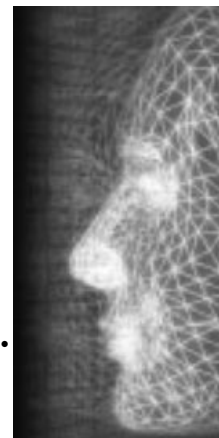


Natural head motion synthesis driven by acoustic prosodic features

By Carlos Busso*, Zhigang Deng, Ulrich Neumann
and Shrikanth Narayanan



Natural head motion is important to realistic facial animation and engaging human–computer interactions. In this paper, we present a novel data-driven approach to synthesize appropriate head motion by sampling from trained hidden markov models (HMMs). First, while an actress recited a corpus specifically designed to elicit various emotions, her 3D head motion was captured and further processed to construct a head motion database that included synchronized speech information. Then, an HMM for each discrete head motion representation (derived directly from data using vector quantization) was created by using acoustic prosodic features derived from speech. Finally, first-order Markov models and interpolation techniques were used to smooth the synthesized sequence. Our comparison experiments and novel synthesis results show that synthesized head motions follow the temporal dynamic behavior of real human subjects. Copyright © 2005 John Wiley & Sons, Ltd.

KEY WORDS: Head motion synthesis; prosody; HMM; facial animation; data-driven; spherical cubic interpolation

Introduction

The development of new human–computer interfaces and exciting applications such as video games and animated feature films has motivated the computer graphics community to generate realistic avatars with the ability to replicate and mirror natural human behavior. Since the use of large motion capture datasets is expensive, and can be only applied to delicately planned scenarios, new automatic systems need to be used to generate natural human facial animation. One

useful and practical approach is to synthesize animated human faces driven by speech.

The straightforward use of speech in facial animation is in lip motion synthesis, in which the acoustic phonemes are used to generate visual visemes that match the spoken sentences. Examples of these approaches include.^{1–7} Also, speech has been used to drive human facial expression, under the assumption that the articulation of the mouth and jaw modify facial muscles, producing different faces poses. Examples of these approaches are.^{8,9} Surprisingly, few efforts have focused on natural generation of rigid head motion, which is an important ingredient for realistic facial animations. In fact, Munhall *et al.*¹⁰ reported that head motion is important for auditory speech perception, which suggests that appropriate head motion can significantly enhance human–computer interfaces.

Although human head motion is associated with many factors, such as speaker style, idiosyncrasies and affective states, linguistic aspects of speech play a crucial role. Kuratate *et al.*⁹ presented preliminary results

*Correspondence to: Carlos Busso, Integrated Media Systems Center, Department of Electrical Engineering, Viterbi School of Engineering, University of Southern California, 3740 McClintock Ave., Room 400, Los Angeles, CA 90089-2564, USA. E-mail: busso@usc.edu

Contract/grant sponsor: NSF; contract/grant number: EEC-9529152.

Contract/grant sponsor: Department of the Army; contract/grant number: DAAD 19-99-D-0046.

about the relationship between head motions and acoustic prosodic features. They concluded based on the strong correlation ($r = 0.8$) that these two are somehow correlated, but perhaps under independent control. This suggests that the tone and the intonation of the speech provide important cues about head motion and vice versa.¹⁰ Notice that, here, it is more important how the speech is uttered rather than just what is said. Therefore, prosodic features (e.g., pitch and energy) are more suitable than vocal tract-based features (e.g., LPC and MFCC). The work of Yehia *et al.*¹¹ even reports that about 80% of the variance observed in the pitch can be determined from head motion.

In this paper, an innovative technique is presented to generate natural head motion directly from acoustic prosodic features. First, vector quantization is used to produce a discrete representation of head poses. Then, a hidden Markov model (HMM) is trained for each cluster, which models the temporal relation between the prosodic features and the head motion sequence. Given that the mapping is not one to one, the observation probability density is modeled with a mixture of Gaussians. The smoothness constraint is imposed by defining a bi-gram model (first order Markov model) on head poses learned from the database. Then, given new speech material, the HMM, working as a sequence generator, produces the most likely head motion sequences. Finally, a smoothing operation based on spherical cubic interpolation is applied to generate the final head motion sequences.

Related Work

Researchers have presented various techniques to model head motion. Pelachaud *et al.*¹² generated head motions from labeled text by predefined rules, based on facial action coding system (FACS) representations.¹³ Cassell *et al.*¹⁴ automatically generated appropriate non-verbal gestures, including head motion, for conversational agents, but their focus was only the 'nod' head motion. Graf *et al.*¹⁵ estimated the conditional probability distribution of major head movements (e.g., nod) given the occurrences of pitch accents, based on their collected head motion data. Costa *et al.*⁸ used *Gaussian mixture model* (GMM) to model the connections between audio features and visual prosody. The connection between eyebrow movements and audio features was specifically studied in their work. Chuang and Bregler¹⁶ presented a data-driven approach to synthesize novel head motion

corresponding to input speech. They first acquired a head motion database indexed by pitch values; then a new head motion sequence was synthesized by choosing and combining best matched recorded head motion segments in the constructed database. Deng *et al.*¹⁷ presented a new audio-driven head motion synthesis technique that synthesized appropriate head motion with keyframing control. After an audio-head motion database was constructed, given novel speech input and user controls (e.g., specified key head poses), a guided dynamic programming technique was used to generate an optimal head motion sequence that maximally satisfies both speech and key frames specifications.

In this paper, we propose to use HMMs to capture the close temporal relation between head motions and acoustic prosodic features. Also, we propose an innovative two-step smoothing technique based on bi-gram models, learned from data, and spherical cubic interpolation.

Data Capture and Processing

Database

The audiovisual database used in this work was recorded from an actress, with 102 markers on her face (left of Figure 1). She was asked to read a custom, phoneme-balanced corpus four times, expressing different emotions (happiness, sadness, anger and neutral state). A VICON motion capture system with three cameras (right of Figure 1) was used to capture her facial expressions and head motions. The sampling frequency was set to 120 Hz. The recording was made in a quiet room using a close talking SHURE microphone at the sampling rate of 48 kHz. The markers' motions and the aligned audio were captured by the system simultaneously. In total, 633 sentences were used in this work. Note that the actress did not receive any instruction about how to move her head.

After the motion data were captured, all the markers were translated to make a nose marker at the local coordinate center of each frame. A neutral pose was chosen as a reference frame, which was used to create a 102×3 matrix, y . For each frame, a matrix x_i was created, using the same marker order as the reference. After that, the singular value decomposition (SVD), UDV^T , of matrix $y^T x_i$ was calculated. Finally, the product of VU^T gave the rotation matrix, R , which defines

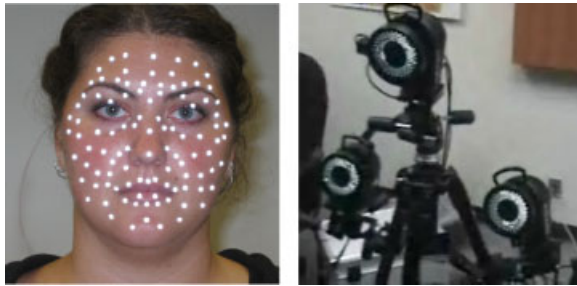


Figure 1. Audio visual database.

the three Euler angles of the head motion of this frame¹⁸ (Figure 2).

$$y^T x_i = UDV^T \quad (1)$$

$$R = VU^T \quad (2)$$

To extract the prosodic features, the acoustic signals were processed by the entropic signal processing system (ESPS), which computes the pitch (F0) and the RMS energy of the audio. The window was set to 25 millisecond with an overlap of 8.3 millisecond. Notice that the pitch takes values only in voiced region of the

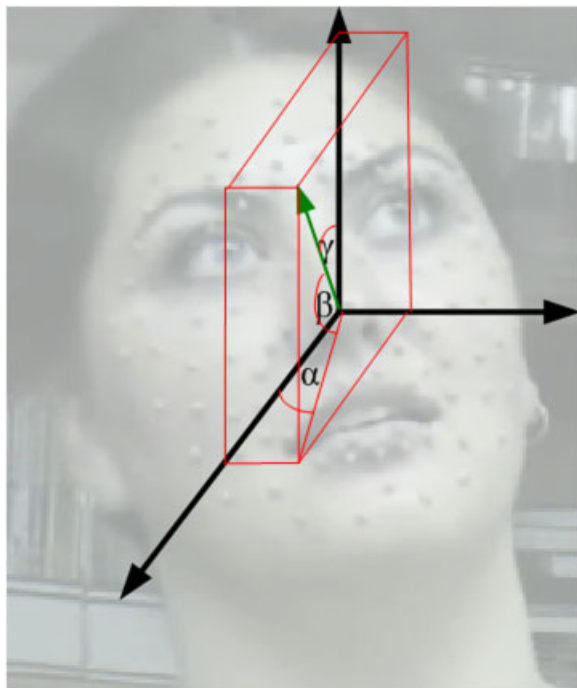


Figure 2. Head poses using Euler angles.

speech. Therefore, to avoid zeros in unvoiced regions, a cubic spline interpolation was applied in those regions. Finally, the first and second derivatives of the pitch and the energy were added to incorporate their temporal dynamics.

Canonical Correlation Analysis

To validate the close relation between head motion and acoustic prosodic features, as suggested by Kuratate *et al.* [9], canonical correlation analysis (CCA)¹⁹ was applied to our audiovisual database. CCA provides a scale-invariant optimum linear framework to measure the correlation between two streams of data with different dimensions. The basic idea is to project both feature vectors into a common dimensional space, in which Pearson's correlation can be computed.

Using pitch, energy and their first and second derivatives (6D feature vector), and the angles that define the head motions (3D feature vector), the average correlation computed from the audiovisual database is $r = 0.7$. This result indicates that useful and meaningful information can be extracted from the prosodic features of speech to synthesize the head motion.

Modeling Head Motion

In this paper, we use HMMs because they provide a suitable and natural framework to model the temporal relation between acoustic prosodic features and head motions. HMMs are used to generate the most likely head motion sequences based on the given observation (prosodic features). The HTK toolkit is used to build the HMMs.²⁰

The output sequences of the HMMs cannot be continuous, so a discrete representation of head motion is needed. For this purpose, the Linde-Buzo-Gray vector quantization (LBG-VQ) algorithm²¹ is used to define K discrete head poses, V_i . The 3D-space defined by the Euler angles is split into K Voronoi regions (Figure 3). For each region, the mean vector U_i and the covariance matrices Σ_i are estimated. The pairs (U, Σ) define the finite and discrete set of code vectors called codebook. In the quantization step, the continuous Euler angles of each frame are approximated with the closest code vector in the codebook.

For each of the clusters, V_i , an HMM model will be created. Consequently, the size of the codebook will determine the number of HMM models.

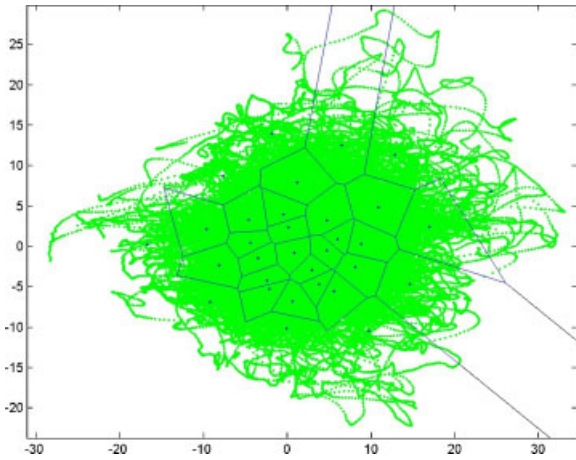


Figure 3. 2D projection of Voronoi regions using 32-size vector quantization.

Learning Natural Head Motion

The posterior probability of being in cluster V_i , given the observation O , is modeled according to Bayes rule as:

$$P(V_i/O) = c \times P(O/V_i) \times P(V_i) \quad (3)$$

where c is a normalization constant. The likelihood distribution, $P(O/V_i)$, is modeled as a Markov process, which is a finite state machine that changes state at every time unit according to the transition probabilities. A first order Markov model is assumed, in which the probabilistic description includes only the current and previous state. The probability density function of the observation is modeled by a mixture of M Gaussian densities which handle, up to some extent, the many-to-many mapping between head motion and prosodic features. Standard algorithms (forward-backward, Baum-Welch re-estimation) are used to train the parameters of the HMMs, using the training data.^{20,22} Notice that the segmentation of the speech according to the head poses clusters is known. Therefore, the HMMs were initialized with this known alignment (force alignment was not needed).

The prior distribution, $P(V_i)$, is used to impose a first smoothing constraint to avoid sudden changes in the synthesized head motion sequence. In this paper, $P(V_i)$ is built using bi-gram models, which are learned from data (similar to standard bi-gram language models used to model word sequence probabilities^{20,23}). The bi-gram model is also a first-order state machine, in which each state models the probability of observing a given output

sequence (in this case, a specific head pose cluster, V_i). The transition probabilities are computed using the frequency of their occurrences. In our training database, the inter-cluster transitions are counted and stored, and the statistic learned is used to reward transitions according to their appearances. Therefore, in the decoding step, this prior distribution will penalize transitions that did not appear in the training data.

Synthesis of Head Motion

Figure 4 describes the procedure to synthesize head motion. For each testing sample, the acoustic prosodic features are extracted and used as input of the HMMs. The model will generate the most likely sequence, $\hat{\mathbf{V}} = (\hat{V}_i^t, \hat{V}_j^{t+1} \dots)$, where \hat{V}_i^t is defined by (\mathbf{U}_i, Σ_i) . The mean vector \mathbf{U}_i will be used to synthesize the head motion sequences.

The transitions between clusters will introduce breaks in the synthesized signal, even if their cluster means are close (see Figure 5). Therefore, a second smoothing step needs to be implemented, to guarantee continuity of the synthesized headpose sequences. A simple solution is to interpolate each Euler angle separately. However, it has been shown that this technique is not optimal, because it introduces jerky movements and other undesired effects such as Gimbal lock.²⁴ As suggested by Shoemake, a better approach is to interpolate in the quaternion unit sphere.²⁴ The basic idea is to transform the Euler angles into quaternions, which are an alternative rotation matrix representation, and then interpolate the frames in this space.

In this paper, we used spherical cubic interpolation,²⁵ squad, which is based on spherical linear interpolation, slerp. For two quaternions q_1 and q_2 , the slerp function is defined as:

$$\text{slerp}(q_1, q_2, \mu) = \frac{\sin(1 - \mu)\theta}{\sin \theta} q_1 + \frac{\sin \mu\theta}{\sin \theta} q_2 \quad (4)$$

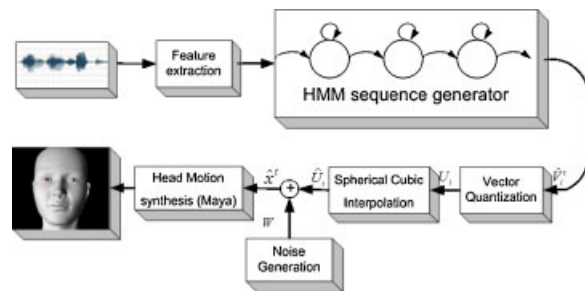


Figure 4. Head motion synthesis framework.

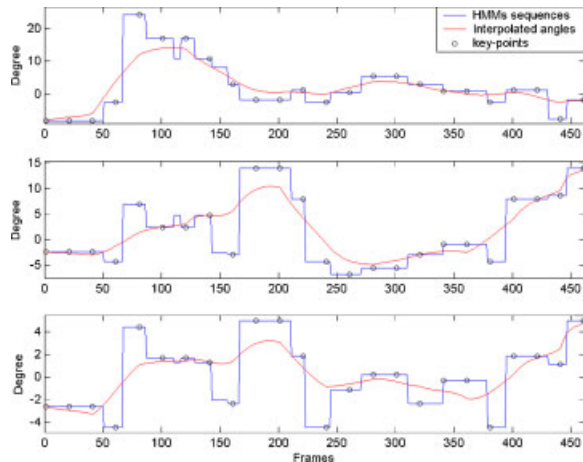


Figure 5. Spherical cubic interpolation.

where $\cos\theta = q_1 \times q_2$ and μ are variables that range from 0 to 1 and determine the frame position of the interpolated quaternion. Given four quaternions, the squad function is defined as:

$$\text{squad}(q_1, q_2, q_3, q_4, \mu) = \text{slerp}(\text{slerp}(q_1, q_4, \mu), \text{slerp}(q_2, q_3, \mu), 2\mu(1 - \mu)) \quad (5)$$

After the Euler angles are transformed into quaternions, key-points are selected by down-sampling the quaternions at a rate of 6 frames per second (this value was empirically chosen). Then, spherical cubic interpolation is used in those key-points by using the squad function. After interpolation, the frame rate of the quaternions is 120 frames per second, as the original data. The last step in this smoothing technique is to transform the interpolated quaternions into Euler angles. Figure 5 shows the interpolation result for one of the sentences. The resulting vectors are denoted $\hat{\mathbf{U}}_i$. The readers are referred to²⁵ for further details about spherical cubic interpolation.

Finally, the synthesized head pose, $\hat{\mathbf{x}}^t$, at time t will be estimated as:

$$\hat{\mathbf{x}}^t = (\alpha, \beta, \gamma)^T = \hat{\mathbf{U}}_i + W_i \quad (6)$$

where W_i is a zero-mean uniformly distributed random white noise. Notice that $\hat{\mathbf{x}}^t$ is a blurred version of V_i 's mean.

If the size of the codebook is large enough, the quantization error will be insignificant. However, the number of HMMs needed will increase and the discrimination

between classes will decrease. Also, more data will be needed to train the models. Therefore, there is a tradeoff between the quantization error and the inter-cluster discrimination.

From Euler Angles to Talking Avatars

A blend shape face model composed of 46 blend shapes is used in this work (eyeball is controlled separately). To create a realistic avatar, lip and eye motion techniques are also included. For novel text/audio input, the speech animation approach presented in References [6,26] was used to generate synchronized visual speech. Then, eye motion is automatically synthesized by a texture-synthesis based approach.²⁷ Hence, appropriate blend shape weights are calculated for each frame. After that, the same audio is used to generate corresponding natural head motion using the proposed approach. The generated head motion Euler angles, $\hat{\mathbf{x}}^t$, are directly applied to the angle control parameters of the face model. This animation procedure is applied to each frame. Besides the animation, the face modeling and rendering are also done in Maya.

Results and Discussion

HMM Configuration

The topology of the HMM is defined by the number and the interconnection of the states. The most popular configurations are the left-to-right topology (LR), in which only transitions in forward direction between adjacent states are allowed, and the ergodic topology (EG), in which transitions between all the states are allowed. The LR topology is simple and needs less data to train its parameters. The EG topology is less restricted, so it can learn a larger set of state transitions from the data. In this particular problem, it is not clear which topology gives better description of the head motion dynamics. Therefore, eight HMM configurations, described in Table 1, with different topologies, number of models, K , number of states, S , and number of mixtures, M , were trained. Notice that the size of the database is not big enough to train more complex HMMs with more states, mixtures or models than those described in Table 1.

HMM configurations	\bar{D}		CCA	
	Mean	Std	Mean	Std
K = 16 S = 5 M = 2 LR	10.2	3.4	0.88	0.11
K = 16 S = 5 M = 4 LR	9.3	3.4	0.87	0.11
K = 16 S = 3 M = 2 LR	9.1	3.6	0.87	0.12
K = 16 S = 3 M = 2 EG	9.1	3.4	0.87	0.10
K = 16 S = 3 M = 4 EG	9.5	3.4	0.83	0.12
K = 32 S = 5 M = 1 LR	12.8	4.0	0.83	0.14
K = 32 S = 3 M = 2 LR	10.7	3.3	0.86	0.12
K = 32 S = 3 M = 1 EG	10.4	3.1	0.86	0.11

Table 1. Results for different configurations

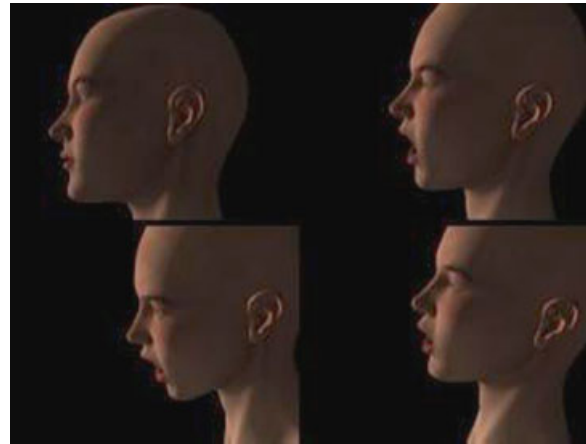
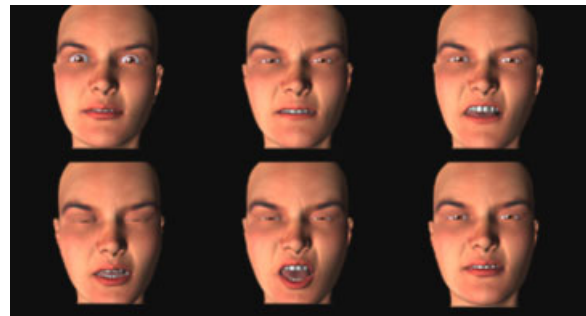
As mentioned before, the pitch, the energy, and their first and second derivatives are used as acoustic prosodic features to train each of the proposed HMM configuration. Eighty per cent of the database was used for training and twenty per cent for testing.

Objectives Evaluation

To evaluate the performance of our approach, the prosodic features from the test data were used to generate head motion sequences, as described in previous section. For those samples, the Euclidian distance, d_{euc} , between the Euler angles of the original frames and the Euler angles of the synthesized data, \hat{x}^f , was calculated. The average and the standard deviation of all the frames of the testing data, is shown in Table 1 (\bar{D}). Notice that the synthesized head motions are directly compared with the original data (not its quantized version), so the quantization error is included in the values of the table. As mentioned in the Introduction, head motion does not depend only on prosodic features, so this level of mismatch is expected.

Table 1 also shows the CCA between the synthesized and original data. As can be observed the correlation was around $r=0.85$ for all the topologies. This result strongly suggests that the synthesized data follow the behavior of the real data, which validates our approach.

As can be seen from Table 1, the performances of the different HMM topologies are similar. The left-to-right HMM, with a 16-size codebook, 3 states, and 2 mixtures, achieves the best result. However, if the database were big enough, an EG with more states and mixture could perhaps give better results. The next experiments were implemented using this topology ($K=16$, $S=3$, $M=2$ LR).

*Figure 6. Synthesized data, side view.**Figure 7. Synthesized data, front view.*

Head Motion Animation Results

We also recorded sentences, not included in the corpus mentioned before, to synthesize novel head motion using our approach. For each recorded audio, the procedure described in this paper was applied. Figures 6 and 7 show frames of the synthesized data. For animation results, please refer to the accompanying video.

Conclusions

This paper presents a novel approach to synthesize natural human head motions driven by speech prosody. HMMs are used to capture the temporal relation between the acoustic prosodic features and head motions. The use of bi-gram models in the sequence generation step guarantees smooth transitions from the discrete representations of head movement configurations. Furthermore, spherical cubic interpolation is used to avoid breaks in the synthesized signal.

The results show that the synthesized sequences follow the temporal dynamic behavior of real data. This proves that the HMMs are able to capture the close

relation between speech and head motion. The results also show that the smoothing techniques used in this work can produce continuous head motion sequences, even when only a 16-word-sized codebook is used to represent head motion poses.

In this paper we show that natural head motion animation can be synthesized by using just speech. In future work, we will add more components to our system. For example, if the emotion of the subject is known, as is usually the case in most of the applications, suitable models that capture the emotional head motion pattern can be used, instead of a general model.

ACKNOWLEDGEMENTS

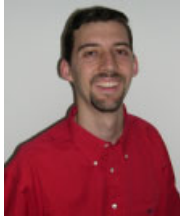
This research was supported in part by funds from the NSF through the Integrated Media Systems Center, a National Science Foundation Engineering Research Center, Cooperative Agreement under Contract number EEC-9529152, a CAREER award (to S. Narayanan), and the Department of the Army under Contract number DAAD 19-99-D-0046. Special thanks go to J.P. Lewis and Murtaza Bulut for helping data capture, Hiroki Itokazu, Bret St. Clair, and Pamela Fox for face model preparation, and anonymous reviewers for useful comments and suggestions.

References

1. Cohen MM, Massaro DW. Modeling coarticulation in synthetic visual speech. *Models and Techniques in Computer Animation*: Magnenat-Thalmann N, Thalmann D (eds). Springer Verlag: Tokyo, Japan, 1993; 139–156.
2. Bregler C, Covell M, Slaney M. Video rewrite: driving visual speech with audio. *Proceedings of ACM SIGGRAPH'97*, 1997; pp. 353–360.
3. Brand M. Voice puppetry. *Proceedings of ACM SIGGRAPH'99*, 1999; pp. 21–28.
4. Kshirsagar S, Thalmann NM. Visyllable based speech animation. *Computer Graphics Forum (Proceedings of Eurographics'03)* 2003; **22**(3): 631–939.
5. Ezzat T, Geiger G, Poggio T. Trainable videorealistic speech animation. *ACM Transaction on Graphics (Proceedings of ACM SIGGRAPH'02)* 2002; 388–398.
6. Deng Z, Lewis JP, Neumann U. Synthesizing speech animation by learning compact speech co-articulation models. In *Computer Graphics International (CGI 2005)*. Stony Brook: NY, USA, June 2005; pp. 19–25.
7. Liu W, Yin B, Jia X. Audio to visual signal mappings with hmm. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP' 4)*. Quebec, Canada, May 2004; pp. 885–888.
8. Costa M, Chen T, Lavagetto F. Visual prosody analysis for realistic motion synthesis of 3d head models. In *International Conference On Augmented, Virtual Environments and Three Dimensional Imaging, ICAV3D*, Ornos, Mykonos, Greece, May–June 2001.
9. Kuratate T, Munhall KG, Rubin PE, Bateson EV, Yehia H. Audio-visual synthesis of talking faces from speech production correlates. In *Sixth European Conference on Speech Communication and Technology, Eurospeech 1999*. Budapest, Hungary, September 1999; pp. 1279–1282.
10. Munhall KG, Jones JA, Callan DE, Kuratate T, Bateson EV. Visual prosody and speech intelligibility: head movement improves auditory speech perception. *Psychological Science* 2004; **15**(2): 133–137.
11. Yehia H, Kuratate T, Bateson EV. Facial animation and head motion driven by speech acoustics. In *5th Seminar on Speech Production: Models and Data*, 2000; pp. 265–268.
12. Pelachaud C, Badler N, Steedman M. Generating facial expressions for speech. *Cognitive Science* 1996; **20**(1): 1–46.
13. Ekman P, Friesen WV. *Unmasking the Face: A Guide to Recognizing Emotions from Facial Clues*. Prentice-Hall: Englewood Cliffs, New Jersey, USA, 1975.
14. Cassell J, Pelachaud C, Badler N, Steedman M, Achorn B, Bechet T, Douville B, Prevost S, Stone M. Animated conversation: ruled-based generation of facial expression gesture and spoken intonation for multiple conversational agents. In *Computer Graphics (Proceedings of ACM SIGGRAPH'94)*, 1994; pp. 413–420.
15. Graf HP, Cosatto E, Strom V, Huang FJ. Visual prosody: facial movements accompanying speech. In *Proceedings of IEEE International Conference on Automatic Faces and Gesture Recognition*, 2002.
16. Chuang E, Bregler C. Head emotion. *Stanford University Computer Science Technical Report, CS-TR-2003-02*, 2003.
17. Deng Z, Busso C, Narayanan S, Neumann U. Audio-based head motion synthesis for avatar-based telepresence systems. In *ACM SIGMM 2004 Workshop on Effective Telepresence (ETP 2004)*. ACM Press: New York, NY, 2004; 24–30.
18. Stegmann MB, Gomez DD. A brief introduction to statistical shape analysis, Informatics and Mathematical Modeling, Technical University of Denmark, 2002.
19. Dehon C, Filzmoser P, Croux C. Robust methods for canonical correlation analysis. In *Data Analysis, Classification, and Related Methods*. Springer-Verlag: Berlin, 2000; pp. 321–326.
20. Young S, Evermann G, Hain T, et al. *The HTK Book*. Entropic Cambridge Research Laboratory: Cambridge, England, 2002.
21. Linde Y, Buzo A, Gray R. An algorithm for vector quantizer design. *IEEE Transactions on Communications* 1980; **28**(1): 84–95.
22. Rabiner LR. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE* 1989; **77**(2): 257–286.
23. Huang X, Acero A, Hon H-W. *Spoken Language Processing: A guide to theory, Algorithm and System Development*. Prentice-Hall: Upper Saddle River, New Jersey, USA, 2001.
24. Shoemake K. Animating rotation with quaternion curves. *Computer Graphics (Proceedings of SIGGRAPH85)* 1985; **19**(3): 245–254.
25. Eberly D. *3D Game Engine Design: A Practical Approach to Real-Time Computer Graphics*. Morgan Kaufmann Publishers: San Francisco, California, USA, 2000.
26. Deng Z, Bulut M, Neumann U, Narayanan S. Automatic dynamic expression synthesis for speech animation. In *IEEE 17th International Conference on Computer Animation and Social Agents (CASA 2004)*, Geneva, Switzerland, 2004; pp. 267–274.

27. Deng Z, Lewis PJ, Neumann U. Automated eye motion using texture synthesis. *IEEE Computer Graphics and Applications* 2005; 25(2): 24–30.

Authors' biographies:



Carlos Busso received the B.Sc. (2000) and M.S. (2003) degrees with high honors in Electrical Engineering from University of Chile, Santiago, Chile. He is currently pursuing his Ph.D. in Electrical Engineering at the University of Southern California (USC), Los Angeles, U.S.A. He is a student member of IEEE since 2001. From 2003, he is a student member of the Speech Analysis and Interpretation Laboratory (SAIL) at USC. His research interests are in digital signal processing, speech and video signal processing, and multimodal interfaces. His currently researches include modeling and understanding human communication and interaction, with application in recognition and synthesis.



Zhigang Deng is a Ph.D. candidate in the Computer Graphics and Immersive Technologies Lab of Integrated Media System Center and Computer Science Department at University of Southern California. His research interests include Computer Graphics, Computer Animation, Statistical Learning for Animation, Human Computer Interaction, and Information Visualization. Deng received a B.S. degree in mathematics from Xiamen University, and M.S. degree in Computer Science from Peking University and USC, respectively. He is a member of IEEE Computer Society, ACM, and ACM SIGGRAPH. (E-mail: zdeng@graphics.usc.edu).



Ulrich Neumann is an associate professor of Computer Science, with a joint appointment in Electrical Engineer-

ing, at the University of Southern California. He completed an MSEE from SUNY at Buffalo in 1980 and his Computer Science Ph.D. at the University of North Carolina at Chapel Hill in 1993, where his focus was on parallel algorithms for interactive volume-visualization. His current research relates to immersive environments and virtual humans. He won an NSF CAREER award in 1995 and the Junior Faculty Research award at USC in 1999. Dr Neumann held the Charles Lee Powell Chair of Computer Science and Electrical Engineering and was the director of the Integrated Media Systems Center (IMSC), an NSF Engineering Research Center (ERC) from 2000 to 2004. He directs the Computer Graphics and Immersive Technologies (CGIT) Laboratory at USC. In his commercial career, he designed multiprocessor graphics and DSP systems, cofounded a video game corporation, and independently developed and licensed electronic products.



Shrikanth Narayanan (Ph.D.'95, UCLA), was with AT&T Research (originally AT&T Bell Labs) first as a senior member, and later as a principal member, of its technical staff from 1995 to 2000. Currently he is an associate professor of Electrical Engineering, Linguistics and Computer Science at the University of Southern California (USC). He is a member of the Signal and Image Processing Institute and a research area director of the Integrated Media Systems Center, an NSF Engineering Research Center, at USC. He was an associate editor of the *IEEE Transactions of Speech and Audio Processing* (2000–2004) and serves on the Speech Processing and Multimedia Signal Processing technical committee of the IEEE Signal Processing Society and the Speech Communication committee of the Acoustical Society of America. Shri Narayanan is a member of Tau-Beta-Pi, Phi Kappa Phi and Eta-Kappa-Nu, a senior member of IEEE and a recipient of an NSF CAREER award, USC Engineering Junior Research Award, USC Electrical Engineering Northrop Grumman Research Award and a faculty fellowship from the USC Center for Interdisciplinary research. His research interests are in signals and systems modeling with applications to speech, language, multimodal and biomedical problems. He has published over 140 papers and holds three U.S. patents.