

Numbers that are highlighted were directly verified against R code.

Uncanny but not confusing

Uncanny but not confusing: Multisite study of perceptual category confusion in the Uncanny Valley

Lead authors: Maya B. Mathur^{1,2*} & David B. Reichling³

Contributing authors: XXX, XXX, XXX, ...

¹Department of Epidemiology, Harvard University, Boston, MA, USA

²Quantitative Sciences Unit, Stanford University, Palo Alto, CA, USA

³University of California at San Francisco (ret.), San Francisco, CA, USA

Abstract (149/150 words)

Android robots that are close, but imperfect, likenesses of humans can provoke negative feelings of dislike and eeriness in humans (“Uncanny Valley” effect). We investigated whether category confusion between the perceptual categories of “robot” and “human” contributes to Uncanny Valley aversion. Using a novel, validated corpus of 182 images of real robot and human faces, we precisely estimated the shape of the Uncanny Valley and the location of the perceived robot/human boundary. To implicitly measure confusion, we tracked 358 subjects’ mouse trajectories as they categorized the faces. We observed a clear Uncanny Valley and a pattern of categorization supporting a perceived categorical boundary. Yet, in contrast to predictions of the category confusion mechanism hypothesis, the Uncanny Valley and category boundary locations did not coincide, and mediation analyses further failed to support a causal role of category confusion. These results suggest category confusion does not explain the Uncanny Valley effect.

1. INTRODUCTION

Android robots have rapidly entered our social sphere. We now entrust them with providing therapy to children with autism and older adults, coaching patients on health behavior change, and collaborating with astronauts in space stations (Rabbitt et al., 2015; Weisberger, 2018). Yet human-robot interactions can be fraught with social peril. In particular, robots that closely resemble humans but are not perfectly human-like can elicit unexpectedly negative emotional reactions in human viewers, jeopardizing the success of social robots. This “Uncanny Valley” effect (Mori, 1970) has dominated discussion of human reactions to anthropomorphic robots in both popular culture and research literature. Specifically, the Uncanny Valley theory posits that as android robots increasingly resemble humans, their likability increases until a point at which it abruptly drops to a negative value because the robots become disliked and eerie (Figure 1). Then, as the robots’ human-likeness continues to increase past this “Uncanny Valley”, they again become likable and eventually reach maximum likability as they become indistinguishable from humans. Recent studies have strongly suggested that the Uncanny Valley does occur in real android robots that were intentionally designed to interact with humans (Mathur & Reichling (2016); Slijkhuis (2017); Lischetzke et al. (2017); Jung & Cho (2018)).

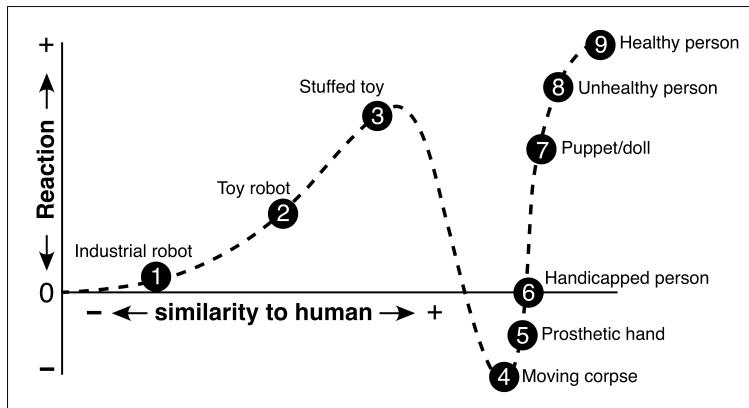


Figure 1: *The hypothesized Uncanny Valley, adapted from Mori (1970).*

With android robots increasingly becoming everyday technology, there is a pressing need

to understand the underlying causes and psychological mechanisms of the Uncanny Valley effect. Some hypotheses posit relatively high-level social and affective mechanisms: robots in the Uncanny Valley might prime awareness of one’s own mortality (Ho et al., 2008), prompt dehumanization responses similar to that directed at human targets of prejudice (Wang et al., 2015), or trigger impulses to avoid pathogens (MacDorman et al., 2009). In contrast, according to a longstanding lower-level hypothesis known as “category confusion”, there is a deep-seated perceptual distinction between the categories “human” and “non-human”, and androids that are difficult to categorize as one or the other cause aversion as a direct result of the categorization difficulty itself (Jentsch, 1997). That is, negative reactions to robots in the Uncanny Valley may be a special case of general aversion to cognitive inhibition caused by competing perceptual representations (Ferrey et al., 2015; Freeman & Johnson, 2016). Similar mechanisms may underlie confusion occurring between categories of gender, race, age, sexual orientation, attitude, attractiveness, skin tone, accent, and relationship type (Fiske & Taylor, 2008).

Supporting some predictions of the category confusion explanation for the Uncanny Valley, human viewers may indeed perceive a categorical boundary between human and various near-human stimuli (Loosser & Wheatley, 2010; Weis & Wiese, 2017), and midrange stimuli seem to elicit the most confusion (Cheetham et al., 2013, 2015; MacDorman & Chattpadhyay, 2016; Mathur & Reichling, 2016; Yamada et al., 2013). However, it remains unclear whether confusion at the category boundary actually causes negative emotional responses to ambiguous faces or is merely an epiphenomenon. One falsifiable prediction of the mechanistic account is that the point on the nonhuman-human spectrum eliciting maximum confusion should coincide with the point eliciting the most negative social and affective responses. Some studies have preliminarily supported this prediction (Mathur & Reichling, 2016; Yamada et al., 2013), but others have not (Cheetham et al., 2014, 2015; Loosser & Wheatley, 2010; MacDorman & Chattpadhyay, 2016).

The apparent conflict between these findings may partly reflect methodological limitations

(Kätsyri et al., 2015; Lay et al., 2016). First, many studies to date have used as few as three stimuli spanning only a limited, very human-like region of the nonhuman-human spectrum (approximately points 4-9 in Figure 1) because, for example, the least human-like stimuli were highly realistic computer renderings of humans. If all stimuli are more human-like than the point eliciting the most negative Uncanny Valley reactions, then affective patterns to these stimuli might increase monotonically with human-likeness rather than showing a characteristic Uncanny Valley curve (as seen in Cheetham et al. (2014, 2015); Looser & Wheatley (2010); MacDorman & Chattpadhyay (2016)). Stimuli whose truncated range precludes observing the Uncanny Valley itself may also disrupt assessment of category confusion in the Uncanny Valley. Second, most studies have generated stimuli through digital image morphing, a method that risks producing unrealistic transitional images exhibiting, for example, partially transparent facial features or incompatible combinations of features that would be avoided in real-world robot design (Kätsyri et al., 2015; Kawabe et al., 2017). It is plausible that these unrealistic features themselves might produce artifactual confusion or aversion. Statistical methods have rarely assessed a mechanistic role of category confusion, instead focusing on certain downstream predictions of the category confusion mechanism hypothesis, for example by testing for associations across stimuli between confusion and emotional responses. Two studies more directly assessed whether reaction time, a coarse measure of confusion, statistically mediated the relationship between human-likeness and ratings of eeriness or weirdness, but with differing findings (Carr et al., 2017; Mathur & Reichling, 2016). Carr et al. (2017) reported mediation, but that study used only three stimuli; Mathur & Reichling (2016) did not detect mediation, but this was a secondary analysis of reaction times in a mechano-humanness rating task rather than a categorization task.

The present study aims to: (1) provide the most precise estimate to date of the shape of the Uncanny Valley curve in real-life robots and humans; (2) estimate the degree of human-likeness marking the perceptual category boundary between “non-human” and “human”; and (3) rigorously assess whether category confusion is a mechanism for Uncanny Valley effects.

We first assembled a large corpus of face images using Internet searches to identify images of socially interactive robots that have actually been built. We used stringent inclusion criteria and validation studies to minimize variation on potential confounders such as the face’s perceived emotion (Lay et al., 2016). We ensured that the faces were well-distributed across the full spectrum of human-likeness, enabling precise estimates of the Uncanny Valley curve and of the location of the category boundary between “robot” and “human”. Subjects recruited at six collaborating sites in four countries rated each face on human-likeness and likability, and they attempted to rapidly categorize each face as “robot” or “human” while we collected measures of category confusion. We assessed confusion using validated measures based on mouse-tracking to supplement coarser existing measures based on reaction time (Freeman & Johnson, 2016; Mathur & Reichling, in press). We assessed for statistical mediation as predicted by the mechanistic account of category confusion in a manner that accommodated the expected nonlinear relationships between human-likeness, confusion, and likability. Subsequent sections of this paper are structured as follows: we will describe methods for stimulus validation, for measurement of category confusion, and for subject recruitment, then describe statistical methods and results for each of the three aims in turn, and conclude with a general discussion.

2. DATA COLLECTION METHODS

All methods and statistical analyses were preregistered in detail; the Supplement describes and justifies some deviations from this protocol. All measures and experiments are reported, and we determined sample sizes in advance. All data, materials, analysis code, and the preregistration are publicly available and documented (<https://osf.io/mu5xj/>).

2.1. Face stimuli

We selected stimuli depicting the faces of real robots designed for social interaction, as well as faces of real humans. To support assessing a potential role of category confusion in mediating the relationship between human-likeness and likability, we attempted to select stimuli in a manner that would minimize confounding of the relationships between human-likeness, confusion, and likability (VanderWeele, 2015). For example, if displaying more positive emotion causes a face to be perceived as more human-like and also causes the face to be perceived as more likable, then perceived emotion could act as a confounder that would compromise causal conclusions from mediation analysis (VanderWeele, 2015). We identified face images using an objective Internet search process similar to that of Mathur & Reichling (2016). In addition to Mathur & Reichling (2016)'s inclusion and exclusion criteria (reproduced in the Supplement), we applied the following inclusion criteria to further minimize variation on such potential confounders. First, the faces had to be photographed in frontal view. Second, the faces had to be perceived as displaying low emotion, defined as having a mean rating between -20 and $+20$ on a visual analog scale ranging from -100 to $+100$. Third, the faces had to be densely spread over the entire spectrum from extremely mechanical to extremely human-like, rather than concentrated in only certain parts of the spectrum. Specifically, on a continuous scale of "mechano-humanness" (MH) score ranging from -100 ("extremely mechanical") to $+100$ ("extremely human-like"), we required that there be no 50-point span of MH score occupied by fewer than 20 faces.

To collect face stimuli, we first reviewed an existing set of 80 robot faces from Mathur & Reichling (2016); these faces had existing estimates of MH score. We discarded images that failed the more stringent inclusion criteria used here (e.g., because the photo was a 3/4 view or the face displayed too much emotion). For failed images, we searched the Internet to try to identify alternative photos of the excluded robot that did meet the criteria. The stimuli in Mathur & Reichling (2016) were somewhat limited by their sparse coverage of the

very human-like range (i.e., MH score above about +90); to resolve this limitation in the present stimuli, we also obtained images of actual human faces from Shutterstock, an online image bank. To obtain these human faces, we performed a broad search using terms such as “portrait”, “face”, and “unemotional”, attempting to select an assortment of faces meeting the same visual criteria as used for the robot faces and spanning a range from unmistakably human to somewhat artificial in appearance. The latter included faces that seemed unusually “perfect” due to marked symmetry, lack of flaws, or heavy use of makeup, and conversely others that seemed unusually “imperfect” due, for example, to particularly prominent features.

To ensure that the resulting robot and human faces were adequately densely distributed throughout the MH spectrum, we iterated between finding additional candidate images that appeared to meet the inclusion criteria and testing the candidate images for perceived emotion and MH score using groups of pilot subjects. At least 26 pilot subjects (mean: 40) rated each candidate face on perceived emotion and MH score. We ultimately included 182 faces, comprising 122 robots and 60 humans (Figure 2). The final validated corpus of images, along with summary measures of their ratings on all analyzed variables in this study, is publicly available for by-attribution use in future research (<https://osf.io/mu5xj/>). The faces had mean MH score -12.3 (-53.5 for the robots and 71.5 for the humans) and had mean likability -5.4 (-32.9 for the robots and 50.5 for the humans).

wrong - should be
-6.1, -47, 78

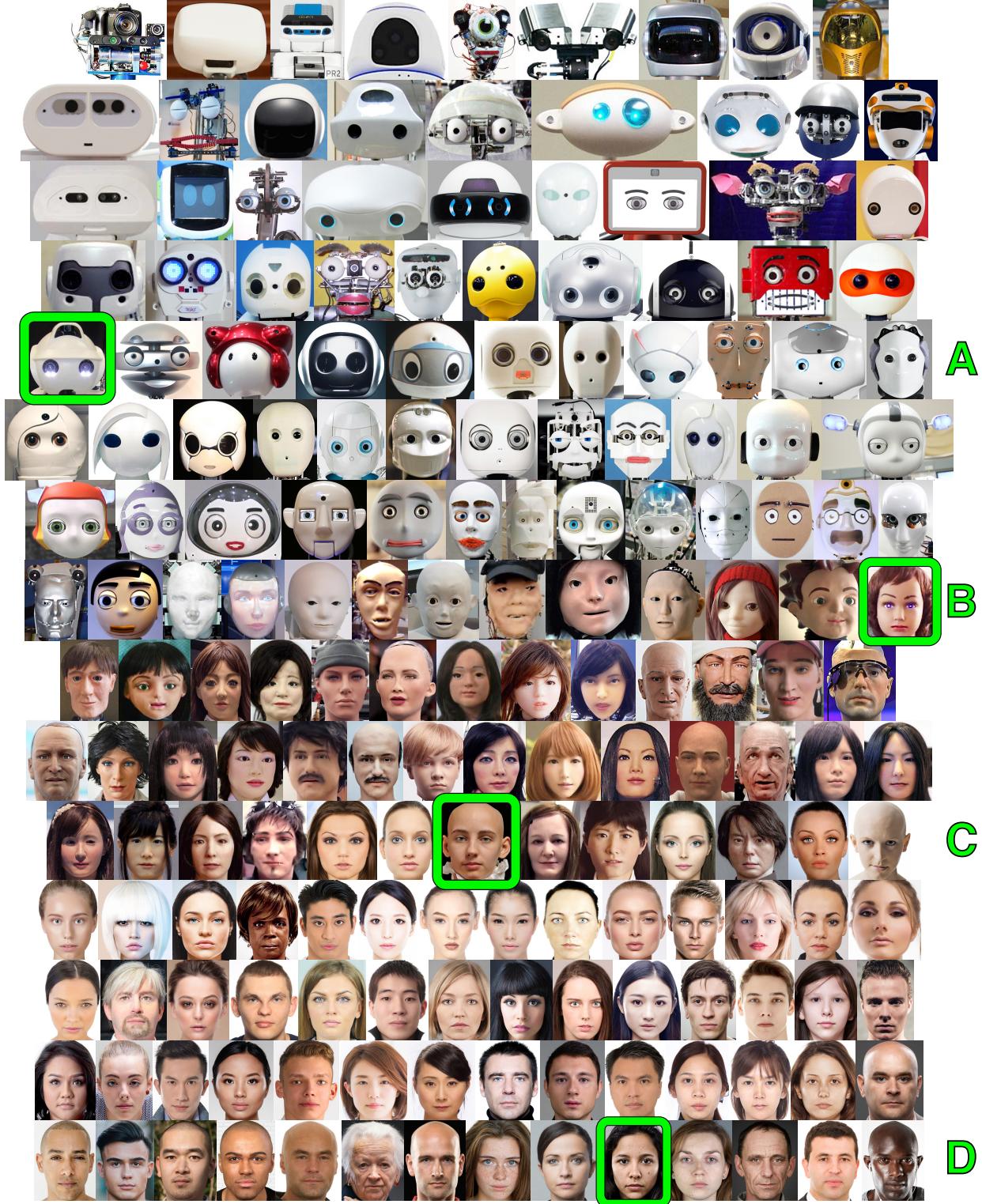


Figure 2: Robot and human face stimuli displayed in ascending order of mean mechano-humanness (MH) score. Boxed faces are those with MH scores closest to the MH scores associated with: (A) the initial likability apex of Uncanny Valley curve (estimation described in Section 3); (B) the likability nadir of Uncanny Valley; (C) the robot/human category boundary (estimation described in Section 4); and (D) the final apex of likability.

2.2. Measures of human-likeness, confusion, and likability

We measured human-likeness by asking, “How mechanical versus human-like does this face look?”; subjects responded using a bipolar visual analog scale ranging from -100 to $+100$ with the endpoints labeled “extremely mechanical” and “extremely human-like” (Mathur & Reichling, 2016). We refer to this measure as “mechano-humanness” or “MH” score. We measured likability by asking subjects to “Estimate how friendly and enjoyable (or creepy) it might be to interact with the robot in some everyday situation, such as asking a question at a museum’s information booth”; subjects responded on a similar visual analog scale with the endpoints labeled “Less friendly; more unpleasant and creepy” and “More friendly and pleasant; less creepy” (Mathur & Reichling, 2016).

We used validated open-source software (Mathur & Reichling, in press) to collect five established measures of category confusion (e.g., Freeman et al. (2008)). Subjects viewed the faces sequentially and were asked to rapidly categorize each face as “robot” or “human” by clicking on one of two buttons presented on the left and right sides of the window (Figure 3). (Methodological details of the categorization task are available in Mathur & Reichling (in press).) Ambiguous stimuli are thought to activate mental representation of both categories simultaneously, leading to dynamic competition that manifests in real time as unstable mouse trajectories (Freeman & Johnson, 2016). That is, because the subject is continuously or alternately attracted to both categories, the mouse trajectory may contain frequent direction changes and may diverge substantially from a direct path from the start position to the location of the category button ultimately chosen.

Therefore, as primary measures of confusion, we collected (Freeman et al., 2008): (1) the number of times the subject’s mouse changed directions horizontally during categorization (*x-flips*); (2) the maximum horizontal deviation between the subject’s mouse trajectory and an ideal trajectory consisting of a straight line from the subject’s initial cursor position to the finally chosen category button (*maximum x-deviation*; red solid line in Figure 3); and (3) the

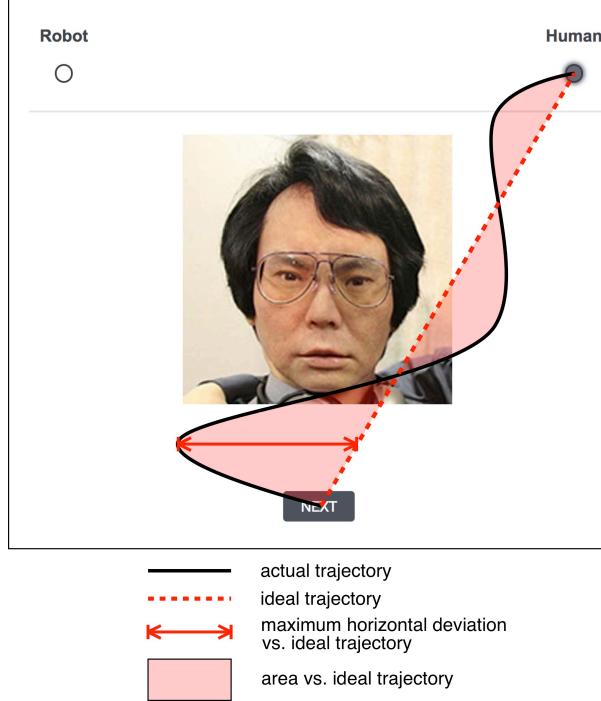


Figure 3: Mouse-tracking measures of category confusion (reproduced from Mathur & Reichling (in press)).

area between the ideal and actual trajectories (pink shading in Figure 3). We additionally measured: (4) the *peak speed* of the subject’s cursor (ambiguous stimuli tend to produce higher peak speeds, reflecting abrupt category shifts (Freeman et al., 2016)); and (5) the total *reaction time* for the trial (ambiguous stimuli tend to produce longer reaction times). Because the latter two measures have limitations as measures of category confusion (Freeman et al., 2016), we made an *a priori* decision to designate them as secondary measures. Each subject provided measures of MH score, confusion, and likability in that order for all 182 faces. To minimize the possibility of task interference or memory effects due to a subject’s repeated exposure to each face, we collected each measure in a separate wave of data collection spaced by approximately one week. Within each wave, the order of the faces was randomized for each subject. At the end of the first wave, subjects also completed basic demographic measures of age, sex, education level, and race/ethnicity. At the end of each wave, subjects reported any technical or comprehension problems.

2.3. Subjects

We collected data at six colleges and universities in the United States, Hungary, the Netherlands, and Italy; we recruited the data collection sites through the authors' previous collaborations and the online platform StudySwap. Supplement Table S1 describes characteristics of the sites. Aggregating across sites, we analyzed data from 358 subjects, who were 72% female with mean age 21.5 years; further demographic characteristics are described in Supplement Table S2. Each site aimed to collect data on at least 50 English-speaking subjects in a quiet lab or classroom on lab-provided computers that were pre-tested for accurate collection of mouse-tracking data. Labs incentivized participation using various monetary compensation, course credit, or volunteer schemes. All subjects completed the study using the same Qualtrics questionnaires provided by the lead authors. Each lab secured its own ethics approval or waiver as appropriate to its location.

Based on *a priori* criteria, we excluded subjects who did not complete all three waves of the questionnaire, whose data indicated technical problems (e.g., reflecting rare, idiosyncratic timing issues that caused no times to be recorded for a subject, or caused timing to stop prematurely; Mathur & Reichling (in press)), or for whom there were known data collection errors (e.g., the waves of data collection were run in the wrong order). Supplement Figure S1 details the number of subjects excluded for each reason. Additionally, we excluded individual trials in which the value of any confusion measure was larger than its 75th percentile plus 1.5 times its interquartile range or smaller its 25th percentile minus 1.5 times the interquartile range; we made this decision because the mediators could potentially take on extreme values if, for example, a subject made an uncontrolled cursor movement. We did not exclude trials in which the subject chose the wrong category for the face (e.g., selected “human” when presented with a robot) because we expected many faces to be quite hard to judge. After all exclusions, the analysis dataset comprised 358 subjects, totalling 55430 ratings of the 182 faces.

3. ESTIMATING THE SHAPE OF THE UNCANNY VALLEY

3.1. Statistical methods

Throughout, one author (MBM) performed all statistical analyses in R (Version 3.5.1). All analyses described in the main text were conducted with means by face as the unit of analysis¹.

In all analyses, we adjusted for a face's mean perceived emotion rating, as estimated during stimulus validation, because we suspected that emotion might statistically confound the relationship between MH score and likability. We first estimated the Uncanny Valley curve by fitting ordinary least squares models that regressed likability on polynomial terms for MH score (e.g., MH, MH^2 , MH^3 , etc.). We mean-centered MH score in analysis, but we report and plot results on the uncentered scale for interpretability, except where otherwise noted. We used Akaike's Information Criterion (AIC) to select the lowest-order and best-fitting polynomial model (Akaike [1974]). We weighted each data point by its inverse-variance of likability to account for the fact that some faces were rated with more precision than others, although unweighted analyses yielded nearly identical results (Supplement).

3.2. Results

Figure 4 shows the best-fitting and most parsimonious model for the relationship between MH score and likability, which was a six-degree polynomial in MH score. As predicted by the Uncanny Valley theory, estimates from this model indicated that as faces progressed from extremely mechanical (MH score near -100) to somewhat less mechanical, likability tended to increase to a point, reaching an initial apex of -18.0 for faces with an MH score of -80.9.

¹For statistical efficiency, we had planned to analyze data at the individual trial level rather than aggregating by face. However, a comparison of individual-level versus face-level estimates of the Uncanny Valley (Supplement) suggested substantial attenuation of the relationship between MH score and likability in the individual-level data. This strongly suggested that subjects' individual ratings of MH score are essentially noisy measurements of a face's "true" MH score and hence that conducting analyses at the individual level would result in downward-biased estimates of the relationship between MH score and likability, a phenomenon that is well-characterized in the literature on nondifferential exposure measurement error (Thomas et al., 1993). This bias is largely mitigated through aggregation as presented in the main text (Prentice & Sheppard, 1995).

After this initial apex, as faces continued to become more human-like, likability began to decrease, dropping to its overall nadir of -67.4 for faces with an MH score of -23.6 . Beyond this Uncanny Valley, as faces continued to become more human-like, their likability once again tended to increase monotonically, ultimately reaching a maximum of 59.6 for faces nearly indistinguishable from humans (i.e., those with an MH score of 93.2). Thus, all key features of the theorized Uncanny Valley were apparent in these stimuli.

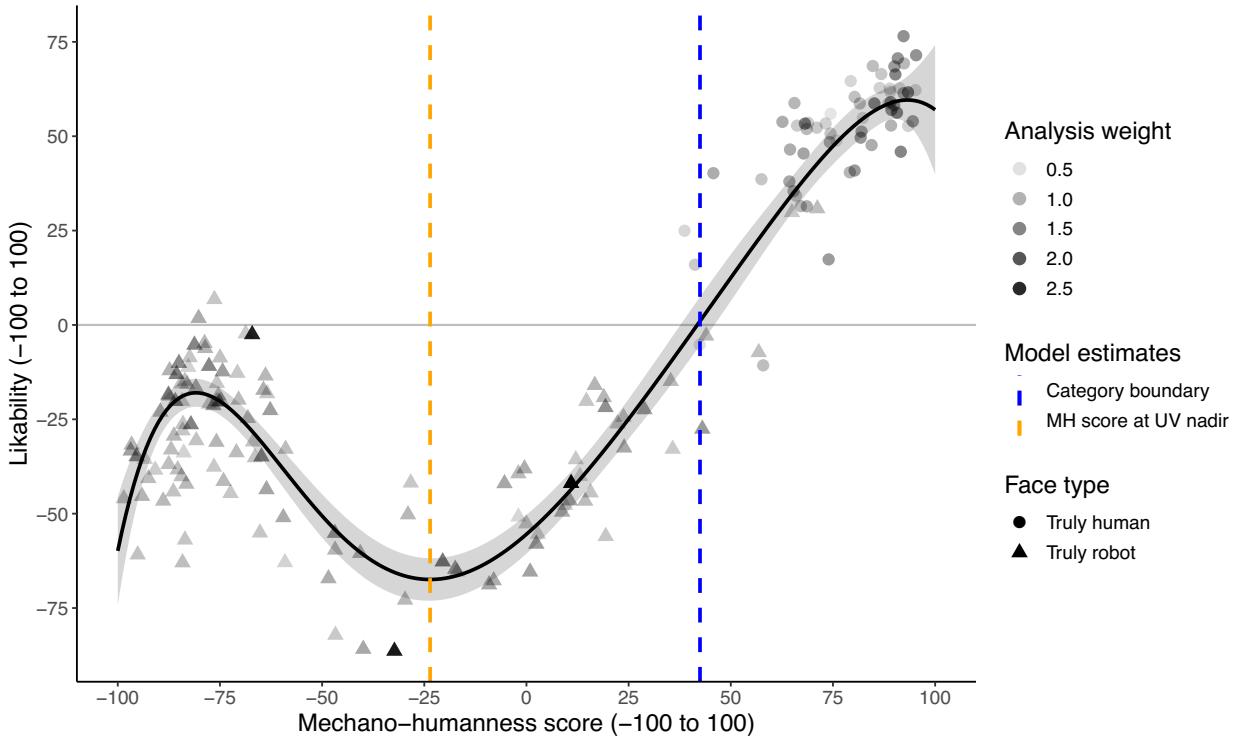


Figure 4: Uncanny Valley curve in 182 faces (points) with empirically estimated curve based on polynomial regression (solid black curve). Triangles indicate faces that were actually robots; circles indicate faces that were actually humans. Point opacity is proportional to the face's inverse-variance weight in analysis. The shaded band is a 95% pointwise confidence interval for the fitted likability values when setting emotion (the adjusted covariate) to its mean. The vertical dashed lines mark the estimated MH scores at which likability reached its nadir (orange) and where the category boundary occurred (blue; estimation described in Section 4).

4. ESTIMATING THE CATEGORY BOUNDARY LOCATION

4.1. Statistical methods

We next estimated the location of the category boundary, defined as the MH score at which the proportion of subjects categorizing the face as “human” is closest to 50%. To do so, we used unweighted ordinary least squares regression to model the proportion of subjects categorizing each face as “human” as a polynomial function of MH score, again choosing the best-fitting and most parsimonious polynomial using the AIC. For this analysis, we made a post hoc decision to exclude the 54 (30%) of faces that were never categorized as “human” because this large mass of faces with a 0% probability would have been challenging to fit accurately using a smooth polynomial model, and regardless, faces so distant from the category boundary would have contributed little to statistically estimating the boundary location². (No faces had a 100% probability of being categorized as human.) We used the estimated coefficients from this model to estimate the category boundary location.

4.2. Results

Figure 5 shows the best-fitting model for the relationship between a face’s MH score and its probability of being categorized as “human”. As expected, the estimated probability of a face’s being categorized as “human” increased monotonically with increasing MH score. We estimated that the category boundary occurred at an MH score of 42.5.

²A sensitivity analysis in which we did not exclude these faces yielded a very similar estimate of the boundary location.

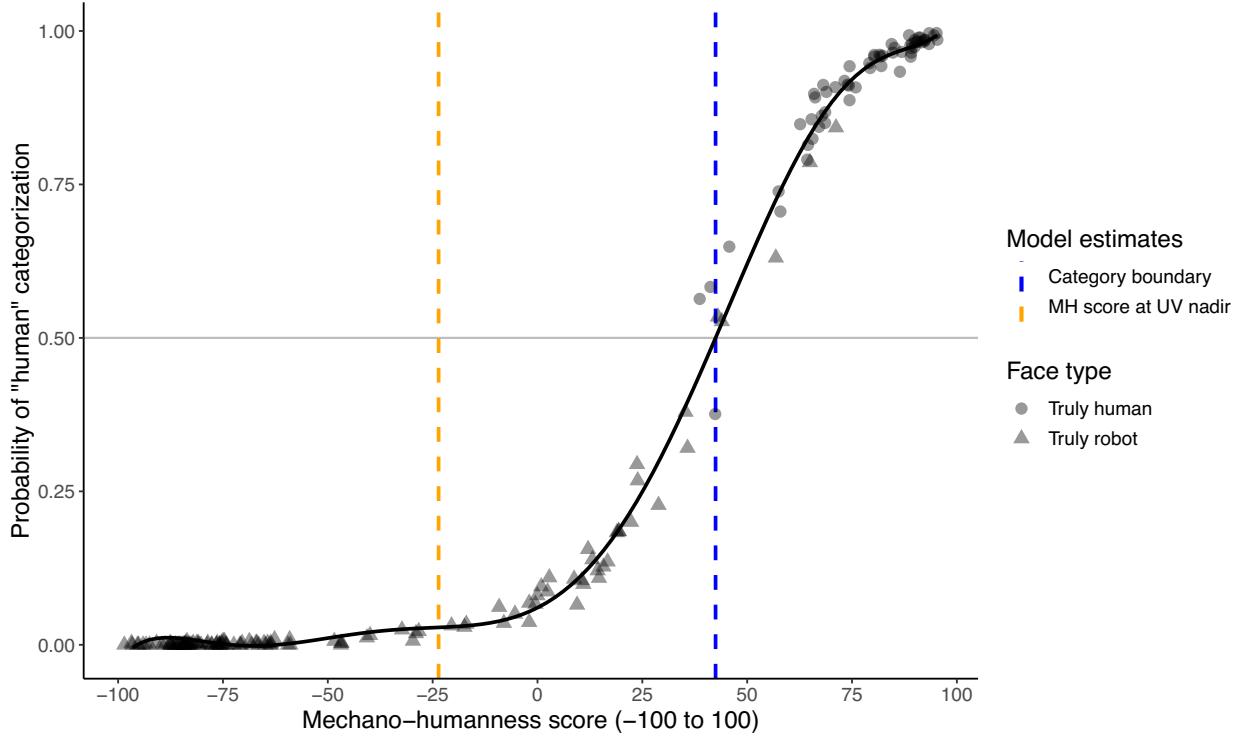


Figure 5: *MH score versus proportion of subjects categorizing the face as “human” with empirically estimated curve based on polynomial regression (solid black curve). Triangles indicate faces that were actually robots; circles indicate faces that were actually humans. The vertical dashed lines mark the estimated MH scores at which likability reached its nadir (orange) and where the category boundary occurred (blue).*

5. ASSESSING CATEGORY CONFUSION AS A MEDIATOR

5.1. Statistical methods

To investigate whether category confusion might be a mechanism of the observed Uncanny Valley effect, we conducted separate mediation analyses for each of the primary and secondary measures of category confusion, treating MH score as the exposure, the confusion measures as the mediators, and likability as the outcome. To improve the interpretability of the direct and indirect effect estimates, we standardized MH score and all mediators for these analyses. In addition, as a simple method to consider in aggregate the three mediators pre-specified as primary measures of confusion (i.e., x -flips, area, and maximum x -deviation), we conducted a final mediation analysis on a composite measure of confusion constructed by summing

the z -scores for these three primary confusion measures, which were very highly correlated (Pearson's r from 0.91 to 0.99).

Estimating causal mediation effects relies on certain no-confounding assumptions (e.g., VanderWeele (2015)). To this end, during stimulus validation, we had attempted to eliminate many sources of confounding by selecting stimuli that were comparable on graphical features and that were almost emotionally neutral; additionally, in analysis, we controlled for a face's mean emotion as rated during stimulus validation. We conducted mediation analyses using a simulation-based method that involves fitting a model for the mediator as a function of the exposure and a model for the outcome as a function of the mediator and the exposure (Imai et al., 2011; Tingley et al., 2014). For each mediator model, we used generalized additive models (GAM) with the identity link to regress the measure of category confusion on a spline basis for MH score. For the outcome models, we similarly used GAM to model likability as a function of each mediator and MH score. The outcome models additionally allowed for nonlinear interactions between MH score and the candidate mediator via a tensor product term, which we dropped from the model if its inclusion worsened the model's AIC. We chose these models in order to flexibly accommodate the expected nonlinearities that characterize the Uncanny Valley, as well as the possibly interactive relationship between MH score and confusion. Because the category confusion measures were often skewed or bimodal, suggesting non-normal errors, we estimated all confidence intervals and p -values using nonparametric bootstrapping.

5.2. Results

Figures 6 and 7 respectively plot MH score versus each confusion measure, and each confusion measure versus likability, along with GAM fits for each relationship. MH score appeared to have nonlinear and non-monotonic relationships with most of the confusion measures; the GAM models for the three primary mediators and their composite estimated that the confusion measures peaked at unstandardized MH scores of 34.3 for x -flips, 95.3 for area,

95.3 for maximum x -deviation, and 95.3 for their composite (Figure 6, dashed vertical lines).

Results were similar for the two secondary confusion measures, namely speed and reaction time. Considering the outcome models, likability did not appear to increase monotonically as the confusion measures increased; rather, the relationships appeared nonlinear and variable across confusion measures (Figure 7).

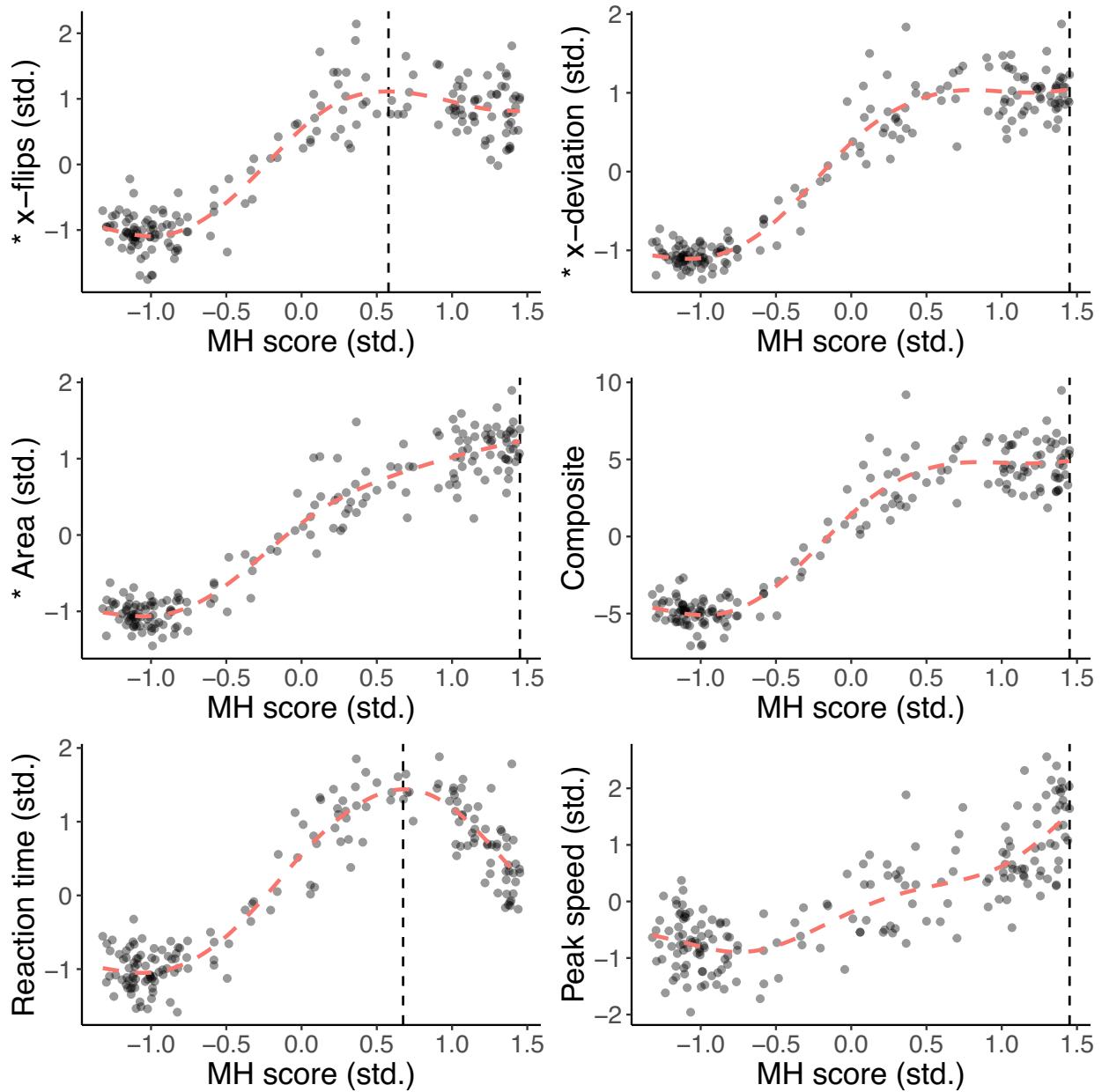


Figure 6: *MH score vs. confusion relationships. The red dashed line represents fitted values from a GAM model as used as in mediation analysis. The vertical dashed line marks the MH score associated with maximum confusion, as estimated by GAM. *: primary confusion measure.*

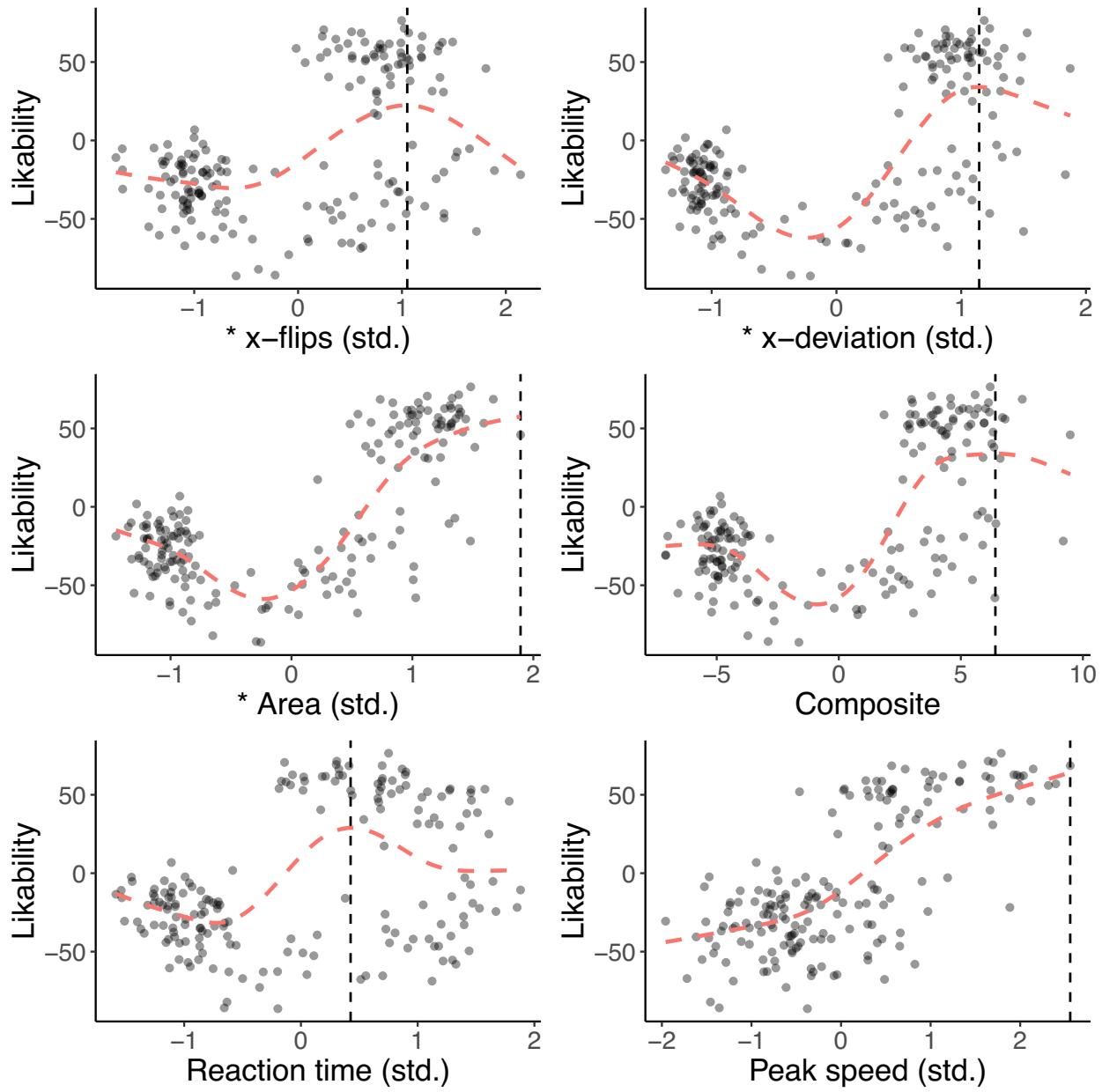


Figure 7: Confusion vs. likability relationships. The red dashed line represents fitted values from a GAM model without interactions between MH score and mediators. The vertical dashed line marks the value of the confusion measure associated with maximum likability, as estimated by GAM. *: primary confusion measure.

Table 1: Analysis of category confusion measures as mediators of the relationship between MH score and likability. Direct and indirect effects are presented for a contrast representing an increase in MH score from its mean to half a standard deviation above its mean. All mediators and MH score were standardized.

Confusion variable	Statistic	Estimate [95% CI]	p-value
Primary measures			
x-flips	Direct effect	32.7 [19.7, 44.9]	< 0.00001
	Indirect effect	-0.6 [-3.6, 2.1]	0.68
	% mediated	-2 [-12, 7]	0.68
x-deviation	Direct effect	31.3 [15.8, 44.7]	0.006
	Indirect effect	0.9 [-5.9, 13.0]	0.62
	% mediated	3 [-22, 41]	0.62
Area	Direct effect	30 [12.6, 39.5]	0.004
	Indirect effect	2.8 [-1.9, 15.4]	0.17
	% mediated	8 [-6, 48]	0.17
Composite	Direct effect	30 [8.9, 45.0]	0.03
	Indirect effect	2.2 [-5.8, 20.7]	0.33
	% mediated	7 [-19, 65]	0.33
Secondary measures			
Reaction time	Direct effect	36.7 [16.3, 49.2]	0.004
	Indirect effect	-5 [-12.1, 10.0]	0.49
	% mediated	-16 [-53, 33]	0.49
Peak speed	Direct effect	29.1 [15.4, 43.1]	0.002
	Indirect effect	3.4 [-3.3, 10.6]	0.13
	% mediated	10 [-11, 34]	0.13

Table 1 presents results of the mediation analyses, in which coefficient estimates represent estimated differences in likability on its original scale ranging from -100 to $+100$. The direct effects represent effects of MH score on likability that occurred independently of each mediating confusion measure; these ranged from 30.0 (95% CI: [12.6, 39.5]) to 32.7 (95% CI: [19.7, 44.9]) for the primary raw measures (x -flips, maximum x -deviation, and area) and 30.0 (95% CI: [8.9, 45.0]) for the composite measure. The indirect effects represent effects of MH score

on likability occurring because of mediation by each confusion measure, and the percent mediated represents the percent of the total effect of MH score on likability that is due to mediation by each confusion measure. The estimated indirect effects were approximately an order of magnitude smaller than the estimated direct effects, ranging from -0.6 (95% CI: [-3.6, 2.1]; estimated percent mediated: -2%) to 2.8 (95% CI: [-1.9, 15.4]; estimated percent mediated: 8%) for the three primary confusion measures and 2.2 (95% CI: [-5.8, 20.7]; estimated percent mediated: 7%) for their composite. (Note that negative estimates for the percent mediated occur when the direct and indirect effects are in different directions.) Results for the two secondary confusion measures were qualitatively similar.

In summary, the confusion measures varied somewhat in their relationships with MH score. For two of the three primary confusion measures, as well as their composite, faces nearly identical to humans (MH scores near +100) produced the most confusion. For the final primary confusion measure, faces at a lower MH score of 34.3 produced the most confusion. Despite these variations across confusion measures, all indicated that confusion peaked for robots that were considerably more human-like than those that were most dislikable (i.e., those occupying the nadir of the Uncanny Valley, estimated to occur at an MH score of -23.6). The relationships between the confusion measures and likability suggested that increased confusion was not clearly and monotonically associated with decreased likability, as the category confusion hypothesis might have predicted. Indeed, the mediation analyses suggested that any mediation by confusion was likely very minimal, with results quite consistent across confusion measures.

6. DISCUSSION

The first two aims of our study were to precisely estimate the shape of the Uncanny Valley curve and the location of the boundary between the categories “robot” and “human”. To this end, we developed a publicly available corpus of 182 images of real, socially interactive robots as well as humans; these images have closely controlled graphical features, are perceived to

be emotionally almost neutral, and are densely distributed throughout a broad spectrum of human-likeness. These faces showed a relationship between human-likeness and likability in which all key features of the theorized Uncanny Valley were apparent: namely, the initial increase in likability as faces progressed from extremely mechanical to somewhat less mechanical, followed by a classic “Uncanny Valley” nadir as faces became considerably more human-like and markedly dislikeable, followed by a gradual increase in likability to its eventual, overall apex as faces became nearly indistinguishable from humans.

Our estimated Uncanny Valley curve also showed interesting differences from traditional predictions. The initial apex of likability for unmistakably mechanical faces was negative (likability = -18.0 ; Figure 4) rather than markedly positive as Mori (1970) originally postulated (c.f. Figure 1), indicating that these robots were still somewhat disliked, and very few individual faces in this region did achieve positive likability ratings. This suggests that the strategy of attempting to avoid the Uncanny Valley by deliberately designing android robots that are unmistakably mechanical (Duffy, 2003; Mori, 1970) might severely stunt the robots’ likability and ultimate social success. Additionally, we found that the nadir of the Uncanny Valley occurred not for faces nearly indistinguishable from humans, as (Mori, 1970) originally theorized, but rather for faces perceived to be more mechanical than human-like (i.e., MH score = -23.6).

In the categorization task, the relationship between the faces’ human-likeness and their probabilities of being categorized as “human” suggested that humans do perceive a perceptual categorical boundary, as opposed to a smooth continuum, between the properties of “robot” and “human”. That is, typical of categorical perception (e.g., de Gelder et al. (1997); Etcoff & Magee (1992)), faces in most regions of the human-likeness spectrum were reliably classified as either “robot” or “human”, but faces within a steeply sloped region around the category boundary elicited much less stable categorizations (Figure 4). This boundary zone is where category confusion is thought to occur. It is interesting that our estimated category boundary occurred at an MH score of 42.5, corresponding to faces perceived to be about 71% human

on the 200-point MH scale, rather than 50% as one might predict by analogy with classic psychophysical studies of category confusion (Harnad, 1987). Conceivably, this off-center position of the category boundary could be related to the fact that the positions of our stimuli on the MH scale were necessarily determined based on subjects' subjective ratings rather than objective metrics, since the faces represented real robots rather than constructed stimuli. However, casting some doubt on this interpretation, previous studies that used morphing methods to generate objectively quantified mixtures of robot and human faces reported similarly located category boundaries in the range of 60-70% human (Cheetham et al., 2011; Weis & Wiese, 2017).

Critically, the Uncanny Valley and the estimated category boundary (at MH score = 42.5) did not coincide; they in fact occurred in two quite distinct regions of the human-likeness spectrum. That is, faces that were maximally disliked (at the nadir of the Uncanny Valley; MH score = -23.6) were almost always categorized as “robots”; in contrast, the zone in which unstable categorization occurred occupied higher MH scores from approximately 0 to 75. Conversely, maximally ambiguous faces (near the category boundary) were not, on average, disliked. This discrepancy casts doubt on the category confusion hypothesis, which would predict that the most ambiguous faces would be most disliked. For the third aim of our study, we more formally assessed category confusion as a possible mechanism for Uncanny Valley effects by conducting mediation analyses treating human-likeness as the exposure, confusion as the mediator, and likability as the outcome. We measured confusion using fine-grained, validated mouse-tracking measures to supplement the coarser measures, such as reaction time, used in previous literature. These analyses did not support mediation by confusion; estimated indirect effects (representing mediation through confusion) were typically an order of magnitude weaker than direct effects.

These analyses have some limitations. Ideally, we would have measured perceptions of human-likeness, confusion, and likability as they occurred in real time, likely within a span of milliseconds. This hypothetical design would have clarified, for example, whether

perceptions of human-likeness indeed influence perceptions of likability, or whether perhaps viewers first make judgments about likability, which in turn affect their perceptions of a face’s human-likeness. However, such a design would be logically unfeasible, requiring subjects to perform three different rating tasks all within the span of milliseconds. Furthermore, designs in which subjects rate all three characteristics in quick succession may introduce task interference or demand characteristics. Although our three-wave design was intended to minimize these types of bias, it cannot rule out reverse causation. Additionally, as discussed, mediation analyses inherently rely on no-confounding assumptions. We tried to minimize confounding by selecting closely matched stimuli and controlling for perceived emotion in analyses. Given the very small size of the estimated indirect effects, we believe it unlikely that any residual confounding would have masked meaningfully strong mediation by confusion.

Ultimately, these findings suggest that although humans do perceive a category boundary between “robot” and “human”, this boundary does not coincide with the Uncanny Valley itself, and category confusion produced by this boundary does not seem to explain Uncanny Valley aversions. It is striking that, despite the decades-long prominence of the Uncanny Valley theory and robot designers’ sophisticated attempts to circumvent it, the robots we sampled — which were purposefully designed for social interaction — nevertheless were dislikeable on average (mean MH score = -32.9) and showed a prominent Uncanny Valley. These findings point to the continued importance of attempting to elucidate the mechanisms underlying the effect.

REPRODUCIBILITY

All materials, data, and code required to reproduce this research are publicly available and documented (<https://osf.io/mu5xj/>). The preregistrations are publicly available (<https://osf.io/mu5xj/registrations/>).

AUTHOR CONTRIBUTIONS

MBM and DBR conceptualized and designed the study. MBM planned and conducted statistical analyses and led writing. The remaining authors collected data and revised the manuscript.

ACKNOWLEDGMENTS

This research was supported by a Harvard University Mind, Brain, & Behavior grant. MBM was supported by NIH grant R01 CA222147. The funders had no role in the design, conduct, or reporting of this research.

REFERENCES

- Akaike, H. (1974). A new look at the statistical model identification. In *Selected Papers of Hirotugu Akaike* (pp. 215–222). Springer.
- Carr, E. W., Hofree, G., Sheldon, K., Saygin, A. P., & Winkielman, P. (2017). Is that a human? Categorization (dis)fluency drives evaluations of agents ambiguous on human-likeness. *Journal of Experimental Psychology: Human Perception and Performance*, 43(4), 651.
- Cheetham, M., Pavlovic, I., Jordan, N., Suter, P., & Jancke, L. (2013). Category processing and the human likeness dimension of the uncanny valley hypothesis: eye-tracking data. *Frontiers in Psychology*, 4, 108.
- Cheetham, M., Suter, P., & Jäncke, L. (2011). The human likeness dimension of the “uncanny valley hypothesis”: behavioral and functional mri findings. *Frontiers in Human Neuroscience*, 5, 126.
- Cheetham, M., Suter, P., & Jancke, L. (2014). Perceptual discrimination difficulty and familiarity in the uncanny valley: more like a “happy valley”. *Frontiers in Psychology*, 5, 1219.
- Cheetham, M., Wu, L., Pauli, P., & Jancke, L. (2015). Arousal, valence, and the uncanny valley: Psychophysiological and self-report findings. *Frontiers in Psychology*, 6, 981.
- de Gelder, B., Teunisse, J.-P., & Benson, P. J. (1997). Categorical perception of facial expressions: Categories and their internal structure. *Cognition & Emotion*, 11(1), 1–23.
- Duffy, B. R. (2003). Anthropomorphism and the social robot. *Robotics and Autonomous Systems*, 42(3-4), 177–190.

- Etcoff, N. L., & Magee, J. J. (1992). Categorical perception of facial expressions. *Cognition*, 44(3), 227–240.
- Ferrey, A. E., Burleigh, T. J., & Fenske, M. J. (2015). Stimulus-category competition, inhibition, and affective devaluation: a novel account of the uncanny valley. *Frontiers in Psychology*, 6, 249.
- Fiske, S. T., & Taylor, S. E. (2008). *Social cognition: From brains to culture*. McGraw Hill New York.
- Freeman, J. B., Ambady, N., Rule, N. O., & Johnson, K. L. (2008). Will a category cue attract you? Motor output reveals dynamic competition across person construal. *Journal of Experimental Psychology: General*, 137(4), 673.
- Freeman, J. B., & Johnson, K. L. (2016). More than meets the eye: Split-second social perception. *Trends in Cognitive Sciences*, 20(5), 362–374.
- Freeman, J. B., Pauker, K., & Sanchez, D. T. (2016). A perceptual pathway to bias: Interracial exposure reduces abrupt shifts in real-time race perception that predict mixed-race bias. *Psychological Science*, 27(4), 502–517.
- Harnad, S. (1987). A new look at the statistical model identification. In S. Harnad (Ed.), *Categorical Perception: The Groundwork of Cognition* (pp. 1–25). Cambridge University Press.
- Ho, C.-C., MacDorman, K. F., & Pramono, Z. D. (2008). Human emotion and the uncanny valley: a GLM, MDS, and Isomap analysis of robot video ratings. In *3rd ACM/IEEE International Conference on Human-Robot Interaction (HRI)* (pp. 169–176).
- Imai, K., Keele, L., Tingley, D., & Yamamoto, T. (2011). Unpacking the black box of causality: Learning about causal mechanisms from experimental and observational studies. *American Political Science Review*, 105(4), 765–789.
- Jentsch, E. (1997). On the psychology of the uncanny (1906). *Angelaki: Journal of the Theoretical Humanities*, 2(1), 7–16.
- Jung, Y., & Cho, E. (2018). Context-specific affective and cognitive responses to humanoid robots. In *The 22nd biennial conference of the international telecommunications society: “beyond the boundaries: Challenges for business, policy and society”*. Seoul: International Telecommunications Society (ITS).
- Kätsyri, J., Förger, K., Mäkäräinen, M., & Takala, T. (2015). A review of empirical evidence on different uncanny valley hypotheses: Support for perceptual mismatch as one road to the valley of eeriness. *Frontiers in Psychology*, 6.
- Kawabe, T., Sasaki, K., Ihaya, K., & Yamada, Y. (2017). When categorization-based stranger avoidance explains the uncanny valley: A comment on MacDorman and Chattopadhyay (2016). *Cognition*, 161, 129–131.

- Lay, S., Brace, N., Pike, G., & Pollick, F. (2016). Circling around the uncanny valley: Design principles for research into the relation between human likeness and eeriness. *i-Perception*, 7(6), 2041669516681309.
- Lischetzke, T., Izquierdo, D., Hüller, C., & Appel, M. (2017). The topography of the uncanny valley and individuals' need for structure: A nonlinear mixed effects analysis. *Journal of Research in Personality*, 68, 96–113.
- Looser, C. E., & Wheatley, T. (2010). The tipping point of animacy: How, when, and where we perceive life in a face. *Psychological Science*, 21(12), 1854–1862.
- MacDorman, K. F., & Chattopadhyay, D. (2016). Reducing consistency in human realism increases the uncanny valley effect; increasing category uncertainty does not. *Cognition*, 146, 190–205.
- MacDorman, K. F., Green, R. D., Ho, C.-C., & Koch, C. T. (2009). Too real for comfort? uncanny responses to computer generated faces. *Computers in Human Behavior*, 25(3), 695–710.
- Mathur, M. B., & Reichling, D. B. (2016). Navigating a social world with robot partners: A quantitative cartography of the uncanny valley. *Cognition*, 146, 22–32.
- Mathur, M. B., & Reichling, D. B. (in press). Open-source software for mouse-tracking in qualtrics to measure category competition. *Behavior Research Methods*. (Preprint: <https://osf.io/ymxau/>.)
- Mori, M. (1970). The uncanny valley. *Energy*, 7(4), 33–35.
- Prentice, R. L., & Sheppard, L. (1995). Aggregate data studies of disease risk factors. *Biometrika*, 82(1), 113–125.
- Rabbitt, S. M., Kazdin, A. E., & Scassellati, B. (2015). Integrating socially assistive robotics into mental healthcare interventions: Applications and recommendations for expanded use. *Clinical Psychology Review*, 35, 35–46.
- Slijkhuis, P. J. (2017). *The uncanny valley phenomenon: A replication with short presentation times* (Unpublished master's thesis). University of Twente.
- Thomas, D., Stram, D., & Dwyer, J. (1993). Exposure measurement error: influence on exposure-disease relationships and methods of correction. *Annual Review of Public Health*, 14(1), 69–93.
- Tingley, D., Yamamoto, T., Hirose, K., Keele, L., & Imai, K. (2014). Mediation: R package for causal mediation analysis.
- VanderWeele, T. (2015). *Explanation in causal inference: methods for mediation and interaction*. Oxford University Press.
- Wang, S., Lilienfeld, S. O., & Rochat, P. (2015). The uncanny valley: Existence and explanations. *Review of General Psychology*, 19(4), 393.

Weis, P. P., & Wiese, E. (2017). Cognitive conflict as possible origin of the uncanny valley. In *Proceedings of the human factors and ergonomics society annual meeting* (Vol. 61, pp. 1599–1603).

Weisberger, M. (2018). *A floating “brain” will assist astronauts aboard the space station.* <https://www.livescience.com/61909-cimon-floating-brain-iss.html>. (Accessed: 2019-4-20)

Yamada, Y., Kawabe, T., & Ihaya, K. (2013). Categorization difficulty is associated with negative evaluation in the “uncanny valley” phenomenon. *Japanese Psychological Research*, 55(1), 20–32.

Supplement

Uncanny but not confusing: Multisite study of perceptual category confusion in the Uncanny Valley

CONTENTS

1	Supplementary methods	2
1.1	Additional inclusion criteria for face stimuli	2
1.2	Changes and additions to preregistered protocol	3
2	Supplementary results	4
2.1	Site and subject characteristics	4
2.2	Data exclusions	6
2.3	Alternative model specifications	6

1. SUPPLEMENTARY METHODS

1.1. Additional inclusion criteria for face stimuli

This section is reproduced directly from [Mathur & Reichling \(2016\)](#).

Inclusion:

1. The full face is shown from top of head to chin.
2. The robot is intended to interact socially with humans.
3. The robot has actually been built.
4. The robot is capable of physical movement (e.g., not a sculpture or purely CGI representation that lacks a three-dimensional body structure).
5. The robot is shown as it is meant to interact with users (e.g., not missing any hair, facial parts, skin, or clothing, if these are intended).
6. The robot represents an android that is plausibly capable of adult social interaction (e.g., not a baby or an animal).
7. The resolution of the original image is sufficient to yield a final cropped image at 100 d.p.i. and 3 in. tall.

Exclusion:

1. The robot represents a well-known character or a famous person (e.g., Einstein).
2. The image includes other faces or human body parts that would appear in the final cropped image.
3. Objects or text overlap the face.
4. The robot is marketed as a toy.

1.2. Changes and additions to preregistered protocol

We used an iterative approach to preregistration; each version of the preregistration is publicly available (<https://osf.io/mu5xj/registrations>). Here, we describe and justify deviations from our preregistered protocol. During stimulus validation, we had intended to include only faces with mean emotion ratings between -10 and $+10$. Because doing so did not yield enough eligible faces to meet the denseness criterion, we slightly relaxed the emotion criterion to allow faces with mean emotion rating between 20 and $+20$. Heuristically, these faces still appeared quite neutral. After stimulus validation but before data collection, we updated the preregistration to stipulate that we would adjust for mean emotion throughout all analyses to adjust for possible residual confounding by emotion. Also, when acquiring stimuli, per the preregistration, we contacted Slijkhuis (2017), who conducted a replication of our previous work that included images of an additional 16 wild-type faces with MH scores near the nadir of the UV. We received a response only after we had successfully acquired and validated the stimulus set, so we did not pursue these 16 faces. An early iteration of the preregistration described a validation study in which we would test for task interference between the MH score, confusion, and likability-rating tasks, but in a later preregistration, we replaced this with a stronger study design with washout periods between each task. Also, we validated the confusion measures using robot stimuli (Mathur & Reichling, *in press*) not used in main analyses rather than color stimuli as described in an early preregistration. For reasons described in the main text and in Section below, we conducted analyses on an aggregated dataset (means by face) rather than at the individual trial level as originally planned. We used different link functions as appropriate due to the change in unit of analysis (e.g., x -flips was modeled with the identity link instead of the log link because mean x -flips is a continuous rather than count variable). The change from trial-level to face-level analyses also reduced the effective sample size, but we found in practice that confidence intervals remained reasonably precise, even for the mediation analyses (Rothman & Greenland, 2018). Last, we added a composite measure of the three primary confusion measures given their very high correlation.

2. SUPPLEMENTARY RESULTS

2.1. Site and subject characteristics

Table S1: Site characteristics and sample sizes.

Site	Location	Analyzed n	Subjects	Compensation	Physical setting
University of Pennsylvania	Philadelphia, PA, USA	140	Undergraduates recruited from Wharton Behavioral Lab's subject pool	Cash	In a group in lab with private cubicles
Eotvos Lorand University	Budapest, Hungary	73	Undergraduates recruited from college courses	Course credit	In a group in lab with private cubicles
Politecnico di Milano	Milano, Italy	55	Volunteers recruited among undergraduates and experimenters' colleagues	None	Individually in quiet lab or room
Eindhoven University of Technology	Eindhoven, Netherlands	35	Undergraduate volunteers	None or cash (CHECK)	In a group in lab with private cubicles
Ithaca College	Ithaca, NY, USA	31	Undergraduates recruited from college courses	Extra credit in course	In a group in lab
Occidental College	Los Angeles, CA, USA	24	Undergraduates recruited from college subject pool	Course credit (CHECK)	In a group in lab

Table S2: Demographic characteristics of all analyzed subjects..

Characteristic	
Female (n (%))	256 (71.5)
Age (mean (sd))	21.54 (4.69)
Education (n (%))	
Did not graduate high school	2 (0.6)
Graduated 2-year college	9 (2.5)
Graduated 4-year college	46 (12.8)
Graduated high school	269 (74.9)
Post-graduate degree	33 (9.2)
Race/ethnicity (n (%))	
Black/African-American	48 (13.4)
Caucasian	208 (57.9)
Native American	5 (1.4)
East Asian	51 (14.2)
Hispanic	31 (8.6)
Middle Eastern	22 (6.1)
Southeast Asian	12 (3.3)
South Asian	17 (4.7)

2.2. Data exclusions

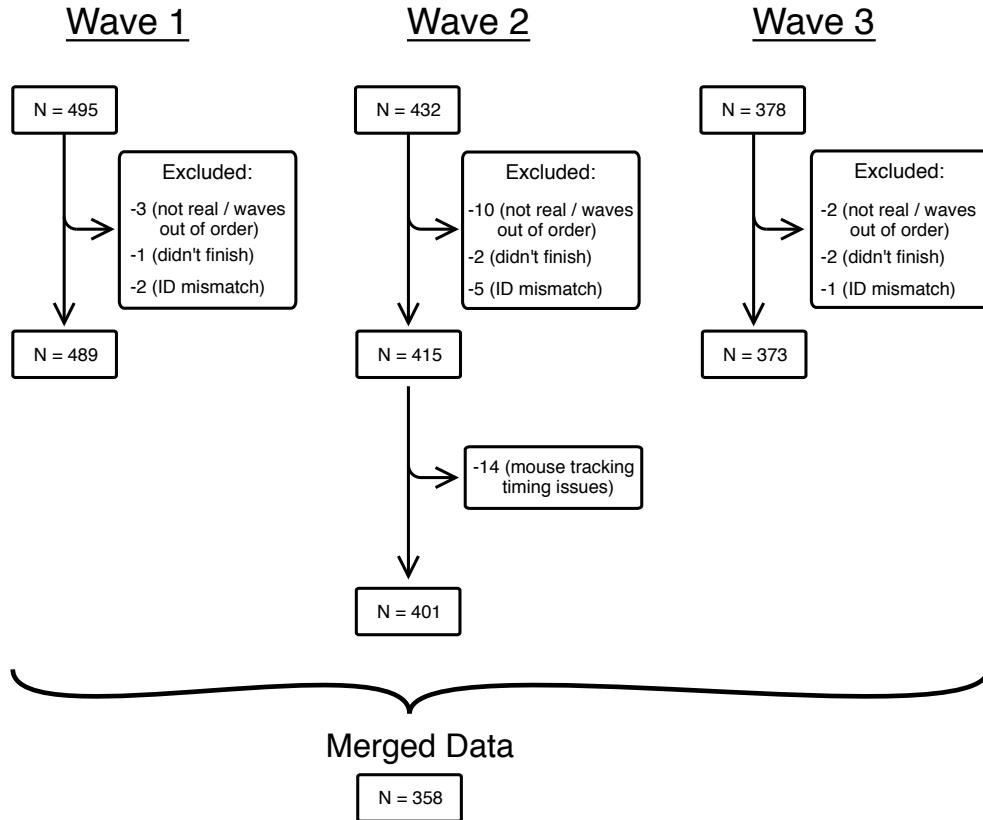


Figure S1: Data exclusions. “Not real”: the data represented pilot subjects or tests. “Waves out of order”: the site accidentally ran the 3 waves in the wrong order for some subjects. “ID mismatch”: the subject entered an invalid experiment ID.

2.3. Alternative model specifications

This section presents results for the initially planned analyses conducted at the individual trial level (55,753 data points) rather than conducted in aggregate at the face level. For the trial-level analysis, we used generalized estimation equations (GEE) with a working exchangeable correlation structure to flexibly account for correlation of observations within subjects and within faces. As a separate sensitivity analysis, we refit the main face-level analysis model without weighting faces by their inverse-variance of likability. Figure S1 compares the three model fits, suggesting that the weighted and unweighted face-level analyses were comparable. However, as described in the main text, the individual-level estimates of the relationship between MH score and likability appeared substantially attenuated compared to the face-level estimates. This is strongly suggestive that subjects’ individual ratings of MH score

are essentially noisy measurements of a face’s “true” MH score and hence that conducting analyses at the individual level would result in downward-biased estimates of the relationship between MH score and likability (Thomas et al., 1993).

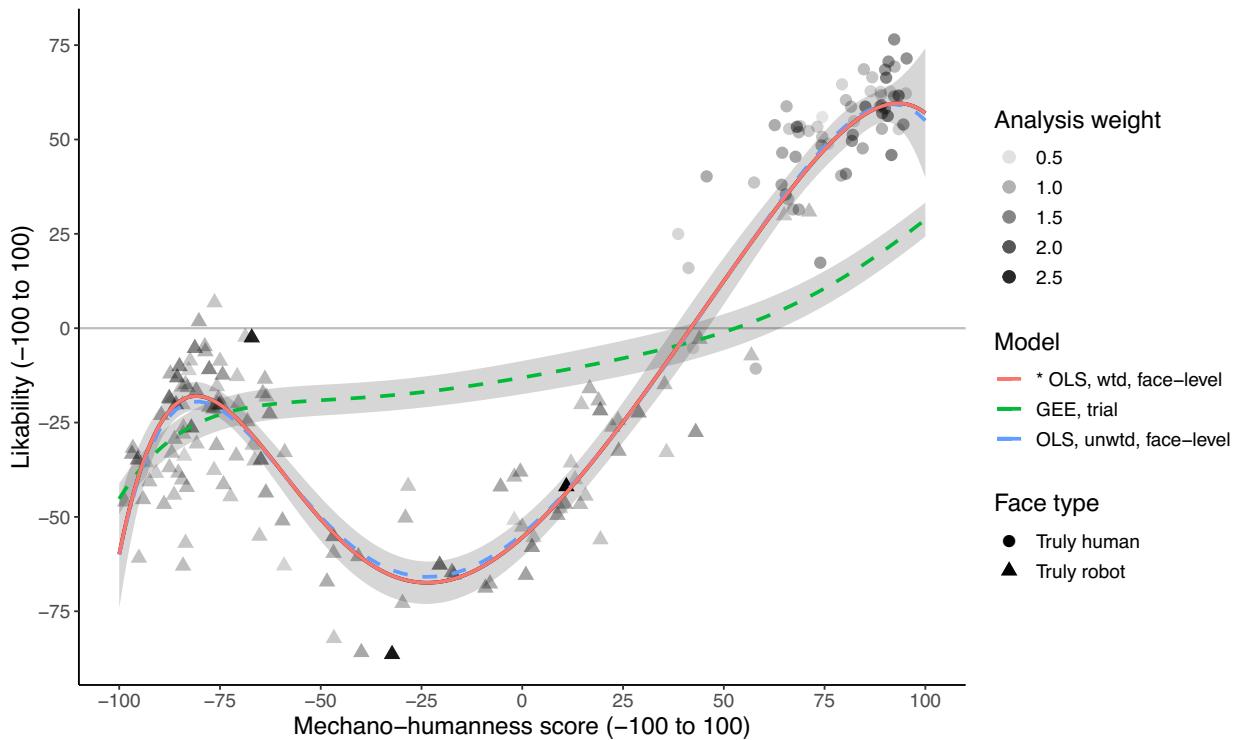


Figure S2: Uncanny Valley curves estimated in face-level analyses (red and blue lines) versus in trial-level analyses (green line). For visual clarity, data points are displayed only at the face level (183 points) and not at the trial level (55,753 points). *: Main face-level analysis model presented in main text. Triangles indicate faces that were actually robots; circles indicate faces that were actually humans. Point opacity is proportional to the face’s inverse-variance weight in the face-level analysis. Shaded bands are 95% pointwise confidence intervals for the fitted likability values when setting emotion to its mean.

REFERENCES

- Mathur, M. B., & Reichling, D. B. (2016). Navigating a social world with robot partners: A quantitative cartography of the uncanny valley. *Cognition*, 146, 22–32.
- Mathur, M. B., & Reichling, D. B. (in press). Open-source software for mouse-tracking in qualtrics to measure category competition. *Behavior Research Methods*. (Preprint: <https://osf.io/ymxau/>)

- Rothman, K. J., & Greenland, S. (2018). Planning study size based on precision rather than power. *Epidemiology*, 29(5), 599–603.
- Slijkhuis, P. J. (2017). *The uncanny valley phenomenon: A replication with short presentation times* (Unpublished master's thesis). University of Twente.
- Thomas, D., Stram, D., & Dwyer, J. (1993). Exposure measurement error: influence on exposure-disease relationships and methods of correction. *Annual Review of Public Health*, 14(1), 69–93.