

---

# The Union of Manifolds Hypothesis and its Implications for Deep Generative Modelling

---

**Bradley C.A. Brown**

University of Waterloo & Layer 6 AI

bcabrown@uwaterloo.ca

**Anthony L. Caterini**

Layer 6 AI

anthony@layer6.ai

**Brendan Leigh Ross**

Layer 6 AI

brendan@layer6.ai

**Jesse C. Cresswell**

Layer 6 AI

jesse@layer6.ai

**Gabriel Loaiza-Ganem**

Layer 6 AI

gabriel@layer6.ai

## Abstract

Deep learning has had tremendous success at learning low-dimensional representations of high-dimensional data. This success would be impossible if there was no hidden low-dimensional structure in data of interest; this existence is posited by the *manifold hypothesis*, which states that the data lies on an unknown manifold of low intrinsic dimension. In this paper, we argue that this hypothesis does not properly capture the low-dimensional structure typically present in data. Assuming the data lies on a single manifold implies intrinsic dimension is identical across the entire data space, and does not allow for subregions of this space to have a different number of factors of variation. To address this deficiency, we put forth the *union of manifolds hypothesis*, which accommodates the existence of non-constant intrinsic dimensions. We empirically verify this hypothesis on commonly-used image datasets, finding that indeed, intrinsic dimension should be allowed to vary. We also show that classes with higher intrinsic dimensions are harder to classify, and how this insight can be used to improve classification accuracy. We then turn our attention to the impact of this hypothesis in the context of deep generative models (DGMs). Most current DGMs struggle to model datasets with several connected components and/or varying intrinsic dimensions. To tackle these shortcomings, we propose *clustered DGMs*, where we first cluster the data and then train a DGM on each cluster. We show that clustered DGMs can model multiple connected components with different intrinsic dimensions, and empirically outperform their non-clustered counterparts without increasing computational requirements.

## 1 Introduction

The manifold hypothesis [9] states that high-dimensional data of interest often lives in an unknown lower-dimensional manifold embedded in ambient space. Recently, there has been renewed interest in studying and exploiting the manifold hypothesis in the context of deep probabilistic modelling [32, 21, 58, 13, 49, 2, 40, 17, 60, 22, 46], as there is strong evidence supporting this phenomenon. From a theoretical perspective, it is known that both manifold learning and density estimation scale exponentially with the (low) *intrinsic* dimension when such structure exists [51, 53], while scaling exponentially with the (high) *ambient* dimension otherwise [15]. Thus, the most plausible explanation for the success of machine learning methods on high-dimensional data is the existence of far lower intrinsic dimension, which facilitates learning on datasets of fairly reasonable size. This is verified empirically by Pope et al. [55], in which a comprehensive study estimating the intrinsic dimension of commonly-used datasets is performed, clearly finding low-dimensional structure.

However, thinking of observed data as lying on a single unknown low-dimensional manifold is quite limiting, as this immediately implies that the intrinsic dimension throughout the dataset is constant. If we consider the intrinsic dimensionality to be the number of factors of variation generating the data, we can see that this formulation prevents distinct regions of the data’s support from having differing factors of variation. Yet this seems to be unrealistic: for example, we should not expect the number of factors needed to describe 8s and 1s in the MNIST dataset [43] to be equal.

To accommodate this intuition, in this paper we propose the *union of manifolds* hypothesis. We posit that high-dimensional data often lies not on a single manifold, but on a disjoint union of manifolds of *different intrinsic dimensions*.<sup>1</sup> While this hypothesis has been mentioned in the clustering literature [29–31, 1], to the best of our knowledge it has never been empirically explored analogously to the way that Pope et al. [55] probe the manifold hypothesis. In Section 3, we carry out this verification, confirming that commonly-used image datasets obey the union of manifolds hypothesis. We then explore connections between the intrinsic dimension of different classes, finding that classes with higher intrinsic dimension are harder to classify. Guided by this insight, we show that classification accuracy can be improved by more heavily weighting the terms corresponding to classes of higher intrinsic dimension in the cross entropy loss.

The union of manifolds hypothesis also provides a convenient basis for reasoning about datasets with multiple connected components: by modelling the data as a union of manifolds, we automatically gain the capacity to represent distinct regions of the data space separately. We expect several connected components to emerge in image-supported datasets with multiple classes; going back to the MNIST example, we would expect ten connected components there – one per digit – irrespective of their intrinsic dimensionalities. However, many deep generative models (DGMs) [12], including Gaussian variational autoencoders (VAEs) [39, 57], normalizing flows (NFs) [24, 38, 6, 18, 28], generative adversarial networks (GANs) [33], and Wasserstein autoencoders (WAEs) [69], are *pushforward models* [61, 63], i.e. they generate samples by transforming noise through a neural network  $G$ . We show in Section 4 that as a consequence of the continuity of  $G$ , pushforward models cannot properly model data whose support has multiple connected components. We then propose to remedy this situation via a mixture-based approach, which is perfectly suited to model such topological structure.

Yet we remain tasked with finding a class of base models which can handle the union of manifolds structure, rather than just multiple connected components. Recently, Loaiza-Ganem et al. [46] showed that likelihood-based models such as the aforementioned VAEs and NFs, along with energy-based models (EBMs) [26] – all of which construct high-dimensional densities through neural networks – can suffer from *manifold overfitting*, a phenomenon where the likelihood can become arbitrarily large while not recovering the data-generating distribution. Manifold overfitting is seen as a direct consequence of using high-dimensional models when the data has lower intrinsic dimension. Loaiza-Ganem et al. [46] provide *two-step models* as a solution, which work superbly in the single-manifold regime, but *not* under the union of manifolds hypothesis: two-step models could in principle still be affected by manifold overfitting if the true intrinsic dimensions vary, and also remain pushforward models susceptible to the issues outlined in the previous paragraph. We can however use these approaches as one of the base models of a larger mixture; this provides the foundation for our approach which we describe in Section 4.

We dub our proposed procedure *clustered DGMs*, which allow for properly modelling data whose support has multiple connected components of potentially varying intrinsic dimensions. Clustered DGMs proceed by first clustering the data, with the hope that each cluster corresponds to a connected component, and then training a pushforward DGM on each cluster. Varying the intrinsic dimension across these DGMs also allows our models to account for the union of manifolds hypothesis. While clustered DGMs require more disk space to be stored, we show that they do *not* require increased training nor memory budgets to outperform their non-clustered counterparts. **Our contributions are:** (1) formulating and empirically verifying the union of manifolds hypothesis, (2) showing that classes with higher intrinsic dimension are more challenging to classify correctly, and that this knowledge can be used to enable increased classification accuracy, (3) highlighting that pushforward DGMs are ill-suited to model distributions whose supports have multiple connected components, and (4) proposing clustered DGMs, both as a way to account for multiple connected components, and for varying intrinsic dimensions.

---

<sup>1</sup>The disjoint union of  $d$ -dimensional manifolds is itself a  $d$ -dimensional manifold: the possibility of having different intrinsic dimensions differentiates the union of manifolds hypothesis from the manifold hypothesis.

## 2 Background and related work

Throughout our paper, we assume access to a dataset  $\mathcal{D} = \{x_i\}_{i=1}^n$  in a high-dimensional ambient space  $\mathcal{X} = \mathbb{R}^D$ . Our first goal will be to better understand the low-dimensional structure underlying the data, and our second goal will be to use this understanding to better learn the distribution  $\mathbb{P}^*$  generating  $\mathcal{D}$ . We also call a *pushforward model* any model whose samples  $X$  are given by:

$$Z \sim \mathbb{P}_Z \quad \text{and} \quad X = G(Z), \quad (1)$$

where  $\mathbb{P}_Z$  is a base distribution in some latent space  $\mathcal{Z}$ , and  $G : \mathcal{Z} \rightarrow \mathcal{X}$  is a neural network. We highlight once again that Gaussian VAEs, GANs, WAEs, injective NFs [13, 40, 17, 60], and NFs are all pushforward models.

**Intrinsic dimension estimation** If we assume that  $\mathbb{P}^*$  is supported on a  $d$ -dimensional manifold embedded in  $\mathcal{X}$  for some unknown  $d < D$ , a natural question is how to estimate  $d$  from  $\mathcal{D}$ . Pope et al. [55, Section 3] provide an excellent overview of the topic and cover the main estimator of intrinsic dimension that we will use in our paper: the Levina and Bickel [44] estimator with the MacKay and Ghahramani [48] extension, given by:

$$\hat{d}_k := \left( \frac{1}{n(k-1)} \sum_{i=1}^n \sum_{j=1}^{k-1} \log \frac{T_k(x_i)}{T_j(x_i)} \right)^{-1}, \quad (2)$$

where  $T_j(x)$  is the Euclidean distance from  $x$  to its  $j^{\text{th}}$ -nearest neighbour, and  $k$  is a hyperparameter specifying the maximum number of nearest neighbours to consider. There also exist recent works estimating intrinsic dimension via the Wasserstein distance [10, 45], but these methods are experimental and have not yet been proven to be as practically reliable as (2). Concurrently, Tempczyk et al. [68] have proposed a new estimator of local intrinsic dimension which may coalesce nicely with our approach; we leave a study of this to future work. As we will see in the next section, popular image datasets exhibit different intrinsic dimensions in separate regions of data space.

**Accounting for low intrinsic dimension in DGMs** As mentioned in the introduction, there is a significant body of work attempting to properly account for  $d$ -dimensional manifold structure in likelihood-based DGMs, including VAEs [21], NFs [32, 58, 13, 49, 40, 17, 60], and EBMs [2]. Not accounting for low intrinsic dimension and using a full-dimensional model can result in several issues, such as numerical instabilities in NFs [7, 19] and *manifold overfitting* [21, 46], the phenomenon wherein maximum-likelihood based modelling of high-dimensional data with low-dimensional manifold structure results in the manifold being learned, but not the density on it [46, Theorem 1]. Loaiza-Ganem et al. [46] propose two-step models as a solution to manifold overfitting, where  $d$ -dimensional representations of the data are recovered through an autoencoding-type method in the first step, a DGM  $\mathbb{P}_Z$  is then trained on these representations in the second step, and the final model is given by pushing forward the second-step DGM through the decoder  $G$  from the first step. However, as mentioned earlier, two-step models provide no mechanism for learning data with multiple disjoint connected components nor varying intrinsic dimension. In other words, two-step models can model a single connected manifold, but are not well-suited to handle the *union of manifolds* hypothesis. Clustered DGMs first identify the connected components, and by modelling them separately with two-step models avoid the aforementioned issues. The only work we are aware of actively considering varying intrinsic dimensions in the context of DGMs is that of Cunningham et al. [20]. While an elegant analysis is provided there, the proposed method is computationally intensive and remains a pushforward model, which is ill-suited for modelling multiple connected components.

**Mixture models and soft clustering in DGMs** Some pushforward models attempt to account for different clusters, i.e. connected components, by mapping a multimodal distribution (e.g. mixture of Gaussians) on  $\mathcal{Z}$ , rather than the more common Gaussian or uniform distributions, through  $G$ . This idea has been used in different contexts in VAEs [50, 23, 37, 16, 56], NFs [36], GANs [8], and WAEs [67]. Due to the connectedness of  $\mathcal{Z}$ , these models will nonetheless still struggle to model distributions with multiple connected components, as we show in Section 4. A more closely related line of work uses mixtures of pushforward models, i.e. they mix multiple  $G$ 's rather than push a mixture on  $\mathcal{Z}$  forward through a single  $G$  [3, 5, 47, 25, 19, 54, 27, 73]. While these models do not all lack the ability to model multiple connected components, they cannot model varying intrinsic dimension,

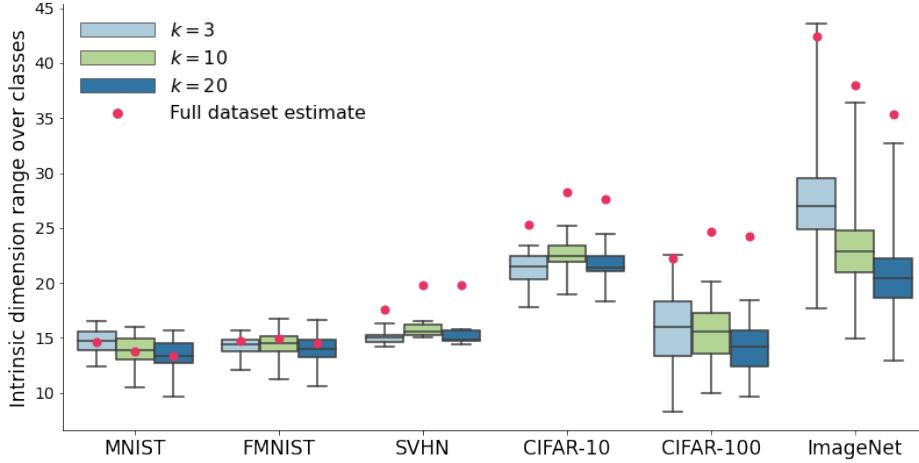


Figure 1: Boxplots showing variability of intrinsic dimension estimates across classes.

often require more convoluted training procedures than their non-mixed counterparts, and are each designed to work for a specific type of DGM. Clustered DGMs can account for varying intrinsic dimension, can be applied to any DGM, and are easy to train – although they could effortlessly be combined with some of these more involved training procedures.

### 3 The union of manifolds hypothesis

We will assume that samples from  $\mathbb{P}^*$ , including  $\mathcal{D}$ , always fall in a set which can be written as:

$$\bigsqcup_{\ell=1}^L \mathcal{M}_\ell \subset \mathcal{X}, \quad (3)$$

where  $\sqcup$  denotes disjoint union, and each  $\mathcal{M}_\ell$  is a connected manifold of dimension  $d^{(\ell)}$ . We call this assumption the *union of manifolds hypothesis*. Note that we make the assumption that each  $\mathcal{M}_\ell$  is connected merely for notational simplicity, as we could collapse each union of manifolds of matching dimension into a single disconnected manifold, but this notation allows us to reason about different  $\mathcal{M}_\ell$  as connected components of interest. Note also that the standard manifold hypothesis is equivalent to adding the assumption that all  $d^{(\ell)}$ s are identical.

#### 3.1 Verifying the hypothesis

In this section we show that commonly-used image datasets are better described by the *union of manifolds hypothesis* than just the manifold hypothesis. In order to achieve this, we first attempt to identify the connected components  $\mathcal{M}_\ell$  of the data, and then obtain estimates  $\hat{d}_k^{(\ell)}$  of the intrinsic dimension of each connected component through (2) for each  $\ell = 1, \dots, L$ . Observing varying estimates across different  $\ell$ s would strongly support the union of manifolds hypothesis, whereas observing similar estimates would support the manifold hypothesis.

Most image datasets are accompanied by class labels. We argue that class labels are a natural proxy for identifying the connected components  $\mathcal{M}_\ell$ , as one can easily conceive of continuously deforming any natural image from a given class into another image from the same class without any intermediate image not belonging to the class. For example, on MNIST we can likely continuously transform any 9 into any other 9 without leaving the manifold of 9s, whereas attempting to similarly transform a 9 into a 3 would likely result in intermediate images that are neither 9s nor 3s.

Figure 1 shows, for different values of the hyperparameter  $k$ , boxplots of the values  $\hat{d}_k^{(\ell)}$  for class labels  $\ell = 1, \dots, L$  in our considered datasets: MNIST, FMNIST [72], SVHN [52], CIFAR-10, CIFAR-100 [41], and ImageNet [62]. Two relevant patterns emerge. First, within each dataset,

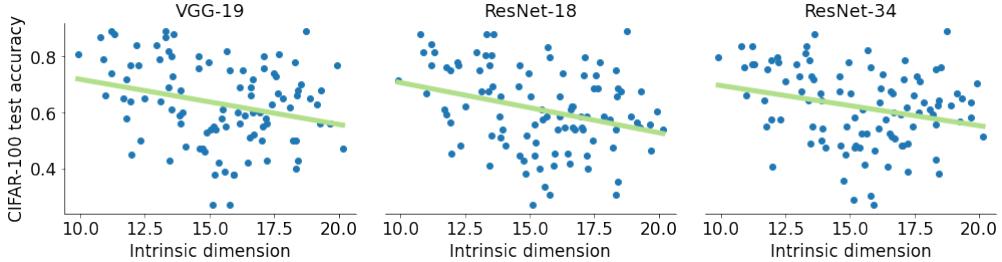


Figure 2: Class intrinsic dimension versus test accuracy on CIFAR-100, along with a least-squares regression line. Correlation coefficients for each model are  $-0.243 \pm 0.015$ ,  $-0.269 \pm 0.012$ , and  $-0.274 \pm 0.007$ , respectively (means and standard errors over 5 runs). Corresponding  $p$ -values,  $0.008 \pm 0.002$ ,  $0.019 \pm 0.008$ , and  $0.006 \pm 0.001$ , show the relationship is significant.

results are consistent across different choices of  $k$ , so that any conclusions we draw are not caused by a specific choice of this hyperparameter. ImageNet is a partial exception to this consistency across  $k$ , likely due to the fact that it has fewer datapoints per class, which can also result in (2) underestimating the true intrinsic dimension [55]. While CIFAR-100 also has fewer datapoints per class than other datasets, consistency across  $k$  suggests this is less of a concern in this case. Second, for all datasets except SVHN, we can observe a wide range of estimated intrinsic dimensions. These results empirically verify that the union of manifolds hypothesis is a more appropriate way to think about images than the manifold hypothesis. Note that even on ImageNet, the results show large ranges of intrinsic dimensions and do not vary wildly with  $k$ , so even if the specific estimated intrinsic dimensions should be taken with a grain of salt, these results still strongly support the union of manifolds hypothesis. We take a more granular look into these estimates in Appendix D.1, where we show further choices of  $k$  as well as the intrinsic dimensions of individual classes.

### 3.2 On the relationship between class intrinsic dimension and classification accuracy

Before moving on to the implications of the union of manifolds hypothesis for DGMs in the next sections, we study how it impacts the performance of classifiers. We train 3 classifiers on CIFAR-100 (see Appendix C.1 for details): a VGG-19 [66], a ResNet-18, and a ResNet-34 [34]. We focused on CIFAR-100 here as the other datasets considered in this work were either too simple to classify (MNIST, FMNIST, SVHN, CIFAR-10), or produced less-reliable intrinsic dimension estimates (ImageNet). For each class, we compute its classification accuracy (i.e. the percentage of times an instance of the class is correctly classified) and plot this against its estimated intrinsic dimension in Figure 2. We can see that, consistently across classifiers, there is an inverse relationship between estimated intrinsic dimension and classification accuracy. We also compute the correlation between these two quantities for a quantitative comparison, and compute a  $p$ -value for independence with a  $t$ -test: we find that the negative correlation is significant (see figure caption). In other words, the higher the intrinsic dimension of a class, the harder it is to classify. While clearly intrinsic dimension does not fully explain accuracy, this correlation suggests the following intuition: learning useful representations for classes with higher intrinsic dimension requires learning more factors of variation and is thus more challenging.

In order to check if this insight can help improve classifiers, we train the ResNet-18 in two different ways (see Appendix C.2 for experimental setup): in the first, we use the standard cross entropy loss, and in the second, we weight the terms corresponding to each class in the cross entropy loss in a manner proportional to their intrinsic dimension. The idea behind this reweighting of the loss is simply to focus more on classes of higher intrinsic dimension, as we have just shown these are harder to classify. Results are shown in Table 1, where we can see that this very simple change to the cross entropy loss marginally (but significantly, as error bars do not overlap) improves the accuracy of

Table 1: Means and standard errors of ResNet-18 accuracy on CIFAR-100 across 5 runs.

Weights	Test accuracy
Standard	$61.38\% \pm 0.0017\%$
Proportional to intrinsic dimension	$61.77\% \pm 0.0020\%$

the network, providing an essentially “free” improvement, given the low computational overhead of estimating intrinsic dimension.

## 4 Accounting for the union of manifolds hypothesis in DGMs

### 4.1 Topological mismatch of pushforward models

As previously mentioned, Gaussian VAEs, NFs, injective NFs, GANs, WAEs, and two-step models are *pushforward* models. In measure-theoretic language, the model from (1) is given by the pushforward measure  $G_{\#}\mathbb{P}_Z$ . We have previously mentioned that pushforward models struggle to model data whose support has multiple connected components. Intuitively this can be understood as follows: if the support  $\text{supp}(\mathbb{P}_Z)$  of  $\mathbb{P}_Z$  is connected – as is the case in the aforementioned DGMs – and  $G$  is continuous, then we should expect the support of the model,  $\text{supp}(G_{\#}\mathbb{P}_Z)$ , to be connected since the continuous image of connected sets is itself connected. The following proposition formalizes this intuition, thus showing that pushforward models whose base distribution  $\mathbb{P}_Z$  has connected support are inappropriate to model multiple connected components.

**Proposition 1:** Let  $\mathcal{Z}$  and  $\mathcal{X}$  be topological spaces and  $G : \mathcal{Z} \rightarrow \mathcal{X}$  be continuous. Considering  $\mathcal{Z}$  and  $\mathcal{X}$  as measurable spaces with their respective Borel  $\sigma$ -algebras, let  $\mathbb{P}_Z$  be a probability measure on  $\mathcal{Z}$  such that  $\text{supp}(\mathbb{P}_Z)$  is connected and  $\mathbb{P}_Z(\text{supp}(\mathbb{P}_Z)) = 1$ .<sup>2</sup> Then  $\text{supp}(G_{\#}\mathbb{P}_Z)$  is connected.

**Proof sketch:** We begin by pointing out that from the formal definition of support,  $\text{supp}(G_{\#}\mathbb{P}_Z)$  need not be equal to  $G(\text{supp}(\mathbb{P}_Z))$ . In other words, the proof does not immediately follow from the fact that continuous functions map connected sets to connected sets. In our proof we first show that  $\text{supp}(G_{\#}\mathbb{P}_Z) = \text{cl}(G(\text{supp}(\mathbb{P}_Z)))$ , where  $\text{cl}(\cdot)$  denotes closure in  $\mathcal{X}$ , and the result follows from the closure of connected sets being connected. The formal proof along with the definition of support is included in Appendix A.

The implication of Proposition 1 is simple: if the support of the data is not connected (i.e. has multiple connected components), then pushforward models are insufficient to properly model the data. Intuitively, this result rules out the possibility of  $G_{\#}\mathbb{P}_Z$  assigning 0 probability to regions of  $G(\text{supp}(\mathbb{P}_Z))$  which connect different disconnected components of the data support  $\text{supp}(\mathbb{P}^*)$ , as illustrated in Figure 3. We highlight two related concurrent works: Salmona et al. [63] show that pushforward models cannot recover multimodal distributions without exploding Lipschitz constants, and Ross et al. [61] argue that pushforward models cannot learn manifolds with complicated topological structures. The former is highly theoretical and proposes no empirical methodology, while the latter’s method could be easily incorporated into our clustered DGMs framework, which we will describe in the next section.

Note that Proposition 1 is orthogonal to the union of manifolds hypothesis, as it is a statement about the connected components of the data support, not about their intrinsic dimension. In other words, this insufficiency of pushforward models is also an issue under the (standard) manifold hypothesis, as a single manifold can have several connected components. Nonetheless, properly accounting for the union of manifolds hypothesis requires being able to model its multiple connected components. In the next section we show how to enable pushforward models to handle multiple connected components of varying intrinsic dimensions.

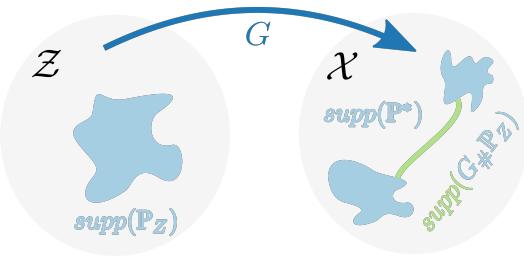


Figure 3: Depiction of the inability of pushforward models to represent multiple connected components. If  $\text{supp}(\mathbb{P}^*)$  is given by the two blue connected components on  $\mathcal{X}$ , then by continuity,  $G$  cannot map  $\text{supp}(\mathbb{P}_Z)$  to this set, and so its image must include the green region connecting them. Proposition 1 shows that the model  $G_{\#}\mathbb{P}_Z$  cannot learn to assign 0 probability to this green region, as the support of a distribution can be intuitively understood as the “smallest” set to which it assigns probability 1.

<sup>2</sup>The condition  $\mathbb{P}_Z(\text{supp}(\mathbb{P}_Z)) = 1$  is highly technical and holds in all settings of practical interest, but can be violated in contrived counterexamples. See Appendix A for a discussion.

## 4.2 Clustered DGMs

In order to properly account for multiple connected components, we must first identify them. We do this in two different ways, one supervised and one unsupervised. The supervised approach assumes access to class labels and uses them to index connected components as in Section 3. In order to not always assume access to labels, the unsupervised way runs a clustering algorithm, whose output partitions the dataset  $\mathcal{D}$  into  $L$  pairwise disjoint clusters  $\mathcal{D}_\ell$ , for  $\ell = 1, \dots, L$ . The intuition behind this step is that clustering algorithms are expected to recover clusters which correspond to the connected components we are trying to recover. In other words, we expect  $\mathcal{D}_\ell$  to be contained in  $\mathcal{M}_\ell$ .

Then, we train a pushforward DGM, with components  $\mathbb{P}_Z^{(\ell)}$  (which can be trained or fixed, depending on the type of DGM) and  $G_\ell : \mathcal{Z}_\ell \rightarrow \mathcal{X}$ , on each cluster  $\mathcal{D}_\ell$ . This procedure is summarized in Algorithm 1. Note that if  $N$  epochs of training are deemed necessary to train a baseline single pushforward model, training our  $L$  models for  $N$  epochs each requires the exact same number of gradient updates, as  $\mathcal{D}_\ell$  is on average  $L$  times smaller than  $\mathcal{D}$ . Therefore, clustered DGMs do not incur added training costs.

Importantly, clustered DGMs are not subject to the topological pathologies explained in the previous section, as each connected component is modelled separately. Our final model, which we refer to as a *clustered DGM*, is given by the mixture of all these models, where we take the mixture weights to be proportional to the size of the clusters. In other words, to obtain a sample  $X$  from a trained clustered DGM, we first sample  $\ell$  with probability proportional to  $|\mathcal{D}_\ell|$ , and then sample from the corresponding DGM:

$$\ell \sim p(\ell) = \frac{|\mathcal{D}_\ell|}{\sum_{\ell'=1}^L |\mathcal{D}_{\ell'}|}, \quad Z \mid \ell \sim \mathbb{P}_Z^{(\ell)}, \quad \text{and} \quad X = G_\ell(Z). \quad (4)$$

While saving a trained clustered DGM requires more storage than a single pushforward model of the same capacity, the sequential nature of Algorithm 1 implies the  $\ell^{\text{th}}$  model need not be in memory while training the  $(\ell + 1)^{\text{th}}$  model, resulting in no added memory costs during training. Alternatively, memory costs could be increased by parallelizing the `for` loop in Algorithm 1, thus reducing training time compared to the baseline non-clustered DGMs. Clustered DGMs can also be easily sampled from without having to load  $L$  models in memory; see Appendix B for a discussion.

While clustered DGMs properly account for disconnected supports, this only covers ‘‘half’’ of the union of manifolds hypothesis. In order to enable clustered DGMs to model varying intrinsic dimensions, it suffices to have each DGM take into consideration the intrinsic dimension of its corresponding cluster. We do this by first obtaining an estimate  $\hat{d}_k^{(\ell)}$  of the intrinsic dimension of each cluster  $\mathcal{D}_\ell$  through (2). Note that this estimator matches the one we used in Section 3 in the supervised case (as the clusters are the same), but might be different in the unsupervised case (simply because clusters need not exactly match labels). We then use a two-step model [46], where we set the latent dimension of the  $\ell^{\text{th}}$  DGM to  $\hat{d}_k^{(\ell)}$ , i.e.  $\mathcal{Z}_\ell = \mathbb{R}^{\hat{d}_k^{(\ell)}}$ . We refer to these models as *clustered two-step models*.

## 5 Experiments with DGMs

In this section we carry out experiments to verify our claims on the implications of the union of manifolds hypothesis for DGMs. In a first set of experiments, we focus exclusively on the modelling of disconnected data supports, without considering varying intrinsic dimensions. This choice allows us to empirically corroborate the theoretical limitations of pushforward models we exhibited in Section 4.1, while avoiding any potential confounding from varying intrinsic dimensionality in our models. Indeed, we observe clustered DGMs outperform their non-clustered counterparts, strongly supporting our thesis. In our second set of experiments, we aim to show that accounting for varying intrinsic dimensions in likelihood-based DGMs helps prevent manifold overfitting, particularly when the difference of intrinsic dimensions between separate connected components of the data support is

Table 2: FID scores. We show means and standard errors across 3 runs. Best models, and those whose error bars overlap with those of the best models, are bolded.

Model	MNIST	FMNIST	SVHN	CIFAR-10	CIFAR-100
NF	$10.3 \pm 0.1$	$11.1 \pm 0.1$	<b><math>5.0 \pm 0.0</math></b>	$11.3 \pm 0.0$	$11.2 \pm 0.0$
NF (GMM)	$10.0 \pm 0.1$	$10.9 \pm 0.1$	<b><math>5.0 \pm 0.0</math></b>	$11.3 \pm 0.0$	<b><math>11.1 \pm 0.0</math></b>
C-NF (labels)	<b><math>8.1 \pm 0.0</math></b>	<b><math>9.1 \pm 0.0</math></b>	$8.2 \pm 0.1$	<b><math>10.8 \pm 0.0</math></b>	$52.3 \pm 2.5$
C-NF	$9.2 \pm 0.0$	$9.7 \pm 0.0$	$29.4 \pm 12.1$	$12.8 \pm 0.1$	$12.7 \pm 0.1$
VAE	$22.4 \pm 0.1$	$18.5 \pm 0.2$	$10.1 \pm 0.4$	$16.2 \pm 1.0$	$15.0 \pm 0.2$
VAE (GMM)	$20.7 \pm 0.1$	$18.7 \pm 0.1$	$11.2 \pm 0.2$	<b><math>15.0 \pm 0.3</math></b>	<b><math>14.4 \pm 0.6</math></b>
C-VAE (labels)	<b><math>19.3 \pm 0.0</math></b>	<b><math>14.4 \pm 0.0</math></b>	<b><math>9.0 \pm 0.2</math></b>	<b><math>14.5 \pm 0.3</math></b>	$15.1 \pm 0.2$
C-VAE	$19.6 \pm 0.1$	$15.5 \pm 0.2$	$9.7 \pm 0.1$	<b><math>14.6 \pm 0.2</math></b>	<b><math>14.0 \pm 0.1</math></b>
WAE	<b><math>5.3 \pm 0.2</math></b>	$8.9 \pm 0.4$	$6.9 \pm 0.2$	$10.6 \pm 0.0$	$9.5 \pm 0.1$
WAE (GMM)	$8.5 \pm 0.4$	$9.1 \pm 0.1$	<b><math>6.1 \pm 0.1</math></b>	$13.1 \pm 0.2$	$11.7 \pm 0.0$
C-WAE (labels)	<b><math>6.3 \pm 1.5</math></b>	<b><math>7.5 \pm 0.2</math></b>	<b><math>6.4 \pm 0.6</math></b>	<b><math>9.9 \pm 0.6</math></b>	$10.3 \pm 0.1$
C-WAE	<b><math>6.1 \pm 2.4</math></b>	<b><math>7.3 \pm 0.2</math></b>	<b><math>6.9 \pm 1.6</math></b>	<b><math>10.6 \pm 0.8</math></b>	<b><math>8.9 \pm 0.0</math></b>

large. We show this through empirical gains obtained by clustered two-step models over clustered two-step models with fixed (non-varying) intrinsic dimension. All of our experimental details are included in Appendix C, and our code is available at <https://github.com/layer6ai-labs/UoMH>.

### 5.1 Better modelling of disconnected data supports

In Table 2 we compare different pushforward models against their clustered versions. We use the FID score [35] (lower is better), a commonly-used sample quality metric, to measure performance. We carry out comparisons both with the commonly-used Gaussian  $\mathbb{P}_Z$ , and with a mixture of Gaussians  $\mathbb{P}_Z$  (indicated as “GMM”). For the clustered models (indicated with the prefix “C-”), we use both the supervised version which models each class separately (indicated as “labels”), and the unsupervised version, for which we use agglomerative clustering [59] with Ward’s criterion [71] (see Appendix C.3 for details on this clustering method).

We can see that for all datasets and models except NFs on SVHN and CIFAR-100, clustered models are always either the best performing ones, or within error bars. Note that CIFAR-100 has 100 different labels, so that the clustered models trained with labels are trained with significantly less data on each cluster, and thus some performance drop is to be expected. We also point out that the fully unsupervised clustered DGMs use 10 clusters, even for CIFAR-100, precisely to avoid the aforementioned issue of each cluster having too few datapoints. Additionally, we can see that both versions of clustered models often outperform both their standard and non-clustered GMM counterparts, strongly suggesting that the clustering scheme is successfully separating the data support into regions which can be modelled as connected components. These results firmly support the conclusion that clustered DGMs are better equipped to handle disconnected data supports. Also, the GMM-based pushforward models achieve only mild or no improvement whatsoever over their baseline Gaussian-based models, and in some cases performance is even decreased. These outcomes indicate that, as shown in Section 4.1, changing  $\mathbb{P}_Z$  from a Gaussian to a mixture of Gaussians is indeed not sufficient to allow pushforward models to better model multiple connected components. We highlight once again that SVHN is an exception to some of these conclusions, which we hypothesize might be due to connected components in this dataset being closer together than in other datasets – making clustering harder – similarly to how SVHN also exhibited less variability of intrinsic dimensions in Section 3.1. We include some additional discussions and samples in Appendix D.2. We believe that the improvements obtained by clustered DGMs over their non-clustered counterparts, in spite of the individual models within a clustered DGM being trained with much less data, strongly emphasize the importance of enabling commonly-used DGMs to better deal with disconnected data supports.

### 5.2 Better modelling of varying intrinsic dimensions

In Appendix D.3, we carry out a comparison between the following: non-clustered two-step models; clustered two-step models where we set the intrinsic dimension across all clusters to its estimate

over the entire dataset; and clustered two-step models, setting intrinsic dimensions as described in Section 4.2. We found that, while both versions of clustered models outperformed non-clustered two-step models – once again confirming the conclusions from the previous section – setting intrinsic dimensions differently on every cluster did not improve performance over keeping intrinsic dimension constant across clusters. At a first glance this might appear to imply that there is no benefit to accounting for varying intrinsic dimensions, seemingly contradicting the theoretical result of Loaiza-Ganem et al. [46, Theorem 1], which states that manifold overfitting will occur whenever intrinsic dimension is overspecified. However, upon closer inspection of the proof of Theorem 1 in [46], we can see that the rate at which manifold overfitting happens depends on the difference between ambient and intrinsic dimension.<sup>3</sup> This observation explains both the good empirical results of Loaiza-Ganem et al. [46] and our aforementioned results: the former are driven by the difference between ambient and intrinsic dimension (on the order of hundreds/thousands), while the latter by the error between the true and estimated intrinsic dimension (on the order of ones/tens).

In order to verify that properly setting varying intrinsic dimensions in DGMs is advantageous, provided the true intrinsic dimensions are different enough, we generate a synthetic dataset using a pretrained BigGAN [14]. We use this model, which has  $\mathcal{Z} = \mathbb{R}^{120}$ , to generate 61,024 samples from the golden retriever class.<sup>4</sup> These samples have true intrinsic dimension at most 120, and their estimated intrinsic dimension through (2) is 22 (we rounded the real-valued estimate up). We also generate another 61,024 samples from this model, except we zero out all but  $m = 1$  latent coordinates before passing them through the generator, resulting in samples of true intrinsic dimension of at most  $m$ , and estimated intrinsic dimension 3. We then combine these two types of samples into a single dataset of size 122,048, with each type of sample corresponding to a class. This dataset has the property that its classes have a large difference in intrinsic dimension, and estimating the intrinsic dimension of the entire dataset while ignoring classes gives a value of 5. We then train two models on this dataset: a (supervised) clustered version C-WAE+NF (constant  $\hat{d}$ ) where both clusters have their latent dimension set to 5 (i.e. the single estimate of intrinsic dimension over the entire dataset), and a (supervised) clustered version C-WAE+NF where clusters have their latent dimensions set to 22 and 3 (i.e. class intrinsic dimension estimates), respectively.

We then repeat this entire process for  $m = 3$  and  $m = 5$ , and show the resulting FID scores in Figure 4 for 2 runs (see caption for the intrinsic dimension estimates we obtained and used for the  $m = 3$  and  $m = 5$  cases). We can see that setting the latent dimension of each cluster to its estimated intrinsic dimension results in a large performance gap versus setting it without accounting for varying intrinsic dimensions. We can also see that this gap in performance tightens as the difference in true intrinsic dimensions decreases. These results show that considering varying intrinsic dimension is practically relevant in DGMs, particularly when the difference in true intrinsic dimensions is large.

## 6 Conclusions, limitations, and future work

**Conclusion and broader impact** In this work, we have proposed the *union of manifolds hypothesis*, which extends the standard manifold hypothesis [9] to accommodate data of varying intrinsic dimension. We first support this hypothesis by showing that intrinsic dimensionality can vary widely

<sup>3</sup>For positive  $\sigma \rightarrow 0$ , Loaiza-Ganem et al. [46] show that the likelihood of a  $D$ -dimensional model can go to infinity at a rate of  $\sigma^{d-D}$ , while not recovering  $\mathbb{P}^*$ , if the data is supported on a  $d$ -dimensional manifold.

<sup>4</sup>This odd-looking number is simply a by-product of the default batch size of the pretrained model.

across classes of standard image datasets. We also establish that classes with higher intrinsic dimension are harder to classify, refining the dataset-level result of Pope et al. [55]. We anticipate that the broader deep learning community will further unlock the potential of the union of manifolds hypothesis for both understanding natural data and discovering empirical improvements.

We then propose our own empirical innovation – the *clustered DGMs* modelling framework – to improve the performance of DGMs on image datasets satisfying the union of manifolds hypothesis. By first clustering the data and then learning a two-step manifold-based density estimator [46] on each cluster, we show that we can not only better model the data by varying the latent dimensions, but also by accounting for the possibility of multiple *disjoint* manifolds. We show theoretically that baseline pushforward models do not have this capacity, and empirically that our models outperform these baselines on image generation. We lastly design an experiment where we have some control over the intrinsic dimension of the observed data, and show that our methods provide further outperformance when the variability of intrinsic dimension over the dataset is higher. Due to the nature of our contributions, we do not foresee any potential negative societal consequences.

**Limitations and future work** We have established that our proposed clustered DGMs approach performs admirably on high-dimensional image datasets with multiple distinct connected components of varying intrinsic dimension. However, it is tough for our proposed approach to outperform baselines on a dataset like SVHN where, although we still assume low-dimensional manifold structure, the inter-class intrinsic dimension variability is incredibly low. We conjecture as well for SVHN that the components comprising the dataset may be “closer together” than for the other datasets, which could hamstring the initial clustering algorithm and make the overall learning procedure more challenging. Furthermore, we have not tested on anything besides image data; verification of the union of manifolds hypothesis and the efficacy of clustered DGMs on other types of data remain to be proven. We have reason to believe our improvements will generalize though, as has often been the case with other DGM innovations developed on image-like data.

Another point to note is that we have assumed that we know the number of classes *a priori*, and have used this information to select the number of clusters in our model. While this may often be a reasonable assumption, a *fully* unsupervised approach would not do this. We might be able to relax this assumption by instead first learning the number of clusters [64, 42, 70]; we leave this for future work as well. Additionally, the current hard clustering approach we follow essentially splits the full learning task into separate smaller tasks. While these tasks are easier, there is less data to perform each of them, and we believe there is room to incorporate some of the “mixture of  $G_s$ ” strategies outlined in Section 2, or parameter sharing, while retaining the benefits outlined in Section 4.2.

Lastly, while we have addressed the modelling of data containing several connected components, we have not addressed any other more complicated topological structure which may be present in data. Nonetheless, it is worth noting that all pushforward methods also have this limitation, and although there remains work to be done, we have expanded the topological structures expressible by DGMs.

## References

- [1] M. Abdolali and N. Gillis. Beyond linear subspace clustering: A comparative study of nonlinear manifold clustering algorithms. *Computer Science Review*, 42:100435, nov 2021. doi: 10.1016/j.cosrev.2021.100435. URL <https://doi.org/10.1016%2Fj.cosrev.2021.100435>.
- [2] M. Arbel, L. Zhou, and A. Gretton. Generalized energy based models. *ICLR*, 2021.
- [3] S. Arora, R. Ge, Y. Liang, T. Ma, and Y. Zhang. Generalization and equilibrium in generative adversarial nets (gans). In *International Conference on Machine Learning*, pages 224–232. PMLR, 2017.
- [4] D. Arthur and S. Vassilvitskii. k-means++: The advantages of careful seeding. In *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pages 1027–1035, 2007.
- [5] E. Banijamali, A. Ghodsi, and P. Popuart. Generative mixture of networks. In *2017 International Joint Conference on Neural Networks (IJCNN)*, pages 3753–3760. IEEE, 2017.
- [6] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J.-H. Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, pages 573–582. PMLR, 2019.

- [7] J. Behrmann, P. Vicol, K.-C. Wang, R. Grosse, and J.-H. Jacobsen. Understanding and mitigating exploding inverses in invertible neural networks. In *International Conference on Artificial Intelligence and Statistics*, pages 1792–1800. PMLR, 2021.
- [8] M. Ben-Yosef and D. Weinshall. Gaussian mixture generative adversarial networks for diverse datasets, and the unsupervised clustering of images. *arXiv preprint arXiv:1808.10356*, 2018.
- [9] Y. Bengio, A. Courville, and P. Vincent. Representation learning: A review and new perspectives. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1798–1828, 2013.
- [10] A. Block, Z. Jia, Y. Polyanskiy, and A. Rakhlin. Intrinsic dimension estimation. *arXiv preprint arXiv:2106.04018*, 2021.
- [11] V. I. Bogachev. *Measure theory*, volume 2. Springer, 2007.
- [12] S. Bond-Taylor, A. Leach, Y. Long, and C. 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.
- [13] J. Brehmer and K. Cranmer. Flows for simultaneous manifold learning and density estimation. *Advances in Neural Information Processing Systems*, 33:442–453, 2020.
- [14] A. Brock, J. Donahue, and K. Simonyan. Large scale GAN training for high fidelity natural image synthesis. *ICLR*, 2019.
- [15] T. Cacoullos. Estimation of a multivariate density. *Annals of the Institute of Statistical Mathematics*, 18(1):179–189, 1966.
- [16] L. Cao, S. Asadi, W. Zhu, C. Schmidli, and M. Sjöberg. Simple, scalable, and stable variational deep clustering. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 108–124. Springer, 2020.
- [17] A. L. Caterini, G. Loaiza-Ganem, G. Pleiss, and J. P. Cunningham. Rectangular flows for manifold learning. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [18] R. T. Q. Chen, J. Behrmann, D. K. Duvenaud, and J.-H. Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [19] R. Cornish, A. Caterini, G. Deligiannidis, and A. Doucet. Relaxing bijectivity constraints with continuously indexed normalising flows. In *International conference on machine learning*, pages 2133–2143. PMLR, 2020.
- [20] E. Cunningham, A. Cobb, and S. Jha. Principal manifold flows. *arXiv preprint arXiv:2202.07037*, 2022.
- [21] B. Dai and D. Wipf. Diagnosing and enhancing VAE models. *ICLR*, 2019.
- [22] V. De Bortoli, E. Mathieu, M. Hutchinson, J. Thornton, Y. W. Teh, and A. Doucet. Riemannian score-based generative modeling. *arXiv preprint arXiv:2202.02763*, 2022.
- [23] N. Dilokthanakul, P. A. Mediano, M. Garnelo, M. C. Lee, H. Salimbeni, K. Arulkumaran, and M. Shanahan. Deep unsupervised clustering with gaussian mixture variational autoencoders. *ICLR*, 2017.
- [24] L. Dinh, J. Sohl-Dickstein, and S. Bengio. Density estimation using Real NVP. *ICLR*, 2017.
- [25] L. Dinh, J. Sohl-Dickstein, H. Larochelle, and R. Pascanu. A RAD approach to deep mixture models. *arXiv preprint arXiv:1903.07714*, 2019.
- [26] Y. Du and I. Mordatch. Implicit generation and modeling with energy based models. *Advances in Neural Information Processing Systems*, 32:3608–3618, 2019.
- [27] L. L. Duan. Transport Monte Carlo: High-accuracy posterior approximation via random transport. *Journal of the American Statistical Association*, (just-accepted):1–30, 2021.
- [28] C. Durkan, A. Bekasov, I. Murray, and G. Papamakarios. Neural Spline Flows. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- [29] E. Elhamifar and R. Vidal. Sparse manifold clustering and embedding. In J. Shawe-Taylor, R. Zemel, P. Bartlett, F. Pereira, and K. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc., 2011. URL <https://proceedings.neurips.cc/paper/2011/file/fc490ca45c00b1249bbe3554a4fdf6fb-Paper.pdf>.

- [30] E. Elhamifar and R. Vidal. Sparse subspace clustering: Algorithm, theory, and applications, 2012. URL <https://arxiv.org/abs/1203.1005>.
- [31] E. Elhamifar and R. Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence*, 35(11):2765–2781, 2013.
- [32] M. C. Gemicci, D. Rezende, and S. Mohamed. Normalizing flows on riemannian manifolds. *arXiv preprint arXiv:1611.02304*, 2016.
- [33] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.
- [34] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778. IEEE, 2016.
- [35] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. 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.
- [36] P. Izmailov, P. Kirichenko, M. Finzi, and A. G. Wilson. Semi-supervised learning with normalizing flows. In *International Conference on Machine Learning*, pages 4615–4630. PMLR, 2020.
- [37] Z. Jiang, Y. Zheng, H. Tan, B. Tang, and H. Zhou. Variational deep embedding: an unsupervised and generative approach to clustering. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence*, pages 1965–1972, 2017.
- [38] D. P. Kingma and P. Dhariwal. Glow: Generative Flow with Invertible  $1 \times 1$  Convolutions. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- [39] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. *ICLR*, 2014.
- [40] K. Kothari, A. Khorashadizadeh, M. de Hoop, and I. Dokmanić. Trumpets: Injective flows for inference and inverse problems. In *Proceedings of the Thirty-Seventh Conference on Uncertainty in Artificial Intelligence*, volume 161, pages 1269–1278, 2021.
- [41] A. Krizhevsky, G. Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- [42] B. Kulis and M. I. Jordan. Revisiting k-means: new algorithms via bayesian nonparametrics. In *Proceedings of the 29th International Conference on International Conference on Machine Learning*, pages 1131–1138, 2012.
- [43] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [44] E. Levina and P. Bickel. Maximum likelihood estimation of intrinsic dimension. *Advances in neural information processing systems*, 17, 2004.
- [45] U. Lim, V. Nanda, and H. Oberhauser. Tangent space and dimension estimation with the wasserstein distance. *arXiv preprint arXiv:2110.06357*, 2021.
- [46] G. Loaiza-Ganem, B. L. Ross, J. C. Cresswell, and A. L. Caterini. Diagnosing and fixing manifold overfitting in deep generative models. *arXiv preprint arXiv:2204.07172*, 2022.
- [47] F. Locatello, D. Vincent, I. Tolstikhin, G. Rätsch, S. Gelly, and B. Schölkopf. Competitive training of mixtures of independent deep generative models. *arXiv preprint arXiv:1804.11130*, 2018.
- [48] D. J. MacKay and Z. Ghahramani. Comments on ‘maximum likelihood estimation of intrinsic dimension’ by e. levina and p. bickel (2004). *The Inference Group Website, Cavendish Laboratory, Cambridge University*, 2005.
- [49] E. Mathieu and M. Nickel. Riemannian continuous normalizing flows. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- [50] E. Nalisnick, L. Hertel, and P. Smyth. Approximate inference for deep latent gaussian mixtures. In *NIPS Workshop on Bayesian Deep Learning*, volume 2, page 131, 2016.
- [51] H. Narayanan and S. Mitter. Sample complexity of testing the manifold hypothesis. *Advances in neural information processing systems*, 23, 2010.

- [52] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [53] A. Ozakin and A. Gray. Submanifold density estimation. *Advances in Neural Information Processing Systems*, 22, 2009.
- [54] G. G. F. Pires and M. A. Figueiredo. Variational mixture of normalizing flows. In *ESANN*, 2020.
- [55] P. Pope, C. Zhu, A. Abdelkader, M. Goldblum, and T. Goldstein. The intrinsic dimension of images and its impact on learning. In *International Conference on Learning Representations*, 2021.
- [56] A. Potapczynski, G. Loaiza-Ganem, and J. P. Cunningham. Invertible gaussian reparameterization: Revisiting the gumbel-softmax. *Advances in Neural Information Processing Systems*, 33:12311–12321, 2020.
- [57] D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.
- [58] D. J. Rezende, G. Papamakarios, S. Racaniere, M. Albergo, G. Kanwar, P. Shanahan, and K. Cranmer. Normalizing flows on tori and spheres. In *International Conference on Machine Learning*, pages 8083–8092. PMLR, 2020.
- [59] L. Rokach and O. Maimon. Clustering methods. In *Data mining and knowledge discovery handbook*, pages 321–352. Springer, 2005.
- [60] B. L. Ross and J. C. Cresswell. Tractable density estimation on learned manifolds with conformal embedding flows. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- [61] B. L. Ross, G. Loaiza-Ganem, A. L. Caterini, and J. C. Cresswell. Neural implicit manifold learning for topology-aware generative modelling. *arXiv preprint arXiv:2206.11267*, 2022.
- [62] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [63] A. Salmona, V. de Bortoli, J. Delon, and A. Desolneux. Can push-forward generative models fit multimodal distributions? *arXiv preprint arXiv:2206.14476*, 2022.
- [64] S. Salvador and P. Chan. Determining the number of clusters/segments in hierarchical clustering/segmentation algorithms. In *16th IEEE international conference on tools with artificial intelligence*, pages 576–584. IEEE, 2004.
- [65] R. L. Schilling and F. Kühn. *Counterexamples in Measure and Integration*. Cambridge University Press, 2021.
- [66] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *ICLR*, 2015.
- [67] M. Śmieja, M. Wołczyk, J. Tabor, and B. C. Geiger. Segma: Semi-supervised gaussian mixture autoencoder. *IEEE transactions on neural networks and learning systems*, 32(9):3930–3941, 2020.
- [68] P. Tempczyk, R. Michaluk, Ł. Garncarek, P. Spurek, J. Tabor, and A. Goliński. Lidl: Local intrinsic dimension estimation using approximate likelihood. *arXiv preprint arXiv:2206.14882*, 2022.
- [69] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders. *ICLR*, 2018.
- [70] V. Villegas, H. J. Braviner, P. Naderian, C. J. Maddison, and G. Loaiza-Ganem. Bayesian nonparametrics for offline skill discovery. *arXiv preprint arXiv:2202.04675*, 2022.
- [71] J. H. Ward Jr. Hierarchical grouping to optimize an objective function. *Journal of the American statistical association*, 58(301):236–244, 1963.
- [72] H. Xiao, K. Rasul, and R. Vollgraf. Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.
- [73] F. Ye and A. G. Bors. Lifelong mixture of variational autoencoders. *IEEE Transactions on Neural Networks and Learning Systems*, 2021.

## A Proof of Proposition 1

Before giving the standard definition of support, we motivate it. Intuitively, the support of a distribution  $\mathbb{P}$  is “the smallest set” to which  $\mathbb{P}$  assigns probability 1. The idea of “the smallest set” satisfying a property is often formalized by taking the intersection of all sets satisfying the property. Naïve application of this principle however results in a “definition” of support with undesirable properties. Consider for example the case where  $\mathbb{P}$  is a standard Gaussian distribution on  $\mathbb{R}$ . Clearly any definition of support that we use should obey that the support of  $\mathbb{P}$  is  $\mathbb{R}$ . However, if we simply follow the “definition” of support as the intersection of all sets of probability 1, we would not obtain this: for any  $x \in \mathbb{R}$ , clearly  $\mathbb{P}(\mathbb{R} \setminus \{x\}) = 1$ , yet  $\bigcap_{x \in \mathbb{R}} \mathbb{R} \setminus \{x\} = \emptyset$ , where  $\emptyset$  denotes the empty set. In other words, arbitrary intersections of sets of probability 1 need not have probability 1, and in fact can have probability 0. More abstractly, sets of probability 1 are not closed under arbitrary intersections, which is desirable when defining “the smallest set” satisfying a given property as the intersection of all such sets. In order to solve this problem, the support is defined as the intersection of all closed sets having probability 1 (closed sets are indeed closed under intersections).

**Definition 1 (Support of probability distributions):** Let  $\mathcal{Z}$  be a topological space. Equipping  $\mathcal{Z}$  with its Borel  $\sigma$ -algebra to make it a measurable space, let  $\mathbb{P}$  be a probability distribution on  $\mathcal{Z}$ . Let  $\mathcal{C}(\mathbb{P})$  be the collection of all closed sets  $C$  in  $\mathcal{Z}$  such that  $\mathbb{P}(C) = 1$ . The support of  $\mathbb{P}$ , denoted as  $\text{supp}(\mathbb{P})$ , is defined as:

$$\text{supp}(\mathbb{P}) = \bigcap_{C \in \mathcal{C}(\mathbb{P})} C. \quad (5)$$

Before continuing, we point out that there exist counterexamples showing that the support of a distribution need not have probability 1 [65, Examples 6.2 and 6.3]. In other words, even though closed sets are closed under intersections, closed sets of probability 1 need not share this property. These examples are however much more contrived than the one above and of no practical interest (if  $\mathcal{Z}$  is a separable metric space, as is the case of  $\mathbb{R}^d$ , then any Borel measure  $\mathbb{P}$  on  $\mathcal{Z}$  satisfies  $\mathbb{P}(\text{supp}(\mathbb{P})) = 1$ , see Bogachev [11, Proposition 7.2.9]), and thus support is commonly defined as above. Note that the assumption that  $\mathbb{P}_Z(\text{supp}(\mathbb{P}_Z)) = 1$  could be replaced by the requirement that  $\mathcal{Z}$  be a separable metric space and  $\mathbb{P}_Z$  a Borel measure on it (which covers all settings of practical interest), although this would yield a slightly less general result. We restate Proposition 1 below for convenience:

**Proposition 1:** Let  $\mathcal{Z}$  and  $\mathcal{X}$  be topological spaces and  $G : \mathcal{Z} \rightarrow \mathcal{X}$  be continuous. Considering  $\mathcal{Z}$  and  $\mathcal{X}$  as measurable spaces with their respective Borel  $\sigma$ -algebras, let  $\mathbb{P}_Z$  be a probability measure on  $\mathcal{Z}$  such that  $\text{supp}(\mathbb{P}_Z)$  is connected and  $\mathbb{P}_Z(\text{supp}(\mathbb{P}_Z)) = 1$ . Then  $\text{supp}(G_\# \mathbb{P}_Z)$  is connected.

**Proof:** As mentioned in the proof sketch in the main manuscript, showing that  $\text{supp}(G_\# \mathbb{P}_Z) = \text{cl}(G(\text{supp}(\mathbb{P}_Z)))$ , where  $\text{cl}(\cdot)$  denotes closure in  $\mathcal{X}$ , is enough, since by continuity of  $G$ ,  $G(\text{supp}(\mathbb{P}_Z))$  is connected, and the closure of connected sets is itself connected.

We begin by verifying that  $\text{cl}(G(\text{supp}(\mathbb{P}_Z))) \subseteq \text{supp}(G_\# \mathbb{P}_Z)$ . Let  $C \in \mathcal{C}(G_\# \mathbb{P}_Z)$ . By definition of  $\mathcal{C}(G_\# \mathbb{P}_Z)$ ,  $C$  is closed in  $\mathcal{X}$  and is such that  $G_\# \mathbb{P}_Z(C) = 1$ . By definition of pushforward measure, it follows that:

$$\mathbb{P}_Z(G^{-1}(C)) = G_\# \mathbb{P}_Z(C) = 1, \quad (6)$$

and by continuity of  $G$ , it also follows that  $G^{-1}(C)$  is closed in  $\mathcal{Z}$ . Then,  $G^{-1}(C) \in \mathcal{C}(\mathbb{P}_Z)$ , and by definition of support,  $\text{supp}(\mathbb{P}_Z) \subseteq G^{-1}(C)$ . It then follows that  $G(\text{supp}(\mathbb{P}_Z)) \subseteq G(G^{-1}(C)) \subseteq C$ . Since this holds for every  $C \in \mathcal{C}(G_\# \mathbb{P}_Z)$ , we have that:

$$G(\text{supp}(\mathbb{P}_Z)) \subseteq \bigcap_{C \in \mathcal{C}(G_\# \mathbb{P}_Z)} C = \text{supp}(G_\# \mathbb{P}_Z). \quad (7)$$

Applying  $\text{cl}(\cdot)$  on both sides, and using the fact that  $\text{supp}(G_\# \mathbb{P}_Z)$  is closed in  $\mathcal{X}$  (since it is an intersection of closed sets), yields  $\text{cl}(G(\text{supp}(\mathbb{P}_Z))) \subseteq \text{supp}(G_\# \mathbb{P}_Z)$ .

It now only remains to show that  $\text{supp}(G_\# \mathbb{P}_Z) \subseteq \text{cl}(G(\text{supp}(\mathbb{P}_Z)))$ . By definition of pushforward measure, we have that:

$$G_\# \mathbb{P}_Z(G(\text{supp}(\mathbb{P}_Z))) = \mathbb{P}_Z(G^{-1}(G(\text{supp}(\mathbb{P}_Z)))) \geq \mathbb{P}_Z(\text{supp}(\mathbb{P}_Z)) = 1, \quad (8)$$

where the inequality follows from  $\text{supp}(\mathbb{P}_Z) \subseteq G^{-1}(G(\text{supp}(\mathbb{P}_Z)))$ . Since  $G(\text{supp}(\mathbb{P}_Z)) \subseteq \text{cl}(G(\text{supp}(\mathbb{P}_Z)))$ , then  $G_\# \mathbb{P}_Z(\text{cl}(G(\text{supp}(\mathbb{P}_Z)))) = 1$ . It follows that  $\text{cl}(G(\text{supp}(\mathbb{P}_Z))) \in \mathcal{C}(G_\# \mathbb{P}_Z)$ , and thus  $\text{supp}(G_\# \mathbb{P}_Z) \subseteq \text{cl}(G(\text{supp}(\mathbb{P}_Z)))$  by definition of support.  $\square$

## B Memory-efficient sampling of clustered DGMs

**Algorithm 2:** Sampling of clustered DGMs

---

**Input:**  $m$ , trained clustered DGM  $\{(\mathbb{P}_Z^{(\ell)}, G_\ell)\}_{\ell=1}^L$ , and corresponding cluster sizes  $|\mathcal{D}_1|, \dots, |\mathcal{D}_L|$ .  
**Output:**  $\mathcal{S}$

```

1  $\mathcal{S} \leftarrow \emptyset$ 
2  $(m_1, \dots, m_L) \sim \text{Multinomial}\left(m, \left(\frac{|\mathcal{D}_1|}{\sum_{\ell'=1}^L |\mathcal{D}_{\ell'}|}, \dots, \frac{|\mathcal{D}_L|}{\sum_{\ell'=1}^L |\mathcal{D}_{\ell'}|}\right)\right)$ 
3 for  $\ell = 1$  to  $L$  do
4   for  $t = 1$  to  $m_\ell$  do
5      $Z \sim \mathbb{P}_Z^{(\ell)}$ 
6      $X = G_\ell(Z)$ 
7      $\mathcal{S} \leftarrow \mathcal{S} \cup \{X\}$ 
8   end
9 end

```

---

As mentioned in the main manuscript, Algorithm 1 can be trivially implemented in a memory-efficient way. Analogously, the sampling process described in (4) can be made memory-efficient. We now describe how to achieve this. Assume we want to generate  $m$  samples from a trained clustered DGM model. The idea is to make sure that each of the  $L$  DGMs are only loaded into memory once, which requires first knowing exactly how many samples will be required from each model. Note that, if we denote the number of samples coming from the  $\ell^{\text{th}}$  model as  $m_\ell$ , then clearly  $(m_1, \dots, m_L)$  follows a multinomial distribution with parameters  $m$  and  $(|\mathcal{D}_1|, \dots, |\mathcal{D}_L|) / \sum_{\ell'=1}^L |\mathcal{D}_{\ell'}|$ . We can thus first sample from this multinomial, and then sample the appropriate number of times from each model, as described in Algorithm 2. Note that, similarly to Algorithm 1, the outer **for** loop in Algorithm 2 can be trivially made memory-efficient by loading the required model  $(\mathbb{P}_Z^{(\ell)}, G_\ell)$  into memory and then removing it from memory once it has been sampled from  $m_\ell$  times (and the inner **for** loop can be trivially parallelized).

## C Experimental details

For this project, we trained all models discussed using internal cluster GPUs and a small number of runs on Tesla T4 GPUs rented on Google Cloud Platform. The internal GPUs were a combination of Tesla V100, Quadro RTX 6000 and TITAN V GPUs. We trained 6 different DGMs (including 2 two-step models) over 5 datasets in 3-4 configurations (e.g. baseline, GMM, and clustered) 3 times each. Each run, took between 1-48 GPU hours with a mean time of approximately 8 hours. We also trained 4 different variants of classification models on the CIFAR-100 dataset 5 times, each consuming around 3 GPU hours. In total, this project took on the order of 3,000 GPU hours to complete.

### C.1 Verifying the union of manifolds hypothesis

To estimate the intrinsic dimension of each class, we run the MLE estimator from (2) on all of the datapoints for every class separately. Note that our estimator slightly differs from the one used by Pope et al. [55], who use  $k - 2$  instead of  $k - 1$  in the denominator (while Pope et al. [55] do not explicitly mention this, it is the version of the estimator implemented in their code<sup>5</sup>) as an asymptotic bias correction [44]. Using  $k - 1$ , we observed more consistent values for the intrinsic dimension estimates across  $k$ s and in both versions we observed the same trends of different classes having a large variance in intrinsic dimension estimates. We also note that the asymptotic correction was designed for large values of  $k$ , which we are not using. For a more detailed view of the results presented in the main manuscript, we show the intrinsic dimension estimates for classes in each dataset for five different  $k$  values in Appendix D.1.

### C.2 On the relationship between class intrinsic dimension and classification accuracy

In order to investigate if there is a correlation between intrinsic dimension and the difficulty of classification for a given class, we train three different classifiers on CIFAR-100 and measure the accuracy for each class. We use the same hyperparameters and training strategy for each of the three models we consider (VGG-19 [66], ResNet-18, and ResNet-34 [34]). We do not modify the model architecture except for changing the dimension of the last linear layer to project into a 100-dimensional vector to match the number of classes in CIFAR-100. We use the CIFAR-100 dataset test split (10,000 images) then randomly select 10% of the training dataset as our validation set (5,000 images) and use the remaining 90% as our train split. Note that this validation split is

<sup>5</sup><https://github.com/ppope/dimensions>

constant across all runs. The classification accuracy is reported as the accuracy on the test set on the same epoch where the best validation accuracy is achieved. We do not add any data augmentations as doing so could modify intrinsic dimension estimates and confound our results. Studying the relationship between intrinsic dimension and data augmentations is an interesting direction for future work. The combination of (1) no data augmentations, and (2) the presence of a validation set causes our results to be lower than the state-of-the-art for the models, although the results remain self-consistent here. We train each model using the cross entropy loss for a total of 200 epochs starting with a learning rate of 0.1 and decrease this to  $10^{-2}$ ,  $4 \times 10^{-3}$ , and  $8 \times 10^{-4}$  after 60, 120, and 160 epochs respectively. We use the standard SGD optimizer with momentum set to 0.9 and train with a weight decay of  $5 \times 10^{-4}$ . For all datasets, we normalize by the channel-wise mean and standard deviation of all training images. We use a batch size of 128 for training, and a batch size of 100 for both validation and test.

We now detail the exact loss we used for our experiments with the re-weighted cross entropy loss. A weighted version of the categorical cross entropy loss is given by:

$$-\sum_{i=1}^n \sum_{\ell=1}^L \omega_\ell \cdot y_{i,\ell} \cdot \log f_\theta(x_i)_\ell, \quad (9)$$

where  $y_i$  is a one-hot vector of length  $L$  corresponding to the label of  $x_i$ ,  $f_\theta(x_i)$  is the  $L$ -dimensional output of the classifier containing assigned class probabilities, and  $\omega_\ell$  is the scalar weight given to the  $\ell^{th}$  class. The standard categorical cross entropy uses the weights  $\omega_\ell = 1$  for  $\ell = 1, \dots, L$ . Our modified weights are given by:

$$\omega_\ell = L \cdot \frac{\hat{d}_k^{(\ell)}}{\sum_{\ell'=1}^L \hat{d}_k^{(\ell')}}. \quad (10)$$

Note that when the intrinsic dimension estimates of all classes match, our proposed weights recover the standard ones, but more heavily weight classes of higher intrinsic dimension otherwise.

### C.3 Agglomerative clustering and Ward’s criterion

In order to generate the dataset clusters to train each component in our clustered models (with the goal of approximating the connected components of the full dataset) without labels, we perform agglomerative clustering on the dataset. At the start of the algorithm, each datapoint is in its own cluster. Then, at every step, the two clusters with the smallest linkage distance are merged, decreasing the number of clusters by 1. This occurs until a pre-specified number of clusters,  $L$ , is reached. For all of our datasets, we set  $L = 10$  which is the number of classes for the MNIST, FMNIST, SVHN, and CIFAR-10 datasets. We experimented with ranges of  $L$  between 7 and 15 and did not find a large variation in performance. Automatically inferring  $L$  from a dataset is the subject of future work and would make our clustering models completely unsupervised. For all of our experiments involving agglomerative clustering, we used Ward’s linkage criterion [71]. The distance between two clusters with this criterion is the variance of the Euclidean distance between all datapoints in the clusters being merged. Therefore, at each step, the two clusters with the smallest variance will be combined. We experimented with using single linkage, which defines the distance between two clusters as the minimum distance between a datapoint in both clusters, but found that this led to very unbalanced cluster sizes and had very poor performance when training our DGMs. We also tried a custom linkage metric that merges the two clusters at each step to maximize the inter-cluster intrinsic dimension estimate variance, with the intuition that building clusters with a large variance in intrinsic dimension would better model the submanifolds of the dataset. In practice, we found this method very sensitive to initialization and saw the same unbalanced cluster failure mode as the single linkage criterion.

For all of our experiments, we also tried using the k-means++ [4] clustering algorithm. This also worked well but gave slightly worse results than agglomerative clustering with Ward’s criterion across the board.

### C.4 Better modelling of disconnected components

For all models, we randomly select 10% of the training dataset to be used for validation and train on the remaining 90%. We use a separate test dataset to report all metrics for our models. A batch size of 128 is used for all datasets. Unless otherwise noted, at the beginning of training, we scale all the data to between 0 and 1. For all experiments, we use the Adam optimizer, typically with learning rate 0.001 and cosine annealing for a maximum of 100 epochs. We also use gradient norm clipping with a value of 10.

**VAEs** 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 512 units. For SVHN, CIFAR-100 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. The latent dimension of the model is set to 20. We do not use early stopping and always train for 100 epochs. For the GMM, we set the

prior to have 10 modes with learnable mixture weights, means and standard deviations. The mixture components of the prior were initialized with fixed means spaced 3 units apart centred at 0, standard deviations of 1 and uniform mixture weights.

**WAEs** For all datasets, we use an MLP with 2 hidden layers of size 256 each for the discriminator. The encoder and decoder for all datasets are the same convolutional networks used for VAEs except we increase the channels to (64, 64, 32, 32) and (32, 32, 64, 64) respectively. The latent dimension of the model is also set to 20. We train every model for a total of 300 epochs and perform early stopping on the reconstruction error with a patience of 30 epochs only for the GMM baselines. We note that when training for the full 300 epochs (without early stopping) the GMM baselines perform significantly worse and the other models perform better. We use a learning rate of 0.0005 for the discriminator and do not use any learning rate schedule for the model. We weight the discriminator loss with a lambda value of 10. For the GMM, we set the prior to have 10 modes with fixed means spaced 3 units apart centred at 0, fixed standard deviations of 1 and uniform mixture weights. We attempted to learn these parameters similar to the VAE but observed worse results.

**NFs** We use a rational quadratic spline flow [28] with 64 hidden units, 4 layers, and 3 blocks per layer. For the SVHN, CIFAR-10 and CIFAR-100 datasets, we use a whitening transform at the start of training to make the data zero-mean and marginally unit-variance; note that this affine transformation does not affect the manifold structure of the data. Otherwise we scale the data followed by a logit transform with epsilon set to 0.001. We configure the mixture base distribution in the same way as the VAE prior. We set the initial learning rate to 0.0005.

## C.5 Better modelling of varying intrinsic dimensions

**VAE+NF** For the VAE, we followed the same setup as the single VAE except we keep the learning rate fixed to 0.001 throughout training. The NF density estimator is also identical to the single NF with the following differences: (1) for all datasets, we standardize the data (by subtracting the mean and dividing by the standard deviation) and do not scale the data or apply a logit transform, and (2) we set the learning rate to 0.001.

**WAE+NF** For the WAE, we followed the same setup as the single WAE except we use residual networks to handle the larger and more complex golden retriever datasets. The encoder first applies a convolution (increasing the channel dimension to 16), batch normalization, and ReLU activation to the input before a max pool operation with stride 2 and kernel size of 3. The main body of the encoder is a set of 4 residual layers with channel dimensions (16, 32, 64, 128). Each layer consists of two basic blocks which are two consecutive convolution, BatchNorm and ReLU layers where the input of the block is added to the activation before the final ReLU. The stride of the first convolution of each layer is set to 2 to halve the spatial resolution. The skip connection is downsampled by a convolution with a stride of 2. The activations are then average-pooled across each channel and projected to the latent space dimension by a linear layer. The decoder is identical to the encoder main body except there are no downsampling operations, the spatial resolution of the activation is doubled at each layer by bilinear upsampling and the channel dimensions are (128, 64, 32, 16). We perform early stopping on validation loss with a patience of 30 epochs, for a maximum of 300 epochs. The NF density estimator is also identical to the single NF with the following differences: (1) for all datasets, we standardize the data (by subtracting the mean and dividing by the standard deviation) and do not scale the data or apply a logit transform, (2) we set the learning rate to 0.001, and (3) we perform early stopping on validation loss with a patience of 30 epochs, for a maximum of 100 epochs.

## D Additional experiments

### D.1 Verifying the union of manifolds hypothesis

We show a more granular breakdown of Figure 1, with intrinsic dimension estimate values for each class in Figure 5 for MNIST, Figure 6 for FMNIST, Figure 7 for SVHN, and Figure 8 for CIFAR-10. Since CIFAR-100 and ImageNet have too many classes to show, we only include the top 5 highest and lowest intrinsic dimension classes in Figure 9 for CIFAR-100, and Figure 10 for ImageNet.

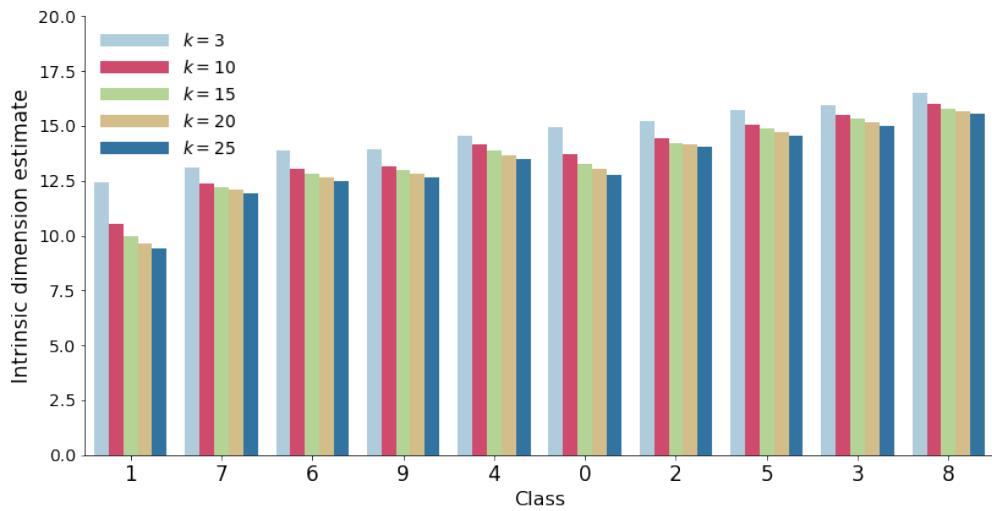


Figure 5: Intrinsic dimension estimates for classes in the MNIST dataset.

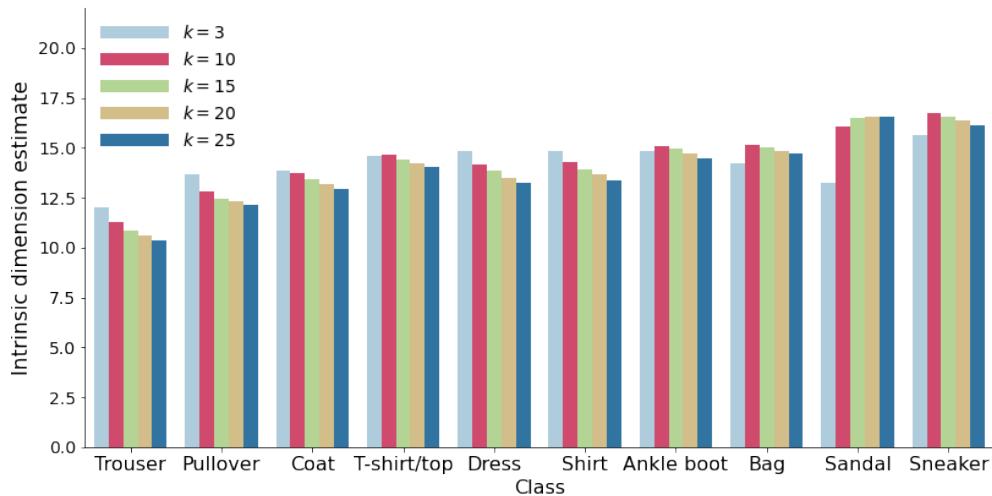


Figure 6: Intrinsic dimension estimates for classes in the FMINST dataset.

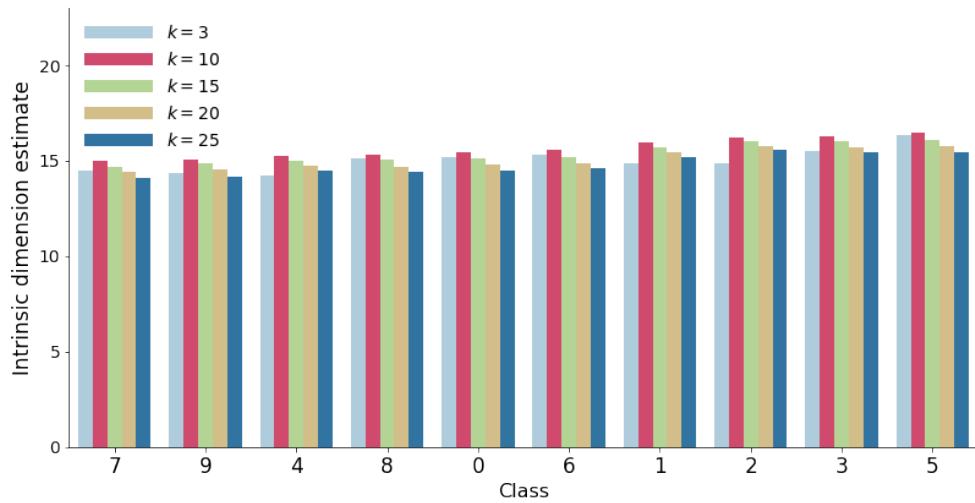


Figure 7: Intrinsic dimension estimates for classes in the SVHN dataset.

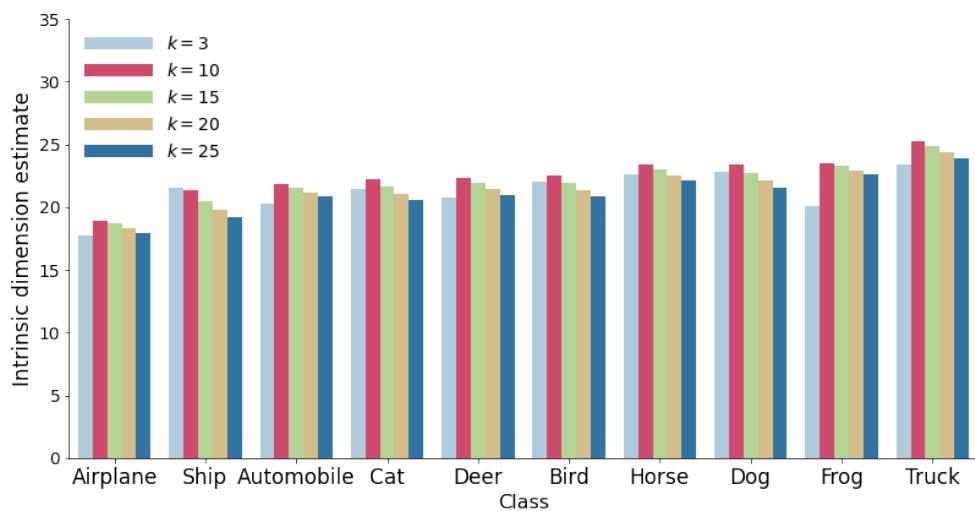


Figure 8: Intrinsic dimension estimates for classes in the CIFAR-10 dataset.

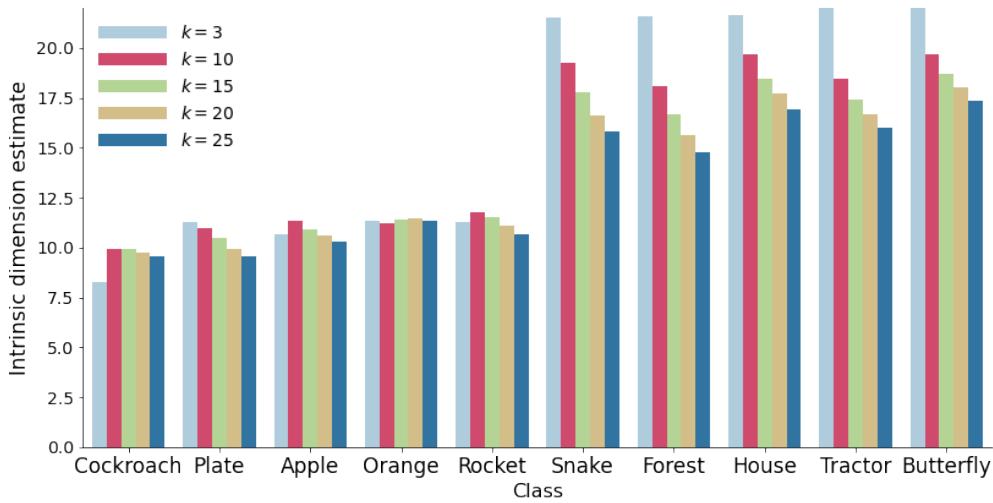


Figure 9: Intrinsic dimension estimates for classes in the CIFAR-100 dataset.

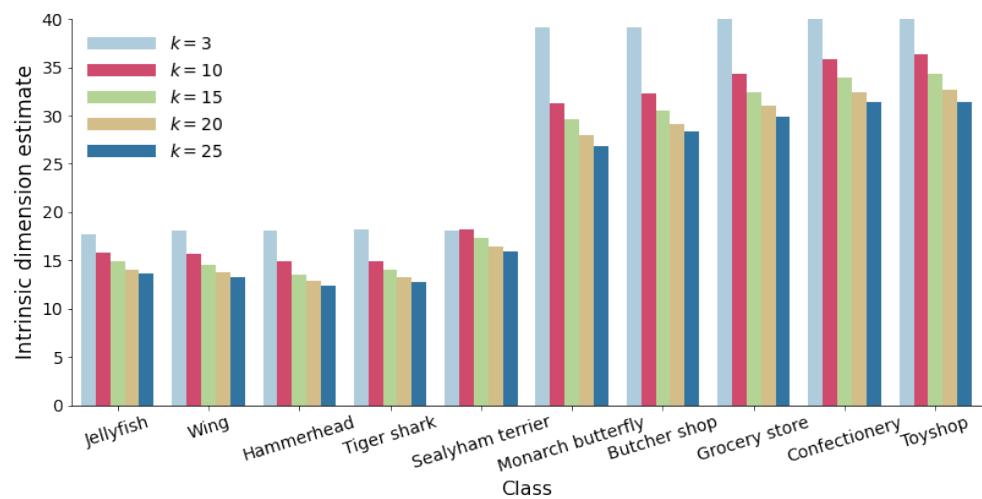


Figure 10: Intrinsic dimension estimates for classes in the ImageNet dataset.

## D.2 Better modelling of disconnected components

Figures 11, 12, and 13 show samples of our VAE, WAE, and NF models, respectively. We note that all models except VAE-based ones, both baseline and clustered, experienced “failure runs” which were of evidently lower quality than regular runs. We suspect that additional hyperparameter tuning would likely remove these instabilities, although we still used these runs to fill Table 2, as they provide valid comparisons between clustered and non-clustered models, and our hyperparameter choices are relatively standard for vanilla models. We reiterate that the performance of the C-VAE (labels) and C-WAE (labels) models on image generation is poor on CIFAR-100 because there are not enough datapoints per class. Similarly, we did not do dequantizing, which is commonly done for many of these models, particularly NFs, so as to not make the data support mathematically connected. Thus, while our comparisons remain fair, our trained NFs have poorer sample quality than is usual. We also note that the C-NF (labels) model on CIFAR-100 was too large to fit in memory, and training and sampling from it were only enabled thanks to our memory-efficient implementations. This highlights the importance of the memory savings previously discussed. Finally, we also note that we did some preliminary experiments with GANs, and did not observe the same benefits of clustered models. Although these results were preliminary and might be fixed through more tuning, we observed clustered GANs were even more susceptible to training instabilities than their non-clustered baselines.

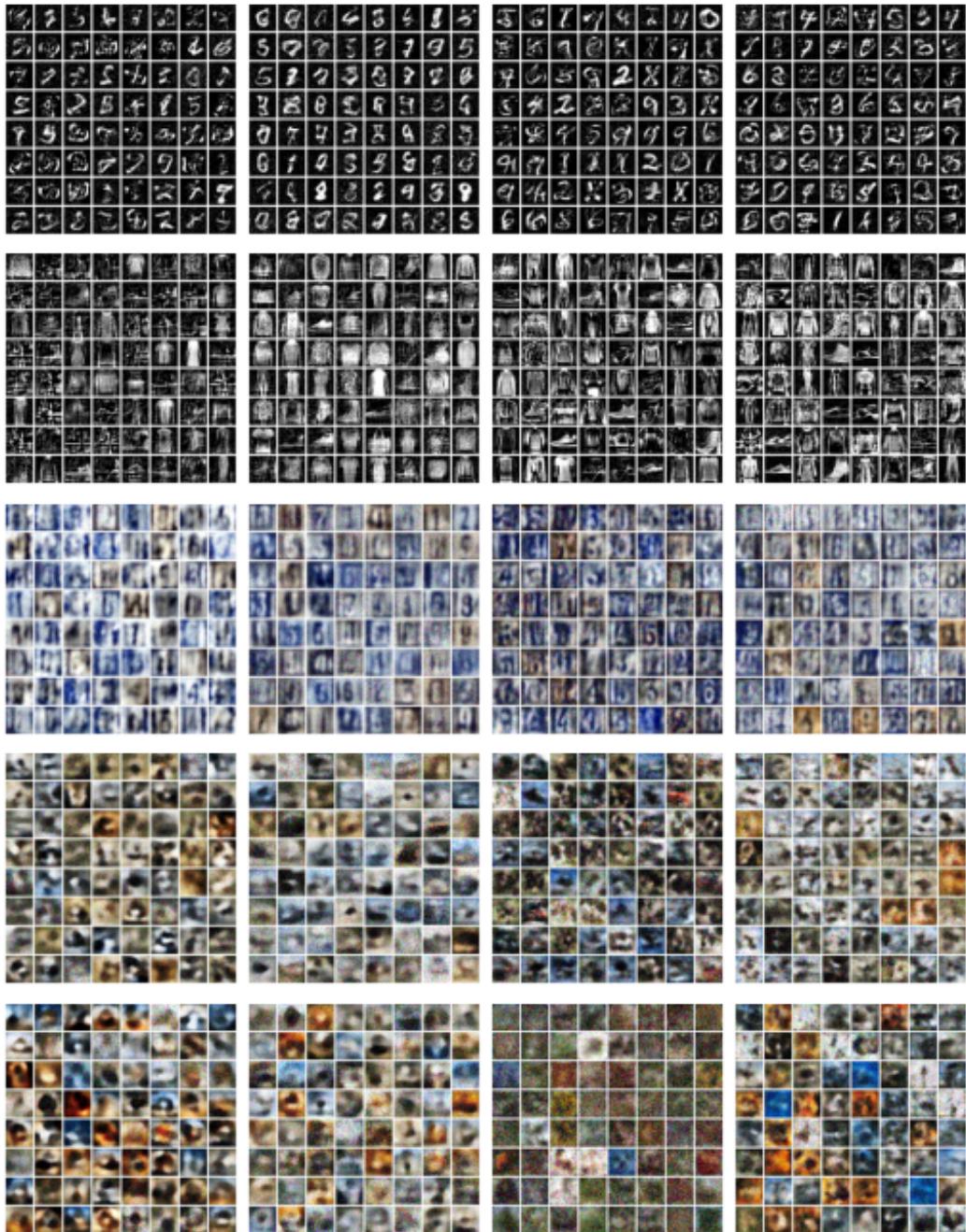


Figure 11: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), CIFAR-10 (**fourth row**), and CIFAR-100 (**fifth row**). Models are VAE (**first column**), VAE (GMM) (**second column**), C-VAE (labels) (**third column**) and C-VAE (**fourth column**).



Figure 12: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), CIFAR-10 (**fourth row**), and CIFAR-100 (**fifth row**). Models are WAE (**first column**), WAE (GMM) (**second column**), C-WAE (labels) (**third column**) and C-WAE (**fourth column**).

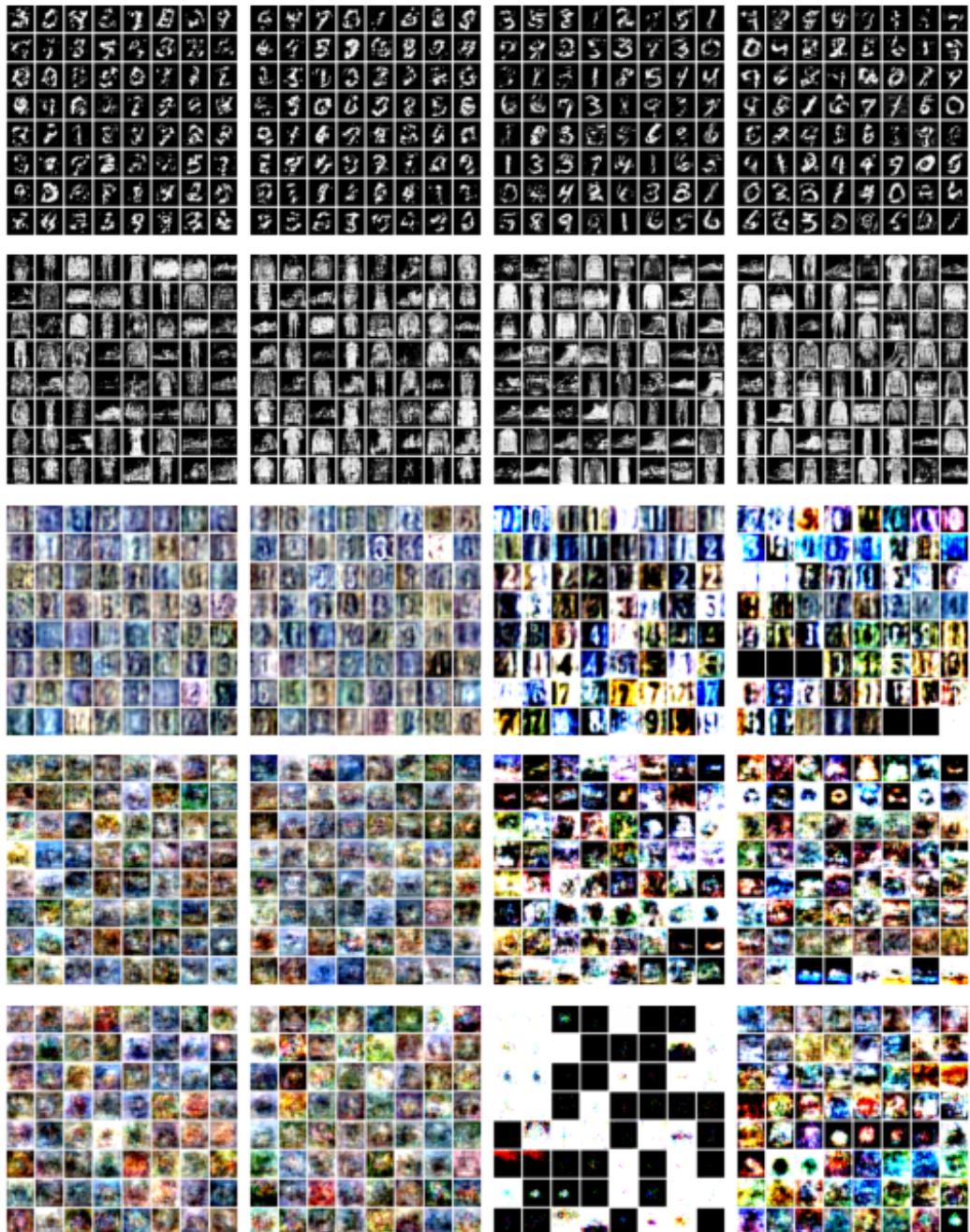


Figure 13: Uncurated samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), CIFAR-10 (**fourth row**), and CIFAR-100 (**fifth row**). Models are NF (**first column**), NF (GMM) (**second column**), C-NF (labels) (**third column**) and C-NF (**fourth column**).

### D.3 Better modelling of varying intrinsic dimensions

Table 3 shows results of running a VAE+NF two-step model. We can see that the best performing clustered versions still mostly outperform the non-clustered version. This is consistent with our previous results and highlights once again the relevance of accounting for multiple connected components. As mentioned in the main manuscript though, there is no significant improvement between fixing the latent dimension of every cluster to its estimated intrinsic dimension, and simply fixing the intrinsic dimension of every cluster to the same value (the estimated intrinsic dimension of the entire dataset). Samples from these models are shown in Figure 14.

As mentioned in the main manuscript, we generate a synthetic dataset where we have control over the difference of intrinsic dimensions between different components of the data. Samples of the resulting models are shown in Figure 15.

Table 3: FID scores. We show means and standard errors across 3 runs. Best models, and those whose error bars overlap with those of the best models, are bolded.

Model	MNIST	FMNIST	SVHN	CIFAR-10	CIFAR-100
VAE+NF	$10.5 \pm 0.1$	$12.4 \pm 0.2$	<b><math>5.5 \pm 0.0</math></b>	$11.2 \pm 0.1$	$10.9 \pm 0.1$
C-VAE+NF constant $\hat{d}$ (labels)	<b><math>9.4 \pm 0.0</math></b>	<b><math>8.7 \pm 0.0</math></b>	<b><math>5.5 \pm 0.1</math></b>	<b><math>10.2 \pm 0.0</math></b>	$11.1 \pm 0.0$
C-VAE+NF constant $\hat{d}$	<b><math>9.4 \pm 0.1</math></b>	<b><math>8.8 \pm 0.1</math></b>	$9.0 \pm 3.3$	$11.5 \pm 0.9$	<b><math>10.0 \pm 0.1</math></b>
C-VAE+NF (labels)	$9.5 \pm 0.0$	$8.8 \pm 0.1$	<b><math>5.6 \pm 0.1</math></b>	$10.5 \pm 0.1$	$11.4 \pm 0.1$
C-VAE+NF	<b><math>9.4 \pm 0.0</math></b>	<b><math>8.7 \pm 0.1</math></b>	$5.7 \pm 0.1$	$13.2 \pm 2.4$	<b><math>9.9 \pm 0.0</math></b>

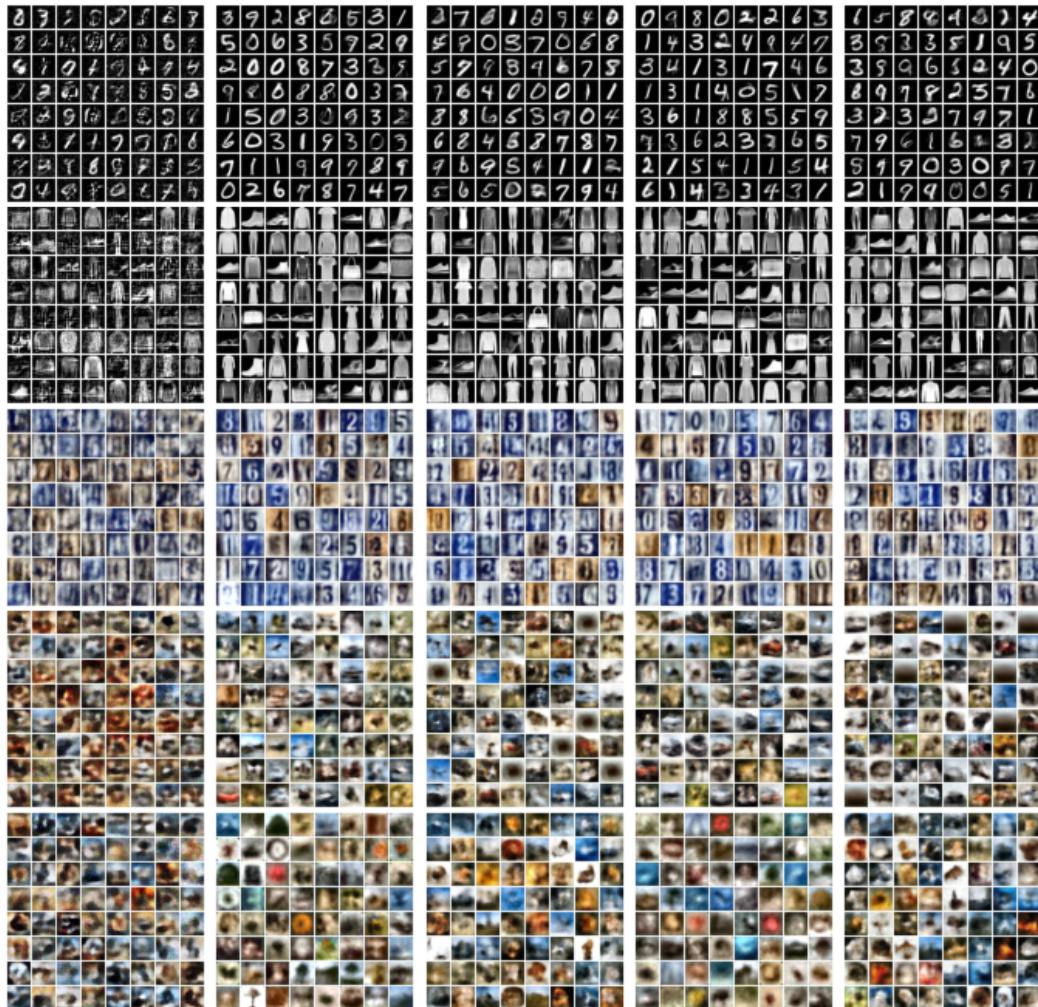


Figure 14: Uncurred samples from models trained on MNIST (**first row**), FMNIST (**second row**), SVHN (**third row**), CIFAR-10 (**fourth row**), and CIFAR-100 (**fifth row**). Models are VAE+NF (**first column**), C-VAE+NF constant  $\hat{d}$  (labels) (**second column**), C-VAE+NF constant  $\hat{d}$  (**third column**), C-VAE+NF (labels) (**fourth column**) and C-VAE+NF (**fifth column**).

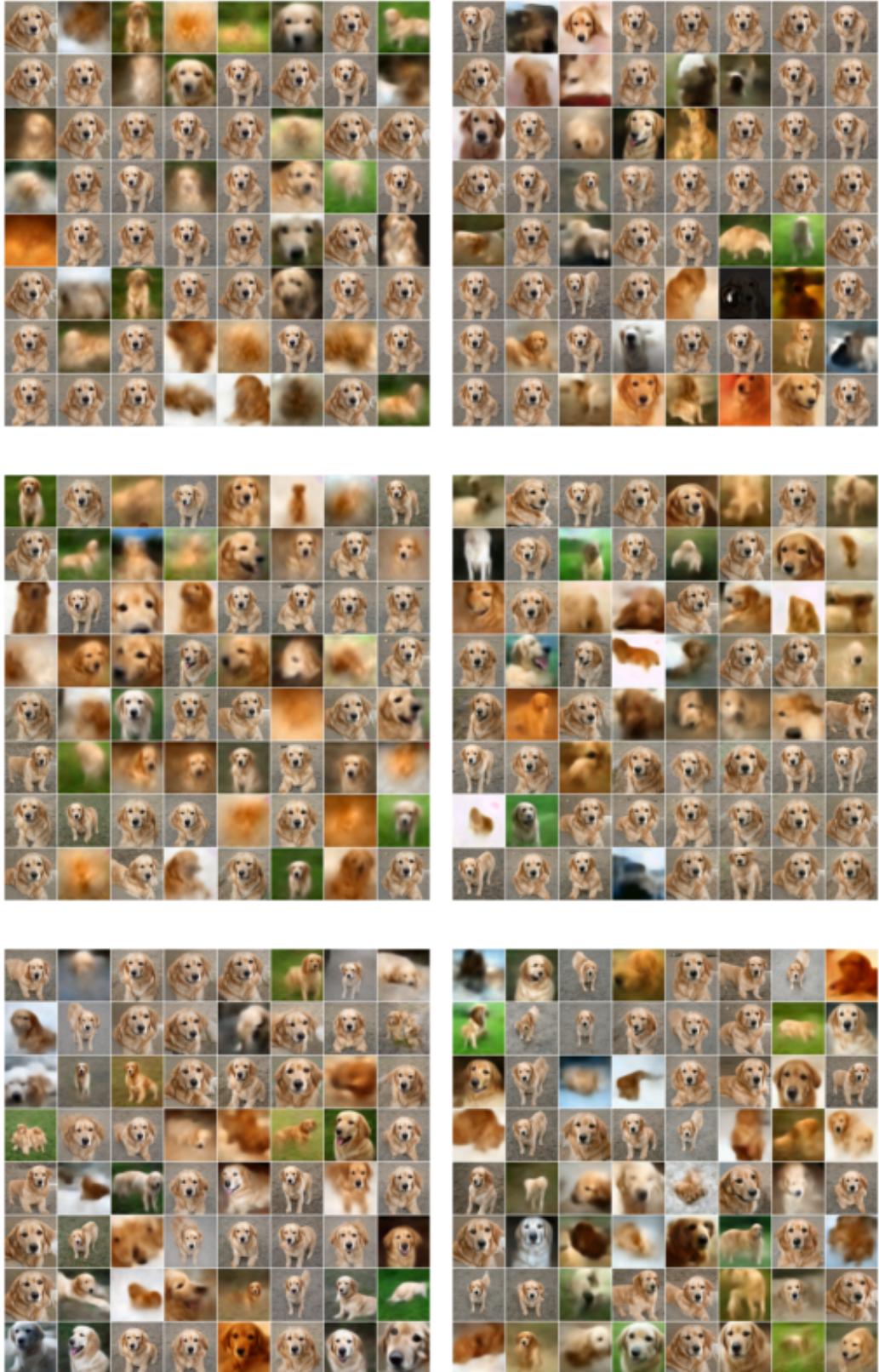


Figure 15: Uncurated samples from models trained on synthetic golden retriever dataset from BigGAN model as outlined in section 5.2. Datasets generated with  $m = 1, 3$  and  $5$  for the first, second and third rows respectively. Models are C-WAE+NF constant  $\hat{d}$  (labels) (**first column**) and C-WAE+NF (labels) (**second column**).