

Pick and Place Robot Using Visual Feedback Control and Transfer Learning-Based CNN

Fusaomi Nagata, Kohei Miki, Akimasa Otsuka, Kazushi Yoshida

Graduate School of Science & Engineering,

Sanyo-Onoda City University

Sanyo-Onoda, Japan

nagata@rs.socu.ac.jp, f120613@ed.socu.ac.jp,

otsuka_a@rs.socu.ac.jp, kazushi.yoshida@rs.socu.ac.jp

Keigo Watanabe

Department of Intelligent Mechanical Systems

Division of Industrial Innovation Sciences

Graduate School of Natural Science and Technology

Okayama University

Okayama, Japan

watanabe@sys.okayama-u.ac.jp

Maki K. Habib

Mechanical Engineering Department

School of Sciences & Engineering

American University in Cairo

Cairo, Egypt

maki@aucegypt.edu

Abstract—Artificial neural network (ANN) which has four or more layers structure is called deep NN (DNN) and it is recognized as one of promising machine learning techniques. Convolutional neural network (CNN) is widely used and powerful structure for image recognition and/or defect inspection. It is also known that support vector machine (SVM) has a superior ability for binary classification in spite of only having two layers. The authors already have developed a CNN&SVM design and training tool for defect detection of resin molded articles, while the effectiveness and the validity have been proved through several CNNs design, training and evaluation. The tool further enables to facilitate the design of a CNN model based on transfer learning concept. In this paper, a pick and place robot is introduced while implementing a visual feedback control and a transfer learning-based CNN. The visual feedback control enables to omit the complicated calibration between image and robot coordinate systems, also the transfer learning-based CNN allows the robot to estimate the orientation of target objects for dexterous picking operation. The usefulness and validity of the system is confirmed through pick and place experiments using a small articulated robot named DOBOT.

Index Terms—convolutional neural network, transfer learning, pick and place, robot

I. INTRODUCTION

Artificial neural network (ANN) with four or more layers structure is called deep NN (DNN). The DNN is recognized as one of promising machine learning techniques. Convolutional neural network (CNN) has the most used and powerful structure for image recognition. It is also known that support vector machine (SVM) has a superior ability for binary classification in spite of only two layers. Nagi et al. proposed the max-pooling convolutional neural networks (MPCNN) for a vision-based hand gesture recognition problem [1]. It is reported that the MPCNN classified six kinds of different gestures with 96% accuracy and allowed mobile robots to perform real-time

gesture recognition. Weimer et al. also designed a deep CNN architectures for automated feature extraction in industrial defect inspection process [2]. The CNN automatically generates target features from massive amount of training image data and demonstrates successful defect detection results with lower false alarm rates. Faghih-Roohi et al. developed a different type of deep CNN for automatic rail surface defect detection [3]. It was concluded that the large CNN model yielded a better classification result than the small and medium CNNs, although the training process required a longer time. Zhou et al. designed and trained a CNN model to classify the surface defects of steel sheets [4]. The CNN could directly acquire better typical features from labeled images including surface defects. Further, Ferguson et al. presented a novel system to identify casting defects in X-ray images based on the Mask Region-based CNN architecture [5], [6]. It is reported that the proposed system simultaneously realized defect detection and segmentation on input test images making it suitable for a range of defect detection tasks. The authors have also developed a CNN&SVM design and training tool for defect detection of resin molded articles and the usefulness and validity have been proved through several design, training and evaluation experiment of CNNs and SVMs [7–9]. The tool further enables to easily design powerful CNN models based on transfer learning concept.

Taryudi and Wang presented an application of stereo vision system to estimate the position and orientation of objects, so that a robot arm equipped with a gripper could well pick up the objects and place them at desired positions in the workspace [10]. When industrial robots are applied to pick and place tasks of resin molded articles, information of each object's position and orientation is essential. Recognition and extraction of the object position in an image are not so difficult if some image

processing technique is used, however, that of orientation is not easy due to the variety in shape. In this paper, a pick and place robot is introduced while implementing a visual feedback control and a transfer learning-based CNN. The visual feedback control enables to omit the complicated calibration between image and robot coordinate systems, also the transfer learning-based CNN allows the robot to estimate the orientation of target objects. The effectiveness and validity of the system is demonstrated through pick and place experiments using an small articulated robot named DOBOT.

II. DESIGN & TRAINING TOOL FOR CNN AND SVM

The authors have already proposed a software to be able to easily design CNNs and SVMs for defect detection. In training of CNN, pre-training using randomly initialize weights and additional (successive) training with once trained weights can be selected. As for SVM, one-class unsupervised learning and two class supervised learning can be selectively executed. Also, favorite CNN, which is used for a feature extractor, and Kernel function are selected. The tool has another promising function to design original CNNs based on transfer learning. Figure 1 shows the part in main dialogue developed for efficiently designing and training transfer learning-based CNNs. For example, the following main items can be set for the operation of transfer learning through the dialogue.

- Folders for training and test images.
- Base CNNs used for transfer learning such as AlexNet, VGG16, VGG19, GoogleNet, Inception-V3 and IncResNetV2.
- Training parameters such as mini batch size, desired accuracy and loss, learning rates for convolution layers and fully connected layers, max epochs, and so on.

The software shown in Fig. 1 is developed on MATLAB system optionally installed with Neural Network Toolbox, Parallel Computing Toolbox for GPU, Deep Learning Toolbox, Statistics and Machine Learning Toolbox.

III. IMAGES FOR TRAINING AND TEST

Training image generator was already proposed to efficiently augment limited number of training images [7]. By using the generator, images for training are prepared considering typical twelve orientations, i.e., 0° , 15° , 30° , 45° , 60° , 75° , 90° , 105° , 120° , 135° , 150° and 165° . Figures 2 and 3 show examples of the training images for the categories of 45° and 165° , respectively. The resolution and channel are 200×200 and 1, respectively.

IV. TRANSFER LEARNING-BASED CNN

A. Design and Training

In this section, a transfer learning-based CNN is designed to learn the feature of orientation included in images as shown in Figs. 2 and 3. Figure 4 illustrates the structure of the original AlexNet consisting of 25 layers, which can classify input images into one of 1,000 categories. In order to make the CNN have an ability to classify input images into 12 categories as 0° , 15° , 30° , 45° , 60° , 75° , 90° , 105° , 120° , 135° , 150° and

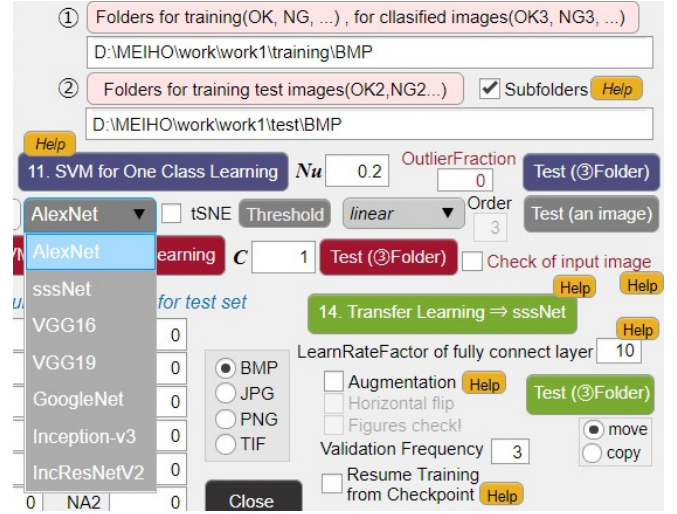


Fig. 1. A part of the main dialogue developed for efficiently designing and training transfer learning-based CNNs.

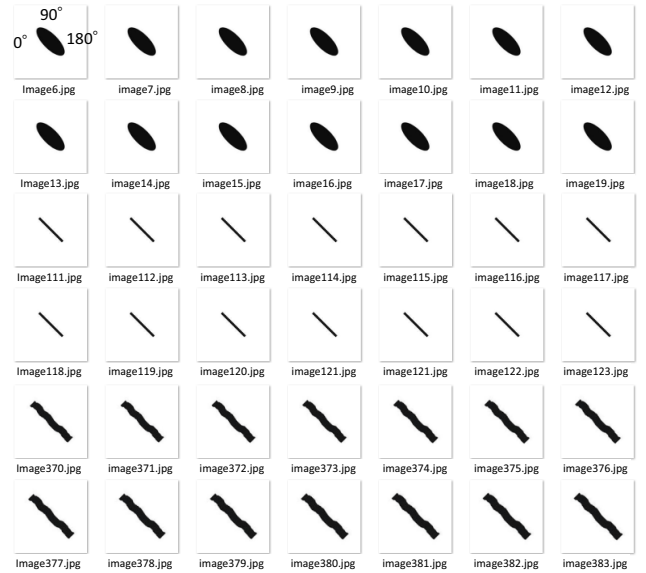


Fig. 2. Examples of training images for the orientation of 45° .

165° , the fully connected layers are replaced as shown in Fig. 5 before executing transfer learning. 6889 images consisting of 12 categories are used for the transfer learning. As for training parameters, mini batch size is given 50. Iteration is the number of mini batches needed to complete one epoch, so that one epoch in this transfer learning is composed of $6889/50=137$ iterations. Desired accuracy and loss are set to 1 and 0, respectively. Besides, learning rates of convolutional layers and fully connected ones are set to 0.0001 and 0.001, respectively. It is important for fast and stable convergence in transfer learning to set the learning rate in convolutional layers smaller than that of fully connected layers.

If n th training image is given to the input layer of the transferred CNN, then the softmax layer yields the probability

Well-known CNN named *AlexNet* trained for classification of 1000 categories

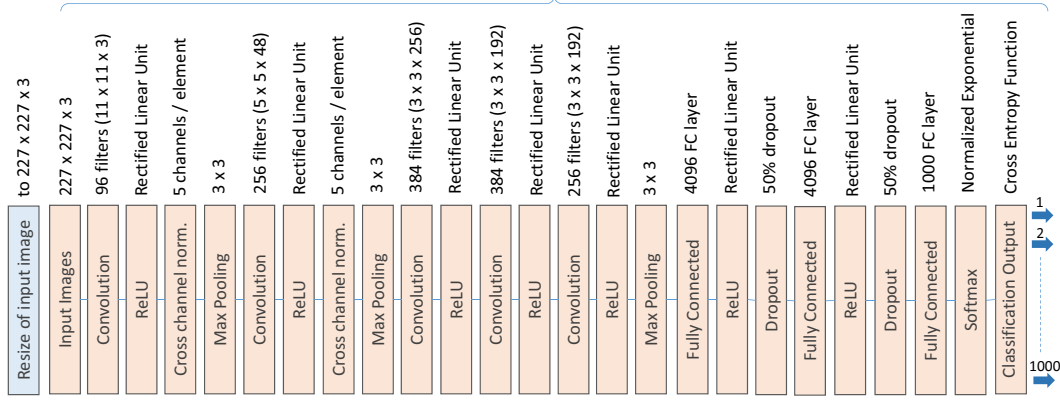


Fig. 4. Network structure of well-known CNN named AlexNet which can classify input images into one of 1000 kinds of categories.

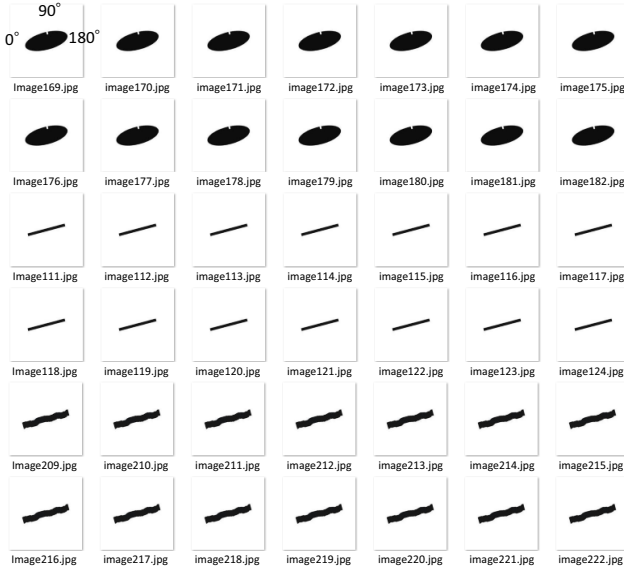


Fig. 3. Examples of training images for the orientation of 165°.

p_{ni} ($i = 1, 2, \dots, 12$) as the score for twelve types of categories, which is written by

$$p_{ni} = \frac{e^{y_{ni}}}{\sum_{k=1}^{12} e^{y_{nk}}} \quad (1)$$

where $\mathbf{y}_n = [y_{n1} \ y_{n2} \ \dots \ y_{n12}]^T$ is the output from the last fully connected layer corresponding to the n th image. The transferred CNN is trained based on the back propagation algorithm using the loss function called cross entropy given by

$$\bar{E} = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^{12} t_{nk} \log(y_{nk}) \quad (2)$$

where $\mathbf{t}_n = [t_{n1} \ t_{n2} \ \dots \ t_{n12}]^T$ denotes the n th desired output for classification, i.e., only one element in \mathbf{t}_n has 1, remained

elements have 0. N is the total number of samples in the training image data set.

The training to optimize the transfer learning-based CNN was processed using a single PC with a Core i7 CPU and a GPU (NVIDIA GeForce GTX 1060, 6GB). The training progressed as shown in Fig. 6, in which both the training accuracy and loss seem to well converge to desired values. Note that the accuracy A_c at each iteration is given by

$$A_c = \frac{50 - N_m}{50} \times 100 \quad (3)$$

where 50 is the mini batch size, i.e., the number of sampled images processed during one iteration, N_m is the number of mis-classified images within the 50 sampled images. It actually took about 40 minutes until the learning was stopped since both the accuracy and loss had not been improved during 10 consecutive iterations or more. Note that this training could be completed within one epoch by severally giving different learning rates in convolutional layers and fully connected layers. Through the process explained above, an original CNN model acquired by transfer learning of AlexNet, which is the winner of ImageNet LSVRC2012, is presented to recognize the orientation of objects.

B. Generalization Ability

After the training, the generalization ability of the transfer learning-based CNN was checked using 15 test images including objects imitating resin molded articles which had not been in the training data set. Figure 7 shows the classified photos and their results, i.e., the angles written in the JPEG images are the outputs from the CNN. It is observed from the results that the obtained CNN has a promising generalization ability that can recognize the orientations of objects in the images. However, some visual inconsistencies, e.g., between “test7.jpg” (75°) and “test12.jpg” (60°); “test.jpg” (150°) and “test3.jpg” (150°) are observed. As can be clearly seen, some images in Fig. 7 are not complete square. That is the reason why the main cause of these results seems to be the

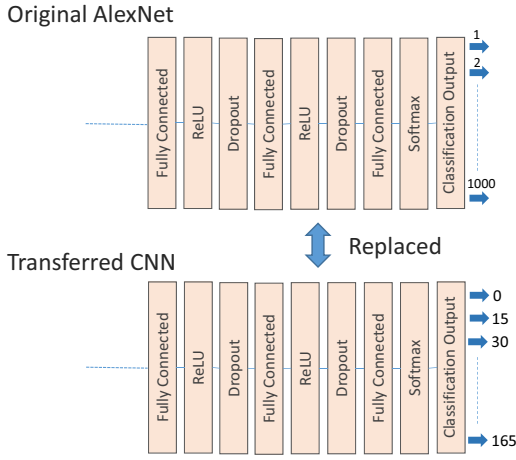


Fig. 5. Replacement of fully connected layers for dealing with target classification task, i.e., 12 categories.

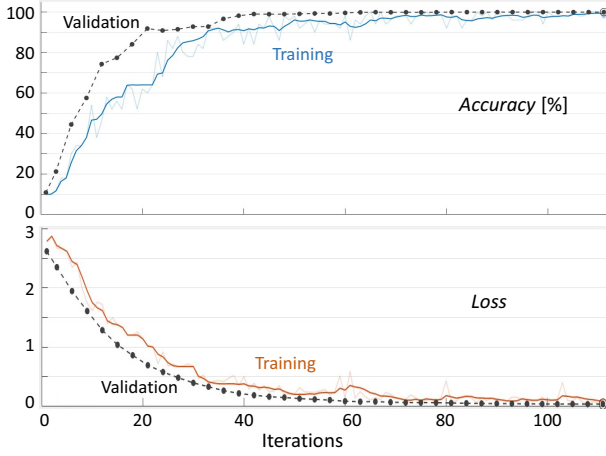


Fig. 6. Training progress of transfer learning shown in Fig. 5.

conversion of resolution before classification. The resolution of images given to the input layer is forced to be converted to $227 \times 227 \times 3$ specified by the input layer of the AlexNet, which brings out some undesirable deformation of images and the resultant ambiguities in classification.

Figure 8 shows another classification result of different types of test images. As can be seen, although shapes in the images are quit different from those in Fig. 7, desirable generalization ability can be observed.

V. EXPERIMENT OF PICK AND PLACE

A. Without Visual Feedback Control

An actual pick and place experiment is conducted using a small articulated robot named DOBOT. The experimental setup is shown in Fig. 9. Position $[x \ y \ z]^T$ and yaw angle R of the gripper in robot coordinate system can be controlled by an API function $\text{SetPTPCmd}(x_d, y_d, z_d, R_d)$. Note that the yaw angle R is dealt with the orientation of a target object in this experiment. Figure 10 shows the developed control

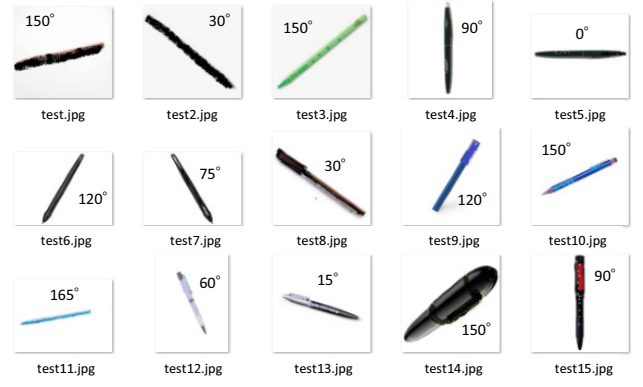


Fig. 7. Classification results of test images using the transfer learning-based CNN.



Fig. 8. Another classification results of different shapes of test images using the transfer learning-based CNN.

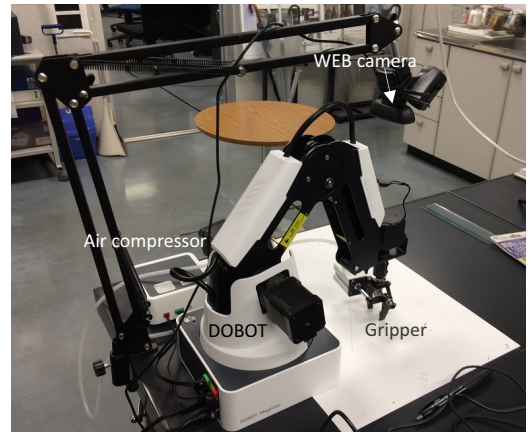


Fig. 9. Experimental setup based on an articulated robot.

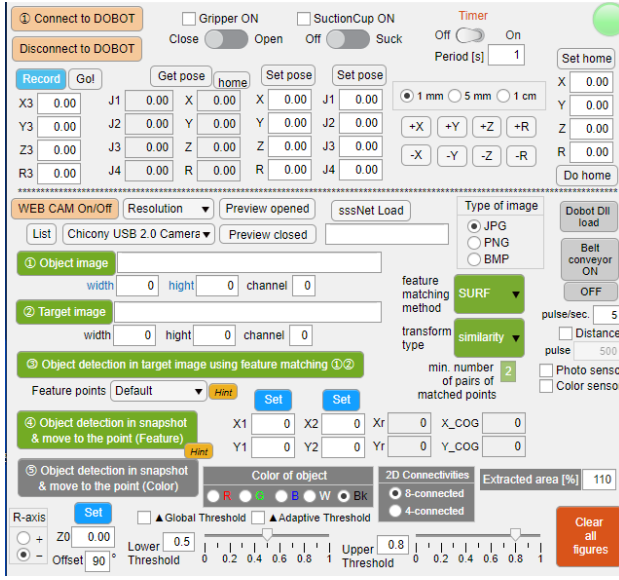


Fig. 10. Developed control dialogue for DOBOT.

dialogue for the robot. Pick and place task while recognizing the orientations of target objects can be executed through the dialogue. Figure 11 illustrates the flowchart of the pick and place task, which is implemented in a timer interrupt routine.

In the timer interrupt, first of all, a snapshot of 1600×1200 resolution is captured. After binarized into black and white, a connected component with the largest area is found as a target object and then the COG position $[I_x \ I_y]^T$ ($1 \leq I_x \leq 1600, 1 \leq I_y \leq 1200$) in image coordinate system is extracted by

$$I_x = \frac{\sum_{x=1}^{1600} \sum_{y=1}^{1200} x f(x, y)}{S} \quad (4)$$

$$I_y = \frac{\sum_{x=1}^{1600} \sum_{y=1}^{1200} y f(x, y)}{S} \quad (5)$$

where the position (x, y) are the variables of column and row in the image. Also, $f(x, y)$ is the binary value, i.e., 1 or 0, at the coordinate (x, y) . S is the area of the identified and extracted object, which is obtained by calculating the total number of pixels with a value of 1 forming the object. Consequently, desired position $[x_d \ y_d]^T$ in robot coordinate system to move the gripper to the COG position can be obtained by

$$x_d = X_1 + I_x \frac{X_2 - X_1}{1600} \quad (6)$$

$$y_d = Y_1 + I_y \frac{Y_2 - Y_1}{1200} \quad (7)$$

where $[X_1 \ Y_1]^T$ and $[X_2 \ Y_2]^T$ in robot coordinate system are the positions of left upper and right bottom of the snapshot as shown in Fig. 12, i.e., they are corresponding to pixels of $(0, 0)$ and $(1600, 1200)$, respectively. The part of the connected component is further cropped centering the COG from the

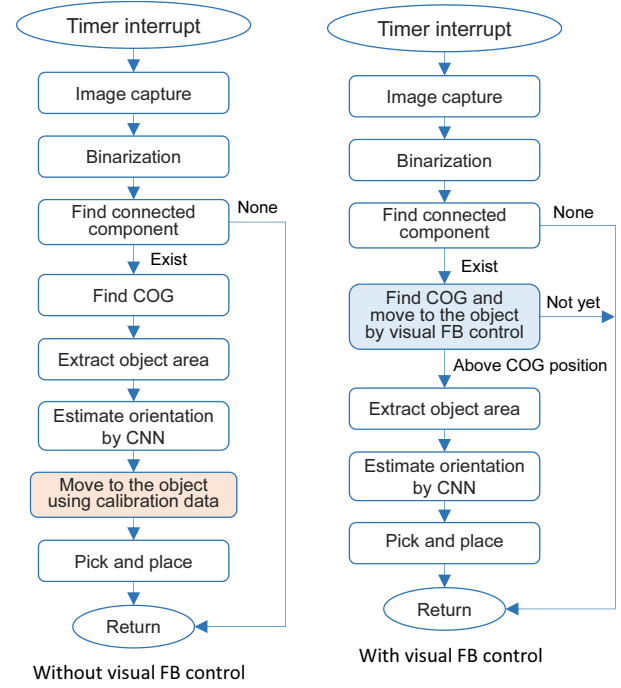


Fig. 11. Flowcharts to realize pick and place task using transfer learning-based CNN, in which left and right figures show the process flow diagrams without and with a visual feedback control, respectively.

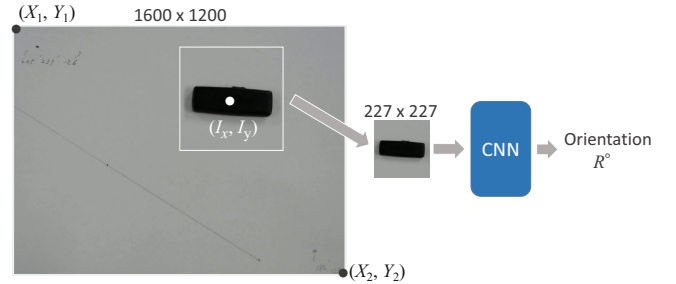


Fig. 12. Procedure to extract the orientation of a workpiece from a captured image.

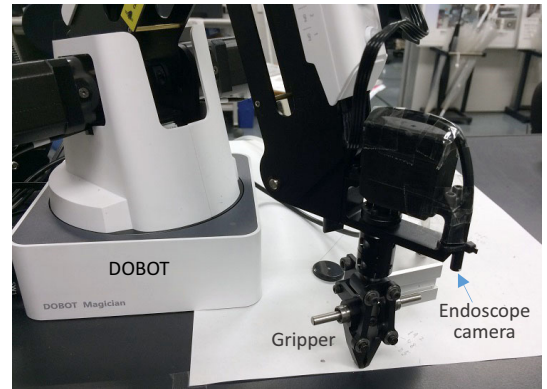


Fig. 13. Experimental setup based on an articulated robot with an endoscope camera.

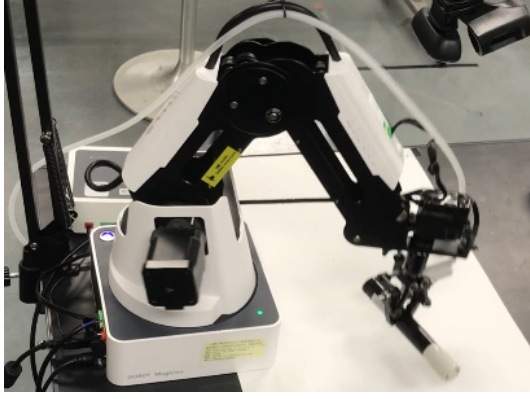


Fig. 14. Experimental scene just after picking.

original snapshot as shown in Fig. 12. The cropped image is resized into 227×227 and given to the input layer of the transfer learning-based CNN designed in the previous section. Finally, the orientation of the object can be estimated by the CNN, which is used for the desired yaw angle R_d so that the robot can successfully grasp the object with a long-axis shape.

Figure 14 shows the scene just after picking, in which the gripper cloud well identify the COG position of the black part and the orientation, then pick the point up.

B. With Visual Feedback Control

Camera configuration in visual feedback control is not required, which is different from that in the previous subsection. A lightweight endoscope camera is attached close to the gripper as shown in Fig. 13. Manipulated variable $v(k) = [v_x(k) \ v_y(k)]^T$ for visual feedback is generated by a PI-action given by

$$v(k) = K_p e(k) + K_i \sum_{n=1}^k e(n) \quad (8)$$

where k is the discrete time. K_p and K_i are the gains for proportional and integral actions, respectively. $e(k) = [e_x(k) \ e_y(k)]^T$ is the error vector in image coordinate system measured by

$$e(k) = X_d - I(k) \quad (9)$$

where $X_d = [\frac{X_r}{2} \ \frac{Y_r}{2}]^T$ and $I(k) = [I_x(k) \ I_y(k)]^T$ are the desired position and the measured object's COG position in image coordinate system, respectively. $[X_r \ Y_r]$ is the resolution of captured image, so that the desired position $X_d = [\frac{X_r}{2} \ \frac{Y_r}{2}]^T$ is the center position of the image.

It was confirmed from a similar experiment as shown in Fig. 14 that the visual feedback control allows the robot to execute the pick and place task without setting $[X_1 \ Y_1]^T$, $[X_2 \ Y_2]^T$ and without using Eqs. (6) and (7).

VI. CONCLUSIONS

In this paper, a CNN acquired by transfer learning of AlexNet, which had been trained using about 1.2 million

images of 1000 categories in ImageNet database, was introduced to recognize the orientations of objects. Originally, the AlexNet had been able to classify input images into one of 1,000 kinds of objects, however, the transferred CNN has been able to recognize the orientation of an object in images with 12 kinds of degrees. Then, a visual feedback control has been implemented so that the gripper of the pick and place robot can move to the position nearly just above a target object. Due to the visual feedback control, the complicated calibration between image and robot coordinate systems can be omitted. The effectiveness of the system is evaluated through experimental pick and place tests using an articulated robot named DOBOT.

ACKNOWLEDGMENT

This work was partially supported by Mitsubishi Pencil Co., Ltd. and Meiho Co., Ltd.

REFERENCES

- [1] J. Nagi, F. Ducatelle, G.A.D. Caro, D. Ciresan, U. Meier, A. Giusti, F. Nagi, J. Schmidhuber, and L.M. Gambardella, "Max-pooling convolutional neural networks for vision-based hand gesture recognition," *2011 IEEE International Conference on Signal and Image Processing Applications (ICSIPA2011)*, pp. 342–347, 2011.
- [2] D. Weimer, B. Scholz-Reiter, and M. Shpitalni, "Design of deep convolutional neural network architectures for automated feature extraction in industrial inspection," *CIRP Annals – Manufacturing Technology*, Vol. 65, No. 1, pp. 417–420, 2016.
- [3] S. Faghih-Roohi, S. Hajizadeh, A. Nunez, R. Babuska, and B.D. Schutter, "Deep convolutional neural networks for detection of rail surface defects," *Procs. of the 2016 International Joint Conference on Neural Networks (IJCNN2016)*, Vancouver, Canada, pp. 2584–2589, 2016.
- [4] S. Zhou, Y. Chen, D. Zhang, J. Xie, and Y. Zhou, "Classification of surface defects on steel sheet using convolutional neural networks," *Materials and Technology*, Vol. 51, No. 1, pp. 123–131, 2017.
- [5] M. Ferguson, R. Ak, Y. Lee, and K. Law, "Detection and segmentation of manufacturing defects with convolutional neural networks and transfer learning," *Smart and Sustainable Manufacturing Systems*, Vol. 2, No. 1, pp. 137–164, 2018.
- [6] K. He, G. Gkioxari, P. Dollr, and R. Girshick, "Mask R-CNN," *Procs. of 2017 IEEE International Conference on Computer Vision (ICCV)*, pp. 2980–2988, 2017.
- [7] F. Nagata, K. Tokuno, H. Tamano, H. Nakamura, M. Tamura, K. Kato, A. Otsuka, T. Ikeda, K. Watanabe, and M.K. Habib, "Basic application of deep convolutional neural network to visual inspection," *Procs. of International Conference on Industrial Application Engineering (ICIAE2018)*, pp. 4–8, Okinawa, 2018.
- [8] F. Nagata, K. Tokuno, K. Nakashima, A. Otsuka, T. Ikeda, H. Ochi, K. Watanabe, and M. K. Habib, "Fusion method of convolutional neural network and support vector machine for high accuracy anomaly detection," *Procs. of the 2019 IEEE International Conference on Mechatronics and Automation (ICMA 2019)*, pp. 970–975, Tianjin, China, 2019.
- [9] F. Nagata, K. Tokuno, K. Mitarai, A. Otsuka, T. Ikeda, H. Ochi, K. Watanabe, and M. K. Habib, "Defect detection method using deep convolutional neural network, support vector machine and template matching techniques," *Artificial Life and Robotics*, Vol. 24, No. 4, pp. 512–519, 2019.
- [10] Taryudi and M.S. Wang, "3D object pose estimation using stereo vision for object manipulation system," *Procs. of 2017 International Conference on Applied System Innovation (ICASI)*, pp. 1532–1535, 2017.