

Learning Novel Objects for Domestic Service Robots

Muhammad Attamimi¹, Tomoaki Nakamura¹, Takayuki Nagai¹,
Komei Sugiura² and Naoto Iwahashi²

¹*The University of Electro-Communications*

²*National Institute of Information and Communications Technology
Japan*

1. Introduction

It is fair to say that robots which can interact with and serve humans especially in the domestic environment will spread widely in near future. A fundamental task called mobile manipulation is required for such domestic service robots. Therefore, many humanoid robots have been developed with the ability of mobile manipulation (1–5). Recently, competitions such as RoboCup@Home (6), Mobile Manipulation Challenge (7), and Semantic Robot Vision Challenge (8), have been proposed to evaluate such robots.

Since the tasks are implemented on domestic service robots, it stands to reason that natural interaction such as speech instruction should be used for the mobile manipulation. Here, we focus on the mobile manipulation using natural speech instruction such as “Bring me X” (X is an out-of-vocabulary (OOV) word). In order to realize this task, the integration of navigation, manipulation, speech recognition, and image recognition is required.

Image and speech recognition are difficult especially when novel objects are involved in the system. For example, there are objects specific to each home and new products can be brought into the home. It is impossible to register the names and images of all these objects with the robot in advance. Hence, we propose a method for learning novel objects with a simple procedure.

The robot, on which the proposed learning method is implemented, is intended to be used in a private domestic environment. Therefore, the procedure of teaching objects to the robot must be simple. For example, the user says, “This object is X” (X is the name of the object) and shows the object to the robot (Fig.1: Left). It is easy for a user to teach a robot many objects with this procedure. Then the user orders the robot to bring him/her something. For example, the user says, “Bring me X” (Fig.1: Right). As we mentioned earlier, such extended manipulation tasks are necessary for domestic service robots. However, there are three problems in teaching novel objects to the robots. The first problem is speech recognition of an object’s name. In usual methods, phonemes of names must be registered in an internal dictionary. However, it is impossible to register all objects in advance. The second problem is the speech synthesis. A robot must utter the name of the recognized object for interaction with humans such as “Is it X?” However, conventional robot utterance systems cannot utter a word which is not registered in the dictionary. Even if the phoneme sequence of an OOV word can be recognized,

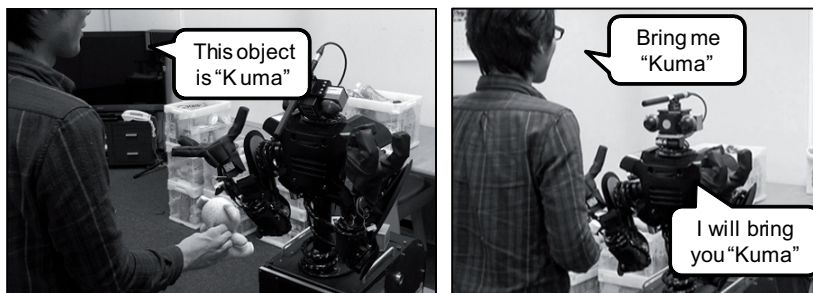


Fig. 1. Left: The user teaches the object to the robot. Right: The robot recognizes and utters the OOV word.

it cannot be used for speech synthesis since accuracy of phoneme recognition is less than 90%. The third problem is segmentation of the object region from a scene in the learning phase. When a robot learns an object, it must find where the object region is in the scene and segment it.

Methods for extracting OOV words in a speech have been proposed (9) for solving the first problem. Phonemes of OOV words can be obtained with these methods but they are not always correct. To solve the second problem, the user is required to restate the OOV word again and again so that correct phonemes are obtained (10). The robot can utter the word correctly but it is best that the robot learns the word from one user's utterance. There are also methods for situations in which the correct phonemes are not obtained. With such methods, the user utters the spelling of the OOV word to correct the phonemes (11). However, this requires a long time for the robot to learn OOV words by recognizing their spelling in Japanese or Chinese.

Considering these problems, we propose a system shown in Fig.2. We solve the first problem by extracting OOV words from a template sentence. The second problem may be solved by uttering phonemes of OOV words using a text-to-speech (TTS) system. However, it is difficult to recognize phonemes correctly. In the proposed method, the OOV part of the user's speech is converted to the robot's voice by voice conversion using Eigenvoice Gaussian mixture models (EGMMs) (12).

There has been research on the segmentation of images (13–15) for solving the third problem. The method developed by Rother *et al.* (13) requires a rough hand-drawn region of an object. Shi and Malik's method (14) can segment images automatically, but it cannot determine which segment is the object region. Mishra and Aloimonos's method (15) can segment the object accurately using color, 3D information, and motion. However, an initial point that locates inside the object region must be specified.

On the other hand, an object, which can be assumed as a human movements (e.g. hand movements), can be extracted from complicated scene because the proposed method is designed for a human to teach a robot an object. A color histogram and scale invariant feature transform (SIFT) are computed from extracted objects and registered in a database. This information is used for object recognition.

We implement the proposed method on a robot called "DiGORO". We believe it is important to evaluate the robot in a realistic domestic environment with a realistic task. When the robot moves to the object, it does not always arrive at an ideal position nor angle, and the illumination changes according to the position. A system is needed to work well in such

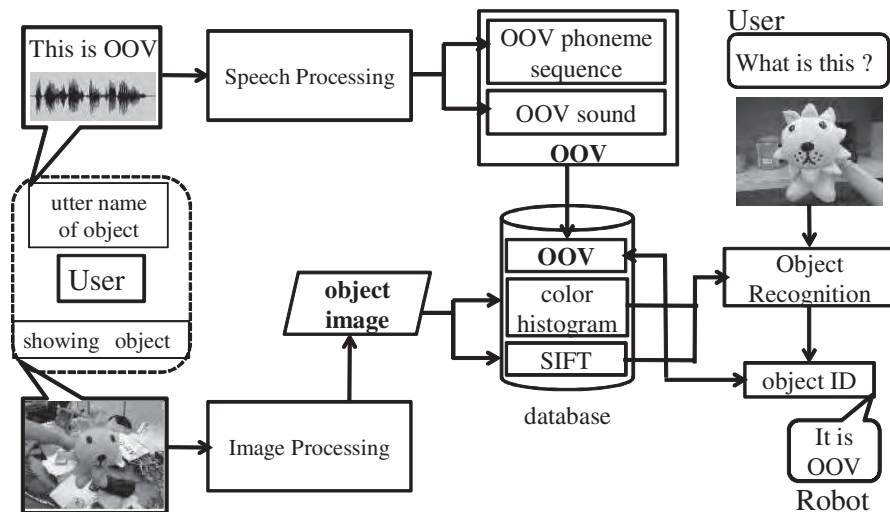


Fig. 2. Overview of learning novel objects.

an environment. In this research, we used the “Supermarket” task of RoboCup@Home (6) as the extended mobile manipulation task. RoboCup@Home is a competition that tests the ability of robots in a domestic environment. Supermarket is a standardized task based on the fetch-and-carry operation. There are also other tasks which can be used for evaluation (7; 8). The “Semantic Robot Vision Challenge” (8) evaluates the ability of a robot to find an object in a real environment. However, only three teams participated in the 2009 competition. Furthermore, Semantic Robot Vision Challenge is not for evaluating manipulation. The “Mobile Manipulation Challenge” was held at the 2010 International Conference on Robotics and Automation. Even this competition evaluates the mobile manipulation ability of robots, only four teams participated. It is difficult to determine what task should be used for evaluating robots, even though there are tasks (6–8) for it. We used one of the tasks of RoboCup@Home, which we believe is the most standard. RoboCup@Home has the largest number of participants¹ and has clearly-stated rules, which are open to the public. Besides, the rules are improved every year. From these reasons, such tasks are better than self-defined ones.

This chapter is organized as follows: the following section discusses a method for finding novel objects in cluttered scene. Then, the idea of pronouncing out-of-vocabulary words using voice conversion will be discussed in section 3. In section 4, the procedure of extended mobile manipulation task is described. Next section will discuss some experimental results to

¹ 24 teams participated in 2010 RoboCup@Home competition (6). On the other hand, a few teams participated in Mobile Manipulation Challenge (7), and Semantic Robot Vision Challenge (8).

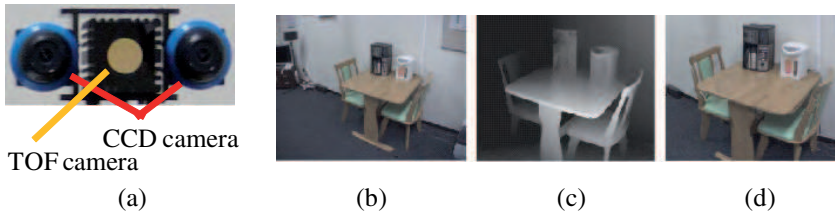


Fig. 3. (a)3D visual sensor. (b)color image (1024×768). (c)depth image (176×144). (d)mapped color image (176×144).

validate the proposed system. Discussion about proposed method will be done in section 6 followed by conclusion of this chapter in 7.

2. Finding novel objects in cluttered scene

2.1 3D visual sensor

Figure 3 shows the visual sensor used in this chapter. This sensor is able to acquire color and accurate depth information in real time by calibrating a TOF and two CCD cameras.

The distance measurement capability of TOF camera is based on the TOF principle. In TOF systems, the time taken for light to travel from an active illumination source to the objects in the field of view and return to the sensor is measured. A TOF camera SwissRanger SR4000 (23) is used as a part of 3D visual sensor. It emits a modulated near-infrared (NIR) and the CMOS/CCD imaging sensor measures the phase delay of the returned modulated signal at each pixel. These measurements in the sensor result in a 176×144 pixel depth map.

In the geometric camera calibration, the parameters that express camera pose and properties can be classified into extrinsic parameters (i.e. rotation and translation) and intrinsic ones (i.e. focal length, coefficient of lens distortion, optical center and pixel size). The extrinsic parameters represent camera position and pose in 3D space, while the intrinsic parameters are needed to project a 3D scene onto the 2D image plane. We use Zhang's calibration method in our proposed system, since the technique only requires the camera to observe a checkerboard pattern shown at a few different orientations. For the calibration of TOF camera, the reflected signal amplitude can be used to observe the checkerboard pattern. Therefore, it is straightforward to apply the same calibration method. Figure 3 (b), (c) and (d) show images captured from the visual sensor.

2.2 Motion attention based object segmentation

Assuming a user shows a target object to the robot, there may be people, objects, or furniture behind that object. The problem is object segmentation in such a complex background. Because the user has the object at hand, the object can be segmented out by taking into account the motion cue. This fact motivates us to use object segmentation based on motion attention. Figure 4 shows an overview of motion attention. A motion detector first extracts the initial object region $M(x,y)$. Then, object information, such as color (hue) image $H(x,y)$ and depth image $D(x,y)$, is taken from the region. In particular, a hue histogram $f_H(h)$ and depth histogram $f_D(d)$ are taken from the region and normalized. Here, h and d represent the quantized value of hue and depth, respectively. Since these two histograms can be considered as probability density functions of the target object, the object probability map

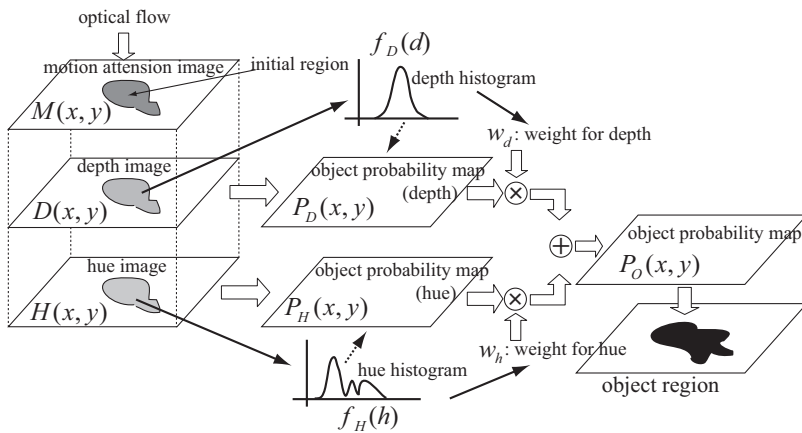


Fig. 4. Segmentation of object region using motion attention.

of each component ($P_D(x, y)$ and $P_H(x, y)$) at each pixel location can be easily obtained.

$$P_D(x, y) = f_D(D(x, y)), \tag{1}$$

$$P_H(x, y) = f_H(H(x, y)). \tag{2}$$

The weighted sum of these two object probability maps results in the object probability map $P_O(x, y)$.

$$P_O(x, y) = LPF[w_d P_D(x, y) + w_h P_H(x, y)], \tag{3}$$

The weights w_d and w_h are automatically assigned inversely proportional to the variance of each histogram. If the variance of the histogram is larger, its information is considered as inaccurate and the weight decreases. LPF represents a low pass filter, and we use a simple 3×3 averaging filter for it. The map is binarized, and then a final object mask is obtained using the connected component analysis. In the learning phase, object images are simply collected, then color histograms and SIFT features are extracted. These are used for object detection and recognition.

2.3 Object detection and identification in recognition phase

When the robot recognizes an object, the target object should be extracted from the scene. However, the same method in the learning phase is not applicable because the object is placed somewhere and it is not held by the user. Therefore, if objects are on the table, the plane detection technique is beneficial for detecting the objects. The 3D randomized Hough transform (24) is used for fast and accurate plane detection. This plane detection method is summarized below.

1. 3D information is captured in the scene.
2. Maximum plane is detected as table top using randomized Hough transform (24).
3. The plane is removed from 3D information.
4. The remaining point is projected on the plane.

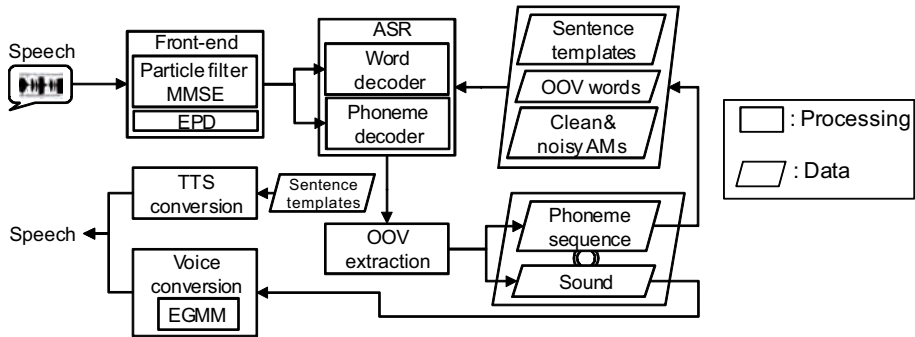


Fig. 5. Overview of the speech processing.

5. Connected components analysis is performed on the plane and each object is segmented out.

A SIFT descriptors is used for recognition. First, the candidates are narrowed down by using color information followed by the matching of SIFT descriptors, which are collected during the learning phase. It should be noted that the SIFT descriptors are extracted from multiple images taken from different viewpoints. Moreover, the number of object images is reduced for speeding up the SIFT matching process by matching among within-class object images and discarding similar ones. This process is also useful for deciding the threshold on the SIFT matching score.

3. Pronouncing out-of-vocabulary words using voice conversion

Figure 5 shows a schematic of the speech processing of the method, which uses automatic speech recognition (ASR) system called ATRASR (25). ATRASR is a hidden Markov model (HMM)-based speech recognition system, and it is used as a front-end and word/phoneme decoder. The phoneme decoder is used for obtaining the phoneme sequence of OOV words. Therefore, word- and phoneme-level speech recognition is possible.

To suppress noise, a particle filter is first applied to the online estimation of non-stationary noise, and then minimum mean square error (MMSE) estimation is used for noise reduction (26). Voice activity detection is conducted using endpoint detection (EPD) based on the frame's energy. This noise reduction part is of critical importance in RoboCup@Home tasks since the noise condition is severe.

Acoustic models (AMs) for the speech recognizer consist of "clean AMs" (male and female voices), which are trained using only clean voices, and "noisy AMs" (male and female voices), which are trained clean voices mixed with noise. This makes the speech recognition system robust in a noisy environment.

We use a template-based segmentation of words. To teach a robot an OOV word, the user is supposed to say template sentences such as "This is X". In terms of practical use, using a standard template sentence is reasonable since it is easy for users to understand how to teach a robot a word. A set of segmented voice and phoneme sequences is registered in a database. The phoneme sequence is used for utterance recognition of an OOV word.

For generating an utterance with an OOV word, the proposed method first converts the segmented voice recorded when the OOV word is learnt. The other part of the utterance is synthesized using XIMERA (27), which is a TTS conversion system. The OOV word part

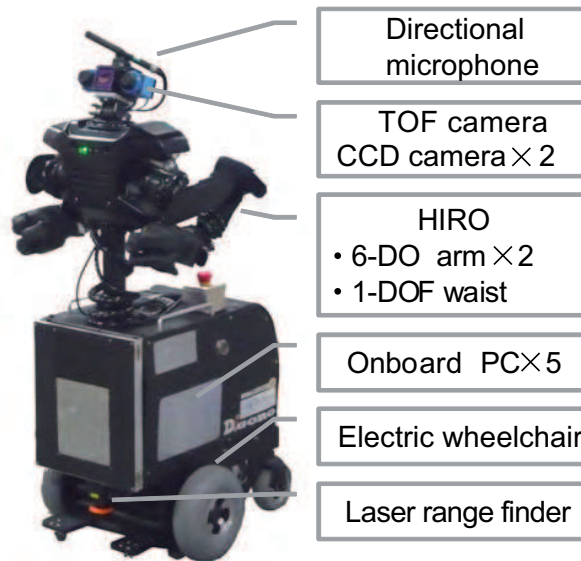


Fig. 6. The robot platform “DiGORO”.

is converted into the robot’s voice since the original sound is the user’s voice, which is not naturally concatenated with a synthesized voice. The voice conversion is based on Eigenvoice Gaussian mixture models (EGMMs) (12). The recognized phoneme sequence of the OOV word is not used for synthesis since phoneme recognition accuracy is less than 90%, and the number of utterances for teaching an OOV word is virtually constrained to one owing to the time constraint of RoboCup@Home.

4. Procedure of extended mobile manipulation task

In this section, we describe the procedure of the mobile manipulation task called “Supermarket”.

4.1 Robot platform: DiGORO

Figure 6 shows the robot “DiGORO” we previously developed (28). It is composed of the following hardware:

- Electric wheelchair
- HOKUYO laser range finder UTM-30LX
- KAWADA upper body humanoid robot
- Onboard PC (Intel Core2Duo processor) × 5
- Sanken directional microphone CS-3e
- YAMAHA loudspeaker NX-U10
- Mesa infrared TOF camera Swissranger
- Imaging Source CCD camera × 2

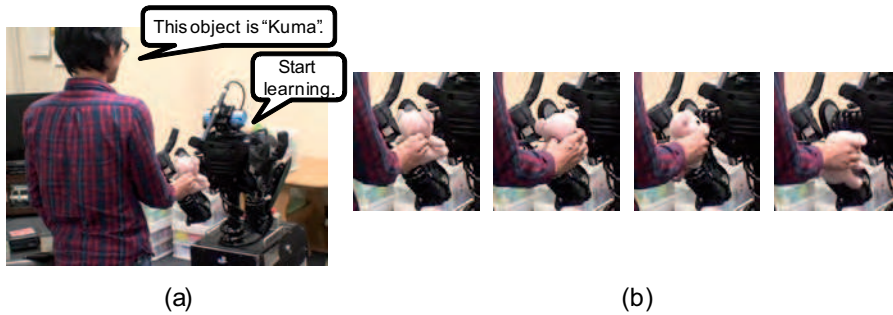


Fig. 7. Scenery of the learning phase. (a)The user gives a command for the robot to learn a new object. (b)The user shows the object from various directions.

4.2 Learning of novel object

Before the task, we need to teach objects to the robot. The procedure of teaching objects is summarized below.

1. The user shows the object to the robot and say “DiGORO, this object is X (X is an OOV word)” in Japanese ². The phoneme sequence and sound of the OOV word are extracted from the user’s speech (Fig.7(a)).
2. The robot says “I start learning of X.”
3. The user moves the object. Then 40 images of the object are segmented out and captured (Fig.7(b)).
4. Visual features (SIFT descriptors and color histogram) are calculated.
5. The phoneme sequence of the OOV word, sound of the OOV word and visual features of the object are registered in the object database. Then the robot says “I’ve memorized X.”

If the user moves the object incorrectly, the motion attention cannot segment out the object. For example, if the user moves the object from outside the observed frame into it, the unintentional region is segmented out. To avoid such a situation, we instructed the user to say “DiGORO, this is X”, while showing the object to the robot and to keep on showing it until the robot says, “I’ve memorized X.”

4.3 Supermarket task

Supermarket is a task that the robot brings three specified objects. The detailed procedure of this task is summarized below.

1. The user says “DiGORO, bring me X (X is an OOV word)” (Fig.8(a)) .
2. The robot says “I’ll bring you X. Is that correct?” Then if the user says “DiGORO, no”, go to 1. If the user says “DiGORO, yes” go to next.
3. The robot says “Where is X?”
4. The user says “It is P (P is a name of the place.)”.
5. The robot says “I’ll go to P. Is that correct?” Then if the user says “DiGORO, no”, go to 4. If the user says “DiGORO, yes” go to next.

² In this paper, the utterances are translated into English.

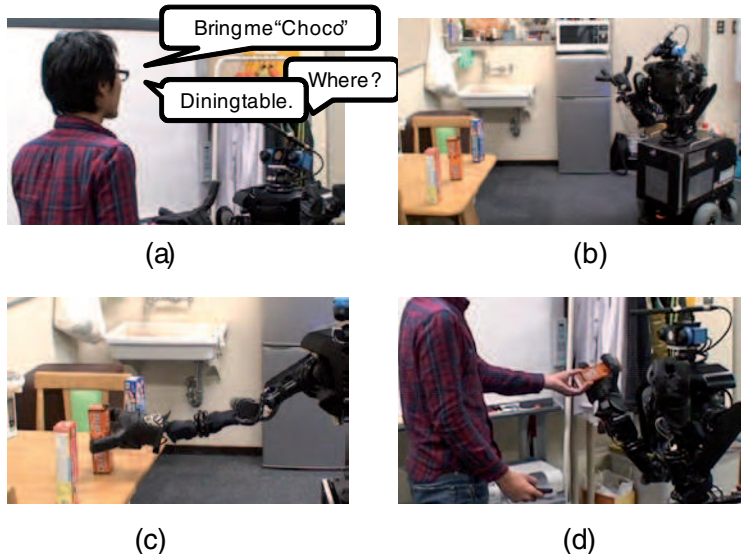


Fig. 8. Scenery of the supermarket task. (a)The user ordered the robot to bring “choco”, and the robot asked the user for needed information. (b)The robot navigated to the place where the user ordered. (c)The robot found the object and grasped it. (d)The robot returned to the start position and handed over the object to the user.

6. The robot goes to P (Fig.8(b)).
7. The robot looks around to find X using plane detection. If the robot can find X, the robot turns toward the object and says “I found X.”
8. The robot grasps the object (Fig.8(c)) and returns to the start position (Fig.8(d)).
9. The task is completed when the robot brings three objects in total. Otherwise go to 1.

5. Experiments

In this section, experiments have been conducted to evaluate the proposed method and the applied system through the task called supermarket task.

5.1 Experiment 1: Evaluation of proposed method

5.1.1 Segmentation accuracy

In this experiment, we evaluated segmentation accuracy. It is necessary to segment the object region from a complex background which generally exists in the domestic environment. We used motion attention discussed in Section 2 because the object is held by a user in the learning phase.

The experiment was carried out in an ordinary living room, shown in Fig.9. We used 120 ordinary objects, as shown in Fig.10. A user taught the robot each object by showing and telling its name to the robot. The robot acquired 40 consecutive frames for each object and extracted the target object region from each image. Figure 11 shows examples of object segmentation.



Fig. 9. Experimental Environment.



Fig. 10. 120 Objects used for experiments.

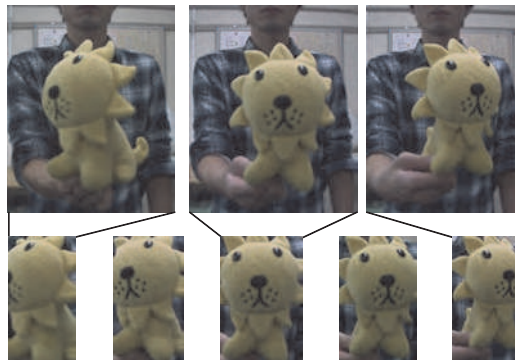


Fig. 11. Examples of object segmentation.

We extracted 10 out of the 40 frames for evaluating the segmentation accuracy of each object. Detection accuracy was measured using recall and precision rates, which are generically used for evaluation of classification, as shown in Fig.12(a), because it can be considered that pixels are classified into two classes, object region and non-object region. Here, "Object region"

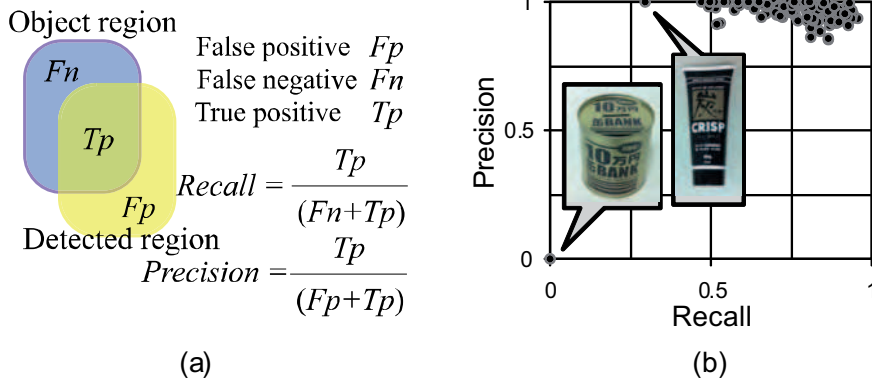


Fig. 12. (a) Definitions of recall and precision. (b) Results of object detection.



Fig. 13. Examples of object segmentation failure. (a) Object which could not be extracted. (b)Left: Object Right: Results of segmentation. (c)Left: Object Right: Results of segmentation.

	Place 1	Place 2	Place 3
Recognition rate	91%	89%	89%

Table 1. Object recognition rates.

indicates the manually labeled object region. Figure 12(b) shows a 2D plot of recall vs. precision. Each point represents an average of a single object (10 frames). As a result, the averages of all objects were 76.2% for recall and 95.8% for precision. This result indicates that the inside of object regions is extracted correctly because the precision was high. Therefore, it will not negatively affect object recognition.

Figure 13 shows examples of segmentation failure. The object in Fig.13(a) was not segmented at all because object’s entire surface reflected near infrared rays which lead to fail on measuring 3D information. A part of object in Fig.13(b) was segmented because the black region absorbed near infrared rays. Moreover, a part of the object in Fig.13(c) was segmented because near infrared rays are reflected partially. We can see that black and metallic objects tend to cause low recall.

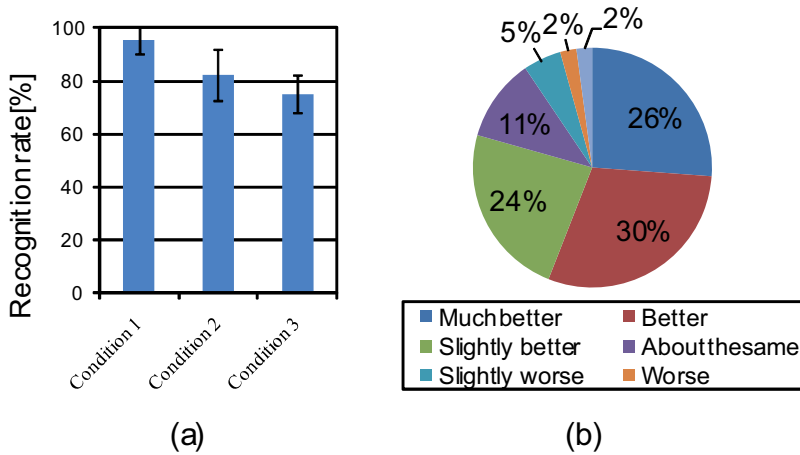


Fig. 14. (a) Recognition results. (Condition 1: Recognition with correct phonemes. Condition 2: Teacher and user are same person. Condition 3: Teachers taught the OOV words.) (b) Evaluation of voice conversion. The CMOS of VC was 1.45.

Quality	Score
Much better	3
Better	2
Slightly better	1
About the same	0
Slightly worse	-1
Worse	-2
Much worse	-3

Table 2. CMOS evaluation and scores.

5.1.2 Object recognition accuracy

We used 120 common objects, which had been learnt by the robot, as mentioned in the previous subsection. Three different locations with different lighting conditions in the living room were selected, and each object, which was segmented out using motion attention, was recognized. The results are listed in Table 1. The average recognition rate was about 90%. A major problem was false recognition between similar kinds of object such as cup noodles with different taste, because those objects have similar texture.

Next, we evaluated the proposed recognition method with the COIL100 database (29). COIL100 consist of 100 objects and 72 images per object. 36 images of each object were used for learning and the other 36 images were used for recognition. The recognition rate was 97.6%.

5.1.3 Recognition accuracy of out-of-vocabulary words

We evaluated the recognition accuracy of OOV words. The experimental procedure is described as follows. The teacher taught the robot OOV words such as “This is X”. In a domestic environment, the teacher may not be only one person but also his/her family or friends. Therefore, we conducted the experiment under the condition that OOV words are

taught by several teachers including the user who asked the robot to bring something. For comparison, we conducted the experiment under simpler conditions. In each condition, volunteers uttered sentences "Bring me X" which consist of 120 words and the robot recognized X. There were eight volunteers and 960 utterances were recognized in total. The distance between the volunteer and the microphone was 50 cm. The ambient noise level in the experiment was set as 55dBA, which simulated the standard noise level in the RoboCup@Home competition when there is no other noise source such as announcement. If the speech recognition system can work in 55dBA noise, it can also work in a domestic environment. The recognition rate was calculated from these utterances.

Figure 14(a) shows the recognition rate in each condition, and the details of each condition are as follows:

- 1. Recognition with correct phonemes:** Correct phonemes of the 120 words were manually registered in the dictionary. Each volunteer uttered "Bring me X" (X is the object name) and the robot recognized the object name.
- 2. Teacher and user is the same person:** Each volunteer uttered 120 sentences "This is X" (X is the object name) and the robot learned the 120 OOV words. The robot recognized the 120 OOV words spoken by the volunteer who was the same as the teacher.
- 3. Teachers taught OOV words:** First, 120 words were randomly assigned to eight teachers and these words were taught to the robot by them. Then, the robot recognized the 120 OOV words spoken by the user who is one of teachers. Therefore, the words were not always taught by the user. 118 out of the 960 were spoken by the teacher, i.e. teacher was the same as the user, and 842 out of the 960 utterances were spoken by others, i.e. teacher was not the same as the user.

The recognition rate was 95.2% in Condition 1, as shown in Fig.14(a). On the other hand, the accuracy of phonemes was 69.3% and the recognition rate was 82.4% in Condition 2. This indicates that the recognition rate was over 80%, which is satisfactory in a practical situation. In Condition 3, the recognition rate was 75.2%, as shown in Fig.14(a). The recognition rate was 83.4% when the teacher was the same as the user and 74.1% when the teacher was not the same as the user. Note that the speech files used in the training and those used in the test were different, even if the trainer and the tester was the same person. We can see that the recognition rate was lower than that in Condition 2. However, this is not a problem if restating is allowed.

5.1.4 Quality evaluation of robot's utterances

The objective of this experiment is to evaluate the quality of the robot's utterances. The experimental procedure is described below.

First, we made a database that included 960 utterances. It had 120 unique words and each word was uttered by eight volunteers. The ambient noise level was 55 dBA and the distance between the volunteer and microphone was 50 cm. Next, robot's utterances were generated using the proposed method. Utterances were also generated using a baseline method for comparison. These two methods are summarized as follows:

Voice Conversion (VC) (proposed): The utterances in the database are converted to robot voice by using EGMMs (12) (details of the proposed method were explained in Section 3).

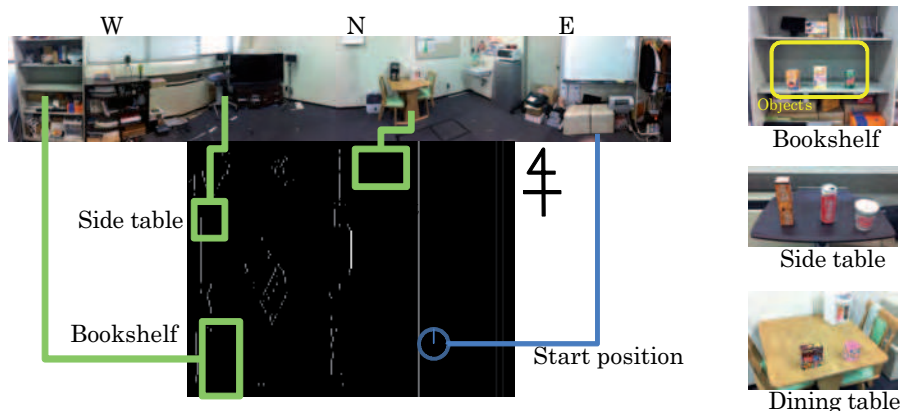


Fig. 15. The map and location of the tables/shelves.

Text-To-Speech (TTS) (baseline): The phoneme sequences obtained by phoneme recognition were used for generating robot utterances.

We then formed another group of six volunteers to evaluate the quality of generated utterances. Each volunteer listened to the utterances generated using TTS and VC. These utterances were composed of 120 unique words. The order of words was chosen at random. The order of TTS and VC samples was also chosen at random for each trial.

The comparison mean opinion score (CMOS) was used for evaluation. CMOS is specified by ITU-T recommendation P.800 (30). In the field of speech synthesis, CMOS is used for comparing voices synthesized with two methods. Specifically, the evaluation was conducted using the following questionnaire.

(Volunteer listens to two robot's utterances.) Do you think the former is more accurate than the latter in terms of pronunciation?

The evaluation and its scores are listed in Table 2.

The evaluation results are shown in Fig.14(b). The CMOS of VC was 1.45, which suggests that VC is preferred. We can see that the proposed method, which utilizes VC, is efficient if the word which has been learnt is uttered once.

5.2 Experiment 2: Evaluation of applied system in mobile manipulation

We implemented an integrated audio-visual processing system on DiGORO and performed an experiment in a living room. The purpose of this experiment is to evaluate the robot in which the proposed method has been applied in mobile manipulation. We chose a task called "Supermarket" in the RoboCup@Home league. RoboCup@Home has several advantages, that is competition that has large number of participants, and clearly-stated rules, which are open to the public. In addition, improvements on the rules are done annually.

5.2.1 Experimental setup

Figure 15 illustrates a map generated from DiGORO's own on board SLAM mapping module. The location of the tables/shelf is also shown.

We designed the task module according to the flow in Section 4.3. A volunteer first interacted with the robot at the start position. Then the robot navigated to a table/shelf, recognized the

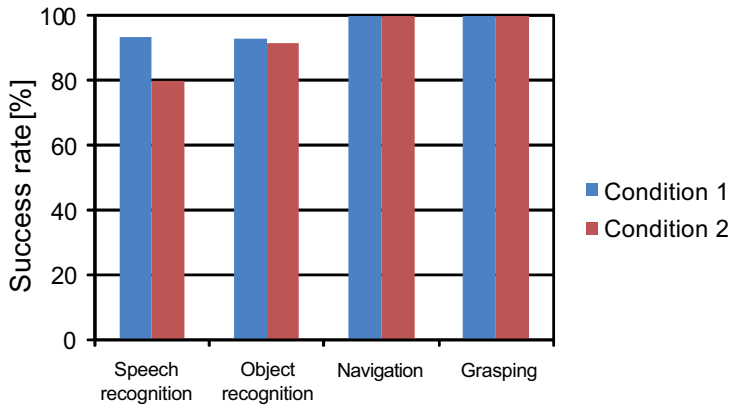


Fig. 16. Success rates. (Condition 1: words are taught by the same as requester. Condition 2: words are taught by different volunteers.)

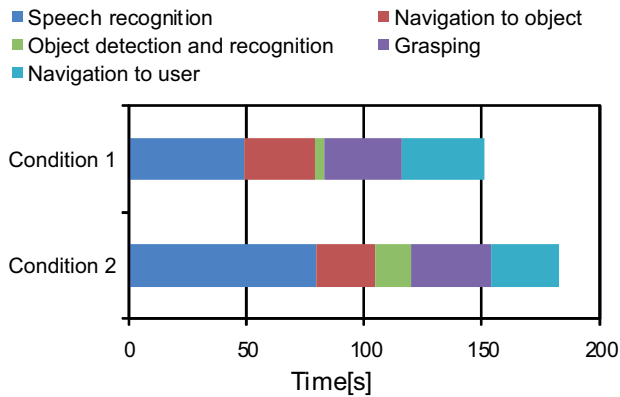


Fig. 17. Elapsed time of each process. (Condition 1: words are taught by the same as requester. Condition 2: words are taught by different volunteers.)

specified object, grasped it, and came back to the volunteer. This process was repeated for three objects.

We conducted the task under two conditions. One was similar to a real competition and the other was a more difficult condition. In each condition, we changed the dictionary of speech recognition because the user who teaches the object to the robot may not be the only person. The details of the conditions are as follows.

Condition 1: In the learning phase, each volunteer taught the robot the objects’ names. The same volunteer asked the robot to bring the objects in the execution phase.

Condition 2: In the learning phase, 120 words were randomly assigned to eight volunteers and they taught these words to the robot. Each volunteer asked the robot to bring objects in the execution phase. Therefore, the names of the objects to bring were not always taught by the same volunteer who commanded the robot in the execution phase.

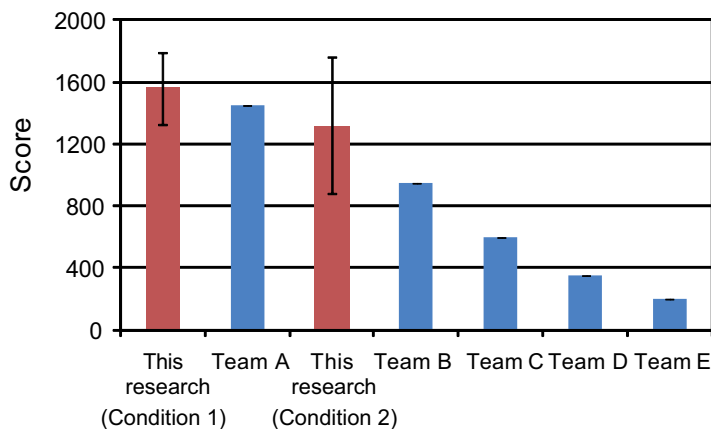


Fig. 18. Score comparison. (Condition 1: words are taught by the same as requester. Condition 2: words are taught by different volunteers.)

In the two different experimental setups, five volunteers who don't have prior knowledge of the robot conducted the task.

Therefore, the robot was supposed to bring 30 objects throughout this experiment. In each task, 30 out of the 120 objects were randomly chosen. The training data for these objects were obtained in Experiment 1.

5.2.2 Experimental results

We evaluated the results from three view points, success rate of each process, process elapsed time, and the score as a total performance.

Figure 16 shows the success rate of each process. We can see that high success rates over 90% were obtained, except for speech recognition. The speech recognition rate was 93% in Condition 1. On the other hand, it was 80% in Condition 2. This is because the phoneme sequences in the lexicon were not accurate.

Figure 17 depicts the average elapsed time for each process (per object). The results suggest that the trial can be completed within 10 min (elapsed time should be tripled and added 60 sec for the robot's instruction). The phase of instruction to the robot took a long time. The confirmation from the robot such as "I will bring X. Is this correct?" or restating the instruction to the robot such as "Bring me X" by the volunteer when the robot could not recognize the object name, were responsible for it. The instruction phase in Condition 1 was shorter than that in Condition 2 because false recognition in Condition 1 was less than that in Condition 2. This figure also shows that the time of the object recognition phase in Condition 2 was longer than that in Condition 1 because the object location was chosen randomly in both conditions. It accidentally took a long time to find objects in Condition 2, depending on their location.

Next, we evaluated the task scores as a reference. Note that the comparison of the scores may be unfair because there are differences between a laboratory and competition environments. However, we used the scores since they can be the only source for comparison among different robots through the same realistic task.

Figure 18 shows a comparison of scores among teams that participated in an actual competition in 2009. The average score in Condition 1 was 1560. From this score, DiGORO would outperform the best team in the competition. Furthermore, the average score in

Condition 2 was 1320, which was comparable to Team A. It should be noted that we used the average scores. Although, a team could performed a task only once in the competition. In that respect, this comparison may be unfair.

In an actual competition, three objects which the robot brings were selected from ten common objects whose names are listed and was given to the teams. Therefore, it was possible to manually register the names of all the objects in the dictionary. On the other hand, objects which the robot brings were chosen from 120 objects in our experiment. Moreover, no manual process was included in the learning process. Considering these conditions, we can see that our robot obtained promising results even though the environment was different from the competition.

6. Discussion

6.1 Image processing

In this section, we will discuss the results from the evaluation of segmentation accuracy. Precision was 95.8%, which indicates that the inside of the object region was extracted correctly. On the other hand, recall was 76.2%, which was less than the precision. This indicates that sometimes only part of the object region was segmented out. This is because the TOF camera could not capture 3D information due to the material of the object. For example, 3D information cannot be captured from black or metallic objects because these objects reflect or absorb near infrared rays. We believe this will be improved by using a stereo vision. DiGORO (Fig.6) has two CCD cameras and can compute stereo disparity with them.

We now discuss the results of object recognition. The object recognition rate was about 90%. We used color and SIFT features for object recognition. Generally, it is difficult to recognize objects that have the same color and/or with no textures. For future work, we plan to use an object recognition method that integrates 3D shape information (31), which can significantly improve object recognition performance.

6.2 Learning and recognition of OOV words

For this research, the robot learned OOV words from one user's utterance and it is possible for the robot to recognize and utter them. The recognition rate was 82.4% and utterance was judged as better than the baseline method, which means a practical system is constructed. Failure in recognition was because false phonemes were learnt in the learning phase. The recognition rate can be improved by a user confirmation which phonemes were learnt correctly or not after learning. For example, a user utters "This is X" and the robot learns the object. Then the user confirms which "X" can be recognized or not by asking "Did you memorize X?" If the robot utters "Yes, I memorized X'" ($X = X'$), then the OOV word is registered correctly. Otherwise, the OOV word may not be registered correctly and the user can teach the object name again to the robot.

6.3 Evaluation in domestic environment

We evaluated the system in a domestic environment using the Supermarket task, which is one of the tasks in the RoboCup@Home league. Here, let us briefly discuss the evaluation task. As we mentioned earlier, it is difficult to determine what task should be used for evaluation, and there is no global standardized tasks for this. This situation makes it very difficult to evaluate robots, which were developed by different groups, through a same realistic task. We cannot compare our robot with others by using a self-defined task, since it is almost impossible to build their robots from scratch. Therefore, we think global standardized tasks are needed.

In this research, we propose to utilize the format of the task of RoboCup@Home, since we strongly believe that the tasks are the most standard tasks for evaluating robots for the following reasons:

1. The rules are open to the public.
2. Many teams from around the world participate, i.e. the task has already been performed by many robots.
3. The rules have been improved by many robotics researchers.

Unfortunately, the comparison of the scores in the current form is not fair enough. Hence the score should be treated as reference. Although the score is just for a reference, DiGORO outperforms the best team who participated in the competition, and it shows DiGORO can perform well in a domestic environment. Any deduction in points was a result of the robot not recognizing what a user wanted it to bring. This can be improved by user confirmation in the learning phase, as mentioned above.

The learning and recognition of OOV words can be applied to other tasks. For example the "Who is who?" task, which is one of the tasks in RoboCup@Home, involves the learning of human faces and names. In this task, a user utters "My name is X" and the robot learns "X" as his/her name. With this method, we can deal with a vast number of names.

Furthermore, DiGORO has many other abilities, and it can carry out eight other tasks. For example the robot can carry out the command "Follow Me", which is for following humans, and "Shopping Mall", which is for learning the location in an unknown place. These advanced features led our team to the 1st place at RoboCup@Home 2010. This suggests that DiGORO can stably work in a domestic environment.

7. Conclusion

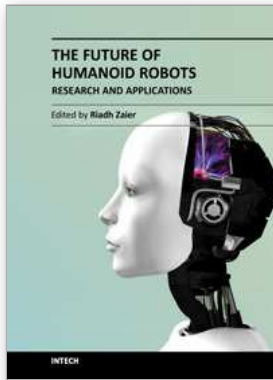
We proposed a practical learning method of novel objects. With this method a robot can learn a word from one utterance. It is possible to utter an OOV word using the segmentation of the word from a template sentence and voice conversion. The object region is extracted from a complicated scene through a user moving the object. We implemented them all in a robot as an object learning system and evaluate it by conducting the Supermarket task. The experimental results show that our robot, DiGORO, can stably work in a real environment.

8. References

- [1] T. Inamura, K. Okada, S. Tokutsu, N. Hatao, M. Inaba, and H. Inoue, "HRP-2W: A humanoid platform for research on support behavior in daily life environments," *Robotics and Autonomous Systems*, vol.57, no.2, pp.145–154, 2009.
- [2] K. WYROBEK, E. BERGER, H. VAN DER LOOS, and J. SALISBURY, "Towards a personal robotics development platform: Rationale and design of an intrinsically safe personal robot," *IEEE International Conference on Robotics and Automation*, pp.2165–2170, 2008.
- [3] F. WEISSHARDT, U. REISER, C. PARLITZ, and A. VERL, "Making High-Tech Service Robot Platforms Available," *Proceedings-ISR/ROBOTIK 2010*, pp.1115–1120, 2010.
- [4] J. STÜCKLER, and S. BEHNKE, "Integrating indoor mobility, object manipulation, and intuitive interaction for domestic service tasks," *IEEE-RAS International Conference on Humanoid Robots*, pp.506–513, 2009.

- [5] D. Holz, J. Paulus, T. Breuer, G. Giorgana, M. Reckhaus, F. Hegger, C. Müller, Z. Jin, R. Hartanto, P. Ploeger, et al., "The b-it-bots RoboCup@ Home 2009 team description paper," RoboCup 2009@ Home League Team Descriptions, Graz, Austria, 2009.
- [6] "RoboCup@Home," <http://www.ai.rug.nl/robocupathome/>, 2010.
- [7] "2010 Mobile Manipulation Challenge," <http://www.willowgarage.com/mmc10>, 2010.
- [8] "Semantic Robot Vision Challenge," <http://www.semantic-robot-vision-challenge.org/>, 2009.
- [9] I. Bazzi, and J. Glass, "A multi-class approach for modelling out-of-vocabulary words," Seventh International Conference on Spoken Language Processing, pp.1613–1616, 2002.
- [10] M. Nakano, N. Iwahashi, T. Nagai, T. Sumii, X. Zuo, R. Taguchi, T. Nose, A. Mizutani, T. Nakamura, M. Attamim, et al., "Grounding New Words on the Physical World in Multi-Domain Human-Robot Dialogues," 2010 AAAI Fall Symposium Series, pp.74–79, 2010.
- [11] H. Holzapfel, D. Neubig, and A. Waibel, "A dialogue approach to learning object descriptions and semantic categories," Robotics and Autonomous Systems, vol.56, no.11, pp.1004–1013, 2008.
- [12] T. Toda, Y. Ohtani, and K. Shikano, "One-to-Many and Many-to-One Voice Conversion Based on Eigenvoices," IEEE International Conference on Acoustics, Speech and Signal Processing, vol.4, pp.1249–1252, 2007.
- [13] C. Rother, V. Kolmogorov, and A. Blake, "Grabcut: Interactive foreground extraction using iterated graph cuts," ACM Transactions on Graphics (TOG), vol.23, no.3, pp.309–314, 2004.
- [14] J. Shi, and J. Malik, "Normalized cuts and image segmentation," IEEE Transactions on Pattern Analysis and Machine Intelligence, vol.22, no.8, pp.888–905, 2002.
- [15] A.K. Mishra, and Y. Aloimonos, "Active Segmentation," International Journal of Humanoid Robotics, vol.6, pp.361–386, 2009.
- [16] S. Hasler, H. Wersing, S. Kirstein, and E. Körner, "Large-Scale Real-Time Object Identification Based on Analytic Features," Artificial Neural Networks–ICANN 2009, pp.663–672, 2009.
- [17] H. Kim, E. Murphy-Chutorian, and J. Triesch, "Semi-autonomous learning of objects," Computer Vision and Pattern Recognition Workshop, p.145, 2006.
- [18] H. Wersing, S. Kirstein, M. Gotting, H. Brandl, M. Dunn, I. Mikhailova, C. Goerick, J. Steil, H. Ritter, and E. Korner, "Online learning of objects in a biologically motivated visual architecture," International Journal of Neural Systems, vol.17, no.4, pp.219–230, 2007.
- [19] N. Iwahashi, "Robots that learn language: Developmental approach to human-machine conversations," Symbol Grounding and Beyond, pp.143–167, 2006.
- [20] D. Roy, "Grounding words in perception and action: computational insights," Trends in Cognitive Sciences, vol.9, no.8, pp.389–396, 2005.
- [21] M. Fujita, R. Hasegawa, T. Takagi, J. Yokono, and H. Shimomura, "An autonomous robot that eats information via interaction with humans and environments," IEEE International Workshop on Robot and Human Interactive Communication, pp.383–389, 2002.
- [22] M. Johnson-Roberson, G. Skantze, J. Bohg, J. Gustafson, R. Carlson, and D. Kragic, "Enhanced Visual Scene Understanding through Human-Robot Dialog," 2010 AAAI Fall Symposium on Dialog with Robots, pp.143–144, 2010.

- [23] "Mesa imaging," <http://www.mesa-imaging.ch/index.php>.
- [24] K. Okada, S. Kagami, M. Inaba, and H. Inoue, "Plane segment finder: Algorithm, implementation and applications," IEEE International Conference on Robotics and Automation, vol.2, pp.2120–2125, 2005.
- [25] S. Nakamura, K. Markov, H. Nakaiwa, G. Kikui, H. Kawai, T. Jitsuhiro, J. Zhang, H. Yamamoto, E. Sumita, and S. Yamamoto, "The ATR multilingual speech-to-speech translation system," IEEE Transactions on Audio, Speech, and Language Processing, vol.14, no.2, pp.365–376, 2006.
- [26] M. Fujimoto, and S. Nakamura, "Sequential non-stationary noise tracking using particle filtering with switching dynamical system," IEEE International Conference on Acoustics, Speech and Signal Processing, vol.1, pp.769–772, 2006.
- [27] H. Kawai, T. Toda, J. Ni, M. Tsuzaki, and K. Tokuda, "XIMERA: A new TTS from ATR based on corpus-based technologies," Fifth ISCA Workshop on Speech Synthesis, pp.179–184, 2004.
- [28] H. Okada, T. Omori, N. Iwahashi, K. Sugiura, T. Nagai, N. Watanabe, A. Mizutani, T. Nakamura, and M. Attamimi, "Team eR@sers 2009 in the @home league team description paper," , 2009.
- [29] S.A. Nene, S.K. Nayar, and H. Murase, "Columbia Object Image Library (COIL-100)," Technical report, Feb. 1996.
- [30] International Telecommunication Union, "ITU-T P.800," <http://www.itu.int/rec/T-REC-P.800/en>.
- [31] M. Attamimi, A. Mizutani, T. Nakamura, T. Nagai, K. Funakoshi, , and M. Nakano, "Real-Time 3D Visual Sensor for Robust Object Recognition," IEEE/RSJ International Conference on Intelligent Robots and Systems, pp.4560–4565, 2010.
- [32] RoboCup@Home league committee, "RoboCup@Home Rules & Regulations," http://www.ai.rug.nl/robocupathome/documents/rulebook2009_FINAL.pdf, 2009.



The Future of Humanoid Robots - Research and Applications

Edited by Dr. Riadh Zaier

ISBN 978-953-307-951-6

Hard cover, 300 pages

Publisher InTech

Published online 20, January, 2012

Published in print edition January, 2012

This book provides state of the art scientific and engineering research findings and developments in the field of humanoid robotics and its applications. It is expected that humanoids will change the way we interact with machines, and will have the ability to blend perfectly into an environment already designed for humans. The book contains chapters that aim to discover the future abilities of humanoid robots by presenting a variety of integrated research in various scientific and engineering fields, such as locomotion, perception, adaptive behavior, human-robot interaction, neuroscience and machine learning. The book is designed to be accessible and practical, with an emphasis on useful information to those working in the fields of robotics, cognitive science, artificial intelligence, computational methods and other fields of science directly or indirectly related to the development and usage of future humanoid robots. The editor of the book has extensive R&D experience, patents, and publications in the area of humanoid robotics, and his experience is reflected in editing the content of the book.

How to reference

In order to correctly reference this scholarly work, feel free to copy and paste the following:

Muhammad Attamimi, Tomoaki Nakamura, Takayuki Nagai, Komei Sugiura and Naoto Iwahashi (2012). Learning Novel Objects for Domestic Service Robots, *The Future of Humanoid Robots - Research and Applications*, Dr. Riadh Zaier (Ed.), ISBN: 978-953-307-951-6, InTech, Available from: <http://www.intechopen.com/books/the-future-of-humanoid-robots-research-and-applications/learning-novel-objects-for-domestic-service-robots>

INTECH
open science | open minds

InTech Europe

University Campus STeP Ri
Slavka Krautzeka 83/A
51000 Rijeka, Croatia
Phone: +385 (51) 770 447
Fax: +385 (51) 686 166
www.intechopen.com

InTech China

Unit 405, Office Block, Hotel Equatorial Shanghai
No.65, Yan An Road (West), Shanghai, 200040, China
中国上海市延安西路65号上海国际贵都大饭店办公楼405单元
Phone: +86-21-62489820
Fax: +86-21-62489821

© 2012 The Author(s). Licensee IntechOpen. This is an open access article distributed under the terms of the [Creative Commons Attribution 3.0 License](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.