# Analyses of Optimization Performance in Longitudinal Pose Tracking with Eye-Vergence System

Yejun Kou[1], Hongzhi Tian[1], Mamoru Minami[1], and Takayuki Matsuno[1]

[1]Okayama University, Japan
(Tel: 81-86-251-8233, Fax: 81-86-251-8233)
[1]ptlg9dvi@s.okayama-u.ac.jp

**Abstract:** In this research, Genetic algorithm(GA) is used a pose-tracking method, which is called"Real-Time Multi-Step GA"(RT-MS GA), to solve on-line optimization problems for 3-D visual servoing. Real-time object tracking has been shortened for on-line pose estimation by using RT-MS GA and Eye-Vergence System, and the pose-tracking accuracy has been verified through fitness function distribution which is a correlation function between the target object projected in camera frame and the model defined in control computer. In this research, the performance of the RT-Ms GA was confirmed, and the error in recognition was be verified.

**Keywords:** Visual servoing, Eye-vergence, Real-Time Multi-Step GA, Longitudinal

## 1 INTRODUCTION

Visual Servoing is a control method to control the motion of the robot. By incorporating visual information obtained from visual sensor [1]-[4]with the feed-back loop, visual servoing is expected to be able to allow the robot adapt the changing or unknown environment. Some methods have been proposed already to improve the observation abilities of the robot, for instance by using stereo cameras [5], multiple cameras [6], and two cameras; with one fixed on the end-effector, and the other one fixed in the workspace [7]. However, these methods obtain different views to observe an object by increasing the number of cameras, leaving the system less adaptive for changing environment.

The problem to find position/orientation , i. e, pose of an object relatively based on hand-eye camera flame can be transposed to the optimization problem of correlation function. In this research, we use Genetic Algorithm(GA) to get the maximum correlation value within less than video rate, "Real-Time Multi-Step GA"(RT-MS GA) algorithm that is on-line estimation method [8].

The pose of an actual object are expressed with the pose of the coordinate system $\Sigma_M$ being fixed to the object, and the object's pose estimated by RT-MS GA is expressed as $\Sigma_{\hat{M}}$. It is natural there should always exist an error, then ${}^M\boldsymbol{T}_{\hat{M}}$ $4 \times 4$ homogeneous Matrix representing the pose relation between the actual object $\Sigma_M$ and the detected one $\Sigma_{\hat{M}}$ is usually not identity matrix. It is necessary to decrease this estimation error ${}^M\boldsymbol{T}_{\hat{M}}$. Then, we compared the pose of the object calculated by full search of fitness value with the pose estimated by 1-step GA. The fitness function means correlation between projected target object to left and right camera image and 3-D target model predefined in the computer.

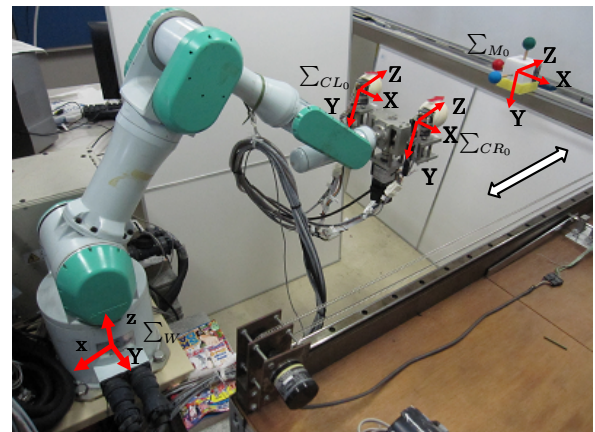By a previous work, it has been shown that the tracking



Fig. 1. Object and the visual-servoing system

performance could be improved by using eye-vergence system. However, the argument is not based on what kind of preferable influence the "RT-MS GA" can affect on-line optimization process during recognition. In this research, visual servoing experiment was performed on the point of frequency response as shown in Fig. 1, where target object represented by $\Sigma_M$ was oscillated by sinusoidal function and end-effector and eye-vergence systems are both controlled to keep desired pose relations that it constant desired pose relations.By comparing the results from full-search, we clarify how the dynamical performance can be improved by "RT-MS GA" in tracking the moving target object.

## 2 3D POSE TRACKING METHOD

Visual servoing system used in this research needs real-time estimation. Dynamic image given by video camera constitutes of still pictures input successively in order by a time series. Tracking an object in video is realized by estimating

the object in a still picture continuously within a video rate (33 [ms]). Therefore, this section explains the outline of the estimation technique for the still picture of one sheet.

The pose of 3-D model $\phi = (x, y, z, \epsilon_1, \epsilon_2, \epsilon_3)$ ($\epsilon$ is orientation variable of quaternion) is determined by the gene of GA. The 2-D models are obtained by projecting the 3-D model onto either side of left and right camera flame. This 2-D model is evaluated by calculating a fitness function defined by correlation between the 2-D model and the input image. And when the pose $\phi$ of the solid model coincides with the pose of the object, the correlation function has been designed to have maximum value. Therefore, the pose estimation problem of the object is turned into optimization problem to find $\phi$ maximizing the correlation function, i.e., fitness function [8]. There are various methods to solve optimization problem that the maximum of a given distribution is searched for and discovered. The most easy and understandable method is full search method. Although full search method can discover the maximum by calculating all the values of the function and can certainly discover the maximum value, this method spends too long time. It is important for real-time pose tracking to complete calculation time in a very short–within video rate, 33[ms]–in this research, GA has been applied for the pose tracking in a form of "RT-MS GA" [8].

## 2.1 Model-based Matching Method

In this section, a model-based matching method was presented. The images input from a right-and-left video cameras are composed by hue value ranging from 0 to 255. The searching model is shown as Fig. 2. The model constituted inside spaces $S_{R,in}(\phi)$ and $S_{L,in}(\phi)$, and outsize spaces $S_{R,out}(\phi)$ and $S_{L,out}(\phi)$, in order to evaluate difference of hue value between the object and the circumference. The hue value of right image at the position $^{IR}\boldsymbol{r}_i$ is expressed as $p(^{IR}\boldsymbol{r}_i)$, and the hue value of left image at the position $^{IL}\boldsymbol{r}_i$ is expressed as $p(^{IL}\boldsymbol{r}_i)$. Equation (1) shows the fitness function that calculate the correlation function between the search model and image.

$$
\begin{aligned}
F(\phi) = & \left\{ \left( \sum_{^{IR}\boldsymbol{r}_i \in S_{R,in}(\phi)} p(^{IR}\boldsymbol{r}_i) - \sum_{^{IR}\boldsymbol{r}_i \in S_{R,out}(\phi)} p(^{IR}\boldsymbol{r}_i) \right) \right. \\
& \left. + \left( \sum_{^{IL}\boldsymbol{r}_i \in S_{L,in}(\phi)} p(^{IL}\boldsymbol{r}_i) - \sum_{^{IL}\boldsymbol{r}_i \in S_{L,out}(\phi)} p(^{IL}\boldsymbol{r}_i) \right) \right\} / 2 \\
= & \ \{ F_R(\phi) + F_L(\phi) \} / 2 \qquad (1)
\end{aligned}
$$

The projected model area, $S_{p,q}(p = L, R; q = in, out)$ are all depending on the model's assumed pose that in designated by gene of GA's evolutionary processes. In the right imaging range, Eq. (1) deducts the total value in $S_{R,out}(\phi)$,
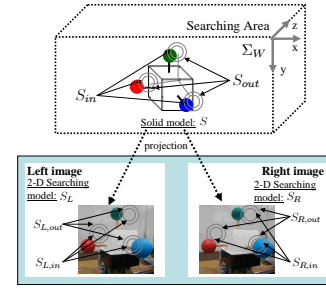


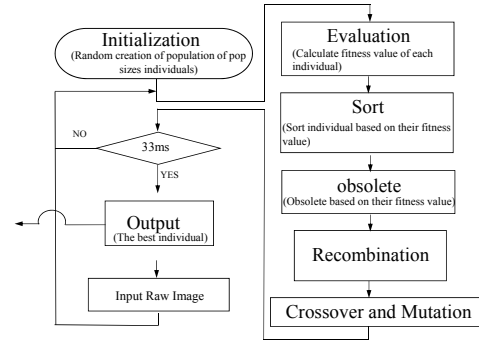Fig. 2. Definition of a solid model and left/right searching models



Fig. 3. Real-Time Multi-Step GA

from $S_{R,in}(\phi)$ and obtains the fitness value of the right image from the total value of the hue value $p(^{IR}\boldsymbol{r}_i)$ of an input image. The left imaging range is also the same. The fitness value of the left and the right is added and an average is taken. The image of the left and the right is simultaneously evaluated using this fitness function. This fitness function composed of "in" and "out" area aimed at to determine correlation value more emphasized then calculating correlation merely by solid model.

When a solid model $S_{R,in}(\phi)$ and $S_{L,in}(\phi)$ are cprrectly in agreement with the object, both in the right-and-left image, the object and the search model must be in agreement in 3-D pose. Though there is no gurantee that the variables to give the highest peek of the correlation function, i .e. , fitness function coincides with the true pose of object, but we can make effects to realize such conditions; righting condition, shape and color of the target, simple backdrop, and so on. Then we assume in this paper that highest peak of the fitness function indicates that the variables to give the peak represents true pose of the target object. It is defined as $F_R(\phi) = 0$ if $F_R(\phi) \leq 0$ and $F_L(\phi) = 0$ if $F_L(\phi) \leq 0$.

## 2.2 The Optimal Solution Searching Method using GA

By using a fitness function, the problem searches for the pose of an object can be transposed to the problem which searches for the maximum of a fitness function $F(\phi)$. In this research, we use GA to get the maximum fitness value within
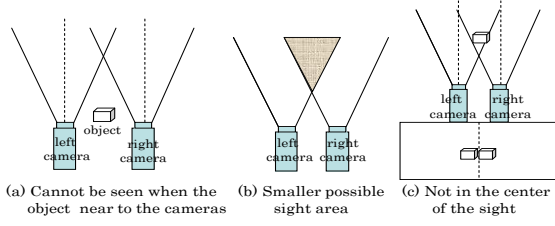
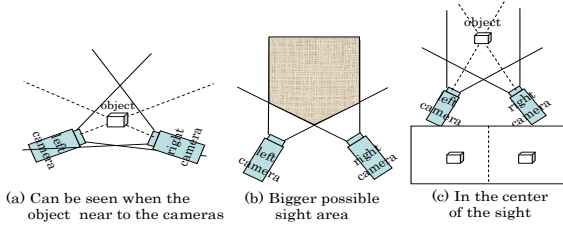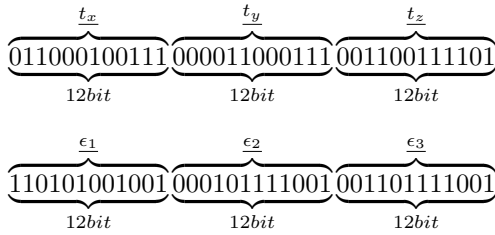Fig. 4. Disadvantage of fixed camera system



Fig. 5. Advantage of Eye-vergence system

less that video rate. Moreover, the gene information showing the position/orientation on the individual in this research is shown below.



The position/orientation of the individual shows the pose of the solid model in the Model-based Matching method. Top 36 bits with every 12 bits of this gene express the position coordinate of a solid model, and remainder 36 bits with every 12 bits expresses the orientation of the solid model, where the orientation is defined by quaternion. Bit used at this time may be reduced for searching time shortening.

Next, each individual gene get fitness value from the fitness function $F(\phi)$ using its pose information. Evolution processing is performed based on the superiority or inferiority of this value, and a set of the next generation is generated through GA's process. At this time, the pose in which fitness value was high in former generation, that is, it approaches toward the maximum neighborhood of the fitness function showing object. By repeating this processing (change of generation), GA discovers the maximum value showing the true pose of the target object.

However, normal GA needs to wait for convergence for a definite period of time. When a fitness function shows a value high enough and estimation of object is completed often more time has passed by, then, there is a possibility that
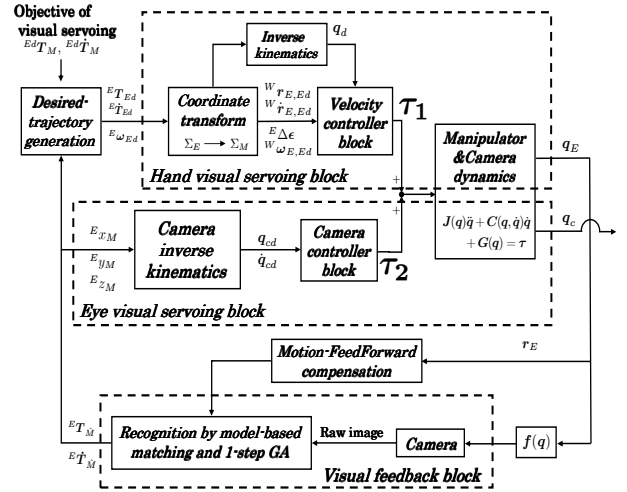


Fig. 6. Hand & Eye-Vergence Visual Servo System

the surrounding situation has changed a lot, that means target object may turned into a very different pose. Therefore we use RT-MS GA (Fig. 3). RT-MS GA is on-line estimation method [8]. To RT-MS GA, its evolving speed to optimize the fitness function should be faster than the target object's moving speed, then we can obtain the best gene in each time point, and by using the best gene, we can realize the recognition of target's pose and posture.

Utilizing RT-MS GA with reasonable performance in one loop and increasing accuracy with repeatable ability within real-time video rate is our approach strategy comparing to others that may provide powerful accuracy but also with computational burden and time-consuming.

## 3 SERVOING SYSTEM

### 3.1 Eye-vergence System

Although hand-eye composition has a shortcoming that servo operation may become unstable easily by vibration of the hand or time-delay in visual pose detection compared with fixed-camera system, but the merits of eye-vergence system is that the viewpoint can be adjusted in order to find a suitable viewpoint. In this paper, the visual servoing of hand-eye composition with two cameras is considered. If it assumes that object form is known, it is possible to measure six variables of a position/orientation also by a single eye. However, it is well known that there is a problem to measure the distance between camera and objects in a high precision, but the compound eye composition can handle this problem well.

A fixed-hand-eye system has some disadvantages, making the observing ability deteriorated because of the relative geometry of the camera and the target. Such as: the robot cannot observe the object well when it is near the cameras (Fig. 4 (a)), small intersection of the possible sight space of
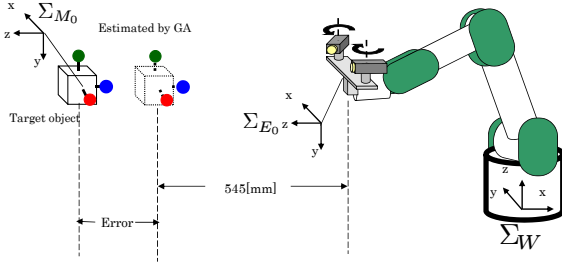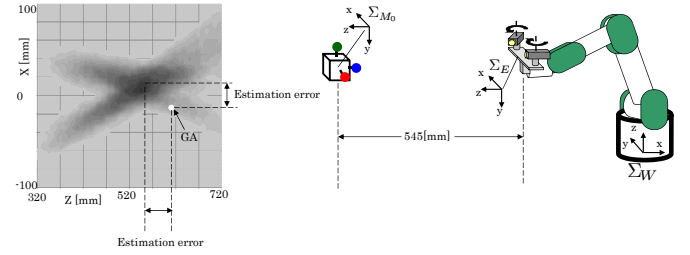
Fig. 7. Coordinate system



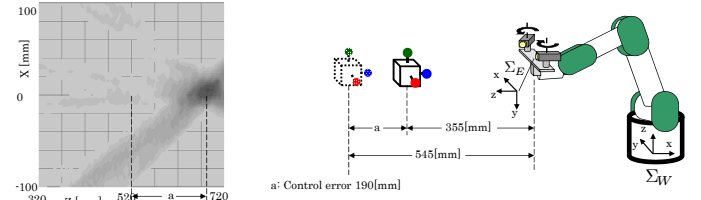Fig. 8. The relation of the position of target object and hand



Fig. 9. The relation of the position of target object and hand (estimated position is further than real target)
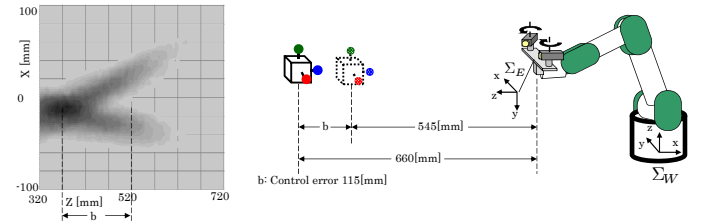


Fig. 10. The relation of the position of target object and hand (estimated position is nearer than real target)

the two cameras (Fig. 4 (b)), and the image of the object cannot appear in the center of both cameras, so we could not get clear image information of target and its periphery, reducing the pose measurement accuracy (Fig. 4 (c)).

To solve the problems above, in this research, we have chosen Eye-Vergence system that gives the cameras an ability to rotate themselves to project a target at center of the images.

Thus it is possible to change the pose of the cameras in order to observe the object better, as it is shown in Fig. 5, enhancing the measurement accuracy in trigonometric calculation and avoiding peripheral distortion of camera lens by observing target at the center of lens.

### 3.2 Hand & Eye Visual Servoing Controller

The block diagram of our proposed hand & eye-vergence visual servoing controller is shown in Fig. 6. Each joint angle of a manipulator is set to $\boldsymbol{q}_E = [q_1, \cdots, q_7]$, pan tilt angle of a camera is set to $\boldsymbol{q}_c = [q_8, q_9, q_{10}]$, the desired angle of each link is set to $\boldsymbol{q_d}$. The hardware control system of the velocity-based servo system of Mitsubishi Heavy Industries, Ltd PA10 is expressed as

$$\boldsymbol{\tau} = \boldsymbol{K}_{SP}(\boldsymbol{q}_d - \boldsymbol{q}) + \boldsymbol{K}_{SD}(\dot{\boldsymbol{q}}_d - \dot{\boldsymbol{q}}) \qquad (2)$$

where $\boldsymbol{K}_{SP}$ and $\boldsymbol{K}_{SD}$ are symmetric positive definite matrices to determine PD gain. Moreover, the target angle of a camera is set with $\boldsymbol{q}_{cd}$ and the error of an angle is defined as

$$\Delta\boldsymbol{q}_c = \boldsymbol{q}_{cd} - \boldsymbol{q}_c. \qquad (3)$$

The controller of eye-visual servoing is given by

$$\dot{\boldsymbol{q}}_{cd} = \boldsymbol{K}_{P_c}\Delta\boldsymbol{q}_c \qquad (4)$$

where $\boldsymbol{K}_{P_c}$ are positive control gain. The visual servoing is performed using these controllers [10].

## 4 EXPERIMENT

### 4.1 Experiment Condition

The initial hand pose is defined as $\Sigma_{E_0}$, and the initial object pose is defined as $\Sigma_{M_0}$. The homogeneous transformation matrix from $\Sigma_W$ to $\Sigma_{E_0}$ and from $\Sigma_W$ to $\Sigma_{M_0}$ are:

$$^W\boldsymbol{T}_{M_0} = \begin{bmatrix} 0 & 0 & -1 & -1435[mm] \\ 1 & 0 & 0 & 0[mm] \\ 0 & -1 & 0 & 420[mm] \\ 0 & 0 & 0 & 1 \end{bmatrix}, \qquad (5)$$

$$^W\boldsymbol{T}_{E_0} = \begin{bmatrix} 0 & 0 & -1 & -890[mm] \\ 1 & 0 & 0 & 0[mm] \\ 0 & -1 & 0 & 420[mm] \\ 0 & 0 & 0 & 1 \end{bmatrix}, \qquad (6)$$

where the above relations between $\Sigma_{M_0}$, $\Sigma_{E_0}$, $\Sigma_W$ are depicted in Fig.7.

The target object moves according to the following time function as:

$$^{M_0}z_M(t) = -150 + 150\cos(\omega t)[mm]. \qquad (7)$$

The relation between the object and the desired end-effector is set as:

$$^{Ed}\boldsymbol{\psi}_M = [0, -100[mm], 545[mm], 0, 0, 0]. \qquad (8)$$

Since it is reasonable that the target object is incorrectly estimated when fitness value was low, visual servoing motion will be stopped in consideration of the safety during experiments on condition that fitness value decreased to 0.1 or below. The fitness function is obtained from correlation function of target object in the camera scene and the model pre-defined in the control computer, then the correlation function value presents convicting level of the detected peak by RT-MS GA expressing the real target moving arround the robot in the environment. Table.1 shows information set on GA used in this experiment.

Table 1. Parameters of GA

| Number of individual | 20 |
|---|---|
| Length of gene | 72[bit] |
| Selection rate | 0.40 |
| Mutation | 0.10 |
| Search range($\Sigma_E$) | $-100 \leq t_x \leq 100$ |
| | $320 \leq t_z \leq 720$ |

## 4.2 Experiment Results

In this paper, in order to examine reliability of MT-RS GA, the position($x, z$ coordinates)was unknown, the posture($\varepsilon_1, \varepsilon_2, \varepsilon_3$) and the $y$-direction coordinate to preform a full search experiment were obtained from the visual-servoing experiment, furthermore, different from the previous experiments, we added two more factors to increase the accuracy of the full which is the angle of the two cameras($\theta_1, \theta_2$), with the visual servoing experiments were performed in $\omega = 0.419$ rad/s (Period: $T = 15[s]$ ) and $\omega = 0.209$ rad/s (Period: $T = 30[s]$ ), the distribution of the fitness value is calculated against all possible positions of $x$-$z$ plane, where the value of $x$-$z$ position that gives the highest peak means the target's $x$-$z$ position.

In Fig. 8, the time varying fitness function distribution combined by both left and right camera is depicted, and the position of the white circle with a notation "GA" represents that the RT-MS GA has found the indicated position represents most possible position. The distance between the white circle and the crosspoint combined by both left and right camera means the error of GA estimated position/posture.

Therefore the deviation between the peak (target real position) and the GA's position means on-line tracking error of object in 3-D space as shown in Fig. 8.

Table 2. Parameters in experiment of $\omega = 0.209$rad/s

| Time points[s] | Fitness value estimated by GA | Maximum fitness value calculated by full search | Position of maximum fitness value in GA($x, y, z$)[mm] | Position of maximum fitness value in full search($x, y, z$)[mm] |
|---|---|---|---|---|
| 6.376 | 0.694 | 0.7685 | -19.336, -128.398, 485.918 | -17, -128.398, 506 |
| 12.774 | 0.796 | 0.8241 | -4.59, -127.373, 519.707 | -6, -127.373, 506 |
| 19.15 | 0.852 | 0.8704 | 5.176, -120.244, 587.676 | 2, -120.244, 620 |
| 25.536 | 0.731 | 0.787 | 2.93, -117.363, 575.859 | 2, -117.363, 612 |
| 31.91 | 0.796 | 0.8611 | -8.691, -123.906, 500.859 | -3, -123.906, 524 |
| 44.659 | 0.796 | 0.8056 | 1.27, -122.93, 552.227 | 0, -122.93, 586 |
| 51.027 | 0.815 | 0.8333 | -8.789, -122.002, 594.023 | -13, 122.002, 624 |
| 57.385 | 0.782 | 0.787 | -0.586, -125.908, 563.75 | -1, -125.908, 596 |

Table 3. Parameters in experiment of $\omega = 0.419$rad/s

| Time points[s] | Fitness value estimated by GA | Maximum fitness value calculated by full search | Position of maximum fitness value in GA($x, y, z$)[mm] | Position of maximum fitness value in full search($x, y, z$)[mm] |
|---|---|---|---|---|
| 5.734 | 0.509 | 0.7685 | 17.676, -126.055, 430.352 | 28, -126.055 , 472 |
| 11.968 | 0.63 | 0.6926 | 0.977, -104.814, 620.781 | 4, -104.814, 648 |
| 18.484 | 0.593 | 0.6944 | 2.539, -132.109, 485.82 | 9, -132.109, 520 |
| 25.515 | 0.667 | 0.713 | -3.32, -107.598, 589.727 | -10, -107.598, 620 |
| 32.547 | 0.685 | 0.7407 | 23.438, -119.365, 501.543 | 29, -119.365, 540 |
| 40.64 | 0.639 | 0.713 | -1.758, -77.031, 607.598 | -1, -77.031, 636 |
| 48.937 | 0.722 | 0.7593 | -2.734, -125.273, 468.242 | -2, -125.273, 495 |
| 58.047 | 0.63 | 0.6574 | 7.129, -92.949, 652.715 | 5, -92.949, 682 |

Fig. 9 depicts the situation that the object place nearer than prescribed desired position relation of the object and the hand, and Fig. 10 shows the opposite situation. Fig. 11 $\omega = 0.209$ rad/s, and Fig. 12 is $\omega = 0.419$ rad/s shows the experiments result of circular frequency. It is the result obtained by full search in the plane of $z$-$x$ in $\Sigma_W$ for fitness value which performed 8 times randomly. The Table.2 and Table.3 are showing the time points of full search, the fitness value in GA recognition and full search. Furthermore, the position of the fitness value in GA and full search are also been shown. GA's gene gotten highest fitness value is displayed on these results by the white circle. The position of the depth direction ( $z$-direction) of the hand in $\Sigma_W$, the estimation results and the relation of the position of the depth direction of an actual object are shown in the graph. When period of the target's motion is $T = 30[s]$, it can be seen that the end-effector tracks the target with fewer place delay, and the RT-MS GA tracks correctly the highest peak even the target was changing its position in real time. Even though the condiction $T = 15[s]$ in Fig. 12, to make the correct tracking of end-effector is difficult, but RT-MS GA maintains real-time tracking of target through real-time optimization of time-varying fitness distributions. Shown as the Fig.11and Fig.12, the relation of full-search results can be easily under-stood. When fitness value is high, the color is deep; Inversely, the color is light when fitness value is low.

## 4.3 Discussion

Full search results of fitness value show that fitness value is distributed in the shape of x or shape of y centering on a deep-colored point (fitness value is high). This visual servoing system is composed by two eyes, and it is asking for fitness value, using simultaneously the image obtained from the right-and-left camera. Thereby, GA becomes easy to discover an object. As the experiments results, we can figure out that the RT-MS GA can obtain a reliable fitness value in a very short time, and according to the Table.2 and Table.3, the position to get the highest fitness value in GA is near to the position of maximum fitness value when calculated by full search, comparing with the long-time-spent full search method, the superiority of RT-MS GA can been easily under-stood.

## 5 CONCLUSION

In this paper, the visual servoing was performed in the experiment of frequency response, and the action of GA gene in a visual servoing was shown. The real-time estimation tracking error has been grasped by clarifying relationship of GA and an object. RT-MS GA became easy to calculate the optimal solution with eye-vergence system. Comparing with the full search method, the superiority of RT-MS GA can been confirmed. During the experiments, there was a case that the result of full search was out of the search range, which means there are maybe a possibility that the best fitness value is still not been found. As a future subject, it is necessary to consider change of the searching area of GA.

## REFERENCES

[1] S.Hutchinson, G.Hager, and P.Corke, "A Tutorial on Visual Servo Control", IEEE Trans. on Robotics and Automation, vol. 12, no. 5, pp. 651-670, 1996.

[2] P.Y.Oh, and P.K.Allen, "Visual Servoing by Partitioning Degrees of Freedom", IEEE Trans. on Robotics and Automation, vol. 17, no. 1, pp. 1-17, 2001.

[3] E.Malis, F.Chaumentte and S.Boudet, "2-1/2-D Visual Servoing", IEEE Trans. on Robotics and Automation, vol. 15, no. 2, pp. 238-250, 1999.

[4] P.K.Allen, A.Timchenko, B.Yoshimi, and P.Michelman, "Automated Tracking and Grasping of a Moving object with a Robotic Hand-Eye System", IEEE Trans. on Robotics and Automation, vol. 9, no. 2, pp. 152-165, 1993.

[5] W. Song, M. Minami, Y. Mae and S. Aoyagi, " On-line Evolutionary Head Pose Measurement by Feedforward Stereo Model Matching ", IEEE Int. Conf. on Robotics and Automation (ICRA), pp.4394-4400, 2007.

[6] J. Stavnitzky, D. Capson, "Mutiple Camera Model-Based 3-D Visual Servoing", IEEE Trans. on Robotics and Automation, vol. 16, no. 6, December 2000.c

[7] C. Dune, E. Marchand, C. leroux, " One Click Focus with Eye-inhand/Eye-to hand Cooperation ", IEEE Int. Conf. on Robotics and Automation (ICRA), pp.2471-2476, 2007.

[8] H. Suzuki, M. Minami, "Visual Servoing to catch fish Using Global/local GA Search", IEEE/ASME Transactions on Mechatronics, Vol.10, Issue 3, 352-357 (2005.6).

[9] M.Minami, W.Song, "Hand-eye-motion Invariant Pose Estimation with On-line 1-step GA -3D Pose Tracking Accuracy Evaluation in Dynamic Hand-eye Oscillation", Journal of Robotics and Mechatronics, Vol.21, No.6, pp.709-719 (2009.12)

[10] Wei. Song, M. Minami, Fujia Yu, Yanan Zhang and Akira Yanou "3-D Hand & Eye-Vergence Approaching Visual Servoing with Lyapunouv-Stable Pose Tracking ", IEEE Int. Conf. on Robotics and Automation (ICRA), pp.11, 2011.
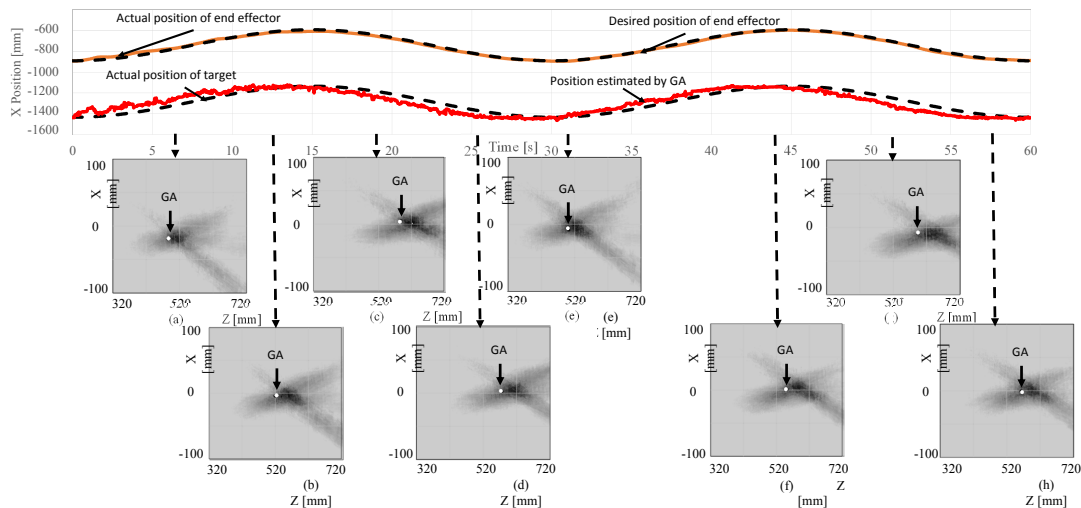
Fig. 11. Relation between fitness value, position of end effector ,actual position of target object and position estimated in z-x plane by GA in $\omega = 0.209$ rad/s
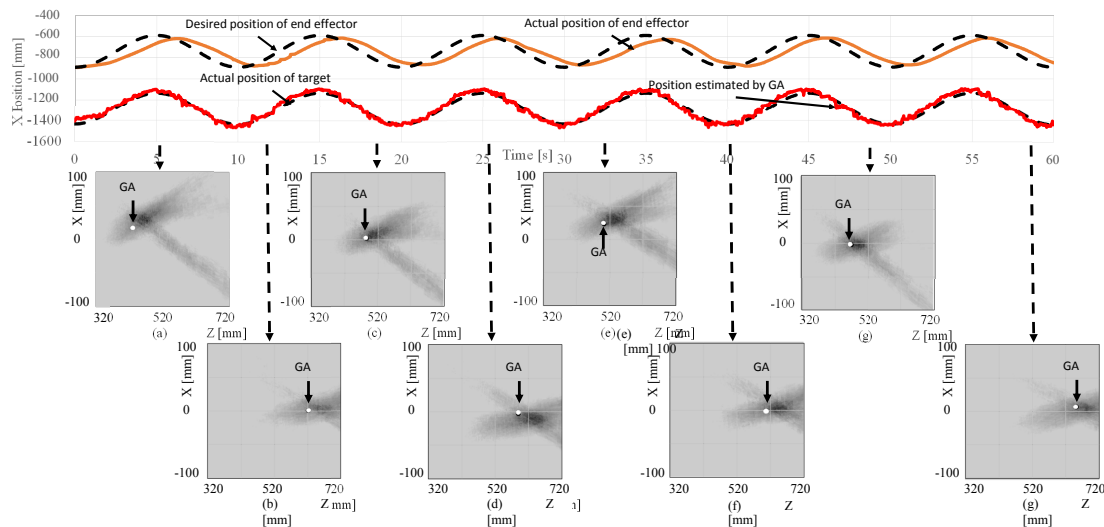


Fig. 12. Relation between fitness value, position of end effector ,actual position of target object and position estimated in z-x plane by GA in $\omega = 0.419$ rad/s