

Quality assessment of anatomical MRI images from Generative Adversarial Networks: human assessment and image quality metrics

Matthias S. Treder¹, Ryan Codrai¹, and Kamen A. Tsvetanov^{2, 3}

¹ *School of Computer Science and Informatics, Cardiff University, Cardiff CF24 3AA, UK*

² *Department of Clinical Neurosciences, University of Cambridge, CB2 0SZ, UK*

³ *Department of Psychology, University of Cambridge, Cambridge, CB2 3EB, UK*

Abstract

Background: Generative Adversarial Networks (GANs) can synthesize brain images from image or noise input. So far, the gold standard for assessing the quality of the generated images has been human expert ratings. However, due to limitations of human assessment in terms of cost, scalability, and the limited sensitivity of the human eye to more subtle statistical relationships, a more automated approach towards evaluating GANs is required.

New method: We investigated to what extent visual quality can be assessed using image quality metrics and we used group analysis and spatial independent components analysis to verify that the GAN reproduces multivariate statistical relationships found in real data. Reference human data was obtained by recruiting neuroimaging experts to assess real Magnetic Resonance (MR) images and images generated by a GAN. Image quality was manipulated by exporting images at different stages of GAN training.

Results: Experts were sensitive to changes in image quality as evidenced by ratings and reaction times, and the generated images reproduced group effects (age, gender) and spatial correlations moderately well. We also surveyed a number of image quality metrics. **Overall, Fréchet Inception Distance (FID), Maximum Mean Discrepancy (MMD) and Naturalness Image Quality Evaluator (NIQE) showed sensitivity to image quality and good correspondence with the human data, especially for lower-quality images (i.e. images from early stages of GAN training). However, only a Deep Quality Assessment (QA) model trained on human ratings was able to reproduce the subtle**

differences between higher-quality images.

Conclusions: We recommend a combination of group analyses, spatial correlation analyses, and both distortion metrics (FID, MMD, NIQE) and perceptual models (Deep QA) for a comprehensive evaluation and comparison of brain images produced by GANs.

Keywords: Generative Adversarial Network, GAN, MRI, machine learning, generative models, ageing, quality assessment, deep learning

1. Introduction

Magnetic resonance imaging (MRI) provided a means to understanding the structural and functional heterogeneity of the human brain in health and disease. The recent surge in computing horsepower together with large international collaborative initiatives advanced neuroimaging into a big data science, opening the field to deep learning with a promise for new discoveries [1]. One of the most successful deep learning models has been the Convolutional Neural Network (CNN). Using brain images (e.g. T1-weighted or grey-matter density maps) as input, it has been used in brain age regression [2], brain tumor classification [3], tumor segmentation [4], and a plethora of other applications such as image enhancement, image modality translation, and data augmentation (see [5] for a review).

Many of the latter applications build on generative models which take as input either an image (*image-to-image* approach) or a random vector (*noise-to-image*) and produce an artificially generated brain image as output. The two main architectures for generative models are Generative Adversarial Networks (GANs; [6]) and Variational Autoencoders (VAEs; [7]). GANs have been shown to produce more crisp images than VAEs in diffusion-weighted [8] and T1-weighted images [9] and will be our focus in the rest of the paper. In GANs, the generator creates brain images as outputs using transposed convolution layers, the discriminator is a CNN that tries to classify images as generated or real. Both networks act as adversaries towards each other, with the generator using feedback from the discriminator to create increasingly realistic brain images. In the *image-to-image* approach, the generator receives an input image, typically from a different modality or lower resolution, and produces an output image in a different modality or higher resolution. For instance, GANs have been used to translate T1-weighted images into T2-weighted images [10], Computed Tomography [11], functional

29 MRI [12], and diffusion weighted images [13]. Furthermore, [14, 15] recov-
30 ered high-resolution images from images reduced in k-space. A potential
31 limitation of the image-to-image approach is that it requires input images
32 from two imaging domains. In the *noise-to-image* approach, realistic MR
33 images are synthesized *de novo*, starting from a random noise vector [16].
34 Their success rests on the fact that MRIs form a low-dimensional mani-
35 fold and the generator acts as a forward model that maps from this mani-
36 fold to image space. The primary application of noise-to-image GANs has
37 been data augmentation. Several authors successfully used noise-to-image
38 GANs to create realistic T1-weighted 2D image slices [17, 18] or 3D brain
39 images [19, 9, 20]. The noise-to-image approach has also been used in tandem
40 with the image-to-image approach in order to improve tumor detection [21].
41 *Noise-to-image* problems are arguably harder than *image-to-image* problems,
42 since all anatomical features have to be learned from scratch, including the
43 position and shape of the brain volume and its internal 3D structures.

44 While much effort has gone into improving generative models, no compre-
45 hensive framework for assessing the image quality and biological plausibility
46 of the generated images has emerged yet. The current gold standard for
47 judging visual quality is assessments by neuroimaging experts [8, 22]. In [17]
48 two neuroimaging experts scored images on a 5-points Likert scale with real
49 images scoring higher than generated ones. Several studies used a binary
50 detection task wherein the experts were presented either real vs generated
51 images [22, 8, 18]. Results indicate that generated MRIs are able to mislead
52 experts into believing they are real to a high extent, albeit not perfectly so.
53 Participants reported abnormalities in the shape of landmarks, changes in
54 image contrast, or unusually high symmetry between the two hemispheres as
55 giving away whether or not an image was real.

56 However, relying on human assessment alone has significant shortcom-
57 ings. Recruiting human experts is costly, time consuming, and the number
58 of images that can be rated is limited. To avoid this bottleneck and en-
59 able fast development cycles, reliable *image quality metrics* are required that
60 can serve as a proxy for human assessment. They can be readily applied
61 to datasets of any size. In image-to-image applications output images can
62 be compared against ground truth reference images using a simple distance
63 metric (see [16] for an overview). A popular metric is L2 loss and quan-
64 tities derived from it such as mean squared error and peak signal-to-noise
65 ratio (PSNR) have been used for quality assessment [10, 23, 24, 25], as well
66 as Structural Similarity Index Measure (SSIM) [26] and kernel density es-

67 timates [22]. A different approach is to use metrics based on intermediate
68 layers of deep learning models. A widely adopted metric in the GAN litera-
69 ture is the Inception Score (IS). However, in [22] IS did not agree well with
70 human perception, since generated images yielded a higher score than real
71 ones.

72 Additionally, while human assessment and image quality metrics can
73 quantify the extent to which brain images 'look' realistic, we believe that they
74 provide limited information on their *biological plausibility*. Human observers
75 are sensitive to image artifacts such as checkerboard patterns but other dis-
76 tortions of brain images involve subtle statistical relationships across images
77 that may not be visible to the human eye. For instance, the discovery of
78 structural networks such as the Default Mode Network requires multivariate
79 statistical analyses such as spatial Independent Component Analysis (ICA)
80 performed across dozens or hundreds of MRIs [27]. Furthermore, establishing
81 group differences (e.g. young vs old, male vs female) requires group statisti-
82 cal analysis. It is not clear whether the GAN learns to reproduce such group
83 differences and large-scale structural networks because it is not explicitly
84 trained to do so. We therefore believe that biological plausibility should be
85 investigated as an additional dimension of quality assessment using a com-
86 bination of group analysis and spatial ICA. To summarize, the goal of this
87 study was to provide ingredients for a systematic and efficient evaluation of
88 generated brain images. Our contributions are the following:

89 *1. Behavioral experiment.* We conducted the hitherto largest behavioral
90 study on generated MRIs with 26 neuroimaging experts assessing real and
91 generated grey-matter (GM) density maps from a noise-to-image Wasserstein
92 GAN. Participants performed a detection task (real vs generated images) and
93 a subjective quality rating task using a 5-points Likert scale. To determine
94 the experts' sensitivity to objective changes in image quality, we exported
95 images at five different stages (iterations) of GAN training, from early stages
96 (where image quality was supposed to be poor) to later stages. **We hypoth-**
97 **esized that the proportion of images labeled as 'real' increases with training**
98 **iteration in the detection task. Analogously, we hypothesized that quality**
99 **ratings increase with iteration in the rating task. Ultimately both quantities**
100 **were expected to approach the responses participants gave for real images.**

101 *2. Biological plausibility.* For the first time, we performed an analysis
102 of the structural properties of the generated 3D MRIs by performing spatial
103 Independent Component Analysis (ICA) to investigate structural networks
104 and group comparisons that show that GANs are sensitive to gender and age

105 differences. We hypothesized that structural networks and group differences
106 show a high degree of correspondence between real and generated images.

107 *3. Image quality metrics.* We performed a comparison of image quality
108 metrics, ranging from metrics popular in the GAN literature (e.g. Inception
109 Score, Fréchet Inception Distance) to metrics used in image quality assess-
110 ment, some of which have not been used in the context of brain images
111 before. Additionally, we trained a Deep Quality Assessment (QA) model on
112 the human data with the goal to provide an automated metric that mimicks
113 human assessment. Since our ultimate goal was to identify a metric that can
114 serve as proxy for human perception, all metrics were tested for their consis-
115 tency with the behavioral data. We hypothesized that the Deep QA model
116 would outperform standard image quality metrics, especially for high-quality
117 generated images at the end of training.

118 2. Method

119 The Method section is split into four parts. In Section 2.1, we introduce
120 our generative modeling approach. In Section 2.2, we introduce the behav-
121 ioral experiment. In Section 2.3 we introduce spatial ICA and group analysis
122 performed in order to investigate the images’ biological plausibility. In Sec-
123 tion 2.4, we introduce various image quality metrics and a Deep QA as an
124 objective way of measuring the visual quality of the generated images, **as**
125 **well as control analyses using a StyleGAN [28]**. For brevity, we refer to any
126 type of brain image derived from an MR sequence (e.g. T1, T2, grey-matter
127 density maps derived from T1) as *MRI* in the rest of the paper.

128 2.1. Generative modeling of MRIs

129 2.1.1. Data

130 Models were trained using the Cambridge Centre for Ageing and Neu-
131 roscience (Cam-CAN) data [29, 30]. A T1-weighted 3D-structural MRI was
132 acquired with the following parameters: repetition time (TR) = 2,250 ms;
133 echo time (TE) = 2.99 ms; inversion time (TI) = 900 ms; flip angle $\alpha = 9^\circ$;
134 field of view (FOV) = $256 \times 240 \times 192$ mm³; resolution = 1 mm isotropic;
135 accelerated factor = 2; acquisition time, 4 min and 32 s [30]. The T1 image
136 was initially coregistered to the MNI template, and the T2 image was then
137 coregistered to the T1 image using a rigid-body linear transformation. The
138 coregistered T1 and T2 images were used in a multichannel segmentation
139 to extract probabilistic maps of six tissue classes: gray matter (GM), white

140 matter (WM), cerebrospinal fluid (CSF), bone, soft tissue, and residual noise.
141 The native space GM and WM images were submitted to diffeomorphic reg-
142 istration (DARTEL; [31]) to create group template images. Each template
143 was normalized to the MNI template using a 12-parameter affine transfor-
144 mation. After applying the combined normalization parameters (native to
145 group template and group template to MNI template) to each individual
146 participant’s GM images, the normalized images were smoothed using an 8
147 mm Gaussian kernel. In total, GM images from 653 participants, aged 18-88
148 years, were considered as input to noise-to-brain GAN.

149 It is worth noting that GM maps were chosen over more common modal-
150 ities such as T1 since our main focus was on the ability of GANs to generate
151 brain structures. In contrast, non-segmented anatomical MR images also
152 involve signals from non-neural structures such as skull, face, and ears. Con-
153 sequently, it would be impossible to unambiguously attribute human rating
154 and image quality results to brain or non-brain structures. Moreover, spatial
155 brain networks (see Section 2.3) are typically defined using using GM-based
156 morphometry images rather than non-segmented T1 images [32].

157 Before being fed into the GAN, the following additional preprocessing
158 steps were performed. In the first step, we applied a cuboidal crop on the
159 MRIs thereby removing the empty lateral space that surrounds the brain.
160 Due to the four up-scaling layers in the generator that each up-scale by a
161 factor of 2, the generator was restricted to producing images whose length
162 of any given dimension is a multiple of 2^4 . To bring the size of the real
163 images in line with the generated ones, we resized the cropped MRIs to a
164 final resolution of $96 \times 112 \times 96$. The resizing was necessary in order to fit the
165 MRIs in GPU memory, since the intermittent tensors in the GAN were 5D
166 (input channels x output channels x spatial dimensions) and could consume
167 several gigabytes of video RAM each.

168 2.1.2. Generative Adversarial Network

169 In a Generative Adversarial Network (GAN) two neural networks are
170 pitted against each other in a zero-sum game [6]. In this framework, the
171 generator G acts as the adversary of the discriminator D . Whereas G aims
172 to learn the generative distribution of the training data p_{real} by approxim-
173 ating it with a generative distribution p_{gen} , D seeks to accurately predict the
174 probability that an image is generated rather than real. The discriminator
175 is a function $D : \mathbb{R}^{96 \times 112 \times 96} \rightarrow [0, 1]$, $x \mapsto D(x)$ that takes an image as input
176 and returns the probability that this image is real. The generator is a func-

177 tion $G : \mathbb{R}^{100} \rightarrow \mathbb{R}^{96 \times 112 \times 96}$, $z \mapsto G(z)$ that takes a random vector z sampled
 178 from a probability distribution p_z over \mathbb{R}^{100} and returns an MRI. Training
 179 is designed to maximize the quality of the generated images by learning the
 180 generative distribution p_{real} . The model training reaches optimality when D
 181 is unable to discriminate between real and generated samples. The objective
 182 of the discriminator is governed by the equation

$$\max_D V(D) = \underbrace{\mathbb{E}_{x \sim p_{\text{real}}} [\log D(x)]}_{\text{Ability to label real samples}} + \underbrace{\mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]}_{\text{Ability to label generated samples}}. \quad (1)$$

183 D seeks to maximize the probability for samples x from the training data
 184 and minimize the probability for samples generated from a random vector z .
 185 This formulation encourages the discriminator to become sensitive to image
 186 features that tell the difference between real and generated images. The
 187 generator is governed by the equation

$$\min_G V(G) = \underbrace{\mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]}_{\text{Ability to generate quality images}}. \quad (2)$$

188 G seeks to maximize the probability that its generated images are labeled
 189 as 'real' by D . Both objectives can be combined into a joint loss function that
 190 is used for GAN training and optimized using alternating gradient descent

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))] . \quad (3)$$

191 Since G learns a mapping from a low dimensional vector space to a high
 192 dimensional output, the space from which the input vector originates is a rep-
 193 resentation of the low-dimensional manifold of MRIs. This quantity is a lower
 194 bound of the Jensen-Shannon divergence between the real and generated dis-
 195 tributions. GAN training using this objective function can fail to converge
 196 due to near-zero gradients when the true and generated distribution do not
 197 overlap and also. A more stable alternative that provides useful gradients
 198 for non-overlapping distributions is given by Wasserstein GANs (WGANs)
 199 [33, 34]. To this end, the objective function is modified to the Earth Mover
 200 or Wasserstein-1 distance between p_{real} and p_{gen} given by

$$W(p_{\text{real}}, p_{\text{gen}}) = \inf_{\gamma \in \Pi(p_{\text{real}}, p_{\text{gen}})} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|] \quad (4)$$

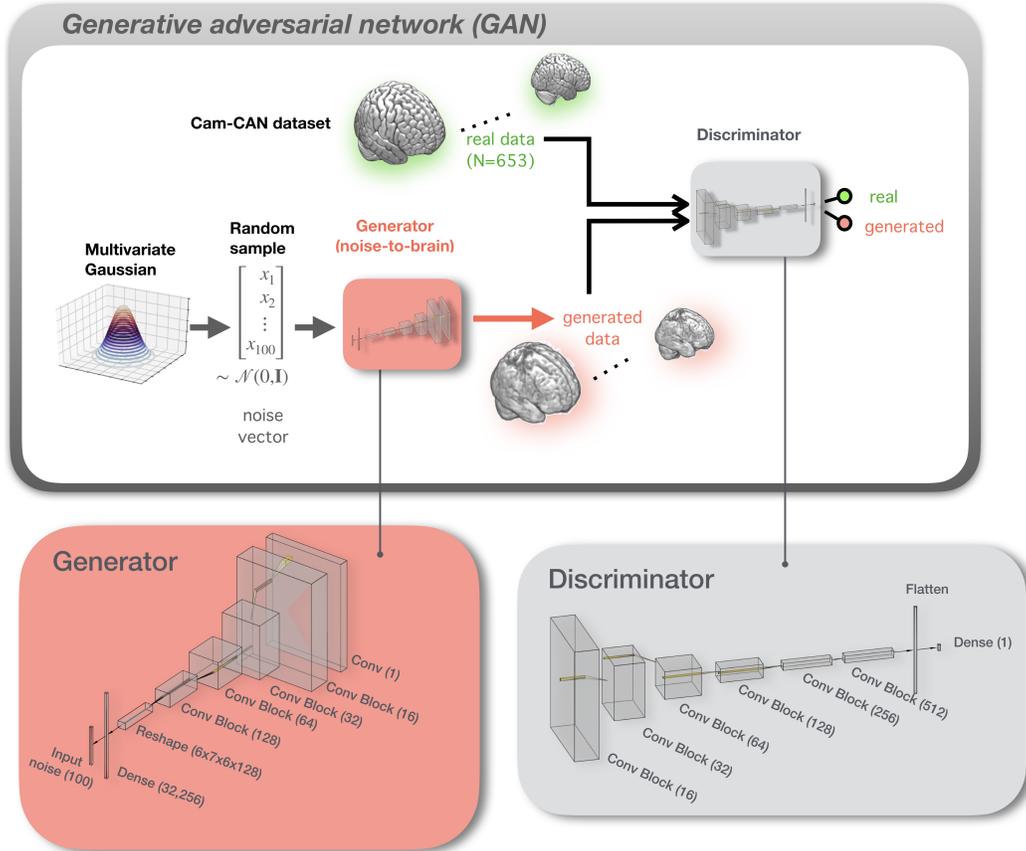


Figure 1: Structure of the noise-to-brain GAN. For the generator and discriminator, the number of filters is given in brackets, whereas for a dense layer, the number of units is shown in brackets.

201 where $\Pi(p_{\text{real}}, p_{\text{gen}})$ is the set of joint distributions whose marginals are
 202 given by p_{real} and p_{gen} . Practical implementations of this loss function are
 203 given in [33, 34].

204 2.1.3. Generating MR images using GANs

205 The structure of the noise-to-brain is shown in Figure 1. Note that the
 206 output of each convolutional layer is 4D (3 spatial dimensions and 1 output
 207 channels dimension). For displaying purposes, it is shown as 3D in the figure.

208 *2.2. Behavioral experiment*

209 *2.2.1. Participants*

210 Forty-five participants volunteered in the online study. Out of these, 26
211 participants (10 female, 16 male) completed all parts of the study (question-
212 naire, detection task, and subjective rating task). Only complete datasets
213 were used for further analysis. Participants were aged 18 to 62 years ($\mu =$
214 $33.12, \sigma = 8.48$), with the majority of participants (16) in their 30s. All
215 participants were neuroimaging experts (e.g. neuroscientists, neurologists,
216 radiologists) at various levels of seniority (from undergraduate student to
217 professor). They did not receive remuneration for their participation. Due
218 to COVID-19 related restrictions (most of the data was collected during UK
219 lockdowns) and the difficulty of recruiting a sufficient number of MR experts
220 locally, the study was conducted online. Participants were mainly recruited
221 through the authors' research networks, although no contact details were
222 recorded. Participants gave written informed consent using an online form.
223 Ethical approval was obtained from the COMSC Research Ethics Group at
224 Cardiff University (COMSC/Ethics/2020/034).

225 *2.2.2. Behavioral experiment*

226 As stimuli, we used middle slides in the sagittal, coronal, and transverse
227 plane extracted from the 3D images. These slices were horizontally concate-
228 nated to 432 x 288 pixels images. An image was either based on the real
229 MRI data or generated by the GAN. To sample the generated images at dif-
230 ferent stages of GAN training, 'fake' images were drawn from five different
231 iterations: 344, 1055, 7954, 24440, and 60000 (final iteration).

232 The experiment was conducted online using PsyToolkit [35, 36] (see Fig-
233 ure S1, Supplementary Material, for a visual depiction). Participants carried
234 it out on their own computer by following a weblink. After filling out a brief
235 questionnaire to indicate their expertise and familiarity with MRI data, the
236 experiment started in fullscreen mode. It took about 15 minutes and was
237 split into three phases: In the training phase, participants performed a prac-
238 tice detection task to familiarize themselves with the stimuli. The next two
239 phases were the detection task and the subjective rating task. The order
240 of these two phases was randomized across participants. Each phase was
241 preceded by a 3 s countdown.

242 Training consisted of 30 randomized trials. In each trial, a real or gen-
243 erated ('fake') MR image was presented at the center of the screen. Partic-
244 ipants were instructed to quickly and accurately indicate their choice with

245 the arrow keys on their keyboard, using the left index finger for the left ar-
246 row key ('real') and the middle finger for the right arrow key ('fake'). The
247 presentation time of the image was unlimited, although for practical reasons
248 the trial ended after 20 s and it was marked as timed out. After the trial,
249 feedback was presented ('correct', 'wrong', 'time out'). Trials were separated
250 by an inter-trial interval of 400-600 ms. Button press and reaction time were
251 recorded. The proportion of stimuli was 1/3 real (10 images) and 2/3 fake
252 (20 images) with equal proportions for each of the five fake levels (4 images
253 each).

254 The detection task was structured just like the training task but no feed-
255 back was presented. There were 240 trials split into three blocks of 80 trials.
256 The proportion of stimuli was 1/3 real (80 images) and 2/3 fake (160 im-
257 ages) with equal proportions for each of the five fake levels (32 images each).
258 Participants were asked to take a break after every block. The experiment
259 resumed upon key press.

260 The subjective rating task consisted of 30 trials comprising 5 real images
261 plus 5 images from each iteration. Participants were rated their visual quality
262 by selecting an item from a 5-points Likert scale with the options "very real",
263 "relatively real", "neutral", "relatively fake", and "very fake". The options
264 were presented below the image and chosen via a mouse click. Presentation
265 time was unlimited although for practical reasons the trial ended after at
266 most 10 s. Across the experiment, all images were unique. Once it was
267 finished, participants had the option to submit a comment in a text box.

268 *2.2.3. Behavioral data processing*

269 For the detection task, the proportion of 'real' responses, proportion of
270 correct responses, and reaction time (RT) were investigated. Trials with
271 timeouts or <150 ms RT were discarded. One participant was removed due
272 to a large amount of timed out trials (>10%). For every participant, mean
273 'real' responses and correct responses were calculated by averaging across
274 all images corresponding to a given iteration. For the RT analysis, mean
275 RTs were calculated using a trimmed mean approach wherein the 5% trials
276 with the largest RTs were excluded and mean RT was calculated across the
277 remaining trials [37]. For the subjective rating task, Mean Opinion Scores
278 (MOS) were calculated by averaging across all images in an iteration. RTs
279 were processed in the same way as for the detection task.

280 *2.3. Group analysis and spatial Independent Components Analysis*

281 To test whether the GAN generates biologically plausible data we com-
282 pared the similarity between generated and real data using data-driven and
283 model-based analyses. For the model-based approach, we sought to repro-
284 duce subtle differences between groups wherein we trained separate GANs on
285 different subsets of the data. Training separate models for different groups
286 is an approach that has been explored before [38]. To investigate age effects,
287 we split the data into elderly (> 70 years; 172 MRIs) and young 287 (< 40
288 years, 185 MRIs). To test for effects of sex, we contrasted male (323 MRIs)
289 and female (330 MRIs) MRIs. Once trained, we randomly generated as many
290 MRIs as were in each of the subsets. The hypothesis was that if the GAN is
291 sensitive to group characteristics, the generated MRIs should display similar
292 age and gender effects as the real MRIs. For the data-driven approach, we
293 used Independent Component Analysis (ICA) to decompose the synthetic
294 data in a small set of spatially independent maps that correspond sensibly
295 with known neurobiological relationships in the real data [32]. In particular,
296 spatial ICA was implemented on real and generated GM maps separately.
297 For each dataset, data were decomposed to a small number of spatially in-
298 dependent sources using the Source-Based Morphometry toolbox [39] in the
299 Group ICA for fMRI Toolbox (GIFT)¹. By combining PCA and ICA, one
300 can decompose a participants-by-voxels into a source matrix that maps in-
301 dependent components (ICs) to voxels (here referred to as “IC maps”), and
302 a mixing matrix that maps ICs to participants. The voxels that carried sim-
303 ilar information across participants would have higher values and group to
304 a set of regions. The spatial IC maps were then converted to z-scores. The
305 correspondence between real and generated MRIs was based on the spatial
306 correlation between thresholded IC maps, where the threshold was set to
307 z-value > 3 . To further explore that the results were independent of the
308 selected number of components, we repeated the procedure for a range of
309 different components (5, 10, 15, 20, 25, 30).

310 *2.4. Image Quality Metrics and Models*

311 *2.4.1. Image quality metrics*

312 In pursuit of an automated metric that can serve as a proxy for human
313 assessment, we surveyed a range of image metrics. For the sake of compara-

¹<http://mialab.mrn.org/software/gift>

314 bility we applied these metrics to the same 2D images that were used in the
 315 behavioral experiment.

316 We first calculated four metrics widely adopted in the GAN literature,
 317 namely Inception Score (IS), Modified Inception Score (MIS), Fréchet Incep-
 318 tion Distance (FID), Maximum Mean Discrepancy (MMD) [40]. **Additional-**
 319 **ly, we tested two reference-free image quality metrics namely Natural Image**
 320 **Quality Evaluator (NIQE) and Blind/Referenceless Image Spatial Quality**
 321 **Evaluator (BRISQUE) [41]. Other popular image quality metrics such as**
 322 **Peak signal to noise ratio (PSNR), Structural Similarity Index (SSIM), and**
 323 **Video Multimethod Assessment Fusion (VMAF), require reference images**
 324 **and thus do not strictly apply in this setting. However, since they are rel-**
 325 **evant to image-to-image GANs, they are investigated in the Supplementary**
 326 **Material.** In the following, each of these metrics is introduced in more detail.
 327 As before, p_{real} and p_{gen} are used to denote the distributions of the real and
 328 generated images. Furthermore, $p(y)$ denotes the distribution of the class
 329 labels and $p(x|y)$ the distribution of images in class y .

- *IS* is a popular GAN metric using a CNN pre-trained on ImageNet [42]. We used the InceptionV3 model trained with 1000 classes. Its final layer is a softmax layer that represents the conditional distribution $p(y|x)$ of class labels y given an input image x . IS is then given by the formula

$$\text{IS} = \exp \mathbb{E}_{x \sim p_{\text{gen}}} (\text{KL}(p(y|x) || p(y)))$$

330 where \mathbb{E} is the expectation, $x \sim p_{\text{gen}}$ are the images sampled from the
 331 generator, KL refers to the Kullback-Leibler divergence, and $p(y)$ is
 332 the marginal distribution of class labels. IS favours generated images
 333 that show a clear class membership, characterized by conditional class
 334 labels $p(y|x)$ with a 'peaky' (low entropy) distribution, at the same
 335 time favouring a broad coverage of multiple classes characterized by a
 336 'flat' (high entropy) marginal distribution $p(y)$. A detailed discussion
 337 of IS is provided in [43].

- *MIS*, unlike IS, takes into account the desire for diversity of images within a given class. To achieve this, [44] incorporated a cross-entropy term $-p(y|x_i) \log p(y|x_j)$ into IS that represents within-class diversity yielding the modified Inception Score

$$\text{MIS} = \exp \mathbb{E}_{x_i} (\mathbb{E}_{x_j} (\text{KL}(p(y|x_i) || p(y|x_j)))).$$

- *FID* uses the feature embeddings of a CNN wherein it models the distributions of real and generated images as multivariate Gaussians [45]. FID is then given as the Fréchet distance between these two Gaussians,

$$\text{FID} = \|\mu_{\text{real}} - \mu_{\text{gen}}\|_2^2 + \text{trace}(\Sigma_{\text{real}} + \Sigma_{\text{gen}} - 2(\Sigma_{\text{real}}\Sigma_{\text{gen}})^{\frac{1}{2}})$$

338 where μ_{real} and μ_{gen} are the means for real and generated images, and
 339 Σ_{real} and Σ_{gen} are their respective covariance matrices. We calculated
 340 FID using InceptionV3’s penultimate layer. The activation maps were
 341 condensed into 2,048 features by applying global average pooling.

- *MMD* measures the distance between two distributions, operationalized as the distance between its mean embeddings in feature space. Gretton et al. [46] implemented MMD using kernels, yielding

$$\text{MMD} = \mathbb{E}_{x, x' \sim p_r} [k(x, x')] + \mathbb{E}_{y, y' \sim p_g} [k(y, y')] - 2\mathbb{E}_{x \sim p_r, y \sim p_g} [k(x, y)]$$

342 where k is a kernel function. We calculated MMD by again using the
 343 penultimate layer of InceptionV3 and the Radial Basis Function (RBF)
 344 kernel $k(x, x') = \exp(-\frac{1}{2}\|x - x'\|^2)$.

- *NIQE* computes image statistics based on normalized luminance values. A multivariate Gaussian distribution is fit to the image statistics and its Mahalanobis distance to a fit obtained on a corpus of natural images is calculated

$$\text{NIQE}(\nu_{\text{real}}, \nu_{\text{gen}}) = \sqrt{(\nu_{\text{real}} - \nu_{\text{gen}})^\top \left(\frac{\Sigma_{\text{real}} + \Sigma_{\text{gen}}}{2} \right)^{-1} (\nu_{\text{real}} - \nu_{\text{gen}})}$$

345 where ν_{real} and Σ_{real} are the mean feature vector and covariance ma-
 346 trix obtained on the corpus and ν_{gen} and Σ_{gen} are the corresponding
 347 quantities computed on the image that is being assessed. We also used
 348 a second version of NIQE referred to as NIQE-MRI wherein the model
 349 was pretrained on an independent set of MRIs. Matlab’s `niqe` and
 350 `fitniqe` functions were used.

- *BRISQUE* is a metric that, in contrast to NIQE, requires reference images with distortions as well as MOS to train a SVR model on the image statistics [47]. Just as NIQE, we also pretrained a second version

354 of the model on MRIs. Since BRISQUE requires MOS it was trained
355 on the subjective rating data and cross-validation was used to obtain
356 unbiased scores for the images. Matlab’s `brisque` and `fitbrisque`
357 functions were used.

358 IS, MIS, NIQE and BRISQUE are reference free metrics that provide
359 scores for both generated and real images. For the other metrics, we used
360 generated images from each of the iterations as target images and the real
361 MRIs as reference images. FID and MMD only require a distribution of
362 generated images. Another distinction between the metrics is that IS and
363 MIS are quality metrics (better quality = higher value) whereas FID, MMD,
364 NIQE and BRISQUE are distance metrics (better quality = lower value).

365 2.4.2. Deep QA model

366 An alternative to image-based metrics is to train a model that tries to
367 mimick human perceptual assessments. Such perceptual models have been
368 used both for MRIs and in the wider image/video quality literature [48, 49].
369 We implemented two CNNs that essentially performed the same task as the
370 human experts, namely classifying images as ‘real’ or ‘fake’ (detection task)
371 and assigning a subjective rating to an image on a 5-points Likert scale
372 (rating task). The model architecture is depicted in Figure 6a.

373 *Detection task model.* The first part of the model consisted of a Incep-
374 tionV3 model pretrained on ImageNet with the classification layer removed.
375 The InceptionV3 output was fed into two dense layers (32 and 16 units,
376 dropout rate 0.1, leaky ReLU activation with $\alpha = 0.1$), denoted as MRI
377 features. The final layer was a single sigmoid unit. The model was trained
378 on a set of 3611 real and generated images not used in the behavioral ex-
379 periment. Images from all iterations were pooled and labeled as ‘fake’. The
380 model was initially trained for 5 epochs with the InceptionV3 model frozen.
381 Subsequently, the last two convolutional layers of InceptionV3 were unfrozen
382 and training continued for another 20 epochs. Binary cross-entropy was used
383 as loss function with an Adam optimizer, a learning rate of 10^{-4} (reduced to
384 10^{-5} after 5 epochs), and a batch size of 16. Finally, the model was tested
385 on the 240 images used in the detection task and the predicted logits were
386 recorded for each image. No behavioral data was used.

387 *Rating task model.* The challenge of training a model on the rating task
388 was the small number of only 30 images. Therefore, we explored a transfer
389 learning approach wherein we first trained the detection task model and then

390 fine-tuned it on the rating task data. To this end, we took the detection task
391 model and replaced the sigmoid by a softmax layer with 5 units represent-
392 ing the Likert scale ratings. The 30 rating task images were split into train
393 and test sets using 5-fold cross-validation. For the training images, class la-
394 bels (representing ratings) could not be assigned unequivocally since there
395 were disagreements between the human raters. To encode this uncertainty in
396 the model we used probabilistic class labels: For each image, the empirical
397 distribution of ratings was determined across raters and used to initialize a
398 probability distribution. In each batch, class labels were randomly sampled
399 from these distributions. The model was initially trained for 20 epochs with
400 the InceptionV3 model and the MRI features frozen. After this, the MRI
401 features were unfrozen and model was trained for another 200 epochs. Cat-
402 egorical cross-entropy was used as loss function with an Adam optimizer,
403 a learning rate of 10^{-4} (reduced to 10^{-5} after 20 epochs), and batch size
404 6. Both models were trained 100 times with randomly initialized weights.
405 Results were averaged across runs.

406 *2.4.3. Hardware and software*

407 GANs were trained on NVIDIA Tesla P100 GPUs with 16 GB VRAM
408 using Supercomputing Wales with TensorFlow 1. Statistical analyses were
409 performed on a standard Desktop computer using Python 3.6 with the pack-
410 ages Statsmodels 0.13, Scipy 1.5.2, and TensorFlow 2.4, as well as Matlab
411 R2018b.

412 *2.4.4. Control analyses using a 3D StyleGAN*

413 To verify that our results are not contingent on a specific choice of a
414 generative model or imaging contrast we performed two control analyses.
415 We used a different generative model called StyleGAN which uses a latent
416 space combined with dense layers to modify activation maps at different
417 stages of the generator [50]. StyleGANs have been successfully applied to 3D
418 structural neuroimaging data [28]. First, to verify that our results generalize
419 to a different model using the same data, we trained the StyleGAN on the
420 GM maps used in our noise-to-image GAN. To this end, we used code from
421 a publicly available repository². The arrangement of 2D slices was very
422 similar but not identical to the arrangement in the main experiment, due

²<https://github.com/sh4174/3DStyleGAN>

423 to differences in the processing pipeline. Second, to verify that our results
424 also hold for a different imaging modality, we trained the model again using
425 T1-weighted images from the same cohort. No behavioral data has been
426 collected for these additional analyses. Therefore our analysis focused on
427 applying the image quality metrics and investigating their consistency across
428 models and imaging contrasts.

429 **3. Results**

430 *3.1. Qualitative analysis of generated MRIs*

431 Figure 2 shows a 3D rendering of the generated images using MRICroGL
432 [51]. To indicate how the images evolve across the different stages of training,
433 Figure 2a shows images for six logarithmically spaced iterations ranging from
434 an early iteration (iteration 11) to the final iteration (iteration 60000). Each
435 of the six renderings is based on the same random vector. The GAN initially
436 outputs a block of noise (iterations 11 and 90) from which a brain is eventually
437 carved out (iterations 344 and 843). After this, GAN training seems to focus
438 on improving details and removing noise voxels outside the brain volume
439 (compare iterations 3,240 and 60000). More incidental evidence for the GAN
440 devoting a lot of computational resources to cancelling noise outside the
441 brain is provided in Figure S3 (Supplementary Material). Figure 2b shows
442 an 'exploded' version of a generated MRI followed by two different random
443 examples with a 3D view and various planar views. The model is able to
444 convincingly mimic the overall shape of the brain and prominent structures
445 such as lateral fissure and the cerebellum. The 'exploded' view indicates that
446 the internal grey-matter structure is reproduced as well. The shortcomings
447 of the GAN become more evident when viewed side by side with a real MRI
448 example in Figure 2c. Here, large-scale structures such as the lateral fissure
449 are more regular and the image looks less noisy overall.

450 *3.2. Correlation analysis and spatial ICA*

451 To ground these observations in a more quantitative analysis, we per-
452 formed correlation analyses, an analysis of age and sex group effects, and
453 a spatial ICA analysis, depicted in Figure 3. Correlations were calculated
454 across all brain voxels between pairs of MRIs. To account for possible effects
455 of small spatial misalignments and high-frequency noise, correlation analyses
456 were repeated for images post-processed with a 2D Gaussian smoothing filter
457 ($\sigma = 3$ voxels).

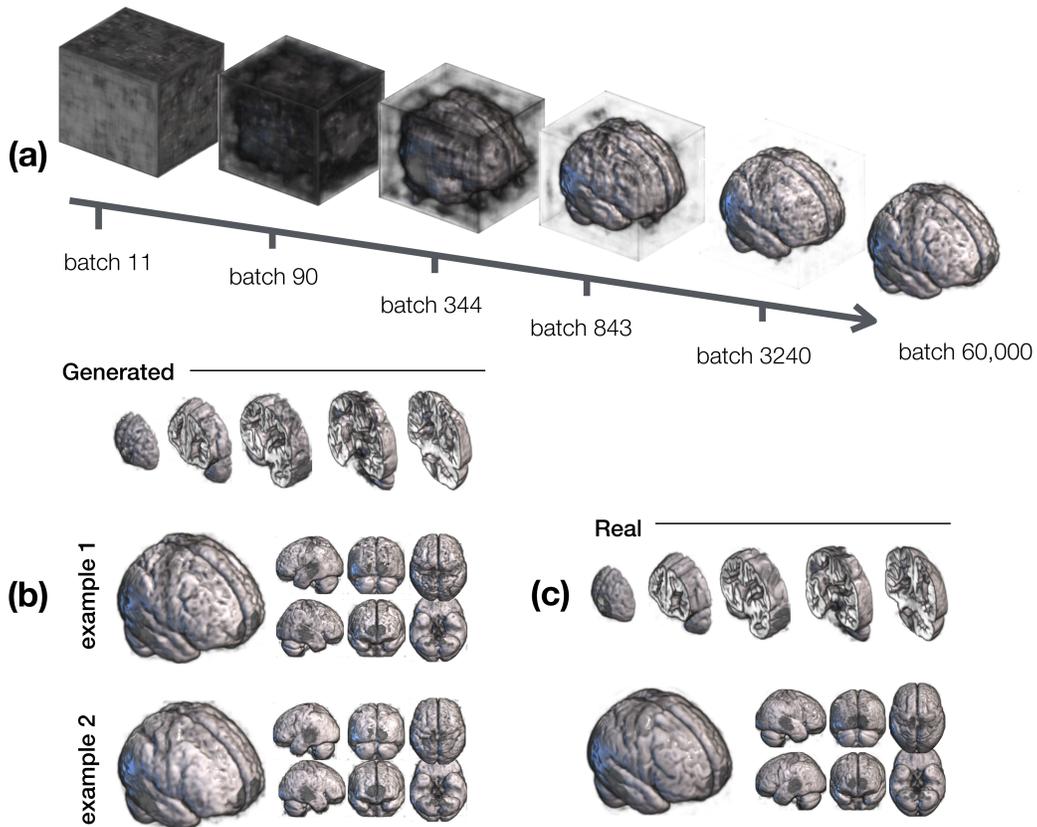


Figure 2: Qualitative results showing a MRIcroGL rendering of the generated 3D images. (a) Generated images at six different times (iterations) during GAN training. Whereas only a volume of noise is generated in the initial stages (iterations 11 and 90), the basic brain structure is evident from iteration 843. Much of the subsequent training appears devoted to refining details and removing the noise outside the brain. Iteration 60000 corresponds to the final model. (b) 'Exploded' view of a generated MRI in the first row indicates that the GAN correctly reproduces the internal structure of the brain, followed by two examples for generated MRIs. (c) For comparative purposes, the same rendering for a real MRI is shown. Comparing real and generated MRI, we observe that the generated image looks slightly more noisy and irregular, especially for large elongated structures such as the lateral fissure.

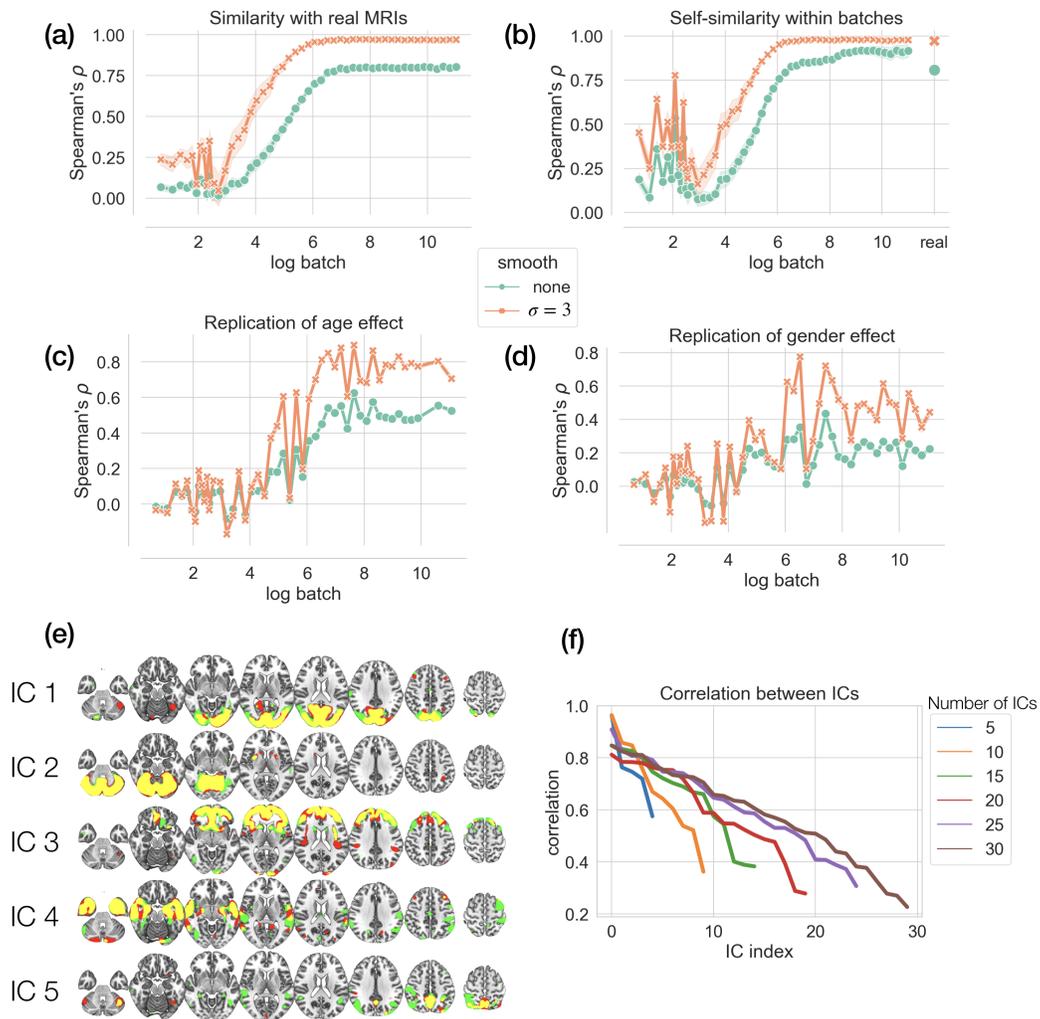


Figure 3: Correlation analyses and spatial ICA. (a) Similarity between generated and real 3D MRIs measured as their average correlation across voxels as a function of log iteration. The shaded area corresponds to 1 standard deviation. Results shown for original data without smoothing (green) and smoothed data (Gaussian kernel, $\sigma = 3$, orange). (b) Within-iteration similarity, average correlation between pairs of images within an iteration. Correlation within real MRIs is shown as an additional data point on the right. (c) Replication of the age effect ('old minus young'). (d) Replication of sex effect ('male minus female'). Correlation starts around zero and increases significantly with training, settling at a significant but moderately low value of 0.17. Smoothing significantly boosts correlation (0.42). (e) Spatial ICA maps for ICA with 5 components. Components for both real and generated MRIs have been matched and overlaid. Color coding indicates whether a voxel belongs to the real (red) or generated (green) MRIs, with the intersection shown in yellow. The large overlap suggests significant spatial correspondence between real and generated MRIs. (f) Correlations between real and generated ICs for ICAs with different numbers of components (between 5 and 30). High degree of correspondence for the first few components is consistently found.

458 Figure 3a depicts Spearman’s rank correlation ρ between generated and
459 real MRIs. Every generated image was paired and correlated with every real
460 image, then mean and standard deviation were calculated across all pairs.
461 Correlation increased significantly then saturated relatively early in training.
462 Since we masked out non-brain voxels, this early saturation nicely dovetails
463 with our earlier observation that the late stages of GAN training seem largely
464 devoted to removing noise outside the brain. The correlation for the final
465 model was 0.8. Smoothing with a Gaussian kernel with a standard deviation
466 of $\sigma = 3$ voxels yielded a correlation of 0.97, a statistically significant increase
467 (Wilcoxon Rank Sum test, $z = 799.76, p < 0.0001$). The significant effect of
468 smoothing suggests a high degree of correspondence between the large-scale
469 structure of generated and real MRIs.

470 Figure 3b depicts the correlation among MRIs within iterations. Corre-
471 lation among real MRIs has been added as additional data points (separate
472 rightmost data points). Correlation started low and quickly saturated at a
473 high level. For the final model the correlation amounted to 0.92, whereas
474 it was 0.98 with smoothing, a significant increase ($z = 563.47, p < 0.0001$).
475 Corresponding correlations for real MRIs were 0.81 without and 0.97 with
476 smoothing, again a significant difference ($z = 565.08, p < 0.0001$). Directly
477 comparing generated and real MRIs, correlation was significantly higher for
478 generated MRIs both for unsmoothed ($z = 549.74, p < 0.0001$) and smoothed
479 ($z = 350.98, p < 0.0001$) data.

480 Figure 3c depicts how well the GAN replicates the age effect found in
481 real data. In real data, the age effect was operationalized as the average
482 MRI in the elderly group (> 70 years) minus the average MRI in the young
483 group (< 40 years). To generate MRIs specifically from these groups, two
484 different GANs were trained, one on the elderly MRIs and one on the young
485 ones. Since the age effect was a group effect, correlations for individual MRIs
486 were not available. Therefore, to calculate errorbars, 100 different bootstrap
487 samples of the GAN data were created. Correlation between real age effects
488 and the generated age effects increased with iteration, suggesting that the
489 GAN was sensitive to these group characteristics. In the final iteration,
490 correlation was 0.52 for the unsmoothed and 0.71 for the smoothed data, a
491 significant difference ($z = 12.22, p < 0.0001$).

492 Figure 3d depicts the replication of a gender effect, operationalized as the
493 average ‘male MRI’ minus the average ‘female MRI’. Again, two different
494 GANs were trained on male and female MRIs separately and bootstrapping
495 was used to calculate errorbars. Correlation increased with iteration, again

496 suggesting that the GAN is sensitive to these group characteristics. However,
497 it settled at a lower value than for the age effect. In the final iteration,
498 correlation was 0.22 for the unsmoothed and 0.44 for the smoothed data, a
499 significant difference ($z = 12.22, p < 0.0001$).

500 Figure 3e depicts 5 ICs for a spatial ICA with 5 components. ICs in real
501 and generated MRIs were matched by their correlation. To binarize the maps,
502 one-sample t-values were calculated for all participants and thresholded at a
503 t-value of 3. A large amount of overlap (shown in yellow) is evident especially
504 for the first few ICs. Non-overlapping voxels are shown in red for real and
505 green for generated data. Figure 3f depicts the correlations between IC maps
506 for different spatial ICAs using a different number of components (5, 10, 15,
507 20, 25, 30). For all analyses, we found correlation values > 0.57 for the
508 first five components, suggesting a moderate to high degree of consistency of
509 components for real and generated MRIs.

510 3.3. Behavioral results

511 Figure 4 depicts the behavioral results. For both tasks, reaction times
512 (RTs) increased with iteration, approaching the RT for real images (Figure 4a
513 and d). For the detection task, the percentage of 'real' responses increased
514 with the iteration, although it stayed short of the corresponding number for
515 real MRIs (Figure 4b). For completeness, the correct score is also depicted
516 in Figure 4c. For the subjective rating task, Mean Opinion Scores (MOS)
517 are shown, that is, averages when ratings are encoded as integer numbers 1
518 through 5. MOS increased with iteration but did not reach the MOS obtained
519 for real MRIs.

520 To substantiate these observations statistically, we used linear mixed-
521 effects models (LMMs) with a fixed effect (log iteration) and a random
522 effect (participant). LMMs are useful in the case of correlated responses
523 arising from multiple measurements per participant and are more flexible
524 than repeated-measures Analysis of Variance [52]. LMMs deal with pseudo-
525 replication and account for the fact that participants display individual dif-
526 ferences in overall reaction time and accuracy.

527 In the detection task, RT increased significantly with log iteration as ev-
528 idenced by the positive regression slope (LMM, $\beta = 224.705, z = 15.837, p <$
529 0.001). Analogously, there was significant increase in 'real' responses with
530 log iteration ($\beta = 6.184, z = 9.343, p < 0.001$) and a significant decrease
531 in correct responses ($\beta = -6.195, z = -9.405, p < 0.001$). We then used
532 Wilcoxon Signed-Rank tests to compare RTs and performance for iteration

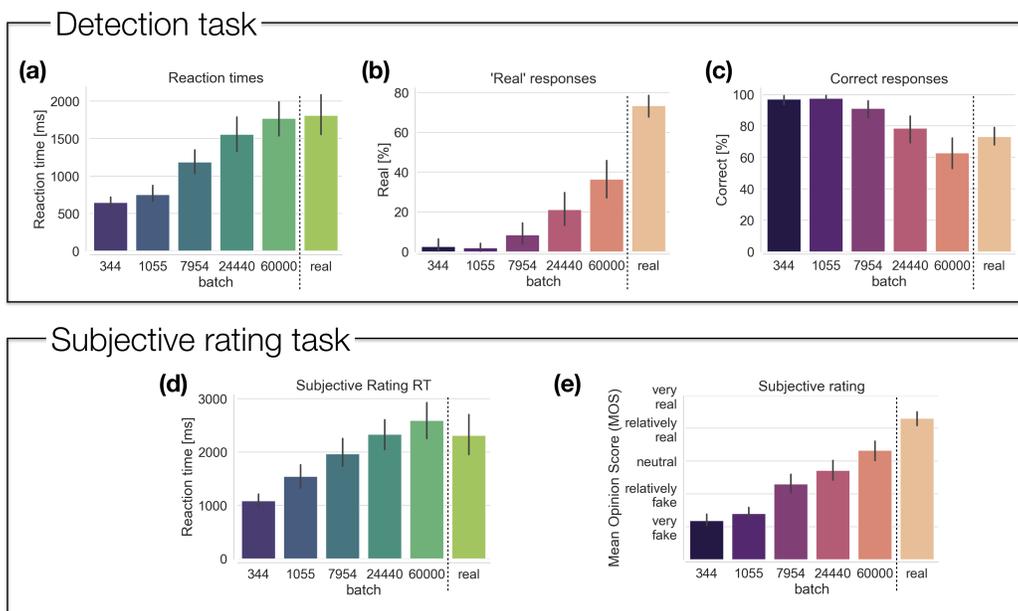


Figure 4: Behavioral experiment results for the detection task and subjective rating task. Results are shown for generated MRIs from different training iterations (indicated by the iterations) and for real MRI data. (a) In the detection task, reaction times (RTs) increased with iteration, approaching the RT obtained for real MRIs. (b) Percentage 'real' responses increased with iteration but stayed short of the proportion obtained for real data. (c) Percentage correct responses decreased with iteration, indicating that it becomes harder to distinguish real and generated MRIs. (d) In the subjective rating task, RTs increased with iteration. For iteration 60000 RT slightly overshoot the RT found for real MRI. (e) Subjective ratings, here summarized as Mean Opinion Scores, increased with iteration albeit staying below the rating for real MRIs for iteration 60000.

533 60000 vs real MRIs across participants. No difference was found for RTs
534 ($w = 146, p = 0.672$) and correct responses ($w = 119, p = 0.252$) but the
535 proportion 'real' was significantly higher for real MRIs than for generated
536 ones ($w = 0, p < 0.0001$).

537 In the subjective rating task, RT increased significantly with log itera-
538 tion ($\beta = 279.875, z = 11.377, p < 0.001$). For the analysis of the ratings,
539 a linear mixed model approach would not be not sufficient since the tar-
540 get variable was ordinal and the 'distances' between neighbouring categories
541 were unknown [53]. The Mean Opinion Scores displayed in Figure 4e were
542 used for illustrative purposes. In line with the recommendations in [53] we
543 used a multi-level proportional odds model with a logit link function which
544 simultaneously deals with mixed effects and ordinal responses. To this end,
545 we fit a Cumulative Link Mixed Model with Laplace approximation [54],
546 wherein rate was used as the ordinal target variable, log iteration as fixed
547 effect and participant as a random effect with random intercept. We found
548 a significant positive slope for log iteration ($\beta = 1.234, z = 17.94, p < 0.001$)
549 signifying higher ratings for later iterations. The same statistical approach
550 was applied to compare ratings for iteration 60000 vs real MRIs, showing
551 significantly higher ratings for real MRIs as compared to generated ones
552 ($\beta = 2.133, z = 7.748, p < 0.0001$).

553 3.4. Image quality metrics and Deep QA model

554 Figure 5 shows the metrics calculated for the detection task images. An
555 analogous analysis applied to the rating task images **and additional full refer-**
556 **ence metrics are** depicted in Figure S2 (Supplementary Material). Figure 6b
557 shows the corresponding result for the Deep QA model. Linear regression
558 analyses of log iteration on the metric were conducted to investigate whether
559 the metric increases/decreases with iteration.

560 For the detection data, all regression slopes were significant (all p-values
561 < 0.0001). IS and MIS both decreased with iteration. Since they are qual-
562 ity metrics, the opposite pattern was expected. FID and MMD both de-
563 creased with iteration but images at iteration 24400 had a lower value than
564 images at iteration 60000 (Wilcoxon Rank Sum, FID $z = 12.217, p < 0.001$;
565 MMD $z = 12.217, p < 0.001$), conflicting with the behavioral data. Although
566 NIQE, NIQE-MRI and BRISQUE-MRI showed overall regression trends in
567 line with behavioral data, the metrics did not change monotonically with it-
568 eration. For instance, NIQE decreased from iteration 344 to 24440 but then
569 increased again for iteration 60000. BRISQUE showed an inverted U-shaped

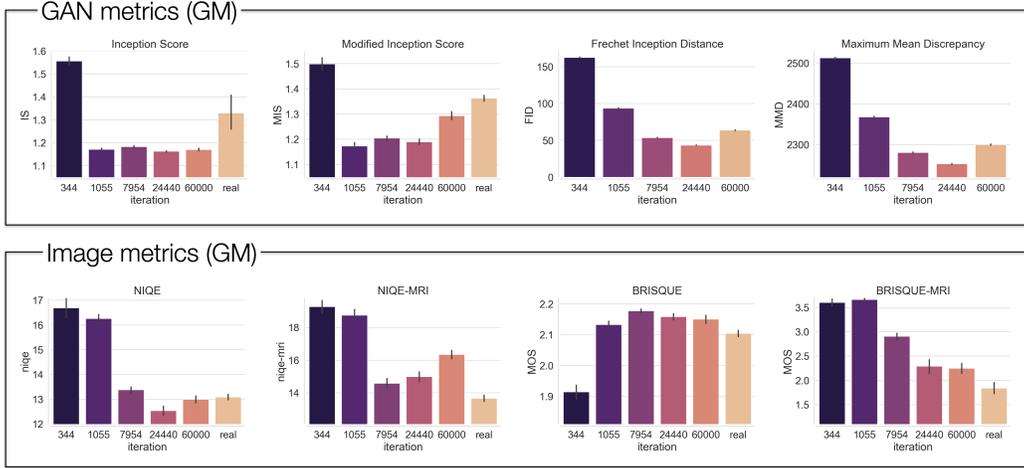


Figure 5: Image quality metrics applied to the 240 images used in the detection experiment. Different GAN and image quality metrics are shown, applied to generated images from different iterations. If the metric does not require a reference, the value for real images is shown as well.

570 pattern. Contrary to expectations, iteration 344 had the lowest score making
 571 it the best in terms of image quality. The Deep QA model was the only
 572 model for which predictions on the detection data changed monotonically
 573 with iteration, although the difference between iterations 24400 and 60000
 574 was not significant ($p = 0.51$).

575 We repeated the same analysis on the rating task data (Figure S2 and
 576 Figure 6b). All regression slopes were significant ($p < 0.001$) except for IS
 577 ($p = 0.44$). Deep QA was the only metric that showed a significant difference
 578 between iterations 24400 and 60000 ($z = -2.193, p = 0.0282$).

579 To statistically compare the image quality metrics to the behavioral data,
 580 we performed Spearman’s rank correlation analyses for both tasks, shown in
 581 Table 1. For the detection task, we used the proportion ‘real’ responses for
 582 each of the 240 images, averaged across participants, and correlated them
 583 with the metrics obtained on the same data (Figure 5 and Figure 6b). **FID**
 584 **and MMD estimate the distance between distributions and hence do not provide**
 585 **values for individual images. To nevertheless estimate correlation, we**
 586 **used the value obtained for each iteration for all images belonging to this**
 587 **iteration.** Furthermore, since some metrics did not provide separate predic-
 588 tions for real images, we conducted two correlation analyses, one analysis
 589 that covered only generated images and another one that included real im-

590 ages (denoted as 'include real' in the table). For the rating task, we used
 591 the Mean Opinion Scores for each of the 30 images, averaged across partici-
 592 pants, and correlated them with the metrics obtained on the rating task data
 593 (Figure 6c and Figure S2).

594 Except for IS and MIS all metrics showed moderate to high correlation
 595 with the behavioral data. For the detection task, FID ($\rho = -0.602$), MMD
 596 ($\rho = -0.602$), NIQE ($\rho = -0.706$), and BRISQUE-MRI ($\rho = -0.859$) were
 597 the best performing metrics when the real data was excluded. When real
 598 images were included, the best performing metrics were NIQE ($\rho = -0.519$),
 599 NIQE-MRI ($\rho = -0.761$), BRISQUE-MRI ($\rho = -0.858$), and Deep QA ($\rho =$
 600 0.831). Despite the fact that the Deep QA model performed the classification
 601 task better than humans (100% accuracy), it assigned higher probabilities to
 602 later iterations which led to a ranking consistent with humans (Figure 6b).
 603 However, when real images were removed from the correlation analysis, the
 604 correlation for the Deep QA model plummeted to $\rho = 0.516$, indicating that
 605 the effect was largely driven by the difference between real and generated
 606 images.

607 For the rating task, NIQE ($\rho = -0.767$), NIQE-MRI ($\rho = -0.575$),
 608 BRISQUE-MRI ($\rho = -0.840$), and Deep QA ($\rho = -0.879$) were the best
 609 performing metrics when the real data was excluded. The same qualitative
 610 pattern was found when real images were included (NIQE $\rho = -0.614$, NIQE-
 611 MRI $\rho = -0.722$, BRISQUE-MRI $\rho = -0.879$, and Deep QA $\rho = -0.912$).

Table 1: Correlation between behavioral data and image quality metrics reported as Spearman’s ρ and corresponding p-value, for the two behavioral tasks separately. Each analysis was performed twice, one time using only the data on the generated images, and a second time including the real data, too. For each analysis, the four metrics with the largest absolute correlation are highlighted in bold.

Task		IS	MIS	FID	MMD	NIQE	NIQE-MRI	BRISQUE	BRISQUE-MRI	DeepQA
Detection	ρ	-0.255	0.11	-0.602	-0.602	-0.706	-0.536	0.465	-0.859	0.516
	p-value	0.001	0.11	0.000	0.000	0.000	0.000	0.000	0.000	0.000
Detection (include real)	ρ	0.008	0.344	-	-	-0.519	-0.761	-0.013	-0.858	0.831
	p-value	0.9	0.000	-	-	0.000	0.000	0.840	0.000	0.000
Rating	ρ	-0.092	0.367	-0.422	-0.422	-0.767	-0.575	0.521	-0.840	0.879
	p-value	0.66	0.071	0.035	0.035	0.000	0.003	0.008	0.000	0.000
Rating (include real)	ρ	0.003	0.232	-	-	-0.614	-0.722	0.294	-0.879	0.912
	p-value	0.98	0.213	-	-	0.000	0.000	0.114	0.000	0.000

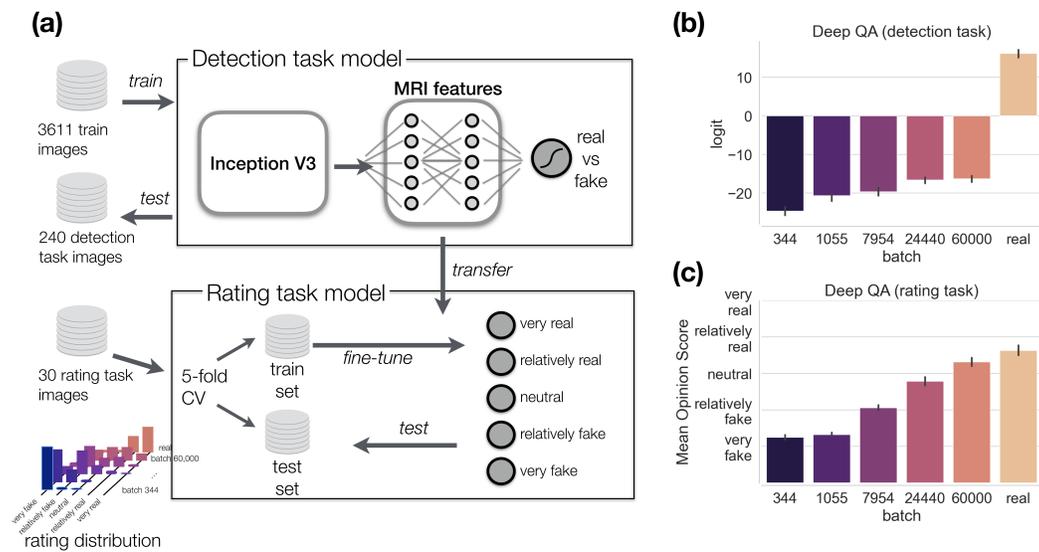


Figure 6: Deep QA model. (a) Model architecture. A pretrained InceptionV3 model was used followed by two dense layers (MRI features) trained on independent data. It was then fine-tuned on the rating data using cross-validation. (b) Logits of the predicted probabilities on the detection task images. (c) Mean Opinion Scores on the rating task images. Means and confidence intervals were calculated across 100 runs.

612 3.5. Control analyses

613 Figure 7 shows the metrics calculated for the GM and T1 images us-
614 ing a StyleGAN. Linear regression analyses of log iteration on the metric
615 were conducted on both GM and T1 to investigate whether the metric in-
616 creases/decreases with iteration. For IS, the regression slope was not signif-
617 icantly different from zero (GM $p = 0.12$; T1 $p = 0.13$). For MIS, results
618 were inconsistent, with a small negative slope for GM data ($\beta = -0.0685, p <$
619 0.0001) but a positive slope for T1 ($\beta = 0.0066, p = 0.03$). For both FID and
620 MMD, regression slopes were significantly negative for both GM and T1 (all
621 $p < 0.0001$). Furthermore, these two metrics were the only to have a strictly
622 monotonic relationship with iteration. NIQE, NIQE-MRI and BRISQUE
623 decreased significantly with iteration for both GM and T1 (all $p < 0.0001$),
624 although the only metric that did so strictly monotonically was BRISQUE
625 applied on GM data. The BRISQUE-MRI transfer used the BRISQUE model
626 trained on the GM data used in the behavioral experiment. The model did
627 not transfer well to T1. It did not even yield consistent results on the GM
628 data used in the StyleGAN when compared to the GM data used in the
629 noise-to-image GAN. Possibly, this was due to differences in the way the
630 slices were placed next to each other to produce a 2D image. The Deep QA
631 model was tested in an out-of-distribution setting on the StyleGAN images.
632 For GM, it yielded higher scores for real images than for early iterations, but
633 iterations 128 and 200 were ranked as higher than real images. For T1 data,
634 it yielded a constant output.

635 3.6. Data sharing

636 All scripts, real and generated images, experimental data, behavioral re-
637 sults and metrics, and the Deep QA model are available at github.com/treder/MRI-GAN-QA.

638 4. Discussion

639 Human assessment is the current gold standard for judging image quality
640 of images generated by GANs [8, 22]. However, the utility of human assess-
641 ment is limited by economical factors, limited scalability of human raters to
642 large datasets, and the limited sensitivity of the human eye to subtle statis-
643 tical relationships across dozens or hundreds of images. By collecting both
644 behavioral data, performing group analysis and spatial ICA, and surveying
645 a whole range of image quality metrics, our goal was to identify a more au-
646 tomated approach for the assessment of generated images. These points will



Figure 7: Image quality metrics applied to the StyleGAN images for (a) GM and (b) T1-weighted images, and separately for (c) Deep QA model applied to the GM and T1 images. Iterations given in units of thousands.

647 be expanded on in the next two subsections, followed by a broader discussion
648 of our findings, the limits of GANs for image generation, and the ideal image
649 quality metric.

650 *4.1. Biological plausibility: group analysis and spatial ICA*

651 Prior to performing spatial ICA, we investigated the degree of correspon-
652 dence between real and generated MRIs at a voxel level. The correlation
653 between generated and real images was 0.80 (Figure 3a), which was very
654 similar to the correlation of 0.81 of real images with themselves (Figure 3b).
655 However, the correlation of generated images with each other was 0.92. This
656 suggests that while generated MRIs adequately approximate real MRIs, the
657 latter have a larger amount of diversity, in line with a previous research [8].
658 One explanation for this is that the GAN lacks *capacity*, that is, it is simply
659 unable to replicate the diversity found in real data because the data does
660 not lie within the range of possible outputs of the generator. Alternatively,
661 the GAN might simply have failed to learn the generative distribution of the
662 MRIs, a possibility discussed in Section 4.3.

663 We then performed group analysis to investigate whether the GAN re-
664 produces group differences. This is not a trivial question since the GAN is
665 only trained to produce images that 'look like brains'. It is not trained to re-
666 produce statistical regularities and the discriminator that guides the training
667 process does not have access to all images simultaneously to facilitate such a
668 learning. Analysis of group differences considered training separate GANs to
669 characterize chronological age (young vs old) and sex (male vs female). The
670 groups of generated images were then subjected to group analysis. We found
671 that the age effect is fairly well reproduced (Figure 3c) while the replication
672 of the gender effect was rather low (Figure 3d). The low reproduction of the
673 gender effect might be due to the fact that gender differences contribute less
674 variance to the grey matter signal than aging which involves global atrophy.

675 Furthermore, we performed spatial ICA to investigate whether the GAN
676 reproduces structural brain networks known to support cognitive function.
677 Again, the model is not explicitly trained to perform this. Performing spatial
678 ICA on real and generated images separately and comparing the results we
679 found a large degree of correspondence between the spatial maps (Figure 3e).
680 We found spatial patterns reminiscent of the networks identified in other
681 studies of ageing with source-based morphometry [55, 32]. This correlation
682 between the spatial maps was high irrespective of the exact choice of the
683 number of ICs.

684 *4.2. Image quality metrics as a proxy for human assessment*

685 To identify whether image quality metrics correspond well with human
686 assessment, we collected human data in two different tasks (detection and
687 rating) and calculated a series of image quality metrics for the same images.
688 For the human data, we found a significant increase in 'real' responses with
689 iteration. Since better performance was associated with lower RT, this could
690 not be explained in terms of a speed-accuracy trade-off. A similar pattern
691 of results was found in the subjective rating task, with Mean Opinion Scores
692 (MOS) increasing significantly with iteration and approaching the ratings for
693 real images. As hypothesized, these results confirmed that changes in objec-
694 tive image quality operationalized by iteration are perceptually relevant and
695 measurable. Behavioral measures converged towards the responses obtained
696 for real images, but there remained a gap, indicating that our model was not
697 consistently able to fool human experts.

698 To relate behavioral data to the image quality metrics, we calculated
699 Spearman's rank correlation between quality metrics and the proportion of
700 'real' responses in the behavioral data. The standard GAN metrics IS and
701 MIS showed poor correspondence with human assessment, in line with ear-
702 lier research [22]. Both FID and MMD showed a good correspondence with
703 behavior in the detection task (both $\rho = -0.602$). In the rating task, the
704 correlation decreased (both $\rho = -0.422$) but this can possibly explained by
705 the fact that there were only 5 images in each iteration. Since FID and
706 MMD are distribution-level metrics, the estimation of distribution paramet-
707 ers could be noisy or unreliable with so few examples. The control analyses
708 using StyleGAN gave further credence to FID and MMD as useful image
709 quality metrics for MRIs. Both metrics were the only to decrease strictly
710 monotonically with iteration.

711 Unlike FID and MMD, NIQE and BRISQUE do not require a reference
712 and provide ratings for individual images. NIQE showed a moderate correla-
713 tion with behavior in the detection task ($\rho = -0.706$, dropping to $\rho = -0.519$
714 when including real images) and the rating task ($\rho = -0.767$, $\rho = -0.614$
715 when including real images). Upon closer inspection, the metric decreased
716 monotonically until iteration 24440, yielding larger values for iteration 60000
717 and real images. A roughly similar pattern for was found for the StyleGAN
718 T1 images. A pattern contradicting this was found for the StyleGAN GM
719 images, with the real images receiving a higher score than all generated im-
720 ages. NIQE can be adapted to the target domain in an unsupervised way.
721 The resultant NIQE-MRI also showed a moderate correlation with behavior

722 (detection task $\rho = -0.536$, $\rho = -0.761$ including real images; rating task
723 $\rho = -0.575$, $\rho = -0.722$ including real images). NIQE-MRI also showed a
724 slightly more consistent pattern for StyleGAN GM and T1 images, with the
725 lowest score given to real images for all but one iteration and generally higher
726 scores given to earlier iterations. Overall, NIQE-MRI provided an accurate
727 tripartite pattern, which high values for early iterations, medium values for
728 later iterations, and lower values for real images. However, both NIQE and
729 NIQE-MRI fail to consistently reproduce monotonically decreasing values for
730 later iterations.

731 BRISQUE showed an inconsistent, inverted U-shaped pattern in the main
732 experiment. Correlation with behavior was moderate (detection task $\rho =$
733 0.465 , rating task $\rho = 0.521$; correlations not significant when including real
734 images) but it was positive whereas a negative correlation was expected. In
735 the StyleGAN, it provided a nearly monotonic pattern for GM images but
736 for not for T1 images. BRISQUE can be adapted to the target domain in
737 a supervised way by training on target MOS. The resultant BRISQUE-MRI
738 showed an adequate monotonic pattern which however did not transfer well
739 to the StyleGAN images, suggesting possible overfitting to the images from
740 the main experiment. We conclude that neither BRISQUE nor BRISQUE-
741 MRI are reliable image quality metrics for the MRI domain.

742 The two Deep QA models showed a moderate correlation with human
743 data in the detection task ($\rho = 0.516$, rising to $\rho = 0.831$ when including
744 real images) and a high correlation in the rating task ($\rho = 0.879$ and $\rho =$
745 0.912 when including real images). The Deep QA models were the only
746 metrics that changed monotonically with iteration. It is worth noting that
747 only the rating model was explicitly trained on human responses, whereas
748 the detection model was trained to discriminate between real and generated
749 images. Crucially, the Deep QA rating model was the only metric that
750 showed a significant increase in quality ratings from the 24400 to the 60000
751 iteration.

752 However, the Deep QA rating model did not readily transfer to the Style-
753 GAN images, with a near constant output for most iterations in both GM and
754 T1 images (Figure 7c). This suggests that the model overfit to the training
755 dataset. To develop a more versatile model that we envision several possible
756 avenues. Firstly, training the model using data augmentation (e.g. skewing
757 and rotations of the images) might increase generalizability. Secondly, our
758 data could be combined with ratings obtained in other studies. A single
759 model could then be trained to predict ratings on different datasets simulta-

760 neously. To achieve this, the model would need to learn some generalization
761 capabilities and hence be more likely to transfer to new datasets. To facili-
762 tate this, our human ratings, corresponding images and training scripts have
763 been published in a GitHub repository³. Thirdly, a semi-supervised learning
764 approach could be considered wherein a labeled dataset is combined with an
765 additional unlabeled dataset (e.g., [56]). For instance, the Deep QA could be
766 trained to simultaneously predict MOS on our data and predict whether im-
767 ages in the unlabeled dataset are real or fake by adding a separate real/fake
768 output unit to the model. At inference time, this output unit is dropped and
769 ratings are obtained via the other units. Exploring these avenues in more
770 detail is left for future work.

771 Across many of the surveyed metrics, we found that they correlate with
772 human behavior reasonably well for lower iterations (iterations 344, 1055,
773 7954) but are less sensitive to image differences later in training (iterations
774 24400 vs 60000). We believe this discrepancy might be due to most met-
775 rics being sensitive to *distortion* rather than *perceptual quality*, two different
776 dimensions of image quality that have been explored theoretically in image
777 restoration tasks [57]. *Distortion* refers to the deviation of a distorted im-
778 age from a noise-free reference whereas *perceptual quality* measures the more
779 elusive 'naturalness' of an image regardless of a reference. In our experi-
780 ment, informal evidence in line with this is a participant stating⁴ that they
781 focused on image artifacts such as checkerboard patterns for low-quality im-
782 ages (i.e., distortion), but that it was rather subtle differences in luminance
783 across the image would that made higher-quality generated images look fake
784 (i.e., perceptual quality).

785 If human assessment of image quality was indeed governed by distortion
786 and artifacts for earlier iterations and a more intuitive understanding of what
787 an MRI should look like for later iterations, it is reasonable to believe that
788 most of the metrics were sensitive to artifacts, hence explaining their good
789 performance for early iterations and bad performance for the last iteration.
790 The only model that consistently agreed with human data across all quality
791 levels was the Deep QA model. The model was directly trained on the images
792 and human responses which gave it the opportunity to learn the more subtle

³<https://github.com/treder/MRI-GAN-QA>

⁴the participant disclosed their participation in the experiment to the author in a personal email

793 features that humans use for more high-quality generated images.

794 4.3. *The limits of GANs*

795 The promise of GANs is somewhat uncanny: After being exposed to a few
796 hundred images it supposedly learns the brain image manifold in image space.
797 A key question therefore is whether the GAN indeed learns (a) the generative
798 distribution of MRIs, whether it learns (b) a different distribution, or (c)
799 simply memorizes and reproduces the training data. Regarding (c), it turns
800 out that memorization of data requires a large number of different samples
801 [58]. Other indicators contra memorization are the fact that interpolation
802 between noise vectors can produce novel and meaningful image variations
803 [59] and the clear disparity seen between real and generated images [60].

804 In accordance with (a), GANs can closely approximate the underlying
805 generative distribution for sufficiently large datasets [6]. Generalization be-
806 yond the training data, indicated by the emergence of novel combinations of
807 features, has been shown in simple datasets [61]. However, there is doubt
808 that the same holds for limited sample sizes [60]. For instance, the *generative*
809 *support* (crudely speaking, the number of individually different images the
810 GAN can produce) found for faces strongly depends on the model architec-
811 ture, with estimates ranging from 160,000 to over a million faces. This casts
812 doubt on option (a). Further support for (b) is given by the phenomenon of
813 mode collapse, i.e. the tendency of GANs to concentrate its probability mass
814 to a few modes. [62] showed that the relative frequency of attributes such
815 as hair style for faces or room type for indoor scenes is distorted from the
816 real proportion in the data. Furthermore, classifier performance deteriorates
817 when trained on GAN images rather than real data, and [62] estimated that
818 the effective dataset size of the GAN is 100x smaller than the training data.

819 The debate on the generative support of GANs is ongoing, but results
820 so far favour alternative (b), that is, GANs probably approximate the true
821 distribution rather than just memorizing the data, but the support of the
822 GAN distribution is limited. This limitation can probably be partially coun-
823 teracted by using a large and diverse training set and a sufficiently capable
824 CNN architecture. For limited size datasets, data augmentation could be a
825 viable tool for increasing set size and diversity.

826 4.4. *The 'ideal' quality metric*

827 Although the focus of this paper was to assess how well existing metrics
828 agree with human data, it seems expedient to list other favourable properties

829 that an ideal MRI quality metric should possess:

- 830 1. *Model agnostic*. The metric should be independent of the specific model
831 architecture in order to enable quantitative comparison across models.
832 This precludes e.g. the discriminator’s loss function to act as a metric.
- 833 2. *Reference free*. It should determine the quality of an image per se,
834 without resorting to a reference image. This would make it applicable
835 to both noise-to-image and image-to-image applications and also in
836 situations wherein the ground truth is not available.
- 837 3. *Individual image ratings*. A measure of image quality should be pro-
838 vided for individual images rather than the whole dataset.
- 839 4. *Multi-dimensional*. Although MRIs are often depicted as a series of 2D
840 slices, the slices together actually form a 3D image and an accurate
841 assessment should take into account all spatial dimensions.
- 842 5. *Human performance*. The metric has to reproduce the perceptual as-
843 sessment of expert viewers. Furthermore, there is evidence that image
844 manipulations that are imperceptible to humans can have a signifi-
845 cant effect on CNN predictions [63]. Therefore, we believe that human
846 assessment should serve as a *lower sensitivity bound* when evaluating
847 image metrics. In other words, if humans are sensitive to the difference
848 between two types of images, so should be the metric – but the metric
849 could be sensitive to differences that humans are not.

850 None of the surveyed metrics meets all of these criteria. All of the sur-
851 veyed metrics are model agnostic but only IS, MIS, NIQE, BRISQUE, and
852 the Deep QA model are reference free. Individual image ratings are not
853 available for FID and MMD because they operate on the level of distribu-
854 tions rather than individual images. Although all metrics were applied to
855 2D image slices, IS, MIS, MMD and MMD can be seamlessly extended to
856 3D by deriving a feature embedding from a 3D CNN. The Deep QA model
857 can be trained on 3D images using 3D convolutions. A difficulty in extend-
858 ing metrics to 3D images is the collection of experimental data. A bespoke
859 experimental protocol for assessing 3D images wherein participants can e.g.
860 rotate, zoom into and slice 3D MRI images would need to be developed.

861 4.5. Limitations

862 Behavioral results were collected using an online experiment rather than
863 a controlled study in a lab environment. Participants performed the exper-
864 iments in different environments using different display devices (e.g. laptop

865 screen vs monitor). This potentially reduces the internal validity of our
866 results, although it may increase their external validity. We believe that
867 this limitation mostly affects group analyses. For instance, the seniority of
868 a participant may be correlated with the display device they use, thereby
869 confounding group differences between junior and senior experts. However,
870 we performed only within-subject analyses. If anything, the fact that our
871 findings are consistent across a variety of environments and display devices
872 strengthens rather than weakens their validity. For instance, using Amazon
873 Mechanical Turk, [64] replicated various visual effects such as Stroop effect
874 and attentional blink that have been studied extensively in laboratory set-
875 tings. Another limitation is that we only considered a noise-to-image model,
876 no image-to-image model. In the latter, reference images are available which
877 might improve the quality of metrics that rely on reference images.

878 *4.6. Conclusion*

879 As a practical recommendation for evaluating brain images produced by
880 GANs, we propose that researchers assess both local properties (image qual-
881 ity of individual images or groups of images) and global properties (statistical
882 regularities across the whole dataset) of the generated data.

883 With regard to global properties, the generated dataset should reproduce
884 relevant group differences such as young vs elderly or male vs female, with
885 the specific choice of relevant groups depending on the nature of the dataset
886 and the research question. Sensitivity of the GAN to group differences can
887 be investigated by either re-training it on separate groups, as done in this
888 paper, or using a conditional GAN with the demographics as separate in-
889 put variables. Furthermore, a realistic GAN should be able to reproduce
890 the large-scale structural networks exposed by multivariate analyses across
891 the dataset. We performed spatial ICA to verify that the GAN reproduces
892 components such as the Default Mode Network. Other approaches such as
893 Principal Component Analysis (PCA) or non-negative matrix factorization
894 [65] could be used for the same purpose.

895 With regard to local properties, the generated brain image should 'look
896 like' a real one. This can be verified using both distortion and perceptual
897 metrics, two dimensions of image quality [57]. **As distortion metrics, FID,
898 MMD, and NIQE stood out as versatile, widely used image quality met-
899 rics that do not require tuning to the MRI data. NIQE has complementary
900 properties compared to FID and MMD in that it yields ratings for individ-
901 ual images without requiring a reference, whereas FID and MMD provide**

902 a distribution-level metric for the difference between the set of generated
903 images vs the set of real images. Therefore, we recommend the use of all
904 three metrics. As a perceptual metric, our Deep QA rating model was the
905 only metric that reproduced human data for higher quality images. Unfor-
906 tunately, the model required human data to be trained. This defeats its
907 purpose to some extent since we set out to escape the need to collect human
908 data. Although speculative, a possible way out of this predicament is to pre-
909 train a Deep QA model on a large, diverse dataset that integrates multiple
910 data modalities [66, 67]. Ideally, this model could be used out-of-the-box for
911 any new dataset. Optionally, it could be fine-tuned to the statistics of a new
912 dataset without requiring human data.

913 Concluding, we propose a combination of local and global analyses for
914 assessing the quality of generated images. For the time being, human as-
915 sessment remains the gold standard for assessing individual images. In the
916 future, we believe that a Deep QA model that has been trained to mimick
917 human perceptual assessments can pave the way for quick and cost-effective
918 image quality assessments that will accelerate GAN research and improve its
919 validation.

920 Acknowledgements

921 We acknowledge the support of the Supercomputing Wales project, which
922 is part-funded by the European Regional Development Fund (ERDF) via
923 Welsh Government. KAT was supported by the Guarantors of Brain (G101149).

924 Author contributions

925 MT conceptualized the study and collected human data. RC trained
926 the GAN. MT and KT performed the statistical analyses and drafted the
927 manuscript. All authors reviewed and approved the manuscript.

928 References

- 929 [1] R. A. Poldrack, K. J. Gorgolewski, Making big data open: data sharing
930 in neuroimaging, *Nature Neuroscience* 17 (2014) 1510–1517.
- 931 [2] H. Peng, W. Gong, C. F. Beckmann, A. Vedaldi, S. M. Smith, Accurate
932 brain age prediction with lightweight deep neural networks, *Medical*
933 *Image Analysis* 68 (2021) 101871.

- 934 [3] M. Sajjad, S. Khan, K. Muhammad, W. Wu, A. Ullah, S. W. Baik,
935 Multi-grade brain tumor classification using deep CNN with extensive
936 data augmentation, *Journal of Computational Science* 30 (2019) 174–
937 182.
- 938 [4] G. Mohan, M. M. Subashini, MRI based medical image analysis: Survey
939 on brain tumor grade classification, *Biomedical Signal Processing and*
940 *Control* 39 (2018) 139–161.
- 941 [5] V. Sorin, Y. Barash, E. Konen, E. Klang, Creating Artificial Images for
942 Radiology Applications Using Generative Adversarial Networks (GANs)
943 – A Systematic Review, *Academic Radiology* 27 (2020) 1175–1185.
- 944 [6] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley,
945 S. Ozair, A. Courville, Y. Bengio, Generative Adversarial Nets, in: Z. G.
946 Weinberger, M. Welling, C. Cortes, N. D. Lawrence, K. Q. (Eds.), *Ad-*
947 *vances in Neural Information Processing Systems 27*, Curran Associates,
948 Inc., 2014, pp. 2672–2680.
- 949 [7] D. P. Kingma, M. Welling, Auto-encoding variational bayes, in: 2nd In-
950 ternational Conference on Learning Representations, ICLR 2014 - Con-
951 ference Track Proceedings, International Conference on Learning Rep-
952 resentations, ICLR, 2014.
- 953 [8] A. U. Hirte, M. Platscher, T. Joyce, J. J. Heit, E. Tranvinh, C. Fed-
954 erau, Diffusion-Weighted Magnetic Resonance Brain Images Generation
955 with Generative Adversarial Networks and Variational Autoencoders: A
956 Comparison Study, *arXiv* (2020) 2006.13944.
- 957 [9] C. K. Chong, E. T. W. Ho, Synthesis of 3D MRI Brain Images with
958 Shape and Texture Generative Adversarial Deep Neural Networks, *IEEE*
959 *Access* 9 (2021) 64747–64760.
- 960 [10] S. U. H. Dar, M. Yurt, L. Karacan, A. Erdem, E. Erdem, T. Çukur,
961 Image Synthesis in Multi-Contrast MRI with Conditional Generative
962 Adversarial Networks, *IEEE Transactions on Medical Imaging* 38 (2019)
963 2375–2388.
- 964 [11] D. Nie, R. Trullo, J. Lian, C. Petitjean, S. Ruan, Q. Wang, D. Shen,
965 Medical Image Synthesis with Context-Aware Generative Adversarial

- 966 Networks, in: Medical Image Computing and Computer Assisted In-
967 tervention - MICCAI 2017. MICCAI 2017. Lecture Notes in Computer
968 Science, Springer, Cham, 2017, pp. 417–425.
- 969 [12] D. Abramian, A. Eklund, Generating fMRI volumes from T1-weighted
970 volumes using 3D CycleGAN, arXiv abs/1907.0 (2019).
- 971 [13] X. Gu, H. Knutsson, M. Nilsson, A. Eklund, Generating Diffusion MRI
972 Scalar Maps from T1 Weighted Images Using Generative Adversarial
973 Networks, Lecture Notes in Computer Science 11482 (2019) 489–498.
- 974 [14] Y. Chen, F. Shi, A. G. Christodoulou, Y. Xie, Z. Zhou, D. Li, Efficient
975 and Accurate MRI Super-Resolution Using a Generative Adversarial
976 Network and 3D Multi-level Densely Connected Network, in: Medical
977 Image Computing and Computer Assisted Intervention – MICCAI 2018,
978 Springer, Cham, 2018, pp. 91–99.
- 979 [15] T. M. Quan, T. Nguyen-Duc, W.-K. Jeong, Compressed Sensing MRI
980 Reconstruction Using a Generative Adversarial Network With a Cyclic
981 Loss, IEEE Transactions on Medical Imaging 37 (2018) 1488–1497.
- 982 [16] X. Yi, E. Walia, P. Babyn, Generative adversarial network in medical
983 imaging: A review, arXiv e-prints 1809.07294 (2018).
- 984 [17] C. Bermudez, A. J. Plassard, T. L. Davis, A. T. Newton, S. M. Resnick,
985 B. A. Landman, Learning Implicit Brain MRI Manifolds with Deep
986 Learning, Proceedings of SPIE—the International Society for Optical
987 Engineering 10574 (2018).
- 988 [18] K. Kazuhiro, R. A. Werner, F. Toriumi, M. S. Javadi, M. G. Pomper,
989 L. B. Solnes, F. Verde, T. Higuchi, S. P. Rowe, Generative Adversarial
990 Networks for the Creation of Realistic Artificial Brain Magnetic Reso-
991 nance Images, Tomography (Ann Arbor, Mich.) 4 (2018) 159–163.
- 992 [19] A. Volokitin, E. Erdil, N. Karani, K. C. Tezcan, X. Chen, L. Van Gool,
993 E. Konukoglu, Modelling the Distribution of 3D Brain MRI Using a
994 2D Slice VAE, in: Lecture Notes in Computer Science, volume 12267
995 LNCS, Springer Science and Business Media Deutschland GmbH, 2020,
996 pp. 657–666.

- 997 [20] G. Kwon, C. Han, D.-s. Kim, Generation of 3D Brain MRI Using Auto-
998 Encoding Generative Adversarial Networks, in: MICCAI 2019: Medical
999 Image Computing and Computer Assisted Intervention, pp. 118–126.
- 1000 [21] C. Han, L. Rundo, K. Murao, T. Noguchi, Y. Shimahara, Z. A. Milacski,
1001 S. Koshino, E. Sala, H. Nakayama, S. Satoh, MADGAN: unsupervised
1002 medical anomaly detection GAN using multiple adjacent brain MRI slice
1003 reconstruction, *BMC Bioinformatics* 2021 22:2 22 (2021) 1–20.
- 1004 [22] F. Calimeri, A. Marzullo, C. Stamile, G. Terracina, Biomedical Data
1005 Augmentation Using Generative Adversarial Neural Networks, in:
1006 P. Verschure, A. Villa, A. Lintas, S. Rovetta (Eds.), *Artificial Neural
1007 Networks and Machine Learning - ICANN 2017. ICANN 2017. Lecture
1008 Notes in Computer Science*, Springer, Cham, 2017, pp. 626–634.
- 1009 [23] C. Han, L. Rundo, K. Murao, Z. A. Milacski, K. Umemoto,
1010 H. Nakayama, S. Satoh, GAN-based Multiple Adjacent Brain MRI
1011 Slice Reconstruction for Unsupervised Alzheimer’s Disease Diagnosis,
1012 in: *In International Meeting on Computational Intelligence Methods for
1013 Bioinformatics and Biostatistics. Lecture Notes in Computer Science*,
1014 Springer, Cham, 2019, pp. 44–54.
- 1015 [24] K. H. Kim, W. Do, S. Park, Improving resolution of MR images with an
1016 adversarial network incorporating images with different contrast, *Med-
1017 ical Physics* 45 (2018) 3120–3131.
- 1018 [25] D. Yang, B. Liu, L. Axel, D. Metaxas, 3D LV Probabilistic Segmentation
1019 in Cardiac MRI Using Generative Adversarial Network, in: *Statistical
1020 Atlases and Computational Models of the Heart. Atrial Segmentation
1021 and LV Quantification Challenges*, Springer, Cham, 2019, pp. 181–190.
- 1022 [26] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, Image quality
1023 assessment: From error visibility to structural similarity, *IEEE Trans-
1024 actions on Image Processing* 13 (2004) 600–612.
- 1025 [27] N. Luo, J. Sui, A. Abrol, J. Chen, J. A. Turner, E. Damaraju, Z. Fu,
1026 L. Fan, D. Lin, C. Zhuo, Y. Xu, D. C. Glahn, A. L. Rodrigue, M. T.
1027 Banich, G. D. Pearlson, V. D. Calhoun, Structural brain networks match
1028 intrinsic functional networks and vary across domains: a study from
1029 15000+ individuals, *Cerebral Cortex* 30 (2020) 5460–5470.

- 1030 [28] S. Hong, R. Marinescu, A. V. Dalca, A. K. Bonkhoff, M. Bretzner, N. S.
1031 Rost, P. Golland, 3D-StyleGAN: A Style-Based Generative Adversarial
1032 Network for Generative Modeling of Three-Dimensional Medical Images,
1033 in: 24th International Conference on Medical Image Computing and
1034 Computer Assisted Intervention (MICCAI).
- 1035 [29] M. Shafto, L. K. Tyler, M. Dixon, J. R. Taylor, J. B. Rowe, R. Cusack,
1036 A. J. Calder, W. D. Marslen-Wilson, J. Duncan, T. Dalgleish, R. N.
1037 Henson, C. Brayne, F. E. Matthews, The Cambridge Centre for Ageing
1038 and Neuroscience (Cam-CAN) study protocol: a cross-sectional, lifes-
1039 pan, multidisciplinary examination of healthy cognitive ageing, *BMC*
1040 *neurology* 14 (2014) 204.
- 1041 [30] J. R. Taylor, N. Williams, R. Cusack, T. Auer, M. A. Shafto, M. Dixon,
1042 L. K. Tyler, Cam-CAN, R. N. Henson, The Cambridge Centre for Age-
1043 ing and Neuroscience (Cam-CAN) data repository: Structural and func-
1044 tional MRI, MEG, and cognitive data from a cross-sectional adult lifes-
1045 pan sample, *NeuroImage* 144 (2017) 262–269.
- 1046 [31] J. Ashburner, A fast diffeomorphic image registration algorithm, *Neu-
1047 roImage* 38 (2007) 95–113.
- 1048 [32] K. A. Tsvetanov, R. N. A. Henson, P. S. Jones, H. Mutsaerts,
1049 D. Fuhrmann, L. K. Tyler, J. B. Rowe, The effects of age on resting-state
1050 BOLD signal variability is explained by cardiovascular and cerebrovas-
1051 cular factors, *Psychophysiology* 58 (2021) e13714.
- 1052 [33] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein Generative Adversar-
1053 ial Networks, in: *Proceedings of the 34th International Conference on*
1054 *Machine Learning*, pp. 214–223.
- 1055 [34] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, A. C. Courville, Im-
1056 proved Training of Wasserstein GANs, in: *Proceedings of the 31st Inter-*
1057 *national Conference on Neural Information Processing Systems*, Curran
1058 Associates Inc., Long Beach, California, USA, 2017, p. 5769–5779.
- 1059 [35] G. Stoet, PsyToolkit: A software package for programming psycholog-
1060 ical experiments using Linux, *Behavior Research Methods* 42 (2010)
1061 1096–1104.

- 1062 [36] G. Stoet, PsyToolkit: A Novel Web-Based Method for Running Online
1063 Questionnaires and Reaction-Time Experiments, *Teaching of Psychol-*
1064 *ogy* 44 (2017) 24–31.
- 1065 [37] R. Whelan, Effective analysis of reaction time data, *Psychological*
1066 *Record* 58 (2008) 475–482.
- 1067 [38] M. Frid-Adar, I. Diamant, E. Klang, M. Amitai, J. Goldberger,
1068 H. Greenspan, GAN-based synthetic medical image augmentation for
1069 increased CNN performance in liver lesion classification, *Neurocomput-*
1070 *ing* 321 (2018) 321–331.
- 1071 [39] L. Xu, K. M. Groth, G. Pearlson, D. J. Schretlen, V. D. Calhoun, Source-
1072 based morphometry: The use of independent component analysis to
1073 identify gray matter differences with application to schizophrenia, *Hum-*
1074 *an Brain Mapping* 30 (2009) 711–724.
- 1075 [40] A. Borji, Pros and cons of GAN evaluation measures, *Computer Vision*
1076 *and Image Understanding* 179 (2019) 41–65.
- 1077 [41] C. Lee, S. Woo, S. Baek, J. Han, J. Chae, J. Rim, Comparison of
1078 objective quality models for adaptive bit-streaming services, in: *2017*
1079 *8th International Conference on Information, Intelligence, Systems &*
1080 *Applications (IISA)*, volume 2017, Institute of Electrical and Electronics
1081 Engineers Inc., 2017, pp. 1–4.
- 1082 [42] J. Jia Deng, W. Wei Dong, R. Socher, L.-J. Li-Jia Li, K. Kai Li, L. Li Fei-
1083 Fei, ImageNet: A large-scale hierarchical image database, in: *2009 IEEE*
1084 *Conference on Computer Vision and Pattern Recognition*, IEEE, 2009,
1085 pp. 248–255.
- 1086 [43] S. Barratt, R. Sharma, A Note on the Inception Score, *arXiv* (2018)
1087 1801.01973.
- 1088 [44] S. Gurumurthy, R. K. Sarvadevabhatla, V. B. Radhakrishnan, DeLi-
1089 GAN : Generative Adversarial Networks for Diverse and Limited Data,
1090 *Proceedings - 30th IEEE Conference on Computer Vision and Pattern*
1091 *Recognition, CVPR 2017 2017-Janua* (2017) 4941–4949.
- 1092 [45] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, S. Hochreiter,
1093 GANs Trained by a Two Time-Scale Update Rule Converge to a Local

- 1094 Nash Equilibrium, *Advances in Neural Information Processing Systems*
1095 2017-Decem (2017) 6627–6638.
- 1096 [46] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, A. Smola,
1097 A Kernel Two-Sample Test, *Journal of Machine Learning Research* 13
1098 (2012) 723–773.
- 1099 [47] A. Mittal, A. Krishna Moorthy, A. Conrad Bovik, No-Reference Image
1100 Quality Assessment in the Spatial Domain, *IEEE Transactions on Image*
1101 *Processing* 21 (2012) 4695–4708.
- 1102 [48] G. Yang, Y. Cao, X. Xing, M. Wei, Perceptual Loss Based Super-
1103 Resolution Reconstruction from Single Magnetic Resonance Imaging,
1104 in: *Lecture Notes in Computer Science*, volume 11632, Springer Verlag,
1105 2019, pp. 411–424.
- 1106 [49] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, O. Wang, The Un-
1107 reasonable Effectiveness of Deep Features as a Perceptual Metric, in:
1108 *Proceedings of the IEEE Conference on Computer Vision and Pattern*
1109 *Recognition (CVPR)*, pp. 586–595.
- 1110 [50] T. Karras, S. Laine, T. Aila, A Style-Based Generator Architecture for
1111 Generative Adversarial Networks, *Proceedings of the IEEE Computer*
1112 *Society Conference on Computer Vision and Pattern Recognition 2019-*
1113 *June (2018)* 4396–4405.
- 1114 [51] C. Rorden, M. Brett, Stereotaxic display of brain lesions, *Behavioural*
1115 *Neurology* 12 (2000) 191–200.
- 1116 [52] D. A. Magezi, Linear mixed-effects models for within-participant psy-
1117 chology experiments: an introductory tutorial and free, graphical user
1118 interface (LMMgui), *Frontiers in Psychology* 6 (2015) 2.
- 1119 [53] C. Keeble, P. D. Baxter, A. J. Gislason-Lee, L. A. Treadgold, A. G.
1120 Davies, Methods for the analysis of ordinal response data inmedical
1121 image quality assessment, *British Journal of Radiology* 89 (2016).
- 1122 [54] R. H. B. Christensen, ordinal—Regression Models for Ordinal Data,
1123 2019.

- 1124 [55] K. Liu, S. Yao, K. Chen, J. Zhang, L. Yao, K. Li, Z. Jin, X. Guo,
1125 Structural Brain Network Changes across the Adult Lifespan, *Frontiers*
1126 in *Aging Neuroscience* 9 (2017) 275.
- 1127 [56] C. Ge, I. Y. H. Gu, A. S. Jakola, J. Yang, Deep semi-supervised learning
1128 for brain tumor classification, *BMC Medical Imaging* 20 (2020).
- 1129 [57] Y. Blau, T. Michaeli, The Perception-Distortion Tradeoff, in: *Proceed-*
1130 *ings of the IEEE Conference on Computer Vision and Pattern Recogni-*
1131 *tion (CVPR)*, pp. 6228–6237.
- 1132 [58] V. Nagarajan, C. Raffel, G. Brain, I. J. Goodfellow Google Brain, The-
1133 oretical Insights into Memorization in GANs, in: *32nd Conference on*
1134 *Neural Information Processing Systems*.
- 1135 [59] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning
1136 with deep convolutional generative adversarial networks, in: *4th Inter-*
1137 *national Conference on Learning Representations, ICLR 2016 - Confer-*
1138 *ence Track Proceedings, International Conference on Learning Repre-*
1139 *sentations, ICLR, 2016*.
- 1140 [60] S. Arora, Y. Zhang, Do GANs actually learn the distribution? An
1141 empirical study, *arXiv* (2017) 1706.08224.
- 1142 [61] S. Zhao, H. Ren, A. Yuan, J. Song, N. Goodman, S. Ermon, Bias
1143 and Generalization in Deep Generative Models: An Empirical Study,
1144 *Advances in Neural Information Processing Systems 2018-Decem* (2018)
1145 10792–10801.
- 1146 [62] S. Santurkar, L. Schmidt, A. Madry, A Classification-Based Perspec-
1147 tive on GAN Distributions, in: *International Conference on Learning*
1148 *Representations (ICLR)*.
- 1149 [63] J. Zhang, C. Li, Adversarial Examples: Opportunities and Challenges,
1150 *IEEE Transactions on Neural Networks and Learning Systems* 31 (2020)
1151 2578–2593.
- 1152 [64] M. J. C. Crump, J. V. McDonnell, T. M. Gureckis, Evaluating Amazon’s
1153 Mechanical Turk as a Tool for Experimental Behavioral Research, *PLoS*
1154 *ONE* 8 (2013) e57410.

- 1155 [65] A. Anderson, P. K. Douglas, W. T. Kerr, V. S. Haynes, A. L. Yuille,
1156 J. Xie, Y. N. Wu, J. A. Brown, M. S. Cohen, Non-negative matrix
1157 factorization of multimodal MRI, fMRI and phenotypic data reveals
1158 differential changes in default mode subnetworks in ADHD, *NeuroImage*
1159 102 (2014) 207–219.
- 1160 [66] K. Aderghal, A. Khvostikov, A. Krylov, J. Benois-Pineau, K. Afdel,
1161 G. Catheline, Classification of Alzheimer Disease on Imaging Modalities
1162 with Deep CNNs Using Cross-Modal Transfer Learning, in: *Proceedings*
1163 *- IEEE Symposium on Computer-Based Medical Systems*, volume 2018-
1164 June, Institute of Electrical and Electronics Engineers Inc., 2018, pp.
1165 345–350.
- 1166 [67] K. Wang, A. Mamidipalli, T. Retson, N. Bahrami, K. Hasenstab,
1167 K. Blansit, E. Bass, T. Delgado, G. Cunha, M. S. Middleton, R. Loomba,
1168 B. A. Neuschwander-Tetri, C. B. Sirlin, A. Hsiao, Automated CT and
1169 MRI Liver Segmentation and Biometry Using a Generalized Convo-
1170 lutional Neural Network, *Radiology: Artificial Intelligence* 1 (2019)
1171 180022.