

Spatial Evolutionary Generative Adversarial Networks

Jamal Toutouh
Massachusetts Institute of
Technology, CSAIL
toutouh@mit.edu

Erik Hemberg
Massachusetts Institute of
Technology, CSAIL
hembergerik@csail.mit.edu

Una-May O'Reilly
Massachusetts Institute of
Technology, CSAIL
unamay@csail.mit.edu

ABSTRACT

Generative adversary networks (GANs) suffer from training pathologies such as instability and mode collapse. These pathologies mainly arise from a lack of diversity in their adversarial interactions. Evolutionary generative adversarial networks apply the principles of evolutionary computation to mitigate these problems. We hybridize two of these approaches that promote training diversity. One, **E-GAN**, at each batch, injects mutation diversity by training the (replicated) generator with three independent objective functions then selecting the resulting best performing generator for the next batch. The other, **Lipizzaner**, injects population diversity by training a two-dimensional grid of GANs with a distributed evolutionary algorithm that includes neighbor exchanges of additional training adversaries, performance based selection and population-based hyper-parameter tuning. We propose to combine mutation and population approaches to diversity improvement. We contribute a superior evolutionary GANs training method, **Mustangs**, that eliminates the single loss function used across **Lipizzaner**'s grid. Instead, each training round, a loss function is selected with equal probability, from among the three **E-GAN** uses. Experimental analyses on standard benchmarks, MNIST and CelebA, demonstrate that **Mustangs** provides a statistically faster training method resulting in more accurate networks.

CCS CONCEPTS

• **Computing methodologies** → **Unsupervised learning**; **Neural networks**; *Distributed algorithms*;

KEYWORDS

Generative adversarial networks, coevolution, diversity

ACM Reference Format:

Jamal Toutouh, Erik Hemberg, and Una-May O'Reilly. 2019. Spatial Evolutionary Generative Adversarial Networks. In *Genetic and Evolutionary Computation Conference (GECCO '19)*, July 13–17, 2019, Prague, Czech Republic. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3321707.3321860>

1 INTRODUCTION

Generative adversarial networks (GANs) have emerged as a powerful machine learning paradigm. They were first introduced for the task of estimating a distribution function underlying a given set of samples [6]. A GAN consists of two neural networks, one a generator and the other a discriminator. The discriminator is trained to correctly discern the “natural/real” samples from “artificial/fake” samples produced by the generator. The generator, given a latent random space, is trained to transform its inputs into samples that

fool the discriminator. Formulated as a minmax optimization problem through the definitions of discriminator and generator loss, training can converge on an optimal generator, one that approximates the latent true distribution so well that the discriminator can only provide a label at random for any sample.

The early successes of GANs in generating realistic, complex, multivariate distributions motivated a growing body of applications, such as image generation [5], video prediction [11], image in-painting [25], and text to image synthesis [19]. However, while the competitive juxtaposition of the generator and discriminator is a compelling design, GANs are notoriously hard to train. Frequently training dynamics show pathologies. Since the generator and the discriminator are differentiable networks, optimizing the minmax GAN objective is generally performed by (variants of) simultaneous gradient-based updates to their parameters [6]. This type of gradient-based training rarely converges to an equilibrium. GAN training thus exhibits degenerate behaviors, such as *mode collapse* [4], *discriminator collapse* [10], and *vanishing gradients* [2].

Different objectives impact the gradient information used to update parameters weights of the networks. Therefore, changing the objective impacts the search trajectory and could eliminate or decrease the frequency of pathological trajectories. A set of recent studies by members of the machine learning community proposed different objective functions. Generally, these functions compute loss as the distance between the fake data and real data distributions according to different measures. The original GAN [6] applies the *Jensen-Shannon divergence (JSD)*. Other measures include: 1) *Kullback-Leibler divergence (KLD)* [15], 2) the *Wasserstein distance* [3], 3) the *least-squares (LS)* [12], and 4) the *absolute deviation* [26].

Each of these objective functions improves training but none entirely solves all of its challenges. An evolutionary computation project investigated an evolutionary generative adversarial network (**E-GAN**), a different approach [22]. **E-GAN**, batch after batch, is able to guide its trajectory with gradient information from a population of three different objectives, which defines the gradient-based mutation to be applied. As we will describe in more detail in Section 2, each batch, **E-GAN** trains each of three copies of the GAN with one of the three objectives in the population. After this independent training, **E-GAN** selects the best GAN according to a given fitness function to start the next batch and train further. This process splices batch-length trajectories from different gradient information together. Essentially, **E-GAN** injects *mutation diversity* into training. As a result, on some benchmarks, **E-GAN** improves and provides comparable performance on a baseline using a single objective.

Another evolutionary computation idea for addressing training pathologies comes from competitive coevolutionary algorithms.

With one population adversarially posed against another, they optimize with a minmax objective like GANs. Pathologies similar to what is reported in GAN training have been observed in co-evolutionary algorithms, such as *focusing*, *relativism*, and *lost of gradient* [17]. They have been attributed to a lack of diversity. Each population converges or the coupled population dynamics lock into a tit-for-tat pattern of ineffective signaling. Spatial distributed populations have been demonstrated to be particularly effective at resolving this. This approach has been transferred to GANs with **Lipizzaner** [1, 20]. **Lipizzaner** uses a spatial distributed competitive coevolutionary algorithm. It places the individuals of the generator and discriminator populations on a grid (each cell contains a pair of generator-discriminators). Each generator is evaluated against all the discriminators of its neighborhood and the same happen with each discriminator. **Lipizzaner** takes advantage of neighborhood communication to propagate models. **Lipizzaner**, in effect, provides diversity in the *genome space*.

In this paper we ask whether a method that capitalizes on ideas from both **E-GAN** and **Lipizzaner** is better than either one of them. Specifically, can a combination of diversity in mutation and genome space train GANs faster, more accurately and more reliably? Thus, we present the *Mutation SpaTial gANs* training method, **Mustangs**. For each cell of the grid, **Mustangs** selects randomly with equal probability a given loss function from among the set of three that **E-GAN** introduced, which is applied for the current training batch. This process is repeated for each batch during the whole GAN training. We experimentally evaluate **Mustangs** on standard benchmarks, MNIST and CelebA, to determine whether it provides more accurate results and/or requires shorter execution times. The main contributions of this paper are: 1) **Mustangs**, a training method of GANs that provides both mutation and genome diversity. 2) A open source software implementation of **Mustangs**¹, 3) A demonstration of **Mustangs**'s higher accuracy and performance on MNIST and CelebA. 4) A deployment of **Mustangs** on cloud computing infrastructure that optimizes the GAN grid in parallel.

The rest of the paper is organized as follows. Section 2 presents related work. Section 3 describes the method. The experimental setup and the results are in sections 4 and 5. Finally, conclusions are drawn and future work is outlined in Section 6.

2 RELATED WORK

Recent work have focused on improving the robustness of GAN training and the overall quality of generated samples [2, 4, 15]. Prior approaches tried to mitigate degenerate GAN dynamics by using heuristics, such as decreasing the optimizer's learning rate along the iterations [18]. Other authors have proposed changing generator's or discriminator's objectives [3, 12, 15, 26]. More advanced methods apply ensemble approaches [23].

A different category of studies employ multiple generators and/or discriminators. Some remarkable examples analyze training a cascade of GANs [23]; sequentially training and adding new generators with boosting techniques [21]; training in parallel multiple generators and discriminators [8]; and training an array of discriminators specialized in a different low-dimensional projection of the data [14].

Recent work by Yao and co-authors proposed **E-GAN**, whose main idea is to evolve a population of three independent loss functions defined according to three distance metrics (JSD, LS, and a metric based on JSD and KL) [22]. One at a time, independently, the loss functions are used to train a generator from some starting condition, over a batch. The generators produced by the loss functions are evaluated by a single discriminator (considered optimal) that returns a fitness value for each generator. The best generator of the three options is then selected and training continues, with the next training batch, and the three different loss functions. The use of different objective functions (mutations) overcomes the limitations of a single individual training objective and a better adapts the population to the evolution of the discriminator. **E-GAN** defines a specific fitness function that evaluates the generators in terms of the quality and the diversity of the generated samples. The results shown that **E-GAN** is able to get higher inception scores, while showing comparable stability when it goes to convergence. **E-GAN** works because the evolutionary population injects diversity into the training. Over one training run, different loss functions inform the best generator of a batch.

Another evolutionary way to improve training is also motivated by diversity [1, 20]. **Lipizzaner** simultaneously trains a spatially distributed population of GANs (pairs of generators and discriminators) that allows neighbors to communicate and share information. Gradient learning is used for GAN training and evolutionary selection and variation is used for hyperparameter learning. Overlapping neighborhoods and local communication allow efficient propagation of improving models. Besides, this strategy has the ability to distribute the training process on parallel computation architectures, and therefore, it can efficiently scale. **Lipizzaner** co-evolutionary dynamics are able to escape degenerate GAN training behaviors, e.g. mode collapse and vanishing gradient, and resulting generators provide accurate and diverse samples.

In this paper we ask whether an advance that capitalizes on ideas from both **E-GAN** (i.e., diversity in the mutation space) and **Lipizzaner** (i.e., diversity in genomes space) is better than either one of them.

3 MUSTANGS METHOD

This section presents **Mustangs** devised in this work. First, we introduce the general optimization problem of GAN training. Then, we describe a method for spatial coevolution GANs training. Finally, we present the multiple mutations applied to produce the generators offspring.

3.1 General GAN Training

In this paper we adopt a mix of notation used in [10]. Let $\mathcal{G} = \{G_u, u \in \mathcal{U}\}$ and $\mathcal{D} = \{D_v, v \in \mathcal{V}\}$ denote the class of generators and discriminators, respectively, where G_u and D_v are functions parameterized by u and v . $\mathcal{U}, \mathcal{V} \subseteq \mathbb{R}^p$ represent the parameters space of the generators and discriminators, respectively. Further, let G_* be the target unknown distribution that we would like to fit our generative model to.

¹**Mustangs** source code - <https://github.com/mustang-gan/mustang>

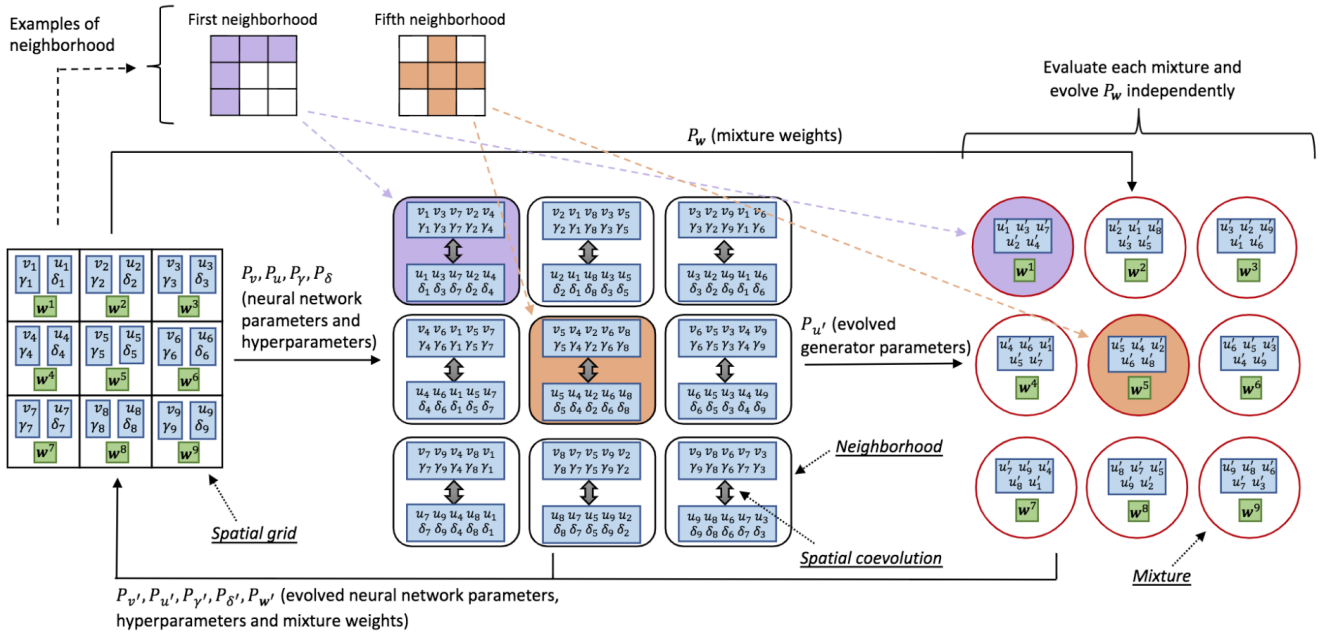


Figure 1: Spatial coevolution training on a 3×3 grid [1]. $P_v = \{v_1, \dots, v_9\}$ and $P_u = \{u_1, \dots, u_9\}$ denote neural network weights of discriminator and generator population respectively. $P_\gamma = \{\gamma_1, \dots, \gamma_9\}$ and $P_\delta = \{\delta_1, \dots, \delta_9\}$ denote the hyperparameters (e.g., learning rate) of discriminator and generator population, respectively. $P_w = \{w_1, \dots, w_9\}$ denote the mixture weights. The $(\cdot)'$ notation denotes the value of (\cdot) after one iteration of coevolution.

Formally, the goal of GAN training is to find parameters u and v so as to optimize the objective function

$$\min_{hu \in \mathcal{U}} \max_{v \in \mathcal{V}} \mathcal{L}(u, v), \text{ where}$$

$$\mathcal{L}(u, v) = \mathbb{E}_{x \sim G_*} [\phi(D_v(x))] + \mathbb{E}_{x \sim G_u} [\phi(1 - D_v(x))], \quad (1)$$

and $\phi : [0, 1] \rightarrow \mathbb{R}$, is a concave function, commonly referred to as the *measuring function*. In practice, we have access to a finite number of training samples $x_1, \dots, x_m \sim G_*$. Therefore, an empirical version $\frac{1}{m} \sum_{i=1}^m \phi(D_v(x_i))$ is used to estimate $\mathbb{E}_{x \sim G_*} [\phi(D_v(x))]$. The same holds for G_u .

3.2 Spatial Coevolution for GAN Training

Evolutionary computation implements mechanisms inspired by biological evolution such as reproduction, diversity generation, and survival of the fittest to address optimization problems. In this case, we apply the competitive coevolutionary algorithm outlined in Algorithm 1 to optimize GANs. It evolves two populations, $P_u = \{u_1, \dots, u_T\}$ a population of generators and $P_v = \{v_1, \dots, v_T\}$ a population of discriminators to create diversity in genomes spaces, where T is the population size. The fitness \mathcal{L} of each generator $u_i \in P_u$ and discriminator $v_j \in P_v$ are assessed according to their interactions with a set of discriminators from P_v and generators from P_u , respectively (Lines 2 to 7). The fittest individuals are used to generate the new of individuals (generators and discriminators) by applying mutation (see Section 3.3). The new individuals replace the ones in the current population if they perform better (better fitness) to produce the next generation.

Mustangs applies the spatially distributed coevolution summarized in Algorithm 1, the individuals of both populations (generators of P_u and discriminators of P_v) are distributed on the cells of a two dimensional toroidal grid [1]. Spatial coevolution has shown a considerable ability in maintaining diversity in the populations and fostering continuous arms races [13, 24].

The cell's *neighborhood* defines the subset of individuals of P_u and P_v to interact with and it is specified by its size s_n . Given a $m \times n$ -grid, there are $m \times n$ neighborhoods. Without losing generality, we consider m^2 square grids to simplify the notation. In our study, we use a five-cell neighborhood, i.e., one center and four adjacent cells (see Figure 1). We apply the same notation used in [1]. For the k -th neighborhood in the grid, we refer to the generator in its center cell by $P_u^{k,1} \subset P_u$ and the set of generators in the rest of the neighborhood cells by $P_u^{k,2}, \dots, P_u^{k,s_n}$, respectively. Furthermore, we denote the union of these sets by $P_u^k = \cup_{i=1}^{s_n} P_u^{k,i} \subseteq P_u$, which represents the k th generator neighborhood.

In the spatial coevolution applied here, each neighborhood performs an instance of Algorithm 1 with the populations P_u^k and P_v^k to update its center cell, i.e. $P_u^{k,1}, P_v^{k,1}$, with the returned values (Line 15 of Algorithm 1).

Given the m^2 neighborhoods, all the individuals of P_u and P_v will get updates as $P_u = \cup_{k=1}^{m^2} P_u^k, P_v = \cup_{k=1}^{m^2} P_v^k$. These m^2 instances of Algorithm 1 run in parallel in an asynchronous fashion when dealing with reading/writing from/to the populations. This implementation scales with lower communication overhead, allows the cells to run its instances without waiting each other, increases the diversity by mixing individuals computed during different stages

Algorithm 1 GeneralCoevGANs($P_u, P_v, \mathcal{L}, \{\alpha_i\}, \{\beta_i\}, p_\mu$)

Input:
 P_u : generator population P_v : discriminator population
 $\{\alpha_i\}$: selection probability $\{\beta_i\}$: mutation probability
 I : number of generations \mathcal{L} : GAN objective function
 p_μ : objective-selection probabilities vector

Return:
 P_u : evolved generator population
 P_v : evolved discriminator population

```

1: for i in [1 . . . I] do
    // Evaluate  $P_u$  and  $P_v$ 
2:  $f_{u_1 \dots u_T} \leftarrow 0$ 
3:  $f_{v_1 \dots v_T} \leftarrow 0$ 
4: for each  $u_i$  in  $P_u$ , each  $v_j$  in  $P_v$  do
5:      $f_{u_i} -= \mathcal{L}(u_i, v_j)$ 
6:      $f_{v_j} += \mathcal{L}(u_i, v_j)$ 
7: end for
    // Sort  $P_u$  and  $P_v$ 
8:  $u_{1 \dots T} \leftarrow u_{s(1) \dots s(T)}$  with  $s(i) = \text{argsort}(f_{u_1 \dots u_T}, i)$ 
9:  $v_{1 \dots T} \leftarrow v_{s(1) \dots s(T)}$  with  $s(j) = \text{argsort}(f_{v_1 \dots v_T}, j)$ 
    // Selection
10:  $u_{1 \dots T} \leftarrow u_{s(1) \dots s(T)}$  with  $s(i) = \text{argselect}(u_{1 \dots T}, i, \{\alpha_i\})$ 
11:  $v_{1 \dots T} \leftarrow v_{s(1) \dots s(T)}$  with  $s(j) = \text{argselect}(v_{1 \dots T}, j, \{\alpha_j\})$ 
    // Mutation & Replacement
12:  $u_{1 \dots T} \leftarrow \text{replace}(\{u_i\}, \{u'_i\})$  with  $u'_i = \text{mutateP}_u(u_i, \{\beta_i\})$ 
13:  $v_{1 \dots T} \leftarrow \text{replace}(\{v_j\}, \{v'_j\})$  with  $v'_j = \text{mutateP}_v(v_j, \{\beta_j\})$ 
14: end for
15: return  $P_u, P_v$ 
    
```

of the training process, and performs better with limited number of function evaluations [16].

Taking advantage of the population of $|P_u|$ generators trained, the spatial coevolution method selects one of the generator neighborhoods $\{P_u^k\}_{1 \leq k \leq m^2}$ as a mixture of generators according to a given performance metric $f : \mathcal{U}^{s_n} \times \mathbb{R}^{s_n} \rightarrow \mathbb{R}$. Thus, it is chosen the best generator mixture $P_u^* \in \mathcal{U}^{s_n}$ according to the mixture weights $\mathbf{w}^* \in [0, 1]^{s_n}$. Hence, the s_n -dimensional mixture weight vector \mathbf{w} is defined as follows

$$P_u^*, \mathbf{w}^* = \arg \max_{P_u^k, \mathbf{w}^k : 1 \leq k \leq m^2} f\left(\sum_{u_i \in P_u^k} w_i G_{u_i}\right), \quad (2)$$

where w_i represents the mixture weight of (or the probability that a data point comes from) the i th generator in the neighborhood, with $\sum_{w_i \in \mathbf{w}^k} w_i = 1$. These hyperparameters $\{\mathbf{w}^k\}_{1 \leq k \leq m^2}$ are optimized during the training process after each step of spatial coevolution by applying an evolution strategy (1+1)-ES [1].

3.3 Mustangs Gradient-based Mutation

Mustangs coevolutionary algorithm generates the offspring of both populations P_h and P_q by by applying *asexual reproduction*, i.e. next generation's of individuals are produced by applying mutation (Lines 12 and 13 in Algorithm 1). These mutation operators are defined according to a giving training objective, which generally attempts to minimize the distance between the generated fake data and real data distributions according to a given measure. **Lipizzaner** applies the same gradient-based mutation for both populations during the coevolutionary learning [1].

In this study, we add mutation diversity to the genome diversity provided by **Lipizzaner**. Thus, we use the mutations used by **E-GAN** to generate the offspring of generators [22]. E-GAN applies three different mutations corresponding with three different minimization objectives w.r.t. the generator: 1) *Minmax mutation*, which

objective is to minimize the JSD between the real and fake data distributions, i.e., $JSD(p_{real} \parallel p_{fake})$ (see Equation (3)). 2) *Least-square mutation*, which is inspired in the least-square GAN [12] that applies this criterion to adapt both, the generator and the discriminator. The objective function is formulated as shown in Equation (4). 3) *Heuristic mutation*, which maximizes the probability of the discriminator being mistaken by minimizing the objective function in Equation (5). This objective is equal to minimizing $[KL(p_{real} \parallel p_{fake}) - 2JSD(p_{real} \parallel p_{fake})]$

$$M_G^{minmax} = \frac{1}{2} \mathbb{E}_{x \sim G_u} [\log(1 - D_v(x))] \quad (3)$$

$$M_G^{least-square} = \mathbb{E}_{x \sim G_u} [\log(D_v(x) - 1)^2] \quad (4)$$

$$M_G^{heuristic} = \frac{1}{2} \mathbb{E}_{x \sim G_u} [\log(D_v(x))] \quad (5)$$

Thus, the mutation applied to the generators (mutateP_u) in Line 12 of Algorithm 1 is defined in the Algorithm 2. The new generator is produced by using a loss function (mutation) to optimize one of the three objectives functions introduced above. i.e., M_G^{minmax} , $M_G^{least-square}$, and $M_G^{heuristic}$, which are *binary cross entropy (BCE) loss*, *mean square error (MSE) loss*, and *heuristic loss*, respectively.

E-GAN applies the three mutations to the generator (ancestor) and it selects the individual that provides the best fitness when is evaluated against the discriminator [22]. In contrast, **Mustangs** picks at random with same the probability ($\frac{1}{3}$) one of the mutations (loss functions), as it is shown in Figure 2, and then the gradient descent method is applied accordingly (see Algorithm 2). This enables diversity in the mutation space without adding noticeable overhead over the spatial coevolutionary training method presented before, since **Mustangs** evaluates only the mutated generator instead of three as **E-GAN** does.

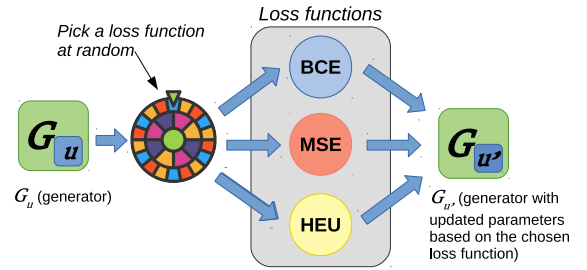


Figure 2: Graphical representation of the mutation used in Mustangs. The generator (ancestor) G_u is mutated to produce the new generator $G_{u'}$, by using one of the loss functions chosen at random.

E-GAN, **Lipizzaner**, and **Mustangs** apply the same mutation (loss function) to update the discriminators, the one defined to address the GAN minmax optimization problem described in Equation (1).

Algorithm 2 mutate $P_u(u)$ **Input:** u : generator parameters**Return:** u' : mutated generator parameters

```

1: operator ← pick_random[1...3]
2: if operator == 1 then
3:   u' ← applyGradientDescent( $M_G^{minmax}$ , u)
4: else if operator == 2 then
5:   u' ← applyGradientDescent( $M_G^{least-square}$ , u)
6: else
7:   u' ← applyGradientDescent( $M_G^{heuristic}$ , u)
8: end if
9: return u'
```

4 EXPERIMENTAL SETUP

Mustangs is evaluated on two common image data sets: MNIST² and CelebA³. The experiments compare the following methods:

GAN-BCE a standard GAN which uses M_G^{minmax} objective

E-GAN the E-GAN method

Lip-BCE **Lipizzaner** with M_G^{minmax} objective

Lip-MSE **Lipizzaner** with $M_G^{least-square}$ objective

Lip-HEU **Lipizzaner** with $M_G^{heuristic}$ objective

Mustangs the **Mustangs** method

The settings used for the experiments are summarized in Table 1. The four spatial covevolutionary GANs use the parameters presented in Table 1. **E-GAN** and **GAN-BCE** both use the Adam optimizer with an initial learning rate (0.0002). The other parameters of **E-GAN** use the same configuration as used in [22].

For MNIST data set experiments, all methods use the same stop condition: a computational budget of nine hours (9h). The distributed methods of **Mustangs**, **Lip-BCE**, **Lip-MSE**, and **Lip-HEU** use a grid size of 3×3 , and are able to train nine networks in parallel. Thus, they are executed during one hour to comply with the computational budget of nine hours. Regarding CelebA experiments, the four spatial covevolutionary GANs are analyzed. Thus, they stop after performing 20 training epochs, since they require similar computational budget to perform them.

All methods have been implemented in Python3 and pytorch⁴. The spatial covevolutionary ones have extended the **Lipizzaner** framework [20].

The experimental analysis is performed on a cloud computation platform that provides 8 Intel Xeon cores 2.2GHz with 32 GB RAM and a NVIDIA Tesla T4 GPU with 16 GB RAM. We run multiple independent runs for each method.

For quantitative assessment of the accuracy of the generated fake data the Frechet inception distance (FID) is evaluated [7]. We analyze the computational performance of each method. Finally, we evaluate the diversity of the data samples generated.

5 RESULTS AND DISCUSSION

This section presents the results and the analyses of the studied optimization methods. The first three subsections evaluate the MNIST

Table 1: Network topology of the GANs trained.

Parameter	MNIST	CelebA
<i>Network topology</i>		
Network type	MLP	DCGAN
Input neurons	64	100
Number of hidden layers	2	4
Neurons per hidden layer	256	16,384 - 131,072
Output neurons	784	64×64×64
Activation function	<i>tanh</i>	<i>tanh</i>
<i>Training settings</i>		
Batch size	100	128
Skip N disc. steps	1	-
<i>Coevolutionary settings</i>		
Stop condition	9 hours comp.	20 training epochs
Population size per cell	1	1
Tournament size	2	2
Grid size	3×3	2×2
Performance metric (m)	FID	FID
Mixture mutation scale	0.01	0.01
<i>Hyperparameter mutation</i>		
Optimizer	Adam	Adam
Initial learning rate	0.0002	0.00005
Mutation rate	0.0001	0.0001
Mutation probability	0.5	0.5

experiments in terms of the FID score, the computational performance, and the diversity of the generated samples, respectively. The last one analyzes the CelebA results in terms of FID score.

5.1 Quality of the Generated Data

Table 2 shows the best FID value from each of the 30 independent runs performed for each method. **Mustangs** has the lowest median (see Figure 3 for a boxplot). All the methods that used **Lipizzaner** are better than **E-GAN** and **GAN-BCE**. However, there is quite a significant increase in FID value between the **Lip-MSE** and the other **Lipizzaner** based methods. The results indicate that **Mustangs** is robust to the varying performance of the individual loss functions and can still find a high performing mixture of generators. This helps to strengthen the idea that diversity, both in genome and mutation space, provides robust GAN training.

Table 2: FID MNIST results in terms of best mean, normalized standard deviation, median and interquartile range (IQR). (Low FID indicates good performance)

Algorithm	Mean	Std%	Median	IQR
Mustangs	42.235	12.863%	43.181	7.586
Lip-BCE	48.958	20.080%	46.068	4.663
Lip-MSE	371.603	20.108%	381.768	104.625
Lip-HEU	52.525	17.230%	52.732	9.767
E-GAN	466.111	10.312%	481.610	69.329
GAN-BCE	457.723	2.648%	459.629	17.865

The results provided by the methods that generates diversity in genome space only (**Lip-BCE**, **Lip-HEU**, and **Lip-MSE**) are significantly more competitive than **E-GAN**, which provides diversity in mutation space only (see Figure 3). Therefore, the spatial distributed coevolution provides an efficient tool to optimize GANs.

²The MNIST Database - <http://yann.lecun.com/exdb/mnist/>

³The CelebA Database - <http://mmlab.ie.cuhk.edu.hk/projects/CelebA.html>

⁴Pytorch Website - <https://pytorch.es/>

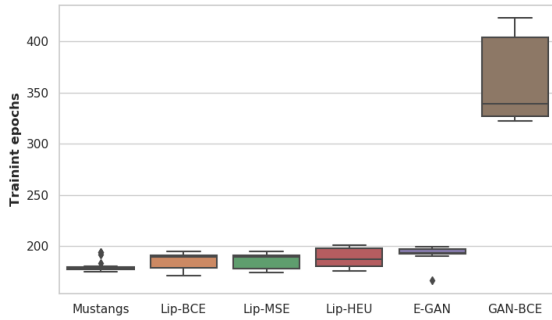


Figure 3: Results on MNIST dataset. Boxplot that shows the best FIDs computed for each independent run.

Surprisingly the median FID score of **GAN-BCE** is better than **E-GAN**. **E-GAN** has a larger variance of FID scores compared to **GAN-BCE**, and in the original paper it was shown that **E-GAN** performance improved with more epochs (by using a computational budget of 30h, compared to the 9h we use here).

The results in Table 2 indicate that spatially distributed coevolutionary training is the best choice to train GANs, even when there is no knowledge about the best loss function to the problem. However, the choice of loss function (mutation) may impact the final results. In summary, the combination of both mutation and genome diversity significantly provides the most best result. A ranksum test with Holm correction confirms that the difference between **Mustangs** and the other methods is significant at confidence levels of $\alpha=0.01$ and $\alpha=0.001$ (see Figure 4).

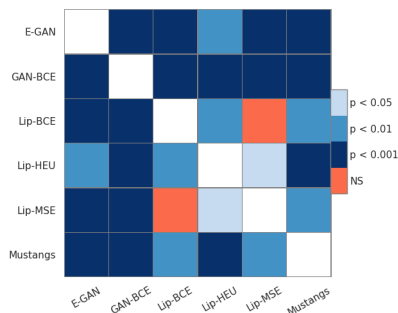


Figure 4: Holm statistical post-hoc analysis on MNIST dataset. It illustrates the p -values computed by the statistical tests.

Next, we evaluate the FID score through out the GAN training process, see Figure 5 illustrates the FID changes during the entire training process. In addition, we zoom in on the first 50 training epochs in Figure 6. None of the evolutionary GAN training methods seem to have converged after 9h of computation. This implies that longer runs can further improve the results.

According to Figure 5, the robustness of the three most competitive methods (**Mustangs**, **Lip-BCE**, and **Lip-HEU**) indicates that the FID almost behaves like a monotonically decreasing function

with small oscillations. The other three methods have larger oscillations and does not seem to have a FID trend that decreases with the same rate.

Focusing on the two methods that apply the same unique loss function in Equation (3), **GAN-BCE** and **Lip-BCE**, we can clearly state the benefits of the distributed spatial evolution. Even the two methods provide comparable FID during the first 30 training epochs, **Lip-BCE** converges faster to better FID values (see Figure 6). This difference is even more noticeable when both algorithms consume the 9h of computational cost (see Figure 5).

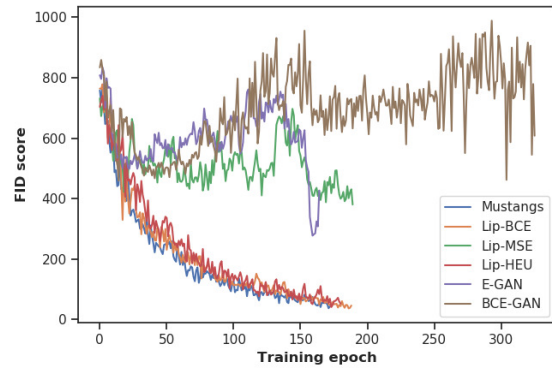


Figure 5: Results on MNIST dataset. FID evolution through the training process during 9h of computational cost.

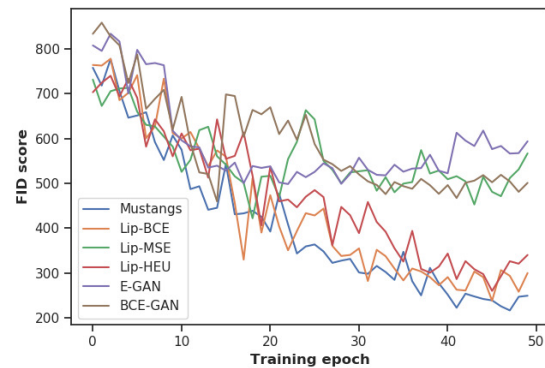


Figure 6: Results on MNIST dataset. FID evolution through the first 50 epochs of the training process.

Notice that the spatial coevolutionary methods use FID as the objective function to select the best mixture of generators during the optimization of the GANs. In contrast, **E-GAN** applies a specific objective function based on the losses [22] and **GAN-BCE** optimizes just one network.

5.2 Computational Performance

In this section, we analyze the computational performance of the GAN training methods, all used the same computational budget (9h). We start analyzing the number of training epochs. As **Mustangs** and **Lipizzaner** variations apply asynchronous parallelism, the

number of training epochs performed by each cell of the grid in the same run varies. Thus, for these methods, we consider that the number of training epochs of a given run is the mean of the epochs performed by each cell.

Table 3 shows the mean, normalized standard deviation, minimum, and maximum values of number of training epochs. The number of epochs are normally distributed for all the algorithms. Please note that, first, all the analyzed methods have been executed on a cloud architecture, which could generate some differences in terms of computational efficiency results; and second, everything is implemented by using the same Python libraries and versions to limit computational differences between them due to the technologies used.

Table 3: Results on MNIST dataset. Mean, normalized standard deviation, minimum, median and maximum of training epochs.

Algorithm	Mean	Std%	Minimum	Median	Maximum
Mustangs	179.456	2.708%	174.778	177.944	194.000
Lip-BCE	185.919	4.072%	171.222	189.444	194.333
Lip-MSE	185.963	3.768%	173.667	189.222	194.778
Lip-HEU	188.644	4.657%	175.667	186.667	200.556
E-GAN	193.167	2.992%	166.000	193.000	199.000
GAN-BCE	365.267	10.570%	322.000	339.500	423.000

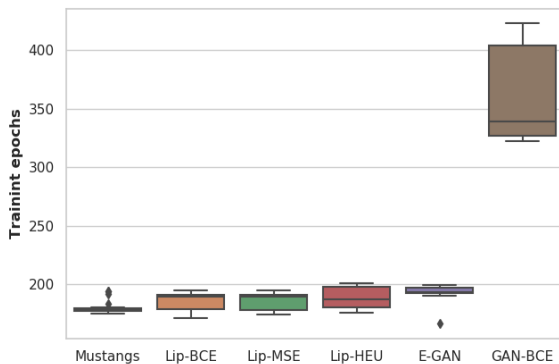


Figure 7: Results on MNIST dataset. Boxplot of the number of training epochs.

The method that is able to train each network for the most epochs is **GAN-BCE**. However, this method trains only one network, in contrast to the other evaluated methods. **GAN-BCE** performs about two times the number of iterations of **E-GAN**, which evaluates three networks.

The spatially distributed coevolutionary algorithms performed significantly fewer training epochs than **E-GAN**. However, during each epoch these methods evaluate 45 GANs, i.e., neighborhood size of 5×9 cells, which is 15 times more networks than **E-GAN**.

One of the most important features of the spatial coevolutionary algorithm is that it is executed asynchronously and in parallel for all the cells [20]. Thus, there is no bottleneck for each cells performance since it operates without waiting for the others. In future work

E-GAN could take advantage of parallelism and optimizing the three discriminators at the same time. However, it has an important synchronization bottleneck because they are evaluated over the same discriminator, which is trained and evaluated sequentially after that operation.

5.3 Diversity of the Generated Outputs

In this section, we evaluate the diversity of the generated samples by the networks that had the best FID score. We report the the total variation distance (TVD) for each algorithm [9] (see Table 4).

Table 4: Results on MNIST dataset. TVD results. (Low TVD indicates more diversity)

Alg.	Mustangs	Lip-BCE	Lip-HEU	Lip-MSE	E-GAN	GAN-BCE
TVD	0.180	0.171	0.115	0.365	0.534	0.516

The methods that provide genome diversity generate more diverse data samples than the other two analyzed methods. This shows that genomic diversity introduces a capability to avoid mode collapse situations as the one shown in Figure 9(a). The three algorithms with the lowest FID score (**Mustangs**, **Lip-BCE**, and **Lip-HEU**) also provide the lowest TVD values. The best TVD result is obtained by **Lip-HEU**.

The distribution of each class of digits for generated images is shown in Figure 8. The diagram supports the TVD results, e.g. **Lip-HEU** and **Mustangs** produce more diverse set of samples spanning across different classes. We can observe that the two methods that do not apply diversity in terms of genome display a possible mode collapse, since about half of the samples are of class 3 and they not generate samples of class 4 and 7.

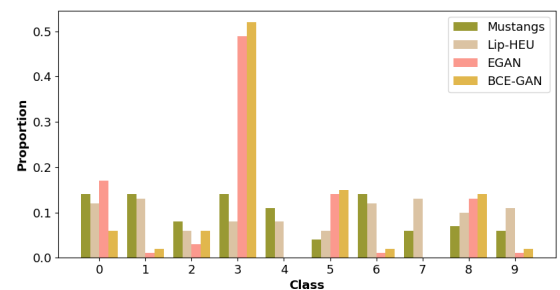


Figure 8: Classes distribution of the samples generated of MNIST by Mustangs, Lip-HEU, E-GAN, and GAN-BCE.

Figure 9 illustrates how spatially distributed coevolutionary algorithms are able to produce robust generators that provide with accurate samples across all the classes.

5.4 CelebA Experimental Results

The spatially distributed coevolutionary methods are applied to perform the CelebA experiments. Table 5 summarizes the results over multiple independent runs.

Mustangs provides the lowest median FID and **Lip-MSE** the highest one. **Lip-BCE** and **Lip-HEU** provide median and average

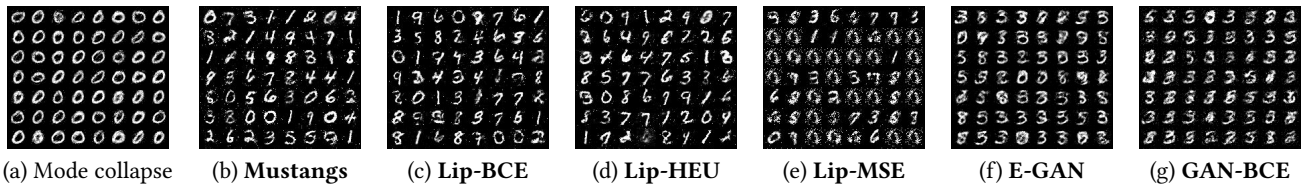


Figure 9: Sequence of samples generated of MNIST dataset. (a) contains for mode collapse, the generator is focused on the character 0. It illustrates samples generated by the best generator (in terms of FID) by Mustangs (b), Lip-BCE (c), Lip-HEU (d), Lip-MSE (e), E-GAN (f), and GAN-BCE (g).

Table 5: Results on CelebA dataset. FID results in terms of best mean, normalized standard deviation, median and interquartile range (IQR). (Low FID indicates good performance)

Algorithm	Mean	Std%	Median	IQR
Mustangs	36.148	0.571%	36.087	0.237
Lip-BCE	36.250	5.676%	36.385	2.858
Lip-MSE	158.740	40.226%	160.642	47.216
Lip-HEU	37.872	5.751%	37.681	2.455

FID scores close to the Mustangs ones. However, Mustangs is more robust to the varying performance of the methods that apply a unique loss function (see deviations in Table 5).

The robustness of the training provided by Mustangs makes it an efficient tool when the computation budget is limited (i.e., performing a limited number of independent runs), since it shows low variability in its competitive results.

Next, we evaluate the FID score through the training process. Figure 10 shows the FID changes during the entire training process. All the analyzed methods behave like a monotonically decreasing function. However, the FID evolution of Lip-BCE presents oscillations. Mustangs, Lip-BCE, and Lip-HEU FIDs show a similar evolution.

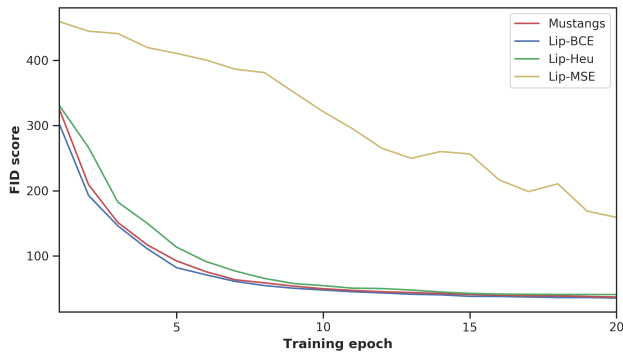


Figure 10: Results on CelebA dataset. FID evolution through the 20 epochs of the training process.

Figure 11 illustrates a sequence of samples generated by the best mixture of generators in terms of FID score of the most competitive two training methods, i.e., Mustangs and Lip-BCE. As it can be seen in these two sets images generated, both methods present similar capacity of generating human faces.



(a) Mustangs (b) Lip-BCE

Figure 11: Sequence of samples generated of CelebA. It illustrates samples generated by the best mixture of generators (in terms of FID) by Mustangs (a) and Lip-BCE (b).

6 CONCLUSIONS AND FUTURE WORK

We have empirically showed that GAN training can be improved by boosting diversity. We enhanced an existing spatial evolutionary GAN training framework that promoted genomic diversity by probabilistically choosing one of three loss functions. This new method, called Mustangs was tested on the MNIST and CelebA data sets showed the best performance and high diversity in label space, as well as on the TVD measure. The high performance of Mustangs is due to its inherent robustness. This allows it to overcome often observed training pathologies, e.g. mode collapse. Furthermore, the Mustangs method can be executed asynchronously and the computation is easy to distribute with low overhead. We extended the Lipizzaner open source implementation to demonstrate this.

Future work will include the evaluation of Mustangs on more data sets and longer training epochs. We can also include other loss functions. In addition, we are exploring the diversity of the networks over their neighborhoods and the whole grid when applying the different methods studied here. This study will be a first step in devising an algorithm that self-adapts the probabilities of the different mutations dynamically. Finally, other advancements in evolutionary algorithms that can improve the robustness of GAN training, e.g. temporal evolutionary training, need to be considered.

ACKNOWLEDGMENTS

This research was partially funded by European Union’s Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 799078, and by MINECO and FEDER projects TIN2016-81766-REDT and TIN2017-88213-R. The Systems that learn initiative at MIT CSAIL.

REFERENCES

- [1] Abdullah Al-Dujaili, Tom Schmedlechner, Erik Hemberg, and Una-May O'Reilly. 2018. Towards distributed coevolutionary GANs. In *AAAI 2018 Fall Symposium*.
- [2] Martin Arjovsky and Léon Bottou. 2017. Towards principled methods for training generative adversarial networks. *arXiv preprint arXiv:1701.04862* (2017).
- [3] Martin Arjovsky, Soumith Chintala, and Léon Bottou. 2017. Wasserstein GAN. *arXiv preprint arXiv:1701.07875* (2017).
- [4] Sanjeev Arora, Andrej Risteski, and Yi Zhang. 2018. Do GANs learn the distribution? Some Theory and Empirics. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=BJehNfW0->
- [5] Zhe Gan, Liqun Chen, Weiyao Wang, Yuchen Pu, Yizhe Zhang, Hao Liu, Chunyuan Li, and Lawrence Carin. 2017. Triangle generative adversarial networks. In *Advances in Neural Information Processing Systems*. 5247–5256.
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative adversarial nets. In *Advances in neural information processing systems*. 2672–2680.
- [7] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Günter Klambauer, and Sepp Hochreiter. 2017. GANs trained by a two time-scale update rule converge to a nash equilibrium. *arXiv preprint arXiv:1706.08500* 12, 1 (2017).
- [8] Daniel Jiwoong Im, He Ma, Chris Dongjoo Kim, and Graham Taylor. 2016. Generative Adversarial Parallelization. *arXiv preprint arXiv:1612.04021* (2016).
- [9] Chengtao Li, David Alvarez-Melis, Keyulu Xu, Stefanie Jegelka, and Suvrit Sra. 2017. Distributional Adversarial Networks. *arXiv preprint arXiv:1706.09549* (2017).
- [10] Jerry Li, Aleksander Madry, John Peebles, and Ludwig Schmidt. 2017. Towards Understanding the Dynamics of Generative Adversarial Networks. *arXiv preprint arXiv:1706.09884* (2017).
- [11] Xiaodan Liang, Lisa Lee, Wei Dai, and Eric P Xing. 2017. Dual motion GAN for future-flow embedded video prediction. In *IEEE International Conference on Computer Vision (ICCV)*, Vol. 1.
- [12] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. 2017. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 2794–2802.
- [13] Melanie Mitchell, Michael D. Thomure, and Nathan L. Williams. 2006. The role of space in the success of coevolutionary learning. In *Artificial Life X: Proceedings of the Tenth International Conference on the Simulation and Synthesis of Living Systems*. MIT Press, 118–124.
- [14] Behnam Neyshabur, Srinadh Bhojanapalli, and Ayan Chakrabarti. 2017. Stabilizing GAN training with multiple random projections. *arXiv preprint arXiv:1705.07831* (2017).
- [15] Tu Nguyen, Trung Le, Hung Vu, and Dinh Phung. 2017. Dual discriminator generative adversarial nets. In *Advances in Neural Information Processing Systems*. 2670–2680.
- [16] Sune S Nielsen, Bernabé Dorronsoro, Grégoire Danoy, and Pascal Bouvry. 2012. Novel efficient asynchronous cooperative co-evolutionary multi-objective algorithms. In *Evolutionary Computation (CEC), 2012 IEEE Congress on*. IEEE, 1–7.
- [17] Elena Popovici, Anthony Bucci, R Paul Wiegand, and Edwin D De Jong. 2012. Coevolutionary principles. In *Handbook of natural computing*. Springer, 987–1033.
- [18] Alec Radford, Luke Metz, and Soumith Chintala. 2015. Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks. *arXiv preprint arXiv:1511.06434* (2015).
- [19] Scott E Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee. 2016. Learning what and where to draw. In *Advances in Neural Information Processing Systems*. 217–225.
- [20] Tom Schmedlechner, Ignavier Ng Zhi Yong, Abdullah Al-Dujaili, Erik Hemberg, and Una-May O'Reilly. 2018. Lipizzaner: A System That Scales Robust Generative Adversarial Network Training. In *the 32nd Conference on Neural Information Processing Systems (NeurIPS 2018) Workshop on Systems for ML and Open Source Software*.
- [21] Ilya O Tolstikhin, Sylvain Gelly, Olivier Bousquet, Carl-Johann Simon-Gabriel, and Bernhard Schölkopf. 2017. Adagan: Boosting generative models. In *Advances in Neural Information Processing Systems*. 5430–5439.
- [22] Chaoyue Wang, Chang Xu, Xin Yao, and Dacheng Tao. 2018. Evolutionary Generative Adversarial Networks. *arXiv preprint arXiv:1803.00657* (2018).
- [23] Yaxing Wang, Lichao Zhang, and Joost van de Weijer. 2016. Ensembles of generative adversarial networks. *arXiv preprint arXiv:1612.00991* (2016).
- [24] Nathan Williams and Melanie Mitchell. 2005. Investigating the Success of Spatial Coevolution. In *Proceedings of the 7th Annual Conference on Genetic and Evolutionary Computation (GECCO '05)*. ACM, New York, NY, USA, 523–530. <https://doi.org/10.1145/1068009.1068096>
- [25] Raymond A Yeh, Chen Chen, Teck Yian Lim, Alexander G Schwing, Mark Hasegawa-Johnson, and Minh N Do. 2017. Semantic image inpainting with deep generative models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 5485–5493.
- [26] Junbo Zhao, Michael Mathieu, and Yann LeCun. 2016. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126* (2016).