# Text-driven object affordance for guiding grasp-type recognition in multimodal robot teaching

Naoki Wake[1], Daichi Saito[2], Kazuhiro Sasabuchi[1], Hideki Koike[1], and Katsushi Ikeuchi[1]

*Abstract*—This study investigates how text-driven object affordance, which provides prior knowledge about grasp types for each object, affects image-based grasp-type recognition in robot teaching. The researchers created labeled datasets of first-person hand images to examine the impact of object affordance on recognition performance. They evaluated scenarios with real and illusory objects, considering mixed reality teaching conditions where visual object information may be limited. The results demonstrate that object affordance improves image-based recognition by filtering out unlikely grasp types and emphasizing likely ones. The effectiveness of object affordance was more pronounced when there was a stronger bias towards specific grasp types for each object. These findings highlight the significance of object affordance in multimodal robot teaching, regardless of whether real objects are present in the images. Sample code is available on GitHub.

## I. INTRODUCTION

Robot grasping has been a major issue in robot teaching for decades [1], [2]. Because robot grasping determines the positional relationship between a robot's hand and an object, grasping objects suitable for the given environment is critical for efficient and successful manipulations after grasping. Recent studies have focused on learning-based end-to-end robot grasping [3]–[8], where contact points or motor commands are estimated from visual input. However, a desired grasp differs depending on the type of manipulation to be achieved, even for the same target object. While such grasp uncertainty can be addressed in an automatic manner using an advanced robot control method (e.g., [9]), a simpler approach can be employed in the context of robot teaching, where a human teaches the robot how to grasp through a demonstration.

We have been developing a platform to teach a robot "how to grasp and manipulate an object" through multimodal human demonstrations [10]–[14] (Fig. 1). The demonstration is accompanied by verbal instructions and captured by a head-mounted device (HMD). The user demonstrates object manipulation using either physical objects or illusory objects superimposed by the HMD on the demonstrator's hands. According to the definition of Milgram et al. [15], we refer to the later setup as mixed reality (MR). Assuming such a multimodal teaching system, this paper focuses on recognizing grasp types based on the name of the object and the first-person image at the time of grasping.

The problem of recognizing grasp types from human grasping images is not new. However, because grasp-type recognition has developed in the context of computer vision, most existing research is image-based (e.g., [17]–[19]). When
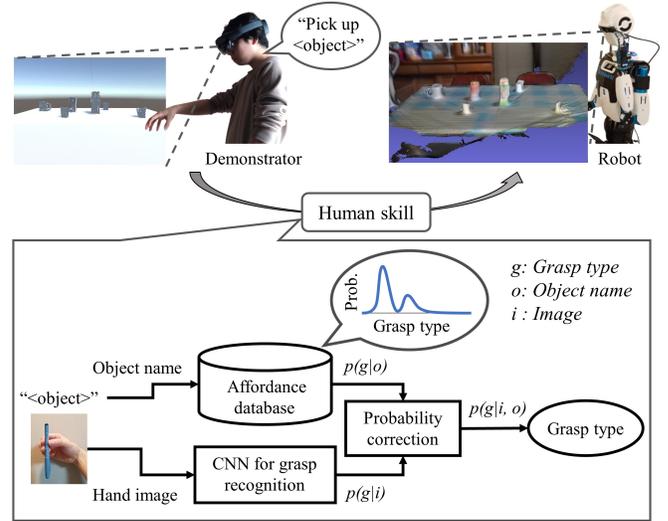


Fig. 1. Conceptual diagram of robot teaching. (Top) Head-mounted device provides first-person images during a demonstration with verbal instructions (modified version of image from [10]). The demonstrations are transferred to a robot in the form of a skill set, which includes a grasp type. (Bottom) Proposed pipeline for grasp-type recognition leveraging object affordance. The pipeline estimates the grasp type from the pairing of an object name and an image of a hand grasping that object. Object affordance is searched from an affordance database using text matching (modified image in [16]).

considered in terms of multimodal robot teaching, further questions arise. 1) How can linguistic input be utilized? 2) In what situations is linguistic input more advantageous? 3) Is linguistic input useful even in a challenging situation, such as MR, where images are not projected? Although these are practical and important questions in robot teaching, to the best of our knowledge, no previous studies have addressed these issues.

In many cases, an object name is known to be associated with the possible grasp types [18], [20]–[22]. Based on this association, we have previously proposed a pipeline that leverages a prior distribution of grasp types to improve a convolutional neural network (CNN)-based image recognition [16](Fig. 1 bottom). We refer to the prior distribution as object affordance, a concept proposed by Gibson [23]. In the pipeline, appropriate object affordance was searched from an affordance database using text matching. Although our preliminary experiments suggested that the object affordance is a promising solution to leverage a user's linguistic input for grasp-type recognition, its effectiveness has not been fully understood due to the lack of a dataset.

This study aimed to investigate the role of object affordance for multimodal grasp-type recognition. To this end, we prepared a large first-person grasping image dataset containing a wider range of labeled grasp types and household objects. We

[1]Applied Robotics Research, Microsoft, Redmond, WA 98052, USA naoki.wake@microsoft.com [2]Department of Computer Science, Tokyo Institute of Technology, Meguro, Tokyo 1528550, Japan

tested the pipeline with two types of affordances, which reflect one or both of the likeliness and unlikeliness of each grasp type. The experiments showed that object affordance guides CNN recognition in two ways: 1) excluding unlikely grasp types from the candidates and 2) enhancing likely grasp types among the candidates. Additionally, the "enhancing effect" was more pronounced with a greater grasp-type bias for each object in a test dataset. Furthermore, we tested the pipeline for recognizing mimed grasping images (i.e., images of a hand grasping an illusory object), assuming that a real object may be absent in some situations (e.g., teaching in MR). Similar to the experiment with real grasping images, object affordance proved to be effective for mimed grasping images. Additionally, the CNN recognition for the mimed images exhibited lower performance compared to its recognition for real grasping, indicating the importance of the presence of real objects in image-based recognition.

The contributions of this study are 1) demonstrating the effectiveness of the object affordance in guiding grasp-type recognition both with and without the real objects in images, 2) demonstrating the conditions under which the merits of object affordance are pronounced, and 3) providing a dataset of first-person grasping images labeled with the possible grasp types for each object.

The remainder of this paper is organized as follows. Section II provides an overview of the proposed pipeline and related works. Section III describes the experiments conducted with and without real objects. Finally, Section IV summarizes the results of the study and describes future work.

## II. System overview

### A. Grasp taxonomy

There are two main approaches to analyzing human grasping from a single image: 1) using hand poses of grasping [24]–[26] and 2) using a specific grasp taxonomy [17], [19], [21], [27]–[31]. Each approach has its own advantages. Hand pose analysis in 3D space enables measurement of object states, such as posture [25] and grasping area [24]. Meanwhile, taxonomy analysis enables human grasps to be represented as discrete intermediate states that focus on the pattern of the fingers in contact. This study aimed to recognize grasp types from human behavior as an extension of taxonomy-based studies. We employed the taxonomy by Feix et al., which contains 33 grasp types [32].

### B. Dataset of human grasps

Building a realistic dataset of human hand shapes while manipulating objects will contribute to the study of human grasping. Some studies collected joint positions using wired sensors [33], a data glove [34], and model fittings [26], [35]. Another study created a dataset of hand–object contact maps obtained using thermography [36]. Additionally, taxonomy-based studies have created datasets annotated with grasp types [17], [19], [29], [37]. For example, Bullock et al. collected a dataset containing first-person images of four workers [37].

Despite the variety of datasets available for grasp-type recognition, they could not be directly applied to our study because they do not aim to cover the possible grasp types associated with an object. Although there exists a pseudo-image dataset focusing on object-specific affordance [18], there is no dataset that provides the actual grasping images. In contrast, the uniqueness of the dataset in this study is that it aimed to cover the possible grasp types for each object while providing RGB images of real human grasps. Additionally, the objects were selected from common household objects (see Section III-A1 for details).

### C. Object affordance

The originality of this research is that we introduce object affordance obtained by searching a database by an object name. Although several studies have reported the effectiveness of using multi-modal cues for grasp-type recognition [17], [38], the understanding of the effectiveness of linguistically-driven object affordance is still limited in the context of multimodal robot teaching.

In concrete implementation, object affordance was represented by a dictionary with object names as keys. When an object name was input to the pipeline, the object affordance corresponding to the object name was retrieved by searching the dictionary (Fig. 1). Note that this affordance database was not acquired automatically but was assumed to be added and modified by the user according to each application.

*1) Definition of object affordance:* Prediction of affordance has become an active research topic in the cross-domain of robotics and computer vision. Affordance, which is generally regarded as an opportunity for interaction in a scene, has been defined in different ways depending on the problem to be solved. For example, in the computer vision research using deep learning, affordances have been formulated as a type of label in semantic segmentation tasks [36], [39]–[42]. In robotics research, affordance is a topic of the task-dependent object grasping problem, which is referred to as task-oriented grasping (TOG) [24]. In the context of TOG, affordance is defined as the possible tasks (e.g., cut and poke) allowed for an object [43]–[46].

In this study, object affordance was defined for each object as "a distribution of the possible grasp types associated with the object's name." This definition is similar to TOG in that it considers affordance to be object-specific. However, our definition focuses on the grasp types and does not scope the information on the possible tasks following the grasps.

*2) Types of object affordance evaluated:* The experiments in Section III evaluate the role of object affordance using sub-datasets that were sampled from the created dataset, which was labeled with the possible grasp types for each object (see Section III-A1 for details). While testing the proposed pipeline (Fig. 1), an affordance database was created for each sub-dataset based on the grasp-type labels found in the sub-dataset. We prepared two types of affordances for each object (Fig. 2):

- *Varied affordance* was calculated as a normalized histogram of the labeled grasp types for each object.
- *Uniform affordance* was calculated by flattening the non-zero values in the histogram.

While the varied affordance contains information regarding the likeliness and unlikeliness of grasping, the uniform affordance only contains information regarding the unlikeliness of grasping.

## D. Convolutional neural network with object affordance

We formulated grasp detection by fusing a CNN with object affordance (Fig. 1) as follows. The image, object name, and grasp type are denoted as $i$, $o$, and $g$, respectively. Assuming the output of the CNN to be a probability of each grasp type $g$ given an image $i$, we represent the output of the CNN as $p(g \mid i)$. Additionally, based on the definition of object affordance (i.e., a distribution of the possible grasp types associated with the object's name), we represent the object affordance as $p(g \mid o)$. Herein, we focused on deriving the probability of each grasp type given both an image and an object name (i.e., $p(g \mid i,o)$) from these conditional probabilities, $p(g \mid i)$ and $p(g \mid o)$. Assuming that $p(i)$ and $p(o)$ are independent, the following equation holds based on mathematical formulas:

$$p(g \mid i,o) = \frac{p(i,o \mid g)\,p(g)}{p(i,o)}$$
$$= \frac{p(i \mid g)\,p(o \mid g)p(g)}{p(i)\,p(o)} \qquad (1)$$
$$= \frac{p(g \mid i)\,p(g \mid o)}{p(g)}$$

Hence, the conditional probability distribution $p(g \mid i,o)$ can be estimated from the available distributions $p(g \mid i)$, $p(g \mid o)$, and $p(g)$. Finally, the grasp type can be determined as that which maximizes $p(g \mid i,o)$. A reasonable interpretation of this equation is that the grasp-type recognition based on object name and image can be approximated by a measure that considers the predictions based on the object name and image respectively, and the rarity of the grasp type (i.e., $1/p(g)$).

A CNN network was obtained by fine-tuning ResNet-101 [47]. To avoid overfitting, we applied random reflection and translation to images, and randomly shifted the image color in the hue, saturation, value space after every training epoch. The learning was conducted using the Adam optimizer and continued until the validation accuracy ceased increasing.

## III. Experiments

### A. Scenario 1: with real objects

In this scenario, the demonstration of grasping a real object was provided as a first-person image using an HMD. We assumed that the system could retrieve object affordance from the affordance database using the name of the object mentioned through verbal instructions (e.g., "Pick up the **apple**.").

*1) Data preparation:* Demonstrations are frequently recorded by an HMD in MR-based robot teaching (e.g., [10], [48]). Even for robot teaching in the physical world, first-person images provided by the demonstrator are preferred over third-person images owing to their ability to avoid self-occlusion. Therefore, we required a dataset of first-person images labeled with the possible grasp types and object names.
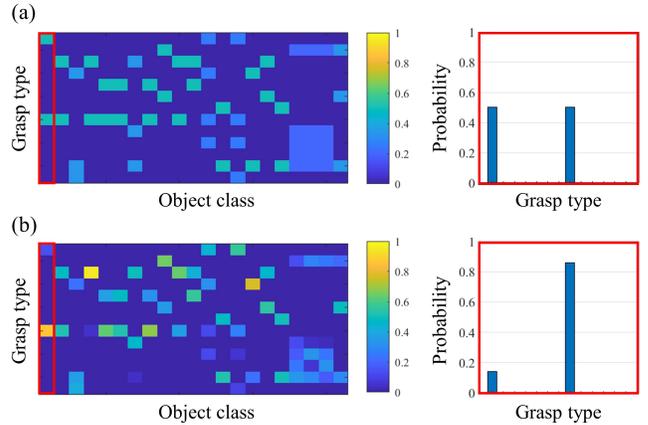


Fig. 2. Examples of object affordance calculated from a sub-dataset: (a) example of uniform affordance and (b) example of varied affordance. Refer to Fig. 3 for the order of grasp types and object classes.

Because we were not able to find any existing dataset that met these requirements, we created one.

The images were captured by a HoloLens2 sensor [49]. We used this sensor because it is a commercially available sensor that can capture first-person images without the need for hand-made attachments. The target object was chosen from the Yale-CMU-Berkeley (YCB) object set [50], which covers common household items. We employed this object set because it has been used as a benchmark for many robotic studies. We selected eight items from the food category and 13 items from the kitchen category: chip can, cracker box, gelatin box, potted meat can, apple, banana, peach, and pear; and pitcher, bleach cleanser, glass cleaner, wine glass, metal bowl, mug, abrasive sponge, cooking skillet, plate, fork, spoon, knife, and spatula, respectively. We selected these items to encompass a variety of sizes. We prepared two datasets to avoid the overestimation of the performance of the network due to CNN overfitting:

- *YCB dataset*: Training dataset containing exactly the same items as the YCB object set.
- *YCB-mimic dataset*: Testing dataset containing objects that are the same as those in the YCB dataset but differ in terms of color, texture, or shape (e.g., a cracker box from another manufacturer).

The datasets were prepared through the following pipeline. Before collecting the images, we manually assigned a set of plausible grasps according to the taxonomy in [32] (Fig. 3). Based on a previous study [51], we focused on 13 grasp types that we believed were achievable for common robot hands. For each object and grasp type, we captured images of a human grasping the object using their right hand. We captured more than 1500 grasping images by varying the arm orientation and rotation as much as possible. Furthermore, to crop the hand regions from the captured images, we applied a third-party hand detector [52] in offline. After manually filtering out the detection errors, 1000 images were randomly collected for each object and grasp type. The following experiments were conducted with sub-datasets that were sampled from the YCB

| | Chips can | Cracker box | Gelatin box | Potted meat can | Apple | Pear | Banana | Peach | Pitcher | Bleach cleanser | Glass cleaner | Wine glass | Metal bowl | Mug | Abrasive sponge | Cooking skillet | Plate | Knife | Spoon | Fork | spatula |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Large diameter | ■ | | | | | | | | | | | ■ | | ■ | | | | | | | |
| Small diameter | | | | | | | | | ■ | | | | | | ■ | | | ■ | ■ | ■ | ■ |
| Medium wrap | | ■ | | ■ | | | ■ | | ■ | | ■ | | | | | | | | | | |
| Power disk | | | ■ | | | | | | | | | ■ | | ■ | | | | | | | |
| Power sphere | | | | | ■ | ■ | | ■ | | | | ■ | | | | | | | | | |
| Adducted thumb | | | | | | | | | ■ | | | | | | ■ | | | | | | ■ |
| Extension type | | | | | | | | | | | | | ■ | | | | ■ | | | | |
| Sphere 4-finger | ■ | ■ | | ■ | ■ | ■ | | ■ | | ■ | | ■ | | ■ | | | | | | | |
| Tripod | | | | | | | ■ | | | | | | | | | | | ■ | ■ | ■ | |
| Prismatic 2-finger | | | | | | | | | | | | ■ | | ■ | | | | ■ | ■ | ■ | |
| Prismatic 3-finger | | | | | | | | | | | | | | ■ | | | | ■ | ■ | | |
| Prismatic 4-finger | | | ■ | | | | ■ | | | | | ■ | | ■ | | | ■ | ■ | ■ | ■ | ■ |
| Precision sphere | | | ■ | | | | | | | | | | | ■ | | | | | | | |

Fig. 3. Grasp types assigned to Yale-CMU-Berkeley (YCB) objects. Images were selected from the database to demonstrate examples of grasping.

or YCB-mimic dataset.

*2) Evaluation of dataset size:* Because small datasets lead to underestimation in CNN recognition, we validated the performances of the CNNs trained with different sized sub-datasets of the YCB dataset. We prepared five sub-datasets containing 10, 50, 100, 500, and 1000 images per grasp type. The images were randomly sampled such that a sub-dataset included all the images from the other smaller sub-datasets. The CNNs were tested with sub-datasets of the YCB-mimic dataset. We refer to these sub-datasets as the test datasets. The test datasets were created by randomly sampling 100 images per grasp type. The performances of the CNNs were validated ten times using different test datasets.

Fig. 4 shows the result. The CNN performance tended to increase with the size of the dataset and converged above 500 images per grasp type. This result indicates that the YCB dataset is sufficiently large to avoid underestimation due to insufficient images.

*3) Effect of affordance on recognition:* We evaluated the effectiveness of the proposed pipeline by comparing five methods: the proposed pipeline using 1) varied affordance (i.e., $p(g \mid i,o)$), 2) uniform affordance, 3) only varied affordance (i.e., $p(g \mid o)$), 4) only uniform affordance, and 5) only the CNN (i.e., $p(g \mid i)$). For fair comparison, the same CNN was used for each method. The grasp type that maximizes the probability distribution was chosen. In the case of using only the uniform affordance, the grasp type was randomly selected from the possible grasp types.

The CNN was trained with a sub-dataset of the YCB dataset. Based on the evaluation of the dataset size in Section III-A2, the sub-dataset was prepared by randomly sampling 1000 images per grasp type. We compared the performances of five methods applied to a set of 100 test datasets. Each test dataset was created by randomly sampling 100 images per object from the YCB-mimic dataset.

Fig. 5 shows the result. The pipelines combining the CNN and affordance performed better than the CNN-only and affordance-only pipelines. While the proposed affordance exhibited the best performance, the proposed pipeline using uniform affordance was comparable. These increased performances indicate the effectiveness of using affordance for guiding grasp-type recognition.

To elucidate the role of affordance, we examined the cases where the CNN failed whereas the use of varied affordance succeeded (Fig. 6). In such cases, the CNN did not output the correct grasp as the best candidate, possibly due to finger occlusion; however, it had a small affordance value, resulting in a small likelihood to be the output of the proposed method. As a result, the correct grasping was chosen as the final output. Therefore, evidently, object affordance contributed to excluding the unlikely grasp types from the candidates of the CNN.

To investigate the advantage of varied affordance over uniform affordance, we examined cases where the use of uniform affordance failed whereas the use of varied affordance succeeded (Fig. 7). In these cases, the pipeline outputted the correct grasping by employing varied affordance. Therefore, evidently, varied affordance contributed to enhancing the grasp types that were likely for an object.

*4) Enhancing effect of varied affordance:* After observing the enhancing effect of the object affordance, based on information theory, we hypothesized that the effect would be stronger with greater grasp-type bias for each object in a test dataset (i.e., grasp-type heterogeneity). To test this hypothesis, we evaluated the effect of grasp-type heterogeneity on recognition.

We used the same 100 test datasets that were prepared for the comparison experiment. The degree of grasp-type heterogeneity, $h$, was defined for each test dataset by the following equation:

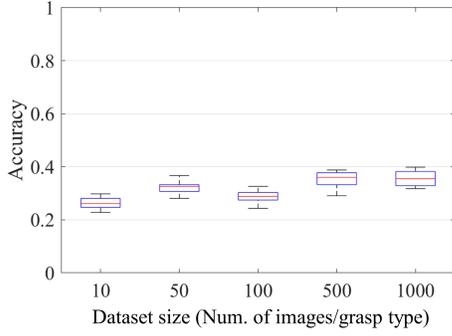$$h = \frac{1}{N} \sum_{i=1}^{N} std\left(a_i\right), \tag{2}$$

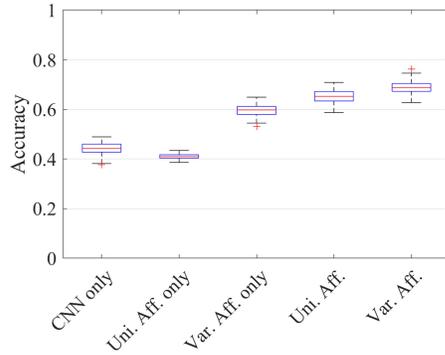Fig. 4. Performances of CNNs trained with different dataset sizes.



Fig. 5. Performance of grasp-type recognition with different pipelines: CNN only (only the CNN), Uni. Aff. only (only uniform affordance), Var. Aff. only (only varied affordance), Uni. Aff. (proposed pipeline using uniform affordance), and Var. Aff. (proposed pipeline using varied affordance).

where $N$, $a_i$, and *std* represent the number of object classes, vector of varied affordance of an object (i.e., each column in Fig. 2 (b)), and an operation to calculate the standard deviation of the non-zero values of a vector, respectively.

The sub-datasets were tested with the proposed pipeline using varied affordance and uniform affordance. The same CNN as in the comparison experiments was used. Fig. 8 shows the difference in performance between the two affordance types plotted against the grasp-type heterogeneity. As hypoth-
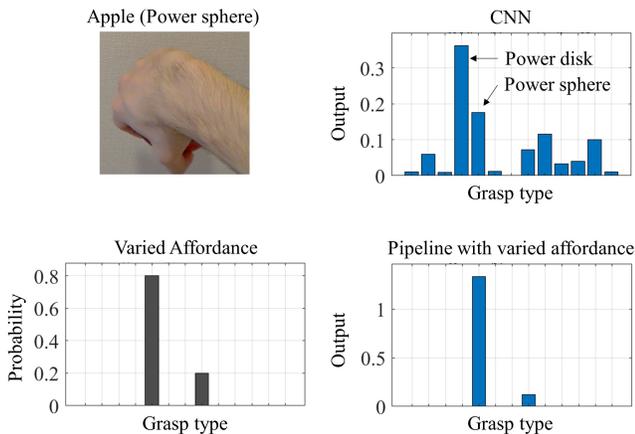


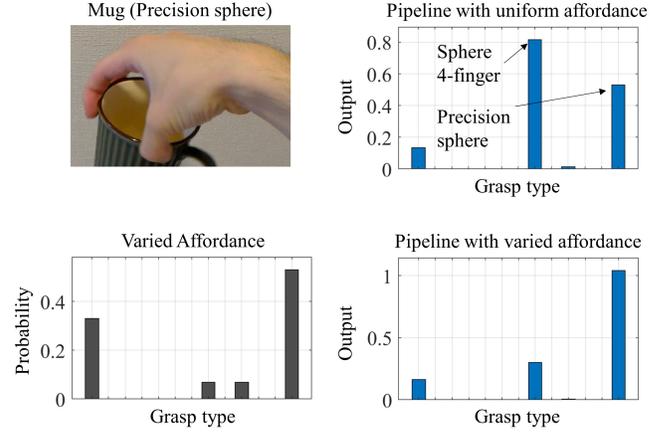Fig. 6. Example of where the CNN failed. The order of grasp types is the same as in Fig. 3.



Fig. 7. Example of where the proposed pipeline using uniform affordance failed. The order of grasp types is the same as in Fig. 3.
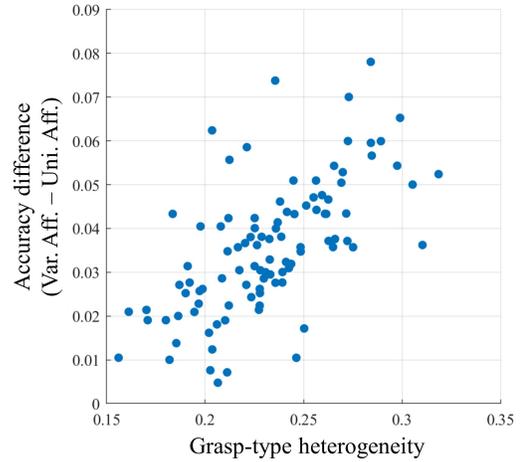


Fig. 8. Performance difference between the pipelines plotted against grasp-type heterogeneity. Each plot represents a test dataset.

esized, the difference increased with the increasing grasp-type heterogeneity. This result indicates that the enhancing effect of the varied affordance is more pronounced when the degree of grasp-type heterogeneity is higher.

### B. Scenario 2: without real objects

In the previous section, we evaluated the pipeline for images of grasping real objects. On the contrary, robot teaching may not require real objects to be grasped in some situations (e.g., teaching in MR). In such situations, the captured images do not include real objects; however, a user can interact with an illusory object in MR (i.e., an MR object). Since such "mimed" images lack visual object information, image-based grasp-type recognition can become challenging. This section provides an evaluation of the performance of the proposed pipeline when mimed images and object affordance are available.

*1) Data preparation:* To obtain the CNN for recognizing the grasp types, we prepared a dataset of the mimed images captured by a HoloLens2 sensor [49]. We used the texture-mapped 3D mesh models of the YCB objects described in Section III-A1 as MR objects. Grasp achievement was
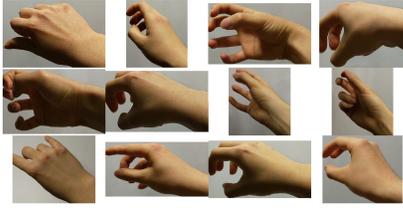
Fig. 9. Examples of the mimed images captured by the HoloLens2 sensor. Although the grasped YCB objects were not captured, they were presented to the user in MR.
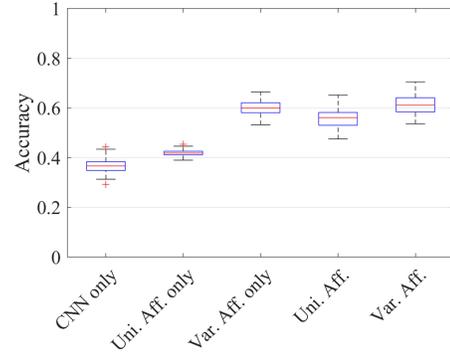


Fig. 10. Performances of grasp-type recognition with different pipelines. The contractions are the same as in Fig. 5.
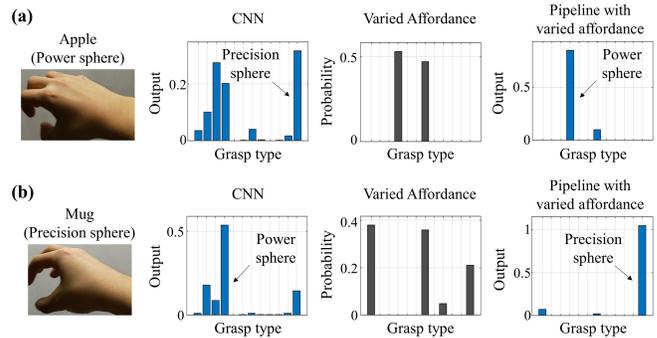


Fig. 11. Examples of the CNN failed in recognizing the mimed images. (a) Recognition of "power sphere" grasping of an illusory apple. (b) Recognition of "precision sphere" grasping of an illusory mug. The order of grasp types is the same as in Fig. 3, excluding "small diameter" grasping.

determined by the type and number of the fingers in contact, following the definition in [32]. The positions of the hand joints were estimated via the HoloLens2 API. During the collection of the images, a user grasped one of the rendered MR objects guided by visual cues that represent the contact state between the user's hand and MR object [10]. Among the object list in Fig. 3, the glass cleaner and the wine glass were ignored due to the lack of 3D models provided by [50], and the abrasive sponge was ignored due to the inability to express soft materials in MR. Furthermore, "small diameter" grasping was ignored because of the difficulty in measuring the joint positions with the corresponding accuracy (i.e., within 1 cm [32]). As the result of eliminating the "small diameter" grasping, we excluded objects with only one type of grasp (i.e., the pitcher and cooking skillet).

Following the same recording and post-processing protocol described in Section III-A1, we collected a dataset containing 1000 mimed images for each object and grasp type (note that the MR objects were not captured in the images). We created two datasets under different lighting conditions and used one for training the CNN and the other for testing the pipeline. Fig. 9 shows examples of the images.

*2) Effect of affordance on recognition:* We compared the same five methods as in Section III-A. The protocols to obtain the CNN and affordance database were the same as those described in Section II-D. That is, we compared the performances of the five methods applied to a set of 100 test datasets. Fig. 10 shows the comparison results. Similar to the results in Section III-A, the proposed pipeline exhibited the highest performance. Although the CNN recognition for mimed images was inferior to recognition for real grasping images (see Fig. 5), the use of affordance proved to be effective.

Additionally, we observed the two functions of object affordance (i.e., excluding the unlikely grasp types from the candidates and enhancing the likely grasp types among the candidates), similar to Section III-A. For example, Fig. 11 shows cases where the CNN failed to discriminate between the "power sphere" and "precision sphere," which appeared similarly in mimed grasping. Despite such similarity, the proposed pipeline using varied affordance succeeded by excluding either of them as candidates.

## IV. DISCUSSION AND CONCLUSION

### A. Summary of the experiments

This study investigated the role of object affordance in guiding grasp-type recognition. To this end, we created two first-person image datasets containing images with and without the grasped objects, respectively. The results revealed the effects of object affordance in guiding CNN recognition: 1) it excludes the unlikely grasp types from the candidates and 2) enhances the likely grasp types among the candidates. The enhancing effect was stronger when there was more heterogeneity between the likely grasp types. These findings suggest that object affordance can be effective in improving grasp-type recognition.

The advantage of our proposed pipeline (Fig. 1) is that it can be updated independently of the CNN. For example, if a user experiences a grasp type that is not assigned for an object, the pipeline can be updated by simply modifying the object affordance according to the user's feedback. As another example, if a user wants to interact with objects that are not registered in the affordance database, the pipeline can be updated by manually adding the object affordances. In the case of using uniform affordances, which showed promising results (Fig. 5), object affordances can be readily added by manually assigning the possible grasp types. Such an approach is less expensive than updating a CNN by collecting a large number of grasping images depending on the use case.

Recognition from the mimed images appears to be more difficult than that from the images of grasping real objects (Fig. 5 and 10), indicating the importance of the presence of real objects in image-based recognition. The inferior performance with mimed images is reasonable because a grasp type depends on the shape of the hand and the fingers that are in contact with the object. Recognition from the mimed images may benefit from the findings of previous studies. For example, a study proposed combining other information, such as contact points and contact normals, which can be easily calculated for MR objects [53]. It may also be possible to utilize techniques developed in other research areas, such as sign language recognition for mimed images [54], [55]. However, most importantly, object affordances can be applied to any recognition method that outputs a probability distribution (see equation 1). As long as visual ambiguities are inherently present (e.g., finger occlusion or absence of object), the proposed pipeline should be beneficial for grasp-type recognition.

### B. Methodological considerations

The proposed pipeline used a text-based database that can be added and modified by the user according to each application (see Section II-C). This approach has two limitations. First, multiple object affordances cannot be associated with one text label. This becomes a problem when a user wants to register different object affordances for objects with the same name (e.g., grasp-type A for a *cup* while grasp-type B for another *cup*). Another limitation is that manual work is required to register object affordances. However, in practical use, we do not consider these characteristics to be a critical problem. Users can address the former issue by assigning different text labels to objects with different affordances. For the latter issue, we believe that the number of objects in the home environment is finite and falls within an acceptable range.

Regarding the system input, this study assumed that the pipeline can access the name of the grasped object and retrieve the affordance using the object name. For practical robot teaching applications, separate solutions to these requirements are required. To access the name of the grasped object, general object recognition or user input information can be used. For example, our robot teaching platform is designed to extract the name of the grasped object from human instructions [12]. While this study used text matching to retrieve the affordance using object names, we could also employ a thesaurus or word embedding methods to cover the word variations.

The image datasets prepared in this study were collected against plain backgrounds under simple lighting conditions. Training images under a controlled environment often results in overfitting of a CNN and reduce its generalization performance. However, we consider the effect of using the controlled images to be limited for the following reasons. First, the CNNs were trained with cropped images with minimum background reflections (see Fig. 3 and 9). Second, to mitigate the lighting biases, we randomly shifted the colors of the training images. The effect of these pre-processing steps could be supported by the fact that the recognition performances of a single CNN were much higher than the chance rate (see Fig. 5 and 10). This study aimed to investigate the role of object affordance and did not scope the improvement of the generalization performance of CNNs; given the reasonable performance of the CNNs, the environmental condition was not critical to the paper's argument.

### C. Future studies

As future research, the proposed pipeline could be employed in a learning-from-observation (LfO) framework, where the object names can be estimated from verbal instructions. We are currently testing this hypothesis by integrating the pipeline with an LfO system that we developed in-house [11], [12].

## REFERENCES

[1] M. R. Cutkosky and R. D. Howe, "Human grasp choice and robotic grasp analysis," in *Dextrous robot hands*, pp. 5–31, Springer, 1990.

[2] S. B. Kang and K. Ikeuchi, "Toward automatic robot instruction from perception-mapping human grasps to manipulator grasps," *IEEE Transactions on Robotics and Automation*, vol. 13, no. 1, pp. 81–95, 1997.

[3] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *The International Journal of Robotics Research*, vol. 34, no. 4-5, pp. 705–724, 2015.

[4] D. Morrison, P. Corke, and J. Leitner, "Closing the loop for robotic grasping: A real-time, generative grasp synthesis approach," in *Robotics: Science and Systems Conference (RSS)*, pp. 1–10, 2018.

[5] Y. Jiang, S. Moseson, and A. Saxena, "Efficient grasping from rgbd images: Learning using a new rectangle representation," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3304–3311, IEEE, 2011.

[6] K. Bousmalis, A. Irpan, P. Wohlhart, Y. Bai, M. Kelcey, M. Kalakrishnan, L. Downs, J. Ibarz, P. Pastor, K. Konolige, *et al.*, "Using simulation and domain adaptation to improve efficiency of deep robotic grasping," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4243–4250, IEEE, 2018.

[7] J. Redmon and A. Angelova, "Real-time grasp detection using convolutional neural networks," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1316–1322, IEEE, 2015.

[8] Q. Yu, W. Shang, Z. Zhao, S. Cong, and Y. Lou, "Robotic grasping of novel objects from rgb-d images by using multi-level convolutional neural networks," in *Proceedings of the IEEE International Conference on Information and Automation (ICIA)*, pp. 341–346, IEEE, 2018.

[9] O. Tutsoy, D. E. Barkana, and K. Balikci, "A novel exploration-exploitation-based adaptive law for intelligent model-free control approaches," *IEEE Transactions on Cybernetics*, 2021.

[10] D. Saito, N. Wake, K. Sasabuchi, H. Koike, and K. Ikeuchi, "Contact web status presentation for freehand grasping in mr-based robot-teaching," in *Companion of the 2021 ACM/IEEE International Conference on Human-Robot Interaction*, pp. 167–171, 2021.

[11] N. Wake, R. Arakawa, I. Yanokura, T. Kiyokawa, K. Sasabuchi, J. Takamatsu, and K. Ikeuchi, "A learning-from-observation framework: One-shot robot teaching for grasp-manipulation-release household operations," in *Proceedings of the IEEE/SICE International Symposium on System Integration (SII)*, pp. 461–466, IEEE, 2021.

[12] N. Wake, I. Yanokura, K. Sasabuchi, and K. Ikeuchi, "Verbal focus-of-attention system for learning-from-demonstration," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2021. Accepted.

[13] K. Sasabuchi, N. Wake, and K. Ikeuchi, "Task-oriented motion mapping on robots of various configuration using body role division," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 413–420, 2020.

[14] D. Saito, K. Sasabuchi, N. Wake, J. Takamatsu, H. Koike, and K. Ikeuchi, "Task-grasping from human demonstration," in *Proceedings of the IEEE-RAS International Conference on Humanoid Robotics (Humanoids)*, IEEE, 2022.

[15] P. Milgram and F. Kishino, "A taxonomy of mixed reality visual displays," *IEICE TRANSACTIONS on Information and Systems*, vol. 77, no. 12, pp. 1321–1329, 1994.

[16] N. Wake, K. Sasabuchi, and K. Ikeuchi, "Grasp-type recognition leveraging object affordance," in *HOBI Workshop, IEEE International Symposium on Robot and Human Interactive Communication (RO-MAN)*, 2020.

[17] M. Cai, K. Kitani, and Y. Sato, "Understanding hand-object manipulation by modeling the contextual relationship between actions, grasp types and object attributes," in *Robotics: Science and Systems Conference (RSS)*, pp. 1–10, 2016.

[18] E. Corona, A. Pumarola, G. Alenya, F. Moreno-Noguer, and G. Rogez, "Ganhand: Predicting human grasp affordances in multi-object scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 5031–5041, 2020.

[19] G. Rogez, J. S. Supancic, and D. Ramanan, "Understanding everyday hands in action from rgb-d images," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 3889–3897, 2015.

[20] H. B. Helbig, J. Steinwender, M. Graf, and M. Kiefer, "Action observation can prime visual object recognition," *Experimental brain research*, vol. 200, no. 3, pp. 251–258, 2010.

[21] T. Feix, I. M. Bullock, and A. M. Dollar, "Analysis of human grasping behavior: Correlating tasks, objects and grasps," *IEEE Transactions on Haptics*, vol. 7, no. 4, pp. 430–441, 2014.

[22] F. Cini, V. Ortenzi, P. Corke, and M. Controzzi, "On the choice of grasp type and location when handing over an object," *Science Robotics*, vol. 4, no. 27, 2019.

[23] J. J. Gibson and L. Carmichael, *The senses considered as perceptual systems*, vol. 2. Houghton Mifflin Boston, 1966.

[24] M. Kokic, D. Kragic, and J. Bohg, "Learning task-oriented grasping from human activity datasets," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 3352–3359, 2020.

[25] M. Kokic, D. Kragic, and J. Bohg, "Learning to estimate pose and shape of hand-held objects from rgb images," in *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 3980–3987, IEEE, 2019.

[26] Y. Hasson, G. Varol, D. Tzionas, I. Kalevatykh, M. J. Black, I. Laptev, and C. Schmid, "Learning joint reconstruction of hands and manipulated objects," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11807–11816, 2019.

[27] M. Cai, K. M. Kitani, and Y. Sato, "A scalable approach for understanding the visual structures of hand grasps," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1360–1366, IEEE, 2015.

[28] D.-A. Huang, M. Ma, W.-C. Ma, and K. M. Kitani, "How do we use our hands? discovering a diverse set of common grasps," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 666–675, 2015.

[29] A. Saudabayev, Z. Rysbek, R. Khassenova, and H. A. Varol, "Human grasping database for activities of daily living with depth, color and kinematic data streams," *Scientific data*, vol. 5, no. 1, pp. 1–13, 2018.

[30] T. Feix, I. M. Bullock, and A. M. Dollar, "Analysis of human grasping behavior: Object characteristics and grasp type," *IEEE Transactions on Haptics*, vol. 7, no. 3, pp. 311–323, 2014.

[31] V. Arapi, C. Della Santina, G. Averta, A. Bicchi, and M. Bianchi, "Understanding human manipulation with the environment: a novel taxonomy for video labelling," *IEEE Robotics and Automation Letters*, vol. 6, no. 4, pp. 6537–6544, 2021.

[32] T. Feix, J. Romero, H.-B. Schmiedmayer, A. M. Dollar, and D. Kragic, "The grasp taxonomy of human grasp types," *IEEE Transactions on Human-Machine Systems*, vol. 46, no. 1, pp. 66–77, 2015.

[33] G. Garcia-Hernando, S. Yuan, S. Baek, and T.-K. Kim, "First-person hand action benchmark with rgb-d videos and 3d hand pose annotations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 409–419, 2018.

[34] Y. Lin and Y. Sun, "Grasp planning based on strategy extracted from demonstration," in *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4458–4463, IEEE, 2014.

[35] S. Hampali, M. Rad, M. Oberweger, and V. Lepetit, "Honnotate: A method for 3d annotation of hand and object poses," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3196–3206, 2020.

[36] S. Brahmbhatt, C. Ham, C. C. Kemp, and J. Hays, "Contactdb: Analyzing and predicting grasp contact via thermal imaging," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8709–8719, 2019.

[37] I. M. Bullock, T. Feix, and A. M. Dollar, "The yale human grasping dataset: Grasp, object, and task data in household and machine shop environments," *The International Journal of Robotics Research*, vol. 34, no. 3, pp. 251–255, 2015.

[38] Y. Yang, Y. Li, C. Fermuller, and Y. Aloimonos, "Robot learning manipulation action plans by" watching" unconstrained videos from the world wide web," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 29, 2015.

[39] T.-T. Do, A. Nguyen, and I. Reid, "Affordancenet: An end-to-end deep learning approach for object affordance detection," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5882–5889, IEEE, 2018.

[40] L. Porzi, S. R. Bulo, A. Penate-Sanchez, E. Ricci, and F. Moreno-Noguer, "Learning depth-aware deep representations for robotic perception," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 468–475, 2016.

[41] M. Lau, K. Dev, W. Shi, J. Dorsey, and H. Rushmeier, "Tactile mesh saliency," *ACM Transactions on Graphics (TOG)*, vol. 35, no. 4, pp. 1–11, 2016.

[42] A. Roy and S. Todorovic, "A multi-scale cnn for affordance segmentation in rgb images," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 186–201, Springer, 2016.

[43] M. Kokic, J. A. Stork, J. A. Haustein, and D. Kragic, "Affordance detection for task-specific grasping using deep learning," in *Proceedings of the IEEE-RAS International Conference on Humanoid Robotics (Humanoids)*, pp. 91–98, IEEE, 2017.

[44] K. Fang, Y. Zhu, A. Garg, A. Kurenkov, V. Mehta, L. Fei-Fei, and S. Savarese, "Learning task-oriented grasping for tool manipulation from simulated self-supervision," *The International Journal of Robotics Research*, vol. 39, no. 2-3, pp. 202–216, 2020.

[45] D. Song, C. H. Ek, K. Huebner, and D. Kragic, "Task-based robot grasp planning using probabilistic inference," *IEEE Transactions on Robotics*, vol. 31, no. 3, pp. 546–561, 2015.

[46] J. Bohg, A. Morales, T. Asfour, and D. Kragic, "Data-driven grasp synthesis—a survey," *IEEE Transactions on Robotics*, vol. 30, no. 2, pp. 289–309, 2013.

[47] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.

[48] Z. Guan, Z. Liu, Y. Li, X. Hong, B. Hu, and C. Xu, "A novel robot teaching system based on augmented reality," in *2019 International Conference on Image and Video Processing, and Artificial Intelligence*, vol. 11321, pp. 304–309, SPIE, 2019.

[49] Microsoft, "Microsoft hololens." https://www.microsoft.com/en-us/hololens, Accessed: 2021-Sep-06.

[50] B. Calli, A. Singh, A. Walsman, S. Srinivasa, P. Abbeel, and A. M. Dollar, "The ycb object and model set: Towards common benchmarks for manipulation research," in *Proceedings of the International Conference on Advanced Robotics (ICAR)*, pp. 510–517, IEEE, 2015.

[51] Y. Lin and Y. Sun, "Robot grasp planning based on demonstrated grasp strategies," *The International Journal of Robotics Research*, vol. 34, no. 1, pp. 26–42, 2015.

[52] Jsk-ros-pkg, "ssd_object_detector." https://github.com/jsk-ros-pkg/jsk_recognition, Accessed: 2021-Sep-06.

[53] J. Aleotti and S. Caselli, "Grasp recognition in virtual reality for robot pregrasp planning by demonstration," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2801–2806, IEEE, 2006.

[54] M. Al-Qurishi, T. Khalid, and R. Souissi, "Deep learning for sign language recognition: Current techniques, benchmarks, and open issues," *IEEE Access*, 2021.

[55] A. Wadhawan and P. Kumar, "Sign language recognition systems: A decade systematic literature review," *Archives of Computational Methods in Engineering*, vol. 28, no. 3, pp. 785–813, 2021.