# Synthetic-Neuroscore: Using A Neuro-AI Interface for Evaluating
# Generative Adversarial Networks

Zhengwei Wang[a], Qi She[b], Alan F. Smeaton[c], Tomás E. Ward[c], Graham Healy[c]

[a]*V-SENSE, School of Computer Science and Statistics, Trinity College Dublin, Dublin 1, Ireland*
[b]*Intel Labs, Beijing, China*
[c]*Insight Centre for Data Analytics, Dublin City University, Dublin 9, Ireland*

## Abstract

Generative adversarial networks (GANs) are increasingly attracting attention in the computer vision, natural language processing, speech synthesis and similar domains. Arguably the most striking results have been in the area of image synthesis. However, evaluating the performance of GANs is still an open and challenging problem. Existing evaluation metrics primarily measure the dissimilarity between real and generated images using automated statistical methods. They often require large sample sizes for evaluation and do not directly reflect human perception of image quality.

In this work, we describe an evaluation metric we call **Neuroscore**, for evaluating the performance of GANs, that more directly reflects psychoperceptual image quality through the utilization of brain signals. Our results show that Neuroscore has superior performance to the current evaluation metrics in that: (1) It is more consistent with human judgment; (2) The evaluation process needs much smaller numbers of samples; and (3) It is able to rank the quality of images on a per GAN basis.

A convolutional neural network (CNN) based **neuro-AI interface** is proposed to predict Neuroscore from GAN-generated images directly without the need for neural responses. Importantly, we show that including neural responses during the training phase of the network can significantly improve the prediction capability

---

*Email addresses:* `zhengwei.wang@tcd.ie`, Work done in the Insight Centre for Data Analytics, Dublin City University (Zhengwei Wang), `qi.she@intel.com` (Qi She), `alan.smeaton@dcu.ie` (Alan F. Smeaton), `tomas.ward@dcu.ie` (Tomás E. Ward), `graham.healy@dcu.ie` (Graham Healy)

of the proposed model. Materials related to this work are provided at *https://github.com/villawang/Neuro-AI-Interface*.

*Keywords:* Neuroscore, Generative adversarial networks, Neuro-AI interface, Brain-computer interface.

## 1. Introduction

There is a growing interest in studying generative adversarial networks (GANs) in the deep learning community [1, 2]. Specifically, GANs have been widely applied to various domains such as computer vision [3], natural language processing [4], speech synthesis [5] and time series synthesis [6]. Compared with other deep generative models (e.g. variational autoencoders (VAEs)), GANs are favored for effectively handling sharp estimated density functions, efficiently generating desired samples and eliminating deterministic bias. Due to these properties GANs have successfully contributed to plausible image generation [3], image to image translation [7], image super-resolution [8], image completion [9] etc.

However, three main challenges currently in research into GANs could be: (1) Mode collapse – the model cannot learn the distribution of the full dataset well, which leads to poor generalization ability; (2) Difficult to train – it is non-trivial for the discriminator and generator in a GAN to achieve Nash equilibrium [10] during training; (3) Hard to evaluate – the evaluation of GANs can be considered as an effort to measure the dissimilarity between the real distribution $p_r$ and the generated distribution $p_g$. Unfortunately, the accurate estimation of $p_r$ is intractable. Thus, it is challenging to have a good estimation of the correspondence between $p_r$ and $p_g$. Challenges (1) and (2) are more concerned with computational aspects where much research has been carried out to mitigate these issues [11, 12, 13]. Challenge (3) is similarly fundamental, however limited literature is available and most of the current metrics only focus on measuring the dissimilarity between training and generated images. A more meaningful GAN evaluation metric that is consistent with human perceptions is paramount in helping researchers to further refine and design better GANs.

Although some evaluation metrics, e.g., Inception Score (IS), Kernel Maximum Mean Discrepancy (MMD) and Fréchet Inception Distance (FID), have already been proposed [12, 10, 14], they have a number of limitations. Firstly, these metrics do not agree with human perceptual judgments and human rankings of GAN models. A small artifact in images can have a large effect on the decision made by a machine learning system [15], whilst the intrinsic image content

2

does not change. In this aspect, we consider human perception to be more robust to adversarial images samples when compared to a machine learning system. Secondly, these metrics require large sample sizes for evaluation [16, 12] and acquiring large-scale samples for evaluation sometimes is not realistic in real-world applications since generating them may be time-consuming. Finally, the existing metrics are not able to rank individual GAN-generated images by their quality i.e., metrics are generated on a collection of images rather than on a single image basis. The within-GAN variances are crucial because they can provide an insight into the variability of that GAN.

The current literature demonstrates that a CNN is able to predict neural responses in the inferior temporal cortex in an image recognition task [17, 18] via invasive BCI techniques [19]. The ways in which a CNN can be used to predict a neural response with a non-invasive BCI aspect is still an open question. Figure 1 illustrates a schematic of different mesoscopic and macroscopic neural measurement techniques using invasive and non-invasive approaches. In this schematic, only EEG (Electroencephalography) is non-invasively measured from the human scalp [20]. Other types of neural dynamics such as ECoG and LFP are measured invasively, which requires electrodes to be implanted. Compared to invasively measured neural dynamics, EEG pros are that it is a simple measurement, a non-painful experience during recording, easier to get ethics approval for and more easily generalized to real-world applications. However, EEG suffers challenges such as low signal quality (i.e., low SNR), low spatial resolution (interesting neural activities can span all of the scalp and are thus difficult to localise), all of which makes predicting EEG responses challenging.

With the success achieved by deep neural networks (DNNs) in areas including computer vision and natural language processing, the operation and functionality of DNNs and its connection with the human brain has been extensively studied and investigated in the literature [22, 23, 24, 25, 26, 27, 18, 28, 29]. In this research area, the convolutional neural network (CNN) is widely studied and compared with the visual system in the human brain because both are hierarchical systems and the processing steps are similar. For example in an object recognition task, both CNNs and humans recognize an object by progressively extracting higher-level representations of the visual input through a hierarchy where successive layers operate on the inputs of the proceeding layers e.g., certain patterns of basic shapes, edges and colors as input can be determined at higher levels of the hierarchy to be a particular complex object composed of the inputs. Work reported in [18] outlines a CNN approach to delving even more deeply into understanding the development and organization of sensory cortical processing. It
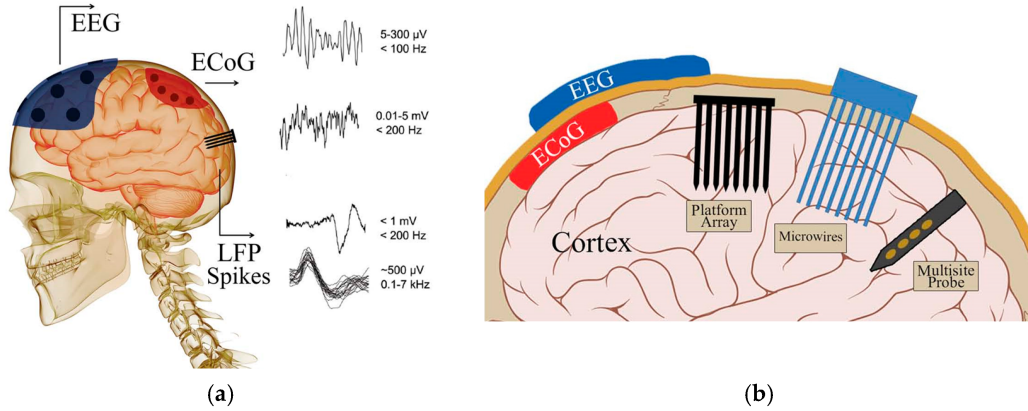
Figure 1: Schematic of different types of recorded neural signals (illustrated in (a)) via invasive and non-invasive measurements (illustrated in (b)). Figure from [21].

has recently been demonstrated that a CNN is able to reflect the spatio-temporal neural dynamics in the human brain visual processing area [22, 26, 25]. Despite much work carried out to reveal the similarity between CNNs and brain systems, research on interactions between CNNs and neural dynamics is limited.

In [17] the authors demonstrate that a CNN matched with neural data recorded from the inferior temporal cortex of a human subject [30] has high performance in an object recognition task. Given the evidence above that a CNN is able to predict neural responses in the brain, we explore the use of CNNs to predict P300 [31, 32] amplitudes in this paper. This type of model can then produce (synthetic) EEG feedback for different types of GANs.

By applying advanced statistical and machine learning techniques to non-invasive EEG, better source localization and reconstruction becomes possible. Our previous work [33, 34] demonstrated the effectiveness of using spatial filtering approaches for reconstructing P300 source ERP signals. Remaining low SNR issues can be further remedied by averaging EEG trials. Based on this evidence, we explore the use of DNNs to predict a metric we call Neuroscore [35], when neural information is available via EEG.

In this work, we describe a metric called Neuroscore to evaluate the performance of GANs, which is derived from a neurophysiological response recorded via non-invasive electroencephalography (EEG). We demonstrate and validate a neural-AI interface (as seen in Figure 2), which uses neural responses as supervisory information to train a CNN. The trained CNN model is then able to predict Neuroscore for images without requiring the corresponding neural responses. We

4

test this framework via three models: Shallow convolutional neural network, Mobilenet V2 [36] and Inception V3 [37].
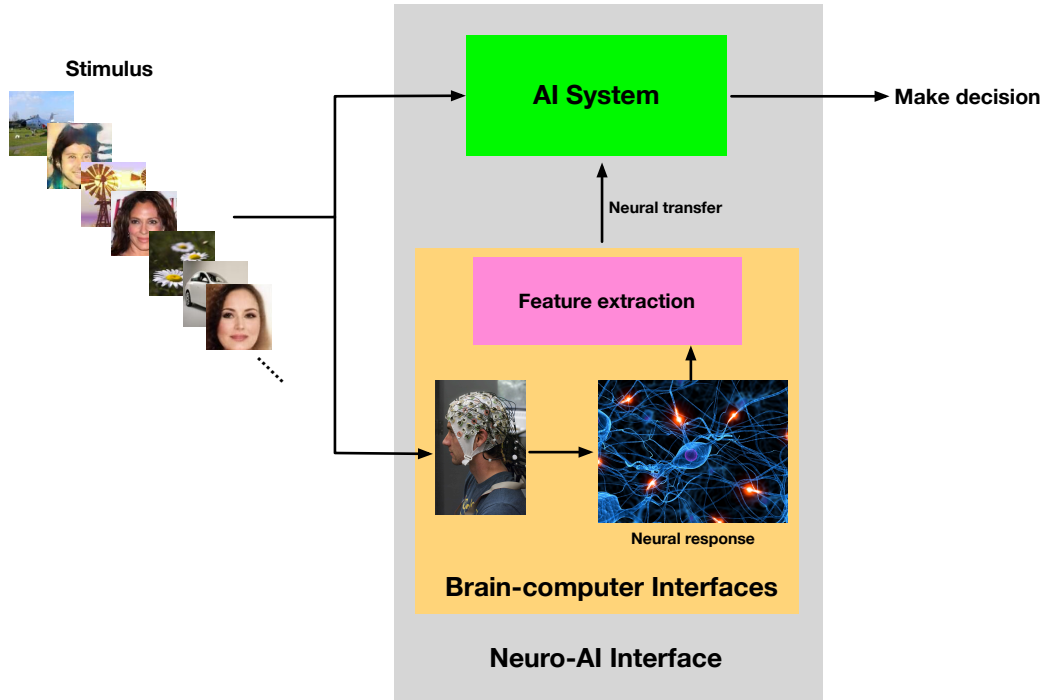


Figure 2: Schematic of a neuro-AI interface. Stimuli (image stimuli used in this work) are simultaneously presented to an AI system and to participants. Participants' neural responses are transferred to the AI system as supervised information for assisting the AI system to make decision.

In outline, Neuroscore is calculated via measurement of the P300, an event-related potential (ERP) present in EEG, via a rapid serial visual presentation (RSVP) paradigm. The P300 and RSVP paradigm are mature techniques in the brain-computer interface (BCI) community and have been applied in a wide variety of tasks such as image search [38], information retrieval [39], and others. The unique benefit of Neuroscore is that it more directly reflects human perceptual judgment of images, which is intuitively more reliable compared to conventional metrics in the literature [14]. In summary, our contributions are two-fold:

- We combine human perception research with GANs and deep learning research. This represents a new avenue of investigation in the development of better GANs technologies.

5

- We propose a type of neuro-AI interface and training strategy to generalize the use of Neuroscore, which can be directly used for GAN evaluations without recording EEG. This enables our Neuroscore to be more widely applied to real-world scenarios, with a new measure we name synthetic-Neuroscore.

## 2. Related Work

Three well-known metrics are compared with Neuroscore in this paper.

### 2.1. Inception Score (IS)

Inception Score is the most widely used metric in the literature [12, 16, 14]. It uses a pre-trained Inception network [37] as an image classification model $\mathcal{M}$ to compute

$$\text{IS} = \exp\left(\mathbb{E}_{\mathbf{x} \sim p_g}\left[\text{KL}\left(p_{\mathcal{M}}(\mathbf{y}|\mathbf{x}) \| p_{\mathcal{M}}(\mathbf{y})\right)\right]\right)$$

where $p_{\mathcal{M}}(\mathbf{y}|\mathbf{x})$ is the label distribution of $\mathbf{x}$ that is predicted by the model $\mathcal{M}$ and $p_{\mathcal{M}}(\mathbf{y})$ is the marginal probability of $p_{\mathcal{M}}(\mathbf{y}|\mathbf{x})$ over the probability $p_g$. A larger inception score will have $p_{\mathcal{M}}(\mathbf{y}|\mathbf{x})$ close to a point mass and $p_{\mathcal{M}}(\mathbf{y})$ close to uniform, which indicates that the Inception network is very confident that the image belongs to a particular ImageNet category [40] where all categories are equally represented. This suggests the generative model has both high quality and diversity.

### 2.2. Kernel Maximum Mean Discrepancy (MMD)

MMD is a method for comparing two distributions, in which the test statistic is the largest difference in expectations over functions in the unit ball of a reproducing kernel Hilbert space [41]. MMD is computed as

$$\text{MMD}^2(p_r, p_g) = \mathbb{E}_{\substack{\mathbf{x}_r, \mathbf{x}_r^\top \sim p_r \\ \mathbf{x}_g, \mathbf{x}_g^\top \sim p_g}}\left[k(\mathbf{x}_r, \mathbf{x}_r^\top) - 2k(\mathbf{x}_r, \mathbf{x}_g) + k(\mathbf{x}_g, \mathbf{x}_g^\top)\right]$$

It measures the dissimilarity between $p_r$ and $p_g$ for some fixed kernel function $k$, such as a Gaussian kernel [11]. A lower MMD indicates that $p_g$ is closer to $p_r$, showing the GAN has better performance.

## 2.3. Fréchet Inception Distance (FID)

FID uses a feature space extracted from a set of generated image samples by a specific layer of the Inception network [10]. The feature space is modelled via a multivariate Gaussian by the mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$. FID is computed as

$$\text{FID}(p_r, p_g) = \|\boldsymbol{\mu}_r - \boldsymbol{\mu}_g\|_2^2 + \text{Tr}\left(\boldsymbol{\Sigma}_r + \boldsymbol{\Sigma}_g - 2(\boldsymbol{\Sigma}_r\boldsymbol{\Sigma}_g)^{\frac{1}{2}}\right)$$

Similar to MMD, lower FID is better, corresponding to more similar real and generated samples as measured by the distance between their activation distributions.

For Inception Score, the score is calculated through the Inception model [37]. It has been shown that Inception Score is very sensitive to the model parameters [42]. Even the score produced by the same model trained using different libraries (e.g., Tensorflow, Keras, PyTorch) differ a lot from each other. It also requires a large sample size for the accurate estimation for $p_{\mathcal{M}}(y)$. FID and MMD both measure the similarity between training images and generated images based on the feature space [16], since the pixel representations of images do not naturally support for meaningful Euclidean distances to be computed [43]. The main concern about these two methods is whether the distributional characteristics of the feature space exactly reflect the distribution for the images [15].

We list the supported features of Neuroscore and traditional metrics in Table 1. Neuroscore can not only evaluate image quality as can the other metrics, but also have 3 unique characteristics, which will be demonstrated in Section 5.

| Feature | IS | MMD | FID | **Neuroscore** |
|---|---|---|---|---|
| Evaluate image quality | ✓ | ✗ | ✓ | ✓ |
| Consistent with human | ✗ | ✗ | ✗ | ✓ |
| Small sample size | ✗ | ✗ | ✗ | ✓ |
| Rank images | ✗ | ✗ | ✗ | ✓ |

Table 1: Comparison between Neuroscore and other metrics.

## 3. Preliminaries

### 3.1. Generative Adversarial Networks

A generative adversarial network (GAN) has two components, the discriminator $D$ and the generator $G$. Given a distribution $\boldsymbol{z} \sim p_{\boldsymbol{z}}$, $G$ defines a probability distribution $p_g$ as the distribution of the samples $G(\boldsymbol{z})$. The objective of a GAN is

to learn the generator's distribution $p_g$ that approximates the real data distribution $p_r$. Optimization of a GAN is performed with respect to a joint loss for $D$ and $G$ as

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim p_r} \log[D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} \log[1 - D(G(\mathbf{z}))]$$

### 3.2. P300 (or P3) Component and Preprocessing

In neuroscience, the P300 ERP component refers to a voltage change measured on the scalp which arises from current flow changes in the brain in response to a target stimulus [31], that can be measured with EEG. It reflects a participant's attention, which can be modulated by the specific instruction given to a participant. It has been shown in previous studies that real face stimuli generate larger P300/LPP potentials than non-real face stimuli such as cartoon face images [44, 45, 46]. Furthermore, the P300/LPP increases linearly with face realism, reflecting increased activity in visual and parietal cortex for more realistic faces[44]. The P300 response elicited by a target stimulus is typically evident between 300 – 600 ms post stimulus presentation depending on the type of task. EEG is normally recorded by using multiple channels e.g.. 32 channels, which makes it difficult to estimate the P300 source amplitude. We use an LDA beamformer [47, 33] to reconstruct the P300 source signal from the recorded raw EEG epochs.

Briefly, given a target EEG epoch $\mathbf{X}_i \in \mathbb{R}^{C \times T}$ and a standard EEG epoch[1] $\mathbf{K}_i \in \mathbb{R}^{C \times T}$ ($C$ is the number of channels and $T$ is time points in each EEG epoch). The optimization problem for the LDA beamformer is to find a projection vector $\mathbf{w} \in \mathbb{R}^{C \times 1}$ that solves the optimization problem:

$$\min_{\mathbf{w}} \mathbf{w}^\top \mathbf{\Sigma} \mathbf{w} \ \text{ s.t.} \mathbf{w}^\top \mathbf{p} = 1 \tag{1}$$

where $\mathbf{\Sigma} \in \mathbb{R}^{C \times C}$ is the EEG epoch covariance matrix ($\mathbf{\Sigma} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{X}_i \mathbf{X}_i^\top$, $N$ is number of trials) and $\mathbf{p} \in \mathbb{R}^{C \times 1}$ is the spatial pattern difference between target and standard condition [47]. The closed-form solution is

$$\mathbf{w} = \mathbf{\Sigma}^{-1} \mathbf{p} (\mathbf{p}^\top \mathbf{\Sigma}^{-1} \mathbf{p})^{-1} \tag{2}$$

The source signal of each single-trial $\mathbf{s}$ can be obtained as

$$\mathbf{s} = \mathbf{w}^\top \mathbf{X}_i = (\mathbf{p}^\top \mathbf{\Sigma}^{-1} \mathbf{p})^{-1} \mathbf{p}^\top \mathbf{\Sigma}^{-1} \mathbf{X}_i \tag{3}$$

---

[1]A target EEG epoch is an EEG trial (time duration 0 – 1 s) which corresponds to a target stimulus i.e., DCGAN, BEGAN, PROGAN and RFACE images in this study. A standard/non-target EEG epoch is an EEG trial which corresponds to a non-target images i.e., non-face image in this work.

where $\mathbf{s} \in \mathbb{R}^{1 \times T}$. Hence, LDA beamformer enables transformation of multi-channel EEG epochs to single-channel EEG epochs facilitating more robust measurement of the P300 and its amplitude.

## 4. Methodology

### 4.1. Data Acquisition and Experiment

We used three GAN models to generate synthetic images of faces: DCGAN [48], BEGAN [49] and progressive growing of GANs (PROGAN) [50] with sample outputs shown in Figure 3. Image streams in the experiment contain generated



Figure 3: Face image examples used in the experiments. From left to right: DCGAN, BEGAN, PROGAN, and real face (RFACE).

images from DCGAN, BEGAN and PROGAN, as well as real face (RFACE) images and non-face category images. RFACE images were sampled from the CelebA dataset [51]. Non-face category (standard images) were sampled from the ImageNet dataset [40], similar to those used in other RSVP experiments such as [52, 53]. EEG data for 12 participants was gathered. Data collection was carried out with approval from Dublin City University Research Ethics Committee (REC/2018/115). Each participant completed two types of task which we call the behavioral experiment (BE) task and the rapid serial visual presentation (RSVP) task. The sequence of blocks presented in the experiment was: BE → RSVP → BE → RSVP → BE. The presented images were randomly shuffled (across and within blocks), meaning the appearance of face images could not be predicted ahead of time by a participant i.e., they occurred at random times but always in the same quantity.

The objective of the BE task was to record participants' responses to each type of image category while the RSVP task was to record EEG when participants were viewing the rapid presentation of images. The ultimate goal of this study was to compare whether the EEG responses in the RSVP task were consistent with participants' responses in the BE task.

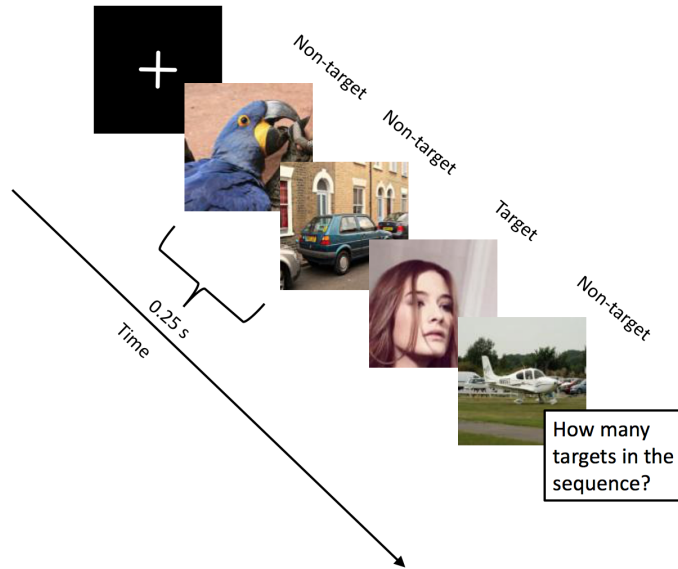An example of the RSVP experimental protocol is shown in Figure 4. The



Figure 4: An example of RSVP experimental protocol used in this work. A rapid image stream containing targets and standards (non-target) is presented to participants at 4 Hz (4 images per second) presentation rate.

RSVP task contained 26 blocks. Each RSVP block contained 240 images (6 images for each face category thus 24 face targets in total and 216 non-face images), thus there were 6,240 images (624 face targets / 5,616 non-face images) available for each participant. In the RSVP task, image streams were presented to participants at a 4 Hz presentation rate. Participants in RSVP blocks were asked to search for real face (RFACE) images[2]. This instruction to participants was constructed so that they would maintain attention to detect face images (from all GAN types), and furthermore focus their attention to what they perceived as real face images [32]. Details of the experiment can be found in [35].

EEG was recorded from participants in both the BE and RSVP tasks along with timestamp information for image presentation and behavioural responses (via a photodiode and hardware trigger) to allow for precise epoching of the EEG signals for each trial [54]. EEG data was acquired using a 32-channel BrainVision

[2]P300 responses were elicited for all GAN image categories e.g., while DCGAN had almost perfect behavioral accuracy labelled as being 'fake', DCGAN images still elicited a P300.

actiCHamp at 1,000 Hz sampling frequency, using electrode locations as defined by the International 10-20 system. To enhance the low signal-to-noise ratio of the acquired EEG, pre-processing is required. Pre-processing typically involves re-referencing, filtering the signal (by applying a bandpass filter to remove environmental noise or to remove activity in non-relevant frequencies), epoching (extracting a time epoch typically surrounding the stimulus onset) and trial/channel rejection (to remove those containing artifacts). In this work, a common average reference (CAR) was utilized and a bandpass filter (i.e., 0.5-20 Hz) was applied prior to epoching. EEG data was then downsampled to 250 Hz. Only trials where behavioral responses occurred between 0 and 1 second after the presentation of a stimulus were used. Trial rejection was carried out to remove those trials containing noise such as eye-related artifacts (via a peak-to-peak amplitude threshold across all electrodes).

*4.2. Neuroscore*

We used a rapid serial visual presentation (RSVP) paradigm [54, 34, 55] to elicit the P300 ERP. Our experimental procedure is illustrated in our previous published work [35]. We average the single-trial P300 amplitude (as Neuroscore) to mitigate the background EEG noise [31], which renders a stable measurement of the EEG response to a typical type of stimulus. In general, our Neuroscore is calculated via two steps: (1) Reconstruct the P300 source signal from the raw EEG; (2) Average the P300 amplitude of each reconstructed single-trial source signal across trials (see Algorithm 1).

The proposed Neuroscore reflects a human's perceptual response to different GANs via EEG measurements, thus it is consistent with human perceptual judgment on GANs. Figure 5 demonstrates the performance of Neuroscore calculated from a human neural response. In Figure 5(a), it can be seen that different image categories activate different P300 responses. Figure 5(b) illustrates a strong correlation between Neuroscore and human judgment (BE accuracy)[3]. These results demonstrate that Neuroscore reflects human judgment perception. More details can be found in our previous work [35].

---

[3]BE accuracy is the recorded accuracy (calculated as the number of correctly labeled images divided by the total number of images) in the behavioral experiment. Normalized BE accuracy is calculated by subtracting the average accuracy (across GAN types for that participant) from BE accuracy.

---

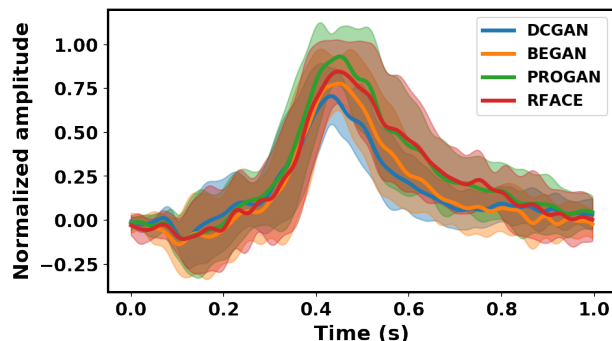**Algorithm 1** Calculation of Neuroscore

---

**Input:**

- $\mathbf{X} \in \mathbb{R}^{N \times C \times T}$ is the EEG signal corresponding to the target stimulus, where $N$ is the number of target trials, $C$ is the number of channels, and $T$ is the number of time points.

- $\mathbf{K} \in \mathbb{R}^{M \times C \times T}$ is the EEG signal corresponding to the standard stimulus, $M$ is number of standard trials, $C$ is number of channels, $T$ is number of time points. The target and standard EEG trials are already explained in section 3.2.
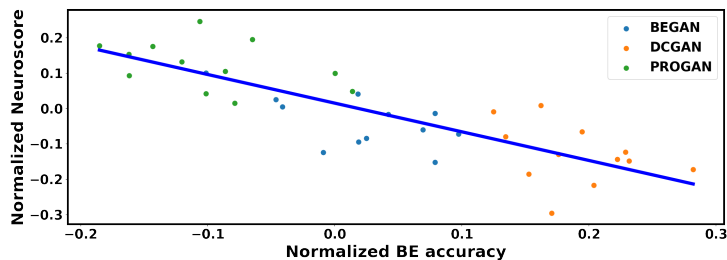
**Output:** Neuroscore

1: $\boldsymbol{\Sigma} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{X}_i \mathbf{X}_i^{\top} + \frac{1}{M} \sum_{i=1}^{M} \mathbf{K}_i \mathbf{K}_i^{\top}$
2: **for** $t_i$ in [400 ms, 600 ms] **do**
3: $\quad \mathbf{p} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{X}_{i,t_i} - \frac{1}{M} \sum_{i=1}^{M} \mathbf{K}_{i,t_i}$
4: $\quad \mathbf{w} = \boldsymbol{\Sigma}^{-1} \mathbf{p} (\mathbf{p}^{\top} \boldsymbol{\Sigma}^{-1} \mathbf{p})^{-1}$
5: $\quad J_{t_i} \leftarrow \mathbf{w}^{\top} \boldsymbol{\Sigma} \mathbf{w}$
6: $\quad W_{t_i} \leftarrow \mathbf{w}$
7: **end for**
8: $t_{optimal} = \text{argmin}_{t_i} J$
9: $\mathbf{w}_{optimal} = W_{t_{optimal}}$
10: $\text{t}_{P300} = [\text{t}_{optimal} - 100 \text{ ms}, \text{t}_{optimal} + 100 \text{ ms}]$     ▷ *This is time window being detected for P300.*
11: **for** $i = 1 : N$ **do**
12: $\quad \mathbf{s} = \mathbf{w}_{optimal}^{\top} \mathbf{X}_i$
13: $\quad a = \max(\mathbf{s}_{t_{p300}})$
14: $\quad A_i \leftarrow a$
15: **end for**
16: $\text{Neuroscore} = \dfrac{1}{N} \sum_{i=1}^{N} A_i$

---

### 4.3. Neuro-AI Interface

We propose a neuro-AI interface in order to generalize the use of Neuroscore. This kind of framework interfaces between neural responses and AI systems (a CNN is used in this study), which use neural responses as supervised information to train a CNN. The trained CNN is then used for generating a **synthetic-Neuroscore** given images generated by one of the popular GAN models i.e., average the outputs of corresponding images. Figure 6. demonstrates the schematic

(a) Reconstructed source P300 signals for each type of image category by using LDA beamformer across 12 participants. P300 component appears in 400 ms – 600 ms. Solid lines are averaged responses across participants while shadow areas represent the standard deviation.



(b) Correlation between Neuroscore and human judgment (i.e., participants' behavioural accuracy) across participants. **Pearson correlation statistics is** $r(36) = -0.828, p = 4.766 \times 10^{-10}$.

Figure 5: Performance of real Neuroscore, calculated from participants' neural responses.

of the neuro-AI interface used in this work[4]. Flow 1 shows that the image processed by a human being's brain produces a single-trial P300 source signal for each input image. Flow 2 in Figure. 6 demonstrates a CNN with included EEG signals during the training stage. The convolutional and pooling layers process the image similarly as retina has done [58]. *It should be noted that a CNN model is trained by using single images with their corresponding* **single-trial** *EEG informa-*

---

[4]We understand that a human being's brain system is much more complex than demonstrated in this work and that information flow in the brain is not one-directional [56, 57]. Our framework can be further extended to be more biologically plausible.
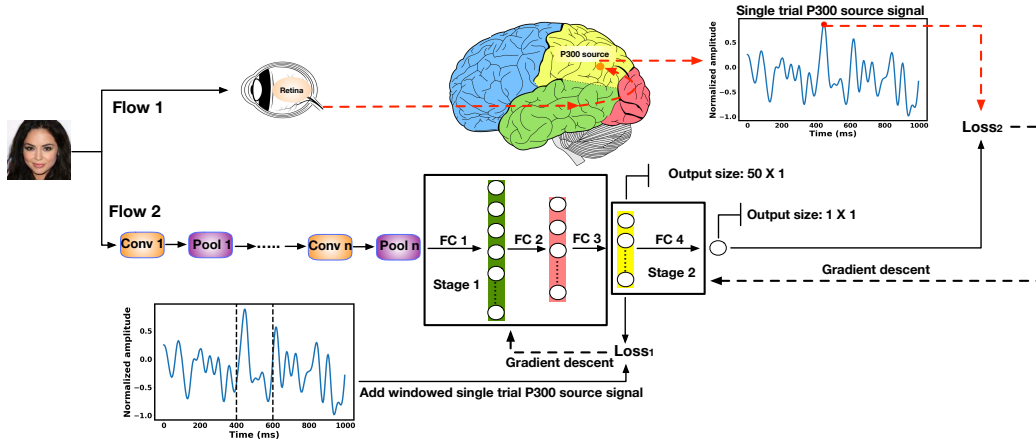
Figure 6: A neuro-AI interface and training details with added EEG information. Our training strategy includes two stages: (1) Learning from image to single-trial P300 source signal; and (2) Learning from single-trial P300 source signal to single-trial P300 amplitude. $loss_1$ is the $L_2$ distance between the yellow layer and the single-trial P300 source signal in the 400 – 600 ms corresponding to the single input image. $loss_2$ is the mean square error between model prediction and the single-trial P300 amplitude. $loss_1$ and $loss_2$ will be introduced in Section 4.4.

*tion (including single-trial P300 signal and single-trial P300 amplitude*[5]. Fully connected layers (FC) 1 – 3 aim to emulate the brain's functionality that produces the EEG signal. The yellow dense layer in the architecture aims to predict the single-trial P300 source signal at 400 – 600 ms in response to each image input. In order to help the model make a more accurate prediction for the single-trial P300 amplitude for the output, the single-trial P300 source signal at 400 – 600 ms is fed to the yellow dense layer to learn parameters for the previous layers in the training step. The model was then trained to predict the single-trial P300 source amplitude (the red point shown in signal-trail P300 source signal of Figure 6).

### 4.4. Training Details

Mobilenet V2, Inception V3 and Shallow network (architecture of Shallow network refers to Figure 7) were explored in this work, where in flow 2 we use

---

[5]Single-trial P300 amplitude refers the maximum value in the 400 ms – 600 ms time window of a single-trial EEG signal. Details can be referred to our previous work [35].). The averaged output of a trained model in terms of one image category can be represented as the **synthesized Neuroscore** (we refer to it as **synthetic-Neuroscore** in this paper)
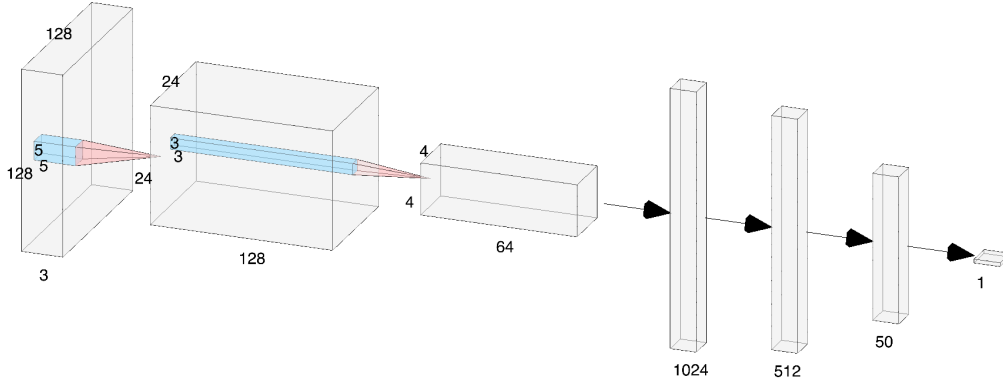
Figure 7: Shallow network architecture used in this work.

these three network bones such as Conv1-pooling layers. For Mobilenet V2 and Inception V3, we used ImageNet pre-trained parameters up to the FC 1 (as shown in Figure 6). Table 2 shows the FC layers details of three networks. Due to no pretrained parameters in the Shallow net, only three FC layers are contained in order to avoid overfitting. We trained parameters from FC 1 to FC 4 for Mobilenet V2 and Inception V3. $\theta_1$ is used to denote the parameters from FC 1 to FC 3 and $\theta_2$ indicates the parameters in FC 4. For the Shallow model, parameters up to FC 2 represent $\theta_1$ and parameters in FC 3 indicate $\theta_2$.

| Model | FC 1 | FC 2 | FC 3 | FC 4 |
|---|---|---|---|---|
| Shallow net | (1024, 512) | (512, 50) | (50, 1) | NA |
| Mobilenet | (1792, 896) | (896, 448) | (448, 50) | (50, 1) |
| Inception | (2048, 1024) | (1024, 512) | (512, 50) | (50, 1) |

Table 2: FC layers details of three networks investigated in this study.

We added EEG to the model because we first want to find a function $f(\chi) \to \mathbf{s}$ that maps the images space $\chi$ to the corresponding single-trial P300 source signal $\mathbf{s}$. This prior knowledge can help us to predict the single-trial P300 amplitude in the second learning stage.

We compared the performance of the models with, without EEG signals and with randomized EEG signals for training. We defined two stage loss function (loss$_1$ for a single-trial P300 source signal in the $400 - 600$ ms time window and

15

loss$_2$ for single-trial P300 amplitude) as

$$\text{loss}_1(\boldsymbol{\theta}_1) = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{S}_i^{true} - \mathbf{S}_i^{pred}(\boldsymbol{\theta}_1)\|_2^2$$

$$\text{loss}_2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \frac{1}{N} \sum_{i=1}^{N} (\text{y}_i^{true} - \text{y}_i^{pred}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2))^2 \tag{4}$$

where $\mathbf{S}_i^{true} \in \mathbb{R}^{1 \times T}$ is the single-trial P300 signal in the 400 - 600 ms time window to the presented image, $y_i$ refers to the single-trial P300 amplitude for each image, and $N$ refers to the batch size. In this case, we trained 20 epochs with batch size equaling to 256. An Adam optimizer with default hyperparameters was used and learning rate is 0.001.

The training of the models without using EEG signal is straightforward, models were trained directly to minimize loss$_2(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)$ by feeding images and the corresponding single-trial P300 amplitude. In this case, training is an end-to-end process i.e., from an image to single-trial a P300 amplitude without considering stage 1. The reason that we do this is to investigate the significance of adding single-trial P300 signal as supervisory information to the network. Training with EEG information is explained in Algorithm 2 and visualized in the "Flow 2" of

---

**Algorithm 2** Two training stages with EEG information.

---

**Stage 1:** Training parameters $\boldsymbol{\theta}_1$.

    **Input:** Images and averaged P300 signal $\mathbf{S}_i^{true}$.

1: **for** number of training iterations **do**

2:     Update $\boldsymbol{\theta}_1$ by descending its stochastic gradient:   $\nabla_{\boldsymbol{\theta}_1} \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{S}_i^{true} - \mathbf{S}_i^{pred}(\boldsymbol{\theta}_1)\|_2^2$

3: **end for**

**Stage 2:** Freezing $\boldsymbol{\theta}_1$, training parameters $\boldsymbol{\theta}_2$.

    **Input:** Images and single-trial P300 amplitude $\text{y}_i^{true}$.

4: **for** number of training iterations **do**

5:     Update $\boldsymbol{\theta}_2$ by descending its stochastic gradient:   $\nabla_{\boldsymbol{\theta}_2} \frac{1}{N} \sum_{i=1}^{N} (\text{y}_i^{true} - \text{y}_i^{pred}(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2))^2$

6: **end for**

---

Figure 6 with two stages. Stage 1 learns parameters $\boldsymbol{\theta}_1$ to predict P300 source signal while stage 2 learns parameters $\boldsymbol{\theta}_2$ to predict single-trial P300 amplitude with $\boldsymbol{\theta}_1$ fixed.

# 5. Results

## 5.1. EEG and Model Performance

*Individual Participant Performance.* Three models have been validated for each individual participant as shown in Figure 8. It can be seen that all three models
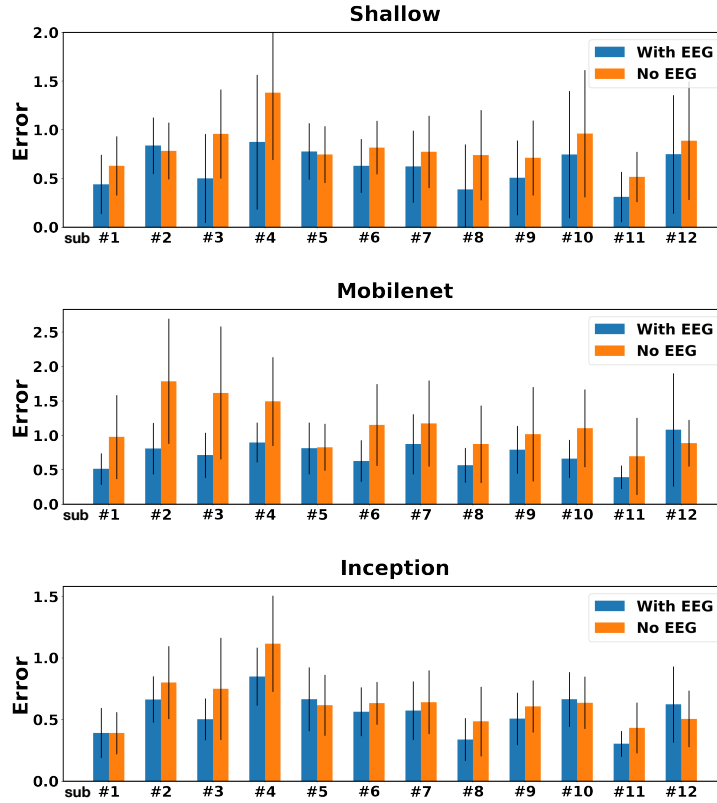


Figure 8: Error of 3 models with and without EEG signals. Error is defined as: $\sum_i^m |\text{Neuroscore}_{pred}^{(i)} - \text{Neuroscore}_{true}^{(i)}|$, where m = 3 is the number of GAN categories used (DCGAN, BEGAN, PROGAN, 12 participants) and Neuroscore is obtained by averaging single-trial P300 amplitudes. A smaller value indicates better performance. Details of numeric values can be refered to Table3.

trained with EEG outperform the models trained without EEG. **In other words, we show that including EEG/P300 time series signals as supervisory information to the yellow dense layer yields an improvement in performance as seen in Figure 6**. with smaller error and standard deviation across almost all individual subjects. For those cases where the reverse is true (7 from 36 have better or equal

17

performance without EEG), this might result from the number of EEG trials for an individual participant not being sufficient enough for training of deep networks to learn the mapping function $f(\chi)$ from image to EEG.

| Model | | Error mean (std) |
|---|---|---|
| Shallow net | Shallow-EEG | 0.151 (±0.245) |
| | Shallow | 0.428 (±0.623) |
| Mobilenet | Mobilenet-EEG | 0.155 (±0.235) |
| | Mobilenet | 0.437 (±0.589) |
| Inception | Inception-EEG | 0.157 (±0.487) |
| | Inception | 0.462 (±0.932) |

Table 3: Details of error mean and standard deviation for Figure 8.

*Cross Participant Performance.* We evaluated the cross participant performance of our approach by pooling trials across participants to see if the use of pooled trials produced a smaller error. In this case, the number of EEG trials across participants is 6012. We split data into training and testing as 2:1 in which there are 4008 trials for training and 2004 trials for testing. All trials are randomly shuffled and we repeat this process for 20 times in order to get a more robust result.

Table 4 shows the error for each model with the EEG signal, with a randomized EEG signal **within each type of GAN** and without an EEG signal. All models with EEG signals perform better than models without EEG signals, with much smaller errors and standard deviation.

Adding the EEG signal to the intermediate layer reduces error in all three models (as the same error is shown in Figure 8), namely 0.151, 0.168 and **0.171** for Shallow, Mobilenet, and Inception respectively. This indicates that the Inception model benefits most when adding EEG signal into the training stage. The performance of models with the EEG signal is ranked as Inception-EEG followwd by Mobilenet-EEG, and Shallow-EEG, which indicates that deeper neural networks may achieve better performance in this task. We used the randomized EEG signal here as a baseline to determine the efficacy of adding the EEG signal to produce better Neuroscore output. When randomizing the EEG signal, it shows that the error for each three model increases significantly. For Mobilenet and Inception, the error with the randomized EEG signal is even higher than those without the EEG signal in the training stage, demonstrating that EEG information in the training stage is crucial to each model.

| Model | | Error mean(std) |
|---|---|---|
| | Shallow-EEG | **0.209 (±0.102)** |
| Shallow net | Shallow-EEG$_{random}$ | 0.348 (±0.114) |
| | Shallow | 0.360 (±0.183) |
| | Mobilenet-EEG | **0.198 (±0.087)** |
| Mobilenet | Mobilenet-EEG$_{random}$ | 0.404 (±0.162) |
| | Mobilenet | 0.366 (±0.261) |
| | Inception-EEG | **0.173 (±0.069)** |
| Inception | Inception-EEG$_{random}$ | 0.392 (±0.057) |
| | Inception | 0.344 (±0.149) |

Table 4: Errors for 9 models across the 12 participants ("*-EEG" indicates models are trained with paired EEG, "*-EEG$_{random}$" refers to EEG trials which are randomized in the loss$_1$ **within each type of GAN**). Results are averaged by shuffling training/testing sets 20 times.

Figure 9 shows that the models with EEG information have a stronger correlation between synthetic-Neuroscore and Neuroscore. The cluster (blue, orange, and green circles) for each category of the model trained with EEG (left column) is more separable than the cluster produced by model without EEG (right column). This indicates that when with EEG is used in training models Neuroscore is more accurate and that Neuroscore is able to rank the performances of different GANs, which cannot be achieved with other metrics [14].

*5.2. Neuroscore Aligns with Human Perception*

Figure 5(b) shows the correlation between Neuroscore and human judgment (BE accuracy) according to three GANs: BEGAN, DCGAN, and PROGAN. The statistical test demonstrates the strong correlation between those two variables. This indicates that Neuroscore can be used to evaluate GANs as it reflects human perceptual judgment. A number of previous studies have noted that increasing task difficulty reduces the amplitude of the P300 [59, 60, 61, 62]. It may be the case that the larger P300 amplitudes observed for the PROGAN images indicate that these face images were easier to detect compared to the images from the other GANs. For example, DCGAN images tended to contain far more visual aberrations and other inherent artefacts that would impede their detection [63]. It has also been noted in another prior study that increased sensory evidence results in shorter reaction times and larger component amplitudes in temporal and spatial regions coinciding with those examined in our work [64]. Another prior study
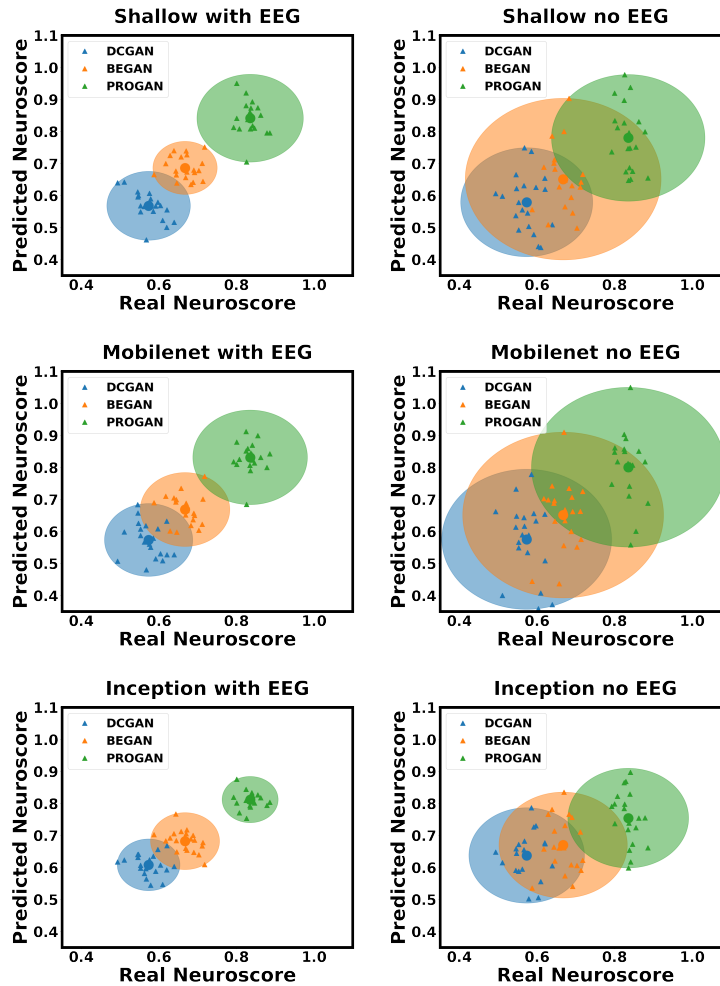
Figure 9: Scatter plot of synthetic-Neuroscore (vertical axis) and Neuroscore (horizontal) for 6 models (Shallow, Mobilenet, Inception with and without EEG signals for training) across participants, with 20 times repeated shuffling training and testing set. Each circle represents the cluster for a specific category. Small triangle markers inside each cluster correspond to each shuffling process. The dot at the center of each cluster is the mean.

explains larger P300 amplitudes for real face images resulting from enhanced perceptual processing [44]. In effect, larger average P300/LPP amplitudes for a particular GAN type are indicative that its images are perceived as being real faces.

We have already demonstrated that the Neuroscore derived from raw EEG is consistent with human perception [35]. We will now demonstrate the same prop-

erty of synthetic-Neuroscore predicted from the neuro-AI interface. We compare the synthetic-Neuroscore with three widely used evaluation metrics. The ultimate goal of GANs is to generate images that are indistinguishable from real images by human beings. Therefore, consistency between an evaluation metric and human perception is a critical requirement for the metric to be considered good. Table 5 shows the comparison between synthetic-Neuroscore and three traditional scores. To be consistent with all the scores (smaller score indicates better GAN), we used 1/IS and 1/synthetic-Neuroscore for comparisons in Table 5. It can be seen that people rank the GAN performance as PROGAN > BEGAN > DCGAN. All three synthetic-Neuroscores produced by the three models with EEG are consistent with human judgment while the other three conventional scores are not (they all indicate that DCGAN outperforms BEGAN).

| Metrics | | DCGAN | BEGAN | PROGAN |
|---|---|---|---|---|
| 1/IS | | 0.44 | 0.57 | 0.42 |
| MMD | | 0.22 | 0.29 | 0.12 |
| FID | | 63.29 | 83.38 | 34.10 |
| **Proposed Methods** | 1/Shallow-EEG | **1.60** | **1.39** | **1.14** |
| | 1/Mobilenet-EEG | **1.71** | **1.29** | **1.20** |
| | 1/Inception-EEG | **1.51** | **1.34** | **1.24** |
| Human (BE accuracy) | | **0.995** | **0.824** | **0.705** |

Table 5: Three conventional scores: Inception Score (IS), Maximum Mean Discrepancy (MMD), Fréchet Inception Distance (FID), and synthetic-Neuroscore produced by three models with EEG for each GAN category. A lower score indicates better performance of the GAN. Neuroscore is consistent with human judgments. Bold text indicates the consistency with human judgment (BE) accuracy.

### 5.3. Synthetic-Neuroscore Needs Far Fewer Samples

The number of samples needed for evaluation of a GAN is crucial in real-world applications considering computational efficiency and efforts needed for labeling and annotation. Traditional metrics need a large sample size to capture the underlying statistical properties of the real and generated images [12, 16]. In practice, we should prefer a metric that is robust when dealing with small sample sizes i.e., where small sample sizes can produce good estimates. Figure 10(b) shows that synthetic-Neuroscore converges stably at around 20 presentations of a specific image (for signal-enhancement purposes), which is far fewer than the thousands of images required by traditional methods [14, 16]. This is due to the

fact that the LDA-beamformed single-trial P300 amplitude becomes stable when as few as dozens of EEG trials corresponding to one category are available [20].
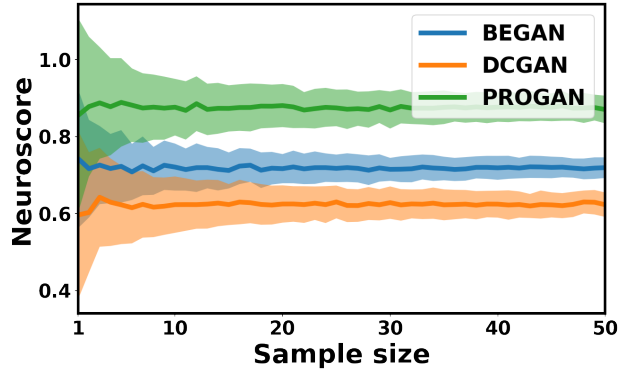


Figure 10: Synthetic-Neuroscore for different evaluated sample sizes for each type of GAN. 200 repeated measurements have been made by randomly shuffling image samples.

## 5.4. Synthetic-Neuroscore Can Rank Images

Another property of using synthetic-Neuroscore is the ability to indicate the quality of an individual image. Traditional evaluation metrics are unable to score each individual image for two reasons. Firstly they need large-scale samples for evaluation and secondly most methods (e.g., MMD and FID) evaluate GANs based on the dissimilarity between real images and generated images so they are not able to score the generated images individually. For our proposed method, the score of each single image can also be evaluated as a synthetic single-trial P300 amplitude measurement. We demonstrate in Figure 11 how the predicted single-trial P300 amplitude conveys perceptual quality at the level of individual images. This property provides synthetic-Neuroscore with a novel capability for tracking variations in image output quality within a typical GAN. Although synthetic-Neuroscore and IS are both generated from deep neural networks, synthetic-Neuroscore is more suitable than IS for evaluating GANs as it is a direct reflection of human perception and fewer sample images are required for evaluation. This has benefits in terms of improved explanation of output than that offered by IS. For example low ranked images can be selected at evaluation time to illustrate cases where the GAN under evaluation is performing poorly.
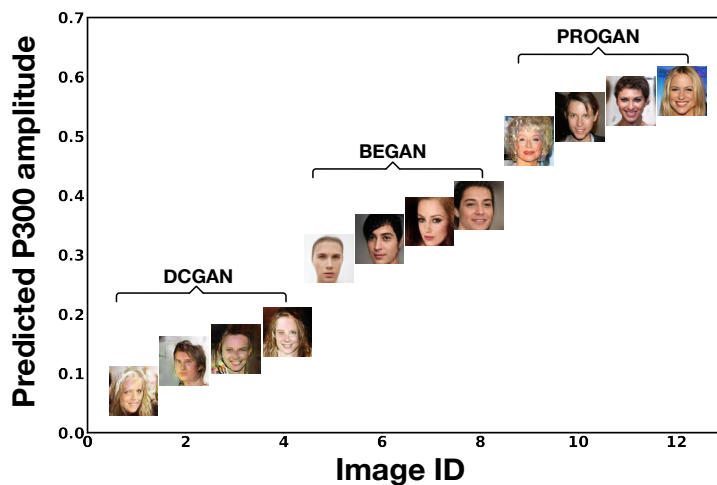
Figure 11: P300 for each single image predicted by the proposed neuro-AI interface in our paper. Higher predicted P300 indicates better image quality.

## 6. Conclusions

In this paper, we outline a metric for evaluating the quality of the outputs from GANs called Neuroscore. Furthermore, we describe a neuro-AI interface to calculate a synthetic-Neuroscore for evaluating GAN performance that only requires EEG signals as supervisory information during model training. Three deep network architectures are explored and our results demonstrate that including neural responses during the training phase of the neuro-AI interface improves its accuracy even though neural measurements are absent when evaluating on a test set. This means that human subjects are not actually needed to evaluate the output from a test GAN, their neural responses are needed only when training the model that produces a synthetic-Neuroscore.

We compared our synthetic-Neuroscore measure to three traditional evaluation metrics and demonstrated the unique advantages of synthetic-Neuroscore, that it is consistent with human perception, that it requires far fewer image samples for calculation and that it can rank individual images in terms of quality, within a specific GAN.

In this work, we demonstrated the use of CNNs to synthesize the neural response. More complicated neural architectures such as mixture of CNNs and recurrent neural networks can be investigated in future work when more EEG data is available.

23

## Acknowledgements

# References

[1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, Y. Bengio, Generative adversarial nets, in: Advances in Neural Information Processing Systems, 2014, pp. 2672–2680.

[2] Z. Wang, Q. She, T. E. Ward, Generative adversarial networks: A survey and taxonomy, arXiv preprint arXiv:1906.01529 (2019).

[3] T. Karras, S. Laine, T. Aila, A style-based generator architecture for generative adversarial networks, arXiv preprint:1812.04948 (2018).

[4] W. Fedus, I. Goodfellow, A. M. Dai, MaskGAN: Better text generation via filling in the _ ., arXiv preprint:1801.07736 (2018).

[5] C. Donahue, J. McAuley, M. Puckette, Synthesizing audio with generative adversarial networks, arXiv preprint:1802.04208 (2018).

[6] E. Brophy, Z. Wang, T. E. Ward, Quick and easy time series generation with established image-based GANs, arXiv preprint arXiv:1902.05624 (2019).

[7] J.-Y. Zhu, T. Park, P. Isola, A. A. Efros, Unpaired image-to-image translation using cycle-consistent adversarial networks, arXiv preprint:1703.10593v6 (2017).

[8] C. Ledig, L. Theis, F. Huszár, J. Caballero, A. Cunningham, A. Acosta, A. Aitken, A. Tejani, J. Totz, Z. Wang, et al., Photo-realistic single image super-resolution using a generative adversarial network, in: 2017 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2017, pp. 105–114.

[9] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, T. S. Huang, Generative image inpainting with contextual attention, arXiv preprint:1801.07892 (2018).

[10] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, S. Hochreiter, GANs trained by a two time-scale update rule converge to a local nash equilibrium, in: Advances in Neural Information Processing Systems, 2017, pp. 6626–6637.

[11] Y. Li, K. Swersky, R. Zemel, Generative moment matching networks, in: International Conference on Machine Learning, 2015, pp. 1718–1727.

[12] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, X. Chen, Improved techniques for training GANs, in: Advances in Neural Information Processing Systems, 2016, pp. 2234–2242.

[13] M. Arjovsky, S. Chintala, L. Bottou, Wasserstein GAN, arXiv preprint:1701.07875 (2017).

[14] A. Borji, Pros and cons of GAN evaluation measures, arXiv preprint:1802.03446 (2018).

[15] P. W. Koh, P. Liang, Understanding black-box predictions via influence functions, in: International Conference on Machine Learning, 2017, pp. 1885–1894.

[16] Q. Xu, G. Huang, Y. Yuan, C. Guo, Y. Sun, F. Wu, K. Q. Weinberger, An empirical study on evaluation metrics of generative adversarial networks, arXiv preprint:1806.07755 (2018).

[17] D. L. Yamins, H. Hong, C. F. Cadieu, E. A. Solomon, D. Seibert, J. J. DiCarlo, Performance-optimized hierarchical models predict neural responses in higher visual cortex, Proceedings of the National Academy of Sciences 111 (23) (2014) 8619–8624.

[18] D. L. Yamins, J. J. DiCarlo, Using goal-driven deep learning models to understand sensory cortex, Nature neuroscience 19 (3) (2016) 356.

[19] S. Waldert, Invasive vs. non-invasive neuronal signals for brain-machine interfaces: Will one prevail?, Frontiers in neuroscience 10 (2016) 295.

[20] A. Mouraux, G. D. Iannetti, Across-trial averaging of event-related EEG responses and beyond, Magnetic resonance imaging 26 (7) (2008) 1041–1054.

[21] N. Lago, A. Cester, Flexible and organic neural interfaces: A review, Applied Sciences 7 (12) (2017) 1292.

[22] R. M. Cichy, A. Khosla, D. Pantazis, A. Torralba, A. Oliva, Comparison of deep neural networks to spatio-temporal cortical dynamics of human visual object recognition reveals hierarchical correspondence, Scientific reports 6 (2016) 27755.

[23] R. M. Cichy, D. Kaiser, Deep neural networks as scientific models, Trends in cognitive sciences (2019).

[24] I. I. Groen, M. R. Greene, C. Baldassano, L. Fei-Fei, D. M. Beck, C. I. Baker, Distinct contributions of functional and deep neural network features to representational similarity of scenes in human brain and behavior, Elife 7 (2018) e32962.

[25] I. Kuzovkin, R. Vicente, M. Petton, J.-P. Lachaux, M. Baciu, P. Kahane, S. Rheims, J. R. Vidal, J. Aru, Activations of deep convolutional neural networks are aligned with gamma band activity of human visual cortex, Communications biology 1 (1) (2018) 107.

[26] T. Tu, J. Koss, P. Sajda, Relating deep neural network representations to EEG-fMRI spatiotemporal dynamics in a perceptual decision-making task, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 2018, pp. 1985–1991.

[27] A. P. Batista, J. J. DiCarlo, Deep learning reaches the motor system, Nature methods 15 (10) (2018) 772.

[28] N. Kriegeskorte, Deep neural networks: a new framework for modeling biological vision and brain information processing, Annual review of vision science 1 (2015) 417–446.

[29] S. R. Kheradpisheh, M. Ghodrati, M. Ganjtabesh, T. Masquelier, Deep networks can resemble human feed-forward vision in invariant object recognition, Scientific reports 6 (2016) 32672.

[30] L. Chelazzi, E. K. Miller, J. Duncan, R. Desimone, A neural basis for visual search in inferior temporal cortex, Nature 363 (6427) (1993) 345.

[31] J. Polich, Updating P300: an integrative theory of P3a and P3b, Clinical Neurophysiology 118 (10) (2007) 2128–2148.

[32] M. Carrillo-De-La-Pena, F. Cadaveira, The effect of motivational instructions on P300 amplitude, Neurophysiologie Clinique/Clinical Neurophysiology 30 (4) (2000) 232–239.

[33] Z. Wang, G. Healy, A. F. Smeaton, T. E. Ward, Spatial filtering pipeline evaluation of cortically coupled computer vision system for rapid serial visual presentation, Brain-Computer Interfaces 5 (4) (2018) 132–145.

[34] Z. Wang, G. Healy, A. F. Smeaton, T. E. Ward, A review of feature extraction and classification algorithms for image RSVP based BCI, Signal Processing and Machine Learning for Brain-machine Interfaces (2018) 243–270.

[35] Z. Wang, G. Healy, A. F. Smeaton, T. E. Ward, Use of neural signals to evaluate the quality of generative adversarial network performance in facial image generation, Cognitive Computation (2019). `doi:https://doi.org/10.1007/s12559-019-09670-y`.

[36] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, L.-C. Chen, Mobilenetv2: Inverted residuals and linear bottlenecks, in: 2018 IEEE Conference on Computer Vision and Pattern Recognition, IEEE, 2018, pp. 4510–4520.

[37] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, Z. Wojna, Rethinking the inception architecture for computer vision, in: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 2818–2826.

[38] A. D. Gerson, L. C. Parra, P. Sajda, Cortically coupled computer vision for rapid image search, IEEE Transactions on Neural Systems and Rehabilitation Engineering 14 (2) (2006) 174–179.

[39] E. Mohedano, K. McGuinness, G. Healy, N. E. O'Connor, A. F. Smeaton, A. Salvador, S. Porta, X. Giró-i Nieto, Exploring EEG for object detection and retrieval, in: Proceedings of the 5th ACM on International Conference on Multimedia Retrieval, ACM, 2015, pp. 591–594.

[40] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, Imagenet: A large-scale hierarchical image database, in: Proceedings of the Conference on Computer Vision and Pattern Recognition, IEEE, 2009, pp. 248–255.

[41] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, A. Smola, A kernel two-sample test, Journal of Machine Learning Research 13 (Mar) (2012) 723–773.

[42] S. Barratt, R. Sharma, A note on the Inception Score, arXiv preprint:1801.01973 (2018).

[43] D. A. Forsyth, J. Ponce, A modern approach, Computer vision: A Modern Approach (2003) 88–101.

[44] S. Schindler, E. Zell, M. Botsch, J. Kissler, Differential effects of face-realism and emotion on event-related brain potentials and their implications for the uncanny valley theory, Scientific reports 7 (2017) 45003.

[45] W. Ling, W. Jingmei, W. Junli, L. Yingjun, A comparative event-related potential study on recognition of cartoon face and real face, Psychological Research 5 (2012).

[46] J. Zhao, Q. Meng, L. An, Y. Wang, An event-related potential comparison of facial expression processing between cartoon and real faces, PloS one 14 (1) (2019) e0198868.

[47] M. S. Treder, A. K. Porbadnigk, F. S. Avarvand, K.-R. Müller, B. Blankertz, The LDA beamformer: Optimal estimation of ERP source time series using linear discriminant analysis, Neuroimage 129 (2016) 279–291.

[48] A. Radford, L. Metz, S. Chintala, Unsupervised representation learning with deep convolutional generative adversarial networks, arXiv preprint arXiv:1511.06434 (2015).

[49] D. Berthelot, T. Schumm, L. Metz, BEGAN: Boundary equilibrium generative adversarial networks, arXiv preprint arXiv:1703.10717 (2017).

[50] T. Karras, T. Aila, S. Laine, J. Lehtinen, Progressive growing of GANs for improved quality, stability, and variation, arXiv preprint arXiv:1710.10196 (2017).

[51] Z. Liu, P. Luo, X. Wang, X. Tang, Deep learning face attributes in the wild, in: IEEE International Conference on Computer Vision (ICCV), 2015.

[52] G. Healy, Z. Wang, C. Gurrin, T. Ward, A. F. Smeaton, An EEG image-search dataset: A first-of-its-kind in IR/IIR. NAILS: Neurally augmented image labelling strategies (2017).

[53] G. Healy, T. E. Ward, C. Gurrin, A. F. Smeaton, Overview of NTCIR-13 NAILS task, in: The 13th NTCIR (2016-2017) Evaluation of Information Access Technologies Conference, Tokyo, Japan, 2017.

[54] Z. Wang, G. Healy, A. F. Smeaton, T. E. Ward, An investigation of triggering approaches for the rapid serial visual presentation paradigm in brain computer interfacing, in: 2016 27th Irish Signals and Systems Conference, IEEE, 2016, pp. 1–6.

[55] G. Healy, Z. Wang, T. Ward, A. Smeaton, C. Gurrin, Experiences and insights from the collection of a novel multimedia EEG dataset, in: International Conference on Multimedia Modeling, Springer, 2020, pp. 475–486.

[56] Q. She, G. Chen, R. H. Chan, Evaluating the small-world-ness of a sampled network: Functional connectivity of entorhinal-hippocampal circuitry, Scientific reports 6 (2016) 21468.

[57] Q. She, Y. Gao, K. Xu, R. H. Chan, Reduced-rank linear dynamical systems, in: Thirty-Second AAAI Conference on Artificial Intelligence, 2018.

[58] L. McIntosh, N. Maheswaranathan, A. Nayebi, S. Ganguli, S. Baccus, Deep learning models of the retinal response to natural scenes, in: Advances in Neural Information Processing Systems, 2016, pp. 1369–1377.

[59] K. H. Kim, J. H. Kim, J. Yoon, K.-Y. Jung, Influence of task difficulty on the features of event-related potential during visual oddball task, Neuroscience letters 445 (2) (2008) 179–183.

[60] A. R. Marathe, A. J. Ries, K. McDowell, A novel method for single-trial classification in the face of temporal variability, in: International Conference on Augmented Cognition, Springer, 2013, pp. 345–352.

[61] D. Senkowski, C. S. Herrmann, Effects of task difficulty on evoked gamma activity and ERPs in a visual discrimination task, Clinical Neurophysiology 113 (11) (2002) 1742–1753.

[62] C. Scharinger, A. Soutschek, T. Schubert, P. Gerjets, Comparison of the working memory load in n-back and working memory span tasks by means of EEG frequency band power and P300 amplitude, Frontiers in human neuroscience 11 (2017) 6.

[63] J. M. Wolfe, E. M. Palmer, T. S. Horowitz, Reaction time distributions constrain models of visual search, Vision research 50 (14) (2010) 1304–1311.

[64] M. G. Philiastides, H. R. Heekeren, P. Sajda, Human scalp potentials reflect a mixture of decision-related signals during perceptual choices, Journal of Neuroscience 34 (50) (2014) 16877–16889.