

Bottom-up Social development through reproducing contingency with sensorimotor clustering

Hidenobu Sumioka¹

Yuji Takeuchi¹

Yuichiro Yoshikawa²

Minoru Asada^{1,2}

¹ Graduate School of Eng., Osaka Univ.

² JST ERATO Asada Synergistic Intelligence Project

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

{hidenobu.sumioka,yuji.takeuchi,asada}@ams.eng.osaka-u.ac.jp
yoshikawa@jeap.org

Abstract

This paper presents a learning mechanism that finds the reasonable segmentation to achieve social behavior as well as incrementally acquires social behavior by reproducing the contingency in interaction with a caregiver. The robot autonomously categorizes sensorimotor activity according to a contingency measure based on transfer entropy. The advantage of adaptive categorization is tested in a task of acquiring joint attention behaviors. The results of computer simulations of human-robot interaction indicate that a robot acquires a series of joint attention behavior such as gaze following and alternation and it finds suitable segmentation that improves the performance of gaze following over time.

1. Introduction

Human infants acquire a variety of social behavior through interaction with others. In particular, joint visual attention is one of the building blocks for social capabilities such as language communication and mind-reading (Moore and Dunham, 1995). Therefore, understanding how infants acquire a variety of joint attention behavior such as gaze following, gaze alternation, and pointing is a central topic in developmental psychology. However, it remains a mystery how infants acquire such behavior.

In robotics, joint attention studies have been recently receiving increasing attention not only from the viewpoint of building communicative robots (Imai et al., 2001) but also from synthetic approaches for modeling and understanding human developmental processes (Nagai et al., 2003, Triesch et al., 2006) as argued in surveys (Kaplan and Hafner, 2004, Asada et al., 2009). One of them has addressed how a robot can acquire different joint attention behavior (Sumioka et al., 2008). Sumioka *et al*

emphasized a statistical structure based on the fact that infants can often attain consistent consequences when they respond adequately to a preceding stimulus including behavior of their caregivers. Such a structure of the relationship among a preceding stimulus, one's own action, and its consequence, called contingency, was utilized to find more contingent sets including sensory and motor variables that provides consistent consequences to a robot among several candidates, and to construct sensorimotor maps based on the found sets. The results of computer simulations of human-robot interaction indicated that finding the contingency and taking actions to reproduce it enable a robot to acquire a series of joint attention behavior such as gaze following and alternation in an order that is almost the same order of infant development.

In their study, each random variable was quantized in advance into the reasonable segments sufficient to reproduce contingencies of interaction between a robot and a caregiver so that a robot can acquire social behavior. However, it is not trivial for a robot to quantize a variable adequately since the most reasonable segmentation depends on the sensor resolution of the robot, the control resolutions of the caregiver's behavior and the robot's one, and an observed object size and its location. If the robot found the contingency based on rough segmentation, the contingency is expected to include stronger contingency based on more sophisticated segmentation. Therefore, we utilize the contingency measure to obtain the most reasonable segmentation, that is, we segment a variable so that a robot can reproduce the contingency. We hypothesize that quantizing variables to experience more contingent consequences leads the robot to acquire social behavior that enables it to interact with its caregiver adequately.

This paper presents a learning mechanism to quantize each variable adaptively as well as to find the contingency and reproduce the contingent relationship. The contingency measure based on information

theory (Sumioka et al., 2008) is utilized for adaptive quantizing. The advantage of adaptive quantizing is tested in a task of acquiring joint attention behavior. The results of computer simulations of human-robot interaction indicate that a robot acquires a series of joint attention behavior such as gaze following and alternation and it finds suitable segmentation that improves the performance of gaze following over time.

2. Face-to-face interaction to develop joint attention behavior

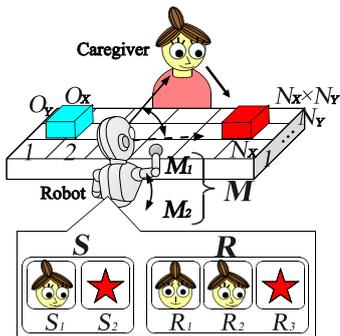


Figure 1: An environmental setting

To examine whether a robot can acquire a variety of joint attention behavior with quantizing sensory variables and motor ones, we start from a rough model of the caregiver's gaze shift. We simulate almost the same interaction as in previous studies (Sumioka et al., 2008).

Figure 1 shows an environmental setting of the interaction. The robot sits across from the caregiver at a fixed distance. An interaction where each of the caregiver and the robot take an action at once is defined as a time step. There is a table having $N_X \times N_Y$ sections where two same objects each of which occupies $O_X \times O_Y$ sections are randomly placed. The positions of objects are determined randomly every ten steps.

In an interaction, the caregiver observes her environment and then shifts her gaze to the robot or an object according to a few policies as described in section 4.1.2. Next, the robot observes its environment and obtains the information about the direction of the caregiver's face (S_1) and the presence of an object (S_2) as sensory variables. It also stores the information about what it is looking at as the result of its actions called resultant sensory variables: caregiver's frontal face (R_1), caregiver's profile (R_2), and the presence of an object (R_3). Finally, it shifts its gaze to the caregiver or a table section and shows a hand gesture, and stores a motor command for gaze shift (M_1) and one for hand gesture (M_2) as motor variables.

Here, a contingency inherent in the interaction appears as a dependency of a resultant sensory variable on a sensory variable and a motor variable. We call a triplet of variables (S_i, M_j, R_k) an event variable. Moreover, an event variable that involves strong dependency is called a contingent event variable. The task of the robot is performed by finding a contingent event variable and then acquiring a sensorimotor mapping based on the found event variable. Moreover, the robot has to determine how it should quantize sensory variables and motor variables.

3. Proposed mechanism to successfully develop social behavior with adaptive partitioning

Instead of a designer quantizing a random variable into several segments in advance, we make the robot quantize them autonomously. A contingency measure proposed by Sumioka *et al.* (Sumioka et al., 2008) is utilized for quantizing variables as well as constructing sensorimotor mappings. The proposed architecture shown in Figure 2 consists of three features: (1) a contingency monitor that sends commands to quantize sensory variables and motor ones, (2) State/Motor categorizer to output one of components in sensory variables and motor ones according to the observed features and motor commands, and (3) sequential contingency learning module to enable the robot to acquire several actions by finding contingency of interaction and its reproduction.

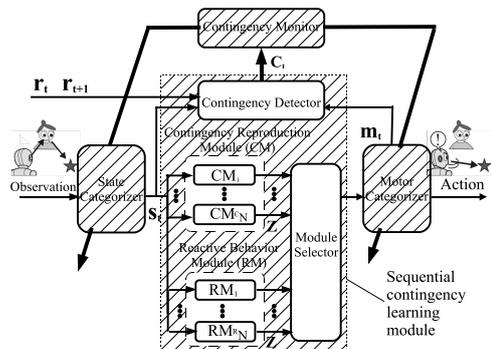


Figure 2: A proposed mechanism

Observed sensory features are categorized by the state categorizer as one of the components in each sensory variable. The selected components are sent to the sequential contingency learning module. The sequential contingency learning module decides one of the components in each motor variable according to the acquired sensorimotor mapping or innate behavior policies as described in section 3.2.3. Finally, the motor categorizer selects motor commands based on the selected components. During this process,

contingency of an event variable is evaluated in the sequential contingency learning module by calculating the contingency measure (Sumioka et al., 2008) as described in next section. According to the measure, the contingency monitor commands the state categorizer and the motor one to update the segmentation in sensory variable and motor one.

3.1 Contingency measure

Sumioka *et al.* proposed an information theoretic measure of contingency, called as saliency of contingency (C-saliency), based on transfer entropy (Schreiber, 2000) to quantify contingency of an event variable (Sumioka et al., 2008).

Suppose that two time series variables X and Y may be approximated by first-order Markov processes and that they form the following contingency: x^{t+1} , i.e., the value of X at time $t + 1$, is only influenced by x^t and y^t , i.e., the values of X and Y at the previous time t . Here, the transfer entropy that indicates the influence of Y on X is defined by

$$T_{Y \rightarrow X} = \sum p(x^{t+1}, x^t, y^t) \log \frac{p(x^{t+1}|x^t, y^t)}{p(x^{t+1}|x^t)}. \quad (1)$$

C-saliency $C_{i,k}^j$ to quantify the joint effect of sensory variable S_i and motor variable M_j on resultant sensory variable R_k is defined as

$$\begin{aligned} C_{i,k}^j &= T_{(S_i, M_j) \rightarrow R_k} - (T_{S_i \rightarrow R_k} + T_{M_j \rightarrow R_k}) \\ &= \sum_{s_i^t, r_k^t} p(r_k^t, s_i^t) \sum_{r_k^{t+1}, m_j^t} e(r_k^{t+1}, m_j^t | r_k^t, s_i^t), \end{aligned} \quad (2)$$

where $e(r_k^{t+1}, m_j^t | r_k^t, s_i^t)$ is called an element of C-saliency under a pair of observed values (r_k^t, s_i^t) and is given by:

$$\begin{aligned} e(r_k^{t+1}, m_j^t | r_k^t, s_i^t) &= \\ & p(r_k^{t+1}, m_j^t | r_k^t, s_i^t) \log \frac{p(r_k^{t+1} | r_k^t, s_i^t, m_j^t)}{p(r_k^{t+1} | r_k^t, s_i^t)} \\ & - p(r_k^{t+1}, m_j^t | r_k^t) \log \frac{p(r_k^{t+1} | r_k^t, m_j^t)}{p(r_k^{t+1} | r_k^t)}. \end{aligned} \quad (3)$$

The element of C-saliency represents strength of the dependency of the state transition from r_k^t to r_k^{t+1} on pair (s_i^t, m_j^t) . If triplet (r_k^t, s_i^t, m_j^t) causes r_k^{t+1} , the difference between $p(r_k^{t+1} | r_k^t, s_i^t, m_j^t)$ and $p(r_k^{t+1} | r_k^t, s_i^t)$ becomes larger. An event variable with the highest C-saliency is regarded as a contingent event variable.

Additionally, C-saliency has an interesting feature to evaluate the performance of an acquired sensorimotor map. The sensorimotor map that usually causes contingent consequences enables a robot to predict the state transitions of a resultant sensory

variable only by states of a sensory variable. C-saliency related to such a map gets lower because the value of the first term in Eq. (3) reduces. Therefore, the derivative of C-saliency is useful to evaluate the accuracy of the acquired sensorimotor map: if the derivative is negative, the robot has acquired a sensorimotor map sufficient to reproduce the contingency while the robot needs to quantize variables related to the map if the derivative is not negative.

3.2 Components in proposed mechanism

C-saliencies of event variables is utilized by the proposed mechanism not only to find contingency and to reproduce it but also to quantize variables based on more reasonable segmentation. Here, the roles of the components in the mechanism are described.

3.2.1 Contingency monitor

The contingency monitor modulates the quantization of sensory variables and motor ones. The quantization of a variable consists of two processes: how a variable should be quantized by the existing segments (arrangement process) and how many segments a variable should be quantized into (insertion process). In each process, we use a derivative of C-saliency for an event variable. Here, the derivative of C-saliency for an event variable (S_i, M_j, R_k) at t time steps is indicated as $\Delta C_{i,k}^j(t) = C_{i,k}^j(t) - C_{i,k}^j(t-1)$, where $C_{i,k}^j(t)$ indicates C-saliency for (S_i, M_j, R_k) at t time steps.

The contingency monitor basically uses the arrangement process where segments in a variable are modulated so that C-saliency of an event variable including the variable becomes higher. In the process, the contingency monitor sends the state categorizer and the motor one the values that determine how much state categorizer and motor one should update segment arrangement. The values ΔC_{max}^S and ΔC_{max}^M for a sensory variable S_i and a motor variable M_j are given by $\Delta C_{max}^S (= \max_{j,k} \Delta C_{i,k}^j)$ and $\Delta C_{max}^M (= \max_{i,k} \Delta C_{i,k}^j)$, respectively. To avoid modulation during a period when C-saliencies are overestimated due to few samples, these values are sent when variance for moving average of each C-saliency during T^A time steps, $\sigma_{A_{i,k}^j}$, is lower than ε_A .

The contingency monitor has a possibility of using the insertion process, which decides whether it should insert new segments into a sensory variable or a motor one, after the contingency detector selects a contingent event variable and generates a CM as described in section 3.2.3. Let ${}^c C_{i,k}^j$ C-saliency for a contingent event variable (S_i, M_j, R_k) . New segments are inserted to S_i and M_j when the variance

$\sigma_{I_{i,k}^j}$ for the moving average of its derivative $\Delta^c C_{i,k}^j$ keeps a lower value than ε_S during T^I time steps. Once new segments are inserted into a sensory variable and a motor one, the insertion process is not applied for those variables during T^D time steps at least.

3.2.2 Sensory/Motor Categorizer

State/Motor categorizer outputs one of segments in each of sensory variable or motor one for given inputs. Suppose that a variable V is quantized into N_v codebook vectors and vector ${}^v\mathbf{a}_\ell$ represents segment ${}^v c_\ell$ ($\ell = 1, 2, \dots, N_v$). When vector ${}^v\mathbf{x}$ related to V is input to state(motor) categorizer, the state(motor) categorizer selects segment ${}^v c_\ell$ with probability $P(V = {}^v c_\ell)$:

$$P(V = {}^v c_\ell) = \frac{\exp\{1/(\tau_v \|{}^v\mathbf{x} - {}^v\mathbf{a}_\ell\|)\}}{\sum_{q=1}^{N_v} \exp\{1/(\tau_v \|{}^v\mathbf{x} - {}^v\mathbf{a}_q\|)\}}, \quad (4)$$

where τ_v is a positive constant.

The selected codebook vector ${}^v\mathbf{a}_\ell$ is updated according to ΔC (which stands for ΔC_{max}^S or ΔC_{max}^M described in previous section) related to an event variable including V :

$${}^v\mathbf{a}_\ell^{t+1} = {}^v\mathbf{a}_\ell^t + \eta \Phi_{v_x, v_{a_\ell}} \Psi_{\Delta C} [{}^v\mathbf{x}^t - {}^v\mathbf{a}_\ell^t] \quad (if \quad {}^v\mathbf{x}^t \in {}^v c_\ell), \quad (5)$$

where, η is a learning rate. $\Phi_{v_x, v_{a_\ell}}$ is given by $\Phi_{v_x, v_{a_\ell}} = \exp\left(\frac{-\|{}^v\mathbf{x}^t - {}^v\mathbf{a}_\ell^t\|}{\zeta}\right)$, where ζ is a constant value. $\Psi_{\Delta C}$ is defined as $\Psi_{\Delta C} = \xi \tanh(\Delta C)$ and ξ is constant.

In addition, the sensory categorizer or the motor one inserts new codebook vectors for a variable when the insertion process is applied for the variable by the contingency monitor. Each categorizer decides where to insert the vectors according to the policy described in section 4.

3.2.3 Sequential contingency learning module

We used the learning module proposed by Sumioka *et al* (Sumioka *et al.*, 2008) that consists of a contingency detector to calculate C-saliencies for all event variables, contingency reproduction modules (CMs) that construct sensorimotor maps to reproduce found contingency, reactive behavior modules (RMs) to output a motor command based on a fixed behavior policy, and a module selector to select motor commands among several outputs from CMs and RMs.

The sequential contingency learning module keeps acquiring different sensorimotor mappings as follows. At the beginning of learning, there are no CMs.

Therefore, the module selector selects the outputs of RMs. As interaction between a caregiver and the robot is iterated, the contingency detector finds a contingent event variable and then generates a new CM that constructs a sensorimotor map to reproduce the found contingency. Once a CM is generated, The module selector starts to select an output from the CM as well as ones from RMs. This iteration of finding contingency and its reproduction, the robot acquire several actions.

Note that whenever a new CM is generated, a new sensory variable S^Π and a new motor one M^Π are added to their sets to indicate whether the new CM was used and is going to be used, respectively. The contingency detector also starts to evaluate new event variables including them. Such event variables may be selected as a next contingent event variable if the found contingency leads novel contingency. Therefore, the robot is expected to find a series of contingent events.

The sensorimotor map in CM is modulated every 200 time steps to utilize more reasonable segmentation. Hereafter i -th CM that is constituted for event variable (S_i, M_j, R_k) is defined as $\Pi_i(R_k | S_i, M_j)$.

4. Experiment

We conducted computer simulations to test whether the proposed mechanism can acquire joint attention actions in different environments. We first examine whether a robot acquires a series of actions related to joint attention such as gaze following and alternation, i.e., successive looking between a caregiver and an object, in simple environmental setting. The size of objects is then changed to show that the mechanism can find the more reasonable segmentation. After that, the performance of the mechanism is tested in more complex situations where a robot has to deal with high-dimensional information or it has a field of view as bias inherent in human. In all experiments, policies for RMs and parameters were set so that a robot can find the contingency related to gaze following at least.

4.1 Experimental setting

4.1.1 Environment and infant model

The initial set of variables is listed in Table 1. The sensory variable for the caregiver's face is denoted by S_1 which consists of two segments (${}^{S_1}c_1$ and ${}^{S_1}c_2$) and two additional components indicating that an infant model (hereafter a robot) is looking at her frontal face (f_r) and that it does not look at the caregiver (f_ϕ), respectively. In the experiments, we fixed the number of components in the sensory variable for an object S_2 . Each member of S_2 indicates whether the robot is looking at an object (o) or at

something else (o_ϕ).

The resultant variables R_1 , R_2 , and R_3 are designed as binary variables indicating whether the robot is looking at its preferred face or an object ("1") or not ("0"). The robot's gaze shift denoted by M_1 consists of two segments ($^{M_1}c_1, ^{M_1}c_2$) and an additional state g_c indicating shifting its gaze to the caregiver's face. Likewise, the gesture denoted by M_2 consists of two segments ($^{M_2}c_1, ^{M_2}c_2$) and h_c indicating pointing its hands to caregiver's face.

Two RMs are used to determine gaze movements and hand by selecting a component of M_1 and M_2 . The RM for M_1 is designed to select either g_c with probability 0.1 or one segment with probability 0.9 while the RM for M_2 selects a component of M_2 randomly. The parameters in the proposed mechanism are set as $(T^A, T^I, T^D, \varepsilon_A, \varepsilon_I, \tau_v) = (2.0 \times 10^3, 5.0 \times 10^3, 2.0 \times 10^5, 1.0 \times 10^{-10}, 1.0 \times 10^{-12}, 2.0 \times 10^{-2})$. The joint and conditional probabilities to calculate C-saliencies were estimated using the histograms of the values of event variables.

Table 1: Initial variables in robot

Type	Name	Elements
S	caregiver's face	$S_1 = \{S_1 c_1, S_1 c_2, f_r, f_\phi\}$
	object	$S_2 = \{o, o_\phi\}$
M	gaze shift	$M_1 = \{^{M_1}c_1, ^{M_1}c_2, g_c\}$
	hand gesture	$M_2 = \{^{M_2}c_1, ^{M_2}c_2, h_c\}$
R	frontal face of caregiver	$R_1 = \{0, 1\}$
	profile of caregiver	$R_2 = \{0, 1\}$
	object	$R_3 = \{0, 1\}$

4.1.2 Behavior rules for caregivers

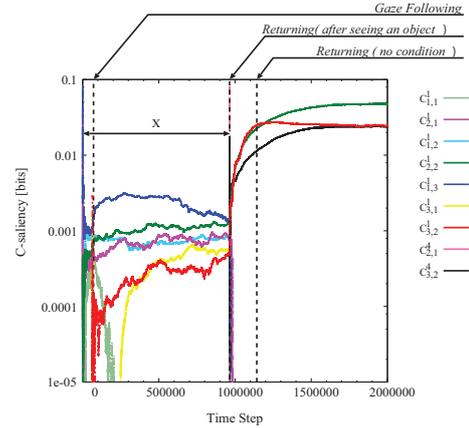
We used the caregiver model described in the previous study (Sumioka et al., 2008). The caregiver, who always looks at the robot's face or an object on the table, not only randomly selects a target but also shows joint attention behavior.

She usually selects a target randomly. If she is looking at the robot's face, she follows the robot's gaze with probability p_{rja}^c . If she is looking at an object, she shifts her gaze between the robot and an object with probability p_{ija}^c . In addition, the caregiver shifts the gaze to the robot's face with probability p_{aja}^c if she and the robot successfully look at the same object. In the following experiments, we used $(p_{rja}^c, p_{ija}^c, p_{aja}^c) = (0.5, 0.5, 1.0)$.

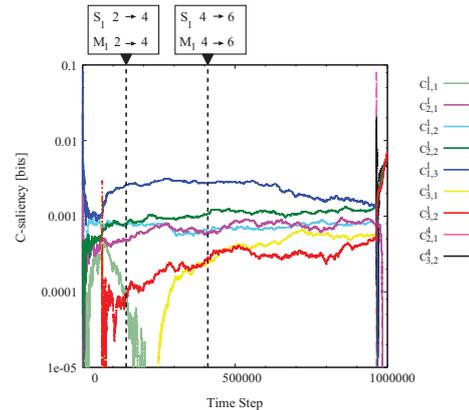
4.2 Development of joint attention with adaptive segmentation

We first confirmed that the proposed mechanism enables a robot to acquire a variety of actions related to

joint attention with quantizing sensory variables and motor variables. We ran 2,000,000 time step simulations five times where two objects each of which occupies 5×1 sections were arranged on a table having 50×1 sections. We set $(\eta, \zeta, \xi) = (0.1, 10, 6.0 \times 10^4)$. In each simulation, two codebook vectors were added on the positions of two existing codebook vectors selected randomly when contingency monitor selected insertion process for a sensory variable and a motor one.



(a) The timing of generating CMs



(b) The timing of increasing codebook vectors

Figure 3: Time courses of saliency of contingency of event variables in simulation face-to-face interactions between caregiver and robot

An average of 2.8 CMs was obtained. In the 80 percent of the simulations, a particular set of CMs was generated in the following fixed order: $\Pi_1(R_3|S_1, M_1)$, $\Pi_2(R_2|S_2, M_1)$, and then $\Pi_3(R_2|S_3^{\Pi_1}, M_1)$. Each of these CMs allowed the robot to achieve social behavior: following the caregiver's gaze ($\Pi_1(R_3|S_1, M_1)$; hereafter called *following-gaze* module), shifting its gaze to the caregiver after seeing an object ($\Pi_2(R_2|S_2, M_1)$; hereafter called *returning (seeing-object)* module), and shifting its gaze to the caregiver regardless of the achievement of gaze following ($\Pi_3(R_2|S_3^{\Pi_1}, M_1)$;

hereafter called *returning (no-condition)* module).

Figure 3 shows examples of the time courses of C-saliencies for nine event variables of which C-saliencies are higher than others. The vertical axis indicates the logarithmic value of the C-saliencies. We also show the timing of generating new CMs as arrows at the top of the graph in Figure 3(a). After sufficient interaction data was collected, $C_{1,3}^1$ became the highest among all C-saliencies (blue curve in Figure 3(a)). As a result, a new CM ($\Pi_1(R_3|S_1, M_1)$) corresponding to the *following-gaze* module was generated, and $S_3^{\Pi_1}$ and $M_3^{\Pi_1}$ were added as sensory and motor variables, respectively. The robot then began to follow the caregiver’s gaze by the *following-gaze* module. However, the success rate of gaze following is not so high at many areas on a table (see Figure 4(a)) because S_1 and M_1 have only two segments, that is, the robot classifies the face of the caregiver that is looking on the table as only two different patterns. In this case, $C_{1,3}^1$ does not decrease since segmentation is not reasonable to achieve gaze following. Therefore, new segments are inserted into S_1 and M_1 (see Figure 3(b)). Finally, S_1 and M_1 had an average of 6.4 segments. The codebook vectors in each variable were arranged at almost equal distance at the end of the simulation (Figure 5). We can see that the found segment arrangement enables a robot to achieve gaze following for the caregiver successfully (see Figure 4(b)).

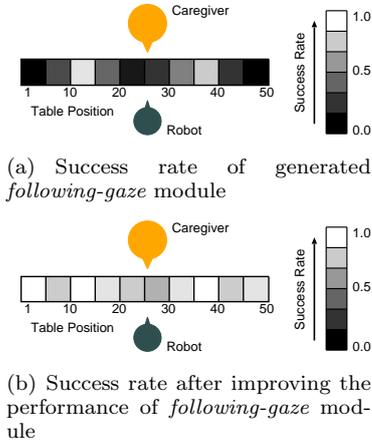
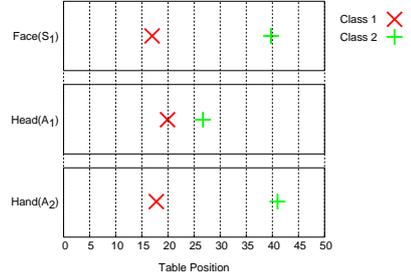


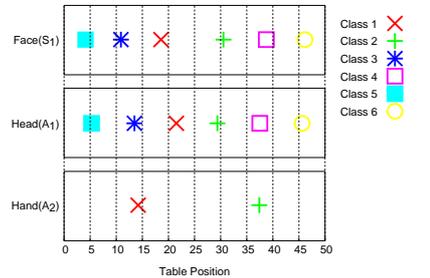
Figure 4: Changes in success rate of gaze following

The increase of success rate of gaze following led $C_{1,3}^1$ to decrease gradually. This decrease made $C_{2,2}^1$ the next highest value, and the *returning (seeing-object)* module was generated. Using output from this module changed the contingency in interaction and promoted an increase of $C_{3,2}^1$ (red curve in Figure 3(a)). This caused the generation of the *returning (no-condition)* module. This enabled the robot to shift its gaze to the caregiver regardless of following the caregiver’s gaze or not. As a result, the

robot alternately shifted its gaze between the caregiver and an object. This indicates that the robot acquired gaze following and alternation with finding the reasonable segment arrangement.



(a) The displacement of codebook vectors before learning



(b) The displacement of codebook vectors after learning

Figure 5: Transition of codebook vectors

4.3 Performance of adaptive segmentation

The environmental features such as the size of a table or an object affect how many segments are needed to achieve gaze following. We examined to what extent the robot can maintain the high performance of gaze following in several different situations where objects with the different size are arranged.

We ran simulations where the size of objects is different. To show the advantage of the proposed mechanism, we also tested mechanisms without arrangement process and insertion process: S_1 , M_1 , and M_2 were quantized into the fixed segments (four, eight, or twelve segments) that were arranged at equal distance in advance.

Figure 6 shows the average of success rate of gaze following in utilizing *following-gaze* module in the cases of the different object size. We can see that the proposed mechanism can achieve high success rate in every case while the mechanisms without arrangement process and insertion one has low success rate except for the case where the number of segments is sufficient to achieve gaze following.

We also checked how many segments S_1 or M_1 is quantized into after learning. The results shown in

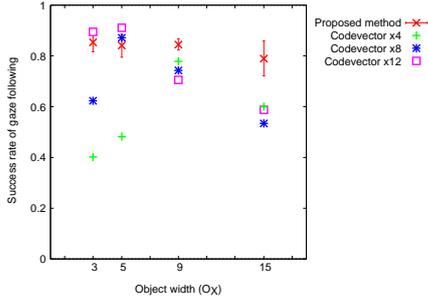


Figure 6: Change in success rate of gaze following for different object size

Figure 7 indicates that the proposed mechanism can find reasonable segmentation to achieve gaze following.

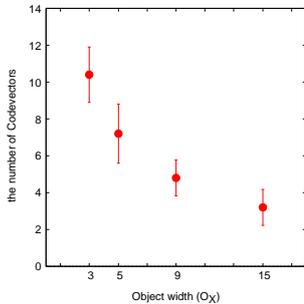


Figure 7: The number of segments in S_1 and M_1 for different object size after learning

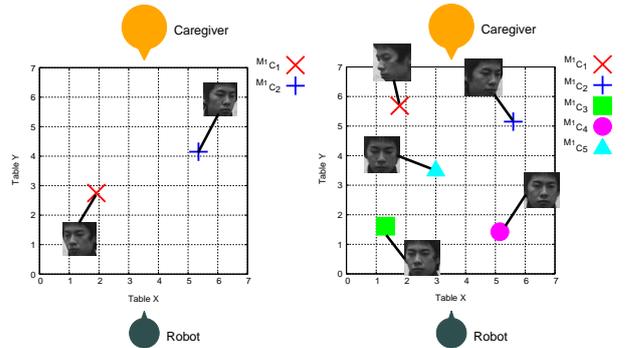
4.4 Segmentation in presence of high-dimensional information

In natural interaction with a caregiver, an infant should deal with high-dimensional information. In this case, a designer has difficulty quantizing variables in advance. We tested whether a robot acquires gaze following when it obtains a camera image of human face and takes actions on a square table.

We ran simulations where two objects each of which is a square having 4 sections were arranged on a square table having 49 sections. The robot observes one of 18 40×40 pixel grayscale images indicating different directions of the human face. As codebook vectors for S_1 , 1600-dimensional vectors were used. The robot's actions were represented as 2-dimensional vectors indicating a position on the table. We set $(\eta, \zeta, \xi) = (1.0, 500, 8.0 \times 10^4)$. A new codebook vector was added on a point through two vectors of the existing vectors in insertion process.

The average of 1.8 CMs was obtained. In over 80 percent of the simulations, *following-gaze* module and *returning(after seeing an object)* module were generated. In the simulations, S_1 and M_1 were quan-

tized into the average of 5.2 segments. Figure 8 shows changes in the sensorimotor map from S_1 to M_1 constructed by *following-gaze* module during a simulation. When *following-gaze* module was generated, the robot can coarsely shift its gaze to where the caregiver is looking (Figure 8(a)). Through the iteration of arrangement process and insertion one, however, it acquired sensorimotor map sufficient to follow the caregiver's gaze.



(a) A sensorimotor map in following-gaze module when the module was generated. (b) A sensorimotor map in following-gaze module after a simulation

Figure 8: Changes in sensorimotor map from S_1 to M_1

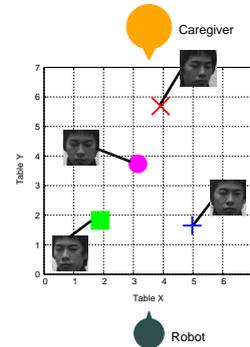


Figure 9: A sensorimotor map of following-gaze module in a robot with realistic visual field

4.5 Segmentation in presence of perspective correct visual field

In the experiments described above, we assumed that the robot can see an area of equal size despite the position on the table. However, this assumption is not reasonable in a real world. The size of the area that a human sees depends on the distance from the human: the closer the area is to the human, the smaller area the human sees. Therefore, we investigated whether codebook vectors are arranged according to the distance from a robot when the robot has a visual field depending on the distance.

We assumed to have a robot that is 0.5 meters tall has 2.5-degree field of view based on the area that a human fovea covers (Fairchild, 2005). This means that the robot can see about three or two sections on the table at the same time when shifting its gaze around the caregiver while it can see about one section when looking at an area around the robot itself. We used the same experimental setting as ones in the previous section except for the view field.

The average number of acquired CMs was 1.6. In about 60 percent of the simulations, *following-gaze* module and *returning(after seeing an object)* module were generated. The number of segments in S_1 and M_1 was the average of 3.8 segments. Figure 9 shows changes in the sensorimotor map from S_1 to M_1 constructed by *following-gaze* module after a simulation. S_1 and M_1 were quantized so that codebook vectors are arranged depending on the distance. This result indicates that the proposed mechanism enables the robot to find reasonable segmentation even when it has to segment a variable with depending on the distance.

5. Conclusion and discussion

We proposed a learning mechanism that finds reasonable segmentation to achieve joint attention behavior as well as incrementally acquires it by reproducing the contingency in interaction with a caregiver. The robot autonomously categorizes sensorimotor activity according to a contingency measure based on transfer entropy. We confirmed that a robot acquires gaze following and alternation and it finds suitable segmentation to reproduce the contingency in several conditions including several kinds of difficulty.

Developmental psychologists have suggested that human infants develop the ability of gaze following gradually (Moore and Dunham, 1995): they utilize only the information of head orientation of another person to achieve joint attention, and then begin to realize that the person’s eyes also direct his/her attention. In the experiment, the robot quantized sensory (and motor) variables to find stronger contingency and, as a result, gradually quantized S_1 representing the caregiver’s direction of gaze at higher resolution. Finding stronger contingency may explain infant development of gaze following.

In the proposed mechanism, the derivative of C-saliency ΔC was utilized to modulate the codebook vectors. We investigated how much ΔC influences this modulation. We ran another simulation using the same experimental setting as the ones reported in section 4.5 except that ΔC was replaced by a constant value, expressly $\Psi_{\Delta C} = 1.0$. Compared to the result with adaptive $\Psi_{\Delta C}$ (Figure 9), codebook vectors of M_1 were distributed evenly on a table although the segmentation in M_1 with adaptive $\Psi_{\Delta C}$

is optimized reflecting visual field. This illustrates that the modulation based on the derivative of C-saliency promotes finding segmentation sufficient to reproduce contingency.

In the experiments, a few components in a variable such as f_r in S_1 were given in advance. However, a robot should quantize each variable without such a priori knowledge. The segmentation to reproduce contingency of interaction with others may generate the components given in the experiments. As a future work, we will investigate whether a robot can autonomously find suitable segmentation including f_r and f_ϕ in S_1 .

Acknowledgements

This work was supported by Grant-in-Aid for JSPS Fellows(20-5227).

References

- Asada, M., Hosoda, K., Kuniyoshi, Y., Ishiguro, H., Inui, T., Yoshikawa, Y., Ogino, M., and Yoshida, C. (2009). Cognitive developmental robotics: a survey. *IEEE Trans. on Autonomous Mental Development*, 1(1):12–34.
- Fairchild, M. (2005). *Color appearance models*. Wiley.
- Imai, M., Ono, T., and Ishiguro, H. (2001). Physical relation and expression: Joint attention for human-robot interaction. In *Proc. of 10th IEEE International Workshop on Robot and Human Communication*.
- Kaplan, F. and Hafner, V. (2004). The challenges of joint attention. *Interaction Studies*, 5:67–74.
- Moore, C. and Dunham, P., (Eds.) (1995). *Joint attention: It’s origins and role in development*. Lawrence Erlbaum Associates.
- Nagai, Y., Hosoda, K., Morita, A., and Asada, M. (2003). A constructive model for the development of joint attention. *Connection Science*, 15(4):211–229.
- Schreiber, T. (2000). Measuring information transfer. *Physical review letters*, 85(2):461–464.
- Sumioka, H., Yoshikawa, Y., and Asada, M. (2008). Development of joint attention related actions based on reproducing interaction contingency. In *Proc. of the 7th Int. Conf. on Developmental and Learning*.
- Triesch, J., Teuscher, C., Deák, G., and Carlson, E. (2006). Gaze following: why (not) learn it? *Developmental Science*, 9(2):125–157.