# Learning Implicit Generative Models by Matching Perceptual Features

Cicero Nogueira dos Santos*, Youssef Mroueh*, Inkit Padhi*, Pierre Dognin*
IBM Research, T.J. Watson Research Center, NY

{cicerons,mroueh,pdognin}@us.ibm.com, inkit.padhi@ibm.com

## Abstract

*Perceptual features (PFs) have been used with great success in tasks such as transfer learning, style transfer, and super-resolution. However, the efficacy of PFs as key source of information for learning generative models is not well studied. We investigate here the use of PFs in the context of learning implicit generative models through moment matching (MM). More specifically, we propose a new effective MM approach that learns implicit generative models by performing mean and covariance matching of features extracted from pretrained ConvNets. Our proposed approach improves upon existing MM methods by: (1) breaking away from the problematic min/max game of adversarial learning; (2) avoiding online learning of kernel functions; and (3) being efficient with respect to both number of used moments and required minibatch size. Our experimental results demonstrate that, due to the expressiveness of PFs from pretrained deep ConvNets, our method achieves state-of-the-art results for challenging benchmarks.*

## 1. Introduction

The use of features from deep convolutional neural networks (DCNNs) pretrained on ImageNet [35] has led to important advances in computer vision. DCNN features, usually called *perceptual features (PFs)*, have been used in tasks such as transfer learning [40, 16], style transfer [9] and super-resolution [17]. While there have been previous works on the use of PFs in the context of image generation and transformation [7, 17], exploration of PFs as key source of information for learning generative models is not well studied. Particularly, the efficacy of PFs for implicit generative models trained through *moment matching* is an open question.

Moment matching approaches for generative modeling are based on the assumption that one can learn the data distribution by matching the moments of the model distribution to the empirical data distribution. Two representative meth-

ods of this family are based on maximum mean discrepancy (MMD) [11, 12, 22] and the method of moments (MoM) [33]. While MoM based methods embed a probability distribution into a finite-dimensional vector (i.e., matching of a finite number of moments), MMD based methods embed a distribution into an infinite-dimensional vector [33]. A challenge for MMD methods is to define a kernel function that is statistically efficient and can be used with small minibatch sizes [21]. A solution comes by using adversarial learning for the online training of kernel functions [21, 3]. However, this solution inherits the problematic min/max game of adversarial learning. The main challenges of using MoM for training deep generative networks consist in defining millions of sufficiently distinct moments and specifying an objective function to learn the desirable moments. Ravuri *et al.* [33] addressed these two issues by defining the moments as features and derivatives from a *moment network* that is trained online (together with the generator) by using a specially designed objective function.

In this work we demonstrate that, by using PFs to perform moment matching, one can overcome some of the difficulties found in current moment matching approaches. More specifically, we propose a simple but effective moment matching method that: (1) breaks away from the problematic min/max game completely; (2) does not use online learning of kernel functions; and (3) is very efficient with regard to both number of used moments and required minibatch size. Our proposed approach, named Generative Feature Matching Networks (GFMN), learns implicit generative models by performing mean and covariance matching of features extracted from all convolutional layers of pretrained deep ConvNets. Some interesting properties of GFMNs include: (a) the loss function is directly correlated to the generated image quality; (b) mode collapsing is not an issue; and (c) the same pretrained feature extractor can be used across different datasets.

We perform an extensive number of experiments with different challenging datasets: CIFAR10, STL10, CelebA and LSUN. We demonstrate that our approach can achieve state-of-the-art results for challenging benchmarks such as CIFAR10 and STL10. Moreover, we show that the same
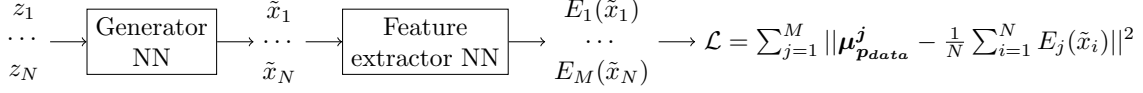
---

* Equal contribution.

Figure 1: **GFMN Training**: From $z_1, \ldots z_N$ noise signals, generator $G$ creates $N$ images $\tilde{x}_1, \ldots \tilde{x}_N$. The fixed pretrained feature extractor $E$ is used to obtain $E_j(\tilde{x}_i)$ features. $\mathcal{L}$ is the $L_2$-norm of the difference between extracted features means of generated and real data, $\boldsymbol{\mu}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}}$. We precompute $\boldsymbol{\mu}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}}$ on the entire real dataset (it does not change during training); the mean of generated data is estimated on a minibatch of size $N$. The same strategy is used for variance terms in $\mathcal{L}$.

feature extractor is effective across different datasets. The main contributions of this work can be summarized as follows: (1) We propose a new effective moment matching-based approach to train implicit generative models that does not use adversarial or online learning of kernel functions, provides stable training, and achieves state-of-the-art results; (2) We show theoretical results that demonstrate GFMN convergence under the assumption of the universality of perceptual features; (3) We propose an ADAM-based moving average method that allows effective training with small minibatches; (4) Our extensive quantitative and qualitative experimental results demonstrate that pretrained autoencoders and DCNN classifiers can be effectively used as (cross-domain) feature extractors for GFMN training.

## 2. Generative Feature Matching Networks

### 2.1. The method

Let $G$ be the generator implemented as a neural network with parameters $\theta$, and let $E$ be a pretrained neural network with $L$ hidden layers. Our proposed approach consists in training $G$ by minimizing the following loss function:

$$\min_\theta \sum_{j=1}^{M} ||\mu_{p_{data}}^j - \mu_{p_G}^j(\theta)||^2 + ||\sigma_{p_{data}}^j - \sigma_{p_G}^j(\theta)||^2 \quad (1)$$

where:

$$\mu_{p_{data}}^j = \mathbb{E}_{x \sim p_{data}} E_j(x) \in \mathbb{R}^{d_j}$$
$$\mu_{p_G}^j(\theta) = \mathbb{E}_{z \sim \mathcal{N}(0, I_{n_z})} E_j(G(z; \theta)) \in \mathbb{R}^{d_j}$$
$$\sigma_{p_{data}, \ell}^j = \mathbb{E}_{x \sim p_{data}} E_{j,\ell}(x)^2 - [\mu_{p_{data}}^{j,\ell}]^2, \ell = 1 \ldots d_j$$
$$\sigma_{p_G, \ell}^j(\theta) = \mathbb{E}_{z \sim \mathcal{N}(0, I_{n_z})} E_{j,\ell}(G(z; \theta))^2 - [\mu_{p_G}^{j,\ell}]^2, \ell = 1 \ldots d_j$$

and $||.||^2$ is the $L_2$ loss; $x$ is a real data point sampled from the data generating distribution $p_{data}$; $z \in \mathbb{R}^{n_z}$ is a noise vector sampled from the normal distribution $\mathcal{N}(0, I_{n_z})$; $E_j(x)$, denotes the output vector/feature map of the hidden layer $j$ from $E$; $M \leq L$ is the number of hidden layers used to perform feature matching. Note that $\sigma_{p_{data}}^2$ and $\sigma_{p_G}^2$ denote the variances of the features from real data and generated data, respectively. We use diagonal covariance matrices as computing full covariance matrices is impractical for large numbers of features.

In practice, we train $G$ by first precomputing estimates of $\boldsymbol{\mu}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}}$ and $\boldsymbol{\sigma}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}}$ on the training data, then running multiple training iterations where we sample a minibatch of generated (*fake*) data and optimize the parameters $\theta$ using stochastic gradient descent (SGD) with backpropagation. The network $E$ is used for the purpose of feature extraction only and is kept fixed during the training of $G$. Fig. 1 presents GFMN training pipeline.

**Autoencoder Features**: A natural choice of unsupervised method to train a feature extractor is the autoencoder (AE) framework. The decoder part of an AE consists exactly of an image generator that uses features extracted by the encoder. Therefore, by design, the encoder network should be a good feature extractor for the purpose of generation.

**Classifier Features:** We experiment with different DCNN architectures pretrained on ImageNet to play the role of the feature extractor $E$. Our hypothesis is that ImageNet-based PFs are informative enough to allow the training of (cross-domain) generators by feature matching.

### 2.2. Matching Feat. with ADAM Moving Average

**From feature matching loss to moving averages.** In order to train with a mean and covariance feature matching loss, one needs large minibatches to obtain good mean and covariance estimates. With images larger than $32 \times 32$, DCNNs produce millions of features, resulting easily in memory issues. We propose to alleviate this problem by using *moving averages of the difference of means (covariances) of real and generated data*. Instead of computing the (memory) expensive feature matching loss in Eq. 1, we keep moving averages $v_j$ of the difference of feature means (covariances) at layer $j$ between real and generated data. We detail our moving average strategy for the mean features only, but the same approach applies for the covariances. The mean features from the first term of Eq. 1, $||\boldsymbol{\mu}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}} - \mathbb{E}_{z \sim \mathcal{N}(0, I_{n_z})} E_j(G(z; \theta))||^2$ can be approximated by:

$$v_j^\top \left( \boldsymbol{\mu}_{\boldsymbol{p_{data}}}^{\boldsymbol{j}} - \frac{1}{N} \sum_{k=1}^{N} E_j(G(z_k; \theta)) \right),$$

where $N$ is the minibatch size and $v_j$ is a moving average on $\Delta_j$, the difference of the means of the features extracted

by the $j$-th layer of $E$:

$$\Delta_j = \boldsymbol{\mu_{p_{data}}^j} - \frac{1}{N}\sum_{k=1}^{N} E_j(G(z_k;\theta)). \quad (2)$$

Using these moving averages we replace the first term of the loss given in Eq. 1 by

$$\min_{\theta}\sum_{j=1}^{M} v_j^{\top}\left(\boldsymbol{\mu_{p_{data}}^j} - \frac{1}{N}\sum_{k=1}^{N} E_j(G(z_k;\theta))\right). \quad (3)$$

The moving average formulation of features matching above has a major advantage on the naive formulation of Eq. 1 since we can now rely on $v_j$ to get better estimates of the population feature means of real and generated data while using a small minibatch of size $N$. For a similar result using the feature matching loss given in Eq. 1, one would need a minibatch with large size $N$, which is problematic for large number of features.

**ADAM moving average: from SGD to ADAM updates.** Note that for a rate $\alpha$, the moving average $v_j$ has the following update:

$$v_{j,\text{new}} = (1-\alpha)*v_{j,\text{old}} + \alpha*\Delta_j, \forall j = 1\dots M$$

It is easy to see that the moving average is a gradient descent update on the following loss:

$$\min_{v_j}\frac{1}{2}||v_j - \Delta_j||^2. \quad (4)$$

Hence, writing the gradient update with learning rate $\alpha$ we have equivalently:

$$v_{j,\text{new}} = v_{j,\text{old}} - \alpha*(v_{j,\text{old}} - \Delta_j) = (1-\alpha)*v_{j,\text{old}} + \alpha*\Delta_j.$$

With this interpretation of the moving average, we propose to get a better moving average estimate by using the ADAM optimizer [18] on the loss of the moving average given in Eq. 4, such that

$$v_{j,\text{new}} = v_{j,\text{old}} - \alpha ADAM(v_{j,\text{old}} - \Delta_j).$$

$ADAM(x)$ function is computed as follows:

$$m_t = \beta_1 * m_{t-1} + (1-\beta_1)*x \quad \hat{m}_t = m_t/(1-\beta_1^t)$$
$$u_t = \beta_2 * u_{t-1} + (1-\beta_2)*x^2 \quad \hat{u}_t = u_t/(1-\beta_2^t)$$
$$ADAM(x) = \hat{m}_t/(\sqrt{\hat{u}_t}+\epsilon),$$

where $x$ is the gradient for the loss function in Eq. 4, $t$ is the iteration number, $m_t$ and $u_t$ are the first and second moment vectors at iteration $t$, $\beta_1 = .9$, $\beta_2 = .999$ and $\epsilon = 10^{-8}$ are constants. $m_0$ and $u_0$ are initialized as proposed by [18]. We refer to [18] for a detailed ADAM optimizer description.

This moving average formulation, which we call *ADAM Moving Average* (AMA) promotes stable training when using small minibatches. Although we detail AMA using

mean feature matching only, we use this approach for both mean and covariance matching. The main advantage of AMA over simple moving average (MA) is in its adaptive first and second order moments that ensure stable estimation of the moving averages $v_j$. In fact, this is a non-stationary estimation since the mean of the generated data changes in the training, and it is well known that ADAM works well for such online non-stationary losses [18].

In Section 5.3 we provide experimental results supporting: (1) The memory advantage that the AMA formulation of feature matching offers over the naive implementation; (2) The stability advantage and improved generation results that AMA allows compared to the naive implementation. We discuss in Appendix 2 the advantage of AMA on MA from a regret bounds point of view [34].

# 3. Universality of PFs and GFMN Convergence

Our proposed approach is related to the recent body of work on MMD or MM based generative models [22, 8, 21, 3, 33]. We highlight the main differences between MMD-GANs and GFMN in terms of requirements on the kernel for MMD-GAN and on the feature map (Extractor) for GFMN, that ensure convergence of the generator to the data distribution. See Tab. 1 for a summary.

**GMMN, MMD-GAN Convergence: MMD Matching with Universal Kernels.** We start by reviewing known results on MMD. Let $\mathcal{H}_k$ be a Reproducing Kernel Hilbert Space (RKHS) defined with a continuous kernel $k$. Informally, *$k$ is universal if any bounded continuous function can be approximated to an arbitrary precision in $\mathcal{H}_k$* (formal definition in Appendix ). Theorem 1 [12] shows that the MMD is a well defined metric for universal kernels.

**Theorem 1** ([12]). *Given a kernel $k$, let $p,q$ be two distributions, their MMD is: $MMD^2(k,p,q) = ||\mu_p - \mu_q||^2_{\mathcal{H}_k}$, where $\mu_p = \mathbb{E}_{x\sim p}k_x$ is the mean embedding. If $k$ is universal then $MMD^2(k,p,q) = 0$ if and only if $p = q$.*

Given a Universal kernel such as a Gaussian Kernel as outlined in GMMN [22, 8], one can learn implicit Generative models $G_\theta$ that defines a family of distribution $\{q_\theta\}$ by minimizing the MMD distance:

$$\inf_{\theta} MMD(k, p_{data}, q_\theta) \quad (5)$$

Assuming $p_{data}$ is in the family $\{q_\theta\}$ ($\exists \theta^*, q_{\theta^*} = p_{data}$), the infimum of MMD minimization for a universal kernel is achieved for $q_\theta = p_{data}$ (immediate consequence of Theorem 1). This elegant setup for MMD matching with universal kernels, while avoiding the difficult min/max game in GAN, does not translate into good results in image generation. To remedy that, other discrepancies introduced in [21, 3, 33] compose universal kernels $k$ with a feature map $\phi \in \Psi$ as follows:

| | Metric | Kernel/Feature Map Convergence Conditions | Generative M. Optimization |
|---|---|---|---|
| GMMN [22, 8] | $\mathrm{MMD}(k, p, q)$ | Universal $k$ | min prob. |
| MMD-GAN [21, 3, 33] | $D_{\mathrm{MMD}}(p, q)$ | $k \circ \phi$<br>$k$ Fixed universal & lipschitz<br>$\phi$ lipschitz learned | min / max prob. |
| GFMN (This work) | $\mathrm{MMD}(K_\Phi, p, q)$ | Universal Feature Map $\Phi$<br>$\Phi$ fixed | min prob.<br>$\mu_{p_{data}}$ precomputed |

Table 1: Comparison of different approaches using MMD matching for implicit generative modeling. GFMN has two practical computational advantages it avoids the min/max game and allows to use a pre-computed mean embedding on the real data. Theoretically GFMN converges to the real data distribution if the feature extractor used was universal (See text for definition of Universal features as given in [26] ) .

$D_{\mathrm{MMD}}(p, q) = \sup_{\phi \in \Psi} \mathrm{MMD}(k \circ \phi, p, q)$. For learning implicit generative models [21] replaces MMD in Eq. (5) by $D_{\mathrm{MMD}}$. Under conditions on the kernel and the learned feature map this discrepancy is continuous in the weak topology (Prop. 2 in [1, 21]). Nevertheless, learning generative models remains challenging with it as it boils down to a min/max game as in original GAN [10].

**GFMN Convergence: MMD Matching with Universal Features.** While universality is usually thought on the kernel level, it is not straightforward to define universality for kernels defined by feature maps. Micchelli *et al.* [26] define universality of feature maps and how it connects to their corresponding kernels. Specifically for a fixed feature set on a space $\mathcal{X}$ (space of images) $S = \{\phi_j, j \in I, \phi_j : \mathcal{X} \to \mathbb{R}\}$, where $I$ is a countable set of indices, define the kernel $K_\phi(x, y) = \sum_{j \in I} \phi_j(x) \phi_j(y)$. Micchelli *et al.* [26] in Thm. 7 show that this Kernel is universal if the set $S$ is universal. Informally speaking, *a feature set $S$ is universal if linear functions in this feature space $(\sum_{j \in I} u_j \phi_j(x))$, are dense in the set of continuous bounded functions* (formal definition in Appendix 1).

This is of interest since GFMN corresponds to MMD matching with a kernel $K_\Phi$ defined on a fixed feature map $\Phi(x) = \{\phi_j(x)\}_{j \in I}$, where $I$ is finite. We have $K_\Phi(x, y) = \langle \Phi(x), \Phi(y) \rangle = \sum_{j \in I} \phi_j(x) \phi_j(y)$ and

$$\mathrm{MMD}^2(K_\Phi, p, q) = ||\mathbb{E}_{x \sim p} \Phi(x) - \mathbb{E}_{x \sim q} \Phi(x)||^2 .$$

For $\mathrm{MMD}^2(K_\Phi, p, q)$ to be a metric it is enough to have the set features $S$ be universal (by Thm.1 and Thm. 7 in [26]). Prop. 1 gives conditions for GFMN convergence:

**Proposition 1.** *Assume $p_{data}$ belongs to the family defined by the generator $\{q_\theta\}_\theta$. GFMN converges to the real distribution by matching in a feature space $S = \{\phi_j, j \in I\}$, where $I$ is a countable set, if the features set $S$ is universal (informally means that any continuous functions can be written as linear combination in the span of $S$) .*

*Proof.* $S$ is universal $\implies k_\Phi$ is universal [26]. Hence $\mathrm{MMD}(k_\Phi, p_{data}, q_\theta) = 0$ iff $q_\theta = p_{data}$. GFMN solves

$\inf_\theta \mathrm{MMD}^2(K_\Phi, p_{data}, q_\theta)$, and the infimum is achieved for $\theta$ such that $q_\theta = p_{data}$ ( $p_{data} \in \{q_\theta\}_\theta$ ). □

**Remark 1.** *The analysis covers here mean matching but the same applies to covariance matching considering $S = \{\phi_j, \phi_j \phi_k, j, k \in I\}$.*

**Universality of Perceptual Features in Computer Vision.** From Prop. 1 we see that for GFMN to be convergent with pretrained feature extractors $E_j$ that are **perceptual features** (such as features from VGG or ResNet pretrained on ImageNet), we need to assume universality of those features in the image domain. We know from transfer learning that features from ImageNet pretrained VGG/ResNet can express any functions for a downstream task by finding a linear weight in their span. Note that this is the definition of universal feature as given in [26]: continuous functions can be approximated in the linear span of those features. Hence, assuming universality of PFs defined by ImageNet pretrained VGG or ResNet, GFMN is guaranteed to converge to the data distribution by Prop. 1. Our results complement the common wisdom on "universality" of PFs in transfer learning and style transfer by showing that they are sufficient for learning implicit generative models.

## 4. Related work

GFMN is related to the recent body of work on MMD and moment matching based generative models [22, 8, 21, 3, 33]. The closest to our method is the Generative Moment Matching Network + Autoencoder (GMMN+AE) proposed in [22]. In GMMN+AE, the objective is to train a generator $G$ that maps from a prior uniform distribution to the latent code learned by a pretrained AE, and then uses the frozen pretrained decoder to map back to image space. As discussed in Section 3 one key difference in our approach is that, while GMMN+AE uses a Gaussian kernel to perform moment matching using the AE low dimensional latent code, GFMN performs mean and covariance matching in a PF space induced by a non-linear kernel function (a DCNN) that is orders of magnitude larger than the AE latent code, and that we argued is universal in the image domain.

Li *et al*. [21] demonstrate that GMMN+AE is not competitive with GANs for challenging datasets such as CIFAR10. MMD-GANs, discussed in Section 3, demonstrated competitive results with the use of adversarial learning by learning a feature map in conjuction with a Gaussian kernel [21, 3]. Finally, Ravuri *et al*. [33] recently proposed a method to perform online learning of the moments while training the generator. Our proposed method differs by using fixed pretrained PF extractors for moment matching.

Bojanowski *et al*. [4] proposed the Generative Latent Optimization (GLO) model that jointly optimizes model parameters and noise input vectors $z$, while avoiding adversarial training. Our work relates also to plug and play generative models of [30] where a pretrained classifier is used to sample new images, using MCMC sampling methods.

Our work is also related to AE-based generative models variational AE (VAE) [19], adversarial AE (AAE) [25] and Wasserstein AE (WAE) [38]. However, GFMN is quite distinct from these methods because it uses pretrained AEs to play the role of feature extractors only, while these methods aim to impose a prior distrib. on the latent space of AEs. Another recent line of work that involves the use of AEs in generative models consists in applying AEs to improve GANs stability [42, 39]. Finally, our objective function is related to the McGan loss function [29], where authors match first and second order moments.

# 5. Experiments

## 5.1. Experimental Setup

**Datasets:** We evaluate our proposed approach on images from CIFAR10 [20] (50k train., 10k test, 10 classes), STL10 [6] (5k train., 8k test, 100k unlabeled, 10 classes), CelebA [24] (200k) and LSUN bedrooms [41] datasets. STL10 images are rescaled to 32×32, while CelebA and LSUN images are rescaled to either 64×64 or 128×128, depending on the experiment. CelebA images are center-cropped to 160×160 before rescaling.

**GFMN Generator:** In most of our experiments the generator $G$ uses a DCGAN-like architecture [32]. For CIFAR10, STL10, LSUN and CelebA$_{64\times64}$, we use two extra layers as commonly used in previous works [28, 13]. For CelebA$_{128\times128}$ and some experiments with CIFAR10 and STL10, we use a ResNet-based generator such as the one in [13]. Architecture details are in the supplementary material.

**Autoencoder Features**: For most AE experiments, we use an encoder network whose architecture is similar to the discriminator in DCGAN (strided convolutions). We use batch normalization and ReLU non-linearity after each convolution. We set the latent code size to 128, 128, and 512 for CIFAR10, STL10 and CelebA, respectively. To perform feature extraction, we get the output of each ReLU in the network. Additionally, we also perform some experiments where the encoder uses a VGG19 architecture. The decoder

network $D$ uses a network architecture similar to our generator $G$. More details in the supplementary material.

**Classifier Features:** We perform our experiments on classifier features with VGG19 [37] and Resnet18 networks [14] which we pretrained using the whole ImageNet dataset [35] with 1000 classes. Pretrained ImageNet classifiers details can be found in the supplementary material.

**GFMN Training:** GFMNs are trained with an ADAM optimizer; most hyperparameters are kept fixed across datasets. We use $n_z = 100$ and minibatch of 64. Dataset dependent learning rates are used for updating $G$ ($10^{-4}$ or $5 \times 10^{-5}$) and AMA ($5\times10^{-5}$ or $10^{-5}$). We use AMA moving average (Sec. 2.2) in all reported experiments.

## 5.2. Autoencoder Features vs. (Cross-domain) Classifier Features

This section presents a comparative study on the use of pretrained autoencoders and cross-domain classifiers as feature extractors in GFMN. Tab. 2 shows the Inception Score (IS) [36] and Fréchet Inception Distance (FID) [15] for GFMN trained on CIFAR10 using different feature extractors $E$. The two first rows in Tab. 2 correspond to GFMN models that use pretrained encoders as $E$, while the last four rows use pretrained VGG19/Resnet18 ImageNet classifiers. We can see in Tab. 2 that there is a large boost in performance when ImageNet classifiers are used as feature extractors instead of encoders. Despite the classifiers being trained on a different domain (ImageNet vs. CIFAR10), the classifier features are significantly more effective. While the best IS with encoders is 4.95, the lowest IS with ImageNet classifier is 7.88. Additionally, when using simultaneously VGG19 and Resnet18 as feature extractors (two last rows), which increases the number of features to 832K, we get even better performance. Finally, we achieve the best performance in terms of both IS and FID (last row[1]) when using a generator architecture that contains residual blocks, similar to the one propose in [13].

Random samples from GFMN$^{\text{VGG19+Resnet18}}$ trained with CIFAR10 and STL10 are shown in Figs. 2a and 2b respectively. Fig. 2c shows random samples from GFMN$^{\text{VGG19}}$ trained with LSUN bedrooms dataset (resolution $64 \times 64$). Fig. 3 presents samples from GFMN$^{\text{VGG19}}$ trained with CelebA dataset with resolution 128×128, which shows that GFMN can achieve good performance with image resolutions larger than 32×32. These results also demonstrate that: (1) the same classifier (VGG19 trained on ImageNet) can be successfully applied to train GFMN models across different domains; (2) perceptual features from DCNNs encapsulate enough statistics to allow the learning of good generative models through moment matching.

Tab. 3 shows IS and FID for increasing number of layers (i.e. number of features) in our extractor VGG19. We se-

---

[1]Average result of five runs with different random seeds.

Table 2: CIFAR10 results for GFMN with different feature extractors.

| E Type | E Arch. | Pre-trained On | # features | G Arch. | IS | FID (5K/50K) |
|--------|---------|----------------|------------|---------|-----|--------------|
| Encoder | DCGAN | CIFAR10 | 60K | DCGAN | $4.51 \pm 0.06$ | 82.8 / 78.3 |
| Encoder | VGG19 | ImageNet | 296K | DCGAN | $4.95 \pm 0.06$ | 61.6 / 57.2 |
| Classifier | Resnet18 | ImageNet | 544K | DCGAN | $7.92 \pm 0.10$ | 29.1 / 24.3 |
| Classifier | VGG19 | ImageNet | 296K | DCGAN | $7.88 \pm 0.08$ | 25.5 / 20.8 |
| Classifier | VGG19 + Resnet18 | ImageNet | 832K | DCGAN | $8.08 \pm 0.08$ | 25.5 / 20.9 |
| Classifier | VGG19 + Resnet18 | ImageNet | 832K | Resnet | $\mathbf{8.27 \pm 0.09}$ | **18.1 / 13.5** |



(a) CIFAR10      (b) STL10      (c) LSUN

Figure 2: Generated samples from GFMN that uses as feature extractor VGG-19+Resnet18 (2a, 2b) and VGG-19 net (2c).
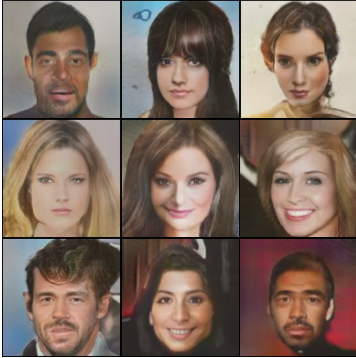


Figure 3: Samples from GFMN$^{\text{VGG19}}$ trained on CelebA with image size $128 \times 128$.

lect up to 16 layers, excluding the output of fully connected layers. Using more layers dramatically improves the performance of the feature extractor, reaching IS and FID peak performance when the maximum number of layers is used. Note that the features are ReLU activation outputs, meaning the encodings may be quite sparse. In Appendix 7 we show qualitative results that corroborate these results.

To verify whether the number of features is the main factor for performance, we conducted an experiment where we train an AE with an encoder using a VGG19 architecture. This encoder is pretrained on ImageNet and produces a total of 296K features. The second row in Tab. 2 shows the

results for this experiment. Although there is improvement in both IS and FID compared to the DCGAN encoder (first row), the boost is not comparable to the one obtained with a VGG19 classifier. In other words, features from classifiers are significantly more informative than AEs features for the purpose of training generators by feature matching.

Table 3: Impact of the number of layers/features used for feature matching in GFMN ($1K=2^{10}$).

| # layers | # features | IS | FID (5K / 50K) |
|----------|-----------|-----|----------------|
| 1 | 64K | $4.68 \pm 0.05$ | 118.6 / 114.8 |
| 3 | 160K | $5.59 \pm 0.08$ | 83.2 / 78.2 |
| 5 | 208K | $6.12 \pm 0.05$ | 53.8 / 49.3 |
| 7 | 240K | $6.99 \pm 0.06$ | 39.4 / 34.9 |
| 9 | 264K | $7.26 \pm 0.06$ | 32.3 / 27.7 |
| 11 | 280K | $7.72 \pm 0.08$ | 29.6 / 25.0 |
| 13 | 290K | $7.49 \pm 0.09$ | 29.2 / 24.8 |
| 15 | 294K | $7.62 \pm 0.04$ | 27.6 / 22.7 |
| 16 | 296K | $\mathbf{7.88 \pm 0.08}$ | **25.5 / 20.8** |

## 5.3. AMA and Training Stability

This section presents experimental results that evidence the advantage of our proposed ADAM moving average (AMA) over the simple moving average (MA). The main benefit of AMA is the promotion of stable training when using small minibatches. The ability to train with small minibatches is *essential* due to GFMN's need for large number of features from DCNNs, which becomes a challenge in

(a) MA - mbs 64         (b) MA - mbs 512         (c) AMA - mbs 64

Figure 4: Generated images from GFMN trained with either simple Moving Average (MA) (4a and 4b) or Adam Moving Average (AMA) (4c), and various minibatch sizes (mbs). While small minibatch sizes have a big negative effect for MA, it is not an issue for AMA.

terms of GPU memory usage. Our Pytorch [31] implementation of GFMN can only handle minibatches of size up to 160 when using VGG19 as a feature extractor and image size $64 \times 64$ on a Tesla K40 GPU w/ 12GB of memory. A more optimized implementation minimizing memory overhead could, in principle, handle somewhat larger minibatch sizes (as could a more recent Tesla V100 w/ 16 GB). However, increase image size or feature extractor size and the memory footprint increases quickly. We will always run out of memory when using larger minibatches, regardless of implementation or hardware.

All experiments in this section use CelebA training set, and a feature extractor using the encoder from an AE following a DCGAN-like architecture. This feature extractor is smaller than VGG19/Resnet18 allowing for minibatches of size up to 512 for image size $64 \times 64$. Fig. 4 shows generated images from GFMN trained with either MA or our proposed AMA. For MA, generated images from GFMN trained with 64 and 512 minibatch size are presented in Figs. 4a and 4b respectively. For AMA, Fig. 4c shows results for minibatch size 64. In MA training, the minibatch size has a *tremendous* impact on the quality of generated images: with minibatches smaller than 512, almost all images generated are quite distorted. On the other hand, when using AMA, GFMN generates much better images with minibatch size 64 (Fig. 4c). For AMA, increasing the minibatch size from 64 to 512 does not improve the quality of generated images for the given dataset and feature extractor. In the supplementary material, we show a comparison between MA and AMA with VGG19 ImageNet classifier as feature extractor for a minibatch size of 64. AMA also displays a very positive effect on the quality of generated images when a stronger feature extractor is used. An alternative for training with larger minibatches would be the use of multi-GPU, multi-node setups. However, performing large scale experiments is beyond the scope of the current work. Moreover,

many practitioners do not have access to a GPU cluster, and the development of methods that can also work on a single GPU with small memory footprint is essential.
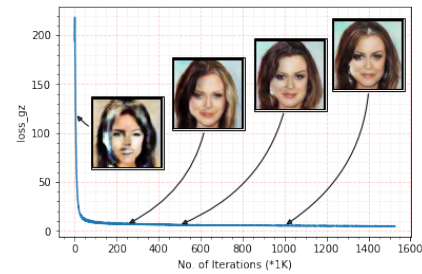


Figure 5: Loss as a function of training epochs with example of generated faces.

An important advantage of GFMN over adversarial methods is its training stability. Fig. 5 shows the evolution of the generator loss per epoch and generated examples when using AMA. There is a clear correlation between the quality of generated images and the loss. Moreover, mode collapsing was not observed in our experiments with AMA.

### 5.4. Comparison to the State-of-the-art

In Tab. 4, we compare GFMN results with different adversarial and non-adversarial approaches for CIFAR10 and STL10. In the middle part of the table, we report results for recent unsupervised models that use a DCGAN-like architecture in the generator. Despite using a frozen cross-domain feature extractor, GFMN outperforms the unsupervised systems in IS and FID for both datasets. The bottom part of Tab. 4 includes results for supervised approaches. Some of these models use a Resnet architecture in the generator as indicated in parenthesis. Note that GAN-based methods that perform conditional generation use direct feedback from the labels in the form of log likelihoods

Table 4: Inception Score and FID of different generative models for CIFAR10 and STL10.

| Model | CIFAR 10 | | STL 10 | |
|---|---|---|---|---|
| | IS | FID (5K / 50K) | IS | FID (5K / 50K) |
| Real data | 11.24±.12 | 7.8 / 3.2 | 26.08±.26 | 8.08 / 4.0 |
| **No Adversarial Training** | | | | |
| GMMN [21] | 3.47±.03 | | | |
| GMMN+AE [21] | 3.94±.04 | | | |
| (ours) GFMN$^{VGG+Resnet}$ | 8.08 ± 0.08 | 25.5 / 20.9 | 8.57 ± 0.08 | 34.2 / 17.2 |
| (ours) GFMN$^{VGG+Resnet}$ (Resnet G) | **8.27 ± 0.09** | **18.1 / 13.5** | **9.12 ± 0.09** | **31.6 / 13.9** |
| **Adversarial Training & Online Moment Learning Methods (Unsupervised)** | | | | |
| MMD GAN [21] | 6.17±.07 | | | |
| MMD$_{rq}$ GAN [3] | 6.51±.03 | 39.9 / - | | |
| WGAN-GP [27] | 6.68±.06 | 40.2 / - | 8.42±.13 | 55.1 / - |
| SN-GANs [27] | 7.58±.12 | 25.5 / - | 8.79±.14 | 43.2 / - |
| MoLM-1024 [33] | 7.55±.08 | 25.0 / 20.3 | | |
| GAN-DFM [39] | 7.72±.13 | | | |
| MoLM-1536 [33] | 7.90±.10 | 23.3 / 18.9 | | |
| **Adversarial Training (Supervised)** | | | | |
| Impr. GAN [36] | 8.09±.07 | | | |
| FisherGAN (Resnet G) [28] | 8.16±.12 | | | |
| WGAN-GP (Resnet G) [13] | 8.42±.10 | | | |

from the discriminator (e.g. using the $k+1$ trick from [36]). In contrast, our generator is trained with a loss function that *only* performs feature matching. Our generator is agnostic to the labels and there is no feedback in the form of a log likelihood from the labeled data. Despite that, GFMN produces results that are at the same level of supervised GAN models that use labels from the target dataset.

We performed additional experiments with a WGAN-GP architecture where: (1) the discriminator is a VGG19 or a Resnet18; (2) the discriminator is pretrained on ImageNet. The goal was to evaluate if WGAN-GP can benefit from DCNN classifiers pretrained on ImageNet. Although we tried different hyperparameter combinations, we were not able to successfully train WGAN-GP with VGG19 or Resnet18 discriminators (details in Appendix 8).

## 6. Discussion & Concluding Remarks

We achieve successful non-adversarial training of implicit generative models by introducing different key ingredients: (1) moment matching on perceptual features from all layers of pretrained neural networks; (2) a more robust way to compute the moving average of the mean features by using ADAM optimizer, which allows us to use small mini-batches; and (3) the use of perceptual features from multiple neural networks at the same time (VGG19 + Resnet18).

Our quantitative results in Tab. 4 show that GFMN achieves better or similar results compared to the state-of-the-art Spectral GAN (SN-GAN) [27] for both CIFAR10 and STL10. This is an impressive result for a non-adversarial feature matching-based approach that uses pretrained cross-domain feature extractors and has stable train-

ing. When compared to MMD approaches [22, 8, 21, 3, 33], GFMN presents important distinctions (some of them already listed in Secs. 3 and 4) which make it an attractive alternative. Compared to GMMN and GMMN+AE [22], we can see in Tab. 4 that GFMN achieves far better results. In the supplementary material, we also show a qualitative comparison between GFMN and GMMN results. Compared to recent adversarial MMD methods (MMD GAN) [21, 3] GFMN also presents significantly better results while avoiding the problematic min/max game. GFMN achieves better results than the Method of Learned Moments (MoLM) [33], while using a much smaller number of features to perform matching. The best performing model from [33], MoLM-1536, uses around 42 million moments to train the CIFAR10 generator, while our best GFMN model uses around 850K moments/features only, almost 50x less.

One may argue that the best GFMN results are obtained with feature extractors trained with classifiers. However, there are two important points to note: (1) we use a cross domain feature extractor and do not use labels from the target datasets (CIFAR10, STL10, LSUN, CelebA); (2) classifier accuracy does not seem to be the most important factor for generating good features: VGG19 classifier produces features as good as the ones from Resnet18, although the former is less accurate (more details in supplementary material). We are confident that GFMN can achieve state-of-the-art results with features from classifiers trained with unsupervised methods such as [5].

In conclusion, this work presents important theoretical and practical contributions that shed light on the effectiveness of perceptual features for training implicit generative models through moment matching.

# References

[1] M. Arbel, D. J. Sutherland, M. Bińkowski, and A. Gretton. On gradient regularizers for mmd gans. In *NIPS*, 2018. 4

[2] M. Arjovsky, S. Chintala, and L. Bottou. Wasserstein generative adversarial networks. In *Proc. of ICML*, pages 214–223, 2017. 13

[3] M. Bikowski, D. J. Sutherland, M. Arbel, and A. Gretton. Demystifying MMD GANs. In *International Conference on Learning Representations*, 2018. 1, 3, 4, 5, 8

[4] P. Bojanowski, A. Joulin, D. Lopez-Paz, and A. Szlam. Optimizing the latent space of generative networks, 2018. 5, 12

[5] M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision*, 2018. 8

[6] A. Coates, A. Ng, and H. Lee. An analysis of single-layer networks in unsupervised feature learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 215–223, 2011. 5

[7] A. Dosovitskiy and T. Brox. Generating images with perceptual similarity metrics based on deep networks. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 658–666. Curran Associates, Inc., 2016. 1

[8] G. K. Dziugaite, D. M. Roy, and Z. Ghahramani. Training generative neural networks via maximum mean discrepancy optimization. In *Proceedings of the Thirty-First Conference on Uncertainty in Artificial Intelligence*, pages 258–267, 2015. 3, 4, 8

[9] L. A. Gatys, A. S. Ecker, and M. Bethge. Image style transfer using convolutional neural networks. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016. 1

[10] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Proc. of NIPS*, page 2672, 2014. 4, 13

[11] A. Gretton, K. M. Borgwardt, M. Rasch, B. Schölkopf, and A. J. Smola. A kernel method for the two-sample-problem. In *Proceedings of the 19th International Conference on Neural Information Processing Systems*, pages 513–520, 2006. 1

[12] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola. A kernel two-sample test. *Journal of Machine Learning Research*, 13:723–773, 2012. 1, 3

[13] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. C. Courville. Improved training of wasserstein gans. *CoRR*, 2017. 5, 8

[14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 5

[15] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, G. Klambauer, and S. Hochreiter. Gans trained by a two time-scale update rule converge to a nash equilibrium. *CoRR*, abs/1706.08500, 2017. 5, 12

[16] M. Huh, P. Agrawal, and A. A. Efros. What makes imagenet good for transfer learning? *CoRR*, abs/1608.08614, 2016. 1

[17] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *European Conference on Computer Vision*, 2016. 1

[18] D. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 3

[19] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 5

[20] A. Krizhevsky. Learning multiple layers of features from tiny images. page 60, 2009. 5

[21] C.-L. Li, W.-C. Chang, Y. Cheng, Y. Yang, and B. Poczos. MMD GAN: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, pages 2203–2213. 2017. 1, 3, 4, 5, 8, 14, 15

[22] Y. Li, K. Swersky, and R. Zemel. Generative moment matching networks. In *Proceedings of the International Conference on International Conference on Machine Learning*, pages 1718–1727, 2015. 1, 3, 4, 8, 14

[23] H. Ling and K. Okada. Diffusion distance for histogram comparison. In *Computer Vision and Pattern Recognition*, 2006. 12

[24] 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)*, 2015. 5

[25] A. Makhzani, J. Shlens, N. Jaitly, and I. Goodfellow. Adversarial autoencoders. In *International Conference on Learning Representations*, 2016. 5

[26] C. A. Micchelli, Y. Xu, and H. Zhang. Universal kernels. *J. Mach. Learn. Res.*, 7:2651–2667, Dec. 2006. 4, 11

[27] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018. 8, 12, 13

[28] Y. Mroueh and T. Sercu. Fisher GAN. In *Proceedings of NIPS*, 2017. 5, 8

[29] Y. Mroueh, T. Sercu, and V. Goel. McGan: Mean and covariance feature matching GAN. In *Proceedings of the 34th International Conference on Machine Learning*, pages 2527–2535, 2017. 5, 12

[30] A. Nguyen, J. Clune, Y. Bengio, A. Dosovitskiy, and J. Yosinski. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Conference on Computer Vision and Pattern Recognition*, pages 3510–3520, 2017. 5

[31] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS-W*, 2017. 7

[32] A. Radford, L. Metz, and S. Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. In *ICLR*, 2016. 5

[33] S. V. Ravuri, S. Mohamed, M. Rosca, and O. Vinyals. Learning implicit generative models with the method of learned moments. In *Proceedings of the 35th International Conference on Machine Learning*, pages 4311–4320, 2018. 1, 3, 4, 5, 8, 12

[34] S. J. Reddi, S. Kale, and S. Kumar. On the convergence of adam and beyond. In *ICLR*, page 23, 2018. 3

[35] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *Int. J. Comput. Vision*, 115(3):211–252, Dec. 2015. 1, 5

[36] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen. Improved techniques for training gans. In *Proc. of NIPS*, pages 2226–2234, 2016. 5, 8, 12

[37] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *CoRR*, abs/1409.1556, 2014. 5

[38] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations*, 2018. 5

[39] D. Warde-Farley and Y. Bengio. Improving generative adversarial networks with denoising feature matching. In *Proceedings of ICLR*, 2017. 5, 8

[40] J. Yosinski, J. Clune, Y. Bengio, and H. Lipson. How transferable are features in deep neural networks? In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'14, 2014. 1

[41] F. Yu, Y. Zhang, S. Song, A. Seff, and J. Xiao. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015. 5

[42] J. J. Zhao, M. Mathieu, and Y. LeCun. Energy-based generative adversarial network. In *Proceedings of ICLR*, 2017. 5

# 7. Appendix

## 7.1. Continuation of Universality of PFs and GFMN Convergence

We summarize here the main definitions and theorems from [26] regarding universality of kernels and feature maps.

**Universal Kernels.** The following defines a universal kernel

**Definition 1** (Universal Kernel). *Given a kernel $K$ defined on $\mathcal{X} \times \mathcal{X}$. Let $\mathcal{Z}$ be any compact subset of $\mathcal{X}$. Define the space of kernel sections:*

$$\mathcal{K}(\mathcal{Z}) = \overline{span}\{K_y, y \in \mathcal{Z}\},$$

*where $K_y : \mathcal{X} \to \mathbb{R}$, $K_y(x) = K(x,y)$. Let $\mathcal{C}(Z)$ be the space of all continuous real valued functions defined on $\mathcal{Z}$. A kernel is said **universal** if for any choice of $\mathcal{Z}$ (compact subset of $\mathcal{X}$) $\mathcal{K}(\mathcal{Z})$ is dense in $\mathcal{C}(\mathcal{Z})$.*

In other words a kernel is universal if $\mathcal{C}(\mathcal{Z}) = K(\mathcal{Z})$. Meaning if any continuous function can be expressed in the span of $K_y$.

**Universal Feature Maps.** We turn now for kernels defined by feature maps and how to characterize their universality. Consider a continuous feature map $\Phi : \mathcal{X} \to \mathcal{W}$, where $(\mathcal{W}, \langle, \rangle_{\mathcal{W}})$ is a Hilbert space; the kernel $K$ has the following form:

$$K(x,y) = \langle \Phi(x), \Phi(y) \rangle_{\mathcal{W}}. \tag{6}$$

Let $\mathcal{Y}$ be an orthonormal basis of $\mathcal{W}$ define the following continuous function $F_y \in \mathcal{C}(\mathcal{Z})$ defined at $x \in \mathcal{Z}$:

$$F_y(x) = \langle \Phi(x), y \rangle_{\mathcal{W}},$$

and let:
$$\Phi(\mathcal{Y}) = \overline{span}\{F_y, y \in \mathcal{Y}\}$$

**Definition 2** (Universal feature Map). *A feature map is universal if $\Phi(\mathcal{Y})$ is dense in $\mathcal{C}(\mathcal{Z})$, for all $\mathcal{Z}$ compact subsets of $\mathcal{X}$.i.e A feature map is universal if $\Phi(\mathcal{Y}) = \mathcal{C}(\mathcal{Z})$.*

The following Theorem shows the relation between universality of a kernel defined by feature map and the universality of the feature map:

**Theorem 2** ([26], Thm 4, Relation between $\mathcal{K}(\mathcal{Z})$ and $\Phi(\mathcal{Y})$ ). *For kernel defined by feature maps in (6) we have $\mathcal{K}(\mathcal{Z}) = \Phi(\mathcal{Y})$. A kernel of form (6) is universal if and only if its feature map is universal.*

Hence the following Theorem 7 from [26]:

**Theorem 3** ([26]). *Let $S = \{\phi_j, j \in I\}$, where $I$ is a countable set and $\phi_j : \mathcal{X} \to \mathbb{R}$ continuous function. Define the following kernel*

$$K(x,y) = \sum_{j \in I} \phi_j(x)\phi_j(y).$$

*$K$ is universal if and only if the set of features $S$ is universal.*

## 7.2. Discussion of AMA versus MA

As we already discussed the moving average of $v$ of the difference of features means

$$\Delta_t = \frac{1}{N}\sum_{i=1}^{N} E(x_i) - \frac{1}{N}\sum_{i=1}^{N} E(G(z_i, \theta_t))$$

between real and generated data at each time step $t$ in the gradient descent up to time $T$, can be seen as a gradient descent in an online setting on the following cost :

$$f^* = \min_v \sum_{t=1}^{T} f_t(v) = \sum_{t=1}^{T} ||v - \Delta_t||_2^2$$

Note that we are in the online setting since $\Delta_t$ is only known when $\theta_t$ of the generator is updated. The sequence $v_t$ generated by MA (moving average) and by AMA (ADAM moving average) is the SGD updates and ADAM updates respectively applied to the cost function $f_t$. Hence we can bound the regret of the sequence $\{v_t^{\text{MA}}\}$ and $\{v_t^{\text{AMA}}\}$ using known results on SGD and ADAM. Let $d$ be the dimension of the encoding $E$. For MA, using classic regret bounds for gradient descents we obtain:

$$R_T^{\text{MA}} = \sum_{t=1}^{T} ||v_t^{\text{MA}} - \Delta_t||_2^2 - f^* \leq O(\sqrt{dT}).$$

For AMA, using ADAM regrets bounds from (Reddi et al., 2018). Let us define

$$R_T^{\text{AMA}} = \sum_{t=1}^{T} ||v_t^{\text{AMA}} - \Delta_t||_2^2 - f^*.$$

We have:

$$R_T^{\text{AMA}} \leq O(\sqrt{T}\sum_{i=1}^{d} \hat{u}_i^{T, \frac{1}{2}}) + \cdots$$

$$O\left(\sum_{i=1}^{d}\sqrt{\sum_{t=1}^{T}(\Delta_{t,i} - v_{t,i}^{\text{AMA}})^2}\right) + C$$

where $\hat{u}$ are defined in the ADAM updates as moving averages of second order moments of the gradients. The

regret bound of AMA is better than MA especially if $\sum_{i=1}^{d} \hat{u}_i^{T,\frac{1}{2}} \ll d$ and

$$\sum_{i=1}^{d} \sqrt{\sum_{t=1}^{T} (\Delta_{t,i} - v_{t,i}^{\text{AMA}})^2} \ll \sqrt{Td}.$$

### 7.3. Mean Matching vs. Mean + Covariance Matching in GFMN

In this Appendix, we present comparative results between GFMN with mean feature matching vs. GFMN with mean + covariance feature matching. Using the first and second moments to perform feature matching gives statistical advantage over using the first moment only. In Table 5, we can see that for different feature extractors, performing mean + covariance feature matching produces significantly better results in terms of both IS and FID. Mroueh *et al.* [29] have also demonstrated the advantages of using mean + covariance matching in the context of GANs.

### 7.4. Neural Network Architectures

In Tables 6 and 7, and Figure 6 we detail the neural net architectures used in our experiments. In both DCGAN-like generator and discriminator, an extra layer is added when using images of size 64×64. In VGG19 architecture, after each convolution, we apply batch normalization and ReLU. The Resnet generator is used for CelebA$_{128\times128}$ experiments and also for some experiments with CIFAR10 and STL10. For these two last datasets, the Resnet generator has 3 ResBlocks only, and the output size of the $DENSE$ layer is $4 \times 4 \times 512$.
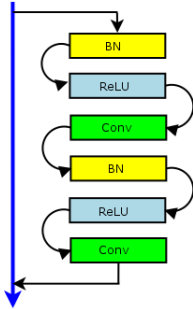


Figure 6: ResBlock

### 7.5. Pretraining of ImageNet Classifiers and Autoencoders

Both VGG19 and Resnet18 networks are trained with SGD with fixed $10^{-1}$ learning rate, 0.9 momentum term, and weight decay set to $5 \times 10^{-4}$. We pick models with best *top-1* accuracy on the validation set over 100 epochs of training; 29.14% for VGG19 (image size 32×32), and 39.63% for Resnet18 (image size 32×32). When training the classifiers we use random cropping and random horizontal flipping for data augmentation. When using VGG19 and Resnet18 as feature extractors in GFMN, we use features from the output of each ReLU that follows a conv. layer, for a total of 16 layers for VGG and 17 for Resnet18.

In our experiments with autoencoders (AE) we pretrained them using either mean squared error (MSE) or the Laplacian pyramid loss [23, 4]. Let $E$ and $D$ be the encoder and the decoder networks with parameters $\phi$ and $\psi$, respectively.

$$\min_{\phi,\psi} \mathbb{E}_{p_{data}} ||x - D(E(x;\phi);\psi)||^2$$

or the Laplacian pyramid loss [23]

$$\text{Lap}_1(x, x') = \sum_j 2^{-2j} |L^j(x) - L^j(x')|_1$$

where $L^j(x)$ is the $j$-th level of the Laplacian pyramid representation of $x$. The Laplacian pyramid loss provides better signal for learning high frequencies of images and overcome some of the blurriness issue known from using a simple MSE loss. [4] recently demonstrated that the Lap$_1$ loss produces better results than $L_2$ loss for both autoencoders and generative models.

### 7.6. Quantitative Evaluation Metrics

We evaluate our models using two quantitative metrics: Inception Score (IS) [36] and Fréchet Inception Distance (FID) [15]. We followed the same procedure used in previous work to calculate IS [36, 27, 33]. For each trained generator, we calculate the IS for randomly generated 5000 images and repeat this procedure 10 times (for a total of 50K generated images) and report the average and the standard deviation of the IS.

We compute FID using two sample sizes of generated images: 5K and 50K. In order to be consistent with previous works [27, 33] and be able to directly compare our quantitative results with theirs, the FID is computed as follows:

- CIFAR10: the statistics for the real data are computed using the 50K training images. This (real data) statistics are used in the FID computation of both 5K and 50K samples of generated images. This is consistent with both Miyato *et al.* [27] and Ravuri *et al.* [33] procedure to compute FID for CIFAR10 experiments.

- STL10: when using 5K generated images, the statistics for the real data are computed using the set of 5K (labeled) training images. This is consistent with the FID

Table 5: CIFAR10 results for GFMN with Mean Feature Matching vs. GFMN with Mean + Covariance Feature Matching.

| Feature Extractor | Mean Matching | | Mean + Covar. Matching | |
|---|---|---|---|---|
| | IS | FID (5K / 50K) | IS | FID (5K / 50K) |
| DCGAN (Encoder) | $3.76 \pm 0.04$ | 96.5 / 92.5 | $4.51 \pm 0.06$ | 82.8 / 78.3 |
| Resnet18 | $7.03 \pm 0.11$ | 35.7 / 31.1 | $7.92 \pm 0.10$ | 29.1 / 24.3 |
| VGG19 | $7.42 \pm 0.09$ | 27.5 / 22.8 | $7.88 \pm 0.08$ | 25.5 / 20.8 |

Table 6: DCGAN like Generator

| $z \in \mathbb{R}^{100} \sim \mathcal{N}(0, I)$ |
|---|
| DENSE $\rightarrow 4 \times 4 \times 512$ |
| $4 \times 4$, STRIDE=2 DECONV BN 256 RELU |
| $4 \times 4$, STRIDE=2 DECONV BN 128 RELU |
| $4 \times 4$, STRIDE=2 DECONV BN 64 RELU |
| $3 \times 3$, STRIDE=1 CONV 3 BN 64 RELU |
| $3 \times 3$, STRIDE=1 CONV 3 BN 64 RELU |
| $3 \times 3$, STRIDE=1 CONV 3 TANH |

Table 7: Resnet Generator

| $z \in \mathbb{R}^{100} \sim \mathcal{N}(0, I)$ |
|---|
| DENSE, $4 \times 4 \times 2048$ |
| RESBLOCK UP 1024 |
| RESBLOCK UP 512 |
| RESBLOCK UP 256 |
| RESBLOCK UP 128 |
| RESBLOCK UP 164 |
| BN, RELU, $3 \times 3$ CONV 3 |
| TANH |

computation of Miyato *et al.* [27]. When using 50K generated images, the statistics for the real data are computed using a set of 50K images randomly sampled from the unlabeled STL10 dataset.

FID computation is repeated 3 times and the average is reported. There is very small variance in the FID results.

### 7.7. Impact of the number of layers used for feature extraction

Figure 7 shows generated images from generators that were trained with a different number of layers employed to feature matching. In all the results in Fig.7, the VGG19 network was used to perform feature extraction. We can see a significant improvement in image quality when more layers are used. Better results are achieved when 11 or more layers are used, which corroborates the quantitative results in Sec. 5.2.

### 7.8. Pretrained Generator/Discriminator in WGAN-GP

The objective of the experiments presented in this section is to evaluate if WGAN-GP can benefit from DCNN classifiers pretrained on ImageNet. In the experiments, we used a WGAN-GP architecture where: (1) the discriminator is a VGG19 or a Resnet18; (2) the discriminator is pretrained on ImageNet; (3) the generator is pretrained on CIFAR10 through autoencoding. Although we tried different hyperparameter combinations, we were not able to successfully train WGAN-GP with VGG19 or Resnet18 discriminators. Indeed, the discriminator, being pretrained on ImageNet, can quickly learn to distinguish between real and fake images. This limits the reliability of the gradient information from the discriminator, which in turn renders the training of a proper generator extremely challenging or even impossible. This is a well-known issue with GAN training [10] where the training of the generator and discriminator must strike a balance. This phenomenon is covered in [2] Section 3 (illustrated in their Figure 2) as one motivation for work like Wassertein GANs. If a discriminator can distinguish perfectly between real and fake early on, the generator cannot learn properly and the min/max game becomes unbalanced, having no good discriminator gradients for the generator to learn from, producing degenerate models. Figure 8 shows some examples of images generated by the unsuccessfully trained models.

### 7.9. Impact of Adam Moving Average for VGG19 feature extractor.

In this appendix, we present a comparison between the simple moving average (MA) and ADAM moving average (AMA) for the case where VGG19 ImageNet classifier is used as a feature extractor. This experiment uses a minibatch size of 64. We can see in Fig. 9 that AMA has a very positive effect in the quality of generated images. GFMN trained with MA produces various images with some sort of crossing line artifacts.

### 7.10. Visual Comparison between GFMN and GMMN Generated Images.

Figure 10 shows a visual comparison between images generated by GFMN (Figs. 10a and 10b) and Generative Moment Matching Networks (GMMN) (Figs. 10c and 10d).

(a) 1 Layer     (b) 3 Layers     (c) 5 Layers

(d) 7 Layers     (e) 9 Layers     (f) 11 Layers

(g) 13 Layers     (h) 15 Layers     (i) 16 Layers

Figure 7: Generated images from GFMN trained with a different number of VGG19 layers for feature extraction.

GMMN [22] generated images were obtained from Li *et al.* [21]. In this experiment, both GMMN and GFMN use a DCGAN-like architecture in the generator. Images generated by GFMN have significantly better quality compared to the ones generated by GMMN, which corroborates the quantitative results in Sec. 5.4.

## 7.11. Autoencoder features vs. VGG19 features for CelebA.

In this appendix, we present a comparison in image quality for autoencoder features vs. VGG19 features for the CelebA dataset. We show results for both simple moving

Figure 8: Generated images by WGAN-GP with pretrained VGG19 as a discriminator.
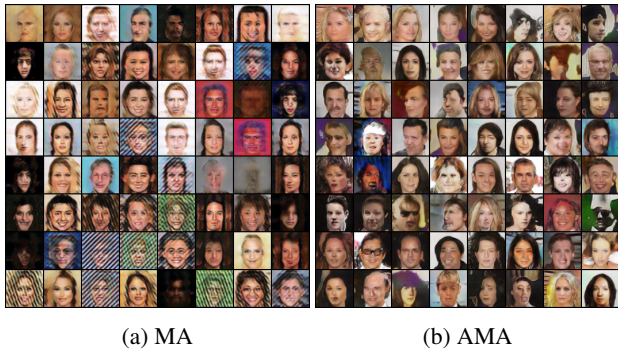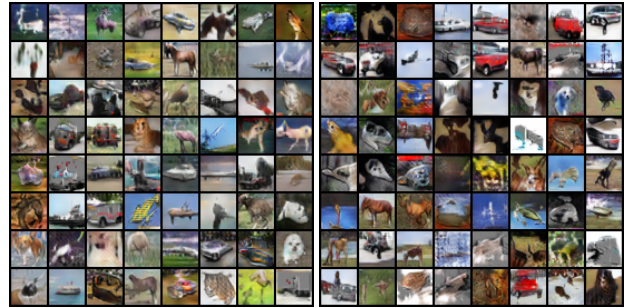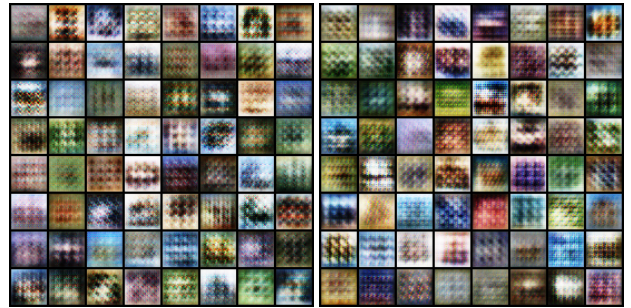


(a) MA          (b) AMA

Figure 9: Generated images from GFMN trained with either simple moving average (MA) or Adam moving average (AMA). VGG19 ImageNet classifier is used as feature extractor.

average (MA) and ADAM moving average (AMA), for both cases we use a minibatch size of 64. In Fig. 11, we show generated images from GFMN trained with either VGG19 features (top row) or autoencoder (AE) features (bottom row). We show images generated by GFMN models trained with simple moving average (MA) and Adam moving average (AMA). We can note in the images that, although VGG19 features are from a cross-domain classifier, they lead to much better generation quality than AE features, specially for the MA case.
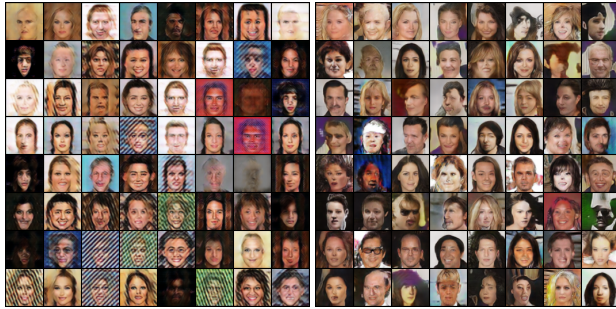


(a) GFMN with VGG19 features    (b) GFMN with Resnet18 features
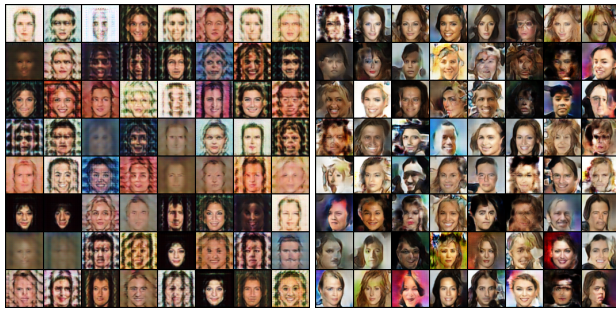


(c) GMMN - Matching on data space    (d) GMMN+AE - Matching on AE space

Figure 10: Generated images from GFMN (10a and 10b) and GMMN (10c and 10d). GMMN images were obtained from Li *et al.* [21].

(a) MA - VGG19 Features      (b) AMA - VGG19 Features

(c) MA - AE Features      (d) AMA - AE Features

Figure 11: Generated images from GFMN trained with either VGG19 features (top row) or autoencoder (AE) features (bottom row). We show images generated by GFMN models trained with simple moving average (MA) and Adam moving average (AMA). Although VGG19 features are from a cross-domain classifier, they perform much better than AE features, specially for the MA case.