# Effects of damping head movement and facial expression in dyadic conversation using real–time facial expression tracking and synthesized avatars

**Steven M. Boker**[1], **Jeffrey F. Cohn**[2,3,*], **Barry-John Theobald**[4],
**Iain Matthews**[3,5], **Timothy R. Brick**[1] and **Jeffrey R. Spies**[1]

[1]*Department of Psychology, Gilmer Hall Room 102, University of Virginia, Charlottesville, VA 22903, USA*
[2]*Department of Psychology, 3137 Sennott Square, 210 S. Bouquet Street, Pittsburgh, PA 15260, USA*
[3]*Robotics Institute, Carnegie Mellon University, 5000 Fifth Avenue, Pittsburgh, PA 15213, USA*
[4]*School of Computing Sciences, University of East Anglia, Earlham Road, Norwich NR4 7TJ, UK*
[5]*Disney Research, Pittsburgh, 4615 Forbes Avenue, Pittsburgh, PA 15213, USA*

When people speak with one another, they tend to adapt their head movements and facial expressions in response to each others' head movements and facial expressions. We present an experiment in which confederates' head movements and facial expressions were motion tracked during videoconference conversations, an avatar face was reconstructed in real time, and naive participants spoke with the avatar face. No naive participant guessed that the computer generated face was not video. Confederates' facial expressions, vocal inflections and head movements were attenuated at 1 min intervals in a fully crossed experimental design. Attenuated head movements led to increased head nods and lateral head turns, and attenuated facial expressions led to increased head nodding in both naive participants and confederates. Together, these results are consistent with a hypothesis that the dynamics of head movements in dyadicconversation include a shared equilibrium. Although both conversational partners were blind to the manipulation, when apparent head movement of one conversant was attenuated, both partners responded by increasing the velocity of their head movements.

**Keywords:** facial expression; dynamics; symmetry

## 1. INTRODUCTION

When people converse, they adapt their movements, facial expressions and vocal cadence to one another. This multi-modal adaptation allows the communication of information that either reinforces or is in addition to the information that is contained in the semantic verbal stream. For instance, back-channel information such as direction of gaze, head nods and 'uh-huh's allow the conversants to better segment speaker–listener turn taking. Affective displays such as smiles, frowns, expressions of puzzlement or surprise, shoulder movements, head nods and gaze shifts are components of the multi-modal conversational dialogue.

When two people adopt similar poses, this could be considered a form of spatial symmetry (Boker & Rotondo 2002). Interpersonal symmetry has been reported in many contexts and across sensory modalities: for instance, patterns of speech (Cappella & Panalp 1981, Neumann & Strack 2000), facial expression (Hsee *et al*. 1990) and laughter (Young & Frye 1966). Increased symmetry is associated with increased rapport and affinity between conversants

(LaFrance 1982; Bernieri 1988). Intrapersonal and cross-modal symmetry may also be expressed. Smile intensity is correlated with cheek raising in smiles of enjoyment (Messinger *et al*. 2009) and with head pitch and yaw in embarrassment (Cohn *et al*. 2004; Ambadar *et al*. 2009). The structure of intrapersonal symmetry may be complex: self-affine multi-fractal dimension in head movements change based on conversational context (Ashenfelter *et al*. 2009).

Symmetry in movements implies redundancy in movements, which can be defined as negative Shannon information (Shannon & Weaver 1949; Redlich 1993). As symmetry is formed between conversants, the ability to predict the actions of one based on the actions of the other increases. When symmetry is broken by one conversant, the other is likely to be surprised or experience change in attention. The conversant's previously good predictions would now be much less accurate. Breaking symmetry may be a method for increasing the transmission of non-verbal information by reducing the redundancy in a conversation.

This view of an ever-evolving symmetry between two conversants may be conceptualized as a dynamical system with feedback as shown in figure 1. Motor activity (e.g. gestures, facial expression or speech) is produced by one conversant and perceived by the other. These perceptions contribute to some system that functions to map the perceived actions of the

* Author and address for correspondence: Department of Psychology, 3137 SQ, 210S. Bouquet Street, Pittsburg, PA 15260, USA (jeffcohn@cs.cmu.edu).
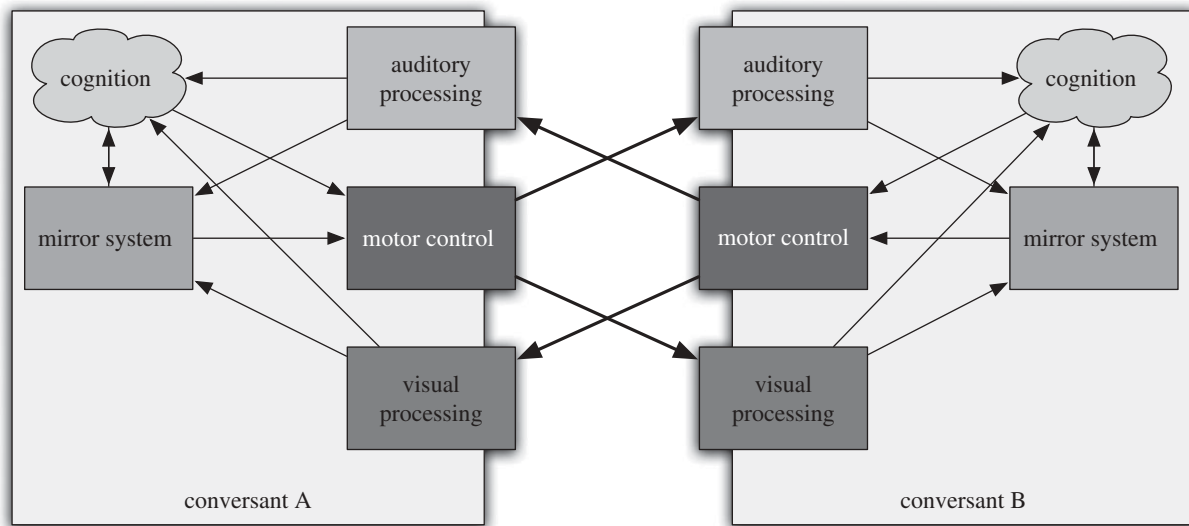
Figure 1. Dyadic conversation involves a dynamical system with adaptive feedback control resulting in complex, non-stationary behaviour.

interlocutor onto potential action: a mirror system. Possible neurological candidates for such a mirror system have been advanced by Rizzolati and colleagues (Iacoboni *et al*. 1999; Rizzolatti & Craighero 2004; Rizzolatti & Fadiga 2007) who argued that such a system is fundamental to communication.

Conversational movements are likely to be non-stationary (Boker *et al*. 2002; Ashenfelter *et al*. in press) and involve both symmetry formation and symmetry breaking (Boker & Rotondo 2002). One technique that is used in the study of non-stationary dynamical systems is to induce a known perturbation into a free running system and measure how the system adapts to the perturbation. In the case of facial expressions and head movements, one would need to manipulate conversant A's perceptions of the facial expressions and head movements of conversant B, while conversant B remained blind to these manipulations as illustrated in figure 2.

Recent advances in active appearance models (AAMs) (Cootes *et al*. 2002) have allowed the tracking and re-synthesis of faces in real-time (Matthews & Baker 2004). Placing two conversants into a videoconference setting provides a context in which a real-time AAM can be applied, since each conversant is facing a video camera and each conversant only sees a video image of the other person. One conversant could be tracked and the desired manipulations of head movements and facial expressions could be applied prior to re-synthesizing an avatar that would be shown to the other conversant. In this way, a perturbation could be introduced as shown in figure 2.

To test the feasibility of this paradigm and to investigate the dynamics of symmetry formation and breaking, we present the results of an experiment in which we implemented a mechanism for manipulating head movement and facial expression in real time during a face-to-face conversation using a computer-enhanced videoconference system. The experimental manipulation was not noticed by naive participants, who were informed that they would be in a videoconference and that we had 'cut out' the face of the person

with whom they were speaking. No participant guessed that he or she was actually speaking with a synthesized avatar. This manipulation revealed the co-regulation of symmetry formation and breaking in two-person conversations.

## 2. MATERIAL AND METHODS
### (a) *Apparatus*
Videoconference booths were constructed in two adjacent rooms. Each $1.5\,m \times 1.2\,m$ footprint booth consisted of a $1.5\,m \times 1.2\,m$ backprojection screen, two $1.2\,m \times 2.4\,m$ non-ferrous side walls covered with white fabric, and a white fabric ceiling. Each participant sat on a stool approximately 1.1 m from the backprojection screen as shown in figure 3. Audio was recorded using Earthworks directional microphones through a Yamaha 01V96 multi-channel digital audio mixer. National Television System Committee format video was captured using Panasonic IK-M44H 'lipstick' colour video cameras and recorded to two JVC BR-DV600U digital video decks. Society of Motion Picture and Television Engineers time stamps generated by an ESE 185-U master clock were used to maintain a synchronized record on the two video recorders and to synchronize the data from a magnetic motion capture device. Head movements were tracked and recorded using an Ascension Technologies MotionStar magnetic motion tracker sampling at 81.6 Hz from a sensor attached to the back of the head using an elastic headband. Each room had an extended range transmitter whose fields overlapped through the non-ferrous wall separating the two video booth rooms.

To track and re-synthesize the avatar, video was captured by an AJA Kona card in an Apple 2-core 2.5 GHz G5 PowerMac with 3 Gb of RAM and 160 Gb of storage. The PowerMac ran software described below and output the resulting video frames to an InFocus IN34 DLP Projector. Thus, the total delay time from the camera in booth 1 through the avatar synthesis process and projected to
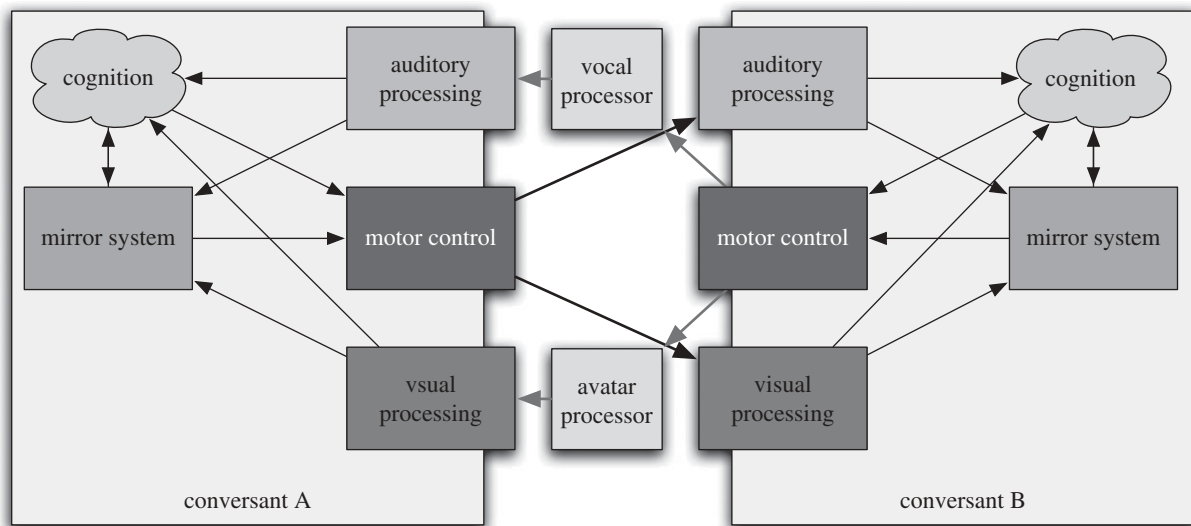
Figure 2. By tracking rigid and non-rigid head movements in real time and re-synthesizing an avatar face, controlled pertubations can be introduced into the shared dynamical system between two conversants.
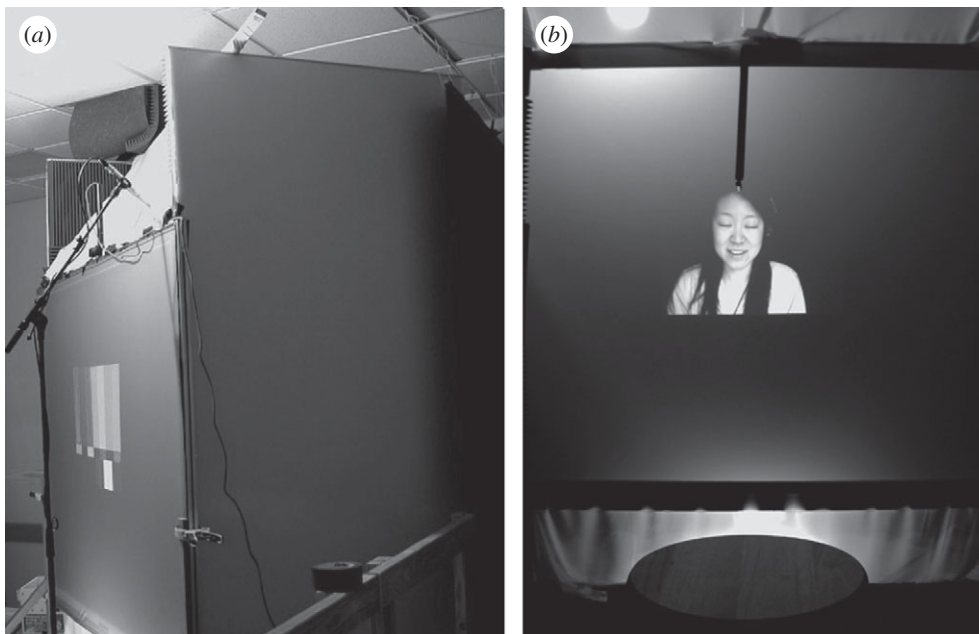


Figure 3. Videoconference booth. (*a*) Exterior of booth showing backprojection screen, side walls, fabric ceiling, and microphone. (*b*) Interior of booth from just behind participant's stool showing projected video image and lipstick videocamera.

booth 2 was 165 ms. The total delay time from the camera in booth 2 to the projector in booth 1 was 66 ms, because the video signal was passed directly from booth 2 to booth 1 and did not need to go through a video analogue/digital and avatar synthesis. For the audio manipulations described below, we reduced vocal pitch inflection using a TC-Electronics VoiceOne Pro. Audio–video synchronization was maintained using digital delay lines built into the Yamaha 01V96 mixer.

### (b) *Active appearance models*

AAMs (Cootes *et al.* 2001) are generative, parametric models commonly used to track and synthesize faces in video sequences. Recent improvements in both the fitting algorithms and the hardware on which they

run allow tracking (Matthews & Baker 2004) and synthesis (Theobald *et al.* 2007) of faces in real-time.

The AAM is formed of two compact models: one describes variation in shape and the other variation in appearance. AAMs are typically constructed by first defining the topological structure of the shape (the number of landmarks and their interconnectivity to form a two-dimensional triangulated mesh), then annotating with this mesh a collection of images that exhibit the characteristic forms of the variation of interest. For this experiment, we label a subset of 40–50 images (less than 0.2% of the images in a single session) that are representative of the variability in facial expression. An individual shape is formed by concatenating the coordinates of the corresponding mesh vertices, $\mathbf{s} = (x_1, y_1, \ldots, x_n, y_n)^{\mathrm{T}}$, so the

collection of training shapes can be represented in matrix form as $\mathbf{S} = [\mathbf{s}_1, \mathbf{s}_2, \ldots, \mathbf{s}_N]$. Applying principal component analysis (PCA) to these shapes, typically aligned to remove in-plane pose variation, provides a compact model of the form

$$\mathbf{s} = \mathbf{s}_0 + \sum_{i=1}^{m} \mathbf{s}_i p_i, \qquad (2.1)$$

where $\mathbf{s}_0$ is the mean shape and the vectors $\mathbf{s}_i$ are the eigenvectors corresponding to the $m$ largest eigenvalues. These eigenvectors are the basis vectors that span the shape space and describe variation in the shape about the mean. The coefficients $p_i$ are the shape parameters, which define the contribution of each basis in the reconstruction of $\mathbf{s}$. An alternative interpretation is that the shape parameters are the coordinates of $\mathbf{s}$ in shape space, thus each coefficient is a measure of the distance from $\mathbf{s}_0$ to $\mathbf{s}$ along the corresponding basis vector.

The appearance of the AAM is a description of the variation estimated from a shape-free representation of the training images. Each training image is first warped from the manually annotated mesh location to the base shape, so the appearance comprises the pixels that lie inside the base mesh, $\mathbf{x} = (x, y)^{\mathrm{T}} \in \mathbf{s}_0$. PCA is applied to these images to provide a compact model of appearance variation of the form

$$A(\mathbf{x}) = A_0(\mathbf{x}) + \sum_{i=1}^{l} \lambda_i A_i(\mathbf{x}), \quad \forall\, \mathbf{x} \in \mathbf{s}_0, \qquad (2.2)$$

where the coefficients $\lambda_i$ are the appearance parameters, $A_0$ is the base appearance and the appearance images, $A_i$, are the eigenvectors corresponding to the $l$ largest eigenvalues. As with shape, the eigenvectors are the basis vectors that span appearance space and describe variation in the appearance about the mean. The coefficients $\lambda_i$ are the appearance parameters, which define the contribution of each basis in the reconstruction of $A(\mathbf{x})$. Because the model is invertible, it may be used to synthesize new face images (see figure 4).

### (c) Manipulating facial displays using AAMs
To manipulate the head movement and facial expression of a person during a face-to-face conversation such that they remain blind to the manipulation, an avatar is placed in the feedback loop, as shown in figure 2. Conversants speak via a videoconference and an AAM is used to track and parameterize the face of one conversant.

As outlined, the parameters of the AAM represent displacements from the origin in the shape and appearance space. Thus, scaling the parameters has the effect of either exaggerating or attenuating the overall facial expression encoded as AAM parameters

$$\mathbf{s} = \mathbf{s}_0 + \sum_{i=1}^{m} \mathbf{s}_i p_i \beta, \qquad (2.3)$$

where $\beta$ is a scalar, which when greater than unity exaggerates the expression and when less than unity attenuates the expression. An advantage of

using an AAM to conduct this manipulation is that a separate scaling can be applied to the shape and appearance to create some desired effect. We stress here that in these experiments, we are not interested in manipulating individual actions on the face (e.g. inducing an eye-brow raise), rather we wish to manipulate, in real time, the overall facial expression produced by one conversant during the conversation.

The second conversant does not see the video of the person to whom they are speaking. Rather, they see a re-rendering of the video from the manipulated AAM parameters as shown in figure 5. To re-render the video using the AAM, the shape parameters $\mathbf{p} = (p_1, \ldots, p_m)^{\mathrm{T}}$, are first applied to the model, equation (2.3), to generate the shape, $\mathbf{s}$, of the AAM, followed by the appearance parameters $\lambda = (\lambda_1, \ldots, \lambda_l)^{\mathrm{T}}$ to generate the AAM image, $A(\mathbf{x})$. Finally, a piece-wise affine warp is used to warp $A(\mathbf{x})$ from $\mathbf{s}_0$ to $\mathbf{s}$, and the result is transferred into image coordinates using a similarity transform (i.e. movement in the $x$–$y$ plane, rotation and scale). This can be achieved efficiently, at video frame rate, using standard graphics hardware.

Typical example video frames synthesized using an AAM before and after damping are shown in figure 6. Note that the effect of the damping is to reduce the expressiveness. Our interest here is to estimate the extent to which manipulating expressiveness in this way can affect the behaviour during conversation.

### (d) Participants
Naive participants ($n = 27$, 15 male, 12 female) were recruited from the psychology department participant pool at a midwestern university. Confederates ($n = 6$, three male and three female) were undergraduate research assistants. AAM models were trained for the confederates so that the confederates could act as one conversant in the dyad. Confederates were informed of the purpose of the experiment and the nature of the manipulations, but were blind to the order and timing of the manipulations. All confederates and naive participants read and signed informed consent forms approved by the Institutional Review Board.

### (e) Procedure
We attenuated three variables: (i) head pitch and turn: translation and rotation in image coordinates from their canonical values by either 1.0 or 0.5; (ii) facial expression: the vector distance of the AAM shape parameters from the canonical expression (by multiplying the AAM shape parameters by either 1.0 or 0.5); and (iii) audio: the range of frequency variability in the fundamental frequency of the voice (by using the VoicePro to either restrict or not restrict the range of the fundamental frequency of the voice) in a fully crossed design. Naive participants were given a cover story that video was 'cut out' around the face and then participated in two 8 min conversations, one with a male and another with a female confederate. Prior to debrief, the naive participants were asked if they 'noticed anything unusual about the experiment'. None mentioned that they thought they were speaking with a computer generated face or noted the experimental manipulations.
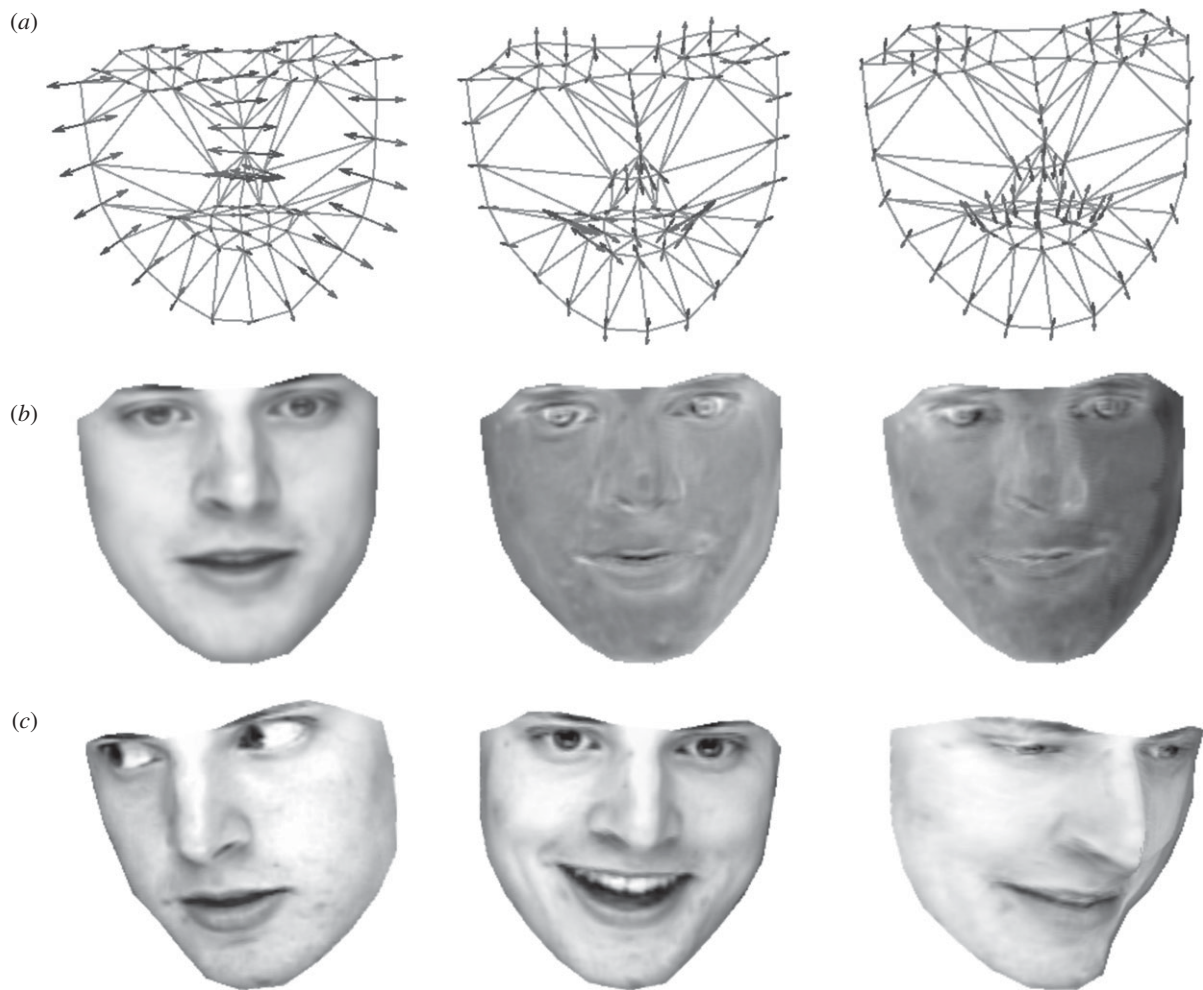
Figure 4. Illustration of AAM re-synthesis. Row (*a*) shows the mean face shape on the left and first shape modes. Row (*b*) shows the mean appearance and the first three appearance modes. The AAM is invertible and can synthesize new faces, four of which are shown in row (*c*) (From Boker & Cohn in press).

## (f) *Data reduction and analysis*

Angles of the Ascension Technologies head sensor in the anterior–posterior (A–P) and lateral directions (i.e. pitch and yaw, respectively) were selected for analysis. These directions correspond to the meaningful motion of a head nod and a head turn, respectively. We focus on angular velocity since this variable can be thought of as how animated a participant was during an interval of time.

To compute angular velocity, we first converted the head angles into angular displacement by subtracting the mean overall head angle across a whole conversation from each head angle sample. We used the overall mean head angle since this provided an estimate of the overall equilibrium head position for each participant independent of the trial conditions. Second, we low-pass filtered the angular displacement time series and calculated angular velocity using a quadratic filtering technique (generalized local linear approximation; Boker *et al*. in press), saving both the estimated displacement and the velocity for each sample. The root mean square (RMS) of the lateral and A–P angular velocity was then calculated for each 1 min condition of each conversation for each naive participant and confederate.

Because the head movements of each conversant both influence and are influenced by the movements of the other, we seek an analytic strategy that models bidirectional effects (Kenny & Judd 1986). Specifically, each conversant's head movements are both a predictor variable and an outcome variable. Neither can be considered to be an independent variable. In addition, each naive participant was engaged in two conversations, one with each of the two confederates. Each of these sources of non-independence in dyadic data needs to be accounted for in a statistical analysis.

To put both conversants in a dyad into the same analysis we used a variant of Actor–Partner analysis (Kashy & Kenny 2000; Kenny *et al*. 2006). Suppose we are analysing RMS-V angular velocity. We place both the naive participants' and the confederates' RMS-V angular velocity into the same column in the data matrix and use a second column as a dummy code labelled 'confederate' to identify whether the data in the angular velocity column came from a naive participant or a confederate. In a third column, we place the RMS-V angular velocity from the other participant in the conversation. We then use the terminology 'actor' and 'partner' to distinguish which variable is the predictor and which is the outcome for a selected row in the data matrix. If confederate = 1, then the confederate is the 'actor' and the naive participant is the 'partner' in that row of the data matrix. If confederate = 0, then the naive participant
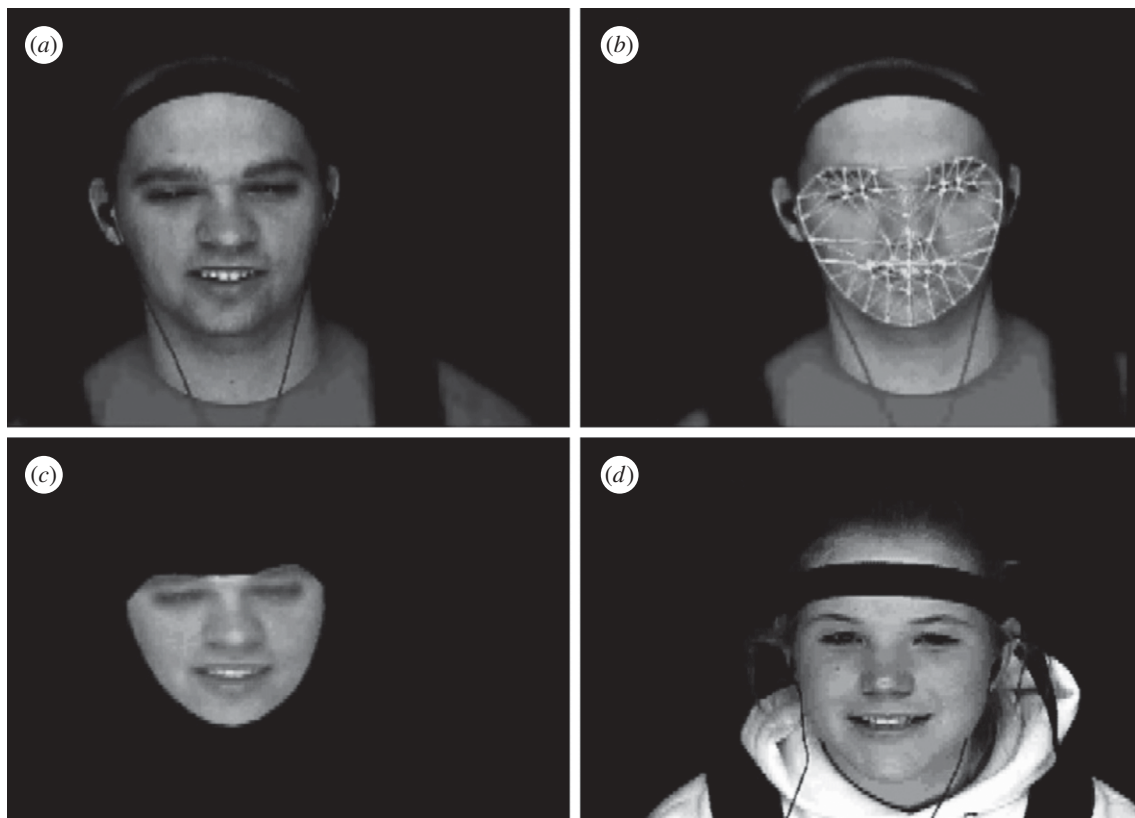
Figure 5. Illustration of the videoconference paradigm. A movie clip can be viewed at http://people.virginia.edu/smb3u/Clip1.avi. (*a*) Video of the confederate. (*b*) AAM tracking of confederate's expression. (*c*) AAM reconstruction that is viewed by the naive participant. (*d*) Video of the naive participant.

is the 'actor' and the confederate is the 'partner.' We coded the sex of the 'actor' and the 'partner' as a binary variables (0 = female, 1 = male). The RMS angular velocity of the 'partner' was used as a continuous predictor variable.

Binary variables were coded for each manipulated condition: attenuated head pitch and turn (0 = normal, 1 = 50% attenuation), and attenuated expression (0 = normal, 1 = 50% attenuation). Because only the naive participant sees the manipulated conditions, we also added interaction variables (confederate × delay condition and confederate × sex of partner), centering each binary variable prior to multiplying. The manipulated condition may affect the naive participant directly, but may also affect the confederate indirectly through changes in the behaviour of the naive participant. The interaction variables allow us to account for an overall effect of the manipulation as well as possible differences between the reactions of the naive participant and of the confederate.

We then fit mixed effect models using restricted maximum likelihood. Because there is non-independence of rows in this data matrix, we need to account for this non-independence. An additional column is added to the data matrix that is coded by experimental session and then the mixed effects model of the data is grouped by the experimental session column (both conversations in which the naive participant engaged). Each session was allowed a random intercept to account for individual differences between experimental sessions in the overall RMS

velocity. This mixed effects model can be written as

$$
\begin{aligned}
y_{ij} = b_{j0} &+ b_1 A_{ij} + b_2 P_{ij} + b_3 C_{ij} + b_4 H_{ij} + b_5 F_{ij} + b_6 V_{ij} + \\
&= b_7 Z_{ij} + b_8 C_{ij} P_{ij} + b_9 C_{ij} H_{ij} + b_{10} C_{ij} F_{ij} \\
&\quad + b_{11} C_{ij} V_{ij} + e_{ij},
\end{aligned} \tag{2.4}
$$

$$
b_{j0} = c_{00} + u_{j0}, \tag{2.5}
$$

where $y_{ij}$ is the outcome variable (lateral or A–P RMS velocity) for condition $i$ and session $j$. The other predictor variables are the sex of the actor $A_{ij}$, the sex of the partner $P_{ij}$, whether the actor is the confederate $C_{ij}$, the head pitch and turn attenuation condition $H_{ij}$, the facial expression attenuation condition $F_{ij}$, the vocal inflection attenuation condition $V_{ij}$ and the lateral or A–P RMS velocity of the partner $Z_{ij}$. As each session was allowed to have its own intercept, the predictions are relative to the overall angular velocity associated with each naive participant's session.

## 3. RESULTS

The results of a mixed effects random intercept model grouped by session predicting A–P RMS angular velocity of the head are displayed in table 1. As expected from previous reports, males exhibited lower A–P RMS angular velocity than females, and when the conversational partner was male there was lower A–P RMS angular velocity than when the conversational partner was female. Confederates exhibited lower A–P RMS velocity than naive participants, although this effect only just reached significance at the $\alpha = 0.05$ level. Both attenuated head pitch and turn, and facial
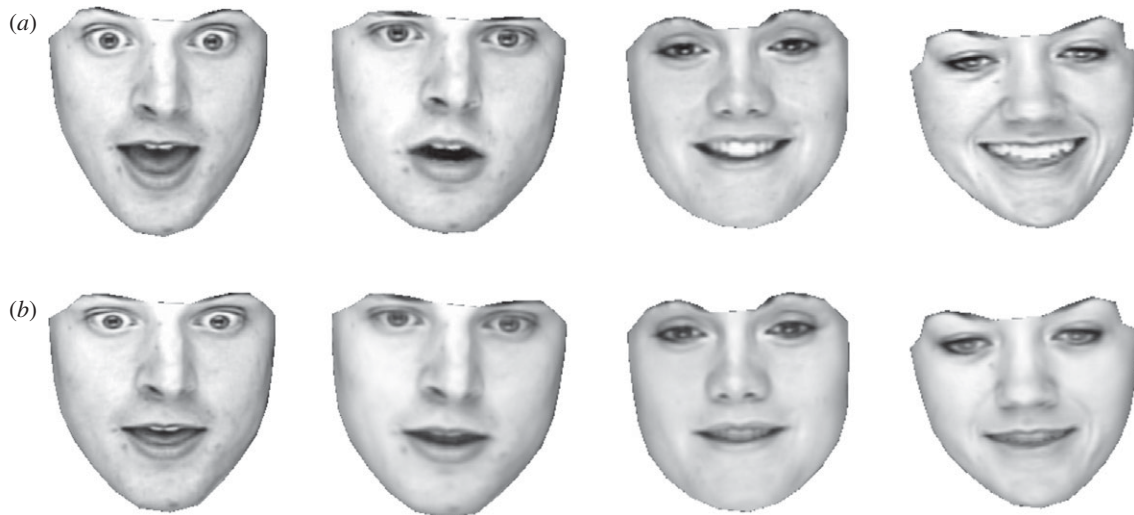
Figure 6. Facial expression attenuation using an AAM. (*a*) Four faces re-synthesized from their respective AAM models showing expressions from tracked video frames. (*b*) The same video frames displayed at 25 per cent of their AAM parameter difference from each individual's mean facial expression (i.e. $\beta = 0.25$).

Table 1. Head A–P RMS angular velocity predicted using a mixed effects random intercept model grouped by session. ('Actor' refers to the member of the dyad whose data are being predicted and 'partner' refers to the other member of the dyad. Akaike information criterion (AIC) = 3985.4, Bayesian information criterion (BIC) = 4051.1, groups = 27, random effects intercept s.d. = 1.641.)

|  | value | s.e. | d.f. | *t*-value | *p* |
|---|---|---|---|---|---|
| intercept | 10.009 | 0.5205 | 780 | 19.229 | <0.0001 |
| actor is male | −3.926 | 0.2525 | 780 | −15.549 | <0.0001 |
| partner is male | −1.773 | 0.2698 | 780 | −6.572 | <0.0001 |
| actor is confederate | −0.364 | 0.1828 | 780 | −1.991 | 0.0469 |
| attenuated head pitch and turn | 0.570 | 0.1857 | 780 | 3.070 | 0.0022 |
| attenuated expression | 0.451 | 0.1858 | 780 | 2.428 | 0.0154 |
| attenuated inflection | −0.037 | 0.1848 | 780 | −0.200 | 0.8414 |
| partner A–P RMS velocity | −0.014 | 0.0356 | 780 | −0.389 | 0.6971 |
| confederate × partner is male | −2.397 | 0.5066 | 780 | −4.732 | <0.0001 |
| confederate × attenuated head pitch and turn | −0.043 | 0.3688 | 780 | −0.116 | 0.9080 |
| confederate × attenuated expression | 0.389 | 0.3701 | 780 | 1.051 | 0.2935 |
| confederate × attenuated inflection | 0.346 | 0.3694 | 780 | 0.937 | 0.3490 |

expression was associated with greater A–P angular velocity: both conversants nodded with greater vigour when either the avatar's rigid head movement or the facial expression was attenuated. Thus, the naive participant reacted to the attenuated movement of the avatar by increasing her or his head movements. Also, the confederate (who was blind to the manipulation) reacted to the increased head movements of the naive participant by increasing his or her head movements. When the avatar attenuation was in effect, both conversational partners adapted by increasing the vigour of their head movements. There were no effects of either the attenuated vocal inflection or the A–P RMS velocity of the conversational partner. Only one interaction reached significance—confederates had a greater reduction in A–P RMS angular velocity when speaking to a male naive participant than the naive participants had when speaking to a male confederate.

The results for RMS lateral angular velocity of the head are displayed in table 2. As was true in the

A–P direction, males exhibited less lateral RMS angular velocity than females, and conversants exhibited less lateral RMS angular velocity when speaking to a male partner. Confederates again exhibited less velocity than naive participants. Attenuated head pitch and turn was again associated with greater lateral angular velocity: participants turned away or shook their heads either more often or with greater angular velocity when the avatar's head pitch and turn variation was attenuated. However, in the lateral direction, we found no effect of the facial expression or the vocal inflection attenuation. There was an independent effect such that lateral head movements were negatively coupled. That is to say in 1 min blocks when one conversant's lateral angular movement was more vigorous, their conversational partner's lateral movement was reduced. Again, only one interaction reached significance—confederates had a greater reduction in A–P RMS angular velocity when speaking to a male naive participant than the

Table 2. Head lateral RMS angular velocity predicted using a mixed effects random intercept model grouped by dyad. (AIC = 9818.5, BIC = 9884.2, groups = 27, random effects intercept s.d. = 103.20.)

|  | value | s.e. | d.f. | *t*-value | *p* |
|---|---|---|---|---|---|
| intercept | 176.37 | 22.946 | 780 | 7.686 | <0.0001 |
| actor is male | −60.91 | 9.636 | 780 | −6.321 | <0.0001 |
| partner is male | −31.86 | 9.674 | 780 | −3.293 | 0.0010 |
| actor is confederate | −21.02 | 6.732 | 780 | −3.122 | 0.0019 |
| attenuated head pitch and turn | 14.19 | 6.749 | 780 | 2.102 | 0.0358 |
| attenuated expression | 8.21 | 6.760 | 780 | 1.215 | 0.2249 |
| attenuated inflection | 4.40 | 6.749 | 780 | 0.652 | 0.5147 |
| partner A–P RMS velocity | −0.30 | 0.034 | 780 | −8.781 | <0.0001 |
| confederate × partner is male | −49.65 | 18.979 | 780 | −2.616 | 0.0091 |
| confederate × attenuated head pitch and turn | −4.81 | 13.467 | 780 | −0.357 | 0.7213 |
| confederate × attenuated expression | 6.30 | 13.504 | 780 | 0.467 | 0.6408 |
| confederate × attenuated inflection | 10.89 | 13.488 | 780 | 0.807 | 0.4197 |

naive participants had when speaking to a male confederate. There are at least three differences between the confederates and the naive participants that might account for this effect: (i) the confederates have more experience in the video booth than the naive participants and may thus be more sensitive to the context provided by the partner as the overall context of the video booth is familiar; (ii) the naive participants are seeing an avatar and it may be that there is an additional partner sex effect of seeing a full body video over seeing a 'floating head'; and (iii) the reconstructed avatars have a reduced number of eye blinks than the video since some eye blinks are not caught by the motion tracking.

## 4. DISCUSSION

Automated facial tracking was successfully applied to create real-time re-synthesized avatars that were accepted as being video by naive participants. No participant guessed that we were manipulating the apparent video in their videoconference conversations. This technological advance presents the opportunity for studying adaptive facial behaviour in natural conversation while still being able to introduce experimental manipulations of rigid and non-rigid head movements without either participant knowing the extent or timing of these manipulations.

The damping of head movements was associated with increased A–P and lateral angular velocity. The damping of facial expressions was associated with increased A–P angular velocity. There are several possible explanations for these effects. During the head movement attenuation condition, naive participants might perceive the confederate as looking more directly at him or her, prompting more incidents of gaze avoidance. A conversant might not have received the expected feedback from an A–P or lateral angular movement of a small velocity and adapted by increasing her or his head angle relative to the conversational partner in order to elicit the expected response. Naive participants may have perceived the attenuated facial expressions of the confederate as being non-responsive and attempted to increase the velocity of their head

nods in order to elicit greater response from their conversational partners.

As none of the interaction effects for the attenuated conditions were significant, the confederates exhibited the same degree of response to the manipulations as the naive participants. Thus, when the avatar's head pitch and turn variation was attenuated, both the naive participant and the confederate responded with increased velocity head movements. This suggests that there is an expected degree of matching between the head velocities of the two conversational partners. Our findings provide evidence in support of a hypothesis that the dynamics of head movement in dyadic conversation include a shared equilibrium: both conversational partners were blind to the manipulation and when we perturbed one conversant's perceptions, both conversational partners responded in a way that compensated for the perturbation. It is as if there were an equilibrium energy in the conversation and when we removed energy by attenuation, and thus changed the value of the equilibrium, the conversational partners supplied more energy in response and thus returned the equilibrium towards its former value.

These results can also be interpreted in terms of symmetry formation and symmetry breaking. The dyadic nature of the conversants' responses to the asymmetric attenuation conditions is evidence of symmetry formation. But head turns have an independent effect of negative coupling, where greater lateral angular velocity in one conversant was related to reduced angular velocity in the other: evidence of symmetry breaking. Our results are consistent with symmetry formation being exhibited in both head nods and head turns, while symmetry breaking being more related to head turns. In other words, head nods may help form symmetry between conversants while head turns, contribute to both symmetry formation and symmetry breaking. One argument for why these relationships would be observed is that head nods may be more related to acknowledgement or attempts to elicit expressivity from the partner, whereas head turns may be more related to new semantic information in the conversational stream (e.g. floor changes) or to signals of disagreement or withdrawal.

With the exception of some specific expressions (e.g. Keltner 1995; Ambadar *et al.* 2009), previous research has ignored the relationship between head movements and facial expressions. Our findings suggest that facial expression and head movement may be closely related. These results also indicate that the coupling between one conversant's facial expressions and the other conversant's head movements should be taken into account. Future research should inquire into these within-person and between-person cross-modal relationships.

The attenuation of facial expression created an effect that appeared to the research team as being that of someone who was mildly depressed. Decreased movement is a common feature of psychomotor retardation in depression, and depression is associated with decreased reactivity to a wide range of positive and negative stimuli (Rottenberg 2005). Individuals with depression or dysphoria, in comparison with non-depressed individuals, are less likely to smile in response to pictures or movies of smiling faces and affectively positive social imagery (Gehricke & Shapiro 2000; Sloan *et al.* 2002). When they do smile, they are more likely to damp their facial expression (Reed *et al.* 2007).

Attenuation of facial expression can also be related to cognitive states or social context. For instance, if one's attention is internally focused, the attenuation of facial expression may result. Interlocutors might interpret damped facial expression of their conversational partner as reflecting a lack of attention to the conversation.

Naive participants responded to damped facial expression and head turns by increasing their own head nods and head turns, respectively. These effects may have been efforts to elicit more responsive behaviour in the partner. In response to simulated maternal depression by their mother, infants attempt to elicit a change in their mother's behaviour by smiling, turning away and then turning again towards her and smiling. When they fail to elicit a change in their mothers' behaviour, they become withdrawn and distressed (Cohn & Tronick 1983). Similarly, adults find exposure to prolonged depressed behaviour increasingly aversive and withdraw (Coyne 1976). Had we attenuated facial expression and head motion for more than a minute at a time, naive participants might have become less active following their failed efforts to elicit a change in the confederate's behaviour. This hypothesis remains to be tested.

There are a number of limitations of this methodology that could be improved with further development. For instance, while we can manipulate the degree of expressiveness as well as the identity of the avatar (Boker & Cohn in press), we cannot yet manipulate specific facial expressions in real time. Depression not only attenuates expression, but also makes some facial actions, such as contempt, more likely (Ekman *et al.* 2005; Cohn *et al.* 2009). As an analogue for depression, it would be important to manipulate specific expressions in real time. In other contexts, cheek raising (AU 6 in the facial action coding system) (Ekman *et al.* 2002) is believed to covary with communicative intent and felt emotion (Coyne 1976). In the past, it has not been possible to

experimentally manipulate discrete facial actions in real time without the source person's awareness. If this capability could be implemented in the videoconference paradigm, it would make possible a wide range of experimental tests of emotion signalling.

Other limitations include the need for person-specific models, restrictions on head rotation and limited face views. The current approach requires manual training of face models, which involves hand labelling about 30–50 video frames. Because this process requires several hours of pre-processing, avatars could be constructed for confederates but not for unknown persons, such as naive participants. It would be useful to have the capability of generating real-time avatars for both conversation partners. Recent efforts have made progress towards this goal (Lucey *et al.* in press; Saragih *et al.* 2009). Another limitation is that if the speaker turns more than about $20°$ from the camera, parts of the face become obscured and the model can no longer track the remainder of the face. Algorithms have been proposed that address this issue (Gross *et al.* 2004), but it remains a research question. Another limitation is that the current system has modelled the face only from the eyebrows to the chin. A better system would include the forehead, and some model of the head, neck, shoulders and background in order to give a better sense of the placement of the speaker in context. Adding forehead features is relatively straight-forward and has been implemented. Tracking of neck and shoulders is well advanced (Sheikh *et al.* 2008). The videoconference avatar paradigm has motivated new work in computer vision and graphics and made possible new methodology to experimentally investigate social interaction in a way not possible before. The timing and identity of social behaviour in real time can now be rigorously manipulated outside of participants' awareness.

## 5. CONCLUSION

We presented an experiment that used automated facial and head tracking to perturb the bidirectionally coupled dynamical system formed by two individuals speaking with one another over a videoconference link. The automated tracking system allowed us to create re-synthesized avatars that were convincing to naive participants and, in real time, to attenuate head movements and facial expressions formed during natural dyadic conversation. The effect of these manipulations exposed some of the complexity of multi-modal coupling of movements during face to face interactions. The experimental paradigm presented here has the potential to transform social psychological research in dyadic and small group interactions owing to an unprecedented ability to control the real-time appearance of facial structure and expression.

are those of the authors and do not necessarily reflect the views of the National Science Foundation. We gratefully acknowledge the help of Kathy Ashenfelter, Tamara Buretz, Eric Covey, Pascal Deboeck, Katie Jackson, Jen Koltiska, Sean McGowan, Sagar Navare, Stacey Tiberio, Michael Villano and Chris Wagner.

## REFERENCES

Ambadar, Z., Cohn, J. F. & Reed, L. I. 2009 All smiles are not created equal: morphology and timing of smiles perceived as amused, polite, and embarrassed/nervous. *J. Nonverbal Behav.* **33**, 17–34. (doi:10.1007/s10919-008-0059-5)

Ashenfelter, K. T., Boker, S. M., Waddell, J. R., Vitanov, N. & Abadjieva, E. 2009 Spatiotemporal symmetry and multifractal structure of head movements during dyadic conversation. *J. Exp. Psychol. Hum. Percept. Perform.* **35**(4), 1072–1091.

Bernieri, F. J. 1988 Coordinated movement and rapport in teacher–student interactions. *J. Nonverbal Behav.* **12**, 120–138. (doi:10.1007/BF00986930)

Boker, S. M. & Cohn, J. F. In press. Real time dissociation of facial appearance and dynamics during natural conversation. In *Dynamic faces: insights from experiments and computation* (eds M. Giese, C. Curio & H. Bültoff). Cambridge, MA: MIT Press.

Boker, S. M. & Rotondo, J. L. 2002 Symmetry building and symmetry breaking in synchronized movement. In *Mirror neurons and the evolution of brain and language* (eds M. Stamenov & V. Gallese), pp. 163–171. Amsterdam, The Netherlands: John Benjamins.

Boker, S. M., Xu, M., Rotondo, J. L. & King, K. 2002 Windowed cross-correlation and peak picking for the analysis of variability in the association between behavioral time series. *Psychol. Meth.* **7**, 338–355. (doi:10.1037/1082-989X.7.3.338)

Boker, S. M., Deboeck, P. R., Edler, C. & Keel, P. K. In press. Generalized local linear approximation of derivatives from time series. In *Statistical methods for modeling human dynamics: an interdisciplinary dialogue* (eds S.-M. Chow & E. Ferrar). Boca Raton, FL: Taylor & Francis.

Cappella, J. N. & Panalp, S. 1981 Talk and silence sequences in informal conversations: III: interspeaker influence. *Hum. Commun. Res.* **7**, 117–132. (doi:10.1111/j.1468-2958.1981.tb00564.x)

Cohn, J. F. & Tronick, E. Z. 1983 Three month old infants' reaction to simulated maternal depression. *Child Dev.* **54**, 185–193. (doi:10.2307/1129876)

Cohn, J. F., Kreuze, T. S., Yang, Y., Gnuyen, M. H., Padilla, M. T. & Zhou, F. 2009 Detecting depression from facial actions and vocal prosody. *Affective computing and intelligent interaction (ACII 2009)*. Amsterdam, The Netherlands: IEEE.

Cohn, J. F., Reed, L. I., Moriyama, T., Xiao, J., Schmidt, K. L. & Ambadar, Z. 2004 Multimodal coordination of facial action, head rotation, and eye motion. *6th IEEE International Conference on Automatic Face and Gesture Recognition*, pp. 645–650. Seoul, Korea: IEEE.

Cootes, T. F., Edwards, G. & Taylor, C. J. 2001 Active appearance models. *IEEE Trans. Pattern Anal. Machine Intell.* **23**, 681–685. (doi:10.1109/34.927467)

Cootes, T. F., Wheeler, G. V., Walker, K. N. & Taylor, C. J. 2002 View-based active appearance models. *Image Vis. Comput.* **20**, 657–664. (doi:10.1016/S0262-8856(02)00055-0)

Coyne, J. C. 1976 Depression and the response of others. *J. Abnorm. Psychol.* **85**, 186–193. (doi:10.1037/0021-843X.85.2.186)

Ekman, P., Friesen, W. & Hager, J. 2002 *Facial action coding system*. Salt Lake City, UT: Research Nexus.

Ekman, P., Matsumoto, D. & Friesen, W. V. 2005 Facial expression in affective disorders. In *What the face reveals* (eds P. Ekman & E. Rosenberg), pp. 331–341. New York, NY: Oxford University Press.

Gehricke, J.-G. & Shapiro, D. 2000 Reduced facial expression and social context in major depression: discrepancies between facial muscle activity and self-reported emotion. *Psychiatry Res.* **95**, 157–167. (doi:10.1016/S0165-1781(00)00168-2)

Gross, R., Matthews, I. & Baker, S. 2004 Constructing and fitting active appearance models with occlusion. *1st IEEE Workshop on Face Processing in Video*. Washington, DC: IEEE.

Hsee, C. K., Hatfield, E., Carlson, J. G. & Chemtob, C. 1990 The effect of power on susceptibility to emotional contagion. *Cogn. Emotion* **4**, 327–340. (doi:10.1080/02699939008408081)

Iacoboni, M., Woods, R. P., Brass, M., Bekkering, H., Mazziotta, J. C. & Rizzolatti, G. 1999 Cortical mechanisms of human imitation. *Science*, **286**, 2526–2528.

Kashy, D. A. & Kenny, D. A. 2000 The analysis of data from dyads and groups. In *Handbook of research methods in social psychology* (eds H. Reis & C. M. Judd), pp. 451–477. New York, NY: Cambridge University Press.

Keltner, D. 1995 Signs of appeasement: evidence for the distinct displays of embarrassment, amusement, and shame. *J. Pers. Soc. Psychol.* **68**, 441–454. (doi:10.1037/0022-3514.68.3.441)

Kenny, D. A. & Judd, C. M. 1986 Consequences of violating the independence assumption in analysis of variance. *Psychol. Bull.* **99**, 422–431. (doi:10.1037/0033-2909.99.3.422)

Kenny, D. A., Kashy, D. A. & Cook, W. L. 2006 *Dyadic data analysis*. New York, NY: Guilford.

LaFrance, M. 1982 Posture mirroring and rapport. In *Interaction rhythms: periodicity in communicative behavior* (ed. M. Davis), pp. 279–298. New York, NY: Human Sciences Press.

Lucey, S., Wang, Y., Cox, M., Sridharan, S. & Cohn, J. F. In press. Efficient constrained local model fitting for non-rigid face alignment. *Image Vis. Comp. J.*

Matthews, I. & Baker, S. 2004 Active appearance models revisited. *Int. J. Comput. Vis.* **60**, 135–164. (doi:10.1023/B:VISI.0000029666.37597.d3)

Messinger, D. S., Chow, S. M. & Cohn, J. F. 2009 Automated measurement of smile dynamics in mother–infant interaction: a pilot study. *Infancy* **14**, 285–305. (doi:10.1080/15250000902839963)

Neumann, R. & Strack, F. 2000 'Mood contagion': the automatic transfer of mood between persons. *J. Pers. Soc. Psychol.* **79**, 158–163. (doi:10.1037/0022-3514.79.2.211)

Redlich, N. A. 1993 Redundancy reduction as a strategy for unsupervised learning. *Neural Comput.* **5**, 289–304. (doi:10.1162/neco.1993.5.2.289)

Reed, L. I., Sayette, M. A. & Cohn, J. F. 2007 Impact of depression on response to comedy: a dynamic facial coding analysis. *J. Abnorm. Psychol.* **116**, 804–809. (doi:10.1037/0021-843X.116.4.804)

Rizzolatti, G. & Craighero, L. 2004 The mirror–neuron system. *Ann. Rev. Neurosci.* **27**, 169–192. (doi:10.1146/annurev.neuro.27.070203.144230)

Rizzolatti, G. & Fadiga, L. 2007 Grasping objects and grasping action meanings: the dual role of monkey rostroventral premotor cortex. In *Novartis Foundation Symposium 218—Sensory Guidance of Movement* (eds P. Ekman &

E. Rosenberg), pp. 81–108. New York, NY: Novartis Foundation.

Rottenberg, J. 2005 Mood and emotion in major depression. *Curr. Dir. Psychol. Sci.* **14**, 167–170. (doi:10.1111/j.0963-7214.2005.00354.x)

Saragih, J., Lucey, S. & Cohn, J. F. 2009 Probabilistic constrained adaptive local displacement experts. Workshop at the *Int. Conf. on Computer Vision*, Kyoto, Japan.

Shannon, C. E. & Weaver, W. 1949 *The mathematical theory of communication*. Urbana, IL: The University of Illinois Press.

Sheikh, Y. A., Datta, A. & Kanade, T. 2008 On the sustained tracking of human motion. *Proc. 8th Int. Conf. on Automatic Face and Gesture Recognition*, Amsterdam, The Netherlands.

Sloan, D. M., Bradley, M. M., Dimoulas, E. & Lang, P. J. 2002 Looking at facial expressions: dysphoria and facial EMG. *Biol. Psychol.* **60**, 79–90. (doi:10.1016/S0301-0511(02)00044-3)

Theobald, B., Matthews, I., Cohn, J. F. & Boker, S. 2007 Real-time expression cloning using appearance models. *Proc. 9th Int. Conf. on Multimodal Interfaces*, pp. 134–139. New York, NY: Association for Computing Machinery.

Young, R. D. & Frye, M. 1966 Some are laughing; some are not—why? *Psychol. Rep.* **18**, 747–752.