# Adversarial Out-domain Examples for Generative Models

1st Dario Pasquini
*Department of Computer Science*
*Sapienza University*
Rome, Italy
pasquini@di.uniroma1.it

2nd Marco Mingione
*Department of Statistics*
*Sapienza University*
Rome, Italy
marco.mingione@uniroma1.it

3rd Massimo Bernaschi
*Institute for Applied Computing (IAC)*
*CNR*
Rome, Italy
massimo.bernaschi@cnr.it

*Abstract*—Deep generative models are rapidly becoming a common tool for researchers and developers. However, as exhaustively shown for the family of discriminative models, the test-time inference of deep neural networks cannot be fully controlled and erroneous behaviors can be induced by an attacker. In the present work, we show how a malicious user can force a pre-trained generator to reproduce arbitrary data instances by feeding it suitable adversarial inputs. Moreover, we show that these adversarial latent vectors can be shaped so as to be statistically indistinguishable from the set of genuine inputs. The proposed attack technique is evaluated with respect to various *GAN* images generators using different architectures, training processes and for both conditional and not-conditional setups.

*Index Terms*—Generative adversarial models, Attacks against machine learning, Adversarial input

## I. INTRODUCTION

The existence of adversarial inputs has been demonstrated for a quite large set of deep learning architectures [9], [10], [35]. An adversarial input is a carefully forged data instance that aims at leading the model to behave in an incorrect or unexpected way. Moreover, the adversarial setup requires that those instances must be indistinguishable from genuine inputs.

In the present work, motivated by the extensive studies carried out on adversarial inputs for discriminative models, we extend the adversarial context into the increasingly popular generative models field. In particular, we focus on the most promising class of architectures, called *Generative Adversarial Networks* (*GANs*) [8]. *GANs* perform generative modeling of a target data distribution by training a deep neural network architecture. This is composed by two neural networks, a generator and a discriminator that are trained simultaneously in a zero-sum game. In the end, the generator learns a deterministic mapping between a latent representation and an approximation of the target data distribution. What we show with the present work is that a pre-trained generator can be forced to reproduce an arbitrary output if fed by a suitable adversarial input. In particular, our findings show that the data space, defined by the generator, contains data instances having very low probability of lying in the space of the expected outputs (i.e., the target data distribution). We will refer to those outputs as *out-domain examples* and to the relative adversarial inputs as *out-domain latent vectors* or *OLV* in short. Figure 1 shows a set of *out-domain examples* for a generator trained
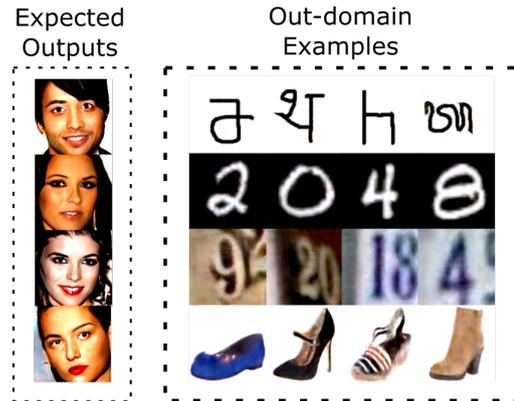


Fig. 1. Comparison between expected generator outputs (left column) and generated out-domain examples (right column) for a Progressive *GAN* generator trained on the *CelebA* dataset.

by using a Progressive GAN [14]. In that example, we found a set of inputs capable to force the generative model[1] to produce images completely different from those belonging to the generator training dataset, i.e., the *CelebA* dataset [22]. Moreover, we show that those *OLV* can be forged in order to be *statistically* indistinguishable from known-to-be trusted inputs. The existence of such adversarial inputs raises new practical questions about the use of *GAN* generators in an untrustworthy environment, as a web application. The main contributions of the present paper may be summarized as follows:

1) We show that a generator may be forced to produce *out-domain* data instances which are arbitrarily different from those for which the generator is trained. Our experiments refer to three common image datasets and for standard and conditional *GAN* architectures: *Deep Convolutional GAN* (*DCGAN*) [29] and *Auxiliary Classifier GAN* (*ACGAN*) [28].
2) We propose a first type of adversarial input for not encoder-based generative models.
3) We investigate the nature of out-domain examples showing that their quality strongly depends on the dimension of both the latent and the data space.

[1] A pre-trained *Progressive GAN* available at https://tfhub.dev/google/progan-128/1.

## II. BACKGROUND

The objective of a generative model is to learn a probability distribution $\tilde{p}_X$ that approximates a target data probability distribution $p_X$. Actually, in general, $p_X$ is unknown and it can only be inferred by a limited set of samples. One of the most powerful approaches to train a generative model is the recently proposed *Generative Adversarial Networks* (*GANs*) framework. *GANs* require the simultaneous training of two neural networks: a generator $G$ and a discriminator $D$. $D$'s objective is to maximize the probability of discriminating between $\tilde{p}_X$ and $p_X$, whereas $G$'s objective is to make $\tilde{p}_X$ and $p_X$ indistinguishable in order to mislead $D$. This kind of unsupervised training process is renamed *Adversarial Training* in this context.

From the mathematical point of view, both the generator and the discriminator can be intended as functions: $G : \mathcal{Z} \to \mathcal{X}$ and $D : \mathcal{X} \to [0, 1]$. In other words, during the training process, the neural network $G$ receives as input a vector $z = (z_1, \ldots, z_n)$ and produces $x = G(z)$. Each element of $z$ is a realization of a random vector $Z_1, \ldots, Z_n$, with $Z_i \overset{i.i.d.}{\sim} p_Z$, where $p_Z$ is an arbitrary density function. Then, the optimization problem can be easily summarized by the following formulation:

$$\min_{\theta_G} \max_{\theta_D} \langle \mathbb{E}[\log(D(x))] + \mathbb{E}[\log(D(1 - G(z)))] \rangle \quad (1)$$

where $\theta_G$ and $\theta_D$ are, respectively, the parameters of the generator and the discriminator network and $x$ is an instance from the training set $X$.

The use of random latent instances as input of $G$ makes possible to explore the data space by generating new data instances, not necessarily available in the training set.

Many extensions to the original *GAN* framework have been successfully developed [2], [7], [28], [29]. One of the most influential work is [29], in which the *DCGAN* architecture was proposed. *DCGAN* is capable of exploiting the potential of *Convolutional Neural Networks (CNNs)* in both the generator and discriminator perspective. The *GAN* framework can be easily extended to train conditional generators [25] using the *ACGAN* architecture [28]. In this case, a supervised training approach is used to the purpose of learning a probability distribution conditioned to a set of classes $Y$. During this training process, the class $y$ is chosen randomly (e.g., with uniform probability) from the set $Y$ including all the possible classes. Then, the generator is modeled as a bivariate function that receives as input an instance of the latent space $z$ labelled with its related class $y \in Y$. *ACGAN* architecture improves the performance of the training process by adding an auxiliary classification task to the discriminator. The latter outputs two probability distributions, the first over the source (i.e., the probability that the instance comes from $p_X$) and the second over the class labels (i.e., the probability that the instance belongs to class $y$). In this case, the optimization problem can be parametrized by extending (1) as follows:

$$L_{\text{source}} = \mathbb{E}[\log(D^{\text{source}}(x))] + \mathbb{E}[\log(D^{\text{source}}(1 - G(z)))]$$

$$L_{\text{class}} = \mathbb{E}[\log(D^{\text{class}}(y|x))] + \mathbb{E}[\log(D^{\text{class}}(y|G(z, y)))]$$

$$L_D = L_{\text{class}} + L_{\text{source}} \qquad L_G = L_{\text{class}} - L_{\text{source}} \quad (2)$$

where $L_D$ and $L_G$ are the loss functions of the discriminator and the generator, respectively.

## III. RELATED WORKS

Some examples of adversarial inputs for encoder-based generative models like *Variational AutoEncoder* (*VAE*) and *VAE-GAN* are analyzed in [16]. In the proposed scenario, given a data instance $x$, an attacker aims at producing an instance $\hat{x}$ that differs in a limited way from $x$ but which is capable of driving the *VAE* to reconstruct a $\hat{x}$ far from the original $x$. The reconstructed $\hat{x}$ can be an approximation of an arbitrary data instance chosen from the attacker. All the results described by the authors refer to $\hat{x}$ as if it came from the same data distribution on which the *VAE* is trained. A similar attack scenario has been investigated also in [34]. At the best of our knowledge no other form of adversarial input targeting *GAN* generators has been proposed.

Many works investigate the possibility of finding or exploiting an inverse mapping from the data space to the latent space of a pre-trained generator [4], [20], [27]. In [27] a pre-trained generator $G$ is used to invert a discriminative model $C$ to the purpose of synthesizing novel images. It is noteworthy that the authors report a first case of partial out-domain example. In particular, a generator trained on *ImageNet* was able to reproduce images belonging to classes known to $C$ but unknown to $G$.

Additionally, a technique to map images into a latent representation with a pre-trained generator is proposed in [4] and in [20]. In those works, the authors mention the possibility of mapping images, which are not present in the training set, into the generator latent space. That is shown only for images coming from the same distribution on which the model is trained. The proposed inversion technique is essentially the same used in the present work and it is based on the direct optimization of the generator input by a gradient descent based approach. In [4] during the optimization process, the latent vectors are encouraged to be *similar* to those of the latent prior distribution by adding a penalty term in the loss function. This term is a weighted sum of the discrepancy between the mean and standard deviation of the latent vector and the latent prior distribution. We exploited the possibility of extending this penalty term beyond the second moment, given that just two moments might not be sufficient to correctly identify the latent vectors.

The same generator inversion technique is used in [31] as a defense against adversarial examples. In this work, the adversarial examples are mapped to unperturbed data instances by inverting a generator trained on the data distribution of clear data. The proposed model inversion aims at finding the closest generator codomain element for each input of a discriminative model $f$. At the end, these codomain elements are taken as input by $f$ instead of the original untrusted data.

A similar technique is also used in [21] to perform a data

membership attack against generative models. That kind of attack aims at inferring the presence of a data instance $x$ in the training-set used during the training of a generative model $G$. In this case, the generator inversion is carried out by a neural network attacker called $\mathcal{A}$. This attacker is trained as an encoder for $G$. Given a data instance $x$, the latent vector $z = \mathcal{A}(x)$ is used to estimate the chances of $x$ of being in the $G$'s dataset by calculating the distance between $G(z)$ and $x$.

## IV. OUT-DOMAIN EXAMPLES

Let $\mathcal{X}$ be the set of all possible data that can be generated by $G$ and let $\mathcal{Z}$ be the set of all the possible latent vectors coming from $p_Z$. Then, an out-domain example for a generator $G$ is defined to be an element $\hat{x} \in \mathcal{X}$ such that:

$$G(\hat{z}) = \hat{x}$$
$$p_X(\hat{x} \in \mathcal{X}) \leq \epsilon_x \quad \text{and} \quad p_Z(\hat{z} \notin \mathcal{Z}) \leq \epsilon_z \quad (3)$$

where $\hat{z}$ is the $OLV$ used to generate an out-domain example $\hat{x}$ and both $\epsilon_x$ and $\epsilon_z$ are negligible probabilities. Hence, the underlying assumption in (3) is that the probability of $\hat{x}$ of belonging to the set of expected outputs is very low. We will refer to the set of expected outputs of a generator with the term *domain*. The domain of a generator can also be intended as the semantic contents defined by $p_X$[2]. To the purpose of finding a suitable out-domain latent vector $\hat{z}$, we choose a target instance $\dot{x}$ and then, we look for the $\hat{z} \in \mathcal{Z}$ such that there is the minimum distance between $\dot{x}$ and $\hat{x} = G(\hat{z})$. We refer to this process with the expression *latent search*. Coherently with (3), the target instance $\dot{x}$ is chosen *ad hoc* to be out of the generator's domain. This scenario, is depicted in Figure 2 where a set of out-domain examples (panel **c**) from a *DCGAN* generator trained on *CIFAR10* is reported. In this case, 25 target instances have been randomly chosen from five datasets different from *CIFAR10*.

In addition, as required by (3), an out-domain latent vector $\hat{z}$ is considered a valid input for the generator $G$ if it lies in a *dense* region of the latent space. This implies that $\hat{z}$ must be statistically indistinguishable from a valid latent vector sampled from the latent probability distribution $p_Z$. In the adversarial perspective, this means that a defender is unable to tell apart a valid latent vector from an out-domain latent vector before the generation of $\hat{x}$.

### A. Motivating adversarial scenario

The recent success of generative models in the scientific [5], [6] and in the entertainment field [12], [36], inspired the development of many *GANs* based software applications. These are often in the form of a web service with an interactive interface by which the user provides direct or indirect input to the model[3]. Assuming a white-box access to the generator model, an attacker can find out-domain latent vectors capable of driving the service to produce inadequate contents such as pornographic or offensive material. The attacker can use these

out-domain examples in order to perform a very effective and straightforward defacing attack direct to the generator owner. Indeed, this type of web application and software allows to share and save internal copies of the images created by the users. This scenario resembles a reflected or stored *Cross-site Scripting (XSS)* attack where the attacker is able to arbitrary modify an image in the web page. The white-box assumption is supported by the observation that often these applications, in order to reduce the server load, run the generative model in the client-side and additionally pre-trained versions of open-source generator are frequently used.

We assume that the owner (referred as defender) performs a validation process on $\hat{z}$ before the calculation of $\hat{x} = G(\hat{z})$. This validation can be intended as a function $\upsilon : \mathcal{Z} \to \{0, 1\}$. Therefore, the defender accepts to calculate $\hat{x}$ if and only if $\upsilon(\hat{z})$ is equal to 1. In our attack scenario, this function $\upsilon$ is represented by a distributional hypothesis test. The null hypothesis $(H_0)$ is that the vector $\hat{z}$ is sampled from $p_Z$. Thus, given a test statistic $t$ and for fixed *type I error* $\alpha$, the decision rule can be formalized as follows:

$$Pr(T \geq t | H_0) \geq \alpha \Rightarrow \upsilon(\hat{z}) = 1 \quad (4)$$

where $T$ is the distribution of the test statistics under $H_0$ and $Pr(T \geq t | H_0)$ corresponds to the classic $p$-value of confirmatory data analysis. The same scenario can be easily extended to conditional generators. In that case, we assume that the defender is able to arbitrary choose and fix a data class $y \in Y$. The attacker aims at finding a suitable out-domain latent vector for the conditioned generator $G(\cdot | Y = y)$.

## V. METHODOLOGY

As mentioned in Section IV, out-domain examples can be found by looking for the closest representation of an arbitrary chosen target instance in the data space defined by $G$. Actually, by leveraging the differentiable nature of the generator and the structure of a well formed latent-data mapping [29], we can transform this searching problem in an efficient optimization process as follows:

$$L(\dot{x}, z) = d(\dot{x}, \ G(z)) + \rho(z)$$
$$\hat{z} = \operatorname*{argmin}_{z \in \mathcal{Z}} L(\dot{x}, z) \quad (5)$$

Where $\dot{x}$ is a given target instance, $d$ is a distance function and $\rho$ is a penalty term applied to $z$. The purpose of a penalty term is to force the solution $\hat{z}$ to be consistent with $p_Z$. More precisely, $\rho$ is defined as the weighted sum of the squared difference of the first $k$ sample moments of $z$ and the theoretical moments of a random variable $Z \sim p_Z$.

$$\rho(z) = \sum_{i=1}^{k} \omega_i \| \mu_Z(i) - \tilde{\mu}_z(i) \|_2^2 \quad (6)$$

Where $\mu_Z(i)$ is the $i^{th}$ moment of $Z$ and $\tilde{\mu}_z(i)$ is the $i^{th}$ sample moment of the latent vector $z$. The parameter $\omega_i$ is the weight assigned to the $i^{th}$ moment difference. In the case of conditional generators, the searching process is performed by

---

[2]For instance, the domain of the *MNIST* dataset is the set of digits representation and the domain of *CelebA* dataset is a set of human faces.

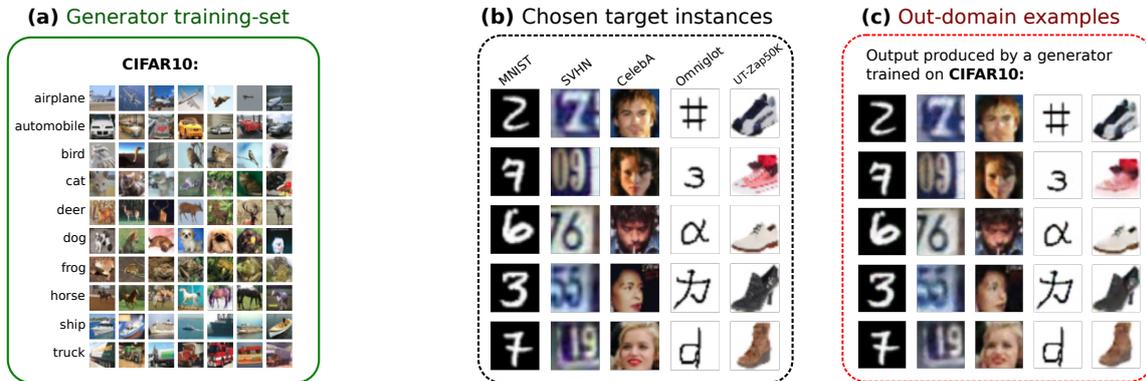[3]One example can be found here: https://make.girls.moe.

Fig. 2. Target images (panel **b**) and reproduced out-domain examples (panel **c**) generated by a DCGAN generator trained on the *CIFAR10* dataset (panel **a**). Target images have been randomly chosen from five common image datasets different from *CIFAR10*. The generated out-domain examples are visually close to their respective targets even if there is not intersection between *CIFAR10* and the datasets of the chosen target instances.

fixing a class $y$ as input of the generator function. This implies that the optimization process acts only on the latent vector $z$ and cannot modify the class representation $y$. More formally, in the conditional case, the problem can be reformulated as:

$$L(\dot{x}, z) = d(\dot{x}, \ G(z, \ y)) + \rho(z)$$
$$\hat{z} = \underset{z \in \mathcal{Z}}{\operatorname{argmin}} L(\dot{x}, z) \tag{7}$$

assuming that $y$ is randomly chosen from $Y$ by the defender.

Starting from a random initialization of $z$, say $z_0$ obtained by sampling from $p_Z$, we iteratively update the current latent vector according to the following rule:

$$z_{i+1} = z_i + \eta \nabla L(\dot{x}, z_i), \quad i = 1, \dots, N \tag{8}$$

where $\eta$ is the learning rate. At each iteration of the optimization process, the distance function $d$ is computed between the target $\dot{x}$ and $G(z_i) = \hat{x}_i$. We tested and compared two distance functions: the mean squared error (*MSE*) and the cross entropy (*XE*). It is noteworthy that in the cross entropy case, the *softmax* function is used in order to ensure the unitary sum in both $\dot{x}$ and $G(z)$. However, given its not-bijectivity, we force the comparison between the target and generated image to be scale invariant. Nonetheless, although this approach diverts from the original objective of founding the closest codomain instance to $\dot{x}$, *XE* is able to provide a very good approximation (at least in the visual form) of $\dot{x}$ with fewer training iterations than *MSE*.

The penalty term is used to ensure the *indistinguishability* of the out-domain vector from a trusted input. The main objective is to find a $\hat{z}$ such that the probability of $\upsilon(\hat{z}) = 1$ is maximized. This can be obtained by forcing the out-domain latent vector to have moments equal to those of a random variable distributed as $p_Z$. Indeed, in probability theory, Moment Generating Functions (MGFs) have great practical relevance not only because they can be used to easily derive moments, but also because they uniquely identify a probability distribution function, a feature that makes them a handy tool

to solve several problems. The MGF (if it exists) can be seen as an alternative specification useful to characterize a random variable. On one hand, the MGF can be used to compute the $n^{th}$ moment of a distribution as the $n^{th}$ derivative of the MGF evaluated in 0 On the other hand, a MGF is useful to compare two different random variables, studying their behavior under limit conditions. Given a random variable $X$, its MGF is defined as the expected value of $e^{tX}$:

$$M_X(t) = \mathbb{E}(e^{tX}), \quad t \in \mathbb{R} \tag{9}$$

If (9) holds, then the $n^{th}$ moment of $X$, denoted by $\mu_X(n)$, exists and it is finite for any $n \in \mathbb{N}$:

$$\mu_X(n) = \mathbb{E}(X^n) = \left. \frac{\partial^n M_X(t)}{\partial t^n} \right|_{t=0} \tag{10}$$

## VI. RESULTS

In our experiments, we tested and compared two common prior distributions, i.e., the standard normal and the continuous uniform distribution in $[-1, 1]$. Given the constraint imposed by the latter, we perform a hard clipping on $z$ values in order to force the latent vector to lie in the allowed hypercube. As proposed in [20], we tested the *stochastic clipping* method but results showed no substantial improvement.

We did not apply any clipping method for the normal prior distribution. Empirically, it has been observed that the penalty on the moments is sufficient to guarantee that $z$ assumes values in an acceptable range.

The quality of the out-domain examples is evaluated on different *DCGAN* generators and on conditional generators trained within an *ACGAN* framework. For the sake of exposition, we will refer to each trained generator with the following compact notation:

*[Architecture]-[Training-dataset]-[Latent_prior]-*
*[Latent_space_dimension]*

In particular, a generator is trained for any combination of the followings:

- **Architecture:** *DCGAN, ACGAN*

| Dataset | $Z$ dim. | (a) **Normal Latent distribution** | | | (b) **Uniform Latent distribution** | | |
|---|---|---|---|---|---|---|---|
| | | **Avg MSE** | **Test Succ.** | Avg MSE* | **Avg MSE** | **Test Succ.** | Avg MSE* |
| CIFAR10 | 100 | 0.010646 | 100% | 0.008881 | 0.012131 | 100% | 0.009354 |
| CIFAR10 | 256 | 0.005094 | 100% | 0.003902 | 0.012131 | 100% | 0.009354 |
| CIFAR10 | 512 | 0.003693 | 100% | 0.002710 | 0.003603 | 100% | 0.002826 |
| SVHN | 100 | 0.018569 | 100% | 0.011541 | 0.016730 | 100% | 0.011359 |
| SVHN | 256 | 0.011374 | 100% | 0.006970 | 0.012121 | 100% | 0.007213 |
| SVHN | 512 | 0.009474 | 100% | 0.005314 | 0.008323 | 100% | 0.005594 |
| C.MNIST | 100 | 0.063097 | 100% | 0.040453 | 0.059342 | 100% | 0.037457 |
| C.MNIST | 256 | 0.052926 | 100% | 0.029160 | 0.045851 | 100% | 0.027178 |
| C.MNIST | 512 | 0.043685 | 100% | 0.025855 | 0.037946 | 99% | 0.022415 |

- **Training Dataset:** *CIFAR10* [17], *SVHN* [26] and a simple variation of *MNIST* [19], called *ColorMNIST*
- **Latent space dimension:** $\{100, 256, 512\}$
- **Latent prior distribution:** $N(0, 1)$ and $U[-1, 1]$

For instance, *DCGAN-CIFAR10-Normal*-100 defines a *DCGAN* generator trained on *CIFAR10* with a normal latent prior distribution and latent space dimension equal to 100. The *ColorMNIST* dataset is obtained by applying a random background color to the original *MNIST*. The reason of that modification is to offer to the generator the chance of representing a larger set of outputs by letting the generator learn a larger number of *RGB* triplets, while keeping virtually unaltered the complexity of *MNIST*. All the generators and discriminators' architectures as well as the hyper-parameters and the training process are the same proposed in [29]. We tested three values of the latent space dimension that are commonly used in literature.

The validation process of the out-domain vectors is performed by fixing $\alpha = 0.05$, $k = 4$ for the normal prior and $k = 6$ for the uniform prior. We performed three different distributional tests, i.e., Kolmogorov-Smirnov, Shapiro-Wilk and Anderson-Darling [32]. Results showed that, given the penalty term $\rho$, all the distributional tests bring to the same decision. The following results refer to the Anderson-Darling test [1], which was finally chosen since its test statistics is based on the *Cumulative Distribution Function CDF* [30] and, compared to other tests, it places more weight on the values in the tails of the distribution.

We defined a test-set of target instances to the purpose of evaluating the capability of different generators to reproduce out-domain examples. This test-set contains randomly chosen instances from four image datasets i.e, *Omniglot* [23], *CelebA* [22], *UT-Zap50K* [37] and *Tiny ImageNet* [33]. A random sample of 32 images is selected for each dataset for a total of 128 target instances. In the case of *Tiny ImageNet*, the images are sampled from classes which are different from those of *CIFAR10*. All images are forced to share the same dimension of $32 \times 32$ pixels and to be normalized in the interval $[0, 1]$. To simplify the understanding of the results, the target distance function used for all the experiments is the *Mean Squared Error* (*MSE*). The average *Mean Squared Error MSE* and the percentage of successfully passed statistical tests are computed on the test-set and used as main evaluation scores. In the latent search process, the *Adam* [15] optimizer is used with a learning rate equal to 0.01. All the experiments are performed using the *TensorFlow* framework [3]. The most relevant codes used for the present work along with an interactive proof of concept are available on: https://github.com/pasquini-dario/OutDomainExamples.

### A. DCGAN

Table I shows the results related to the *DCGAN* generators. All the forged *OLV* pass successfully the distributional test, regardless of the chosen prior distribution. Several checks on the biases of the *OLV* have been carried out providing pretty good results. In particular, the estimation of the odd moments for the normal is precise; the second and the fourth moment are instead slightly overestimated and underestimated, respectively. For the uniform prior, the bias for the second and the sixth moment is slightly positive and for the remainings the estimation is precise, with a quite large variance in the estimation of the third moment. No evident patterns are worth to notice when looking at the training dataset or at changes of the latent space dimension for both the prior distributions. These results are due to the fact that the penalty term $\rho$ strongly constraints the values of the latent vectors in a well defined range. It is worth to notice that relaxing the moments penalty during the optimization process (5) would reduce further the *MSE*. In contrast to in-domain inversion [4], we can state that out-domain examples do not take any significant advantage from latent vectors statistically close to those used during the training process.

Figure 3 shows a set of target instances and related out-domain examples for a total of six generators. The upper panel reports the out-domain examples produced by three generators trained on different training sets but with same latent space dimension and prior. When the training set of the generator is *ColorMNIST*, the method fails in finding suitable *OLV* capable of reproducing the target images. For the other two, the generator is able to provide a valid reconstruction for all
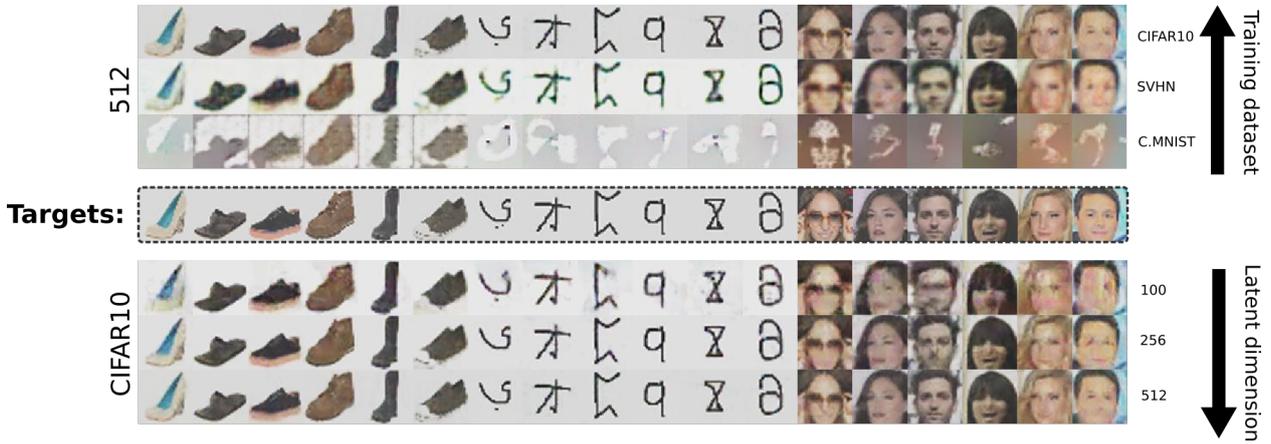
Fig. 3. Out-domain examples for a set of *DCGAN* architectures trained with uniform latent prior. The central row shows several different randomly chosen targets from the test-set. The upper panel shows the variability in the out-domain generation when the latent space dimension is fixed to 512 and the training dataset of the generator varies. The lower panel shows the variability in the out-domain generation when the training dataset of the generator is fixed to *CIFAR10* and the latent space dimension varies.
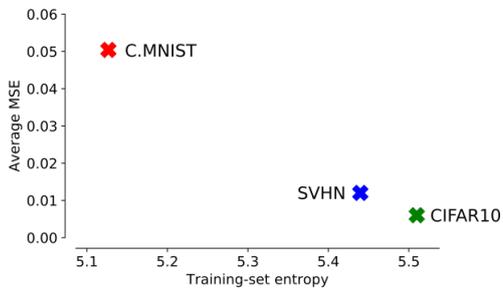


Fig. 4. Average *MSE* compared to the estimated entropy of each training set.
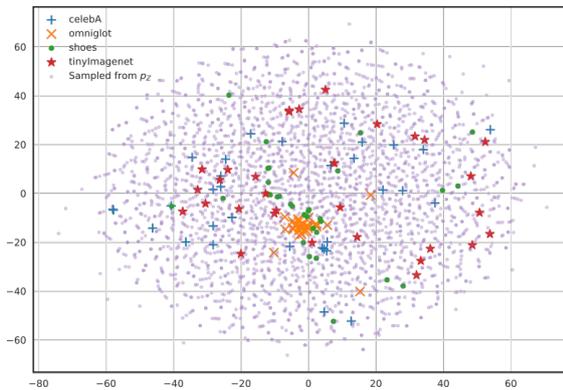


Fig. 5. Two-dimensional representation of out-domain latent vectors for the *DCGAN-CIFAR10-Uniform*-100. Out-domain latent vectors with target instances coming from the same dataset are represented with the same marker.

the targets. The failure of *ColorMNIST* may be connected to the fact that it is less *heterogeneous* with respect to *SVHN* and *CIFAR10*. By heterogeneity we intend the actual number of different pixels which are necessary in order to reproduce the same heterogeneity of the whole set (i.e., the entropy). It is reasonable to expect that the larger the variety of images in the training set, the larger will be the set of potential out-domain examples reproduced by the generator. As an estimator of that variety, we computed the Shannon entropy [13] for a sample of $2^{10}$ images from each training set. Results are depicted in Figure 4 and show that there is a strong dependence between the average *MSE* (i.e., the reconstruction error) and the entropy of the training set.

Figure 5 depicts a two-dimensional projection of a set of out-domain latent vectors and latent vectors directly sampled by $p_Z$. This representation is obtained by applying the dimension reduction algorithm called *t-distributed Stochastic Neighborhood Embedding* (*t-SNE*) [24] on vectors of size 100. It is possible to note how the out-domain latent vectors tend to be uniformly distributed in the space. In the case of the *Omniglot*

dataset, the *OLV* tend to cluster in a specific region and this may be due to its intrinsic homogeneity.

Even if the entropy is a sort of predictor of the success of our method, it is still possible, given a target image, a latent prior and a training set, to enhance the quality of the generated image by increasing the dimension of the latent space. As a matter of fact, we can observe, by looking at the lower panel of Figure 3, that an increase of the latent space dimension makes the generated image more similar to the target one. An additional motivation can be that the latent space acts as an information bottleneck for the target instance during the latent search process All these possibilities are evaluated in terms of *MSE*. Figure 6 shows the average *MSE* for each latent space dimension, latent prior and training set confirming that the more is the entropy of the training set, the higher the probability of success in the generation and, at the same time, the larger the dimension of the latent space, the higher the quality of the reconstruction. Instead, there is no relevant difference in the quality of the out-domain examples when the latent prior distribution varies.
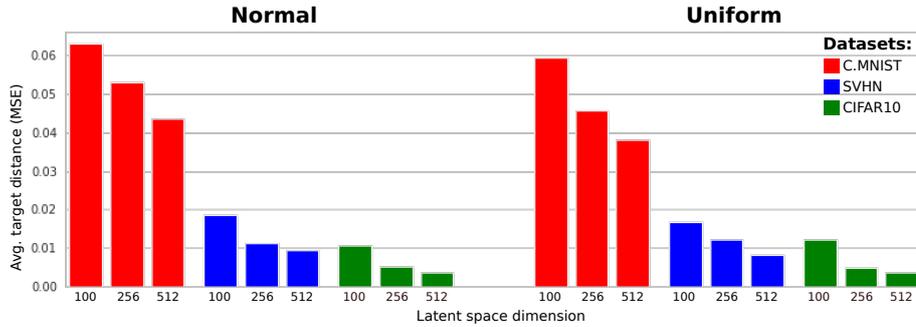
Fig. 6. Average *MSE* for all the *DCGAN* generators trained on all combinations of dataset, latent space dimension and latent prior distribution.
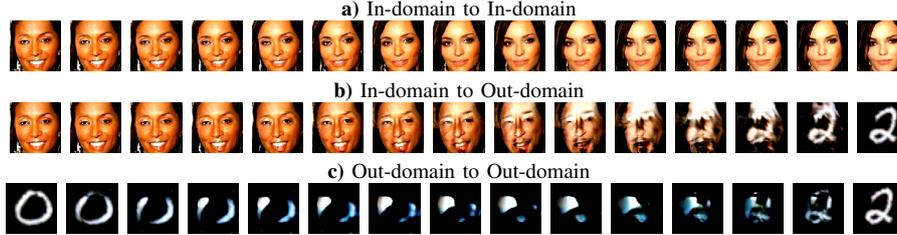


Fig. 7. Three examples of linear interpolation between two latent vectors for a *ProGAN* trained on the *CelebA* dataset. The row (a) depicts the interpolation process between two randomly chosen latent vectors. Row (b) depicts interpolation from a randomly chosen latent vector and an out-domain latent vector. Row (c) depicts interpolation between two out-domain latent vectors.

Figure 7 shows three examples of linear interpolation between latent vectors [29]. The first row depicts a smooth and semantic meaningful transaction between two random vectors sampled from $p_Z$, referred as *in-domain* latent vectors. By semantic meaningful transaction, we mean that each image between the two interpolation points remains coherent with $p_X$. The second row depicts the interpolation between an in-domain vector and an out-domain vector. In contrast with the first case, the transaction is unbalanced and not particularly smooth. From the sequence, it can be noticed that the semantic valid attributes of the starting image, i.e. the black of the hair and the reflection on the forehead, are deformed to recreate the final *MNIST* digit. The last row shows the extreme case of interpolation between two out-domain latent vectors. In this case, all the intermediate data instances never cross the in-domain set.

### B. ACGAN

Conditional generators are trained to the purpose of enforcing their outputs to be part of a meaningful, from the semantic viewpoint, data class. Typically, this implies a better global coherence and quality in the definition of the generator's data space [28]. This is especially true for models trained with the *ACGAN* framework in which the generators are encouraged to produce images that are correctly classified from the discriminator as genuine and belonging to its class. The experiments described below aim at finding out if the conditional extension is sufficient to the purpose of preventing the generation of out-domain examples. We trained different *ACGAN* generators using the same set of parameters reported in Section VI-A.

Also the architecture used for the generator and the discriminator is the same used for the previous *DCGAN* experiments[4]. Training hyper-parameters are the same proposed in [28]. As aforementioned, in the conditional setup, the hypothesis is that the class $y$ is randomly chosen by the defender and the attacker can not modify its representation during the latent search process. Tests and validation are performed as in Section VI-A but they are evaluated conditionally to each class $y$ in the generators' classes set. All the tested training sets are composed by 10 classes. In this case, the *MSE* is calculated as the average over the classes.

We are able to find an out-domain example for each image in the test-set, conditionally to each class. We do not report the results when the training set is *ColorMNIST*, since it already failed in the less severe *DCGAN* experiment. As an example, Figure 8 shows the generated out-domain examples, conditionally to each class of *CIFAR10*, for four randomly chosen target instances in the test-set. It can be noticed that the class has no relevant impact on the quality of the out-domain examples: the attack succeeds regardless of the class. The same happens when attacking the generators trained on *SVHN*. Results in terms of *MSE* are summarized in Table II. It is possible to observe that the average *MSE* is uniformly larger compared to the *DCGAN* experiments due to the conditional setup. In Figure 9, we also report the distribution of the *MSE*, conditionally to each class, for each training set. No specific patterns are registered: the *MSE* distribution is approximately

---

[4]The only difference is the number of neurons in the generator's input layer and in the discriminator's output layer due to the conditional setup

| Dataset | $Z$ dim. | (a) **Normal Latent distribution** | | (b) **Uniform Latent distribution** | |
|---|---|---|---|---|---|
| | | **Avg MSE** | **Test Succ.** | **Avg MSE** | **Test Succ.** |
| CIFAR10 | 100 | 0.023457 | 100% | 0.019615 | 100% |
| CIFAR10 | 256 | 0.013144 | 100% | 0.009547 | 100% |
| CIFAR10 | 512 | 0.009075 | 100% | 0.005944 | 99% |
| SVHN | 100 | 0.026686 | 100% | 0.024732 | 100% |
| SVHN | 256 | 0.016879 | 100% | 0.016268 | 100% |
| SVHN | 512 | 0.016291 | 100% | 0.013707 | 99% |



Fig. 8. Comparison between out-domain examples produced by an *ACGAN-CIFAR10-Normal*-512. Each (but the first) column depicts the out-domain example produced by the generator, conditionally to each class.



Fig. 9. *MSE* distribution conditionally to each class for *CIFAR10* and *SVHN* datasets.

the same for each class and training set. However, there is a slight variability for *CIFAR10* given the larger heterogeneity among its classes. The validation of the out-domain latent vectors is the same described in Section VI-A. Also in this case, all the latent vectors pass successfully the Anderson-Darling test. Moments distributions for each dataset, latent prior and latent space dimension are also checked. No relevant difference has been observed with respect to the not-conditional generators experiments.

## VII. CONCLUSION AND FURTHER DEVELOPMENTS

We showed how to forge suitable adversarial inputs capable of driving a trained generator to produce arbitrary data instances in output. This is possible for both conditional and not-conditional generators. Additionally, we showed that an adversarial input can be shaped in order to be statistically indistinguishable from the set of trusted inputs. We also showed that the success of our method strongly depends on two main factors: the heterogeneity of the set on which the generator is trained and the latent space dimension.
In additional experiments we found a set of generators showing a greater resilience to the generation of out-domain examples. In particular, the *Non-saturating GAN*s with *ResNet* architecture analyzed in [18] shows an inherent difficulty to produce out-domain examples even when the generator is trained on high entropic datasets such as *CIFAR10*. We conjecture that this property is strongly related to the generator architecture.
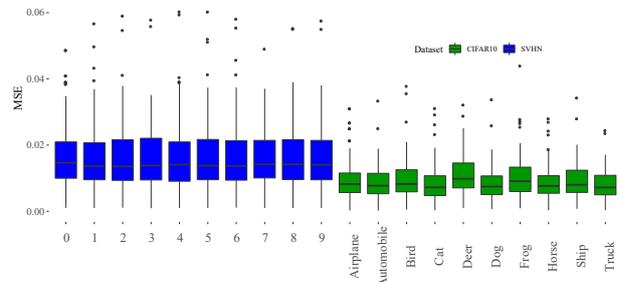In the described adversarial scenario, we supposed that an

aware defender can just test the validity of the model's input in order to evaluate the genuineness of the latent vectors. However, it is possible to imagine a more powerful defender able to verify the generator's output in order to spot unexpected generation.
As future directions of activity, we expect to *i)* investigate the generation of out-domain examples for other *GAN* architectures; *ii)* study the generation of out-domain examples in contexts other than those of images; *iii)* investigate the possibility of training an arbitrary complex generator which is resilient to the generation of out-domain examples; *iv)* evaluate the possibility of extending the attack to a black-box scenario using an approach inspired by [11].

## REFERENCES

[1] T. W. Anderson and D. A. Darling. A test of goodness of fit. *Journal of the American statistical association*, 49(268):765–769, 1954.
[2] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein GAN. *ArXiv e-prints*, Jan. 2017.
[3] M. A. at al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
[4] A. Creswell and A. A. Bharath. Inverting the generator of a generative adversarial network (ii). *arXiv preprint arXiv:1802.05701*, 2018.
[5] N. De Cao and T. Kipf. MolGAN: An implicit generative model for small molecular graphs. *ArXiv e-prints*, May 2018.
[6] L. de Oliveira, M. Paganini, and B. Nachman. Learning Particle Physics by Example: Location-Aware Generative Adversarial Networks for Physics Synthesis. *ArXiv e-prints*, Jan. 2017.
[7] J. Donahue, P. Krähenbühl, and T. Darrell. Adversarial feature learning. *arXiv preprint arXiv:1605.09782*, 2016.

[8] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[9] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and Harnessing Adversarial Examples. *ArXiv e-prints*, Dec. 2014.

[10] S. H. Huang, N. Papernot, I. J. Goodfellow, Y. Duan, and P. Abbeel. Adversarial attacks on neural network policies. *CoRR*, abs/1702.02284, 2017.

[11] A. Ilyas, L. Engstrom, A. Athalye, and J. Lin. Black-box adversarial attacks with limited queries and information. *CoRR*, abs/1804.08598, 2018.

[12] Y. Jin, J. Zhang, M. Li, Y. Tian, H. Zhu, and Z. Fang. Towards the Automatic Anime Characters Creation with Generative Adversarial Networks. *ArXiv e-prints*, Aug. 2017.

[13] L. Jost. Entropy and diversity. *Oikos*, 113(2):363–375, 2006.

[14] T. Karras, T. Aila, S. Laine, and J. Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *CoRR*, abs/1710.10196, 2017.

[15] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

[16] J. Kos, I. Fischer, and D. Song. Adversarial examples for generative models. *CoRR*, abs/1702.06832, 2017.

[17] A. Krizhevsky. Learning multiple layers of features from tiny images. 05 2012.

[18] K. Kurach, M. Lucic, X. Zhai, M. Michalski, and S. Gelly. The GAN landscape: Losses, architectures, regularization, and normalization. *CoRR*, abs/1807.04720, 2018.

[19] Y. LeCun and C. Cortes. MNIST handwritten digit database. 2010.

[20] Z. C. Lipton and S. Tripathi. Precise recovery of latent vectors from generative adversarial networks. *CoRR*, abs/1702.04782, 2017.

[21] K. S. Liu, B. Li, and J. Gao. Generative model: Membership attack, generalization and diversity. *CoRR*, abs/1805.09898, 2018.

[22] Z. Liu, P. Luo, X. Wang, and X. Tang. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

[23] B. M Lake, R. Salakhutdinov, and J. B Tenenbaum. Human-level concept learning through probabilistic program induction. 350:1332–1338, 12 2015.

[24] L. v. d. Maaten and G. Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(Nov):2579–2605, 2008.

[25] M. Mirza and S. Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[26] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011.

[27] A. Nguyen, J. Clune, Y. Bengio, A. Dosovitskiy, and J. Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. In *CVPR*, page 7, 2017.

[28] A. Odena, C. Olah, and J. Shlens. Conditional Image Synthesis With Auxiliary Classifier GANs. *ArXiv e-prints*, Oct. 2016.

[29] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *CoRR*, abs/1511.06434, 2015.

[30] S. Ross. *A first course in probability*. Pearson, 2014.

[31] P. Samangouei, M. Kabkab, and R. Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *CoRR*, abs/1805.06605, 2018.

[32] S. S. Shapiro. *How to test normality and other distributional assumptions*, volume 3. ASQC Milwaukee, WI, 1990.

[33] Stanford. Tiny imagenet visual recognition challenge.

[34] P. Tabacof, J. Tavares, and E. Valle. Adversarial images for variational autoencoders. *CoRR*, abs/1612.00155, 2016.

[35] C. Xie, J. Wang, Z. Zhang, Y. Zhou, L. Xie, and A. L. Yuille. Adversarial examples for semantic segmentation and object detection. *CoRR*, abs/1703.08603, 2017.

[36] L. Yang, S. Chou, and Y. Yang. Midinet: A convolutional generative adversarial network for symbolic-domain music generation using 1d and 2d conditions. *CoRR*, abs/1703.10847, 2017.

[37] A. Yu and K. Grauman. Fine-grained visual comparisons with local learning. In *Computer Vision and Pattern Recognition (CVPR)*, Jun 2014.