



# The human likeness dimension of the “uncanny valley hypothesis”: behavioral and functional MRI findings

Marcus Cheetham\*, Pascal Suter and Lutz Jäncke

Department of Neuropsychology, University of Zürich, Zürich, Switzerland

**Edited by:**

Hans-Jochen Heinze, University of  
Magdeburg, Germany

**Reviewed by:**

Bernd Weber,  
Rheinische-Friedrich-Wilhelms  
Universität, Germany  
Beatrice De Gelder, Donders  
Institute, Netherlands

**\*Correspondence:**

Marcus Cheetham, Department of  
Neuropsychology, University of  
Zurich, Binzmühlestrasse 14/Box 25,  
CH-8050 Zürich, Switzerland.  
e-mail: m.cheetham@psychologie.  
uzh.ch

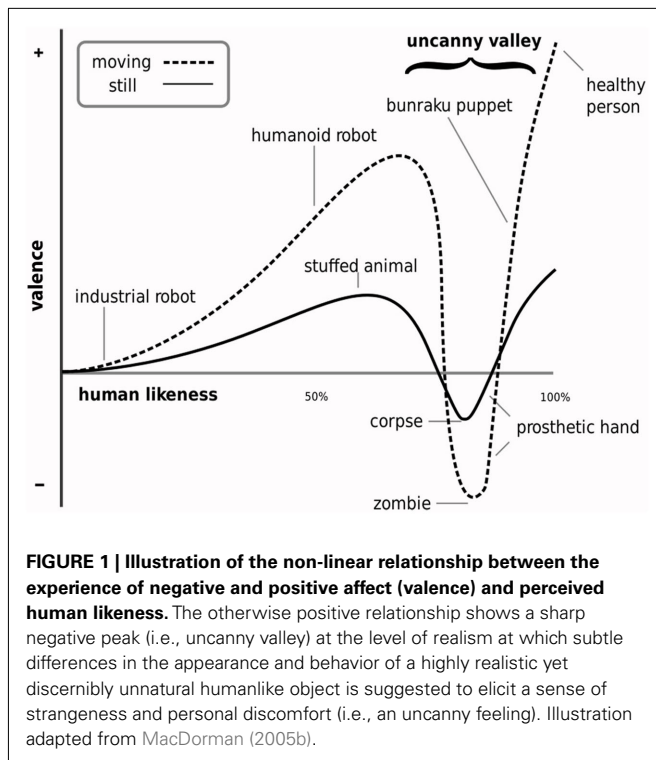
The uncanny valley hypothesis (Mori, 1970) predicts differential experience of negative and positive affect as a function of human likeness. Affective experience of humanlike robots and computer-generated characters (avatars) dominates “uncanny” research, but findings are inconsistent. Importantly, it is unknown how objects are actually perceived along the hypothesis’ dimension of human likeness (DOH), defined in terms of human physical similarity. To examine whether the DOH can also be defined in terms of effects of categorical perception (CP), stimuli from morph continua with controlled differences in physical human likeness between avatar and human faces as endpoints were presented. Two behavioral studies found a sharp category boundary along the DOH and enhanced visual discrimination (i.e., CP) of fine-grained differences between pairs of faces at the category boundary. Discrimination was better for face pairs presenting category change in the human-to-avatar than avatar-to-human direction along the DOH. To investigate brain representation of physical change and category change along the DOH, an event-related functional magnetic resonance imaging study used the same stimuli in a pair-repetition priming paradigm. Bilateral mid-fusiform areas and a different right mid-fusiform area were sensitive to physical change within the human and avatar categories, respectively, whereas entirely different regions were sensitive to the human-to-avatar (caudate head, putamen, thalamus, red nucleus) and avatar-to-human (hippocampus, amygdala, mid-insula) direction of category change. These findings show that Mori’s DOH definition does not reflect subjective perception of human likeness and suggest that future “uncanny” studies consider CP and the DOH’s category structure in guiding experience of non-human objects.

**Keywords:** uncanny valley, fMRI, categorical perception, virtual reality, avatar, human likeness

## INTRODUCTION

The uncanny valley hypothesis (Mori, 1970) proposes that observation of a humanlike object (e.g., industrial robot, lifelike prosthetic hand, human corpse, doll, mannequin) can evoke positive or negative feelings and cognitions (referred to here as *valence*) depending on the object’s degree of physical similarity to human appearance (for recent overviews, see, e.g., Pollick, 2009; MacDorman et al., 2009). Importantly, the relationship between valence and human likeness is suggested to be non-linear. As illustrated in **Figure 1**, valence increases positively with greater human likeness up to the point of relatively high realism at the first peak of the curve along the *dimension of human likeness* (DOH). Positive valence reflects the experience of emotional engagement with and feelings of empathy for the humanlike object. At greater degrees of realism, the observer encounters difficulty in distinguishing an object from its natural human counterpart and experiences personal discomfort. Mori characterizes this discomfort as an uncanny feeling marked by a sense of strangeness, eeriness, and disquiet that can extend to feelings of disgust and revulsion. This uncanny feeling is reflected in a sharp negative peak or valley (i.e., uncanny valley) in the slope of the depicted valence–human likeness relationship. When an object’s appearance is so realistic that it is perceived to be human the valence of associated affect and cognition is thought to reach a second positive peak.

The hypothesis’ central prediction that individuals can feel less emotionally engaged or even distracted by relatively realistic humanlike objects has been very influential in guiding animators, video game designers, and roboticists in the design of virtual characters (e.g., Fabri et al., 2004) and robots (e.g., Minato et al., 2006). This has led to research into how characters should be designed to avoid “falling” into the uncanny valley (e.g., Walters et al., 2008; MacDorman et al., 2009) and into tool development to aid such design decision making (Ho and MacDorman, 2010). The idea of the uncanny valley has also been used in empirical studies to account for unexpected findings and negative responses of participants to the appearance and behavior of avatars or robots (e.g., Aylett, 2004; Wages et al., 2004; Yamamoto et al., 2009). In view of this interest, it is noteworthy that Mori did not subject his working hypothesis to empirical examination. In fact, research of the valence–human likeness relationship is in its infancy, largely focusing on the valence of subjective feelings and evaluations in response to variously realistic non-human characters. But findings have been inconsistent (e.g., Hanson, 2006; MacDorman, 2006; Tinwell and Grimshaw, 2009; Tinwell et al., 2010), this being in part attributable to the uncertainty surrounding Mori’s vague definition of the valence dimension (e.g., Seyama and Nagayama, 2007; MacDorman et al., 2009). In contrast, the way in which human likeness



along the DOH is actually perceived has not been scrutinized at all.

While Mori defined the DOH in terms of a linear change in the degree of physical humanlike similarity, we proposed that the subjective perception of objects along the DOH might be described differently, namely, in terms of the effects of categorical perception (for CP see, e.g., Harnad, 1987). Applied to the DOH, CP means that, irrespective of physical differences in humanlike appearance, objects along the DOH are treated as conceptually equivalent members of either the category “non-human” or the category “human,” except at those levels of physical realism at the boundary between these two categories. At this category boundary, the available sensory evidence does not allow rapid and effortless discrimination of an object on the basis of the observer’s category representations of non-human and human. Consistent with this, ambiguity in discriminating a humanlike object from its natural human counterpart lies at the heart of Mori’s hypothesis, and, when no such ambiguity is experienced, Mori implicitly assumes that the observer assigns his or her sensory impressions of an object to the non-human or the human category. Importantly, the defining feature of CP is considered to be enhanced discrimination of pairs of different stimuli that are perceptually adjacent (along a dimension such as the DOH) but straddle opposite sides of the category boundary (such as between the categories human and non-human) and poorer discrimination of pairs of different stimuli from within a given category (Pastore, 1987). In other words, the ability to discriminate between physically different stimuli might not be the same at all points along the DOH, with enhanced sensitivity for objects closest to the category boundary (for CP criteria, see, e.g., Studdert-Kennedy et al., 1970).

It is important to note that CP does not occur along any human likeness dimension. Campbell et al. (1987) found CP for faces of humans and cows but not humans and monkeys. CP for the non-human and human categories along Mori’s DOH has not been investigated. In fact, a prominent feature of Mori’s hypothesis is that it does not consider the possibility that objects assigned to the human category might also differ in the degree of humanlike similarity along the DOH. Our interest in understanding the categorical structure of the DOH is based on the assumption that negatively valenced or uncanny experience is likely to occur in association with categorization ambiguity for stimuli at or closest to either side of the non-human–human category boundary. For this, the human category and the potential impact of CP on processing objects along the full length of the DOH need to be first considered. The present investigation focused therefore on the subjective perception of human likeness rather than on valence in order to provide a clear theoretical–psychological framework for further studies of uncanny experience.

Two behavioral and one neuroimaging study was conducted. The behavioral studies aimed to demonstrate that faces assigned to the human category can also vary in humanlike appearance along the DOH and to test the proposal that subjective perception of human likeness along the DOH shows effects of CP. To represent the DOH, linear continua of morphed faces were generated using avatar and human faces as endpoints to create a controlled transition of physical similarity between them. Using stimuli drawn from these continua, the first study entailed a forced choice classification task to determine the presence and location of the avatar–human category boundary, and the second used a perceptual discrimination task to determine the presence and location of the discrimination boundary (i.e., the point of enhanced perceptual discrimination) and to verify CP. The subsequent neuroimaging study used functional magnetic resonance imaging (fMRI) to investigate non-human–human category processing in the brain. The first aim of this was to show that distinctly different brain regions are responsive to processing the physical similarity of faces along the DOH and to processing the DOH’s discrete human and non-human categories. For this, we used a pair-repetition priming paradigm (for reviews see, Grill-Spector and Malach, 2001; Henson, 2003). This paradigm entailed the presentation of stimulus trials, each comprising a pair of faces (i.e., prime and target) displayed in quick succession, and the measurement of the physiological BOLD response to the primed targets as an indicator of neural activity. Repetition in the target of information presented in the prime has been shown to induce neural adaptation (i.e., attenuation of neural activity) in brain areas selectively responsive to processing the repeated information (e.g., Grill-Spector et al., 2006; Jiang et al., 2006). This neural adaptation effect and its regional localization have been successfully used to disentangle regions responsive to physical and category change (e.g., Rotshtein et al., 2005; Jiang et al., 2007).

We focused in the fMRI study on category change, hypothesizing that implicit processing of the discrete change in category between our prime and target stimuli would modulate neural activity in regions associated with category learning and

categorization uncertainty, that is, basal ganglia, medial temporal lobe (MTL), thalamus, medial frontal gyrus, and the anterior insula (e.g., Ashby and Maddox, 2005; Grinband et al., 2006; Li et al., 2009; Fleming et al., 2010; Seger and Miller, 2010). Mori focused in his informal description of the hypothesis on the situation in which a non-human object is initially mistaken for human. The second aim of this fMRI study was to explore the possibility that the anticipated effects of category change in these brain regions might be different depending on the actual direction of category change between the stimuli used as prime and target, that is, in the human-to-avatar direction along the DOH (i.e., human as prime and avatar as target) or the avatar-to-human direction (i.e., avatar as prime and human as target). The basis for this idea was that different brain regions (e.g., in the basal ganglia and MTL) are known to be differentially modulated depending on categorization experience with a given category (e.g., Poldrack et al., 1999). Given the asymmetrical category knowledge of our participants with human and novel non-human faces, we assumed that this differential effect might similarly apply for category processing along the DOH in individuals with little or no experience of our non-human stimuli. Category-related processing of highly realistic and potentially biologically salient humanlike faces might also modulate neural responses in further regions specifically associated with appraisal of affective valence (e.g., Vuilleumier, 2005) or processing under conditions of valence ambiguity (e.g., Herwig et al., 2007). We explored this, anticipating involvement of regions commonly associated with affective processing, especially the amygdala (e.g., Herwig et al., 2007; Todorov and Engell, 2008).

## MATERIALS AND METHODS

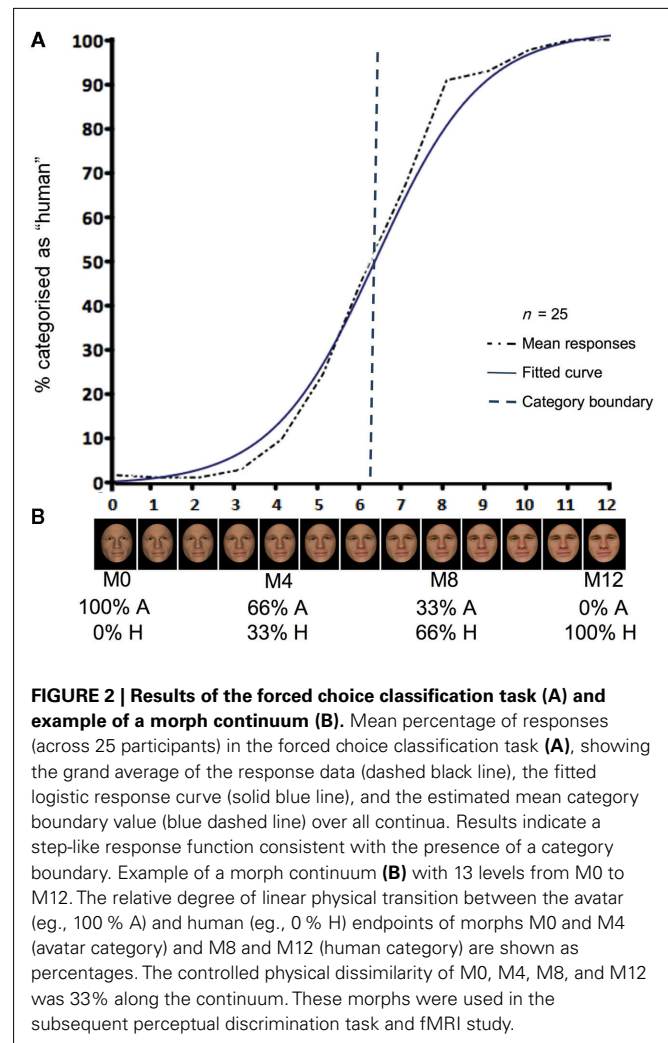
### PARTICIPANTS

Healthy male and female, consistently right-handed (Annett, 1970) adult volunteers with no record of neurological or psychiatric illness and no current medication use volunteered for one of the three studies. All participants were students of the University of Zurich, native speakers of Swiss-German or Standard German, reported having no explicit experience with avatars such as those used in video games, virtual role playing games, or second life, and had normal or corrected-to-normal vision. Written informed consent was obtained before participation according to the guidelines of the Declaration of Helsinki. Each volunteer received 20 Swiss Francs for participation. The study and all procedures and consent forms were approved by the Ethics Committee of the University of Zurich.

### STIMULI

Avatar-human morph continua were generated to represent Mori's DOH. The selected stimuli for use as endpoints in the morphing procedure were 32 photographic images of unknown people and 32 facial images of avatars, together forming 32 morph continua (see Figure 2B). All faces were male, presented with full face, frontal view, and neutral expression. The avatars were generated with the modeling suite Poser 7 (Smith Micro Software)<sup>1</sup>,

<sup>1</sup><http://www.smithmicro.com>



which permits considerable texture detail and control over the facial mesh. The facial geometry and texture of the avatars were modeled (age and configural cues) to closely match the human counterpart to minimize perception of biological motion during quick successive presentation of faces (e.g., Schultz and Pilz, 2009). The images were edited in Adobe Photoshop CS3, masking external features with an elliptic form and black background (72 dpi and 560 × 650 pixels). Contrast levels, overall brightness, and skin tone of each pair of color images of faces used as the endpoints of each continuum were adjusted to match before morphing. Morpher 3.3 (Zealsoft Inc., Eden Prairie, MN, USA) was used to generate the linear morph continua. Each continuum comprised 13 different morphed images. These images or morphs were labeled M0 (beginning with the avatar endpoint) through to M12 (at the human endpoint), each morph representing a difference in physical appearance at increments of 8.33%.

All morphs (i.e., M0–M12) were presented in the forced choice classification task of study 1. As described in the following, only those morphs representing increments of 33.33% in physical difference along each continuum (i.e., M0, M4, M8, M12) were used

for the subsequent perceptual discrimination task of study 2 and the target monitoring task in the (fMRI) study 3.

### STUDY 1: FORCED CHOICE CLASSIFICATION TASK

A two-alternative forced choice classification task was conducted to determine the presence of a sigmoid shaped response function along the avatar–human morph continua. This response function is considered to indicate a category boundary (Harnad, 1987).

#### Participants

$N = 25$  volunteers (13 female, mean age 21.8 years; range 19–28 years) participated.

#### Materials and Procedure

All participants were tested individually. Each participant received written instructions presented on the screen before commencement of the experiment. A practice pre-test of five trials using stimuli from continua not included in the main test was then performed to ensure that the instructions had been understood. It was again ensured that participants understood the meaning of the word avatar. All morph stimuli were presented individually in random order for a total presentation of 416 trials. Each trial began with the presentation of a fixation point for 500 ms (participants were required to maintain fixation), followed by a morph image for 750 ms. The participant was asked to identify the stimulus as either an avatar or human as quickly and precisely as possible by pressing one of two response keys. A black screen with fixation point remained after morph image presentation until the participant pressed a key, after which a blank black screen without fixation cross remained for 1500 ms until the next trial began. The task was conducted in a sound attenuated and light-dimmed room, and morph stimuli were presented on a LCD monitor (1280 × 1024 resolution, 60 Hz refresh rate, at eye-to-monitor viewing distance of 62 cm), using Presentation<sup>®</sup> software (Version 14.1)<sup>2</sup>.

#### Analyses

The classification data were summarized by the shape of the avatar–human classification curve as described by the slope of the response function. The slope was determined by fitting logistic function models to the response data of each participant and continuum, and the parameter estimates were derived. Individual continua were analyzed across participants to ensure best fit of logistic functions, these reflecting a step-like shape in the avatar–human classification curve and thus a categorical component (Harnad, 1987). To test for a step-like shape across all continua, the derived parameter estimates for the logistic function of each continuum, averaged across participants, were tested against zero in a one-sample  $t$ -test. This and the average boundary value (i.e., the morph position associated with greatest decision uncertainty) of the fitted logistic curves are reported. SPSS 16 was used for data analysis.

### STUDY 2: PERCEPTUAL DISCRIMINATION TASK

A variant of the same-different discrimination task (e.g., Angeli et al., 2008) was performed to determine CP along the morph

continua. Participants judged whether the presented faces of a face pair were both the same or different in appearance. Greater performance accuracy in “different” judgments for face pairs that are perceptually adjacent (and therefore physically different along the avatar–human continua) but straddle opposite sides of a category boundary (as determined in the preceding task) is taken as evidence of CP (Pastore, 1987).

#### Participants

$N = 20$  participants (nine female, mean age 25.1 years; range 18–30 years) were tested.

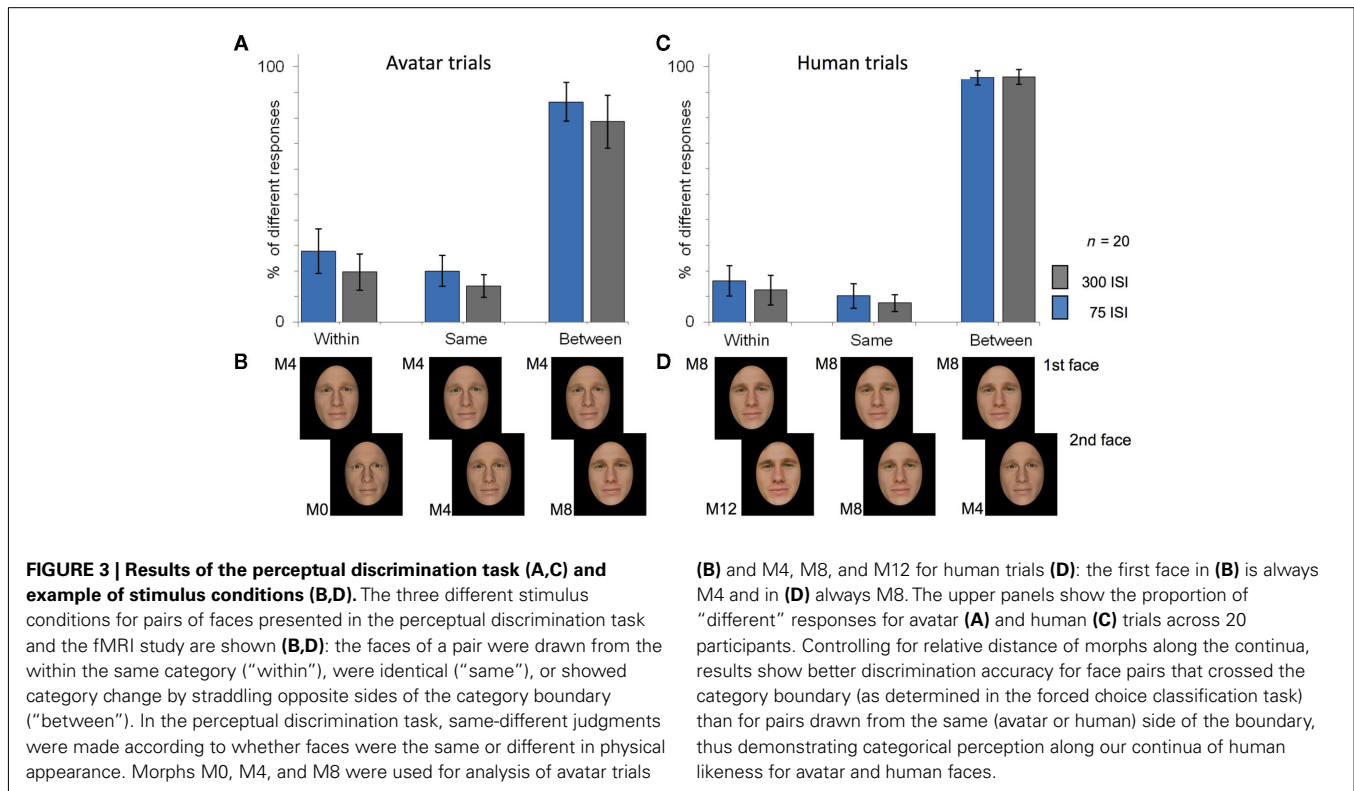
#### Materials and Procedure

The choice of continua for use in this task was determined on the basis of the results (see Forced Choice Classification) of the preceding forced choice classification task. Four morphs M0 and M4 (categorized as avatars) and M8 and M12 (categorized as human) were selected from each of the continua (see **Figure 2B**), with M4 and M8 straddling the category boundary. M0, M4, M8, and M12 represent increments of 33% along the length of each continuum to ensuring control of physical dissimilarity.

For the perceptual discrimination task, stimulus trials were presented that each comprised a pair of faces (see **Figure 3**). The first face of a face pair was either an avatar or a human. Trials in which the first face was an avatar are referred to as “avatar” trials (see **Figure 3C**) and those in which a human face was presented first are referred to as “human” trials (see **Figure 3D**). In the avatar trials, the first face was always the avatar morph M4, and the first face in the human trials was always the human morph M8. The choice of the second face of each face pair was determined according to three different face pair conditions: “within” category face pairs (both faces drawn from within a category), “between” category face pairs (first and second face morphs lying either side of the category boundary), and the “same” stimulus face pairs (faces of a pair are identical). In the within category condition, the second face was morph M0 for avatar trials and M12 for human trials. In the between category condition, the second face was morph M8 for avatar trials and M4 for human trials. In the “same” stimulus condition, M4 was used as the first and second face for avatar trials and M8 as the first and second face for human trials. Based on these conditions, the presentation of face pairs was therefore as follows: Avatar trials comprised the face pair morphs M4–M0 for the within, M4–M4 for the same, and M4–M8 for the between conditions, and humans trials comprised the face pair morphs M8–M12 for the within, M8–M8 for the same, and M8–M4 for the between conditions.

Both faces of each face pair were always drawn from the same continuum in which they were originally morphed. The presentation of face pairs was pseudo-randomized so that no trials of face pairs from within the same continuum were shown in close sequence. The presentation of avatar or human trials from a given continuum was determined randomly but counterbalanced across all participants in such a way as to ensure that each participant viewed either avatar or human trials from any given continuum, and that in total an equal number of avatar or human trials were viewed. Each face of a face pair was presented for 500 ms with

<sup>2</sup>www.neurobs.com



an inter stimulus interval (ISI) of either 75 or 300 ms between the faces of a pair. This was to determine whether ISI duration would differentially impact judgments (see, e.g., Rotshtein et al., 2005). Stimulus onset asynchrony (SOA) was 2,500 ms between trials of face pairs. Please note that these features of the experimental design of stimulus presentation were also carried over to the fMRI investigation of Study 3 (and are therefore not described again in the Materials and Procedure of Study 3). This experimental design allowed within category effects for faces from the avatar category and for faces from the human category to be examined, and the effects of category change between face pairs comprising an avatar and a human face to be examined in both directions along the DOH, that is, in the avatar-to-human and human-to-avatar directions.

The task was conducted in a sound attenuated and light-dimmed room, and morph stimuli were presented on a LCD monitor (1280 × 1024 resolution, 60 Hz refresh rate, at eye-to-monitor viewing distance of 62 cm), using Presentation® software (Version 14.1, see text footnote 2)

### Analyses

To examine discrimination accuracy for face pairs that crossed the category boundary compared with face pairs from the same side of the boundary, the “different” responses (indicating the judgment that both faces of a pair were of different physical appearance) were computed as proportions of the total number of face trials and subjected to  $3 \times 2$  ANOVA, with 3 face-pair conditions (within, between, same) and 2 ISI (short and long). The data for avatar trials and human trials were treated separately in analysis. Greenhouse–Geisser adjustment was applied to correct the degrees

of freedom whenever the assumption of sphericity was violated. SPSS 16 was used for data analysis.

### STUDY 3: TARGET MONITORING TASK AND FMRI

#### Participants

$N = 22$  volunteers (10 female; mean age, 21.9 years; range, 19–27 years) participated in the fMRI study.

#### Materials and Procedure

A pair-repetition paradigm with a target monitoring task was applied. The stimulus conditions (i.e., the morph stimuli for the face pairs in the within, same and between conditions in the avatar and human trials, and the presentation times for the stimuli, ISI, and SOA) were the same as described in the preceding perceptual discrimination task.

The target monitoring task required participants to press a response button upon detection of one of four possible up-turned rare target faces (M0, M4, M8, or M12) selected at random from an unused morph continuum. This was to ensure attention of the participants to the stimuli of interest. This task requires no explicit judgments of human or avatar category or physical similarity of faces, thus eliminating the confounding effects on BOLD signal modulation by motor response requirements (Henson, 2003). Differential attention to prime and target was avoided in that the up-turned target face was presented as the first or second face of a face pair trial. There were 192 trials of face pairs, 15% of this trial number being in addition rare target faces and 25% in addition null events with fixation cross but no face stimuli. Participants were required to maintain fixation. Each scanning session consisted of two experimental runs. The visual stimuli were presented using



the “VisuaStim – Digital” MRI-compatible head-mounted display (Resonance Technology Inc.), with a visual mono display, resolution of  $800 \times 600$ , and  $30^\circ$  field of view. The total scanning time was approximately 26 min.

### **fMRI data acquisition**

Structural and functional images were acquired using a 3-T whole-body MR unit (Philips Medical Systems, Best, The Netherlands) and eight-channel Philips SENSE head coil. Structural images of the entire brain were registered using a T1-weighted three-dimensional, spoiled gradient echo pulse sequence (180 slices, TR = 20 ms, TE = 2.3 ms, flip angle =  $20^\circ$ , FOV =  $220 \text{ mm} \times 220 \text{ mm} \times 135 \text{ mm}$ , matrix size =  $224 \times 187$ , voxel size =  $0.98 \text{ mm} \times 1.18 \text{ mm} \times 0.75 \text{ mm}$ , resliced to  $0.86 \text{ mm} \times 0.86 \text{ mm} \times 0.75 \text{ mm}$ ). Functional images were acquired from 225 whole-head scans per run using a Sensitivity Encoded (SENSE; Pruessmann et al., 1999) single-shot echo planar imaging technique (repetition time, TR = 2.6 s; echo time, TE = 35 ms; field of view =  $220 \text{ mm} \times 220 \text{ mm} \times 132 \text{ mm}$ ; flip angle =  $78^\circ$ ; matrix size =  $80 \times 80$ ; voxel size =  $2.75 \text{ mm} \times 2.75 \text{ mm} \times 4 \text{ mm}$ , resliced to  $1.72 \text{ mm} \times 1.72 \text{ mm} \times 4 \text{ mm}$ ). Susceptibility artifacts in temporal cortices were reduced by adjusting the slice tilt to  $30^\circ$  from the transverse plane (Weiskopf et al., 2007). Three dummy scans at the beginning of each run were acquired and discarded in order to establish a steady state in T1 relaxation for all functional scans.

### **fMRI data analysis**

Preprocessing and MRI data analysis were performed using MATLAB 2006b (Mathworks Inc., Natick, MA, USA) and the SPM5 software package<sup>3</sup>. All images were realigned to the first recorded volume, normalized into standard stereotactical space (using the EPI-template provided by the Montreal Neurological Institute, MNI brain), resliced to  $2 \text{ mm} \times 2 \text{ mm} \times 2 \text{ mm}$  voxel size and smoothed using a 6-mm full-width-at-half-maximum Gaussian kernel. Activated voxels were identified by a general linear model (Friston et al., 1995) implemented in SPM5. High-pass filtering (cut-off 128 s) was applied to the time series. For each subject, the fMRI responses were modeled with a design matrix using the onset of the second face of each face-pair in each of the six face pair conditions (i.e., within, same, and between for avatar and human trials) and the onset of the target events (for target monitoring) as regressors convolved with SPM’s standard canonical hemodynamic response function (HRF). The parameter estimates of the HRF for each regressor were calculated for each voxel, and linear contrasts were computed for each subject (Friston et al., 1995). These contrasts were entered into a second-level model. All contrast images were first smoothed using a 8-mm FWHM Gaussian kernel to allow inter-subject localization differences to be accounted for, with a final estimated overall smoothing of 10 mm FWHM  $[(6^2 + 8^2)^{1/2}]$ . One-sample random effects *t*-statistics across subjects was performed to allow for population inferences.

The pair-repetition paradigm used the same face pairs and face-pair conditions as described for the preceding perceptual discrimination task. Briefly, the “within,” “same,” and “between” face pair conditions used the face pair morphs M4–M0 (i.e., within), M4–M4 (i.e., same), and M4–M8 (i.e., between) in the avatar trials and M8–M12 (i.e., within), M8–M8 (i.e., same), and M8–M4 (i.e., between) in the human trials (see **Figure 3B**). Contrasts were defined on the basis of these face-pair conditions in order to identify brain regions sensitive to physical and category-related change between the presented prime and target stimuli. The rationale behind this is that (further to the description in the introduction) the neural adaptation effect (Grill-Spector and Malach, 2001) evokes a smaller BOLD signal in response to the second of a pair of identical stimuli than to the second of a pair of dissimilar stimuli. As the targets in the same, within, and between conditions represent dissimilar points along the DOH in terms of physical attributes and/or category, relative differences in signal decrease between the face-pair conditions permits localization of neuron populations responsive to change in physical attributes and/or category along the DOH (see e.g., Murray and Wojciulik, 2004; Grill-Spector et al., 2006; Jiang et al., 2009). As the prime was always M4 in avatar and M8 in human trials, the following contrasts refer to the targets of each face pair condition in the avatar (M0, M4, M8) and human trials (M12, M8, M4). Sensitivity to physical change was detected using the contrast of conditions M0 (i.e., within) plus M8 (i.e., between) > M4 (i.e., same) for avatar trials and M12 (i.e., within) plus M4 (i.e., between) > M8 (i.e., same) for human trials. To detect brain regions selectively responsive to category change across the boundary in the direction avatar-to-human and human-to-avatar, the contrasts M8 (i.e., between) > M4 (i.e., same) plus M0 (i.e., within) for avatar trials and M4 (i.e., between) > M8 (i.e., same) plus M12 (i.e., within) for human trials were used.

Whole-brain voxel-wise analyses was performed. In examining the sensitivity of brain regions to variation in physical attributes, we expected and focused on activations in bilateral regions of the mid-fusiform gyrus (e.g., Jiang et al., 2009; Xu et al., 2009), but for purposes of comparison with other studies (e.g., Rotshstein et al., 2005) all voxels are reported that survived significance thresholding at  $p < 0.001$ , uncorrected for multiple comparisons with a spatial extent threshold of  $k = 5$  voxels.

The hypothesized regions of interest (ROIs) sensitive to category change are described in the introduction. The exploration of further regions associated with affective processing and affective ambiguity focused on the amygdala. Because of this dual approach with predefined ROIs and exploration of additional regions, all clusters of voxels ( $k = 20$  voxels) responsive to category change that survived a threshold of  $p < 0.005$  (uncorrected) are reported and discussed. The more lenient threshold (but more stringent level of contiguous voxels) reflected the exploratory interest in reducing type-2 errors (e.g., Wager et al., 2003; Phelps et al., 2008) coupled with the view that hemodynamic responses in sub-cortical structures of the affect processing network are considered more difficult to detect (e.g., Phelps, 2004; Etkin et al., 2006; Herwig et al., 2007). But, those voxels surviving significance thresholding at  $p < 0.001$ , uncorrected for multiple comparisons with a spatial

<sup>3</sup><http://fil.ion.ucl.ac.uk/spm>

extent threshold of  $k = 20$  voxels, are also reported in the Section “Results.”

## RESULTS

### FORCED CHOICE CLASSIFICATION

The logistic slope value of the fitted regression curve of each individual continuum was highly significant at  $p > 0.001$ . Given the large number of values for all continua, we report the test of this response function across all continua and against zero in a one-sample  $t$ -test. This showed that the response function has a highly significant logistic component ( $t_{31} = 29.28$ ,  $p > 0.001$ ) reflecting a sigmoid step-like function consistent with the presence of a category boundary (Harnad, 1987). The grand mean (see **Figure 2A**) of the fitted logistic curves (and of the response data) across continua shows the sigmoid shaped curve with lower and upper asymptotes of avatar or human categorization responses (as percentages of “different” responses) nearing 100% for avatars and 100% for humans, respectively. Across continua, the mean category boundary value ( $M = 6.30$ ,  $SD = 0.85$ ) corresponds with morph M6 (i.e., the midpoint along the morph continua). This value is derived from the fitted logistic curve and the ordinate midpoint between the lower and upper asymptotes of the categorization responses, indicating therefore that the maximum uncertainty of 50% in categorization judgments was associated with the morph M6.

Morphs M0, M4, M8, and M12 were selected for presentation in the subsequent two studies. The mean percentage of avatar or human identifications across all continua (reported here in terms of the response “human”, as shown in **Figure 2A**) for the morphs M0, M4, M8, and M12 was 2.25 ( $SD = 2.5$ ), 10.25 ( $SD = 4.8$ ), 89.27 ( $SD = 8.4$ ), 98 ( $SD = 2.9$ ), respectively (see **Figure 2B**).

In the interest of gaining an overall picture of categorization response times (RT) along the DOH, an RM-ANOVA with factors continuum (32 continua)  $\times$  morph position (13 levels) and RT as dependent variable was conducted. The analysis revealed no effect for continuum, but there was a main effect for morph position,  $F(2.75, 66.01) = 27.04$ ,  $p < 0.001$ . Consistent with the preceding result of maximum uncertainty in decision responses at morph M6, the mean RT across participants indicated longest RTs for M6. To characterize this more clearly, the mean RT values at M6 were compared with the mean RT values at all other morph positions. A one-way RM-ANOVA analysis with morph position (two levels: M6 versus all other morphs) and RT in seconds as dependent variable collapsed across continua showed that RT for M6 ( $M = 1.42$ ,  $SD = 0.26$ ) differed highly significantly from RT for the other morph positions ( $M = 0.99$ ,  $SD = 0.46$ ),  $F(1,24) = 62.04$ ,  $p < 0.001$ .

As the morphs M0, M4, M8, and M12 were to be used in the subsequent two studies, we tested also for differences in RT between these morphs. A one-way RM-ANOVA analysis with morph position (four levels: M0, M4, M8, and M12) and RT as dependent variable collapsed across continua was conducted. The analysis revealed an effect for morph position,  $F(1.77, 41.32) = 19.19$ ,  $p < 0.001$ . Tests of planned within-subject contrasts showed significant differences in RT within each category such that RT for M0 ( $M = 0.67$  s,  $SD = 0.27$ ) was shorter than for M4 ( $M = 1.05$  s,  $SD = 0.71$ ) [ $F(1,24) = 48.2$ ,  $p < 0.001$ ] and RT for M12 ( $M = 0.91$ ,

$SD = 0.56$ ) was shorter than RT for M8 ( $M = 1.19$  s,  $SD = 0.56$ ) [ $F(1,24) = 17.2$ ,  $p < 0.001$ ], thus indicating quicker responses for morphs at either end of the continua than for the morphs M4 and M8 that straddled the category boundary. But, there was no significant difference in RT between M4 and M8.

### PERCEPTUAL DISCRIMINATION TASK

#### Avatar trials

For the “avatar” trials (see **Figure 3A**), the ANOVA analyses showed a significant main effect on discrimination accuracy of face pair condition, that is, within, same, and between [ $F(2,38) = 149.74$ ,  $p < 0.001$ ] and of ISI [ $F(1,19) = 10.28$ ,  $p = 0.006$ ]. The ISI effect appears to reflect a general bias irrespective of face pair condition toward more “different” responses in the short ISI condition than in the long ISI condition. Pre-planned tests of within-subject contrasts revealed that face pairs that crossed the category boundary (“between” trial type) were significantly more often indicated as different than were those pairs from within a category [ $F(1,19) = 142.89$ ,  $p < 0.001$ ]. There was also a significant effect of discrimination accuracy within the avatar category because pairs from within a category (i.e., “within” trial type) were more frequently indicated to be different than those pairs of the “same” trial type,  $F(1,19) = 6.09$ ,  $p < 0.026$ .

#### Human trials

Similarly, for the “human” trials (see **Figure 3B**), the ANOVA analyses showed a significant effect on discrimination accuracy of face pair condition [ $F(2,38) = 876.46$ ,  $p < 0.001$ ] but no effect of ISI. Pre-planned tests of within-subject contrasts revealed that face pairs that crossed the category boundary (“between” trial type) were significantly more often indicated as different than were those pairs from within a category (“within” trial type),  $F(1,19) = 932.03$ ,  $p < 0.001$ . There was also an effect of discrimination accuracy within the human category because pairs from within a category were significantly more frequently indicated to be different than those pairs of the “same” trial type,  $F(1,19) = 478.52$ ,  $p = 0.28$ .

In summary, the results demonstrate that processes of CP are very likely to influence processing of objects along the DOH. The effect shown in this task was one of better discrimination accuracy for pairs that crossed the category boundary than for equidistant pairs drawn from within a category on one or other side of the boundary, this effect applying similarly for both avatar and human trials.

### TARGET MONITORING TASK AND EVENT-RELATED FMRI

On this basis of the perceptual discrimination task, brain activity was expected to be differentially affected by the within and between category conditions. The ISI condition (short versus long) of the perceptual discrimination task was also carried over to the fMRI study.

#### Sensitivity to physical change

To detect brain regions sensitive to differences in physical features, face pairs in which there was a change in physical features between prime and target were compared with those face pairs in which there was no change in physical features (see **Table 1**).

**Avatar trials.** The contrast “within plus between versus same” was applied to face pairs in which an avatar was the prime, and this revealed right-lateralized increased activations in the fusiform gyrus (BA37) as shown in **Figure 4A**, superior temporal gyrus (BA22), and middle frontal gyrus (BA6). No differences between short and long ISI were found.

**Human trials.** Similarly, the contrast “within plus between versus same” for face pairs in which a human was the prime revealed modulation of activation bilaterally in the fusiform gyrus (BA37) as shown in **Figure 4B**, left precuneus (BA4), right mid-cingulum (BA32), and left insula (BA13). No differences between short and long ISI were found.

### Sensitivity to category change

To detect brain regions selectively responsive to category change across the boundary in the avatar-to-human and human-to-avatar directions along the continuum, face pairs in which there was a change in category between prime and target were compared with the conditions in which there was no such change in category (see **Table 2**).

**Avatar-to-human direction along DOH.** For the avatar-to-human direction, the contrast “between versus same plus within”

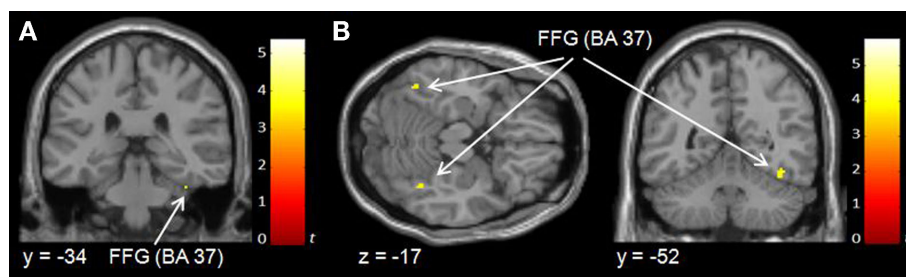
was applied to face pairs in which an avatar was the prime. This showed (at  $p < 0.005$ ) a large cluster that included the left hippocampus (BA28), entorhinal area (BA34), perirhinal area (BA35), and amygdala, and a further cluster in the right mid-insula (BA13) (see **Figure 5A**). The activation in the left hippocampus (BA28) was found also at the higher significance threshold of  $p < 0.001$  [max  $t$  value 4.12 ( $k = 12$ ) at  $-24, -14, -14$ ]. No differences between short and long ISI were found.

**Human-to-avatar direction along DOH.** To establish the pattern of activation associated with brain regions sensitive to crossing the category boundary in the direction of human-to-avatar, the same contrast (“between versus same plus within”) was applied for face pairs in which a human was the prime. This showed a different pattern of brain areas sensitive to category change (at  $p < 0.005$ ), including left and right putamen, left caudate head, right lateral posterior nucleus of thalamus, and left red nucleus (see **Figure 5B**). Both the right thalamus (lateral posterior nucleus) [max  $t$  value 3.88 ( $k = 13$ ) at  $14, -22, 12$ ], left caudate head [max  $t$  value 3.62 ( $k = 5$ ) at  $12, 14, -6$ ] and left putamen [max  $t$  value 3.88 ( $k = 5$ ) at  $-20, 6, 14$ ] were found also at the higher significance threshold of  $p < 0.001$ . No differences between short and long ISI were found.

**Table 1 | Brain areas sensitive to physical change.**

Change in physical qualities versus no change in physical qualities:							
Region of activation	BA	L/R	x	y	z	Max $t$ value	No. of voxels
<b>AVATAR</b>							
Fusiform gyrus	37	R	38	-36	-30	4.01	8
Superior temporal gyrus	22	R	50	-60	14	4.85	16
Middle frontal gyrus	6	R	36	-2	42	4.24	6
<b>HUMAN</b>							
Cerebellum		R	42	-58	-26	4.01	31
*Fusiform gyrus	37	R	42	-52	-18	3.84	
Fusiform gyrus	37	L	-46	-58	-18	3.75	5
Precuneus	4	L	-14	-42	58	5.73	33
Mid-cingulate cortex	32	R	10	18	38	4.70	31
Insula	13	L	-38	20	8	4.01	5

Coordinates in MNI space and maximum  $t$  values are shown for local voxel maxima in each cluster (SPM( $t$ ) maps thresholded at  $p < 0.001$  uncorrected for multiple comparisons, with a cluster extent threshold of 5 voxels). \*, Subpeaks of a cluster; BA, approximate Brodmann area; L, left hemisphere; R, right hemisphere; MNI voxel coordinates; listed brain regions exceeding  $p$  (uncorr.)  $< 0.001$ .



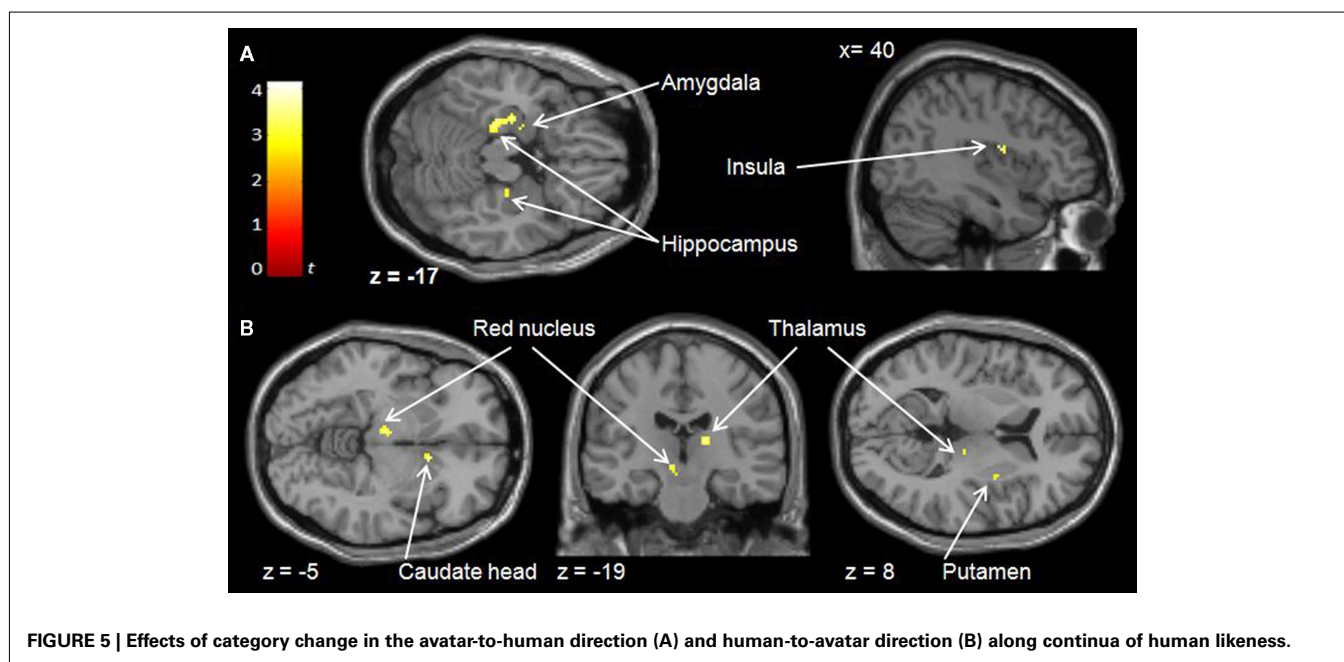
**FIGURE 4 | Effects of physical change on hemodynamic activity during observation of avatar (A) and human (B) trials.**



**Table 2 | Brain areas sensitive to category change.**

Change in category versus no change in category							
Region of activation	BA	L/R	x	y	z	Max t value	No. of voxels
<b>AVATAR PRIME AND HUMAN TARGET</b>							
Hippocampus/entorhinal area	28/34	L	-24	-14	-14	4.12	97
*Perirhinal cortex	35	L	-22	-26	-16	3.45	
*Amygdala		L	-20	-8	20	3.12	
Mid-insula	13	R	40	-5	16	3.54	29
<b>HUMAN PRIME AND AVATAR TARGET</b>							
Putamen		L	20	6	14	3.88	31
Thalamus/lateral posterior nucleus		R	14	-22	12	3.86	54
Caudate head		R	12	14	-6	3.62	23
Red nucleus		R	-6	-16	-6	3.36	24

Coordinates in MNI space and maximum t values are shown for local voxel maxima in each cluster (SPM(t) maps thresholded at  $p < 0.005$  uncorrected for multiple comparisons, with a cluster extent threshold of 20 voxels). \*, Subpeaks of a cluster; BA, approximate Brodmann area; L, left hemisphere; R, right hemisphere; MNI voxel coordinates; listed brain regions exceeding  $p$  (uncorr.)  $< 0.005$ .



## DISCUSSION

The forced choice classification task showed that faces drawn from continua representing the DOH are subjectively assigned to the discrete avatar and human categories. Because morphed faces are computer-generated and modeled, Seyama and Nagayama (2009) suggest that they are artificial and might therefore be processed differently than are natural human faces. But this view does not apply for the outcome of our perceptual category judgments: Morph stimuli regarded as artificial in a technical sense can be explicitly judged to be human. Importantly, the faces assigned to the human category were drawn from different points along the morph continua, thus demonstrating that there is variation in humanlike appearance within the human category. Mori did not consider

this in his hypothesis, focusing instead on the wide variation of potential physical forms of non-human objects along the DOH as often found in uncanny research (e.g., Minato et al., 2006; Walters et al., 2008). Ramey (2005) reflects that these objects are typically designed and modeled after the human image. The implication of this design and research approach to the uncanny valley is that the human image is often treated as a general point of reference irrespective of the fact, as shown in the present study, that there are differences in human likeness within the human category.

Mori was particularly concerned with the threshold between non-human and human object perception and associated uncertainty and discomfort. Defining the non-human–human category boundary is therefore critical to examining Mori’s hypothesis.

The classification response function of all our continua showed a sigmoid shape consistent with the presence of an avatar–human category boundary (Harnad, 1987). The reaction time data indicated significant increases in and a similar degree of decision uncertainty for morphs close to each side of the category boundary. One could reframe Mori's hypothesis in terms of category processing and suggest that any personal discomfort associated with processing human likeness is most likely to occur in response to stimuli at or near the category boundary where there is greatest categorization ambiguity. It is possible that this prediction applies to stimuli either side of the avatar–human category boundary. This awaits investigation.

The perceptual discrimination task showed that a pair of faces located either side of the avatar–human category boundary is easier to distinguish than a pair of faces (with the same degree of physical difference along the DOH) drawn from within a category, thus indicating CP for the DOH. In other words, the cognitive representation of the category structure of the DOH influences the observer's sensitivity for perceptual information relating to human likeness. Evidently, this sensitivity is greatest at or near to the threshold of non-human and human object perception at which small changes in realism (i.e., concerning human or non-human-specifying perceptual features) might evoke uncertainty and discomfort. Taken together, these behavioral findings demonstrate that the DOH may be defined both in terms of a gradual change in the degree of physical humanlike similarity, as described by Mori, as well as in terms of the effects of CP and the dimension's underlying category structure, as proposed in this investigation.

The stimulus conditions described by Mori (i.e., observing novel non-human objects subtly different in physical appearance from that of the human counterpart) were simulated within the constraints of fMRI methodology. The imaging data confirmed that different regions of the brain are responsive to processing change in physical humanlike similarity and processing change in category along the DOH. We consider the findings for physical change first.

The right and left mid-fusiform areas were sensitive to fine-grained change in physical human likeness for human trials and a different right mid-fusiform area was sensitive to physical change for avatar trials. This is consistent with reports in other studies of sensitivity in these areas to facial physical similarity (Xu et al., 2009), to similarity of facial geometry (Jiang et al., 2006, 2009) and to similarity of surface properties of facial texture (e.g., Jiang et al., 2006). Our fMRI data give no indication as to the relative importance of shape and texture in processing faces along the DOH (for face recognition see, e.g., O'Toole et al., 1999; Jiang et al., 2006; Russell et al., 2007). But participants consistently reported attending to the skin texture of avatars in the forced choice classification task. Surface cues are an important diagnostic aid for judging human likeness (MacDorman et al., 2009) and for distinguishing synthetic and natural objects and faces (e.g., Biederman and Ju, 1998; Russell and Sinha, 2007), especially when the objects are highly similar in structure (Price and Humphreys, 1989). The reported attention to skin texture may reflect greater difficulty with or preferential use of specific facial information for forced choice classification of avatars, but this is not clear. As performance in

visual discriminations of facial surface properties and shape does depend on experience (Vuong et al., 2005; Balas and Nelson, 2010), there should at least be differences between the processing of perceptual features of human and novel avatar faces. An additional analysis comparing the avatar and human “within” conditions in our perceptual discrimination task supports this: Fine visual discrimination of avatars was less accurate than that of human faces,  $F(1,19) = 6,31$ ,  $p = 0.02$ . Whether and how the relative importance of structural and textural information changes along the course of the DOH has not been investigated (but see MacDorman et al., 2009).

A change in the category of sequentially presented faces (i.e., prime and target of paired faces) was found to evoke modulations of neural activity in brain regions previously associated with category learning and category uncertainty. The pattern of regions was entirely different depending on whether the target was an avatar in the human–avatar pairs or human in the avatar–human pairs. This shows that the direction of change along the DOH does influence the way in which humanlike faces are processed, at least within our paradigm. The dorsal and ventral striatum (putamen and head of caudate), thalamus, and red nucleus were responsive to the avatar target and the MTL (hippocampus, entorhinal, and perirhinal areas), mid-insula, and amygdala were responsive to the human target.

The results for the striatum and MTL were the most prominent. These regions are thought to have different roles in and an antagonistic relationship during category processing and learning (e.g., Seger and Cincotta, 2005; for overview see Poldrack and Rodriguez, 2004). For example, relative activity in the MTL memory system subsides and that of the striatal memory system increases during categorization training (Poldrack et al., 1999, 2001). Alternatively, these and other memory systems (e.g., Poldrack and Foerde, 2008) might actually work in concert during category processing (Cincotta and Seger, 2007). Irrespective of their possible interplay, the differential response of these regions to the direction of category change between prime and target (i.e., human–avatar or avatar–human) suggests that the avatar and human faces represent different categorization problems that require, at least in part, dissimilar processes or strategies to resolve them. There is some support for this in the behavioral data from the perceptual discrimination task. An additional analysis compared the avatar and human “between” conditions and found a significantly higher proportion of judgment errors [ $F(1,19) = 8,69$ ,  $p < 0.01$ ] for avatar–human pairs than for human–avatar pairs.

In the pair-repetition paradigm, processing of the target stimulus is biased by the implicit memory of the preceding prime that represents a different “expected” or “predicted” category than is actually shown in the target. The caudate head (responsive to avatar targets) is sensitive to prediction error in various tasks (e.g., O'Doherty et al., 2003; King-Casas et al., 2005; Bray and O'Doherty, 2007; Jensen et al., 2007). Besides contributing to learning stimulus–category associations (Seger and Cincotta, 2005, 2006) and selecting and executing motor responses for signaling category membership (Williams and Eskandar, 2006; Seger et al., 2010), the putamen (also responsive to avatar targets) is associated with processing prediction error (den Ouden et al.,

2010). It is possible that these striatal regions contribute to processing the deviation of the novel avatar target from the human prime and do so on the basis of the prime and well defined category representations of human faces. If this is true, one would expect the “between” condition of the perceptual discrimination task to be easier for the human–avatar pairs than for the avatar–human pairs for which corresponding category representations are less well defined. The preceding additional analysis of the “between” condition supports this. Predictive processing has been investigated explicitly under conditions of decision uncertainty and shown to modulate activity in the ventral striatum (and the thalamus and red nucleus as also found in response to our targets of the avatar–human face pairs; Volz et al., 2003). Filoteo et al. (2005) suggest that the caudate head might be involved in switching between potential category rules used to establish category membership, while Grinband et al. (2006) found activations highly similar to ours in their study of categorization ambiguity, and they suggest that the ventral striatum signals adjustment of the represented categorical boundary in order to minimize errors.

Medial temporal lobe structures (responsive to human targets) are involved in processing novel and familiar stimuli, novelty detection, and uncertainty coding (e.g., Stern et al., 1996; Stark and Squire, 2000; Rutishauser et al., 2006; Vanni-Mercier et al., 2009). MTL is thought to encode and retain individual category instances of novel stimuli (Poldrack et al., 1999; Seger and Cincotta, 2005) and to facilitate further category-related processing of features by other systems including the basal ganglia (Meeter et al., 2008). Given our paired presentation of similar human and novel avatar faces, it is noteworthy that processes of novelty detection in the hippocampus are considered to entail processing of the deviation of actual sensory input from expected input (e.g., Kumaran and Maguire, 2006). This would be consistent with the notion that the processing of our human targets was strongly guided by representations of the preceding input from the avatar primes. The hippocampus is involved in visual categorization and perceptual learning, while impaired processing of faces is found only when both hippocampus and perirhinal cortex are damaged (Graham et al., 2006). The perirhinal cortex played a clear role in response to category change between human primes and avatar targets. The perirhinal cortex appears to be critical for object memory, contributing to resolving complex visual discriminations and those of high ambiguity (Bussey et al., 2003) when both novel and familiar stimuli share visual features (Barens et al., 2005).

Consistent with its role in decision making under conditions of category processing and uncertainty (e.g., Volz et al., 2003; Hsu et al., 2005; Grinband et al., 2006; Heekeren et al., 2008; Fleming et al., 2010), we expected the anterior rather than the mid-insular cortex to be responsive to category change. The mid-insular cortex was responsive to human targets. The role of this region in the present study and generally (Wager and Feldman Barrett, 2004) is not clear. The mid-insular is responsive to affective ambiguity, anticipation of emotionally aversive visual stimuli, and to physical pain (Simmons et al., 2008), this reflecting the processing of negative expectation particularly in the context of pain (e.g., Petrovic et al., 2000). Our participants were not

exposed to noxious stimulation, but similar neural mechanisms are thought to mediate the experience of pain and feelings of personal discomfort (Price, 2000). Right mid (Lutz et al., 2009) and anterior insula activity (Critchley et al., 2002) is also associated with measures of arousal (heart rate and galvanic skin responses). Given the suggested role of prediction error processing, Paulus and Stein (2006) proposed anxiety mechanism is interesting. Applied here, the deviation between the expected arousal state on the basis of the avatar prime and the actual arousal associated with the unexpected human target might be interpreted as signaling the presence of uncertainty, threat, or potential threat. Uncanny-related arousal and discomfort is suggested to be rooted in mechanisms of threat avoidant behavior (MacDorman, 2005a; Green et al., 2008). Interestingly, Gray and Critchley (2007) and Wager et al. (2003) associate mid-insular representations with a threat-related component and avoidant behavior. Alternatively, arousal and mid-insula responsiveness to category change may be more closely related to the role of the right anterior insula during category processing in marshaling attentional resources to enhance performance during categorization under conditions of uncertainty (Heekeren et al., 2008).

In exploring additional brain areas otherwise associated with affective processing, the amygdala was found to be responsive to category change in avatar–human face pairs. The amygdala is responsive to natural and computer-generated human faces (e.g., Todorov and Engell, 2008), humanlike but unnatural faces (Rotshtein et al., 2001) novelty, uncertainty, unclear predictive value, and ambiguous valence (e.g., Phelps and LeDoux, 2005; Herwig et al., 2007; Levy et al., 2010; Neta and Whalen, 2010). Whether and how appraisal of affective valence (and prospective outcomes in terms of potential threat or reward) contributes to category processing is not clear, though prospective outcomes are thought to interact with and bias processes of perceptual categorization (Heekeren et al., 2008; see also Gupta et al., 2011). In view of the different responses of the striatum and MTL in this study, Seger and Miller’s (2010) suggestion is interesting. They propose that the amygdala might play a role during category processing in altering the balance between the memory systems of these two regions depending on the affective meaning of a situation. For example, modulation of the amygdala might reflect its involvement in processing of novel category information in conjunction with the MTL and enhancement of attentional and memory processing of the novel avatar primes (Kleinbans et al., 2007). Alternatively, the amygdala might have been responsive to the “unexpected” category of the targets (i.e., the human faces; Roesch et al., 2010).

This investigation demonstrates that the definition of the DOH in terms of the degree of object similarity to human appearance does not reflect the way in which human likeness is subjectively perceived along this dimension. The forced choice categorization task showed that there is variation in human likeness within the human category, and the perceptual discrimination task confirmed that processes of CP and associated representations of the DOH’s category structure influence the perception of human likeness. Given the inconsistent findings in “uncanny” research to date, careful definition of the category boundary especially in studies

using morph continua is therefore important. It is likely that mechanisms associated with category processing and the representation of human likeness also influence affective experience. For example, increased uncertainty in forced choice category decisions might be associated with discomfort (assuming that Mori's hypothesis is correct), though any discomfort need not be limited to the non-human side of the category boundary along the DOH. The present behavioral findings suggest that the general conception of the DOH should be revised. As Mori did not embed his ideas in any psychological framework, we suggest that his hypothesis should be considered in terms of the well-established psychological empirical-theoretical framework of category perception and learning. This could prove useful in examining where along the DOH negative or uncanny experience is most likely to occur for a given set of stimuli and in throwing light on the potential role of categorization ambiguity in evoking negative affect. It is of course possible that any association between category ambiguity and negative affect is not specific to humanlike stimuli along the DOH but is a more general feature of category ambiguity itself. The framework set out in this study could also be applied to other fields of related research to understand how variously realistic humanlike characters influence in association with category ambiguity for example the experience of presence (i.e., immersive experience) in virtual reality (e.g., Brenton et al., 2008), audience persuasion and identification with fictive characters in narrative

(e.g., Cohen, 2001), communication in educational virtual environments (e.g., Fabri et al., 2004), or consumer trust in electronic commerce (e.g., Bauer and Neumann, 2005). For example, greater category ambiguity might enhance attentional processing of the human or non-human-specifying perceptual features of digital human representations, and this might render it less likely that a person responds to the character as if it were in some way real, experiences identification, processes the media content as the designers intend, or develops a sense of trust. Using this framework, the fMRI study showed that different brain regions are responsive to the direction along the DOH in which there is a change in the category of observed objects (i.e., avatar-to-human or human-to-avatar). This appears to reflect the impact of differences in category knowledge for avatar and human faces. Re-examination of this effect in the context of category training and experience with avatars might throw more light on this. Replicating the behavioral findings of this study with a larger number of participants and with stimuli presenting biological motion (see Chaminade et al., 2007; Saygin et al., 2011) would reinforce our findings. Whether there are gender differences remains to be investigated.

## ACKNOWLEDGMENTS

This work was supported by the European Union FET Integrated Project PRESENCIA (Contract number 27731).

## REFERENCES

- Angeli, A., Davidoff, J., and Valentine, T. (2008). Face familiarity, distinctiveness, and categorical perception. *Q. J. Exp. Psychol. (Hove)* 61, 690–707.
- Annett, M. (1970). A classification of hand preference by association analysis. *Br. J. Psychol.* 61, 303–321.
- Ashby, F. G., and Maddox, W. T. (2005). Human category learning. *Annu. Rev. Psychol.* 56, 149–178.
- Aylett, R. S. (2004). Agents and affect: why embodied agents need affective systems. *Lect. Notes Comput. Sci.* 3025, 496–504.
- Balas, B., and Nelson, C. A. (2010). The role of face shape and pigmentation in other-race face perception: an electrophysiological study. *Neuropsychologia* 48, 498–506.
- Barense, M., Bussey, T., Lee, A., Rogers, T. T., Hodges, J. R., Saksida, L., Murray, E., and Graham, K. (2005). Feature ambiguity influences performance on novel object discriminations in patients with damage to perirhinal cortex. *Cognitive Neuroscience Society Annual Meeting Program 2005*, New York, NY, 125–116.
- Bauer, H. H., and Neumann, M. (2005). "Investigating the Effects of Avatars as Virtual Representatives in Electronic Commerce," in *Proceedings of the Australian-New Zealand Marketing Academy Conference (ANZMAC)*, Perth, 18.
- Biederman, I., and Ju, G. (1998). Surface vs. edge-based determinants of visual recognition. *Cogn. Psychol.* 20, 38–64.
- Bray, S., and O'Doherty, J. (2007). Neural coding of reward-prediction error signals during classical conditioning with attractive faces. *J. Neurophysiol.* 97, 3036–3045.
- Brenton, H., Gillies, M., Ballin, D., and Chatting, D. (2008). The uncanny valley: does it exist? *Paper presented at the Proceedings of the Human-Animated Characters Interaction, HCI 2005: The Bigger Picture*, Edinburgh.
- Bussey, T. J., Saksida, L. M., and Murray, E. A. (2003). Impairments in visual discrimination after perirhinal cortex lesions: testing "declarative" versus "perceptual-mnemonic" views of perirhinal cortex function. *Eur. J. Neurosci.* 17, 649–660.
- Campbell, R., Pascalis, O., Coleman, M., Wallace, S. B., and Benson, P. J. (1987). Are faces of different species perceived categorically by human observers? *Proc. R. Soc. Lond. B Biol. Sci.* 264, 1429–1434.
- Chaminade, T., Hodgins, J., and Kawato, M. (2007). Anthropomorphism influences perception of computer-animated characters' actions. *Soc. Cogn. Affect. Neurosci.* 2, 206–216.
- Cincotta, C. M., and Seger, C. A. (2007). Dissociation between striatal regions while learning to categorize via feedback and via observation. *J. Cogn. Neurosci.* 19, 249–265.
- Cohen, J. (2001). Defining identification: a theoretical look at the identification of audiences with media characters. *Mass Commun. Soc.* 4, 245–264.
- Critchley, H. D., Mathias, C. J., and Dolan, R. J. (2002). Fear conditioning in humans: the influence of awareness and autonomic arousal on functional neuroanatomy. *Neuron* 33, 653–663.
- den Ouden, H. E., Daunizeau, J., Roiser, J., Friston, K. J., and Stephan, K. E. (2010). Striatal prediction error modulates cortical coupling. *J. Neurosci.* 30, 3210–3219.
- Etkin, A., Egner, T., Peraza, D. M., Kandel, E. R., and Hirsch, J. (2006). Resolving emotional conflict: a role for the rostral anterior cingulate cortex in modulating activity in the amygdala. *Neuron* 51, 871–882.
- Fabri, M., Moor, D., and Hobbs, D. (2004). Mediating the expression of emotion in educational collaborative virtual environments: an experimental study. *Int. J. Virtual Reality* 7, 66–81.
- Filoteo, J. V., Maddox, W. T., Salmon, D. P., and Song, D. D. (2005). Information-integration category learning in patients with striatal dysfunction. *Neuropsychology* 19, 212–222.
- Fleming, S. M., Whiteley, L., Hulme, O. J., Sahani, M., and Dolan, R. J. (2010). Effects of category-specific costs on neural systems for perceptual decision-making. *J. Neurophysiol.* 103, 3238–3247.
- Friston, K. J., Holmes, A. P., Worsley, K. J., Poline, J. P., Frith, C. D., and Frackowiak, R. S. J. (1995). Statistical parametric maps in functional imaging: a general linear approach. *Hum. Brain Mapp.* 2, 189–210.
- Graham, K. S., Scahill, V. L., Hornberger, M., Barense, M. D., Lee, A. C., Bussey, T. J., and Saksida, L. M. (2006). Abnormal categorization and perceptual learning in patients with hippocampal damage. *J. Neurosci.* 26, 7547–7554.
- Gray, M. A., and Critchley, H. D. (2007). Interoceptive basis to craving. *Neuron* 54, 183–186.
- Green, R. D., MacDorman, K. F., Ho, C.-C., and Vasudevan, S. K. (2008). Sensitivity to the proportions of faces that vary in human likeness. *Comput. Hum. Behav.* 24, 2456–2474.
- Grill-Spector, K., Henson, R., and Martin, A. (2006). Repetition and the brain: neural models of stimulus-specific effects. *Trends Cogn. Sci. (Regul. Ed.)* 10, 14–23.

- Grill-Spector, K., and Malach, R. (2001). fMR-adaptation: a tool for studying the functional properties of human cortical neurons. *Acta Psychol. (Amst.)* 107, 293–321.
- Grinband, J., Hirsch, J., and Ferrera, V. P. (2006). A neural representation of categorization uncertainty in the human brain. *Neuron* 49, 757–763.
- Gupta, R., Kosciak, T. R., Bechara, A., and Tranel, D. (2011). The amygdala and decision-making. *Neuropsychologia* 49, 760–766.
- Hanson, D. (2006). Exploring the aesthetic range for humanoid robots. *Paper Presented at the Proceedings of the ICCS/CogSci-2006 Long Symposium: Toward Social Mechanisms of Android Science*, Vancouver.
- Harnad, S. R. (1987). “Introduction: psychological and cognitive aspects of categorical perception: a critical overview,” in *Categorical Perception: The Groundwork of Cognition*, ed. S. Harnad (New York: Cambridge University Press), 1–25.
- Heekeren, H. R., Marrett, S., and Ungerleider, L. G. (2008). The neural systems that mediate human perceptual decision making. *Nat. Rev. Neurosci.* 9, 467–479.
- Henson, R. N. (2003). Neuroimaging studies of priming. *Prog. Neurobiol.* 70, 53–81.
- Herwig, U., Kaffenberger, T., Baumgartner, T., and Jancke, L. (2007). Neural correlates of a “pessimistic” attitude when anticipating events of unknown emotional valence. *Neuroimage* 34, 848–858.
- Ho, C.-C., and MacDorman, K. (2010). Revisiting the uncanny valley theory: developing and validating an alternative to the Godspeed indices. *Comput. Hum. Behav.* 26, 1508–1518.
- Hsu, M., Bhatt, M., Adolphs, R., Tranel, D., and Camerer, C. F. (2005). Neural systems responding to degrees of uncertainty in human decision-making. *Science* 310, 1680–1683.
- Jensen, J., Smith, A. J., Willeit, M., Crawley, A. P., Mikulis, D. J., Vitcu, I., and Kapur, S. (2007). Separate brain regions code for salience vs. valence during reward prediction in humans. *Hum. Brain Mapp.* 28, 294–302.
- Jiang, F., Dricot, L., Blanz, V., Goebel, R., and Rossion, B. (2009). Neural correlates of shape and surface reflectance information in individual faces. *Neuroscience* 163, 1078–1091.
- Jiang, X., Bradley, E., Rini, R. A., Zeffiro, T., Vanmeter, J., and Riesenhuber, M. (2007). Categorization training results in shape- and category-selective human neural plasticity. *Neuron* 53, 891–903.
- Jiang, X., Rosen, E., Zeffiro, T., Vanmeter, J., Blanz, V., and Riesenhuber, M. (2006). Evaluation of a shape-based model of human face discrimination using fMRI and behavioral techniques. *Neuron* 50, 159–172.
- King-Casas, B., Tomlin, D., Anen, C., Camerer, C. F., Quartz, S. R., and Montague, P. R. (2005). Getting to know you: reputation and trust in a two-person economic exchange. *Science* 308, 78–83.
- Kleinhans, N. M., Johnson, L. C., Mahurin, R., Richards, T., Stegbauer, K. C., Greenon, J., Dawson, G., and Aylward, E. (2007). Increased amygdala activation to neutral faces is associated with better face memory performance. *Neuroreport* 18, 987–991.
- Kumaran, D., and Maguire, E. A. (2006). An unexpected sequence of events: mismatch detection in the human hippocampus. *PLoS Biol.* 4, e424. doi:10.1371/journal.pbio.0040424
- Levy, I., Snell, J., Nelson, A. J., Rustichini, A., and Glimcher, P. W. (2010). Neural representation of subjective value under risk and ambiguity. *J. Neurophysiol.* 103, 1036–1047.
- Li, S., Mayhew, S. D., and Kourtzi, Z. (2009). Learning shapes the representation of behavioral choice in the human brain. *Neuron* 62, 441–452.
- Lutz, A., Greischar, L. L., Perlman, D. M., and Davidson, R. J. (2009). BOLD signal in insula is differentially related to cardiac function during compassion meditation in experts vs. novices. *Neuroimage* 47, 1038–1046.
- MacDorman, K. (2005a). “Mortality salience and the uncanny valley,” in *IEEE-RAS International Conference on Humanoid Robots* (Tsukuba: IEEE), 339–405.
- MacDorman, K. (2005b). “Androids as an Experimental Apparatus: Why is There an Uncanny Valley and Can We Exploit it?” in *Toward Social Mechanisms of Android Science: A CogSci 2005 Workshop*, Stresa, 106–118.
- MacDorman, K. (2006). Subjective ratings of robot video clips for human likeness, familiarity, and eeriness: an exploration of the uncanny valley. *Paper presented at the ICCS/CogSci-2006 Long Symposium: Toward Social Mechanisms of Android Science*, Vancouver.
- MacDorman, K., Green, R., Ho, C.-C., and Koch, C. (2009). Too real for comfort? Uncanny responses to computer generated faces. *Comput. Hum. Behav.* 25, 695–710.
- MacDorman, K., and Ishiguro, H. (2006). The uncanny advantage of using androids in cognitive science research. *Interact. Stud.* 7, 297–337.
- Maddox, W. T., and Ashby, F. G. (2004). Dissociating explicit and procedural-learning based systems of perceptual category learning. *Behav. Processes* 66, 309–332.
- Meeter, M., Radics, G., Myers, C. E., Gluck, M. A., and Hopkins, R. O. (2008). Probabilistic categorization: how do normal participants and amnesic patients do it? *Neurosci. Biobehav. Rev.* 32, 237–248.
- Minato, T., Shimada, M., Itakura, S., Lee, K., and Ishiguro, H. (2006). Evaluating the human likeness of an android by comparing gaze behaviors elicited by the android and a person. *Adv. Robot.* 20, 1147–1163.
- Mori, M. (1970). Bukimi no tani [the uncanny valley]. *Energy* 7, 33–35.
- Murray, S. O., and Wojciulik, E. (2004). Attention increases neural selectivity in the human lateral occipital complex. *Nat. Neurosci.* 7, 70–74.
- Neta, M., and Whalen, P. J. (2010). The primacy of negative interpretations when resolving the valence of ambiguous facial expressions. *Psychol. Sci.* 21, 901–907.
- O’Doherty, J. P., Dayan, P., Friston, K., Critchley, H., and Dolan, R. J. (2003). Temporal difference models and reward-related learning in the human brain. *Neuron* 38, 329–337.
- O’Toole, A. J., Vetter, T., and Blanz, V. (1999). Three-dimensional shape and two-dimensional surface reflectance contributions to face recognition: an application of three-dimensional morphing. *Vision Res.* 39, 3145–3155.
- Pastore, R. E. (1987). “Categorical perception: some psychophysical models,” in *Categorical perception: The groundwork of Cognition*, ed. S. Harnad (New York: Cambridge University Press), 29–52.
- Paulus, M. P., and Stein, M. B. (2006). An insular view of anxiety. *Biol. Psychiatry* 60, 383–387.
- Petrovic, P., Pettersson, K. M., Ghatan, P. H., Stone-Elander, S., and Ingvar, M. (2000). Pain-related cerebral activation is altered by a distracting cognitive task. *Pain* 85, 19–30.
- Phelps, E. A. (2004). Human emotion and memory: Interactions of the amygdala and hippocampal complex. *Curr. Opin. Neurobiol.* 14, 198–202.
- Phelps, E. A., and LeDoux, J. E. (2005). Contributions of the amygdala to emotion processing: from animal models to human behavior. *Neuron* 48, 175–187.
- Phelps, E. A., O’Connor, K. J., Gatenby, J. C., Gore, J. C., Grillon, C., and Davis, M. (2008). Activation of the left amygdala to a cognitive representation of fear. *Neurosciences* 4, 37–41.
- Poldrack, R. A., Clark, J., Pare-Blagoev, E. J., Shohamy, D., Creso Moyano, J., Myers, C., and Gluck, M. A. (2001). Interactive memory systems in the human brain. *Nature* 414, 546–550.
- Poldrack, R. A., and Forde, K. (2008). Category learning and the memory systems debate. *Neurosci. Biobehav. Rev.* 32, 197–205.
- Poldrack, R. A., Prabhakaran, V., Seger, C. A., and Gabrieli, J. D. (1999). Striatal activation during acquisition of a cognitive skill. *Neuropsychology* 13, 564–574.
- Poldrack, R. A., and Rodriguez, P. (2004). How do memory systems interact? Evidence from human classification learning. *Neurobiol. Learn. Mem.* 82, 324–332.
- Pollick, F. (2009). “In search of the Uncanny Valley,” in *Lecture notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering, UC Media 2009*, eds P. Daras and O. M. Ibarra (Venice: Springer), 69–78.
- Price, D. D. (2000). Psychological and neural mechanisms of the affective dimension of pain. *Science* 288, 1769–1772.
- Price, C. J., and Humphreys, G. W. (1989). The effects of surface detail on object categorization and naming. *Q. J. Exp. Psychol. A.* 41, 797–827.
- Pruessmann, K. P., Weiger, M., Scheidegger, M. B., and Boesinger, P. (1999). SENSE: Sensitivity encoding for fast MRI. *Magn. Reson. Med.* 42, 952–962.
- Ramey, C. H. (2005). The uncanny valley of similarities concerning abortion, baldness, heaps of sand, and humanlike robots. *Paper Presented at the Proceedings of Views of the Uncanny Valley Workshop: IEEE-RAS International Conference on Humanoid Robots*, Tsukuba.
- Roesch, M. R., Calu, D. J., Esber, G. R., and Schoenbaum, G. (2010). All that glitters. Dissociating attention and outcome expectancy from prediction errors signals. *J. Neurophysiol.* 104, 587–595.
- Rotshtein, P., Henson, R. N., Treves, A., Driver, J., and Dolan, R. J. (2005). Morphing Marilyn into Maggie dissociates physical and identity face representations in the brain. *Nat. Neurosci.* 8, 107–113.



- Rotshtein, P., Malach, R., Hadar, U., Graif, M., and Hendler, T. (2001). Feeling or features: different sensitivity to emotion in high-order visual cortex and amygdala. *Neuron* 32, 747–757.
- Russell, R., Biederman, I., Nederhouser, M., and Sinha, P. (2007). The utility of surface reflectance for the recognition of upright and inverted faces. *Vision Res.* 47, 157–165.
- Russell, R., and Sinha, P. (2007). Real-world face recognition: the importance of surface reflectance properties. *Perception* 36, 1368–1374.
- Rutishauser, U., Mamelak, A. N., and Schuman, E. M. (2006). Single-trial learning of novel stimuli by individual neurons of the human hippocampus-amygdala complex. *Neuron* 49, 805–813.
- Saygin, A. P., Chaminade, T., Ishiguro, H., Driver, J., and Frith, C. (2011). The thing that should not be: predictive coding and the uncanny valley in perceiving human and humanoid robot actions. *Soc. Cogn. Affect. Neurosci.* doi: 10.1093/scan/nsr025
- Schultz, J., and Pilz, K. S. (2009). Natural facial motion enhances the cortical responses to faces. *Exp. Brain Res.* 194, 465–475.
- Seger, C. A., and Cincotta, C. M. (2005). The roles of the caudate nucleus in human classification learning. *J. Neurosci.* 25, 2941–2951.
- Seger, C. A., and Cincotta, C. M. (2006). Dynamics of frontal, striatal, and hippocampal systems during rule learning. *Cereb. Cortex* 16, 1546–1555.
- Seger, C. A., and Miller, E. K. (2010). Category learning in the brain. *Annu. Rev. Neurosci.* 33, 203–219.
- Seger, C. A., Peterson, E. J., Cincotta, C. M., Lopez-Paniagua, D., and Anderson, C. W. (2010). Dissociating the contributions of independent corticostriatal systems to visual categorization learning through the use of reinforcement learning modeling and Granger causality modeling. *Neuroimage* 50, 644–656.
- Seyama, J., and Nagayama, R. S. (2007). The uncanny valley: effect of realism on the impression of artificial human faces. *Presence (Camb.)* 16, 337–351.
- Seyama, J., and Nagayama, R. S. (2009). Probing the uncanny valley with the eye size aftereffect. *Presence-Teleop. Virt.* 18, 321–339.
- Simmons, A., Matthews, S. C., Paulus, M. P., and Stein, M. B. (2008). Intolerance of uncertainty correlates with insula activation during affective ambiguity. *Neurosci. Lett.* 430, 92–97.
- Stark, C. E., and Squire, L. R. (2000). Functional magnetic resonance imaging (fMRI) activity in the hippocampal region during recognition memory. *J. Neurosci.* 20, 7776–7781.
- Stern, C. E., Corkin, S., Gonzalez, R. G., Guimaraes, A. R., Baker, J. R., Jennings, P. J., Carr, C. A., Sugiura, R. M., Vedantham, V., and Rosen, B. R. (1996). The hippocampal formation participates in novel picture encoding: evidence from functional magnetic resonance imaging. *Proc. Natl. Acad. Sci. U.S.A.* 93, 8660–8665.
- Studdert-Kennedy, M., Lineman, A. M., Hassis, K. S., and Cooper, F. R. (1970). Motor theory of speech perception: a reply to Lane's critical review. *Psychol. Rev.* 77, 234–249.
- Tinwell, A., and Grimshaw, M. (2009). "Bridging the uncanny: an impossible traverse?" in *Proceedings of the 13th International MindTrek Conference: Everyday Life in the Ubiquitous Era*, eds. O. Sotamaa, A. Lugmayr, H. Franssila, P. Näränen, and J. Vanhala (Tampere: ACM), 66–73.
- Tinwell, A., Grimshaw, M., and Williams, A. (2010). Uncanny behaviour in survival horror games. *J. Gaming Virtual Worlds* 2, 3–25.
- Todorov, A., and Engell, A. (2008). The role of the amygdala in implicit evaluation of emotionally neutral faces. *Soc. Cogn. Affect. Neurosci.* 3, 303–312.
- Vanni-Mercier, G., Mauguire, F., Isnard, J., and Dreher, J. C. (2009). The hippocampus codes the uncertainty of cue-outcome associations: an intracranial electrophysiological study in humans. *J. Neurosci.* 29, 5287–5294.
- Volz, K. G., Schubotz, R. I., and von Cramon, D. Y. (2003). Predicting events of varying probability: uncertainty investigated by fMRI. *Neuroimage* 19(Pt 1), 271–280.
- Vuilleumier, P. (2005). How brains beware: neural mechanisms of emotional attention. *Trends Cogn. Sci. (Regul. Ed.)* 9, 585–594.
- Vuong, Q. C., Peissig, J. J., Harrison, M. C., and Tarr, M. J. (2005). The role of surface pigmentation for recognition revealed by contrast reversal in faces and Greebles. *Vision Res.* 45, 1213–1223.
- Wager, T. D., and Feldman Barrett, L. (2004). From affect to control: functional specialization of the insula in motivation and regulation. Available at: [http://www.affective-science.org/pubs/2004/Wager\\_Edfest\\_submitted\\_copy.pdf](http://www.affective-science.org/pubs/2004/Wager_Edfest_submitted_copy.pdf) [Published online at PsycExtra].
- Wager, T. D., Phan, K. L., Liberzon, I., and Taylor, S. F. (2003). Valence, gender, and lateralization of functional brain anatomy in emotion: a meta-analysis of findings from neuroimaging. *Neuroimage* 19, 513–531.
- Wages, R., Grünvogel, S. M., and Grützmaker, B. (2004). How realistic is realism? Considerations on the aesthetics of computer games. *Lect. Notes Comput. Sci.* 3166, 216–225.
- Walters, M. L., Syrdal, D. S., Dautenhahn, K., te Boekhorst, R., and Koay, K. L. (2008). Avoiding the uncanny valley: robot appearance, personality and consistency of behavior in an attention-seeking home scenario for a robot companion. *Auton. Robots* 24, 159–178.
- Weiskopf, N., Hutton, C., Joseph, O., Turner, R., and Deichmann, R. (2007). Optimized EPI for fMRI studies of the orbitofrontal cortex: compensation of susceptibility-induced gradients in the readout direction. *MAGMA* 20, 39–49.
- Williams, Z. M., and Eskandar, E. N. (2006). Selective enhancement of associative learning by microstimulation of the anterior caudate. *Nat. Neurosci.* 9, 562–568.
- Xu, X., Yue, X., Lescroart, M. D., Biederman, I., and Kim, J. G. (2009). Adaptation in the fusiform face area (FFA): image or person? *Vision Res.* 49, 2800–2807.
- Yamamoto, K., Tanaka, S., Kobayashi, H., Kozima, H., and Hashiya, K. (2009). A non-humanoid robot in the "uncanny valley": experimental analysis of the reaction to behavioral contingency in 2-3 year old children. *PLoS ONE* 4, e6974. doi:10.1371/journal.pone.0006974

**Conflict of Interest Statement:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Received: 18 March 2011; paper pending published: 01 July 2011; accepted: 13 October 2011; published online: 24 November 2011.

Citation: Cheetham M, Suter P and Jäncke L (2011) The human likeness dimension of the "uncanny valley hypothesis": behavioral and functional MRI findings. *Front. Hum. Neurosci.* 5:126. doi: 10.3389/fnhum.2011.00126

Copyright © 2011 Cheetham, Suter and Jäncke. This is an open-access article subject to a non-exclusive license between the authors and Frontiers Media SA, which permits use, distribution and reproduction in other forums, provided the original authors and source are credited and other Frontiers conditions are complied with.