

Pros and Cons of GAN Evaluation Measures: New Developments

Ali Borji
aliborji@gmail.com

Abstract

This work is an update of a previous paper on the same topic published a few years ago (Borji, 2019). With the dramatic progress in generative modeling, a suite of new quantitative and qualitative techniques to evaluate models has emerged. Although some measures such as Inception Score, Fréchet Inception Distance, Precision-Recall, and Perceptual Path Length are relatively more popular, GAN evaluation is not a settled issue and there is still room for improvement. Here, I describe new dimensions that are becoming important in assessing models (*e.g.* bias and fairness) and discuss the connection between GAN evaluation and deepfakes. These are important areas of concern in the machine learning community today and progress in GAN evaluation can help mitigate them.

Keywords: Generative Models, Generative Adversarial Networks, Neural Networks, Deep Learning, Deepfakes

1. Introduction

Generative models have revolutionized AI and are now being widely adopted¹. They are remarkably effective at synthesizing strikingly realistic and diverse images, and learn to do so in an unsupervised manner (Wang et al., 2019; Liu et al., 2021; Bond-Taylor et al., 2021). Models such as generative adversarial networks (GANs) (Goodfellow et al., 2014), autoregressive models such as PixelCNNs (Oord et al., 2016), variational autoencoders (VAEs) (Kingma and Welling, 2013), and recently transformer GANs (Hudson and Zitnick, 2021; Jiang et al., 2021) have been constantly improving the state of the art in image generation. See for example Brock et al. (2018); Karras et al. (2019); Ramesh et al. (2021) and Razavi et al. (2019).

A key difficulty in generative modeling is evaluating performance, *i.e.* how good is a model in approximating a data distribution? Quantitative evaluation²

¹For a list of GAN applications, please refer to [here](#).

²Generative models are commonly evaluated in terms of fidelity (how realistic a generated image is) and diversity (how well generated samples capture the variations in real data) of the learned distribution. Generative models, however, can also be tested against other criteria such

of generative models of images and video is an open problem. Here, I give a summary and discussion of the latest progress, since [Borji \(2019\)](#), in this field covering newly proposed measures, benchmarks, and GAN visualization and error diagnosis techniques.

1.1. Background

Evaluation of generative models is an active area of research. Several quantitative and qualitative measures have been proposed so far. In this section, I review some of the influential measures that have inspired the researchers.

A classic approach is to compare the log-likelihoods of models. This approach, however, has several shortcomings. A model can achieve high likelihood, but low image quality, and conversely, low likelihood and high image quality. Also, kernel density estimation in high-dimensional spaces is very challenging ([Theis et al., 2015](#)).

The two most common GAN evaluation measures are Inception Score (IS) and Fréchet Inception Distance (FID). They rely on a pre-existing classifier (InceptionNet) trained on ImageNet. IS ([Salimans et al., 2016](#)) computes the KL divergence between the conditional class distribution and the marginal class distribution over the generated data. FID ([Heusel et al., 2017](#)) calculates the Wasserstein-2 (*a.k.a* Fréchet) distance between multivariate Gaussians fitted to the embedding space of the Inception-v3 network of generated and real images³.

IS does not capture intra-class diversity, is insensitive to the prior distribution over labels (hence is biased towards ImageNet dataset and Inception model⁴), and is very sensitive to model parameters and implementations ([Barratt and Sharma, 2018](#); [Odena, 2019](#)). It requires a large sample size to be reliable. A simple class-conditional model that memorizes one example per ImageNet class achieves a high IS. Using a 1D example, [Barratt and Sharma \(2018\)](#) show that the true underlying data distribution may achieve a lower IS than other distributions.

FID has been widely adopted because of its consistency with human inspection and sensitivity to small changes in the real distribution (*e.g.* slight blurring or small artifacts in synthesized images). FID, unlike IS, can detect intra-class mode collapse⁵. The Gaussian assumption in FID calculation, however, might not hold in practice. A major drawback with FID is its high bias. The sample size to calculate FID has to be large enough (usually above 50K). Smaller sample sizes can lead to over-estimation of the actual FID ([Chong and Forsyth, 2020](#)).

While single-value metrics like IS and FID successfully capture many aspects of the generator, they lump quality and diversity evaluation and are therefore not ideal for diagnostic purposes (Precision-Recall addresses this shortcoming).

as linear separability, mode collapse, etc. For a complete list of desired properties of GANs, please see [Borji \(2019\)](#).

³To learn how to implement FID, please refer to [here](#).

⁴See also [here](#).

⁵Mode collapse happens when a generator produces data from only a few modes of the target distribution, thus failing to generate diverse enough outputs. Please see [here](#).

Further, they continue to have a blind spot for image quality (Karras et al., 2019).

Another popular approach to evaluate GANs is manual inspection (*e.g.* Denton et al. (2015); Zhou et al. (2019)). While this approach can directly tell us about the image quality, it is limited in assessing sample diversity. It is also subjective as the reviewers may incorporate their own biases and opinions. It also requires knowledge of what is realistic and what is not for the target domain which might be hard to learn in some domains (*e.g.* medical domain⁶). Further, it is constrained by the number of images that can be reviewed in a reasonable time prohibiting its usage during model development. Finally, due to the subtle differences in experimental protocols (*e.g.* user interface design, fees, duration), it is often difficult to replicate the results across different publications.

1.2. Organization of the Paper

This paper is organized as follows. Recently proposed evaluation measures are categorized under quantitative and qualitative measures and are covered in sections 2 and 3, respectively⁷. In section 4, I then discuss some benchmark and model comparison studies as well as analysis works. Finally, section 5 highlights some major points and offers suggestions for future research in GAN evaluation. Note that some of the works reviewed here have not been accepted in prior published work and may have not been peer reviewed. Therefore, they should be treated with additional caution in the absence of widespread usage.

2. New Quantitative GAN Evaluation Measures

2.1. Specialized Variants of Fréchet Inception Distance and Inception Score

2.1.1. Spatial FID (sFID)

Nash et al. (2021) propose a variant of FID called sFID that uses spatial features rather than the standard pooled features. They compute FID using both the standard pool3 inception features and the first 7 channels from the intermediate mixed 6/conv feature maps. The reason to include pool3 is because it compresses spatial information to a large extent, making it less sensitive to spatial variability. The reason to include the intermediate mixed 6/conv is because it provides a sense of spatial distributional similarity between models.

2.1.2. Class-aware FID (CAFD) and Conditional FID

Liu et al. (2018) argue against the single-manifold Gaussian assumption in FID, and employ a Gaussian mixture model (GMM) to better fit the feature

⁶For example, to tell whether a GAN can generate X-rays with pneumonia in them, it is better to ask doctors rather than Amazon Mechanical Turk workers!

⁷Some methods may fall in both categories.

distribution and to include class information. They compute Fréchet distance in each of the K classes and average the results to obtain CAFD:

$$CAFD(P_r, P_g) = \frac{1}{K} \sum_{i=1}^K \|\mu_i^r - \mu_i^g\| + \text{Tr}(C_i^r + C_i^g - 2(C_i^r C_i^g)^{\frac{1}{2}}),$$

where μ_i and C_i are the mean and the covariance matrix of class i , respectively.

Soloveitchik et al. (2021) introduce a conditional variant of FID for evaluating conditional generative models. Instead of having two distributions, they consider two classes of distributions conditioned on a common input.

2.1.3. Fast FID

Mathiasen and Hvilshøj (2020) propose a method to speed up the FID computation in order to make it computationally tractable such that it can be used as a loss function for training GANs. The real data samples do not change during training, so their Inception encodings need to be computed once. Therefore, reducing the number of fake samples lowers the time to compute the Inception encodings. The bottleneck, however, is computing $\text{tr}(\sqrt{\Sigma_1 \Sigma_2})$ in the FID formula. Σ_1 and Σ_2 are the covariance matrices of Gaussians fitted to the real and the generated data. FastFID circumvents this issue without explicitly computing $\sqrt{\Sigma_1 \Sigma_2}$. Depending on the number of fake samples, FastFID can speed up FID computation 25 to 500 times. Notice that FastFID makes backpropagating over FID faster, and as such it is not a metric, rather an optimization method.

2.1.4. Memorization-informed FID (MiFID)

Bai et al. (2021) conduct the first generative model competition⁸, where participants were invited to generate realistic dog images given 20,579 images of dogs from ImageNet. They modified the FID score to penalizes models producing images too similar to the training set as

$$\text{MiFID}(S_g, S_t) = m_\tau(S_g, S_t) \cdot \text{FID}(S_g, S_t)$$

where S_g is the generated set and S_t is the original training set. m_τ is the memorization penalty which is based on thresholding the memorization distance s of generated and true distribution defined as:

$$s(S_g, S_t) = \frac{1}{|S_g|} \sum_{x_g \in S_g} \min_{x_t \in S_t} \left(1 - \frac{|\langle x_g, x_t \rangle|}{|x_g| \cdot |x_t|} \right)$$

$$m_\tau(S_g, S_t) = \begin{cases} \frac{1}{s(S_g, S_t) + \epsilon} & (\epsilon \ll 1), \text{ if } s(S_g, S_t) < \tau \\ 1, & \text{otherwise} \end{cases}$$

Intuitively, lower memorization distance is associated with more severe training sample memorization.

⁸<https://www.kaggle.com/c/generative-dog-images>

2.1.5. Unbiased FID and IS

Chong and Forsyth (2020) show that FID and IS are biased in the sense that their expected value, computed using a finite sample set, is not their true value. They find that bias depends on the number of images used for calculating the score and the generators themselves, making objective comparisons between models difficult. To mitigate the issue, they propose an extrapolation approach to obtain a bias-free estimate of the scores, called \overline{FID}_∞ and \overline{IS}_∞ , computed with an infinite number of samples. These extrapolated scores are simple, drop-in replacements for the finite sample scores⁹.

2.1.6. Clean FID

Parmar et al. (2021) examine the sensitivity of FID score (and also other scores such as KID) to inconsistent and often incorrect implementations across different image processing libraries when training and evaluating generative models. FID score is widely used to evaluate generative models, but each FID implementation uses a different low-level image processing method. Image resizing functions in commonly-used deep learning libraries often introduce aliasing artifacts. They observe that numerous subtle choices need to be made for FID calculation and a lack of consistencies in these choices can lead to vastly different FID scores. To address these challenges, they introduce a standardized protocol, and provide an easy-to-use FID evaluation library called cleanfid¹⁰.

2.1.7. Fréchet Video Distance (FVD)

Unterthiner et al. (2019) propose an extension of the FID, called FVD, to evaluate generative models of video. In addition to the quality of each frame, FVD captures the temporal coherence in a video. To obtain a suitable feature video representation, they used a pre-trained network (the Inflated 3D Convnet (Carreira and Zisserman, 2017) pre-trained on Kinetics-400 and Kinetics-600 datasets) that considers the temporal coherence of the visual content across a sequence of frames. The I3D network generalizes the Inception architecture to sequential data, and is trained to perform action-recognition on the Kinetics dataset consisting of human-centered YouTube videos. FVD considers a distribution over videos, thereby avoiding the drawbacks of frame-level measures.

Some other measures for video synthesis evaluation include Average Content Distance (Tulyakov et al., 2018) (the average L2 distance among all consecutive frames in a video), Cumulative Probability Blur Detection (CPBD) (Narvekar and Karam, 2009), Frequency Domain Blurriness Measure (FDBM) (De and Masilamani, 2013), Pose Error (Yang et al., 2018), and Landmark Distance (LMD) (Chen et al., 2018). For more details on video generation and evaluation, please consult Oprea et al. (2020).

In addition to above, Fréchet Audio Distance (FAD) (Roblek et al., 2019) and Fréchet ChemNet Distance (FCD) (Preuer et al., 2018) have been suggested

⁹https://github.com/mchong6/FID_IS_infinity

¹⁰github.com/GaParmar/clean-fid

for evaluation of music enhancement and molecule generation algorithms.

2.2. Methods based on Self-supervised Learned Representations

FID scores computed with classification-pretrained embeddings have been shown to correlate well with human evaluations (Heusel et al., 2017). Meanwhile, they can be misleading as they are biased towards ImageNet. On non-ImageNet datasets, FID can result in inadequate evaluation. Morozov et al. (2020) advocate for using self-supervised representations to evaluate GANs on the established non-ImageNet benchmarks¹¹. They show that representations, typically obtained via contrastive or clustering-based approaches, provide better transfer to new tasks and domains, and thus can serve as more universal embeddings of natural images. Further, they demonstrate that self-supervised representations produce a more reasonable ranking of models and often improve the sample efficiency of FID.

2.3. Methods based on Analysing Data Manifold

These methods measure disentanglement in the representations learned by generative models, and are useful for improving the generalization, robustness, and interpretability of models.

2.3.1. Local Intrinsic Dimensionality (LID)

Barua et al. (2019) introduce a new evaluation measure called “CrossLID” to assess the local intrinsic dimensionality (LID) of real-world data with respect to neighborhoods found in GAN-generated samples. Intuitively, CrossLID measures the degree to which the manifolds of two data distributions coincide with each other. The idea behind CrossLID is depicted in Fig. 1. Barua *et al.* compare CrossLID with other measures and show that it a) is strongly correlated with the progress of GAN training, b) is sensitive to mode collapse, c) is robust to small-scale noise and image transformations, and 4) is robust to sample size. It is not clear whether this measure can be applied to complex and high dimensional data where it is hard to define local dimensionality.

2.3.2. Intrinsic Multi-scale Distance (IMD)

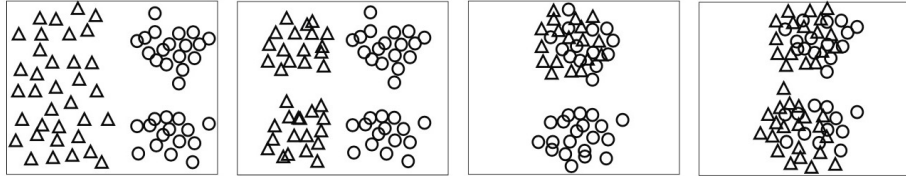
Tsitsulin et al. (2019) argue that current evaluation measures only reflect the first two (mean and covariance) or three moments of distributions. As illustrated in Fig. 2, FID and KID (Kernel Inception Distance (Bińkowski et al., 2018))¹² are insensitive to the global structure of the data distribution. They propose a measure, called IMD, to take all data moments into account, and claim that IMD is intrinsic¹³ and multi-scale¹⁴. Tsitsulin *et al.* experimentally show that their method is effective in discerning the structure of data manifolds even on

¹¹<https://github.com/stanis-morozov/self-supervised-gan-eval>

¹²KID is referred to as MMD in Borji (2019).

¹³It is invariant to isometric transformations of the manifold such as translation or rotation.

¹⁴It captures both local and global information.



(a) CrossLID=15.16 (b) CrossLID=7.33 (c) CrossLID=4.78 (d) CrossLID=2.10

Figure 1: Four 2D examples illustrating how generated samples by a GAN (triangles) relate to real data samples from a bimodal Gaussian (circles), together with CrossLID scores. (a) generated data distributed uniformly, spatially far from the real data. (b) generated data with two modes, spatially far from the real data (c) generated data associated with only one mode of the real data. (d) generated data associated with both modes of the real data (the desired situation). The lower the CrossLID, the better. Figure compiled from [Barua et al. \(2019\)](#).

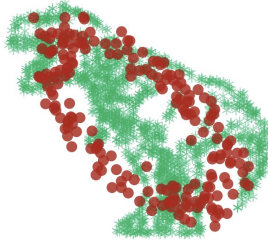


Figure 2: The motivation behind the Intrinsic Multi-scale Distance (IMD) score. Here, the two distributions have the same first 3 moments. Since FID and KID are insensitive to higher moments, they can not distinguish the two distributions, whereas the IMD score can. Figure from [Tsitsulin et al. \(2019\)](#).

unaligned data. IMD compares data distributions based on their geometry and is somewhat similar to the Geometry Score ([Khruikov and Oseledets, 2018](#)).

[Barannikov et al. \(2021\)](#) introduce Cross-Barcode(P,Q), given a pair of distributions in a high-dimensional space, tracks multiscale topology discrepancies between manifolds on which the distributions are concentrated. Based on the Cross-Barcode, they introduce the Manifold Topology Divergence score (MTop-Divergence) and apply it to assess the performance of deep generative models in various domains: images, 3D-shapes, and time-series. Their proposed method is domain agnostic and does not rely on pre-trained networks

2.3.3. Perceptual Path Length (PPL)

First introduced in StyleGAN3 ([Karras et al., 2019](#)), PPL measures whether and how much the latent space of a generator is entangled (or if it is smooth and

the factors of variation are properly separated¹⁵). Intuitively, a less curved latent space should result in perceptually smoother transition than a highly curved latent space (See Fig. 3). Formally, PPL is the empirical mean of the perceptual difference between consecutive images in the latent space \mathcal{Z} , over all possible endpoints:

$$l_{\mathcal{Z}} = \mathbb{E}\left[\frac{1}{\epsilon^2}d(G(\text{slerp}(\mathbf{z}_1, \mathbf{z}_2; t)), G(\text{slerp}(\mathbf{z}_1, \mathbf{z}_2; t + \epsilon)))\right],$$

where $\mathbf{z}_1, \mathbf{z}_2 \sim P(\mathbf{z}), t \sim U(0, 1)$, G is the generator, and $d(\cdot, \cdot)$ is the perceptual distance between the resulting images. `slerp` denotes the spherical interpolation (Shoemake, 1985), and $\epsilon = 10^{-4}$ is the step size. The LPIPS (learned perceptual image patch similarity) (Zhang et al., 2018) can be used to measure the perceptual distance between two images. LPIPS is a weighted L2 difference between two VGG16 (Simonyan and Zisserman, 2014) embeddings, where the weights are learned to make the metric agree with human perceptual similarity judgments. PPL captures semantics and image quality as shown in Fig. 4. Fig. 5 illustrates the superiority of PPL over FID and Precision-Recall scores in comparing two models.

2.3.4. Linear Separability in Latent Space

Karras et al. (2019) propose another measure in StyleGAN3 to quantify the latent space disentanglement by measuring how well the latent-space points can be separated into two distinct sets via a linear hyperplane. They argue that if a latent space is sufficiently disentangled, it should be possible to find directions that consistently correspond to individual factors of variation.

First, some auxiliary classification networks for a number of binary attributes are trained (*e.g.* to distinguish male and female faces). To measure the separability of an attribute:

1. Generate 200K images with $\mathbf{z} \sim P(\mathbf{z})$ and classify them using an auxiliary classification network with label Y ,
2. Sort the samples according to classifier confidence and remove the least confident half, yielding 100K labeled latent-space vectors,
3. Fit a linear SVM to predict the label X based only on the latent-space points and classify the points by this plane,
4. Compute the conditional entropy $H(Y|X)$ where X represents the classes predicted by the SVM and Y represents the classes determined by the classifier. A low value suggests consistent latent space directions for the corresponding factor(s) of variation.

The final separability score is $\exp(\sum_i H(Y_i|X_i))$, where i enumerates over a set of attributes (*e.g.* $i = 40$ over CelebA dataset).

¹⁵For example, features that are absent in either endpoint may appear in the middle of a linear interpolation path between two random inputs.

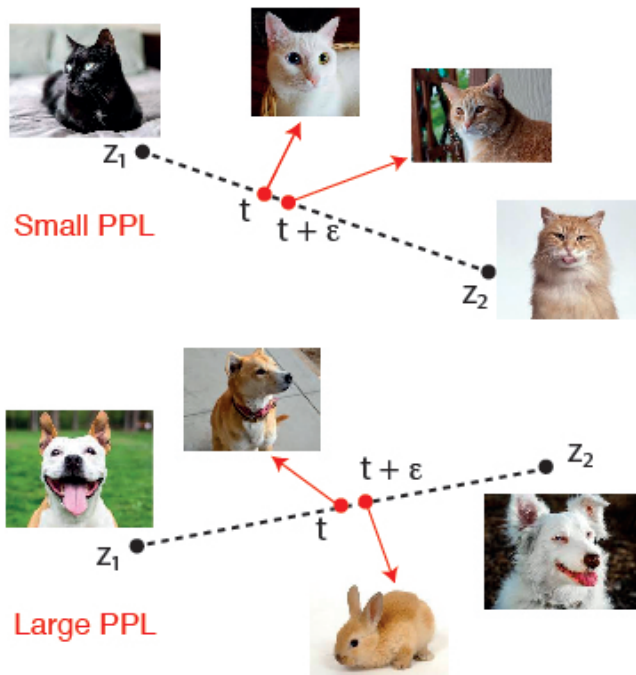


Figure 3: Illustration of PPL. Small perturbations around a point result in smaller average perceptual distance for a model that has learned a disentangled representation.

2.4. Classification Accuracy Score (CAS)

Ravuri and Vinyals (2019) propose a measure that is in essence similar to the “Classification Performance” discussed in Borji (2019). They argue that if a generative model is learning the data distribution in a perceptually meaningful space then it should perform well in downstream tasks. They use class-conditional generative models from a number of generative models such as GANs and VAEs to infer the class labels of real data. They then train an image classifier using only synthetic data and use it to predict labels of real images in the test set. Their findings are as follows: a) when using a state-of-the-art GAN (BigGAN-deep by Brock et al. (2018)), Top-1 and Top-5 accuracies decrease by 27.9% and 41.6%, respectively, compared to the original data, b) CAS automatically identifies particular classes for which generative models fail to capture the data distribution (Fig. 6), and c) IS and FID are neither predictive of CAS, nor useful when evaluating non-GAN models.

2.5. Non-Parametric Tests to Detect Data-Copying

Meehan et al. (2020) formalize a notion of overfitting called data-copying where a generative model memorizes the training samples or their small variations (Fig. 7). They provide a three sample test for detecting data-copying that uses

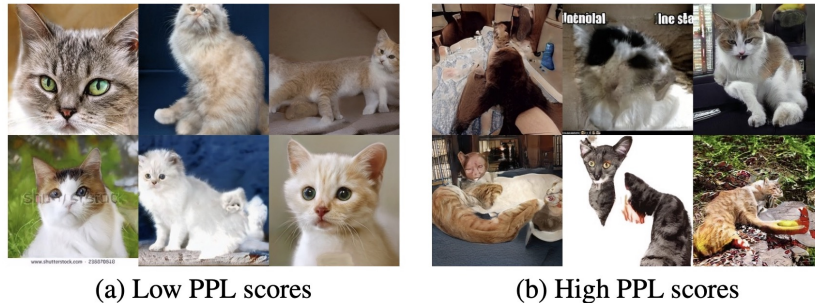


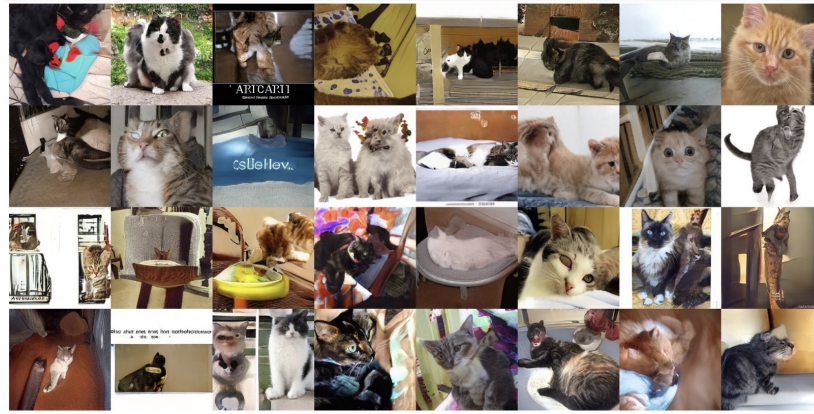
Figure 4: Correlation between PPL and image quality (generated by StyleGAN). Panels (a) and (b) show random examples with low and high (per-image) PPL. PPL score is capable of capturing the consistency of the images. The lower the PPL, the better. Figure from [Karras et al. \(2020\)](#).

the training set, a separate held-out sample from the target distribution, and a generated sample from the model. The key insight is that an overfitted GAN generates samples that are on average closer to the training samples, therefore the average distance of generated samples to training samples is smaller than the corresponding distances between held-out test samples and the training samples. They also divide the instance space into cells and conduct their test separately in each cell. This is because, as they argue, generative models tend to behave differently in different regions of space.

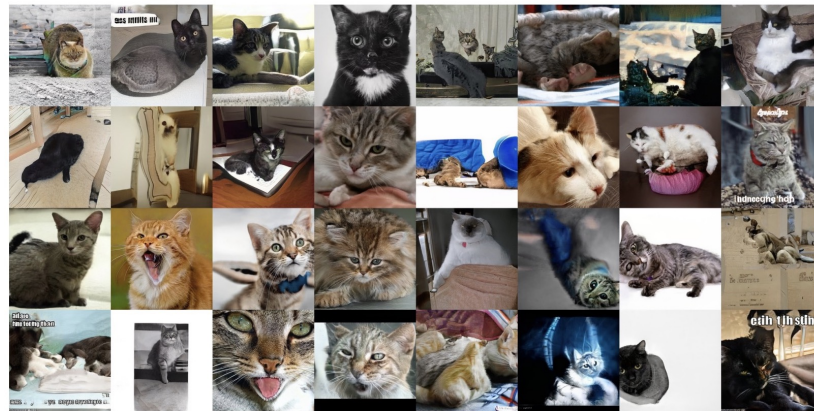
2.6. Measures that Probe Generalization in GANs

[Zhao et al. \(2018\)](#)¹⁶ utilize carefully designed training datasets to characterize how existing models generate novel attributes and their combinations. Some of their findings are as follows. When presented with a training set with all images having exactly 3 objects, both GANs and VAEs typically generate 2-5 objects (Fig. 8.B). Over a multi-modal training distribution (*e.g.* images with either 2 or 10 objects), the model acts as if it is trained separately on each mode, and it averages the two modes. When the modes are close to each other (*e.g.* 2 and 4 objects), the learned distribution assigns higher probability to the mean number of objects (3 in this example), even though there was no image with 3 objects in the training set (Fig. 8.C). When the training set contains certain combinations (*e.g.* red cubes but not yellow cubes; Fig. 8.D), the model memorizes the combinations in the training set when it contains a small number of them (*e.g.* 20), and generates novel combinations when there is more variety (*e.g.* 80). [Xuan et al. \(2019\)](#) made a similar observation shown in Fig. 9. On a geometric-object training dataset where the number of objects for each training image is fixed, they found a variable number of objects in generated images.

¹⁶Please see also the subsection on “Evaluating Mode Drop and Mode Collapse” in [Borji \(2019\)](#).



Model 1: FID=8.53, P=0.64, R=0.28, PPL=924



Model 2: FID=8.53, P=0.62, R=0.29, PPL=387

Figure 5: Synthesized samples from two generative models trained on LSUN (Yu et al., 2015) CAT without truncation (*i.e.* drawing latent vectors from a truncated or shrunk sampling space to improve average image quality, at the expense of losing diversity (Kingma and Welling, 2013; Brock et al., 2018; Karras et al., 2020)). FID, precision (P), and recall (R) are similar for the two models, even though the latter produces cat-shaped objects more often. Perceptual path length (PPL) shows a clear preference for model 2. Figure from Karras et al. (2020).

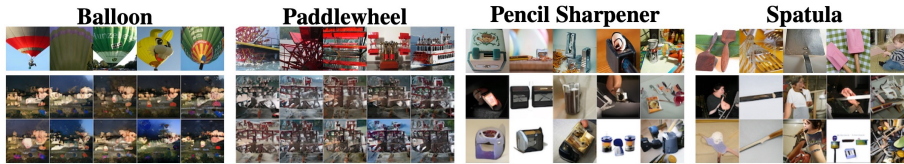


Figure 6: CAS can identify classes for which a generative model, here BigGAN-deep (Brock et al., 2018), fails to capture the data distribution (top row: real images, bottom two rows: generated samples). Figure compiled from Ravuri and Vinyals (2019).

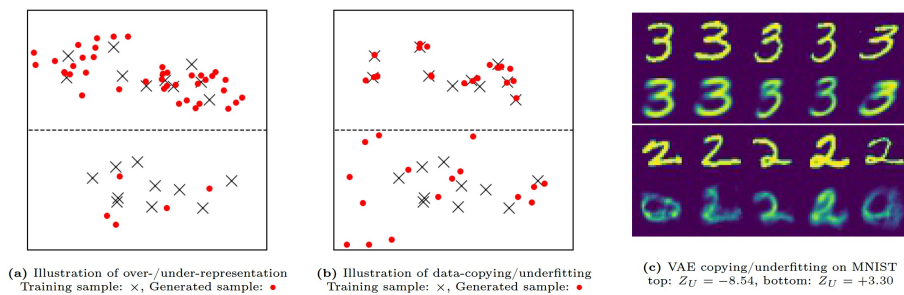


Figure 7: Illustration of the data-copying concept introduced in Meehan et al. (2020). Each panel depicts a single instance space partitioned into two regions. Panel (a) shows an over-represented region (top) and an under-represented region (bottom). This is the kind of overfitting evaluated by measures such as FID and Precision-Recall. Panel (b) shows a data-copied region (top) and an underfitted region (bottom). Panel (c) shows VAE-generated and training samples from a data-copied (top) and underfitted (bottom) region over MNIST. In each 10-image strip, the bottom row provides random generated samples from the region and the top row shows their training nearest neighbors. Samples in the bottom region are on average further to their training nearest neighbor than held-out test samples in the region, and samples in the top region are closer, and thus “data-copying”. Figure from Meehan et al. (2020).

Some other studies that have investigated GAN generalization include O’Brien et al. (2018); van Steenkiste et al. (2020).

2.7. New Ideas based on Precision and Recall ($P\&R$)

Sajjadi et al. (2018) propose to use precision and recall to explicitly quantify the trade off between quality (precision) and coverage (recall). Precision measures the similarity of generated instances to the real ones and recall measures the ability of a generator to synthesize all instances found in the training set (Fig. 10). P&R curves can distinguish mode-collapse (poor recall) and bad quality (poor precision). Some new ideas based on P&R are summarized below.

2.7.1. Density and Coverage

Naeem et al. (2020) argue that even the latest version of the precision and recall metrics are still not reliable as they a) fail to detect the match between two

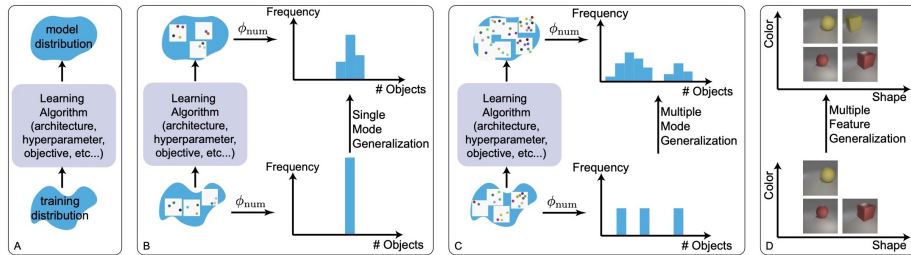


Figure 8: A) A generative model can be probed with carefully designed training data. Examining the learned distribution when training data B) takes a single value for a feature (*e.g.* all training images have 3 objects), C) has multiple modes for a feature (*e.g.* all training images have 2, 4 or 10 objects), or D) has multiple modes over multiple features. Figure compiled from Zhao et al. (2018).

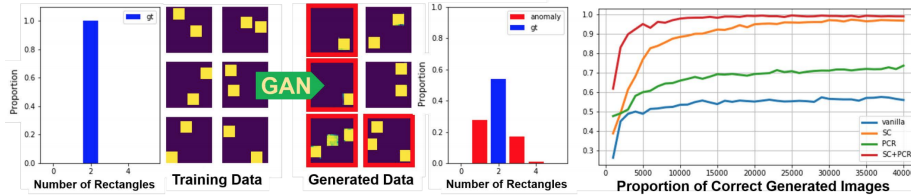


Figure 9: Xuan et al. (2019) show that training a GAN over images with exactly two rectangles results in a model that generates one, two, or three rectangles (anomalous ones shown in red). They also propose a model that generates a high fraction of correct images (the right-most panel).

identical distributions, b) are not robust against outliers, and c) have arbitrarily selected evaluation hyperparameters. To solve these issues, they propose density and coverage metrics. Precision counts the binary decision of whether the fake data Y_j is contained in any neighbourhood sphere. Density, instead, counts how many real-sample neighbourhood spheres contain Y_j . See Fig. 11. They analytically and experimentally show that density and coverage provide more interpretable and reliable signals for practitioners than the existing measures¹⁷.

2.7.2. Alpha Precision and Recall

Alaa et al. (2021) introduce a 3-dimensional evaluation metric, (α -Precision, β -Recall, Authenticity), to characterizes the fidelity, diversity and generalization power of generative models (Fig. 12). The first two assume that a fraction $1 - \alpha$ (or $1 - \beta$) of the real (and synthetic) data are “outliers”, and α (or β) are “typical”. α -Precision is the fraction of synthetic samples that resemble the “most typical” α real samples, whereas β -Recall is the fraction of real samples

¹⁷<https://github.com/clovaai/generative-evaluation-prdc>

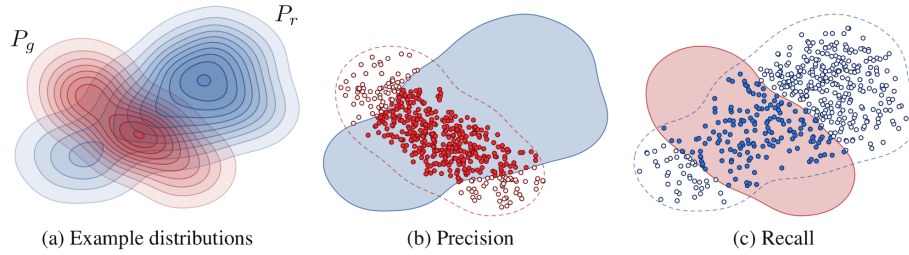


Figure 10: (a) Illustration of precision-recall for distribution of real images P_r (blue) and the distribution of generated images P_g (red). (b) Precision is the probability that a random image from P_g falls within the support of P_r . (c) Recall is the probability that a random image from P_r falls within the support of P_g . Figure compiled from [Kynkäänniemi et al. \(2019\)](#).

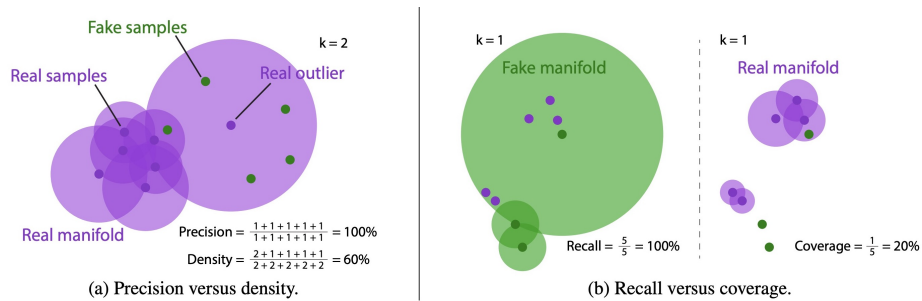


Figure 11: Pictorial depiction of density and coverage measures ([Naeem et al., 2020](#)) (See text for definitions). Note that in the recall versus coverage figure (panel b), the real and fake samples are identical across left and right. a) Here, the real manifold is overestimated due to the real outlier sample. Generating many fake samples around the real outlier increases the precision measure. b) Here, although the fake samples are far from the modes in real samples, the recall is perfect because of the overestimated fake manifold. Figure compiled from [Naeem et al. \(2020\)](#).

covered by the most typical β synthetic samples. α -Precision and β -Recall are evaluated for all $\alpha, \beta \in [0, 1]$, providing entire precision and recall curves instead of a single number. To compute both metrics, the (real and synthetic) data are embedded into hyperspheres with most samples concentrated around the centers. Typical samples are located near the centers whereas outliers are close to the boundaries. To quantify Generalization they introduce the Authenticity metric to measure the probability that a synthetic sample is copied from the training data. Some other works that have proposed extensions to P&R include [Djologua et al. \(2020\)](#); [Simon et al. \(2019\)](#); [Kynkäänniemi et al. \(2019\)](#).

2.8. Duality GAP Metric

[Grnarova et al. \(2019\)](#) leverage the notion of duality gap from game theory to propose a domain-agnostic and computationally-efficient measure that can be used to assess different models as well as monitoring the progress of a single model

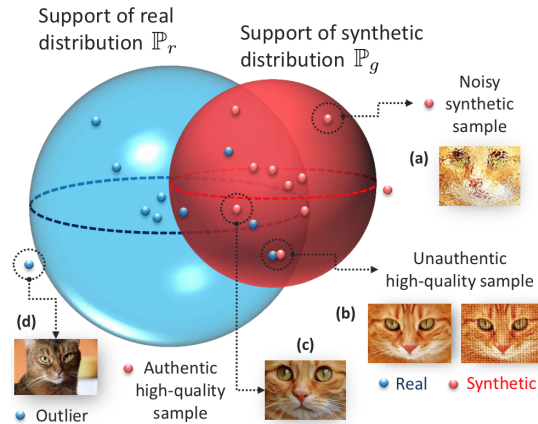


Figure 12: Illustration of α -Precision, β -Recall and Authenticity metrics (Alaa et al., 2021). Blue and red spheres correspond to the α and β -supports of real and generative distributions, respectively. Blue and red points correspond to real and synthetic data. (a) Generated samples falling outside the blue sphere look unrealistic or noisy. (b) Overfitted models can generate high-quality samples that are “unauthentic” because they are copied from training data. (c) High-quality samples should reside in the blue sphere. (d) Outliers do not count in the β -Recall metric. (Here, $\alpha=\beta=0.9$, α -Precision=8/9, β -Recall=4/9, Authenticity=9/10). Figure from Alaa et al. (2021).

throughout training. Intuitively, duality gap measures the sub-optimality (*w.r.t.* an equilibrium) of a given solution (G,D) where G and D are the generator and the discriminator, respectively. They also show that their measure highly correlates with FID on natural image datasets, and can also be used in other modalities such as text and sound. Further, their measure requires no labels or a pretrained classifier, making it domain agnostic. In a follow-up work, Sidheekh et al. (2021) extend the notion of duality gap to proximal duality gap that is applicable to the general context of training GANs.

2.9. Spectral Methods

A number of studies (*e.g.* Durall et al. (2020); Frank et al. (2020); Dzanic et al. (2019); Zeng et al. (2017)) have shown that current generators are unable to correctly approximate the spectral distributions of real data. Durall et al. (2020) showed that up-scaling operations commonly used in GANs (*e.g.* up-convolutions) alter the spectral properties of the images causing high frequency distortions in the output. This can be seen in the average azimuthal integration of the power spectrum shown in the top panel of Fig. 13. Based on this observation, they proposed a novel spectral regularization term to compensate spectral distortions. They also proposed a very simple but highly accurate detector for generated images and videos, *i.e.* a DeepFake detector.

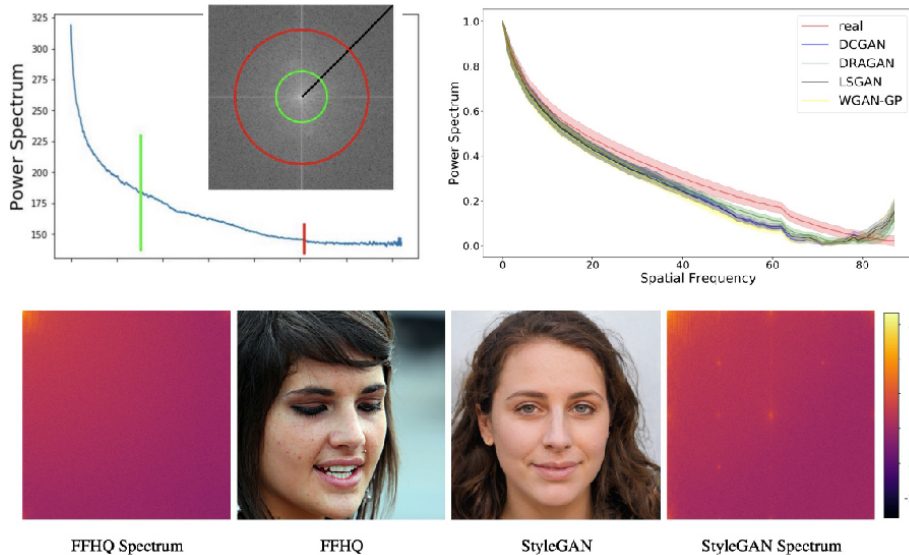


Figure 13: Top-left) An example of azimuthal integral for an image (inset) taken from [Durrall et al. \(2020\)](#). In the 1D Power Spectrum, each frequency component is the radial integral over the 2D spectrum (shown in red and green). Top-right) Statistics (mean and variance) after azimuthal integration over the power-spectrum of real and GAN generated images over CelebA dataset ([Liu et al., 2015](#)). Notice in particular the sharp increase at higher frequencies. Bottom) A side-by-side comparison of real and generated faces in spatial and frequency domains. The left-most panel shows that mean DCT spectrum of the FFHQ with a sample from this dataset next to it. The right-most side shows the mean DCT spectrum of a dataset sampled from StyleGAN trained on FFHQ, with a generated face to its left. Results are averaged over 10K images. Image from [Frank et al. \(2020\)](#).

2.10. Caption Score (CapS)

Almost all existing metrics are solely for evaluating models that generative image. They do not take into account the corresponding text in the context of text-to-image generation. Thus, they may be fooled by a network that ignores the textual input and only focuses on generating realistic looking images. To remedy this shortcoming, [Ding et al. \(2021\)](#) introduce the caption score to evaluate the correspondence between images and text. This score measures the quality and accuracy for text-image generation at a finer granularity than FID and Inception Score (IS), and is defined as:

$$\text{CapS}(x, t) = \sqrt{|t| \prod_{i=0}^{|t|} p(t_i | x, t_{0:i-1})},$$

where t is a sequence of text tokens and x is the image. The higher the CapS, the better. In addition to CapS, they also report results using other measures such as FID and IS.

2.11. Perplexity

A commonly used metric to evaluate generative models of text is the perplexity. It measures the probability for a sentence to be produced by a language model that has been trained on a dataset:

$$\text{perplexity} = \prod_{t=1}^T \left(\frac{1}{P_{\text{LM}}(\mathbf{x}^{(t+1)}|\mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T}$$

where \mathbf{x} represent a single word and T is the sentence length. The lower the perplexity value, the better the model (averaged over a set of test sentences). See [here](#), [here](#), [Iqbal and Qureshi \(2020\)](#), and [Tevet et al. \(2018\)](#) for some other measures for evaluating text generation models.

3. New Qualitative GAN Evaluation Measures

A number of new qualitative measures have also been proposed. These measures typically focus on how convincing a generated image is from human perception perspective.

3.1. Human Eye Perceptual Evaluation (HYPE)

[Zhou et al. \(2019\)](#) propose a human-in-the-loop evaluation approach that is grounded in psychophysics research on visual perception (Fig. 14). They introduce two variants. The first one measures visual perception under adaptive time constraints to determine the threshold at which a model’s outputs appear real (*e.g.* 250 ms). The second variant is less expensive and measures human error rate on fake and real images without time constraints¹⁸. [Kolchinski et al. \(2019\)](#) propose an approach to approximate HYPE and report 66% accuracy in predicting human scores of image realism, matching the human inter-rater agreement rate. A major drawback with human evaluation approaches such as HYPE is scaling them. One way to remedy this is to train a model from human judgments and interact with a person only when the model is not certain. HYPE is more reliable compared to automated ones but cannot be used for monitoring the training process.

3.2. Neuroscore

[Wang et al. \(2020b\)](#) outline a method called Neuroscore using neural signals and rapid serial visual presentation (RSVP) to directly measure human perceptual response to generated stimuli. Participants are instructed to attend to target images (real and generated) amongst a larger set of non-target images (left panel in Fig. 15). This paradigm, known as the oddball paradigm, is commonly used to elicit the P300 event-related potential (ERP), which is a positive voltage deflection typically occurring between 300 ms and 600 ms after the appearance

¹⁸An important point when conducting experiments of this sort (*e.g.* [Denton et al. \(2015\)](#)) is to make sure that real and generated images are shuffled during the presentation. Otherwise, subjects maybe able to guess what type of images are shown in a session from the first few images, if all the images in the session in either real or fake.

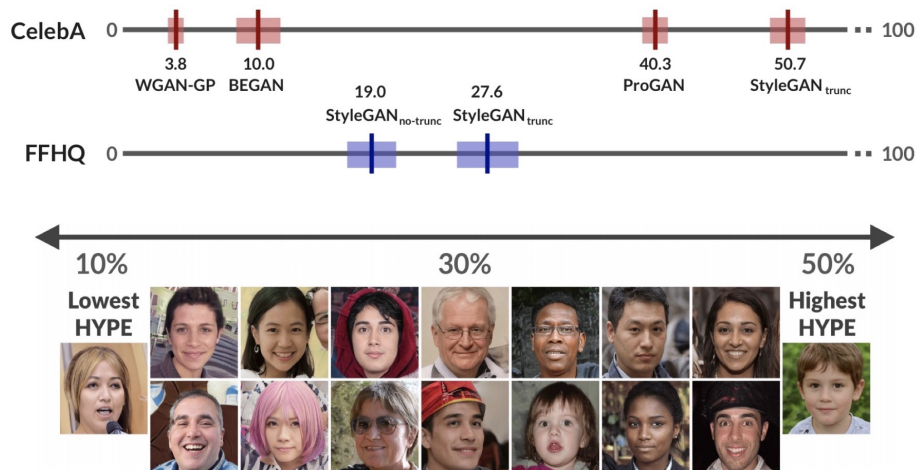


Figure 14: HYPE tests generative models for how realistic their images look to the human eye (Zhou et al., 2019). Top) HYPE scores of different models over CelebA and FFHQ datasets. A score of 50% represents indistinguishable results from real, while a score above 50% represents hyper-realism. Bottom) Example images sampled with the truncation trick from StyleGAN trained on FFHQ dataset. Images on the right have the highest HYPE scores (*i.e.* exhibit the highest perceptual fidelity). Figure compiled from Zhou et al. (2019).

of a rare visual target. Fig. 15 shows a depiction of this approach. Wang *et al.* show that Neuroscore is more consistent with human judgments compared to the conventional metrics. They also trained a convolutional neural network to predict Neuroscore from GAN-generated images directly without the need for neural responses.

3.3. Seeing What a GAN Can Not Generate

Bau et al. (2019) visualize mode collapse at both distribution level and instance level. They employ a semantic segmentation network to compare the distribution of segmented objects in the generated images versus the real images. Differences in statistics reveal object classes that are omitted by a GAN. Their approach allows to visualize the GAN’s omissions for an omitted class and to compare differences between individual photos and their approximate inversions by a GAN. Fig. 16 illustrates this approach.

3.4. Measuring GAN Steerability

Jahanian et al. (2019) propose a method to quantify the degree to which basic visual transformations are achieved by navigating the latent space of a GAN (See Fig. 17). They first learn a N -dimensional vector representing the optimal path in the latent space for a given transformation. Formally, the task is to learn the walk w by minimizing the following objective function:

$$w^* = \operatorname{argmin}_w \mathbb{E}_{z, \alpha} [\mathcal{L}(G(z + \alpha w), \operatorname{edit}(G(z), \alpha))],$$

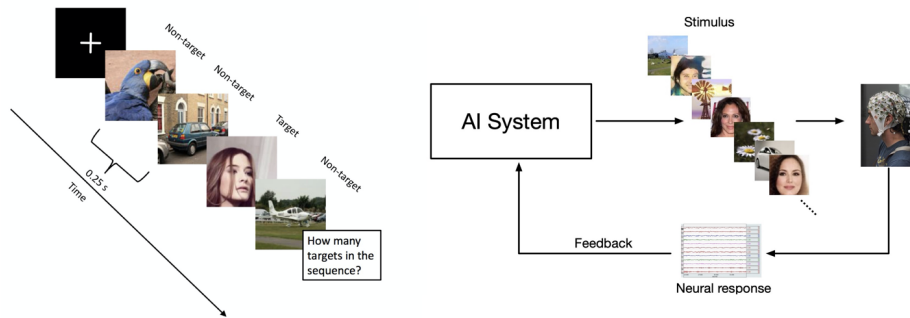


Figure 15: Left) An example of RSVP experimental protocol in which a rapid image stream (4 images per second) containing target and non-target images is presented to participants. Participants are instructed to search for real face images. The idea is that GAN generated faces will elicit a different response than real faces. Right) Schematic diagram of neuro-AI interface for computing the Neuroscore. Generated images by a GAN are shown to the participants and the corresponding recorded neural responses are used to evaluate the performance. Figure compiled from Wang et al. (2020b).

where α is the step size, and \mathcal{L} is the distance between the generated image $G(z + \alpha w)$ after taking an α -step in the latent direction and the target image $\text{edit}(G(z), \alpha)$. The latter is the new image derived from the source image $G(z)$. To quantify how well a desired image manipulation under each transformation is achieved, the distributions of a given attribute (*e.g.* “luminance”) in the real data and generated images (after walking in latent space) are compared.

3.5. GAN Dissection

Bau et al. (2018) propose a technique to dissect and visualize the inner workings of an image generator¹⁹. The main idea is to identify GAN units (*i.e.* generator neurons) that are responsible for semantic concepts and objects (such as tree, sky, and clouds) in the generated images. Having this level of granularity into the neurons allows editing existing images (*e.g.* to add or remove trees as shown in Fig. 18) by forcefully activating and deactivating (ablating) the corresponding units for the desired objects. Their technique also allows finding artifacts in the generated images and hence can be used to evaluate and improve GANs. A similar approach has been proposed in Park et al. (2019) for semantic manipulation and editing of GAN generated images²⁰.

3.6. A Universal Fake vs. Real Detector

Wang et al. (2020a) ask whether it is possible to create a “universal” detector to distinguish real images from synthetic ones using a dataset of synthetic images generated by 11 CNN-based generative models. See Fig. 19. With careful pre-

¹⁹A demo of this work is available at [here](#).

²⁰See [here](#) for an illustration.

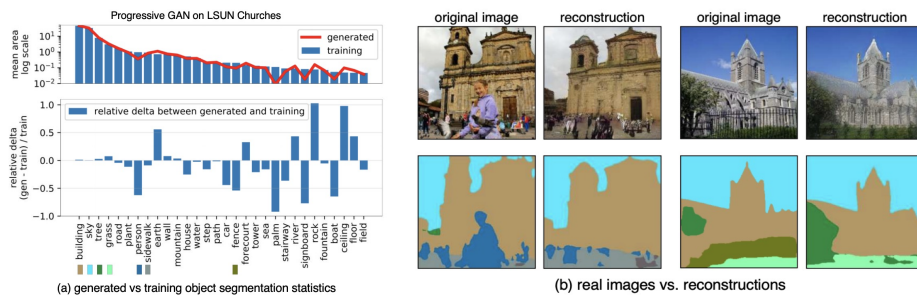


Figure 16: Seeing what a GAN cannot generate (Bau et al., 2019). (a) Distribution of object segmentations in the training set of LSUN churches vs. the corresponding distribution over the generated images. Objects such as people, cars, and fences are dropped by the generator. (b) Pairs of real images and their reconstructions in which individual instances of a person and a fence cannot be generated. Figure compiled from Bau et al. (2019).

and post-processing and data augmentation, they show that an image classifier trained on only one specific CNN generator is able to generalize well to unseen architectures, datasets, and training methods. They highlight that today’s CNN-generated images share common systematic flaws, preventing them from achieving realistic image generation. Similar works have been reported in Chai et al. (2020); Yu et al. (2019). In alignment with these results, Gragnaniello et al. (2021) also conclude that we are still far from having reliable tools for GAN image detection.

4. Discussion

4.1. GAN Benchmarks and Analysis Studies

Following previous works (*e.g.* Lucic et al. (2017); Shmelkov et al. (2018)), new studies have investigated formulating good criteria for GAN evaluation, or have conducted systematic GAN benchmarks. Gulrajani et al. (2020) argue that a good evaluation measure should not have a trivial solution (*e.g.* memorizing the dataset) and show that many scores such as IS and FID can be won by simply memorizing the training data (Fig. 20). They suggest that a necessary condition for a metric not to behave this way is to have a large number of samples. They also propose a measure based on neural network divergences (NND). NND works by training a discriminative model to discriminate between samples of the generative model and samples from a held-out test set. The poorer the discriminative model performs, the better the generative model is. Through experimental validation, they show that NND can effectively measure diversity, sample quality, and generalization. Lee and Town (2020) introduce Mimicry²¹, a lightweight PyTorch library that provides implementations of popular GANs and

²¹<https://github.com/kwotsin/mimicry>

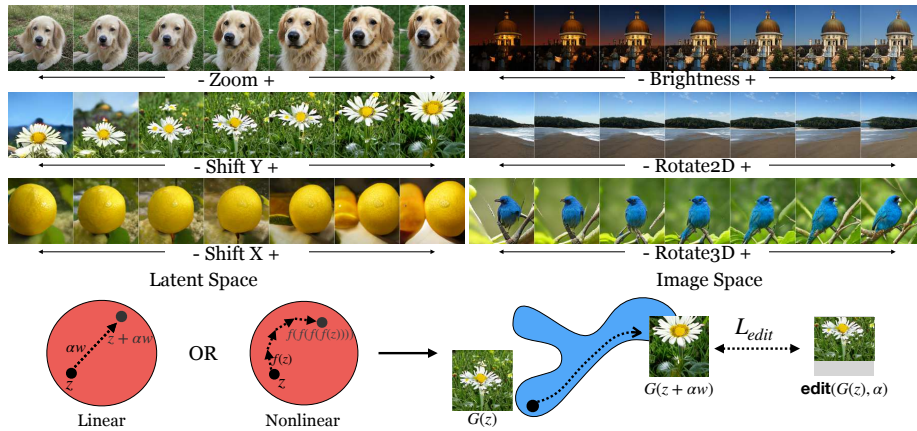


Figure 17: Top) A walk in the latent space of a GAN corresponds to visual transformations such as zoom and camera shift, Bottom) The goal is to find a path in z space (linear or non-linear) to transform the generated image $G(z)$ to its edited version $\text{edit}(G(z), \alpha)$, e.g. an $\alpha \times$ zoom. To measure steerability, the distributions of a given attribute in real images and generated images (after walking in the latent space) are compared. Figure from [Jahanian et al. \(2019\)](#).

evaluation metrics to closely reproduce reported scores in the literature. They also compare several GANs on seven widely-used datasets by training them under the same conditions, and evaluating them using three popular GAN metrics. Some tutorials and interaction tools have also been developed to understand, visualize, and evaluate GANs²².

Also mention that deepfakes might be detected because they have different spectral properties than real images. See for example [Durall et al. \(2020\)](#)

4.2. Assessing Fairness and Bias of Generative Models

Fairness and bias in ML algorithms, datasets, and commercial products have become growing concerns recently (e.g. [Buolamwini and Gebu \(2018\)](#)) and have attracted widespread attention from public and media²³. Even without a single precise definition for fairness ([Verma and Rubin, 2018](#)), it is still possible to observe the lack of fairness across many domains and models, with GANs being no exception. Some recent works (e.g. [Xu et al. \(2018\)](#); [Yu et al. \(2020\)](#)) have tried to mitigate the bias in GANs or use GANs to reduce bias in other ML algorithms (e.g. [Sattigeri et al. \(2019\)](#); [McDuff et al. \(2019\)](#)). Bias can enter a model during training (through data, labeling, architecture), evaluation (such as who created the evaluation measure), as well as deployment (how is the model being distributed and for what purposes). Therefore it is critical to address these aspects when evaluating and comparing generative models.

²²See [this link](#) and <https://poloclub.github.io/ganlab/>.

²³<https://thegradiant.pub/pulse-lessons/>

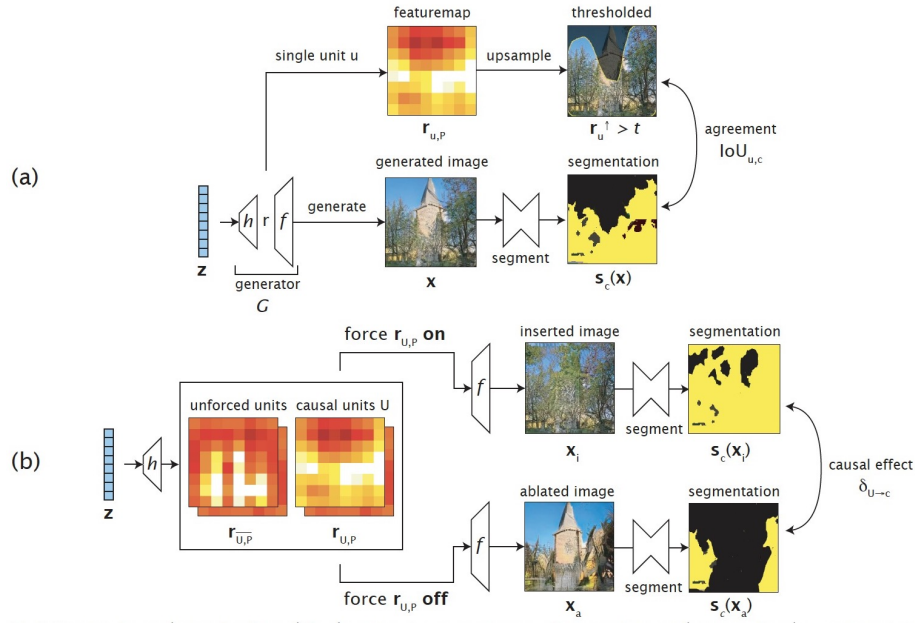


Figure 18: An overview of the GAN dissection approach (Bau et al., 2018). Top) Measuring the relationship between representation units and trees in the output using (a) dissection and (b) intervention. Dissection measures agreement between a unit u and a concept c by comparing its thresholded upsampled heatmap with a semantic segmentation of the generated image $s_c(x)$. Intervention measures the causal effect of a set of units U on a concept c by comparing the effect of forcing these units on (unit insertion) or off (unit ablation). The segmentation s_c reveals that trees increase after insertion and decrease after ablation. The average difference in the tree pixels measures the average causal effect. Bottom) Applying the dissection method to a generated outdoor church image. Dissection method can also be used to diagnose and improve GANs by identifying and ablating the artifact-causing units (panels e to g). Figure compiled from Bau et al. (2018).



Figure 19: Wang et al. (2020a) show that a classifier trained to distinguish images generated by only one GAN (ProGAN, the left-most column) from real ones can detect the images generated by other generative models (remaining columns) as well. Please see <https://peterwang512.github.io/CNNDetection/>.

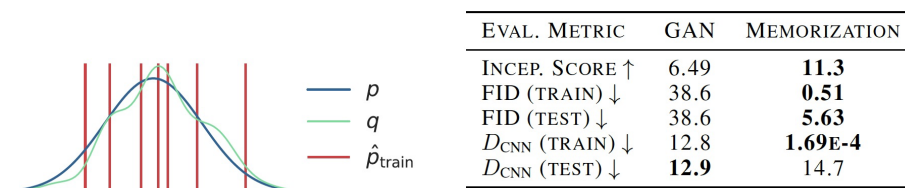


Figure 20: Left) Gulrajani et al. (2020) show that common evaluation measures such as IS and FID prefer a model that memorizes the dataset (p_{train} , red) to a model (q , green) which imperfectly fits the true distribution (p , blue) but covers more of p 's support. Right) Neural network divergence, DCNN, discounts memorization and prefers GANs that generalize beyond the training set. Figure compiled from Gulrajani et al. (2020)

4.3. Connection to Deepfakes

A alarming application of generative models is fabricating fake content²⁴. It is thus crucial to develop tools and techniques to detect and limit their use (See Tolosana et al. (2020) for a survey on deepfakes). There is a natural connection between deepfake detection and GAN evaluation. Obviously, as generative models improve, it becomes increasingly harder for a human to distinguish between what is real and what is fake (manipulated images, videos, text, and audio). Telling the degree of fakeness of an image or video directly tells us about the performance of a generator. Even over faces, where generative models excel, it is still possible to detect fake images (Fig. 21), although in some cases it can be very daunting (Fig. 22). Difficulty of deepfake detection for humans is also category dependent, with some categories such as faces, cats, and dogs being easier than bedrooms or cluttered scenes (Fig. 23). Fortunately, or unfortunately, it is still possible to build deep networks that can detect the subtle artifacts in the doctored images (*e.g.* using the universal detectors mentioned above (Wang et al., 2020a; Chai et al., 2020)). Moving forward, it is important to study whether and how GAN evaluation measures can help us mitigate the threat from deepfakes.

²⁴Some concerns can be found at these links 1, 2, 3, and 4



Figure 21: Sample generated/fake faces and cues to tell them apart from real ones. See also <https://chail.github.io/patch-forensics/>. Image courtesy of Twitter.

5. Summary and Conclusion

Here, I reviewed a number of recently proposed GAN evaluation measures. A summary of the discussed measures is shown in Figure 24. Similar to generative models, evaluation measures are also evolving and improving over time. Although some measures such as IS, FID, P&R, and PPL are relatively more popular than others, objective and comprehensive evaluation of generative models is still an open problem. Some directions for future research in GAN evaluation are as follows.

1. Prior research has been focused on examining generative models of faces or scenes containing one or few objects. Relatively, less effort has been devoted to assess how good GANs are in generating more complex scenes such as bedrooms, street scenes, and nature scenes. As an example, [Casanova et al. \(2020\)](#) is an early effort in this direction.
2. Generative models have been primarily evaluated in terms of the quality and diversity of their generated images. Other dimensions such as generalization and fairness have been less explored. Generalization assessment can provide a deeper look into what generative models learn. For example, it can tell how and to what degree models capture compositionality (*e.g.* does a model generate the right number of paws, nose, eyes, etc for dogs?) and logic (does a model properly capture physical properties such as gravity, light direction, reflection, and shadow?). Evaluating models in terms of fairness is critical for mitigating the potential risks that may arise at the deployment time and to ensure that a model has the right societal impact.
3. An important matter in GAN evaluation is task dependency. In other words, how well a generative model works depends on its intended use. Sometimes,

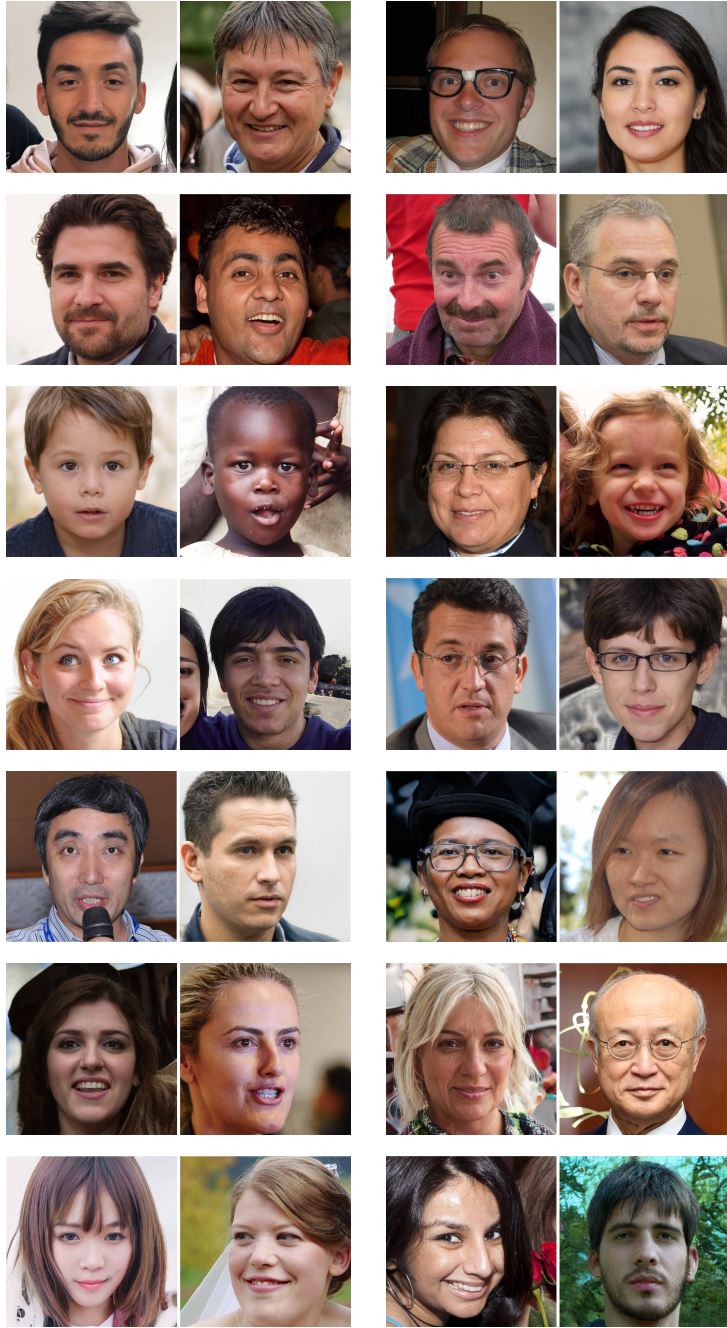


Figure 22: Can you determine which face in each pair is real? Key (L for left and R for right) row wise: L, L, R, L, R, R, L, L, L, R, R, R, R, R, L. Images taken from <https://www.whichfaceisreal.com/>.

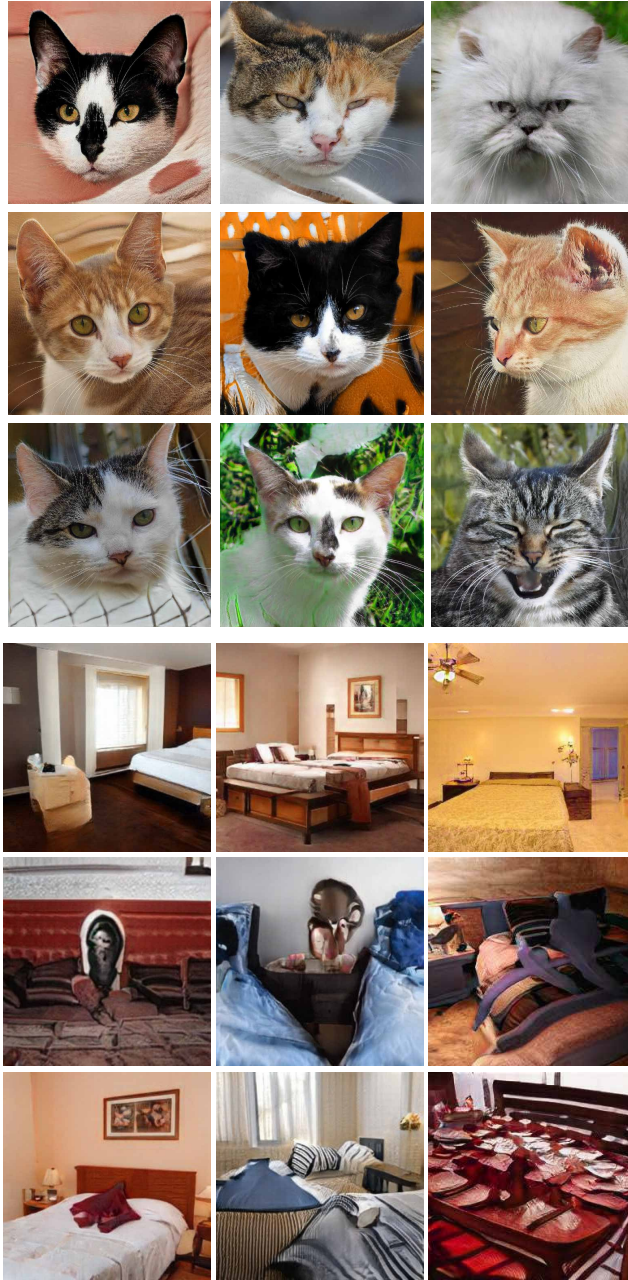


Figure 23: Some sample generated cats (top) and beds (bottom), generated from <http://thiscatdoesnotexist.com/> and <https://thisrentaldoesnotexist.com/>, respectively. See also <https://thisxdoesnotexist.com/>. Although images look realistic in the first glance, a closer examination reveals the artifacts (See it for yourself!).

	Quantitative /Analysis /Optimization	Qualitative	Overfitting /Memorization	Latent Space Disentanglement	Deepfake Detection
FID & IS Variants					
Spatial FID (sFID)	Analysis				
Class-aware FID (CAFD)	✓		✓		
Conditional FID	✓		✓		
Fast FID	Optimization				
Memorization-informed FID (MIFID)	✓		✓		
Unbiased FID and IS	✓				
Clean FID	Analysis				
Fre'chet Video Distance (FVD)	✓				
Methods based on Self-supervised Learned Representations	Analysis				
Methods based on Analysing Data Manifold					
Local Intrinsic Dimensionality (LID)	✓				
Intrinsic Multi-scale Distance (IMD)	✓				
Perceptual Path Length (PPL)	✓	✓		✓	
Linear Separability in Latent Space	✓			✓	
Classification Accuracy Score (CAS)	✓				
Non-Parametric Tests to Detect Data-Copying	✓		✓		
Measures that Probe Generalization	Analysis		✓		
New Ideas based on Precision and Recall (P&R)					
Density and Coverage	✓		✓		
Alpha Precision and Recall	✓		✓		
Duality GAP Metric	✓				
Spectral Methods	✓	✓			✓
Caption Score (CapS)	✓				
Human Eye Perceptual Evaluation (HYPE)		✓			
Neuroscore		✓			
GAN Steerability & Dissection		✓		✓	
A Universal Fake vs. Real Detector		✓			✓

Figure 24: A summary of evaluation measures covered in this work.

it might be easier to evaluate a model on downstream tasks when those tasks have a clear target for a given input. In some tasks (*e.g.* graphics applications such as image synthesis, image translation, image inpainting, and attribute manipulation) image quality is more important, whereas in some other tasks (*e.g.* generating synthetic data for data augmentation) diversity may weigh more. Thus, evaluation metrics should be tailored to the target task. It should be noted that the representation power of the discriminator or encoder of a GAN does not necessarily reflect its sample quality and diversity.

4. Having good evaluation measures is important not only for ranking generative models but also for diagnosing their errors. Reliable GAN evaluation is specially important in domains where humans are less attuned to discern the quality of samples (*e.g.* medical images). An important question thus is whether current evaluation measures generalize across different domains (*i.e.* are domain-agnostic).
5. The degree to which generative models memorize the training data is still unclear. A common technique to assess memorization is to look for nearest neighbors. It is known that this approach has several shortcomings (Theis et al., 2015; Borji, 2019). Motivated by works that study memorization in supervised learning works, van den Burg and Williams (2021) proposed new methods to understand and quantify memorization in generative models.
6. Ultimately, future works should also study how research in evaluating performance of generative models can help mitigate the threat from fabricated content. In this regard, extending works such as Durall et al. (2020); Wang et al. (2020a), conducting benchmarks on standard datasets containing fabricated content, as well as cross-talks between research in GAN evaluation and deep fake detection will lead to advances in these areas.

References

- Alaa, A.M., van Breugel, B., Saveliev, E., van der Schaar, M., 2021. How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models. arXiv preprint arXiv:2102.08921 .
- Bai, C.Y., Lin, H.T., Raffel, C., Kan, W.C.w., 2021. On training sample memorization: Lessons from benchmarking generative modeling with a large-scale competition. arXiv preprint arXiv:2106.03062 .
- Barannikov, S., Trofimov, I., Sotnikov, G., Trimbach, E., Korotin, A., Filipov, A., Burnaev, E., 2021. Manifold topology divergence: a framework for comparing data manifolds. arXiv preprint arXiv:2106.04024 .
- Barratt, S., Sharma, R., 2018. A note on the inception score. arXiv preprint arXiv:1801.01973 .
- Barua, S., Ma, X., Erfani, S.M., Houle, M.E., Bailey, J., 2019. Quality evaluation of gans using cross local intrinsic dimensionality. arXiv preprint arXiv:1905.00643 .

- Bau, D., Zhu, J.Y., Strobel, H., Zhou, B., Tenenbaum, J.B., Freeman, W.T., Torralba, A., 2018. Gan dissection: Visualizing and understanding generative adversarial networks. arXiv preprint arXiv:1811.10597 .
- Bau, D., Zhu, J.Y., Wulff, J., Peebles, W., Strobel, H., Zhou, B., Torralba, A., 2019. Seeing what a gan cannot generate, in: Proceedings of the IEEE International Conference on Computer Vision, pp. 4502–4511.
- Bińkowski, M., Sutherland, D.J., Arbel, M., Gretton, A., 2018. Demystifying mmd gans. arXiv preprint arXiv:1801.01401 .
- Bond-Taylor, S., Leach, A., Long, Y., Willcocks, C.G., 2021. Deep generative modelling: A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models. arXiv preprint arXiv:2103.04922 .
- Borji, A., 2019. Pros and cons of gan evaluation measures. *Computer Vision and Image Understanding* 179, 41–65.
- Brock, A., Donahue, J., Simonyan, K., 2018. Large scale gan training for high fidelity natural image synthesis. arXiv preprint arXiv:1809.11096 .
- Buolamwini, J., Gebru, T., 2018. Gender shades: Intersectional accuracy disparities in commercial gender classification, in: Conference on fairness, accountability and transparency, PMLR. pp. 77–91.
- van den Burg, G.J., Williams, C.K., 2021. On memorization in probabilistic deep generative models. arXiv preprint arXiv:2106.03216 .
- Carreira, J., Zisserman, A., 2017. Quo vadis, action recognition? a new model and the kinetics dataset, in: proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 6299–6308.
- Casanova, A., Drozdal, M., Romero-Soriano, A., 2020. Generating unseen complex scenes: are we there yet? arXiv preprint arXiv:2012.04027 .
- Chai, L., Bau, D., Lim, S.N., Isola, P., 2020. What makes fake images detectable? understanding properties that generalize, in: European Conference on Computer Vision, Springer. pp. 103–120.
- Chen, L., Li, Z., Maddox, R.K., Duan, Z., Xu, C., 2018. Lip movements generation at a glance, in: Proceedings of the European Conference on Computer Vision (ECCV), pp. 520–535.
- Chong, M.J., Forsyth, D., 2020. Effectively unbiased fid and inception score and where to find them, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 6070–6079.
- De, K., Masilamani, V., 2013. Image sharpness measure for blurred images in frequency domain. *Procedia Engineering* 64, 149–158.

- Denton, E., Chintala, S., Szlam, A., Fergus, R., 2015. Deep generative image models using a laplacian pyramid of adversarial networks. arXiv preprint arXiv:1506.05751 .
- Ding, M., Yang, Z., Hong, W., Zheng, W., Zhou, C., Yin, D., Lin, J., Zou, X., Shao, Z., Yang, H., et al., 2021. Cogview: Mastering text-to-image generation via transformers. arXiv preprint arXiv:2105.13290 .
- Djolonga, J., Lucic, M., Cuturi, M., Bachem, O., Bousquet, O., Gelly, S., 2020. Precision-recall curves using information divergence frontiers, in: International Conference on Artificial Intelligence and Statistics, pp. 2550–2559.
- Durall, R., Keuper, M., Keuper, J., 2020. Watch your up-convolution: Cnn based generative deep neural networks are failing to reproduce spectral distributions, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 7890–7899.
- Dzanic, T., Shah, K., Witherden, F., 2019. Fourier spectrum discrepancies in deep network generated images. arXiv preprint arXiv:1911.06465 .
- Frank, J., Eisenhofer, T., Schönherr, L., Fischer, A., Kolossa, D., Holz, T., 2020. Leveraging frequency analysis for deep fake image recognition, in: International Conference on Machine Learning, PMLR. pp. 3247–3258.
- Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial networks. arXiv preprint arXiv:1406.2661 .
- Gragnaniello, D., Cozzolino, D., Marra, F., Poggi, G., Verdoliva, L., 2021. Are gan generated images easy to detect? a critical analysis of the state-of-the-art. arXiv preprint arXiv:2104.02617 .
- Grnarova, P., Levy, K.Y., Lucchi, A., Perraudin, N., Goodfellow, I., Hofmann, T., Krause, A., 2019. A domain agnostic measure for monitoring and evaluating gans, in: Advances in Neural Information Processing Systems, pp. 12092–12102.
- Gulrajani, I., Raffel, C., Metz, L., 2020. Towards gan benchmarks which require generalization. arXiv preprint arXiv:2001.03653 .
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S., 2017. Gans trained by a two time-scale update rule converge to a local nash equilibrium. arXiv preprint arXiv:1706.08500 .
- Hudson, D.A., Zitnick, C.L., 2021. Generative adversarial transformers. arXiv preprint arXiv:2103.01209 .
- Iqbal, T., Qureshi, S., 2020. The survey: Text generation models in deep learning. Journal of King Saud University-Computer and Information Sciences .

- Jahanian, A., Chai, L., Isola, P., 2019. On the "steerability" of generative adversarial networks. arXiv preprint arXiv:1907.07171 .
- Jiang, Y., Chang, S., Wang, Z., 2021. Transgan: Two transformers can make one strong gan. arXiv preprint arXiv:2102.07074 .
- Karras, T., Laine, S., Aila, T., 2019. A style-based generator architecture for generative adversarial networks, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 4401–4410.
- Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T., 2020. Analyzing and improving the image quality of stylegan, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 8110–8119.
- Khrulkov, V., Oseledets, I., 2018. Geometry score: A method for comparing generative adversarial networks, in: International Conference on Machine Learning, PMLR. pp. 2621–2629.
- Kingma, D.P., Welling, M., 2013. Auto-encoding variational bayes. arXiv preprint arXiv:1312.6114 .
- Kolchinski, Y.A., Zhou, S., Zhao, S., Gordon, M., Ermon, S., 2019. Approximating human judgment of generated image quality. arXiv preprint arXiv:1912.12121 .
- Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., Aila, T., 2019. Improved precision and recall metric for assessing generative models. arXiv preprint arXiv:1904.06991 .
- Lee, K.S., Town, C., 2020. Mimicry: Towards the reproducibility of gan research. arXiv preprint arXiv:2005.02494 .
- Liu, M.Y., Huang, X., Yu, J., Wang, T.C., Mallya, A., 2021. Generative adversarial networks for image and video synthesis: Algorithms and applications. Proceedings of the IEEE .
- Liu, S., Wei, Y., Lu, J., Zhou, J., 2018. An improved evaluation framework for generative adversarial networks. arXiv preprint arXiv:1803.07474 .
- Liu, Z., Luo, P., Wang, X., Tang, X., 2015. Deep learning face attributes in the wild, in: Proceedings of the IEEE international conference on computer vision, pp. 3730–3738.
- Lucic, M., Kurach, K., Michalski, M., Gelly, S., Bousquet, O., 2017. Are gans created equal? a large-scale study. arXiv preprint arXiv:1711.10337 .
- Mathiasen, A., Hvilshøj, F., 2020. Fast fr\'echet inception distance. arXiv preprint arXiv:2009.14075 .

- McDuff, D., Ma, S., Song, Y., Kapoor, A., 2019. Characterizing bias in classifiers using generative models. arXiv preprint arXiv:1906.11891 .
- Meehan, C., Chaudhuri, K., Dasgupta, S., 2020. A non-parametric test to detect data-copying in generative models. arXiv preprint arXiv:2004.05675 .
- Morozov, S., Voynov, A., Babenko, A., 2020. On self-supervised image representations for gan evaluation, in: International Conference on Learning Representations.
- Naeem, M.F., Oh, S.J., Uh, Y., Choi, Y., Yoo, J., 2020. Reliable fidelity and diversity metrics for generative models. arXiv preprint arXiv:2002.09797 .
- Narvekar, N.D., Karam, L.J., 2009. A no-reference perceptual image sharpness metric based on a cumulative probability of blur detection, in: 2009 International Workshop on Quality of Multimedia Experience, IEEE. pp. 87–91.
- Nash, C., Menick, J., Dieleman, S., Battaglia, P.W., 2021. Generating images with sparse representations. arXiv preprint arXiv:2103.03841 .
- O’Brien, S., Groh, M., Dubey, A., 2018. Evaluating generative adversarial networks on explicitly parameterized distributions. arXiv preprint arXiv:1812.10782 .
- Odena, A., 2019. Open questions about generative adversarial networks. Distill 4, e18.
- Oord, A.v.d., Kalchbrenner, N., Vinyals, O., Espeholt, L., Graves, A., Kavukcuoglu, K., 2016. Conditional image generation with pixelcnn decoders. arXiv preprint arXiv:1606.05328 .
- Oprea, S., Martinez-Gonzalez, P., Garcia-Garcia, A., Castro-Vargas, J.A., Orts-Escolano, S., Garcia-Rodriguez, J., Argyros, A., 2020. A review on deep learning techniques for video prediction. IEEE Transactions on Pattern Analysis and Machine Intelligence .
- Park, T., Liu, M.Y., Wang, T.C., Zhu, J.Y., 2019. Semantic image synthesis with spatially-adaptive normalization, in: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2337–2346.
- Parmar, G., Zhang, R., Zhu, J.Y., 2021. On buggy resizing libraries and surprising subtleties in fid calculation. arXiv preprint arXiv:2104.11222 .
- Preuer, K., Renz, P., Unterthiner, T., Hochreiter, S., Klambauer, G., 2018. Fréchet chemnet distance: a metric for generative models for molecules in drug discovery. Journal of chemical information and modeling 58, 1736–1741.
- Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., Chen, M., Sutskever, I., 2021. Zero-shot text-to-image generation. arXiv preprint arXiv:2102.12092 .

- Ravuri, S., Vinyals, O., 2019. Classification accuracy score for conditional generative models, in: *Advances in Neural Information Processing Systems*, pp. 12268–12279.
- Razavi, A., Oord, A.v.d., Vinyals, O., 2019. Generating diverse high-fidelity images with vq-vae-2. *arXiv preprint arXiv:1906.00446* .
- Roblek, D., Kilgour, K., Sharifi, M., Zuluaga, M., 2019. Fréchet audio distance: A reference-free metric for evaluating music enhancement algorithms .
- Sajjadi, M.S., Bachem, O., Lucic, M., Bousquet, O., Gelly, S., 2018. Assessing generative models via precision and recall, in: *Advances in Neural Information Processing Systems*, pp. 5228–5237.
- Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Radford, A., Chen, X., 2016. Improved techniques for training gans. *arXiv preprint arXiv:1606.03498* .
- Sattigeri, P., Hoffman, S.C., Chenthamarakshan, V., Varshney, K.R., 2019. Fairness gan: Generating datasets with fairness properties using a generative adversarial network. *IBM Journal of Research and Development* 63, 3–1.
- Shmelkov, K., Schmid, C., Alahari, K., 2018. How good is my gan?, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 213–229.
- Shoemake, K., 1985. Animating rotation with quaternion curves, in: *Proceedings of the 12th annual conference on Computer graphics and interactive techniques*, pp. 245–254.
- Sidheekh, S., Aimen, A., Krishnan, N.C., 2021. Characterizing gan convergence through proximal duality gap. *arXiv preprint arXiv:2105.04801* .
- Simon, L., Webster, R., Rabin, J., 2019. Revisiting precision and recall definition for generative model evaluation. *arXiv preprint arXiv:1905.05441* .
- Simonyan, K., Zisserman, A., 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* .
- Soloveitchik, M., Diskin, T., Morin, E., Wiesel, A., 2021. Conditional frechet inception distance. *arXiv preprint arXiv:2103.11521* .
- van Steenkiste, S., Kurach, K., Schmidhuber, J., Gelly, S., 2020. Investigating object compositionality in generative adversarial networks. *Neural Networks* 130, 309–325.
- Tevet, G., Habib, G., Shwartz, V., Berant, J., 2018. Evaluating text gans as language models. *arXiv preprint arXiv:1810.12686* .
- Theis, L., Oord, A.v.d., Bethge, M., 2015. A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844* .

- Tolosana, R., Vera-Rodriguez, R., Fierrez, J., Morales, A., Ortega-Garcia, J., 2020. Deepfakes and beyond: A survey of face manipulation and fake detection. *Information Fusion* 64, 131–148.
- Tsitsulin, A., Munkhoeva, M., Mottin, D., Karras, P., Bronstein, A., Oseledets, I., Müller, E., 2019. The shape of data: Intrinsic distance for data distributions. *arXiv preprint arXiv:1905.11141* .
- Tulyakov, S., Liu, M.Y., Yang, X., Kautz, J., 2018. Mocogan: Decomposing motion and content for video generation, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1526–1535.
- Unterthiner, T., van Steenkiste, S., Kurach, K., Marinier, R., Michalski, M., Gelly, S., 2019. Fvd: A new metric for video generation .
- Verma, S., Rubin, J., 2018. Fairness definitions explained, in: *2018 IEEE/ACM international workshop on software fairness (fairware)*, IEEE. pp. 1–7.
- Wang, S.Y., Wang, O., Zhang, R., Owens, A., Efros, A.A., 2020a. Cnn-generated images are surprisingly easy to spot... for now, in: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Wang, Z., Healy, G., Smeaton, A.F., Ward, T.E., 2020b. Use of neural signals to evaluate the quality of generative adversarial network performance in facial image generation. *Cognitive Computation* 12, 13–24.
- Wang, Z., She, Q., Ward, T.E., 2019. Generative adversarial networks in computer vision: A survey and taxonomy. *arXiv preprint arXiv:1906.01529* .
- Xu, D., Yuan, S., Zhang, L., Wu, X., 2018. Fairgan: Fairness-aware generative adversarial networks, in: *2018 IEEE International Conference on Big Data (Big Data)*, IEEE. pp. 570–575.
- Xuan, J., Yang, Y., Yang, Z., He, D., Wang, L., 2019. On the anomalous generalization of gans. *arXiv preprint arXiv:1909.12638* .
- Yang, C., Wang, Z., Zhu, X., Huang, C., Shi, J., Lin, D., 2018. Pose guided human video generation, in: *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 201–216.
- Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J., 2015. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365* .
- Yu, N., Davis, L.S., Fritz, M., 2019. Attributing fake images to gans: Learning and analyzing gan fingerprints, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 7556–7566.
- Yu, N., Li, K., Zhou, P., Malik, J., Davis, L., Fritz, M., 2020. Inclusive gan: Improving data and minority coverage in generative models, in: *European Conference on Computer Vision*, Springer. pp. 377–393.

- Zeng, Y., Lu, H., Borji, A., 2017. Statistics of deep generated images. arXiv preprint arXiv:1708.02688 .
- Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O., 2018. The unreasonable effectiveness of deep features as a perceptual metric, in: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 586–595.
- Zhao, S., Ren, H., Yuan, A., Song, J., Goodman, N., Ermon, S., 2018. Bias and generalization in deep generative models: An empirical study, in: Advances in Neural Information Processing Systems, pp. 10792–10801.
- Zhou, S., Gordon, M., Krishna, R., Narcomey, A., Fei-Fei, L.F., Bernstein, M., 2019. Hype: A benchmark for human eye perceptual evaluation of generative models, in: Advances in Neural Information Processing Systems, pp. 3449–3461.