



September 1994

Sight and Sound: Generating Facial Expressions and Spoken Intonation from Context

Scott Prevost
University of Pennsylvania

Catherine Pelachaud
University of Pennsylvania

Follow this and additional works at: <http://repository.upenn.edu/hms>

Recommended Citation

Prevost, S., & Pelachaud, C. (1994). Sight and Sound: Generating Facial Expressions and Spoken Intonation from Context. Retrieved from <http://repository.upenn.edu/hms/40>

Copyright 1994 IEEE. Reprinted from *Proceedings of the second ESCA/AAAI/IEEE Workshop on Speech Synthesis*, 5 pages.

This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the University of Pennsylvania's products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to pubs-permissions@ieee.org. By choosing to view this document, you agree to all provisions of the copyright laws protecting it.

This paper is posted at ScholarlyCommons. <http://repository.upenn.edu/hms/40>
For more information, please contact libraryrepository@pobox.upenn.edu.

Sight and Sound: Generating Facial Expressions and Spoken Intonation from Context

Abstract

This paper presents a model for automatically producing prosodically appropriate speech and corresponding facial expression for agents that respond to simple database queries in a 3D graphical representation of the world. This work addresses two major issues in human-machine interaction. First, proper intonation is necessary for conveying information structure, including important distinctions of contrast and focus. Second, facial expressions and lip movements often provide additional information about discourse structure, turn-taking protocols and speaker attitudes.

Comments

Copyright 1994 IEEE. Reprinted from *Proceedings of the second ESCA/AAAI/IEEE Workshop on Speech Synthesis*, 5 pages.

This material is posted here with permission of the IEEE. Such permission of the IEEE does not in any way imply IEEE endorsement of any of the University of Pennsylvania's products or services. Internal or personal use of this material is permitted. However, permission to reprint/republish this material for advertising or promotional purposes or for creating new collective works for resale or redistribution must be obtained from the IEEE by writing to pubs-permissions@ieee.org. By choosing to view this document, you agree to all provisions of the copyright laws protecting it.

Sight and Sound: Generating Facial Expressions and Spoken Intonation from Context

Catherine Pelachaud and Scott Prevost
Computer and Information Science, University of Pennsylvania
200 South 33rd Street, Philadelphia, PA 19104-6389, USA
Internet contact: prevost@linc.cis.upenn.edu

Preliminary draft accepted for presentation at the Second
ESCA/IEEE Workshop on Speech Synthesis

1 Introduction

This paper presents a model for automatically producing prosodically appropriate speech and corresponding facial expression for agents that respond to simple database queries in a 3D graphical representation of the world. This work addresses two major issues in human-machine interaction. First, proper intonation is necessary for conveying information structure, including important distinctions of contrast and focus. Second, facial expressions and lip movements often provide additional information about discourse structure, turn-taking protocols and speaker attitudes ([7], [8], [14], [15]).

The intonation generation model is based on Combinatory Categorical Grammar (CCG – cf. [20]), a formalism which easily integrates the notions of syntactic constituency, prosodic phrasing and information structure. Based on the CCG grammar, a simple discourse model and a domain-independent knowledge base, the system produces spoken responses to database queries with appropriate intonation. Given the timings for phonemes and intonational phenomena in the speech wave, we produce precise specifications for generating the lip movements and facial expressions for a graphical model of a human head. Results from our current implementation demonstrate the system’s ability to generate a variety of intonational possibilities and facial animations for a given sentence depending on the discourse context.

Previous work in the area of intonation generation includes studies by Terken

- (4) URINALYSIS addresses HEMATURIA
L+H* LH% H* LL\$

The final aspect of speech generation involves translating such a string into a form usable by a suitable speech synthesizer. The current implementation uses the Bell Laboratories TTS system [13] as a post-processor to synthesize the speech wave and produce precise timing specifications for phonemes. The duration specifications are then annotated with pitch accent peaks and intonational boundaries before being sent to the animation system for processing ([3]).

Starting from a functional group (lip shapes, conversational signal, punctuator, regulator or manipulator), we offer algorithms which incorporate synchrony [5], create coarticulation effects, emotional signals, and eye and head movements ([16], [17]). Our rules generate automatically the facial actions corresponding to an input utterance. A conversational signal (movements occurring on accents, like raising of eyebrow) starts and ends with the accented word; while punctuator signals (such as smiling) coincide with pauses. Blinking is synchronized at the phoneme level. Head nods and shakes appear on accent and pause. The head of the speaker turns away from the listener at the beginning of a speaking turn and turns toward to the listener at end of a speaking turn to signal a change of turn.

The computation of the lip shape is done in three passes. Phonemes are characterized by their degree of deformability. For each deformable segment, the program looks for the nearby segment whose associated lip shapes influence it using the look-ahead model; The properties of muscle contractions are taken into account in two ways: spatially, by adjusting the sequence of contracting muscles if antagonist movements (i.e. movements which show very different lip positions like pucker movements versus extension of the lips) succeed each other; and temporally by noticing if a muscle has enough time to contract (respectively relax) before (respectively after) the surrounding lip shape. Those two muscle constraints act on the final computation of the lip shapes.

3 Conclusions

The system described above produces quite sharp and natural-sounding distinctions of intonation contour as well as visually distinct facial animations for minimal pairs of queries. The examples in the full paper illustrate the system's capability for producing appropriately different audio and visual output for a single string of words under the control of differing discourse contexts. We believe the system provides a sound basis for exploring the role of prosody and facial expressions in human-machine interactions, particularly those involving autonomous virtual human agents.

References

- [1] C. Benoit. Why synthesize talking faces? In *Proceedings of the ESCA Workshop on Speech Synthesis*, pages 253–256, Autrans, 1990. ESCA.
- [2] N.M. Brooke. Computer graphics synthesis of talking faces. In *Proceedings of the ESCA Workshop on Speech Synthesis*, Autrans, 1990. ESCA.
- [3] Justine Cassell, Catherine Pelachaud, Norman Badler, Mark Steedman, Brett Achorn, Tripp Becket, Brett Douville, Scott Prevost, Chin Seah, and Matthew Stone. Animated conversation: rule-based generation of facial expression, gesture and spoken intonation for multiple conversational agents. *ACM SIGGRAPH 94*, 1994. submitted.
- [4] Michael M. Cohen and Dominic W. Massaro. Modeling coarticulation in synthetic visual speech. In D. Thalmann N. Magnenat-Thalmann, editor, *Computer Animation '93*. Springer-Verlag, 1993.
- [5] W.S. Condon and W.D. Osgton. Speech and body motion synchrony of the speaker-hearer. In D.H. Horton and J.J. Jenkins, editors, *The perception of Language*, pages 150–184. Academic Press, 1971.
- [6] James Davis and Julia Hirschberg. Assigning intonational features in synthesized spoken discourse. In *Proceedings of the 26th Annual Meeting of the Association for Computational Linguistics*, pages 187–193, Buffalo, 1988.
- [7] S. Duncan. Some signals and rules for taking speaking turns in conversations. In Weitz, editor, *Nonverbal Communication*. Oxford University Press, 1974.
- [8] P. Ekman. About brows: emotional and conversational signals. In M. von Cranach, K. Foppa, W. Lepenies, and D. Ploog, editors, *Human ethology: claims and limits of a new discipline: contributions to the Colloquium*, pages 169–248. Cambridge University Press, Cambridge, England; New-York, 1979.
- [9] D.R. Hill, A. Pearce, and B. Wyvill. Animating speech: an automated approach using speech synthesised by rules. *The Visual Computer*, 3:277–289, 1988.
- [10] Julia Hirschberg. Accent and discourse context: Assigning pitch accent in synthetic speech. In *Proceedings of AAAI: 1990*, page ??, 1990.
- [11] George Houghton. *The Production of Language in Dialogue: a Computational Model*. PhD thesis, University of Sussex, 1986.
- [12] J.P. Lewis and F.I. Parke. Automated lip-synch and speech synthesis for character animation. *CHI + GI*, pages 143–147, 1987.

- [13] Mark Liberman and A. L. Buchsbaum. Structure and usage of current Bell Labs text to speech programs. Technical Memorandum TM 11225-850731-11, AT&T Bell Laboratories, 1985.
- [14] D.W. Massaro. *Speech Perception by Ear and Eye: A Paradigm for Psychological Inquiry*. Cambridge University Press, 1989.
- [15] H. McGurk and J. MacDonald. Hearing lips and seeing voices. *Nature*, 264:746–748, 1976.
- [16] C. Pelachaud, N.I. Badler, and M. Steedman. Linguistic issues in facial animation. In N. Magnenat-Thalmann and D. Thalmann, editors, *Computer Animation '91*, pages 15–30. Springer-Verlag, 1991.
- [17] C. Pelachaud, M.L. Viaud, and H. Yahia. Rule-structured facial animation system. In *IJCAI 93*, 1993.
- [18] Scott Prevost and Mark Steedman. Generating contextually appropriate intonation. In *Proceedings of the 6th Conference of the European Chapter of the Association for Computational Linguistics*, pages 332–340, Utrecht, 1993.
- [19] Scott Prevost and Mark Steedman. Using context to specify intonation in speech synthesis. In *Proceedings of EuroSpeech 93*, Berlin, 1993. To appear.
- [20] Mark Steedman. Structure and intonation. *Language*, pages 260–296, 1991.
- [21] Jacques Terken. The distribution of accents in instructions as a function of discourse structure. *Language and Structure*, 27:269–289, 1984.
- [22] D. Terzopoulos and K. Waters. Techniques for realistic facial modelling and animation. In N. Magnenat-Thalmann and D. Thalmann, editors, *Computer Animation '91*, pages 45–58. Springer-Verlag, 1991.
- [23] R. Zacharski, A.I.C. Monaghan, D.R. Ladd, and J. Delin. BRIDGE: Basic research on intonation in dialogue generation. Technical report, HCRC: University of Edinburgh, 1993. Unpublished manuscript.