

Realizing Personality in Audio-Visually Triggered Non-verbal Behaviors

Hiroshi G. Okuno^{†,*}, Kazuhiro Nakadai^{*} Hiroaki Kitano^{*,‡}

[†] Graduate School of Informatics, Kyoto University, Kyoto, Japan

^{*} Kitano Symbiotic Systems Project, ERATO, Japan Science and Tech. Corp., Tokyo, Japan

[‡] Sony Computer Science Laboratories, Inc., Tokyo, Japan

okuno@i.kyoto-u.ac.jp, nakadai@nakadai.com, kitano@csl.sony.co.jp

Abstract—Controlling robot behaviors becomes more important recently as active perception for robot, in particular active audition in addition to active vision, has made remarkable progress. We are studying how to create social humanoids that perform actions empowered by real-time audio-visual tracking of multiple talkers. In this paper, we present *personality* as a means of controlling non-verbal behaviors. It consists of two dimensions, dominance vs. submissiveness and friendliness vs. hostility, based on the Interpersonal Theory in psychology. The upper-torso humanoid *SIG* equipped with real-time audio-visual multiple-talker tracking system is used as a testbed for social interaction. As a companion robot, with friendly personality, it turns toward a new sound source in order to show its attention, while with hostile personality, it turns away from a new sound source. As a receptionist robot with dominant personality, it focuses its attention on the current customer, while with submissive personality, its attention to the current customer is interrupted by a new one.

Keywords—robot interaction, active audition, personality, focus-of-attention, social interaction

I. INTRODUCTION

Social interaction is essential for humanoid robots, because they are getting more common in social and home environments, such as a pet robot in a living room, a service robot at office, or a robot serving people at a party [1]. Social skills of such robots require robust complex perceptual abilities; for example, it identifies people in the room, pays attention to their voice and looks at them to identify, and associates voice and visual images. Intelligent behavior of social interaction should emerge from rich channels of input sensors; vision, audition, tactile, and others.

Perception of various kinds of sensory inputs should be *active* because we hear and see things and events that are important to us as individuals, not sound waves or light rays [2]. Selective attention of sensors such as looking against seeing or listening against hearing plays an important role in social interaction. Other important factors in social interaction are recognition and synthesis of emotion in facial expression and verbal tones [3], [4].

Selectivity and *capacity limitation* are two main factors in attention control [5]. A humanoid does some perception intentionally based on selectivity [6]. It also has some limitation in the number of sensors or processing capabilities, and thus only a limited number of sensory information is processed. Since selectivity and capacity limitation are the

flip side of the same coin, only selectivity is argued in this paper. Selective attention of auditory processing called the *cocktail party effect* was reported by Cherry in 1953 [7]. At a crowded party, one can attend to one conversation and then change to another. But the question is to which one pays one's attention and how one changes one's attention.

Personality in selective attention consists in answers of these questions. Reeves and Nass use the *Five-Factor Model* in analyzing the personality of media including software agents [8]. The *big five* dimensions of personality are *Dominance/Submissiveness*, *Friendliness*, *Conscientiousness*, *Emotional Stability*, and *Openness*. Although these five dimensions generally define a human's basic personality, they are not appropriate to define humanoid's one, because the latter three dimensions cannot be applied to current capabilities of current humanoids.

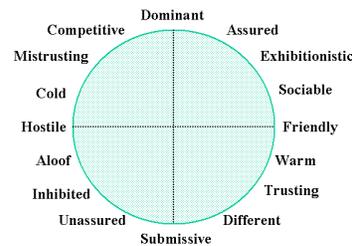


Fig. 1. Interpersonal Circumplex: variation of personality

We use the *Interpersonal Theory* instead for defining personality in selective attention. It deals with people's characteristic interaction patterns, as is shown in Figure 1, varying along the *Dominance/Submissiveness* and *Friendliness/Hostility*. The variation is represented by the *interpersonal circumplex*, which is a circular model of the interpersonal domain of personality [9].

Physically embodied agents, or humanoid robots have no explicit personality as far as we know. Usually personality is emphasized in language generation, whether verbal or textual. Although the most important human communication means is language, non-verbal sensori-motor

based behavior is non-the-less important. In this paper, we use personality to define attention control and report some observations of non-verbal interactions between humanoid and human.

A. Related Work

Personality for software agents has been studied extensively. Bates and his group propose *believable agents* that can express emotion clearly in appropriately timed manner [10]. Cassell developed conversational agents that integrate face and gesture [11]. She also argues that implementation of conversational agents should be based on actual study of human-human interaction. Hayes-Roth organizes the Virtual Theater project, which studies the creation of intelligent, automated characters that can act either in well-defined stories or in improvisational environments [12].

Personality for robots has also been investigated to widen communication channels in human-robot interaction, although not all works mention personality explicitly. Miwa *et al* have developed human-like head robots and implement personality to attain smooth and effective communication with human [13]. In their system, personality consists of the sensing and expression personality. The sensing personality determines how a stimulus works for a robot's mental state. Seven emotions were mapped out in the 3D mental space based on the *Five Factor Model*. Once the robot determines its emotion, it expresses its emotion based on the expression personality. They realized six kinds of personality on their robot.

Not a few works mention focus of attention. Ono *et al*. use the robot called *Robovie* to make common attention between human and robot by using gestures [14]. Breazeal incorporates the capabilities of recognition and synthesis of emotion in facial expression and verbal tones into the robot called *Kismet* [3], [4]. Waldherr *et al*. have developed the robot called *AMELLA* that can recognize pose and motion gestures [15]. Matsusaka *et al*. have built the robot called *Hadaly* that can localize the talker as well as recognize speeches by speech-recognition system so that it can interact with multiple people [16]. Nakadai *et al* developed *real-time* auditory and visual multiple-tracking system for the upper-torso humanoid called *SIG* [17], [18], [19]. They extended the system to attain in-face interaction by incorporating *auditory fovea* that is the azimuth dependency in performance of sound source localization [20].

Usually personality is emphasized in language generation, whether verbal or not. Although the most important human communication means is language, non-verbal sensori-motor based behavior is non-the-less important. In this paper, we use personality to define focus-of-attention control and report some observations of non-verbal interactions between humanoid and human.

II. HUMANOID HARDWARE

As a testbed of integration of perceptual information to control motor of high degree of freedom (DOF), we used the humanoid robot (hereafter, referred as *SIG*) with the following components:

4 DOFs of body driven by 4 DC motors — Each DC motor has a potentiometer to measure the direction.

A pair of CCD cameras of Sony EVI-G20 for stereo vision input.

Two pairs of omni-directional microphones (Sony ECM-77S). One pair of microphones are installed at the ear position of the head to collect sounds from the external world. Each microphone is shielded by the cover to prevent from capturing internal noises. The other pair are to collect sounds within the cover.

A cover of the body (Figure 2) reduces sounds to be emitted to external environments, which is expected to reduce the complexity of sound processing. This cover, made of FRP, is designed by Mr. Tatsuya Matsui for making human robot interaction smoother [21].

III. PERCEPTUAL SYSTEMS IN REAL-TIME MULTIPLE-TALKER TRACKING

The real-time multiple-talker tracking system is designed based on the client/server model (Figure 3). Each server or client executes the following logical modules:

1. Audition module extracts auditory events by pitch extraction, sound source separation and localization, and sends those events to Association module.
2. Vision module with a pair of cameras, extracts visual events by face detection, identification and localization, and then sends visual events to Association module.
3. Motor module generates PWM (Pulse Width Modulation) signals to DC motors and sends motor events to Association.
4. Association module groups various events into a stream and maintains association and deassociation between streams.
5. Attention module selects some stream on which it focuses its attention and makes a plan of motor control.
6. Dialog module communicates with people according to its attention by speech synthesis and speech recognition. The “Julian” automatic speech recognition [22] is used.

The status of each module is displayed on each node. *SIG* server displays the radar chart of objects and the stream chart. Motion module displays the radar chart of the body direction. Audition module displays the spectrogram of input sound and pitch (frequency) vs sound source direction chart. Vision module displays the image of the camera and the status of face identification and tracking.

To attain real-time tracking, the above modules are physically distributed to five Linux nodes connected by TCP/IP

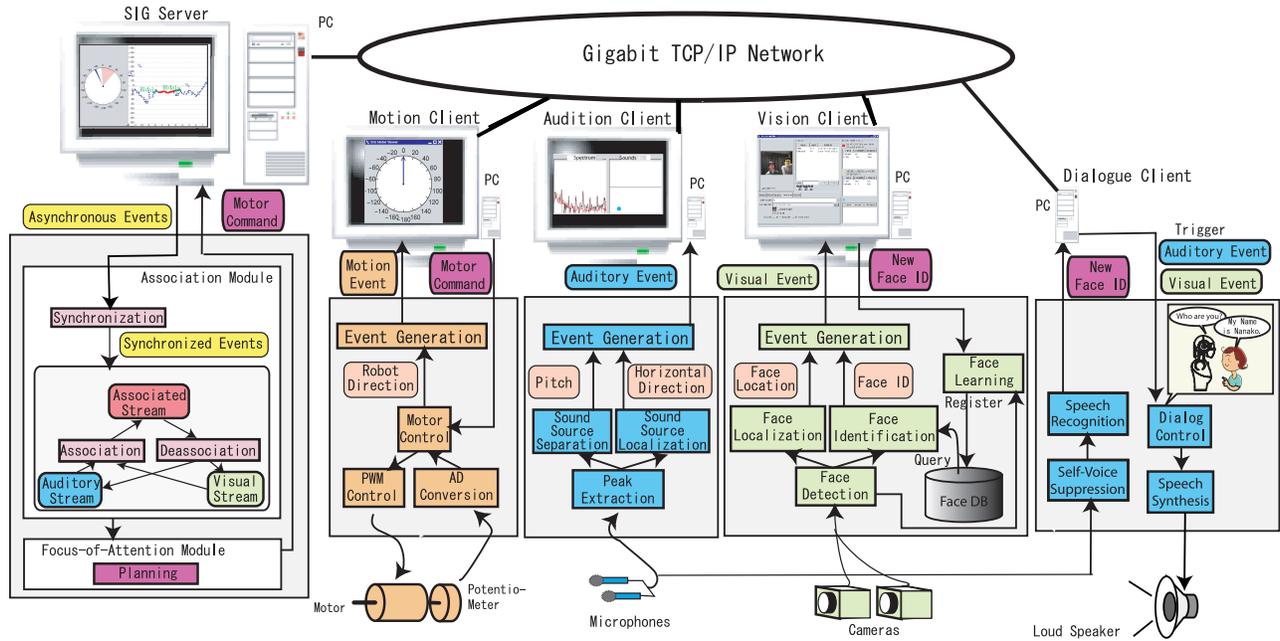


Fig. 3. Hierarchical architecture of real-time audio and visual tracking system

over Gigabit Ethernet network and run asynchronously. The system is implemented by distributed processing of five nodes with Pentium-IV 1.8 GHz. Each node serves Vision, Audition, Motion and Dialogue modules, and SIG server.

A. Active audition module

To localize sound sources with two microphones, first a set of peaks are extracted for left and right channels, respectively. Then, the same or similar peaks of left and right channels are identified as a pair and each pair is used to calculate interaural phase difference (IPD) and interaural intensity difference (IID). IPD is calculated from frequencies of less than 1500 Hz because of uniqueness of the phase.

Since auditory and visual tracking involves motor movements, which cause motor and mechanical noises, audition should suppress or at least reduce such noises. In human robot interaction, when a robot is talking, it should suppress its own speeches. Nakadai *et al* presented the *active audition* for humanoids to improve sound source tracking by integrating audition, vision, and motor controls [23]. We also use their heuristics to reduce internal burst noises caused by motor movements.

From IPD and IID, the epipolar geometry is used to obtain the direction of sound source [23]. The key ideas of their real-time active audition system are twofold; one is to exploit the property of the harmonic structure (fundamental frequency, F_0 , and its overtones) to find a more accurate pair of peaks in left and right channels. The other is to search the sound source direction by combining the belief factors of IPD and IID based on Dempster-Shafer theory.

Finally, audition module sends an auditory event consisting of pitch (F_0) and a list of 20-best direction (θ) with reliability for each harmonics.

B. Face recognition and identification module

Vision extracts lengthwise objects such as persons from a disparity map to localize them by using a pair of cameras. First a disparity map is generated by an intensity based area-correlation technique. This is processed in real-time on a PC by a recursive correlation technique and optimization peculiar to Intel architecture [24]. In addition, left and right images are calibrated by affine transformation in advance.

An object is extracted from a 2-D disparity map by assuming that a human body is lengthwise. A 2-D disparity map is defined by

$$DM_{2D} = \{D(i, j) | i = 1, 2, \dots, W, j = 1, 2, \dots, H\} \quad (1)$$

where W and H are width and height, respectively and D is a disparity value.

To extract lengthwise objects, the median of DM_{2D} along the direction of height is extracted. A 1-D disparity map DM_{1D} as a sequence of $D_l(i)$ is created.

$$DM_{1D} = \{D_l(i) | i = 1, 2, \dots, W\} \quad (2)$$

Next, a lengthwise object such as a human body is extracted by segmentation of a region with similar disparity in DM_{1D} . Then, for object localization, epipolar geometry is applied to the center of gravity of the extracted region. Finally, vision module sends a visual event consisting of a list

of 5-best Face ID (Name) with its reliability and position (distance r , azimuth θ and elevation ϕ) for each face.

C. Stream formation and association

Association synchronizes events given by other modules. It forms an auditory, visual or associated stream by their proximity. Events are stored in the short-term memory only for 2 seconds. Synchronization process runs with the delay of 200msec, which is forced by the largest delay of the vision module.

An auditory event is connected to the nearest auditory stream within $\pm 10^\circ$ and with common or harmonic pitch. A visual event is connected to the nearest visual stream within 40cm and with common face ID. In either case, if there are plural candidates, the most reliable one is selected. Unless any appropriate stream is found, such an event becomes a new stream. In case that no event is connected to an existing stream, such a stream remains alive for up to 500 msec. After 500 msec of keep-alive state, the stream terminates.

An auditory and a visual streams are associated if their direction difference is within $\pm 10^\circ$ and this situation continues for more than 50% of the 1 sec period. If either auditory or visual event has not been found for more than 3 sec, such an associated stream is deassociated and only existing auditory or visual stream remains. If the auditory and visual direction have been apart more than 30° for 3 sec, such an associated stream is deassociated to two separate streams.

IV. ATTENTION SYSTEM WITH PERSONALITY

Attention control focuses on one of auditory, visual, or associated streams. This selective attention is basically performed at two levels, that is personality and task. To define personality, the interpersonal circumplex of the Interpersonal Theory is used. With its two mutually independent axes, dominant/submissive and friendly/hostile, variations of personality are *Dominant*, *Assured*, *Exhibitionistic*, *Sociable*, *Friendly*, *Warm*, *Trustaing*, *Different*, *Submissive*, *Unassured*, *Inhibited*, *Aloof*, *Hostile*, *Cold*, *Mistrusting*, and *Competitive* (Figure 1) [9].

Since these variations are represented as a circle (circumplex), each variation of personality is represented as a point, (r, θ) , inside the interpersonal circumplex, where $0 \leq r \leq 1$ and $0 \leq \theta \leq 2\pi$. Therefore, the value of *Friendly/Hostile* axis and that of *Dominant/Submissive* axis are represented as $r \cos \theta$ and $r \sin \theta$, respectively. Each variation occupies a pie of $\pi/8$. For example, *Friendly* is specified as a pie section of $-\frac{\pi}{16} \sim \frac{\pi}{16}$, and *Dominant* as that of $\frac{3\pi}{16} \sim \frac{5\pi}{16}$.

To what the system attend is called “*interested*”. The total amount of interest in the system keeps the same and a newly focused stream takes all the amount of *interest* in winner-take-all competition between streams. attention con-

trol module selects the stream of the largest *interest*. Three mental factors are defined.

1. *interest* in a new stream — When a new stream is generated, the stream gets *interest* according to its status, auditory, visual or associated. The initial value of *interest* for a new stream is given at the time of stream generation.
2. decay of *interest* — The *interest* of a focused stream is reduced at the rate of e^{-kT} every minute, where k is $\{1.5 - r \sin \theta\}/3$. The lost *interest* is distributed to other streams.
3. decay of belief — Disappeared stream still remains in the system, because a unseen talker resumes to talk after a short period of silence. If disappeared stream is deleted immediately, the continuity of stream is difficult to maintain. In this paper, the constant value is used for the decay factor of belief.

The initial value of interest for a new stream is determined by what kind of interaction the robot will attend. For task-oriented manner, an associated stream has the highest initial value, while for socially-oriented manner, any new stream has the equal opportunity.

Task-oriented attention control forces Attention to behave according to a specific script. For example, a receptionist robot should focus on the user for whom an associated stream is generated. Therefore, the initial values of interest for auditory, visual and associated stream are 1, 1 and 2, respectively, in this paper. The essence of assignment is that the value for associated streams is highest.

Socially-oriented attention control forces Attention to show the interest of the robot. As an example of socially-oriented control, we implement a companion robot. It should pay attention to a new auditory or visual event, and thus all initial values of interest for any kind stream is the same, say 1 in this paper.

V. EXPERIMENTS AND OBSERVATION

Experiments was conducted in a small room of a normal residential apartment. The width, length and height of the room is about 3 m, 3 m, and 2 m, respectively. The room has 6 down-lights embedded on the ceiling. Two kinds of experiments are conducted in this section.

A. Task-oriented interaction: as a receptionist robot

One scenario to evaluate the above control is specified as follows: (1) A known participant comes to the receptionist robot. His face has been registered in the face database. (2) He says Hello to *SIG*. (3) *SIG* replies “Hello. You are XXX-san, aren’t you?” (4) He says “yes”. (5) *SIG* says “XXX-san, Welcome to the party. Please enter the room.”.

Figure 4 depicts two snapshots of this script. Figure 4 a) shows the initial state. The loud speaker on the stand is the mouth of *SIG*’s. When a participant comes to the receptionist, *SIG* has not noticed him yet, because he is out of

SIG's sight. When he speaks to *SIG*, Audition generates an auditory event with sound source direction, and sends it to Association, which creates an auditory stream. This stream triggers Attention to make a plan that *SIG* should turn to him, and *SIG* does it (Figure 4 b)).

This experiment demonstrates *SIG*'s two interesting behaviors. One is voice-triggered tracking, and the other is that *SIG* does not pay attention to its own speech. As a receptionist robot, once an association is established, *SIG* keeps its face fixed to the direction of the talker of the associated stream. Therefore, even when *SIG* utters via a loud speaker on the left, *SIG* does not pay an attention to the sound source, that is, its own speech.

Another script is that a hostile *SIG* with $\alpha = 1$ $\theta =$ turns away from an associated stream. In Figure 5, when a participant says "Hello" to *SIG*, *SIG* turns away from him.

B. Socially-oriented interaction: as a companion robot

When four talkers actually talk spontaneously in attendance of *SIG*, *SIG* tracks some talker and then changes focus-of-attention to others. The observed behavior is evaluated by checking the internal states of *SIG*; that is, auditory and visual localization shown in the radar chart, auditory, visual, and associated streams shown in the stream chart, and peak extraction as shown in Figure 6.

The top-right image consists of the radar chart (left) and the stream chart (right) updated in real-time. The former shows the environment recognized by *SIG* at the moment of the snapshot. A pink sector indicates a visual field of *SIG*. Because of using the absolute coordinate, the pink sector rotates as *SIG* turns. A green point with a label is the direction and the face ID of a visual stream. A blue sector is the direction of an auditory stream. Green, blue and red lines indicate the direction of visual, auditory and associated stream, respectively. Blue and green *thin* lines indicate auditory and visual streams, respectively. Blue, green and red *thick* lines indicate associated streams with only auditory, only visual, and both information, respectively.

The bottom-left image shows the auditory viewer consisting of the power spectrum and auditory event viewer. The latter shows an auditory event as a filled circle with its pitch in X axis and its direction in Y axis. The bottom-right image shows the visual viewer captured by the *SIG*'s left eye. A detected face is displayed with a red rectangle. The top-left image shows the scene of this experiment recorded by a video camera.

The temporal sequence of *SIG*'s recognition and actions shows that the design of companion robot works well and pays its attention to a new talker. The current system has attained a passive companion. To design and develop an active companion may be important future work.



a) When a participant comes and says "Hello", *SIG* turns toward him. b) *SIG* asks his name and he introduces himself to it.

Fig. 4. Task-oriented Control of Friendly *SIG*



a) A participant says "Hello". b) *SIG* turns away from him.

Fig. 5. Task-oriented Control of Hostile *SIG*



Fig. 6. Socially-oriented Control of Friendly *SIG*. Scene (upper-left), radar and sequence chart (upper-right), spectrogram and pitch-vs-direction chart (lower-left), and face-tracking chart (lower-right).

C. Observation: *SIG* as a non-verbal Eliza

As socially-oriented attention control, interesting human behaviors are observed. The mechanism of associating auditory and visual streams and that of socially-oriented attention control are explained in advance to the user.

1. Some people walk around talking with their hand converging *SIG*'s eyes in order to confirm the performance of auditory tracking.
2. Some people creep on the floor with talking in order to confirm the performance of auditory tracking.
3. Some people play hide-and-seek games with *SIG*.
4. Some people play sounds from a pair of loud speakers with changing the balance control of pre-amplifier in order to confirm the performance of auditory tracking.

5. When one person reads loud a book and then another person starts to read loud a book, *SIG* with *Dominant* personality turns its head to the second talker for a short time and then is back to the first talker and keeps its attention on him/her. On the contrary, *SIG* with *Submissive* personality often turns its head to each talker. In either case, the value of α is set to 1.

Above observations remind us of Eliza [25], although *SIG* does not say anything except a receptionist robot. When the user says something to *SIG*, it turns to him/her, which invites the participation of the user into interaction. *SIG* also invites exploration of the principles of its functioning, that is, the user is drawn in to see how *SIG* will respond to variations in behavior. Since *SIG* takes only passive behaviors, it does not arouse higher expectations of verisimilitude that it can deliver on.

Needless to say, there are lots of work remaining to validate the proposed approach for personality of artifacts. We are currently working to incorporate active social interaction by developing the capability of listening to simultaneous speeches.

VI. CONCLUSIONS

In this paper, we demonstrate that auditory and visual multiple-talker tracking subsystem can improve social aspects of human robot interaction. Although a simple scheme of behavior is implemented, human robot interaction is drastically improved by real-time multiple-talker tracking system. We can pleasantly spend an hour with *SIG* as a companion robot even if its behavior is quite passive.

Since the Interpersonal Theory research community provides software for analysing circumplex correlation matrices, we have plan to gather the data of user interaction to evaluate whether the presented architecture of selective attention based on personality realizes the target variation of personality. This pursuit may lead to a general theory of personality for software agents and humanoid robots. This research was partially supported by the Ministry of Education, Culture, Sports, Science and Technology, Grant-in-Aid "Informatics" No.14019051.

REFERENCES

- [1] R. A. Brooks, C. Breazeal, R. Irie, C. C. Kemp, M. Marjanovic, B. Scassellati, and M. M. Williamson, "Alternative essences of intelligence," in *Proceedings of 15th National Conference on Artificial Intelligence (AAAI-98)*. 1998, pp. 961–968, AAAI.
- [2] S. Handel, *Listening*, The MIT Press, MA., 1989.
- [3] C. Breazeal and B. Scassellati, "A context-dependent attention system for a social robot," in *Proceedings of the 16th International Joint Conference on Artificial Intelligence (IJCAI-99)*, 1999, pp. 1146–1151.
- [4] C. Breazeal, "Emotive qualities in robot speech," in *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS-2001)*. 2001, pp. 1389–1394, IEEE.
- [5] H.E. Pashler, *The Psychology of Attention*, The MIT Press, MA., 1997.
- [6] J.M. Wolfe, K. R. Cave, and S.L. Franzel, "Guided search: An alternative to the feature integration model for visual search.," *Journal of Experimental Psychology: Human Perception and Performance*, vol. 15, no. 3, pp. 419–433, 1989.
- [7] E. C. Cherry, "Some experiments on the recognition of speech, with one and with two ears.," *Journal of Acoustic Society of America*, vol. 25, pp. 975–979, 1953.
- [8] B. Reeves and C. Nass, *The Media Equation: How People Treat Computers, Television, and New Media Like Real People and Places*, Cambridge University Press, Cambridge, UK, 1996.
- [9] D.J. Kiesler, "The 1982 interpersonal circle: A taxonomy for complementarity in human transactions.," *Psychological Review*, vol. 90, pp. 185–214, 1993.
- [10] J. Bates, "The role of emotion in believable agents," *Comm. of the ACM*, vol. 37, no. 7, pp. 122–125, 1994.
- [11] J. Cassell, "More than just another pretty face: Embodied conversational interface agents," *Comm. of the ACM*, vol. 43, no. 4, pp. 70–78, 2000.
- [12] B. Hayes-Roth, G. Ball, C. Lisetti, R. Picard, and A. Stern, "Affect and emotion in the user interface.," in *Proceedings of 1998 International Conference on Intelligent User Interfaces*. 1998, pp. 91–96, ACM.
- [13] H. Miwa, A. Takanishi, and H. Takanobu, "Experimental study on robot personality for humanoid head robot.," in *Proceedings of 2001 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS 2001)*. 2001, pp. 1183–1188, IEEE.
- [14] T. Ono, M. Imai, and H. Ishiguro, "A model of embodied communications with gestures between humans and robots.," in *Proceedings of Twenty-third Annual Meeting of the Cognitive Science Society (CogSci2001)*. 2000, pp. 732–737, AAAI.
- [15] S. Waldherr, S. Thrun, R. Romero, and D. Margaritis, "Template-based recognition of pose and motion gestures on a mobile robot," in *Proceedings of 15th National Conference on Artificial Intelligence (AAAI-98)*. 1998, pp. 977–982, AAAI.
- [16] Y. Matsusaka, T. Tojo, S. Kuota, K. Furukawa, D. Tamiya, K. Hayata, Y. Nakano, and T. Kobayashi, "Multi-person conversation via multi-modal interface — a robot who communicates with multi-user," in *Proceedings of 6th European Conference on Speech Communication Technology (EUROSPEECH-99)*. 1999, pp. 1723–1726, ESCA.
- [17] K. Nakadai, K. Hidai, H. Mizoguchi, H. G. Okuno, and H. Kitano, "Real-time auditory and visual multiple-object tracking for robots," in *Proceedings of the 17th International Joint Conference on Artificial Intelligence (IJCAI-01)*. 2001, pp. 1425–1432, IJCAI.
- [18] H.G. Okuno, K. Nakadai, K. Hidai, H. Mizoguchi, and H. Kitano, "Human-robot interaction through real-time auditory and visual multiple-talker tracking," in *Proceedings of IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS-2001)*. 2001, pp. 1402–1409, IEEE.
- [19] H.G. Okuno, K. Nakadai, K. Hidai, H. Mizoguchi, and H. Kitano, "Human-robot non-verbal interaction empowered by real-time auditory and visual multiple-talker tracking," *Advanced Robotics*, vol. 17, no. 2, pp. in print, 2003.
- [20] K. Nakadai, H. G. Okuno, and H. Kitano, "Exploiting auditory fovea in humanoid-human interaction," in *Proceedings of the 18th National Conference on Artificial Intelligence (AAAI-2002)*. 2002, pp. 431–438, AAAI.
- [21] K. Nakadai, T. Matsui, H. G. Okuno, and H. Kitano, "Active audition system and humanoid exterior design," in *Proceedings of IEEE/RAS International Conference on Intelligent Robots and Systems (IROS 2000)*. 2000, pp. 1453–1461, IEEE.
- [22] T. Kawahara, A. Lee, T. Kobayashi, K. Takeda, N. Minematsu, K. Itou, A. Ito, M. Yamamoto, A. Yamada, T. Utsuro, and K. Shikano, "Japanese dictation toolkit – 1997 version –," *Journal of Acoustic Society Japan (E)*, vol. 20, no. 3, pp. 233–239, 1999.
- [23] K. Nakadai, T. Lourens, H. G. Okuno, and H. Kitano, "Active audition for humanoid," in *Proceedings of 17th National Conference on Artificial Intelligence (AAAI-2000)*. 2000, pp. 832–839, AAAI.
- [24] S. Kagami, K. Okada, M. Inaba, and H. Inoue, "Real-time 3d optical flow generation system," in *Proc. of International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI'99)*, 1999, pp. 237–242.
- [25] J. Weizenbaum, "Eliza – a computer program for the study of natural language communication between man and machine," *Comm. of the ACM*, vol. 9, no. 1, pp. 36–45, 1966.