# Visual servoing to arbitrary target with photo-model-based recognition method

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

[1]Okayama University, Japan
(Tel: 81-86-251-8233, Fax: 81-86-251-8233)
[1]psnc8ytd@s.okayama-u.ac.jp

**Abstract:** Visual servoing is one of the methods to control robots. By incorporating visual information obtained from the installed vision sensor into the feedback loop, it is desirable for the robot to be able to operate in a changing environment or an unknown environment. For visual servo systems, the authors proposed a photo-model-based recognition method for processing an arbitrary object with a still picture. In the previous work, a flat shape clothes handling robotic system has been proposed to handle deformable and unique clothes. In this paper, we extend the photo-model-based robot handling system (pick and place) to a real-time pose tracking system. And the recognition method is used for attitude tracking of different objects of dynamic pictures. Furthermore, to verify the ability to track with the photo-model-based recognition method. The authors design some frequency response experiments with arbitrary aquatic creature toys to keep the relative pose between a sea animal and hand-eye. The results of visual servoing experiments show that the proposed identification method is feasible, flexible and effective.

**Keywords:** visual servoing, photo-model-based recognition, genetic algorithm(GA)

## 1 INTRODUCTION

Since robots have higher reliability and accuracy than humans, they have been used extensively in production factories to perform a wide variety of tasks instead of human workers.

However, until now, robots cannot entirely replace humans. While human beings can conduct intended tasks in pending circumstances, an automated robot is not adept at being similarly adaptable. Therefore, the researchers have tried to improve the abilities of automated robots.

About automated robots, the visual servoing, a robot control technology using visual information obtained from a vision sensor (camera) in the feedback loop, is expected to be able to allow the robot to adapt to changing or unknown environments [1, 2, 3].

In the previous works [4] and [5], a photo-model-based matching method has been proposed. With this method, we developed a clothes-handling robot. In [6], 3D recognition accuracy has been confirmed experimentally by using 12 different samples cloths. Except for static clothes handling, there is a need for tracking arbitrary moving target object. With photo-model-based recognition method, we developed a visual servoing system. As shown in Fig. 1. The dual-eye cameras that are fixed at the end-effector of a PA-10 robot perform the object recognition and pose estimation process based on the photograph model.

For verifying the ability of tracking an arbitrary object, in this paper, the authors conduct some frequency response experiments with arbitrary aquatic creature toys to keep the relative pose between a sea animal and hand-eye. The results of visual servoing experiments show that the proposed identification method is feasible, flexible and effective.
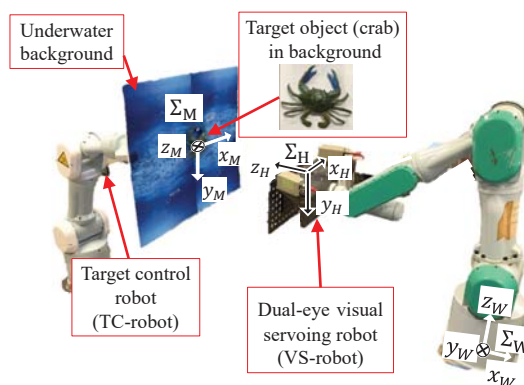


Fig. 1. The visual servoing system with dual eye-in-hand cameras. World coordinate system ($\Sigma_W$) and target coordinate system ($\Sigma_M$)

## 2 PHOTO-MODEL-BASED RECOGNITION

To make it easier to understand the photo-model-based recognition method, we will describe the kinematics of stereo-vision before an explanation of the proposed system in details.

### 2.1 Kinematics of stereo-vision

Figure 2 shows a perspective projection of the dual-eyes vision system. Each coordinate system are as follows:

- $\Sigma_W$: world coordinate system,

- $\Sigma_H$: end-effector (hand) coordinate system, as

- $\Sigma_{CL}, \Sigma_{CR}$: left and right camera coordinate systems,

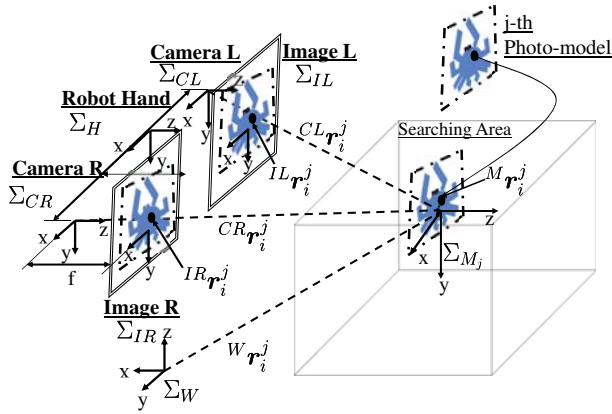- $\Sigma_{IL}, \Sigma_{IR}$: left and right image coordinate systems,

Fig. 2. Perspective projection of dual-eye vision-system: In the searching area, a 3D solid model is represented by the picture of crab with black point (j-th photo model)

- $\Sigma_{M_j}$: j-th model coordinate system,

- $\Sigma_M$: object coordinate system.

- ${}^M r_i^j$: position of an arbitrary i-th point on j-th 3D model in $\Sigma_{M_j}$, where ${}^M r_i^j$ is a constant vector

- ${}^{CR} r_i^j$ and ${}^{CL} r_i^j$: position of an arbitrary i-th point on j-th 3D model based on $\Sigma_{CR}$ and $\Sigma_{CL}$

- ${}^{IL} r_i^j$ and ${}^{IR} r_i^j$: projected position on $\Sigma_{IL}$ and $\Sigma_{IR}$ of an arbitrary i-th point on j-th 3D model

The homogeneous transformation matrix from $\Sigma_{CR}$ to $\Sigma_M$ is defined as ${}^{CR}T_M(\phi_M^j, q)$, where $\phi_M^j$ is j-th model's pose and $q$ means robot's joint angle vector. Then, ${}^{CR} r_i^j$ can be calculated by using Eq. (1),

$$ {}^{CR} r_i^j = {}^{CR}T_M(\phi_M^j, q) \, {}^M r_i^j. \qquad (1) $$

The position vector of the i-th point in the right and left camera image coordinates ${}^{IR} r_i^j$ can be described by using projective transformation matrix $P_k$ as,

$$ {}^{IR} r_i^j = P_k \, {}^{CR} r_i^j = P_k {}^{CR}T_M(\phi_M^j)^M r_i^j \qquad (2) $$

Then, ${}^{IR} r_i^j$ can be described as,

$$ \begin{cases} {}^{IR} r_i^j(\phi_M^j) = f_R(\phi_M^j, {}^M r_i^j) \\ {}^{IL} r_i^j(\phi_M^j) = f_L(\phi_M^j, {}^M r_i^j) \end{cases} \qquad (3) $$

where ${}^{IL} r_i^j$ can also be described as the same manner like ${}^{IR} r_i^j$.

## 2.2 Model generation

The model generation process is represented as Figure 3. It should be noted that the model is only part of a picture and the picture is not the model. Firstly, a background image is captured by the first camera and the averaged hue value of



(a) Background     (b) Target object in background

(c) $S_{in}$ space of model is shown by black points group     (d) Enveloping space of $S_{in}$ is shown by points group $S_{out}$
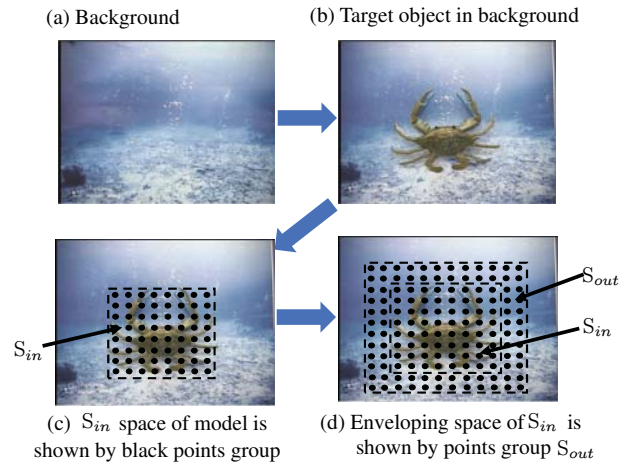
Fig. 3. (a) shows a photograph of background image, (b) shows a photograph of the target object (the blue crab) in background, (c) represents a photograph of surface space model $S_{in}$ by inner points group and (d) represents a photograph of outside space of model $S_{out}$ that enveloping $S_{in}$.

the background image is calculated as shown in Fig. 3 (a). Then, the crab is put on the background. Take a $640 \times 480$ pixels picture at a distance of 400[mm] from the object as shown in Fig. 3 (b). As shown in Fig. 3 (c), scan from the four corners at the same time in the arrow direction compare with the averaged hue value of the background, and generate the surface space $S_{in}$ of the model. Finally, the outside space $S_{out}$ of the model is generated by enveloping $S_{in}$ as shown in Fig. 3 (d). It is gotten that the size of the tangential plane of the object in the real 3D space.

## 2.3 3D photo-model-based matching

In Fig. 4, a generated solid model is projected from the 3D space onto the left and right 2D searching planes. The sub figure on the top of Fig. 4 shows a generated 3D solid model with its pose $S_{in}(\phi_M^j)$ (inner dotted points) and the outside space enveloping $S_{in}(\phi_M)$ denoted as outer dotted line $(S_{out}(\phi_M^j))$. The sub figure on the left/right bottom of Fig. 4 show the left/right 2D searching models $S_L(\phi_M^j)$ and $S_R(\phi_M^j)$ respectively. Both $S_L(\phi_M^j)$ and $S_R(\phi_M^j)$ consist of $S_{L,in}(\phi_M^j)$ and $S_{L,out}(\phi_M^j)$ and $S_{R,in}(\phi_M^j)$ and $S_{R,out}(\phi_M^j)$. The evaluation of the correlation between the projected model and the images from the dual-eye cameras attached at the end-effector defined as a fitness function.

## 2.4 Definition of the fitness function

The concept of the fitness function in this study can be said to be an extension of the work in [8] in which different models including a rectangular shape surface-strips model was evaluated using images from a single camera. The correlation between the projected model $\phi$ and captured images on the left and right 2D searching areas is calculated by the
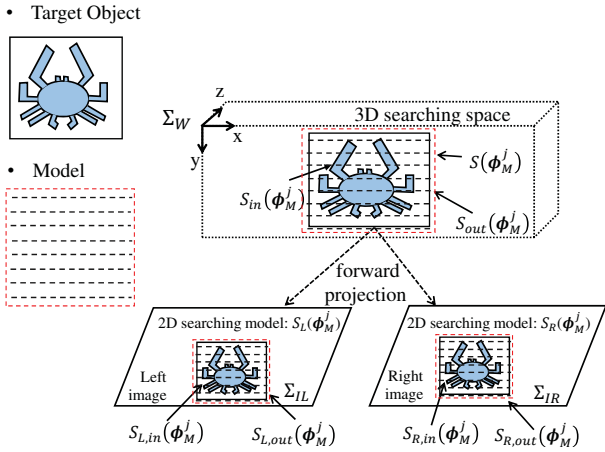
Fig. 4. Through projection transformation, a 3D solid model in the 3D searching space is projected to 2D left and right images. Searching models represented as $S_L(\phi_M^j)$ and $S_R(\phi_M^j)$

equations Eq. (4) to Eq. (6).

$$
F(\phi_M^j) = \left\{ \left( \sum_{\substack{IR r_i^j \in \\ S_{R,in}(CR \phi_M^j)}} p(IR r_i^j) + \sum_{\substack{IR r_i^j \in \\ S_{R,out}(CR \phi_M^j)}} p(IR r_i^j) \right) \right.
$$
$$
\left. + \left( \sum_{\substack{IL r_i^j \in \\ S_{L,in}(CL \phi_M^j)}} p(IL r_i^j) + \sum_{\substack{IL r_i^j \in \\ S_{L,out}(CL \phi_M^j)}} p(IL r_i^j) \right) \right\} / (2N) \tag{4}
$$

The evaluation of every point in the input image that lie inside the surface model frame and outside area of the model frame is represented as $^{IL}r_i^j \in S_{L,in}(\phi_M^j)$ and $^{IL}r_i^j \in S_{L,out}(\phi_M)$ respectively. $N$ is the total number of sampling points Eqs. (5) and (6) is used for calculating $p_{L,in}(^{IL}r_i^j)$ and $p_{L,out}(^{IL}r_i^j)$.

$$
p_{L,in}(^{IL}r_i^j) = \begin{cases} 2, & \text{if}(|H_{IL}(^{IL}r_i^j) - H_{ML}(^{IL}r_i^j)| \le 30); \\ -1, & \text{if}(|H_{IL}(^{IL}r_i^j) - H_{ML}(^{IL}r_i^j)| \ge 50); \\ -0.005, & \text{if}(|\bar{H}_B - H_{ML}(^{IL}r_i^j)| \le 30); \\ 0, & \text{otherwise.} \end{cases} \tag{5}
$$

$$
p_{L,out}(^{IL}r_i^j) = \begin{cases} 0.1, & \text{if}(|\bar{H}_B - H_{IL}(^{IL}r_i^j)| \le 20); \\ -0.5, & \text{otherwise.} \end{cases} \tag{6}
$$

where

- $H_{IL}(^{IL}r_i^j)$: the hue value of the left camera image at the point $^{IL}r_i^j$ (i-th point in $S_{L,in}$),

- $H_{ML}(^{IL}r_i^j)$: the hue value of the point $^{IL}r_i^j$ (i-th point

in $S_{L,in}$) on the model ,

- $\bar{H}_B$: the average hue value of the background image

The evaluation values are tuned experimentally. In Eq. (5), if the hue value of each point of captured images, which lies inside the surface model frame $S_{L,in}$, is same to the hue value of each point in a model, the fitness value will increase with the voting value of "+2." The fitness value will decrease with the value of "−0.005" for every point of crabs in the left camera image that are similar to the average hue value of the background. Similarly, in Eq. (6), if the hue value of each point in the left camera image, which are in $S_{L,out}$, is same to the hue value of the background, with the tolerance of 20, the fitness value will increase with the value of "0.1." Otherwise, the fitness value will be decreased with the value of "−0.5." Similarly, a function $p_{R,in}(^{IR}r_i^j)$ and $p_{R,out}(^{IR}r_i^j)$ are represented for the right camera image.

### 2.5 Fitness distribution of position

To evaluate the feasibility of fitness function Eq. (4), a good way is a brute-force search or an exhaustive search. For still pictures of a moment, the fitness values of all candidates that represent different poses of models are calculated. We call it "fitness distribution."

In fact, for the measurement of position or orientation, it is impossible to exhaust all possibilities. For fitness distribution, the accuracy of position $^H x_M - ^H y_M$ is 5[mm]. We prepared 3D toys of 6 marine creatures as shown in Fig. 5 (a). The four labels corresponding to each model are number, English name, size and Japanese name. To recognize, all the objects were taken pictures and saved in a database as shown in Fig. 5 (b). The true value of all objects is

$$
^H\psi_M = [0, 0, 500, 0, 0, 0][mm]. \tag{7}
$$

Fig. 6 are experimental results. In each frame, the lower right two pictures are the images captured by the left and right camera. And target object is explained on the left side. From the results, each peak appears near the true value. The error is small. The last experiment searches C02 from several objects. Since its peak is near the true value, it can be said that this method can search for an object from a complex background.

### 2.6 Real-time Multi-step Genetic Algorithm (RM-GA)

As shown in Fig. 8, searching for all possible models is time-consuming for real-time recognition. Therefore, we convert the problem of finding/recognizing the target object's pose into an optimization problem with a multi-peak distribution. GA is a simple and effective searching method. For real-time recognition through one frame image in 33[ms], we have proposed a Real-time Multi-step Genetic Algorithm

(a) Details of 6 target toys (Unit: cm)



(b) 12 target images in database

Fig. 5. (a) 6 marine biological models. Code name is from C01 to C06. The size of each 3D toy is shown in left down of each frame. The four labels corresponding to each model are number, English name, size and Japanese name. (b) pictures of marine biological models with blue sea background. The size of each picture is $640 \times 480$ pixels. It should be noted that the model is only part of a picture and the picture is not the model.
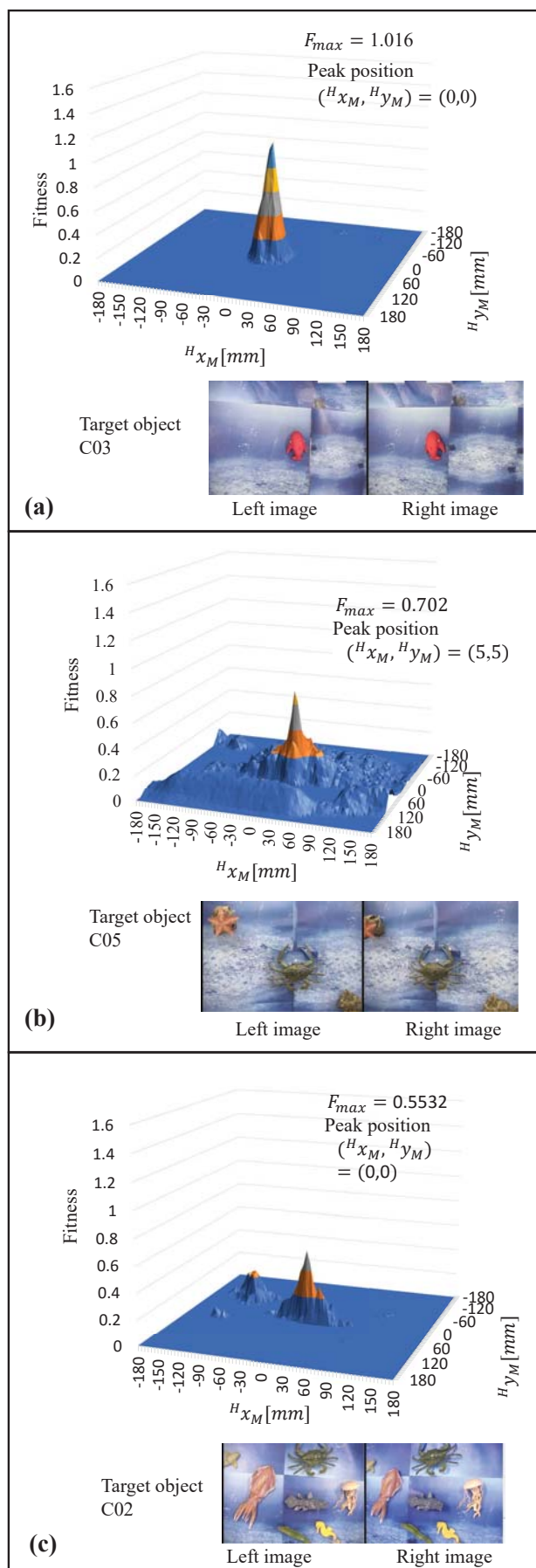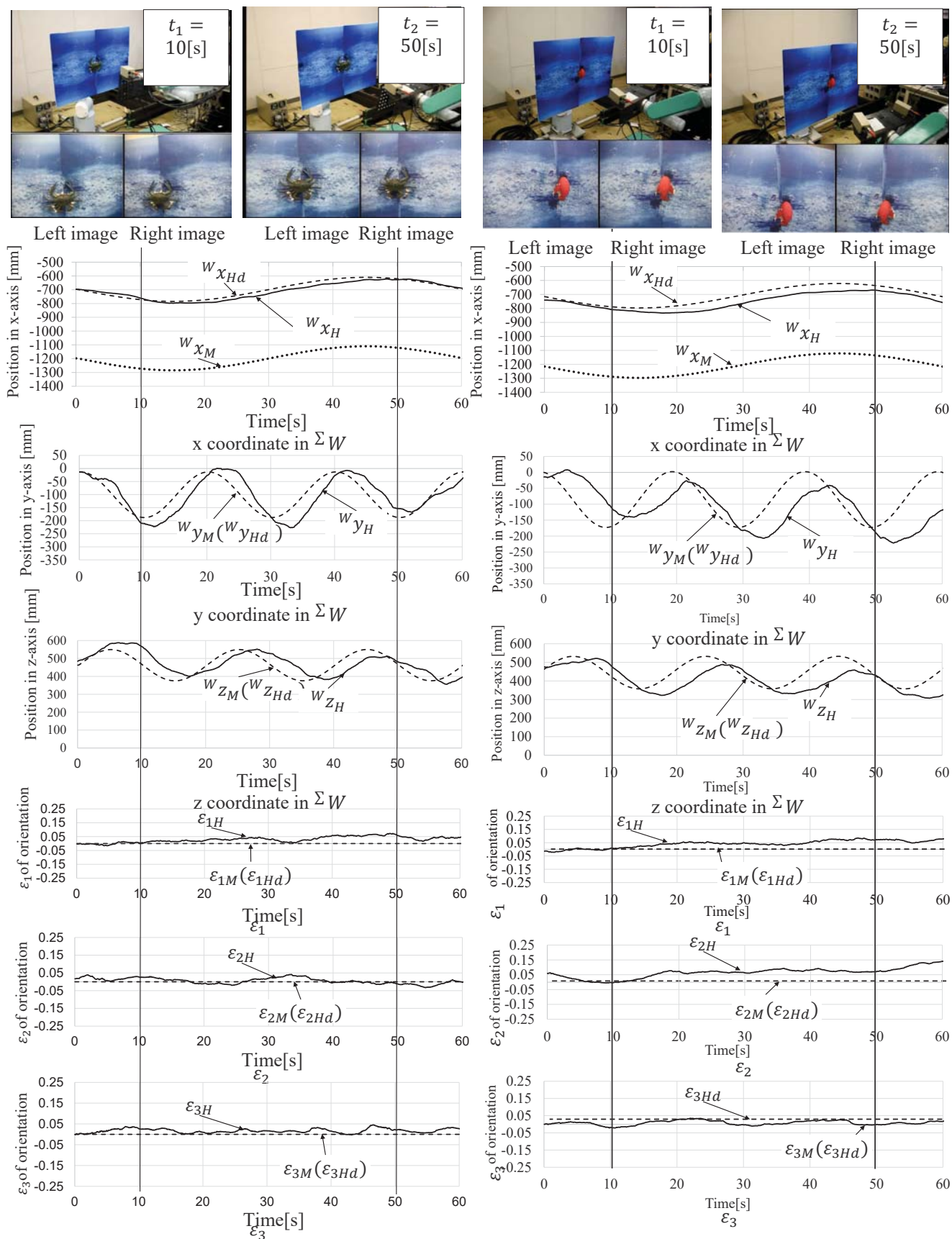






Fig. 6. Due to the limited space, this figure shows fitness distribution results of object C02, C03, and C05 on $^H x_M - ^H y_M$ plane. C03 and C05 are chosen randomly for later visual servoing experiment. In frame (c), target object is C02. We verify the robustness of this recognition method to the background with (c) experiment.

Fig. 7. Pose estimation and tracking performance experiments of the 3D solid target object. Positions of crab and dolphin change in sine wave. Object's orientation does not change.

Fitness distribution: Exploring every possible location takes a lot of time.

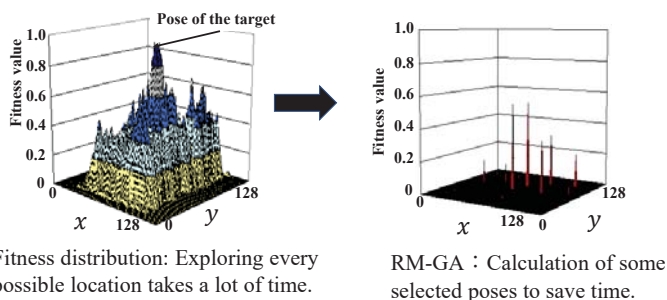RM-GA：Calculation of some selected poses to save time.

Fig. 8. we have proposed Real-time Multi-step Genetic Algorithm (RM-GA) for searching the pose of target object in real-time.
(RM-GA)[7] although it may not be the best GA in comparison to other optimization methods. We did not compare GA with other optimization methods in this study. With this algorithm, Each chromosome consists of six variables. Each variable is coded by 12bits that can provide sufficient accuracy to get the optimal solution. The first three variables of a model in 3D space $(t_x, t_y, t_z)$ are represented as the position and the last three variables $(\varepsilon_1, \varepsilon_2, \varepsilon_3)$ are represented as the orientation.

$$\underbrace{01\cdots01}_{12bits}\underbrace{00\cdots01}_{12bits}\underbrace{11\cdots01}_{12bits}\underbrace{01\cdots01}_{12bits}\underbrace{01\cdots11}_{12bits}\underbrace{01\cdots10}_{12bits}.$$

Readers can refer to [7], which has a more detailed explanation.

## 3   EXPERIMENTAL ENVIRONMENT

Two frequency response experiments with crab and dolphin toy are conducted to confirm the tracking ability of the developed photo-model-based visual servoing system. As shown in Fig. 1, during the experiments the desired position and orientation is the same as Eq. (7). The target trajectories of the two experiments are the same. Both of them are sine curves with an amplitude of 100 [mm], a period of 20 [s] in the x and y-axis directions and amplitude of 100 [mm], a period of 60 [s] in the z-axis direction.

## 4   EXPERIMENTAL RESULTS AND DISCUSSION

In Fig.7, $^W\boldsymbol{r}_H$ and $^W\boldsymbol{\varepsilon}_H$ are the tracking results of the end-effector. $^W\boldsymbol{r}_M$ and $^W\boldsymbol{\varepsilon}_M$ are the motion of target object. And the desired positions of end-effector $^W\boldsymbol{r}_{Hd}$ and $^W\boldsymbol{\varepsilon}_{Hd}$ are ideal positions during the experiments. It can be seen that even though the tracking curves delay somewhat in phase, the visual servoing system with photo-model-based recognition method can track the object in time. And different objects have different effects on the tracking results. On this point, we will conduct further research and discussion.

## 5   CONCLUSION

The visual servoing experiments were conducted to confirm the performance of the photo-model-based recognition method. According to the experimental results, this system can recognize and track the 3D crab and dolphin toys with the prepared pictures.

We conclude that if we have prepared the pictures of objects the system can recognize them and track them. However, different objects seem to have different effects on the tracking performance of the system. In the future, we would like to discuss this problem.

## 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, 1996, pp. 651-670.

[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, 2001, pp. 1-17.

[3]  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.

[4]  Funakubo R, Phyu KW, Tian H, Minami M, Recognition and handling of clothes with different pattern by dual hand-eyes robotic system, IEEE/SICE International Symposium, 2016, pp. 742-747.

[5]  Phyu KW, Cui Y, Tian H, Hagiwara R, Funakubo R, Yanou A, Minami M, Accuracy on Photo-Model-Based Clothes Recognition, SICE Annual Conference, Tsukuba, Japan, 2016, September 20-23.

[6]  Phyu, K.W., Funakubo, R., Fumiya, I., Shinichiro, Y. and Minami, M., Verification of recognition performance of cloth handling robot with photo-model-based matching, IEEE International Conference on Mechatronics and Automation (ICMA), 2017, pp. 1750-1756.

[7]  Myint, M., Yonemori, K., Lwin, K.N., Yanou, A. and Minami, M., Dual-eyes vision-based docking system for autonomous underwater vehicle: an approach and experiments, Journal of Intelligent & Robotic Systems, 92(1), 2018, pp.159-186.

[8]  Minami M, Agbanhan J, Asakura T, Evolutionary scene recognition and simultaneous position/orientation detection, in Soft Computing in Measurement and Information Acquisition, Springer Berlin Heidelberg, 2003, pp. 178-207.