gauss Newton step. With these variables the dog leg step  $u_n = u_{dl_n}$  is calculated as follows:

$$u_{dl_n} := \begin{cases} u_{gn_n} & \text{if } \ell_{gn_n} \leq \alpha_n \\ \alpha_n \frac{u_{sd_n}}{\|u_{sd_n}\|_2} & \text{if } \ell_{gn_n} > \alpha_n \text{ and } \ell_{sd_n} \geq \alpha_n \\ \lambda_n u_{sd_n} + \beta_n (u_{gn_n} - \lambda_n u_{sd_n}) & \text{else} \end{cases} \quad (52)$$

where in the third case  $\beta_n$  is chosen such that the maximum movement on the sensor has length  $\alpha_n$ .

The complete dogleg algorithm for visual servoing is shown in Figure 10.

## 4. Experimental Evaluation

### 4.1 Experimental Setup and Test Methods

The robot setup used in the experimental validation of the presented controllers is shown in Figure 11. Again an eye-in-hand configuration and an object with 4 identifiable markings are used. Experiments were carried out both on a Thermo CRS F3 (pictured here) and on a Unimation Stäubli RX-90 (Figure 2 at the beginning of the chapter). In the following only

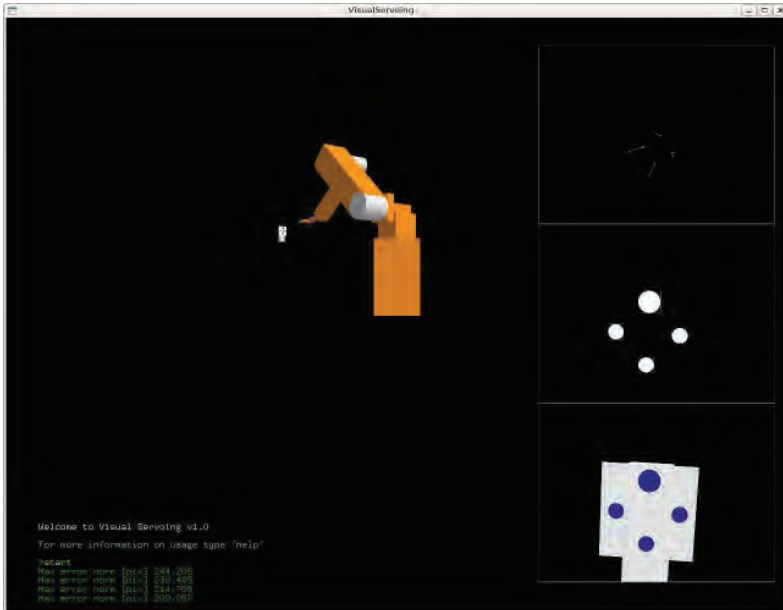


Fig. 12. OpenGL Simulation of camera-robot system with simulated camera image (bottom right), extracted features (centre right) and trace of objects markings on the sensor (top right)

the CRS F3 experiments are considered; the results with the Stäubli RX-90 were found to be equivalent. The camera was a Sony DFW-X710 with IEEE1394 interface,  $1024 \times 768$  pixel resolution and an  $f = 6.5$  mm lens.

In addition to the experiments with a real robot two types of simulations were used to study the behaviour of controllers and models in detail. In our *OpenGL Simulation*<sup>3</sup>, see Figure 12, the complete camera-robot system is modelled. This includes the complete robot arm with inverse kinematics, rendering of the camera image in a realistic resolution and application of the same image processing algorithms as in the real experiments to obtain the image features. Arbitrary robots can be defined by their Denavit-Hartenberg parameters (cf. Spong et al., 2005) and geometry in an XML file. The screenshot above shows an approximation of the Stäubli RX-90.

The second simulation we use is the *Multi-Pose Test*. It is a system that uses the exact model as derived in Section 2.2, without the image generation and digitisation steps as in the OpenGL Simulation. Instead, image coordinates of objects points as seen by the camera are calculated directly with the pinhole camera model. Noise can be added to these measurements in order to examine how methods react to these errors. Due to the small computational complexity of the Multi-Pose Test it can be, and has been applied to many start and teach pose combinations (in our experiments, 69,463 start poses and 29 teach poses). For a given algorithm and parameter set the convergence behaviour (success rate and speed) can thus be studied on a statistically relevant amount of data.

<sup>3</sup> The main parts of simulator were developed by Andreas Jordt and Falko Kellner when they were students in the Cognitive Systems Group.

## 4.2 List of Models and Controllers Tested

In order to test the advantages and disadvantages of the models and controllers presented above we combine them in the following way:

Short Name	Controller	Model	Parameters
Trad const	Traditional	$\Delta y_n \approx J^* u$	$\lambda = 0.2$
Trad dyn	Traditional	$\Delta y_n \approx J_n u$	$\lambda = 0.1$ , sometimes $\lambda = 0.07$
Trad PMJ	Traditional	$\Delta y_n \approx \frac{1}{2}(J_n + J^*) u$	$\lambda = 0.25$
Trad MPJ	Traditional	$u \approx \frac{1}{2}(J_n^+ + J^{*+}) \Delta y_n$	$\lambda = 0.15$
Trad cyl	Traditional	$\Delta y_n \approx \bar{J}_n u$ (cylindrical)	$\lambda = 0.1$
TR const	Trust-Region	$\Delta y_n \approx J^* u$	$\alpha_0 = 0.09, e_{des} = 0.18$
TR dyn	Trust-Region	$\Delta y_n \approx J_n u$	$\alpha_0 = 0.07, e_{des} = 0.04$
TR PMJ	Trust-Region	$\Delta y_n \approx \frac{1}{2}(J_n + J^*) u$	$\alpha_0 = 0.07, e_{des} = 0.09$
TR MPJ	Trust-Region	$u \approx \frac{1}{2}(J_n^+ + J^{*+}) \Delta y_n$	$\alpha_0 = 0.05, e_{des} = 0.1$
TR cyl	Trust-Region	$\Delta y_n \approx \bar{J}_n u$ (cylindrical)	$\alpha_0 = 0.04, e_{des} = 0.1$
Dogleg const	Dogleg	$u \approx J^{*+} \Delta y_n$ and $u \approx J_n^T \Delta y_n$	$\alpha_0 = 0.22, e_{des} = 0.16, \lambda = 0.5$
Dogleg dyn	Dogleg	$u \approx J_n^+ \Delta y_n$ and $u \approx J_n^T \Delta y_n$	$\alpha_0 = 0.11, e_{des} = 0.28, \lambda = 0.5$
Dogleg PMJ	Dogleg	$\Delta y_n \approx \frac{1}{2}(J_n + J^*) u$ and $u \approx J_n^T \Delta y_n$	$\alpha_0 = 0.29, e_{des} = 0.03, \lambda = 0.5$
Dogleg MPJ	Dogleg	$u \approx \frac{1}{2}(J_n^+ + J^{*+}) \Delta y_n$ and $u \approx J_n^T \Delta y_n$	$\alpha_0 = 0.3, e_{des} = 0.02, \lambda = 0.5$

Here we use the definitions as before. In particular,  $J_n$  is the dynamical Image Jacobian as defined in the current pose, calculated using the current distances to the object,  $z_i$  for marking  $i$ , and the current image features in its entries. The distance to the object is estimated in the real experiments using the known relative distances of the object markings, which yields a fairly precise estimate in practice.  $J^*$  is the constant Image Jacobian, defined in the teach (goal) pose  $x^*$ , with the image data  $y^*$  and distances at that pose.  $\Delta y_n = y_{n+1} - y_n$  is the change in the image predicted by the model with the robot command  $u$ .

The values of the parameters detailed above were found to be useful parameters in the Multi-Pose Test. They were therefore used in the experiments with the real robot and the OpenGL Simulator. See below for details on how these values were obtained.

$\lambda$  is the constant dampening factor applied as the last step of the controller output calculation. The Dogleg controller did not converge in our experiments without such an additional dampening which we set to 0.5. The Trust-Region controller works without additional dampening.  $\alpha_0$  is the start and minimum value of  $\alpha_n$ . These, as well as the desired model error  $e_{des}$  are given in mm on the sensor. The sensor measures  $4.8 \times 3.6$  mm which means that at its  $1024 \times 768$  pixel resolution  $0.1 \text{ mm} \approx 22$  pixels after digitisation.

## 4.3 Experiments and Results

The Multi-Pose Test was run first in order to find out which values of parameters are useful for which controller/model combination. 69,463 start poses and 29 teach poses were combined randomly into 69,463 fixed pairs of tasks that make up the training data. We studied the following two properties and their dependence on the algorithm parameters:

1. *Speed*: The number of iterations (steps/robot movements) needed for the algorithm to reach its goal. The mean number of iterations over all successful trials is measured.
2. *Success rate*: The percentage of experiments that reached the goal. Those runs where an object marking was lost from the camera view by a movement that was too large and/or mis-directed were considered not successful, as were those that did not reach the goal within 100 iterations.

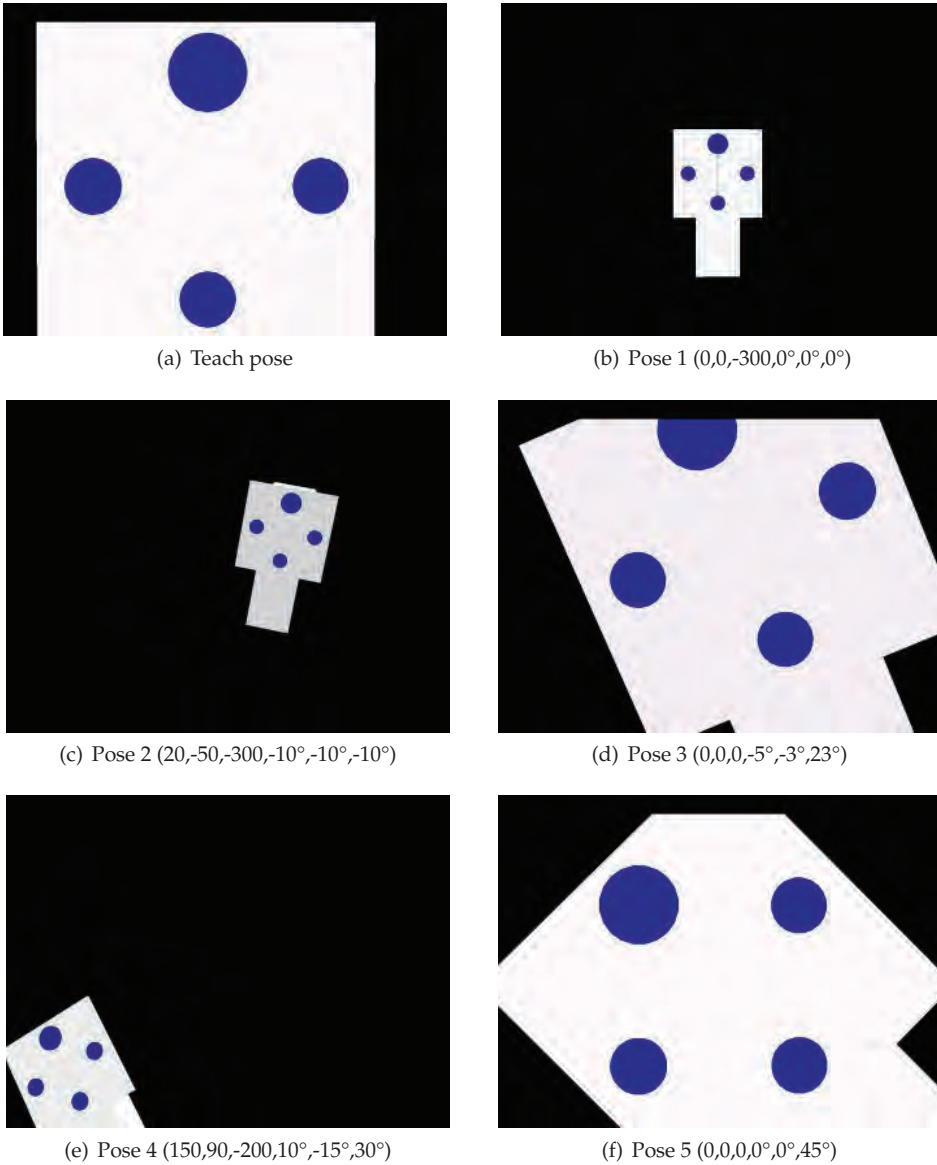


Fig. 13. Teach and start poses used in the experiments; shown here are simulated camera images in the OpenGL Simulator. Given for each pose is the relative movement in {C} from the teach pose to the start pose. Start pose 4 is particularly difficult since it requires both a far reach and a significant rotation by the robot. Effects of the linearisation of the model or errors in its parameters are likely to cause a movement after which an object has been lost from the camera's field of view. Pose 5 is a pure rotation, chosen to test for the retreat-advance problem.

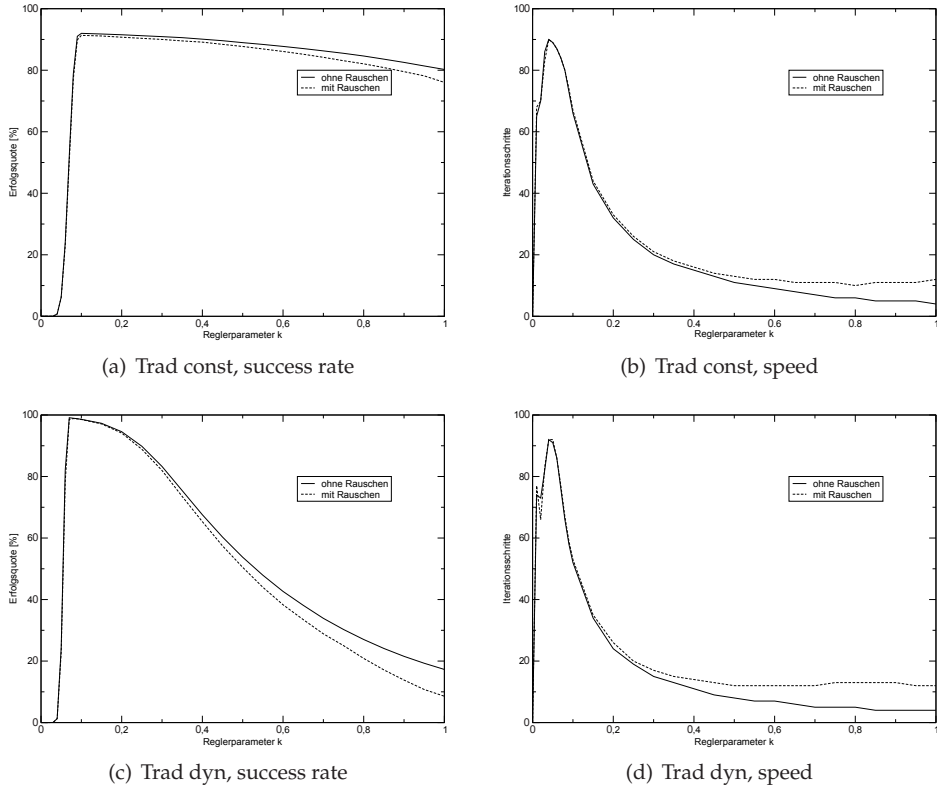
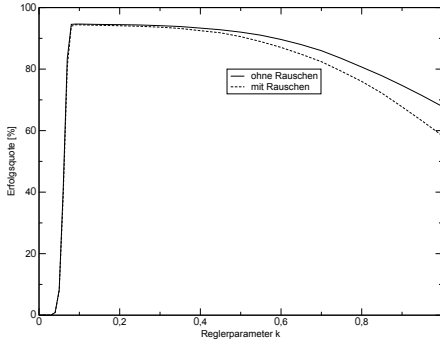


Fig. 14. Multi-Pose Test: Traditional Controller with const. and dyn. Jacobian. Success rate and average speed (number of iterations) are plotted as a function of the damping parameter  $\lambda$ .

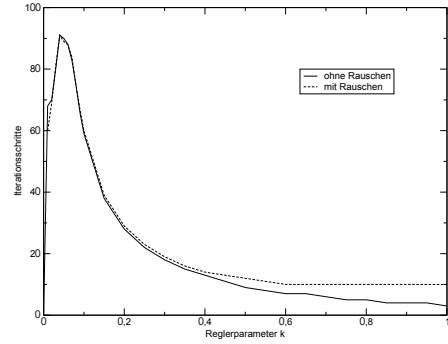
Using the optimal parameters found by the Multi-Pose Test we ran experiments on the real robot. Figure 13 shows the camera images (from the OpenGL simulation) in the teach pose and five start poses chosen such that they cover the most important problems in visual servoing. The OpenGL simulator served as an additional useful tool to analyse why some controllers with some parameters would not perform well in a few cases.

#### 4.4 Results with Non-Adaptive Controllers

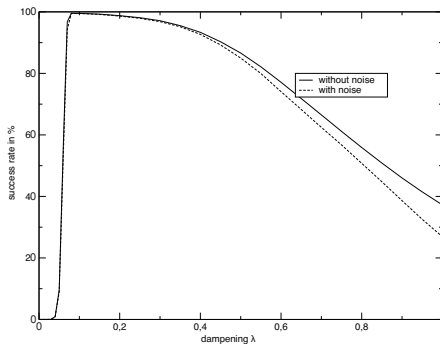
Figures 14 and 15 show the results of the Multi-Pose Test with the Traditional Controller using different models. For the *success rates* it can be seen that with  $\lambda$ -values below a certain value  $\approx 0.06$ – $0.07$  the percentages are very low. On the other hand, raising  $\lambda$  above  $\approx 0.08$ – $0.1$  also significantly decreases success rates. The reason is the proportionality of image error and (length of the) robot movement inherent in the control law with its constant factor  $\lambda$ . During the course of the servoing process the norm of the image error may vary by as much as a factor of 400. The controller output varies proportionally. This means that at the beginning of the control process very large movements are carried out, and very small movements at the end.



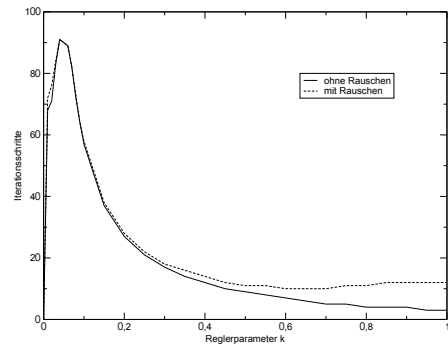
(a) Trad PMJ, success rate



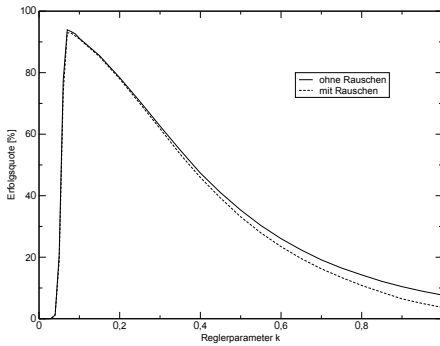
(b) Trad PMJ, speed



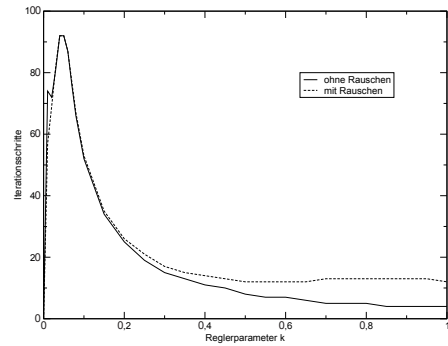
(c) Trad MJP, success rate



(d) Trad MJP, speed



(e) Trad cyl, success rate



(f) Trad cyl, speed

Fig. 15. Multi-Pose Test: Traditional Controller with PMJ, MPJ and cylindrical models. Shown here are again the success rate and speed (average number of iterations of successful runs) depending on the constant dampening factor  $\lambda$ . As before, runs that did not converge in the first 100 steps were considered unsuccessful.



Controller	param. $\lambda$	Real Robot start pose					OpenGL Sim. start pose					Multi-Pose speed (iter.)		success (%)
		1	2	3	4	5	1	2	3	4	5			
Trad const	0.2	49	55	21	46	31	44	44	23	44	23	32	91.53	
Trad dyn	0.1	63	70	48	$\infty$	58	46	52	45	$\infty$	47	52	98.59	
	0.07				121					81		76	99.11	
Trad MJP	0.15	41	51	33	46	37	35	39	31	41	32	37	99.27	
Trad PMJ	0.25	29	29	17	$\infty$	35	26	26	18	$\infty$	32	38	94.52	
Trad cyl	0.1	59	$\infty$	50	70	38	46	49	49	58	49	52	91.18	

Table 1. All results, Traditional Controller, optimal value of  $\lambda$ . “ $\infty$ ” means no convergence

The movements at the beginning need strong dampening (small  $\lambda$ ) in order to avoid large mis-directed movements (Jacobians usually do not have enough validity for 400 mm movements), those at the end need little or no dampening ( $\lambda$  near 1) when only a few mm are left to move. The version with the constant image Jacobian has a better behaviour for larger ( $\geq 0.3$ ) values of  $\lambda$ , although even the optimum value of  $\lambda = 0.1$  only gives a success rate of 91.99%. The behaviour for large  $\lambda$  can be explained by  $J^*$ 's smaller validity away from the teach pose; when the robot is far away it suggests smaller movements than  $J_n$  would. In practise this acts like an additional dampening factor that is stronger further away from the object.

The adaptive Jacobian gives the controller a significant advantage if  $\lambda$  is set well. For  $\lambda = 0.07$  the success rate is 99.11%, albeit with a speed penalty, at as many as 76 iterations. With  $\lambda = 0.1$  this decreases to 52 at 98.59% success rate.

The use of the PMJ and MJP models show again a more graceful degradation of performance with increasing  $\lambda$  than  $J_n$ . The behaviour with PMJ is comparable to that with  $J^*$ , with a maximum of 94.65% success at  $\lambda = 0.1$ ; here the speed is 59 iterations. Faster larger  $\lambda$ , e.g. 0.15 which gives 38 iterations, the success rate is still at 94.52%. With MJP a success rate of 99.53% can be achieved at  $\lambda = 0.08$ , however, the speed is slow at 72 iterations. At  $\lambda = 0.15$  the controller still holds up well with 99.27% success and significantly less iterations: on average 37.

Using the cylindrical model the traditional controller's success is very much dependant on  $\lambda$ . The success rate peaks at  $\lambda = 0.07$  with 93.94% success and 76 iterations; a speed 52 can be achieved at  $\lambda = 0.1$  with 91.18% success. Overall the cylindrical model does not show an overall advantage in this test.

Table 1 shows all results for the traditional controller, including real robot and OpenGL results. It can be seen that even the most simple pose takes at least 29 steps to solve. The Trad MJP method is the clearly the winner in this comparison, with a 99.27% success rate and on average 37 iterations. Pose 4 holds the most difficulties, both in the real world and in the OpenGL simulation. In the first few steps a movement is calculated that makes the robot lose the object from the camera's field of view. The Traditional Controller with the dynamical Jacobian achieves convergence only when  $\lambda$  is reduced from 0.1 to 0.07. Even then the object marking comes close to the image border during the movement. This can be seen in Figure 16 where the trace of the centre of the object markings on the sensor is plotted. With the cylindrical model the controller moves the robot in a way which avoids this problem. Figure 16(b) shows that there is no movement towards the edge of the image whatsoever.

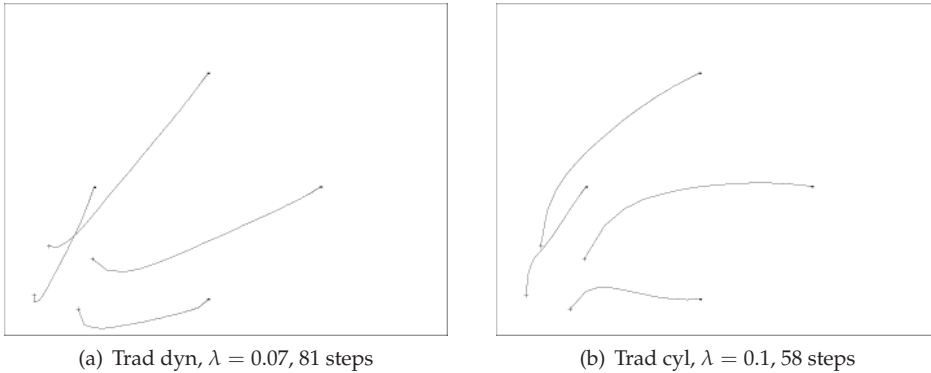


Fig. 16. Trad. Controller, dyn. and cyl. model, trace of markings on sensor, pose 4 (OpenGL).

#### 4.5 Results with Adaptive Controllers

In this section we wish to find out whether the use of dynamical dampening by a limitation of the movement on the sensor (image-based trust region methods) can speed up the slow convergence of the traditional controller. We will examine the Trust-Region controller first, then the Dogleg controller.

Figure 17 shows the behaviour for the constant and dynamical Jacobians as a function of the main parameter, the desired maximum model error  $e_{des}$ . The success rate for both variants is only slightly dependent on  $e_{des}$ , with rates over 91 % (Trust const) and 99 % (Trust dyn) for the whole range of values from 0.01 to 0.13 mm when run without noise. The speed is significantly faster than with the Traditional Controller at 13 iterations ( $e_{des} = 0.18$ , 91.46 % success) and 8 iterations ( $e_{des} = 0.04$ , 99.37 % success), respectively. By limiting the step size dynamically the Trust Region methods calculate smaller movements than the Traditional Controller at the beginning of the experiment but significantly larger movements near the end. This explains the success rate (no problems at beginning) and speed advantage (no active dampening towards the end). The use of the mathematically more meaningful dynamical model  $J_n$  helps here since the Trust Region method avoids the large mis-directed movements far away from the target without the need of the artificial dampening through  $J^*$ . The Trust/dyn. combination shows a strong sensitivity to noise; this is mainly due to the amplitude of the noise (standard deviation 1 pixel) which exceeds the measurement errors in practice when the camera is close to the object. This results in convergence problems and problems detecting convergence when the robot is very close to its goal pose. In practise (see e.g. Table 2 below) the controller tends to have fewer problems. In all five test poses, even the difficult pose 4 the controller converges with both models without special adjustment (real world and OpenGL), with a significant speed advantage of the dynamical model. In pose 5 both are delayed by the retreat-advance problem but manage to reach the goal successfully.

The use of the MJP model helps the Trust-Region Controller to further improve its results. Success rates (see Figure 18) are as high as 99.68 % at  $e_{des} = 0.01$  (on average 16 iterations), with a slightly decreasing value when  $e_{des}$  is increased: still 99.58 % at  $e_{des} = 0.1$  (7 iterations, which makes it the fastest controller/model combination in our tests).

As with the Traditional Controller the use of the PMJ and cylindrical model do not show overall improvements for visual servoing over the dynamical method. The results, are also

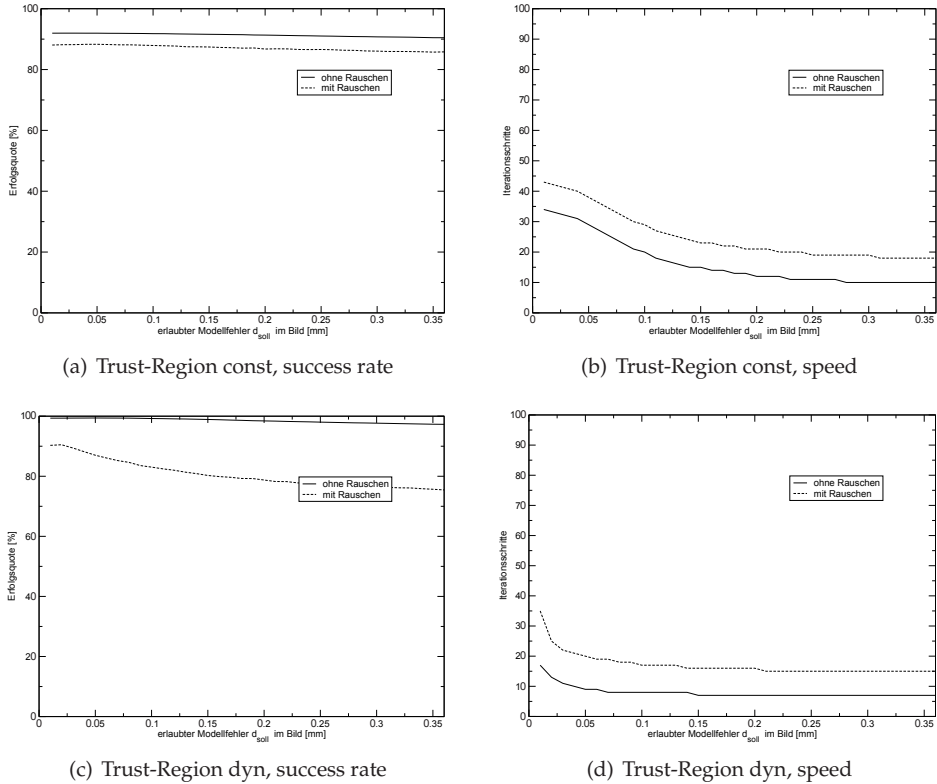
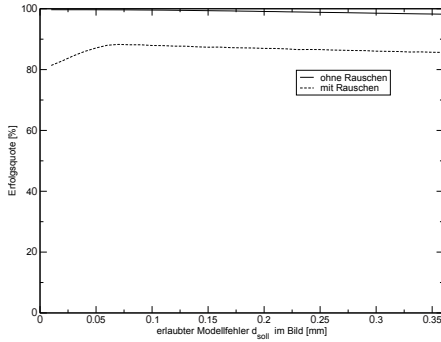


Fig. 17. Multi-Pose Test: Trust-Region Controller with const. and dyn. Jacobian

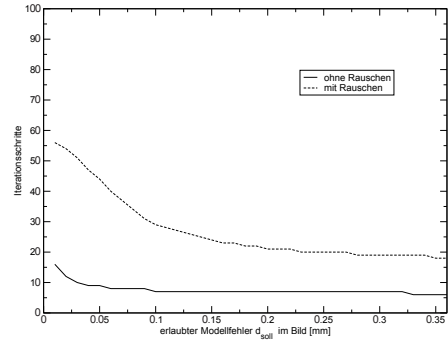
shown also in Figure 18. Table 2 details the results for all three types of tests. It can be seen that while both models have on average better results than with the constant Jacobian they do have convergence problems that show in the real world. In pose 2 (real robot) the cylindrical model causes the controller to calculate an unreachable pose for the robot at the beginning, which is why the experiment was terminated and counted as unsuccessful.

The Dogleg Controller shows difficulties irrespective of the model used. Without an additional dampening with a constant  $\lambda = 0.5$  no good convergence could be achieved. Even with dampening its maximum success rate is only 85%, with  $J^*$  (at an average of 10 iterations). Details for this combination are shown in Figure 19 where we see that the results cannot be improved by adjusting the parameter  $e_{des}$ . With other models only less than one in three poses can be solved, see results in Table 2.

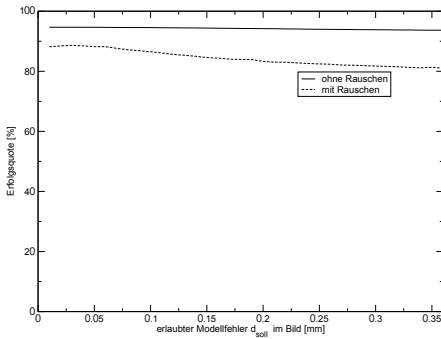
A thorough analysis showed that the switching between gradient descent and Gauss-Newton steps causes the problems for the Dogleg controller. This change in strategy can be seen in Figure 20 where again the trace of projected object markings on the sensor is shown (from the real robot system). The controller first tries to move the object markings towards the centre of the image, by applying gradient descent steps. This is achieved by changing yaw and pitch angles only. Then the Dogleg step, i.e. a combination of gradient descent and Gauss-Newton



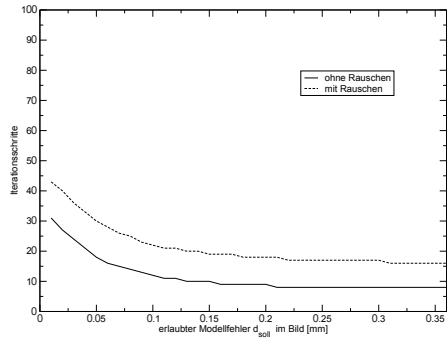
(a) Trust-Region MJP, success rate



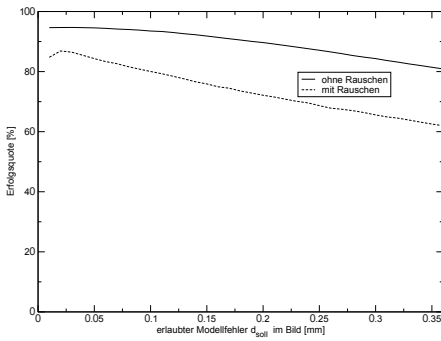
(b) Trust-Region MJP, speed



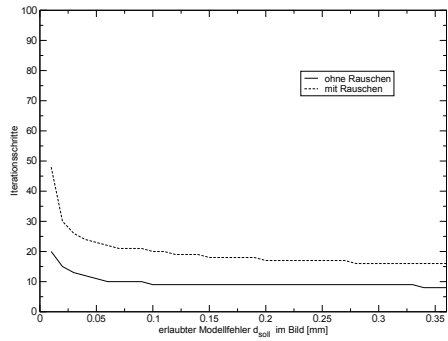
(c) Trust-Region PMJ, success rate



(d) Trust-Region PMJ, speed



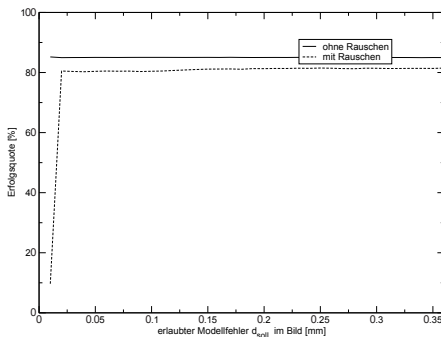
(e) Trust-Region cyl, success rate



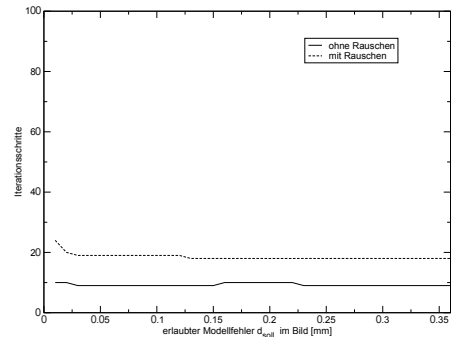
(f) Trust-Region cyl, speed

Fig. 18. Multi-Pose Test: Trust-Region Controller with PMJ, MPJ and cylindrical model. Plotted are the success rate and the speed (average number of iterations of successful runs) depending on the desired (maximum admissible) error,  $e_{des}$ .

Controller	param.		Real Robot start pose					OpenGL Sim. start pose					Multi-Pose	
	$\alpha_{\text{start}}$	$\ell_{\text{des}}$	1	2	3	4	5	1	2	3	4	5	speed (iter.)	success (%)
Trust const	0.09	0.18	22	29	11	39	7	20	26	6	31	7	13	91.46
Trust dyn	0.07	0.04	10	15	9	17	17	9	12	7	14	6	8	99.37
Trust MJP	0.05	0.1	8	9	11	13	7	7	9	6	11	5	7	99.58
Trust PMJ	0.07	0.09	21	28	7	$\infty$	13	20	25	6	$\infty$	5	13	94.57
Trust cyl	0.04	0.1	10	$\infty$	7	11	15	8	18	6	11	6	9	93.5
Dogleg const	0.22	0.16	19	24	8	$\infty$	12	17	25	4	21	9	10	85.05
Dogleg dyn	0.11	0.28	13	$\infty$	$\infty$	$\infty$	13	8	$\infty$	6	$\infty$	16	9	8.4
Dogleg MJP	0.3	0.02	$\infty$	$\infty$	10	$\infty$	13	$\infty$	$\infty$	5	$\infty$	7	8	26.65
Dogleg PMJ	0.29	0.03	14	13	5	$\infty$	12	9	13	5	14	7	8	31.47

Table 2. All results, Trust-Region and Dogleg Controllers. “ $\infty$ ” means no success.

(a) Dogleg const, success rate



(b) Dogleg const, speed

Fig. 19. Multi-Pose Test: Dogleg Controller with constant Image Jacobian

step (with the respective Jacobian), is applied. This causes zigzag movements on the sensor. These are stronger when the controller switches back and forth between the two approaches, which is the case whenever the predicted and actual movements differ by a large amount.

## 5. Analysis and Conclusion

In this chapter we have described and analysed a number of visual servoing controllers and models of the camera-robot system used by these controllers. The inherent problem of the traditional types of controllers is the fact that these controllers do not adapt their controller output to the current state in which the robot is: far away from the object, close to the object, strongly rotated, weakly rotated etc. They also cannot adapt to the strengths and deficiencies of the model, which may also vary with the current system state. In order to guarantee successful robot movements towards the object these controllers need to restrict the steps the robot takes, and they do so by using a constant scale factor (“dampening”). The constancy of this scale factor is a problem when the robot is close to the object as it slows down the movements too much.

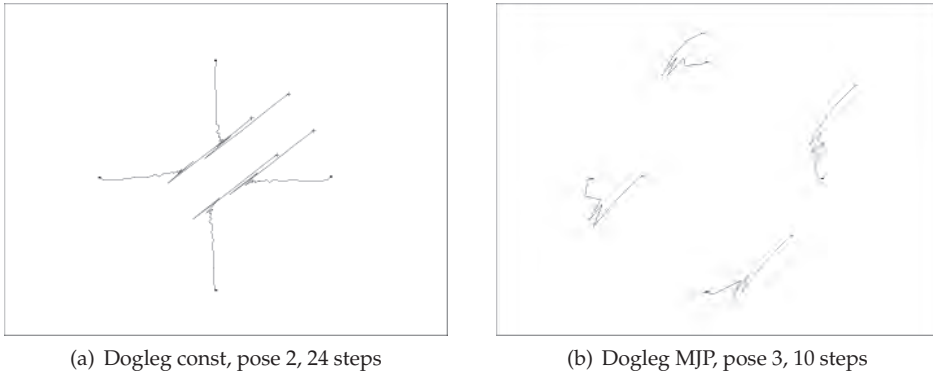


Fig. 20. Dogleg, const and MJP model, trace of markings on sensor, poses 2 and 3 (real robot).

Trust-region based controllers successfully overcome this limitation by adapting the dampening factor in situations where this is necessary, but only in those cases. Therefore they achieve both a better success rate and a significantly higher speed than traditional controllers.

The Dogleg controller which was also tested does work well with some poses, but on average has much more convergence problems than the other two types of controllers.

Overall the Trust-Region controller has shown the best results in our tests, especially when combined with the MJP model, and almost identical results when the dynamical image Jacobian model is used. These models are more powerful than the constant image Jacobian which almost always performs worse.

The use of the cylindrical and PMJ models did not prove to be helpful in most cases, and in those few cases where they have improved the results (usually pure rotations, which is unlikely in most applications) the dynamical and MJP models also achieved good results.

The results found in experiments with a real robot and those carried out in two types of simulation agree on these outcomes.

## Acknowledgements

Part of the visual servoing algorithm using a trust region method presented in this chapter was conceived in 1998–1999 while the first author was at the University of Bremen. The advice of Oliver Lang and Fabian Wirth at that time is gratefully acknowledged.

## 6. References

- François Chaumette. Potential problems of stability and convergence in image-based and position-based visual servoing. In David J Kriegmann, Gregory D Hager, and Stephen Morse, editors, *The Confluence of Vision and Control*, pages 66–78. Springer Verlag, New York, USA, 1998.
- François Chaumette and Seth Hutchinson. Visual servoing and visual tracking. In Bruno Siciliano and Oussama Khatib, editors, *Springer Handbook of Robotics*, pages 563–583. Springer Verlag, Berlin, Germany, 2008.

- Peter I. Corke and Seth A. Hutchinson. A new partitioned approach to image-based visual servo control. *IEEE Transactions on Robotics and Automation*, 237(4):507–515, August 2001.
- Peter Deuffhard and Andreas Hohmann. *Numerical Analysis in Modern Scientific Computing: An Introduction*. Springer Verlag, New York, USA, 2nd edition, 2003.
- Roger Fletcher. *Practical Methods of Optimization*. John Wiley & Sons, New York, Chichester, 2nd edition, 1987.
- Seth Hutchinson, Gregory D Hager, and Peter Corke. A tutorial on visual servo control. Tutorial notes, Yale University, New Haven, USA, May 1996.
- Masami Iwatsuki and Norimitsu Okiyama. A new formulation of visual servoing based on cylindrical coordinate system. *IEEE Transactions on Robotics*, 21(2):266–273, April 2005.
- Martin Jägersand. Visual servoing using trust region methods and estimation of the full coupled visual-motor Jacobian. In *Proceedings of the IASTED Applications of Control and Robotics, Orlando, USA*, pages 105–108, January 1996.
- Kenichi Kanatani. *Statistical Optimization for Geometric Computation: Theory and Practice*. Elsevier Science, Amsterdam, The Netherlands, 1996.
- Kaj Madsen, Hans Bruun Nielsen, and Ole Tingleff. Methods for non-linear least squares problems. Lecture notes, Department of Informatics and Mathematical Modelling, Technical University of Denmark, Lyngby, Denmark, 1999.
- Ezio Malis. Improving vision-based control using efficient second-order minimization techniques. In *Proceedings of 2004 International Conference on Robotics and Automation (ICRA 2004), New Orleans, USA*, pages 1843–1848, April 2004.
- Michael J D Powell. A hybrid method for non-linear equations. In Philip Rabinowitz, editor, *Numerical Methods for Non-linear Algebraic Equations*, pages 87–114. Gordon and Breach, London, 1970.
- Andrew P Sage and Chelsea C White. *Optimum Systems Control*. Prentice-Hall, Englewood Cliffs, USA, 2nd edition, 1977.
- Nils T Siebel, Oliver Lang, Fabian Wirth, and Axel Gräser. Robuste Positionierung eines Roboters mittels Visual Servoing unter Verwendung einer Trust-Region-Methode. In *Forschungsbericht Nr. 99-1 der Deutschen Forschungsvereinigung für Meß-, Regelungs- und Systemtechnik (DFMRS) e.V.*, pages 23–39, Bremen, Germany, November 1999.
- Mark W Spong, Seth Hutchinson, and Mathukumalli Vidyasagar. *Robot Modeling and Control*. John Wiley & Sons, New York, Chichester, 2005.
- Lee E Weiss, Arthur C Sanderson, and Charles P Neuman. Dynamic sensor-based control of robots with visual feedback. *IEEE Journal of Robotics and Automation*, 3(5):404–417, October 1987.

