

# Diagnosing and Fixing Manifold Overfitting in Deep Generative Models

Anonymous authors

Paper under double-blind review

## Abstract

Likelihood-based, or *explicit*, deep generative models use neural networks to construct flexible high-dimensional densities. This formulation directly contradicts the manifold hypothesis, which states that observed data lies on a low-dimensional manifold embedded in high-dimensional ambient space. In this paper we investigate the pathologies of maximum-likelihood training in the presence of this dimensionality mismatch. We formally prove, in a measure-theoretic way, that degenerate optima are achieved wherein the manifold itself is learned but not the distribution on it, a phenomenon we call *manifold overfitting*. We propose a class of two-step procedures consisting of a dimensionality reduction step followed by maximum-likelihood density estimation, and prove that they recover the data-generating distribution in the nonparametric regime, thus avoiding manifold overfitting. We also show that these procedures enable density estimation on the manifolds learned by *implicit* models, such as generative adversarial networks, hence addressing a major shortcoming of these models. Several recently proposed methods are instances of our two-step procedures; we thus unify, extend, and theoretically justify a large class of models.

## 1 Introduction

We consider the standard setting for generative modelling, where samples  $\{x_n\}_{n=1}^N \subset \mathbb{R}^D$  of high-dimensional data from some unknown distribution  $\mathbb{P}^*$  are observed, and the task is to estimate  $\mathbb{P}^*$ . Many deep generative models (DGMs) (Bond-Taylor et al., 2021), including variational autoencoders (VAEs) (Kingma & Welling, 2014; Rezende et al., 2014; Ho et al., 2020; Kingma et al., 2021) and variants such as adversarial variational Bayes (AVB) (Mescheder et al., 2017), normalizing flows (NFs) (Dinh et al., 2017; Kingma & Dhariwal, 2018; Behrmann et al., 2019; Chen et al., 2019; Durkan et al., 2019; Cornish et al., 2020), energy-based models (EBMs) (Du & Mordatch, 2019), and continuous autoregressive models (ARMs) (Uria et al., 2013; Theis & Bethge, 2015), use neural networks to construct a flexible density trained to match  $\mathbb{P}^*$  by maximizing either the likelihood or a lower bound of it. This modelling choice implies the model has  $D$ -dimensional support, thus directly contradicting the manifold hypothesis (Bengio et al., 2013), which states that high-dimensional data is supported on an unknown  $d$ -dimensional manifold  $\mathcal{M} \subset \mathbb{R}^D$ , where  $d < D$ .

There is strong evidence supporting the manifold hypothesis (Pope et al., 2021), and a natural question arises: *how relevant is this modelling mismatch?* We answer this question by proving that when  $\mathbb{P}^*$  is supported on  $\mathcal{M}$ , maximum-likelihood training of a flexible  $D$ -dimensional density results in  $\mathcal{M}$  itself being learned, but not  $\mathbb{P}^*$ . Our result extends that of Dai & Wipf (2019) beyond VAEs to all likelihood-based models and drops the empirically unrealistic assumption that  $\mathcal{M}$  is homeomorphic to  $\mathbb{R}^d$  (e.g. one can imagine the MNIST (LeCun, 1998) manifold as having 10 connected components, one per digit). This phenomenon – which we call *manifold overfitting* – has profound consequences for generative modelling. Maximum-likelihood is indisputably one of the most important concepts in statistics, and enjoys well-studied theoretical properties such as consistency and asymptotic efficiency under seemingly mild regularity conditions (Lehmann & Casella, 2006). These conditions can indeed be reasonably expected to hold in the setting of “classical statistics” under which they were first considered, where models were simpler and available data was of much lower ambient dimension than by modern standards. However, in the presence of  $d$ -dimensional manifold structure, the previously innocuous assumption that there exists a ground truth  $D$ -dimensional density cannot possibly hold. Manifold overfitting thus shows

that DGMs do not enjoy the supposed theoretical benefits of maximum-likelihood, which is often regarded as a principled objective for training DGMs, because they will recover the manifold but not the distribution on it.

In order to address manifold overfitting, we propose a class of two-step procedures, depicted in Fig. 1. The first step, which we call *generalized autoencoding*<sup>1</sup>, reduces the dimension of the data through an encoder  $g : \mathbb{R}^D \rightarrow \mathbb{R}^d$  while also learning how to map back to  $\mathcal{M}$  through a decoder  $G : \mathbb{R}^d \rightarrow \mathbb{R}^D$ . In the second step, maximum-likelihood estimation with a DGM is performed on the low-dimensional representations  $\{g(x_n)\}_{n=1}^N$ . Intuitively, the first step removes the dimensionality mismatch in order to avoid manifold overfitting in the second step. This intuition is confirmed in a second theoretical result where we prove that, given enough capacity, our two-step procedures indeed recover  $\mathbb{P}^*$  in the infinite data limit while retaining density evaluation. We also identify DGMs that are instances of our procedure class. Our methodology thus results in novel models, and provides a unifying perspective and theoretical justification for all these related works.

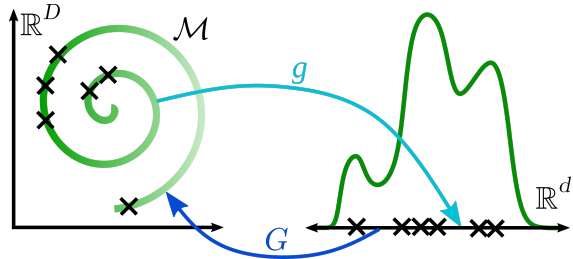


Figure 1: Depiction of our two-step procedures. In the first step, we learn to map from  $\mathcal{M}$  to  $\mathbb{R}^d$  through  $g$ , and to invert this mapping through  $G$ . In the second step, we perform density estimation (green density on the right) on the dataset encoded through  $g$ . Our learned distribution on  $\mathcal{M}$  (shades of green on the spiral) is given by pushing forward the density from the second step through  $G$ .

We also show that some implicit models (Mohamed & Lakshminarayanan, 2016), e.g. generative adversarial networks (GANs) (Goodfellow et al., 2014), can be made into generalized autoencoders. Consequently, in addition to preventing manifold overfitting on explicit models, our two-step procedures enable density evaluation for implicit models, thus addressing one of their main limitations.

Finally, we achieve significant empirical improvements in sample quality over maximum-likelihood, strongly supporting our theoretical findings. We show these improvements persist even when accounting for the additional parameters of the second-step model, or when adding Gaussian noise to the data as an attempt to remove the dimensionality mismatch that causes manifold overfitting. We also obtain very promising results in out-of-distribution (OOD) detection using only likelihoods.

## 2 Related Work and Motivation

**Manifold mismatch** It has been observed in the literature that  $\mathbb{R}^D$ -supported models exhibit undesirable behaviour when the support of the target distribution has complicated topological structure. For example, Cornish et al. (2020) show that the bi-Lipschitz constant of topologically-misspecified NFs must go to infinity, even without dimensionality mismatch, explaining phenomena like the numerical instabilities observed by Behrmann et al. (2021). Mattei & Frellsen (2018) observe VAEs can have unbounded likelihoods and are thus susceptible to similar instabilities. Dai & Wipf (2019) study dimensionality mismatch in VAEs and its effects on posterior collapse. These works motivate the development of models with low-dimensional support. Goodfellow et al. (2014) and Nowozin et al. (2016) model the data as the pushforward of a low-dimensional Gaussian through a neural network, thus making it possible to properly account for the dimension of the support. However, in addition to requiring adversarial training – which is more unstable than maximum-likelihood (Chu et al., 2020) – these models minimize the Jensen-Shannon divergence or  $f$ -divergences, respectively, in the *nonparametric* setting (i.e. infinite data limit with sufficient capacity), which are ill-defined due to dimensionality mismatch. Attempting to minimize Wasserstein distance has also been proposed (Arjovsky et al., 2017; Tolstikhin et al., 2018) as a way to remedy this issue, although estimating this distance is hard in practice Arora et al. (2017) and unbiased gradient estimators are not available. In addition to having a more challenging training objective than maximum-likelihood, these *implicit* models lose a key advantage of *explicit* models: density evaluation. Our work aims to both properly account

<sup>1</sup>Our generalized autoencoders are unrelated to those of Wang et al. (2014).

for the manifold hypothesis in likelihood-based DGMs while retaining density evaluation, and endow implicit models with density evaluation.

**NFs on manifolds** Several recent flow-based methods properly account for the manifold structure of the data. Gemici et al. (2016), Rezende et al. (2020), and Mathieu & Nickel (2020) construct flow models for prespecified manifolds, with the obvious disadvantage that the manifold is unknown for most data of interest. Brehmer & Cranmer (2020) propose injective NFs, which model the data-generating distribution as the pushforward of a  $d$ -dimensional Gaussian through an injective function  $G : \mathbb{R}^d \rightarrow \mathbb{R}^D$ , and avoid the change-of-variable computation through a two-step training procedure; we will see in Sec. 5 that this procedure is an instance of our methodology. Caterini et al. (2021) and Ross & Cresswell (2021) endow injective flows with tractable change-of-variable computations, the former through automatic differentiation and numerical linear algebra methods, and the latter with a specific construction of injective NFs admitting closed-form evaluation. We build a general framework encompassing a broader class of DGMs than NFs alone, giving them low-dimensional support without requiring injective transformations over  $\mathbb{R}^d$ .

**Adding noise** Denoising approaches add Gaussian noise to the data, making the  $D$ -dimensional model appropriate at the cost of recovering a noisy version of  $\mathbb{P}^*$  (Vincent et al., 2008; Vincent, 2011; Alain & Bengio, 2014; Chae et al., 2021; Horvat & Pfister, 2021a;b; Cunningham & Fiterau, 2021). In particular, Horvat & Pfister (2021b) show that recovering the true manifold structure in this case is only guaranteed when adding noise orthogonally to the tangent space of the manifold, which cannot be achieved in practice when the manifold itself is unknown. In the context of score-matching (Hyvärinen, 2005), denoising has led to empirical success (Song & Ermon, 2019; Song et al., 2021). In Sec. 3.2 we show that adding small amounts of Gaussian noise to a distribution supported on a manifold results in highly peaked densities, which can be hard to learn. Zhang et al. (2020a) also make this observation, and propose to add the same amount of noise to the model itself. However, their method requires access to the density of the model after having added noise, which in practice requires a variational approximation and is thus only applicable to VAEs. Our first theoretical result can be seen as a motivation for the method of Zhang et al. (2020a), and our procedures are applicable to all likelihood-based DGMs. We empirically verify that simply adding noise is not enough to avoid manifold overfitting in practice, and that our two-step methodology outperforms this approach.

### 3 Manifold Overfitting

#### 3.1 An Illustrative Example

Consider the simple case where  $D = 1$ ,  $d = 0$ ,  $\mathcal{M} = \{-1, 1\}$ , and  $\mathbb{P}^* = 0.3\delta_{-1} + 0.7\delta_1$ , where  $\delta_x$  denotes a point mass at  $x$ . Suppose the data is modelled with a mixture of Gaussians  $p(x) = \lambda \cdot \mathcal{N}(x; m_1, \sigma^2) + (1 - \lambda) \cdot \mathcal{N}(x; m_2, \sigma^2)$  parameterized by  $\lambda \in [0, 1]$ ,  $m_1, m_2 \in \mathbb{R}$ , and  $\sigma^2 \in \mathbb{R}_{>0}$ , which we will think of as a flexible density. This model can learn the correct distribution in the limit  $\sigma^2 \rightarrow 0$ , as shown on the left panel of Fig. 2 (dashed line). However, arbitrarily large likelihood values can be achieved by other densities – the one shown with a dotted line approximates a distribution  $\mathbb{P}^\dagger$  on  $\mathcal{M}$  which *is not*  $\mathbb{P}^*$  but nonetheless has large likelihoods. The implication is simple: maximum-likelihood estimation will not necessarily recover the data-generating distribution  $\mathbb{P}^*$ . Our choice of  $\mathbb{P}^\dagger$  (see figure caption) was completely arbitrary, hence any distribution on  $\mathcal{M}$  other than  $\delta_{-1}$  or  $\delta_1$  could be recovered with likelihoods diverging to infinity. Recovering  $\mathbb{P}^*$  is then a coincidence which we should not expect to occur when training via maximum-likelihood. In other words, we should expect maximum-likelihood to recover the manifold (i.e.  $m_1 = \pm 1$ ,  $m_2 = \mp 1$  and  $\sigma^2 \rightarrow 0$ ), but not the distribution on it (i.e.  $\lambda \notin \{0.3, 0.7\}$ ). We also plot the density learned by a Gaussian VAE (see App. C.1) in blue to show this issue empirically. While this model assigns some probability outside of  $\{-1, 1\}$  due to limited capacity, the probabilities assigned around  $-1$  and  $1$  are far off from  $0.3$  and  $0.7$ , respectively; even after quantizing with the sign function, the VAE only assigns probability  $0.53$  to  $x = 1$ .

The underlying issue here is that  $\mathcal{M}$  is “too thin in  $\mathbb{R}^D$ ” (it has Lebesgue measure  $0$ ), and thus  $p(x)$  can “spike to infinity” at *every*  $x \in \mathcal{M}$ . If the dimensionalities were correctly matched this could not happen, as the requirement that  $p$  integrate to  $1$  would be violated. We highlight that this issue is not only a problem with data having  $d = 0$  dimensions, and can happen whenever  $d < D$ . The right panel of Fig. 2 shows

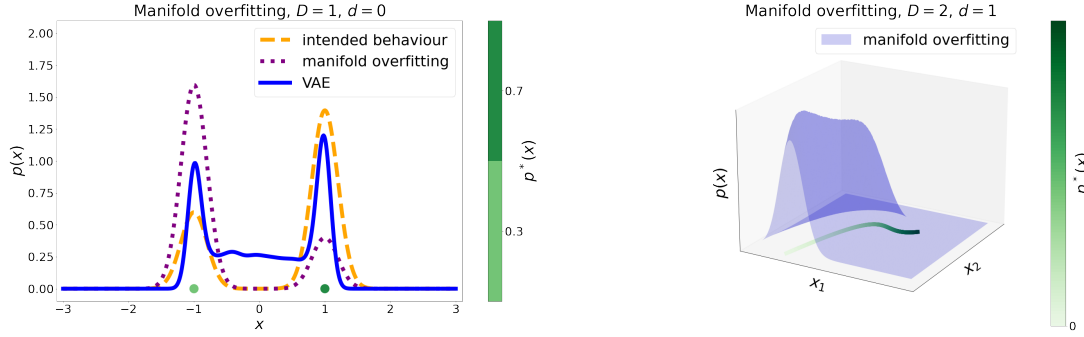


Figure 2: **Left panel:**  $\mathbb{P}^*$  (green);  $p_t(x) = 0.3 \cdot \mathcal{N}(x; -1, 1/t) + 0.7 \cdot \mathcal{N}(x; 1, 1/t)$  (orange, dashed) for  $t = 5$ , which converges weakly to  $\mathbb{P}^*$  as  $t \rightarrow \infty$ ; and  $p'_t(x) = 0.8 \cdot \mathcal{N}(x; -1, 1/t) + 0.2 \cdot \mathcal{N}(x; 1, 1/t)$  (purple, dotted) for  $t = 5$ , which converges weakly to  $\mathbb{P}^\dagger = 0.8\delta_{-1} + 0.2\delta_1$  while getting arbitrarily large likelihoods under  $\mathbb{P}^*$ , i.e.  $p'_t(x) \rightarrow \infty$  as  $t \rightarrow \infty$  for  $x \in \mathcal{M}$ ; Gaussian VAE density (blue, solid). **Right panel:** Analogous phenomenon with  $D = 2$  and  $d = 1$ , with the blue density “spiking” around  $\mathcal{M}$  in a manner unlike  $\mathbb{P}^*$  (green) while achieving large likelihoods.

another example of this phenomenon with  $d = 1$  and  $D = 2$ , where a distribution  $\mathbb{P}^*$  (green curve) is poorly approximated with a density  $p$  (blue surface) which nonetheless would achieve high likelihoods by “spiking around  $\mathcal{M}$ ”. Looking ahead to our experiments, the middle panel of Fig. 4 shows a 2-dimensional EBM suffering from this issue, spiking around the ground truth manifold on the left panel, but not correctly recovering the distribution on it. The intuition provided by these examples is that if a flexible  $D$ -dimensional density  $p$  is trained with maximum-likelihood when  $\mathbb{P}^*$  is supported on a low-dimensional manifold, it is possible to simultaneously achieve large likelihoods while being close to *any*  $\mathbb{P}^\dagger$ , rather than close to  $\mathbb{P}^*$ : this calls into question the validity of maximum-likelihood as a training objective in this setting.

### 3.2 The Manifold Overfitting Theorem

We now formalize the intuition developed so far. We assume some familiarity with measure theory (Billingsley, 2008) and with smooth (Lee, 2013) and Riemannian manifolds (Lee, 2018). Nonetheless, we provide a measure theory primer in App. A, where we informally review relevant concepts such as absolute continuity of measures ( $\ll$ ), weak convergence, and pushforward measures. We also use the concept of Riemannian measure (Pennec, 2006), which plays an analogous role on manifolds to that of the Lebesgue measure on Euclidean spaces. We briefly review Riemannian measures in App. B.1, and refer the reader to Dieudonné (1973) for a thorough treatment.<sup>2</sup> We begin by defining a useful condition on probability distributions for the following theorems, which captures the intuition of “continuously spreading mass all around  $\mathcal{M}$ ”.

**Definition 1 (Smoothness of Probability Measures):** Let  $\mathcal{M}$  be a finite-dimensional  $C^1$  manifold, and let  $\mathbb{P}$  be a probability measure on  $\mathcal{M}$ . Let  $\mathbf{g}$  be a Riemannian metric on  $\mathcal{M}$  and  $\mu_{\mathcal{M}}^{(\mathbf{g})}$  the corresponding Riemannian measure. We say that  $\mathbb{P}$  is *smooth* if  $\mathbb{P} \ll \mu_{\mathcal{M}}^{(\mathbf{g})}$  and it admits a continuous density  $p : \mathcal{M} \rightarrow \mathbb{R}_{>0}$  with respect to  $\mu_{\mathcal{M}}^{(\mathbf{g})}$ .

Note that smoothness of  $\mathbb{P}$  is independent of the choice of Riemannian metric  $\mathbf{g}$  (see App. B.1). We emphasize that this is a weak requirement, corresponding in the Euclidean case to  $\mathbb{P}$  admitting a continuous and positive density with respect to the Lebesgue measure. Denoting the Lebesgue measure on  $\mathbb{R}^D$  as  $\mu_D$ , we now state our first result.

**Theorem 1 (Manifold Overfitting):** Let  $\mathcal{M} \subset \mathbb{R}^D$  be an analytic  $d$ -dimensional embedded submanifold of  $\mathbb{R}^D$  with  $d < D$ , and  $\mathbb{P}^\dagger$  a smooth probability measure on  $\mathcal{M}$ . Then there exists a sequence of probability measures  $(\mathbb{P}_t)_{t=1}^\infty$  on  $\mathbb{R}^D$  such that:

<sup>2</sup>See especially Sec. 22 of Ch. 16. Note Riemannian measures are called Lebesgue measures in this reference.



1.  $\mathbb{P}_t \rightarrow \mathbb{P}^\dagger$  weakly as  $t \rightarrow \infty$ .
2. For every  $t \geq 1$ ,  $\mathbb{P}_t \ll \mu_D$  and  $\mathbb{P}_t$  admits a density  $p_t : \mathbb{R}^D \rightarrow \mathbb{R}_{>0}$  with respect to  $\mu_D$  such that:
  - (a)  $\lim_{t \rightarrow \infty} p_t(x) = \infty$  for every  $x \in \mathcal{M}$ .
  - (b)  $\lim_{t \rightarrow \infty} p_t(x) = 0$  for every  $x \notin \text{cl}(\mathcal{M})$ , where  $\text{cl}(\cdot)$  denotes closure in  $\mathbb{R}^D$ .

**Proof sketch:** We construct  $\mathbb{P}_t$  by convolving  $\mathbb{P}^\dagger$  with 0 mean,  $\sigma_t^2 I_D$  covariance Gaussian noise for a sequence  $(\sigma_t^2)_{t=1}^\infty$  satisfying  $\sigma_t^2 \rightarrow 0$  as  $t \rightarrow \infty$ , and then carefully verify that the stated properties of  $\mathbb{P}_t$  indeed hold. See App. B.2 for the full formal proof.

Informally, part 1 says that  $\mathbb{P}_t$  can get arbitrarily close to  $\mathbb{P}^\dagger$ , and part 2 says that this can be achieved with densities diverging to infinity on all  $\mathcal{M}$ . The relevance of this statement is that large likelihoods of a model do not imply it is adequately learning the target distribution  $\mathbb{P}^*$ , showing that maximum-likelihood is not a valid objective when data has low-dimensional manifold structure. Maximizing  $\frac{1}{N} \sum_{n=1}^N \log p(x_n)$ , or  $\mathbb{E}_{X \sim \mathbb{P}^*} [\log p(X)]$  in the nonparametric regime, over a  $D$ -dimensional density  $p$  need not recover  $\mathbb{P}^*$ : since  $\mathbb{P}^*$  is supported on  $\mathcal{M}$ , it follows by Theorem 1 that not only can the objective be made arbitrarily large, but that this can be done while recovering *any*  $\mathbb{P}^\dagger$ , which need not match  $\mathbb{P}^*$ . The failure to recover  $\mathbb{P}^*$  is caused by the density being able to take arbitrarily large values on all of  $\mathcal{M}$ , thus *overfitting to the manifold*. When  $p$  is a flexible density, as for many DGMs with universal approximation properties (Hornik, 1991; Koehler et al., 2021), manifold overfitting becomes a key deficiency of maximum-likelihood – which we fix in Sec. 4.

Note also that the proof of Theorem 1 applied to the specific case where  $\mathbb{P}^\dagger = \mathbb{P}^*$  formalizes the intuition that adding small amounts of Gaussian noise to  $\mathbb{P}^*$  results in highly peaked densities, suggesting that the resulting distribution, which denoising methods aim to estimate, might be empirically difficult to learn. More generally, even if there exists a ground truth  $D$ -dimensional density which allocates most of its mass around  $\mathcal{M}$ , this density will be highly peaked. In other words, even if Theorem 1 does not technically apply in this setting, it still provides useful intuition as manifold overfitting might still happen in practice. Indeed, we empirically confirm in Sec. 6 that even if  $\mathbb{P}^*$  is only “very close” to  $\mathcal{M}$ , manifold overfitting remains a problem.

**Differences from regular overfitting** Manifold overfitting is fundamentally different from regular overfitting, where the empirical distribution  $\hat{\mathbb{P}}_N = \frac{1}{N} \sum_{n=1}^N \delta_{x_n}$  is recovered.<sup>3</sup> First, regular overfitting requires more model capacity as  $N$  increases, since every new data point needs to be memorized. In contrast, once a model has enough capacity to learn  $\mathcal{M}$  and concentrate mass around it, no extra capacity is needed for manifold overfitting to occur regardless of how much data is observed. Second, a standard result from empirical process theory (Van Der Vaart & Wellner, 1996) states that  $\hat{\mathbb{P}}_N$  converges in distribution to  $\mathbb{P}^*$  as  $N \rightarrow \infty$ , whereas in manifold overfitting  $\mathbb{P}^*$  is not recovered even with infinite data, making the latter a more severe problem. Finally, an unseen test datapoint  $x_{N+1} \in \mathcal{M}$  will still be assigned very high likelihood – in line with the training data – under manifold overfitting, yet very low likelihood under regular overfitting. Manifold overfitting is thus undetectable when comparing train and test likelihoods.

**A note on divergences** Maximum-likelihood is often thought of as minimizing the KL divergence  $\mathbb{KL}(\mathbb{P}^* || \mathbb{P})$  over the model distribution  $\mathbb{P}$ . Naïvely one might believe that this contradicts the manifold overfitting theorem, but this is not the case. In order for  $\mathbb{KL}(\mathbb{P}^* || \mathbb{P}) < \infty$ , it is required that  $\mathbb{P}^* \ll \mathbb{P}$ , which does not happen when  $\mathbb{P}^*$  is a distribution on  $\mathcal{M}$  and  $\mathbb{P} \ll \mu_D$ . For example,  $\mathbb{KL}(\mathbb{P}^* || \mathbb{P}_t) = \infty$  for every  $t \geq 1$  even if  $\mathbb{E}_{X \sim \mathbb{P}^*} [\log p_t(X)]$  varies in  $t$ . In other words, minimizing the KL divergence is not equivalent to maximizing the likelihood in the setting of dimensionality mismatch, and the manifold overfitting theorem elucidates the effect of maximum-likelihood training in this setting. Similarly, other commonly considered divergences – such as  $f$ -divergences – cannot be meaningfully minimized. Arjovsky et al. (2017) propose using the Wasserstein distance as it is well-defined even in the presence of support mismatch, although we highlight once again that estimating and/or minimizing this distance is difficult in practice.

**Non-convergence of maximum-likelihood** The manifold overfitting theorem shows that any smooth distribution  $\mathbb{P}^\dagger$  on  $\mathcal{M}$  can be recovered through maximum-likelihood, even if it does not match  $\mathbb{P}^*$ . It does not,

<sup>3</sup>As an example of regular overfitting, the flexible model  $p(x) = \frac{1}{N} \sum_{n=1}^N \mathcal{N}(x; x_n, \sigma^2 I_D)$  with  $\sigma^2 \rightarrow 0$  recovers  $\hat{\mathbb{P}}_N$  and achieves arbitrarily large likelihoods, but may not generalize.

however, guarantee that *some*  $\mathbb{P}^\dagger$  will even be recovered. It is thus natural to ask whether it is possible to have a sequence of distributions achieving arbitrarily large likelihoods while not converging at all. The result below shows this to be true: in other words, training a  $D$ -dimensional model could result in maximum-likelihood not even converging.

**Corollary 1:** Let  $\mathcal{M} \subset \mathbb{R}^D$  be an analytic  $d$ -dimensional embedded submanifold of  $\mathbb{R}^D$  with more than a single element, and  $d < D$ . Then, there exists a sequence of probability measures  $(\mathbb{P}_t)_{t=1}^\infty$  on  $\mathbb{R}^D$  such that:

1.  $(\mathbb{P}_t)_{t=1}^\infty$  does not converge weakly.
2. For every  $t \geq 1$ ,  $\mathbb{P}_t \ll \mu_D$  and  $\mathbb{P}_t$  admits a density  $p_t : \mathbb{R}^D \rightarrow \mathbb{R}_{>0}$  with respect to  $\mu_D$  such that:
  - (a)  $\lim_{t \rightarrow \infty} p_t(x) = \infty$  for every  $x \in \mathcal{M}$ .
  - (b)  $\lim_{t \rightarrow \infty} p_t(x) = 0$  for every  $x \notin \text{cl}(\mathcal{M})$ .

**Proof:** Let  $\mathbb{P}^{\dagger 1}$  and  $\mathbb{P}^{\dagger 2}$  be two different smooth probability measures on  $\mathcal{M}$ , which exist since  $\mathcal{M}$  has more than a single element. Let  $(\mathbb{P}_t^1)_{t=1}^\infty$  and  $(\mathbb{P}_t^2)_{t=1}^\infty$  be the corresponding sequences from Theorem 1. The sequence  $(\mathbb{P}_t)_{t=1}^\infty$ , given by  $\mathbb{P}_t = \mathbb{P}_t^1$  if  $t$  is even and  $\mathbb{P}_t = \mathbb{P}_t^2$  otherwise, satisfies the above requirements.  $\square$

## 4 Fixing Manifold Overfitting

### 4.1 The Two-Step Correctness Theorem

The previous section motivates the development of likelihood-based methods which work correctly even in the presence of dimensionality mismatch. Intuitively, fixing the mismatch should be enough, which suggests (i) first reducing the dimension of the data to some  $d$ -dimensional representation, and then (ii) applying maximum-likelihood density estimation on the lower-dimensional dataset. The following theorem, where  $\mu_d$  denotes the Lebesgue measure on  $\mathbb{R}^d$ , confirms that this intuition is correct.

**Theorem 2 (Two-Step Correctness):** Let  $\mathcal{M} \subseteq \mathbb{R}^D$  be a  $C^1$   $d$ -dimensional embedded submanifold of  $\mathbb{R}^D$ , and let  $\mathbb{P}^*$  be a distribution on  $\mathcal{M}$ . Assume there exist measurable functions  $G : \mathbb{R}^d \rightarrow \mathbb{R}^D$  and  $g : \mathbb{R}^D \rightarrow \mathbb{R}^d$  such that  $G(g(x)) = x$ ,  $\mathbb{P}^*$ -almost surely. Then:

1.  $G_\#(g_\# \mathbb{P}^*) = \mathbb{P}^*$ , where  $h_\# \mathbb{P}$  denotes the pushforward of measure  $\mathbb{P}$  through the function  $h$ .
2. Moreover, if  $\mathbb{P}^*$  is smooth, and  $G$  and  $g$  are  $C^1$ , then:
  - (a)  $g_\# \mathbb{P}^* \ll \mu_d$ .
  - (b)  $G(g(x)) = x$  for every  $x \in \mathcal{M}$ , and the functions  $\tilde{g} : \mathcal{M} \rightarrow g(\mathcal{M})$  and  $\tilde{G} : g(\mathcal{M}) \rightarrow \mathcal{M}$  given by  $\tilde{g}(x) = g(x)$  and  $\tilde{G}(z) = G(z)$  are diffeomorphisms and inverses of each other.

**Proof:** See App. B.3.

We now discuss the implications of Theorem 2.

**Assumptions and correctness** The condition  $G(g(x)) = x$ ,  $\mathbb{P}^*$ -almost surely, is what one should expect to obtain during the dimensionality reduction step, for example through an autoencoder (AE) (Rumelhart et al., 1985) where  $\mathbb{E}_{X \sim \mathbb{P}^*}[\|G(g(X)) - X\|_2^2]$  is minimized over  $G$  and  $g$ , provided these have enough capacity and that population-level expectations can be minimized. We do highlight however that we allow for a much more general class of procedures than just autoencoders, nonetheless we still refer to  $g$  and  $G$  as the “encoder” and “decoder”, respectively. Part 1,  $G_\#(g_\# \mathbb{P}^*) = \mathbb{P}^*$ , justifies using a first step where  $g$  reduces the dimension of the data, and then having a second step attempting to learn the low-dimensional distribution  $g_\# \mathbb{P}^*$ : if a model  $\mathbb{P}_Z$  on  $\mathbb{R}^d$  matches the encoded data distribution, i.e.  $\mathbb{P}_Z = g_\# \mathbb{P}^*$ , it follows that  $G_\# \mathbb{P}_Z = \mathbb{P}^*$ . In other words, matching the distribution of encoded data and then decoding recovers the target distribution.

Part 2a guarantees that maximum-likelihood can be used to learn  $g_\# \mathbb{P}^*$ : note that if the model  $\mathbb{P}_Z$  is such that  $\mathbb{P}_Z \ll \mu_d$  with density  $p_Z = d\mathbb{P}_Z/d\mu_d$ , and  $g_\# \mathbb{P}^* \ll \mu_d$ , then both distributions are dominated by  $\mu_d$ .

Their KL divergence can then be expressed in terms of their respective densities:

$$\text{KL}(g_{\#}\mathbb{P}^*||\mathbb{P}_Z) = \int_{g(\mathcal{M})} p_Z^* \log \frac{p_Z^*}{p_Z} d\mu_d, \quad (1)$$

where  $p_Z^* = dg_{\#}\mathbb{P}^*/d\mu_d$  is the density of the encoded ground truth distribution. Assuming that  $|\int_{g(\mathcal{M})} p_Z^* \log p_Z^* d\mu_d| < \infty$ , the usual decomposition of KL divergence into expected log-likelihood and entropy applies, and it thus follows that maximum-likelihood over  $p_Z$  is once again equivalent to minimizing  $\text{KL}(g_{\#}\mathbb{P}^*||\mathbb{P}_Z)$  over  $\mathbb{P}_Z$ . In other words, learning the distribution of encoded data through maximum-likelihood with a flexible density approximator such as a VAE, AVB, NF, EBM, or ARM, and then decoding the result is a valid way of learning  $\mathbb{P}^*$  which avoids manifold overfitting.

**Density evaluation** Part 2b of the two-step correctness theorem bears some resemblance to injective NFs. However, note that the theorem does not imply  $G$  is injective, it only implies its restriction to  $g(\mathcal{M})$ ,  $G|_{g(\mathcal{M})}$ , is injective (and similarly for  $g$ ). Fig. 3 exemplifies how this can happen even if  $g$  and  $G$  are not injective. As with injective NFs, the density  $p_X$  of  $G_{\#}\mathbb{P}_Z$  (with respect to the Riemannian measure on  $\mathcal{M}$  corresponding to the Riemannian metric inherited from  $\mathbb{R}^D$ , which can equivalently be understood as the volume form on  $\mathcal{M}$ ) for a model  $\mathbb{P}_Z$  on  $g(\mathcal{M})$  is given by the injective change-of-variable formula:

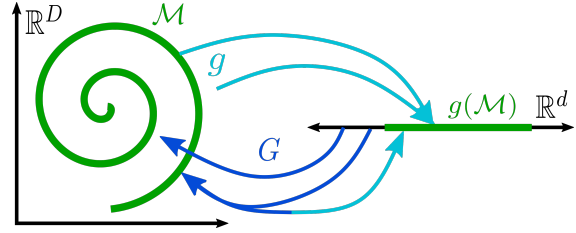


Figure 3: Illustration of how  $g$  and  $G$  can be bijective between  $\mathcal{M}$  (spiral) and  $g(\mathcal{M})$  (line segment) while not being fully bijective.

$$p_X(x) = p_Z(g(x)) \left| \det J_G^\top(g(x)) J_G(g(x)) \right|^{-1/2}, \quad (2)$$

for  $x \in \mathcal{M}$ , where  $J_G(g(x)) \in \mathbb{R}^{D \times d}$  is the Jacobian matrix of  $G$  evaluated at  $g(x)$ . Practically, this observation enables density evaluation of a trained two-step model, for example for OOD detection. Implementation-wise, we can use the approach proposed by Caterini et al. (2021) in the context of injective NFs, which uses forward-mode automatic differentiation to efficiently construct the Jacobian in (2). We highlight that, unlike Caterini et al. (2021), we do not train our models through (2). Furthermore, injectivity is not enforced in  $G$ , but rather achieved at optimality of the encoder/decoder pair, and only on  $g(\mathcal{M})$ .

## 4.2 Generalized Autoencoders

We now explain different approaches for obtaining  $G$  and  $g$ . As previously mentioned, a natural choice would be an AE minimizing  $\mathbb{E}_{X \sim \mathbb{P}^*} [\|G(g(X)) - X\|_2^2]$  over  $G$  and  $g$ . However, many other choices are also valid. We call a *generalized autoencoder* (GAE) any procedure in which both (i) low-dimensional representations  $z_n = g(x_n)$  are recovered for  $n = 1, \dots, N$ , and (ii) a function  $G$  is learned with the intention that  $G(z_n) = x_n$  for  $n = 1, \dots, N$ .

As alternatives to an AE, some DGMs can be used as GAEs, either because they directly provide  $G$  and  $g$  or can be easily modified to do so. These methods alone might obtain a  $G$  which correctly maps to  $\mathcal{M}$ , but might not be correctly recovering  $\mathbb{P}^*$ . From the manifold overfitting theorem, this is what we should expect from likelihood-based models, and we argue it is not unreasonable to expect from other models as well. For example, the high quality of samples generated from adversarial methods (Brock et al., 2019) suggests they are indeed learning  $\mathcal{M}$ , but issues such as mode collapse (Che et al., 2017) suggest they might not be recovering  $\mathbb{P}^*$  (Arbel et al., 2021). Among other options (Wang et al., 2020), we can use the following explicit DGMs as GAEs: (i) VAEs, taking  $G$  as the decoder mean, and  $g$  as the encoder mean, or (ii) AVB, using the encoder as  $g$  and the mean from the decoder as  $G$ . We can also use the following implicit DGMs as GAEs: (iii) Wasserstein autoencoders (WAEs) (Tolstikhin et al., 2018), again using the decoder as  $G$  and the encoder as  $g$ , (iv) bidirectional GANs (BiGANs) (Donahue et al., 2017; Dumoulin et al., 2017), taking  $G$  as the generator and  $g$  as the encoder, or (v) any GAN, by fixing  $G$  as the generator and then learning  $g$  by minimizing reconstruction error  $\mathbb{E}_{X \sim \mathbb{P}^*} [\|G(g(X)) - X\|_2^2]$ .

Note that explicit construction of  $g$  can be avoided as long as the representations  $\{z_n\}_{n=1}^N$  are learned, which could be achieved through non-amortized models (Gershman & Goodman, 2014; Kim et al., 2018), or with optimization-based GAN inversion methods (Xia et al., 2021). While Theorem 2 justifies using GAN inversion or related techniques as GAEs, these methods fall outside of our current scope, and we thus do not explore them empirically in Sec 6.

We summarize our two-step procedure class once again:

1. Learn  $G$  and  $\{z_n\}_{n=1}^N$  from  $\{x_n\}_{n=1}^N$  with a GAE.
2. Learn  $p_Z$  from  $\{z_n\}_{n=1}^N$  with a likelihood-based DGM.

The final model is then given by pushing  $p_Z$  forward through  $G$ . Any choice of GAE and likelihood-based DGM gives a valid instance of a two-step procedure.

## 5 Towards Unifying Deep Generative Models

**Making implicit models explicit** As noted above, some DGMs are themselves GAEs, including some implicit models for which density evaluation is not typically available, such as WAEs, BiGANs, and GANs. Ramesh & LeCun (2018) use (2) to train implicit models, but they do not train a second step DGM and thus have no mechanism to encourage trained models to satisfy the change-of-variable formula. Dieng et al. (2019) aim to provide GANs with density evaluation, but add  $D$ -dimensional Gaussian noise in order to achieve this, resulting in an adversarially-trained explicit model, rather than truly making an implicit model explicit. The two-step correctness theorem not only fixes manifold overfitting for explicit likelihood-based DGMs, but also enables density evaluation for these implicit models through (2) once a low-dimensional likelihood model has been trained on  $g(\mathcal{M})$ . We highlight the relevance of training the second step model  $p_Z$  for (2) to hold: even if  $G$  mapped some base distribution on  $\mathbb{R}^d$ , e.g. a Gaussian, to  $\mathbb{P}^*$ , it need not be injective to achieve this, and could map distinct inputs to the same point on  $\mathcal{M}$  (see Fig. 3). Such a  $G$  could be the result of training an implicit model, e.g. a GAN, which correctly learned its target distribution. Training  $g$ , and  $p_Z$  on  $g(\mathcal{M}) \subseteq \mathbb{R}^d$ , is still required to ensure  $G|_{g(\mathcal{M})}$  is injective and (2) can be applied, even if the end result of this additional training is that the target distribution remains properly learned.

**Two-step procedures** Several methods can be seen through the lens of our two-step approach, and can be interpreted as addressing manifold overfitting thanks to Theorem 2. Dai & Wipf (2019) use a two-step VAE, where both the GAE and DGM are taken as VAEs. Xiao et al. (2019) use a standard AE along with an NF. Brehmer & Cranmer (2020), and Kothari et al. (2021) use an AE as the GAE where  $G$  is an injective NF and  $g$  its left inverse and use an NF as the DGM. Ghosh et al. (2020) use an AE with added regularizers along with a Gaussian mixture model. Rombach et al. (2021) use a VAE along with a diffusion model (Ho et al., 2020) and obtain highly competitive empirical performance, which is justified by our theoretical results.

Other methods, while not exact instances, are philosophically aligned. Razavi et al. (2019) first obtain discrete low-dimensional representations of observed data and then train an ARM on these, which is similar to a discrete version of our own approach. Arbel et al. (2021) propose a model which they show is equivalent to pushing forward a low-dimensional EBM through  $G$ . The design of this model fits squarely into our framework, although a different training procedure is used.

The methods of Zhang et al. (2020b), Caterini et al. (2021), and Ross & Cresswell (2021) simultaneously optimize  $G$ ,  $g$ , and  $p_Z$  rather than using a two-step approach, combining in their loss a reconstruction term with a likelihood term as in (2). The validity of these methods however is not guaranteed by the two-step correctness theorem, and we believe a theoretical understanding of their objectives to be an interesting direction for future work.

## 6 Experiments

We now experimentally validate the advantages of our proposed two-step procedures across a variety of settings. We use the nomenclature A+B to refer to the two-step model with A as its GAE and B as its DGM.

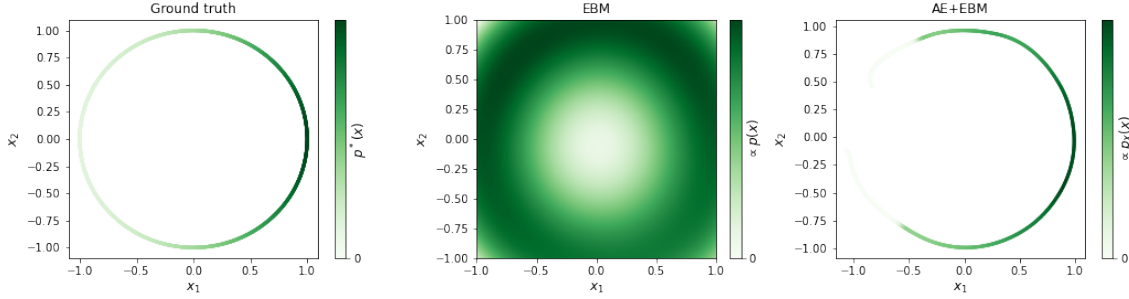


Figure 4: Results on simulated data: von Mises ground truth (left), EBM (middle), and AE+EBM (right).

All experimental details are provided in App. C. Our code<sup>4</sup> provides baseline implementations of all our considered GAEs and DGMs, which we hope will be useful to the community even outside of our proposed two-step methodology.

### 6.1 Simulated Data

We consider a von Mises distribution on the unit circle in Fig. 4. We learn this distribution both with an EBM and a two-step AE+EBM model. While the EBM indeed concentrates mass around the circle, it assigns higher density to an incorrect region of it (the top, rather than the right), corroborating manifold overfitting. The AE+EBM model not only learns the manifold more accurately, it also assigns higher likelihoods to the correct part of it.

### 6.2 Comparisons Against Maximum-Likelihood

We now show that our two-step methods empirically outperform maximum-likelihood training. Conveniently, some likelihood-based DGMs recover low-dimensional representations and hence are GAEs too, providing the opportunity to compare two-step training and maximum-likelihood training directly. In particular, AVB and VAEs both maximize a lower bound of the log-likelihood, so we can train a first model as a GAE, recover low-dimensional representations, and then train a second-step DGM. Any performance difference compared to maximum-likelihood is then due to the second-step DGM rather than the choice of GAE.

We show the results in Table 1 for MNIST, FMNIST (Xiao et al., 2017), SVHN (Netzer et al., 2011), and CIFAR-10 (Krizhevsky et al., 2009). We use Gaussian decoders with learnable scalar variance for both models, even for MNIST and FMNIST, as opposed to Bernoulli or other common choices (Loaiza-Ganem & Cunningham, 2019) in order to properly model the data as continuous and allow for manifold overfitting to happen. While ideally we would compare models based on log-likelihood, this is only sensible for models sharing the same dominating measure, here this is not the case as the single-step models are  $D$ -dimensional,

Table 1: FID scores (lower is better). Means  $\pm$  standard errors across 3 runs are shown. The superscript “+” indicates a larger model, and the subscript “ $\sigma$ ” indicates added Gaussian noise. Unreliable FID scores are highlighted in red (see text for description).

MODEL	MNIST	FMNIST	SVHN	CIFAR-10
AVB	22.1 $\pm$ 0.1	19.2 $\pm$ 0.2	8.5 $\pm$ 0.0	14.4 $\pm$ 0.3
AVB <sup>+</sup>	21.1 $\pm$ 0.2	17.9 $\pm$ 0.4	9.3 $\pm$ 0.4	14.0 $\pm$ 0.0
AVB <sub><math>\sigma</math></sub> <sup>+</sup>	21.0 $\pm$ 0.0	17.8 $\pm$ 0.2	9.1 $\pm$ 0.0	14.1 $\pm$ 0.1
AVB+ARM	10.1 $\pm$ 0.0	12.0 $\pm$ 0.0	5.4 $\pm$ 0.1	12.0 $\pm$ 0.1
AVB+AVB	13.2 $\pm$ 0.1	14.6 $\pm$ 0.3	6.6 $\pm$ 0.1	12.1 $\pm$ 0.2
AVB+EBM	11.0 $\pm$ 0.1	13.9 $\pm$ 0.6	6.2 $\pm$ 0.1	12.4 $\pm$ 0.3
AVB+NF	10.2 $\pm$ 0.1	12.0 $\pm$ 0.1	5.5 $\pm$ 0.0	12.0 $\pm$ 0.0
AVB+VAE	11.0 $\pm$ 0.1	12.4 $\pm$ 0.1	5.6 $\pm$ 0.3	12.2 $\pm$ 0.0
VAE	21.4 $\pm$ 0.1	17.8 $\pm$ 0.2	8.8 $\pm$ 0.1	14.8 $\pm$ 0.3
VAE <sup>+</sup>	21.1 $\pm$ 0.1	18.3 $\pm$ 0.2	8.9 $\pm$ 0.2	14.7 $\pm$ 0.1
VAE <sub><math>\sigma</math></sub> <sup>+</sup>	21.3 $\pm$ 0.1	18.0 $\pm$ 0.2	8.9 $\pm$ 0.1	14.6 $\pm$ 0.2
VAE+ARM	10.1 $\pm$ 0.1	12.0 $\pm$ 0.1	5.3 $\pm$ 0.1	12.0 $\pm$ 0.1
VAE+AVB	14.0 $\pm$ 0.2	14.5 $\pm$ 0.4	6.9 $\pm$ 0.0	12.3 $\pm$ 0.3
VAE+EBM	11.7 $\pm$ 0.2	13.6 $\pm$ 0.2	6.3 $\pm$ 0.1	12.9 $\pm$ 0.4
VAE+NF	10.1 $\pm$ 0.1	11.8 $\pm$ 0.0	5.3 $\pm$ 0.1	11.9 $\pm$ 0.1
ARM <sup>+</sup>	17.9 $\pm$ 2.4	11.8 $\pm$ 1.3	8.2 $\pm$ 0.5	11.4 $\pm$ 0.1
ARM <sub><math>\sigma</math></sub> <sup>+</sup>	5.3 $\pm$ 0.1	4.7 $\pm$ 0.1	6.4 $\pm$ 0.1	10.8 $\pm$ 0.1
AE+ARM	9.8 $\pm$ 0.1	12.2 $\pm$ 0.1	5.8 $\pm$ 0.0	12.2 $\pm$ 0.1
EBM <sup>+</sup>	8.7 $\pm$ 0.3	12.7 $\pm$ 0.3	7.7 $\pm$ 0.2	12.0 $\pm$ 0.1
EBM <sub><math>\sigma</math></sub> <sup>+</sup>	8.2 $\pm$ 0.4	12.1 $\pm$ 0.4	8.4 $\pm$ 0.5	11.9 $\pm$ 0.2
AE+EBM	11.3 $\pm$ 0.2	12.7 $\pm$ 0.1	6.7 $\pm$ 0.2	11.9 $\pm$ 0.1

<sup>4</sup>See supplementary material. We will publicly release our code upon publication.

while our two-step models are not. We thus use the FID score (Heusel et al., 2017) as a measure of how well models recover  $\mathbb{P}^*$ . Table 1 shows that our two-step procedures consistently outperform single-step maximum-likelihood training, even when adding Gaussian noise to the data, thus highlighting that manifold overfitting is still an empirical issue even when the ground truth distribution is  $D$ -dimensional but highly peaked around a manifold. We also note that some of the baseline models are significantly larger, e.g. the VAE<sup>+</sup> on MNIST has approximately 824k parameters, while the VAE model has 412k, and the VAE+EBM only 416k. The parameter efficiency of two-step models highlights that our empirical gains are not due to increasing model capacity but rather from addressing manifold overfitting.

Table 1 also shows comparisons between single and two-step models for ARMs and EBMs, which unlike AVB and VAEs, are not GAEs themselves; we thus use an AE as the GAE for these comparisons. Although FID scores did not consistently improve for these two-step models over their corresponding single-step baselines, we found the visual quality of samples was significantly better for almost all two-step models, as demonstrated in the first two columns of Fig. 5, and by the additional samples shown in App. D.1. We thus highlight with red the corresponding FID scores as unreliable in Table 1. We believe these failures modes of the FID metric itself, wherein the scores do not correlate with visual quality, emphasize the importance of further research on sample-based scalar evaluation metrics for DGMs (Borji, 2022), although developing such metrics falls outside our scope. We also show comparisons using precision and recall (Kynkäänniemi et al., 2019) in App. D.3, and observe that two-step models still outperform single-step ones.

We also point out that one-step EBMs exhibited training difficulties consistent with maximum-likelihood non-convergence (App. D.2). Meanwhile, Langevin dynamics (Welling & Teh, 2011) for AE+EBM exhibit better and faster convergence, yielding good samples even when not initialized from the training buffer (see Fig. 12 in App. D.2), and AE+ARM speeds up sampling over the baseline ARM by a factor of  $\mathcal{O}(D/d)$ , in both cases because there are fewer coordinates in the sample space. Of the 44 two-step models shown in Table 1, only one (AE+EBM on MNIST) did not visually outperform its single-step counterpart (App. D.1), empirically corroborating our theoretical findings.

Finally, we have omitted some comparisons verified in prior work: Dai & Wipf (2019) show VAE+VAE outperforms VAE, and Xiao et al. (2019) that AE+NF outperforms NF.

### 6.3 OOD Detection with Implicit Models

Having verified that, as predicted by Theorem 2, two-step models outperform maximum-likelihood training, we now turn our attention to the other consequence of this theorem, namely endowing implicit models with density evaluation after training a second-step DGM. We demonstrate that our approach advances fully-unsupervised likelihood-based out-of-distribution detection. Nalisnick et al. (2019) discovered the counter-intuitive phenomenon that likelihood-based DGMs sometimes assign higher likelihoods to OOD data than to in-distribution data. In particular, they found models trained on FMNIST and CIFAR-10 assigned higher likelihoods to MNIST and SVHN, respectively. While there has been a significant amount of research trying to remedy and explain this situation (Choi et al., 2018; Ren et al., 2019; Le Lan & Dinh, 2020; Caterini & Loaiza-Ganem, 2021), there is little work achieving good OOD performance using only likelihoods (Caterini et al., 2021).

Table 2: OOD classification accuracy as a percentage (higher is better). Means  $\pm$  standard errors across 3 runs are shown. Arrows point from in-distribution to OOD data.

MODEL	FMNIST $\rightarrow$ MNIST	CIFAR-10 $\rightarrow$ SVHN
ARM <sup>+</sup>	9.9 $\pm$ 0.6	15.5 $\pm$ 0.0
BiGAN+ARM	81.9 $\pm$ 1.4	38.0 $\pm$ 0.2
WAE+ARM	69.8 $\pm$ 13.9	40.1 $\pm$ 0.2
AVB <sup>+</sup>	96.0 $\pm$ 0.5	23.4 $\pm$ 0.1
BiGAN+AVB	59.5 $\pm$ 3.1	36.4 $\pm$ 2.0
WAE+AVB	90.7 $\pm$ 0.7	43.5 $\pm$ 1.9
EBM <sup>+</sup>	32.5 $\pm$ 1.1	46.4 $\pm$ 3.1
BiGAN+EBM	51.2 $\pm$ 0.2	48.8 $\pm$ 0.1
WAE+EBM	57.2 $\pm$ 1.3	49.3 $\pm$ 0.2
NF <sup>+</sup>	36.4 $\pm$ 0.2	18.6 $\pm$ 0.3
BiGAN+NF	84.2 $\pm$ 1.0	40.1 $\pm$ 0.2
WAE+NF	95.4 $\pm$ 1.6	46.1 $\pm$ 1.0
VAE <sup>+</sup>	96.1 $\pm$ 0.1	23.8 $\pm$ 0.2
BiGAN+VAE	59.7 $\pm$ 0.2	38.1 $\pm$ 0.1
WAE+VAE	92.5 $\pm$ 2.7	41.4 $\pm$ 0.2

We train several two-step models where the GAE is either a BiGAN or a WAE, which do not by themselves allow for likelihood evaluation, and then use the resulting log-likelihoods (or lower bounds/negative energy





Figure 5: Uncurated samples from single-step models (**first row**, showing  $\text{ARM}_\sigma^+$ ,  $\text{EBM}^+$ ,  $\text{AVB}_\sigma^+$ , and VAE) and their respective two-step counterparts (**second row**, showing AE+ARM, AE+EBM, AVB+NF, and VAE+AVB), for MNIST (**first column**), FMNIST (**second column**), SVHN (**third column**), and CIFAR-10 (**fourth column**).

functions) for OOD detection. Two-step models allow us to use either the high-dimensional  $\log p_X$  from (2) or low-dimensional  $\log p_Z$  as heuristics for this task. We conjecture that the latter is more reliable, since (i) the base measure is always  $\mu_d$ , and (ii) the encoder-decoder is unlikely to exactly satisfy the conditions of Theorem 2. Hence, we use  $\log p_Z$  here, and show results for  $\log p_X$  in App. D.4.

Table 2 shows the classification accuracy of a decision stump given only the log-likelihood; we show some corresponding histograms in App. D.4. The stump is forced to assign large likelihoods as in-distribution, so that accuracies below 50% indicate it incorrectly assigned higher likelihoods to OOD data. We correct the classification accuracy to account for datasets of different size (details in App. D.4), resulting in an easily interpretable metric which can be understood as the expected classification accuracy if two same-sized samples of in-distribution and OOD data were compared. Not only did we enable implicit models to perform OOD detection, but we also outperformed likelihood-based single-step models in this setting. To the best of our knowledge, no other model achieves nearly 50% accuracy on CIFAR-10→SVHN using *only* likelihoods. Although admittedly the problem is not yet solved, we have certainly made progress on a challenging task for fully-unsupervised methods. For completeness, we show samples from these models in App. D.1 and FID scores in App. D.3. Implicit models see less improvement in FID from adding a second-step DGM than explicit models, suggesting that manifold overfitting is a less dire problem for implicit models. Nonetheless, we do observe some improvements, particularly for BiGANs, hinting that our two-step methodology not only endows these models with density evaluation, but that it can also improve their generative performance. We further show in App. D.4 that OOD improvements obtained by two-step models apply to explicit models too.

Interestingly, whereas the VAEs used in Nalisnick et al. (2019) have Bernoulli likelihoods, we find that our single-step likelihood-based Gaussian-decoder VAE and AVB models perform quite well on distinguishing FMNIST from MNIST, yet still fail on the CIFAR-10 task. Studying this is of future interest but is outside the scope of this work.

## 7 Conclusions, Scope, and Limitations

In this paper we diagnosed manifold overfitting, a fundamental problem of maximum-likelihood training with flexible densities when the data lives in a low-dimensional manifold. We proposed to fix manifold overfitting



with a class of two-step procedures which remedy the issue, theoretically justify a large group of existing methods, and endow implicit models with density evaluation after training a low-dimensional likelihood-based DGM on encoded data.

Our two-step correctness theorem remains nonetheless a nonparametric result. In practice, the reconstruction error will be positive, i.e.  $\mathbb{E}_{X \sim \mathbb{P}^*} [\|G(g(X)) - X\|_2^2] > 0$ . Note that this can happen even when assuming infinite capacity, as  $\mathcal{M}$  needs to be diffeomorphic to  $g(\mathcal{M})$  for some  $C^1$  function  $g : \mathbb{R}^D \rightarrow \mathbb{R}^d$  for the reconstruction error to be 0. We leave a study of learnable topologies of  $\mathcal{M}$  for future work. The density in (2) might then not be valid, either if the reconstruction error is positive, or if  $p_Z$  assigns positive probability outside of  $g(\mathcal{M})$ . However, we note that our approach at least provides a mechanism to encourage our trained encoder-decoder pair to invert each other, suggesting that (2) might not be too far off. We also believe that a finite-sample extension of our result, while challenging, would be a relevant direction for future work. We hope our work will encourage follow-up research exploring different ways of addressing manifold overfitting, or its interaction with the score-matching objective.

Finally, we treated  $d$  as a hyperparameter, but in practice  $d$  is unknown and improvements can likely be had by estimating it (Levina & Bickel, 2004). Still, we observed significant empirical improvements across a variety of tasks and datasets, demonstrating that manifold overfitting is not just a theoretical issue in DGMs, and that two-step methods are an important class of procedures to deal with it.

## Broader Impact Statement

Due to the theoretical and methodological nature of our contributions, we do not foresee our work having any negative societal consequences. We advance the understanding of the distributions learned by our DGMs, which can be important for fairness and reducing biases (Humayun et al., 2022), particularly in the context of natural image and facial generation.

## References

- Guillaume Alain and Yoshua Bengio. What regularized auto-encoders learn from the data-generating distribution. *The Journal of Machine Learning Research*, 15(1):3563–3593, 2014.
- Michael Arbel, Liang Zhou, and Arthur Gretton. Generalized energy based models. *ICLR*, 2021.
- Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein generative adversarial networks. In *International conference on machine learning*, pp. 214–223. PMLR, 2017.
- Sanjeev Arora, Rong Ge, Yingyu Liang, Tengyu Ma, and Yi Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *International Conference on Machine Learning*, pp. 224–232. PMLR, 2017.
- Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, pp. 573–582. PMLR, 2019.
- Jens Behrmann, Paul Vicol, Kuan-Chieh Wang, Roger Grosse, and Jörn-Henrik Jacobsen. Understanding and mitigating exploding inverses in invertible neural networks. In *International Conference on Artificial Intelligence and Statistics*, pp. 1792–1800. PMLR, 2021.
- Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- Patrick Billingsley. *Probability and measure*. John Wiley & Sons, 2008.
- S Bond-Taylor, A Leach, Y Long, and CG Willcocks. Deep generative modelling: A comparative review of VAEs, GANs, normalizing flows, energy-based and autoregressive models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- Ali Borji. Pros and cons of gan evaluation measures: New developments. *Computer Vision and Image Understanding*, 215:103329, 2022.

- Johann Brehmer and Kyle Cranmer. Flows for simultaneous manifold learning and density estimation. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. *ICLR*, 2019.
- Anthony L Caterini and Gabriel Loaiza-Ganem. Entropic Issues in Likelihood-Based OOD Detection. *arXiv preprint arXiv:2109.10794*, 2021.
- Anthony L Caterini, Gabriel Loaiza-Ganem, Geoff Pleiss, and John P Cunningham. Rectangular flows for manifold learning. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Minwoo Chae, Dongha Kim, Yongdai Kim, and Lizhen Lin. A likelihood approach to nonparametric estimation of a singular distribution using deep generative models. *arXiv preprint arXiv:2105.04046*, 2021.
- Tong Che, Yanran Li, Athul Paul Jacob, Yoshua Bengio, and Wenjie Li. Mode regularized generative adversarial networks. *ICLR*, 2017.
- Ricky T. Q. Chen, Jens Behrmann, David K Duvenaud, and Joern-Henrik Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Hyunsun Choi, Eric Jang, and Alexander A Alemi. WAIC, but why? Generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018.
- Casey Chu, Kentaro Minami, and Kenji Fukumizu. Smoothness and stability in GANs. *ICLR*, 2020.
- Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *ICLR*, 2016.
- Rob Cornish, Anthony Caterini, George Deligiannidis, and Arnaud Doucet. Relaxing bijectivity constraints with continuously indexed normalising flows. In *International Conference on Machine Learning*, pp. 2133–2143. PMLR, 2020.
- Edmond Cunningham and Madalina Fiterau. A change of variables method for rectangular matrix-vector products. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, volume 130. PMLR, 2021.
- Bin Dai and David Wipf. Diagnosing and enhancing VAE models. *ICLR*, 2019.
- Adji B Dieng, Francisco JR Ruiz, David M Blei, and Michalis K Titsias. Prescribed generative adversarial networks. *arXiv preprint arXiv:1910.04302*, 2019.
- Jean Dieudonné. *Treatise on Analysis: Volume III*. Academic Press, 1973.
- Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. *ICLR*, 2017.
- Jeff Donahue, Philipp Krähenbühl, and Trevor Darrell. Adversarial feature learning. *ICLR*, 2017.
- Yilun Du and Igor Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 32:3608–3618, 2019.
- Vincent Dumoulin, Ishmael Belghazi, Ben Poole, Olivier Mastropietro, Alex Lamb, Martin Arjovsky, and Aaron Courville. Adversarially learned inference. *ICLR*, 2017.
- Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural Spline Flows. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Mevlana C Gemici, Danilo Rezende, and Shakir Mohamed. Normalizing flows on riemannian manifolds. *arXiv preprint arXiv:1611.02304*, 2016.
- Samuel Gershman and Noah Goodman. Amortized inference in probabilistic reasoning. In *Proceedings of the annual meeting of the cognitive science society*, volume 36, 2014.

- Partha Ghosh, Mehdi SM Sajjadi, Antonio Vergari, Michael Black, and Bernhard Schölkopf. From variational to deterministic autoencoders. *ICLR*, 2020.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- Alfred Gray. The volume of a small geodesic ball of a riemannian manifold. *The Michigan Mathematical Journal*, 20(4):329–344, 1974.
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron C Courville. Improved training of wasserstein gans. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. GANs trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, volume 33, pp. 6840–6851, 2020.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- Christian Horvat and Jean-Pascal Pfister. Denoising normalizing flow. *Advances in Neural Information Processing Systems*, 34, 2021a.
- Christian Horvat and Jean-Pascal Pfister. Density estimation on low-dimensional manifolds: an inflation-deflation approach. 2021b.
- Ahmed Imtiaz Humayun, Randall Balestrierio, and Richard Baraniuk. Magnet: Uniform sampling from deep generative network manifolds without retraining. *ICLR*, 2022.
- Aapo Hyvärinen. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005.
- Yoon Kim, Sam Wiseman, Andrew Miller, David Sontag, and Alexander Rush. Semi-amortized variational autoencoders. In *International Conference on Machine Learning*, pp. 2678–2687. PMLR, 2018.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- Diederik P Kingma and Prafulla Dhariwal. Glow: Generative Flow with Invertible  $1 \times 1$  Convolutions. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Diederik P Kingma and Max Welling. Auto-encoding Variational Bayes. *ICLR*, 2014.
- Diederik P Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. In *Advances in neural information processing systems*, volume 34, 2021.
- Frederic Koehler, Viraj Mehta, and Andrej Risteski. Representational aspects of depth and conditioning in normalizing flows. In *International Conference on Machine Learning*, pp. 5628–5636. PMLR, 2021.
- Konik Kothari, AmirEhsan Khorashadizadeh, Maarten de Hoop, and Ivan Dokmanić. Trumpets: Injective flows for inference and inverse problems. In *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161, pp. 1269–1278, 2021.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila. Improved precision and recall metric for assessing generative models. *Advances in Neural Information Processing Systems*, 32, 2019.

- Charline Le Lan and Laurent Dinh. Perfect density models cannot guarantee anomaly detection. In *"I Can't Believe It's Not Better!" NeurIPS 2020 workshop*, 2020.
- Y LeCun. The MNIST database of handwritten digits, 1998. URL <http://yann.lecun.com/exdb/mnist>, 1998.
- John M Lee. Smooth manifolds. In *Introduction to Smooth Manifolds*, pp. 1–31. Springer, 2013.
- John M Lee. *Introduction to Riemannian manifolds*. Springer, 2018.
- Erich L Lehmann and George Casella. *Theory of point estimation*. Springer Science & Business Media, 2006.
- Elizaveta Levina and Peter Bickel. Maximum likelihood estimation of intrinsic dimension. *Advances in neural information processing systems*, 17, 2004.
- Gabriel Loaiza-Ganem and John P Cunningham. The continuous bernoulli: fixing a pervasive error in variational autoencoders. *Advances in Neural Information Processing Systems*, 32:13287–13297, 2019.
- Emile Mathieu and Maximilian Nickel. Riemannian continuous normalizing flows. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Pierre-Alexandre Mattei and Jes Frellsen. Leveraging the exact likelihood of deep latent variable models. *Advances in Neural Information Processing Systems*, 31, 2018.
- Lars Mescheder, Sebastian Nowozin, and Andreas Geiger. Adversarial Variational Bayes: Unifying variational autoencoders and generative adversarial networks. In *International Conference on Machine Learning*, pp. 2391–2400. PMLR, 2017.
- Shakir Mohamed and Balaji Lakshminarayanan. Learning in implicit generative models. *arXiv preprint arXiv:1610.03483*, 2016.
- Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. 2020.
- Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don't know? *ICLR*, 2019.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- Sebastian Nowozin, Botond Cseke, and Ryota Tomioka. F-GAN: Training generative neural samplers using variational divergence minimization. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, pp. 271–279, 2016.
- Xavier Pennec. Intrinsic statistics on riemannian manifolds: Basic tools for geometric measurements. *Journal of Mathematical Imaging and Vision*, 25(1):127–154, 2006.
- Phillip Pope, Chen Zhu, Ahmed Abdelkader, Micah Goldblum, and Tom Goldstein. The intrinsic dimension of images and its impact on learning. *ICLR*, 2021.
- Prajit Ramachandran, Barret Zoph, and Quoc V Le. Searching for activation functions. *arXiv preprint arXiv:1710.05941*, 2017.
- Aditya Ramesh and Yann LeCun. Backpropagation for implicit spectral densities. *arXiv preprint arXiv:1806.00499*, 2018.
- Ali Razavi, Aaron van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with VQ-VAE-2. In *Advances in neural information processing systems*, pp. 14866–14876, 2019.
- Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Advances in Neural Information Processing Systems*, volume 32, pp. 14707–14718, 2019.

- Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pp. 1278–1286. PMLR, 2014.
- Danilo Jimenez Rezende, George Papamakarios, Sébastien Racaniere, Michael Albergo, Gurtej Kanwar, Phiala Shanahan, and Kyle Cranmer. Normalizing flows on tori and spheres. In *International Conference on Machine Learning*, pp. 8083–8092. PMLR, 2020.
- Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. *arXiv preprint arXiv:2112.10752*, 2021.
- Brendan Leigh Ross and Jesse C Cresswell. Tractable density estimation on learned manifolds with conformal embedding flows. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning internal representations by error propagation. Technical report, California Univ San Diego La Jolla Inst for Cognitive Science, 1985.
- Mehdi SM Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. *Advances in Neural Information Processing Systems*, 31, 2018.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In *Proceedings of the 33rd Annual Conference on Neural Information Processing Systems*, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *ICLR*, 2021.
- Lucas Theis and Matthias Bethge. Generative image modeling using spatial LSTMs. *Advances in Neural Information Processing Systems*, 28:1927–1935, 2015.
- Ilya Tolstikhin, Olivier Bousquet, Sylvain Gelly, and Bernhard Schoelkopf. Wasserstein auto-encoders. *ICLR*, 2018.
- Benigno Uria, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems 26 (NIPS 26)*, 2013.
- Aad W Van Der Vaart and Jon Wellner. *Weak convergence and empirical processes: with applications to statistics*. Springer Science & Business Media, 1996.
- Aaron Van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. In *International conference on machine learning*, pp. 1747–1756. PMLR, 2016.
- Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011.
- Pascal Vincent, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. Extracting and composing robust features with denoising autoencoders. In *Proceedings of the 25th international conference on Machine learning*, pp. 1096–1103, 2008.
- Wei Wang, Yan Huang, Yizhou Wang, and Liang Wang. Generalized autoencoder: A neural network framework for dimensionality reduction. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, pp. 490–497, 2014.
- Yingfan Wang, Haiyang Huang, Cynthia Rudin, and Yaron Shaposhnik. Understanding how dimension reduction tools work: An empirical approach to deciphering t-SNE, UMAP, TriMAP, and PaCMAP for data visualization. *arXiv preprint arXiv:2012.04456*, 2020.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688. Citeseer, 2011.
- Weihao Xia, Yulun Zhang, Yujiu Yang, Jing-Hao Xue, Bolei Zhou, and Ming-Hsuan Yang. GAN inversion: A survey. *arXiv preprint arXiv:2101.05278*, 2021.

- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- Zhisheng Xiao, Qing Yan, and Yali Amit. Generative latent flow. *arXiv preprint arXiv:1905.10485*, 2019.
- Mingtian Zhang, Peter Hayes, Thomas Bird, Raza Habib, and David Barber. Spread divergence. In *International Conference on Machine Learning*, pp. 11106–11116. PMLR, 2020a.
- Zijun Zhang, Ruixiang Zhang, Zongpeng Li, Yoshua Bengio, and Liam Paull. Perceptual generative autoencoders. In *International Conference on Machine Learning*, pp. 11298–11306. PMLR, 2020b.

## A Informal Measure Theory Primer

Before stating Theorems 1 and 2, and studying their implications, we provide a brief tutorial on some aspects of measure theory that are relevant to follow our discussion. This review is not meant to be comprehensive, and we prioritize intuition over formalism. Readers interested in the topic may consult textbooks such as Billingsley (2008).

### A.1 Probability Measures

Let us first motivate the need for measure theory in the first place and consider the question: *what is a density?* Intuitively, the density  $p_X$  of a random variable  $X$  is a function having the property that integrating  $p_X$  over any set  $A$  gives back the probability that  $X \in A$ . This density characterizes the distribution of  $X$ , in that it can be used to answer any probabilistic question about  $X$ . It is common knowledge that discrete random variables are not specified through a density, but rather a probability mass function. Similarly, in our setting, where  $X$  might always take values in  $\mathcal{M}$ , such a density will not exist. To see this, consider the case where  $A = \mathcal{M}$ , so that the integral of  $p_X$  over  $\mathcal{M}$  would have to be 1, which cannot happen since  $\mathcal{M}$  has volume 0 in  $\mathbb{R}^D$  (or more formally, Lebesgue measure 0). Measure theory provides the tools necessary to properly specify *any distribution*, subsuming as special cases probability mass functions, densities of continuous random variables, and distributions on manifolds.

A measure  $\mu$  on  $\mathbb{R}^D$  is a function mapping subsets  $A \subseteq \mathbb{R}^D$  to  $\mathbb{R}_{\geq 0}$ , obeying the following properties: (i)  $\mu(A) \geq 0$  for every  $A$ , (ii)  $\mu(\emptyset) = 0$ , where  $\emptyset$  denotes the empty set, and (iii)  $\mu(\cup_{k=1}^{\infty} A_k) = \sum_{k=1}^{\infty} \mu(A_k)$  for any sequence of pairwise disjoint sets  $A_1, A_2, \dots$  (i.e.  $A_i \cap A_j = \emptyset$  whenever  $i \neq j$ ). Note that most measures of interest are only defined over a large class of subsets of  $\mathbb{R}^D$  (called  $\sigma$ -algebras, the most notable one being the Borel  $\sigma$ -algebra) rather than for every possible subset due to technical reasons, but we omit details in the interest of better conveying intuition. A measure is called a *probability* measure if it also satisfies  $\mu(\mathbb{R}^D) = 1$ . To any random variable  $X$  corresponds a probability measure  $\mu_X$ , having the property that  $\mu_X(A)$  is the probability that  $X \in A$  for any  $A$ . Analogously to probability mass functions or densities of continuous random variables,  $\mu_X$  allows us to answer any probabilistic question about  $X$ . The probability measure  $\mu_X$  is often called the *distribution* or *law* of  $X$ . Throughout our paper,  $\mathbb{P}^*$  is the distribution from which we observe data.

Let us consider two examples to show how probability mass functions and densities of continuous random variables are really just specifying distributions. Given  $a_1, \dots, a_K \in \mathbb{R}^D$ , consider the probability mass function of a random variable  $X$  given by  $p_X(x) = 1/K$  for  $x = a_1, a_2, \dots, a_K$  and 0 otherwise. This probability mass function is simply specifying the distribution  $\mu_X(A) = 1/K \cdot \sum_{k=1}^K \mathbb{1}(a_k \in A)$ , where  $\mathbb{1}(\cdot \in A)$  denotes the indicator function for  $A$ , i.e.  $\mathbb{1}(a \in A)$  is 1 if  $a \in A$ , and 0 otherwise. Now consider a standard Gaussian random variable  $X$  in  $\mathbb{R}^D$  with density  $p_X(x) = \mathcal{N}(x; 0, I_D)$ . Similarly to how the probability mass function from the previous example characterized a distribution, this density does so as well through  $\mu_X(A) = \int_A \mathcal{N}(x; 0, I_D) dx$ . We will see in the next section how these ideas can be extended to distributions on manifolds.

The concept of integrating a function  $h : \mathbb{R}^D \rightarrow \mathbb{R}$  with respect to a measure  $\mu$  on  $\mathbb{R}^D$  is fundamental in measure theory, and can be thought of as “weighting the inputs of  $h$  according to  $\mu$ ”. In the case of the Lebesgue measure  $\mu_D$  (which assigns to subsets  $A$  of  $\mathbb{R}^D$  their “volume”  $\mu_D(A)$ ), integration extends the concept of Riemann integrals commonly taught in calculus courses, and in the case of random variables integration defines expectations, i.e.  $\mathbb{E}_{X \sim \mu_X}[h(X)] = \int h d\mu_X$ . In the next section we will talk about the interplay between integration and densities.

### A.2 Absolute Continuity

So far we have seen that probability measures allow us to talk about distributions in full generality, and that probability mass functions and densities of continuous random variables can be used to specify probability measures. A distribution on a manifold  $\mathcal{M}$  embedded in  $\mathbb{R}^D$  can simply be thought of as a probability measure  $\mu$  such that  $\mu(\mathcal{M}) = 1$ . We would like to define densities on manifolds in an analogous way to probability mass functions and densities of continuous random variables, in such a way that they allow us to



characterize distributions on the manifold. Absolute continuity of measures is a concept that allows us to formalize the concept of *density with respect to a dominating measure*, and encompasses probability mass functions, densities of continuous random variables, and also allows us to define densities on manifolds. We will see that our intuitive definition of a density as a function which, when integrated over a set gives back its probability, is in fact correct, just as long as we specify the measure we integrate with respect to.

Given two measures  $\mu$  and  $\nu$ , we say that  $\mu$  is *absolutely continuous* with respect to  $\nu$  if for every  $A$  such that  $\nu(A) = 0$ , it also holds that  $\mu(A) = 0$ . If  $\mu$  is absolutely continuous with respect to  $\nu$ , we also say that  $\nu$  *dominates*  $\mu$ , and denote this property as  $\mu \ll \nu$ . The Radon-Nikodym theorem states that, under some mild assumptions on  $\mu$  and  $\nu$  which hold for all the measures considered in this paper,  $\mu \ll \nu$  implies the existence of a function  $h$  such that  $\mu(A) = \int_A h d\nu$  for every  $A$ . This result provides the means to formally define densities:  $h$  is called the *density* or *Radon-Nikodym derivative* of  $\mu$  with respect to  $\nu$ , and is often written as  $d\mu/d\nu$ .

Before explaining how this machinery allows us to talk about densities on manifolds, we first continue our examples to show that probability mass functions and densities of continuous random variables are Radon-Nikodym derivatives with respect to appropriate measures. Let us reconsider the example where  $p_X(x) = 1/K$  for  $x = a_1, a_2, \dots, a_K$  and 0 otherwise, and  $\mu_X(A) = 1/K \cdot \sum_{k=1}^K \mathbb{1}(a_k \in A)$ . Consider the measure  $\nu(A) = \sum_{k=1}^K \mathbb{1}(a_k \in A)$ , which essentially just counts the number of  $a_k$ s in  $A$ . Clearly  $\mu_X \ll \nu$ , and so it follows that  $\mu_X$  admits a density with respect to  $\nu$ . This density turns out to be  $p_X$ , since  $\mu_X(A) = \int_A p_X d\nu$ . In other words, the probability mass function  $p_X$  can be thought of as a Radon-Nikodym derivative, i.e.  $p_X = d\mu_X/d\nu$ . Let us now go back to the continuous density example where  $p_X(x) = \mathcal{N}(x; 0, I_D)$  and  $\mu_X$  is given by the Riemann integral  $\mu_X(A) = \int_A \mathcal{N}(x; 0, I_D) dx$ . In this case,  $\nu = \mu_D$ , and since the Lebesgue integral extends the Riemann integral, it follows that  $\mu_X(A) = \int_A p_X d\mu_D$ , so that the density  $p_X$  is actually also a density in the formal sense of being a Radon-Nikodym derivative, so that  $p_X = d\mu_X/d\mu_D$ . We can thus see that the formal concept of density or Radon-Nikodym derivative generalizes both probability mass functions and densities of continuous random variables as we usually think of them, allowing to specify distributions in a general way.

The concept of Radon-Nikodym derivative also allows us to obtain densities on manifolds, the only missing ingredient being a dominating measure on the manifold. Riemannian measures (App. B.1) play this role on manifolds, in the same way that the Lebesgue measure plays the usual role of dominating measure to define densities of continuous random variables on  $\mathbb{R}^D$ .

### A.3 Weak Convergence

A key point in Theorem 1 is weak convergence of the sequence of probability measures  $(\mathbb{P}_t)_{t=1}^\infty$  to  $\mathbb{P}^\dagger$ . The intuitive interpretation that this statement simply means that “ $\mathbb{P}_t$  converges to  $\mathbb{P}^\dagger$ ” is correct, although formally defining convergence of a sequence of measures is still required. Weak convergence provides such a definition, and  $\mathbb{P}_t$  is said to *converge weakly* to  $\mathbb{P}^\dagger$  if the sequence of scalars  $\mathbb{P}_t(A)$  converges to  $\mathbb{P}^\dagger(A)$  for every  $A$  satisfying a technical condition (for intuitive purposes, one can think of this property as holding for every  $A$ ). In this sense weak convergence is a very natural way of defining convergence of measures: in the limit,  $\mathbb{P}_t$  will assign the same probability to every set as  $\mathbb{P}^\dagger$ .

### A.4 Pushforward Measures

We have seen that to a random variable  $X$  in  $\mathbb{R}^D$  corresponds a distribution  $\mu_X$ . Applying a function  $h : \mathbb{R}^D \rightarrow \mathbb{R}^d$  to  $X$  will result in a new random variable,  $h(X)$  in  $\mathbb{R}^d$ , and it is natural to ask what its distribution is. This distribution is called the *pushforward measure* of  $\mu_X$  through  $h$ , which is denoted as  $h_\# \mu_X$ , and is defined as  $h_\# \mu_X(B) = \mu_X(h^{-1}(B))$  for every subset  $B$  of  $\mathbb{R}^d$ . A way to intuitively understand this concept is that if one could sample  $X$  from  $\mu_X$ , then sampling from  $h_\# \mu_X$  can be done by simply applying  $h$  to  $X$ . Note that here  $h_\# \mu_X$  is a measure on  $\mathbb{R}^d$ .

The concept of pushforward measure is relevant in Theorem 2 as it allows us to formally reason about e.g. the distribution of encoded data,  $g_\# \mathbb{P}^*$ . Similarly, for a distribution  $\mathbb{P}_Z$  corresponding to our second-step model, we can reason about the distribution obtained after decoding, i.e.  $G_\# \mathbb{P}_Z$ .

## B Proofs

### B.1 Riemannian Measures

We begin with a quick review of a Riemannian measures. Let  $\mathcal{M}$  be a  $d$ -dimensional Riemannian manifold with Riemannian metric  $\mathbf{g}$ , and let  $(U, \phi)$  be a chart. The local Riemannian measure  $\mu_{\mathcal{M}, \phi}^{(\mathbf{g})}$  on  $\mathcal{M}$  (with its Borel  $\sigma$ -algebra) is given by:

$$\mu_{\mathcal{M}, \phi}^{(\mathbf{g})}(A) = \int_{\phi(A \cap U)} \sqrt{\det \left( \mathbf{g} \left( \frac{\partial}{\partial \phi^i}, \frac{\partial}{\partial \phi^j} \right) \right)} d\mu_d \quad (3)$$

for any measurable  $A \subseteq \mathcal{M}$ . The Riemannian measure  $\mu_{\mathcal{M}}^{(\mathbf{g})}$  on  $\mathcal{M}$  is such that:

$$\mu_{\mathcal{M}}^{(\mathbf{g})}(A \cap U) = \mu_{\mathcal{M}, \phi}^{(\mathbf{g})}(A) \quad (4)$$

for every measurable  $A \subseteq \mathcal{M}$  and every chart  $(U, \phi)$ .

If  $\mathbf{g}_1$  and  $\mathbf{g}_2$  are two Riemannian metrics on  $\mathcal{M}$ , then  $\mu_{\mathcal{M}}^{(\mathbf{g}_1)} \ll \mu_{\mathcal{M}}^{(\mathbf{g}_2)}$  and  $\mu_{\mathcal{M}}^{(\mathbf{g}_1)}$  admits a continuous and positive density with respect to  $\mu_{\mathcal{M}}^{(\mathbf{g}_2)}$ . Thus, as mentioned in the main manuscript, smoothness of probability measures is indeed independent of the choice of Riemannian metric.

Below we prove a lemma which we will later use, showing that much like the Lebesgue measure, Riemannian measures assign positive measure to nonempty open sets. While we are sure this is a known property, we could not find a proof and thus provide one.

**Lemma 1:** Let  $\mathcal{M}$  be a  $d$ -dimensional Riemannian manifold, and  $\mu_{\mathcal{M}}^{(\mathbf{g})}$  a Riemannian measure on it. Let  $A \subseteq \mathcal{M}$  be a nonempty open set in  $\mathcal{M}$ . Then  $\mu_{\mathcal{M}}^{(\mathbf{g})}(A) > 0$ .

**Proof:** Let  $(U, \phi)$  be a chart such that  $U \cap A \neq \emptyset$ , which exists because  $A \neq \emptyset$ . Clearly  $U \cap A$  is open, and since  $\phi$  is a diffeomorphism onto its image, it follows that  $\phi(U \cap A) \subseteq \mathbb{R}^d$  is also open and nonempty, and thus  $\mu_d(\phi(U \cap A)) > 0$ . As a result,

$$\mu_{\mathcal{M}}^{(\mathbf{g})}(A) \geq \mu_{\mathcal{M}}^{(\mathbf{g})}(U \cap A) = \int_{\phi(U \cap A)} \sqrt{\det \left( \mathbf{g} \left( \frac{\partial}{\partial \phi^i}, \frac{\partial}{\partial \phi^j} \right) \right)} d\mu_d > 0, \quad (5)$$

where the last inequality follows since the integrand is positive and the integration set has positive measure.  $\square$

### B.2 Manifold Overfitting Theorem

We restate the manifold overfitting theorem below for convenience:

**Theorem 1 (Manifold Overfitting):** Let  $\mathcal{M} \subset \mathbb{R}^D$  be an analytic  $d$ -dimensional embedded submanifold of  $\mathbb{R}^D$  with  $d < D$ , and  $\mathbb{P}^\dagger$  a smooth probability measure on  $\mathcal{M}$ . Then there exists a sequence of probability measures  $(\mathbb{P}_t)_{t=1}^\infty$  on  $\mathbb{R}^D$  such that:

1.  $\mathbb{P}_t \rightarrow \mathbb{P}^\dagger$  weakly as  $t \rightarrow \infty$ .
2. For every  $t \geq 1$ ,  $\mathbb{P}_t \ll \mu_D$  and  $\mathbb{P}_t$  admits a density  $p_t : \mathbb{R}^D \rightarrow \mathbb{R}_{>0}$  with respect to  $\mu_D$  such that:
  - (a)  $\lim_{t \rightarrow \infty} p_t(x) = \infty$  for every  $x \in \mathcal{M}$ .
  - (b)  $\lim_{t \rightarrow \infty} p_t(x) = 0$  for every  $x \notin \text{cl}(\mathcal{M})$ , where  $\text{cl}(\cdot)$  denotes closure in  $\mathbb{R}^D$ .

Before proving the theorem, note that  $\mathbb{P}^\dagger$  is a distribution on  $\mathcal{M}$  and  $\mathbb{P}_t$  is a distribution on  $\mathbb{R}^D$ , with their respective Borel  $\sigma$ -algebras. Weak convergence is defined for measures on the same probability space, and

so we slightly abuse notation and think of  $\mathbb{P}^\dagger$  as a measure on  $\mathbb{R}^D$  assigning to any measurable set  $A \subseteq \mathbb{R}^D$  the probability  $\mathbb{P}^\dagger(A \cap \mathcal{M})$ , which is well-defined as  $\mathcal{M}$  is an embedded submanifold of  $\mathbb{R}^D$ . We do not differentiate between  $\mathbb{P}^\dagger$  on  $\mathcal{M}$  and  $\mathbb{P}^\dagger$  on  $\mathbb{R}^D$  to avoid cumbersome notation.

**Proof:** Let  $Y$  be a random variable whose law is  $\mathbb{P}^\dagger$ , and let  $(Z_t)_{t=1}^\infty$  be a sequence of i.i.d. standard Gaussians in  $\mathbb{R}^D$ , independent of  $Y$ . We assume all the variables are defined on the same probability space  $(\Omega, \mathcal{F}, \mathbb{P})$ . Let  $X_t = Y + \sigma_t Z_t$  where  $(\sigma_t)_{t=1}^\infty$  is a positive sequence converging to 0. Let  $\mathbb{P}_t$  be the law of  $X_t$ .

First we prove 1. Clearly  $\sigma_t Z_t \rightarrow 0$  in probability and  $Y \rightarrow Y$  in distribution as  $t \rightarrow \infty$ . Since  $\sigma_t Z_t$  converges in probability to a constant, it follows that  $X_t \rightarrow Y$  in distribution, and thus  $\mathbb{P}_t \rightarrow \mathbb{P}^\dagger$  weakly.

Now we prove that  $\mathbb{P}_t \ll \mu_D$ . Let  $A \subseteq \mathbb{R}^D$  be a measurable set such that  $\mu_D(A) = 0$ . We denote the law of  $\sigma_t Z_t$  as  $\mathbb{G}_t$  and the Gaussian density in  $\mathbb{R}^D$  with mean  $m$  and covariance matrix  $\Sigma$  evaluated at  $y$  as  $\mathcal{N}(y; m, \Sigma)$ . Let  $B = \{(w, y) \in \mathbb{R}^D \times \mathcal{M} : y + w \in A\}$ . By Fubini's theorem:

$$\mathbb{P}_t(A) = \mathbb{P}(Y + \sigma_t Z_t \in A) = \int_B d\mathbb{G}_t \times \mathbb{P}^\dagger(w, y) = \int_B \mathcal{N}(w; 0, \sigma_t^2 \mathbf{I}_D) d\mu_D \times \mathbb{P}^\dagger(w, y) \quad (6)$$

$$= \int_{A \times \mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mu_D \times \mathbb{P}^\dagger(x, y) = \int_{\mathcal{M}} \int_A \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mu_D(x) d\mathbb{P}^\dagger(y) \quad (7)$$

$$= \int_{\mathcal{M}} 0 d\mathbb{P}^\dagger(y) = 0. \quad (8)$$

Then,  $\mathbb{P}_t \ll \mu_D$ , proving the first part of 2. Note also that:

$$p_t(x) = \int_{\mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) \quad (9)$$

is a valid density for  $\mathbb{P}_t$  with respect to  $\mu_D$ , once again by Fubini's theorem since, for any measurable set  $A \subseteq \mathbb{R}^D$ :

$$\int_A p_t(x) d\mu_D(x) = \int_A \int_{\mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) d\mu_D(x) \quad (10)$$

$$= \int_{A \times \mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mu_D \times \mathbb{P}^\dagger(x, y) = \mathbb{P}_t(A). \quad (11)$$

We now prove 2a. Since  $\mathbb{P}^\dagger$  being smooth is independent of the choice of Riemannian measure, we can assume without loss of generality that the Riemannian metric  $\mathfrak{g}$  on  $\mathcal{M}$  is the metric inherited from thinking of  $\mathcal{M}$  as a submanifold of  $\mathbb{R}^D$ , and we can then take a continuous and positive density  $p^\dagger$  with respect to the Riemannian measure  $\mu_{\mathcal{M}}^{(\mathfrak{g})}$  associated with this metric.

Take  $x \in \mathcal{M}$  and let  $B_r^{\mathcal{M}}(x) = \{y \in \mathcal{M} : d_{\mathcal{M}}^{(\mathfrak{g})}(x, y) \leq r\}$  denote the geodesic ball on  $\mathcal{M}$  of radius  $r$  centered at  $x$ , where  $d_{\mathcal{M}}^{(\mathfrak{g})}$  is the geodesic distance. We then have:

$$p_t(x) = \int_{\mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) \geq \int_{B_{\sigma_t}^{\mathcal{M}}(x)} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) \quad (12)$$

$$= \int_{B_{\sigma_t}^{\mathcal{M}}(x)} p^\dagger(y) \cdot \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mu_{\mathcal{M}}^{(\mathfrak{g})}(y) \geq \int_{B_{\sigma_t}^{\mathcal{M}}(x)} \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} p^\dagger(y') \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) d\mu_{\mathcal{M}}^{(\mathfrak{g})}(y) \quad (13)$$

$$= \mu_{\mathcal{M}}^{(\mathfrak{g})}(B_{\sigma_t}^{\mathcal{M}}(x)) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} p^\dagger(y') \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) \quad (14)$$

$$\geq \mu_{\mathcal{M}}^{(\mathfrak{g})}(B_{\sigma_t}^{\mathcal{M}}(x)) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} p^\dagger(y'). \quad (15)$$

Since  $B_{\sigma_t}^{\mathcal{M}}(x)$  is compact in  $\mathcal{M}$  for small enough  $\sigma_t$  and  $p^\dagger$  is continuous in  $\mathcal{M}$  and positive, it follows that  $\inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} p^\dagger(y')$  is bounded away from 0 as  $t \rightarrow \infty$ . It is then enough to show that as  $t \rightarrow \infty$ ,

$$\mu_{\mathcal{M}}^{(\mathfrak{g})}(B_{\sigma_t}^{\mathcal{M}}(x)) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) \rightarrow \infty \quad (16)$$

in order to prove that 2a holds. Let  $B_r^d(0)$  denote an  $L_2$  ball of radius  $r$  in  $\mathbb{R}^d$  centered at  $0 \in \mathbb{R}^d$ , and let  $\mu_d$  denote the Lebesgue measure on  $\mathbb{R}^d$ , so that  $\mu_d(B_r^d(0)) = C_d r^d$ , where  $C_d > 0$  is a constant depending only on  $d$ . It is known that  $\mu_{\mathcal{M}}^{(\mathfrak{g})}(B_r^{\mathcal{M}}(x)) = \mu_d(B_r^d(0)) \cdot (1 + \mathcal{O}(r^2))$  for analytic  $d$ -dimensional Riemannian manifolds (Gray, 1974), and thus:

$$\begin{aligned} & \mu_{\mathcal{M}}^{(\mathfrak{g})}(B_{\sigma_t}^{\mathcal{M}}(x)) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) \\ &= C_d \sigma_t^d (1 + \mathcal{O}(\sigma_t^2)) \cdot \inf_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \frac{1}{\sigma_t^D (2\pi)^{D/2}} \exp \left\{ -\frac{\|x - y'\|_2^2}{2\sigma_t^2} \right\} \end{aligned} \quad (17)$$

$$= \frac{C_d}{(2\pi)^{D/2}} \cdot (1 + \mathcal{O}(\sigma_t^2)) \cdot \sigma_t^{d-D} \cdot \exp \left\{ -\frac{\sup_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \|x - y'\|_2^2}{2\sigma_t^2} \right\}. \quad (18)$$

The first term is a positive constant, and the second term converges to 1. The third term goes to infinity since  $d < D$ , which leaves only the last term. Thus, as long as the last term is bounded away from 0 as  $t \rightarrow \infty$ , we can be certain that the product of all four term goes to infinity. In particular, verifying the following equation would be enough:

$$\sup_{y' \in B_{\sigma_t}^{\mathcal{M}}(x)} \|x - y'\|_2^2 \leq \sigma_t^2. \quad (19)$$

This equation holds, since for any  $x, y' \in \mathcal{M}$ , it is the case that  $\|x - y'\|_2 \leq d_{\mathcal{M}}^{(\mathfrak{g})}(x, y')$  as  $\mathfrak{g}$  is inherited from  $\mathcal{M}$  being a submanifold of  $\mathbb{R}^D$ .

Now we prove 2b for  $p_t$ . Let  $x \in \mathbb{R}^D \setminus \text{cl}(\mathcal{M})$ . We have:

$$p_t(x) = \int_{\mathcal{M}} \mathcal{N}(x - y; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) \leq \int_{\mathcal{M}} \sup_{y' \in \mathcal{M}} \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) d\mathbb{P}^\dagger(y) = \sup_{y' \in \mathcal{M}} \mathcal{N}(x - y'; 0, \sigma_t^2 \mathbf{I}_D) \quad (20)$$

$$= \sup_{y' \in \mathcal{M}} \frac{1}{\sigma_t^D (2\pi)^{D/2}} \exp \left\{ -\frac{\|x - y'\|_2^2}{2\sigma_t^2} \right\} = \frac{1}{\sigma_t^D (2\pi)^{D/2}} \cdot \exp \left\{ -\frac{\inf_{y' \in \mathcal{M}} \|x - y'\|_2^2}{2\sigma_t^2} \right\} \xrightarrow{t \rightarrow \infty} 0, \quad (21)$$

where convergence to 0 follows from  $x \notin \text{cl}(\mathcal{M})$  implying that  $\inf_{y' \in \mathcal{M}} \|x - y'\|_2^2 > 0$ . □

### B.3 Two-Step Correctness Theorem

We restate the two-step correctness theorem below for convenience:

**Theorem 2 (Two-Step Correctness):** Let  $\mathcal{M} \subseteq \mathbb{R}^D$  be a  $C^1$   $d$ -dimensional embedded submanifold of  $\mathbb{R}^D$ , and let  $\mathbb{P}^*$  be a distribution on  $\mathcal{M}$ . Assume there exist measurable functions  $G : \mathbb{R}^D \rightarrow \mathbb{R}^D$  and  $g : \mathbb{R}^D \rightarrow \mathbb{R}^d$  such that  $G(g(x)) = x$ ,  $\mathbb{P}^*$ -almost surely. Then:

1.  $G_{\#}(g_{\#}\mathbb{P}^*) = \mathbb{P}^*$ , where  $h_{\#}\mathbb{P}$  denotes the pushforward of measure  $\mathbb{P}$  through the function  $h$ .
2. Moreover, if  $\mathbb{P}^*$  is smooth, and  $G$  and  $g$  are  $C^1$ , then:
  - (a)  $g_{\#}\mathbb{P}^* \ll \mu_d$ .
  - (b)  $G(g(x)) = x$  for every  $x \in \mathcal{M}$ , and the functions  $\tilde{g} : \mathcal{M} \rightarrow g(\mathcal{M})$  and  $\tilde{G} : g(\mathcal{M}) \rightarrow \mathcal{M}$  given by  $\tilde{g}(x) = g(x)$  and  $\tilde{G}(z) = G(z)$  are diffeomorphisms and inverses of each other.

Similarly to the manifold overfitting theorem, we think of  $\mathbb{P}^*$  as a distribution on  $\mathbb{R}^D$ , assigning to any Borel set  $A \subseteq \mathbb{R}^D$  the probability  $\mathbb{P}^*(A \cap \mathcal{M})$ , which once again is well-defined since  $\mathcal{M}$  is an embedded submanifold of  $\mathbb{R}^D$ .

**Proof:** We start with part 1. Let  $A = \{x \in \mathbb{R}^D : G(g(x)) \neq x\}$ , which is a null set under  $\mathbb{P}^*$  by assumption. By applying the definition of pushforward measure twice, for any measurable set  $B \subseteq \mathcal{M}$ :

$$G_{\#}(g_{\#}\mathbb{P}^*)(B) = g_{\#}\mathbb{P}^*(G^{-1}(B)) = \mathbb{P}^*(g^{-1}(G^{-1}(B))) = \mathbb{P}^*(g^{-1}(G^{-1}((B \setminus A) \cup (A \cap B)))) \quad (22)$$

$$= \mathbb{P}^*(g^{-1}(G^{-1}(B \setminus A)) \cup g^{-1}(G^{-1}(A \cap B))) = \mathbb{P}^*(g^{-1}(G^{-1}(B \setminus A))) \quad (23)$$

$$= \mathbb{P}^*(B \setminus A) = \mathbb{P}^*(B), \quad (24)$$

where we used that  $g^{-1}(G^{-1}(A \cap B)) \subseteq A$ , and thus  $G_{\#}(g_{\#}\mathbb{P}^*) = \mathbb{P}^*$ . Note that this derivation requires thinking of  $\mathbb{P}^*$  as a measure on  $\mathbb{R}^D$  to ensure that  $A$  and  $g^{-1}(G^{-1}(A \cap B))$  can be assigned 0 probability.

We now prove 2b. We begin by showing that  $G(g(x)) = x$  for all  $x \in \mathcal{M}$ . Consider  $\mathbb{R}^D \times \mathcal{M}$  endowed with the product topology. Clearly  $\mathbb{R}^D \times \mathcal{M}$  is Hausdorff since both  $\mathbb{R}^D$  and  $\mathcal{M}$  are Hausdorff ( $\mathcal{M}$  is Hausdorff by the definition of a manifold). Let  $E = \{(x, x) \in \mathbb{R}^D \times \mathcal{M} : x \in \mathcal{M}\}$ , which is then closed in  $\mathbb{R}^D \times \mathcal{M}$  (since diagonals of Hausdorff spaces are closed). Consider the function  $H : \mathcal{M} \rightarrow \mathbb{R}^D \times \mathcal{M}$  given by  $H(x) = (G(g(x)), x)$ , which is clearly continuous. It follows that  $H^{-1}(E) = \{x \in \mathcal{M} : G(g(x)) = x\}$  is closed in  $\mathcal{M}$ , and thus  $\mathcal{M} \setminus H^{-1}(E) = \{x \in \mathcal{M} : G(g(x)) \neq x\}$  is open in  $\mathcal{M}$ , and by assumption  $\mathbb{P}^*(\mathcal{M} \setminus H^{-1}(E)) = 0$ . It follows by Lemma 1 in App. B.1 that  $\mathcal{M} \setminus H^{-1}(E) = \emptyset$ , and thus  $G(g(x)) = x$  for all  $x \in \mathcal{M}$ .

We now prove that  $\tilde{g}$  is a diffeomorphism. Clearly  $\tilde{g}$  is surjective, and since it admits a left inverse (namely  $G$ ), it is also injective. Then  $\tilde{g}$  is bijective, and since it is clearly  $C^1$  due to  $g$  being  $C^1$  and  $\mathcal{M}$  being an embedded submanifold of  $\mathbb{R}^D$ , it only remains to show that its inverse is also  $C^1$ . Since  $G(g(x)) = x$  for every  $x \in \mathcal{M}$ , it follows that  $G(g(\mathcal{M})) = \mathcal{M}$ , and thus  $\tilde{G}$  is well-defined (i.e. the image of its domain is indeed contained in its codomain). Clearly  $\tilde{G}$  is a left inverse to  $\tilde{g}$ , and by bijectivity of  $\tilde{g}$ , it follows  $\tilde{G}$  is its inverse. Finally,  $\tilde{G}$  is also  $C^1$  since  $G$  is  $C^1$ , so that  $\tilde{g}$  is indeed a diffeomorphism.

Now, we prove 2a. Let  $K \subset \mathbb{R}^d$  be such that  $\mu_d(K) = 0$ . We need to show that  $g_{\#}\mathbb{P}^*(K) = 0$  in order to complete the proof. We have that:

$$g_{\#}\mathbb{P}^*(K) = \mathbb{P}^*(g^{-1}(K)) = \mathbb{P}^*(g^{-1}(K) \cap \mathcal{M}). \quad (25)$$

Let  $\mathbf{g}$  be a Riemannian metric on  $\mathcal{M}$ . Since  $\mathbb{P}^* \ll \mu_{\mathcal{M}}^{(\mathbf{g})}$  by assumption, it is enough to show that  $\mu_{\mathcal{M}}^{(\mathbf{g})}(g^{-1}(K) \cap \mathcal{M}) = 0$ . Let  $\{U_{\alpha}\}_{\alpha}$  be an open (in  $\mathcal{M}$ ) cover of  $g^{-1}(K) \cap \mathcal{M}$ . Since  $\mathcal{M}$  is second countable by definition, by Lindelöf's lemma there exists a countable subcover  $\{V_{\beta}\}_{\beta \in \mathbb{N}}$ . Since  $g|_{\mathcal{M}}$  is a diffeomorphism onto its image,  $(V_{\beta}, g|_{V_{\beta}})$  is a chart for every  $\beta \in \mathbb{N}$ . We have:

$$\mu_{\mathcal{M}}^{(\mathbf{g})}(g^{-1}(K) \cap \mathcal{M}) = \mu_{\mathcal{M}}^{(\mathbf{g})}\left(g^{-1}(K) \cap \mathcal{M} \cap \bigcup_{\beta \in \mathbb{N}} V_{\beta}\right) = \mu_{\mathcal{M}}^{(\mathbf{g})}\left(\bigcup_{\beta \in \mathbb{N}} g^{-1}(K) \cap \mathcal{M} \cap V_{\beta}\right) \quad (26)$$

$$\leq \sum_{\beta \in \mathbb{N}} \mu_{\mathcal{M}}^{(\mathbf{g})}(g^{-1}(K) \cap \mathcal{M} \cap V_{\beta}) \quad (27)$$

$$= \sum_{\beta \in \mathbb{N}} \int_{g|_{V_{\beta}}(g^{-1}(K) \cap \mathcal{M} \cap V_{\beta})} \sqrt{\det\left(\mathbf{g}\left(\frac{\partial}{\partial g|_{V_{\beta}}^i}, \frac{\partial}{\partial g|_{V_{\beta}}^j}\right)\right)} d\mu_d = 0, \quad (28)$$

where the final equality follows from  $g|_{V_{\beta}}(g^{-1}(K) \cap \mathcal{M} \cap V_{\beta}) \subseteq K$  for every  $\beta \in \mathbb{N}$  and  $\mu_d(K) = 0$ .  $\square$

## C Experimental Details

### C.1 VAE from Fig. 2

We generated  $N = 1000$  samples from  $\mathbb{P}^* = 0.3\delta_{-1} + 0.7\delta_1$ , resulting in a dataset containing 1 a total of 693 times. The Gaussian VAE had  $d = 1$ ,  $D = 1$ , and both the encoder and decoder have a single hidden layer with 25 units and ReLU activations. We use the Adam optimizer (Kingma & Ba, 2015) with learning rate 0.001 and train for 200 epochs. We use gradient norm clipping with a value of 10.

## C.2 Simulated Data

For the ground truth, we use a von Mises distribution with parameter  $\kappa = 1$ , and transform to Cartesian coordinates to obtain a distribution on the unit circle in  $\mathbb{R}^D = \mathbb{R}^2$ . We generate  $N = 1000$  samples from this distribution. For the EBM model, we use an energy function with two hidden layers of 25 units each and Swish activations (Ramachandran et al., 2017). We use the Adam optimizer with learning rate 0.01, and gradient norm clipping with value of 1. We train for 100 epochs. We follow Du & Mordatch (2019) for the training of the EBM, and use 0.1 for the objective regularization value, iterate Langevin dynamics for 60 iterations at every training step, use a step size of 10 within Langevin dynamics, sample new images with probability 0.05 in the buffer, use Gaussian noise with standard deviation 0.005 in Langevin dynamics, and truncate gradients to  $(-0.03, 0.03)$  in Langevin dynamics. For the AE+EBM model, we use an AE with  $d = 1$  and two hidden layers of 20 units each with ELU activations (Clevert et al., 2016). We use the Adam optimizer with learning rate 0.001 and train for 200 epochs. We use gradient norm clipping with a value of 10. For the EBM of this model, we use an energy function with two hidden layers of 15 units each, and all the other parameters are identical to the single step EBM. We observed some variability with respect to the seed for both the EBM and the AE+EBM models; the manuscript shows the best performing versions.

## C.3 Comparisons Against Maximum-Likelihood and OOD Detection with Implicit Models

For all experiments, we use the Adam optimizer, typically with learning rate 0.001. For all experiments we also clip gradient entries larger than 10 during optimization. We also set  $d = 20$  in all experiments.

### C.3.1 Single and First Step Models

For all single and first step models, unless specified otherwise, we pre-process the data by scaling it, i.e. dividing by the maximum absolute value entry. For all versions with added Gaussian noise, we tried standard deviation values  $\sigma \in \{1, 0.1, 0.01, 0.001, 0.0001\}$  and kept the best performing one ( $\sigma = 0.1$ , as measured by FID) unless otherwise specified.

**AEs** For MNIST and FMNIST, we use MLPs for the encoder and decoder, with ReLU activations. The encoder and decoder have each a single hidden layer with 256 units. For SVHN and CIFAR-10, we use convolutional networks. The encoder and decoder have 4 convolutional layers with  $(32, 32, 16, 16)$  and  $(16, 16, 32, 32)$  channels, respectively, followed by a flattening operation and a fully-connected layer. The convolutional networks also use ReLU activations, and have kernel size 3 and stride 1. We perform early stopping on reconstruction error with a patience of 10 epochs, for a maximum of 100 epochs.

**ARMs** We use an updated version of RNADE (Uria et al., 2013), where we use an LSTM (Hochreiter & Schmidhuber, 1997) to improve performance. More specifically, every pixel is processed sequentially through the LSTM, and a given pixel is modelled with a mixture of Gaussians whose parameters are given by transforming the hidden state obtained from all the previous pixels through a linear layer. The dimension of a pixel is given by the number of channels, so that MNIST and FMNIST use mixtures of 1-dimensional Gaussians, whereas SVHN and CIFAR-10 use mixtures of 3-dimensional Gaussians. We also tried a continuous version of the PixelCNN model (Van Oord et al., 2016), where we replaced the discrete distribution over pixels with a mixture of Gaussians, but found this model highly unstable – which is once again consistent with manifold overfitting – and thus opted for the LSTM-based model. We used 10 components for the Gaussian mixtures, and used an LSTM with 2 layers and hidden states of size 256. We train for a maximum of 100 epochs, and use early stopping on log-likelihood with a patience of 10. We also use cosine annealing on the learning rate. For the version with added Gaussian noise, we used  $\sigma = 1.0$ .

**AVB** We use the exact same configuration for the encoder and decoder as in AEs, and use an MLP with 2 hidden layers of size 256 each for the discriminator, which also uses ReLU activations. We train the MLPs for a maximum of 50 epochs, and CNNs for 100 epochs, using cosine annealing on the learning rates. For the large version, AVB<sup>+</sup>, we use two hidden layers of 256 units for the encoder and decoder MLPs, and increase the encoder and decoder number of hidden channels to  $(64, 64, 32, 32)$  and  $(32, 32, 64, 64)$ , respectively, for convolutional networks. In all cases, the encoder takes in 256-dimensional Gaussian noise with covariance

$9 \cdot I_D$ . We also tried having the decoder output per-pixel variances, but found this parameterization to be numerically unstable, which is again consistent with manifold overfitting.

**BiGAN** We used a Wasserstein-GAN (W-GAN) objective Arjovsky et al. (2017) with gradient penalties Gulrajani et al. (2017) where both the data and latents are interpolated between the real and generated samples. The gradient penalty weight was 10. The generator-encoder loss includes the W-GAN loss, and the reconstruction loss (joint latent regressor from Donahue et al. (2017)), equally weighted. For both small and large versions, we use the exact same configuration for the encoder, decoder, and discriminator as for AVB. We used learning rates of 0.0001 with cosine annealing over 200 epochs. The discriminator was trained for two steps for every step taken with the encoder/decoder.

**EBMs** For all datasets, our energy functions match the structure of the large version of the AVB encoders, except with the Swish activation function, spectral normalization, and scalar outputs. We set the energy function’s output regularization coefficient to 1 and the learning rate to 0.0003. Otherwise, we use the same hyperparameters as on the simulated data. At the beginning of training, we scale all the data to between 0 and 1. We train for 100 epochs without early stopping, which tended to halt training too early.

**NFs** We use a rational quadratic spline flow (Durkan et al., 2019) with 128 hidden units, 4 layers, and 3 blocks per layer. We train using early stopping on validation loss with a patience of 30 epochs, up to a maximum of 100 epochs. We use a learning rate of 0.0005, and use a whitening transform at the start of training to make the data zero-mean and marginally unit-variance, whenever possible (some pixels, particularly in MNIST, were only one value throughout the entire training set); note that this affine transformation does not affect the manifold structure of the data.

**VAEs** The setting for VAEs were largely identical to those of AVB, except we did not do early stopping and always trained for 100 epochs, in addition to not needing a discriminator. For large models a single hidden layer of 512 units was used for each of the encoder and decoder MLPs. We also tried the same decoder per-pixel variance parameterization that we attempted with AVB and obtained similar numerical instabilities, once again in line with manifold overfitting.

**WAEs** The setting for WAEs were identical to those of AVB, except (i) we used a patience of 30 epochs, trained for a maximum of 300 epochs, and (ii) we used only convolutional encoders and decoders, with (64, 64, 32, 32) and (32, 32, 64, 64) hidden channels, respectively. For large models the number of hidden channels was increased to (96, 96, 48, 48) and (48, 48, 96, 96) for the encoder and decoder, respectively.

### C.3.2 Second Step Models

All second step models, unless otherwise specified, pre-process the encoded data by standardizing it (i.e. subtracting the mean and dividing by the standard deviation).

**ARMs** We used the same configuration for second step ARMs as for the first step version, except the LSTM has a single hidden layer with hidden states of size 128.

**AVB** We used the same configuration for second step AVB as we did for the first step MLP version of AVB, except that we do not do early stopping and train for 100 epochs. The latent dimension is set to  $d$  (i.e. 20).

**EBMs** We used the same configuration as the EBM used for simulated data, except we use a learning rate of 0.001, clip gradient entries larger than 10, and take the energy function to have two hidden layers with (64, 32) units. We scale the data rather than standardizing as the pre-processing step.

**NFs** We used the same settings for second step NFs as we did for first step NFs, except (i) we use 64 hidden units, (ii) we do not do early stopping, training for a maximum of 100 epochs, and (iii) we use a learning rate of 0.001.



**VAEs** We used the same settings for second step VAEs as we did for first step VAEs. The latent dimension is also set to  $d$  (i.e. 20).

## D Additional Experimental Results

### D.1 Samples

We show samples obtained by the VAE,  $\text{VAE}^+$ ,  $\text{VAE}_\sigma^+$ , and VAE+ARM models in Fig. 6. In addition to the FID improvements shown in the main manuscript, we can see a very noticeable qualitative improvement obtained by the two-step models. Note that the VAE in the VAE+ARM model is the same as the single-step VAE model. Similarly, we show samples from  $\text{AVB}_\sigma^+$ , AVB+NF, AVB+EBM, and AVB+VAE in Fig. 7 where two-step models greatly improve visual quality. We also show samples from the  $\text{ARM}^+$ ,  $\text{ARM}_\sigma^+$ , and AE+ARM from the main manuscript in Fig. 8; and for the  $\text{EBM}^+$ ,  $\text{EBM}_\sigma^+$ , and AE+EBM models in Fig. 9. We can see that FID score is indeed not always indicative of image quality, and that our AE+ARM and AE+EBM models significantly outperform their single-step counterparts (except AE+EBM on MNIST). Finally, the BiGAN and WAE samples shown in Fig. 10 and Fig. 11 respectively are not consistently better for two-step models, but neither BiGANs nor WAEs are trained via maximum likelihood so manifold overfitting is not necessarily implied by Theorem 1. Other two-step combinations not shown gave similar results.

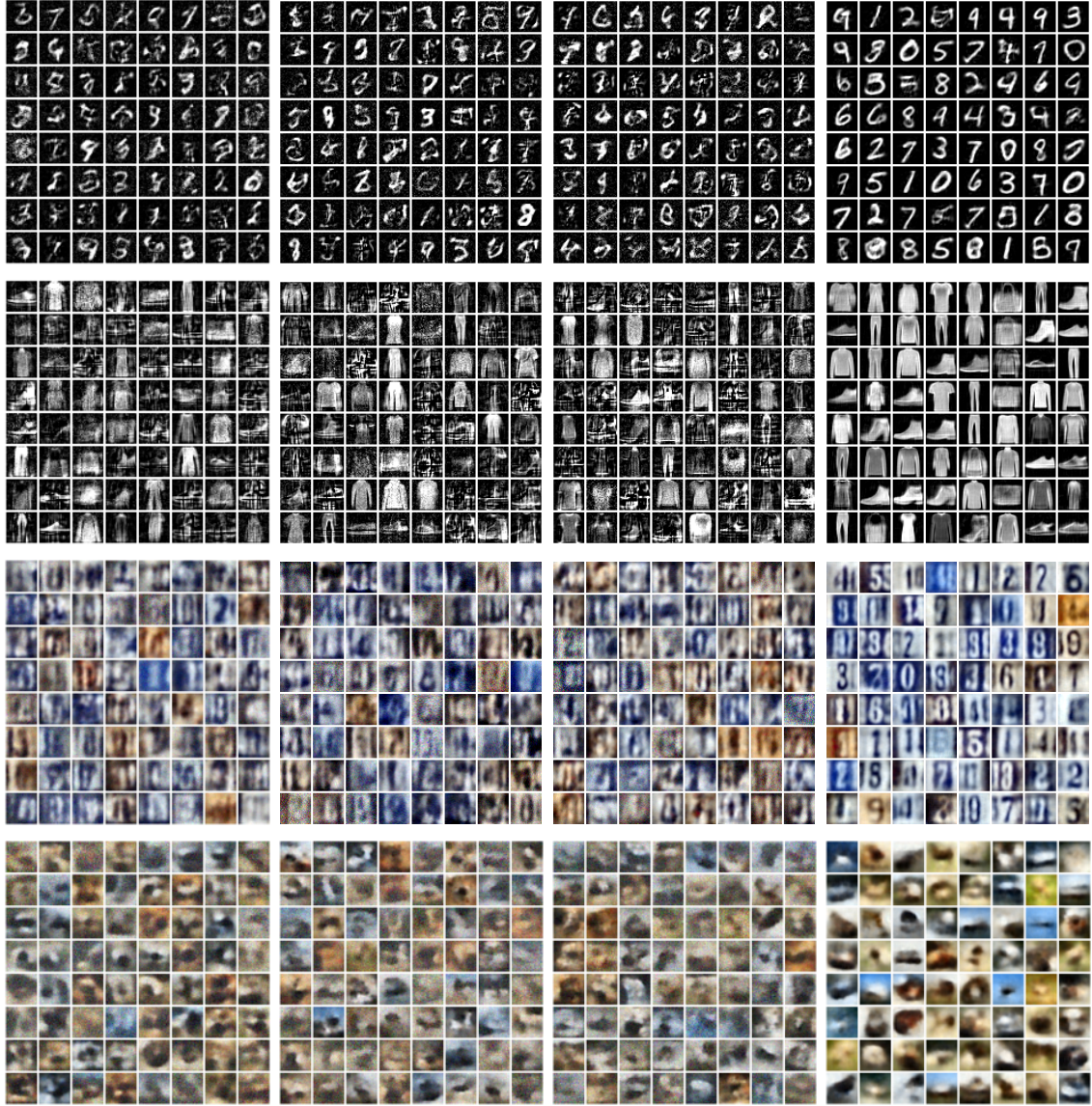


Figure 6: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are VAE (**first column**), VAE<sup>+</sup> (**second column**), VAE<sub>σ</sub><sup>+</sup> (**third column**), and VAE+ARM (**fourth column**).



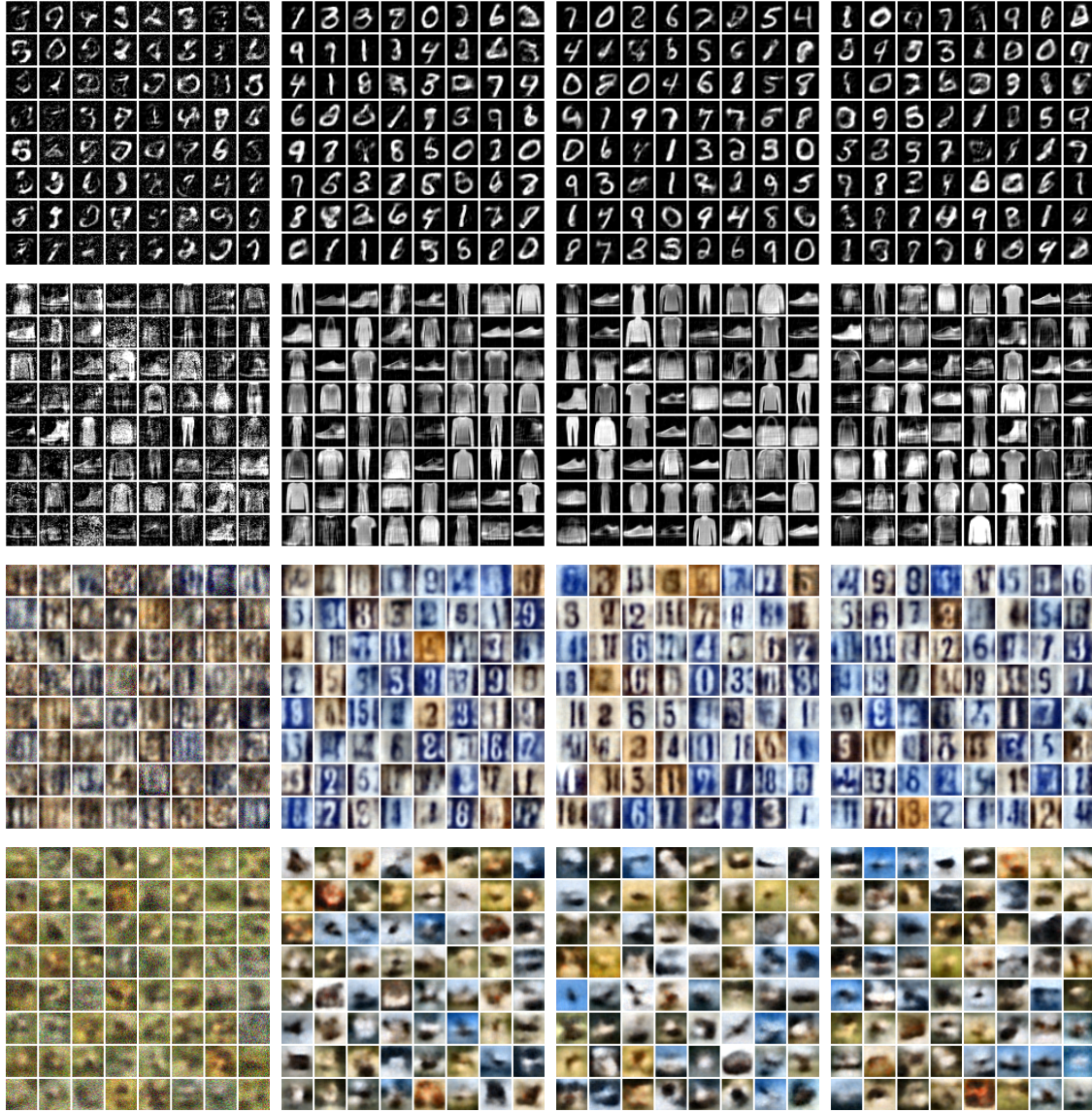


Figure 7: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are AVB+ $\sigma^+$  (**first column**), AVB+EBM (**second column**), AVB+NF (**third column**), and AVB+VAE (**fourth column**).



Figure 8: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are  $ARM^+$  (**first column**),  $ARM^+_\sigma$  (**second column**), and AE+ARM (**third column**).



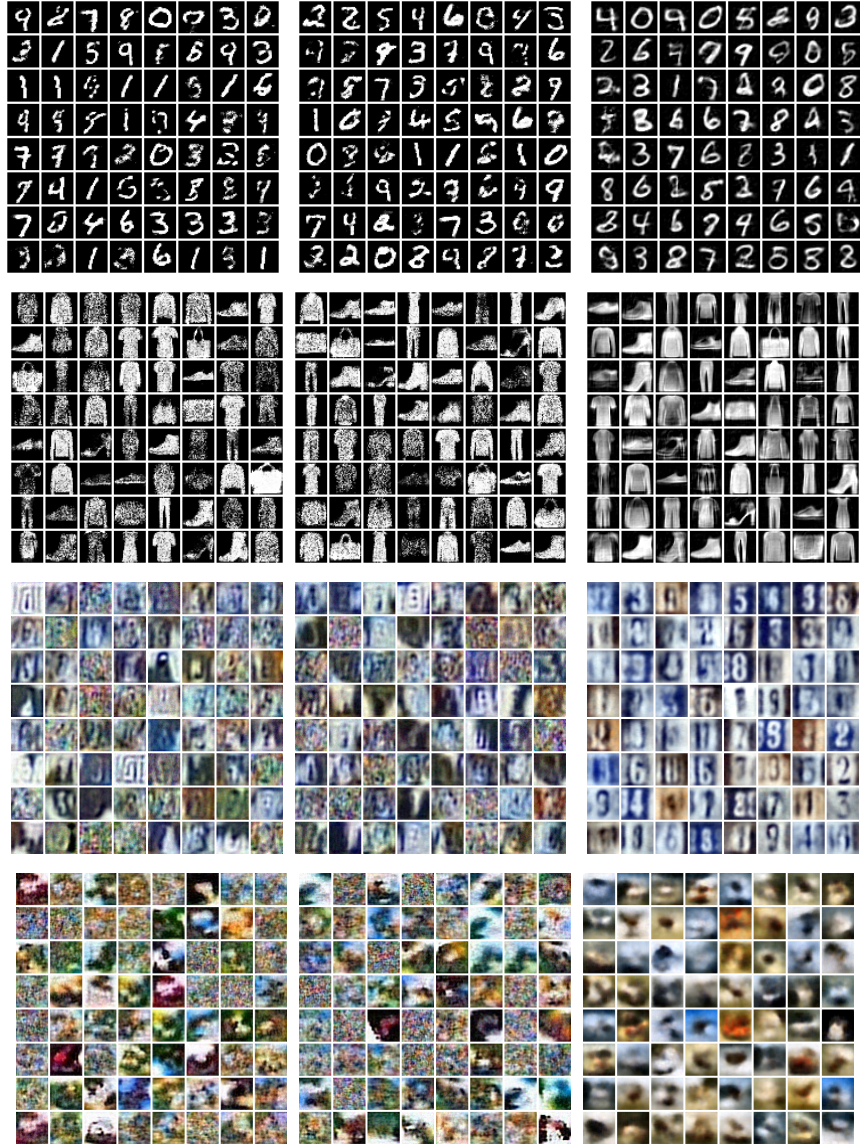


Figure 9: Uncurated samples with Langevin dynamics run for 60 steps initialized from training buffer on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are  $\text{EBM}^+$  (**first column**),  $\text{EBM}^+_{\sigma}$  (**second column**), and AE + EBM (**third column**).

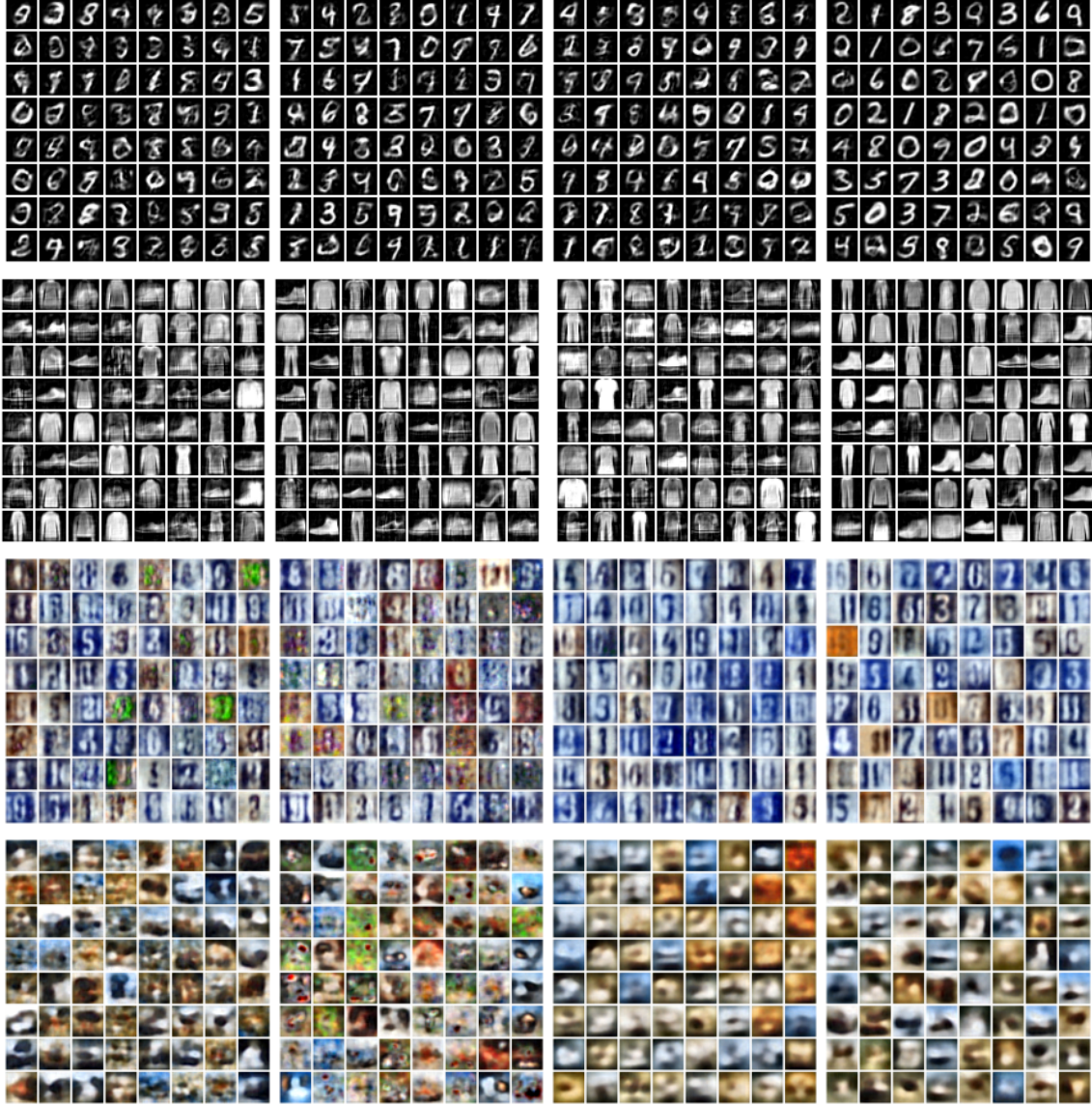


Figure 10: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are BiGAN (**first column**), BiGAN<sup>+</sup> (**second column**), BiGAN+AVB (**third column**), and BiGAN+NF (**fourth column**). BiGANs are not trained via maximum-likelihood, so Theorem 1 does not imply that manifold overfitting should occur.



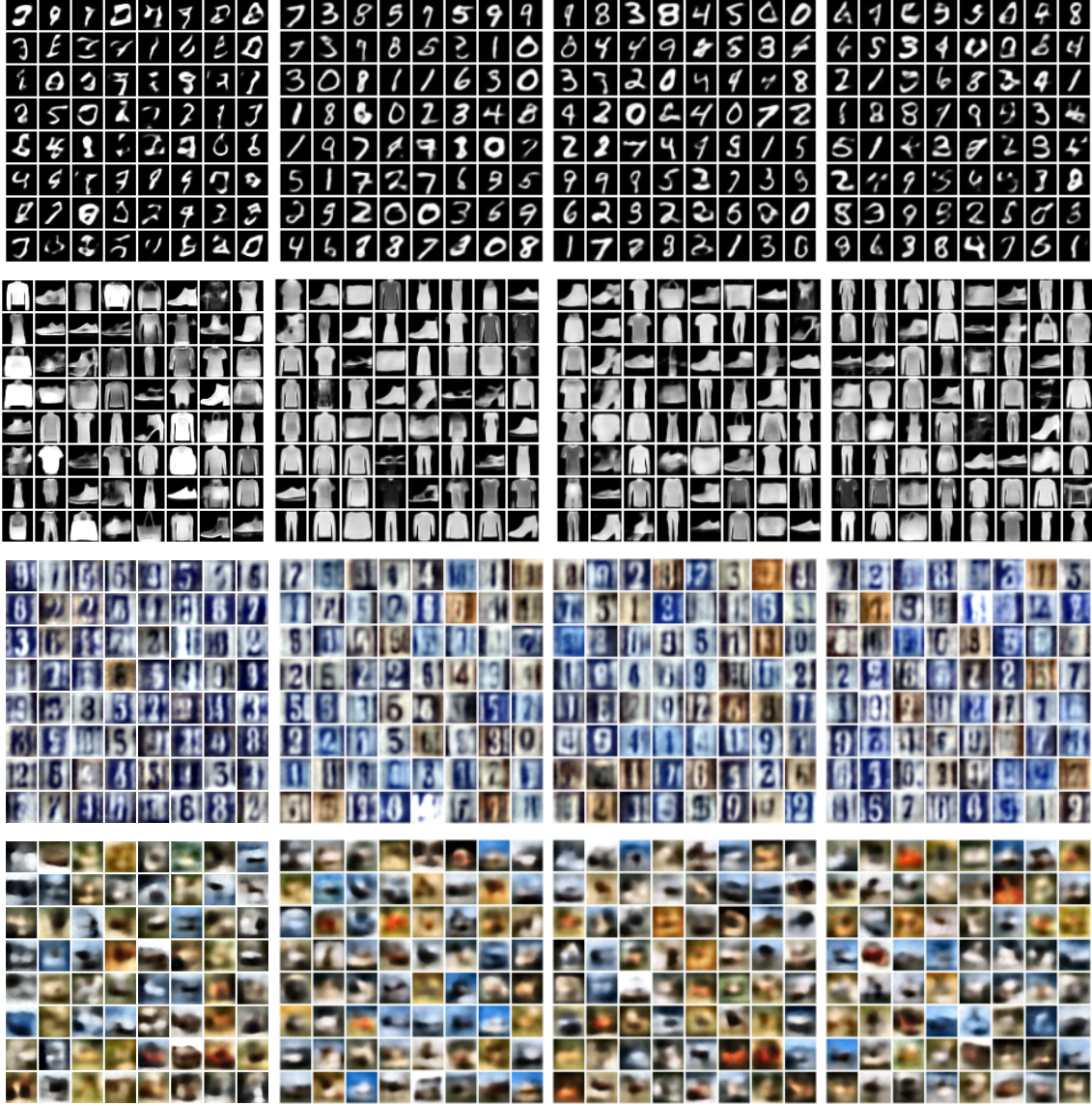


Figure 11: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are WAE<sup>+</sup> (**first column**), WAE+ARM (**second column**), WAE+NF (**third column**), and WAE+VAE (**fourth column**). WAEs are not trained via maximum-likelihood, so Theorem 1 does not imply that manifold overfitting should occur.



## D.2 EBM Improvements

Following Du & Mordatch (2019), we evaluated the single-step EBM’s sample quality on the basis of samples initialized from the training buffer. However, when MCMC samples were initialized from uniform noise, we observed that all samples would converge to a small collection of low-quality modes (see Fig. 12). Moreover, at each training epoch, these modes would change, even as the loss value decreased.

The described non-convergence in the EBM’s model distribution is consistent with Corollary 1. On the other hand, when used as a low-dimensional density estimator in the two-step procedure, this problem vanished: MCMC samples initialized from random noise yielded diverse images. See Fig. 12 for a comparison.

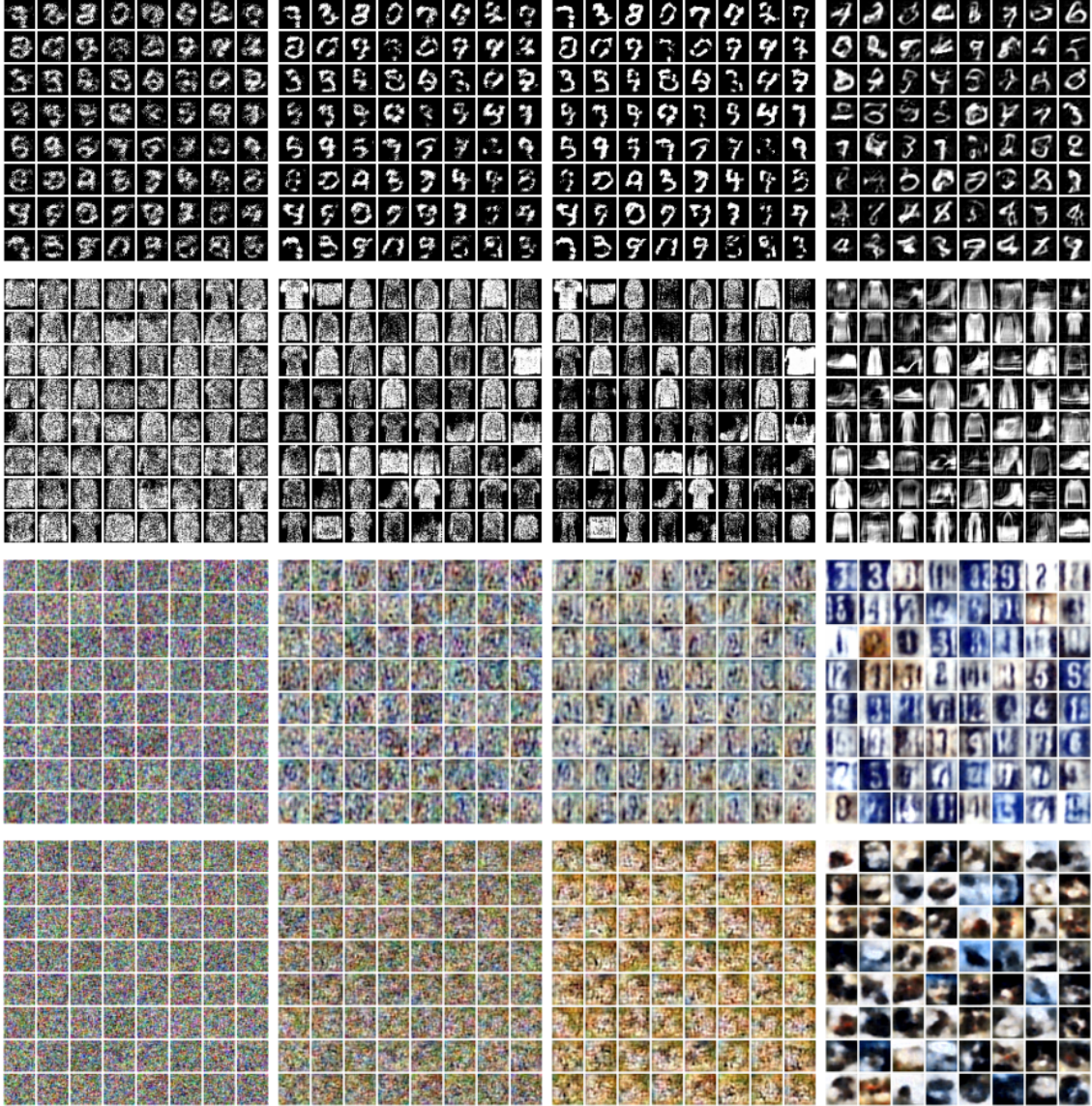


Figure 12: Uncurated samples with Langevin dynamics initialized from random noise (with no buffer) trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), and CIFAR-10 (**fourth row**). Models are EBM<sup>+</sup> with 60 steps (**first column**), EBM<sup>+</sup> with 200 steps (**second column**), EBM<sup>+</sup> with 500 steps (**third column**), and AE + EBM with 60 steps, (**fourth column**).

### D.3 FID, Precision, and Recall Scores

We show in Table 3 precision and recall (along with FID) of all the models used in Sec. 6.2. We opt for the precision and recall scores of Kynkäänniemi et al. (2019) rather than those of Sajjadi et al. (2018) as the former aim to improve on the latter. We also tried the density and coverage metrics proposed by Naeem et al. (2020), but found these metrics to correlate with visual quality less than FID. We can see in Table 3 that two-step models consistently outperform single-step models in recall, while either also outperforming or not underperforming in precision. Much like with FID score, some instances of AE+ARM and AE+EBM have worse scores on both precision and recall than their corresponding single-step model. Given the superior visual quality of those two-step models, we also consider these as failure cases of the evaluation metrics themselves, which we highlight in red in Table 3. We believe that some non-highlighted results do not properly reflect the magnitude by which the two-step models outperformed single-step models, and encourage the reader to see the corresponding samples (as previously mentioned, the only two-step model out of 44<sup>5</sup> that did not visually outperform its single-step counterpart was AE+EBM on MNIST).

We show in Table 4 the FID scores of models involving BiGANs and WAEs. These methods are not trained via maximum likelihood, so Theorem 1 does not apply. In contrast to the likelihood-based models from Table 1, there is no significant improvement in FID for BiGANs and WAEs from using a two-step approach, and sometimes two-step models perform worse. However, for BiGANs we observe similar visual quality in samples (see Fig. 10), once again highlighting a failure of the FID score as a metric. We show these failures with red in Table 4.

Table 3: FID (lower is better) and Precision, and Recall scores (higher is better). Means  $\pm$  standard errors across 3 runs are shown. Unreliable scores are highlighted in red.

MODEL	MNIST			FMNIST			SVHN			CIFAR-10		
	FID	Precision	Recall	FID	Precision	Recall	FID	Precision	Recall	FID	Precision	Recall
AVB	22.1 $\pm$ 0.1	0.08 $\pm$ 0.00	0.12 $\pm$ 0.01	19.2 $\pm$ 0.2	0.21 $\pm$ 0.01	0.06 $\pm$ 0.01	8.5 $\pm$ 0.0	0.79 $\pm$ 0.01	0.06 $\pm$ 0.00	14.4 $\pm$ 0.3	0.93 $\pm$ 0.01	0.04 $\pm$ 0.00
AVB <sup>+</sup>	21.1 $\pm$ 0.2	0.10 $\pm$ 0.00	0.12 $\pm$ 0.00	17.9 $\pm$ 0.4	0.31 $\pm$ 0.01	0.05 $\pm$ 0.01	9.3 $\pm$ 0.4	0.82 $\pm$ 0.01	0.04 $\pm$ 0.01	14.0 $\pm$ 0.0	0.96 $\pm$ 0.00	0.03 $\pm$ 0.00
AVB <sub><math>\sigma</math></sub> <sup>+</sup>	21.0 $\pm$ 0.0	0.10 $\pm$ 0.00	0.11 $\pm$ 0.01	17.8 $\pm$ 0.2	0.31 $\pm$ 0.01	0.05 $\pm$ 0.00	9.1 $\pm$ 0.0	0.77 $\pm$ 0.02	0.06 $\pm$ 0.01	14.1 $\pm$ 0.1	0.96 $\pm$ 0.01	0.03 $\pm$ 0.00
AVB+ARM	10.1 $\pm$ 0.0	0.28 $\pm$ 0.00	0.28 $\pm$ 0.00	12.0 $\pm$ 0.0	0.35 $\pm$ 0.00	0.11 $\pm$ 0.00	5.4 $\pm$ 0.1	0.64 $\pm$ 0.01	0.24 $\pm$ 0.00	12.0 $\pm$ 0.1	0.81 $\pm$ 0.01	0.10 $\pm$ 0.00
AVB+AVB	13.2 $\pm$ 0.1	0.20 $\pm$ 0.00	0.25 $\pm$ 0.01	14.6 $\pm$ 0.3	0.27 $\pm$ 0.01	0.08 $\pm$ 0.00	6.6 $\pm$ 0.1	0.73 $\pm$ 0.01	0.12 $\pm$ 0.00	12.1 $\pm$ 0.2	0.88 $\pm$ 0.00	0.07 $\pm$ 0.00
AVB+EBM	11.0 $\pm$ 0.1	0.27 $\pm$ 0.00	0.21 $\pm$ 0.02	13.9 $\pm$ 0.6	0.36 $\pm$ 0.01	0.08 $\pm$ 0.02	6.2 $\pm$ 0.1	0.67 $\pm$ 0.00	0.18 $\pm$ 0.01	12.4 $\pm$ 0.3	0.85 $\pm$ 0.01	0.08 $\pm$ 0.00
AVB+NF	10.2 $\pm$ 0.1	0.28 $\pm$ 0.00	0.28 $\pm$ 0.00	12.0 $\pm$ 0.1	0.36 $\pm$ 0.00	0.11 $\pm$ 0.00	5.5 $\pm$ 0.0	0.64 $\pm$ 0.00	0.23 $\pm$ 0.00	12.0 $\pm$ 0.0	0.82 $\pm$ 0.01	0.10 $\pm$ 0.00
AVB+VAE	11.0 $\pm$ 0.1	0.26 $\pm$ 0.00	0.28 $\pm$ 0.00	12.4 $\pm$ 0.1	0.32 $\pm$ 0.00	0.12 $\pm$ 0.00	5.6 $\pm$ 0.3	0.65 $\pm$ 0.00	0.21 $\pm$ 0.00	12.2 $\pm$ 0.0	0.82 $\pm$ 0.01	0.10 $\pm$ 0.00
VAE	21.4 $\pm$ 0.1	0.13 $\pm$ 0.00	0.09 $\pm$ 0.00	17.8 $\pm$ 0.2	0.20 $\pm$ 0.01	0.09 $\pm$ 0.00	8.8 $\pm$ 0.1	0.71 $\pm$ 0.01	0.10 $\pm$ 0.01	14.8 $\pm$ 0.3	0.90 $\pm$ 0.00	0.05 $\pm$ 0.00
VAE <sup>+</sup>	21.1 $\pm$ 0.1	0.08 $\pm$ 0.00	0.16 $\pm$ 0.00	18.3 $\pm$ 0.2	0.18 $\pm$ 0.00	0.10 $\pm$ 0.00	8.9 $\pm$ 0.2	0.69 $\pm$ 0.01	0.11 $\pm$ 0.00	14.7 $\pm$ 0.1	0.90 $\pm$ 0.01	0.05 $\pm$ 0.00
VAE <sub><math>\sigma</math></sub> <sup>+</sup>	21.3 $\pm$ 0.1	0.08 $\pm$ 0.00	0.15 $\pm$ 0.00	18.0 $\pm$ 0.2	0.18 $\pm$ 0.01	0.11 $\pm$ 0.00	8.9 $\pm$ 0.1	0.67 $\pm$ 0.02	0.12 $\pm$ 0.01	14.6 $\pm$ 0.2	0.90 $\pm$ 0.01	0.05 $\pm$ 0.00
VAE+ARM	10.1 $\pm$ 0.1	0.28 $\pm$ 0.00	0.29 $\pm$ 0.01	12.0 $\pm$ 0.1	0.38 $\pm$ 0.00	0.10 $\pm$ 0.00	5.3 $\pm$ 0.1	0.64 $\pm$ 0.00	0.23 $\pm$ 0.01	12.0 $\pm$ 0.1	0.81 $\pm$ 0.00	0.10 $\pm$ 0.00
VAE+AVB	14.0 $\pm$ 0.2	0.20 $\pm$ 0.01	0.24 $\pm$ 0.01	14.5 $\pm$ 0.4	0.27 $\pm$ 0.01	0.11 $\pm$ 0.00	6.9 $\pm$ 0.0	0.68 $\pm$ 0.01	0.16 $\pm$ 0.00	12.3 $\pm$ 0.3	0.84 $\pm$ 0.00	0.09 $\pm$ 0.00
VAE+EBM	11.7 $\pm$ 0.2	0.27 $\pm$ 0.00	0.17 $\pm$ 0.00	13.6 $\pm$ 0.2	0.37 $\pm$ 0.00	0.07 $\pm$ 0.00	6.3 $\pm$ 0.1	0.66 $\pm$ 0.00	0.19 $\pm$ 0.01	12.9 $\pm$ 0.4	0.84 $\pm$ 0.00	0.09 $\pm$ 0.00
VAE+NF	10.1 $\pm$ 0.1	0.27 $\pm$ 0.00	0.29 $\pm$ 0.00	11.8 $\pm$ 0.0	0.38 $\pm$ 0.01	0.11 $\pm$ 0.00	5.3 $\pm$ 0.1	0.64 $\pm$ 0.00	0.23 $\pm$ 0.00	11.9 $\pm$ 0.1	0.82 $\pm$ 0.00	0.10 $\pm$ 0.00
ARM <sup>+</sup>	<b>17.9 <math>\pm</math> 2.4</b>	0.53 $\pm$ 0.03	0.11 $\pm$ 0.05	<b>11.8 <math>\pm</math> 1.3</b>	0.31 $\pm$ 0.03	0.23 $\pm$ 0.04	<b>8.2 <math>\pm</math> 0.5</b>	0.85 $\pm$ 0.02	0.01 $\pm$ 0.00	<b>11.4 <math>\pm</math> 0.1</b>	0.92 $\pm$ 0.00	0.07 $\pm$ 0.00
ARM <sub><math>\sigma</math></sub> <sup>+</sup>	<b>5.3 <math>\pm</math> 0.1</b>	<b>0.46 <math>\pm</math> 0.02</b>	<b>0.47 <math>\pm</math> 0.03</b>	<b>4.7 <math>\pm</math> 0.1</b>	<b>0.40 <math>\pm</math> 0.03</b>	<b>0.38 <math>\pm</math> 0.05</b>	<b>6.4 <math>\pm</math> 0.1</b>	0.73 $\pm$ 0.03	0.10 $\pm$ 0.04	<b>10.8 <math>\pm</math> 0.1</b>	<b>0.90 <math>\pm</math> 0.01</b>	<b>0.08 <math>\pm</math> 0.01</b>
AE+ARM	9.8 $\pm$ 0.1	0.28 $\pm$ 0.00	0.29 $\pm$ 0.00	12.2 $\pm$ 0.1	0.38 $\pm$ 0.00	0.09 $\pm$ 0.00	5.8 $\pm$ 0.0	0.64 $\pm$ 0.00	0.21 $\pm$ 0.00	12.2 $\pm$ 0.1	0.82 $\pm$ 0.00	0.09 $\pm$ 0.00
EBM <sup>+</sup>	8.7 $\pm$ 0.3	0.42 $\pm$ 0.01	0.02 $\pm$ 0.00	<b>12.7 <math>\pm</math> 0.3</b>	<b>0.38 <math>\pm</math> 0.01</b>	<b>0.38 <math>\pm</math> 0.01</b>	<b>7.7 <math>\pm</math> 0.2</b>	0.73 $\pm$ 0.01	0.08 $\pm$ 0.01	<b>12.0 <math>\pm</math> 0.1</b>	0.89 $\pm$ 0.00	0.05 $\pm$ 0.00
EBM <sub><math>\sigma</math></sub> <sup>+</sup>	8.2 $\pm$ 0.4	0.42 $\pm$ 0.01	0.02 $\pm$ 0.00	<b>12.1 <math>\pm</math> 0.4</b>	0.38 $\pm$ 0.00	0.01 $\pm$ 0.00	<b>8.4 <math>\pm</math> 0.5</b>	0.73 $\pm$ 0.01	0.08 $\pm$ 0.01	<b>11.9 <math>\pm</math> 0.2</b>	0.90 $\pm$ 0.00	0.05 $\pm$ 0.00
AE+EBM	11.3 $\pm$ 0.2	0.27 $\pm$ 0.01	0.21 $\pm$ 0.01	12.7 $\pm$ 0.1	0.37 $\pm$ 0.01	0.09 $\pm$ 0.01	6.7 $\pm$ 0.2	0.68 $\pm$ 0.00	0.16 $\pm$ 0.00	11.9 $\pm$ 0.1	0.85 $\pm$ 0.01	0.09 $\pm$ 0.00

<sup>5</sup>Among the 44 models comparing against maximum-likelihood in Table 3.

Table 4: FID scores (lower is better) for non-likelihood based GAEs and two-step models. These GAEs are not trained to maximize likelihood, so Theorem 1 does not apply. Means  $\pm$  standard errors across 3 runs are shown. Unreliable scores are shown in red. Samples for unreliable scores are provided in Fig. 10.

MODEL	MNIST	FMNIST	SVHN	CIFAR-10
BiGAN	12.2 $\pm$ 0.1	13.2 $\pm$ 0.0	6.6 $\pm$ 0.0	12.9 $\pm$ 0.2
BiGAN <sup>+</sup>	12.0 $\pm$ 0.1	12.8 $\pm$ 0.1	7.1 $\pm$ 0.1	11.5 $\pm$ 0.1
BiGAN+ARM	10.6 $\pm$ 0.1	12.8 $\pm$ 0.0	6.1 $\pm$ 0.0	13.9 $\pm$ 0.1
BiGAN+AVB	12.5 $\pm$ 0.0	13.3 $\pm$ 0.2	7.4 $\pm$ 0.2	15.4 $\pm$ 0.2
BiGAN+EBM	11.4 $\pm$ 0.1	13.5 $\pm$ 0.0	6.9 $\pm$ 0.2	14.5 $\pm$ 0.2
BiGAN+NF	10.5 $\pm$ 0.0	12.8 $\pm$ 0.0	6.2 $\pm$ 0.2	13.7 $\pm$ 0.1
BiGAN+VAE	11.4 $\pm$ 0.1	13.4 $\pm$ 0.0	6.2 $\pm$ 0.0	14.0 $\pm$ 0.1
WAE	8.0 $\pm$ 0.9	9.7 $\pm$ 0.9	7.5 $\pm$ 0.9	13.0 $\pm$ 0.2
WAE <sup>+</sup>	7.7 $\pm$ 0.8	10.0 $\pm$ 0.5	7.6 $\pm$ 1.1	12.8 $\pm$ 0.0
WAE+ARM	6.0 $\pm$ 0.6	12.2 $\pm$ 0.5	9.2 $\pm$ 3.3	12.7 $\pm$ 0.2
WAE+AVB	8.2 $\pm$ 0.9	10.8 $\pm$ 0.9	9.3 $\pm$ 3.0	12.8 $\pm$ 0.3
WAE+EBM	7.2 $\pm$ 0.8	14.0 $\pm$ 1.3	8.6 $\pm$ 2.0	12.9 $\pm$ 0.1
WAE+NF	6.0 $\pm$ 0.5	10.6 $\pm$ 1.5	9.2 $\pm$ 3.3	12.6 $\pm$ 0.2
WAE+VAE	6.6 $\pm$ 0.7	10.7 $\pm$ 1.1	9.1 $\pm$ 3.1	12.8 $\pm$ 0.2

#### D.4 OOD Detection

**OOD Metric** We now precisely describe our classification metric, which properly accounts for datasets of imbalanced size and ensures correct directionality, in that higher likelihoods are considered to be in-distribution. First, using the in- and out-of-sample training likelihoods, we train a decision stump – i.e. a single-threshold-based classifier. Then, calling that threshold  $T$ , we count the number of in-sample test likelihoods which are greater than  $T$ ,  $n_{I>T}$ , and the number of out-of-sample test likelihoods which are greater than  $T$ ,  $n_{O>T}$ . Then, calling the number of in-sample test points  $n_I$ , and the number of OOD test points  $n_O$ , our final classification rate  $\text{acc}$  is given as

$$\text{acc} = \frac{n_{I>T} + \frac{n_I}{n_O} \cdot (n_O - n_{O>T})}{2n_I}. \quad (29)$$

Intuitively, we can think of this metric as simply the fraction of correctly-classified points (i.e.  $\text{acc}' = \frac{n_{O>T} + (n_O - n_{I>T})}{n_I + n_O}$ ), but with the contributions from the OOD data re-weighted by a factor of  $\frac{n_I}{n_O}$  to ensure both datasets are equally weighted in the metric.

We show further OOD detection results using  $\log p_Z$  in Table 5, and using  $\log p_X$  in Table 6. We also show corresponding histograms in Fig. 13, Fig. 14, and Fig. 15.

Table 5: OOD classification accuracy as a percentage (higher is better), using  $\log p_Z$ . Means  $\pm$  standard errors across 3 runs are shown. Arrows point from in-distribution to OOD data.

MODEL	FMNIST $\rightarrow$ MNIST	CIFAR-10 $\rightarrow$ SVHN
AVB <sup>+</sup>	96.0 $\pm$ 0.5	23.4 $\pm$ 0.1
AVB+ARM	89.9 $\pm$ 2.4	40.6 $\pm$ 0.2
AVB+AVB	74.4 $\pm$ 2.2	45.2 $\pm$ 0.2
AVB+EBM	49.5 $\pm$ 0.1	49.0 $\pm$ 0.0
AVB+NF	89.2 $\pm$ 0.9	46.3 $\pm$ 0.9
AVB+VAE	78.4 $\pm$ 1.5	40.2 $\pm$ 0.1
VAE <sup>+</sup>	96.1 $\pm$ 0.1	23.8 $\pm$ 0.2
VAE+ARM	92.6 $\pm$ 1.0	39.7 $\pm$ 0.4
VAE+AVB	80.6 $\pm$ 2.0	45.4 $\pm$ 1.1
VAE+EBM	54.1 $\pm$ 0.7	49.2 $\pm$ 0.0
VAE+NF	91.7 $\pm$ 0.3	47.1 $\pm$ 0.1
ARM <sup>+</sup>	9.9 $\pm$ 0.6	15.5 $\pm$ 0.0
AE+ARM	86.5 $\pm$ 0.9	37.4 $\pm$ 0.2
EBM <sup>+</sup>	32.5 $\pm$ 1.1	46.4 $\pm$ 3.1
AE+EBM	50.9 $\pm$ 0.2	49.4 $\pm$ 0.6

Table 6: OOD classification accuracy as a percentage (higher is better), using  $\log p_X$ . Means  $\pm$  standard errors across 3 runs are shown. Arrows point from in-distribution to OOD data.

MODEL	FMNIST $\rightarrow$ MNIST	CIFAR-10 $\rightarrow$ SVHN
AVB <sup>+</sup>	96.0 $\pm$ 0.5	23.4 $\pm$ 0.1
AVB+ARM	90.8 $\pm$ 1.8	37.7 $\pm$ 0.5
AVB+AVB	75.0 $\pm$ 2.2	43.7 $\pm$ 2.0
AVB+EBM	53.3 $\pm$ 7.1	39.1 $\pm$ 0.9
AVB+NF	89.2 $\pm$ 0.8	43.9 $\pm$ 1.3
AVB+VAE	78.7 $\pm$ 1.6	40.2 $\pm$ 0.2
VAE <sup>+</sup>	96.1 $\pm$ 0.1	23.8 $\pm$ 0.2
VAE+ARM	93.7 $\pm$ 0.7	37.6 $\pm$ 0.4
VAE+AVB	82.4 $\pm$ 2.4	42.2 $\pm$ 1.0
VAE+EBM	63.7 $\pm$ 1.7	42.4 $\pm$ 0.9
VAE+NF	91.7 $\pm$ 0.3	42.4 $\pm$ 0.3
ARM <sup>+</sup>	9.9 $\pm$ 0.6	15.5 $\pm$ 0.0
AE+ARM	89.5 $\pm$ 0.2	33.8 $\pm$ 0.3
EBM <sup>+</sup>	32.5 $\pm$ 1.1	46.4 $\pm$ 3.1
AE+EBM	56.9 $\pm$ 14.4	34.5 $\pm$ 0.1
BiGAN+ARM	81.5 $\pm$ 1.4	35.7 $\pm$ 0.4
BiGAN+AVB	59.6 $\pm$ 3.2	34.3 $\pm$ 2.3
BiGAN+EBM	57.4 $\pm$ 1.7	47.7 $\pm$ 0.7
BiGAN+NF	83.7 $\pm$ 1.2	39.2 $\pm$ 0.3
BiGAN+VAE	59.3 $\pm$ 2.1	35.6 $\pm$ 0.4
WAE+ARM	68.1 $\pm$ 12.5	37.4 $\pm$ 0.5
WAE+AVB	88.5 $\pm$ 1.4	39.6 $\pm$ 1.3
WAE+EBM	46.3 $\pm$ 4.0	37.8 $\pm$ 0.9
WAE+NF	92.2 $\pm$ 1.9	41.5 $\pm$ 2.4
WAE+VAE	88.4 $\pm$ 3.1	38.3 $\pm$ 0.6

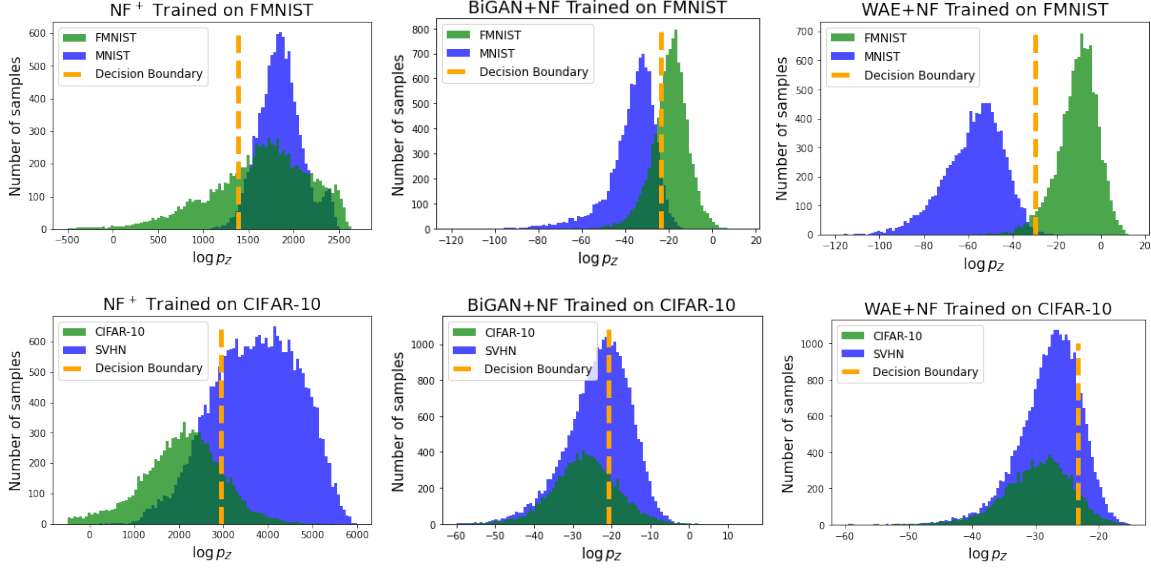


Figure 13: Comparison of the distribution of log-likelihood values between in-distribution (green) and out-of-distribution (blue) data. In both cases, the two-step models push the in-distribution likelihoods further to the right than the NF<sup>+</sup> model alone. *N.B.*: The absolute value of the likelihoods in the NF<sup>+</sup> model on its own are off by a constant factor because of the aforementioned whitening transform used to scale the data before training. However, the *relative* value within a single plot remains correct.

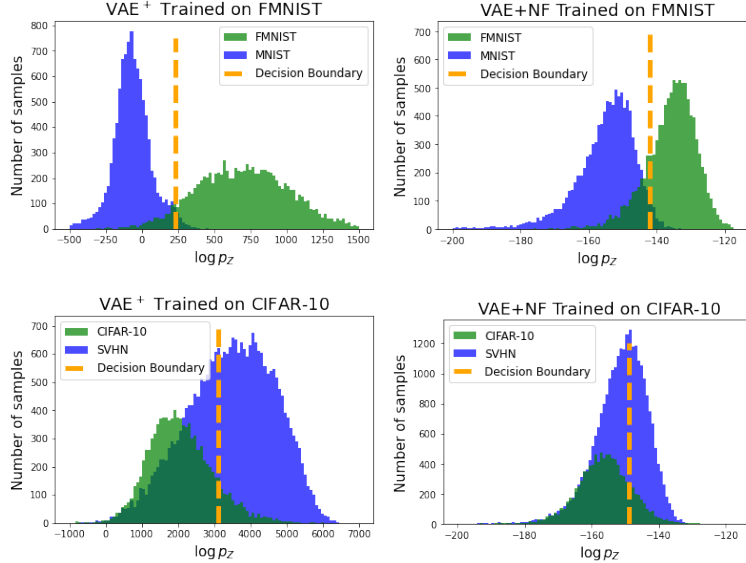


Figure 14: Comparison of the distribution of log-likelihood values between in-distribution (green) and out-of-distribution (blue) data for VAE-based models. While the VAE<sup>+</sup> model does well on FMNIST→MNIST, its performance is poor for CIFAR-10→SVHN. The two-step model VAE+NF improves on the CIFAR-10→SVHN task.

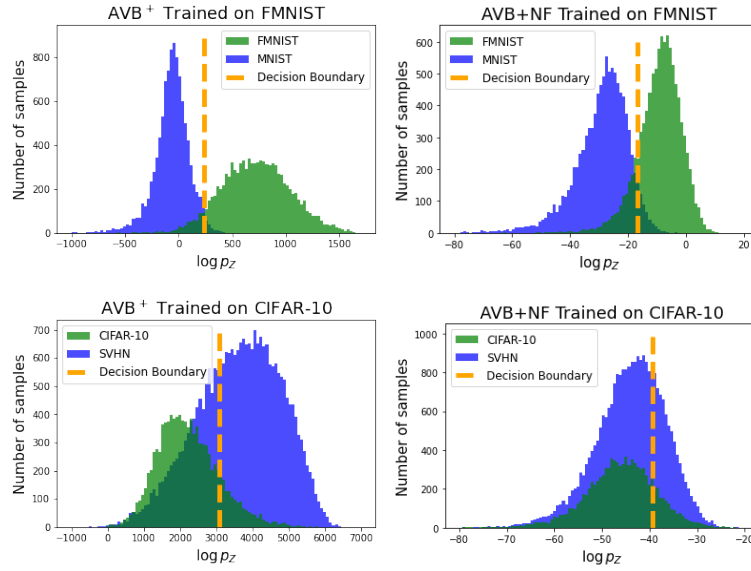


Figure 15: Comparison of the distribution of log-likelihood values between in-distribution (green) and out-of-distribution (blue) data for AVB-based models. While the  $\text{AVB}^+$  model does well on  $\text{FMNIST} \rightarrow \text{MNIST}$ , its performance is poor for  $\text{CIFAR-10} \rightarrow \text{SVHN}$ . The two-step model  $\text{AVB}+\text{NF}$  improves on the  $\text{CIFAR-10} \rightarrow \text{SVHN}$  task.