# Realization and Evaluation of Realistic Nod with Receptionist Robot SAYA

Takuya HASHIMOTO[1], Student Member, IEEE, Sachio HIRAMATSU[1], Toshiaki TSUJI[1], Member, IEEE, and Hiroshi KOBAYASHI[1], Member, IEEE

[1] Department of Mechanical Engineering, Tokyo University of Science, Tokyo, Japan,
e-mail : (tak, hiramatsu, tsuji, hiroshi)@kobalab.com

*Abstract*— Android robots that have a human-like appearance have been developed recently. The purpose of android robots is to realize realistic communication between human and robot by implementation of human-like behaviors. Then knowledge and techniques for generating human-like natural motions and behaviors with android robot are required. In face-to-face communication, we can observe that the speaker nods during his utterance, and we call this behavior "Speaker's nod". Thus, in this study, we tried to realize realistic "Speaker's nod" as a natural interpersonal behavior with receptionist robot SAYA that has human-like appearance. Therefore we experimentally investigated expression-timing, angle and angular velocity of the human "Speaker's nod", and we implemented "Speaker's nod" to SAYA based on the investigation results. Then we have verified that "Speaker's nod" at appropriate timing affected the humanity of SAYA, and it greatly enhanced emotional impression of human by subjective experiment.

## I. INTRODUCTION

These days, research and development of multiple interaction-oriented robots have been frequently observed. These robots are expected to interact with human in daily-lives or office-spaces [1][2]. Since it is said that non-verbal medium are important for realization of smooth communication, most of these robots try to share information with human through not only verbal communication but also non-verbal communication [3]-[5] such as a glance, a nod and gestures.

However, these robots are immediately recognized as robots because of mechanical appearance. On the other hand a few robots that have human-like appearance are developed [6][7], and they are called android-robot. Kobayashi et al. [6] tried to realize natural facial expressions because facial expressions are said to be playing the most important role in face-to-face communication. Ishiguro et al. [7] developed whole body type android robots by ordering to a company. They tried to generate natural behaviors and motions of the robot. In addition, they try to evaluate humanity of android robots in cognitive science perspective. The greatest asset of android robots is that they give us a strong feeling of presence as if we communicate with real human. Therefore it seems that android robots bring human-robot communication close to human-human communication [8], while robots with mechanical appearance lack the ability to express human-like behaviors in particularly non-verbal communication.

We aim for realization of human-like natural behaviors with android receptionist robot SAYA shown in Fig.1. In this paper, a nod of speaking person especially draws attention as a natural behavior in face-to-face communication, and that is called "Speaker's nod". Although it is generally acknowledged that a nod is a listener's reaction to show agreement of a speaker, it can be also observed that a speaker often nods in daily communication. We experimentally investigate expression-timings, amplitude and velocity of such "Speaker's nod". Furthermore, human-like "Speaker's nod" is implemented to SAYA, and then we evaluate effect of nod-timing difference on the humanity of SAYA and also confirm the effect on the communication by subjective experiment.

Chapter 2 describes the mechanical structure and control system of SAYA. Chapter 3 shows the investigation of expression-timings, amplitude and velocity of human "Speaker's nod" in face-to-face conversation. In Chapter 4, "Speaker's nod" is experimentally implemented to SAYA as a natural behavior based on investigation results. Chapter 5 explains our evaluation experiment. Finally, we offer conclusions in Chapter 6.
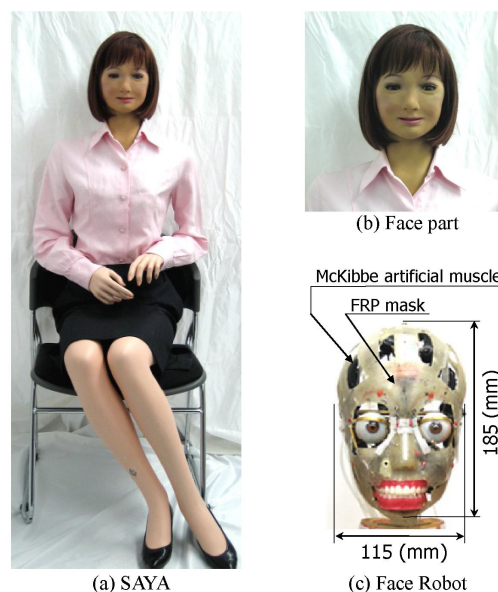


(b) Face part

McKibbe artificial muscle
FRP mask
185 (mm)
115 (mm)

(a) SAYA    (c) Face Robot
Fig.1 Receptionist Robot SAYA

## II. RECEPTIONIST ROBOT SAYA

### A. Overview of receptionist robot SAYA

We have been developing the receptionist robot SAYA that works in our university entrance. The task is to offer information about the university, i.e., (1) introduction of the university and department, (2) information of offices and laboratories (the extension number and the location), (3) introduction of research laboratories, and (4) introduction of SAYA.

Since the anthropomorphism face-robot [6] had been applied to mannequin body (Fig.1(a)), SAYA looks like a real human. In face-to-face communication, the facial expressions are said to play a most important role in non-verbal communication [9]. Thus, facial expressions seem one of the important factor for realizing smooth communication between human and robot. In order to generate human-like facial expressions, we referred to FACS (Facial Action Coding System) proposed by P. Ekman et al. [10]. 19 control points on the face are selected according to FACS. Combining movements of these control points, various facial expressions are realized with the face robot similar to human beings. Especially, a high correct recognition rate of 6 typical facial expressions that are possible to be recognized and express universally was achieved in previous research [11].
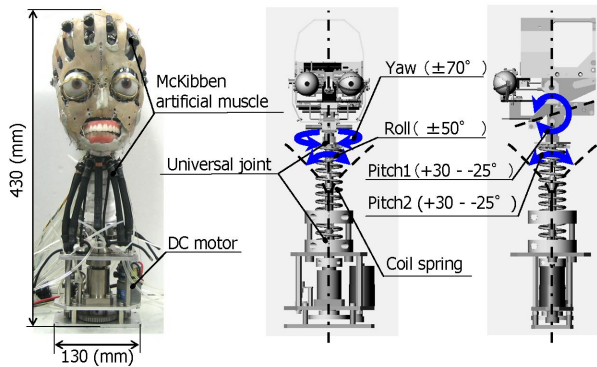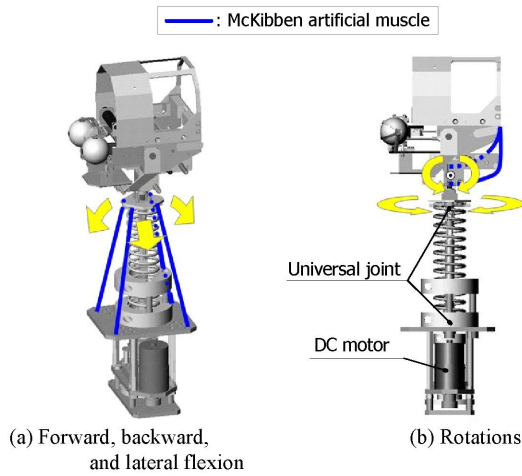


Fig.2 Internal structure of SAYA



Fig.3 Actuator distribution

### B. Hardware configuration

Fig.1(c) and Fig.2 show the structure of the face robot. We use McKibben artificial muscle [12][13] for actuating the control points. Since it is small, light and flexible, it can be distributed to curved surface of the skull such as human muscles. In addition we form the facial skin with soft urethane resin to realize the texture of human facial skin.

The face robot has an oculomotor mechanism. The mechanism controls both pitch and yaw rotations of eyeballs by 2 DC motors. The two eyeballs move together since they are linked to each other. Moreover we mounted a CCD camera inside of left-side eyeball. Then recognizing the human skin color region from image of the CCD camera, the sight line is controlled so that the face robot can pursue the visitor.

In order to realize flexible motion like a cervical spine, we adopted a coil spring for the head motion mechanism. Furthermore, the center of rotations for pitch rotation ("Pitch1") and yaw rotation were set in the base of head. Referring to anatomical knowledge, the movable positions and movable ranges were defined as shown in Fig.2. A forward and a backward motion are flexed by combination of a head rotation and a neck bending. Roll-rotation and both pitch-rotations ("Pitch1" and "Pitch2") are also actuated by McKibben artificial muscle. Yaw-rotation is controlled by a DC motor (Fig.3).

### C. Control system

Fig.4 shows the control system. We use an electro pneumatic regulator for controlling pressurized air from compressor according to commands from PC2. In PC1, skin color regions are extracted by binarization of input image. Then the extracted skin color regions are labeled, and largest area is detected as a face area. Moreover PC1 controls eyes and head motion so that face area always comes to the center of image. SAYA also communicates with visitor using speech recognition and reproduction of recorded female voice.
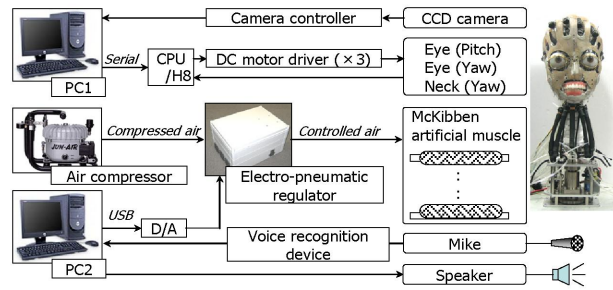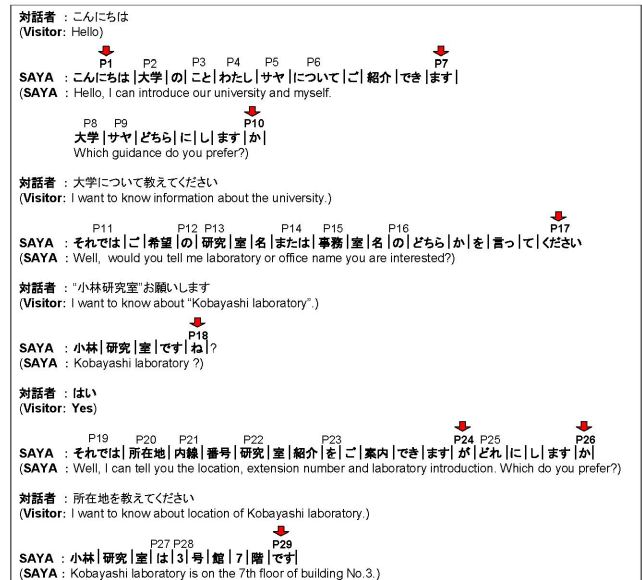


Fig.4 Control system of SAYA



Fig.5 Conversation example (in Japanese)

### D. Flow of conversation scenario

Fig.5 shows a conversation example between SAYA and a visitor. SAYA takes control of conversation and the visitor answers the SAYA's queries one by one. We use speech recognition system that is made in Toshiba IT & Control Systems Corporation. Thus, SAYA can recognize the keyword from visitor's utterance by word-spotting method. In addition we use recorded female voices as a SAYA's

voices, and SAYA's utterances are produced by combination of these voices.

## III. ANALYSIS OF "SPEAKER'S NOD"

Watanabe et al. [5] reported the research on the timing of a reaction and bodily motions. He showed the embodied interaction robot that moves their head and body based on on/off information of voice. However a context, linguistic and emotional information were not considered in this system. On the other hand, it is said that "Speaker's nod" works for not only emphasis of term and meaning but also controlling a rhythm and a flow of conversation [14]. Thus "Speaker's nod" is not expressed randomly. It is also thought that angle and angular velocity are the critical factors to realize human-like "Speaker's nod". We then investigated timings, angle and angular velocity of human's nod in conversation.

### A. Analysis procedure

We investigated expression-timings, angle and angular velocity of human "Speaker's nod" in face-to-face conversation experiment shown in Fig.6. In experiment, an experimenter and a subject sit facing each other, and we asked them to talk mutually. Their conversation was recorded, and 3-dimensional trajectories of subject head motions were measured by a motion capture system (Optotrak Certus, made in Northern Digital Inc.). Marker-positions of motion capture are shown in Fig.7. We asked 5 subjects to participate in this experiment, and experimental duration of each subject was about five minutes.
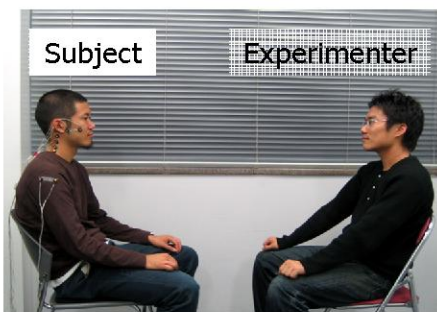


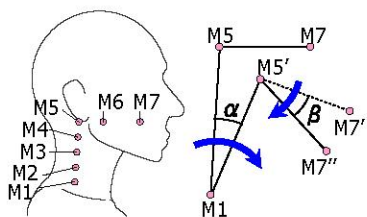Fig.6 Overview of face-to-face conversation



Fig.7 Marker positions of motion capture system

### B. Analysis result of expression-timing of "Speaker's nod" and discussions

The relationship between subject's utterance and expression-timings of "Speaker's nod" were investigated from recorded videos. Then we delimited subject's utterance on the boundary of 200 ms silent pause, and we classified "Speaker's nod" into four categories: "listener's nod (that is, reaction for giving responses to make the conversation go smoothly)", "during of utterance", "end of utterance" and "others". Fig.8 shows analysis result.

It is generally thought that there are a lot of listener's nods as a reaction in conversation. However this result shows that there are a lot of speaker's nods too, and especially the speaker nods at the end of utterance frequently. Maynard [14] assumed that "Speaker's nod" works for emphasis of a term and a meaning, and especially the "Speaker's nod" at the end of utterance clarifies the alternation of the utterance-initiative. Thus, "Speaker's nod" plays an important role in the alternation of utterance.

Then, in this study, we pay attention to the end of utterance as effective nod-timing in face-to-face communication. Since it is easily predicted that the end of sentence is the end of utterance, it is thought that "Speaker's nod" should be expressed at the end of sentence. That is, in the case of SAYA's utterance shown in Fig.5, it is predicted that "Speaker's nod" are expressed at P7, P10, P17, P18, P26 and P29.
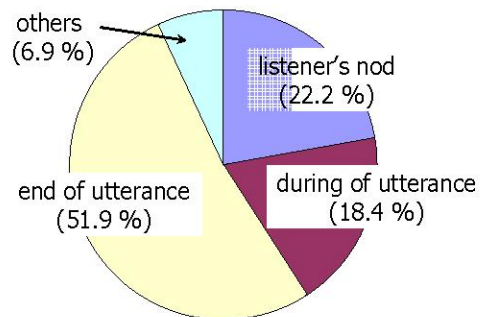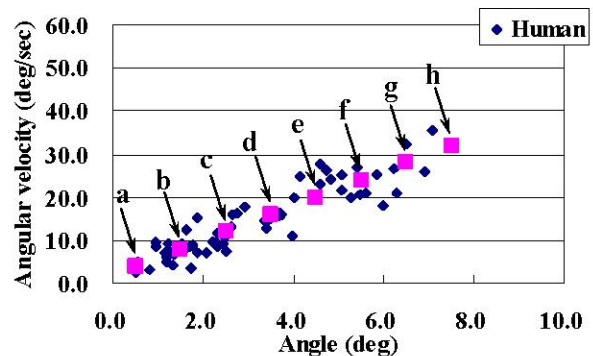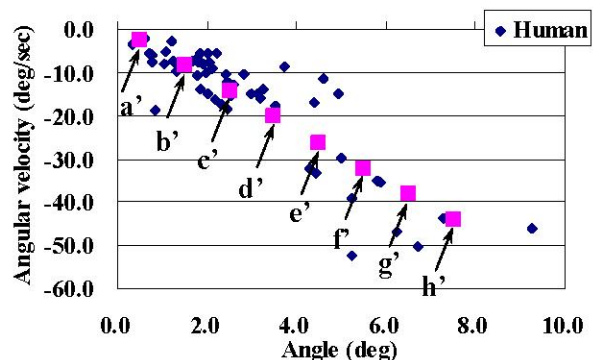


Fig.8 Classification of "Speaker's nod"



(a) Downward angle



(b) Upward angle
Fig.9 Angle and angular velocity of "Speaker's nod"

### C. Analysis result of angle and angular velocity of "Speaker's nod" and discussions

We extracted the total of 58 speaker's nods by extracting 10 to 15 ones from each subject. Moreover we measured the

328

angle and the angular velocity of nod as shown in Fig.9. Since most of nods were done by the head rotation, we defined that the angle-β shown in Fig.7 was the nod-angle. Fig.9(a)(b) show the relations between angle and angular velocity of downward and upward motion of nod.

We found that both angle and angular velocity were changed widely, and angular velocity grew larger according to the growing of angle.

## IV. REALIZATION OF "SPEAKER'S NOD"

### A. Experimental approach of expression-timing

The previous chapter shows that most of "Speaker's nods" are expressed at the end of utterance. However there is also the end of utterance in the middle of sentence, and "Speaker's nod" might be expressed while uttering. Thus, it is difficult to determine the expression-timing definitely. We experimentally investigated where "Speaker's nod" actually appeared during the SAYA's utterance shown in Fig.5.

We put an experimenter and a subject as shown in Fig.10 assuming conversation of reception. The distance between the experimenter and subject was about 1.5 m that is called social distance [15]. Such distance is supposed to be appropriate as the situation of reception. We asked the subject to play the role of SAYA, and we asked the experimenter to play the role of visitor. Then we asked them to speak each sentence shown in Fig.5. The conversation between subject and experimenter was recorded by video. Since spontaneity is required for communication between the experimenter and the subject, we asked them to memorize and practice each dialogue before experiment. Subjects were 20 people and we asked them to attend experiment twice.
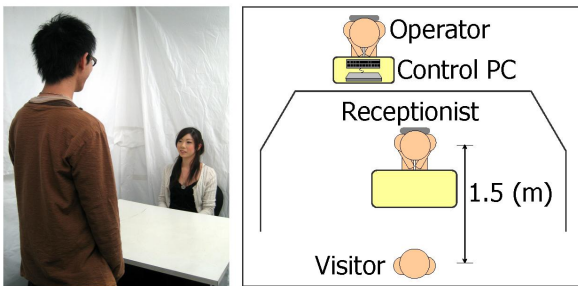


Fig.10 Experimental environment

### B. Experimental result of expression-timing and discussions

40 conversation samples were obtained after experiment mentioned above. SAYA's utterance was parsed as shown in Fig.5 to analyze the relationship between the nod-timing and the linguistic information. Positions at which subjects had nodded in conversation were numbered. We counted the number of subjects who had nodded at each numbered point, and Fig.11 shows the proportion of subjects at each numbered point to all subjects.

The average is 10.5. In order to decide which nod-timing we use, we count numbered points that are included at each 0.1 range. For example, in case of range of 0.6 < <= 0.7, P7 and P24 are extracted and number is 2. It was found that the section 0.5 < <= 0.6 was the threshold as the result of the discriminant analysis. Then eight points (P1, P7, P10, P17,

P18, P24, P26 and P29) were extracted as the nod-timing in SAYA's utterance shown in Fig.5.

P1 is the start of conversation, and that is the scene of greeting in daily life. P1 is the bowing rather than the nod. These words at P7, P10, P17, P18, P26 and P29 are used well at the end of sentence in Japanese, and these points were also the end of utterance. Although P24 is the middle of sentence, the pause of utterance appeared immediately after this point.

These results almost correspond to the estimate results in the preceding chapter. Furthermore since a lot of subjects nod at these points mutually, it is thought that the influence of the individual variation is a little at these points. Therefore we decided to implement "Speaker's nod" to SAYA at these eight points in her utterance.
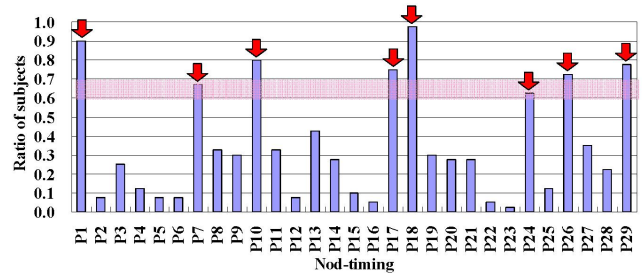


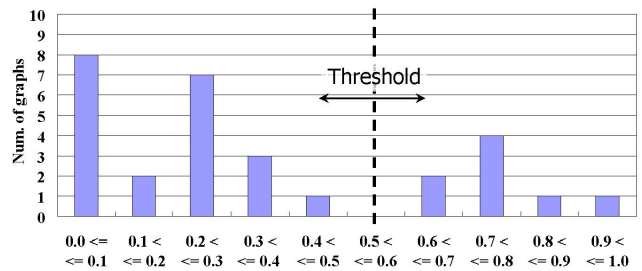Fig.11 Ratio of subjects who nodded at each expression-timing of nod



Fig.12 Result of discriminant analysis

### C. Experimental approach of nod-angle and angular velocity

Nod-angle and angular velocity vary widely as shown in Fig.9. However this range might be different in comparing human and robot. Therefore the range that seems natural behavior should be evaluated with SAYA.

In the first place, we divided the range of human's nod into seven (a-h, a'-h') as shown in Fig.9. We prepared nod-pattern A by combination of "a" and "a'". The remaining nod-patterns B (b-b') - H (h-h') were created as the same manner. Table 1 shows the angle and the angular velocity of each nod pattern, and it shows averages when SAYA had tried each pattern five times. Values in parenthesis show standard deviations. Each pattern was implemented as SAYA's "Speaker's nod" according to the eight nod-timings described in preceding section.

We had SAYA sit as "Receptionist", and asked the subject to stand as "Visitor" as shown in Fig.10. The distance between SAYA and the subject was 1.5 m. Then subject had to observe SAYA's behaviors while SAYA was speaking unilaterally. At this time, the movements of mouth and blink were uniformed through each pattern. After that, subject had to evaluate which nod-pattern seems "natural" and "real" behavior. 25 people participated in this experiment as a subject.

## D. Experimental result of nod-angle and angular velocity and discussions

The number of subjects who estimated that the behavior was "natural" and "real" in each pattern is shown in Fig.13. We found that pattern D got highest score, and the evaluation was high in the order of E, C, F, B, G and A. Since there is no degree of freedom in SAYA's trunk, the whole body can't move. Therefore if the nod-angle is small, SAYA hardly seems to move. On the other hand, if the nod-angle is large, only the head moves greatly without the trunk motion. These cause unnatural behaviors. As a result, intermediate C, D and E were preferred rather than others, and it is thought that natural "Speaker's nod" can be realized within C to E.

In analysis of human motion, the amplitude of nods didn't depend on nod-timing and context, and it varied widely. We then decided to implement pattern C, D and E disorderly as SAYA's "Speaker's nod" for experiment.

Table1 Angle and angular velocity of each nod-pattern

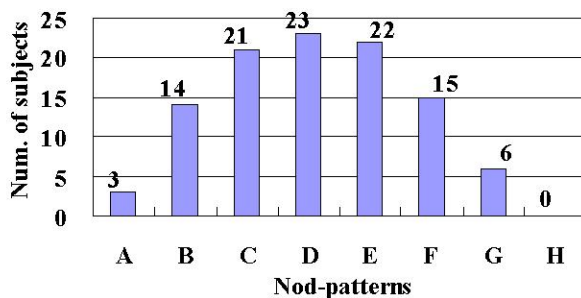|  |  | Angle [deg] | | Angular velocity [deg/sec] | |
|---|---|---|---|---|---|
|  |  | Target | SAYA | Target | SAYA |
| A | a | 0.50 | 0.55 (0.03) | 3.94 | 4.13 (0.24) |
|  | a' | -0.50 | -0.56 (0.03) | -2.26 | -2.80 (0.16) |
| B | b | 1.50 | 1.51 (0.04) | 7.93 | 8.13 (0.68) |
|  | b' | -1.50 | -1.47 (0.11) | -8.25 | -7.94 (0.88) |
| C | c | 2.50 | 2.52 (0.11) | 11.93 | 11.51 (0.70) |
|  | c' | -2.50 | -2.43 (0.17) | -14.25 | -13.51 (0.61) |
| D | d | 3.50 | 3.50 (0.11) | 15.93 | 15.50 (0.68) |
|  | d' | -3.50 | -3.44 (0.05) | -20.24 | -20.59 (0.31) |
| E | e | 4.50 | 4.50 (0.08) | 19.92 | 19.32 (0.35) |
|  | e' | -4.50 | -4.44 (0.14) | -26.23 | -26.58 (0.85) |
| F | f | 5.50 | 5.53 (0.04) | 23.92 | 23.72 (0.16) |
|  | f' | -5.50 | -5.41 (0.15) | -32.23 | -32.40 (0.92) |
| G | g | 6.50 | 6.58 (0.05) | 27.92 | 28.23 (0.21) |
|  | g' | -6.50 | -6.49 (0.07) | -38.22 | -38.86 (0.41) |
| H | h | 7.50 | 7.48 (0.01) | 31.91 | 32.09 (0.05) |
|  | h' | -7.50 | -7.47 (0.10) | -44.22 | -44.72 (0.62) |



Fig.13 Evaluation result about humanity and reality of each nod-pattern

## V. EXPERIMENTAL EVALUATION

We conducted experiment to verify the effects of the implemented nod-pattern for the humanity of SAYA in human-robot communication. The hypothesis for the experiment was "if SAYA nods at appropriate timing, her humanity and reality are improved, and the human can smoothly communicate with her".

### A. Experimental method

In this experiment, we evaluated effects of difference of the nod-timings. The following describes the detailed procedure.

[Experimental conditions]

In order to investigate the effect of the difference of nod-timings, following three nod-timing conditions were prepared.
- T1 : SAYA doesn't nod
- T2 : SAYA nods eight times randomly. The nod frequency in one utterance is same as condition T3.
- T3 : SAYA nods at the eight-points detected in previous chapter.

Nod-pattern C, D and E shown in previous experiment are implemented disorderly in condition T2, T3. The movements of mouth and blink were uniformed through each pattern.

[Experimental environment]

The experimental environment are shown in Fig.10. SAYA sat at "Receptionist", and the subject stood at "Visitor". We adjusted the distance between SAYA and the subject to almost 1.5 m. We also put microphone for speech recognition.

[Experimental procedure]

At the first, we gave a task (that is, "Please say hello to SAYA, and ask her about either the address or the extension number of a laboratory") to the subject. Each subject participated in these three experiments. The order of the three experiments was counterbalanced (that is, the experiment was conducted in the order of either T1-T2-T3 or T3-T2-T1).

[Evaluation method]

In order to obtain subjective evaluation, we managed a questionnaire which is used for investigating how the difference of nod-timings affects communications in terms of conveying the information, smoothness, familiarity and humanity. We used following indicators for subjective evaluation. The subjects answered each item by range of -3 to 3.
- Aspects of conveying information, smoothness
  - Smoothness of conversation
  - Politeness
  - Meaningful
- Aspects of familiarity
  - Friendly
  - Familiarity
  - Amenity
- Aspects of humanity
  - Humanity
  - Nature

[Subjects]

50 university students (43 males, 7 females) participated in the experiments as subjects.

### B. Experimental result and discussions

Fig.4 shows the result of questionare. It shows the mean, the standard deviation and the result of analysis of variance (ANOVA) of the questionnaire.

ANOVA shows that there are significant difference in all questionnaire items. Thus, we applied the multiple comparison to each item by LSD method. As a result, in Q.1 ($p<.05$), Q.4 ($p<.05$), Q.5 ($p<.05$), Q.7 ($p<.05$) and Q.8 ($p<.05$), the value of T3 is significantly better than that of T2, and the T2 is significantly better than T1. The result also
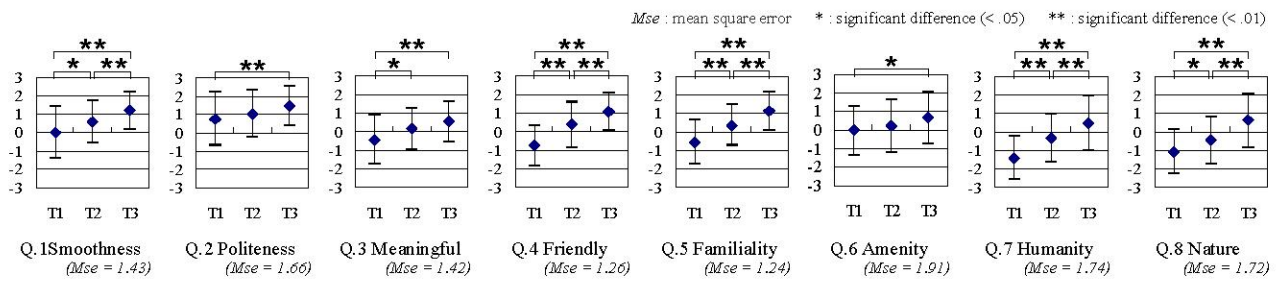
*Mse* : mean square error    \* : significant difference (< .05)    \*\* : significant difference (< .01)

Fig.14 Result of comparison among 3 nod-patterns (T1, T2, T3)

shows that the value of T3 is better than that of T1 in Q.2 (p<.05), and the value of T3 is better than that of T1 in Q.6 (p<.05).

Since it is said that "Speaker's nod"makes the alternation of the utterance smoothly [14], it seems that the evaluation of Q.1 ("Smoothness of conversation") was improved in T2, T3 condition. Especially, participants tended to prefer T3 rather than T2 condition. Furthermore, Q.7 ("Humanity") and Q.8 ("Nature") were improved by nodding at appropriate timing, and Q.4 ("Friendly") and Q.5 ("Familiarity") were also influenced. Therefore, "Speaker's nod" at appropriate timing contributed to aspects of conveying information, smoothness, familiarity and humanity in face-to-face communication. These results confirm the importance of "Speaker's nod" in communication, and these results highlight the positive effect of nodding according to appropriate timing.

In this study, the nod-timings were defined experimentally as shown in chapter IV. Although extracting nod-timings automatically including linguistic and emotional information remains as one of our future challenge, it will be solved by enhancing the research on the reaction-timing [5][15].

## VI. CONCLUSION

In this report, we aimed to examine guideline to generate realistic "Speaker's nod" with the receptionist robot SAYA that has human-like appearance. In order to realize realistic "Speaker's nod" to SAYA, we investigated human "Speaker's nod" paying attention to expression-timing, angle and angular velocity of nod. Then we implemented "Speaker's nod" to SAYA based on investigation and experimental results.

We investigated the effects of implemented "Speaker's nod" on the SAYA's humanity and on the communication by subjective experiment. The experiment results indicated that appropriate "Speaker's nod" greatly enhanced the humanity of SAYA, and it improved the aspects of conveying information, smoothness and familiarity in conversation. We believe that our findings could lead to android robot that has feeling of human-like presence and perform human-like behaviors in communication.

Our future work is to investigate the natural facial expression changes, head motion and involuntary movement of body during communication for realizing human-like behaviors.

## REFERENCES

[1] John. F, H. Asoh and T. Matsui, "Natural dialogue with the Jijo-2 office robot", *Proceeding of the IROS'98*, 1998, pp.1278-1283.

[2] T. Kanda, H. Ishiguro, T. Ono, M. Imai and R. Nakatsu, "Development and evaluation of an interactive humanoid robot "Robovie"", Proceeding of the ICRA'02, 2002, pp.1848-1855.

[3] T. Tojo, Y. Matsusaka, T. Ishi and T. Kobayashi, "A conversational robot utilizing facial and body expressions", *Systems, Man, and Cybernetics, 2000 IEEE International Conference on*, 2000, pp.858-863.

[4] M. Shiomi, T. Kanda, M. Imai, T. Ono, D. Sakamoto, H. Ishiguro and Y. Anzai, "Embodied cooperative behaviors by an autonomous humanoid robot", *Proceeding of the IROS'04*, 2004, pp.2506-2513.

[5] T. Watanabe, M. Okubo and H. Ogawa, "A speech driven embodied interaction robots system for human communication support", *Systems, Man, and Cybernetics, 2000 IEEE International Conference on*, 2000, pp.852 -857.

[6] H. Kobayashi and F. Hara, "Study on Face Robot for Active Human Interface", *Journal of the Robotics Society of Japan*, Vol.12, No.1 (1994), pp.155-163. (in Japanese)

[7] H. Ishiguro, "Android Science -Toward a new cross-interdisciplinary framework-", *Proc. of International Symposium of Robotics Research*, (2005).

[8] D. Matsui, T. Minato, K. F. MacDorman and H. Ishiguro, "Generating natural motion in an android by mapping human motion", *Procedings of IROS'05*, 2005, pp.3301-3308.

[9] J. Cole, "About Face", *MIT Press*, 1998.

[10] P. Ekman and W.V. Friesen, "The Facial Action Coding System", *Consulting Psychologists Press*, 1978.

[11] T. Hashimoto, S. Hiramatsu, T. Tsuji and H. Kobayashi, "Development of the Face Robot SAYA for Rich Facial Expressions", *SICE-ICASE International Joint Conference 2006*, 2006, pp.5423-5428.

[12] C.P. Chou and B. Hannaford, "Measurement and Modeling of McKibben Pneumatic Artificial Muscle", *IEEE Transactions on Robotics and Automation*, vol.12, 1996, pp.90-102.

[13] H. F. Schulte, "The characteristics of the McKibben artificial muscle, In The Application of External Power in Prosthetics and Orthotics", *National Academy of Sciences-National Research Council, Publication874*, 1961, pp.94-115.

[14] S. K. Maynard, "Anlysis of conversation", *Kuroshio Publishers*, 1993, pp.23-179. (in Japanese)

[15] T. Matsuo, "Psychology of communication", *Nakanisiya Publishers*, 1999, pp.57-68. (in Japanese)