# Bimanual Shelf Picking Planner Based on Collapse Prediction

Tomohiro Motoda[1], Damien Petit[1], Weiwei Wan[1,2], and Kensuke Harada[1,2]

*Abstract*— In logistics warehouse, since many objects are randomly stacked on shelves, it becomes difficult for a robot to safely extract one of the objects without other objects falling from the shelf. In previous works, a robot needed to extract the target object after rearranging the neighboring objects. In contrast, humans extract an object from a shelf while supporting other neighboring objects. In this paper, we propose a bimanual manipulation planner based on collapse prediction trained with data generated from a physics simulator, which can safely extract a single object while supporting the other object. We confirmed that the proposed method achieves more than 80% success rate for safe extraction by real-world experiments using a dual-arm manipulator.

## I. INTRODUCTION

In logistics warehouses, we often have to extract a single object that is wedged between other objects on a shelf, which is potentially dangerous for heavy objects to fall and injure human workers. In this case, when a robot tries to extract one of the objects, it has to consider the positional relationship of overlapping objects and manipulate them accordingly. So far, various approaches have been proposed to extract an object from a shelf. In [1]–[3] different methods are proposed but require a series of rearrangement operations. In other cases, extraction and support relations are analyzed between pairs of objects from 3D visual perception [4]. However, in all previous approaches, a robot extracts the target object after rearranging its neighboring objects.

Humans however, extract an object from a shelf while supporting other neighboring objects as shown in Fig. 1 (a). Based on this observation, we propose a bimanual manipulation planner to extract a target object from a shelf while supporting the other object as shown in Fig. 1 (b). To extract an object from a pile without collapse, we need to determine which of the target's neighboring object the robot have to support. We propose a learning-based approach on extracting the target object from the pile while supporting the objects. A network model based on a Fully Convolutional Network (FCN) [5] has been designed to predict the pile state while extracting the target object with a pixel-wise collapse probability map. The inputs of the network are a depth image of the shelf content, and two binary masks corresponding to the two objects selected for extraction and support. The output of the network model is a labeled image predicting the collapsing region while the target object is extracted. If the output includes large collapsing region, we judge that such selection of supported object is not better. Given this output, the robot can select the proper object to support by defining

[1]Graduate School of Engineering Science, Osaka University, Japan
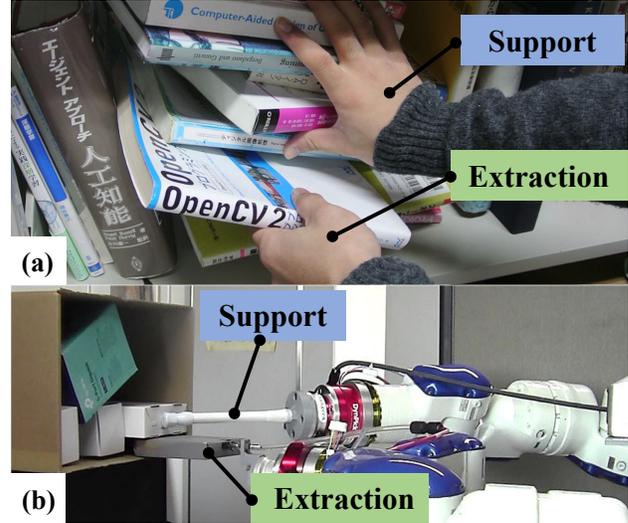[2]National Institute of Advanced Industrial Science and Technology (AIST), Japan

Fig. 1. Extracting the target object while supporting others: (a) a human is extracting a book from the shelf while supporting the neighboring books, (b) robotic bimanual manipulation for safely extracting an object from the shelf.

the ratio of the predicted collapsing region as the safety index to the shelf picking. In addition, to generate a large number of training data of depth images, related binary masks and label images, we use a physics simulation of the piled objects and of the extraction/support action. We experimentally verify the effectiveness of our proposed method by using a real dual-arm manipulator. We show that the robot can safely extract the target object from a shelf with a success rate larger than 80%. By using our proposed method, we do not need to rearrange the objects placed on a shelf to extract the target object and so we increase the picking efficiency.

Our main contributions are:

- A Fully Convolutional Networks to infer the pixel-wise probability map of the collapsing region while extracting a selected object from a shelf (Subsection III-B).
- A physics simulation that generate the necessary training data for the FCN (Subsection III-A).
- A robotic system able to extract a target object from a pile, in a shelf, without rearranging its surrounding objects

This paper is organized as follows. In Section II, we discuss the related work on shelf picking. In Section III our proposed shelf picking method is described including a detailed explanation of the network model and its implementation. In Section IV we describe the experimental setup. In Section V, we discuss the result of our approach and
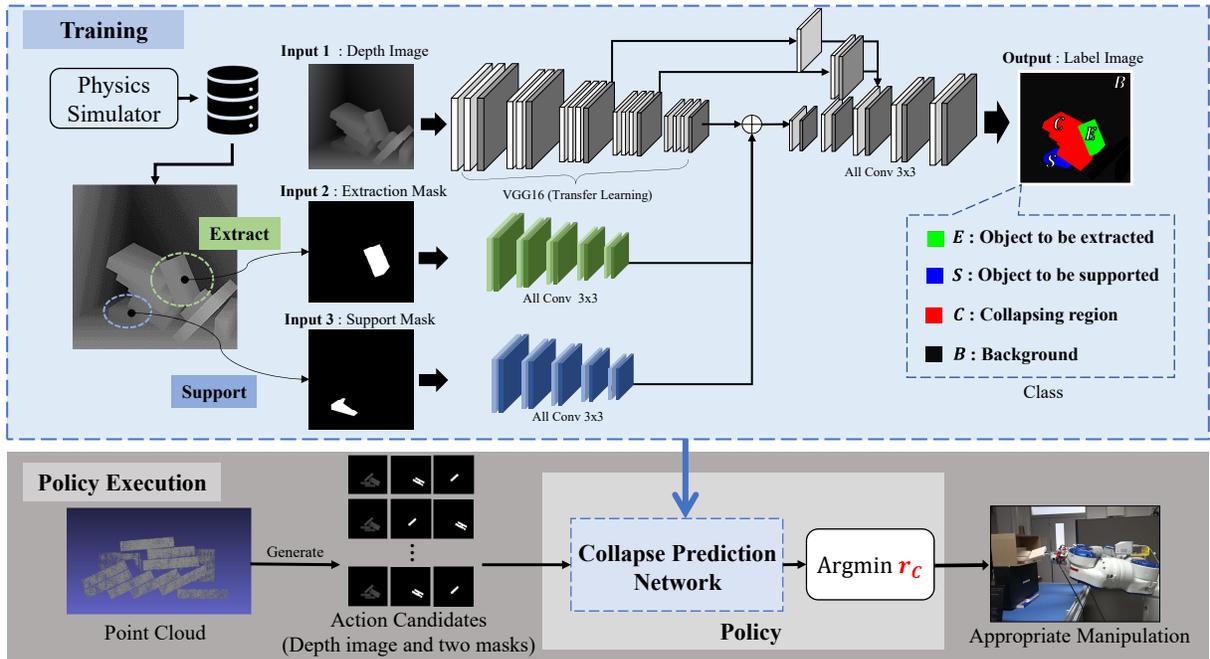
Fig. 2. Schematic overview of the proposed manipulation policy based on Fully Convolutional Network (FCN). The inputs are a depth image and two binary masks. The model encodes the depth image with the VGG-16 network and two masks with a five-layer network, and concatenates these three networks. We focus on the collapsing region, $C$ (highlighted in red), in the output of the network. The method uses the argmin of $r_C$ (ratio of $C$ in the image), to return the appropriate action. The size of the depth image and its related two masks are $256 \times 256$.

experiment. Finally, the paper closes with conclusions and future work.

## II. RELATED WORK

In this section, we introduce some related works on object picking in logistics warehouse. This topic has been extensively researched. There have been some works done on picking an object stored in a box, such as [6], [7]. Among them, we focus on studies extracting an object from a randomly piled objects on a shelf, such as [1]–[4]. Temtsin et al. [8] ranked each object using a measure based on the geometrical relationships of objects and extracted an object with high rank. Mojtahedzadeh et al. [4] and Wu et al. [9] proposed methods to learn the motion of robots in stacked or scattered environments. Some researches achieved manipulation in clutter based on partially observable Markov decision processes (POMDPs). For example. Pajarinen et al. [10] used iterative picking/observation to disassemble the cluttered objects based on the POMDPs. To pick an object from a shelf, Li et al. [2] used POMDP to find the target object by rearranging the objects on the shelf safely and efficiently. Zhang et al. [11], [12] used a Convolutional Neural Networks (CNNs) to estimate the order of extracting the overlapped objects by using the graph representation of the objects' position. Grotz et al. [13] determined the order of objects to be manipulated by taking into account their support relations. However, all these methods used to extract a target object from a shelf need to repeatedly rearrange the overlapping objects of the target before the extraction to avoid a collapse. As far as the authors know, there has been no research on bimanual manipulation planning to extract the target object while supporting its neighboring objects at the same time, in spite of its efficiency.

## III. SHELF PICKING METHOD IMPLEMENTATION

In this paper, we propose a bimanual manipulation method to extract a target object from a pile while supporting the other object. In order to first verify the effectiveness of our new approach we assume that the robot achieves the task by pulling a box-shaped object out horizontally. Assuming a situation in which the insertion of fingers between objects is difficult for the robot, one arm is mounted with a suction gripper to extract the target object. The other arm has a rod-shaped end-effector to support other objects as seen in Fig. 1. We use a depth sensor to provide a 3D point cloud captured from a front point of view of the shelf containing the pile of object.

Fig. 2 illustrates the flow of our overall architecture. The user selects the object subject to extraction, then a FCN is used to predict which objects will be affected during the extraction (I.e. collapsing region).

In the following subsections, the different steps are explained in detail.

### A. Physics Simulator for Data Generation

In this subsection, we describe the setup of the physics simulation system used for data generation.

*1) Scene Generation:* We generate a randomly stacked state of objects in the simulator. In this study, we use PhysX[1],

---

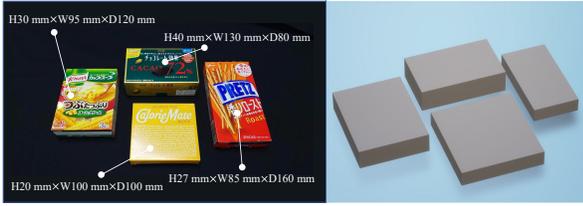[1]https://developer.nvidia.com/gameworks-physx-overview

Fig. 3. Types of objects used in the simulations. The left actual boxes. The right shows the 3D models used in the simulator.
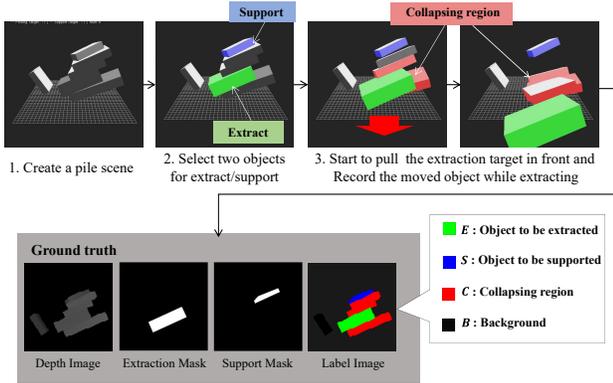


Fig. 4. Dataset generation procedure with physics simulator. The upper row shows a scene of a simulation while the lower row shows the images of ground truth generated from the simulated scene.

i.e., a physics simulator, to configure and simulate the environment. Our simulator is designed with the following settings. Considering a situation where many product boxes are on a shelf, thus we set the simulation parameter referred to the actual movements of them. For both objects and a shelf in our environment, we set the coefficient of static and dynamic friction to be 0.9 and 0.8, respectively. We also set the coefficient of restitution to be 0.1 and the density to be 1.0 $kg/m^3$. We perform the shelf picking simulation by placing six objects from a set of objects. As the number of objects on a shelf increases, the extraction generally becomes more difficult. In our study, we fix the number of the objects to be six, which can generate the successful cases empirically in about 50% even if the target object for extraction/support is randomly selected. Moreover, for the after-mentioned verification, we prepared two sets of objects: One type of object (H 20 mm× W 100 mm× D 100 mm), and four objects of various sizes (the height is 27 – 40 mm, the width is 85 – 130 mm, and the depth is 80 – 160 mm), illustrated in Fig. 3 for the detail. We generate the dataset with either of the sets according to the conditions.

*2) Data generation and Simulation Procedure:* Fig. 4 shows the process of simulation. First, the simulator creates a pile scene. Second, one object moves horizontally in the simulation. Here, we assume that the robot pulls out one object from the shelf horizontally toward the observer. The other object is supported, and we assume that it remains stationary, hence, it is not affected by interference or gravity, and its pose does not change. Finally, in case there is some change in any object pose other than the two targets, we record

these objects as collapsed subject to entangling/collision. In one simulation, we can obtain the tuple, consisting of three images as input data and one labeled image as output data, as shown in Fig. 4.

*B. Collision Prediction Network*

This subsection describes the neural network that predicts the objects affected by the target object extraction and so most likely to fall or collapse. Similar to the fully convolutional network (FCN) used in [5], our model classifies each pixel belonging to the collapsing region.

*1) Ground Truth:* The input data consist of a depth image ($256 \times 256$) and two binary masks ($256 \times 256$). One mask is an object for extraction, and the other mask is an object for support. The output data is a prediction of the classifications for each pixel of the image ($256 \times 256$). We define four classes: Object to be extracted $E$, Object to be supported $S$, Collapsing region $S$, and Background region $B$, as shown in Fig. 4. The collapsing region $C$ expresses the region of objects which move or fall from the shelf while extracting the target object. These input/output data are automatically generated from the simulator.

*2) Network Architecture:* Our network model consists of an encoder for extracting the feature value of the input and a decoder for producing the segmented image at its original resolution. Fig. 2 illustrates the network architecture. First, the encoder part consists of three networks. The model generates the feature maps for a depth image with the VGG-16 network [14] pre-trained by ImageNet [15] and for two masks, each with five convolution layer network. These three outputs are then concatenated into one feature map. Next, the decoder part with five convolution layers up-sample the feature maps to the original resolution with deconvolution. Moreover, our network uses the skip architecture by referring to prior examples [5], [16], which combines the feature maps of the lower layers with those of the upper layers to recover the general location information while preserving the local information.

*C. Manipulation Planning*

This section describes the manipulation procedure to perform the extraction by applying the trained network. The robot acquires point clouds from a depth sensor installed in front of the shelf and generates three input images from this observation. One is a depth image converted from the point clouds to the depth map. The other two images are mask images representing the object to be extracted and the object to be supported. The mask image, $M_c$, is a binary image from each cluster, $c_i$ ($i = 0, 1, 2, ..., N - 1$) of point clouds, which is classified by object segmentation based on the region growing method [17] and the binarization. Fig. 5 shows the process. In the actual experiments, the robot end-effectors approach each object toward the center of gravity in these masks.

We can define action candidates $A$ according to use situations: (1) the other situation is that the robot chooses the safest pair of extraction/support objects (for example,
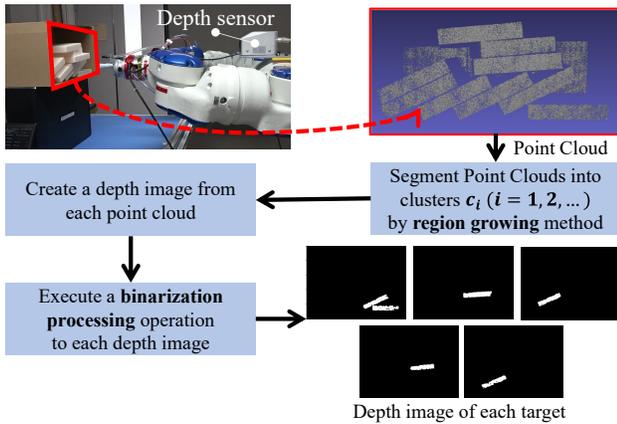
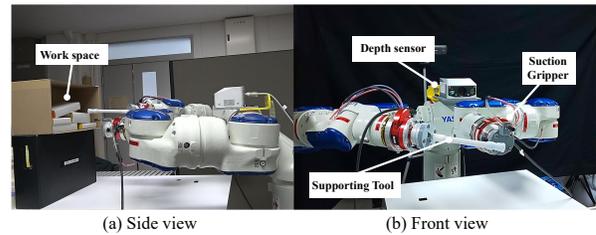Fig. 5. Process of generating candidates with the segmentation method.



Fig. 6. Experimental setup.



Fig. 7. Objects used to evaluate the generalization of the proposed method: (a) experiments for objects of the same size, and experiments for objects of various sizes, and (c) experiments for new objects (not used in simulations).

empty a shelf). In this case, we define action candidates $A$ by preparing all the combinations of two different targets of the extraction/support action. (2) one is that we need to extract a predetermined target object. In this case, in advance the user selects an object subject to extraction arbitrarily, we define action candidates $A$ by choosing another subject to support,

Next, we define the safety index to select the best action from all the candidates. Here, the output shows the region, $R_E$, $R_S$, $R_C$, and $R_B$ that indicates the regions of four different classes ($E$, $S$, $C$, and $B$). If $R_C$ is large, it will increase the risk of collapse. Based on this assumption, we can define the following risk index:

$$r_c\ (a) = \frac{area(R_C^a)}{area(R_E^a \cup R_S^a \cup R_C^a \cup R_B^a)} \quad (1)$$

$R_E^a$, $R_S^a$, $R_C^a$, and $R_B^a$ denote the regions in the output of an action candidate $a$ for two selected objects. $area(\cdot)$ indicates the area of the region. Our algorithm selects the input data that is the smallest for index $r_c$ based on Eq. (1) and determines the best action, $a$, of all the action candidates to be manipulated by the robot.

$$a = \arg\ \min_{a' \in A} r_c\ (a') \quad (2)$$

If the robot's motion is out of the control range, we eliminate it from the candidates and select the next best move.

## IV. EXPERIMENTS

### A. Training Settings

In our experiments, we acquired 15,000 pairs of input and output images from the simulator. From these datasets, we used 90% as training data and 10% as validation data. We augmented our training dataset through left-right inversion and utilized the network using 30,000 pairs. We set the initial learning rate to 0.0001 up to 30 epochs and 0.00001 from 30 epochs onward. The batch size was 1, and we use the Adam [18] as the optimizer. The number of epochs during training was 50, and each epoch required 27,000 iterations. In our training, we used the NVIDIA RTX 2060 super (8 GB VRAM).

### B. Experimental Setup

To verify the effectiveness of the proposed method, we conducted experiments using an actual robot under several conditions. Fig. 6 shows the experimental environment used in the verification. We use the MOTOMAN-SDA5F (Yaskawa Electric Corp.)[2], a bimanual robot with 7 degrees-of-freedom robot arms, which has a suction gripper and a plastic rod-shaped end effector (the bar's length is 20 cm) at the tip of each arm of the robot. The YCAM3D-10L (YOODS Co., Ltd.)[3], a 3D depth sensor, is installed at the front of the shelf.

### C. Results

To evaluate the performance by using the actual robot, we consider experiments on two use situations.

*1) Choosing the safest pair of extraction/support object:* We verify the performance of our prediction network through experiments that the obot always chooses the safest action. In our real-world experiments, we used the target objects as shown in Fig. 7(a), (b), which is the same size as models used in our simulations (Fig. 3). Moreover, in order to evaluate the generalization capability, we separately prepared new objects (Fig. 7(c)). We used the following conditions in experiments:

- 5 objects of same size, as shown in Fig. 7(a).
- 5 objects of various sizes, as shown in Fig. 7(b).

TABLE I

EXPERIMENTAL RESULT

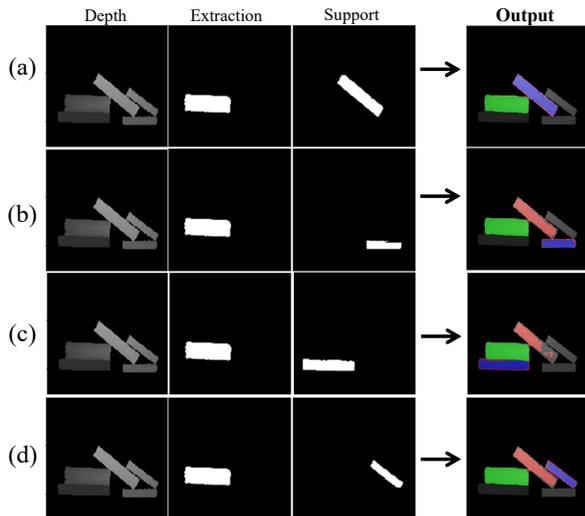| | Objects used in experiments | | | | |
| Proposed model | Same-size (5 objects) | Various-size (5 objects) | New various-size (5 objects) | Same-size (10 objects) | Total |
|---|---|---|---|---|---|
| Trained with same-size objects | 18/20 | 16/20 | 15/20 | 16/20 | **65/80 [81.3%]** |
| Trained with various-size objects | 17/20 | 17/20 | 17/20 | 13/20 | **64/80 [80.0%]** |

Fig. 8. Visualization results from the proposed network model: (a)-(d) Outputs of extraction for a target object while supporting different objects. The red area on the images shows the predicted collapse region ($R_C$). The green region shows the object to be extracted correspond to $R_E$, and the blue region shows the object to be supported correspond to $R_S$.
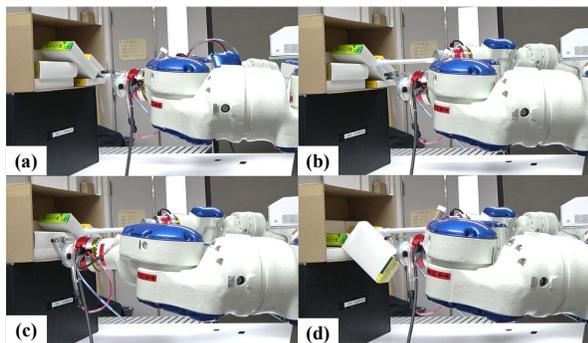


Fig. 9. Description of target manipulation with the proposed method: (a)-(d) a series of scenes in one task.

- 5 new objects of various sizes (not used on the simulator), as shown in Fig. 7(c).
- 10 objects of same size, as shown in Fig. 7(a).

Furthermore, we prepared two network models trained with different datasets generated from the simulations with objects of the same size or various sizes. With each trained model, we conducted 20 trials in every four patterns by changing the size or number of objects. If the robot removes only one object from the shelf, we regard it as a success; otherwise, we consider it failed.

As shown in Table I, we conducted robotic experiments under the above mentioned conditions. The robot achieved the high success rate across all conditions, and the overall extraction success rate were 81.3% (65/80) for the model trained with objects of the same size and 80% (64/80) for the model trained with the dataset of various sizes.

*2) Extracting a predetermined target object:* We assume that a specific target object is needed in a practical situation. The user chooses one object to extract from a shelf in advance, and in that case, our policy determines which is

TABLE II
TRAINED MODEL EVALUATION.

| Proposed Model | Ground-truth | Network performance | | |
| --- | --- | --- | --- | --- |
| | | Avg. IoU | Avg. Recall | Avg. Precision |
| Trained with same-size objects | Same-size target | 0.339 | 0.511 | 0.437 |
| | Various-size target | 0.438 | 0.641 | 0.530 |
| Trained with various-size objects | Same-size target | 0.359 | 0.438 | 0.576 |
| | Various-size target | 0.452 | 0.511 | 0.697 |

the object to support correctly.

We show the output of our network as Fig. 8(a)-(d) in the case. The robot selects the best action from the output that region $R_C$ (highlighted in red) is small as Fig. 8 (a). Fig. 9 shows the experimental scene. When the correct action is selected, the robot first presses down the one object with the stick on the right-hand and then pulls out the target with the left-hand suction gripper. Based on the learning results, we confirmed that the robot selected combinations of objects that are less likely to collapse and execute the safest manipulation.

Moreover, we conduct 20 trials in that case. At each trial, the object to be extracted is not changed. It should be noted that we trained the network with the dataset generated in the simulations using various-size objects (Fig. 3) in this verification. In 20 trials of experiments under this conditions, the robot can extract single target object without collapse in success rate of 85% (17/20). We confirm that our network works well.

## V. EVALUATION

In this section, the performance is evaluated concerning two points. (1) We set a benchmark of the prediction performance based on segmentation metrics and compare our proposed network under different conditions, (2) we acquire the success rate, representing the percentage of the completed when extracting a single target object without collapse by using a real robot.

### A. Prediction Performance

We confirm the network can correctly predict the collapsing regions with ground-truth data, as shown Table II. To evaluate the performance of the collapse prediction, we focus only on the collapsing region $C$ in this study. Our metrics include $precision$, $recall$, and $IoU$ calculated in pixels, between the predicted and ground-truth. We calculate the average values on metrics with a hundred ground-truth data, and compare two networks trained with different training datasets. Moreover, to verify a generalization of the performance, we prepare the ground-truth in two different patterns: target objects of same size or various sizes. We empirically set the threshold of classification for each pixel to 0.4.

As shown in Table II, even when we use networks with different training datasets, there is no significant difference on each metric to the same ground-truth. This result indicates that the size of the object has little effect on learning. In contract, when we use the network trained with the

objects of various sizes , $IoU$ and $precision$ increase in both ground-truth data. By using our method, the collapsing region tends to become a shape similar to the object model. It is assumed that the network trained with objects of same size is relatively sensitive to shape differences. Therefore, training with objects of various sizes works well for correctly predicting the region.

### B. Real-world Manipulation

As shown in Table I, the robot extracted successfully up to 81.3% (65/80) for the same object dataset and 80% (64/80) for the dataset of objects of different sizes. The success rate of each object is not significantly affected in different datasets. Similarly, there is no difference in the success rate when the objects are the same (Fig. 7(a)) and when the size of the objects is randomized (Fig. 7(b)). The success rates of 75% (15/20) and 85% (17/20) were confirmed in the experiments with objects of new various sizes (Fig. 7(c)), indicating that there was no overfitting of our learning results.

In extracting a predetermined target object, our method achieved a high success rate of 85%, indicating that our method can work well in logistics warehouse. Moreover, the success rate is almost equal to other experimental results. We showed that it was possible to make good predictions even when the target object to be extracted was limited.

In some failed cases, the robot executed intuitively incorrect actions, such as supporting an object unrelated to extracting a target. This indicates a problem with simulator settings. For example, there are trials where the robot can remove one target safely without supporting the other object. We also consider that the robot misunderstood some uncertain manipulations as successful trial. It is necessary to reconstruct the dataset or evaluate each action on each successful trial. As shown in Table I, the success rate decreased in ten objects of the same size. For example, if an object is not simply put on another object, the robot needs to support more than two objects. In our method, however, the robot can only support one object, causing a low success rate. In our future work, we will address this issue.

### C. Discussion

We proposed the learning-based approach to predict collapse after extracting one object while supporting the other one. The conventional learning-based approaches [11], [12] predict the support relationship as Section II mentioned. However, considering the complex pile, it will become more difficult to determine the support object by geometry shape and/or physical interaction. In contrast, our proposed method can directly predict whether the selected action is proper or not without checking the complex structure. However, in order to realize the new idea, we focused only on box-shaped objects for the sake of simplicity. Our future work will be extended to more complex-shaped objects and apply them to daily life.

## VI. CONCLUSIONS

This paper described a shelf picking method for safely extracting a single object from a shelf while supporting the other object. By using our proposed network model that predicts the objects that would collapse, a bi-manual robot was able to extract the object without objects falling.

In the future, we plan to make improvements on support actions and our simulator. In particular, we will analyze the trial result on each simulation by adding actions to support and extract in an appropriate way for object types.

## REFERENCES

[1] C. Nam, J. Lee, S. H. Cheong, B. Y. Cho, and C. Kim, "Fast and resilient manipulation planning for target retrieval in clutter," in *Proc. of IEEE Int. Conf. on Robotics and Automation (ICRA)*, pp. 3777–3783, 2020.

[2] J. K. Li, D. Hsu, and W. S. Lee, "Act to see and see to act: POMDP planning for objects search in clutter," in *Proc. of IEEE/RSJ Int. Conf. on Intelligent Robots and System (IROS)*, pp. 5701–5707, 2016.

[3] J. Lee, Y. Cho, C. Nam, J. Park and C. Kim, "Efficient Obstacle Rearrangement for Object Manipulation Tasks in Cluttered Environments," In *Proc. of IEEE Int. Conf. on Robotics and Automation (ICRA)*, pp. 183–189, 2019.

[4] R. Mojtahedzadeh, A. Bouguerra, E. Schaffernicht, and A. J. Lilienthal, "Support relation analysis and decision making for safe robotic manipulation tasks," *Robotics and Autonomous Systems*, vol. 71, pp. 99–117, 2015.

[5] E. Shelhamer, J. Long, and T. Darrell, "Fully Convolutional Networks for Semantic Segmentation," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 640–651, 2017.

[6] J. Mahler, J. Liang, S. Niyaz, M. Laskey, R. Doan , X. Liu, J. A. Ojea, and K. Goldberg, "Dex-Net 2.0: Deep learning to plan robust grasps with synthetic point clouds and analytic grasp metrics," *Robotics: Science and Systems (RSS)*, 2017.

[7] Y. Deng, X. Guo, Y. Wei, K. Lu, B. Fang, D. Guo, H. Liu, and F. Sun, "Deep Reinforcement Learning for Robotic Pushing and Picking in Cluttered Environment," in *Proc. of IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, pp. 619–626, 2019.

[8] S. Temstsin, and A. Degami, "Decision-making algorithms for safe robotic disassembling of randomly piled objects," *Advanced Robotics*, vol. 31(23-24), pp. 1281–1295, 2017.

[9] H. Wu, Z. Zhang, H. Cheng, K. Yang, J. Liu, and Z. Guo, "Learning Affordance Space in Physical World for Vision-based Robotic Object Manipulation," in *Proc. of IEEE Int. Conf. on Robotics and Automation (ICRA)*, pp. 4652–4658, 2020.

[10] J. Pajarinen and V. Kyrki, "Robotic manipulation in object composition space," In *Proc. of IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, pp. 1–6, 2014.

[11] H. Zhang, X. Lan, X. Zhou, Z. Tian, Y. Zhang, and N. Zheng, "Visual Manipulation Relationship Network for Autonomous Robotics," in *Proc. of 2018 IEEE-RAS 18th Int. Conf. on Humanoid Robots (Humanoids)*, pp. 118–125, 2018.

[12] H. Zhang, X. Lan, S. Bai, L. Wan, C. Yang, and N. Zheng, "A Multi-task Convolutional Neural Network for Autonomous Robotic Grasping in Object Stacking Scenes," in *Proc of IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS)*, pp. 6435–6442, 2019.

[13] M. Grotz, D. Sippel, and T. Asfour, "Active Vision for Extraction of Physically Plausible Support Relations," in *Proc. of IEEE-RAS Int. Conf. on Humanoid Robots (Humanoids)*, pp. 439–445, 2019.

[14] K. Simonyan, and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," in *Proc. of Int. Conf. on Learning Representations (ICLR)*, 2015.

[15] J. Deng, W. Dong, R. Socher, L. J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," *IEEE Computer Vision and Pattern Recognition (CVPR)*, 2009.

[16] O. Ronneberger, P. Fischer, and T. Brox, "U-Net: Convolutional Networks for Biomedical Image Segmentation," *Int. Conf. on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, pp. 234–241, 2015.

[17] T. Rabbani, F.A. van den Heuvel, and G. Vosselman, "Segmentation of point clouds using smoothness constraint," in *Proc. of the ISPRS commission V symposium Image Engineering and Vision Metrology*, pp. 248–253, 2006.

[18] D. Kingma, and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint*, arXiv:1412.6980, 2014