

More Real than Real: A Study on Human Visual Perception of Synthetic Faces

Federica Lago¹, Cecilia Pasquini¹, Rainer Böhme², Hélène Dumont³,
Valérie Goffaux³, and Giulia Boato¹

¹Department of Information Engineering and Computer Science, University of Trento

²Department of Computer Science, University of Innsbruck

³Institute of Neuroscience, Université Catholique de Louvain

Deep fakes have become popular recently. The term refers to doctored media contents where one face's is swapped with someone else's face or performs someone else's face movements. In the last couple of years, numerous video clips, often involving celebrities and politicians, have gone viral on social media platforms. This has been enabled by easy to use apps capable to process user-generated content in real-time.

Informally, deep fakes can be defined as realistic digital media (images, videos, or audio tracks) depicting untruthful content, obtained either by manipulating pristine material or generated from scratch. The attribute *deep* refers to the use of algorithms based on deep learning, a subfield of modern Artificial Intelligence (AI), which pushed the boundaries for many applications including media data manipulation and generation. Besides offering exciting opportunities in several fields (such as entertainment, content production, e-learning, and e-health), these advanced creation technologies are now widely recognized as a pressing threat to the reliability of visual information.¹ This may have severe consequences for the digital identity and reputation of individuals.

Many cases of misuse reported in the past months demonstrate the potential impact of manipulated data on disinformation. For example, the last US presidential campaign witnessed the viral diffusion of multiple deceptive visual manipulations: amongst others, an altered video of Joe Biden greeting the wrong state in a public speech, as well as retouched pictures of celebrities untruthfully implying their endorsement for Donald Trump.² Not very surprisingly, deep fakes become increasingly ubiquitous on the web: the estimated number of deep fakes on the web doubled every six months since 2018, reaching a total of 85'000 in 2020,³ although it is not possible to tell which

¹<https://www.theguardian.com/technology/ng-interactive/2019/jun/22/the-rise-of-the-deepfake-and-the-threat-to-democracy>

²<https://www.fastcompany.com/90575763/we-have-the-technology-to-fight-manipulated-images-and-videos-its-time-to-use-it>

³<https://sensity.ai/how-to-detect-a-deepfake-with-sensity/>

fraction is actually misuse.

Images depicting synthetic faces of non-existing people are easy to produce. They can be created at scale using Generative Adversarial Networks (GAN), a variant of deep learning. These images result realistic and hard to recognize as synthetic, and thus can be used for scam or to facilitate fraud.⁴ Turning back to the 2020 election, a fake report on Joe Biden’s son and his connections to China was disseminated by a fabricated digital identity⁵; the alleged author was a Swiss security analyst, portrayed over the web by a synthetic face presumably created with StyleGAN2 [1], the state-of-the-art GAN for still image generation. A few months earlier, a 17-years-old student used a GAN-generated face to impersonate and promote the campaign of a fully fictitious congress candidate.⁶ The hoax so credible that Twitter activated the “verified” icon for the fake account.

Against this backdrop, it is no surprise that different research communities, from multimedia security to computer vision, joined efforts towards the automated detection of AI generated media. A variety of methods have been proposed and refined over the last years [2]. This way accompanied by the development of benchmark datasets (e.g., FaceForensics++ [3]) and global contests (e.g., Facebook Deep fake Detection Challenge).

While forgeries are almost as old as photography itself, the striking features of the latest advancements in image and video manipulations are the higher accessibility combined with the ever increasing level of realism. This is particularly true for still images, which today’s GANs impressively synthesise through easy-to-use interfaces.⁷ By contrast, the creation of realistic manipulated videos of arbitrary subjects and in high resolution still requires significant resources and expertise.

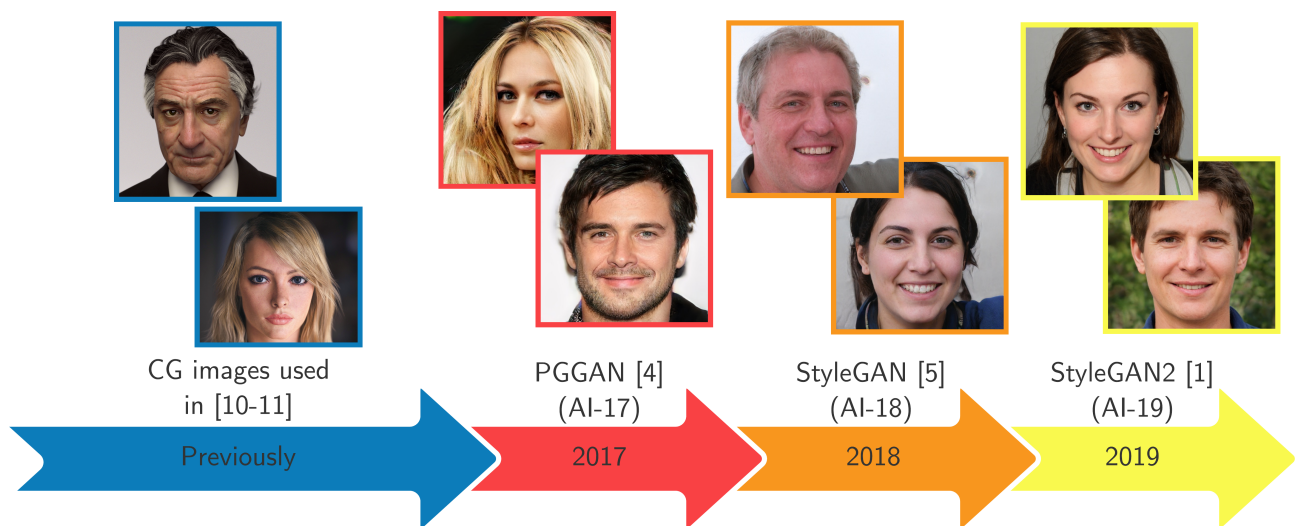


Figure 1: Examples of synthetically generated images over the years. Digital portraits on the left were edited manually, while the following ones are generated through GANs proposed in recent years. We indicate them as AI-17, AI-18, and AI-19 and keep this convention throughout the column.

⁴<https://www.ft.com/content/b50d22ec-db98-4891-86da-af34f06d1cb1>

⁵<https://www.nbcnews.com/tech/security/how-fake-persona-laid-groundwork-hunter-biden-conspiracy-deluge-n1245387>

⁶<https://edition.cnn.com/2020/02/28/tech/fake-twitter-candidate-2020/index.html>

⁷See for instance <https://thispersondoesnotexist.com/>

Figure 1 shows examples of images generated by different GANs: three networks denoted as AI-17 (PGGAN [4]), AI-18 (StyleGAN [5]) and AI-19 (StyleGAN2 [1]) are reported. By looking at the manually-edited computer generated portraits on the left, which was the only way to obtain realistic images until a few years ago, one can appreciate the huge leap forward in terms of human likeness achieved through AIs. Recent progress is best visible at the level of detail. The latest GAN images no longer contain artifacts and inaccuracies, such as mismatches in eye color, ear shape, teeth rendering, or face symmetry. Thus, they evade detection by automated tools relying on exactly these artifacts [2]. Outdated algorithms, however, might not be the only ones deceived by last generation data: we ask if human viewers can still tell the difference.

In this column, we present research that measures human’s ability to distinguish between real and synthetic face images when confronted with cutting-edge AI-based creation technologies.

Human perception of manipulated data has been addressed in previous work from several perspectives. In the field of neuroscience, researchers have investigated how the use of synthetically generated stimuli alters people’s ability in face processing tasks, such as face recognition [6]. In the field of multimedia forensics, researchers have studied to which extent humans perceive face morphing operations, i.e., a specific kind of face manipulation where two faces are blended together to obtain a third hybrid face that carries characteristics of both original subjects. This is particularly relevant for face authentication applications that prevent unauthorized access to places or services, where it is fundamental to study how new morphs affect the ability of humans and algorithms to recognize them as synthetic [7].

Regarding the discrimination of real and synthetic data, the most prominent works are those conducted by Hany Farid and his team, who carried out a series of experiments between 2007 and 2017 on the ability of humans to identify computer generated characters, starting with generic content [8] and then focusing on faces [9, 10, 11]. Their latest studies conclude that humans, on average, are still able to correctly recognize synthetic images.

While those studies are relatively recent, a stringent research question is whether those findings still hold in the light of the striking latest advancements in AI-based synthetic image creation. To the best of our knowledge, no existing studies on the human perception involve last generation synthetic data. To address this gap, we have designed and conducted a perceptual experiment where a wide and diverse group of volunteers has been exposed to synthetic face images produced by state-of-the-art GANs (i.e., AI-17, AI-18, AI-19, in Figure 1). In the remainder of this column, we report the results of this experiment. They reveal how strongly the human ability to discriminate synthetic from real should be called into question.

1 Human ratings of real and synthetic faces

	REAL	AI-17	AI-18	AI-19	
Stimuli sample	150	50	50	50	300
Participant's sample	15	5	5	5	30

Table 1: Distribution of images belonging to each dataset for the whole stimuli sample and for each randomly generated sequence of images seen by a participant (participant's sample).

The perceptual experiment has been carried out through a dedicated, self-hosted web interface which displays a sequence of stimuli varied between participants. The database of stimulus material consisted of 300 images, equally distributed between real and synthetic ones. Real images (from now on referred to as REAL) were selected from the FFHQ dataset [5]; synthetic images were created by the three different GANs, as summarized in Table 1. As mentioned earlier, it is worth observing how the creation technologies suffer from different kind of visual imperfection (see Figure 1): PGGAN [4] images (AI-17) typically present several face artifacts (e.g., different eyes color and size, unnatural hair shape), while StyleGAN [5] (AI-18) often produces visible blobs, especially in the background; StyleGAN2 [1] (AI-19) is the most advanced method and overcomes such limitations almost entirely.

The experiment was split into five parts:

- *Briefing*: participants were informed about the purpose of the study, the expected duration, the target group (18 years and above), the voluntary and anonymous nature of the study as well as their rights concerning the protection of personal data. Active consent was sought before proceeding to the next part.
- *Questionnaire*: participants were asked to provide demographic information, to self-assess their ability to recognize faces, and to report their a priori familiarity with deep fakes;
- *Warm-up phase*: the task of the main experiment was explained and participants had the opportunity to practice with the user interface (see Figure 2). It displayed a face image sized $10 \times 10 \text{ cm}^8$ ($\approx 4 \times 4$ inches) at the center of the screen for 3 seconds. Participants were asked to rate this image by indicating their agreement with the statement “This image is synthetic” on a 7-point rating scale where the leftmost 1 was labelled “Completely disagree”, (i.e., the image is real) and the rightmost 7 with “Completely agree” (i.e., the image is synthetic). The midpoint, labelled “Unsure,” was selected by default.

The response time was unconstrained: participants could select and confirm their rating at any time during or after the 3 seconds in which the image was displayed.

All participants saw the same six images in the warm-up phase. The ratings in this phase were not used in the analysis.

⁸Users whose devices were detected to be unfit for this task were not invited to continue the experiment.

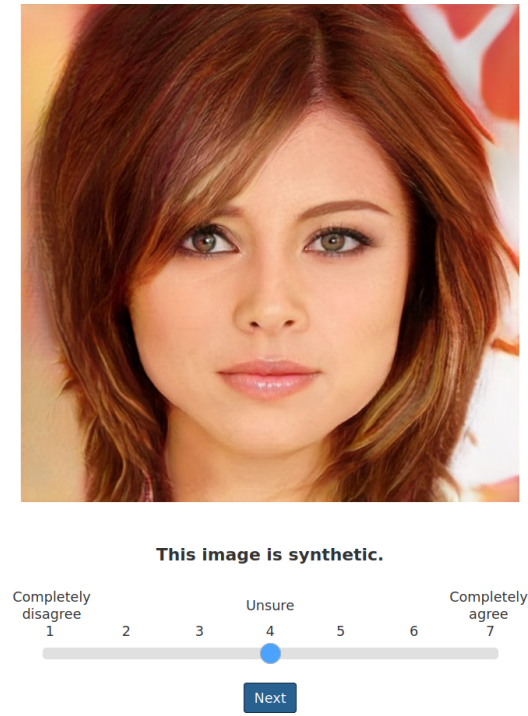


Figure 2: Example screenshot of the user interface for the main task in the English version. The specific face image shown here belongs to AI-17.

- *Main experiment*: participants sequentially complete a total of 34 task instances. A subset of 30 images (the participant’s sample) is randomly selected for each participant from the overall collection of stimuli material. The images in each participant’s sample are distributed as reported in Table 1 (bottom row). Participants were not informed that images were half real and half synthetic.

The remaining four images serve as control a set (two real and two synthetic ones) and were selected to be trivial to rate correctly. The same four images were shown to all participants to check for their compliance with the task.

- *Exit question and debriefing*: finally, participants were invited to describe in an open-ended question their strategy or cues used to rate the images. A response set was stored after confirming consent, followed by a debriefing. Participants did not learn their performance in order to avoid disappointment and to discourage improvement attempts by repeated participation.

Participants were recruited by extending invitations among professional and personal contacts of all researchers involved (spanning four countries). In order to reach a more diverse sample, the webpage has been made available in four languages: English, French, German, and Italian. Participants completed the experiment remotely on their own devices.

We collected response sets of 630 participants from 38 countries over the field time of 4 months (9th July – 13th November 2020). Participants were mostly from Europe (93%), younger than average ($M = 28.15$ years, $SD = 10.56$), and of balanced gender: 45.5% identified as females, 52.9% as males, while the remaining preferred to not declare (1.4%) or identified as non-binary

(0.2%). On average, it took participants 8.43 minutes to complete a response set (median: 7.61). An estimate of 1000 users have accessed the web interface but did not submit a response set.

In order to limit potential biases, we cleaned the data by discarding participants who exhibited indicators of low interest or distraction. To this end, we removed 30 (4.76%) response sets containing outliers in the form of:

- skip rate (how many times a participant clicked on the “Next” button without moving the scale slider);
- median response time to rate task instances; or
- error rate on the control set (wrong rating on the four control images).

We used as cutoffs the 99th percentile of each measure (32.35% for the skip rate, 7.41 seconds for the median response time, and 3 out of 4 for the error rate on the control set). A total of 986 (5.48%) cases where the scale slider was not moved from its default position were treated as item non-responses.

2 Results

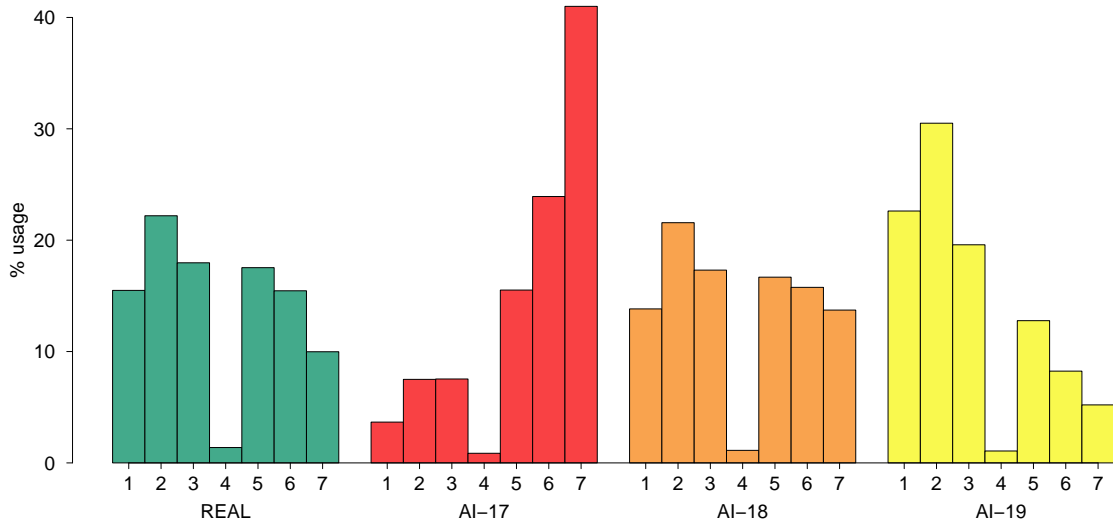


Figure 3: Percentage of scale values (from 1 to 7) used to rate the agreement with the statement “This image is synthetic”. The $N = 17014$ ratings from 600 subjects are broken down by datasets.

Figure 3 shows the distribution of the ratings over all participants on images belonging to different datasets. The central value (marked as “Unsure” in the scale) is low because unchanged default values were treated as non-responses, although, in general, participants avoided it.

Observe that the distributions of answers for REAL and AI-18 are very similar, while for AI-17 and AI-19 they are essentially reversed. In particular for AI-17, participants were able to correctly recognize the images as synthetic 80% of the times, of which 41% with the highest level of agreement

with the statement “This image is synthetic”. However, for the successive AIs, the distribution of responses drifts towards the left end of the scale, suggesting that, **while images generated with earlier AI were still relatively easy to recognize as synthetic, data produced by the newest AIs are increasingly perceived as being real**. To further investigate the results, we drill down into the dataset using different metrics.

Realism rate. We aggregate response values in the range from 1 to 3 as a judgement for real and define the **realism rate** as the percentage of images judged as real. In Figure 4a, the *realism rate* is computed on all the instances of all images belonging to a certain dataset. It is striking to observe that the highest *realism rate* (68%) is achieved by AI-19, while REAL images do not exceed 52%. The difference is statistically significant (Welch’s $t(121) = -6.8$, $p < .001$). In other words, **synthetic faces generated through StyleGAN2 are judged as real more often than real faces**.

Figure 4b provides a richer visualization as it shows the distribution over $[0, 100]$ of *realism rates* of individual images (i.e., computed over multiple task instances involving the same image) of the same dataset. It appears that, **while some real images were almost never considered to be real (realism rate ≈ 0), all images in AI-19 were judged as real at least 40% of the times**. Furthermore, skewness values clearly show that, while for REAL and AI-18 the distribution is almost symmetrical (0.016 and 0.149, respectively), the distributions of AI-17 and AI-19 are skewed in opposite directions (0.613 and -0.445, respectively).

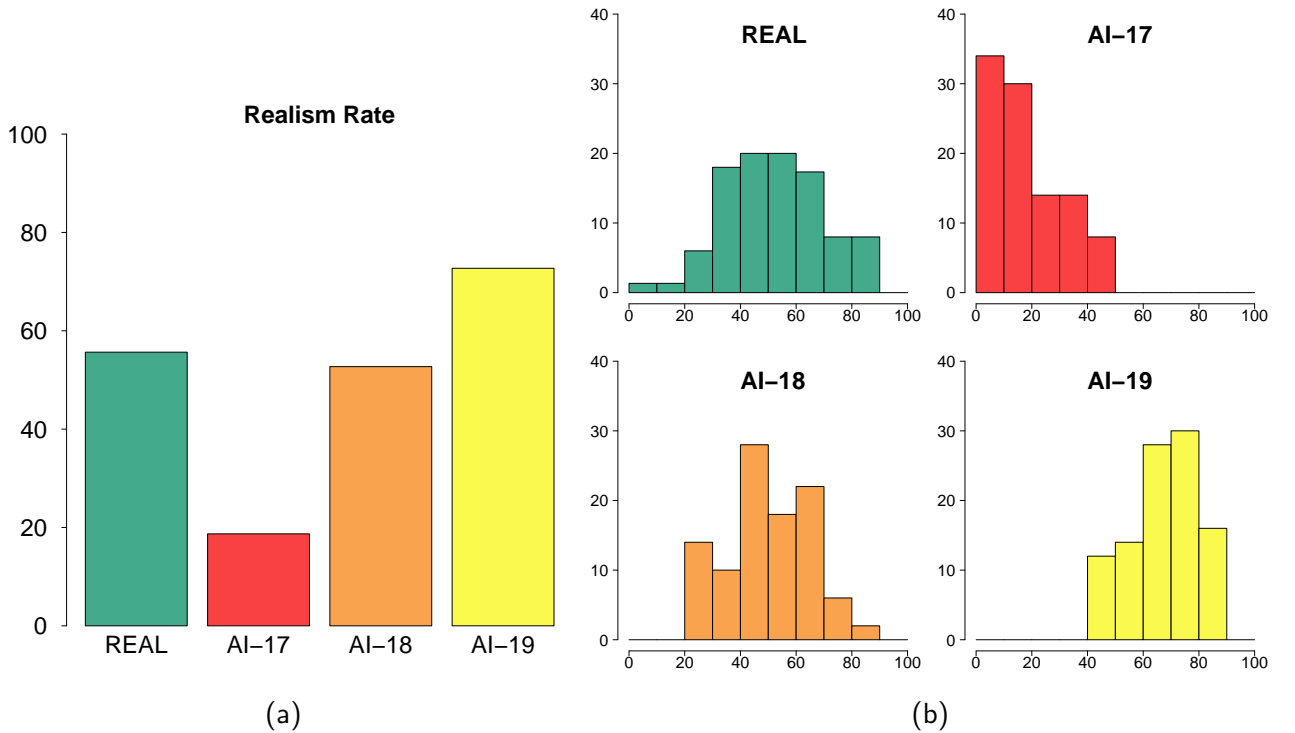


Figure 4: (a) shows the *realism rate* per dataset, while (b) shows the distribution of the *realism rates* over the images in each dataset. In (b) the *realism rate* is on the x-axes, while the y-axes indicates the relative frequency of a *realism rate* bin.

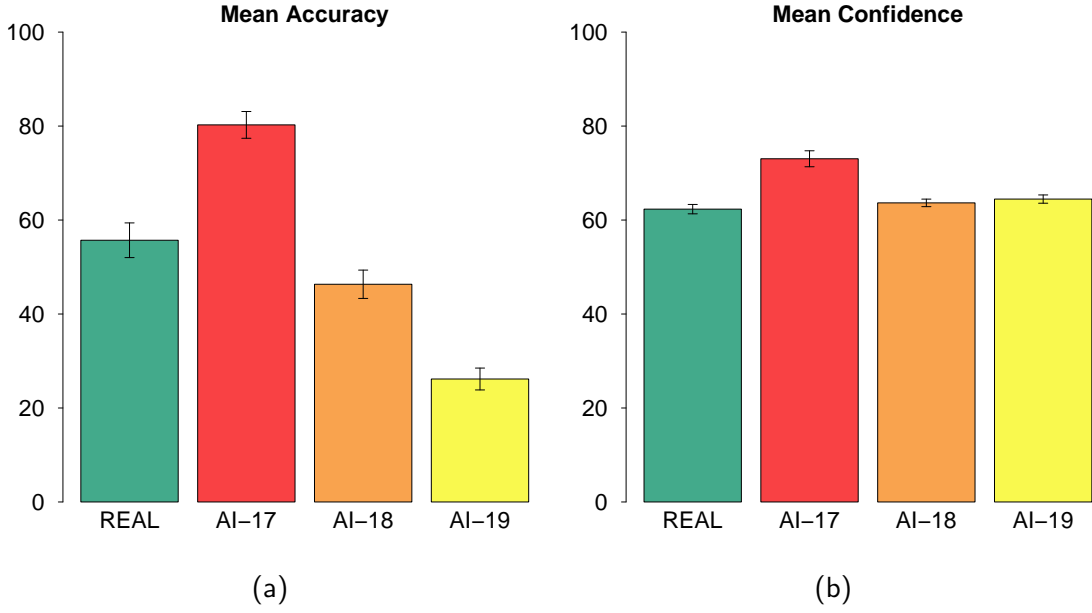


Figure 5: Comparison of the datasets in terms of (a) mean accuracy and (b) mean confidence. The bars show the average, while the error bars indicate the standard deviation.

Accuracy. In relation to the *realism rate*, it is interesting to observe the values of the **accuracy** per image, defined as the frequency of correct responses over all task instances with the given image. The responses are considered correct if they lie between 1 and 3 for real images or between 5 and 7 for synthetic ones. Figure 5a reports the mean values of *accuracy* by dataset, showing that **the accuracy decreases progressively for newer AIs**. The *accuracy* is only slightly above the level of random guessing (56%) for REAL, while it keeps decreasing for the three synthetic datasets, from 80% to as little as 26%.

Confidence. Another indicator is the **confidence**, which is also computed per image. This is obtained by mapping the scale value selected by the participant for a given task instance to the values $[1, \frac{2}{3}, \frac{1}{3}, 0, \frac{1}{3}, \frac{2}{3}, 1]$ and by averaging over image. The mapping serves to make the extremes result in the highest confidence and the central value in the lowest confidence. Figure 5b reports the mean values of the *confidence* by dataset, showing limited variability both across and within datasets, with a maximum of 73% for AI-17 and a minimum of 62% for REAL.

Figure 6 offers a richer visualization of the same indicators. Each dot in the scatter plot represents one image with coordinates defined by its *accuracy* and *confidence*. Color indicates the dataset. Selected examples are also reported visually. From the plot, we can observe a “smile” relationship between *accuracy* and *confidence*. Overall, these two metrics present a weak non-linear correlation

	REAL	AI-17	AI-18	AI-19
p	.006	<.001	.761	<.001
ρ	.222	.813	-.044	-.462

Table 2: Spearman correlation results on accuracy and confidence for each dataset.

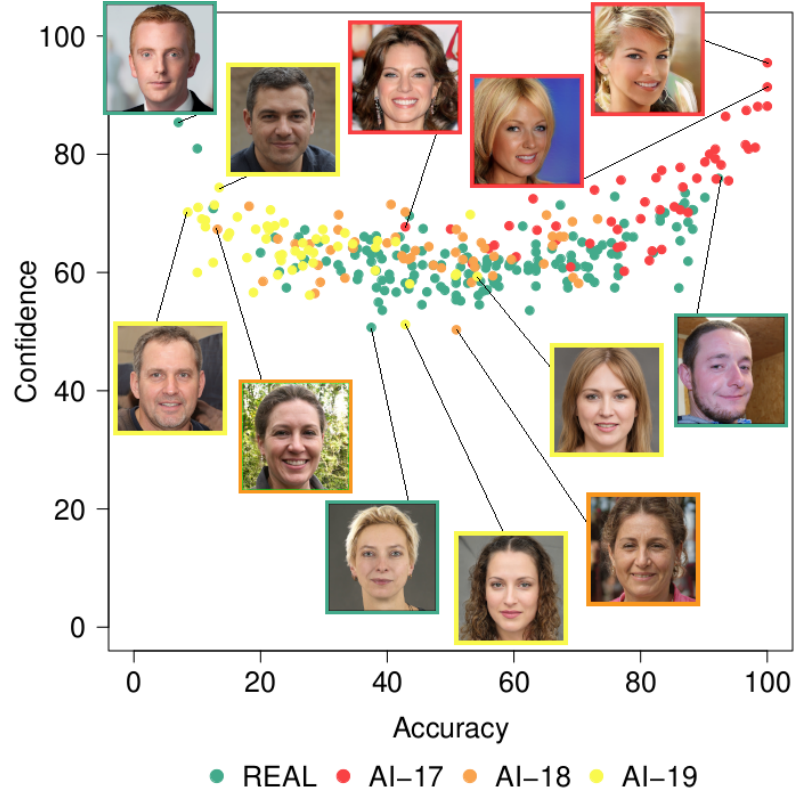


Figure 6: Accuracy (x-axis) and confidence (y-axis) for each image in the stimulus database with selected examples. Color encodes the dataset.

(Spearman $\rho = .228$, $p < .001$). By splitting this analysis by dataset (see Table 2), it is possible to observe how for REAL and AI-18, that in Figure 6 present a higher dispersion of the points, *accuracy* and *confidence* are not correlated. On the other hand, for AI-17 and AI-19's, whose point's are more homogeneously distributed, correlations can be observed. The former has a positive correlation, while the latter present a negative correlation between *accuracy* and *confidence*. In fact, while images in AI-17 are frequently classified correctly with a high confidence, for AI-19, the more a face is inaccurately perceived as real, the higher the observer's confidence.

Indicators of response quality. We computed the group median response time for each instance of the main experiment (i.e., from 1 to 34). We noticed that, **in the course of the experiment, participants tend to speed up their decisions**: for the first task instance, the median response time was 6.0s, while progressively reducing to 4.6s (median) for the last one. This tendency is confirmed by a significant negative correlation between the task index and the median response time (Spearman $\rho = -.968$, $p < .001$). Interestingly, this does not seem to affect the accuracy (Spearman $\rho = -.035$, $p = .843$): **participants neither lose accuracy when speeding up, nor do they learn how to correctly classify images in the course of the experiment**. The latter is quite expected given that participants did not receive any feedback about their performance. In general, this analysis ensures that our experiment was designed correctly and was not contaminated by a speed/accuracy trade-off.

Subject-based analysis. We also analyzed inter-subject variability of the participants by computing the previously defined metrics over the task instances performed by individual subjects.

In general, we found that there is a weak positive correlation between *accuracy* and *confidence* (Spearman $\rho = .201$, $p < .001$), suggesting that participant's confidence is a significant but weak predictor of accuracy.

This result was further investigated by looking for potential causal factors based on the information collected in the questionnaire. Participants have been asked to self-evaluate their face recognition skills and their knowledge on the concept of deep fake, on a scale from 1 to 5.⁹ Through a linear regression analysis, we found a significant effect of self-reported expertise on the participants' accuracy ($p < .001$): **participants who are more familiar with deep fakes tend to classify images more accurately than those who are not**. Such an effect is however not present for the mean confidence ($p = .879$), i.e., **more knowledgeable participants do not seem to provide more confident answers**. Also, self-reported face recognition skills¹⁰ do not significantly influence participants' accuracy nor confidence ($p = .072$ and $p = .100$, respectively).

Participants' strategy. Finally, we analyzed the participants' free-form answers to the request to describe the strategy used to recognize synthetic images. Around 20% of the participants provided an explanation. Figure 7 visualizes the answers as a word cloud.¹¹ It can be seen that the background of the image is very important, similarly to the presence of artifacts in general, which also appear as *blur*, *light*, or *anomaly*. Also, specific face elements such as *symmetry*, *teeth*, *hair*, *eyes*, and *wrinkles* are mentioned often. We also grouped comments according to participants' accuracy ($\geq 70\%$, between 50% and 70%, $\leq 50\%$) to better analyze them.

Artifacts appear to be a key feature for subjects with accuracy ≥ 0.7 , as they were mentioned in 80% of the comments. These well-performing participants also predominantly report their focus on the background of the image (67%) and on the specific face element of the hair (47%). On the contrary, participants with lower performance, and in particular those who performed equal or worse than random guessing, tend to focus less on the background (30% and 14%, for those above or below random chance, respectively) and the hair (26% and 13%, respectively). They concentrate more on the appearance of the face, as roughly 60% of them mention observing either the general structure of the face (symmetry, proportion, expression) or the realism of facial details, such as eyes, teeth, ears, skin, and wrinkles. The poorer performance of participants reporting as strategy the

⁹For the deep fake knowledge the exact question was *How familiar are you with the concept of deep fakes?* 1. You did not know they existed 2. You heard about them but never seen any example 3. You have seen examples but have no intuition of the technology behind them 4. You have seen examples and have some knowledge on the technology behind them 5. You have knowledge on the technology behind them and/or have tried to create some yourself.

¹⁰For the face recognition skills the exact question was *How good do you consider yourself to be at remembering faces?* 1. You tend to forget faces very quickly after seeing them 2. You sometimes mistake people you have met before for strangers 3. You easily recognize faces that you see frequently or occasionally, but you often do not recognize people you have only met briefly before 4. You are better than most people at putting a name to a face 5. You never forget a face.

¹¹Answers reported in different languages have been translated to English through standard web translation tools.

	Mader et al. [11]		AI-17		AI-18		AI-19	
	untrained	trained	all	best	all	best	all	best
Sensitivity d'	1.75	2.04	1.32	1.63	0.02	0.19	-0.68	-0.45
Bias β	2.36	1.45	0.7	0.73	1.05	1.29	3.39*	2.86*

Table 3: SDT analysis in terms of sensitivity d' and bias β reported in [11] – both with and without training – and computed on our data for the three synthetic datasets compared to REAL. For our results, we reported d' and β for all participants (“all”) and a selection of the 25% of participants with the highest confidence (“best”). * indicates that, due to the negative d' , β was computed on the alternative hypothesis (i.e., the signal is the presence of a real face).

of the two populations on the older (although different in nature) kind of manipulations is fairly similar.

A richer analysis can be obtained using Signal Detection Theory (SDT), which is widely used in experimental psychology to model human detection abilities under uncertainty. In particular, the sensitivity index (d') quantifies how hard it is to detect the signal (in this case, the presence of a synthetic face), while the bias (β) measures the extent to which one response is more probable than the other regardless of target presence. In our setting, a value of $\beta > 1$ ($\beta < 1$) indicates that participants are more prone to classify a face as natural (synthetic). The authors of [11] reported that their participants had a clear bias toward classifying an image as real. This bias could only partly be reduced with training (see Table 3).

We decomposed the accuracy computed so far into sensitivity and bias for the three distinct synthetic datasets. In addition, we extracted these metrics from all participants (“all”) and, as a proxy for training, from the subset of the 25% of participants with the highest confidence (denoted as “best”).

As expected, the analysis shows different results among the three technologies. For AI-17 the sensitivity index is rather high (although always lower than in [11]) and the bias indicates that participants are more prone to consider those faces as synthetic. Sensitivity and bias for AI-18 indicate a performance close to random guessing. Finally, the negative sensitivity registered for AI-19, along with the extremely high bias registered for the alternative hypothesis, show to what extent participants are misled into thinking that those faces are real. These results hold consistently for both the “all” and “best” samples.

4 Discussion and future work

This study provides first quantitative evidence on how the quality and realism of face images generated with cutting-edge AIs makes it hard for human viewers to recognize them as synthetic. This trend is rather prominent as, within just two years, we moved from synthetic images that were reliably detected (82% accuracy for AI-17) to synthetic images whose realism rate even surpasses real images (68% for AI-19 versus 52% for REAL). In other words, our participants questioned the authenticity of recent GAN-generated faces less than they did for images of real faces!

The fact that recent AI has jumped this bar so impressively poses a string of follow-up questions. One interesting direction would be to study how these AI-based creation technologies position themselves with regards to the *uncanny valley* effect. This well-known phenomenon in robotics and computer graphics predicts that people's response to human-like avatars quickly shifts from empathy to revulsion as the avatars approach, but do not fully attain, a human appearance [12]. In fact, our results suggest that the recent developments have pushed the synthetic face generation technology away from the uncanny valley, calling for a rise in awareness among users on the existence, threats and opportunities of deep fakes.

On the one hand, there is the need for new automated methods capable to detect synthetic faces in applications where authenticity matters, such as authentication services, the media industry, but also on consumer devices where content is viewed to reduce the dissemination and harmfulness of fake news. However, automatic detectors might not ever work sustainably, given the way GANs work: every improved detector can be used as discriminator and make the next generation GAN evade it altogether.

On the other hand, deep fakes could provide a valuable resource for fields such as cognitive neuroscience, that currently are often limited by the scarcity of extensive sets of controlled face images. The perception of faces by humans is a core research topic in the vast field of the cognitive neuroscience, as a substantial portion of the human visual system is dedicated to perceiving face signals; the cortical territory devoted to the visual perception of other visual categories (e.g., natural scene images) is much less extensive. Human newborns, since the very first minutes of life, preferentially turn their gaze towards faces. In fact, faces emit rich and complex signals that inform on others' intentions, trustworthiness, center of attention, identity, age, emotional expressions, physical health, all cues that are important for an individual's social life quality.

The current results open exciting avenues for the generation of extensive face stimulus sets for human face perception research. As recent research suggests [13], the investigation and manipulation of the most recent deep fakes latent spaces promises to further deepen our understanding of how the human brain perceives faces.

Ethics

The experimental protocol adhered to the Declaration of Helsinki and was approved by the local ethical committee (Psychological Sciences Research Institute, UCLouvain). [no official reference number attributed by this commission]

Acknowledgement

This work was supported by the project PREMIER (PRE-serving Media trustworthiness in the artificial Intelligence ERa), funded by the Italian Ministry of Education, University, and Research (MIUR) within the PRIN 2017 program.

V.G. is a research associate of the National Fund for Scientific Research (F.R.S.-FNRS).

This study has received funding from the FWO and F.R.S.-FNRS under the Excellence of Science (EOS) programme (HUMVISCAT-30991544) awarded to V.G..

References

- [1] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119, 2020.
- [2] Luisa Verdoliva. Media forensics and deepfakes: An overview. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):910–932, 2020.
- [3] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics: A large-scale video dataset for forgery detection in human faces. *arXiv preprint arXiv:1803.09179*, 2018.
- [4] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. In *International Conference on Learning Representations*, 2018.
- [5] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [6] Kate Crookes, Louise Ewing, Ju-dith Gildenhuys, Nadine Kloth, William G Hayward, Matt Oxner, Stephen Pond, and Gillian Rhodes. How well do computer-generated faces tap face expertise? *PloS one*, 10(11):e0141353, 2015.
- [7] Andrey Makrushin, Dennis Siegel, and Jana Dittmann. Simulation of border control in an ongoing web-based experiment for estimating morphing detection performance of humans. In

Proceedings of the 2020 ACM Workshop on Information Hiding and Multimedia Security, pages 91–96, 2020.

- [8] Hany Farid and Mary Bravo. Photorealistic rendering: How realistic is it? *Journal of Vision*, 7(9):766–766, 2007.
- [9] Hany Farid and Mary J Bravo. Perceptual discrimination of computer generated and photographic faces. *Digital Investigation*, 8(3-4):226–235, 2012.
- [10] Olivia Holmes, Martin S Banks, and Hany Farid. Assessing and improving the identification of computer-generated portraits. *ACM Transactions on Applied Perception (TAP)*, 13(2):1–12, 2016.
- [11] Brandon Mader, Martin S Banks, and Hany Farid. Identifying computer-generated portraits: The importance of training and incentives. *Perception*, 46(9):1062–1076, 2017.
- [12] Masahiro Mori, Karl F MacDorman, and Norri Kageki. The uncanny valley [from the field]. *IEEE Robotics & Automation Magazine*, 19(2):98–100, 2012.
- [13] Rufin VanRullen and Leila Reddy. Reconstructing faces from fMRI patterns using deep generative neural networks. *Communications Biology*, 2(1):1–10, 2019.