

# Mutual Development of Behavior Acquisition and Recognition based on Value System

Yasutake Takahashi, Yoshihiro Tamura, and Minoru Asada\*

Dept. of Adaptive Machine Systems, Graduate School of Engineering, Osaka University,

\*JST ERATO Asada Synergistic Intelligence Project

Yamadaoka 2-1, Suita, Osaka, 565-0871, Japan

Email: {yasutake,yoshihiro.tamura,asada}@ams.eng.osaka-u.ac.jp

**Abstract**—Both self-learning architecture (embedded structure) and explicit/implicit teaching from other agents (environmental design issue) are necessary not only for one behavior learning but more seriously for life-time behavior learning. This paper presents a method for a robot to understand unfamiliar behavior shown by others through the collaboration between behavior acquisition and recognition of observed behavior, where the state value has an important role not simply for behavior acquisition (reinforcement learning) but also for behavior recognition (observation). That is, the state value updates can be accelerated by observation without real trial and error while the learned values enrich the recognition system since it is based on estimation of the state value of the observed behavior. The validity of the proposed method is shown by applying it to a dynamic environment where two robots play soccer.

## I. INTRODUCTION

Reinforcement learning has been studied well for motor skill learning and robot behavior acquisition in both single and multi-agent environments. Especially, in the multi-agent environment, observation of others make the behavior learning rapid and therefore much more efficient [1], [2], [3]. Actually, it is desirable to acquire various unfamiliar behavior with some instructions from others, for example, surrounding robots and/or humans in real environment because of huge exploration space and enormous learning time to learn. Therefore, behavior learning through observation has been more important. Understanding observed behavior does not mean simply following the trajectory of an end-effector or joints of demonstrator. It means reading its/his/her intention, that is, the goal of the observed behavior and finding a way how to achieve the goal by oneself regardless of the difference of the trajectory. From a viewpoint of the reinforcement learning framework, this means reading rewards of the observed behavior and estimating sequence of the value through the observation.

Takahashi et al.[4] proposed a method of not only to learn and execute a variety of behavior but also to recognize behavior of others supposing that the observer has already acquired the values of all kinds of behavior the observed agent can do. The recognition means, in this paper, that the robot categorizes the observed behavior to a set of its own behavior acquired beforehand. The method seamlessly combines behavior acquisition and recognition based on “state value” in reinforcement learning scheme. Reinforcement learning generates not only an appropriate behavior (a

map from states to actions) to accomplish a given task but also a utility of the behavior, an estimated discounted sum of rewards that will be received in future while the robot is taking an appropriate policy. This estimated discounted sum of reward is called “state value.” This value roughly indicates closeness to the goal state of the given task if the robot receives a positive reward when it reaches the goal and zero else, that is, if the agent is getting closer to the goal, the value becomes higher. This suggests that the observer may recognize which goal the observed agent likes to achieve if the value of the corresponding task is going higher.

This paper proposes a novel method that enhances behavior acquisition and recognition based on interaction between learning and observation of behavior. A robot learns its behavior through not only trial and error but also reading rewards of the observed behavior of others (including surrounding robots and humans). Fig.1 shows a rough idea of our proposed method.  $V(s)$  and  $\hat{V}(s)$  are the state value updated by oneself and the state value estimated through observation, respectively. Takahashi et al.[4] showed the

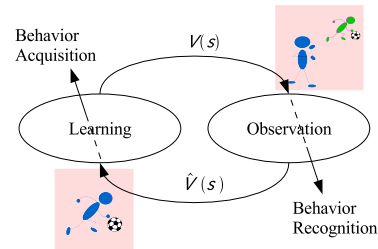


Fig. 1. Interaction between Learning and Observation of Behavior

capability of the proposed method mainly in case that the observer has already acquired a number of behavior to be recognized beforehand. Their case study showed how this system recognizes observed behavior based on the state value functions of self-behavior. This paper shows how the estimated state value of observed behavior,  $\hat{V}(s)$ , gives feedback to learning and understanding unfamiliar observed behavior and this feedback loop enhances the performance of observed behavior recognition. The validity of the proposed method is shown by applying it to a dynamic environment where two robots play soccer.

## II. OUTLINE OF THE MECHANISMS

The reinforcement learning scheme, the state/action value function, the modular learning system for various behavior acquisition/emulation, and behavior recognition based on the state value are explained, here.

### A. Behavior Learning Based on Reinforcement Learning

An agent can discriminate a set  $S$  of distinct world states. The world is modeled as a Markov process, making stochastic transitions based on its current state and the action taken by the agent based on a policy  $\pi$ . The agent receives reward  $r_t$  at each step  $t$ . State value  $V^\pi$ , the discounted sum of the reward received over time under execution of policy  $\pi$ , will be calculated as follows:

$$V^\pi = \sum_{t=0}^{\infty} \gamma^t r_t . \quad (1)$$

In case that the agent receives a positive reward if it stays a specified goal and zero else, then, the state value increases if the agent follows an appropriate policy  $\pi$ . The agent updates its policy through a process of trial and error in order to receive higher positive rewards in future. Analogously, as animals get closer to former action sequences that led to goals, they are more likely to retry it. For further details, please refer to the textbook of Sutton and Barto[5] or a survey of robot learning[6].

Here we introduce a model-based reinforcement learning method. A learning module has a forward model which represents the state transition model and a behavior learner which estimates the state-action value function based on the forward model in a reinforcement learning manner.

Each learning module has its own state transition model. This model estimates the state transition probability  $\hat{P}_{ss'}^a$  for the triplet of state  $s$ , action  $a$ , and next state  $s'$ :

$$\hat{P}_{ss'}^a = Pr\{s_{t+1} = s' | s_t = s, a_t = a\} \quad (2)$$

Each module has a reward model  $\hat{R}_s$ , too:

$$\hat{R}(s) = E\{r_t | s_t = s\} \quad (3)$$

All experiences (sequences of state-action-next state and reward) are simply stored to estimate these models. Now we have the estimated state transition probability  $\hat{P}_{ss'}^a$  and the expected reward  $\hat{R}_s$ , then, an approximated state-action value function  $Q(s, a)$  for a state action pair  $s$  and  $a$  is given by

$$Q(s, a) = \sum_{s'} \hat{P}_{ss'}^a [\hat{R}(s') + \gamma V(s')] \quad (4)$$

$$V(s) = \max_a Q(s, a) , \quad (5)$$

where  $\gamma$  is a discount factor.

### B. Modular Learning System

In order to observe/learn/execute a number of behavior, a modular learning system is adopted. Many modular architectures have been proposed so far (for example [7], [8], [9]). Each module is responsible for learning to achieve a single

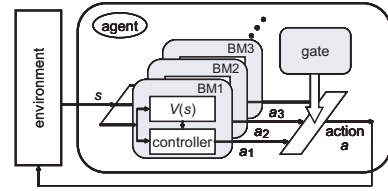


Fig. 2. Modular learning system

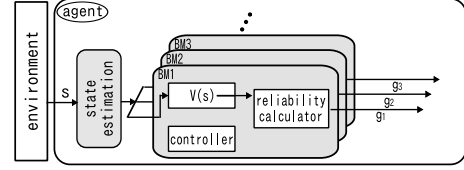


Fig. 3. Behavior inference diagram

goal. One arbiter or a gate module is responsible for merging information from the individual modules in order to derive a single action performed by the robot.

Fig.2 shows a sketch of such a modular learning system. We prepare a number of behavior modules (BMs in the figure) each of which adopts the behavior learning method described in II-A. The module is assigned to one goal-oriented behavior and estimates one action value function  $Q(s, a)$ . A module receives a positive reward when it accomplishes the assigned behavior or zero reward else. The behavior module has a controller that selects the action with the maximum value based on the predictions of next state values. The gating module will then select one output from the inputs of the different behavior modules according to the player's intention.

Fig.3 shows a diagram to recognize an observed behavior. The same behavior modules are used for the behavior recognition. Each behavior module estimates the state value based on the estimated state of the observed demonstrator<sup>1</sup> and calculates reliability of observed behavior, that is, how likely the demonstrator is taking the behavior of the module. The details are described in following sections.

### C. Behavior Recognition based on Estimated Values

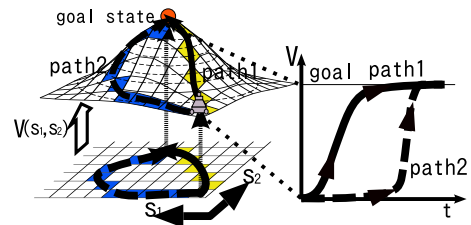


Fig. 4. Behavior recognition based on the change of state value

Fig.4 shows an example task of navigation in a grid world

<sup>1</sup>For reasons of consistency, the term "demonstrator" is used to describe any agent from which an observer can learn, even if the demonstrator does not have an intention to show its behavior to the observer.

and a map of the state value of the task. There is a goal state at the top center of the world. An agent can move one of the neighboring square in the grids every step. It receives a positive reward only when it stays at the goal state while zero else. There are various optimal/suboptimal policies for this task as shown in the figure. If one tries to match the action that the agent took and the one based on a certain policy in order to recognize the observed behavior, it has to maintain various optimal policies and evaluate all of them in the worst case. On the other hand, if the agent follows an appropriate policy, the value is going up even if it is not exactly the optimal one. Likewise, in emulation one is not committed with the optimal policy, as the behavior are the ones available in the portfolio of the agent, which are not necessarily the optimal ones, but the ones that the agent knows to lead to the goal.

While an observer watches a demonstrator's behavior, it uses the same behavior modules for recognition of observed behavior as shown in Fig.3. Each behavior module estimates the state value based on the estimated state of the observed demonstrator and sends it to the selector. The selector watches the sequence of the state values and selects a set of possible behavior modules of which state values are going up as a set of behavior the demonstrator is currently taking. As mentioned before, if the state value goes up during a behavior, it means that the module is valid for explaining the behavior. The observed behavior is recognized by a set of behavior whose modules' values are increasing.

Here we define behavior recognition reliability  $g$  that indicates how much the observed behavior would be reasonable to be recognized as a behavior

$$g = \begin{cases} g + \beta & \text{if } V_t - V_{t-1} > 0 \text{ and } g < 1 \\ g & \text{if } V_t - V_{t-1} = 0 \\ g - \beta & \text{if } V_t - V_{t-1} < 0 \text{ and } g > 0 \end{cases},$$

where  $\beta$  is an update parameter, and 0.1 in this paper. This equation indicates that the reliability  $g$  will become large if the estimated utility rises up and it will become low when the estimated utility goes down. Another condition is to keep  $g$  value from 0 to 1.

#### D. Learning by Observation

In the previous section, behavior recognition system based on state value of its own behavior is described. This system shows robust recognition of observed behavior[10] only when the behavior to be recognized has been well-learned beforehand. If the behavior is under learning, then, the recognition system is not able to show good recognition performance at beginning. The trajectory of the observed behavior can be a bias for learning behavior and might enhance the behavior learning based on the trajectory. The observer cannot watch actions of observed behavior directly and can only estimate the sequence of the state of the observed robot. Let  $s_t^o$  be the estimated state of the observed robot at time  $t$ . Then, the estimated state value  $\hat{V}^o$  of the

observed behavior can be calculated as below:

$$\hat{V}^o(s) = \sum_{s'} \hat{P}_{ss'}^o [\hat{R}(s') + \gamma V^o(s')] \quad (6)$$

where  $\hat{P}_{ss'}^o$  is state transition probability estimated from the behavior observation. This state value function  $\hat{V}^o$  can be used a bias of the state value function of the learner  $V$ . The learner updates its state-action value function  $Q(s, a)$  during the process of trial and error based on the estimated state value of observed behavior  $\hat{V}^o$  as below:

$$Q(s, a) = \sum_{s'} \hat{P}_{ss'}^a [\hat{R}(s') + \gamma V'(s')] \quad (7)$$

while

$$V'(s) = \begin{cases} V(s) & \text{if } V(s) > \eta^n \hat{V}^o(s) \\ \eta^n \hat{V}^o(s) & \text{else} \end{cases}$$

This is a normal update equation as shown in (4) except using  $V'(s)$ . The update system switches the state value of the next state  $s'$  between the state value of own learning behavior  $V(s')$  and the one of the observed behavior  $\hat{V}^o(s')$  discounted by  $\eta^n$ .  $\eta^n$  is the discount factor ( $0 < \eta < 1$ ) based on the time of experiences  $n$ . This means that it receives bigger feedback with  $\hat{V}^o(s')$  if the state-action transition “ $(s, a) \rightarrow s'$ ” is not experienced so far while smaller else. This means the state value update system takes  $\hat{V}^o(s')$  if the learner does not estimate the state value  $V(s')$  because of lack of experience at the state  $s'$  from which it reaches to the goal of the behavior.  $\hat{V}^o(s')$  becomes a bias for reinforcing the action  $a$  from the state  $s$  even though the state value of its own behavior  $V(s')$  is small so that it leads the learner to explore the space near to the goal state of the behavior effectively.

A demonstrator is supposed to show a number of behavior which are not informed directly to the observer. The observer, who learns/recognizes the behavior, assumes that all agents share reward models of the behavior, that is, all agents, including human players, receive a positive reward when the goal of the behavior is achieved. This assumption is very natural as we assume that we share “value” with colleagues, friends, or our family in our daily life. In order to update the estimate values of the behavior the demonstrator is taking, the observer has to estimate which behavior the demonstrator is taking correctly. If the observer waits to learn some specific behavior by observation until it becomes able to recognize the observed behavior well, bootstrap of leaning unfamiliar behavior by observation cannot be expected. Therefore, the observer(learner) maintains a history of the observed trajectories and updates value function of the observed behavior with high reliability or high received reward. The observer estimates the state of the demonstrator every step and the reward received by the demonstrator is estimated as well. If it is estimated that the demonstrator receives a positive reward by reaching to the goal state of the behavior, then, the observer updates the state value of the corresponding behavior even if it has low reliability for the observed behavior. The update strategy enhances to estimate appropriate values of the observed behavior.

### III. BEHAVIOR LEARNING BY OBSERVATION

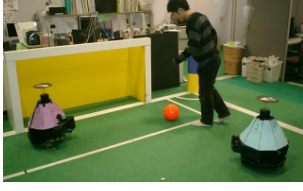


Fig. 5. Robots with a human player in a Soccer Field

Fig.5 shows two robots, a human player and color-coded objects, e.g., an orange ball, and a goal. The robot has an omni-directional camera on top. A simple color image processing is applied in order to detect the color-coded objects and players in real-time. The mobile platform is based on an omni-directional vehicle. These two robots and the human play soccer such as dribbling a ball, kicking it to a goal, passing a ball to the other, and so on. While playing with objects, they watch each other, try to understand observed behavior of the other, and emulate them.

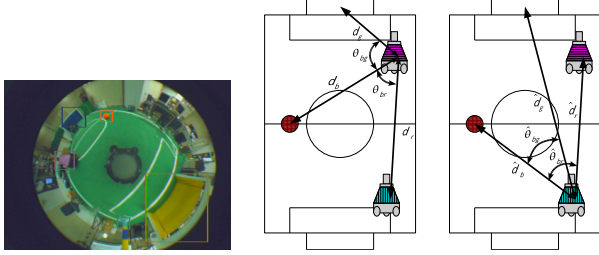


Fig. 6. Estimation of view of the demonstrator. Left : a captured image the of observer, Center : object detection and state variables for self, Right : estimation of view of the demonstrator

Each behavior module estimates a state value of observed behavior at an arbitrary time  $t$  to accomplish the specified task. An observer watches a demonstrator's behavior and maps the sensory information from an observer viewpoint to a demonstrator's one with a simple mapping of state variables. Fig.6 shows a simple example of this transformation. It detects color-coded objects on the omni-directional image, calculates distances and directions of the objects in the world coordinate of the observer, and shifts the axes so that the position of the demonstrator comes to center of the demonstrator's coordinate. Then it roughly estimates the state information in the egocentric coordinate and the state of the demonstrator. Every behavior module estimates a sequence of its state value from the estimated state of the observed demonstrator and the system selects modules which values are increasing. The learner tries to acquire a number of behavior shown in Table II. The table also describes necessary state variables for each behavior. Each state variable is divided into 11 in order to construct a quantized state space. 4 actions are prepared to be selected by the learning modules: Approaching the goal, approaching the teammate, going in front of the ball while watching the goal, and going in front of the ball while watching the teammate.

TABLE I

LIST OF BEHAVIOR LEARNED BY SELF AND STATE VARIABLES FOR EACH

BEHAVIOR	
Behavior	State variables
Approaching a ball	$d_b$
Approaching a goal	$d_g$
Approaching the teammate	$d_r$
Shooting a ball	$d_b, d_g, \theta_{bg}$
Passing a ball	$d_b, d_r, \theta_{br}$

#### A. Comparison Experiment Setup

In order to validate the effect of interaction between acquisition and recognition of behavior through observation, two experiments are set up. One is that the learner does not observe the behavior of other but tries to acquire shooting/passing behavior by itself. The other is that the learner observes the behavior of other and enhances the learning of the behavior based on the estimated state value of the observed behavior. In former experiment, the learner follows the learning procedure:

- 1) 15 episodes for behavior learning by itself
- 2) evaluation of self-behavior performance
- 3) evaluation of behavior recognition performance
- 4) goto 1.

On the other hand, the later experiment, it follows :

- 1) 5 episodes for observation of the behavior of the other
- 2) 10 episodes for behavior learning by self-trials with observed experience
- 3) evaluation of self-behavior performance
- 4) evaluation of behavior recognition performance
- 5) goto 1.

Both learners attempt to acquire behavior listed in Table II. The demonstrator shows one of the behavior one by one and the observer does not know which behavior the demonstrator is taking. In both experiments, the learner follows  $\epsilon$ -greedy method; it follows the greedy policy with 80% probability and takes a random action else. Performance of the behavior execution and recognition of observed behavior during the learning time is evaluated every 15 learning episodes. The performance of the behavior execution is the success rate of the behavior while the learner, the ball, and the teammate are placed at a set of pre-defined positions. The one of the behavior recognition is the average length of period in which the recognition reliability of the right behavior is larger than 70% during the observation. The soccer field area is divided 3 by 3 and the center of the each area is a candidate of the position of the ball, the learner, or the teammate. The performances are evaluated in all possible combinations of the positions.

#### B. Recognition of Observed Behavior

Before evaluating the performance of the behavior execution and behavior recognition of other during learning the behavior, we briefly review how this system estimates the values of behavior and recognizes the observed behavior

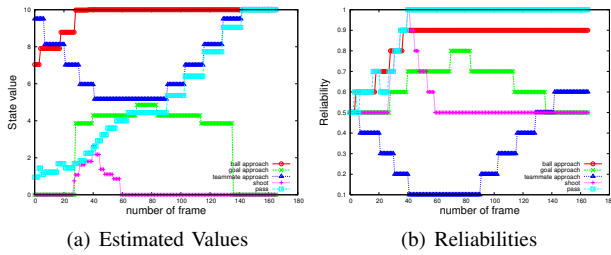


Fig. 7. Sequence of estimated values and reliabilities during a behavior of pushing a ball to the magenta player, red line : approaching a ball, green line : approaching the goal, light blue line : passing, blue line : approaching the other, magenta line : shooting

after the observer has learned behavior. When the observer watches a behavior of the other, it recognizes the observed behavior based on repertoire of its own behavior. Figs.7 (a) and (b) show sequences of estimated values and reliabilities of the behavior, respectively. The line that indicates the passing behavior keeps tendency of increasing value during the behavior in this figures. This behavior is composed of behavior of approaching a ball and approaching the teammate again, then, the line of approaching a ball goes up at the earlier stage and the line of approaching the teammate goes up at the later stage in Fig.7(a). All reliabilities start from 0.5 and increase if the value goes up and decrease else. Even when the value stays low, if it is increasing with small value, the reliability of the behavior increases rapidly. The reliability of the behavior of pushing a ball into the teammate reaches 1.0 at middle stage of the observed behavior. “Recognition period rate” of observed behavior is introduced here to evaluate how long the observer can recognize the observed behavior as a correct one. The recognition period rate is 85% here, that means, the period in which the reliability of passing behavior is over 70% is 85% during the observation.

### C. Performance of Behavior Learning and Recognition

In this section, performances of the behavior execution and behavior recognition during learning the behavior are shown.

Fig.8 shows the success rates of the behavior and their variances during learning in cases of learning with/without value update through observation. The success rates with value update of all kinds of behavior grows more rapidly than the one without observation feedback. Rapid learning is one of the most important aspect for a real robot application. The success rate without value update through observation sometimes could not reach the goal of the behavior at the beginning of the learning because there is no bias to lead the robot to learn appropriate actions. This is the reason why the variances of the rate is big. On the other hand, the system with value update through observation utilizes the observation to bootstrap the learning even though it cannot read exact actions of observed behavior.

Recognition performance and recognition period rate of observed behavior and their variances are shown in Figs.9 and 10, respectively. They indicate a similar aspect with the ones of the success rates. The performance of the behavior

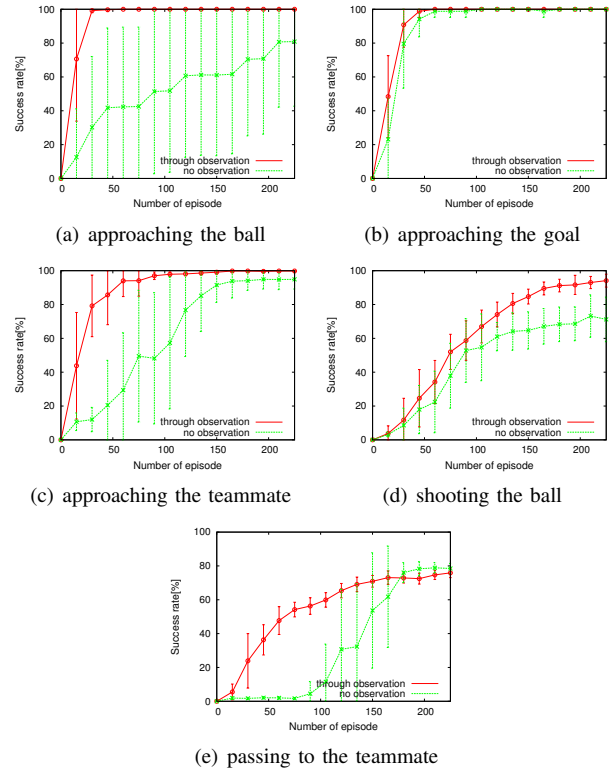


Fig. 8. Success rate of the behavior during learning with/without observation of demonstrator's behavior

recognition depends on the learning performance. If the learning system has not acquired data enough to estimate state value of the behavior, it cannot perform well. The learning system with value update with observed behavior rapidly enables to recognize the behavior while the system without value update based on the observation has to wait to realize a good recognition performance until it estimates good state value of the behavior by its own trial and error. Those figures show the importance of learning through interaction between behavior acquisition and recognition of observed behavior.

### D. Experiments with a real robot and a human demonstrator

The proposed architecture was applied to the real robot system. The robot observed demonstrations of shooting and passing behavior played by a human player first, then, the robot learned the behavior by itself. The human player showed shooting and passing behavior 5 times for each. The real learning time was half an hour for each behavior. Figs. 11 and 12 show one scene of human demonstration of the passing behavior and acquired passing behavior, respectively.

### E. Experiment with Two Humanoid Robots

In order to show wide applicability of the proposed method, another experiment has been done with two simulated humanoid robots. Fig. ?? shows the scenario of our experiment. There are two players in front of a table. Two objects (red and blue boxes) are on the table. A player is independently acquiring behavior of reaching for one of the

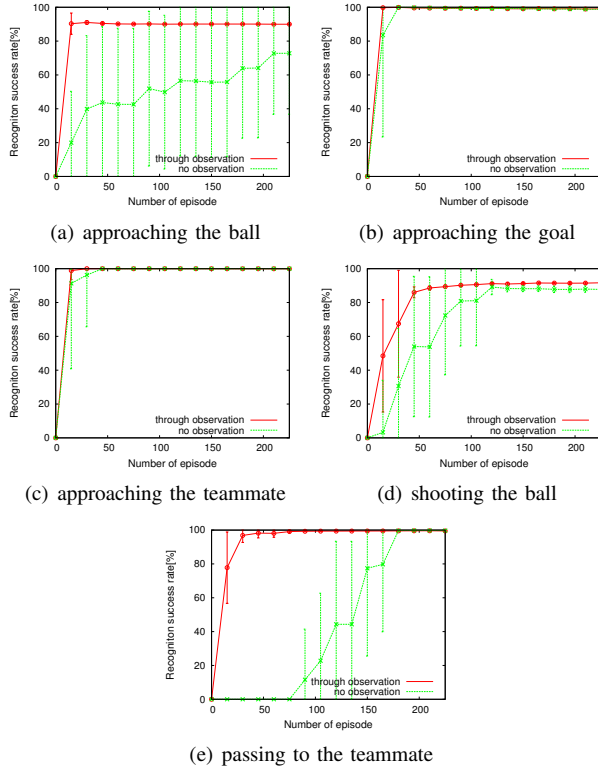


Fig. 9. Recognition performance of the behavior during learning with/without observation of demonstrator's behavior

objects and pulling it to one of the color-coded (magenta and cyan) areas on the table, and update a value system based on reinforcement learning. While playing with objects, they watch each other, try to understand observed behavior of the other, and emulate them.

TABLE II  
LIST OF BEHAVIOR LEARNED BY SELF AND STATE VARIABLES FOR EACH BEHAVIOR

Behavior	State variables
Putting red box on magenta area	$d_r$ and $d_m$
Putting red box on cyan area	$d_r$ and $d_c$
Putting blue box on magenta area	$d_b$ and $d_m$
Putting blue box on cyan area	$d_b$ and $d_c$

The player tries to acquire a number of behavior shown in Table II. The table also describes necessary state variables for each behavior.  $d_r$ ,  $d_b$ ,  $d_m$ , and  $d_c$  indicate distances between the palm of the player and the red box, the blue box, magenta area, and cyan area, respectively. Each state variable is divided into 30 in order to construct a quantized state space. 4 actions are prepared to be selected by the learning modules: reaching the red box, reaching the blue box, pulling the box to the magenta area, and pulling the box to the cyan area.

The success rates of the behavior, the recognition performance and the recognition period rate of observed behavior and their variances are shown in Figs. 14, 15 and 16, respectively. Behavior “Putting red box on magenta

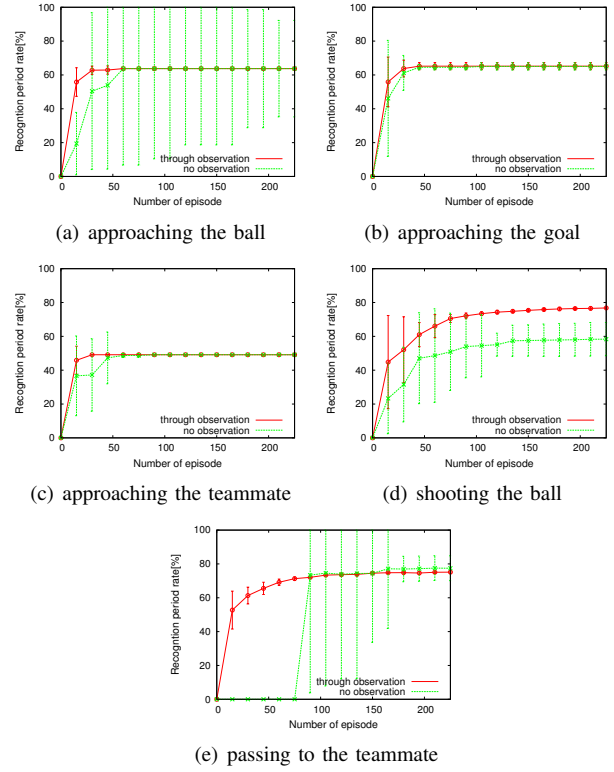


Fig. 10. Recognition period rate of the behavior during learning with/without observation of demonstrator's behavior



Fig. 11. Demonstration of a passing behavior by a human

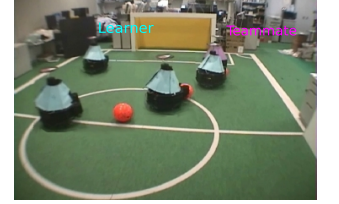


Fig. 12. Acquired passing behavior after the observation and learning

area” and “Putting red box on cyan area” show that the observation feedback makes better performances on learning and recognition of the observed behavior while the case of behavior “Putting blue box on magenta area” and “Putting blue box on cyan area” does not show a big advantage. The reason is that the cyan area is too close to the red and blue boxes then the learner fails to have feedback of the observation of behavior “Putting blue box on areas.” It is no problem because those behavior is easy to learn. Behavior “Putting red box on areas” needs a long period to be achieved and learned and the observation feedback works effectively.

#### IV. CONCLUSION

Above, values are defined as behavior, which are defined by the achieved goals. The observer uses its own value functions to recognize what the demonstrator will do. Preliminary investigations in a similar context have been done by Takahashi et al. [10] and they showed much better

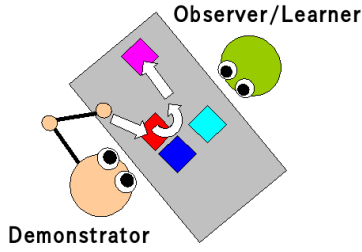


Fig. 13. Scenario of the experiment

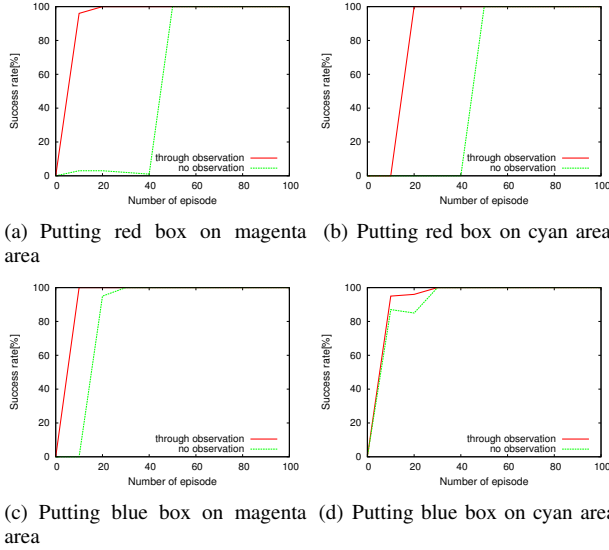


Fig. 14. Success rate of the behavior during learning with/without observation of demonstrator's behavior

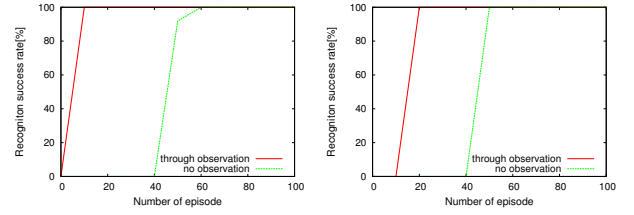
robustness of behavior recognition than a typical method. In this paper, unknown behavior are also understood in term of one's own value function through learning based on the estimated values derived from the observed behavior. Furthermore, value update through the observation enhances not only the performance of behavior learning but also the one of recognition of the observed behavior effectively.

## V. ACKNOWLEDGMENTS

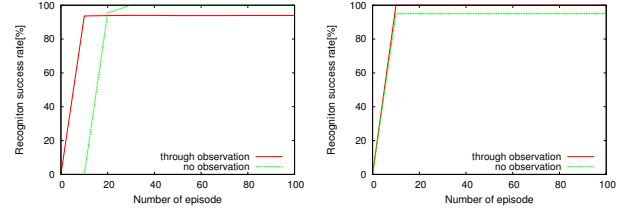
This work is partially supported by Kayamori Foundation of Informational Science Advancement.

## REFERENCES

- [1] S. D. Whitehead, "Complexity and cooperation in q-learning," in *Proceedings Eighth International Workshop on Machine Learning (ML91)*, 1991, pp. 363–367.
- [2] B. Price and C. Boutilier, "Accelerating reinforcement learning through implicit imitation," *Journal of Artificial Intelligence Research*, 2003.
- [3] D. C. Bentivegna, C. G. Atkeson, and G. Chenga, "Learning tasks from observation and practice," *Robotics and Autonomous Systems*, vol. 47, pp. 163–169, 2004.
- [4] Y. Takahashi, T. Kawamata, M. Asada, and M. Negrello, "Emulation and behavior understanding through shared values," in *Proceedings of the 2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, Oct 2007, pp. 3950–3955.

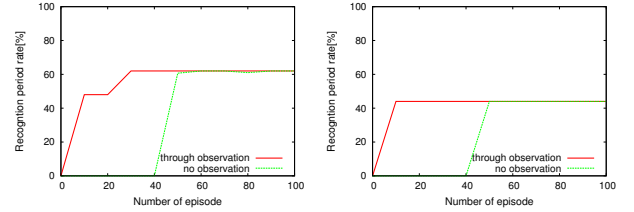


(a) Putting red box on magenta (b) Putting red box on cyan area

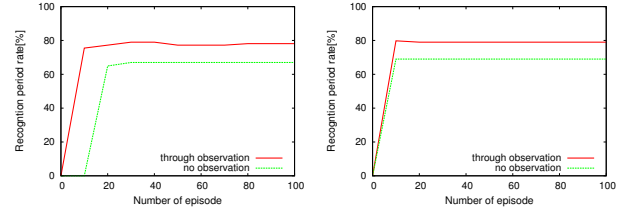


(c) Putting blue box on magenta (d) Putting blue box on cyan area

Fig. 15. Recognition performance of the behavior during learning with/without observation of demonstrator's behavior



(a) Putting red box on magenta (b) Putting red box on cyan area



(c) Putting blue box on magenta (d) Putting blue box on cyan area

Fig. 16. Recognition period rate of the behavior during learning with/without observation of demonstrator's behavior

- [5] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. Cambridge, MA: MIT Press, 1998. [Online]. Available: [citeseer.ist.psu.edu/sutton98reinforcement.html](http://citeseer.ist.psu.edu/sutton98reinforcement.html)
- [6] J. H. Connell and S. Mahadevan, *ROBOT LEARNING*. Kluwer Academic Publishers, 1993.
- [7] R. Jacobs, M. Jordan, N. S., and G. Hinton, "Adaptive mixture of local experts," *Neural Computation*, vol. 3, pp. 79–87, 1991.
- [8] S. P. Singh, "Transfer of learning by composing solutions of elemental sequential tasks," *Machine Learning*, vol. 8, pp. 323–339, 1992. [Online]. Available: [citeseer.nj.nec.com/singh92transfer.html](http://citeseer.nj.nec.com/singh92transfer.html)
- [9] S. Whitehead, J. Karlsson, and J. Tenenbarg, "Learning multiple goal behavior via task decomposition and dynamic policy merging," in *ROBOT LEARNING*, J. H. Connell and S. Mahadevan, Eds. Kluwer Academic Publishers, 1993, ch. 3, pp. 45–78.
- [10] Y. Takahashi, T. Kawamata, and M. Asada, "Learning utility for behavior acquisition and intention inference of other agent," in *Proceedings of the 2006 IEEE/RSJ IROS 2006 Workshop on Multi-objective Robotics*, Oct 2006, pp. 25–31.