# EVOLUTIONARY POSE MEASUREMENT BY STEREO MODEL MATCHING

Wei SONG,[1] Takahiro NATSUME[2], and Mamoru MINAMI[3]

[1] Graduate School of Engineering, Fukui University[2] Toyota Macs Inc.
[3] Dept. Human and Artificial Intelligent System, Fukui University
3-9-1, Bunkyo, Fukui-shi, 910-0017, Japan

## Abstract

This paper presents a pose measurement method of a moving 3-D object. The proposed method utilizes the genetic algorithm (GA) and unprocessed gray-scale image input from stereo-vision, in order to perform recognition of a target being imaged with known target object shape. In fact here, the problem to recognize the target shape and simultaneous detection of the position/orientation, is converted to an optimum problem of a model-based evaluation function, named as surface-strips model-based fitness function that consists in the computation of the brightness difference between an internal surface and a contour-strips. In order to evaluate the proposed 3-D recognition method, experiments to detect position/orientation of a rectangular solid block have been conducted to show its effectiveness of recognizing objects in static image. Furthermore, experiments to recognize a ball on a turning table by a robot manipulator equipped with two hand-eye cameras have also been conducted to show the effectiveness of this method for real-time visual servoing.

## 1 Introduction

In recent years, recognition of an object target in an image has been researched intensively. The object recognition includes a lot of research subjects such as, finding object, recognition of shape, measurement of position/orientation in varieties of environments, and so on. Furthermore, the object recognition for dynamic scene being input by video rate (33[ms]) is important practically to use the recognized results for visual feedback, which requires the system to be processed in real time. For example, object recognition method for dynamic scene using parallel computing hardware, performs recognizing process in a period of 1 [ms] [1]. This method is thought that it is really processing in real time, since the time delay of the detected information is less than 1 [ms] as long as the target being imaged is continuously recognized in the dynamic images. However, such hardware computing system may not comparatively be portable than recognizing system using just software. Contrarily to the above direction, we have proposed a recognition method using soft computing such as genetic algorithm (GA) to recognize a target in dynamic image being input by video rate. Using the proposed concept which we called "1-step GA", we realized the ability

to track a target object by video rate [2], [4]. Furthermore, we have confirmed that this method enabled a hand-eye robot system to catch a swimming fish by a net equipped at hand. However, the cognizable information was limited to 2-D position/orientation, since the system used a single camera. To extend the visual recognition information from previous 2-D to 3-D, we propose here how to extend the recognition space to 3-D while keeping the real time processing ability and utilizing soft computing.

The research direction of a 3-D measurement using images obtained from stereo-vision is divided into main two methods. The one is to measure 3-D information by using geometric characteristics in left and right images, which are corresponding to one characteristic of actual object. In this case, the 3-D measurement is performed by calculating the geometric relation that is based on a principle of triangle surveying, through finding out corresponding points in right and left images. So far many methods being included in this concept were proposed, e.g., [5]. The other main method is to use a model to search a target object in the image whose model is composed based on how the target object can be seen in the input image [6], [7]. An advantage of this method using model-based matching is that it recognizes an object without searching the corresponding points in stereo-vision camera images. Our method is included in this category.

## 2 Evolutionary measurement

In this section, we take 2-D position/orientation measurement as an example to introduce an evolutionary measurement method, in which a GA is used directly to a 2-D gray-scale input image termed here as raw-image.

A searching model, which is composed of an internal surface and contour-strips, is employed for recognition purposes of a target in the raw-image, and such model is designated as surface-strips model. The internal surface approximated the most of the varieties of the 2-D top surface of the target. Consider the 2-D raw-image of a target fish shown in Fig.1(a), and corresponding 3-D plot is shown in Fig.1(b). In this figure, the vertical axis represents the image brightness values, and the horizontal axes, the image plane. To search for such a target fish in the raw-image, a geometrical triangular model of the surface-strips model as shown in Fig.2(a) is used. Let us denote the inside surface of the model as $S_{ss1}$ and
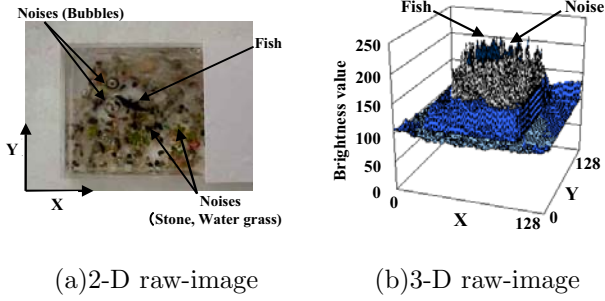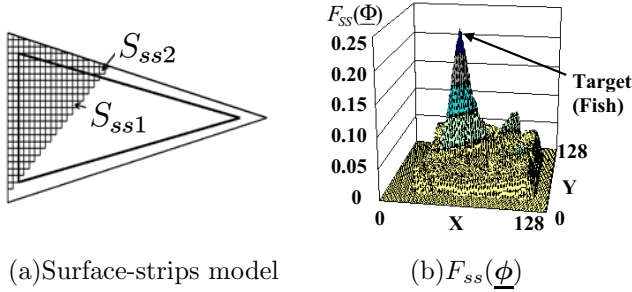
(a)2-D raw-image  (b)3-D raw-image

Figure. 1: Raw-image of swimming fish



(a)Surface-strips model  (b)$F_{ss}(\underline{\phi})$

Figure. 2: Surface-strips Model to search a fish

$$position:\underline{x} \quad position:\underline{y} \quad orientation:\underline{\theta}$$
$$\underbrace{0011101} \quad \underbrace{0110010} \quad \underbrace{100110110}$$
$$an \; individual's \; position/orientation$$

Figure. 3: Individual binary string in GA

the contour-strips as $S_{ss2}$, and also, the combination is designated as $S_{ss}$. Here, $S_{ss1}$ is determined by the size of the target fish, in order to give the following fitness function the highest value. When the position and orientation of surface-strips model $S_{ss}$ is defined as a variable of $\underline{\phi} = [\underline{x}, \underline{y}, \underline{\theta}]^T$, which designates the position and orientation of the origin of the model, then $S_{ss}$ moves in the camera frame and a set of x-y coordinates of the moving model is expressed as $S_{ss}(\underline{\phi})$. The brightness distribution of input raw-image corresponding to the area of the moving model $S_{ss}$ is expressed as $p(\tilde{\boldsymbol{r}}_{x,y}), \tilde{\boldsymbol{r}}_{x,y} \in S_{ss}(\underline{\phi})$, then the evaluation function $F_{ss}(\underline{\phi})$ of the moving surface-strips model is defined here as follows in Eq.(1).

$$F_{ss}(\underline{\phi}) = \sum_{\tilde{\boldsymbol{r}}_{x,y} \in S_{ss1}(\underline{\phi})} p(\tilde{\boldsymbol{r}}_{x,y}) - \sum_{\tilde{\boldsymbol{r}}_{x,y} \in S_{ss2}(\underline{\phi})} p(\tilde{\boldsymbol{r}}_{x,y}) \quad (1)$$

This expression means the integrated brightness difference between the brightness of the internal surface and the one of the contour-strips of the surface-strips model. The filtering result of the surface-strips model-based fitness function, with respect to Fig.1(a) is shown in Fig.2(b). The filtering result with the surface-strips model-based fitness function has a peak corresponding to the target fish in the raw-image. This evaluation using the surface-strips model means that $F_{ss}(\underline{\phi})$ takes into account the differentiation between an object signal and the background. We can make efforts to set such an environment that the highest

value of the $F_{ss}(\underline{\phi})$ is obtained only if $S_{ss1}$ fits to the target object being imaged. Then the problem of recognition of a fish and detection of its position/orientation is converted to a searching problem of $\underline{\phi}$ such that maximizes $F_{ss}(\underline{\phi})$, which is used as a fitness function in GA process.

The positin and orientation of each searching model is characterized as a binary string, shown in Fig.3 which describes an individual in GA. The searching by GA for the solution is performed through an evolution process, from generation to generation. To recognize a target in a dynamical changing image, the recognition system must have been processed in real-time, that is, the searching model must keep the convergence to a fish in the successively input raw-images. We have proposed a new idea of an evolutionary recognition process for dynamic image, in which the GA evolving process is applied only one time to the newly input raw-image. Therefore every input image is evaluated only one time, we named it as "1-step GA", and its flow chart is shown in Fig.4. If the evolutionary calculation is executed in less than 33[ms] and the model keeps to lie on the moving target in the successively input images, it can be said real-time recognition.

The effectiveness of the proposed 2-D recognition method is confirmed by the experiment of catching a fish by visual servoing of a robot with a camera and a net at the hand [4], as shown in Fig.5. In addition, this real time tracking performance can be realized even though the time used for making the model converge to the target be bigger than 33[ms] in step response of recognition experiment[3].

# 3  Extending to 3-D

## 3.1  Kinematics of Stereo-vision

We utilize perspective projection as projection transformation. The coordinate systems of left and right cameras and object in Fig.6 consist of world coordinate system as $\Sigma_W$, model coordinate system as $\Sigma_M$, camera coordinate systems as $\Sigma_{CR}$ and $\Sigma_{CL}$, image coordinate systems as $\Sigma_{IR}$ and $\Sigma_{IL}$. A point $i$ on the target can be described using these coordinates and homogeneous transformation matrices. At first, a homogeneous transformation matrix from $\Sigma_{CR}$ to $\Sigma_M$ is defined as $^{CR}\boldsymbol{T}_M$. And an arbitrary point $i$ on the target object in $\Sigma_{CR}$ and $\Sigma_M$ is defined $^{CR}\boldsymbol{r}_i$ and $^M\boldsymbol{r}_i$. Then $^{CR}\boldsymbol{r}_i$ is,

$$^{CR}\boldsymbol{r}_i = {}^{CR}\boldsymbol{T}_M \, {}^M\boldsymbol{r}_i. \quad (2)$$

Using a homogeneous transformation matrix from $\Sigma_W$ to $\Sigma_{CR}$, i.e., $^W\boldsymbol{T}_{CR}$, then $^W\boldsymbol{r}_i$ is got as,

$$^W\boldsymbol{r}_i = {}^W\boldsymbol{T}_{CR} \, {}^{CR}\boldsymbol{r}_i. \quad (3)$$

The position vector of $i$ point in right image coordinates, $^{IR}\boldsymbol{r}_i$ is described by using projection matrix $\boldsymbol{P}$ of camera as,

$$^{IR}\boldsymbol{r}_i = \boldsymbol{P} \, {}^{CR}\boldsymbol{r}_i. \quad (4)$$

Figure. 4: Flow chart of 1-step GA recognition



Figure. 5: Fish catching by hand-eye visual servo



Figure. 6: Coordinate systems



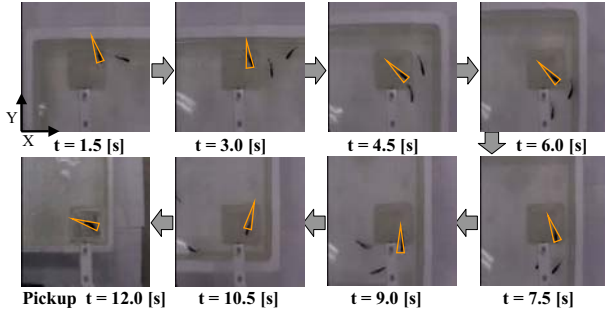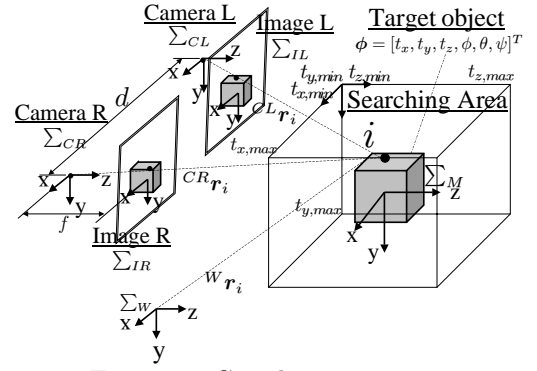Figure. 7: Solid model searching for a block



(a) Left searching model    (b)Right searching model

Figure. 8: Searching model

By the same way as above, using a homogeneous transformation matrix of fixed values defining the kinematical relation from $\Sigma_{CL}$ to $\Sigma_{CR}$, $^{CL}\boldsymbol{T}_{CR}$, $^{CL}\boldsymbol{r}_i$ is,

$$^{CL}\boldsymbol{r}_i = {}^{CL}\boldsymbol{T}_{CR} {}^{CR}\boldsymbol{r}_i. \tag{5}$$

As we have obtained $^{IR}\boldsymbol{r}_i$, $^{IL}\boldsymbol{r}_i$ is described by the following Eq.(6) through projection matrix $\boldsymbol{P}$.

$$^{IL}\boldsymbol{r}_i = \boldsymbol{P} \, {}^{CL}\boldsymbol{r}_i \tag{6}$$

Then position vectors projected in the $\Sigma_{IR}$ and $\Sigma_{IL}$ of arbitrary point $i$ on target object can be described $^{IR}\boldsymbol{r}_i$ and $^{IL}\boldsymbol{r}_i$. Here, position and orientation of the origin of $\Sigma_M$ based on $\Sigma_{CR}$ are represented as $\underline{\boldsymbol{\phi}} = [\underline{t_x}, \underline{t_y}, \underline{t_z}, \underline{\phi}, \underline{\theta}, \underline{\psi}]^T$, in which $\underline{\phi}, \underline{\theta}, \underline{\psi}$ are Euler angles, and then Eq.(4), (6) are rewritten as,

$$\begin{cases} ^{IR}\boldsymbol{r}_i = f_R(\underline{\boldsymbol{\phi}}, \, ^M\boldsymbol{r}_i) \\ ^{IL}\boldsymbol{r}_i = f_L(\underline{\boldsymbol{\phi}}, \, ^M\boldsymbol{r}_i). \end{cases} \tag{7}$$

This relation connects the arbitrary points on the object and projected points on the left and right images with the variables $\underline{\bold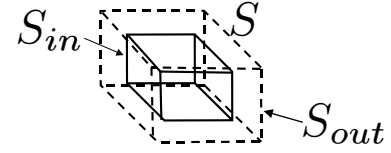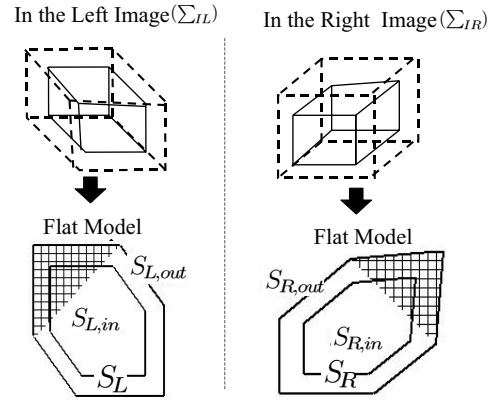symbol{\phi}}$, which is considered to be unknown in this paper. When evaluating each point above, the matching problem of corresponding point in left and right images mentioned in the introduction is arisen. Therefore, to avoid this problem, the 3-D model-based matching that treats the image as a set, is chosen instead of point-based corresponding. The 3-D model for the target object of a rectangular block is shown in Fig.7, in which the set of coordinates inside of the block is depicted as $S_{in}$ and the outside space enveloping $S_{in}$ is denoted as $S_{out}$, and the combination is named as $S$. Then, the set of the points of solid searching model consisted of $S_{in}$ and $S_{out}$, which is projected onto the two dimensional coordinates of left camera are expressed as,

$$\begin{cases} S_{L,in}(\underline{\boldsymbol{\phi}}) = \{^{IL}\boldsymbol{r}_i \in \Re^2 \mid {}^{IL}\boldsymbol{r}_i = f_L(\underline{\boldsymbol{\phi}}, ^M\boldsymbol{r}_i), \\ \qquad\qquad\qquad\qquad ^M\boldsymbol{r}_i \in S_{in} \in \Re^3 \} \\ S_{L,out}(\underline{\boldsymbol{\phi}}) = \{^{IL}\boldsymbol{r}_i \in \Re^2 \mid {}^{IL}\boldsymbol{r}_i = f_L(\underline{\boldsymbol{\phi}}, ^M\boldsymbol{r}_i), \\ \qquad\qquad\qquad\qquad ^M\boldsymbol{r}_i \in S_{out} \in \Re^3 \} \end{cases} \tag{8}$$

(a)Left image          (b)Left birightness

(c)Right image         (d)Right brightness

Figure. 9: Input image and brightness distribution



Figure. 10: $F_{ss}(\phi)$ of searching area

where the left searching model projected to left camera coordinates is shown in Fig.8(a). The area composed of $S_{L,in}$ and $S_{L,out}$ is named as $S_L$. The above defines only the left-image searching model, the right one is defined in the same way and the projected searching model is shown in Fig.8(b).

## 3.2  3-D measurement method

The size of the input images from stereo cameras is $320 \times 240[pixel]$, and both are shown in Fig.9(a) and (c). These images are constructed by gray scale brightness values, and their distributions are shown in Fig.9(b) and (d). In this research, the input images, i.e., unprocessed gray scale images, are directly matched by the projected moving models, $S_L$ and $S_R$, which are located by only $\underline{\phi} = [\underline{t_x}, \underline{t_y}, \underline{t_z}, \underline{\phi}, \underline{\theta}, \underline{\psi}]^T$ as described in Eq.(8) that includes the kinematical relations of the left and right camera coordinates. Therefore, if the camera parameters and kinematical relations are completely accurate, and the solid searching model describes precisely the target object shape, then the $S_{L,in}$ and $S_{R,in}$ will be completely lies on the target reflected on the left and right images, provided that true values of $\underline{\phi}$ is given.

Here, we define evaluation function to estimate how match the moving solid model $S_{in}$ defined by $\underline{\phi}$ lies on the target being imaged on the left and right cameras. In order to search for the target object in the gray scale image, the surface-strips model shown in Fig.8 and its position calculated by Eq.(7) are used. The inside surface of the model in the left and right cameras are $S_{L,in}$ and $S_{R,in}$ and the contour-strips as $S_{L,out}$ and $S_{R,out}$. The brightness distribution of input image lying on the area of searching model is expressed as $p(^{IL}\boldsymbol{r}_i), \boldsymbol{r}_i \in S_L(\underline{\phi})$, and $p(^{IR}\boldsymbol{r}_i), \boldsymbol{r}_i \in S_R(\underline{\phi})$, then the evaluation function $F(\underline{\phi})$
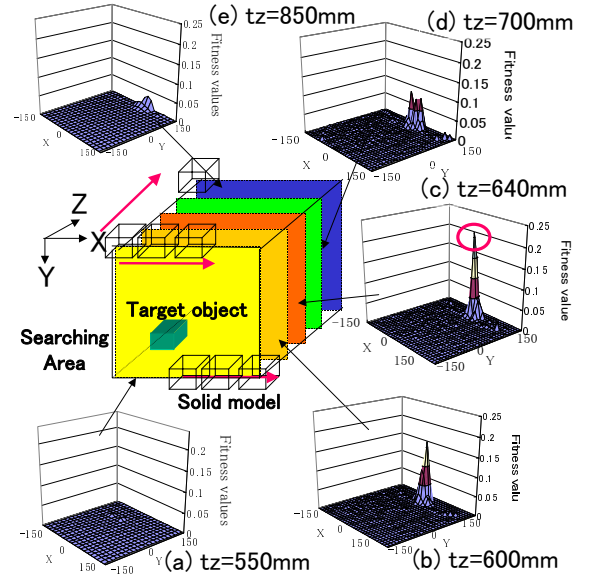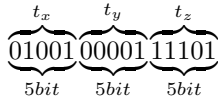
of the moving surface-strips model is given as,

$$
\begin{aligned}
F(\underline{\phi}) &= \left( \sum_{^{IR}\boldsymbol{r}_i \in S_{R,in}(\underline{\phi})} p(^{IR}\boldsymbol{r}_i) - \sum_{^{IR}\boldsymbol{r}_i \in S_{R,out}(\underline{\phi})} p(^{IR}\boldsymbol{r}_i) \right) \\
&\quad \cdot \left( \sum_{^{IL}\boldsymbol{r}_i \in S_{L,in}(\underline{\phi})} p(^{IL}\boldsymbol{r}_i) - \sum_{^{IL}\boldsymbol{r}_i \in S_{L,out}(\underline{\phi})} p(^{IL}\boldsymbol{r}_i) \right) \\
&= F_R(\underline{\phi}) \cdot F_L(\underline{\phi}), \qquad (9)
\end{aligned}
$$

where, in case of $F(\underline{\phi}) \leq 0$, $F(\underline{\phi})$ is given to zero. This function expresses the brightness difference between the one of the internal surface and the one of the contour-strips of the surface-strips model, and used as a fitness function in GA process. When the moving searching model fits to the target object being imaged in the right and left images, then the fitness function $F(\underline{\phi})$ gives maximum value. We confirm this as follows.
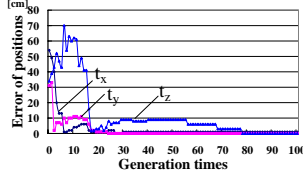
Here, a 3-D plot of fitness distribution that is calculated by Eq.(9) and by scanning the position of the 3-D space of the moving model with fixed orientation coinciding with the true values of the target object, is shown in Fig.10. A position of target object that we measured beforehand is $(\underline{t_x}, \underline{t_y}, \underline{t_z}) = (0, 75, 640)[mm]$. Consider that the position $\underline{t_z}$ is fixed to several values as shown in Fig.10(a) $\sim$ (e), then the fitness distributions on $(\underline{t_x}, \underline{t_y})$ plane are obtained as shown in (a) $\sim$ (e). It can be seen that when the position of the model fits to the true poison $(\underline{t_x}, \underline{t_y}, \underline{t_z}) = (0, 75, 640)[mm]$, the fitness function has the maximum value as shown in Fig.10(c). Therefore the problem of recognition of target object and detection of its position/orientation can be converted to searching problem of $\underline{\phi}$ such that maximizes $F(\underline{\phi})$. To recognize the
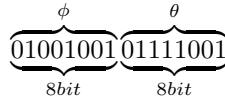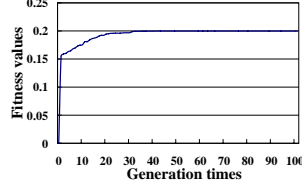
$$\underbrace{01001}_{5bit}\underbrace{00001}_{5bit}\underbrace{11101}_{5bit}$$
$$t_x \quad t_y \quad t_z$$

$$\underbrace{01001001}_{8bit}\underbrace{01111001}_{8bit}$$
$$\phi \quad \theta$$

(a)Individual binary string in GA

(a)Individual binary string in GA



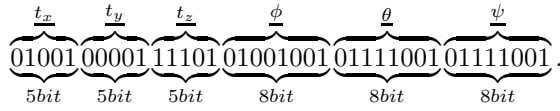(b)Fitness value

(b)Fitness value



(c)Error of position

(c)Error of Orientation

Figure. 11: Position measurement experiment

Figure. 12: Orientation measurement experiment

target object in short time, we solve this optimization problem to search for $\underline{\phi}$ to maximize $F(\underline{\phi})$ by GA whose gene representing $\underline{\phi}$ is defined as,

$$\underbrace{01001}_{5bit}\underbrace{00001}_{5bit}\underbrace{11101}_{5bit}\underbrace{01001001}_{8bit}\underbrace{01111001}_{8bit}\underbrace{01111001}_{8bit}.$$
$$t_x \quad t_y \quad t_z \quad \phi \quad \theta \quad \psi$$

The 39 bits of gene refers to the range of the searching area: $-150 \leq t_x, t_y \leq 150 \quad 550 \leq t_z \leq 850[mm]$, $-90 \leq \phi, \theta, \psi \leq 90[deg]$

# 4 Experiments

## 4.1 Separate Recognition of Block Position and Orientation

Experiments to detect objects position/orientation of the target using the method proposed in static image have been conducted. In these experiments, we used a rectangular solid block $(25 \times 40 \times 80[mm])$ as target object and right and left images shown in Fig.9(a), (c). To appraise the effectiveness, we have performed experiment with the following conditions.

**(1)** Position detection with fixed orientation
The searching space of $(t_x, t_y, t_z)$ is $-150 \leq t_x, t_y \leq 150, 550 \leq t_z \leq 850[mm]$. The orientation is given the true values as $(\phi, \theta, \psi) = (0, 0, 0)$.

**(2)** Orientation detection with fixed position
The searching space of $(\phi, \theta, \psi)$ is $-90 \leq \phi, \theta \leq$

$+90, \psi = 0[deg]$. The position is given the true values as $(t_x, t_y, t_z) = (0, 75, 640)$.

The parameters of GA are, individuals: 30, selection rate: 50%, and mutation rate: 10%. Figures 11(a) and 12(a) give the information of gene in GA process. Figures 11(b) and 12(b) show relations between generation times and maximum fitness values, and Figs.11(c) and 12(c) show relations between generation times and error of position and orientation. These experiments have been repeated for 10 times, and the results in Figs.11 and 12 are the average value of 10 results. Figure 11, 12 show that fitness function converged to the maximum value and the error of searching results of GA decreased gradually. And further, the position detection error is almost $\pm 1[mm]$ after 70 generations, and the orientation detection error is $\pm 3[deg]$ after 30 generations. The one generation in GA process required 300[ms] by computer using Celleron700MHz CPU.

## 4.2 Simultaneous Recognition of Block Position and Orientation

Here, we tried to recognize the Position and Orientation of the solid block (the same target used in the separate recognition experiment) simultaneously. That means searching for all the components of $\underline{\phi}$ (shown at the end of part 3.2)at a time. During the experiment, we keep the target position and orientation as $(\underline{t_x}, \underline{t_y}, \underline{t_z}, \underline{\phi}, \underline{\theta}, \underline{\psi}) = (0, 110, 640, 0, 60, 0)$. Figure 13(a) shows relations between generation times and maximum fitness values, and Figs.13(b) and 13(c) show relations between generation times and error of position and orientation. The detection error of position $t_z$ is almost 30[mm], and that of $\underline{\phi} \, \underline{\psi}$ remains 30[deg] which is quite an obvious contrast to what we receive from the position/orientation separate recognition experiment(in part 4.1). The one generation in GA process required 150[ms] by computer using Pentium4 2.53GHz CPU, which means it cost 150[ms]×300[generations]=45[s] to receive the result of one GA process. So it is necessary to improve the accuracy and speediness in order to recognize the position and orientation simultaneously in real time.

## 4.3 Visual Servoing of ball on Turning table

Furthermore, we have performed a Visual Servoing experiment of a ball on a turning table and the photograph of our experimental system is shown in Fig.14. Here, we keep the hand position from the camera center to the ball center as 600mm. Figure.15 gives the true position of the ball in $\Sigma_W$ as the table turning, and Fig.16 shows the experimental results of the hand position in $\Sigma_W$ under the speed of the turning table was $\pi/60$ rad/sec. It can be seen that the hand position of Y,Z directions in $\Sigma_W$ could be obtained accurately, but in X direction, hand position is not satisfactory. When the speed of the turning table was increasing, it became more difficult for the robot to

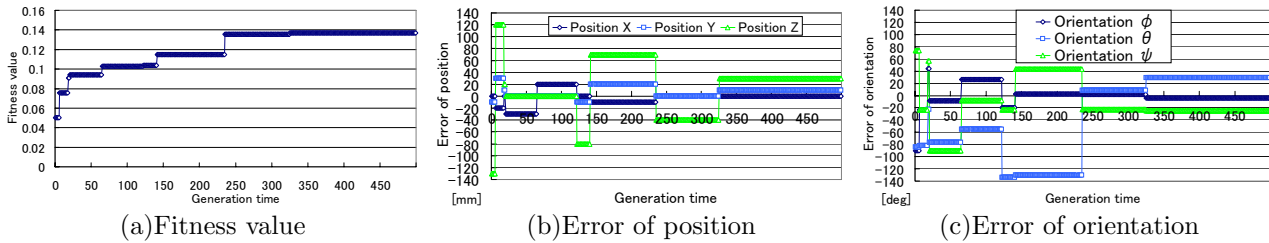(a)Fitness value　　　　　　　(b)Error of position　　　　　　(c)Error of orientation

Figure. 13: GA Searching results of Block (6 degree of freedom, 150[ms/generation](Pentium4 2.53GHz))
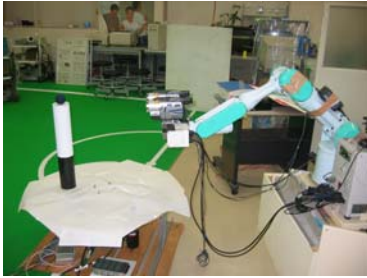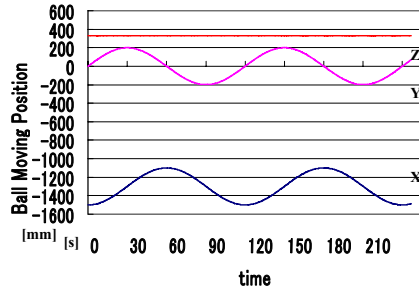


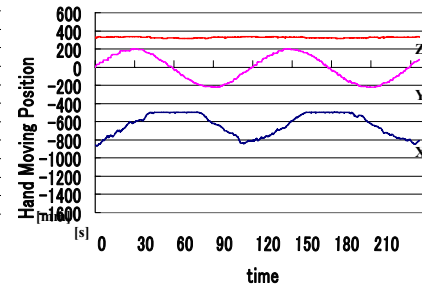Figure. 14: Visual servoing system　　　Figure. 15: Ball moving position　　　Figure. 16: Hand moving position

track the ball because the position recognition in X direction (in $\Sigma_W$) is not so fast as the speed of the turning table.

## 5　Conclusion

We proposed 3-D position/orientation measurement methods which utilizes a genetic algorithm and the unprocessed gray-scale image from the stereo-vision. And experiments to detect position/orientation of a rectangular solid block have been conducted to evaluate the recognition in static image. Furthermore, the experiments of real-time visual servoing by manipulator to a target ball (three degree of freedom) have been performed to appraise the performances of the proposed method. These experimental results have shown the effectiveness of the proposed method. As future work, the experiment "Visual servo to block" will be conduct by reduction of process time to verify the effectiveness of this 3-D position/orientation (six degree of freedom) measurement in real time recognition.

## References

[1] Idaku Ishii, Masatoshi Ishikawa, "High Speed Target Tracking Algorithm for 1ms Visual Feedback System", JRSJ, Vol.17, No.2, pp.39-45, 1999 space-3mm

[2] Minami,M., Agbanhan,J., Asakura,T., "GA-Pattern matching based Manipulator Control System for Real Time Visual Servoing", Advanced Robotics, Vol.12, No.78, pp.711-734, 1999. space-3mm

[3] M.Minami, H.Suzuki, J.Agbanhan, "Fish Catching by Robot Using Gazing GA Visual Servoing ", JSME, Vol.68, No.668, pp.1198-1206, 2002.

[4] M.Minami, H.Suzuki, J.Agbanhan, T.Asakura:"Visual Servoing to Fish and Catching Using Global/Local GA Search", 2001 IEEE/ASME Int. Conf. on Advanced Intelligent Mechatronics Proc., pp.183-188, 2001. space-3mm

[5] Yukie MAEDA, Gang XU, "Smooth Matching of Feature and Recovery of Epipolar Equation by Tabu Search", IEICE, Vol.J83-D-2, No.3, pp.440-448, 1999.

[6] Sadaaki Yamane, Masao Izumi, Kunio Fukunaga, "A Method of Model-Based Pose Estimation", IEICE, Vol.J79-D-2, No.2, pp.165-173, Feb, 1996. space-3mm

[7] Fubito Toyama, Kenji Shoji, Juichi Miyamichi, "Pose Estimation from a Line Drawing Using Genetic Algorithm", IEICE, Vol.J81-D-2, No.7, pp.1584-1590, July, 1998.