# Are you for real? Decoding hyperrealistic AI-generated faces from neural activity

Michoel L. Moshel [a,b]*, Amanda K. Robinson[b], Thomas A. Carlson [b], Tijl Grootswagers[b,c]

[a] School of Psychology, Macquarie University, NSW, Australia

[b] School of Psychology, University of Sydney, NSW, Australia

[c] The MARCS Institute for Brain, Behaviour and Development, Western Sydney University, NSW, Australia

* corresponding author: michoel.moshel@students.mq.edu.au

## Abstract

Can we trust our eyes? Until recently, we rarely had to question whether what we see is indeed what exists, but this is changing. Artificial neural networks can now generate hyperrealistic images that challenge our perception of what is real. This new reality can have significant implications in cybersecurity, counterfeiting, fake news, and border security. We investigated how the human brain encodes and interprets hyperrealistic artificially generated images using behaviour and brain imaging. We found that we could reliably detect AI-generated fake images using neural activity, even though people could not consciously report seeing differences between real and fake images. Understanding this dissociation between brain and

20behaviour may be key in determining the 'real' in our new reality. Stimuli, code,

21and data for this study can be found at https://osf.io/n2z73/.

22Introduction

23The novel and rapidly emerging phenomena of fake multimedia have swept

24through modern culture to the extent that the fake has become the expected norm

25(Adelani et al., 2020; Shen et al., 2019; Shu et al., 2017). The degree to which

26terms like 'fake news' or 'photoshopped' have become common parlance is

27indicative of a general and commonly experienced inability to distinguish between

28what is real and what is not (Fletcher, 2018). Meanwhile, AI technologies, in

29particular Generative Adversarial Networks (GANs), have been making

30increasingly rapid advances in generating realistic images with face generation as

31a major focus (Karras et al., 2019, 2020; Wang et al., 2018; Yu et al., 2020). These

32advances in realism have begun to have real-world consequences including

33undetectable videos of fake events ("Deepfakes": Kietzmann et al., 2020), art and

34audio-visual counterfeits (Farokhmanesh, 2018), and fraudulent social media

35accounts (Gleicher, 2019). For instance, in 2019, Facebook announced that fake

36accounts were being created with profile pictures generated by artificial

37intelligence in an attempt to evade detection (Gleicher, 2019). Crucially,

38understanding how people respond to AI images, in terms of both behaviour and

39neural responses, will inform us about how realistic artificial images and faces are

40perceived differently to real ones, how this dissociation is encoded by the brain,

41and can ultimately aid in the development of future policy and strategies to curb

42the potentially nefarious uses of fake media.

2

43One area in which AI technology has made increasingly rapid and apparent 44progress in is the generation of realistic faces. Until now, fooling observers with 45artificial faces has been a particularly difficult task to achieve given the expertise 46humans have with face perception and recognition (Farid & Bravo, 2007, 2012; 47Gauthier & Tarr, 2002; Sinha et al., 2006). Not only are faces perceived differently 48than objects (Shakeshaft & Plomin, 2015; Sunday et al., 2019) but neuroimaging 49studies highlight distinct brain networks for face processing (Axelrod & Yovel, 502015; Gauthier & Tarr, 2002). The specialized and expert processing of faces 51results in the rapid and automatic detection of artificial face appearance (Wheatley 52et al., 2011). For example, the uncanny valley effect describes how observers 53remain viscerally aware of artificial faces indicated by a steady drop in affinity as 54an artificial face approaches human likeness, despite not being able to identify any 55perceivable defects (MacDorman & Chattopadhyay, 2016). In another example, 56photographs of real faces yield a higher recognition accuracy than computer-57generated equivalents demonstrative of enhanced face expertise for the former 58(Crookes et al., 2015). Likewise, observers have typically performed well at 59discriminating human faces from computer-generated faces depending on image 60resolution, training, and incentives (Holmes et al., 2016). However, more recent 61studies have shown increasingly poorer performances at telling real from fake 62(Mader et al., 2017; Nightingale et al., 2017; Sanders et al., 2019; Zhou et al., 632019). As the capacity for image realism is steadily increasing, identification of 64fake faces will likely be further challenged.

65Neuroimaging has provided useful insight into how face perception unfolds over 66time. Electroencephalography (EEG), which measures electrical activity at the 67scalp with very high temporal resolution, has been used to identify unique neural

3

68responses that reflect the temporal emergence and dynamics of facial processing 69(Bentin et al., 1996; Rossion et al., 2000). Wheatley and colleagues (2011) 70demonstrated the brain's discrimination of real and artificial faces by comparing 71neural responses to real faces with responses to doll faces. The authors found that 72both human and artificial faces elicited an N170, a face-specific neural response 73approximately 170ms after image presentation. However, sustained positivity 74beyond 400ms was associated only with human faces, suggesting that this EEG 75potential could index a process that distinguishes between real and fake faces 76(Wheatley et al., 2011). Indeed, in other studies, sustained positivity, characterised 77by the late positive amplitude (LPP), increased as face realism increased, 78suggesting that real faces, more so than artificial faces, engage high-level 79attentional, semantic and identity evaluations (Schindler et al., 2017). The new 80generation of realistic faces produced by GAN technology, however, is of a far 81superior quality than previously studied artificial faces and often practically 82indistinguishable from real faces. Whether the brain elicits neural indicators 83consistent with artificial fake detection for the new generation of GAN-produced 84images has yet to be seen. Considering that humans remain the gold standard of 85fake image and face detection (Natsume et al., 2019, Marra et al., 2018), 86examining the neural mechanisms in fake face detection is instrumental in 87understanding how to best tackle and understand the new age of fake media. EEG 88remains an ideal method to provide useful insights into the neural processing of 89fake GAN faces. Firstly, it allows for an insight into the sequential stages of face 90processing, from low-level visual features to holistic face perception. Secondly, 91closer examination at the neuronal population level enables us to answer at what 92temporal stages GAN face perception may differ from real face perception. Thirdly,

4

93using newer multivariate methods applied to EEG data enables analysis of signal-
94level information on a trial-by-trial basis and can pinpoint the precise temporal
95emergence of visual processing (Grootswagers, Robinson, & Carlson, 2019;
96Haynes & Rees, 2006; Teichmann et al., 2020).

97With progressive advances in realistic image generation, have we reached a point
98where observers can no longer tell apart real from the fake? Can measuring the
99brain's response reveal how hyper-realistic fake faces are distinguished from real
100faces? We measured whether observers could behaviourally discriminate real faces
101from GAN-generated faces at two levels of face realism; one level of realism similar
102to fake images used in previous work ("unrealistic"), and another level which
103represents the current state-of-the-art hyper-realistic artificial images ("realistic").
104We expected that participants would not be able to discriminate real from realistic
105faces but could for unrealistic faces, consistent with previous research using AI-
106generated faces (Hulzebosch et al., 2020; Zhou et al., 2019). To investigate
107whether we could decode real and fake images from brain activity we used time-
108resolved multivariate pattern analysis (MVPA) and EEG. To ensure the real and
109fake stimuli evoked typical categorical effects that could be decoded in the neural
110signal,  we also included cars and bedrooms stimuli. We presented images upright
111in rapid sequences, which we have previously shown captures low- and high-level
112image processing (Grootswagers, Robinson, & Carlson; Oosterhof et al., 2016). To
113determine the contribution of low-level image properties, we used a much faster
114presentation rate (20Hz; Robinson et al., 2019) and also investigated how real/fake
115face processing is affected by image inversion, which limits high-level expert face-
116processing. Consistent with the brain's sensitivity to artificial face appearance, we
117found it was possible to decode real faces from GAN-generated faces at both levels

118of face realism using the EEG data. However, when asked to behaviourally classify

119faces as either real or fake, a large group of participants could differentiate the

120unrealistic, but not the realistic fake faces. Understanding dissociations between

121observer-reported perceptions of fake images and the brain's response can yield

122important insights into human face perception in general as well as raise

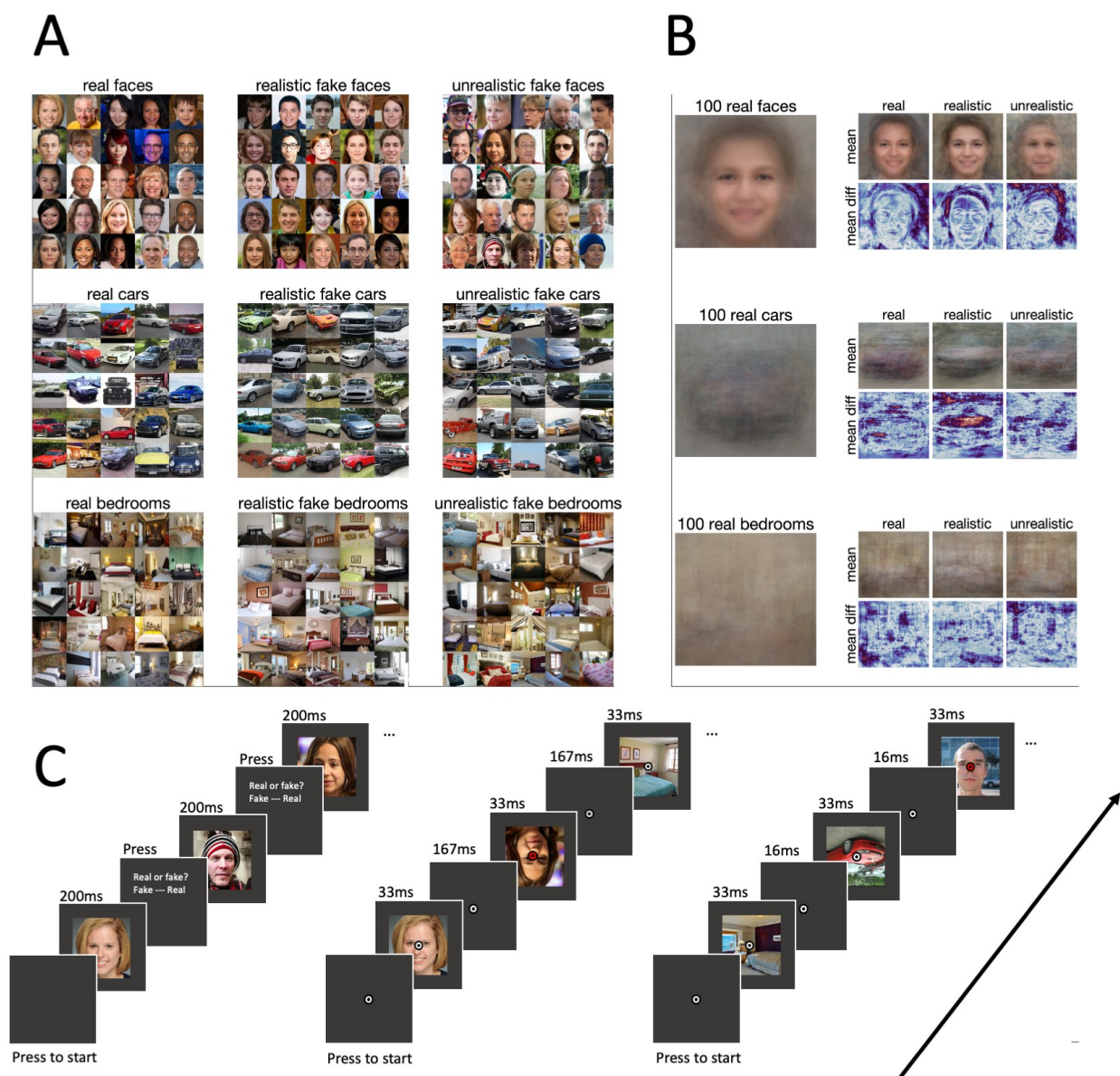123possibilities for training observers to tell apart real from fake.



124

125**Figure 1. Stimuli and design.** Experimental stimuli and design. A) Face, car and 126bedroom stimuli used in the experiment from three conditions (real, realistic fake, 127unrealistic fake), taken from StyleGAN. B) Mean image for each condition and the 128absolute pixel difference between 100 independent real images not used in the 129experiment. Brighter colours (orange) indicate greater absolute differences. C) 130Experimental designs from left to right; behavioural experiment, 5Hz EEG 131experiment and 20Hz EEG experiment.

132Methods

133We performed two experiments that investigated fake versus real image 134identification: one behavioural and one neuroimaging. The stimuli, data, and 135analysis code can be found at https://osf.io/n2z73/.

136Participants

137For behavioural testing, we recruited 200 participants from Amazon Mechanical 138Turk (MTurk) in return for payment. For the EEG component, 22 participants (15 139females, 7 males; mean age 20, range: 18-28) were recruited from the University 140of Sydney in return for course credit. Subjects all had normal or corrected-to-141normal vision and had no reported history of psychiatric or neurological disorders. 142The study was approved by the Human Ethics Committee of the University of 143Sydney. Verbal and written consent was obtained from each participant.

144Stimuli & Design

145GAN-generated stimuli were obtained from StyleGAN output found at 146shorturl.at/josOY (Karras et al., 2019). For a full description of the StyleGAN 147generative procedure and output, see Karras et al. (2019). Fake stimuli consisted 148of 25 faces, cars, and bedrooms at truncation levels of Ψ0.5 (realistic) and Ψ1.0 149(unrealistic), (Figure 1A). To best match image statistics across real and fake

150images, real images were obtained from training images used for GAN output.
151These real training faces were obtained from the Flickr-Faces-HQ dataset (Karras
152et al., 2019). Real cars and bedrooms were randomly selected from the LSUN
153dataset (Yu et al., 2015). To maintain consistent aspect ratios, all images were
154cropped to a square aspect ratio and resized to a 256 × 256 pixel dimension. No
155other filtering or editing was applied to the stimuli in order to provide a
156naturalistic demonstration of visual processing. To reduce obvious surface-level
157inconsistencies between real and fake images, real faces with eyes not facing
158frontward and/or with overly pronounced facial expressions (e.g. crying, laughing)
159were excluded. Upon surface inspection, we found no consistent delineating
160features between the real and fake bedrooms and cars. All images were presented
161in both upright and inverted orientations totalling 450 stimuli overall (Figure 1A).

162Behavioural testing for real versus fake face discrimination was conducted online
163(Grootswagers, 2020). The experiment was programmed in jsPsych (De Leeuw,
1642015) and hosted on Pavlovia.org (Peirce, 2019). Two hundred participants
165performed real or fake face judgements for one of four comparisons (50 in each
166group): 1) upright unrealistic vs upright real, 2) upright realistic vs upright real, 3)
167inverted unrealistic vs inverted real, and 4) inverted realistic vs inverted real. Each
168observer was shown 50 images in total: 25 fake and 25 real. Participants were
169informed that 50% of the images were real photos and 50% were computer-
170generated and were instructed to choose whether each image was real or fake.
171Each image was individually presented on the screen for 200ms, followed by a
172blank screen until the participant pressed a button to indicate if the face was real
173or fake. Stimuli were presented at 256 x 256 pixel dimension against a grey

8

174background. Presentation of images was randomised, and each image was only
175presented once. The experiment took around 3-5 minutes to complete (Figure 1C).

176For the EEG component, the experiment was presented in Psychopy2 (Peirce et al.,
1772019). Participants sat in a dimly lit room approximately 60cm away from a 1920 x
1781080 pixel Asus computer monitor. Stimuli subtended approximately 6.4 degrees
179visual angle on a grey background with a white fixation circle superimposing the
180stimuli at approximately 1.3 degrees. Images were presented in a rapid serial
181visual presentation (RSVP) paradigm, whereby stimuli are presented in rapid
182succession, at 20Hz and 5Hz sequences (33ms image duration and 167ms or 16ms
183gap). There were 20 sequences at each presentation rate comprising 40 in total
184with 18,000 images presented overall (with 20 repeats of each stimulus at each
185presentation rate). A sequence was started with a button press and lasted
186approximately 40 seconds. Subjects were instructed to fixate upon a white circle
187superimposed over each stimulus at the centre of the screen and told to respond
188by pressing any button on a 4-way button box whenever they spotted the fixation
189circle turn red (Figure 1C). Fixation colour changes were randomised to occur
190between 2 and 5 times in each sequence. Length of colour change corresponded to
191the time of one image presentation (33ms). At the conclusion of the experiment,
192participants were debriefed and informed that half the images had been fake.

193EEG recordings and preprocessing

194Continuous EEG data were recorded using a 64-electrode Brain Products EEG cap
195(Standard 64Ch actiCAP; GmbH, Herrsching, Germany) at a sample rate of 1000-
196Hz. Ag/AgCl active electrodes were placed in accordance with a 10/20
197international system (Oostenveld & Praamstra, 2001). Electrode gel was applied to

9

198the scalp under each electrode, aiming to reduce signal impedances to below
19910kΩ. Stimulus onset was synchronised to the EEG using transistor-transistor logic
200(TTL) pulses from the stimulus presentation computer to a separate recording
201computer. Pre-processing of the EEG data was computed offline using EEGLAB
202(Delorme & Makeig, 2004). The continuous EEG data were filtered with a high-
203pass filter of 0.1-Hz and a low-pass filter of 100-Hz and re-referenced to the
204average of all electrodes. No notch filter was applied. The data were then
205separated into epochs corresponding to stimulus presentation ranging from 100ms
206to 1000ms pre and post-stimulus onset. This produced 180,000 pre-processed
207epochs for each participant.

208Decoding analysis

209Time-resolved MVPA decoding analysis of EEG data was implemented in MATLAB
210with the CoSMoMVPA toolbox (Oosterhof, Connolly, & Haxby, 2016). We used
211Linear Discriminant Analysis (LDA) classifiers as implemented in CoSMoMVPA in a
212leave-one-out cross-validation scheme. The LDA classifier estimated the probability
213of EEG data belonging to a certain group (e.g., real or fake) where the higher
214estimate is the predicted class (Grootswagers, Wardle, & Carlson, 2017). This was
215repeated at every time point, for every exemplar, and averaged across subjects to
216generate the mean cross-validation decoding performance at each time point.
217Classification performance was characterized as significant if it produced an
218above-chance accuracy (>50% for real versus fake decoding or 33% for 3-way
219category decoding). An above-chance decoding accuracy informs us that the EEG
220data contains information relevant the contrast of interest (Grootswagers, Wardle,
221& Carlson, 2017; Olivetti et al., 2012; Pereira et al., 2009).

222Category Decoding Analysis

223We performed a category decoding analysis to investigate whether there were
224meaningful differences among the face, car and bedroom stimuli. We used an
225image-by-sequence cross-validation approach (Grootswagers, Robinson, & Carlson,
2262019), which entailed training the classifier on all-but-one image from each of the
227three categories from all-but-one sequence and testing the classifier on left-out
228images from the left-out sequence. This ensured that the classifier had to
229generalize to novel exemplars to successfully decode between faces, cars, and
230bedrooms for each of the real, realistic, and unrealistic conditions (Carlson et al.,
2312013). Decoding accuracy was characterized by an above-chance classifier
232performance (>33%). Contrasts were broken down into presentation rate (5-Hz or

11

23320-Hz), realism level (real, unrealistic, realistic), and configuration (upright, 234inverted).

235Real versus Fake Decoding Analysis

236We investigated whether real and fake image differences could be decoded from 237the EEG data using a leave-one-out cross validation approach. The leave-one-out 238cross-validation approach consists of dividing the data into training and testing 239sets whereby the classifiers are trained on all stimuli but one pair of real and fake 240stimuli from all but one RSVP sequence and then tested on the left-out stimulus 241pair from the remaining sequence. This ensured that the classifier had to 242generalise to the novel stimulus in order to successfully decode the category (i.e. 243real or fake) and could not rely on individual image-specific properties. Real 244stimuli were decoded against fake stimuli. Contrasts were broken down into 245presentation rate (5-Hz or 20-Hz), realism level (unrealistic, realistic), and 246configuration (upright, inverted). Thus, there were 8 decoded contrast 247combinations per image category. Given the large face processing literature and 248our clear hypotheses regarding faces, we were mainly interested in fake versus 249real decoding of faces; results from the car and bedroom categories are included 250for completeness on https://osf.io/n2z73/.

251To map the spatial distribution of the signal, we repeated the real versus fake 252decoding analysis at separate locations on the scalp. For each channel, we selected 253the four closest neighbouring channels and performed the exact same decoding 254analysis described above on just this local cluster of channels, storing the resulting 255accuracies at the centre channel. This results in a channel topography of decoding 256results that provides insight into the spatial origins of the signal.
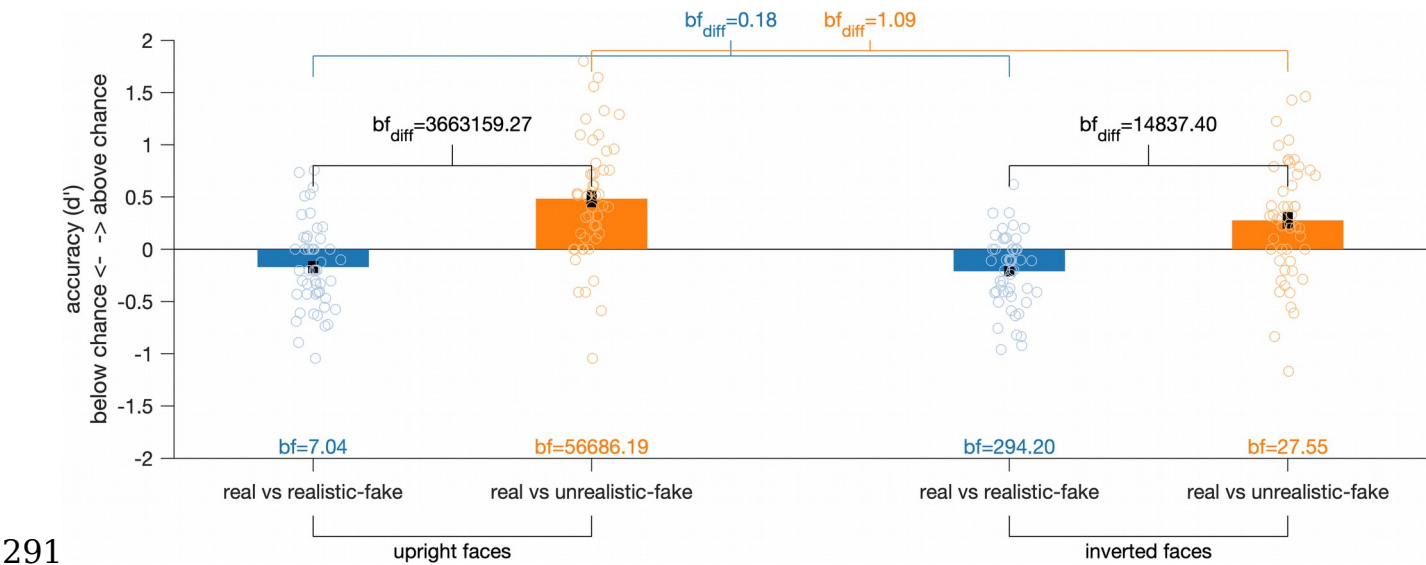
257As an exploratory follow-up analysis, we examined the relationship between real-258fake decoding accuracy and behavioural categorisation accuracy (Grootswagers et 259al., 2018; Ritchie et al., 2019). For each subject and each time point in the real-260fake decoding analysis, we correlated (Spearman's rho) the image-specific average

13

261classifier accuracies with their corresponding behavioural accuracies. We then
262performed group level inference on the resulting subject-wise time-varying brain-
263behaviour correlations. If successful real/fake decoding in EEG reflects the
264real/fake signal that is 'used' by the brain to guide behaviour (Grootswagers et al.,
2652018; Ritchie et al., 2019), then we would expect a positive correlation between
266image-specific EEG-classification accuracy and behavioural accuracy. That is, faces
267identified as real or fake by the classifier would also be identified as real or fake by
268the participants.

269Statistical inference

270For the decoding and behavioural analyses, we used Bayesian statistics to
271characterize evidence arising from the data as either supporting the presence
272(alternative hypothesis) or absence (null hypothesis) of an effect. (Dienes, 2011;
273Jeffreys, 1998; Rouder et al., 2009; Wagenmakers, 2007). We used a standard JZS
274prior to calculate the null and alternative hypotheses (Rouder et al., 2018), which
275is a Cauchy distribution with a scale factor of 0.707 to determine the evidence of
276above-chance performance (e.g., >50% decoding) and a null-hypothesis point prior
277at chance-level (Morey & Rouder, 2011). For ease of interpretation, we
278thresholded Bayes factor (BF) values > 10 for strong evidence for the alternative
279hypothesis and BF values < ⅓ as evidence in favour of the null hypothesis (Morey
280& Rouder, 2011). For the decoding analyses, BFs serve as continuous degrees of
281evidence across multiple time points and not specific hypothesis testing at single
282time points. Thus, isolated BFs at single time points which did not reach threshold
283were not treated as evidence for either hypothesis if the surrounding points did not
284reach threshold or were interspersed with below-threshold values. Rather, BFs

14

285 were treated as evidence if surrounding points were at threshold (Mai et al., 2019).
286 For the decoding analyses, we, in addition, computed corresponding frequentist
287 statistics using sign-permutation tests (1000 permutations) and Monte-Carlo
288 cluster statistics with TFCE as cluster-statistic (Smith & Nichols, 2009), corrected
289 for multiple comparisons across time using the max-statistic method (Maris &
290 Oostenveld, 2007).



291

292 **Figure 2. Behavioural discrimination of real and fake faces.** In an upright
293 (left) and inverted (right) configuration, discriminability for real/realistic (blue)
294 faces was below chance but above chance for real/unrealistic faces (orange).
295 Performance was similar regardless of whether faces were upright or inverted.
296 Bars show mean and standard error. Each circle represents the response of one
297 subject in one condition. The Bayes Factors (displayed above the x-axis) compute
298 the evidence for a difference from chance discriminability (50% accuracy), and
299 difference between conditions (stimulus and orientation).
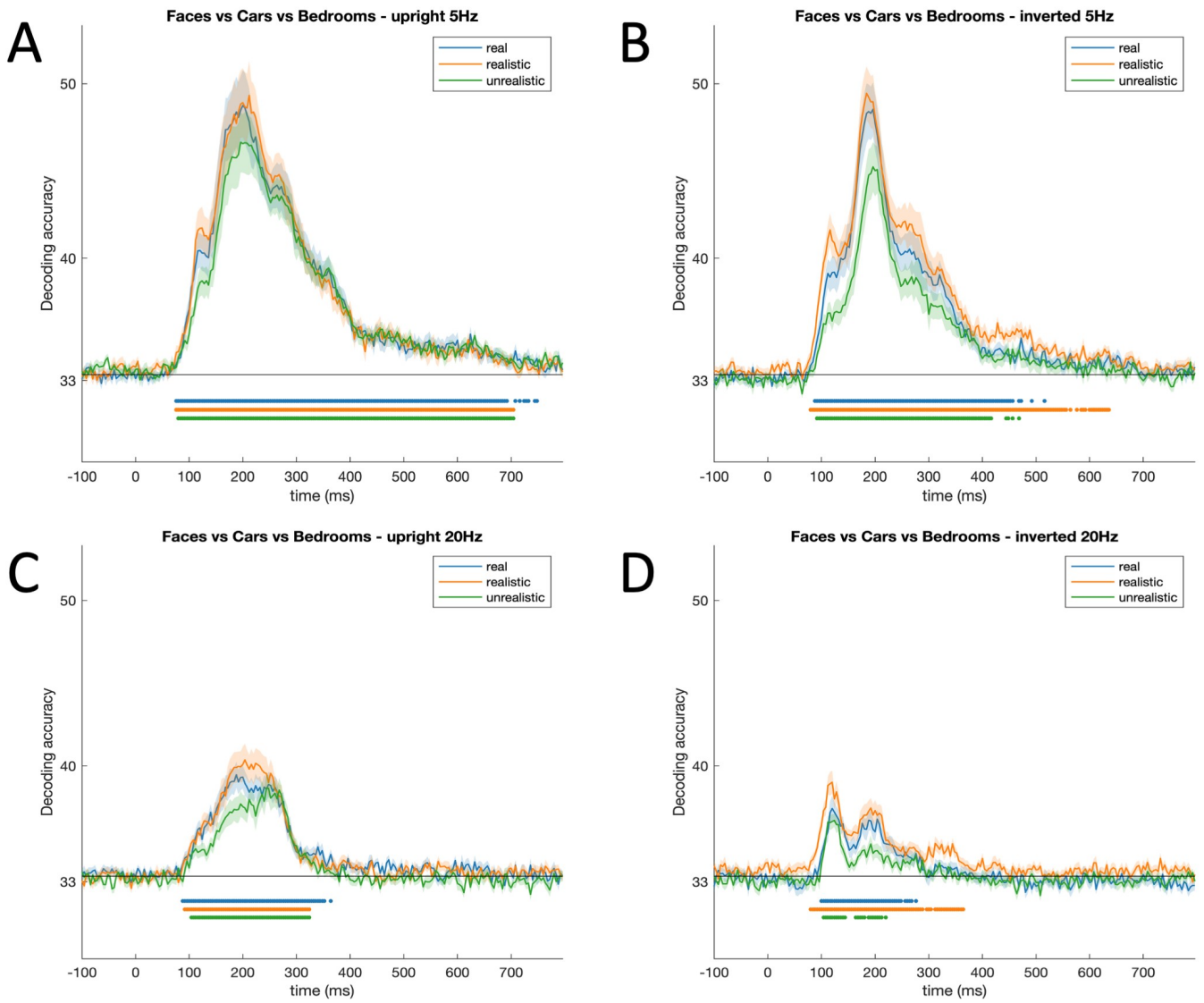
300 Results

301 Behavioural Performance

302 We were interested in whether participants could discriminate between real and
303 fake faces. We calculated the proportion of images that were judged correctly as
304 real or fake for each of the realistic/unrealistic and upright/inverted conditions and

15

305aggregated the judgements over participants. The main findings are presented in
306Figure 2. As indexed by d' discriminability analysis, we found that participants
307could reliably discriminate real from unrealistic fake faces (orange bars) but could
308not discriminate real from realistic fake faces (blue bars). Orientation had little
309effect on discriminability. Interestingly, performance in the real versus realistic
310face condition was below chance. Further inspection of the data revealed a general
311bias for participants to judge faces as real than as fake.  When discriminating
312between upright real and realistic fake faces, observers correctly classified 63%
313(se = 0.026, BF > 100) of real faces and 31% (se = 0.023, BF > 100) of realistic
314fake faces. For discriminating between upright real and unrealistic fake faces,
315observers correctly classified 68% (se = 0.026, BF > 100) of real faces but
316performed at chance (49%, se = 0.027, BF = 0.16) at classifying unrealistic fake
317faces. Classification performances were similar for inverted faces Overall,
318observers could identify real faces (although were more biased to do so) but had
319much more difficulty spotting the fakes.

320Overall, the behavioural results show that observers could not reliably differentiate
321real from realistic fake faces but performed better for real versus unrealistic fake
322faces. Interestingly, observers were more likely to judge artificial faces as being
323more real than fake consistent with Sanders et al. (2019). Inverting the faces had
324little effect on discriminability suggesting that detection was not reliant on
325configural or featural information (Tanaka et al., 2014).

16

327To examine whether real and fake images evoked similar categorical decoding 328effects compared to the previous literature, we decoded image category (cars, 329faces, and bedrooms) at all levels of realism (real, realistic, unrealistic), (Figure 3). 330As expected, we observed similar category-related dynamics for the real, realistic 331and unrealistic images across all conditions. At a 5Hz presentation rate, we 332observed above-chance decoding for all categories at real, realistic, and unrealistic 333(Figure 3A). Decoding emerged and remained above-chance from 100ms until 334700ms post-stimulus onset with an early peak at 120ms, a second peak at 200ms 335and a third peak at 250ms-300ms.

336We then tested how category decoding was affected by our control manipulations 337(inversion and presentation rate). We observed similar above-chance decoding for 338all categorical and realism levels upon inversion (Figure 3B) and at a 20Hz 339presentation rate (Figure 3C), albeit less pronounced with simultaneous stimulus 340inversion and 20Hz presentation (Figure 3D). When upright and inverted, faces, 341cars, and bedrooms could be decoded at all levels of realism with similar temporal 342dynamics reported elsewhere (Grootswagers, Robinson, & Carlson, 2019; 343Grootswagers, Wardle, & Carlson, 2017).

344

**Figure 3**. **Summary of category decoding using orientation and presentation rate manipulation.** A classifier was trained on EEG data from all categories, orientations, and presentation rates. Above-chance distinct category decoding was found for real (blue), realistic (orange), and unrealistic (green) stimuli regardless of orientation, presentation rate or stimuli type. Lines represent decoding accuracy over time with shaded areas displaying standard error across subjects (N = 22). Thresholded p-values below 0.05 are displayed under each pot.
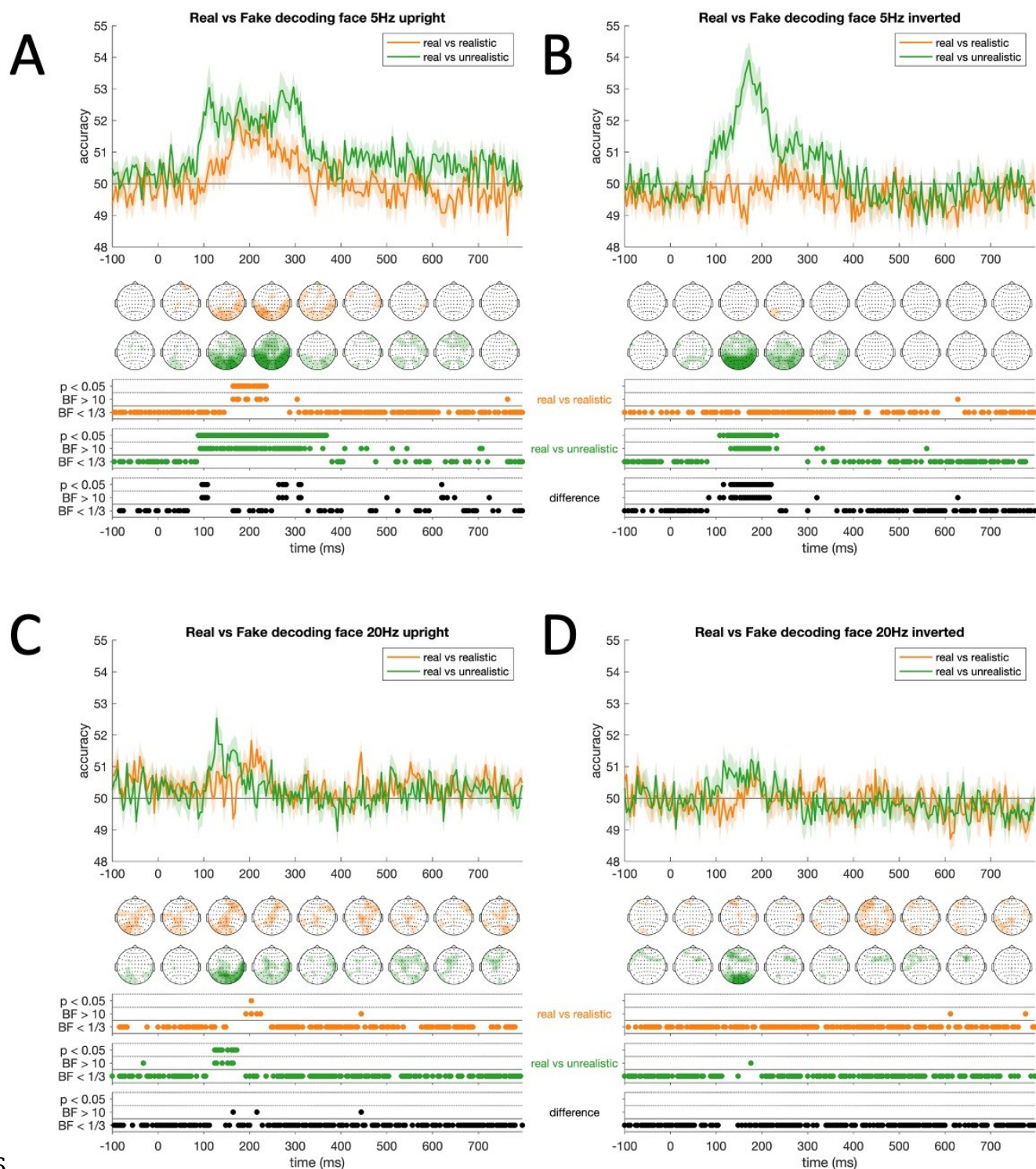
Decoding Realness from EEG: Real vs Fake Faces

To determine if the brain could distinguish real from fake, we then investigated differences in neural patterns evoked from real and fake faces. At 5Hz and upright

18

355(Figure 4A), above-chance decoding emerged and peaked for unrealistic faces at 356around 100ms, 200ms, and 300ms (BF > 10) and fell below-chance at 357approximately 370ms (BF < ⅓). This decodability is reflective of early, rapid, low-358level image perception followed by a later, higher-level, holistic decoding 359consistent with the temporal unfolding of face perception (Dobs et al., 2019; Balas 360& Koldewyn, 2013; Muhlberger et al., 2009). For realistic fake faces, decoding 361emerged at around 170ms and remained above-chance until approximately 240ms 362(BF > 10), suggesting a higher-level basis for discrimination of realistic and real 363faces. Although observers could not reliably tell apart real faces from realistic fake 364faces, the EEG data contains signal information relevant to this distinction which 365meaningfully differs between realistic fakes and unrealistic fakes, and this signal 366appears to be constrained to a relatively short stage of processing.

367If the information that we were decoding at 5Hz was reliant on image features 368rather than a face-processing effect, then we would predict that we could achieve a 369similar decoding result on inverted faces. However, at 5Hz and inverted (Figure 3704B), only unrealistic fake faces were decodable from real faces. Above-chance 371decoding emerged at around 100ms (BF > 10), peaked at around 170ms, and was 372at chance again at approximately 250ms (BF < ⅓). In contrast, realistic faces 373remained at-chance and were not decodable from real faces (BF < ⅓). This 374suggests that inversion, known to disrupt configural processing of faces, is 375similarly disrupting a face-specific mechanism accounting for decoding differences 376between realistic and unrealistic faces (Jacques, d'Arripe, & Rossion, 2007; 377Rossion et al., 2000).

19

378An alternative way to disrupt face-processing is to use faster presentation rates 379(Collins, Robinson, & Behrmann, 2018). At 20Hz and upright (Figure 4C), above-380chance decoding emerged for unrealistic faces at around 100ms and was sustained 381until approximately 170ms (BF >10). Decodability for realistic faces emerged at 382170ms and remained above chance until around 230ms (BF >10), showing very 383similar dynamics to the upright condition. Faster presentation rates have been 384shown to limit the extent and capacity for visual processing (Robinson, 385Grootswagers, & Carlson, 2019), but this result suggests short presentations can 386still yield information informative of real versus fake face distinctions, albeit with 387numerically lower and less sustained decoding accuracy.

388Lastly, at 20Hz and inverted (Figure 4D), decoding performance was at chance for 389realistic and unrealistic fake faces (BF < ⅓). This suggests that inversion plus a 390faster presentation rate is enough for the EEG data to no longer contain any 391relevant information pertaining to real versus fake face distinctions. In other 392words, configural processing has been disrupted to an extent that activity patterns 393evoked from fake faces were not differentiable from activity evoked from real 394faces. As expected, real versus fake bedroom and car decoding was not so evident 395and can be found on https://osf.io/n2z73/.

**Figure 4. Decoding real versus fake faces.** Different effects of orientation and presentation rate on decoding real and fake faces. Plots show decoding performance over time for real and fake (realistic or unrealistic) faces in upright and inverted orientations and at 5Hz and 20Hz presentation rates. The lines in each plot indicate classifier accuracy from time of stimulus onset until 800ms, with shaded areas showing standard errors across each subject (N = 22). Time-varying topographies are presented below each plot averaged across 100ms time bins

21

404where darker shades indicate contribution of channels to real/fake decoding. In the
405lowest panel, thresholded p-values and Bayes Factors indicate above-chance
406decoding or non-zero differences.

407Finally, we examined the relationship between real-fake decoding accuracy and
408behavioural categorisation accuracy. If successful real/fake decoding in EEG
409reflects the real/fake signal that is 'used' by the brain to guide behaviour
410(Grootswagers et al., 2018; Ritchie et al., 2019) then we would predict to observe a
411positive correlation between image-specific EEG-classification accuracy and
412behavioural accuracy. Figure 5 shows the time-varying correlations for the upright
413and inverted 5Hz conditions. We did not perform this analysis for the 20Hz
414conditions due to limited above-chance decoding. We observed evidence for a
415positive brain-behaviour correlation around 170ms for the upright and inverted
416unrealistic faces, which is consistent with time points of above-chance decoding
417(Figure 4A). This result suggests that, at least for the unrealistic faces, the signal
418that is used by the classifier for real/fake distinction could be used by the brain to
419make the real/fake decision (Grootswagers, Cichy, & Carlson., 2018; Ritchie,
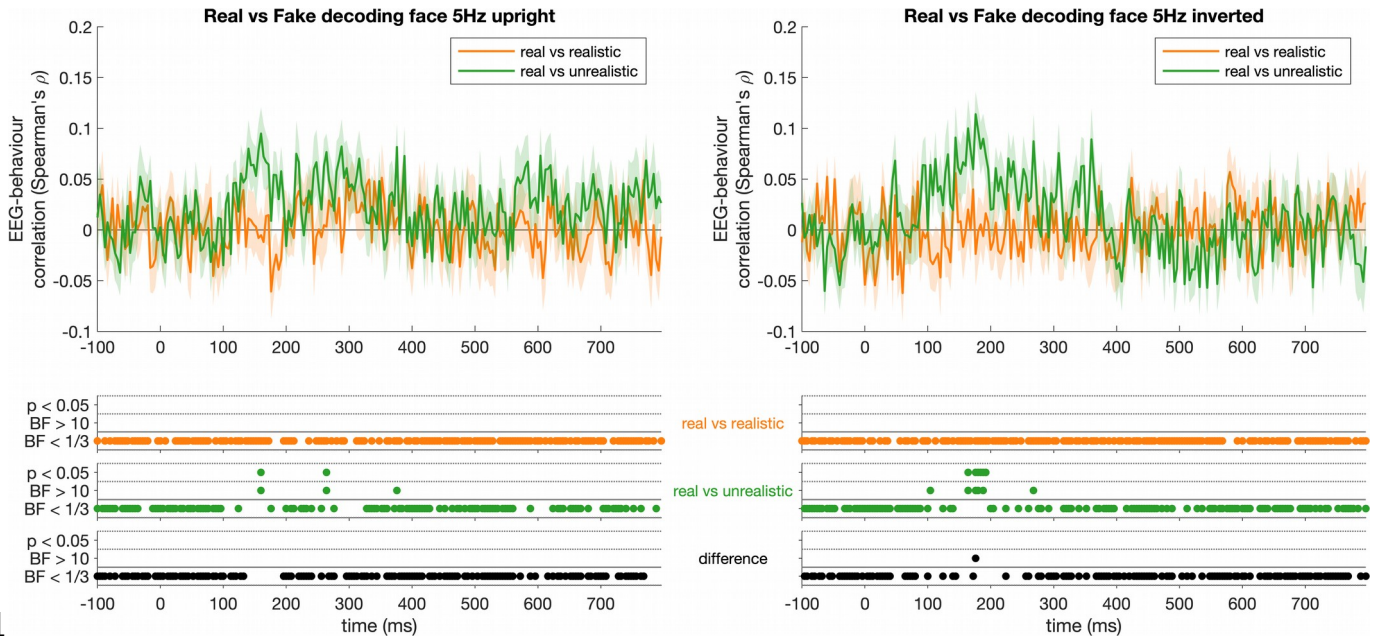420Kaplan, & Klein, 2019).

22

421

**Figure 5. Correlating behavioural accuracy with decoding**. Plots show the relationship between image-specific EEG decoding accuracy and behavioural accuracy over time for the 5Hz upright condition (left) and 5Hz inverted condition (right). The lines indicate correlation from time of stimulus onset until 800ms for realistic versus real faces (orange) and unrealistic versus real faces (green), with shaded areas showing standard errors. In the lowest panel, thresholded p-values and Bayes Factors indicate above-chance correlation or non-zero differences. Positive brain-behaviour correlations can be seen at around 170ms and 270ms for upright unrealistic faces (green) and at around 150ms-200ms for inverted unrealistic faces (BF>10).

## Discussion

There is growing concern that hyperrealism is advancing at such a rate that humans will have difficulty discerning between what is real and what is fake (Fletcher, 2018; Khodabakhsh et al., 2019; Nightingale et al., 2017; Shen et al., 2019). Our results justify these concerns by revealing that observers cannot consciously and reliably identify realistic fake faces amongst real faces. However, using time-resolved EEG and multivariate pattern classification methods, we found that it was possible to decode both unrealistic *and* realistic fake faces from real faces using brain activity. This dissociation between behaviour and neural

23

441responses for realistic faces yields important new evidence about fake face
442perception as well as implications involving the increasingly realistic class of GAN-
443generated faces. Namely, the brain encodes information relevant to artificial face
444appearance even though humans do not consciously perceive any differences
445between GAN-generated faces and real faces.

446Our behavioural results are consistent with previous research that suggests that
447observers typically display difficulties with correctly discriminating between real
448and realistic fake faces despite face expertise (Holmes et al., 2016; Nightingale et
449al., 2017; Sanders et al., 2019; Zhou et al., 2019). For example, in a two-alternative
450forced-choice task, participants would judge realistic artificial faces as being more
451realistic than human faces on a third of all trials (Sanders et al., 2019). Artificial
452faces made by GANs have also recently received attention and have been similarly
453demonstrated to fool observers (Hulzebosch et al., 2020; Isola et al., 2017; Zhou et
454al., 2019; Liu et al., 2020). As expected, we found that it was much harder to
455discriminate fake from real faces in our realistic condition relative to the
456unrealistic condition, confirming that the newer generation of GAN images are
457much more naturalistic.  We presented faces for 200ms, which could be considered
458a brief exposure period, but the images were not masked so processing would have
459continued even after the images had disappeared (Robinson, Grootswagers, &
460Carlson, 2019). Given a long enough time to observe, Liu et al., (2020) found that
461identifying artifacts such as "asymmetrical eyes" and "irregular teeth" in artificial
462faces can assist in spotting fakes. Presumably, assessing such details requires
463more time and eye movements. Indeed, observers can be trained to reliably spot
464fake faces by learning what to look for (Hills & Lewis, 2006; Tanaka & Farah,
4651993). Here, our primary focus was examining first impression responses by

24

466limiting the time spent looking at each face and giving participants unlimited time
467to make a response. Future studies may investigate whether training observers on
468GAN-generated faces enhances detection.

469We found that although observers may be fooled behaviourally by artificial faces,
470they have distinct representations in the human visual system. Given that category
471decoding was most pronounced and sustained in the 5Hz and upright condition,
472enough for each image to reach a high-level representation in the brain
473(Grootswagers, Robinson, & Carlson, 2019), we expected real/fake decoding to be
474most pronounced in this condition too. Above-chance decoding represents the
475classifier successfully distinguishing neural activity evoked from real and fake
476faces, namely, real/fake differences. Critically, a leave-one-out cross validation
477approach (see methods) ensured that the classifier could not learn to categorise
478the EEG data based on visual features or low-level properties belonging to specific
479faces, but rather had to generalize learned category information (real/fake) onto
480novel stimuli (Carlson et al., 2013; Grootswagers, Wardle, & Carlson, 2016;
481Teichmann et al., 2020). This guaranteed that the classifier performance related to
482a group-level distinction rather than to individual image-level properties.

483Indeed, for the 5Hz, upright condition, we found that the classifier successfully
484discriminated between unrealistic/real as well as realistic/real faces (Figure 4A).
485Decoding for unrealistic faces displayed a triple peak pattern, emerging at around
486100ms maintained until around 370ms. Early decoding differences are consistent
487with rapid face detection and face-specific processing (Rossion et al., 2015; Dobs
488et al., 2019; Crouzet, Kirchner, & Thorpe, 2010; Wardle et al., 2020). The latter
489two peaks (at around 170-200ms and 270-320ms) have been similarly

25

490demonstrated to emerge in real versus artificial face perception (Wheatley et al.,
4912011; Balas & Koldewyn, 2013; Sagiv & Bentin, 2001; Schindler et al., 2017,
492Schindler et al., 2019, Wardle et al., 2020). Schindler et al (2017) suggest that
493early-stage N170 processing is related to assessing the structural configuration of
494faces as seen by a greater occipital involvement whilst the later-staged LPP, seen
495to increase linearly with face realism, suggests a deeper person-related, semantic
496involvement (also see Abdel Rahman, 2011, Taylor, Shehzad, & McCarthy, 2016).
497Differences at the triple peak correspond to N250 and P300 components typically
498associated with face familiarity (Collins et al., 2018) and semantic information
499(Tanaka et al., 2006), the latter especially important for behaviour (Hanso et al.,
5002010). In contrast, realistic/real decoding displayed a single-peak emergence
501between around 170ms to 240ms indicating a difference in processing between
502realistic and unrealistic faces. Namely, that differences in perception between real
503and realistic faces were constrained to the 170ms time period. Indeed, in
504comparing human faces to doll faces and artificial faces, others have shown that
505only the human faces typically evoke sustained neural responses beyond the N170
506component necessary for higher-order perception (Balas & Koldewyn, 2013;
507Wheatley et a., 2011). Balas and Koldewyn (2013) found that the N170 was better
508characterised by encoding deviations from facial appearance than it was for
509animacy perception. In other words, realistic faces were perceived as configurally
510different to real faces, but that only unrealistic faces engaged later processing
511necessary for high-order animacy or familiarity perception. Overall, earlier
512decoding for unrealistic faces, consistent with apparent low-level image
513differences (Figure 1B), suggests that early and low to mid-level processing
514differences may account for decodability between real and unrealistic faces. The

26

515decoding for realistic faces, by contrast, emerges later and is constrained to the 516170ms time period, suggesting a face-specific configural process may be 517responsible for this distinction.

518Assessing fake/real decoding for inverted faces allows us to evaluate whether the 519fake/real distinction relies on mechanisms that are responsible for the superiority 520in face recognition for upright faces relative to inverted faces. Inversion disrupts 521the configural processing of faces by making them appear more like objects whilst 522retaining low-level stimulus attributes (Eimer, 2000; Leder & Bruce, 2000; 523Rousselet et al., 2003). Firstly, we found that inversion led to the disruption of 524decoding for realistic faces (Figure 4B). In contrast, we found that decoding for 525unrealistic inverted faces was preserved but less sustained when compared to 526upright. The peak in decoding may be reflective of increased featural processing 527for inverted unrealistic faces, also seen to occur with distorted or 'Thatcherized' 528faces (Carbon et al., 2005; Milivojevic et al., 2003). Lack of above-chance decoding 529for inverted realistic faces may reflect the contribution of high-level, expertise-530driven capabilities for upright fake face detection when face processing 531mechanisms, rather than object processing, were available. Overall, we found that 532upon stimulus inversion our decoding results were consistent with a face-specific 533or expertise response, such that realistic fake faces could not be discriminated 534from real faces when typical face perception was disrupted, even though the same 535visual features were present.

536The presentation of images at a faster presentation rate limits the consolidation of 537each image and build-up of higher-order representation (Grootswagers, Robinson, 538& Carlson, 2019)., allowing an analysis of the contribution of low-level processing.

27

539At a faster presentation rate of 20Hz, we found that upright fake faces could be
540discriminated from real faces for the realistic and unrealistic conditions (Figure
5414C). Indeed, early, low-level visual processing is fairly unaffected by image
542presentation durations (Grootswagers, Robinson, & Carlson, 2019). Observing less
543sustained decoding is consistent with the limited capacity and extent of visual
544processing since each image is masked by every successive image to a greater
545extent and therefore places limits on visual processing compared to a slower
546presentation rate (Collins, Robinson, & Behrmann, 2018; Robinson, Grootswagers,
547& Carlson, 2019). Additionally, higher-level, identity or semantically related face
548information discernible in the slow condition was possibly limited at the faster
549presentation rate consistent with Collins et al. (2018). In sum, we found that
550unrealistic faces could be decoded upon inversion and at a faster presentation rate
551suggesting the contribution of low-level visual differences. By contrast, we could
552not decode realistic faces when inverted, but we could decode at a faster
553presentation rate, indicating that fake/real perception was likely driven by
554expertise and face-specific processing.

555Interestingly, we found that neural differences between real and realistic fake
556faces did not translate into a reliable behavioural decision for realistic face
557discrimination at the population level. We found a brain-behaviour correlation at
558around 150ms-200ms for unrealistic versus real faces, suggesting that this time
559period of processing is important for behaviour. However, the same correlation
560was not observed for the realistic faces. One possibility is that whilst our data
561indicates that a realistic fake/real signal is present, this signal gets 'lost' in the
562visual hierarchy and consequently remains uninformative for behaviour. For
563instance, although animacy categorisation can be decoded throughout the entire

28

564ventral visual stream, this information is most suitably formatted for behaviour in 565higher-level visual areas like the ventral occipital and parahippocampal cortex 566(Grootswagers, Cichy, & Carlson, 2018). Since decoding unrealistic/real faces was 567more sustained than realistic/real faces, associated more with in-depth face 568processing at later stages (i.e., LPP), it is possible that this level of extended 569processing is required for behavioural "readout" (see de-Wit et al., 2016; 570Grootswagers, Cichy, & Carlson, 2018; Ritchie, Kaplan, & Klein 2019). Yet, the 571highest brain-behaviour correlation for unrealistic faces was observed at 150-572200ms, a time when decoding was not reliably different between the realistic and 573unrealistic condition. This has a number of implications. In an applied setting such 574as cyber security or Deepfakes, examining the detection ability for hyper-realistic 575fake faces might be best pursued using machine learning classifiers applied to 576neuroimaging data rather than targeting behavioural performance. As we have 577shown, the former contains discriminative relevance whereas observers may 578actually perform worse than chance given the decision (and a brief glance). A third 579related possibility is that the decodable real/fake face signal is operating below 580conscious access and therefore is not picked up by our behavioural task. This is 581reminiscent of findings that individuals with prosopagnosia who cannot 582behaviourally classify or recognise faces as familiar or unfamiliar nevertheless 583display stronger autonomic responses to familiar faces than unfamiliar faces 584(Tranel & Damasio, 1985). Similarly, what we have shown in this study is that 585participants could not reliably discriminate between real and realistic fake faces 586even though we could accurately decode this difference from their neural activity. 587Still, it is possible that a different behavioural task may have yielded a better 588performance. Forced to respond via a two-alternative forced-choice task or an

29

589implicit task such as face familiarity or trustworthiness may have engaged
590different behavioural processes more conducive for real/fake face discrimination.
591For instance, behaviourally categorising faces as threatening, competent, or
592trustworthy has been shown to occur as quickly as 33- 100ms after onset (Bar et
593al., 2006; Willis & Todorov, 2006). Conversely, real or fake judgments may occur
594as late as 240ms after stimulus presentation (Zhou et al., 2019). Therefore, future
595work could investigate whether judgments about face trustworthiness or threat
596may be a better cue for detection than real or fake.

597In sum, we found that there is a dissociation between the ability of participants to
598categorise faces as real or fake and the decodability of this distinction in the brain.
599In other words, although the brain can 'recognise' the difference between real and
600realistic fake faces, observers cannot consciously tell them apart. Our findings of
601the dissociation between brain response and behaviour has implications for the
602ways in which we study fake face perception, the questions we pose when asking
603about fake image identification, and the possible ways in which we can establish
604protective standards against fake image misuse.

605Future studies may investigate the contribution of face expertise for decoding and
606behaviour.  Expertise influences how deeply and configurally a face is perceived
607allowing for more subtle identification of spatial relations, features, and same-race
608faces (Wong et al., 2009; Tanaka, 2001; Tanaka & Taylor, 1991; Hancock &
609Rhodes, 2008; Meissner & Brigham, 2001). Indeed, individuals with digital
610manipulation training and experience (i.e., photo-editing and photography) are
611more able to identify fake images than non-experienced individuals (Shen et al.,
6122019). Having the same participants participate in both the EEG and behaviour

30

613experiments may be useful in exploring inter-individual differences and the 614influence of expertise.

615In conclusion, we investigated to what extent state-of-the-art GAN faces made by 616AI fool human observers. Using behavioural and neuroimaging methods we found 617that it was possible to reliably detect AI-generated fake images using EEG activity 618given only a brief glance, even though observers could not consciously report 619seeing differences. Given that observers are already struggling with differentiating 620between fake and real faces, it is of immediate and practical concern to further 621investigate the important ways in which the brain is able to tell the two apart. It is 622becoming increasingly possible to rapidly and effortlessly generate hyper-realistic 623fake images, videos, writing, and multimedia that are practically indiscernible from 624real (Radford et al., 2019; Maras & Alexandrou, 2018; Asensio et al., 2014; Ledig 625et al., 2017). This capacity is only going to become more widespread and has 626profound implications for cybersecurity, fake news, detection bypass, and social 627media (Damiani, 2019; Fletcher, 2018; Maddocks, 2020). Already, a newer and 628more realistic set of images and faces have been generated by GANs that might 629challenge human perception more drastically than we have investigated here 630(Karras et al., 2020). Understanding the dissociation between brain and behaviour 631for fake face detection will have practical implications for the way we tackle the 632potentially detrimental and universal spread of artificially generated information.

31

640 References

641 Abdel Rahman, R. (2011). Facing good and evil: Early brain signatures of affective
642      biographical knowledge in face recognition. *Emotion, 11*(6), 1397.

643 Adelani, D. I., Mai, H., Fang, F., Nguyen, H. H., Yamagishi, J., & Echizen, I. (2020,
644      April). Generating sentiment-preserving fake online reviews using neural
645      language models and their human-and machine-based detection. In
646      *International Conference on Advanced Information Networking and*
647      *Applications* (pp. 1341-1354). Springer, Cham.

648 Asensio, J. M. L., Peralta, J., Arrabales, R., Bedia, M. G., Cortez, P., & Peña, A. L.
649      (2014). Artificial intelligence approaches for the generation and assessment
650      of believable human-like behaviour in virtual characters. *Expert Systems*
651      *with Applications, 41*(16), 7281-7290.

652 Axelrod, V., & Yovel, G. (2015). Successful decoding of famous faces in the
653      fusiform face area. *PloS one, 10*(2), e0117126.

654 Balas, B., & Koldewyn, K. (2013). Early visual ERP sensitivity to the species and
655      animacy of faces. *Neuropsychologia, 51*(13), 2876–2881.
656      https://doi.org/10.1016/j.neuropsychologia.2013.09.014

657 Bar, M., Neta, M., & Linz, H. (2006). Very first impressions. *Emotion, 6*(2), 269.

658 Bentin, S., Allison, T., Puce, A., Perez, E., & McCarthy, G. (1996).
659      Electrophysiological studies of face perception in humans. *Journal of*
660      *cognitive neuroscience, 8*(6), 551-565.

661 Carbon, C. C., Schweinberger, S. R., Kaufmann, J. M., & Leder, H. (2005). The
662      Thatcher illusion seen by the brain: an event-related brain potentials study.
663      *Cognitive Brain Research, 24*(3), 544-555.

664 Carlson, T., Tovar, D. A., Alink, A., & Kriegeskorte, N. (2013). Representational
665      dynamics of object vision: the first 1000 ms. *Journal of vision, 13*(10), 1-1.

Carlson, T., Goddard, E., Kaplan, D. M., Klein, C., & Ritchie, J. B. (2018). Ghosts in machine learning for cognitive neuroscience: Moving from data to theory. *NeuroImage, 180*, 88–100. https://doi.org/10.1016/j.neuroimage.2017.08.019

Collins, E., Robinson, A. K., & Behrmann, M. (2018). Distinct neural processes for the perception of familiar versus unfamiliar faces along the visual hierarchy revealed by EEG. *NeuroImage, 181*, 120-131.

Crookes, K., Ewing, L., Gildenhuys, J. D., Kloth, N., Hayward, W. G., Oxner, M., ... & Rhodes, G. (2015). How well do computer-generated faces tap face expertise?. *PloS one, 10*(11), e0141353.

Crouzet, S. M., Kirchner, H., & Thorpe, S. J. (2010). Fast saccades toward faces: face detection in just 100 ms. *Journal of vision, 10*(4), 16-16.

Damiani, J. (2019). A voice deepfake was used to scam a CEO out of $243,000.

De Leeuw, J. R. (2015). jsPsych: A JavaScript library for creating behavioral experiments in a Web browser. *Behavior research methods, 47(*1), 1-12.

de-Wit, L., Alexander, D., Ekroll, V., & Wagemans, J. (2016). Is neuroimaging measuring information in the brain? *Psychonomic Bulletin & Review, 23(*5), 1415–1428.

Dienes, Z. (2011). Bayesian versus orthodox statistics: Which side are you on?. *Perspectives on Psychological Science, 6*(3), 274-290.

Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-trial EEG dynamics including independent component analysis. *Journal of neuroscience methods, 134*(1), 9-21.

Dobs, K., Isik, L., Pantazis, D., & Kanwisher, N. (2019). How face perception unfolds over time. *Nature communications, 10*(1), 1-10.

Eimer, M. (2000). Effects of face inversion on the structural encoding and recognition of faces: Evidence from event-related brain potentials. *Cognitive Brain Research, 10*(1-2), 145-158.

Farid, H., & Bravo, M. (2007). Photorealistic rendering: How realistic is it?. *Journal of Vision, 7*(9), 766-766.

Farid, H., & Bravo, M. J. (2012). Perceptual discrimination of computer generated and photographic faces. *Digital Investigation, 8*(3-4), 226-235.

Farokhmanesh, M. (2018). Deepfakes Are Disappearing from Parts of the Web, But They're Not Going Away. *The Verge*.

Fletcher, J. (2018). Deepfakes, Artificial Intelligence, and Some Kind of Dystopia: The New Faces of Online Post-Fact Performance. *Theatre Journal, 70*(4), 455-471.

33

702 Gauthier, I., & Tarr, M. J. (2002). Unraveling mechanisms for expert object
703     recognition: bridging brain activity and behavior. *Journal of Experimental*
704     *Psychology: Human Perception and Performance, 28(*2), 431.

705 Gleicher, N. (2019). Removing Coordinated Inauthentic Behavior From Georgia,
706     Vietnam and the US. Facebook. *Retrieved from*
707     http://about.fb.com/news/2019/12/removing-coordinated-inauthentic-
708     behavior-from-georgia-vietnam-and-the-us/

709 Grootswagers, T. (2020). A primer on running human behavioural experiments
710     online. *Behavior Research Methods, 52*, 2283–2286.
711     https://doi.org/10.3758/s13428-020-01395-3

712 Grootswagers, T., Cichy, R. M., & Carlson, T. A. (2018). Finding decodable
713     information that can be read out in behaviour. *NeuroImage, 179*, 252–262.
714     https://doi.org/10.1016/j.neuroimage.2018.06.022

715 Grootswagers, T., Robinson, A. K., & Carlson, T. A. (2019). The representational
716     dynamics of visual objects in rapid serial visual processing streams.
717     *NeuroImage, 188*, 668–679.
718     https://doi.org/10.1016/j.neuroimage.2018.12.046

719 Grootswagers, T., Wardle, S. G., & Carlson, T. A. (2016). Decoding Dynamic Brain
720     Patterns from Evoked Responses: A Tutorial on Multivariate Pattern Analysis
721     Applied to Time Series Neuroimaging Data. *Journal of Cognitive*
722     *Neuroscience, 29*(4), 677–697. https://doi.org/10.1162/jocn_a_01068

723 Grootswagers, T., Wardle, S. G., & Carlson, T. A. (2017). Decoding dynamic brain
724     patterns from evoked responses: A tutorial on multivariate pattern analysis
725     applied to time series neuroimaging data. *Journal of cognitive neuroscience,*
726     *29(*4), 677-697.

727 Hancock, K. J., & Rhodes, G. (2008). Contact, configural coding and the other-race
728     effect in face recognition. *British Journal of Psychology, 99*(1), 45-56.

729 Hanso, L., Bachmann, T., & Murd, C. (2010). Tolerance of the ERP signatures of
730     unfamiliar versus familiar face perception to spatial quantization of facial
731     images. *Psychology, 1*(03), 199.

732 Haynes, J. D., & Rees, G. (2006). Decoding mental states from brain activity in
733     humans. Nature Reviews *Neuroscience, 7*(7), 523-534.

734 Hills, P. J., & Lewis, M. B. (2006). Short article: reducing the own-race bias in face
735     recognition by shifting attention. Quarterly Journal of Experimental
736     Psychology, 59(6), 996-1002.

737 Holmes, O., Banks, M. S., & Farid, H. (2016). Assessing and improving the
738     identification of computer-generated portraits. *ACM Transactions on Applied*
739     *Perception (TAP), 13(*2), 1-12.

Hulzebosch, N., Ibrahimi, S., & Worring, M. (2020). Detecting CNN-Generated Facial Images in Real-World Scenarios. 642–643. https://openaccess.thecvf.com/content_CVPRW_2020/html/w39/Hulzebosch_Detecting_CNN-Generated_Facial_Images_in_Real-World_Scenarios_CVPRW_2020_paper.html

Isola, P., Zhu, J.-Y., Zhou, T., & Efros, A. A. (2017). Image-To-Image Translation With Conditional Adversarial Networks. 1125–1134. https://openaccess.thecvf.com/content_cvpr_2017/html/Isola_Image-To-Image_Translation_With_CVPR_2017_paper.html

Jacques, C., d'Arripe, O., & Rossion, B. (2007). The time course of the inversion effect during individual face discrimination. *Journal of Vision, 7*(8), 3-3.

Jeffreys, H. (1998). The theory of probability. *OUP Oxford*.

Karras, T., Laine, S., & Aila, T. (2019). A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition (*pp. 4401-4410).

Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., & Aila, T. (2020). Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 8110-8119).

Kietzmann, J., Lee, L. W., McCarthy, I. P., & Kietzmann, T. C. (2020). Deepfakes: Trick or treat?. *Business Horizons, 63*(2), 135-146.

Khodabakhsh, A., Ramachandra, R., & Busch, C. (2019). Subjective Evaluation of Media Consumer Vulnerability to Fake Audiovisual Content. In *2019 Eleventh International Conference on Quality of Multimedia Experience (QoMEX)* (pp. 1-6). IEEE.

Leder, H., & Bruce, V. (2000). When inverted faces are recognized: The role of configural information in face recognition. *The quarterly journal of experimental psychology Section A, 53*(2), 513-536.

Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Ledig, C., Theis, L., Huszár, F., Caballero, J., Cunningham, A., Acosta, A., Aitken, A., Tejani, A., Totz, J., Wang, Z. & Shi, W. (2017). Photo-realistic single image super-resolution using a generative adversarial network. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 4681-4690).

Liu, Z., Qi, X., & Torr, P. H. (2020). Global texture enhancement for fake face detection in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 8060-8069).

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.

35

Maddocks, S. (2020). 'A Deepfake Porn Plot Intended to Silence Me': exploring continuities between pornographic and 'political'deep fakes. Porn Studies, 1-9.

Mader, B., Banks, M. S., & Farid, H. (2017). Identifying computer-generated portraits: The importance of training and incentives. *Perception, 46*(9), 1062-1076.

Mai, A. T., Grootswagers, T., & Carlson, T. A. (2019). In search of consciousness: Examining the temporal dynamics of conscious visual perception using MEG time-series data. *Neuropsychologia, 129*, 310-317.

Maras, M. H., & Alexandrou, A. (2019). Determining authenticity of video evidence in the age of artificial intelligence and in the wake of Deepfake videos. *The International Journal of Evidence & Proof, 23*(3), 255-262.

Maris, E., & Oostenveld, R. (2007). Nonparametric statistical testing of EEG-and MEG-data. *Journal of neuroscience methods, 164*(1), 177-190.

Marra, F., Gragnaniello, D., Cozzolino, D., & Verdoliva, L. (2018, April). Detection of gan-generated fake images over social networks. In *2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)* (pp. 384-389). IEEE.

Milivojevic, B., Clapp, W. C., Johnson, B. W., & Corballis, M. C. (2003). Turn that frown upside down: ERP effects of thatcherization of misorientated faces. *Psychophysiology, 40(*6), 967-978.

Meissner, C. A., & Brigham, J. C. (2001). Thirty years of investigating the own-race bias in memory for faces: A meta-analytic review. *Psychology, Public Policy, and Law, 7(*1), 3.

Morey, R. D., & Rouder, J. N. (2011). Bayes factor approaches for testing interval null hypotheses. *Psychological methods, 16*(4), 406.

Mühlberger, A., Wieser, M. J., Herrmann, M. J., Weyers, P., Tröger, C., & Pauli, P. (2009). Early cortical processing of natural and artificial emotional faces differs between lower and higher socially anxious persons. *Journal of neural transmission, 116*(6), 735-746.

Natsume, R., Saito, S., Huang, Z., Chen, W., Ma, C., Li, H., & Morishima, S. (2019). Siclope: Silhouette-based clothed people. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition* (pp. 4480-4490).

Nightingale, S. J., Wade, K. A., & Watson, D. G. (2017). Can people identify original and manipulated photos of real-world scenes?. *Cognitive research: principles and implications, 2*(1), 30.

Olivetti, E., Veeramachaneni, S., & Nowakowska, E. (2012). Bayesian hypothesis testing for pattern discrimination in brain decoding. *Pattern Recognition, 45*(6), 2075-2084.

Oostenveld, R., & Praamstra, P. (2001). The five percent electrode system for high-resolution EEG and ERP measurements. *Clinical neurophysiology, 112*(4), 713-719.

Oosterhof, N. N., Connolly, A. C., & Haxby, J. V. (2016). CoSMoMVPA: Multi-Modal Multivariate Pattern Analysis of Neuroimaging Data in Matlab/GNU Octave. *Frontiers in Neuroinformatics, 10*. https://doi.org/10.3389/fninf.2016.00027

Peirce, J., Gray, J. R., Simpson, S., MacAskill, M., Höchenberger, R., Sogo, H., ... & Lindeløv, J. K. (2019). PsychoPy2: Experiments in behavior made easy. *Behavior research methods, 51(*1), 195-203.

Pereira, F., Mitchell, T., & Botvinick, M. (2009). Machine learning classifiers and fMRI: a tutorial overview. *Neuroimage, 45*(1), S199-S209.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). Language models are unsupervised multitask learners. *OpenAI blog, 1*(8), 9.

Ritchie, J. B., Kaplan, D. M., & Klein, C. (2019). Decoding the Brain: Neural Representation and the Limits of Multivariate Pattern Analysis in Cognitive Neuroscience. *British Journal For The Philosophy Of Science, 70*(2), 581–607. https://doi.org/10.1093/bjps/axx023

Robinson, A. K., Grootswagers, T., & Carlson, T. A. (2019). The influence of image masking on object representations during rapid serial visual presentation. *NeuroImage,                         197*,                         224–231. https://doi.org/10.1016/j.neuroimage.2019.04.050

Rossion, B., Gauthier, I., Tarr, M. J., Despland, P., Bruyer, R., Linotte, S., & Crommelinck, M. (2000). The N170 occipito-temporal component is delayed and enhanced to inverted faces but not to inverted objects: an electrophysiological account of face-specific processes in the human brain. *Neuroreport, 11*(1), 69-72.

Rossion, B., Torfs, K., Jacques, C., & Liu-Shuang, J. (2015). Fast periodic presentation of natural images reveals a robust face-selective electrophysiological response in the human brain. *Journal of vision, 15(*1), 18-18.

Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic bulletin & review, 16(*2), 225-237.

Rousselet, G. A., Macé, M. J. M., & Fabre-Thorpe, M. (2003). Is it an animal? Is it a human face? Fast processing in upright and inverted natural scenes. *Journal of vision, 3*(6), 5-5.

Sagiv, N., & Bentin, S. (2001). Structural encoding of human and schematic faces: holistic and part-based processes. *Journal of cognitive neuroscience, 13*(7), 937-951.

37

Sanders, J. G., Ueda, Y., Yoshikawa, S., & Jenkins, R. (2019). More human than human: a Turing test for photographed faces. *Cognitive research: principles and implications, 4*(1), 1-10.

Schindler, S., Bruchmann, M., Bublatzky, F., & Straube, T. (2019). Modulation of face-and emotion-selective ERPs by the three most common types of face image manipulations. *Social cognitive and affective neuroscience, 14*(5), 493-503.

Schindler, S., Zell, E., Botsch, M., & Kissler, J. (2017). Differential effects of face-realism and emotion on event-related brain potentials and their implications for the uncanny valley theory. *Scientific Reports, 7*(1), 45003. https://doi.org/10.1038/srep45003

Shakeshaft, N. G., & Plomin, R. (2015). Genetic specificity of face recognition. Proceedings of the National *Academy of Sciences, 112*(41), 12887-12892.

Shen, C., Kasra, M., Pan, W., Bassett, G. A., Malloch, Y., & O'Brien, J. F. (2019). Fake images: The effects of source, intermediary, and digital media literacy on contextual assessment of image credibility online. *New media & society, 21*(2), 438-463.

Shu, K., Sliva, A., Wang, S., Tang, J., & Liu, H. (2017). Fake news detection on social media: A data mining perspective. *ACM SIGKDD explorations newsletter, 19*(1), 22-36.

Sinha, P., Balas, B., Ostrovsky, Y., & Russell, R. (2006). Face recognition by humans: Nineteen results all computer vision researchers should know about. *Proceedings of the IEEE, 94*(11), 1948-1962.

Smith, S. M., & Nichols, T. E. (2009). Threshold-free cluster enhancement: addressing problems of smoothing, threshold dependence and localisation in cluster inference. *Neuroimage, 44*(1), 83-98.

Sunday, M. A., Dodd, M. D., Tomarken, A. J., & Gauthier, I. (2019). How faces (and cars) may become special. *Vision research, 157*, 202-212.

Tanaka, J. W. (2001). The entry point of face recognition: evidence for face expertise. *Journal of Experimental Psychology: General, 130*(3), 534.

Tanaka, J. W., Curran, T., Porterfield, A. L., & Collins, D. (2006). Activation of preexisting and acquired face representations: the N250 event-related potential as an index of face familiarity. *Journal of Cognitive Neuroscience, 18*(9), 1488-1497.

Tanaka, J. W., & Farah, M. J. (1993). Parts and wholes in face recognition. *The Quarterly journal of experimental psychology, 46*(2), 225-245.

894 Tanaka, J. W., & Taylor, M. (1991). Object categories and expertise: Is the basic
895       level in the eye of the beholder?. *Cognitive psychology, 23*(3), 457-482.

896 Tanaka, J. W., Kaiser, M. D., Hagen, S., & Pierce, L. J. (2014). Losing face:
897       impaired discrimination of featural and configural information in the mouth
898       region of an inverted face. Attention, *Perception, & Psychophysics, 76*(4),
899       1000-1014.

900 Taylor, J., Shehzad, Z., & McCarthy, G. (2016). Electrophysiological correlates of
901       face-evoked person knowledge. *Biological psychology, 118*, 136-146.

902 Teichmann, L., Quek, G. L., Robinson, A. K., Grootswagers, T., Carlson, T. A., &
903       Rich, A. N. (2020). The influence of object-color knowledge on emerging
904       object representations in the brain. *Journal of Neuroscience, 40*(35), 6779-
905       6789.

906 Tranel, D., & Damasio, A. R. (1985). Knowledge without awareness: An autonomic
907       index of facial recognition by prosopagnosics. *Science, 228*(4706), 1453-
908       1454.

909 Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p
910       values. *Psychonomic bulletin & review, 14*(5), 779-804.

911 Wang, T. C., Liu, M. Y., Zhu, J. Y., Tao, A., Kautz, J., & Catanzaro, B. (2018). High-
912       resolution image synthesis and semantic manipulation with conditional gans.
913       In *Proceedings of the IEEE conference on computer vision and pattern
914       recognition* (pp. 8798-8807).

915 Wardle, S. G., Taubert, J., Teichmann, L., & Baker, C. I. (2020). Rapid and dynamic
916       processing of face pareidolia in the human brain. *Nature Communications,
917       11*(1), 4518. https://doi.org/10.1038/s41467-020-18325-8

918 Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a
919       100-ms exposure to a face. *Psychological science, 17*(7), 592-598.

920 Wheatley, T., Weinberg, A., Looser, C., Moran, T., & Hajcak, G. (2011). Mind
921       Perception: Real but Not Artificial Faces Sustain Neural Activity beyond the
922       N170/VPP. *PLOS ONE, 6*(3), e17960.
923       https://doi.org/10.1371/journal.pone.0017960

924 Wong, A. C. N., Palmeri, T. J., & Gauthier, I. (2009). Conditions for facelike
925       expertise with objects: Becoming a Ziggerin expert—but which type?.
926       *Psychological Science, 20*(9), 1108-1117.

927 Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., & Xiao, J. (2015). Lsun:
928       Construction of a large-scale image dataset using deep learning with humans
929       in the loop. *arXiv preprint* arXiv:1506.03365.

930 Yu, J., Xu, X., Gao, F., Shi, S., Wang, M., Tao, D., & Huang, Q. (2020). Toward
931       Realistic Face Photo-Sketch Synthesis via Composition-Aided GANs. *IEEE
932       Transactions on Cybernetics*.

933 Zhou, S., Gordon, M., Krishna, R., Narcomey, A., Fei-Fei, L. F., & Bernstein, M.
934      (2019). Hype: A benchmark for human eye perceptual evaluation of
935      generative models. In *Advances in Neural Information Processing Systems*
936      (pp. 3449-3461).