

# Rapid Behavior Learning in Multi-Agent Environment based on State Value Estimation of Others

Yasutake Takahashi, Kentaro Noma, and Minoru Asada  
Dept. of Adaptive Machine Systems, Graduate School of Engineering,  
Osaka University, Yamadagaoka 2-1, Suita, Osaka, 565-0871, Japan  
Email: {yasutake,kentaro.noma,asada}@ams.eng.osaka-u.ac.jp

**Abstract**—The existing reinforcement learning approaches have been suffering from the curse of dimension problem when they are applied to multiagent dynamic environments. One of the typical examples is a case of RoboCup competitions since other agents and their behaviors easily cause state and action space explosion. This paper presents a method of modular learning in a multiagent environment by which the learning agent can acquire cooperative behaviors with its team mates and competitive ones against its opponents. The key ideas to resolve the issue are as follows. First, a two-layer hierarchical system with multi learning modules is adopted to reduce the size of the sensor and action spaces. The state space of the top layer consists of the state values from the lower level, and the macro actions are used to reduce the size of the physical action space. Second, the state of the other to what extent it is close to its own goal is estimated by observation and used as a state value in the top layer state space to realize the cooperative/competitive behaviors. The method is applied to 4 (defense team) on 5 (offense team) game task, and the learning agent successfully acquired the teamwork plays (pass and shoot) within much shorter learning time (30 times quicker than the earlier work).

## I. INTRODUCTION

Recently, there have been increasing number of studies on cooperative/competitive behavior acquisition in a multiagent environment by using reinforcement learning methods [1], [2], [3], [4], [5]. In such an environment, the state and action spaces for the learning can be easily exploded since not only objects but also other agents should be involved in the state and action spaces, and therefore the sensor and actuator level descriptions may cause information explosion that disables the learning methods to be applied within practical learning time. Shivaram et al. [3] showed that the learning rate can be accelerated by sharing the learned information in the 4 on 5 game task. However, they need still long learning time since they directly use the sensory information as state value to decide the situation. Stefan et al. [1] achieved the cooperative behavior learning task between two robots in real time by introducing the macro action that is abstracted action code predefined by the designer. However, only the macro actions do not seem sufficient to accelerate the learning time in a case that more agents are included in the environment. Therefore, the sensory information should be also abstracted to reduce the size of the state space.

M. Asada is with JST ERATO Asada Synergistic Intelligence Project, Yamadagaoka 2-1, Suita, Osaka, 565-0871, Japan

Jacobs and Jordan [6] proposed a learning system called “Mixture of Experts” that outputs the summation of the weighted outputs of the multiple learning modules. The idea that the weight reflects the fitness of each module to the current situation is generally useful to build an efficient system. Doya et al. [7] proposed a similar system called “MOSAIC” that spatio-temporally divides the nonlinear and unsteady environment into linear and steady segments, and selects the module corresponding to one of these segments that predicts the state transition most correctly. However, they showed a simple example and its scalability to more complicated tasks such as multiagent interactions in a dynamic environment has not been made clear.

The prediction of other’s behavior is important to realize the cooperative (competitive) behaviors with (against) others in general. Takahashi et al. [8] proposed a method to infer the other’s intention by observation based on the idea that the increase of the state value (the larger the state value, the closer to the goal) means the other intends to achieve the corresponding goal regardless of the differences of viewpoint and/or action to achieve the goal. If this prediction capability is incorporated into the learning system, the learner can efficiently acquire the desired behaviors.

This paper presents a method of modular learning in a multiagent environment by which the learning agent can acquire cooperative behaviors with its team mates and competitive ones against its opponents. The key ideas to resolve the issue are as follows. First, a two-layer hierarchical system with multi learning modules is adopted to reduce the size of the sensor and action spaces. The state space of the top layer consists of the state values from the lower level, and the macro actions are used to reduce the size of the physical action space. Second, the state of the other to what extent it is close to its own goal is estimated by observation and used as a state value in the top layer state space to realize the cooperative/competitive behaviors. The method is applied to 4 (defense team) on 5 (offense team) game task, and the learning agent successfully acquired the teamwork plays (pass and shoot) within much shorter learning time (30 times quicker than the earlier work).

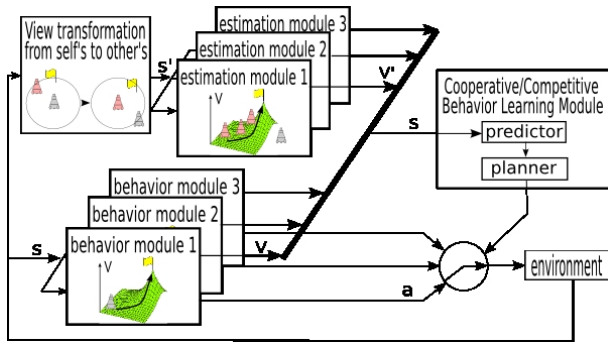


Fig. 1. A multi-module learning system

## II. MULTI MODULE LEARNING SYSTEM WITH OTHER'S STATE INFERENCE MODULES

### A. Architecture

Fig. 1 shows a basic architecture of the proposed system, i.e., a two-layered multi-module reinforcement learning system. The bottom layer consists of two kinds of modules: action modules and other's state inference ones. The top layer consists of a single gate module that learns which action module should be selected according to the current state that consists of state values sent from the modules at the bottom layer. More correctly, it selects one of the action modules which has the best estimation of a state transition sequence by activating a gate signal corresponding to a module while deactivating the gate signals of other modules; the selected module then sends action commands based on its policy.

### B. Action module

An action module has a forward model (predictor) which represents the state transition model and a behavior learner (action planner) which estimates the state-action value function based on the forward model in a reinforcement learning manner.

- (a) **Predictor:** 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  $s'$ :

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

Each module has a reward model  $\hat{R}_{ss'}^a$ , too:

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

We simply store all experiences (sequences of state-action-next state and reward) to estimate these models.

- (b) **Planner:** Now we have the estimated state transition probabilities  $\hat{P}_{ss'}^a$  and the expected rewards  $\hat{R}_{ss'}^a$ , then, an approximated state-action value function  $Q(s, a)$  for a state action pair  $s$  and  $a$  is given by

$$Q(s, a) = \sum_{s'} \hat{P}_{ss'}^a \left[ \hat{R}_{ss'}^a + \gamma \max_{a'} Q(s', a') \right], \quad (3)$$

where  $\gamma$  is a discount rate.

The predictor and planner in each action module have been already acquired before the learning of the gate module here.

### C. Other's state inference module

The role of the other's state inference module is to estimate the state value that indicate the degree of achievement of the other's task by observation, and to send this value to the state space of the gate module at the top layer. In order to estimate the degree of achievement, the following procedure is taken.

- 1) The learner acquires the various kinds of behaviors that the other agent may take, and each behavior corresponds to each action module, and each other's state inference one.
- 2) The learner estimates the sensory information observed by the other through the 3-D reconstruction of its own sensory information.
- 3) Based on the estimated sensory information of the other, each other's state inference module estimates the other's state value by assigning the state value of the corresponding action module of its own.

## III. TASK AND ASSUMPTIONS

The game consists of the offense team (five players and one of them can be the passer) and the defense team (four players attempt to intercept the ball) The offense player nearest to the ball becomes a passer who passes the ball to one of its teammates (receivers) or shoot the ball to the goal if possible while the opposing team tries to intercept it (see Fig. 2).

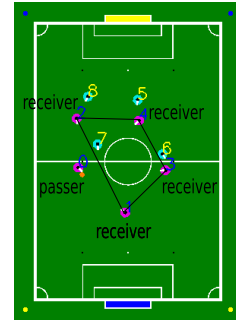


Fig. 2. A passer and the defense formation



Fig. 3. A real robot

Only the passer learns its behavior while the receivers and the defense team members take the fixed control policies. The

receiver becomes the passer after receiving the ball and the passer becomes the receiver after passing the ball. After one episode, the learned information is circulated among team members through communication channel but no communication during one episode. The action and inference modules are given a priori.

The offense (defense) team color is magenta (cyan), and the goal color is blue (yellow) in the following figures.

#### A. Offense team

The passer who is the nearest to the ball passes the ball to one of four receivers or dribble-shoots the ball to the goal. After its passing, the passer shows a pass-and-go behavior that is a motion to the goal during the fixed period of time. The receivers face to the ball and move to the positions so that they can form a rectangle by taking the distance to the nearest teammates (the passer or other receivers) (see Fig. 2). The initial positions of the team members are randomly arranged inside their territory.

#### B. Defense team

The defense team member who is nearest to the passer attempts to intercept the ball, and each of other members attempts to “block” the nearest receiver. “Block” means to move to the position near the offense team member and between the offense and its own goal (see Fig. 2). The offense team member attempts to catch the ball if it is approaching. In order to avoid the disadvantage of the offense team, the defense team members are not allowed inside the penalty area during the fixed period of time. The initial positions of the team members are randomly arranged inside their territory but outside the center circle.

#### C. Robots and the environment

Fig. 3 shows a mobile robot we have designed and built. Fig. 4 shows the viewer of our simulator for our robots and the environment. The robot has an omni-directional camera system. A simple color image processing is applied to detect the ball, the interceptor, and the receivers on the image in real-time (every 33ms.) The left of Fig. 4 shows a situation the agent can encounter while the right images show the simulated ones of the normal and omni vision systems. The mobile platform is an omni-directional vehicle (any translation and rotation on the plane.)

We suppose that the omni directional vision system provides the robot with 3-D construction of the scene. This assumption is needed for the other’s state inference module since it is needed to estimate the sensory information observed by other robots.

### IV. STRUCTURE OF THE STATE AND ACTION SPACES

The passer is only one learner, and the state and action spaces for the lower modules and the gate one are constructed as follows. The action modules are four passing ones for four individual receivers, and one dribble-shoot module. The other’s state inference modules are the ones to estimate the degree of achievement of ball receiving for four individual



Fig. 4. Viewer of simulator

receivers, that is how easily the receiver can receive the ball from the passer. These modules are given in advance before the learning of the gate module.

The action spaces of the lower modules adopt the macro actions that the designer specifies in advance to reduce the size of the exploration space without searching at the physical motor level.

The state space  $S$  for the gate module consists of the following state values from the lower modules:

- four state values of passing action modules corresponding to four receivers,
- one state value of dribble-shoot action module, and
- four state values of receiver’s state inference modules corresponding to four receivers.

In order to reduce the size of the whole state space, these values are binarized, therefore its size is  $2^4 \times 2 \times 2^4 = 512$ .

The rewards are given as follows:

- 10 when the ball is shot into the goal (one episode is over),
- -1 when the ball is intercepted (one episode is over),
- 0.1 when the ball is successfully passed,
- 0.3 when the ball is dribbled.

When the ball is out of the field or the pre-specified time period elapsed, the game is called “draw” and one episode is over.

#### A. State space for the passing module

The state space of the passing module  $S$  is defined on the omni directional camera image as follows (see Fig. 5(a)):

- the smallest angle among angles between the receiver and one the defense players who is nearer to the passer than the receiver ( $\theta_1$ ), and
- the angle between the receiver and one of the defense players who is nearest to the passer ( $\theta_2$ ).

The both angles are quantized into ten levels including an invisible case, therefore the total number of states is 100. An example of the state values of four receivers is shown in Fig. 5(b) where the passer is the robot 3 (hereafter, r3 in short), and the color bars near four robots (r0, r1, r2, and r4) indicate the state values of the pass modules for four receivers, respectively. The higher the bar is, the higher the state value is. Since the pass courses for r1 and r2 are not

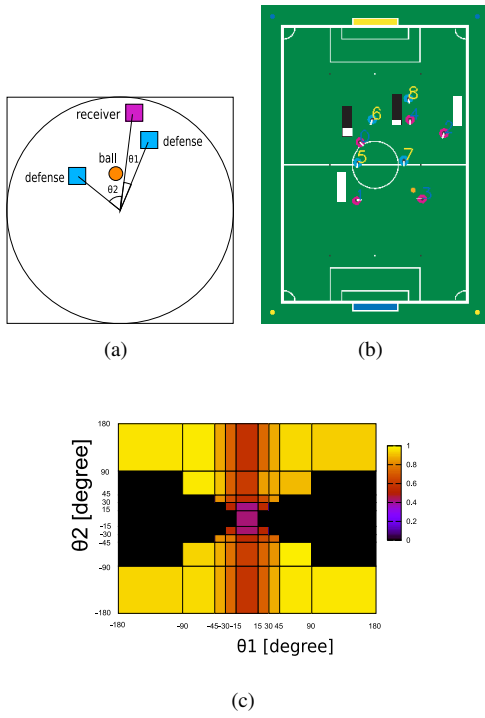


Fig. 5. State variable ((a)), examples of state values ((b)), and state value map of the pass module ((c))

intercepted by the defense players, their state values are high while the state values for  $r_0$  and  $r_4$  are low since their pass courses are intercepted by the defense players.

The state value map is shown in Fig. 5(c) that indicates the smaller the angle between the receiver and the defense player is, the lower the state value is. The black region (one region is separated in the figure) is inexperience area, and the state value sent to the top layer is binarized, that is the value smaller than 1.0 is 0.0.

### B. State space for the dribble-shoot module

The state space of the dribble-shoot module  $S$  is defined on the omni directional camera image as follows (see Fig. 6(a)):

- the angle between the opponent goal and one of the defense players who is nearest to the passer ( $\theta_1$ ),
- the angle between the ball and one of the defense players who is nearest to the passer ( $\theta_2$ ),
- the distance to the nearest defense player ( $r$ ), and
- the angle between the both edges of the opponent goal ( $\theta_3$ ) that represents the distance to the goal.

These state values are quantized into eight, five, eight, and seven, respectively. The total number of states is  $8 \times 8 \times 5 \times 7 = 2240$ .

The state value map of the dribble-shoot module in terms of  $\theta_1$  and  $r$  with fixed values of  $\theta_2$  and  $\theta_3$  is shown in Fig. 6(b) that indicates the nearer the defense player is, the smaller the state value is. The state value sent to the top layer is binarized, that is the value smaller than 1.0 is 0.0.

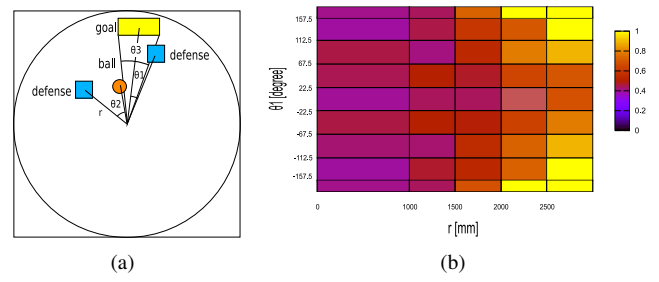


Fig. 6. State variables ((a)) and state values ((b)) for the dribble and shoot module

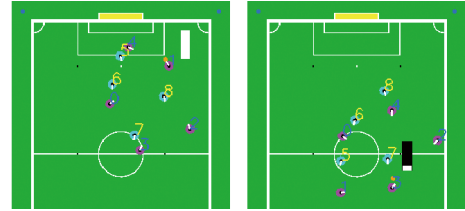


Fig. 7. Two examples of the state values: high (left) and low (right)

Two examples of the state values of the passer expected to take a role of a shooter is shown in Fig. IV-B where the color bars near the passer indicate the state values of the dribble-shoot modules. The higher the bar is, the higher the state value is. Since the passer ( $r_1$ ) is near the goal and no defense players around in Fig. IV-B (left), the state value is high while the state value of the passer ( $r_3$ ) in Fig. IV-B (right) is low since it is located far from the goal and the defense players are around it.

### C. State space for the receiver's state inference module

The passer infers each receiver's state that indicates how easily the receiver can shoot the passed ball to the goal by reconstructing its TV camera view of the scene from the passer's omnidirectional view. Since we suppose that the passer has already learned the shooting behavior, the passer can estimate the receiver's state value by assigning its own experienced state of the shooting behavior.

The state space  $S$  for the receiver's state inference module consists of:

- The distance to the nearest defense player ( $r$ )
- The angle between the both side edged of the opponent goal ( $\theta_1$ ) that represents the distance to the goal (see Fig. (8(a)).).

The both are quantized into five and seven levels, therefore the number of states are  $5 \times 7 = 35$ .

An examples of the state values of the receiver's state inference modules is shown in Fig. 8(b) where the color bars near the four receivers indicate their state values. The higher the bar is, the higher the state value is. Since the receiver ( $r_0$ ) is near the goal and no defense players around, the state value is high while the state values of other receivers ( $r_1$ ,  $r_2$ , and  $r_4$ ) are low since it is located far from the goal and/or the defense players are around it.

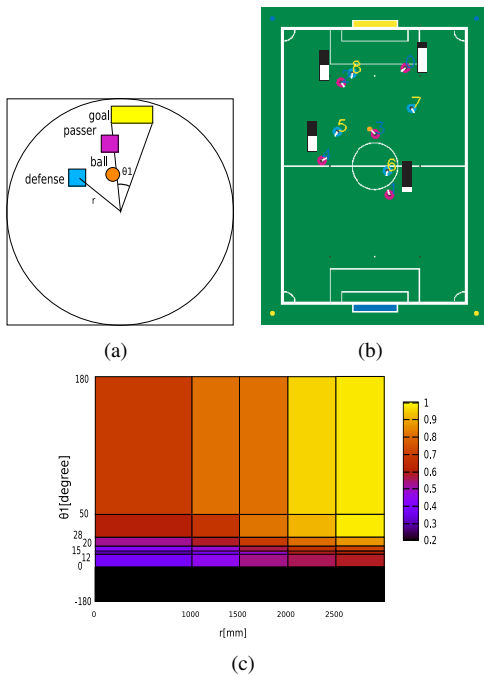


Fig. 8. State variables ((a)), examples of state values ((b)), and state value map ((c)) of the receiver module

The state value map of the receiver's state inference module in terms of  $\theta_1$  and  $r$  is shown in Fig. 8(c) that indicates the nearer (further) the defense player is and the further (nearer) the goal is, the smaller (larger) the state value is. The black region is inexperienced area, and the state value sent to the top layer is binarized, that is the value smaller than 1.0 is 0.0.

## V. EXPERIMENTAL RESULTS

The success rate is shown in Fig. 9(a) where the action selection is 80% greedy and 20% random to cope with new situations. Around the 900th trial, the learning seems to have converged at 30% success, 70% failure, and 10% draw. Compared to the results of [3] that has around 30% success rate with 30,000 trials, the learning time is drastically improved (30 times quicker). Fig. 9(b) indicates the number of passes where it decreases after the 350 trials that means the number of useless passes decreased.

In cases of the success, failure, and draw rates when 100% greedy and 100% random are 55%, 35%, 10%, and 2%, 97%, 1%, respectively. The reason why the success rate in case of 100% greedy is better than in case of 80% greedy seems that the control policies of the receivers and the defense players are fixed, therefore not so new situations happened.

An example of acquired behavior is shown in Fig. 10 where a sequence of twelve top views indicates a successful pass and shoot scene.

## VI. DISCUSSION

We have used the state values instead of the physical sensor values and macro actions instead of the physical motor

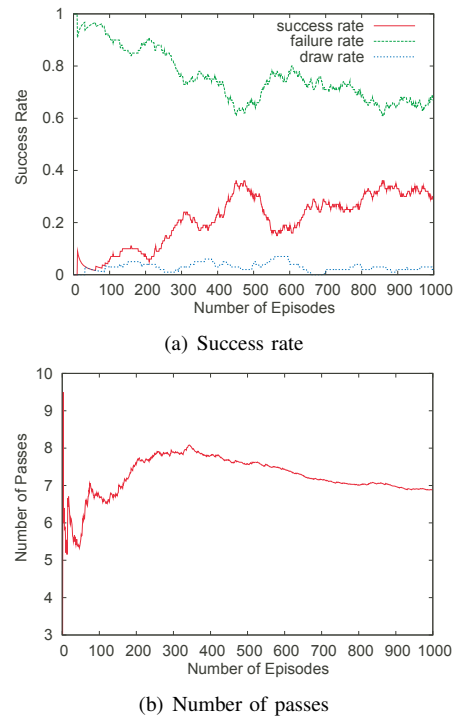


Fig. 9. Success rate and the number of passes

commands, and adopted the receiver's state inference modules that infer how easy for each receiver to receive the ball in order to accelerate the learning. As a result, we have much improved the learning time (30 times quicker!) compared to the result of the existing method [3] that has 32% success with communication and 23% without communication at around the 30,000th trial when the learning seems to have converged.

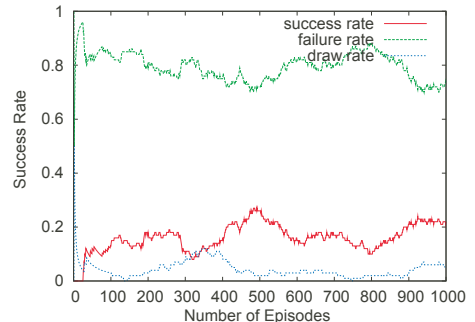


Fig. 11. Success rate without the receiver's state inference modules

Although we have not used the communication between agents during one episode, the receiver's state inference modules seem to take the similar role. Then, we performed the learning without these modules. Fig. 11 shows the success rate, and we can see that the converged success rate is around 21% that is close to 23% of the success rate of the result of the existing method [3]. We may conclude that the state and action space abstraction (the use of state values and macro actions) contributed to the reduction of the

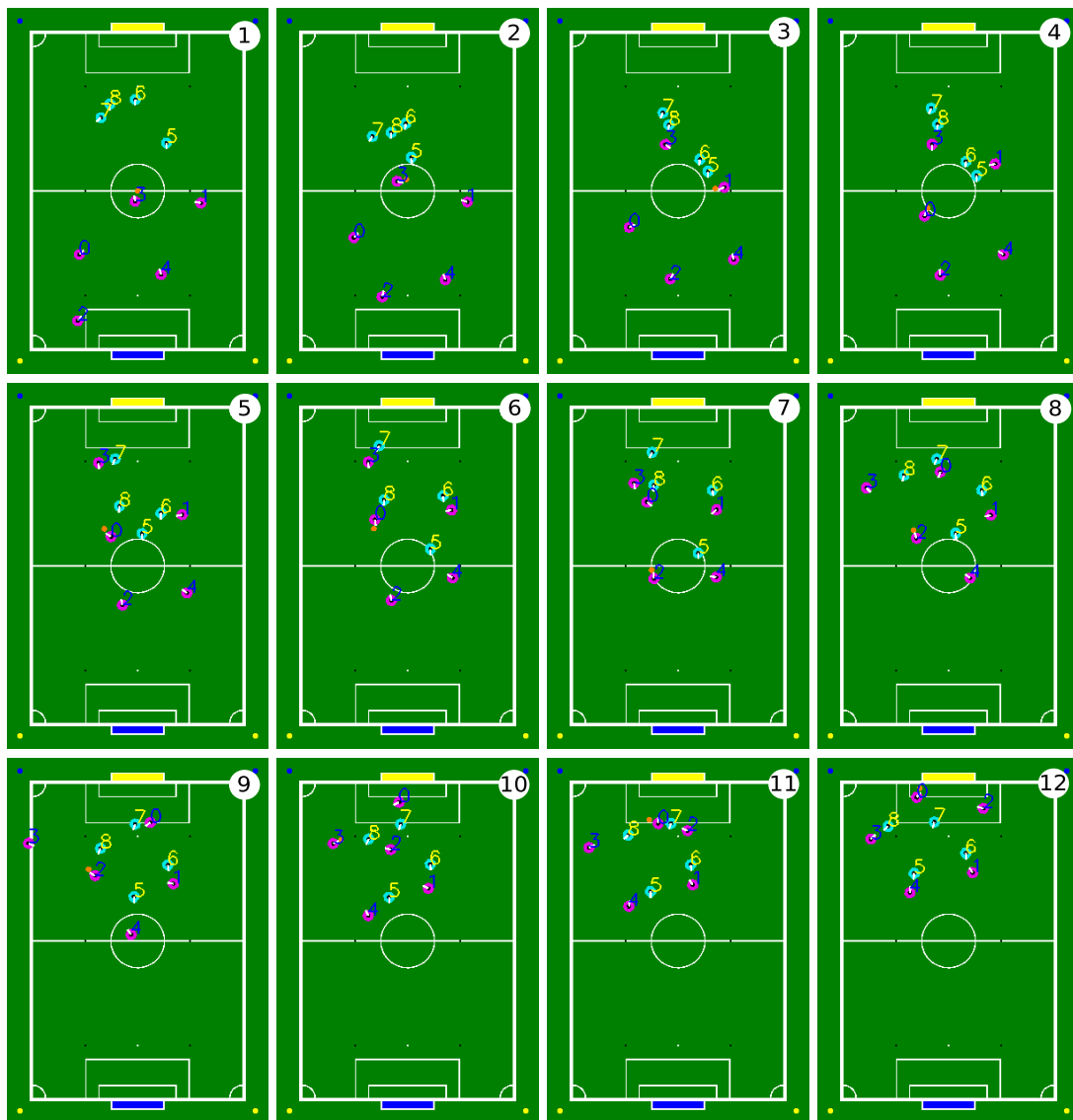


Fig. 10. An example of the acquired behavior

learning time while the use of the receiver's state inference modules contributed to the improvement of the teamwork performance. The real robot implementation is our future work.

#### REFERENCES

- [1] S. Elfving, E. Uchibe, K. Doya, and H. I. Chirstensen, "Multi-agent reinforcement learning: Using macro actions to learn a mating task," *Proceedings of 2004 IEEE/RSJ International Conference on Intelligent Robots and Systems*, vol. 13, pp. 3164–2220, 2004.
- [2] S. Ikenoue, M. Asada, and K. Hosoda, "Cooperative behavior acquisition by asynchronous policy renewal that enables simultaneous learning in multiagent environment," in *Proceedings of the 2002 IEEE/RSJ Intl. Conference on Intelligent Robots and Systems*, 2002, pp. 2728–2734.
- [3] S. Kalyanakrishnan, Y. Liu, and P. Stone, "Half field offense in robocup soccer: A multiagent reinforcement learning case study," in *Proceedings CD RoboCup*, 2006.
- [4] P. Stone, R. S. Sutton, and G. Kuhlmann, "Scaling reinforcement learning toward robocup soccer," *Journal of Machine Learning Research*, vol. 13, pp. 2201–2220, 2003.
- [5] Y. Takahashi, K. Edazawa, and M. Asada, "Multi-module learning system for behavior acquisition in multi-agent environment," in *Proceedings of 2002 IEEE/RSJ International Conference on Intelligent Robots and Systems*, October 2002, pp. CD-ROM 927–931.
- [6] R. Jacobs, M. Jordan, S. Nowlan, and G. Hinton, "Adaptive mixtures of local experts," *Neural Computation*, vol. 3, pp. 79–87, 1991.
- [7] K. Doya, K. Samejima, K. ichi Katagiri, and M. Kawato, "Multiple model-based reinforcement learning," Kawato Dynamic Brain Project Technical Report, KDB-TR-08, Japan Science and Technology Corporation, Tech. Rep., June 2000.
- [8] 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*, 2006, pp. 25–31.