# ADVERSARIAL JAMMING FOR AUTOENCODER DISTRIBUTION MATCHING

*Waleed El-Geresy and Deniz Gündüz*
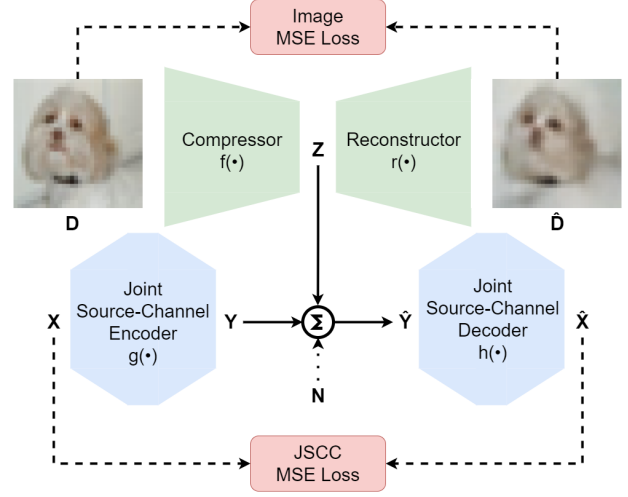
Imperial College London

## ABSTRACT

We propose the use of adversarial wireless jamming to regularise the latent space of an autoencoder to match a diagonal Gaussian distribution. We consider the minimisation of a mean squared error distortion, where a jammer attempts to disrupt the recovery of a Gaussian source encoded and transmitted over the adversarial channel. A straightforward consequence of existing theoretical results is the fact that the saddle point of a minimax game - involving such an encoder, its corresponding decoder, and an adversarial jammer - consists of diagonal Gaussian noise output by the jammer. We use this result as inspiration for a novel approach to distribution matching in the latent space, utilising jamming as an auxiliary objective to encourage the aggregated latent posterior to match a diagonal Gaussian distribution. Using this new technique, we achieve distribution matching comparable to standard variational autoencoders and to Wasserstein autoencoders. This approach can also be generalised to other latent distributions.

***Index Terms***— generative models, variational autoencoders, minimax game, game theory

## 1. INTRODUCTION

Generative modelling can be described as the task of drawing artificial samples from an underlying distribution based on the observation of a real set of samples from that distribution. The distribution can be simple and tractable, or it can be complex and intractable. A variety of methods for generative modelling in the field of machine learning and deep learning have been proposed. These include diffusion models [1], generative adversarial networks (GANs) [2], variational autoencoders (VAEs) [3], normalising flows [4], and energy-based models [5]. Many of these can generate samples from the modelled distribution given a representation (invert the mapping). However, the ability to map a sample (either real or generated) back to an appropriate representation is unique to symmetric techniques, namely VAEs and normalising flows.

Here, we focus on VAEs, [3] which use a variational approximation to model an intractable prior distribution for a particular data distribution. Commonly, a parameterised multivariate Gaussian distribution is used as the prior. This is

**Fig. 1**. Adversarial jamming for prior distribution matching. The game between a compressor/jammer $f$ (compressing the images) and a JSCC autoencoder $(g, h)$ (for source $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$) has a saddle point with $Z^* \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, implicitly imposing this prior on the latent space.

because the Kullback-Leibler (KL) divergence - in the case of a multivariate Gaussian - has an analytic differentiable form, which can be used as part of the loss function. Secondly, using a known distribution such as a multivariate Gaussian can allow us to sample from the data distribution by first sampling this prior and then generating samples from the resultant posterior using the autoencoder's decoder.

Existing work in the field of VAEs involving the matching of the aggregated posterior to a known prior, such as a diagonal Gaussian distribution, has consisted of a number of methods. The initial VAE relied on maximising the so-called evidence lower bound (ELBO) - a lower bound on the likelihood of generated data coming from the true distribution, used as a proxy for direct maximisation of the likelihood. The formulation can also be viewed as a sum of two terms: a reconstruction loss, and an aggregated posterior distribution matching loss term that is the KL divergence from the diagonal Gaussian prior. Wasserstein autoencoders (WAEs) proposed using divergences induced by the optimal transport problem, as an alternative to the KL divergence, for minimising the distance

between the aggregated posterior and an independent Gaussian distribution. [6]. They achieved this in two different ways: by using a GAN to enforce distribution matching in the latent space, or by using a maximum mean discrepancy (MMD) loss term (with a choice of kernel function e.g. the inverse multi-quadratics (IMQ) kernel).

We propose a new adversarial jamming (AJ) approach to distribution matching based on a minimax game between a jammer, $f(\cdot)$, and an auxiliary joint source channel coding (JSCC) encoder and decoder pair, $g(\cdot)$ and $h(\cdot)$, respectively, for a diagonal Gaussian source, $X$, with a mean squared error (MSE) distortion measure for reconstruction under constraints on the power of both the jammer and the JSCC (see Figure 1). Existing theoretical results concerning the optimal (worst-case) additive jamming noise for such a setting [7, 8, 9] imply that under these conditions, the optimal jammer output will be a diagonal Gaussian random vector to maximise the expected value of the reconstruction error. We exploit this fact in order to produce an adversarial jamming regularisation term for an autoencoder that tries to simultaneously minimise the distortion in compressing and reconstructing images from a dataset, while also matching a diagonal Gaussian prior. We show that this technique is effective for distribution matching in the latent space, comparable to a VAE and a WAE. To the best of our knowledge, this is a novel result that exploits a communication theoretic saddle point result to regularise the latent space of an autoencoder for distribution matching, which can lead to similar, potential future extensions.

**Notation.** We use $\mathbf{0}$ and $\mathbf{I}$ to represent the $k$-dimensional zero vector and $(k \times k)$-dimensional identity matrix, respectively. Also, $R_A = Q_A \Lambda_A Q_A^T$ denotes the eigenvalue decomposition of the covariance matrix $R_A$ of random vector $\mathbf{A}$ and $F_A(\omega)$ denotes its characteristic function.

## 2. PROBLEM DEFINITION

Let $\mathbf{X} \in \mathbb{R}^k$, $\mathbf{D} \in \mathbb{R}^n$, and $\mathbf{N} \in \mathbb{R}^k$ be random vectors representing the source distribution, jammer's source of randomness, and channel noise, respectively. Let $f : \mathbb{R}^n \to \mathbb{R}^k$, $g : \mathbb{R}^k \to \mathbb{R}^k$, and $h : \mathbb{R}^k \to \mathbb{R}^k$ be Borel-measurable functions. These respectively represent the jammer, with output $\mathbf{Z} := f(\mathbf{D})$; the encoder (transmitter), with output $\mathbf{Y} := g(\mathbf{X})$; and the decoder (receiver), with output $\hat{\mathbf{X}} := h(\hat{\mathbf{Y}})$, where $\hat{\mathbf{Y}}$ defines the additive noise channel mapping $\hat{\mathbf{Y}} = \mathbf{Y} + \mathbf{N} + \mathbf{Z}$. Then, the zero-sum, minimax game among the transmitter, receiver, and adversarial jammer is shown below with cost function $J(f, g, h)$, and power constraints $P_t$ and $P_a$ for the transmitter and jammer, respectively.

$$\max_f \min_{h,g} J(f, g, h) \qquad (1)$$

where $J(f, g, h) = \mathbb{E}[(\mathbf{X} - \hat{\mathbf{X}})^2]$,

$\mathbb{E}[||g(\mathbf{X})||^2] \le P_t$, and $\mathbb{E}[||f(\mathbf{D})||^2] \le P_a$

The optimal combined policy for the transmitter, receiver, and jammer is the saddle point solution to the zero-sum game, $(f^*, h^*, g^*)$, where the two inequalities in Equation 2 are simultaneously satisfied.

$$J(f, h^*, g^*) \le J(f^*, h^*, g^*) \le J(f^*, h, g) \qquad (2)$$

It was shown in [8] that the problem of optimal (zero-delay) jamming over a scalar additive noise channel (i.e. $k = 1$) is closely connected to the linearity of $g^*$ and $h^*$. Where possible, a jammer will generate additive noise that forces the optimal transmitter and receiver to be linear transformations (the "matching condition"). In the case of a scalar Gaussian source being sent over a Gaussian channel, the output of the jammer will also be Gaussian.

The scalar jamming result was extended and generalised in [7], where it was shown that a vector version of the matching condition can be used to find saddle points for the vector version of the minimax game. An optimal encoding function for the transmitter will be a linear transformation (followed by a random sign change through multiplication by a Bernoulli random variable with $p = 0.5$ over the alphabet $\{0, 1\}$). It was shown that in the vector setting, the power allocation of the transmitter will obey reverse water-filling.

Our focus in this work is on the nature of the optimal jamming function. The necessary and sufficient condition for the linearity of optimal estimation - and a requirement for the optimal jamming strategy - is given in Equation 3, where $\Sigma$ is a diagonal power allocation matrix. The optimal jammer also obeys water-filling, dependent on the eigenvalues of $\mathbf{R_N}$ [7].

$$\frac{\delta \log F_{\Sigma Q_X^T \mathbf{X}}(\omega)}{\delta \omega_i} = S_i \frac{\delta \log F_{Q_Z^T (\mathbf{N} + \mathbf{Z})}(\omega)}{\delta \omega_i}, 1 \le i \le m \quad (3)$$

$$\text{where } S = \Sigma \Lambda_X \Sigma \Lambda_Z^{-1}$$

Such a characterisation of the nature of the optimal jamming solution implies that if $\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\mathbf{N} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then the optimal jammer output will also be an isotropic diagonal Gaussian distribution $\mathbf{Z} \sim \mathcal{N}(\mathbf{0}, \kappa \mathbf{I})$, where $\kappa$ is a scalar constant, to match the source distribution and enforce the transmitter's linearity. This can be easily verified by noting that, in this special case, $I = a\Sigma = b\Lambda_X = c\Lambda_Z = dQ_Z = eQ_X$, where $a, b, c, d, e$ are scalar constants. Thus we have a solution $\mathbf{Z} = \beta(\mathbf{N} - \mathbf{X})$, where $\beta$ is scalar, which means that a saddle point solution for $\mathbf{Z}$ is also a diagonal Gaussian random vector. This saddle point solution to the minimax game is also (almost surely) unique [7].

### 2.1. Jamming as Regularisation for an Autoencoder

We propose the use of this communication theoretic game to regularise the latent space of an autoencoder. Consider a data distribution random vector, $\mathbf{D} \in \mathbb{R}^n$; for example, flattened

image data with $n = l \cdot w \cdot m$, where $l$, $w$, and $m$ are the length, width, and colour channels of the images, respectively. $\mathbf{D}$ shall be the source of randomness for our jammer. We introduce an auxiliary DeepJSCC autoencoder that attempts to transmit and recover a source $X \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ over the additive noise channel. DeepJSCC has been shown to be an effective technique for transmitting complex source distributions, such as natural images, over unknown communication channels [10, 11, 12] and continues to be a topic that is actively researched in the field of semantic communication [13].

Since we deal with two autoencoders, one operating on data $\mathbf{D}$ and one operating on Gaussian source $\mathbf{X}$, we use the terms "data autoencoder" and "DeepJSCC autoencoder" to distinguish between the two. We also refer to the encoder/decoder of each as the compressor/reconstructor and the transmitter/receiver, respectively. In the context of the vector jamming game as defined above, the reconstructor is $f$, the transmitter is $g$, and the receiver is $h$. The reconstructor will be a function $r : \mathbb{R}^k \to \mathbb{R}^n$ with output $\hat{\mathbf{D}} := r(\mathbf{Z})$.

The objective of the DeepJSCC autoencoder is to minimise the reconstruction error of the Gaussian source, $\mathcal{L}_{jscc}$, given in Equation 4. In turn, the objective of the data autoencoder is to minimise an objective function, $\mathcal{L}_{data}$, which is composed of two terms: its own reconstruction error for the data distribution, as well as the negative of $\mathcal{L}_{jscc}$, as shown in Equation 5 below. Here, $\mathcal{L}_{data}$ mirrors the objective of a VAE, with $-\mathcal{L}_{jscc}$ replacing the usual KL divergence term.

$$\mathcal{L}_{jscc}(g, h) = \mathbb{E}[(\mathbf{X} - \hat{\mathbf{X}})^2] \qquad (4)$$
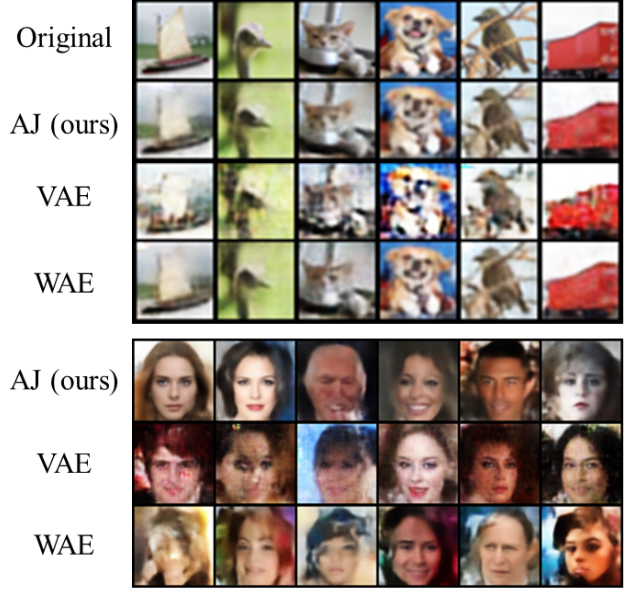
$$\mathcal{L}_{data}(f, r) = \mathbb{E}[(\mathbf{D} - \hat{\mathbf{D}})^2] - \eta \mathbb{E}[(\mathbf{X} - \hat{\mathbf{X}})^2] \qquad (5)$$

where $\eta$ is a positive regularisation constant.

We enforce power constraints $P_t = P_a = 1$ on the outputs $\mathbf{Y}$ and $\mathbf{Z}$ by normalising them (in practice, in mini-batch training with stochastic gradient descent, by using the empirical batch statistics). The competition between the compressor, acting as a jammer, and the DeepJSCC autoencoder, sets up a communication game in the form seen in Equation 1. The saddle point of this communication game is for the output of the jammer $\mathbf{Z}$ to be an isotropic diagonal Gaussian. Our proposed approach makes use of the maximisation term from the game in Equation 1 as a regularisation term that encourages the latent distribution to match this distribution. The proposed scheme is illustrated in Figure 1. For our experiments, we set the power of the fixed noise, $\mathbf{N}$, to 0, using only the adversarial jammer output as noise.

## 3. EXPERIMENTAL RESULTS

We use a DCGAN-like [14] convolutional autoencoder architecture for our data autoencoder, with output normalisation to ensure adherence to constraints on the variance and mean. We use a fully connected architecture that includes a linear
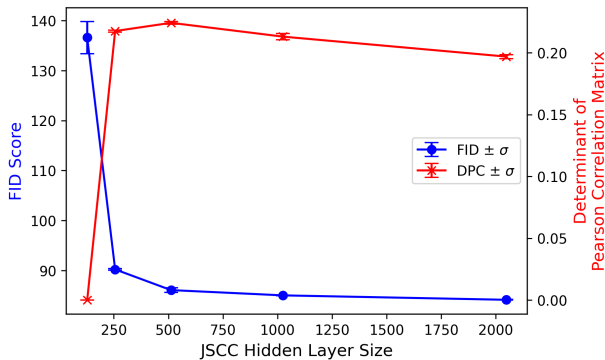


**Fig. 2**. Reconstruction results for the CIFAR-10 dataset (top) and generation results for the CelebA dataset (bottom), for experiments using adversarial jamming, KL divergence (VAE), and maximum mean discrepancy (WAE) as the prior distribution latent regularisation term.

mapping for our transmitter and receiver. The transmitter and receiver respectively have output and input shapes $k$ (which is the transmitter output $\mathbf{Y}$), and $k$ is also the length of the representation, which is the adversarial noise $\mathbf{Z}$, output by the compressor. $k$ is varied for each experiment, depending on the dataset. We test our method on three different datasets: CelebA [15], CIFAR-10 [16], and MNIST [17]. For CelebA we train for 75 epochs with $k = 64$, for CIFAR-10 we train for 250 epochs with $k = 64$, and for MNIST we run two experiments at $k = 2$ and $k = 8$, both for 250 epochs. We compare our method to a VAE trained using an ELBO loss term involving the KL divergence distribution matching, as well as a WAE with an MMD regularisation term using an IMQ Kernel. We use the same objective, parameters and architectures for these ablation experiments, changing only the distribution matching term.

Figure 2 shows the reconstruction results for CIFAR-10 and generation results for CelebA. Table 1 shows a comparison of three metrics for all the experiments. We measure the closeness of the distribution to an isotropic Gaussian indirectly, through two metrics: the Fréchet Inception Distance (FID) score [18], measuring the realism of the generated images; and the determinant of the Pearson correlation matrix (DPC) for the learned features, which is a measure of the degree of statistical independence between the learned features and thus the level of disentanglement of the features. Finally, the MSE is used to measure the fidelity of the image recon-

**Table 1**. FID scores, MSE losses for image reconstruction and the determinants of the Pearson correlation matrices (DPC). Arrows indicate whether higher (↑) or lower (↓) scores are better. These are shown for the three datasets. In each case, we show results for adversarial jamming, a standard VAE, and a WAE, trained using the same parameters and architecture. We see that the performance of adversarial jamming is comparable to WAEs and VAEs. **Boldface** entries denote the best performance.

| Dataset | CelebA | | | CIFAR-10 | | | MNIST | | | MNIST | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Latent Dim | 64 | | | 64 | | | 2 | | | 8 | | |
| Score | FID ↓ | MSE ↓ | DPC ↑ | FID ↓ | MSE ↓ | DPC ↑ | FID ↓ | MSE ↓ | DPC ↑ | FID ↓ | MSE ↓ | DPC ↑ |
| AJ (ours) | 46.9 | 0.0065 | 0.50 | 91.1 | **0.0104** | 0.55 | 23.6 | **0.0131** | 0.96 | **85.8** | **0.0362** | **1.00** |
| VAE | 50.1 | 0.0065 | 0.02 | 105.3 | 0.0116 | 0.17 | 24.1 | 0.0133 | 0.78 | 94.7 | 0.0369 | 0.99 |
| WAE | **44.2** | **0.0061** | **0.74** | **87.1** | 0.0105 | **0.75** | **22.0** | 0.0133 | **0.98** | 86.1 | 0.0370 | 1.00 |



**Fig. 3**. The FID score of a generative model trained using adversarial jamming versus the hidden layer size of the Deep-JSSC autoencoder. Increasing the capacity of the DeepJSCC autoencoder for learning improves the FID score of the generative model.

struction. In Figure 3, we see a plot of the FID score and the determinant of the Pearson correlation matrix against the DeepJSCC autoencoder's hidden layer size for experiments conducted on the CIFAR-10 dataset for $k = 128$.

## 4. DISCUSSION

Visually, the results in Figure 2 demonstrate performance on par with a WAE in terms of reconstruction and generation. The high determinants of the Pearson correlation matrices in Table 1 suggest that our method may achieve a good level of feature independence, comparable to a VAE or WAE. This is reflected in the FID scores of the models, showing that when a diagonal Gaussian prior is used to sample new data, the output samples appear comparatively realistic. However, it should be noted that the choice of hyperparameters, e.g. the regularisation constant, may also affect the results and the trade-off between the MSE and FID.

In Figure 3, the hidden layer size is a way of parameterising the complexity/capacity of the autoencoder. We see that in general, the FID score decreases with increasing au-

toencoder capacity, and the opposite is approximately true for the determinant of the Pearson correlation matrices. This can be interpreted intuitively as the fact that a more competitive DeepJSCC autoencoder is able to compete more effectively with the adversarial jammer and thus improves the prior distribution matching performance.

As an approach to impose a prior distribution on the latent space, adversarial jamming opens up the possibility for novel contexts in which competition can naturally promote the matching of the distribution, with no requirement for the explicit calculation of distribution divergences. A similar approach to distribution matching, where the Jensen-Shannon divergence is implicitly minimised through an adversarial game between a generator and a discriminator network, has already achieved success in the form of GANs [2].

It is also possible for a range of other distributions to be matched based on different non-Gaussian source distributions. We have seen that in the case of input random variables with isotropic covariance matrices, the jamming noise will be a linear combination of the source distribution and the channel noise distribution. It is therefore possible for the jammer to match the source distribution arbitrarily by specifying different source distributions for the DeepJSCC autoencoder. This will be explored as part of future work.

## 5. CONCLUSION

We have introduced adversarial jamming as a novel approach to distribution matching. This technique implicitly imposes distribution matching by encouraging an encoder to learn a mapping from a data distribution to adversarial jamming noise that disrupts the reconstruction task of an auxiliary joint source-channel coder. It increases the scope of possible settings in which distribution matching can be performed in the latent space of autoencoders and other neural networks. The technique has demonstrated good performance in the case of a diagonal Gaussian prior and could be extended to a range of other prior distributions, by considering different reconstruction problems involving alternative, non-Gaussian, source distributions.

## 6. REFERENCES

[1] Jonathan Ho, Ajay Jain, and Pieter Abbeel, "Denoising Diffusion Probabilistic Models," in *Advances in Neural Information Processing Systems*. 2020, vol. 33, pp. 6840–6851, Curran Associates, Inc.

[2] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, "Generative Adversarial Networks," *Advances in neural information processing systems*, June 2014.

[3] Diederik P Kingma and Max Welling, "Auto-Encoding Variational Bayes," in *2nd International Conference on Learning Representations*, 2014.

[4] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio, "Density estimation using Real NVP," in *International Conference on Learning Representations*, Nov. 2016.

[5] Sam Bond-Taylor, Adam Leach, Yang Long, and Chris G Willcocks, "Deep Generative Modelling: A Comparative Review of VAEs, GANs, Normalizing Flows, Energy-Based and Autoregressive Models," *IEEE Trans Pattern Anal Mach Intell.*, p. 20, 2022.

[6] Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf, "Wasserstein Auto-Encoders," in *International Conference on Learning Representations*, 2018.

[7] Emrah Akyol and Kenneth Rose, "Optimal jamming Over additive noise: Vector source-channel case," in *2013 51st Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, Oct. 2013, pp. 1329–1335.

[8] Emrah Akyol, Kenneth Rose, and Tamer Başar, "On Optimal Jamming Over an Additive Noise Channel," in *52nd IEEE Conference on Decision and Control*, Dec. 2013, pp. 3079–3084.

[9] T. Basar and Y.-W. Wu D., "A complete characterization of minimax and maximin encoder- decoder policies for communication channels with incomplete statistical description," *IEEE Transactions on Information Theory*, vol. 31, no. 4, pp. 482–489, July 1985.

[10] Eirina Bourtsoulatze, David Burth Kurka, and Deniz Gunduz, "Deep Joint Source-Channel Coding for Wireless Image Transmission," *IEEE Transactions on Cognitive Communications and Networking*, vol. 5, no. 3, pp. 567–579, Sept. 2019.

[11] Jialong Xu, Bo Ai, Wei Chen, Ang Yang, Peng Sun, and Miguel Rodrigues, "Wireless Image Transmission Using Deep Source Channel Coding With Attention Modules," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.

[12] Hao Ye, Geoffrey Ye Li, Biing-Hwang Fred Juang, and Kathiravetpillai Sivanesan, "Channel Agnostic End-to-End Learning Based Communication Systems with Conditional GAN," in *2018 IEEE Globecom Workshops (GC Wkshps)*, Dec. 2018, pp. 1–5.

[13] Jialong Xu, Tze-Yang Tung, Bo Ai, Wei Chen, Yuxuan Sun, and Deniz Gunduz, "Deep Joint Source-Channel Coding for Semantic Communications," July 2023.

[14] Alec Radford, Luke Metz, and Soumith Chintala, "Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks," Jan. 2016.

[15] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang, "Deep Learning Face Attributes in the Wild," in *2015 IEEE International Conference on Computer Vision (ICCV)*, Dec. 2015, pp. 3730–3738.

[16] Alex Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," *Technical Report, Department of Computer Science, University of Toronto*, 2009.

[17] Li Deng, "The MNIST Database of Handwritten Digit Images for Machine Learning Research," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.

[18] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter, "GANs trained by a two time-scale update rule converge to a local nash equilibrium," in *Advances in Neural Information Processing Systems*, 2017, vol. 30.