

# Vision Research

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

Manuscript Number:	
Article Type:	Full Length Article
Keywords:	Face perceptionDecodingFake facesArtificial intelligenceNeuroimaging
Corresponding Author:	Michoel L Moshel, M.A  AUSTRALIA
First Author:	Michoel L Moshel, M.A
Order of Authors:	Michoel L Moshel, M.A  Amanda K Robinson  Tijl Grootswagers  Thomas A Calrson
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 behaviour may be key in determining the 'real' in our new reality. Stimuli, code, and data for this study can be found at <a href="https://osf.io/n2z73/">https://osf.io/n2z73/</a> .
Suggested Reviewers:	Nancy kanwisher ngk@mit.edu  Bruno Rossion bruno.rossion@uclouvain.be  marcel van gerven m.vangerven@donders.ru.nl

Dear Editor,

We are pleased to submit our manuscript entitled “**Are you for real? Decoding hyperrealistic AI-generated faces from brain activity**” to be considered in *Vision Research* for publication.

This is, to our knowledge, the first study to investigate and analyse how hyperrealistic faces generated by cutting-edge AI are perceived differently to human faces.

Recent advances in AI technology and capabilities have made it possible to generate very realistic images which have, in turn, caused wide-spread concern over the human ability to tell apart real multimedia from fake. This concern is justified as AI-generated content has made its way into journalism, politics, and illegal counterfeiting. One important consideration is that the uniquely and distinctly human specialization for face perception has made it historically difficult to fool observers with artificial faces. However, current AI-generated faces may challenge this ability. The current study used computational analyses applied to behavioural and EEG neuroimaging data to examine in fine-grained detail how real humans and AI-generated faces are perceived.

Consistent with the brain’s sensitivity to artificial face appearance, we found it was possible to decode AI-generated faces from real faces using the EEG data. However, observers could not reliably distinguish fakes faces from real faces. Our EEG data reveal that real versus fake face decoding displayed a characteristic double-peak pattern associated with an initial in-depth processing and a secondary higher-level perception. Control analyses (face inversion and an increased presentation rate) showed that the decoding results were consistent with face-specific neural mechanisms and could not be attributed to low-level visual features. Our result showing that decoding did not translate into above-chance behavioural discrimination between real and fake faces has a number of implications for AI-generated face detection, especially in the possible ways protective standards can be established. The methods used in the current study provides an efficient way to investigate real versus fake distinctions both in the brain and behaviour that can be more generally applied to artificially-generated content in other domains.

We believe that this study is of timely and practical interest and addresses a concern that broadly appeals to researchers and the general public. The findings have important implications for understanding face perception and the ways in which we tackle the potentially detrimental and universal spread of artificially generated information.

We look forward to hearing from you and working towards publication of this manuscript in *Vision Research*.

Michoel Moshel<sup>a,b</sup>, Amanda Robinson<sup>b</sup>, Thomas Carlson<sup>b\*</sup> and Tijl Grootswagers<sup>b,c</sup>

<sup>a</sup>School of Psychology, Macquarie University, NSW 2109, Australia

<sup>b</sup>School of Psychology, University of Sydney, NSW 2006, Australia

<sup>c</sup>The MARCS Institute for Brain, Behaviour and Development, Western Sydney University, NSW 2214, Australia

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

2 Michoel L. Moshel <sup>a,b\*</sup>, Amanda K. Robinson<sup>a</sup>, Thomas A. Carlson <sup>a#</sup>, Tijl Grootswagers<sup>a,c#</sup>

3 <sup>a</sup> School of Psychology, University of Sydney, NSW, Australia

4 <sup>b</sup> School of Psychology, Macquarie University, NSW, Australia

5 <sup>c</sup>The MARCS Institute for Brain, Behaviour and Development, Western Sydney University, NSW, Australia

6 <sup>#</sup> shared authorship

7 \* corresponding author: [michoel.moshel@students.mq.edu.au](mailto:michoel.moshel@students.mq.edu.au)

## 8 Significance Statement

9 The generation of hyper-realistic images and faces has become increasingly possible in the last few years due to  
 10 advances in artificial intelligence (AI) technologies. Given the human expertise and specialization for face  
 11 perception, understanding how the brain is fooled by realistic AI-generated faces may be crucial in navigating  
 12 this new era of hyper-realism. We investigated how humans deal with hyper-realistic faces using behavioural  
 13 data and computational neuroimaging. We found that we could reliably decode AI-generated realistic faces  
 14 using people's neural activity, even though observers could not consciously report seeing differences between  
 15 real and realistic fake faces. Understanding this dissociation has significant implications for fake face detection  
 16 as well as face perception in general.

## 17 Abstract

18 Can we trust our eyes? Until recently, we rarely had to question whether what we see is indeed what exists,  
 19 but this is changing. Artificial neural networks can now generate hyperrealistic images that challenge our  
 20 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 behaviour may be key in determining the 'real' in our new reality. Stimuli, code, and data for this study can be found at <https://osf.io/n2z73/>.

## Keywords

Face perception      Decoding      Fake faces      Artificial intelligence      Neuroimaging

## Introduction

The novel and rapidly emerging phenomena of fake multimedia have swept through modern culture to the extent that the fake has become the expected norm (Adelani et al., 2020; Shen et al., 2019; Shu et al., 2017). The degree to which terms like ‘fake news’ or ‘photoshopped’ have become common parlance is indicative of a general and commonly experienced inability to distinguish between what is real and what is not (Fletcher, 2018). Meanwhile, AI technologies, in particular Generative Adversarial Networks (GANs), have been making increasingly rapid advances in generating realistic images with face generation as a major focus (Karras et al., 2019, 2020; Wang et al., 2018; Yu et al., 2020). These advances in realism have begun to have real-world consequences including undetectable videos of fake events (“Deepfakes”: Kietzmann et al., 2020), art and audio-visual counterfeits (Farokhmanesh, 2018), and fraudulent social media accounts (Gleicher, 2019). For instance, in 2019, Facebook announced that fake accounts were being created with profile pictures generated by artificial intelligence in an attempt to evade detection (Gleicher, 2019). Crucially, understanding how people respond to AI images, in terms of both behaviour and neural responses, will inform us about how realistic artificial images and faces are perceived differently to real ones, how this dissociation is encoded by the brain, and can ultimately aid in the development of future policy and strategies to curb the potentially nefarious uses of fake media.

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

61 Neuroimaging has provided useful insight into how face perception unfolds over time. Electroencephalography  
62 (EEG), which measures electrical activity at the scalp with very high temporal resolution, has been used to  
63 identify unique neural responses that reflect the temporal emergence and dynamics of facial processing (Bentin  
64 et al., 1996; Rossion et al., 2000). Wheatley and colleagues (2011) demonstrated the brain's discrimination of  
65 real and artificial faces by comparing neural responses to real faces with responses to doll faces. The authors  
66 found that both human and artificial faces elicited an N170, a face-specific neural response approximately  
67 170ms after image presentation. However, sustained positivity beyond 400ms was associated only with human  
68 faces, suggesting that this EEG potential could index a process that distinguishes between real and fake faces  
69 (Wheatley et al., 2011). Indeed, in other studies, sustained positivity, characterised by the late positive amplitude  
70 (LPP), increased as face realism increased, suggesting that real faces, more so than artificial faces, engage high-  
71 level attentional, semantic and identity evaluations (Schindler et al., 2017). The new generation of realistic faces

72 produced by GAN technology, however, is of a far superior quality than previously studied artificial faces and  
73 often practically indistinguishable from real faces. Whether the brain elicits neural indicators consistent with  
74 artificial fake detection for the new generation of GAN-produced images has yet to be seen. Considering that  
75 humans remain the gold standard of fake image and face detection (Natsume et al., 2019, Marra et al., 2018),  
76 examining the neural mechanisms in fake face detection is instrumental in understanding how to best tackle  
77 and understand the new age of fake media. EEG remains an ideal method to provide useful insights into the  
78 neural processing of fake GAN faces. Firstly, it allows for an insight into the sequential stages of face  
79 processing, from low-level visual features to holistic face perception. Secondly, closer examination at the  
80 neuronal population level enables us to answer at what temporal stages GAN face perception may differ from  
81 real face perception. Thirdly, using newer multivariate methods applied to EEG data enables analysis of signal-  
82 level information on a trial-by-trial basis and can pinpoint the precise temporal emergence of visual processing  
83 (Grootswagers, Robinson, & Carlson, 2019; Haynes & Rees, 2006; Teichmann et al., 2020).

84 With progressive advances in realistic image generation, have we reached a point where observers can no longer  
85 tell apart real from the fake? Can measuring the brain's response reveal how hyper-realistic fake faces are  
86 distinguished from real faces? We measured whether observers could behaviourally discriminate real faces from  
87 GAN-generated faces at two levels of face realism; one level of realism similar to fake images used in previous  
88 work (“unrealistic”), and another level which represents the current state-of-the-art hyper-realistic artificial  
89 images (“realistic”). We expected that participants would not be able to discriminate real from realistic faces  
90 but could for unrealistic faces, consistent with previous research using AI-generated faces (Hulzebosch et al.,  
91 2020; Zhou et al., 2019). To investigate whether we could decode real and fake images from brain activity we  
92 used time-resolved multivariate pattern analysis (MVPA) and EEG. To ensure the real and fake stimuli evoked  
93 typical categorical effects that could be decoded in the neural signal, we also included cars and bedrooms  
94 stimuli. We presented images upright in rapid sequences, which we have previously shown captures low- and  
95 high-level image processing (Grootswagers, Robinson, & Carlson; Oosterhof et al., 2016). To determine the  
96 contribution of low-level image properties, we used a much faster presentation rate (20Hz; Robinson et al.,  
97 2019) and also investigated how real/fake face processing is affected by image inversion, which limits high-  
98 level expert face-processing. Consistent with the brain’s sensitivity to artificial face appearance, we found it was

99 possible to decode real faces from GAN-generated faces at both levels of face realism using the EEG data.  
 100 However, when asked to behaviourally classify faces as either real or fake, a large group of participants could  
 101 differentiate the unrealistic, but not the realistic fake faces. Understanding dissociations between observer-  
 102 reported perceptions of fake images and the brain's response can yield important insights into human face  
 103 perception in general as well as raise possibilities for training observers to tell apart real from fake.



104

105 **Figure 1. Stimuli and design.** Experimental stimuli and design. A) Face, car and bedroom stimuli used in the  
 106 experiment from three conditions (real, realistic fake, unrealistic fake), taken from StyleGAN. B) Mean image  
 107 for each condition and the absolute pixel difference between 100 independent real images not used in the

experiment. Brighter colours (orange) indicate greater absolute differences. C) Experimental designs from left to right; behavioural experiment, 5Hz EEG experiment and 20Hz EEG experiment.

## Methods

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

## Participants

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

## Stimuli & Design

GAN-generated stimuli were obtained from StyleGAN output found at [shorturl.at/josOY](https://shorturl.at/josOY) (Karras et al., 2019). For a full description of the StyleGAN generative procedure and output, see Karras et al. (2019). Fake stimuli consisted of 25 faces, cars, and bedrooms at truncation levels of  $\Psi 0.5$  (realistic) and  $\Psi 1.0$  (unrealistic), (Figure 1A). To best match image statistics across real and fake images, real images were obtained from training images used for GAN output. These real training faces were obtained from the Flickr-Faces-HQ dataset (Karras et al., 2019). Real cars and bedrooms were randomly selected from the LSUN dataset (Yu et al., 2015). To maintain consistent aspect ratios, all images were cropped to a square aspect ratio and resized to a  $256 \times 256$  pixel dimension. No other filtering or editing was applied to the stimuli in order to provide a naturalistic demonstration of visual processing. To reduce obvious surface-level inconsistencies between real and fake images, real faces with eyes not facing frontward and/or with overly pronounced facial expressions (e.g. crying, laughing) were excluded. Upon surface inspection, we found no consistent delineating features between the



132 real and fake bedrooms and cars. All images were presented in both upright and inverted orientations totalling  
133 450 stimuli overall (Figure 1A).

134 Behavioural testing for real versus fake face discrimination was conducted online (Grootswagers, 2020). The  
135 experiment was programmed in jsPsych (De Leeuw, 2015) and hosted on Pavlovia.org (Peirce, 2019). Two  
136 hundred participants performed real or fake face judgements for one of four comparisons (50 in each group):  
137 1) upright unrealistic vs upright real, 2) upright realistic vs upright real, 3) inverted unrealistic vs inverted real,  
138 and 4) inverted realistic vs inverted real. Each observer was shown 50 images in total: 25 fake and 25 real.  
139 Participants were informed that 50% of the images were real photos and 50% were computer-generated and  
140 were instructed to choose whether each image was real or fake. Each image was individually presented on the  
141 screen for 200ms, followed by a blank screen until the participant pressed a button to indicate if the face was  
142 real or fake. Stimuli were presented at 256 x 256 pixel dimension against a grey background. Presentation of  
143 images was randomised, and each image was only presented once. The experiment took around 3-5 minutes to  
144 complete (Figure 1C).

145 For the EEG component, the experiment was presented in Psychopy2 (Peirce et al., 2019). Participants sat in  
146 a dimly lit room approximately 60cm away from a 1920 x 1080 pixel Asus computer monitor. Stimuli subtended  
147 approximately 6.4 degrees visual angle on a grey background with a white fixation circle superimposing the  
148 stimuli at approximately 1.3 degrees. Images were presented in a rapid serial visual presentation (RSVP)  
149 paradigm, whereby stimuli are presented in rapid succession, at 20Hz and 5Hz sequences (33ms image duration  
150 and 167ms or 16ms gap). There were 20 sequences at each presentation rate comprising 40 in total with 18,000  
151 images presented overall (with 20 repeats of each stimulus at each presentation rate). A sequence was started  
152 with a button press and lasted approximately 40 seconds. Subjects were instructed to fixate upon a white circle  
153 superimposed over each stimulus at the centre of the screen and told to respond by pressing any button on a  
154 4-way button box whenever they spotted the fixation circle turn red (Figure 1C). Fixation colour changes were  
155 randomised to occur between 2 and 5 times in each sequence. Length of colour change corresponded to the  
156 time of one image presentation (33ms). At the conclusion of the experiment, participants were debriefed and  
157 informed that half the images had been fake.

158 EEG recordings and preprocessing

159 Continuous EEG data were recorded using a 64-electrode Brain Products EEG cap (Standard 64Ch actiCAP;  
160 GmbH, Herrsching, Germany) at a sample rate of 1000-Hz. Ag/AgCl active electrodes were placed in  
161 accordance with a 10/20 international system (Oostenveld & Praamstra, 2001). Electrode gel was applied to  
162 the scalp under each electrode, aiming to reduce signal impedances to below 10k $\Omega$ . Stimulus onset was  
163 synchronised to the EEG using transistor-transistor logic (TTL) pulses from the stimulus presentation  
164 computer to a separate recording computer. Pre-processing of the EEG data was computed offline using  
165 EEGLAB (Delorme & Makeig, 2004). The continuous EEG data were filtered with a high-pass filter of 0.1-  
166 Hz and a low-pass filter of 100-Hz and re-referenced to the average of all electrodes. No notch filter was  
167 applied. The data were then separated into epochs corresponding to stimulus presentation ranging from 100ms  
168 to 1000ms pre and post-stimulus onset. This produced 180,000 pre-processed epochs for each participant.

169 Decoding analysis

170 Time-resolved MVPA decoding analysis of EEG data was implemented in MATLAB with the CoSMoMVPA  
171 toolbox (Oosterhof, Connolly, & Haxby, 2016). We used Linear Discriminant Analysis (LDA) classifiers as  
172 implemented in CoSMoMVPA in a leave-one-out cross-validation scheme. The LDA classifier estimated the  
173 probability of EEG data belonging to a certain group (e.g., real or fake) where the higher estimate is the  
174 predicted class (Grootswagers, Wardle, & Carlson, 2017). This was repeated at every time point, for every  
175 exemplar, and averaged across subjects to generate the mean cross-validation decoding performance at each  
176 time point. Classification performance was characterized as significant if it produced an above-chance accuracy  
177 (>50% for real versus fake decoding or 33% for 3-way category decoding). An above-chance decoding accuracy  
178 informs us that the EEG data contains information relevant the contrast of interest (Grootswagers, Wardle, &  
179 Carlson, 2017; Olivetti et al., 2012; Pereira et al., 2009).

180 Category Decoding Analysis

181 We performed a category decoding analysis to investigate whether there were meaningful differences among  
182 the face, car and bedroom stimuli. We used an image-by-sequence cross-validation approach (Grootswagers,

183 Robinson, & Carlson, 2019), which entailed training the classifier on all-but-one image from each of the three  
184 categories from all-but-one sequence and testing the classifier on left-out images from the left-out sequence.  
185 This ensured that the classifier had to generalize to novel exemplars to successfully decode between faces, cars,  
186 and bedrooms for each of the real, realistic, and unrealistic conditions (Carlson et al., 2013). Decoding accuracy  
187 was characterized by an above-chance classifier performance ( $>33\%$ ). Contrasts were broken down into  
188 presentation rate (5-Hz or 20-Hz), realism level (real, unrealistic, realistic), and configuration (upright, inverted).

#### 189 Real versus Fake Decoding Analysis

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

201 To map the spatial distribution of the signal, we repeated the real versus fake decoding analysis at separate  
202 locations on the scalp. For each channel, we selected the four closest neighbouring channels and performed  
203 the exact same decoding analysis described above on just this local cluster of channels, storing the resulting  
204 accuracies at the centre channel. This results in a channel topography of decoding results that provides insight  
205 into the spatial origins of the signal.

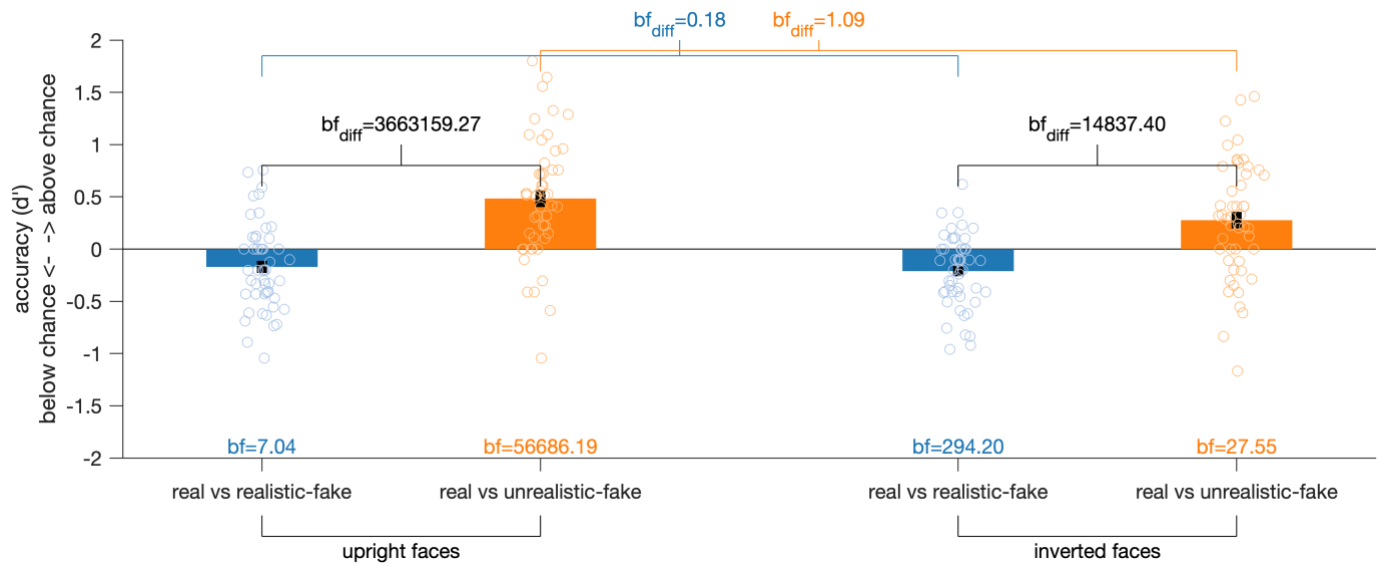
206 As an exploratory follow-up analysis, we examined the relationship between real-fake decoding accuracy and  
207 behavioural categorisation accuracy (Grootswagers et al., 2018; Ritchie et al., 2019). For each subject and each  
208 time point in the real-fake decoding analysis, we correlated (Spearman's  $\rho$ ) the image-specific average

classifier accuracies with their corresponding behavioural accuracies. We then performed group level inference on the resulting subject-wise time-varying brain-behaviour correlations. If successful real/fake decoding in EEG reflects the real/fake signal that is ‘used’ by the brain to guide behaviour (Grootswagers et al., 2018; Ritchie et al., 2019), then we would expect a positive correlation between image-specific EEG-classification accuracy and behavioural accuracy. That is, faces identified as real or fake by the classifier would also be identified as real or fake by the participants.

## Statistical inference

For the decoding and behavioural analyses, we used Bayesian statistics to characterize evidence arising from the data as either supporting the presence (alternative hypothesis) or absence (null hypothesis) of an effect. (Dienes, 2011; Jeffreys, 1998; Rouder et al., 2009; Wagenmakers, 2007). We used a standard JZS prior to calculate the null and alternative hypotheses (Rouder et al., 2018), which is a Cauchy distribution with a scale factor of 0.707 to determine the evidence of above-chance performance (e.g., >50% decoding) and a null-hypothesis point prior at chance-level (Morey & Rouder, 2011). For ease of interpretation, we thresholded Bayes factor (BF) values  $> 10$  for strong evidence for the alternative hypothesis and BF values  $< \frac{1}{3}$  as evidence in favour of the null hypothesis (Morey & Rouder, 2011). For the decoding analyses, BFs serve as continuous degrees of evidence across multiple time points and not specific hypothesis testing at single time points. Thus, isolated BFs at single time points which did not reach threshold were not treated as evidence for either hypothesis if the surrounding points did not reach threshold or were interspersed with below-threshold values. Rather, BFs were treated as evidence if surrounding points were at threshold (Mai et al., 2019). For the decoding analyses, we, in addition, computed corresponding frequentist statistics using sign-permutation tests (1000 permutations) and Monte-Carlo cluster statistics with TFCE as cluster-statistic (Smith & Nichols, 2009), corrected for multiple comparisons across time using the max-statistic method (Maris & Oostenveld, 2007).

231



232

233 **Figure 2. Behavioural discrimination of real and fake faces.** In an upright (left) and inverted (right)  
234 configuration, discriminability for real/realistic (blue) faces was below chance but above chance for  
235 real/unrealistic faces (orange). Performance was similar regardless of whether faces were upright or inverted.  
236 Bars show mean and standard error. Each circle represents the response of one subject in one condition. The  
237 Bayes Factors (displayed above the x-axis) compute the evidence for a difference from chance discriminability  
238 (50% accuracy), and difference between conditions (stimulus and orientation).

239 Results

240 Behavioural Performance

241 We were interested in whether participants could discriminate between real and fake faces. We calculated the  
242 proportion of images that were judged correctly as real or fake for each of the realistic/unrealistic and  
243 upright/inverted conditions and aggregated the judgements over participants. The main findings are presented  
244 in Figure 2. As indexed by d' discriminability analysis, we found that participants could reliably discriminate real  
245 from unrealistic fake faces (orange bars) but could not discriminate real from realistic fake faces (blue bars).  
246 Orientation had little effect on discriminability. Interestingly, performance in the real versus realistic face  
247 condition was below chance. Further inspection of the data revealed a general bias for participants to judge  
248 faces as real than as fake. When discriminating between upright real and realistic fake faces, observers correctly  
249 classified 63% (se = 0.026, BF > 100) of real faces and 31% (se = 0.023, BF > 100) of realistic fake faces. For  
250 discriminating between upright real and unrealistic fake faces, observers correctly classified 68% (se = 0.026,  
251 BF > 100) of real faces but performed at chance (49%, se = 0.027, BF = 0.16) at classifying unrealistic fake

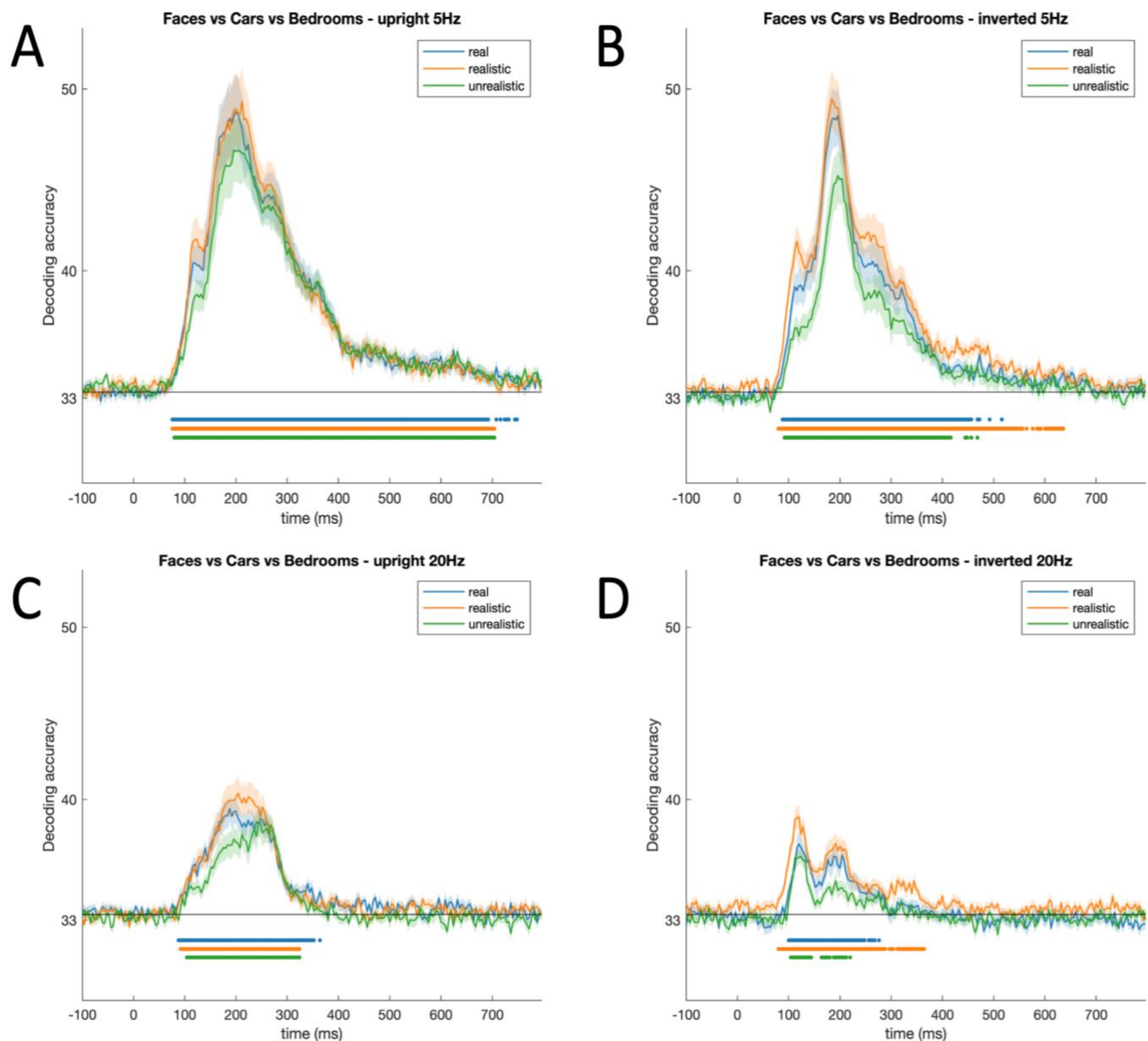
252 faces. Classification performances were similar for inverted faces Overall, observers could identify real faces  
253 (although were more biased to do so) but had much more difficulty spotting the fakes.

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

259 Categorical decoding analysis

260 To examine whether real and fake images evoked similar categorical decoding effects compared to the previous  
261 literature, we decoded image category (cars, faces, and bedrooms) at all levels of realism (real, realistic,  
262 unrealistic), (Figure 3). As expected, we observed similar category-related dynamics for the real, realistic and  
263 unrealistic images across all conditions. At a 5Hz presentation rate, we observed above-chance decoding for all  
264 categories at real, realistic, and unrealistic (Figure 3A). Decoding emerged and remained above-chance from  
265 100ms until 700ms post-stimulus onset with an early peak at 120ms, a second peak at 200ms and a third peak  
266 at 250ms-300ms.

267 We then tested how category decoding was affected by our control manipulations (inversion and presentation  
268 rate). We observed similar above-chance decoding for all categorical and realism levels upon inversion (Figure  
269 3B) and at a 20Hz presentation rate (Figure 3C), albeit less pronounced with simultaneous stimulus inversion  
270 and 20Hz presentation (Figure 3D). When upright and inverted, faces, cars, and bedrooms could be decoded  
271 at all levels of realism with similar temporal dynamics reported elsewhere (Grootswagers, Robinson, & Carlson,  
272 2019; Grootswagers, Wardle, & Carlson, 2017).



**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 (Figure 4A), above-chance decoding emerged and peaked for unrealistic faces at around 100ms, 200ms, and 300ms ( $BF > 10$ ) and fell below-chance at approximately

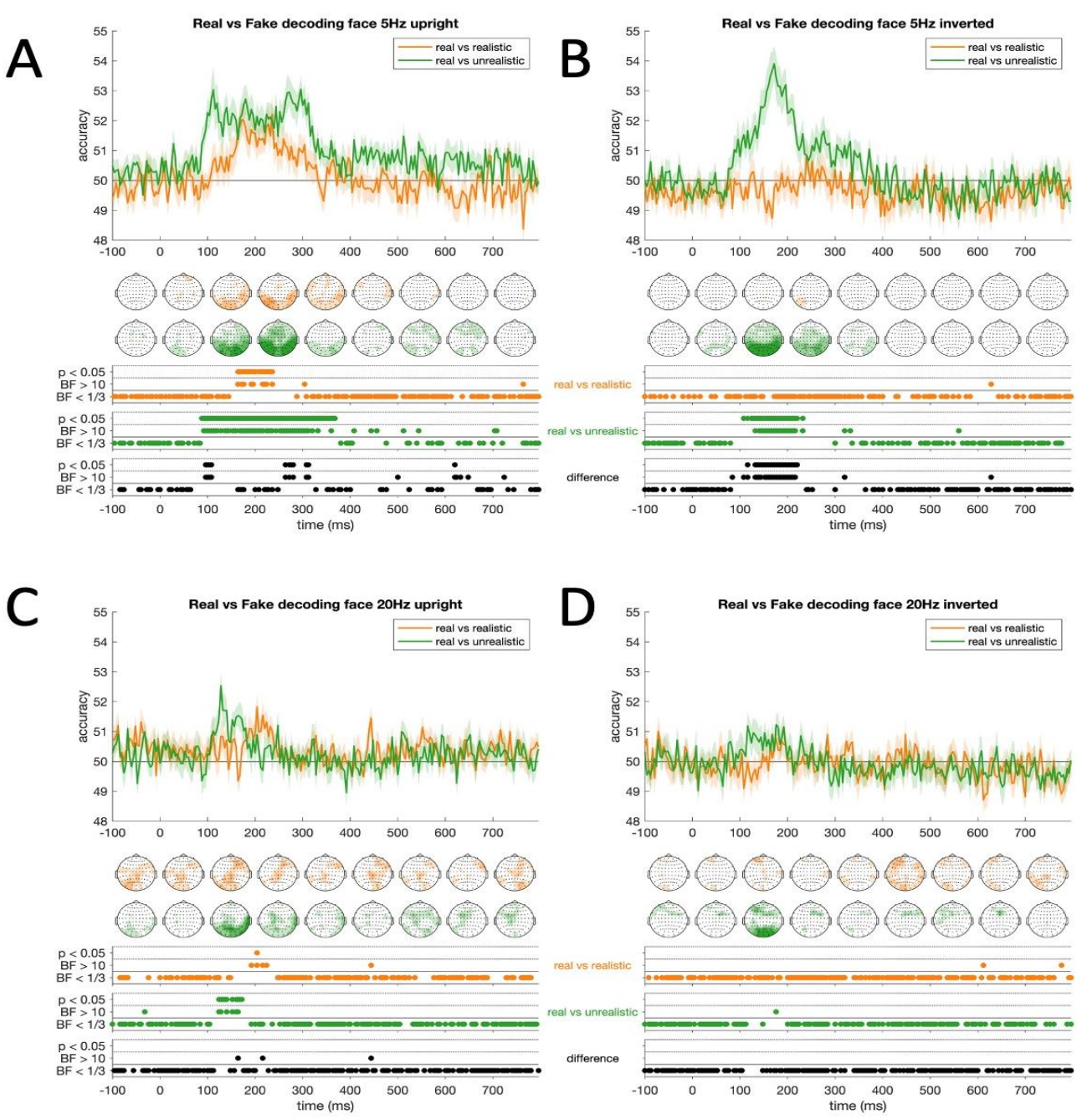
370ms ( $BF < \frac{1}{3}$ ). This decodability is reflective of early, rapid, low-level image perception followed by a later, higher-level, holistic decoding consistent with the temporal unfolding of face perception (Dobs et al., 2019; Balas & Koldewyn, 2013; Mühlberger et al., 2009). For realistic fake faces, decoding emerged at around 170ms and remained above-chance until approximately 240ms ( $BF > 10$ ), suggesting a higher-level basis for discrimination of realistic and real faces. Although observers could not reliably tell apart real faces from realistic fake faces, the EEG data contains signal information relevant to this distinction which meaningfully differs between realistic fakes and unrealistic fakes, and this signal appears to be constrained to a relatively short stage of processing.

If the information that we were decoding at 5Hz was reliant on image features rather than a face-processing effect, then we would predict that we could achieve a similar decoding result on inverted faces. However, at 5Hz and inverted (Figure 4B), only unrealistic fake faces were decodable from real faces. Above-chance decoding emerged at around 100ms ( $BF > 10$ ), peaked at around 170ms, and was at chance again at approximately 250ms ( $BF < \frac{1}{3}$ ). In contrast, realistic faces remained at-chance and were not decodable from real faces ( $BF < \frac{1}{3}$ ). This suggests that inversion, known to disrupt configural processing of faces, is similarly disrupting a face-specific mechanism accounting for decoding differences between realistic and unrealistic faces (Jacques, d'Arripe, & Rossion, 2007; Rossion et al., 2000).

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



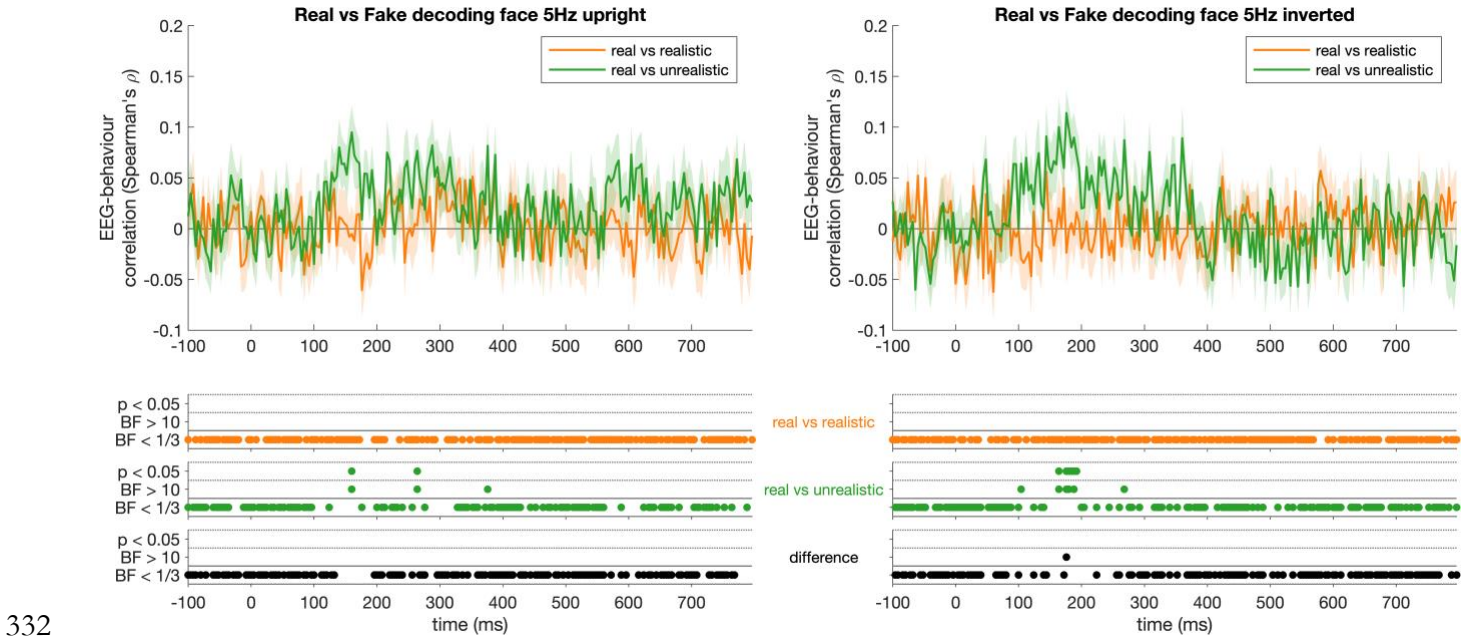
308 Lastly, at 20Hz and inverted (Figure 4D), decoding performance was at chance for realistic and unrealistic fake  
 309 faces ( $BF < 1/3$ ). This suggests that inversion plus a faster presentation rate is enough for the EEG data to no  
 310 longer contain any relevant information pertaining to real versus fake face distinctions. In other words,  
 311 configural processing has been disrupted to an extent that activity patterns evoked from fake faces were not  
 312 differentiable from activity evoked from real faces. As expected, real versus fake bedroom and car decoding  
 313 was not so evident and can be found on <https://osf.io/n2z73/>.



315 **Figure 4. Decoding real versus fake faces.** Different effects of orientation and presentation rate on decoding  
 316 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 where darker shades indicate contribution of channels to real/fake decoding. In the lowest panel, thresholded p-values and Bayes Factors indicate above-chance decoding or non-zero differences.

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



**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-

338 zero differences. Positive brain-behaviour correlations can be seen at around 170ms and 270ms for upright  
339 unrealistic faces (green) and at around 150ms-200ms for inverted unrealistic faces (BF>10).

340 Discussion

341 There is growing concern that hyperrealism is advancing at such a rate that humans will have difficulty  
342 discerning between what is real and what is fake (Fletcher, 2018; Khodabakhsh et al., 2019; Nightingale et al.,  
343 2017; Shen et al., 2019). Our results justify these concerns by revealing that observers cannot consciously and  
344 reliably identify realistic fake faces amongst real faces. However, using time-resolved EEG and multivariate  
345 pattern classification methods, we found that it was possible to decode both unrealistic *and* realistic fake faces  
346 from real faces using brain activity. This dissociation between behaviour and neural responses for realistic faces  
347 yields important new evidence about fake face perception as well as implications involving the increasingly  
348 realistic class of GAN-generated faces. Namely, the brain encodes information relevant to artificial face  
349 appearance even though humans do not consciously perceive any differences between GAN-generated faces  
350 and real faces.

351 Our behavioural results are consistent with previous research that suggests that observers typically display  
352 difficulties with correctly discriminating between real and realistic fake faces despite face expertise (Holmes et  
353 al., 2016; Nightingale et al., 2017; Sanders et al., 2019; Zhou et al., 2019). For example, in a two-alternative  
354 forced-choice task, participants would judge realistic artificial faces as being more realistic than human faces on  
355 a third of all trials (Sanders et al., 2019). Artificial faces made by GANs have also recently received attention  
356 and have been similarly demonstrated to fool observers (Hulzebosch et al., 2020; Isola et al., 2017; Zhou et al.,  
357 2019; Liu et al., 2020). As expected, we found that it was much harder to discriminate fake from real faces in  
358 our realistic condition relative to the unrealistic condition, confirming that the newer generation of GAN  
359 images are much more naturalistic. We presented faces for 200ms, which could be considered a brief exposure  
360 period, but the images were not masked so processing would have continued even after the images had  
361 disappeared (Robinson, Grootswagers, & Carlson, 2019). Given a long enough time to observe, Liu et al.,  
362 (2020) found that identifying artifacts such as “asymmetrical eyes” and “irregular teeth” in artificial faces can  
363 assist in spotting fakes. Presumably, assessing such details requires more time and eye movements. Indeed,  
364 observers can be trained to reliably spot fake faces by learning what to look for (Hills & Lewis, 2006; Tanaka

365 & Farah, 1993). Here, our primary focus was examining first impression responses by limiting the time spent  
366 looking at each face and giving participants unlimited time to make a response. Future studies may investigate  
367 whether training observers on GAN-generated faces enhances detection.

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

379 Indeed, for the 5Hz, upright condition, we found that the classifier successfully discriminated between  
380 unrealistic/real as well as realistic/real faces (Figure 4A). Decoding for unrealistic faces displayed a triple peak  
381 pattern, emerging at around 100ms maintained until around 370ms. Early decoding differences are consistent  
382 with rapid face detection and face-specific processing (Rossion et al., 2015; Dobs et al., 2019; Crouzet, Kirchner,  
383 & Thorpe, 2010; Wardle et al., 2020). The latter two peaks (at around 170-200ms and 270-320ms) have been  
384 similarly demonstrated to emerge in real versus artificial face perception (Wheatley et al., 2011; Balas &  
385 Koldewyn, 2013; Sagiv & Bentin, 2001; Schindler et al., 2017, Schindler et al., 2019, Wardle et al., 2020).  
386 Schindler et al (2017) suggest that early-stage N170 processing is related to assessing the structural configuration  
387 of faces as seen by a greater occipital involvement whilst the later-staged LPP, seen to increase linearly with  
388 face realism, suggests a deeper person-related, semantic involvement (also see Abdel Rahman, 2011, Taylor,  
389 Shehzad, & McCarthy, 2016). Differences at the triple peak correspond to N250 and P300 components typically  
390 associated with face familiarity (Collins et al., 2018) and semantic information (Tanaka et al., 2006), the latter

391 especially important for behaviour (Hanso et al., 2010). In contrast, realistic/real decoding displayed a single-  
392 peak emergence between around 170ms to 240ms indicating a difference in processing between realistic and  
393 unrealistic faces. Namely, that differences in perception between real and realistic faces were constrained to the  
394 170ms time period. Indeed, in comparing human faces to doll faces and artificial faces, others have shown that  
395 only the human faces typically evoke sustained neural responses beyond the N170 component necessary for  
396 higher-order perception (Balas & Koldewyn, 2013; Wheatley et al., 2011). Balas and Koldewyn (2013) found  
397 that the N170 was better characterised by encoding deviations from facial appearance than it was for animacy  
398 perception. In other words, realistic faces were perceived as configurally different to real faces, but that only  
399 unrealistic faces engaged later processing necessary for high-order animacy or familiarity perception. Overall,  
400 earlier decoding for unrealistic faces, consistent with apparent low-level image differences (Figure 1B), suggests  
401 that early and low to mid-level processing differences may account for decodability between real and unrealistic  
402 faces. The decoding for realistic faces, by contrast, emerges later and is constrained to the 170ms time period,  
403 suggesting a face-specific configural process may be responsible for this distinction.

404 Assessing fake/real decoding for inverted faces allows us to evaluate whether the fake/real distinction relies on  
405 mechanisms that are responsible for the superiority in face recognition for upright faces relative to inverted  
406 faces. Inversion disrupts the configural processing of faces by making them appear more like objects whilst  
407 retaining low-level stimulus attributes (Eimer, 2000; Leder & Bruce, 2000; Rousselet et al., 2003). Firstly, we  
408 found that inversion led to the disruption of decoding for realistic faces (Figure 4B). In contrast, we found that  
409 decoding for unrealistic inverted faces was preserved but less sustained when compared to upright. The peak  
410 in decoding may be reflective of increased featural processing for inverted unrealistic faces, also seen to occur  
411 with distorted or ‘Thatcherized’ faces (Carbon et al., 2005; Milivojevic et al., 2003). Lack of above-chance  
412 decoding for inverted realistic faces may reflect the contribution of high-level, expertise-driven capabilities for  
413 upright fake face detection when face processing mechanisms, rather than object processing, were available.  
414 Overall, we found that upon stimulus inversion our decoding results were consistent with a face-specific or  
415 expertise response, such that realistic fake faces could not be discriminated from real faces when typical face  
416 perception was disrupted, even though the same visual features were present.

417 The presentation of images at a faster presentation rate limits the consolidation of each image and build-up of  
 418 higher-order representation (Grootswagers, Robinson, & Carlson, 2019)., allowing an analysis of the  
 419 contribution of low-level processing. At a faster presentation rate of 20Hz, we found that upright fake faces  
 420 could be discriminated from real faces for the realistic and unrealistic conditions (Figure 4C). Indeed, early,  
 421 low-level visual processing is fairly unaffected by image presentation durations (Grootswagers, Robinson, &  
 422 Carlson, 2019). Observing less sustained decoding is consistent with the limited capacity and extent of visual  
 423 processing since each image is masked by every successive image to a greater extent and therefore places limits  
 424 on visual processing compared to a slower presentation rate (Collins, Robinson, & Behrmann, 2018; Robinson,  
 425 Grootswagers, & Carlson, 2019). Additionally, higher-level, identity or semantically related face information  
 426 discernible in the slow condition was possibly limited at the faster presentation rate consistent with Collins et  
 427 al. (2018). In sum, we found that unrealistic faces could be decoded upon inversion and at a faster presentation  
 428 rate suggesting the contribution of low-level visual differences. By contrast, we could not decode realistic faces  
 429 when inverted, but we could decode at a faster presentation rate, indicating that fake/real perception was likely  
 430 driven by expertise and face-specific processing.

431 Interestingly, we found that neural differences between real and realistic fake faces did not translate into a  
 432 reliable behavioural decision for realistic face discrimination at the population level. We found a brain-  
 433 behaviour correlation at around 150ms-200ms for unrealistic versus real faces, suggesting that this time period  
 434 of processing is important for behaviour. However, the same correlation was not observed for the realistic  
 435 faces. One possibility is that whilst our data indicates that a realistic fake/real signal is present, this signal gets  
 436 'lost' in the visual hierarchy and consequently remains uninformative for behaviour. For instance, although  
 437 animacy categorisation can be decoded throughout the entire ventral visual stream, this information is most  
 438 suitably formatted for behaviour in higher-level visual areas like the ventral occipital and parahippocampal  
 439 cortex (Grootswagers, Cichy, & Carlson, 2018). Since decoding unrealistic/real faces was more sustained than  
 440 realistic/real faces, associated more with in-depth face processing at later stages (i.e., LPP), it is possible that  
 441 this level of extended processing is required for behavioural "readout" (see de-Wit et al., 2016; Grootswagers,  
 442 Cichy, & Carlson, 2018; Ritchie, Kaplan, & Klein 2019). Yet, the highest brain-behaviour correlation for  
 443 unrealistic faces was observed at 150-200ms, a time when decoding was not reliably different between the

444 realistic and unrealistic condition. This has a number of implications. In an applied setting such as cyber security  
445 or Deepfakes, examining the detection ability for hyper-realistic fake faces might be best pursued using machine  
446 learning classifiers applied to neuroimaging data rather than targeting behavioural performance. As we have  
447 shown, the former contains discriminative relevance whereas observers may actually perform worse than  
448 chance given the decision (and a brief glance). A third related possibility is that the decodable real/fake face  
449 signal is operating below conscious access and therefore is not picked up by our behavioural task. This is  
450 reminiscent of findings that individuals with prosopagnosia who cannot behaviourally classify or recognise  
451 faces as familiar or unfamiliar nevertheless display stronger autonomic responses to familiar faces than  
452 unfamiliar faces (Tranel & Damasio, 1985). Similarly, what we have shown in this study is that participants  
453 could not reliably discriminate between real and realistic fake faces even though we could accurately decode  
454 this difference from their neural activity. Still, it is possible that a different behavioural task may have yielded a  
455 better performance. Forced to respond via a two-alternative forced-choice task or an implicit task such as face  
456 familiarity or trustworthiness may have engaged different behavioural processes more conducive for real/fake  
457 face discrimination. For instance, behaviourally categorising faces as threatening, competent, or trustworthy  
458 has been shown to occur as quickly as 33- 100ms after onset (Bar et al., 2006; Willis & Todorov, 2006).  
459 Conversely, real or fake judgments may occur as late as 240ms after stimulus presentation (Zhou et al., 2019).  
460 Therefore, future work could investigate whether judgments about face trustworthiness or threat may be a  
461 better cue for detection than real or fake.

462 In sum, we found that there is a dissociation between the ability of participants to categorise faces as real or  
463 fake and the decodability of this distinction in the brain. In other words, although the brain can ‘recognise’ the  
464 difference between real and realistic fake faces, observers cannot consciously tell them apart. Our findings of  
465 the dissociation between brain response and behaviour has implications for the ways in which we study fake  
466 face perception, the questions we pose when asking about fake image identification, and the possible ways in  
467 which we can establish protective standards against fake image misuse.

468 Future studies may investigate the contribution of face expertise for decoding and behaviour. Expertise  
469 influences how deeply and configurally a face is perceived allowing for more subtle identification of spatial

relations, features, and same-race faces (Wong et al., 2009; Tanaka, 2001; Tanaka & Taylor, 1991; Hancock & Rhodes, 2008; Meissner & Brigham, 2001). Indeed, individuals with digital manipulation training and experience (i.e., photo-editing and photography) are more able to identify fake images than non-experienced individuals (Shen et al., 2019). Having the same participants participate in both the EEG and behaviour experiments may be useful in exploring inter-individual differences and the influence of expertise.

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



489 Acknowledgements

490 The authors would like to thank Sophia Shatek for her help with EEG recording and Roy Raskolnikov for his  
491 suggestions and support. This research was supported by ARC DP160101300 (TAC), ARC DP200101787  
492 (TAC), and ARC DE200101159 (AKR). The authors acknowledge the University of Sydney HPC service for  
493 providing High Performance Computing resources. The authors declare no competing financial interests.

494 References

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

497 Adelani, D. I., Mai, H., Fang, F., Nguyen, H. H., Yamagishi, J., & Echizen, I. (2020, April). Generating  
498 sentiment-preserving fake online reviews using neural language models and their human-and machine-  
499 based detection. In *International Conference on Advanced Information Networking and Applications* (pp. 1341-  
500 1354). Springer, Cham.

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

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

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

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

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

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

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

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

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

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

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

524 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.

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

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

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

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

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

Grootswagers, T., Wardle, S. G., & Carlson, T. A. (2016). Decoding Dynamic Brain Patterns from Evoked Responses: A Tutorial on Multivariate Pattern Analysis Applied to Time Series Neuroimaging Data. *Journal of Cognitive Neuroscience*, 29(4), 677–697. [https://doi.org/10.1162/jocn\\_a\\_01068](https://doi.org/10.1162/jocn_a_01068)

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

- Hancock, K. J., & Rhodes, G. (2008). Contact, configural coding and the other-race effect in face recognition. *British Journal of Psychology*, 99(1), 45-56.
- Hanso, L., Bachmann, T., & Murd, C. (2010). Tolerance of the ERP signatures of unfamiliar versus familiar face perception to spatial quantization of facial images. *Psychology*, 1(03), 199.
- Haynes, J. D., & Rees, G. (2006). Decoding mental states from brain activity in humans. *Nature Reviews Neuroscience*, 7(7), 523-534.
- Hills, P. J., & Lewis, M. B. (2006). Short article: reducing the own-race bias in face recognition by shifting attention. *Quarterly Journal of Experimental Psychology*, 59(6), 996-1002.
- Holmes, O., Banks, M. S., & Farid, H. (2016). Assessing and improving the identification of computer-generated portraits. *ACM Transactions on Applied Perception (TAP)*, 13(2), 1-12.
- Hulzebosch, N., Ibrahim, 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](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](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.

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.

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.
- Tanaka, J. W., & Taylor, M. (1991). Object categories and expertise: Is the basic level in the eye of the beholder?. *Cognitive psychology*, 23(3), 457-482.
- Tanaka, J. W., Kaiser, M. D., Hagen, S., & Pierce, L. J. (2014). Losing face: impaired discrimination of featural and configural information in the mouth region of an inverted face. *Attention, Perception, & Psychophysics*, 76(4), 1000-1014.
- Taylor, J., Shehzad, Z., & McCarthy, G. (2016). Electrophysiological correlates of face-evoked person knowledge. *Biological psychology*, 118, 136-146.
- Teichmann, L., Quek, G. L., Robinson, A. K., Grootswagers, T., Carlson, T. A., & Rich, A. N. (2020). The influence of object-color knowledge on emerging object representations in the brain. *Journal of Neuroscience*, 40(35), 6779-6789.
- Tranel, D., & Damasio, A. R. (1985). Knowledge without awareness: An autonomic index of facial recognition by prosopagnosics. *Science*, 228(4706), 1453-1454.
- Wagenmakers, E. J. (2007). A practical solution to the pervasive problems of p values. *Psychonomic bulletin & review*, 14(5), 779-804.
- Wang, T. C., Liu, M. Y., Zhu, J. Y., Tao, A., Kautz, J., & Catanzaro, B. (2018). High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 8798-8807).
- Wardle, S. G., Taubert, J., Teichmann, L., & Baker, C. I. (2020). Rapid and dynamic processing of face pareidolia in the human brain. *Nature Communications*, 11(1), 4518. <https://doi.org/10.1038/s41467-020-18325-8>
- Willis, J., & Todorov, A. (2006). First impressions: Making up your mind after a 100-ms exposure to a face. *Psychological science*, 17(7), 592-598.
- Wheatley, T., Weinberg, A., Looser, C., Moran, T., & Hajcak, G. (2011). Mind Perception: Real but Not Artificial Faces Sustain Neural Activity beyond the N170/VPP. *PLOS ONE*, 6(3), e17960. <https://doi.org/10.1371/journal.pone.0017960>
- Wong, A. C. N., Palmeri, T. J., & Gauthier, I. (2009). Conditions for facelike expertise with objects: Becoming a Ziggerin expert—but which type?. *Psychological Science*, 20(9), 1108-1117.

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

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

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