# BINOCULAR VISUAL SERVOING BASED ON LINEAR TIME-INVARIANT MAPPING

Takashi Mitsuda[1], Noriaki Maru[2], Kazunobu Fujikawa[1] and Fumio Miyazaki[1]

[1]Faculty of Engineering Science, Osaka University

[2]Faculty of System Engineering, Wakayama University

Author Note

Correspondence concerning this article should be addressed to Takashi Mitsuda,

College of Information Science and Engineering, Ritsumeikan University , 1-1-1 Noji-

Higashi, Kusatsu, Shiga, Japan 525-8577.

Phone: +81-77-561-5068

E-mail: mitsuda@is.ritsumei.ac.jp

# BINOCULAR VISUAL SERVOING BASED ON LINEAR TIME-INVARIANT MAPPING

Takashi Mitsuda[1], Noriaki Maru[2], Kazunobu Fujikawa[1] and Fumio Miyazaki[1]

[1]Faculty of Engineering Science, Osaka University
Machikaneyama 1-3, Toyonaka, Osaka 560, Japan
[2]Faculty of System Engineering, Wakayama University
930 Sakaedani, Wakayama 640, Japan
E-mail: mitsuda@robotics.me.es.osaka-u.ac.jp

**Abstract** We propose a simple visual servoing scheme based on the use of binocular visual space. When we use a hand-eye system which has a similar kinematic structure to a human being, we can approximate the transformation from a binocular visual space to a joint space of the manipulator as a linear time-invariant mapping. This relationship makes it possible to generate joint velocities from image observations using a constant linear mapping. This scheme is robust to calibration error, especially to camera turning, because it uses neither camera angles nor joint angles. Some experimental results are also shown to demonstrate the positioning precision remained unchanged despite the calibration error.

## 1   Introduction

Vision is indispensable for the intelligent robot that performs specified tasks in an uncertain and changing environment. In the recent past, the use of visual information in robot control has been an active research area. Various kinds of mechanisms of visual feedback have been proposed and they are called "*visual servoing*" in general terms[1]. It is most important for visual servoing to represent the mapping from robot coordinates to images. This mapping, in general, includes an intermediate stage, Cartesian coordinates. For example, a position-based visual servoing uses the visual image of the scene to represent the 3D environment in Cartesian coordinates[2]. On the other hand, an image-based visual servoing also uses Cartesian-to-image coordinate transformation(the image Jacobian) and the robot Jacobian [3][4][5]. Both approaches require many parameters of the robot and the camera in the transformation from or to Cartesian coordinates. This increases the amount of calculation, moreover, causes the weakness in the calibration error and the disturbance of the parameters. Image-based visual servoing is more robust than position-based visual servoing. However the convergence behavior is sensitive to the change in parameters, and an accurate estimate of the parameters is needed for the stable convergence.

Using free-standing cameras for visual servoing has a great advantage, because both of the end-effector and the target are observable. However the calibration of the hand-eye coordination is difficult, and it causes the problem noted above. If

we use an active vision system whose cameras turn dynamically, the problem will be more serious. Hager et al.[6] describe a stereo-based visual servoing system using free-standing cameras by the use of position estimator. Their system is robust to the static calibration error, but the influence of the changing parameters is not discussed, and active turning of the cameras is not mentioned. Hosoda et al.[7] propose a visual servoing scheme using a Jacobian estimator. It doesn't require any priori knowledge of the kinematic structure of the hand-eye system. However the dynamic change in the parameters affects the behavior of the system, because the system depends on the stable estimation. Hence, both approaches require stable measurement of the system such as camera angles and joint angles, and their disturbance affect the system behavior. Hollinghurst et al.[8] use an affine approximation to the inverse perspective transformation to compute the approximate Cartesian positions. Their system is robust to camera disturbance, but the active turning of the cameras also affects the system behavior.

In this paper, we focus on the hand-eye system whose cameras turn actively. Active vision system is useful to treat a real world for image processing system. However, applying it to a hand-eye system such as image-based visual servoing system described above is difficult, because the hand-eye coordination must be carried out dynamically. In this paper, we overcome this defect by the use of binocular visual space instead of Cartesian space for controlling the manipulator.

There is much evidence from a variety of experiments that the geometry of binocularly-perceived space is not Cartesian[9], and the planning of a reaching movement at visual targets is based on intrinsic coordinate systems[10, 11]. The binocular visual space is a model of this binocularly-perceived space that has been employed by phisiologists and psychologists[9]. In the binocular visual space, the binocular parallax and the horizontal direction serve as coordinates that specify the positions of point in physical eye-level plane.

In this paper, it is shown that the use of binocular visual space makes it possible to approximate the robot kinematics as a linear mapping when we use a hand-eye system which has a similar kinematic structure to a human being. This is a design problem of finding appropriate kinematic structure to allow the hand-eye system to treat the coordinate transformation in a simple and convenient manner. In addition, we propose a simple visual servoing scheme which generates joint velocities from image observations using a constant linear mapping. This scheme is robust to calibration error, especially to camera turning, because it uses neither camera angles nor joint angles. Stability of this scheme is also discussed on the basis of the invariant set theorem.

This paper first defines the hand-eye system which has a similar kinematic structure to a human being, then focuses on the approximation of the kinematics to a linear mapping. Next we show a simple visual servoing scheme based on this linear mapping. Experimental results are also provided to demonstrate the effectiveness of this scheme.

## 2 Approximation of the Kinematics as a Linear mapping

### 2.1 Model of the Hand-eye system

**Fig.1** shows a hand-eye system whose kinematic structure is similar to that of a human being. The manipulator consists of two links and two joints. The elbow joint has 1DOF, and the shoulder joint has 2DOF. The length of upper arm and forward arm is $L_u = 250[\text{mm}], L_f = 380[\text{mm}]$ respectively. The shoulder joint is located at the origin of world coordinate$\Sigma_0$. The two cameras are mounted on pan-tilt heads,

and the heads are mounted on a base frame which turns horizontally round the neck joint. The length of the baseline is 70[mm]. The focal length is 3[mm]. The center of the baseline is located at $(W, K, G) = (-200, -200, 0)$. These parameters are defined to be similar to those of a human being. They are the most suitable values for the approximation of the transformation from binocular visual space to joint space as a linear mapping. We describe in detail the relationship between these parameters and the accuracy of the approximation in Sec.2.5.
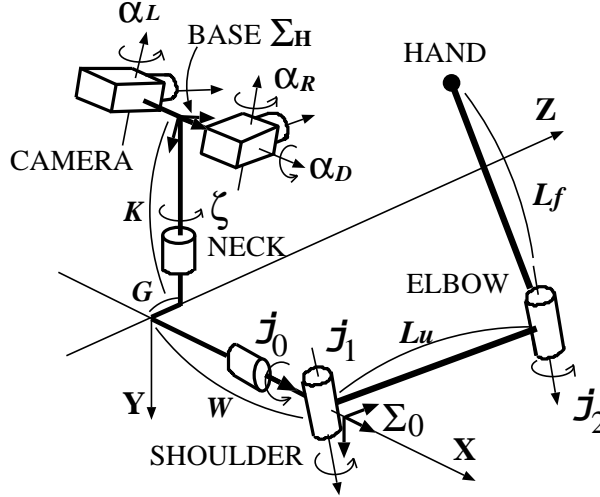


**Fig.1**: Model of the Hand-eye system

## 2.2   Joint Space

The kinematics of the manipulator is

$$
\begin{aligned}
x &= L_u \cos(j_1) + L_f \cos(j_1 + j_2) \\
y &= \tilde{z} \sin(j_0) \\
z &= \tilde{z} \cos(j_0) \\
\tilde{z} &= L_u \sin(j_1) + L_f \sin(j_1 + j_2).
\end{aligned} \tag{1}
$$

The continuous lines in **Fig.2** show the joint space which is projected onto Cartesian space(X-Z plane). It represents the nonlinearity in the transformation between the joint space and the Cartesian space.

## 2.3   Binocular Visual Space

The binocular visual space is defined as the vergence angle $\gamma$ and the viewing directions $\theta, \delta$(see **Fig.3**). This space has been employed by psychologists and physiologists as a model of binocularly-perceived space[9]. The binocular visual coordinate of a fixiation point is described as

$$
\boldsymbol{V} = \begin{bmatrix} \gamma \\ \theta \\ \delta \end{bmatrix} = \begin{bmatrix} \alpha_L - \alpha_R \\ (\alpha_L + \alpha_R)/2 \\ \alpha_D \end{bmatrix} \tag{2}
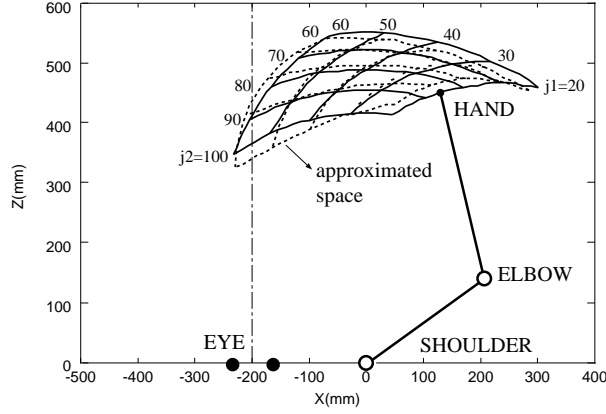$$

**Fig.2**: Joint Space projected onto Cartesian Space

where $\alpha_L, \alpha_R, \alpha_D$ are the camera angles. When the neck is fixed as $\zeta = 0$, the binocular visual space is transformed into Cartesian coordinates by

$$
\begin{aligned}
x &= E \sin(2\theta)/\sin(\gamma) \\
y &= \bar{z} \sin(\delta) \ , \quad z = \bar{z} \cos(\delta) \\
\bar{z} &= E\{\cos(\gamma) + \cos(2\theta)\}/\sin(\gamma),
\end{aligned}
\tag{3}
$$

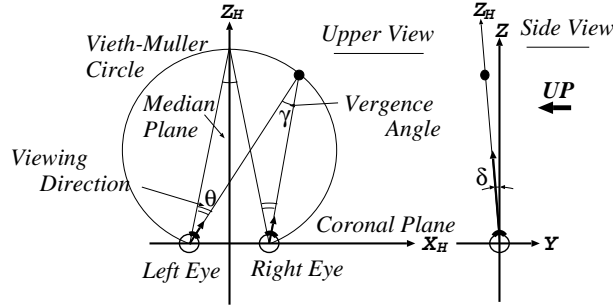where $E$ is half of the baseline length. The continuous lines in **Fig.4** depict the



**Fig.3**: Binocular Visual Space

binocular visual space which is projected onto a Cartesian subspace(X-Z plane). It is obvious that the transformation from the binocular visual space to Cartesian coordinates is a nonlinear mapping.

The binocular visual space has close relation to the camera image. The stereo camera geometry is shown in **Fig.5**. The coordinates of a feature point projected on the camera image planes are transformed into binocular visual coordinates by

$$
\boldsymbol{V} = \begin{bmatrix} \alpha_L - \alpha_R \\ (\alpha_L + \alpha_R)/2 \\ \alpha_D \end{bmatrix} + \begin{bmatrix} (X^L - X^R)/f \\ (X^L + X^R)/2f \\ (Y^L + Y^R)/2f \end{bmatrix}
\tag{4}
$$

where $(X^L, Y^L), (X^R, Y^R)$ are the coordinates of the image planes, and it is assumed that $\tan^{-1}(X^{L,R}/f) \simeq X^{L,R}/f$. Camera angles and image data are trans-
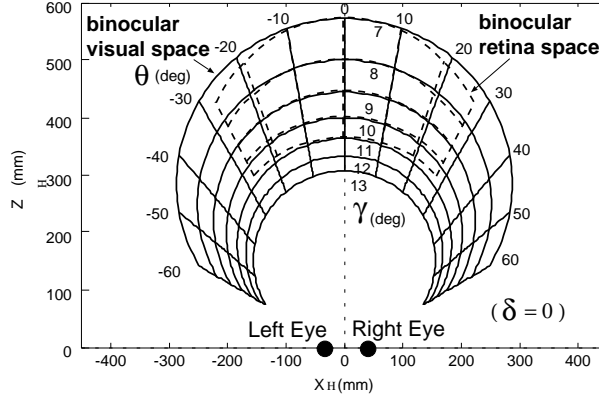
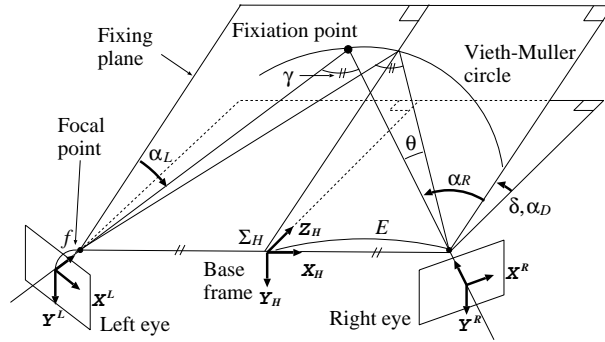**Fig.4**: Binocular Visual Space projected onto Cartesian Space



**Fig.5**: Geometry of stereo camera

formed into binocular visual coordinates by a linear time-invariant sum. The dotted lines in Fig.4 show the space aquired by eq.(4), when the fixiation point is $(\gamma, \theta, \delta) = (9, 0, 0)$. This approximation is available around the fixiation point. We call this space binocular retina space.

## 2.4 Linearity of the transformation between Binocular Visual Space and Joint Space

The continuous lines in **Fig.6** show the joint space which is projected onto a binocular visual space. This figure demonstrates the validiy of linear approximation of the transformation between the binocular visual space and the joint space. We linearize this transformation using the least-squares approximation within a region defined as $j_0 = 0, 20 \le j_1 \le 60, 60 \le j_2 \le 100$. Then the transformation can be represented as a linear equation of the form
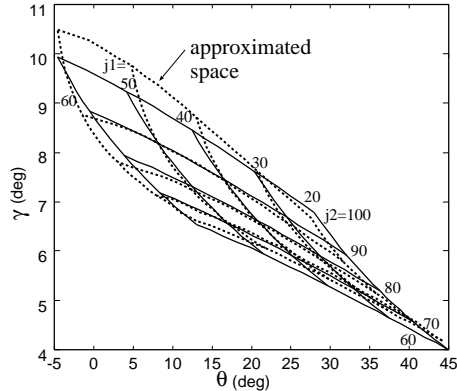
$$\boldsymbol{j} = \boldsymbol{R}\ \boldsymbol{V} + \boldsymbol{C} \tag{5}$$

**Fig.6**: Joint space projected onto Binocular Visual Space

where $\boldsymbol{V} = (\gamma, \theta, \delta)^T$, $\boldsymbol{j} = (j_0, j_1, j_2)^T$. When we restrict the arm motion in a level plane with the shoulder$(y = 0)$, the least-squares approximation results in

$$\boldsymbol{R} = \begin{pmatrix} 0 & 0 & 0 \\ -12.1 & -2.26 & 0 \\ 22.9 & 1.94 & 0 \end{pmatrix}, \boldsymbol{C} = \begin{bmatrix} 0 \\ 2.93 \\ -1.98 \end{bmatrix}. \tag{6}$$

The dotted lines in Fig.6 and Fig.2 show the approximated space by eq.5. This approximation has an accuracy enough to apply to our control scheme proposed in Sec.3. When the arm motion is not restricted in a level plane, $j_0$ only depends on $\delta$, $\gamma$ in the approximated transformation. And when we use the neck joint $\zeta$, it can be combined into an angle of viewing direction $\theta$, because both indicate a direction of a polar coordinate. We are currently investigating these subjects.

## 2.5 Suitable kinematic structure for the approximation

We'll explain that the best approximation of the mapping can be obtained by choosing similar kinematic parameters to a human being. We especially consider those parameters which determine the camera position, i.e.,$(W, K, G)$ given in Fig.1. In case of humans, the parameter values are $(W, K, G) = (-200, -200, 50)$[mm] on average. To optimize the parameters $(W, K, G)$, we evaluate the sum of squared defferences(SSD) over all given data points in the region where the linearization is carried out. **Fig.7** shows SSD–W characteristics under the condition that $K = -200$[mm] and $G = 0$[mm]. From this figure, we can see that the optimal value of $W$ is $-160$[mm]. **Fig.8** shows SSD–K and G characteristics under the condition that $W = -160$[mm]. This figure indicates that the optimal values of $K$ and $G$ are nearly $(K, G) = (-200, 0)$[mm]. In conclusion, $(W, K, G) = (-200, -200, 0)$[mm] are almost optimal in the sense of SSD and similar to those of a human being.

## 3 Linear Visual Servoing

We propose a simple visual servoing scheme based on the use of binocular visual space given by

$$\boldsymbol{u} = -\lambda \boldsymbol{R} (\boldsymbol{V} - \boldsymbol{V}_d) \tag{7}$$

where $\boldsymbol{u}$ are control signals to joint velocity controllers, $\boldsymbol{V}$ is the binocular visual coordinates of the hand, $\boldsymbol{V}_d$ is the binocular visual coordinates of a target and $\lambda$ is
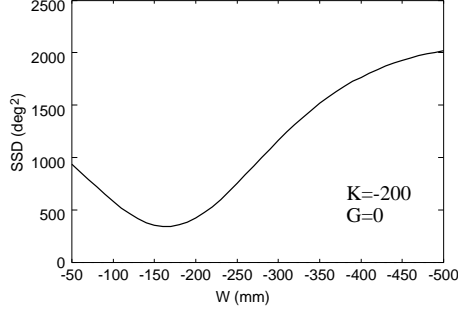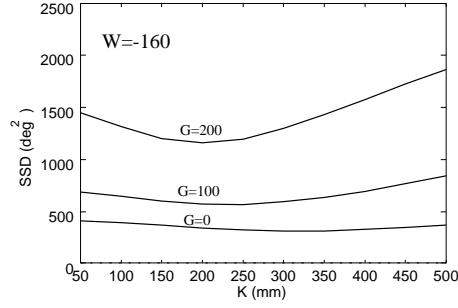
**Fig.7**: SSD-W characteristics



**Fig.8**: SSD-K and G characteristics

a scalar gain, $\boldsymbol{R}$ is the best approximation of the mapping obtained in the previous section. Using eq.(4), We get

$$
\begin{aligned}
\boldsymbol{u} &= -\lambda \boldsymbol{R} \left[ \begin{array}{c} \{(X^L - X^R) - (X_d^L - X_d^R)\}/f \\ \{(X^L + X^R) - (X_d^L + X_d^R)\}/2f \\ \{(Y^L + Y^R) - (Y_d^L + Y_d^R)\}/2f \end{array} \right] \\
&= -\lambda \ \boldsymbol{R} \ \boldsymbol{T} \ (\boldsymbol{I} - \boldsymbol{I}_d) \\
\boldsymbol{T} &= \left( \begin{array}{cccc} 1/f & -1/f & 0 & 0 \\ 1/2f & 1/2f & 0 & 0 \\ 0 & 0 & 1/2f & 1/2f \end{array} \right) \\
\boldsymbol{I} &= \left( X^L, X^R, Y^L, Y^R \right)^T .
\end{aligned} \tag{8}
$$

We call this simple control "*Linear Visual Servoing*". Features of this control scheme are summarized as follows.

- **The control law includes neither camera angles nor joint angles of the manipulator.**
  This system is robust to camera angle errors and joint angle errors. Measuring them is not required. Furthermore, camera angles have little influence on the system. It means that it is possible to turn cameras without considering the control of the manipulator. Hence it is especially suitable for the image processing system using active stereo vision.

- **The control law is very simple.**
  The amount of the calculation is small. Moreover, the estimation of the

trajectory is easy in spite of being an image-based visual servoing, because the trajectory is specified on a binocular visual space.

- **Flexible and global calibration is available.**
  The control scheme doesn't require the parameters of cameras and manipulator (focal length, camera position, length of the link....). The matrix $R$ can be obtained easily from the camera image and joint angles. Hence we can set the suitable parameters for the requisite work space flexibly and globally.

In this paper, we introduce a simple proportional feed-back control. Needless to say, we can apply other traditional control scheme in binocular visual space.

# 4  Effects of Map Approximation

We analyse the effects of a time-invariant linear mapping $\boldsymbol{R}$ in the proposed linear visual servoing scheme from the viewpoint of stability. We assume that the robot perfectly track joint velocity commands using an ideal velocity feedback control. Let $\boldsymbol{V}_d$ be a stationary target point. Define the error $\boldsymbol{e}_V = \boldsymbol{V} - \boldsymbol{V}_d$ and the error system

$$\dot{\boldsymbol{e}_V} = \boldsymbol{M}(\boldsymbol{j})\dot{\boldsymbol{j}} \tag{9}$$

where $\boldsymbol{M}(\boldsymbol{j})$ denotes the Jacobian from joint space to binocular visual space. In the linear visual servoing, joint velocities $\dot{\boldsymbol{j}}$ are given in the form of eq.(7). Then the closed loop system is given by

$$\dot{\boldsymbol{e}_V} = -\lambda \boldsymbol{M}\boldsymbol{R}\,\boldsymbol{e}_V \tag{10}$$

The problem for the moment is to make sure of the fact that a solution starting from an arbitrary point in the region where the linearlization is carried out (see Fig.2) converges to a target point set in the same region. To do so, we define a scalar function

$$\boldsymbol{U}(\boldsymbol{V}) = \frac{1}{2}\boldsymbol{e}_V^T \boldsymbol{e}_V \tag{11}$$

and differentiate it with respect to time along any state trajectory of system eq.(11), i.e.,

$$\frac{d}{dt}\boldsymbol{U}(\boldsymbol{V}) = -\lambda \boldsymbol{e}_V^T \boldsymbol{M}\boldsymbol{R}\,\boldsymbol{e}_V \tag{12}$$

Let $\Omega$ be a bounded region define by

$$\boldsymbol{U}(\boldsymbol{V}) < \boldsymbol{U}_0 \tag{13}$$

where $\boldsymbol{U}_0$ is a positive constant. If $\frac{d}{dt}\boldsymbol{U}(\boldsymbol{V}) \leq 0$ for all $\boldsymbol{V}$ in $\Omega$, we can conclude from the invariant set theorem [12] that any trajectory starting within $\Omega$ converges to the target point. The region $\Omega$ called "*a domain of attraction*" can be numerically determined. **Fig.9** and **Fig.10** show domains of attraction in the neighborhood of the region where the linearization is carried out around a target point(a bullet). To obtain these figures, we restricted the analysis in a level plane of the shoulder ($y = 200$) where the state vector $\boldsymbol{V}$ is two dimensional, i.e., $\boldsymbol{V} = (\gamma, \theta)$ and the linear mapping $\boldsymbol{R}$ in eq.(7) is a $2 \times 2$ matrix whose entries are given in eq.(6) (submatrix with non-zero entries). Several lines in both figures indicate boundaries corresponding to different values of $\boldsymbol{U}_0$. The dark shaded regions of both figures correspond to domains of attraction which were obtained by considering the manipulator's movable area and the slope of the tangent to the state trajectory on its boundaries. The light shade region is presumably included in the domain of attraction because the time derivative of $\boldsymbol{U}(\boldsymbol{V})$ is non-positive, though more careful analysis is required. As expected, the domain of attraction is large enough to include the region where the linearization is carried out.
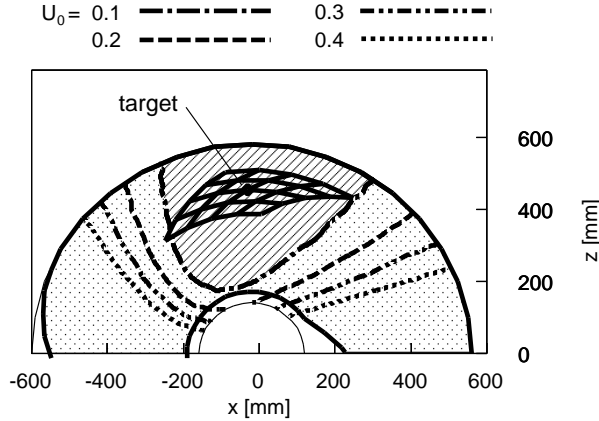
**Fig.9**: Domain of Attraction projected onto Cartesian Space

## 5 Experiments and Simulations

We show the availability of the linear visual servoing through experiments and simulations. To show the robustness to calibration error, we first compare the proposed scheme with a stereo visual servoing using pseudo-inverse matrix[4]. The control law is given by

$$\boldsymbol{u} \;\; = \;\; -\lambda \, \boldsymbol{J}_R^{-1} \, \boldsymbol{J}_I^+ \;\; (\boldsymbol{I} - \boldsymbol{I}_d) \,, \tag{14}$$

where $\boldsymbol{J}_R^{-1}$ is the inverse matrix of the robot Jacabian, and $\boldsymbol{J}_I^+$ is the pseudo-inverse matrix of the image Jacobian. Other parameters are the same as those of eq.(9). In the following, we describes L-VS as the linear visual servoing of eq.(9), and C-VS as the visual servoing using pseudo-inverse matrix of eq.(15). Then, other experimental results are given to illustrate that the active stereo cameras are available to the linear visual sevoing.

### 5.1 Experimental setup

**Fig.11** shows a schematic representation of the experimental setup. Our stereo vision system has a sampling rate of $30Hz$. It should be noted that system calibration was not performed in the experiment. We restricted the robot's motion in a horizontal plane($y = 0$) in the same way as the stability analysis mentioned in the previous section.

### 5.2 Properties of the trajectory

We compared the trajectories of L-VS and C-VS in two situations through experiments and simulations. The initial position and the target position in the first situation are $(x, y)^T = (-300, 200)^T$ and $(100, 300)^T$, respectively. In the second situation, they are exchanged. The trajectories of L-VS are not different between the simulation and the experiment, whereas the trajectories of C-VS are different. This implies the robustness of L-VS to calibration error.

Observing the motion of joint angles, we can see that L-VS's trajectories are almost straight, while C-VS's trajectories are tortuous. Because L-VS scheme is based on the joint space, whereas C-VS scheme is based on Cartesian coordinates. The above-mentioned properties of L-VS are useful to avoid un-reacheable religions.
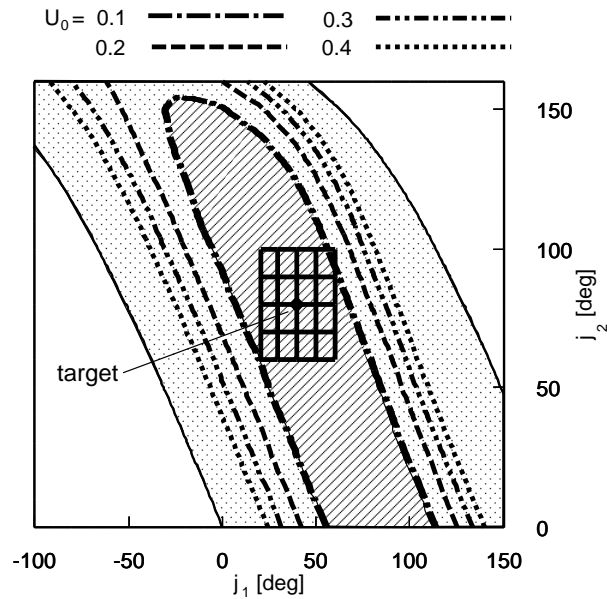
**Fig.10**: Domain of Attraction projected onto Joint Space

## 5.3 Independence of Camera Angles

Camera angles are not required by L-VS control scheme. It means that it is possible to turn cameras without considering the control of manipulator. To confirm this fact, we tried to turn cameras cyclically at the rate of about $1Hz$ by hand, and compared the trajectries between turning cameras and fixing cameras. **Camera angles were not measured, and compensator was not used in this experiment.**

**Fig.13** shows the trajectories of the hand acquired in the experiments. We can see that the influence by camera turning is very small. **Fig.14** shows the time histories of the feature points on the left image plane. The cyclical changes in the target trajectories represent the effect of camera turning. **Fig.15** shows the time histories of errors on the image planes between the hand and the target. It converges zero smoothly, and we can see that the disturbance by camera turning is very small.
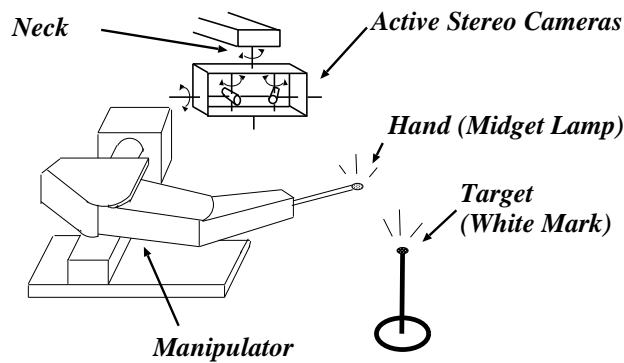


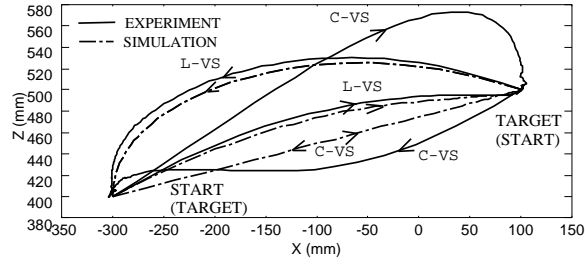**Fig.11**: Schematic representation of the experimental setup

**Fig.12**: Trajectories of the hand on Cartesian Space
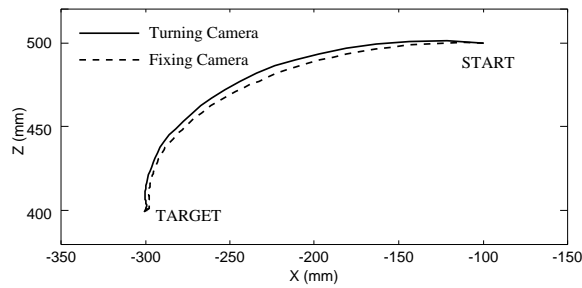


**Fig.13**: Trajectories of the hand

This disturbance is caused by the calibration error of optical model, and the error between binocular retina space and binocular visual space (sec.2.3). We found that the disturbance can be ignored as long as the target and hand is captured around the center of the image plane, although more analysis is required.

## 6 Conclusion

In this paper, we proposed a simple visual servoing scheme based on the use of binocular visual space. The robustness to calibration error, especially to the camera turning, was confirmed through some simulations and experiments. Conven-
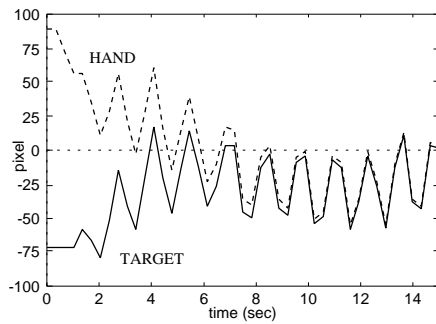


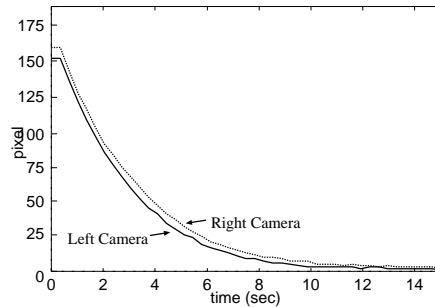**Fig.14**: $X$ coordinates of the feature points in left image

**Fig.15**: Parallaxes between the hand and the target

tional visual servoing schemes require that cameras are fixed or camera angles are measured accurately. The visual servoing without measuring camera angles will be useful to active vision systems.

We also showed that a similar kinematic structure as a human being can approximate the transformation from binocular visual space to joint space of the manipulator as a linear mapping.

This is an interesting design problem of finding appropriate kinematic structure for hand-eye system. Sharma et al.[13] describe an optimized hand-eye configuration as to motion perceptibility for visual servoing. On the other hand, our approach optimized hand-eye configuration as to linearity and the uniformity of hand-eye mapping, which would reduce dynamic effects such as time lag. This approach would be useful to achieve a high-performance no matter what control scheme is taken. We are currently investigating these subjects.

In this paper, the use of the neck joint is not described. We are investigating how to expand the field of cameras by adding the motion of the neck joint. Learning the coefficients of linear mapping is another important subject of our future works.

# References

[1] L. E. Weiss, A. C. Sanderson and C.P.Neuman, "Dynamic sensor-based control of robots with visual feedback," *IEEE Journal of Robotics and Automation*, RA-3(5), pp.404–417, 1987.

[2] W. J. Wilson, "Visual Servo Control of Robots Using Kalman filter Estimates of Relative POSE," *International Federation of Automatic Control*, Sydney, Australia, vol.9, pp.399 – 404, 1993.

[3] B. Espiau, F. Chaumette, and P. Rives, "A New Approach to Visual Servoing in Robotics," *IEEE Trans. Robotics and Automation*, vol.8, no.3, pp.313 – 326, 1992.

[4] N. Maru, H. Kase, S. Yamada, A. Nishikawa and F. Miyazaki, "Manipulator Control by Visual Servoing with the Stereo Vision," In *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol.3, pp.1865–1870, 1993.

[5] K. Hashimoto, T. Ebine and H. Kimura, "Dynamic Visual Feedback Control for a Hand-Eye Manipulator," *IEEE/RSJ Int. Conf. Intelligent Robots and Systems*, Raleigh,NC, pp.1863 – 1868, 1992.

[6] G. Hager, W. Chang, and A. Morse, "Robot Hand-Eye Coordination Based on Stereo Vision," *IEEE Control Systems magazine*, vol.15, no.1, pp.30–38, February 1995.

[7] K. Hosoda and M. Asada, "Versatile Visual Servoing without Knowledge of True Jacobian," *Proceeding of the IEEE International Conference on Intelligent Robots and Systems*, vol.1, pp.186–191, 1994.

[8] N. Hollinghurst and R. Cipolla, "Uncalibrated stereo hand-eye coordination," *Image and Vision Computing*, vol.12, no.3, pp.187–192, April 1994.

[9] R. K. Luneburg, "The Metric of Binocular Visual Space," *Journal of the Optical Society of America, Information and Communication Engineers*, vol.40, no.10, pp.627–642, 1950.

[10] R. Caminiti, P. B. Johnson and A. Urbano, "Making Arm Movements Within Different Parts of Space: Dynamic Aspects in the Primate Motor Cortex," *The journal of Neuroscience*, vol.10(7), pp.2039–2058, 1990.

[11] M. S. A. Graziano and C. G. Gross, "A bitmodal map of space: somatosensory receptive fields in the macaque putamen with corresponding visual receptive fields," *Experimental Brain Research*, vol.97, pp.96–109, 1993.

[12] J.-J. E. Slotine and W. Li, Applied nonlinear control. Englewood Cliffs, N.Z: *Prentice Hall*, 1991.

[13] R. Sharma and S. Hutchinson, "Optimizing Hand/Eye Configuration for Visual-Servo Systems," *Proceeding of the IEEE International Conference on Robotics and Automation*, vol.1, pp.172–177, 1995.