Accuracy on Photo-Model-Based Clothes Recognition

Khaing Win Phyu^{1†}, Yu Cui¹, Hongzhi Tian¹, Ryota Hagiwara¹, Ryuki Funakubo¹, Akira Yanou¹ and

Mamoru Minami¹

¹Department of Intelligent Robotics and Control Laboratory, Okayama University, Okayama, Japan (Tel: +81-86-251-8233; E-mail: minami-m@cc.okayama-u.ac.jp)

Abstract: Recently, robots have been used in clothing industries for mass production with countless merits. However, there remain many challenges for robots in recognition, pose (position and orientation) detection operations, especially when the working object is deformable and every working object has unique shape and color. In this paper, pose detection of clothes through 3D recognition is proposed for the task in which the manipulator recognizes clothes, estimates relative pose and performs pick and place function. In proposed cloth recognition system, a variety of models of different clothes with unique shape and color are generated as BMP (bit map file) format extracted from the camera. Using the photograph model, recognition of cloth is performed by using images from the two cameras that are fixed at the end effector of the robot arm. 12 varieties of different clothes samples are used for this experiment. The pose of individual cloth is estimated by Genetic Algorithm (GA). 1000 times recognition experiment has been executed, having shown the effectiveness of proposed Photo-Model-Based pose recognition system.

Keywords: Visual Servoing, GA, 3D-MoS, Recognition, Pose Detection

1. INTRODUCTION

Nowadays, most of the garment companies especially in Japan are facing with two main inconveniences as follow:

- Growing shortage of labor force because an aging population have been progressing.
- Weak point of human workers are laziness, boring and tiring due to the long working hours.

Recently, IT based technology and Robotics have began to be used in the garment (cloth) companies considering above problems. However, robotics in garment industry are capable of operations only if preconditions are met such as (1) surrounding the operating environment, the light environment is guaranteed not to change with time, (2) the handling object is able to be defined the shape of object in advance, and (3) the design of the robot hand is predefined based on the object shape. A number of research for robot in recognition deformable objects especially clothes have been done[9].

In application to cloth handling, the main tasks are to get the purchasing order of clothes from on line customer and classify these clothes and package and place in different box for storage every day. Since clothes is deformable object, no definition of clothes can be predefined in computer. Consequently it is too difficult to handle a wide variety of clothes that are in irregular shape and size. Therefore, we have developed vision-based robot system as shown in Fig.1 to solve the above problems for mass handling robot of clothes with varieties.

On the other hand, robot control technology using visual information, called as visual servoing, is playing an important role in industry. In our previous works, we have been developing a three-dimensional move on sensing system named 3D-MoS using two cameras as stereo vision sensor. The robot control technology using visual



Fig. 1 System Configuration

information has been already confirmed in guidance and control of underwater robot (3D-MoS / AUV)[6] and soil decontamination robot system (fully automated robotic system to decontaminate the radioactive contamination soil)[8]. In these previous works, model-based method is used to recognize 3D (three dimensional) pose.

However, in cloth handling application in which the manipulator has to pick and place unique-shape-andcolor clothes with random appearance in camera images, it is impossible to use models that are predefined for all clothes. Therefore, as a main contribution in this paper, we introduce a new approach to generate model for every clothes during the operation of pick and place. In proposed "Photo-model 3D-MoS", the model of object is automatically created from a photo by robot itself. Then, relative pose estimation of clothes is performed using generated model through model-based clothes recognition. PA10 with 7-DoF (Degree of Freedom) is being used for recognition and pose detection operations.

The configuration of the system is shown in Fig. 1. In this configuration, three cameras are used as vision sensors. First camera is used for model generation. The other two are used for recognition based on the photo-

[†] Khaing Win Phyu is the presenter of this paper.

graph model which are fixed at the end-effector of the mobile manipulator (PA10 robot). Instead of being done by the system, labeling to the cloth is done by manually. The order of appearance of clothes in view of first camera (first stage) and second two cameras (second stage) are assumed to be the same. However, the existence of other objects may occur between two stages. At that time, the robot will identify clothes among other objects by matching stored photo-models, and perform pose estimation and, pick and place the cloth automatically. Therefore, the tasks of the robot is (1) to recognize the cloth in the second stage to be the same cloth as detected one in the first stage, and (2) to measure the pose of the cloth to pick up it. Identification of cloths can be extended using bar coding method. The merits of this proposed system are to save the cost of staff and to get the better performance and higher accuracy than workers in the garment company.

2. PROPOSED PHOTOGRAPH MODELING-BASED CLOTHES RECOGNITION

There are two main portions in proposed system. The first portion is cloth model template generation and the latter is relative pose estimation using generated model template through model-based matching method. Here is description of kinematics of stereo-vision before explanation of proposed system in detail.

2.1. Kinematics of Stereo-Vision

Fig. 2 shows the relationships between the world coordinate system of the manipulator Σ_W and the hand coordinate system Σ_H . The coordinate system of dual-eyes vision system can be seen in Fig.3. An target object coordinate system is expressed in Σ_M . In the image coordinate system, the coordinate system of the left and right cameras are represented as Σ_{CR} and Σ_{CL} . Σ_{IR} and Σ_{IL} are the coordinate systems of the left and right cameras' images. According to coordinate system, the j-th point of the i-th model can be represented by the following the simultaneous transformation matrix.

- ${}^{CR}T_M$:Homogeneous transformation matrix from right camera coordinate system Σ_{CR} to the object coordinate system Σ_M
- ${}^{CR}r_i$: The object as viewed from the search point i-th on the model in Σ_{CR}
- ${}^{M}\boldsymbol{r}_{i}$: The object as viewed from the search point i-th on the model in Σ_{M}

Therefore, ${}^{CR}r_i$ can be calculated by using Eq. (1)

$$^{CR}\boldsymbol{r}_{i} = {}^{CR}\boldsymbol{T}_{M} {}^{M}\boldsymbol{r}_{i}. \tag{1}$$

The homogeneous transformation matrix ${}^{W}\boldsymbol{T}_{CR}$ from world coordinate system Σ_{W} to the right camera coordinate system Σ_{CR} can be obtained from Eq. (2).

$$^{W}\boldsymbol{r}_{i} = {}^{W}\boldsymbol{T}_{CR} {}^{CR}\boldsymbol{r}_{i}. \tag{2}$$

By using matrix P according to projective transformation, the position vector of the i-th point in the right camera image ${}^{IR}r_i$ can be described as Eq. (3). Eq. (4) is described as matrix P.

1

$$^{IR}\boldsymbol{r}_{i}=\boldsymbol{P}^{CR}\boldsymbol{r}_{i}. \tag{3}$$

$$\boldsymbol{P} = \frac{1}{C_{z_i}} \begin{bmatrix} \frac{f}{\eta_x} & 0 & {}^{I}x_0 & 0\\ 0 & \frac{f}{\eta_y} & {}^{I}y_0 & 0 \end{bmatrix}.$$
 (4)

Using the same method, it is possible to obtained the position vector of the i-th point in the left camera image ${}^{IL}r_i$.

$$^{CL}\boldsymbol{r}_{i} = {}^{CL}\boldsymbol{T}_{M} {}^{M}\boldsymbol{r}_{i}.$$

$${}^{W}\boldsymbol{r}_{i} = {}^{W}\boldsymbol{T}_{CL} {}^{CL}\boldsymbol{r}_{i}. \tag{6}$$

$${}^{L}\boldsymbol{r}_{i} = \boldsymbol{P} \; {}^{CL}\boldsymbol{r}_{i}. \tag{7}$$

According to Eq. (3) and Eq. (7), the relationship of Eq. (8) is connected an arbitrary point on a 3D-model with a pose ${}^{C}\phi_{M}$ – the pose Σ_{M} based on Σ_{CR} – with the projected point on the left camera image ${}^{IL}\mathbf{r}_{i}$ and right camera image ${}^{IR}\mathbf{r}_{i}$ can be written as .

$$\begin{cases} {}^{IR}\boldsymbol{r}_i = \boldsymbol{f}_R({}^{CR}\boldsymbol{\psi}_M, {}^{M}\boldsymbol{r}_i) \\ {}^{IL}\boldsymbol{r}_i = \boldsymbol{f}_L({}^{CL}\boldsymbol{\psi}_M, {}^{M}\boldsymbol{r}_i). \end{cases}$$
(8)



Fig. 2 Coordinate System of 3D-MoS Robot



Fig. 3 Coordinate System of Dual-eyes

2.2. Photo-Model Generation

In our system, three cameras are used as vision sensors. Among these three cameras, the first camera that is fixed in the workspace for capturing photos of clothes is used for cloth model generation. Captured photos are saved as BMP (bit map file) format. We use BMP rather than JPEG because of its full information of RGB. It can be converted into RGB and then HSV images, and also inversely. Comparison of image processing in different image format has not not done in this work. Note that generated model is not for just matching between saved ones and current images. Instead, generated model is used to estimate relative pose with respect to end effector using images from two cameras attached together with end effector. In proposed model generation technique, firstly, the background photo is taken as shown in Fig. 4 (a) and average hue value is calculated. Secondly, the cloth (target object) is put on the background and hue value of each point in the image is calculated as shown in Fig. 4 (b). Thirdly, the individual pixel of captured image is compared by scanning with the average of the hue value of the background image to define model frame based on error. Then, inner surface space of model S_{in} is generated by sampling hue value of each point inside defined frame. Finally, the outside space S_{out} of model is generated as shown in Fig. 4 (c).



Fig. 4 Model Generation Technique

2.3. Model-based Recognition

In geometrically pose estimation methods such as Epipolar Geometry, matching between interested features makes the system performance too much dependent especially on corresponding authenticity of points in the images of plural cameras. The wrong points pairs induce 3D pose estimation errors, resulting in corruption of visual servoing closed loop. To avoid this limitation, we have developed model-based recognition using dual-eye camera and generated photo-model. After generating a model from a bitmap image, the model is used for recognition the cloth (target object). Here, an overview of the recognition method with respect to the camera image is given as a description. 3D pose of the 3D model ${}^{C}\phi_{M} = [{}^{CR}x_{M}, {}^{CR}y_{M}, {}^{CR}z_{M}, {}^{CR}\epsilon_{1M}, {}^{CR}\epsilon_{2M}, {}^{CR}\epsilon_{3M}]^{T}$

is determined using model-based matching method. Generated models with different poses are projected



from 3D-model in searching area onto the left and right 2D images plane as shown in Fig.5. By comparing the projected models with images from two cameras attached at the end effector, relative pose is estimated by using fitness function $F({}^C\phi_M)$ to evaluate. It means the pose of the best model, that is fully matched with captured images from left and right cameras, is selected as estimated pose. The top of Fig. 5 is represented as searching area of a 3D-model named S to search for a cloth(target object). S_{in} is depicted by the space of coordinates on the surface of the model and S_{out} was enveloped the outside space of S_{in} . The left and right 2D searching model are named as S_L and S_R . In order to evaluate, the evaluation and change in hue of the surrounding of the object as shown in the interior region is represented as $S_{R,in}$, $S_{L,in}$ and the outside space enveloping $S_{R,in}$, $S_{L,in}$ is defined as $S_{R,out}$, $S_{L,out}$.

Theoretically, considering clothes to be flat shape, the similar performance can be achieved using only one camera. However, depth information from stereo vision was utilized in proposed system for future work in which depth information will be dominant. For example, the thickness of cloth No.12 in Fig.7 is about 6cm. The thickness varieties are ignored in this work.

2.4. Fitness Function

A function $P(^{IR}r_i)$ represents the matched degree of the i-th point of the model on the right image area, $^{IR}r_i$. Similarly, the left image area, $^{IL}r_i$ represented as a function $P(^{IL}r_i)$.

$$F({}^{C}\phi_{M}) = \left\{ \left(\sum_{\substack{IR_{r_{i}\in}\\S_{R,in}({}^{CR}\phi_{M})}} p({}^{IR}r_{i}) - \sum_{\substack{IR_{r_{i}\in}\\S_{R,out}({}^{CR}\phi_{M})}} p({}^{IR}r_{i}) \right) + \left(\sum_{\substack{IL_{r_{i}\in}\\S_{L,in}({}^{CL}\phi_{M})}} p({}^{IL}r_{i}) - \sum_{\substack{IL_{r_{i}\in}\\S_{L,out}({}^{CL}\phi_{M})}} p({}^{IL}r_{i}) \right) \right\} / 2 = \left\{ F_{R}({}^{CR}\phi_{M}) + F_{L}({}^{CL}\phi_{M}) \right\} / 2$$
(9)

As shown in Eq.(9), the fitness value will be increase

with the voting value of "+2" for every point of clothes in captured images that lies inside of the model frame $S_{R,in}$ and $S_{L,in}$. The fitness value will decrease with the value of "-0.005" for every point of clothes in images that lies in the background area $S_{R,out}$, $S_{L,out}$ and otherwise is "0". The correlation between the i-th point of the model (the search model) and the image having the evaluation value with such a sign is used the following Eq. (9). The whole evaluation function $F({}^{C}\phi_{M})$ is obtained by the average of the fitness function of both left camera image $F_{L}({}^{CL}\phi_{M})$ and right camera image $F_{R}({}^{CR}\phi_{M})$.

2.5. Genetic Algorithm (GA)

Recognition problem of the object can be converted to a searching problem of maximum value $F({}^{C}\phi_{M})$. There are various ways in finding the maximum value of the fitness function. The simplest and easiest way is the full search method. It is intended to find the maximum value by scanning all possible pixels. However, it has inefficient drawback in term of computing time. Even though there are powerful optimization methods, GA with long history is selected in this work because of its simplicity and effectiveness. By applying the GA evaluation process as an optimization solution, the maximum value search processing can be completed efficiently in a short period of time. In this experiment, GA has 60 individuals representing different poses of model. Each individual chromosome has six variables. Each variable are coded by 12 bits. The former 36 bits represent for the position coordinate of the 3D-model and the last 36bits represent for the orientation of the 3D-model. The characteristics of GA individual is defined as



These 60 chromosomes are evaluated by fitness value. Fitter ones are selected to regenerate next generations. Finally the best chromosome that has the most trustful pose is achieved as shown in Fig.6.



Fig. 6 GA Evolution Process

(No.1~No.12)



3. EXPERIMENTAL ENVIRONMENT

There are two units in experimental environment. One is for generation cloth model including one camera as shown in Fig. 8. Another one is end effector equiped with two cameras installed in manipulator's end effector as shown in Fig. 9 and Fig. 10. In Fig. 8, the distance from the camera lens to the model creating plane is 400 mm and the plane color is green. The size of clothes models can be up to 250mm \times 200mm. Each coordinate system of the robot and the cloth used in this experiments are shown in Fig. 9 and Fig. 10. The cloth coordinate system is represented as Σ_M and Σ_H is defined as the hand coordinate system of the robot – end effector –. Σ_M can be viewed from (x=0, y=0, z=685mm). It is centered on the recognition range of the position as a reference of the 510mm \times 390mm. After defining about the position, the recognition range of the angles are from 53° to -53° . The size of the collection box is a 220mm \times 220mm. However, in this experiment, we mainly emphasized to the recognition experiment and handling experiment is our follow-up work.



Fig. 8 Coordinate system of target object (unit is mm in Figure 8)



Fig. 9 Environment of model generation (unit is mm in Figure 9)



Fig. 10 Coordinate system of robot and end-effector (unit is mm in Figure 10)

4. EXPERIMENTAL CONTENT

Even though we conducted experiments using 12 different clothes as shown in Fig. 7, we will discuss in detail about experiments of cloth No.3 according to the following significant characteristics. The reason why we choose No.3 is based on three criterion;

- small size
- colorful pattern
- light weight

In this experiment, we have conducted for 1000 times recognition and analyzed in term of fitness distribution diagram, histogram of frequency distribution of each cloth and table of standard error and distribution. Fig. 11 shows the fitness function distribution in x-y plane of cloth No.3. We can see clearly that the maximum value of peak in fitness function distribution can be searched by proposed GA. Then, for cloth No.3, the number of frequencies (times) caused by position error (position x[mm] and position y[mm]) are shown in Fig. 12 and Fig. 13. Table 1 shows the average error and standard

deviation for cloth No.3. Finally, Fig. 14~16 show the pose estimation error for all 12 different clothes. The error average $\pm 3\sigma$ (standard deviation) is within a 10 mm for position (x-y) and the orientation for angle θ is 10°.



Fig. 11 Fitness function distribution in x-y plane from side view (No.3)

Table 1 Average error and standard deviation

	x[mm]	y[mm]	z[mm]	$\theta[^{\circ}]$
Average error	0.449	0.136	-3.46	0.379
Standard devia-	1.23	1.01	5.94	0.945
tion (σ)				

	x[mm]	y[mm]	z[mm]	$\theta[^{\circ}]$
Average error (-3σ)	-3.24	-2.89	-21.3	-2.46
Average error (-2σ)	-2.01	-1.88	-15.3	-1.51
Average error (-1σ)	-0.782	-0.874	9.39	-0.566
Average error	0.449	0.136	-3.46	0.379
Average error $(+1\sigma)$	1.68	1.15	2.48	1.32
Average error $(+2\sigma)$	2.91	2.16	8.24	2.27
Average error $(+3\sigma)$	4.14	3.17	14.4	3.21



Fig. 12 Histogram of position x[mm] error (No.3)



Fig. 13 Histogram of position y[mm] error (No.3)



(No.1~No.12)



(No.1~No.12)



5. CONCLUSION

In this paper, we introduced new cloth model generation method and pose estimation method using modelbased clothes recognition. The merits of the photomodel-based clothes recognition system are (1) photomodel-based allows the deformable different clothes recognized automatically, (2) 3D-pose measurement is possible, (3) photo-model-based 3D-pose recognition is not limited to clothing, and any object can recognized by using proposed system. We conducted 1000 times of experiment using 12 different clothes with different colorful pattern and multiple size. Recognition accuracy is analyzed in term of fitness distribution, histogram of pose estimation error. According to the experimental results, it can be confirmed that pose estimation of clothes for mass production can be implemented successfully using proposed photo-model-based clothes recognition system.

REFERENCES

- W. Song, M. Minami, Y. Mae, S. Aoyagi, On-line Evolutionary Head Pose Measurement by Feedforward Stereo Model Matching,*International Conference on Robotics and Automation (ICRA).*
- [2] J. Stavnitzky, D. Capson, Mutiple Camera Model-Based 3-D Visual Servoing, *IEEE Trans. on Robotics and Automation*, vol. 16, no. 6, December 2000.
- [3] 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.
- [4] Wei Song, Mamoru Minami, Fujia Yu, Yanan Zhang and Akira Yanou : "3-D Hand and Eye-Vergence Approaching Visual Servoing with Lyapunouv-Stable Pose Tracking", *IEEE International Conference on Robotics and Automation Shanghai International Conference Center*, May 9-13, 2011, Shanghai China.
- [5] S. Nakamura, K. Okumura, A. Yanou and M. Minami : "Visual Servoing of Patient Robot's Face and Eye-Looking Direction to Moving Human," *SICE Annual Conference* pp.1314-1319 2011.
- [6] Myo Myint, Kenta Yonemori, Akira Yanou, Shintaro Ishiyama and Mamoru Minami, "Robustness of Visual-Servo against Air Bubble Disturbance of Underwater Vehicle Using Three-dimensional Marker and Dual-eye Cameras", *International Conference* OCEANS 15 MTS/IEEE, Washington DC, USA, October 19-22, 2015.
- [7] Koichi Maeda, Mamoru Minami, Akira Yanou, Hiroaki Matsumoto, Fujia Yu, Sen Hou: "Frequency Response Experiments of 3-D Full Tracking Visual Servoing with Eye-Vergence Hand-Eye Robot System "SICE Annual Conference pp.101-107, 2012.
- [8] Y. Cui, X. Li, M. Minami, A. Yanou, M. Yamashita and S. ISHIYAMA, "Robot for Decontamination with 3D Move on Sensing", *Nuclear Safety and Simulation*, Vol. 6, No. 2, pp. 142-154, June 2015.
- [9] Yinxiao Li, Chih-Fan Chen, and Peter K.Allen, "Recognition of Deformable Object Category and Pose", *Robotics and Automation (ICRA), 2014 IEEE International Conference*, pp.5558-5564, May 31-June 7, 2014, Hong Kong.