

User Acceptance of Autonomously Generated Social Head Gaze

Vasant Srinivasan
Texas A & M University
College Station, TX, USA
vasant_s18@tamu.edu

Robin Murphy
Texas A & M University
College Station, TX, USA
murphy@cse.tamu.edu

Zachary Henkel
Texas A & M University
College Station, TX, USA
zmhenkel@cse.tamu.edu

ABSTRACT

This paper shows that *autonomous* generation of head gaze using our novel behavioral framework maintains a high level of user acceptance equivalent to that from *manual coding* of head gaze for robotic systems. The current practice for implementing social head gaze is to either use *manual coding*, which requires significant human effort and restricts real-time dialog, or *autonomous* generation through ad hoc implementations addressing only one or two of the five known social contexts. Our behavioral robotics framework approach *autonomously* produces head gaze acts without semantic understanding or modeling of the agents' theory of mind using eight novel affordances from sentence structure and time intervals between exchanges along with coordination functions to trigger behaviors that generate head gaze acts loosely synchronized with speech. The correctness and sufficiency of the *autonomously* generated head gaze acts was validated by a 93 participant study in which participants played the role of trapped victim in a simulated disaster setting and interacted with a rescue robot for 15 minutes. The *autonomous* condition was evaluated more positively ($p < 0.05$) compared to *no gaze* condition for user's affective state and perceptions of the robot. There was no statistical significance between the *autonomous* and *manual* conditions for any of the measures. Our method is real-time, repeatable, extensible, and reduces the work in Wizard of Oz type studies.

Categories and Subject Descriptors

H.1.2 [Models and Principles]: User/Machine Systems—*Human Factors*

General Terms

Gaze, Human Robot Interaction, Model, Behaviors

Keywords

Gaze, Human Robot Interaction, Model

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICMI 2012 Workshop on Speech and Gesture Production in Virtually and Physically Embodied Conversational Agents, October 26, 2012, Santa Monica, CA, USA.

Copyright 2012 ACM 978-1-4503-1514-2/12/10 ...\$15.00.

1. INTRODUCTION

Social head gaze for robots has been shown to be beneficial in applications such as eldercare [17], entertainment [20, 3], search and rescue [2, 27], and telepresence [5, 15]. Social head gaze has been shown to have benefits like increased task performance [20, 3, 26], improved perception of robots physical, social and intellectual characteristics like perceived intelligence [3, 21], understandability [3, 26, 9], reliability [26, 28], likeability [20, 27], and positive emotional response by participant's [26, 27].

Autonomous generation of head gaze has two major benefits in theory. First, it eliminates the human effort and semantic understanding needed to *manually* hardcode the association between what the agents are saying and the appropriate gaze act. Second, *autonomous* generation allows head gaze acts to be produced in real-time with unpredictable dialogs. In practice, *autonomous* generation of head gaze acts has been handicapped by barriers in semantic understanding or modeling of the beliefs, desires, or intentions of the conversant. As a result, *autonomous* systems such as [12, 18, 11, 24, 1, 19] have been able to generate acts for no more than two of the five social contexts for social head gaze [28].

The central tenet of this paper is that "good enough" head gaze acts can be *autonomously* generated without semantic understanding or modeling. The generation process relies on three mechanisms from behavioral robotics: affordances, behaviors and coordination functions. The behavioral framework in Section 3 identifies and implements eight novel affordances in the sentence structure and time intervals that cue a particular behavior which generates the gaze act (e.g., the behavior of back-channeling during a conversation). These affordances and behaviors correspond to the social contexts identified in previous research [28]. The implementation uses two coordination functions for smooth transitioning of behaviors between five social contexts and resolving conflicts, as it is possible for multiple gaze act behaviors to be triggered or active at the same time.

2. RELATED WORK

The *manual coding* of head gaze acts is more precise than *autonomous* head gaze and has been used by 10 studies [16, 20, 14, 30, 21, 22, 29, 13, 25, 10]. However it is undesirable because the speech-gaze synchronization is done a-priori as the cues to generate head gaze are not easily computable, especially those based on sentence understanding like the theme or rheme in a sentence. There is little scope for a more open and richer interaction.

A total of 10 studies implement head gaze *autonomously*

without *manual coding* at any stage [12, 18, 11, 3, 26, 24, 1, 19, 23, 9], however they have two important shortcomings. First, out of the seven studies that operate within the *regulation of an interaction process* context [12, 24, 26, 1, 19, 23, 9] six studies [12, 24, 26, 1, 19, 23] ignore the rules from communication theory for implementation of this context. The two remaining studies [9, 10] use only a smaller subset of the appropriate rules necessary for implementing this context. This is because the cues humans use to generate head gaze are not easily computable, especially those based on semantic understanding like the theme or rheme in a sentence. Second, the majority of these studies (eight out of 10) [12, 18, 11, 24, 1, 19, 23, 9] use just one or two social contexts with only two studies use more than two social contexts [3, 26].

None of studies that generate *autonomous* head gaze explore conflict resolution that may arise because of multiple active gaze behaviors. However one study that used *manual coding* to synchronize gaze acts and pre-recorded speech [10] addressed conflict resolution of head gaze behaviors for two social contexts. The study prioritized gaze behaviors using outcomes, that is rules with task outcomes are preferred outcomes such as building rapport. These priority rules apply for only two out of the five social contexts *regulation of an interaction process* and *manifestation of an interaction*.

3. APPROACH

This section discusses behavioral robotics framework that anchor the affordances that triggers gaze behaviors and the coordination of multiple behaviors. Our approach overcomes the four shortcomings listed in Section 2 and uses affordances from sentence structures and time intervals that are easily computable to generate *autonomous* head gaze acts synched with synthetic speech in real time and coordinate them for five social contexts. Figure 1 lists the eight behaviors identified from previous research [28, 27] (column 1), how the behavior would be triggered from either understanding the semantic context or having a model of the dialog (column 2), how this paper substitutes an affordance of either sentence structure or timing (column 3), which produce a social head gaze act.

3.1 Affordances

The motivation for using affordances from sentence structures is derived from two sources: the superiority of information structure approaches in dialog, mostly notably [4], over turn taking [6], and a novel behavioral model of engagement [28, 27] which maps the social context to a specific gaze act.

The first four affordances (rows 1 to 4) in Figure 1 are the sentence structure approximations of cues used in literature for dialog based on information structure [4] and used to trigger the *fixate(human)* and *avert(human)* gaze acts. The rationale for the approximation is that in the English language the theme always occurs at the beginning of an independent clause/simple sentence and rheme occurs towards the end of the independent clause/simple sentence [8]. The affordance listed in Row 5 is a tag $\langle \text{fixate}, \text{objectname} \rangle$ inserted before an object name. The tag can be inserted either manually when there is an operator present and is typing or dictating a sentence or can be *autonomously* inserted by a reasoning system. This affordance is used to trigger a *fixate(object)* gaze act. Row 6 is an affordance based on timer interval and used to trigger the *concurrency* gaze act.

Behaviors	Cue from Literature	Affordances in Sentence Structure and Timing	Social Gaze Act
1 Start of turn	Start of turn and start of theme (Cassell 98, Mutlu 06, and Mutlu 09a)	Initial word (Start of sentence and start of turn)	Aversion (human)
2 Middle of turn (a)	Start of theme and middle of turn (Cassell 98, Mutlu 06, Mutlu 09a)	Word following punctuation : ! ? (Start of sentence and middle of turn)	Aversion (human) (Probability = 0.73)
3 Middle of turn (b)	Start of rheme and middle of turn (Cassell 98, Mutlu 06, and Mutlu 09)	After 75% of words between punctuations : ! ? (Towards end of sentence and middle of turn)	Fixation (human) (Probability = 0.70)
4 End of turn	Start of rheme and end of turn (Cassell 98, Mutlu 06, and Mutlu 09a)	Carriage return at end of paragraph (Towards end of sentence and end of turn)	Fixation (human)
5 Robot manifesting interaction	Look at object 800ms to 1 sec before utterance of object (Suadee & Crocker 09 a, b.)	The object tag	Fixation (object)
6 Back-channeling	Ad hoc (Breazeal 03, Sidner 05, and Gu 06)	> IERT Elapsed response time	Concurrency (human)
7 Establishing agency	Set time interval t (Breazeal 03, Sidner 05)	> IERT Idle elapsed time	Scan (environment)
8 Projecting mental state	Robot does not understand (Breazeal 03)	Deleting, retying	Confusion (human)

Figure 1: Overview of the behavioral framework

It is based elapsed response time: time since the robot last responded to the human during a dialog. The affordance listed in Row 7 is a timeout on the idleness of the robot and used to trigger the *scan(environment)* gaze act. The last affordance in Row 8 is based on typing and used only when an operator is using the robot. The number of deletes and retying is taken in account and if it exceeds a certain preset value, it triggers the *confusion* gaze act. The affordances from sentence structure and time intervals are close approximations of cues from sentence understanding, can be easily computed and produce *autonomous* head gaze acts that are loosely synchronized with speech. It should be noted that here that Figure 1 only lists affordances from sentence structure for four out of the five social contexts. This is because the behavior *communicating social attention* does not depend on what is being said or time intervals, instead occurring during periods of non-verbal activity. The Figure 1 can be expanded if new studies identify new uses of head gaze or acts.

3.2 Coordination Functions

The coordination of behaviors is done in two stages based on prioritization by timestamps and type. In the first stage a set of primitive behaviors ($b_{e_{primitive_i}}$) within a larger meta or abstract behavior ($b_{e_{abstract}}$) like *Regulation of Interaction* or *Manifesting an Interaction* (shown as dotted lines in Figure 2) is coordinated by assigning priorities on timestamps. This function can be captured as:

$$b_{e_{abstract}} = \text{Max}(\text{timestamp}(b_{e_{primitive_i}}))$$

As seen from the above equation the most recent behavior has the highest priority and is the output of the abstract behaviors. In the second stage all current active behaviors B_e are coordinated for the overall robotic response ρ_g (ORR). This coordination function is based on both timestamps (highest priority for the most recent behavior) and the type of behavior (atomic or non atomic). The atomic behaviors always run through completion, cannot be interrupted. However the non-atomic behaviors are interruptible.

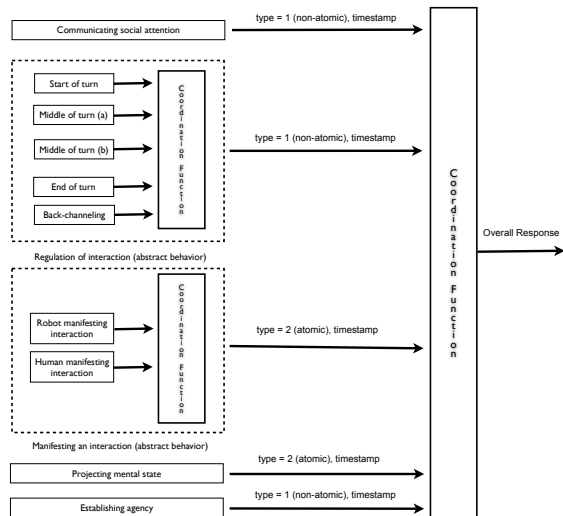


Figure 2: Behaviors and Coordination Functions

The algorithm below describes this function.

```

if ( $Max(type(B_e)) == "atomic"$ ) then
   $\rho_g = (G * Min(B_e(type = "atomic", timestamp)))$ ;
else
   $\rho_g = (G * Max(timestamp(B_e)))$ ;
end if

```

4. IMPLEMENTATION

The implementation of the behavioral framework currently requires a script with typed text and punctuation. This is available transparently to a robot that is a social actor, or can be typed in or dictated by a robot operator if it is tele-operated.

The behaviors were implemented as production rules. The production rules are a set of reactive behaviors, where if a social context perception cue is perceived, the gaze act is called. There are a total of 10 production rules for the five social contexts used in our behavioral framework [28, 27]. An example production rule would be - "If the robot and human(s) are engaged in a conversational interaction, the robot should always *avert* at the start of a theme and start of a turn". This production rule corresponds to the *start of turn* behavior. The affordance in sentence structure for this production rule is the *Initial Word* of the sentence. Hence when a robot is having a conversation with a human, when the head gaze generation system encounters the *Initial Word* in the current buffer, the *start of turn* reactive behavior is fired. The response for this behavior is a *avert* gaze act and this will be the overall response of the robot if there are no conflicts with other behaviors.

The identification of affordances from sentence and timing operates over a textual input. The input from the script fills the input buffer which is a data structure of variable length (n) delimited by the start of a sentence and the end of a sentence. If an affordance from sentence structure is perceived the corresponding production rule is fired and return an output buffer which interleaves the head gaze acts with the words to be spoken. This output buffer is passed through Microsoft speech system, which converts the text to speech and triggers the gaze act when a gaze tag is encountered.

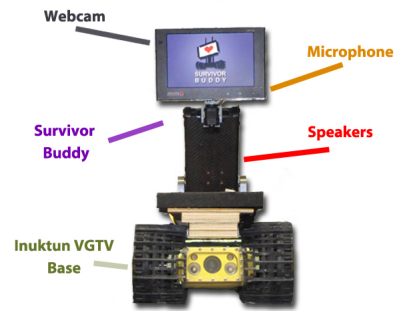


Figure 3: Survivor Buddy

The running time for this implementation is $O(n)$ where n is the length of the text that needs to be parsed.

5. EXPERIMENTS

A 1 x 3 (*no gaze*, *manual gaze*, *autonomous gaze*) between subjects comparative study was conducted to validate the correctness and sufficiency of *autonomous* head gaze. The study investigated two hypotheses:

- H1 - Participants who interact with the robot in the *autonomous* head gaze condition will rate the robot more positively than participants who interact with the robot in the *no gaze* condition. This hypothesis tests the correctness of *autonomous* head gaze.
- H2 - The *autonomous* head gaze condition will perform as well as *manual* head gaze condition when compared to the *no gaze* condition. This hypothesis tests whether the simplifying inferences used to enable *autonomous* head gaze are sufficient for social head gaze.

The study used 93 participants consisting of 53 males and 40 females with an age range of 18-67 ($M = 32.38, SD = 14.84$). On a scale of 1 to 7, the participants familiarity with robots was low ($M = 2.09, SD = 1.67$). The *no gaze* condition represents the state of the practice and thus was the baseline case. Random head gaze acts were expected to produce confounding effects preventing the understanding of whether all of the reported acts worked together.

5.1 Robot, Setting and Dialog

The study used Survivor Buddy, an affective multimedia head mounted on an Inuktun Extreme-VGTV robot (refer to Figure 3). Survivor Buddy has four degrees of freedom and is capable of very agile movements. Survivor Buddy's head is a MIMO 740 touch screen monitor with a webcam and microphone. The speaker system is mounted within its neck.

The setting for the study was a simulated disaster scenario where in Survivor Buddy, a rescue robot shown in Figure 3 has a dialog with a trapped victim using the guidelines for a 911 dispatcher. This setting was chosen because i) it has been shown to elicit affect [2], ii) the dispatcher protocol involves all five social contexts (none of the 21 previous studies reported a dialog that captured all known social head gaze acts), and iii) participants tend to naturally follow the protocol and thus a pre-scripted dialog would not confound the experiments.

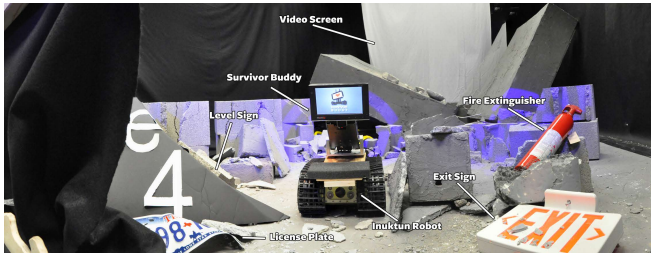


Figure 4: The simulated disaster area from the participants point of view.

The study used purposeful dialog that followed the 911 dispatch protocol focussing on assessing the participants physical health and gaining information as to location and nature of the event. The dialog involves all five social contexts with a total of 162 possible gaze acts. It should be noted here that some gaze acts were produced probabilistically and thus were not guaranteed to be used every time. For the *manually coded* condition, the 162 head gaze acts were *manually coded* following the rules based on semantic understanding (see Figure 1 column 2) by inserting tags into the input buffer and supplying the location of the objects. However for the *autonomous* condition 154 gaze acts were *autonomously* generated. The remaining eight gaze acts belonging to *manifesting an interaction process* is considered *manual*. This is because of two reasons. 1) Object recognition are not consistently repeatable in a dark high fidelity search and rescue environment. Hence the locations of objects was pre-determined to get a more controlled study, 2) The current study also does not use any reasoning system or knowledge base like [26, 3, 10]. It is expected that the *autonomous* head gaze system described in this paper and the behavioral framework from [28, 27] will complement these systems.

5.2 Role of Participant and Robot

This section details the role of the participants and role of the robot in the 93 participant study used to validate the *autonomous* head gaze acts generated by a US&R robot in a simulated disaster setting.

The participants laid down in a simulated confined space wooden box (3 ft x 3 ft x 8 ft) in a dark room and played the role of a victim. A video of a first person view of a parking garage collapse was played with flashes of light to increase emotional arousal. The surroundings became visible, a simulated collapsed parking structure (refer to Figure 4). The Survivor Buddy robot system was placed at a distance of 1.22 meters, within the participant’s personal zone [2], so as to increase the likelihood of social interaction. The screen was a Survivor Buddy logo so that they only social cues were voice and head gaze acts. .

Survivor Buddy then engaged the participant, asking the questions following the pre-scripted dialog. The robot supervisor (hidden) would activate the text for the robot’s turn after the participant stopped speaking from a user interface. The entire interaction last for about 15 minutes. The dialog initially follows the 911 dispatch protocol focussing on assessing the victim’s physical health and gaining information as to location and nature of the event. Then the dialog shifts to questions that assess mental functioning. An ex-

Measure	Main Effect	Post-Hoc Results			Mean, Standard Deviation		
		Manual vs No gaze	Autonomous vs No gaze	Autonomous vs Manual	Autonomous gaze	Manual gaze	No gaze
Creativity	F[2, 90] = 3.301, p < .05, η = .073	p < .05	p < .05	ns	14, 4.49	13.96, 4.43	11.17, 3.5
Person at Ease	F[2, 90] = 3.8, p < .05, η = 0.079	ns	p < .05	ns	6.23, 1.63	6.19, .91	5.61, 1.18
Emotional Response: Positivity	F[2, 90] = 3.97, p < .05, η = .081	ns	p < .05	ns	7.91, 1.49	7.75, 1.36	6.94, 1.45
Emotional Response: Comfort	F[2, 90] = 8.43, p < .001, η = .157	p < .001	p < .001	ns	7.68, 1.42	7.61, 1.7	6.09, 1.97
Robot Likeability	F[2, 90] = 6.75, p < .01, η = .130	p < .01	p < .01	ns	4.93, .601	4.97, .63	4.36, .89
Robot Empathy	F[2, 90] = 4.541, p < .05, η = .092	p < .05	p < .05	ns	5.59, .93	5.65, .89	5.02, .86
Robot Loyalty	F[2, 90] = 4.079, p < .05, η = .083	ns	p < .05	ns	5.43, .85	5.27, .81	4.75, .099
Robot Integrity	F[2, 90] = 3.791, p < .05, η = .078	ns	p < .05	ns	6.02, .911	5.95, .8	5.49, .76
Robot Caring	F[2, 90] = 4.197, p < .05, η = .085	p < .05	p < .05	ns	5.18, .83	5.23, .89	4.56, 1.23
Natural Movement	F[2, 90] = 16.7, p < .001, η = .270	p < .001	p < .001	ns	4.7, 1.62	4.83, 1.6	2.77, 1.38
Gaze-Speech Synchronization	F[2, 90] = 47.9, p < .001, η = .515	p < .001	p < .001	ns	5.9, 1.27	5.78, 1.08	2.93, 1.63
Looking at Objects at Appropriate Times	F[2, 90] = 14.6, p < .001, η = 0.54	p < .001	p < .001	ns	5.72, 1.67	6.1, .99	3.28, 1.43
Robot Engagement	F[2, 90] = 3.301, p < .05, η = .073	ns	p < .05	ns	6.06, .89	5.87, 1.17	5.26, 1.65
Human-Like Behavior	F[2, 90] = 8.9, p < .001, η = 0.165	p < .001	p < .001	ns	5.1, 1.63	4.74, 1.34	3.5, 1.74
Understanding Robot Behavior	F[2, 90] = 18.087, p < .001, η = .56	p < .001	p < .001	ns	5.57, .95	5.63, 1.59	2.83, 1.15

Figure 5: Results

ample question was: “Paper is used commonly for writing. Can you name as many alternative uses as you can?”. If the robot is presenting a cognitively acceptable presence, then the participant will be able to focus more on tasks and be less distracted. The script concluded with the participant being told that rescuers had arrived. The participants then completed a 143 question post interaction assessment, after [2, 20, 7], in order to access participants’ affective state and perceptions of the robot.

6. RESULTS

This section reports the results for 15 measures used in the 93 participant study to validate the correctness and sufficiency of *autonomous* head gaze, and the performance of coordination functions. Figure 5 lists the main effects (column 2) from a ANOVA and tukeys post-hoc results between the conditions (columns 3, 4, and 5) for the each of measures (column 1) and the mean and standard deviation for each of the measure (columns 6, 7, and 8).

The results validate the correctness of the head gaze acts generated by *autonomous* head gaze condition. The *autonomous* head gaze acts generated by the robot is correct if it elicits positive ratings from the participants in key measures of user’s affective state and perceptions of the robot listed in Figure 5 and if the results are similar to or better than those observed in prior studies. This is tested by hypothesis H1. As predicted by H1, the participants who interacted with the robot in the *autonomous* head gaze condition rated the robot more positively in all 15 measures ($p < .05$) when compared to participants who interacted with robot in the *no gaze* condition.

The results validate that the affordances from sentence used to generate *autonomous* head gaze is sufficient. The affordances from the sentence structure used to enable *autonomous* head gaze are sufficient for social head gaze if the *autonomous* head gaze condition performs as well the *manual* head gaze condition. This is tested by hypothesis H2. Consistent with H2, there was no statistically significant re-

sults between the participants who evaluated the robot in the *autonomous* head gaze condition and participants who evaluated the robot in the *manual* head gaze condition for any of the measures.

The results validate the choice and two stage approach for implementing the coordination functions. The choice and approach for implementing coordination function is correct if the generated gaze behaviors are consistent with the robots communication intent and the robot elicits positive ratings from the participants in measures for appropriateness of robot's head movements. Four measures used to test for consistency and appropriateness of robot's head movements: understandability of robot behaviors, natural movement, gaze-speech synchronization and looking at objects at appropriate times showed highly positive results as seen in Figure 5 validating the choice and two stage approach used for implementing coordination functions.

7. DISCUSSION

The results from the study in Section 6 indicate that the head gaze acts generated in the *autonomous* condition are correct, the affordances from the sentence structure and time intervals used to enable *autonomous* head gaze are sufficient for social head gaze, and the coordination function enables seamless transition between five different social contexts naturally and successfully. These results are important because a) Tedious hand coding of every gaze-speech synchronization wizard of oz style is not required, b) Loosely synchronized gaze acts are sufficient for user acceptance and semantic understanding of conversation may not be required.

The high levels of user acceptance observed in the study for *autonomous* condition are surprising because the condition has less precise synchronization. This may be because of two reasons a) The expectations of gaze from a robot with a low degree of anthropomorphism may not be as high as those for human interpersonal interaction, b) The conversation was purposeful and the robot's role as a response professional is typically not emotive. These reasons require further investigation.

The benefits of using our *autonomous* head gaze generation system is that it is real-time having a runtime of $O(n)$, independent of dialogue, repeatable and extensible. These qualities along with the high quality of communication it produces make it a preferred method of implementation for social head gaze.

The autonomous generation of social head gaze described in this paper is also applicable to robots that are tele-operated or puppeteered. Social head gaze becomes very valuable in scenarios where the face of the operator is not visible, like the telenoid robot from ATR [15]. In such a scenario the operator would type the appropriate responses in chat like interface on the remote end and the robot will perform autonomous gaze acts on the local side. If speech was used, dictation systems in combination with prosody could determine the sentence structure.

The study was limited in two ways, but these limitations do not appear to be serious. First, it tested only robot head gaze acts for a single party conversation in one domain. Further explorations of multi-party conversations, other domains, and eye gaze acts is warranted. However, the strong results indicate that the novel inference mechanisms may be of similar value for multi-party dialog and other domains. The second limitation was that of the implementation itself.

The affordances assumed that the robot's speech appeared in punctuated text, which could be parsed from speech and prosody. The implementation could not *autonomously* generate acts for the *manifestation of an interaction* behavior because the robot did not have object recognition capabilities. These limitations can be expected in the future to disappear in the future. Even if the limitation persisted, the implementation is of great practical value because on average *manual coding* was required on average in only 5% of the conversation.

8. CONCLUSIONS

The study shows that *autonomous* generation of head gaze acts, even if loosely synchronized, is sufficient for social interactions and provides the correct act for the context. As seen in Section 6 the robot in the autonomous gaze condition was evaluated more positively than the *no gaze* condition in all 15 measures of user's affective state and perceptions of the robot, and *autonomous* head gaze performed as well as the *manual* head gaze in each of those measures. The dialog used in the study incorporated all five known social contexts providing the most complete coverage found in the literature to date. This contextual coverage provided extensive opportunities for social violations to occur and elicit a negative perception of the robot, but it did not.

While this paper focuses on the user acceptance of the *autonomously* generated loosely synchronized gaze acts, it suggests that the behavioral robotics framework is a useful software mechanism. The novel use of sentence structure and timing as affordances which replace semantic understanding of the conversation along with the coordination functions that allow seamless transition between five different social contexts successfully and naturally are well-suited for the *autonomous* generation of head gaze acts by artificial agents. While the implementation was limited to typed text and did not have access to computer vision, the software has the advantages of modularity, extensibility, and low order complexity $O(n)$. It can immediately replace most, and with object recognition eventually all, of *manual coding* of gaze acts for Wizard of Oz experiments or surrogates.

9. ACKNOWLEDGEMENTS

This work was supported in part by NSF IIS - 0905485 "The Social Medium is the Message", and Microsoft External Research. The author's would like thank Clifford Nass and Jessica Gonzales for their helpful comments and suggestions.

10. REFERENCES

- [1] M. Bennewitz, F. Faber, D. Joho, M. Schreiber, and S. Behnke. Towards a humanoid museum guide robot that interacts with multiple persons. In *in Proc. of the IEEE/RSJ International Conference on Humanoid Robots (Humanoids)*, 2005.
- [2] C. Bethel and R. Murphy. Non-facial and non-verbal affective expression for appearance-constrained robots used in victim management. *Paladyn. Journal of Behavioral Robotics*, 1:219-230, 2010.
- [3] C. Breazeal, C. D. Kidd, A. L. Thomaz, G. Hoffman, and M. Berlin. Effects of nonverbal communication on efficiency and robustness in human-robot teamwork. In *in Proceedings of IROS*, 2003.

- [4] J. Cassell, O. Torres, and S. Prevost. Turn taking vs. discourse structure: How best to model multimodal conversation. In *In Machine Conversations*. Kluwer, 1998.
- [5] M. Cherubini, R. de Oliveira, N. Oliver, and C. Ferran. Gaze and gestures in telepresence: multimodality, embodiment, and roles of collaboration. In *CoRR*, 2010.
- [6] S. Duncan. Some signals and rules for taking speaking turns in conversations. *Journal of Personality and Social Psychology*, 23(2):283 – 292, 1972.
- [7] V. Groom, V. Srinivasan, C. Bethel, R. Murphy, L. Dole, and C. Nass. Responses to robot social roles and social role framing. In *Collaboration Technologies and Systems (CTS), 2011 International Conference on*, pages 194 –203, may 2011.
- [8] M. Halliday. *Intonation and grammar in British English*. Janua linguarum: Series practica. Mouton, 1967.
- [9] A. Holroyd, C. Rich, C. L. Sidner, and B. Ponsler. Generating connection events for human-robot collaboration. In *RO-MAN, 2011 IEEE*, pages 241 –246, 31 2011-aug. 3 2011.
- [10] C. Huang and B. Mutlu. Robot behavior toolkit: Generating effective social behaviors for robots. *Proceedings of to the 7th ACM/IEEE Conference on Human-Robot Interaction (HRI 2012)*, 2012.
- [11] M. Imai, T. Kanda, T. Ono, H. Ishiguro, and K. Mase. Robot mediated round table: Analysis of the effect of robot’s gaze. In *Robot and Human Interactive Communication, 2002. Proceedings. 11th IEEE International Workshop on*, pages 411 – 416, 2002.
- [12] M. Imai, T. Ono, and H. Ishiguro. Physical relation and expression: joint attention for human-robot interaction. In *Robot and Human Interactive Communication, 2001. Proceedings. 10th IEEE International Workshop on*, pages 512 –517, 2001.
- [13] C. T. Ishi, C. Liu, H. Ishiguro, and N. Hagita. Head motions during dialogue speech and nod timing control in humanoid robots. In *Proceedings of the 5th ACM/IEEE international conference on Human-robot interaction, HRI ’10*, pages 293–300, Piscataway, NJ, USA, 2010. IEEE Press.
- [14] Y. Kuno, K. Sadazuka, M. Kawashima, K. Yamazaki, A. Yamazaki, and H. Kuzuoka. Museum guide robot based on sociological interaction analysis. In *Proceedings of the SIGCHI conference on Human factors in computing systems, CHI ’07*, pages 1191–1194, New York, NY, USA, 2007. ACM.
- [15] C. Liu, C. T. Ishi, H. Ishiguro, and N. Hagita. Generation of nodding, head tilting and eye gazing for human-robot dialogue interaction. In *Proceedings of the seventh annual ACM/IEEE international conference on Human-Robot Interaction, HRI ’12*, pages 285–292, New York, NY, USA, 2012. ACM.
- [16] K. F. Macdorman, T. Minato, M. Shimada, S. Itakura, S. Cowley, and H. Ishiguro. Assessing human likeness by eye contact in an android testbed. *Science*, pages 21–23, 2004.
- [17] M. Maja. The role of embodiment in assistive interactive robotics for the elderly, 2005.
- [18] Y. Matsusaka, S. Fujie, and T. Kobayashi. Modeling of conversational strategy for the robot participating in the group conversation. In *INTER_SPEECH’01*, pages 2173–2176, 2001.
- [19] N. Mitsunaga, C. Smith, T. Kanda, H. Ishiguro, and N. Hagita. Adapting robot behavior for human–robot interaction. *Robotics, IEEE Transactions on*, 24(4):911 –916, aug. 2008.
- [20] B. Mutlu, J. Forlizzi, and J. Hodgins. A storytelling robot: Modeling and evaluation of human-like gaze behavior. In *in Proc. International Conference on Humanoid Robots. IEEE*, 2006.
- [21] B. Mutlu, T. Shiwa, T. K, H. Ishiguro, and N. Hagita. Footing in human-robot conversations: How robots might shape participant roles using gaze cues.
- [22] B. Mutlu, F. Yamaoka, T. Kanda, H. Ishiguro, and N. Hagita. Nonverbal leakage in robots: communication of intentions through seemingly unintentional behavior. In *HRI ’09: Proceedings of the 4th ACM/IEEE international conference on Human robot interaction*, pages 69–76, New York, NY, USA, 2009. ACM.
- [23] C. Rich, B. Ponsleur, A. Holroyd, and C. L. Sidner. Recognizing engagement in human-robot interaction. In *Proceeding of the 5th ACM/IEEE international conference on Human-robot interaction, HRI ’10*, pages 375–382, New York, NY, USA, 2010. ACM.
- [24] D. Sakamoto, T. Kanda, T. Ono, M. Kamashima, M. Imai, and H. Ishiguro. Cooperative embodied communication emerged by interactive humanoid robots. In *Robot and Human Interactive Communication, 2004. ROMAN 2004. 13th IEEE International Workshop on*, pages 443 – 448, sept. 2004.
- [25] M. Shimada, Y. Yoshikawa, M. Asada, N. Saiwaki, and H. Ishiguro. Effects of observing eye contact between a robot and another person. *International Journal of Social Robotics*, 3:143–154, 2011. 10.1007/s12369-010-0072-9.
- [26] C. L. Sidner, C. Lee, C. D. Kidd, N. Lesh, and C. Rich. Explorations in engagement for humans and robots. *Artificial Intelligence*, (166):10–1016, 2005.
- [27] V. Srinivasan, C. Bethel, R. Murphy, and C. Nass. Validation of a behavioral robotics framework for head social gaze. In *HRI ’12: Proceedings of the workshop on Gaze in HRI: From Modeling to Communication*, 2012.
- [28] V. Srinivasan and R. Murphy. A survey of social gaze. In *Proceedings of the 6th international conference on Human-robot interaction, HRI ’11*, pages 253–254, New York, NY, USA, 2011. ACM.
- [29] M. Staudte and M. W. Crocker. Visual attention in spoken human-robot interaction. In *HRI ’09: Proceedings of the 4th ACM/IEEE international conference on Human robot interaction*, pages 77–84, New York, NY, USA, 2009. ACM.
- [30] A. Yamazaki, K. Yamazaki, Y. Kuno, M. Burdelski, M. Kawashima, and H. Kuzuoka. Precision timing in human-robot interaction: coordination of head movement and utterance. In *Proceeding of the twenty-sixth annual SIGCHI conference on Human factors in computing systems, CHI ’08*, pages 131–140, New York, NY, USA, 2008. ACM.