

Contact-Rich Manipulation of a Flexible Object based on Deep Predictive Learning using Vision and Tactility

Hideyuki Ichiwara¹, Hiroshi Ito^{1,2}, Kenjiro Yamamoto¹, Hiroki Mori² and Tetsuya Ogata²

Abstract—We achieved contact-rich flexible object manipulation, which was difficult to control with vision alone. In the unzipping task we chose as a validation task, the gripper grasps the puller, which hides the bag state such as the direction and amount of deformation behind it, making it difficult to obtain information to perform the task by vision alone. Additionally, the flexible fabric bag state constantly changes during operation, so the robot needs to dynamically respond to the change. However, the appropriate robot behavior for all bag states is difficult to prepare in advance. To solve this problem, we developed a model that can perform contact-rich flexible object manipulation by real-time prediction of vision with tactility. We introduced a point-based attention mechanism for extracting image features, softmax transformation for predicting motions, and convolutional neural network for extracting tactile features. The results of experiments using a real robot arm revealed that our method can realize motions responding to the deformation of the bag while reducing the load on the zipper. Furthermore, using tactility improved the success rate from 56.7% to 93.3% compared with vision alone, demonstrating the effectiveness and high performance of our method.

I. INTRODUCTION

People can easily perform actions to manipulate flexible objects, such as unzipping a bag, peeling a banana, or changing clothes. However, it is very challenging to make robots manipulate these flexible objects. Flexible object manipulations achieved by robots so far include folding clothes, tying ropes, and folding paper [1][2][3][4]. For these tasks, it is important to predict how the object will be deformed by the robot’s actions, and these tasks have been accomplished mainly using vision. However, people perform many tasks on a daily basis that cannot be accomplished using vision alone. This is because our vision is blocked by our own hands and objects when we work with objects. Robot need to perform such tasks to expand their application range. In this research, we aim to achieve contact-rich manipulation of flexible objects, which is difficult to do with vision alone, using a two-fingered robot equipped with two tactile sensors as a means to supplement vision. We focus on the task of unzipping a fabric bag, which has common problems with these tasks, so as not to lose its generality as a validation task. Common issues are that occlusion occurs with one’s hands and fingers and that it is difficult to predict the state of the target object.

¹Hideyuki Ichiwara, Hiroshi Ito and Kenjiro Yamamoto are with Center for Technology Innovation - Mechanical Engineering, Research & Development Group, Hitachi, Ltd., Ibaraki, 312-0034, Japan hideyuki.ichiwara.bn@hitachi.com

²Hiroshi Ito, Hiroki Mori and Tetsuya Ogata are with Department of Intermedia Art and Science School of Fundamental Science and Engineering, Waseda University, Tokyo, 169-855, Japan ogata@waseda.jp

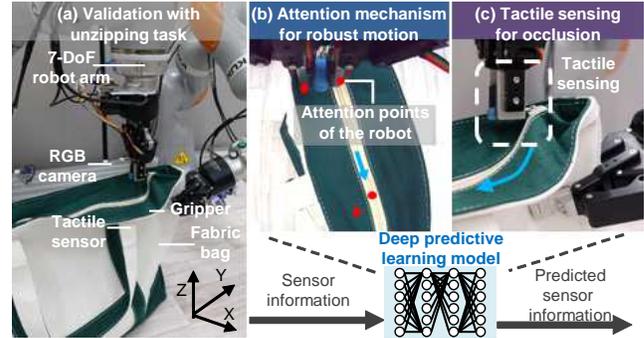


Fig. 1. Overview of our study. (a) The experimental setup for flexible fabric bag unzipping. (b) Obtain attention points of the robot using point-based attention mechanism. (c) Motion generation based on two tactile sensors attached to the gripper of the robot.

One methods to manipulate flexible objects is the physical model-based approach. In this approach, a physical model of the flexible object is built, and the object is manipulated by predicting how the object will be deformed by the robot [5][6][7][8]. However, the dynamics is complex, and these models are very expensive to build. Moreover, even when the models are constructed, ideal trajectories are difficult to define for abstract goals such as “unzipping”.

Therefore, methods based on deep learning with vision have been investigated for flexible object manipulation. Deep predictive learning, which learns to simultaneously predict motions and images to generate appropriate robot motions while predicting the behavior of flexible objects [4][9][10], and reinforcement learning-based methods using real robots [11][12] have been proposed. These methods are promising for contact-rich flexible object manipulation because they do not require a physical model of the target flexible object since they learn from the robot’s motion data. However, since these methods generate motions based mainly on images, it was difficult to deal with background changes and occlusions that were not present during training.

General-purpose object manipulation methods using deep learning have also been proposed, which are supervised learning-based methods with high sampling efficiency using attention mechanisms [13][14]. Such methods are expected to enable flexible object manipulation without the need for a large amount of training data, by paying attention to the part of the flexible object’s state with complex dynamics that is important for the task. In addition, Ichiwara et al. [14] showed that narrowing down the image regions that are important for the task and obtaining attention based on real-

time prediction of images during the task makes the system robust to background changes and other factors.

Moreover, various types of tactile sensors have been developed to improve the manipulation capability of robots [15][16][17][18]. However, most research using these tactile sensors has been done on manipulation by robot hands, including object and material recognition, posture estimation during object grasping, and in-hand manipulation [19][20][21]. On the other hand, there are few examples of applications to systems including a robot arm, and they have been applied to simple tasks using reinforcement learning or verified by simulation [22][23][24].

In this work, we propose a deep predictive learning-based robot control system that performs contact-rich flexible object manipulation, which is difficult to achieve with vision alone. The contributions of this research are as follows: (1) To make the system robust to dynamic visual changes and occlusions caused by flexible object manipulation, we proposed an attention mechanism to narrow down the important regions of the image. (2) Using tactile sensors to compensate for occlusion, we confirmed that the success rate of the unzipping task increased from 56.7 to 93.3 % compared to the case using only vision, demonstrating the effectiveness of tactility in action generation. (3) We proposed a method that integrates vision, tactility, and motion to generate motions using a deep predictive learning model based on real-time state prediction.

II. RELATED WORK

A. Deep Learning-Based Motion Learning

Deep learning is a powerful machine learning method that has been successfully applied in a wide range of fields such as object recognition and language processing [25][26], and is also being applied in robot behavior generation.

Yang et al. [9] used teleoperation to teach a robot to perform a motion and enabled the robot to fold flexible cloth, which was difficult to achieve in the past. To generate an appropriate robot motion while predicting the behavior of a flexible object, deep predictive learning has been proposed, in which learning is performed to predict motion and an image at the same time based on training data. Suzuki et al. [4] enabled robots to manipulate ropes by two arms on the bases of the method of Yang et al. [9]. In addition to the image and the joint angle of the robot, the proximity sensor was used to recognize the rope state. However, since the proximity sensor was used to recognize whether or not there is a rope in the robot hand, high-dimensional contact information such as tactility was not used.

In addition, sampling efficient methods [13][14][27] have been proposed using the attention mechanism. Levine et al. [13][27] proposed using CNN and soft argmax to obtain the position coordinates in the image for image feature extraction. Since positional information is important for the robot's task, the sampling efficiency was improved by compressing the information-rich image to a lower order. Ichiwara et al. [14] proposed a new method based on the method of Levine et al. [13][27] that was made robust to

changes in background and object positions by predicting attention points on the basis of image prediction. By adding top-down attention based on prediction in addition to bottom-up attention from images, it is possible to extract only the location information important for the task. These methods using the attention mechanism have been validated for tasks such as object grasping and pick-and-place, where the object has a fixed shape. On the other hand, they have not been validated for flexible object manipulation, which involves dynamic vision changes due to object manipulation.

B. Manipulation using Tactility

Much of the research on manipulation using tactility has been done on in-hand manipulation with robot hands [19][20][21]. As for the application to systems that include a robot arm, a robot manipulation method based on deep reinforcement learning using tactility has been proposed, and using tactility has been shown to improve search efficiency [24]. However, the effectiveness of this method has not yet been verified on a real robot. In research using real robots, a reinforcement learning-based method using vision and tactility has been proposed, and using tactility has been shown to enable stable learning [23]. This method has been validated with a simple task such as using a five-axis robot to control a single finger attached to a tactile sensor to make a pole rotating in two axes stand upright.

In this study, we propose a deep predictive learning model that integrates and predicts vision, tactility, and motion, using attention mechanisms. Furthermore, as an example of contact-rich flexible object manipulation, we test the effectiveness of our model on a real robot for the unzipping task of a fabric bag. Finally, we show that the tactility is important in the task.

III. METHOD

Fig. 2 shows the proposed model for contact-rich flexible object manipulation using tactility. The network parameters in this study are shown in Fig. 2. For example, the explanation in CNN and Transposed CNN shows (channels, kernel size, strides), and the explanation in FCNN and LSTM shows (number of nodes). The camera image i_t , tactility data p_t , and joint angle a_t of the robot at time t are input to the model, which predicts the camera image \hat{i}_{t+1} , tactility data \hat{p}_{t+1} , and joint angle $\hat{S}(a_{t+1})$ after softmax transformation (SMT) at the next time $t + 1$. Here, it is important to simultaneously predict sensory inputs (images, tactility) other than motion. If the sensory inputs are not predicted, the sensory input are not referenced and only the time t and $t + 1$ motions are learned, and the appropriate motion is not generated [14].

To train the model, we used the time series data when the robot executes the action by human operation as the training data. The model was then trained to output the data at the next time $t + 1$ from the data at time t . When the robot executes a motion using the learned model, each data at the current time is input to the model, and the predicted joint angle data is input to the robot as commands. Our model has

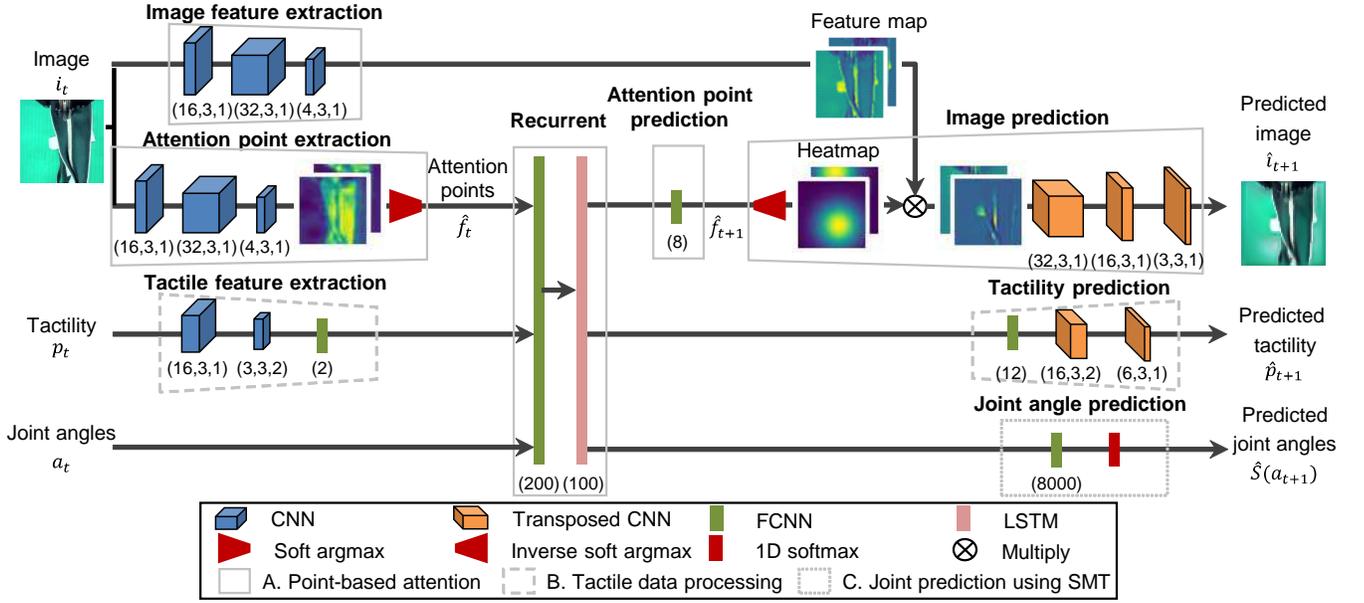


Fig. 2. The proposed model architecture for contact-rich flexible object manipulation using tactility. It consists of point-based attention, which extracts image features and position information, tactility data processing, and joint prediction.

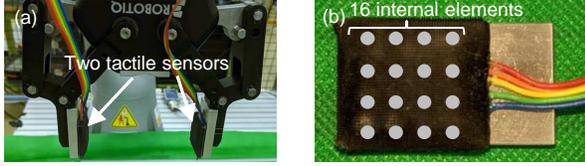


Fig. 3. (a) Two-finger gripper with two tactile sensors. (b) Tactile sensor that has 16 built-in elements.

the following features. To deal with dynamic visual changes, we used the point-based attention mechanism on the basis of prediction, as shown by the solid line in Fig. 2. To make it robust to tactility position shifts caused by shifting of the puller, we processed the tactility data using CNNs, as shown by the dashed lines. To generate finer movements of the robot, SMT was introduced, which is a method of expressing joint angles in a sparse form, as shown by the dashed lines.

A. Point-Based Attention Mechanism

The key components of the attention mechanism are the blocks shown by the solid gray lines in Fig. 2. The attention point extraction block using CNNs and soft argmax outputs the positional coordinates of points in the image (called the attention point in this paper) from the input image. Here, soft argmax obtains the index that takes the maximum value of the array and relaxes the conditions of the argmax function, which has difficulty back-propagating errors [28]. In this block, the feature maps of the input image are obtained using CNNs. In addition, soft argmax is used to extract high-intensity positions in the feature map and treat them as attention points. Next, the recurrent block, attention points prediction block, and joint angle prediction block using the FC layer and LSTM layer predict the attention points, joint

angle and tactility feature at the next time $t + 1$ from those data at time t . In addition, the image prediction block predicts the image at the next time $t + 1$ on the basis of the prediction attention points and the feature maps obtained from the image feature extraction block. Specifically, it generates heat maps with high intensity at the prediction points of interest and predicts the image at the next time $t + 1$ on the basis of the feature maps weighted by the heat maps. This is because the heat map alone is not sufficient to predict the image, but the image can be easily predicted by using the information near the attention points among the image features obtained from the image feature extraction block.

Note that since all blocks are end-to-end connected and the attention points are determined by a self-organized manner, the training data of the attention points is not necessary. The points are obtained in such a way as to minimize the prediction error between the image and the joint angle, and they are obtained at positions important to the task. The predicted attention points output by attention point prediction block are called predicted attention points because they are obtained to minimize the prediction error of the image.

B. Tactility Data Processing using CNNs

In this study, we used the tactile sensor (XR1944 from XELA Robotics) attached to the gripper shown in Fig. 3. This sensor has 16 elements inside and uses magnetism. For each contact point, the digital values in the three axial directions of the pressure direction and the shear direction can be detected, and one sensor outputs $4 \times 4 \times 3$ values. Shear force and normal force resolution are 0.1 gf and 1 gf. Therefore, the output of the sensor can be treated in the same way as the image, and as shown by the broken line in Fig. 2, CNNs were used for extracting tactile features. In image processing,

CNN usage is robust to positional shifts [29]. In this study, we used CNNs for tactile feature extraction to make our system robust to positional shifts of the puller. In addition, for some tasks, it may be better to have a higher sensitivity to position shift. In this case, we can increase the sensitivity by reducing the CNN stride, increasing the number of tactile dimension input to the recurrent block, etc.

C. Joint Prediction using Softmax Transformation(SMT)

To enhance the learning, SMT [30] was used to predict joint angles. SMT converts joint angles into sparse representations instead of representing them as continuous values such as radians. For example, one joint angle a , which can take values from 0 to 1, is transformed using J neurons:

$$n_j^a = \frac{\exp\left(\frac{-\|\frac{j}{J}-a\|^2}{\sigma}\right)}{\sum_{j \in J} \exp\left(\frac{-\|\frac{j}{J}-a\|^2}{\sigma}\right)}, \quad (1)$$

where n_j^a represents j th neuron, and σ is a parameter that determines the distribution after the transformation. In our experiment, J was set to 1000, and σ was set to 0.1. On the other hand, the joint angle prediction block shown in Fig. 2 predicts the joint angle after SMT through a 1-dimensional (1D) softmax layer represented by J neurons per joint angle. When testing after training, the output of the model represented by SMT was converted into the joint commands of the robot:

$$a = \sum_{j \in J} n_j^a \times \frac{j}{J}. \quad (2)$$

In this way, by using multidimensional neurons for each joint, finer movements can be generated.

D. Loss Function

The loss function was defined as follows:

$$g = \sum_{t \in T} \{g_i(\hat{\mathbf{i}}_{t+1}, \hat{\mathbf{i}}_{t+1}) + g_p(\hat{\mathbf{p}}_{t+1}, \hat{\mathbf{p}}_{t+1}) + g_a(S(\mathbf{a}_{t+1}), \hat{S}(\mathbf{a}_{t+1})) + \alpha g_f(\hat{\mathbf{f}}_t, \hat{\mathbf{f}}_{t+1})\}, \quad (3)$$

$$g_i = \frac{1}{H_i \times W_i \times C_i} \|\hat{\mathbf{i}}_{t+1} - \mathbf{i}_{t+1}\|_2^2, \quad (4)$$

$$g_p = \frac{1}{H_p \times W_p \times C_p} \|\hat{\mathbf{p}}_{t+1} - \mathbf{p}_{t+1}\|_2^2, \quad (5)$$

$$g_a = \frac{1}{M \times J} \sum_{l \in M} \sum_{m \in J} S(\mathbf{a}_{t+1})_l^m \log(\hat{S}(\mathbf{a}_{t+1})_l^m), \quad (6)$$

$$g_f = \frac{1}{K} \|\hat{\mathbf{f}}_t - \hat{\mathbf{f}}_{t+1}\|_2^2, \quad (7)$$

The sequence length of the training data is T , the time is t , the image is $\mathbf{i} \in \mathbb{R}^{H_i \times W_i \times C_i}$, the tactility is $\mathbf{p} \in \mathbb{R}^{H_p \times W_p \times C_p}$, the joint angle $\mathbf{a} \in \mathbb{R}^M$ after SMT is $S(\mathbf{a}) \in \mathbb{R}^{M \times J}$, and the coordinates of the attention points are $\mathbf{f} \in \mathbb{R}^K$. g_i and g_p are prediction errors of the image and tactility using the mean square error. g_a is the prediction error of the joint angle after SMT using cross-entropy loss. Also, g_f denotes the error in the Euclidean distance between the attention points $\hat{\mathbf{f}}_t$ of the encoder output and the attention points $\hat{\mathbf{f}}_{t+1}$ of the decoder input. α denotes the weight of the

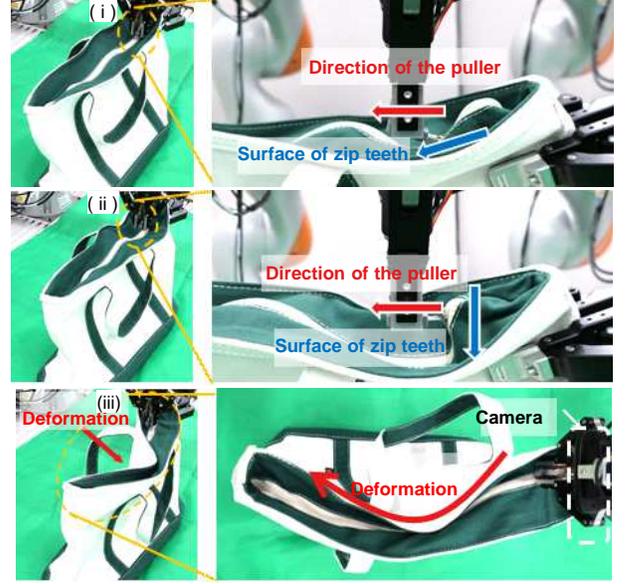


Fig. 4. The initial bag state to verify the task success rate. (i): without visible bag deformation. (ii): without visible bag, but the puller is perpendicular to the zip teeth. (iii): with visible bag deformation.

loss function. By adding g_f , the predicted attention points was efficiently searched near the current attention points. This is because the attention points are obtained for the gripper and the target object, and the position does not change significantly with a one-step time change. In this study, images were RGB color images with a resolution of 64×64 , so $H_i = 64, W_i = 64, C_i = 3$. One tactile sensor consists of 4×4 elements and outputs values in the 3-axis direction, so $H_p = 4, W_p = 4, C_p = 6$. The robot arm used has seven axes and the end-effector was a two-fingered gripper, so $M = 8$ and $J = 1000$. The number of the attention points was four, and considering the two coordinates x and y , $K = 2 \times 4$. Also, $\alpha = 0.0001$ (epoch=0), $\alpha = 0.1$ (epoch>1000), and α was increased in steps.

IV. EXPERIMENTS

A. Hardware

The experimental environment is shown in Fig. 1(a) and Fig. 3. As shown in Fig. 1(a), it consists of two robot arms of KUKA LBR iiwa 14 R820. Robotiq 2F-85 Adaptive Gripper was attached to each arm. An RGB camera was placed on the gripper, and Buffalo's BSW500M Series was used. In addition, tactile sensors in Fig. 3(b) were attached to each of the two fingers of the gripper, as shown in Fig. 3(a). For the fabric bag, a commercially available tote bag manufactured by Captain Stag was used. The puller was made with a 3D printer so that it could be easily grasped with the gripper. In this study, we focused on performing tasks in situations where the task state cannot be recognized only by images. Therefore, the upper part of the bag was fixed with one arm, the puller was grasped by the arm to be controlled, and the state of contact with the object was set as the initial state of the task. At this time, since the initial state of the arm

was the same, the position of the bag was placed within the range where the puller could be grasped. The bag and the puller were not fixed in a strict position, but were set to be in approximately the same position. Then, the puller was pulled to the other side of the bag, and the state where the bag was opened, which was defined as the end of the task.

B. Evaluation

The purpose of the experiment was to verify the effectiveness of our model. In particular, we investigated the effectiveness of using CNNs for tactility feature extraction, SMT, and tactility. Therefore, we compare the success rates of the tasks in the five models.

- (A) Without SMT and tactility, same as [14]
- (B) With SMT without tactility
- (C) With tactility without SMT
- (D) With SMT and tactility using CNN for tactility (ours)
- (E) With SMT and tactility using FC for tactility

Fig.4 shows the initial bag state to verify the task success rate. (i) was when the bag was not deformed significantly as shown in the figure on the left. Also, as shown in the figure on the right, the direction of the puller was almost horizontal to the zip teeth, and pulling the puller in the direction of the teeth end opens the bag, so it was not expected to be difficult to open. (ii) had no visible deformation in the figure like (i). However, as shown in the figure on the right, unlike (i), the puller is perpendicular to the teeth, so if it is pulled directly in the direction of the teeth end, it will become stuck and the zipper cannot be opened. Furthermore, the camera was located above the zipper, a case that we expected to require particularly tactility, since it appeared to be able to be opened by pulling it straight out using only vision. (iii) was a case where the bag was given a visible deformation, as shown in the left figure, and the arm fingertips needed to be moved appropriately to match the direction of the zipper. The deformation was also visible from the direction of the camera, as shown in the right figure. For each case, we performed 30 trials each and evaluated the success rate. The trial time was set to 10 seconds, which is the same as the training data, and if it did not finish within the time, it was defined as a failure.

C. Training Setup

As for the training data, the bag was randomly deformed, and a human operated the robot to acquire 36 patterns. Generally, Robot operation methods for acquiring the training data include arm hand position / posture control method using an operation user interface such as a game controller, bilateral control method [31], and direct teaching method. In this research, we chose the operation by direct teaching from the viewpoint of ease of implementation and the fact that the operator can operate while feeling the reaction force. The operator grabbed the arm and performed the unzipping motion of the bag while keeping the load on the hand small. Joint angle data of the robot arm, image data of the camera, and tactile sensor data were acquired for 10 seconds and 20 Hz. Each model was trained at 10,000 epochs and did not use

TABLE I

TASK SUCCESS RATE FOR EACH MODEL AND INITIAL BAG STATE.

Model	Bag state			
	(i)	(ii)	(iii)	All
(A) Base [14]	23.3%	0.0%	26.7%	16.7%
(B) w/ SMT	83.3%	0.0%	86.7%	56.7%
(C) w/ tac. (CNN)	63.3%	23.3%	36.7%	41.1%
(D) w/ SMT, tac. (CNN)	96.7%	86.7%	96.7%	93.3%
(E) w/ SMT, tac. (FC)	63.3%	63.3%	56.7%	61.1%

TABLE II

LOAD ON THE FINGERTIPS IN EACH MODEL AND INITIAL BAG STATE.

Model	Bag state		Average [N]	
	Maximum [N] (i)	(iii)	(i)	(iii)
(B) w/o tactility	15.8	22.5	8.0	10.1
(D) w/ tactility	12.5	14.7	7.8	8.0

pre-trained models. The batch size was 18, and the input data was scaled to [0,1.0]. The Adam optimizer [32] was used for training. Optimizer parameters were set to $\alpha = 0.001$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. We used TensorFlow as a software library. To train the models, the ABCI system of AIST was used. Due to the large network size of our model, large computational resources with graphics processing unit (GPU) memory were required for training. The ABCI system has 16GB per node, and we used two nodes in this study.

V. RESULTS

Table I shows the task success rates for each model and initial bag state. The comparison of the success rates of each model revealed the following.

- (A) and (B), (C) and (D) differed only in the presence or absence of SMT. In the comparison, the success rates of (B) and (D) using SMT were higher, indicating the effectiveness of SMT.
- (B) and (D) differed only in the presence or absence of tactility. In the comparison, the success rate of the model using tactility was higher, which indicates the effectiveness of tactility usage.
- (D) and (E) differed only in the tactility data processing method. In the comparison, the model that extracted features by CNNs had a higher success rate than the model that extracted features by FC, which indicates the effectiveness of CNNs for tactility data processing.

From the above, we were able to verify the effectiveness of using CNNs for tactility data processing, SMT, and tactility. Furthermore, we were able to achieve the highest success rate of 93.3% for our model with the three methods, demonstrating the effectiveness of our model. In addition, two questions were investigated on the basis of the results.

A. Is the attention mechanism functioning properly?

Fig. 5 shows (a) the input image and attention points, (b) predicted image, and (c) bird's eye view for our model (D) in the initial bag state (iii), where the visible deformation of the bag is large and the effect of attention can be easily determined. Time elapsed with the image on the right.

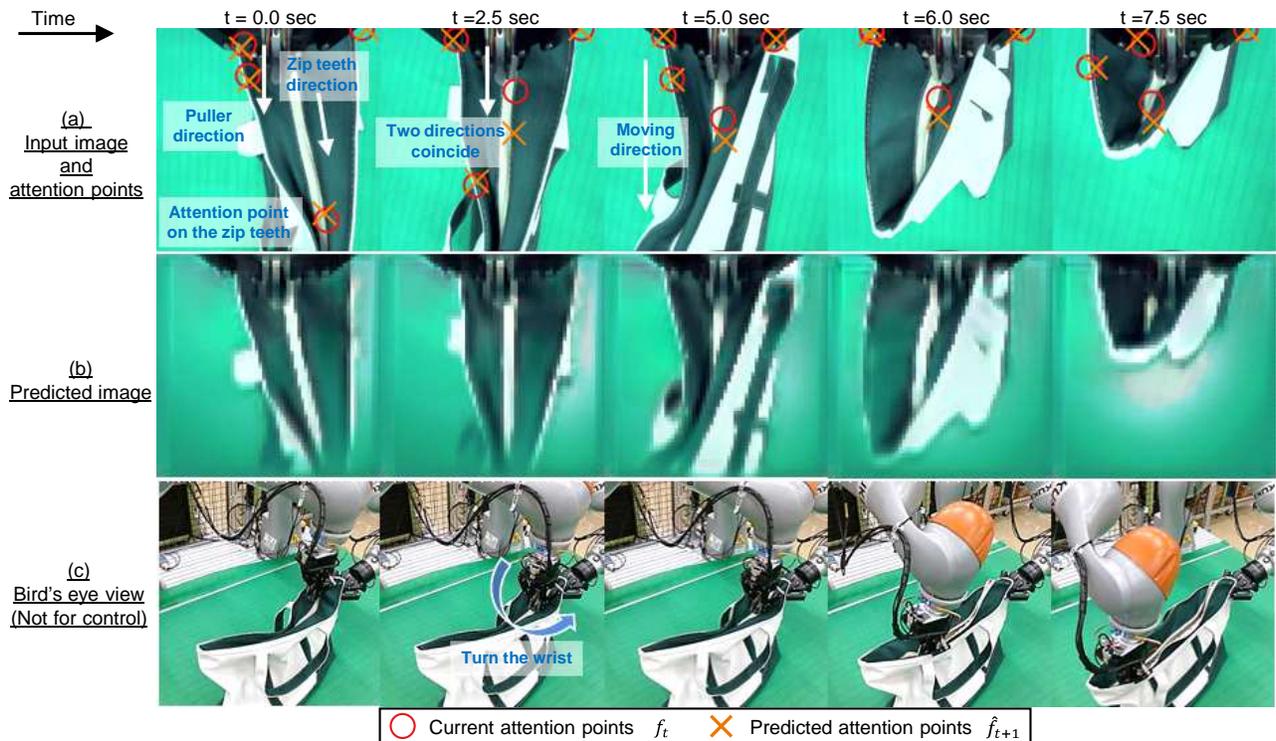


Fig. 5. Snapshots of unzipping task with our model in the initial bag state (iii). (a): images of the robot and attention points. The attention points are the position to which the robot was paying attention. (b): predicted images of time $t + 1$. (c): bird's eye view camera images.

There are two types of attention points: ○ indicates current attention points, and × indicates predicted attention points. Some attention points were on the gripper and zip teeth. These are important positions for the task. In addition, the attention points obtained on the zip teeth were ahead of the current attention points in the direction of motion, and the predicted attention points were also appropriately predicted. From the above, it can be considered that the attention points are obtained appropriately. Note that similar attention points were found for the other bag initial state. In addition, at the initial time 0 sec, the directions of the puller and the zip teeth were misaligned, but the attention points were output on the zip teeth, and at 2.5 sec, the arm wrist was rotated to match the direction of the puller and the zip teeth. After that, the zipper was successfully opened to the end by generating arm movements on the basis of the attention points obtained on the zip teeth. From the above, it can be concluded that the attention mechanism was functioning properly, as the unzipping task was successfully performed by generating motions on the basis of the obtained attention points.

B. Does tactility contribute to action generation?

Table II shows the maximum and average loads on the arm fingertips at the time of success in bag states (i) and (iii), where the success rates were nearly equal for models (B) and (D) without and with tactility. Here, the maximum load was the maximum value of the load for each trial averaged over the number of trials, and the average load was the average of all trials. Both the maximum and average loads were smaller

when using tactility in all bag conditions. In particular, in the initial bag state (iii), where the deformation of the bag was large and the operation based on tactility was necessary comparing to (i), the difference between the two cases was large. These results suggest that the motion was generated on the basis of tactility.

VI. CONCLUSION

We proposed a deep predictive learning robot control system that performs contact-rich manipulation of a flexible object, which was difficult to achieve by vision alone. Furthermore, we demonstrated its effectiveness with unzipping of a fabric bag. Our methods introduced an attention mechanism based on convolutional neural networks (CNNs) and soft argmax for extracting image features, softmax transformation for predicting motions, and CNN for extracting tactile features. In the experiment, the robot was able to generate motions while responding to the bag state, which dynamically deforms when it was unzipped, and showed a high success rate. In addition, the importance of tactility in the task was demonstrated by not only the success rate but also the difference in the load on the zipper. In this study, we considered the unzipping task as an example. Our model is not task-specific but can be widely applied to tasks.

ACKNOWLEDGEMENT

AI Bridging Cloud Infrastructure of National Institute of Advanced Industrial Science and Technology was used.

REFERENCES

- [1] J. Maitin-Shepard, M. Cusumano-Towner, J. Lei, and P. Abbeel, "Cloth grasp point detection based on multiple-view geometric cues with application to robotic towel folding," in *2010 IEEE International Conference on Robotics and Automation*, 2010, pp. 2308–2315.
- [2] S. Miller, J. van den Berg, M. Fritz, T. Darrell, K. Goldberg, and P. Abbeel, "A geometric approach to robotic laundry folding," *International Journal of Robotic Research - IJRR*, vol. 31, pp. 249–267, 02 2012.
- [3] Y. Gao, H. J. Chang, and Y. Demiris, "Iterative path optimisation for personalised dressing assistance using vision and force information," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016, pp. 4398–4403.
- [4] K. Suzuki, M. Kanamura, Y. Suga, H. Mori, and T. Ogata, "In-air knotting of rope using dual-arm robot based on deep learning," *CoRR*, vol. abs/2103.09402, 2021. [Online]. Available: <https://arxiv.org/abs/2103.09402>
- [5] C. Elbrechter, R. Haschke, and H. Ritter, "Folding paper with anthropomorphic robot hands using real-time physics-based modeling," in *2012 12th IEEE-RAS International Conference on Humanoid Robots (Humanoids 2012)*, 2012, pp. 210–215.
- [6] T. M. Caldwell, D. Coleman, and N. Correll, "Optimal parameter identification for discrete mechanical systems with application to flexible object manipulation," in *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2014, pp. 898–905.
- [7] P. Jimenez, "Survey on model-based manipulation planning of deformable objects," *Robotics and Computer-Integrated Manufacturing*, vol. 28, no. 2, pp. 154–163, 2012.
- [8] T. Power and D. Berenson, "Keep it simple: Data-efficient learning for controlling complex systems with simple models," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 1184–1191, 2021.
- [9] P.-C. Yang, K. Sasaki, K. Suzuki, K. Kase, S. Sugano, and T. Ogata, "Repeatable folding task by humanoid robot worker using deep learning," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 397–403, 2016.
- [10] K. Kawaharazuka, T. Ogawa, J. Tamura, and C. Nabeshima, "Dynamic manipulation of flexible objects with torque sequence using a deep neural network," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 2139–2145.
- [11] D. Tanaka, S. Arnold, and K. Yamazaki, "Emd net: An encode-manipulate-decode network for cloth manipulation," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 1771–1778, 2018.
- [12] Y. Wu, W. Yan, T. Kurutach, L. Pinto, and P. Abbeel, "Learning to manipulate deformable objects without demonstrations," *arXiv preprint arXiv:1910.13439*, 2019.
- [13] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *CoRR*, vol. abs/1504.00702, 2015. [Online]. Available: <http://arxiv.org/abs/1504.00702>
- [14] H. Ichiwara, H. Ito, K. Yamamoto, H. Mori, and T. Ogata, "Spatial attention point network for deep-learning-based robust autonomous robot motion generation," *arXiv preprint arXiv:2103.01598*, 2021.
- [15] R. S. Dahiya, P. Mittendorf, M. Valle, G. Cheng, and V. J. Lumelsky, "Directions toward effective utilization of tactile skin: A review," *IEEE Sensors Journal*, vol. 13, no. 11, pp. 4121–4138, 2013.
- [16] A. Yamaguchi and C. G. Atkeson, "Combining finger vision and optical tactile sensing: Reducing and handling errors while cutting vegetables," in *2016 IEEE-RAS 16th International Conference on Humanoid Robots (Humanoids)*. IEEE, 2016, pp. 1045–1051.
- [17] T. P. Tomo, A. Schmitz, W. K. Wong, H. Kristanto, S. Somlor, J. Hwang, L. Jamone, and S. Sugano, "Covering a robot fingertip with uskin: A soft electronic skin with distributed 3-axis force sensitive elements for robot hands," *IEEE Robotics and Automation Letters*, vol. 3, no. 1, pp. 124–131, 2017.
- [18] M. K. Johnson and E. H. Adelson, "Retrographic sensing for the measurement of surface texture and shape," in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 1070–1077.
- [19] M. Li, H. Yin, K. Tahara, and A. Billard, "Learning object-level impedance control for robust grasping and dexterous manipulation," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*, 2014, pp. 6784–6791.
- [20] S. Funabashi, A. Schmitz, T. Sato, S. Somlor, and S. Sugano, "Versatile in-hand manipulation of objects with different sizes and shapes using neural networks," in *2018 IEEE-RAS 18th International Conference on Humanoid Robots (Humanoids)*, 2018, pp. 1–9.
- [21] S. Tian, F. Ebert, D. Jayaraman, M. Mudigonda, C. Finn, R. Calandra, and S. Levine, "Manipulation by feel: Touch-based control with deep predictive models," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 818–824.
- [22] M. A. Lee, Y. Zhu, K. Srinivasan, P. Shah, S. Savarese, L. Fei-Fei, A. Garg, and J. Bohg, "Making sense of vision and touch: Self-supervised learning of multimodal representations for contact-rich tasks," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8943–8950.
- [23] H. van Hoof, N. Chen, M. Karl, P. van der Smagt, and J. Peters, "Stable reinforcement learning with autoencoders for tactile and visual data," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2016, pp. 3928–3934.
- [24] N. Vulin, S. Christen, S. Stevšić, and O. Hilliges, "Improved learning of robot manipulation tasks via tactile intrinsic motivation," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2194–2201, 2021.
- [25] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [26] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [27] C. Finn, X. Y. Tan, Y. Duan, T. Darrell, S. Levine, and P. Abbeel, "Deep spatial autoencoders for visuomotor learning," in *2016 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2016, pp. 512–519.
- [28] D. C. Luvizon, D. Picard, and H. Tabia, "2d/3d pose estimation and action recognition using multitask deep learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5137–5146.
- [29] J. Ngiam, Z. Chen, D. Chia, P. Koh, Q. Le, and A. Ng, "Tiled convolutional neural networks," *Advances in neural information processing systems*, vol. 23, 2010.
- [30] J. Hwang, J. Kim, A. Ahmadi, M. Choi, and J. Tani, "Predictive coding-based deep dynamic neural network for visuomotor learning," in *2017 Joint IEEE International Conference on Development and Learning and Epigenetic Robotics (ICDL-EpiRob)*, 2017, pp. 132–139.
- [31] A. Sasagawa, S. Sakaino, and T. Tsuji, "Motion generation using bilateral control-based imitation learning with autoregressive learning," *IEEE Access*, vol. 9, pp. 20 508–20 520, 2021.
- [32] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.