# Air Conditioning Control System Learning Sensory Scale Based on Reinforcement Learning

Yohei Yamaguchi, Noritaka Shigei, Hiromi Miyajima

*Abstract*—This study proposes the air conditioner (AC) control system based on users' sensations. The purpose is to realize the control system that improves low-performance ACs in terms of energy efficiency and comfortableness performance. The system consists of wireless sensor nodes and user nodes such as PCs and smartphones, and it is applicable to already installed ACs. The users enter their sensation such as *cold*, *good*, *a little hot* and *very hot* through the user node, and the system determines the appropriate control according to the users' sensations. The appropriate control policy is determined based on Q-Learning, which is a reinforcement learning method. For multiple users, several methods for integrating users' sensations are presented. These makes our proposed system applicable to a large number of users. Further, in order to reduce the energy consumption and the number of users' inputs, several types of reward functions are presented. In the simulation, four types of methods are evaluated in terms of the time needed for providing a comfortable environment and the energy consumption. We clarify the effective methods among them.

*Index Terms*—air conditioning, reinforcement learning, wireless sensor network, sensory scale, integration of sensations

## I. INTRODUCTION

In our modern life, the air conditioning system is an essential system. The system needs to be energy efficient and to provide a comfortable room temperature environment[1]. Its state-of-the-art systems, which are generally expensive, equip many high-precision sensors and can achieve high energy efficiency and high comfortableness. Due to budgetary constraints, many facilities often have to use old systems or introduce systems with low introducing cost, which are of low-performance in the control ability. However, such systems are of low energy efficiency and cannot always reasonably provide a comfortable environment. Further, for a large room with many workers, the difference between the ideal performance and the actual one would be large. Because the low-performance systems generally do not equip sufficient number of sensors for sensing the whole room.

Thanks to the recent advances in wireless sensor network (WSN) technologies, WSNs have been effectively employed in various fields such as industry[2] and home automation[3]. The sensing capability of WSN can be utilized for improving the old or low-performance systems. Since wireless sensor devices become much cheaper year by year, it is also a realistic approach. Further, recent years, mobile devices such as tablet PCs and smartphones are to be found everywhere. By incorporating them into WSN system, a powerful system can be realized with a reasonable cost.

For realizing intelligent control systems, reinforcement learning (RL) is a promising technique. RL does not need any teacher signal and it can obtain an appropriate behavior pattern by trial and error according to rewards given from the environment[7], [6]. Therefore, RL can be applied to problems with unknown environment. In RL, in order to obtain an appropriate behavior pattern (policy) for the environment defined by the target problem, the agent repeats trial and error and updates its policy according to some reward obtained from the environment (see Fig.4). RLs have been applied to room environment conditioning. In [4], for air conditioning, PI controller has been combined with RL. The purpose of the control is the minimization of the control error. In [5], for lighting system, an intelligent control system has been proposed by actor-critic algorithm, which is one of RLs. The system controls the lighting system according to the presented user's sensation such as brighter and darker. This approach is gentle to humans and attractive in the case where the control has to be determined according to many users.

In this study, we propose the AC control system based on users' sensations. The purpose is to realize the control system that improves the low-performance AC in terms of energy efficiency and comfortableness. The system consists of wireless sensor nodes and user nodes such as PCs and smartphones, and it is applicable to already installed ACs. The users enter their sensation such as *cold*, *good*, *a little hot* and *very hot* through the user node, and the system determines the appropriate control according to the users' sensations. The appropriate control policy is determined based on Q-Learning, which is one of RLs. For multiple users, several methods for integrating users' sensations are presented. Further, in order to reduce the energy consumption and the number of users' inputs, several types of reward functions are presented. In the simulation, four types of RL methods are evaluated in terms of the time needed for providing a comfortable environment and the energy consumption. We clarify the effective methods among them.

## II. MODEL OF ROOM ENVIRONMENT

In this section, a model of the temperature in a room environment with an air conditioner is described. Let $C > 0$ and $W > 0$ be the cooling strength of the air conditioner and the fan speed, respectively. Let $T(t, d)$ be the air temperature at time $t$ and at distance $d$ away from the air conditioner. Then, we define the relation of $T(t, d)$, $C$, $W$ and $d$, based on the RC electric circuit model. The temperature $T(t, d)$ is defined as follows:

$$T(t, d) = T(t - 1, d) + T_{\mathrm{drp}} \cdot \left( 1 - \exp \left( \frac{-1}{\tau(d, C, W)} \right) \right), \tag{1}$$

$$T_{\mathrm{drp}} = -T(t - 1, d) + \alpha \cdot T_{\mathrm{o}} + c(C), \text{ and} \tag{2}$$

TABLE I
THE USED PARAMETERS IN THE ROOM TEMPERATURE MODEL.

| Param. | The used values |
|---|---|
| AC op. mode $(C, W)$ | (0,0), (1,1), (1,1.6), (1,3), (2,1), (2,1.6), (2,3), (3,1), (3,1.6), (3,3) |
| $\alpha$ | 0.1 |
| $\beta$ | 10.0 |
| $c(0)$ | $T_{\mathrm{o}}$ |
| $c(1)$ | 27.0 |
| $c(2)$ | 23.0 |
| $c(3)$ | 18.0 |



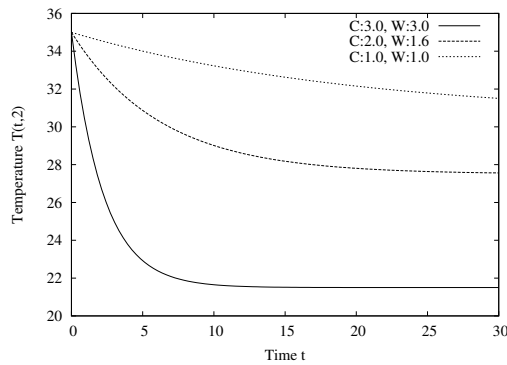Fig. 2.   Changes in the room temperature when changing $d$

.



Fig. 1.   Changes in the room temperature when changing $C$ and $W$.

$$\tau(d, C, W) = \min\left(\frac{\beta \cdot d}{C \cdot W}, \ \beta \cdot d\right) \qquad (3)$$

where $T_{\mathrm{o}}$ is the outer air temperature, $\alpha$ and $\beta$ are constant numbers, and $c(C)$ is a decreasing function with $C$. The factor $T_{\mathrm{drp}}$ determines whether the temperature $T(t, d)$ raises or falls. If $T_{\mathrm{drp}}$ is positive, the temperature increases. Otherwise, it decreases. The coefficient $\alpha$ is the influence rate of $T_{\mathrm{o}}$ and $0 \leq \alpha \leq 1.0$. The function $c(C)$ determines the lower (upper) limit of the temperature in the case of a rise (fall) in temperature. The parameters used in this paper is shown in Table I.

Fig.1 shows the changes of the temperature $T(t, d)$ for $d = 2$ and the different $C$ and $W$. In the figure, it is observed that, as $C$ and $W$ increase, the convergence temperature decreases. Fig.1 shows the changes of the temperature $T(t, d)$ for $C = 3$, $W = 2$ and the different $d$.

In the figure, it is observed that, the convergence temperatures are same but the convergence speed becomes faster with the distance $d$.

### III.   AIR CONDITIONING CONTROL SYSTEM BASED ON SENSUOUS INSTRUCTION

#### A. System Configuration

In this section, we explain our air conditioning system based on sensuous instruction. Our system consists of an air conditioner (AC), two types of sensor devices, user nodes such as personal computer (PC), tablet and smartphone (see Fig.3). The AC equips an infrared remote controller. The first type of sensor device is installed around the receiver of the AC so as to receive the infrared signal from the remote controller. The sensor device sniffs the infrared signal and monitors the operation mode of the AC. The sensed information on the operation mode is sent to the server program run
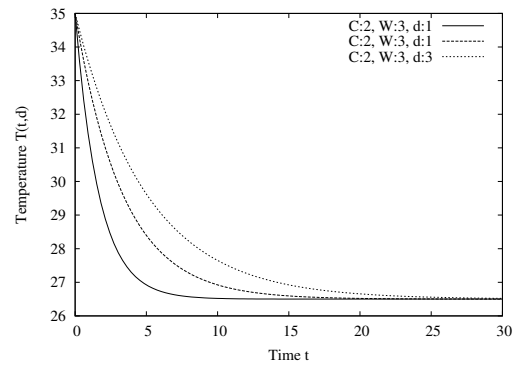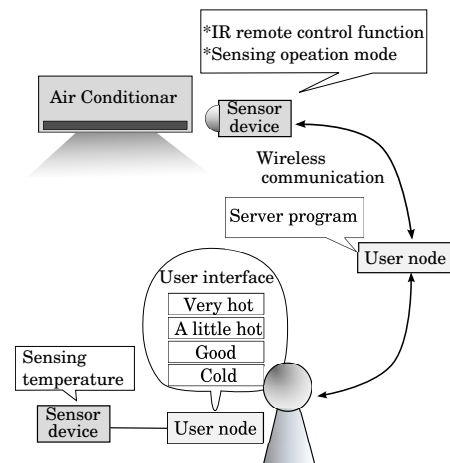


Fig. 3.   Air conditioning control system.

on one of the user nodes by wireless communication. Further, the sensor device also equips the infrared LED, which is used for controlling the AC according to the command from the server program. The second type of the sensor device is installed around the users and each sensor monitors the air temperature around the user. The sensed temperature is also sent to the server program by wireless communication. The primary role of the user nodes is the interface between the user and the system. Through the user interface application run on the node, the user enters his/her sensation such as *cold*, *good*, *a little hot*, *very hot*. The entered sensation is also sent to the server program. The server program is run on one of the user nodes. The program collects the information from the sensor groups and the user nodes, calculates the action for control according to the collected information, and controls the AC by emitting the infrared signal from the sensor device.

#### B. Q-Learning for AC Control Based on User's Sensation

We apply Q-learning[6], which is one of reinforcement learning methods, to air conditioning control based on user's sensory scale. In reinforcement learning, in order to obtain an appropriate behavior pattern (policy) for the environment defined by the target problem, the agent repeats trial and error and updates its policy according to some reward obtained from the environment (see Fig.4). The advantage of reinforcement learning is to not need any teacher signal, which indicates the ideal action at each state. Therefore, it is
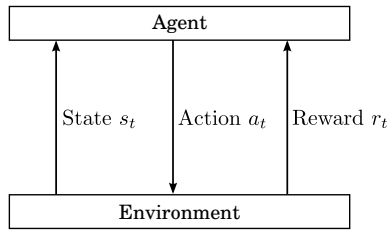
Fig. 4. The concept of reinforcement learning.

applicable to problems with unknown environments. Profit Sharing (PS)[7] and Q-learning are well-known reinforcement learning methods. Compared with PS, it is know that, Q-learning has a higher ability to obtain an optimal policy.

In this subsection, we present the basic form of our Q-learning algorithm for controlling our AC system. In section IV, some key components in the basic form are customized, and four types of the algorithm are presented. Let $N$ be the number of users. Let $S$ be the set of states, in which each member corresponds to the user sensation. $S$ is a key component defined in the next section. Further, the determination of the state involves the integration of $N$ users' sensations, which is also described in the next section. Let $A$ be the set of actions, in which each member corresponds to the AC operation mode $(C, W)$, that is, $A = \{(0,0), (C,W) | C \in \{1,2,3\}, W \in \{1, 1.6, 3\}\}$. Let $Q(s,a)$ be the Q-value of the pair of state $s$ and action $a$, which indicates how much worth is the action $a$ at the state $s$. As the learning progresses, the Q-value of an appropriate pair of state and action increases. At each state $s$, the action with higher Q-value is taken with a higher probability compared with other actions. We assume that the goal of air conditioning is to keep the satisfaction (*good*) of two-thirds users for $T_{\text{sat}}$ time steps. The algorithm is given as follows:

**ALGORITHM Q-Learning for AC Control**
**Initialization:**
  Let $l$ be the current epoch number, and set $l \leftarrow 0$.
  Let $t$ be the current time, and set $t \leftarrow 0$.
  For all $s \in S$ and $a \in A$, $Q(s,a) \leftarrow 0$.
**Step 1 (State observation):** Determine the current state $s_0$ according to the $N$ user inputs.
**Step 2 (Taking action):** For each action $a \in A$, calculate the probability $\pi(s_t, a)$ as follows:

$$\pi(s_t, a) = \frac{\exp(Q(s_t, a_t)/T_{\text{b}})}{\sum_{a' \in A} \exp(Q(s_t, a')/T_{\text{b}})} \qquad (4)$$

$$T_{\text{b}} = T_{\text{b0}} \cdot \left(\frac{T_{\text{b1}}}{T_{\text{b0}}}\right)^{\frac{l}{L_{\max}}} \qquad (5)$$

where $L_{\max}$ is the number of maximum epochs, $T_{\text{b}}$ is the temperature for Boltzmann selection, and $T_{\text{init}}$ and $T_{\text{fin}}$ are maximum and minimum temperatures $T_{\text{b}}$ respectively.

Probabilistically take an action $a_t \in A$ according to the probabilities $\pi(s_t, a)$ $(a \in A)$.
**Step 3 (Reward acquisition and state observation):** Let $r_t$ be the acquired reward, which is a key component described in section IV. Determine the next state $s_{t+1}$ according to $N$ users' inputs.
**Step 4 (Updating Q value):**

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \cdot \delta_t, \qquad (6)$$

where

$$\delta_t = r_t + \gamma \cdot \max_{a \in A} Q(s_{t+1}, a) - Q(s_t, a_t), \qquad (7)$$

$\alpha$ is the learning rate and $\gamma$ is a constant such that $0 \leq \gamma \leq 1$.
**Step 5 (Judgment of goal achievement):** If the goal is reached, that is, two-thirds $N$ users feel *good* for consecutive $T_{\text{sat}}$ time steps, go to the next step. Otherwise, set $t \leftarrow t+1$ and go to Step 2.
**Step 6 (Judgment of termination):** Set $l \leftarrow l + 1$. If $l = L_{\max}$, then terminate the algorithm. Otherwise, for the next epoch, set $T(0,d)$, $T_{\text{o}}$ and the users' sensation scales, which depend on the environment. Set $t \leftarrow 0$ and go to Step 1. $\square$

In section IV, the set of states $S$, the integration of users' sensations and the reward $r_t$ are described in detail.

## IV. INTEGRATION OF USERS' SENSATIONS AND REWARD FUNCTION

In this section, we explain the set of states $S$, the integration of users' sensations and how to give the reward. The naive design of the state set is to represent all the possible combinations of the users' sensations. However, this design needs the number of states that is proportional to the exponential of the number of users $N$. As $N$ increases, the memory size needed by the naive design exponentially increases. Therefore, in this section, we presents three types of integration methods of users' sensations. Further, how to give the reward is an important issue in reinforcement learning. Firstly, we present the reward function using not only the users' sensations and but also the changes of room temperature obtained from the sensor. Then, we also present the reward function using only the users' sensations, which does not need sensing the room temperature.

### A. Naive Design of State Set

The state is coded as a tuple of $N + 1$ elements $(s^{(1)}, s^{(2)}, \cdots, s^{(N)}, s^{\Delta T})$. Each of $N$ elements corresponds to a user sensation, that is, for each $n \in \{1, 2, \cdots, N\}$, $s^{(n)} \in \mathcal{S}$, where $\mathcal{S} = \{cold, \; good, \; a \; little \; hot, \; very \; hot\}$ is the set of user's sensations. The rest one element $s^{\Delta T}$ represents the state of temperature change, which is determined according to the sensing data of a sensor node. Let $\Delta T = T(t, d_{\text{sn}}) - T(t - 1, d_{\text{sn}})$ be the change of the temperature, where $d_{\text{sn}}$ is the distance of the sensor node from AC. The element $s^{\Delta T}$ is defined as follows:

$$s^{\Delta T} = \begin{cases} Down & ; \; \Delta T < -0.5 \\ Unchanged & ; \; |\Delta T| \leq 0.5 \\ Up & ; \; \Delta T > 0.5. \end{cases} \qquad (8)$$

Then, the number of states is $|\mathcal{S}|^N \cdot 3 = 4^N \cdot 3$.

### B. Integration of Users' Sensations

Let us consider to integrate the users' sensations into the collective sensation. We present three types of sensation integration (SI) methods: SI based on Majority Vote (SI-MV), SI based on Averaging (SI-A) and SI based on Likelihood (SI-L).

The first method, SI-MV, integrates $N$ users' sensations into the collective sensation $s_{\text{c}}$ as follows:

$$s_{\text{c}} = \underset{s \in \mathcal{S}}{\arg\max} \, U(s, t), \qquad (9)$$

where $U(s,t)$ is the number of users whose sensations are $s$ at time $t$, $\mathcal{S} = \{$cold, good, a little hot, very hot$\}$ is the set of user's sensations, and if $U(s,t)$s for plural sensations $s$ tie in majority vote then the function $\mathrm{argmax}$ returns the sensation in order of *cold*, *good* and *a little hot*. The state set for SI-MV is $S = \{(s_{\mathrm{c}}, s^{\Delta\mathrm{T}}) | s_{\mathrm{c}} \in \mathcal{S}, s^{\Delta\mathrm{T}} \in \{Down, Unchanged, Up\}\}$. Then, the number of states for SI-MV is $|\mathcal{S}| \cdot 3 = 12$ for any number $N$.

The second method, SI-A, represents each user sensation as a scalar value, calculates the average of the scalar values and returns the integrated one from 12 possible sensations whose segmentation is finer than the one of the original sensation with 4 segments. The mapping from the user sensation $s^{(n)} \in \mathcal{S}$ to a scalar value $v(s^{(n)}) \in \mathcal{V}$ is as follows:

$$v(s^{(n)}) = \begin{cases} 0 & ;\ s^{(n)} = cold \\ 1.0 & ;\ s^{(n)} = good \\ 2.0 & ;\ s^{(n)} = a\ little\ hot \\ 3.0 & ;\ s^{(n)} = very\ hot, \end{cases} \quad (10)$$

where $\mathcal{V} = \{0, 1.0, 2.0, 3.0\}$. The averaging value $\bar{v} = \frac{1}{N}\sum_{n=1}^{N} v(s^{(n)})$ is mapped to the collective value $v^{\mathrm{c}} \in \mathcal{V}^{\mathrm{c}}$ as follows:

$$v^{\mathrm{c}} = \underset{v \in \mathcal{V}^{\mathrm{c}}}{\mathrm{argmax}} |v - \bar{v}|, \quad (11)$$

where $\mathcal{V}^{\mathrm{c}} = \{0.25,\ 0.5,\ 0.75,\ 1.0,\ 1.25,\ 1.5,\ 1.75,\ 2.0, 2.25, 2.5,\ 2.75,\ 3.0\}$. The state set for SI-MV is $S = \{(v^{\mathrm{c}}, s^{\Delta\mathrm{T}}) | v^{\mathrm{c}} \in \mathcal{V}^{\mathrm{c}}, s^{\Delta\mathrm{T}} \in \{Down, Unchanged, Up\}\}$. Then, the number of states for SI-A is $|\mathcal{V}^{\mathrm{c}}| \cdot 3 = 36$ for any $N$.

The last method, SI-L, determines the current state based on the likelihood on the users' sensations at time $t$ and $t-1$. Likewise SI-A, SI-L represents the user sensation as a scalar value. Unlike SI-MV and SI-A, SI-L takes into account not only the user sensations at the current time $t$ but also the one at the previous time $t-1$. With this feature, SI-L does not use the information on the temperature change.

Let $s_t^{(n)}$ and $s_{t-1}^{(n)}$ be the $n$-th user's sensations at time $t$ and $t-1$, respectively. Let $v_t^{(n)}$ and $v_{t-1}^{(n)}$ be the converted scalar values of the $n$-th user's sensations at time $t$ and $t-1$, respectively. The conversion is done by using Eq.(10). Then, the set of states is $S = \{(v_{\mathrm{cur}}^{\mathrm{c}}, v_{\mathrm{prev}}^{\mathrm{c}}) | v_{\mathrm{cur}}^{\mathrm{c}}, v_{\mathrm{prev}}^{\mathrm{c}} \in \mathcal{V}\}$, where $v_{\mathrm{cur}}^{\mathrm{c}}$ and $v_{\mathrm{prev}}^{\mathrm{c}}$ correspond to the collective values at time $t$ and $t-1$, respectively. The number of states for SI-L is $|\mathcal{V}|^2 = 16$.

In the following, the determination of the current state for SI-L is described in detail. The probability density function $p_v$ of the user sensation corresponding to the converted scalar value $v \in \mathcal{V}$ is defined as follows:

$$p_v(v^{(n)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(v^{(n)} - v)}{2\sigma^2}\right), \quad (12)$$

where $\sigma^2 = 1$ is the variance of the distribution. Given the $N$ users' sensation scalar values at time $t$ and $t-1$, $v_t^{(1)}$, $v_t^{(2)}, \cdots, v_t^{(N)}, v_{t-1}^{(1)}, v_{t-1}^{(2)}, \cdots, v_{t-1}^{(N)}$, then the likelihood of each state $(v_{\mathrm{cur}}, v_{\mathrm{prev}}) \in S$ is calculated as follows:

$$L(v_t^{(1)}, \cdots, v_t^{(N)}, v_{t-1}^{(1)}, \cdots, v_{t-1}^{(N)}; v_{\mathrm{cur}}; v_{\mathrm{prev}}) =$$
$$\prod_{n=1}^{N} p_{v_{\mathrm{cur}}}(v_t^{(n)}) \cdot p_{v_{\mathrm{prev}}}(v_{t-1}^{(n)}). \quad (13)$$

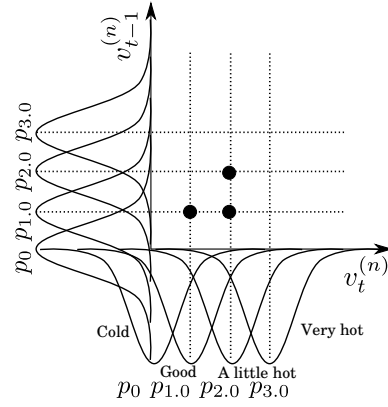

Fig. 5. The probability density function arrangement for SI-L and an input example for 3 users.

The collective state $s^{\mathrm{c}} = (v_{\mathrm{cur}}^{\mathrm{c}}, v_{\mathrm{prev}}^{\mathrm{c}}) \in S$ is calculated as follows:

$$s^{\mathrm{c}} =$$
$$\underset{(v_{\mathrm{cur}}, v_{\mathrm{prev}}) \in S}{\mathrm{argmax}} L(v_t^{(1)}, \cdots, v_t^{(N)}, v_{t-1}^{(1)}, \cdots, v_{t-1}^{(N)}; v_{\mathrm{cur}}; v_{\mathrm{prev}}).$$
$$(14)$$

### C. Reward Function

We consider three types of reward functions. The functions take into account the following issues.

- $R_{\mathrm{L}} > 0$: Reward for goal achievement.
- $R_{\mathrm{S}} > 0$: Reward for keeping a good control.
- $R_{\mathrm{N}} < 0$: Penalty on the degradation on the users' sensations.
- $R_{\mathrm{C}} < 0$, $\rho \cdot R_{\mathrm{L}}$, $\rho \cdot R_{\mathrm{S}}$: Penalty for excess cooling, which means dissipation of energy, where $0.0 < \rho < 1.0$.

The first type is Excess-Cooling-Unaware Reward Function (ECU-RF). ECU-RF does not care about any excess cooling, that is, it does not use $R_{\mathrm{C}} < 0$, $\rho \cdot R_{\mathrm{L}}$ nor $\rho \cdot R_{\mathrm{S}}$. Fig.6 shows how ECU-RF determines the reward $r_t$, where $U(s,t)$ is the number of users whose sensations are $s \in \mathcal{S}$ at time $t$, and $M_{\mathrm{S}} = \lceil \frac{2}{3}N \rceil$ is the minimum number of users to be satisfied.

The second type is Soft-Penalty-for-excess-Cooling Reward Function (SPC-RF). SPC-RF reduces the reward with the reduction rate $\rho$ when the AC control is excess cooling. Fig.7 shows how SPC-RF determines the reward $r_t$, where $M_{\mathrm{C}}$ is the number of users that are not allowed to feel *cold*.

The third type is Hard-Penalty-for-excess-Cooling Reward Function (HPC-RF). HPC-RF gives a negative reward $R_{\mathrm{C}}$ when the AC control is excess cooling. Fig.8 shows how HPC-RF determines the reward $r_t$.

## V. NUMERICAL SIMULATION

### A. Simulation Setting

The virtual users used in the simulation have sensory scales as shown in Fig.9, where thresholds $\theta_1$, $\theta_2$ and $\theta_3$ represent the boundaries between *cold* and *good*, between *good* and *a little hot* and between *a little hot* and *very hot*.

The simulation conditions are summarized in Table II. The number of users is $N = 12$ or 3. The equal number of virtual
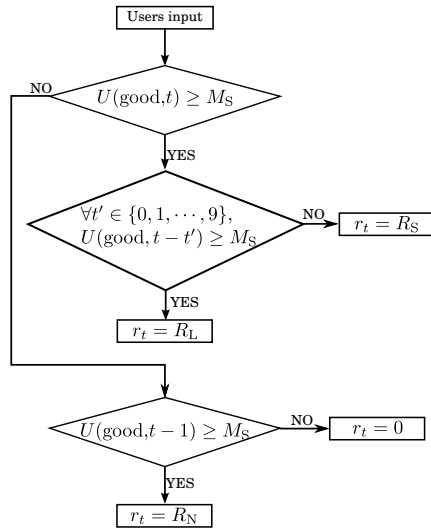
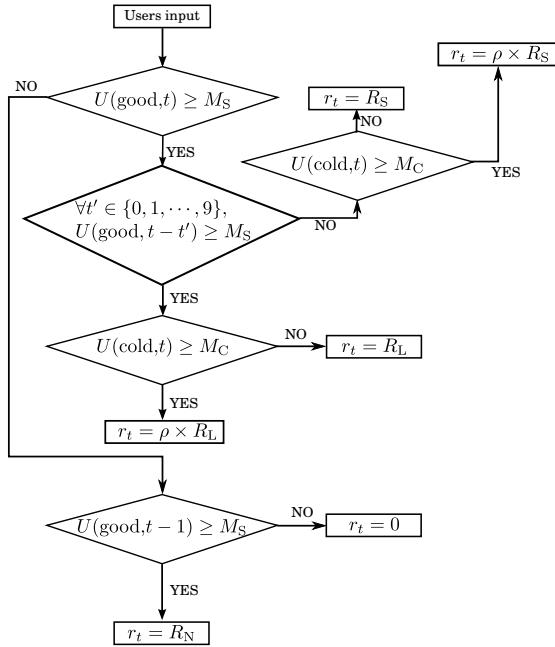Fig. 6. The definition of Excess-Cooling-Unaware Reward Function (ECU-RF).



Fig. 8. The definition of Hard-Penalty-for-excess-Cooling Reward Function (HPC-RF).



Fig. 9. Sensory scale of virtual user.



Fig. 7. The definition of Soft-Penalty-for-excess-Cooling Reward Function (SPC-RF).

In Fig.10, a smaller time means that the users may enter their sensation fewer times. Although all the methods need a large time at the first epoch, the needed time quickly converges to reasonable values. For the needed time, the method (L, HPC) is the best among all the methods. The methods (A, HPC) or (w/o, SPC) are the second best. However, the method (w/o, SPC) is for $N = 3$ users. Therefore, it is notable that the methods (L, HPC) and (A, HPC)

users is arranged at each distance $d \in \{1, 2, 3\}$. At each epoch, set randomly the initial room temperature $T(0, d)$, the outer air temperature $T_o$, and the target temperature and the thresholds $\theta_1$, $\theta_2$ and $\theta_3$ for each user.

We evaluate the four types of methods as shown in Table III. The used parameters for the evaluated methods are summarized in Table IV. These parameters are determined by preliminary simulations.

### B. Simulation Result

We evaluate four types of methods (w/o, SPC), (MV, ECU), (A, HPC) and (L, HPC) in terms of the time needed for reaching goal and the averaging amount of cold air needed for reaching goal. The simulation results are shown in Figs. 10 and 11.
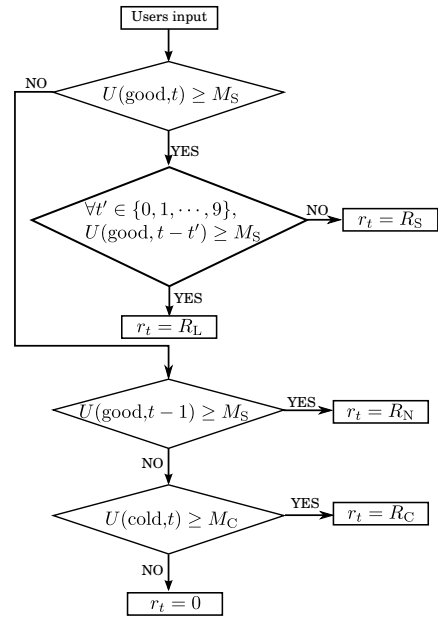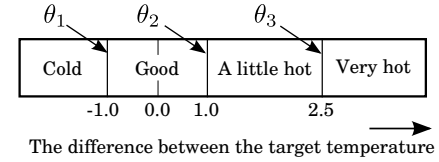
TABLE II
SIMULATION CONDITIONS.

| Item | Used values |
|---|---|
| User's dist. from AC $d$ | 1, 2, 3 |
| Target temperature | 26.0 $\pm$0.5 °C |
| Initial temperature $T(0, d)$ | Same as outer air temperature |
| Outer air temperature $T_o$ | 30$\sim$35 °C |
| $\theta_1$ | $-1.0$ $\pm$0.3 °C |
| $\theta_2$ | 1.0 $\pm$0.3 °C |
| $\theta_3$ | 2.2 $\pm$0.3 °C |
| Number of users to be satisfied $M_S$ | $\lceil \frac{2}{3}N \rceil$ |
| Maximum number of epochs $L_{max}$ | 500 |
| Duration of satisfaction state | 10 minutes |

TABLE III
EVALUATED METHODS.

| Name | Sensation Integration | | | | Reward Function | | |
|---|---|---|---|---|---|---|---|
| | w/o | MV | A | L | ECU | SPC | HPC |
| (w/o, SPC) | o | | | | | o | |
| (MV, ECU) | | o | | | o | | |
| (A, HPC) | | | o | | | | o |
| (L, HPC) | | | | o | | | o |

TABLE IV
PARAMETERS FOR EACH METHOD.

| Item | (w/o,SPC) | (MV, ECU) | (A, HPC) | (L, HPC) |
|---|---|---|---|---|
| Number of users $N$ | 3 | 12 | 12 | 12 |
| Discount rate $\gamma$ | 0.7 | 0.7 | 0.9 | 0.9 |
| Learning rate $\alpha$ | 0.2 | 0.2 | 0.5 | 0.3 |
| Max. temp. in B.S. $T_{b0}$ | 60 | 60 | 5.0 | 3.0 |
| Min. temp. in B.S. $T_{b1}$ | 50 | 50 | 0.3 | 0.3 |
| Reduction rate $\rho$ | 0.3 | 1.0 | – | – |
| $M_S$ | 1 | 8 | 8 | 8 |
| $M_C$ | 1 | – | 5 | 5 |
| $R_L$ | 400 | 400 | 400 | 40 |
| $R_S$ | 30 | 30 | 3 | 3 |
| $R_N$ | −400 | −400 | −400 | −20 |
| $R_C$ | – | – | – | −10 |

for $N = 12$ users outperform the method (w/o, SPC) for $N = 3$ users. This result demonstrates that our methods for integrating sensations, SI-L and SI-A, are effective. Further, the proposed reward function HPC is also effective. The other method (MV, ECU) requires the largest time for reaching goal. Therefore, the majority vote is not good for the needed time.

In Fig.11, a smaller amount of cold air means more energy efficient. We obtain a little bit unexpected results. Among them, the method (MV, ECU) is the best despite that it does not take into account excess cooling. This is because (MV, ECU) prefers lower cooling strength. Therefore, the needed time was largest. It is observed that both of (A, HPC) and (L, HPC) increase with the number of epoch $l$. This is because their needed time decreases with the number of epoch $l$. We can observe that there are tradeoff between the needed time and the energy efficiency. However, the difference on the amount of cool air is not so large, that is, at most approximately 2%.

## VI. CONCLUSION

In this paper, we proposed the air conditioner (AC) control system based on users' sensations. The purpose was to realize the control system that improves the low-performance AC in terms of energy efficiency and comfortableness performance. The system trained the control policy by using Q-learning with entered users' sensations. In order to cope with the large number of users, we proposed three types of integration methods of users' sensations. Especially, the proposed method based on likelihood was the most effective in terms of the time required for the control. Further, in order to reduce the energy consumption in cooling, we proposed several reward functions, which penalize the policy performing excess cooling. However, we could not confirm that the penalties clearly reduce the energy consumption. Instead, we confirmed that there exists a tradeoff between the speed of control and the energy consumption. One of our future works is to overcome the tradeoff.



Fig. 10. Time needed for reaching goal at each epoch $l$.



Fig. 11. Averaging amount of cold air needed for reaching goal at each epoch $l$.

## REFERENCES

[1] K. Sato, M. Samejima, M. Akiyoshi and N. Komada, "A Scheduling Method of Air Conditioner Operation using Workers Daily Action Plan towards Energy Saving and Comfort at Office," *2012 IEEE 17th Conference on Emerging Technologies and Factory Automation*, pp. 1-6, 2012.

[2] M. Bal, "Industrial applications of collaborative Wireless Sensor Networks: A survey," *2014 IEEE 23rd International Symposium on Industrial Electronics*, pp. 1463-1468, 2014.
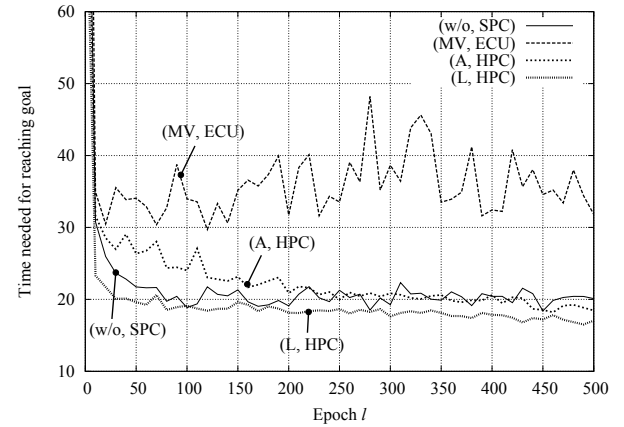
[3] C. Tunca, H. Alemdar, H. Ertan, O.D. Incel and C. Ersoy, "Multimodal Wireless Sensor Network-Based Ambient Assisted Living in Real Homes with Multiple Residents," *Sensors*, Vol. 14, No. 6, pp. 9692-9719, 2014.

[4] J. Si, A. Barto, W. Powell and D. Wunsch, "Robust Reinforcement Learning for Heating, Ventilation, and Air Conditioning Control of Buildings," *Handbook of Learning and Approximate Dynamic Programming*, pp.517-534, Wiley-IEEE Press, 2004.

[5] T. Hiroyasu, A. Nakamura, M. Yoshimi, M. Miki and H. Yokouchi, "Lighting Control System using an Actor-Critic type Learning Algorithm," *2010 Second World Congress on Nature and Biologically Inspired Computing*, pp. 140-145, 2010.

[6] C.J.C.H. Watkins, "Learning from Delayed Rewards," Ph D Thesis, University of Cambridge, England, 1989.

[7] J.J. Grefenstette, "Credit Assignment in Rule Discovery Systems Based on Genetic Algorithms," *Machine Learning*, Vol. 3, pp. 225-245, 1988.