

**Too Beautiful to be Fake: Attractive Faces are Less Likely to be Judged as
Artificially Generated**

Dominique Makowski¹, An Shu Te¹, Stephanie Kirk¹, Ngoi Zi Liang¹, Panagiotis Mavros²,
& S.H. Annabel Chen^{1, 3, 4, 5}

¹ School of Social Sciences, Nanyang Technological University, Singapore

² Singapore-ETH Centre, Future Cities Laboratory, Singapore

³ LKC Medicine, Nanyang Technological University, Singapore

⁴ National Institute of Education, Singapore

⁵ Centre for Research and Development in Learning, Nanyang Technological University,
Singapore

Correspondence concerning this article should be addressed to Dominique Makowski,
HSS 04-18, 48 Nanyang Avenue, Singapore (dom.makowski@gmail.com).

The authors made the following contributions. Dominique Makowski:
Conceptualization, Data curation, Formal Analysis, Funding acquisition, Investigation,
Methodology, Project administration, Resources, Software, Supervision, Validation,
Visualization, Writing – original draft; An Shu Te: Data curation, Project administration,
Resources, Investigation, Writing – original draft; Stephanie Kirk: Project administration,
Resources, Writing – original draft; Ngoi Zi Liang: Project administration, Resources,
Writing – review & editing; Panagiotis Mavros: Supervision, Writing – review & editing;
S.H. Annabel Chen: Project administration, Supervision, Writing – review & editing.

Correspondence concerning this article should be addressed to Dominique Makowski,
HSS 04-18, 48 Nanyang Avenue, Singapore. E-mail: dom.makowski@gmail.com

Abstract

Technological advances render the distinction between artificial (e.g., computer-generated faces) and real stimuli increasingly difficult, yet the factors driving our beliefs regarding the nature of ambiguous stimuli remain largely unknown. In this study, 150 participants rated 109 pictures of faces on 4 characteristics (attractiveness, beauty, trustworthiness, familiarity). The stimuli were then presented again with the new information that some of them were AI-generated, and participants had to rate each image according to whether they believed them to be real or fake. Strikingly, despite all images being pictures of real faces from the same database, most participants rated a large portion of them as “fake”. Moreover, our results suggest a gender-dependent role of attractiveness on reality judgements, with faces rated as more attractive being classified as more real. We also report links between reality beliefs tendencies and dispositional traits such as narcissism and paranoid ideation.

Significance Statement. Computer-generated images of faces are likely to become objectively indistinguishable from real photos in the near future, creating important issues in the context of fake news and misinformation, as well as virtual reality developments. Given the evolutionary importance of perceived attractiveness, we investigated if faces rated as more attractive would be more likely judged as real (vs “fake”, i.e., artificially generated). We indeed found a gender-dependent role of attractiveness on reality judgements, as well as a global influence of personality traits such as narcissism. These results are discussed in the light of consciousness psychology and evolutionary science, and are relevant to AI-researchers and misinformation management agencies.

Keywords: attractiveness, simulation monitoring, fiction, deep fakes, sense of reality

Word count: 4088

Too Beautiful to be Fake: Attractive Faces are Less Likely to be Judged as Artificially Generated

For the first time in human history, technology has enabled the creation of near-perfect simulations indistinguishable from reality. These artificial, yet realistic constructs permeate all areas of life through immersive works of fiction, deep fakes (real-like images and videos generated by deep learning algorithms), virtual and augmented reality (VR and AR), artificial beings (artificial intelligence “bots” with or without a physical form), fake news and skewed narratives, of which ground truth is often hard to access¹. Such developments not only carry important consequences for the technological and entertainment sectors, but also for security and politics - for instance if used for propaganda and disinformation, recruitment into malevolent organizations, or religious indoctrination². This issue is central to what has been coined the “post-truth era”³, in which the distinction (and lack thereof) between authentic and simulated objects will play a critical role.

While not all simulations have achieved perfect realism (e.g., Computer Generated Images - CGI in movies often lack certain key details that makes them visually distinct from real images)⁴, it is fair to assume that these technical limitations will become negligible in the near future, particularly in the field of face generation^{1,5,6}. Such performance, however, leads to a new question: if real and fake stimuli cannot be distinguished based on their objective characteristics, how can we make judgements regarding their nature?

Literature shows that the context surrounding a stimulus often plays an important role in the assessment of its reality (a process henceforth referred to as *simulation monitoring*)^{7,8}. With the extensive search and processing of cues within ambiguous stimuli being an increasingly complex and cognitively effortful strategy^{9,10}, people tend to draw on peripheral contextual cues (**Figure 1**), such as the source of the stimulus (e.g., in what

74 journal has information been published), and its credibility, authority and expertise, to
 75 help facilitate their evaluation⁹⁻¹¹. However, the automization and decontextualization of
 76 information allowed by online social media (where text snippets or video excerpts are
 77 mass-shared with little context) makes this task increasingly difficult^{12,13}. Thus, in the
 78 absence of clear contextual information, what drives our beliefs of reality?

Determinants of Simulation Monitoring

« Is this information *real* or *fake*? »

« *Real* » = genuine, authentic

« *Fake* » = artificial, simulated, deceptive

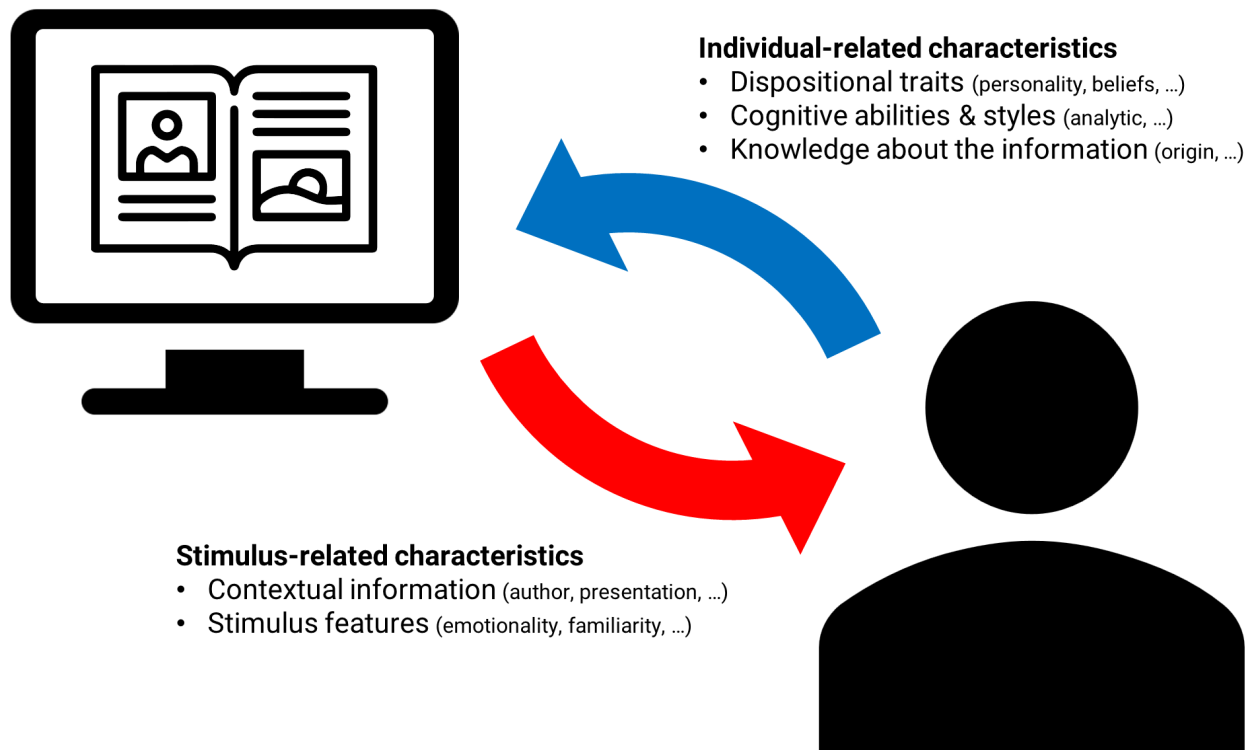


Figure 1. The decision to believe that an ambiguous stimulus (of any form, e.g., images, text, videos, environments, ...) is real or fake depends of individual characteristics (e.g., personality and cognitive styles), stimulus-related features (context, emotionality), and their interaction, which can manifest for instance in our bodily reaction.

79 Evidence suggests that inter-individual characteristics play a crucial role in
 80 simulation monitoring, with factors such as cognitive style, prior beliefs, and personality
 81 traits¹⁴⁻¹⁶. For instance, individuals with stronger analytical reasoning have been found to

82 better discriminate real from fake stimuli^{17,18}, and prior knowledge or beliefs about the
83 stimulus influences one’s perception of it by biasing the attention deployment towards
84 information that is in line with one’s expectations¹⁹. Furthermore, dispositional traits, such
85 as high levels of narcissism and low levels of openness and conscientiousness, have been
86 associated with greater susceptibility to fake news^{16,20}.

87 Beyond stimulus- and individual-related characteristics, evidence suggests that the
88 interaction between the two (i.e., the subjective reaction associated with the experience of
89 a given stimulus), contributes to simulation monitoring decisions. For instance, the
90 intensity of experienced emotions have been shown to increase one’s sense of presence - the
91 extent to which one feels like “being there”, as if the object of experience was real - when
92 engaged in a fictional movie or a VR environment^{21,22}. Conversely, beliefs that emotional
93 stimuli were fake (e.g., that emotional scenes were not authentic but instead involved
94 actors and movie makeup) were found to result in emotion down-regulation^{8,23}. In line with
95 these findings, studies on susceptibility to fake news have also found heightened stimulus
96 emotionality to be associated with greater belief^{24,25}. Additionally, other factors, such as
97 the stimuli’s perceived self-relevance^{26,27}, as well as familiarity²⁸, could also play a role in
98 guiding our appraisal of a stimulus.

99 AI-generated images of faces, due to their popularity as a target of CGI technology
100 and to the possibility of experimentally manipulating facial features, are increasingly used
101 to study face processing in relationship with saliency or emotions, as well as to other
102 important components of face evaluation, such as trustworthiness or attractiveness^{29–32}.
103 Interestingly, artificially created faces rated as more attractive (by an independent group of
104 raters) were perceived as less real⁵. Conversely,³³ reports that attractiveness ratings were
105 significantly lower when participants who were told that the faces were AI-generated were
106 compared to those who had no prior knowledge. Whereas this line of evidence suggests that
107 reality beliefs have an effect on face attractiveness ratings, the opposite question - whether

attractiveness could drive simulation monitoring - has received little attention to date.

This study primarily aims at exploring the effect of facial attractiveness on simulation monitoring, i.e., on the beliefs that an image is real or artificially generated. Based on the embodied reality theory^{7,8}, which suggests that salient and emotional stimuli are perceived to be more real, we hypothesize a quadratic relationship between perceived realness and attractiveness: faces rated as highly attractive or unattractive will more likely be believed to be real. We expect a similar relationship with trustworthiness ratings given its well-established link with attractiveness^{33–36}, and a positive relationship with familiarity (as more familiar faces would appear as more salient, self-relevant and anchored in reality). Additionally, we will further explore the link shared by dispositional traits, such as personality and attitude towards AI, with simulation monitoring tendencies. This study aims beyond the investigation of the discriminative accuracy between “true” photos and “true” artificially-generated images, focusing on the beliefs that a stimulus is real or fake, independently of its true nature.

Methods

Ethics Statement. This study was approved by the NTU Institutional Review Board (NTU IRB-2022-187) and all procedures performed were in accordance with the ethical standards of the institutional board and with the 1964 Helsinki Declaration. All participants provided their informed consent prior to participation and were incentivized after completing the study.

Procedure. In the first part of the study, participants answered a series of personality questionnaires, including the *Mini-IPIP6* (24 items)³⁷ measuring 6 personality traits, the *SIAS-6* and the *SPS-6* (6 items each)³⁸ assessing social anxiety levels, the *FFNI-BF* (30 items)³⁹ measuring 9 facets of narcissism; the *R-GPTS* (18 items)⁴⁰ measuring 2 dimensions related to paranoid thinking; and the *IUS-12* (12 items)⁴¹ measuring intolerance to uncertainty. Self-rated attractiveness was also assessed using 2

items - one measuring general attractiveness (“How attractive would you say you are?”)⁴² and the other measuring physical attractiveness (“How would you rate your own physical attractiveness relative to the average”)⁴³ Finally, we devised 5 items pertaining to expectations about AI-generated image technology (“I think current Artificial Intelligence algorithms can generate very realistic images”). To lower their saliency and the possibility of it priming the subjects about the task, we mixed these items with 5 items from the general attitudes towards AI scale (*GAAIS*)⁴⁴. This scale was presented after the social anxiety questionnaires. 3 attention check questions were also embedded in the surveys.

In the second part of this study, 109 images of neutral-expression faces from the validated American Multiracial Face Database (AMFD)⁴⁵ were presented to the participants for 500ms each, in a randomized order, following a fixation cross display (750 ms). After each stimulus presentation, ratings of *Trustworthiness* (“I find this person trustworthy”) and *Familiarity* (“This person reminds me of someone I know”) were collected using visual analog scales. Notably, as facial attractiveness is a multidimensional construct, encompassing evolutionary, sociocultural, biological as well as cognitive aspects^{46,47}, we assessed attractiveness using 2 visual analog scales, measuring general *Attractiveness* (“I find this person attractive”) and physical *Beauty* (“This face is good-looking”).

In the last part of the study, participants were informed that “about half” of the images previously seen were AI-generated (the instructions used a cover story explaining that the aim of the research was to validate a new face generation algorithm). The same set of stimuli was displayed again for 500 ms in a new randomized order. This time, after each display, participants were asked to express their belief regarding the nature of the stimulus using a visual analog scale (with *Fake* and *Real* as the two extremes). The study was implemented using *jsPsych*⁴⁸, and the exact instructions are available in the experiment code.

Participants. One hundred and fifty participants were recruited via *Prolific*, a crowd-sourcing platform recognized for providing high quality data⁴⁹. The only inclusion criterion was a fluent proficiency in English to ensure that the experiment instructions would be well-understood. Participants were incentivised with a reward of about £7.5 for completing the study, which took about 45 minutes to finish. Demographic variables (age, gender, sexual orientation, education and ethnicity) were self-reported on a voluntary basis.

We excluded 5 participants that either failed 2 ($\geq 66\%$) or more attention check questions, took an implausibly short time to finish the questionnaires or had incomplete responses. The final sample included 145 participants (Mean age = 28.3, SD = 9.0, range: [19, 66]; Sex: 48.3% females, 51.0% males, 0.7% others).

Data Analysis. The real-fake ratings (measured originally on a $[-1, 1]$ analog scale) were converted into two scores, corresponding to two conceptually distinct mechanisms: the dichotomous *belief* (real or fake, based on the sign of the rating) and the *confidence* (the rating's absolute value) associated with that belief. The former was analyzed using logistic mixed models, which modelled the probability of assigning a face to the real (≥ 0) as opposed to fake (< 0). The latter, as well as the other face ratings (attractiveness, beauty, trustworthiness and familiarity), was modelled using mixed beta regressions (suited for outcome variables expressed in percentages). The models included the participants and stimuli as random factors.

We started by investigating the effect of the procedure and instructions to check whether the stimuli (which were all images of real faces) were judged as fake in sufficient proportion to warrant their analysis. Additionally, we assessed the effect of the re-exposure delay, i.e., the time between the first presentation of the image (corresponding to the face ratings) and the second presentation (for the real-fake rating).

The determinants of reality beliefs were modelled separately for attractiveness, beauty, trustworthiness, and familiarity, using second order raw polynomials coefficients to

allow for possible quadratic relationships (**Figure 2**). Aside from attractiveness (conceptualized as a general construct), models for beauty, trustworthiness and familiarity were adjusted for the the two remaining variables *mutatis mutandis*. We took into account the gender of participants and stimuli by retaining the stimuli that were aligned with the participants' sexual preference (e.g., female faces for homosexual females, male faces for heterosexual females, and both for bisexual participants), and modeling the interaction with the participants' gender. For the attractiveness and beauty models, we then added the interaction with the reported self-attractiveness (the average of the two questions pertaining to it) to investigate its potential modulatory effect. Finally, we investigated the inter-individual correlates of simulation monitoring with similar models (but this time, for all items regardless of the participant's gender or sexual orientation) for each questionnaire, with all of the subscales as orthogonal predictors.

The analysis was carried out using *R* 4.2⁵⁰, the *tidyverse*⁵¹, and the *easystats* collection of packages^{52–56}. As all the details, scripts and complimentary analyses are open-access, we will focus in the manuscript on findings that are highly statistically significant ($p < .01$).

Results

Manipulation Check. Only one image file yielded a strong simulation monitoring bias ($> 85\%$), being classified as fake by 88.7% of participants. This image was removed from further analysis, leaving 108 trials per participant. On average, across participants, 44% of images (95%~CI [0.12, 0.64]) were judged as fake and 56% of images (95%~CI [0.36, 0.84]) as real. An intercept-only model with the participants and images as random factors showed that the Intraclass Correlation Coefficient (ICC), which can be interpreted as the proportion of variance explained by the random factors, was of 9.0% for the participants and 9.6% for the stimuli.

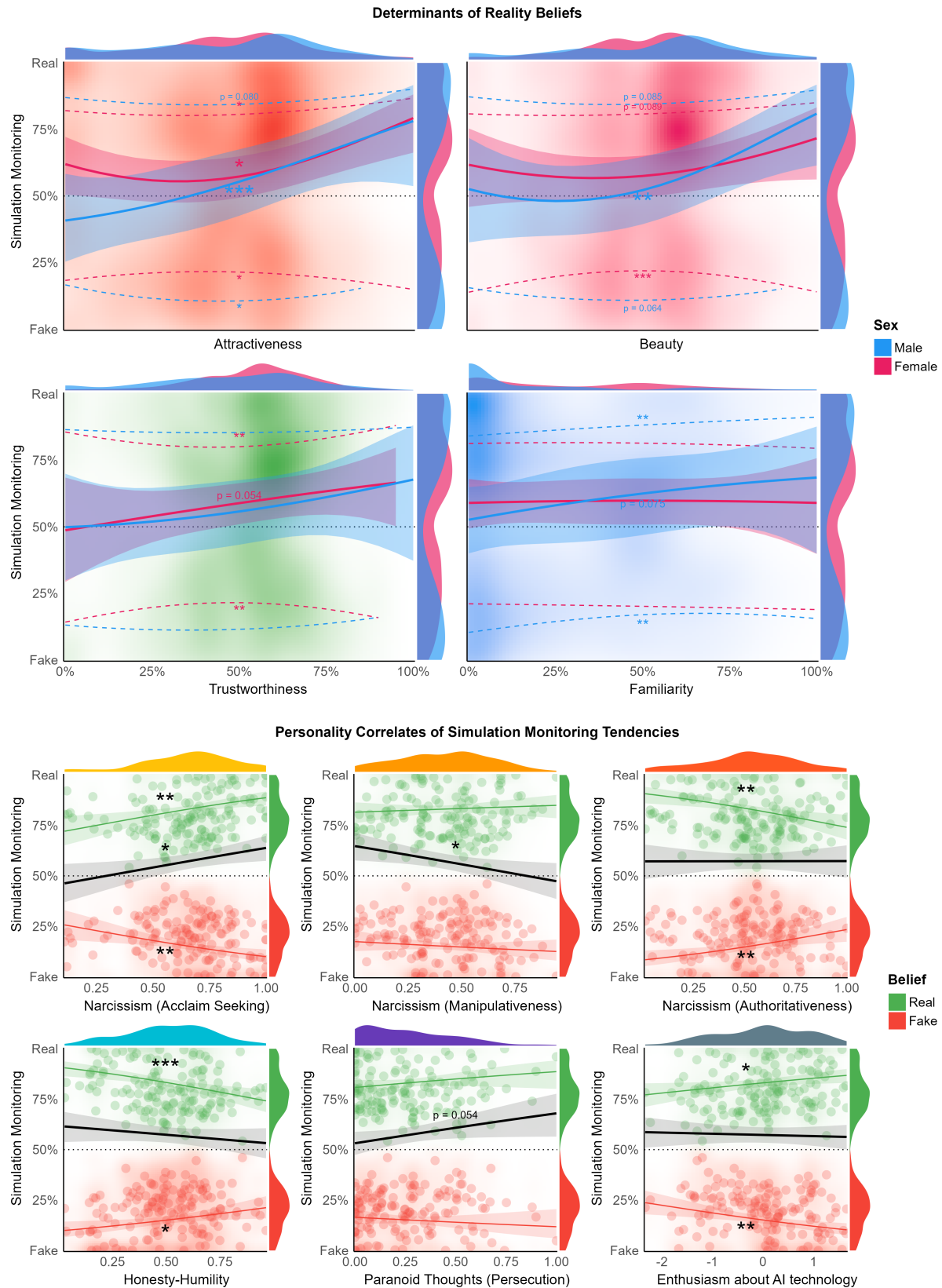


Figure 2. Top part shows the effect of face ratings on 1) the probability of judging a face as real vs. fake (solid line) and 2) on the confidence associated with that judgement (dashed lines) depending on the sex. Bottom part shows the effect of personality traits on the belief (black line) and the confidence associated with it (colored lines). The points are the average per participant confidence for both types of judgements. Stars indicate significance ($p < .001^{***}$, $p < .01^{**}$, $p < .05^{*}$).

While the delay of stimulus re-exposure stimulus did not have a significant effect on participants' beliefs of reality ($OR = 1.00$, 95% $CI = [0.99, 1.00]$), judgement confidence was found to be negatively associated with re-exposure delay when the faces were judged as real ($\beta = -0.006$, 95% $CI = [-0.1, 0.002]$, $p = .004$)

Determinants of Simulation Monitoring. Attractiveness had a significant positive and linear relationship ($R^2_{\text{marginal}} = 2.0\%$) with the belief that a stimulus was real ($\beta_{\text{poly1}} = 16.57$, 95% $CI = [7.33, 25.82]$, $z = 3.51$, $p < .001$) for males, and a quadratic relationship for females ($\beta_{\text{poly2}} = 7.82$, 95% $CI = [1.81, 13.84]$, $z = 2.55$, $p = .011$), with both non-attractive and attractive faces being judged as more real. Attractiveness was also found to have a significant positive and quadratic relationship with confidence in judging faces both as real ($\beta_{\text{poly2}} = 4.30$, 95% $CI = [0.97, 7.64]$, $z = 2.53$, $p = .011$) and as fake ($\beta_{\text{poly2}} = 5.23$, 95% $CI = [0.86, 9.60]$, $z = 2.35$, $p = .019$) for females. For males, however, a significant negative and quadratic relationship was found between attractiveness ratings and belief confidence only for faces judged as fake ($\beta_{\text{poly2}} = -9.92$, 95% $CI = [-18.99, -0.86]$, $z = -2.15$, $p = .032$). There was no interaction with reported self-attractiveness.

Beauty, adjusted for trustworthiness and familiarity, had a significant positive and linear relationship ($R^2_{\text{marginal}} = 2.0\%$) with the belief that a stimulus was real ($\beta_{\text{poly1}} = 11.82$, 95% $CI = [4.28, 20.21]$, $z = 2.76$, $p = .006$) for males only. No effect on confidence was found, aside from a quadratic relationship in females for faces judged as fake, suggesting that non-beautiful and highly beautiful faces were rated as fake with more confidence than average faces ($\beta_{\text{poly2}} = 7.84$, 95% $CI = [3.39, 12.29]$, $z = 3.46$, $p < .001$). There was no interaction with reported self-attractiveness.

Trustworthiness, adjusted for beauty and familiarity, had a predominantly positive and linear relationship ($R^2_{\text{marginal}} = 2.0\%$) with the belief that a stimulus was real ($\beta_{\text{poly1}} = 6.44$, 95% $CI = [-0.11, 13.00]$, $z = 1.93$, $p = .0054$) for females only. No effect on

confidence was found for males, whereas a quadratic relationship was found for females for both faces judged as real ($\beta_{poly2} = 6.14$, 95% $CI = [2.13, 10.14]$, $z = 3.00$, $p = .003$) as well as fake ($\beta_{poly2} = 6.12$, 95% $CI = [1.49, 10.75]$, $z = 2.59$, $p = .001$), suggesting that non-trustworthy and highly trustworthy faces were rated with more confidence than average faces.

We did not find any significant relationships for familiarity adjusted for beauty and trustworthiness ($R^2_{marginal} = 2.0\%$). However, a significant positive and linear relationship was found between familiarity and the confidence judgements of rating faces as real ($\beta_{poly1} = 9.98$, 95% $CI = [3.83, 16.13]$, $z = 3.18$, $p = .001$) whereas a negative linear relationship was found with those judged as fake ($\beta_{poly1} = -12.41$, 95% $CI = [-20.27, -4.54]$, $z = -3.09$, $p = .002$) for males only. This hence suggests that males more confidently judge faces as real with when they are familiar, and as fake when they are unfamiliar.

Inter-Individual Correlates of Simulation Monitoring. The models including the personality traits suggested that *Honesty-Humility* had a significant negative relationship with the confidence associated with real as well as fake judgements ($\beta_{real} = -1.62$, 95% $CI = [-2.55, -0.70]$, $z = -3.43$, $p < .001$; $\beta_{fake} = -1.16$, 95% $CI = [-2.09, -0.23]$, $z = -2.45$, $p = 0.014$).

Significant positive associations were found between the probability of judging faces as real and dimensions of narcissism such as *Acclaim Seeking* ($\beta = 2.24$, 95% $CI = [1.17, 4.27]$, $z = 2.44$, $p = .015$), and *Manipulativeness* ($\beta = 0.47$, 95% $CI = [0.25, 0.87]$, $z = 2.4$, $p = 0.017$). Confidence judgements also shared significant links with narcissism through various facets, such as a positive relationship between the confidence for both real and fake judgements with *Acclaim Seeking* ($\beta_{real} = 1.65$, 95% $CI = [0.59, 2.70]$, $z = 3.07$, $p = .002$; $\beta_{fake} = 1.62$, 95% $CI = [0.56, 2.68]$, $z = 3.00$, $p = .003$), and a negative relationship with *Authoritativeness* ($\beta_{real} = -1.57$,

95% $CI = [-2.58, -0.57]$, $z = -3.08$, $p = .002$; $\beta_{fake} = -1.49$, 95% $CI = [-2.50, -0.48]$,
 $z = -2.89$, $p = .004$).

A positive trend was found in the relationship between the *Persecutory Ideation*
dimension of paranoid thinking and the belief that the faces were real ($\beta = 1.87$,
95% $CI = [0.99, 3.54]$, $z = 1.93$, $p = .054$).

The *Prospective Anxiety* aspect of intolerance to uncertainty shared a negative trend
in its association with confidence ratings ($\beta_{real} = 1.43$, 95% $CI = [0.10, 2.76]$, $z = 2.10$,
 $p = .036$; $\beta_{fake} = -0.91$, 95% $CI = [-1.93, 0.11]$, $z = -1.75$, $p = .081$). No significant
effect was found for social anxiety.

Questions pertaining to the attitude towards AI were reduced to 3 dimensions
through factor analysis, labelled AI-Enthusiasm (loaded by items expressing interest and
excitement in AI development and applications), AI-Realness (loaded by items expressing
positive opinions on the ability of AI to create realistic material), and AI-Danger (loaded
by items expressing concerns on the unethical misuse of AI technology). However, only
AI-Enthusiasm displayed a significant positive relationship with the confidence in both real
and fake judgements ($\beta_{real} = 0.21$, 95% $CI = [0.02, 0.40]$, $z = 2.20$, $p = .028$; $\beta_{fake} = 0.31$,
95% $CI = [0.12, 0.50]$, $z = -8.90$, $p < 0.001$).

Discussion

This study aimed at investigating the effect of facial ratings (attractiveness, beauty,
trustworthiness and familiarity) on simulation monitoring, i.e., on the belief that a stimulus
was artificially generated. Most strikingly, despite all the stimuli being real faces from the
same database, all participants believed (to high degrees of confidence) that a significant
proportion of them were fake. This finding not only attests to the effectiveness of our
instructions, but highlights the current levels of expectation regarding CGI technology.
The strong impact of prior expectations and information on reality beliefs demonstrated

here underlines the volatility of our sense of reality. In fact, stimuli-related and participant-related characteristics accounted together for less than 20% of the beliefs variance, suggesting a large contribution of other subjective processes.

Although attractiveness did not seem to be the primary drive underlying simulation monitoring of face images, we do nonetheless report significant associations, with a different pattern observed depending on the participant's gender. The quadratic relationship found for female participants is aligned with our hypothesis that salient faces (i.e., rated as very attractive or very unattractive) are judged to be more real. The fact that this effect did not reach significance for beauty underlines that attractiveness judgement, and its role in simulation monitoring, is a multidimensional construct that cannot be reduced to physical facial attractiveness, in particular for women^{57,58}. In fact, female participants were more confident in judging faces as fake only when they were rated very high or low on beauty, suggesting that physical beauty and attractiveness are not analogous in their effects on simulation monitoring decisions.

Interestingly, we found a significant positive linear relationship in male participants for both attractiveness and beauty on simulation monitoring that we could interpret under an evolutionary lens. Specifically, males purportedly place more emphasis on facial attractiveness as a sign of reproductive potential, as compared with females, who tend to value characteristics signaling resource acquisition capabilities⁵⁷⁻⁵⁹. It is thus possible that the evolutionary weight associated with attractiveness skewed the perceived saliency of men towards attractive faces, rendering them significantly more salient than unattractive faces, and in turn distorted the relationship with simulation monitoring. However, future studies should test this saliency-based hypothesis by measuring constructs closer to salience and its effects, for instance using neuroimaging^{60,61} or physiological markers (e.g., heart rate deceleration)⁶².

Our results found a positive linear trend between trustworthiness and simulation

monitoring for females only. Given prior evidence that faces presented as computer-generated were rated less trustworthy^{30,33,63}, we expected such a linear association to be more clearly present for both genders. One of the underlying mechanisms that possibly contributed to this dimorphism could be the increased risk-taking aversion reported in females (explained evolutionarily as a compromise to their reproductive potential⁶⁴), to which perceived facial trustworthiness relates⁶⁵. However, if that was the case, faces judged as highly untrustworthy should have appeared as even more salient (representing an evolutionary threat), and hence be judged as more real, leading to a quadratic relationship between trustworthiness and simulation monitoring instead. Further studies are needed to investigate the causes of the increased simulation monitoring sensitivity to trustworthiness in females.

Contrary to our hypothesis, we did not find familiarity to be significantly related to simulation monitoring decisions. Interestingly, there were significant linear relationships between familiarity and confidence judgements for males only, where familiarity increased the confidence of reality beliefs. Although the familiarity measure was not a “recognition” measure, evidence from studies pertaining to the latter could be linked, reporting better face memory for females^{66–68}, as well as an overconfidence in face recall for males^{69,70}. However, it should be noted that the present study’s distribution of familiarity ratings was strongly skewed, and only a low number of pictures was rated as highly familiar. As such, future studies should clarify this point by experimentally manipulating familiarity, for instance by modulating the amount of exposure to items before querying the simulation monitoring judgements.

Regarding the role of inter-individual characteristics in simulation monitoring tendencies, we found higher scores of honesty-humility - a trait related to an increased risk perception and aversion^{71,72} - to be related to a lower confidence in simulation monitoring judgements. Notably, greater narcissistic tendencies in dimensions such as acclaim seeking

and manipulativeness were associated with a higher number of faces judged as real. This is in line with recent research which found people with narcissism to be less likely to engage in analytical reasoning strategies such as reflective thinking^{73,74}, and to be more vigilant and attentive to external stimuli⁷⁵⁻⁷⁷.

Moreover, putting the significant positive links between narcissistic acclaim seeking and confidence judgements in perspective with the negative correlation between honesty-humility and narcissism⁷⁸, we confirm previous evidence regarding the relationship between narcissistic grandiosity and over-confidence in decision-making⁷⁹⁻⁸². Although an inverse effect was found for the narcissistic facet of authoritativeness, we interpret this relationship as related to a higher response assertiveness. Taken together, these results suggest that participants with low humility and high recognition desires are more confident in their judgement regarding the real or fake nature of ambiguous stimuli. Alternatively, participants with opposite traits might perceive a higher risk in the decision-making process and its potential consequences (e.g., being seen as bad at the task at hand), resulting in more conservative confidence ratings.

Our findings suggest - though with weak significance - a positive link between paranoid ideation and the tendency to believe that the stimuli were real. Given previous reports that people with higher levels of paranoia are more sensitive to cues of social threat⁸³⁻⁸⁵, it is plausible that paranoid traits confer greater saliency and emotionality to observed faces, hence increasing perceptions of its realness. This hypothesis, if confirmed by future studies, would be in line with previous findings that persecutory delusions are predicted by a greater sense of presence in VR environments populated with virtual characters⁸⁶.

Despite the ubiquity of AI, the literature pertaining to the influence of people's AI attitudes on simulation monitoring is scarce. Contrary to our expectations, we did not find evidence for the role of participants' expectations regarding the capabilities of AI

technology (in terms of the realism of its productions). Instead, we found only one's enthusiasm about AI technology to be related to an increased confidence in simulation monitoring ratings. This could potentially be because participants with a highly positive attitude towards AI perceive themselves as having greater knowledge about AI and its capabilities⁸⁷, hence permitting themselves to be more confident in their simulation monitoring decisions. In fact, this result is in line with reports that AI attitudes interacts with people's perceived self-knowledge to influence their perception of the opportunities and risks accorded by AI applications⁸⁷.

On a methodological level, although the order of presentation of the facial images was randomized to reduce effects of adaptation, participants were more confident in their judgements for faces perceived as real following a shorter re-exposure delay. Such shorter durations could be associated with the faces being better remembered and appearing more familiar, thereby triggering self-referential and autobiographical memory processing during the repeated display⁸⁸⁻⁹⁰. Indeed, this finding is consistent with studies in which fictional stimuli that were associated with familiarity up-regulated emotions, biasing its salience and perceived realness^{22,26}. However, if that was the case, we would expect shorter re-exposure delays to impact the decision bias as well towards reality, rather than simply the confidence. Future studies should further investigate the modulatory effects of types and degrees of familiarity on perceived realness judgements.

Several limitations have to be noted. The current experimental paradigm required participants to judge the realness of faces they had prior exposure to (which was done to prevent reality judgements from influencing the other ratings). Although the effect of re-exposure delay was negligible, the potential bias induced by face familiarity (as compared to judging completely new items) cannot be discarded. Future studies could examine that by incorporating novel face images or increasing the duration of the re-exposure delay. Moreover, the magnitude of the effects found in the study was relatively

small, suggesting that the facial features measured in the study were not the key determinants of simulation monitoring. Hence, beyond exploring new potential mechanisms, future studies should include a more thorough debriefing to try to capture what conscious strategies (if any) the participants used (e.g., focusing on some features of the stimulus - like hair or eyes in the case of faces) to guide their reality beliefs.

In summary, the aim of the present study was to examine whether a subset of specific characteristics, in particular face attractiveness, significantly influences our simulation monitoring decisions. Notably, we found faces rated as attractive to be perceived as more real, with a possible sexual dimorphism affecting the shape of the relationship. We also found that inter-individual traits, such as narcissistic acclaim-seeking and manipulativeness, as well as persecutory ideation, were related to a systematic bias towards beliefs that the stimuli were real or fake. We believe that these findings provide the foundations to help us understand what drives reality beliefs in an increasingly reality-ambiguous world.

Data Availability

The datasets generated and/or analysed during the current study are available in the GitHub repository <https://github.com/RealityBending/FakeFace>

Funding

This work was supported by the Presidential Postdoctoral Fellowship Grant (NTU-PPF-2020-10014) from Nanyang Technological University (awarded to DM) and the Intra-CREATE Seed Collaboration Grant (NRF2021-ITS008-0010) from the National Research Foundation, Prime Minister's Office, Singapore, under its Campus for Research Excellence and Technological Enterprise (CREATE) programme (awarded to DM and PM).

Acknowledgments

414

415 We would like to thank Taong Ren Qing Malcolm for his contribution to the selection
416 of the materials.

References

1. Nightingale, S. J. & Farid, H. AI-synthesized faces are indistinguishable from real faces and more trustworthy. *Proceedings of the National Academy of Sciences* **119**, e2120481119 (2022).
2. Pantserev, K. The malicious use of AI-based deepfake technology as the new threat to psychological security and political stability. in 37–55 (2020). doi:10.1007/978-3-030-35746-7_3.
3. Lewandowsky, S., Ecker, U. K. & Cook, J. Beyond misinformation: Understanding and coping with the ‘post-truth’ era. *Journal of applied research in memory and cognition* **6**, 353–369 (2017).
4. McDonnell, R. & Breidt, M. Face reality: Investigating the uncanny valley for virtual faces. in *ACM SIGGRAPH ASIA 2010 sketches* 1–2 (2010).
5. Tucciarelli, R., Vehar, N. & Tsakiris, M. *On the realness of people who do not exist: the social processing of artificial faces*. <https://osf.io/dnk9x> (2020) doi:10.31234/osf.io/dnk9x.
6. Moshel, M. L., Robinson, A. K., Carlson, T. A. & Grootswagers, T. Are you for real? Decoding realistic AI-generated faces from neural activity. *Vision Research* **199**, 108079 (2022).
7. Makowski, D. Cognitive neuropsychology of implicit emotion regulation through fictional reappraisal. (Sorbonne Paris Cité, 2018).
8. Makowski, D. *et al.* Phenomenal, bodily and brain correlates of fictional reappraisal as an implicit emotion regulation strategy. *Cognitive, Affective, & Behavioral Neuroscience* **19**, 877–897 (2019).
9. Michael, R. B. & Sanson, M. Source information affects interpretations of the news across multiple age groups in the united states. *Societies* **11**, 119 (2021).

10. Susmann, M. W. *et al.* Persuasion amidst a pandemic: Insights from the elaboration likelihood model. *European Review of Social Psychology* 1–37 (2021).
11. Petty, R. E. & Cacioppo, J. T. The elaboration likelihood model of persuasion. in *Communication and persuasion* 1–24 (Springer, 1986).
12. Berghel, H. Weaponizing twitter litter: Abuse-forming networks and social media. *Computer* **51**, 70–73 (2018).
13. Chen, Y., Conroy, N. K. & Rubin, V. L. News in an online world: The need for an ‘automatic crap detector’. *Proceedings of the Association for Information Science and Technology* **52**, 1–4 (2015).
14. Bryanov, K. & Vziatysheva, V. Determinants of individuals’ belief in fake news: A scoping review determinants of belief in fake news. *PLoS one* **16**, e0253717 (2021).
15. Ecker, U. K. *et al.* The psychological drivers of misinformation belief and its resistance to correction. *Nature Reviews Psychology* **1**, 13–29 (2022).
16. Sindermann, C., Cooper, A. & Montag, C. A short review on susceptibility to falling for fake political news. *Current Opinion in Psychology* **36**, 44–48 (2020).
17. Pehlivanoglu, D. *et al.* The role of analytical reasoning and source credibility on the evaluation of real and fake full-length news articles. *Cognitive research: principles and implications* **6**, 1–12 (2021).
18. Pennycook, G. & Rand, D. G. Lazy, not biased: Susceptibility to partisan fake news is better explained by lack of reasoning than by motivated reasoning. *Cognition* **188**, 39–50 (2019).
19. Britt, M. A., Rouet, J.-F., Blaum, D. & Millis, K. A reasoned approach to dealing with fake news. *Policy Insights from the Behavioral and Brain Sciences* **6**, 94–101 (2019).
20. Piksa, M. *et al.* Cognitive processes and personality traits underlying four phenotypes of susceptibility to (mis) information. *Frontiers in psychiatry* 1142 (2022).

- 457
- 458 21. Sanchez-Vives, M. V. & Slater, M. From presence to consciousness through virtual
459 reality. *Nature Reviews Neuroscience* **6**, 332–339 (2005).
- 460 22. Makowski, D., Sperduti, M., Nicolas, S. & Piolino, P. ‘Being there’ and remembering
461 it: Presence improves memory encoding. *Consciousness and cognition* **53**, 194–202
(2017).
- 462 23. Sperduti, M. *et al.* The distinctive role of executive functions in implicit emotion
463 regulation. *Acta Psychologica* **173**, 13–20 (2017).
- 464 24. Martel, C., Pennycook, G. & Rand, D. G. Reliance on emotion promotes belief in
465 fake news. *Cognitive research: principles and implications* **5**, 1–20 (2020).
- 466 25. Bago, B., Rosenzweig, L. R., Berinsky, A. J. & Rand, D. G. Emotion may predict
467 susceptibility to fake news but emotion regulation does not seem to help. *Cognition
and Emotion* 1–15 (2022).
- 468 26. Sperduti, M. *et al.* The paradox of fiction: Emotional response toward fiction and
469 the modulatory role of self-relevance. *Acta psychologica* **165**, 53–59 (2016).
- 470 27. Goldstein, T. R. The pleasure of unadulterated sadness: Experiencing sorrow in fic-
471 tion, nonfiction, and "in person.". *Psychology of Aesthetics, Creativity, and the Arts*
3, 232 (2009).
- 472 28. Begg, I. M., Anas, A. & Farinacci, S. Dissociation of processes in belief: Source
473 recollection, statement familiarity, and the illusion of truth. *Journal of Experimental
Psychology: General* **121**, 446 (1992).
- 474 29. Sobieraj, S. & Krämer, N. C. What is beautiful in cyberspace? Communication with
475 attractive avatars. in *International conference on social computing and social media*
125–136 (Springer, 2014).
- 476 30. Balas, B. & Pacella, J. Trustworthiness perception is disrupted in artificial faces.
Computers in Human Behavior **77**, (2017).

- 477
478 31. Tsikandilakis, M., Bali, P. & Chapman, P. Beauty is in the eye of the beholder: The
appraisal of facial attractiveness and its relation to conscious awareness. *Perception*
479 **48**, 72–92 (2019).
- 480 32. Calbi, M. *et al.* How context influences our perception of emotional faces: A behav-
481 ioral study on the kuleshov effect. *Frontiers in Psychology* **8**, (2017).
- 482 33. Liefoghe, B. *et al.* Faces merely labelled as artificial are trusted less. (2022).
483
- 484 34. Bartosik, B., Wojcik, G. M., Brzezicka, A. & Kawiak, A. Are you able to trust
me? Analysis of the relationships between personality traits and the assessment of
485 attractiveness and trust. *Frontiers in Human Neuroscience* **15**, 685530 (2021).
- 486 35. Garrido, M. V. & Prada, M. KDEF-PT: Valence, emotional intensity, familiarity and
attractiveness ratings of angry, neutral, and happy faces. *Frontiers in Psychology* **8**,
487 2181 (2017).
- 488 36. Little, A. C., Jones, B. C. & DeBruine, L. M. Facial attractiveness: Evolutionary
based research. *Philosophical Transactions of the Royal Society B: Biological Sciences*
489 **366**, 1638–1659 (2011).
- 490 37. Sibley, C. *et al.* The mini-IPIP6: Validation and extension of a short measure of the
big-six factors of personality in new zealand. *New Zealand Journal of Psychology* **40**,
491 142–159 (2011).
- 492 38. Peters, L., Sunderland, M., Andrews, G., Rapee, R. M. & Mattick, R. P. Develop-
ment of a short form social interaction anxiety (SIAS) and social phobia scale (SPS)
using nonparametric item response theory: The SIAS-6 and the SPS-6. *Psychological*
493 *assessment* **24**, 66 (2012).

39. Jauk, E., Olaru, G., Schürch, E., Back, M. D. & Morf, C. C. Validation of the german five-factor narcissism inventory and construction of a brief form using ant colony optimization. *Assessment* 10731911221075761 (2022).
40. Freeman, D. *et al.* The revised green et al., Paranoid thoughts scale (r-GPTS): Psychometric properties, severity ranges, and clinical cut-offs. *Psychological Medicine* **51**, 244–253 (2021).
41. Carleton, R. N., Norton, M. P. J. & Asmundson, G. J. Fearing the unknown: A short version of the intolerance of uncertainty scale. *Journal of anxiety disorders* **21**, 105–117 (2007).
42. Marcinkowska, U. M., Jones, B. C. & Lee, A. J. Self-rated attractiveness predicts preferences for sexually dimorphic facial characteristics in a culturally diverse sample. *Scientific Reports* **11**, 1–8 (2021).
43. Spielmann, S. S., Maxwell, J. A., MacDonald, G., Peragine, D. & Impett, E. A. The predictive effects of fear of being single on physical attractiveness and less selective partner selection strategies. *Journal of Social and Personal Relationships* **37**, 100–123 (2020).
44. Schepman, A. & Rodway, P. Initial validation of the general attitudes towards artificial intelligence scale. *Computers in human behavior reports* **1**, 100014 (2020).
45. Chen, J. M., Norman, J. B. & Nam, Y. Broadening the stimulus set: Introducing the american multiracial faces database. *Behavior Research Methods* **53**, 371–389 (2021).
46. Rhodes, G. *et al.* The evolutionary psychology of facial beauty. *Annual review of psychology* **57**, 199 (2006).
47. Han, S. *et al.* Beauty is in the eye of the beholder: The halo effect and generalization effect in the facial attractiveness evaluation. *Acta Psychologica Sinica* **50**, 363 (2018).
48. De Leeuw, J. R. jsPsych: A JavaScript library for creating behavioral experiments in a web browser. *Behavior research methods* **47**, 1–12 (2015).

49. Peer, E., Rothschild, D., Gordon, A., Evernden, Z. & Damer, E. Data quality of platforms and panels for online behavioral research. *Behavior Research Methods* **54**, 1643–1662 (2022).
50. R Core Team. *R: A language and environment for statistical computing*. (R Foundation for Statistical Computing, 2022).
51. Wickham, H. *et al.* Welcome to the tidyverse. *Journal of Open Source Software* **4**, 1686 (2019).
52. Lüdecke, D., Waggoner, P. & Makowski, D. Insight: A unified interface to access information from model objects in R. *JOSS* **4**, 1412 (2019).
53. Makowski, D., Ben-Shachar, M. & Lüdecke, D. bayestestR: Describing effects and their uncertainty, existence and significance within the Bayesian framework. *JOSS* **4**, 1541 (2019).
54. Lüdecke, D., Ben-Shachar, M., Patil, I., Waggoner, P. & Makowski, D. performance: An R package for assessment, comparison and testing of statistical models. *JOSS* **6**, 3139 (2021).
55. Lüdecke, D., Ben-Shachar, M., Patil, I. & Makowski, D. Extracting, computing and exploring the parameters of statistical models using R. *JOSS* **5**, 2445 (2020).
56. Makowski, D., Ben-Shachar, M., Patil, I. & Lüdecke, D. Methods and algorithms for correlation analysis in R. *JOSS* **5**, 2306 (2020).
57. Buunk, B. P., Dijkstra, P., Fetchenhauer, D. & Kenrick, D. T. Age and gender differences in mate selection criteria for various involvement levels. *Personal relationships* **9**, 271–278 (2002).
58. Qi, Y. & Ying, J. Gender biases in the accuracy of facial judgments: Facial attractiveness and perceived socioeconomic status. *Frontiers in Psychology* **13**, (2022).

59. Fink, B., Neave, N., Manning, J. T. & Grammer, K. Facial symmetry and judgements of attractiveness, health and personality. *Personality and Individual differences* **41**, 491–499 (2006).
60. Lou, B., Hsu, W.-Y. & Sajda, P. Perceptual salience and reward both influence feedback-related neural activity arising from choice. *Journal of Neuroscience* **35**, 13064–13075 (2015).
61. Indovina, I. & Macaluso, E. Dissociation of stimulus relevance and saliency factors during shifts of visuospatial attention. *Cerebral Cortex* **17**, 1701–1711 (2007).
62. Skora, L., Livermore, J. & Roelofs, K. The functional role of cardiac activity in perception and action. *Neuroscience & Biobehavioral Reviews* 104655 (2022).
63. Hoogers, E. The effect of attitude towards computer generated faces on face perception. (2021).
64. Van Den Akker, O. R., Assen, M. A. van, Van Vugt, M. & Wicherts, J. M. Sex differences in trust and trustworthiness: A meta-analysis of the trust game and the gift-exchange game. *Journal of Economic Psychology* **81**, 102329 (2020).
65. Hou, C. & Liu, Z. The survival processing advantage of face: The memorization of the (un) trustworthy face contributes more to survival adaptation. *Evolutionary Psychology* **17**, 1474704919839726 (2019).
66. Mishra, M. V. *et al.* Gender differences in familiar face recognition and the influence of sociocultural gender inequality. *Scientific reports* **9**, 1–12 (2019).
67. Lewin, C. & Herlitz, A. Sex differences in face recognition—women’s faces make the difference. *Brain and cognition* **50**, 121–128 (2002).
68. Sommer, W., Hildebrandt, A., Kunina-Habenicht, O., Schacht, A. & Wilhelm, O. Sex differences in face cognition. *Acta psychologica* **142**, 62–73 (2013).
69. Bailey, A. A gender in-group effect on facial recall. (University of Tasmania, 2021).

- 555
556 70. Herbst, T. H. Gender differences in self-perception accuracy: The confidence gap and
women leaders' underrepresentation in academia. *SA Journal of Industrial Psychology*
557 **46**, 1–8 (2020).
- 558 71. Levidi, M. D. C., McGrath, A., Kyriakoulis, P. & Sulikowski, D. Understanding crim-
inal decision-making: Links between honesty-humility, perceived risk and negative
559 affect: Psychology, crime & law. *Psychology, Crime and Law* 1–29 (2022).
- 560 72. Weller, J. A. & Thulin, E. W. Do honest people take fewer risks? Personality correlates
of risk-taking to achieve gains and avoid losses in HEXACO space. *Personality and*
561 *individual differences* **53**, 923–926 (2012).
- 562 73. Littrell, S., Fugelsang, J. & Risko, E. F. Overconfidently underthinking: Narcissism
negatively predicts cognitive reflection. *Thinking & Reasoning* **26**, 352–380 (2020).
563
- 564 74. Ahadzadeh, A. S., Ong, F. S. & Wu, S. L. Social media skepticism and belief in
conspiracy theories about COVID-19: The moderating role of the dark triad. *Current*
565 *Psychology* 1–13 (2021).
- 566 75. Carolan, P. L. Searching 'ineffectively': A behavioral, psychometric, and electroen-
cephalographic investigation of psychopathic personality and visual-spatial attention.
567 (Arts & Social Sciences: Department of Psychology, 2017).
- 568 76. Grapsas, S., Brummelman, E., Back, M. D. & Denissen, J. J. The 'why' and 'how' of
narcissism: A process model of narcissistic status pursuit. *Perspectives on Psycho-*
569 *logical Science* **15**, 150–172 (2020).
- 570 77. Eddy, C. M. Self-serving social strategies: A systematic review of social cognition in
narcissism. *Current Psychology* 1–19 (2021).
571
- 572 78. Hodson, G. *et al.* Is the dark triad common factor distinct from low honesty-humility?
573 *Journal of Research in Personality* **73**, 123–129 (2018).

79. Campbell, W. K., Goodie, A. S. & Foster, J. D. Narcissism, confidence, and risk attitude. *Journal of behavioral decision making* **17**, 297–311 (2004).
80. Chatterjee, A. & Pollock, T. G. Master of puppets: How narcissistic CEOs construct their professional worlds. *Academy of Management Review* **42**, 703–725 (2017).
81. Brunell, A. B. & Buelow, M. T. Narcissism and performance on behavioral decision-making tasks. *Journal of Behavioral Decision Making* **30**, 3–14 (2017).
82. O'Reilly, C. A. & Hall, N. Grandiose narcissists and decision making: Impulsive, overconfident, and skeptical of experts—but seldom in doubt. *Personality and Individual Differences* **168**, 110280 (2021).
83. Fornells-Ambrojo, M. *et al.* How do people with persecutory delusions evaluate threat in a controlled social environment? A qualitative study using virtual reality. *Behavioural and Cognitive Psychotherapy* **43**, 89–107 (2015).
84. Freeman, D. *et al.* Can virtual reality be used to investigate persecutory ideation? *The Journal of nervous and mental disease* **191**, 509–514 (2003).
85. King, A. & Dudley, R. Paranoia, worry, cognitive avoidance and intolerance of uncertainty in a student population. *Journal of Applied Psychology and Social Science* **3**, 70–89 (2017).
86. Freeman, D. *et al.* The psychology of persecutory ideation II: A virtual reality experimental study. *The Journal of nervous and mental disease* **193**, 309–315 (2005).
87. Said, N. *et al.* An artificial intelligence perspective: How knowledge and confidence shape risk and opportunity perception. (2022).
88. Abraham, A. & Von Cramon, D. Y. Reality= relevance? Insights from spontaneous modulations of the brain's default network when telling apart reality from fiction. *PloS one* **4**, e4741 (2009).
89. Gobbini, M. I. *et al.* Prioritized detection of personally familiar faces. *PloS one* **8**, e66620 (2013).

595

- 596 90. Taylor, M. J. *et al.* Neural correlates of personally familiar faces: Parents, partner
597 and own faces. *Human brain mapping* **30**, 2008–2020 (2009).