
Diagnosing and Enhancing Gaussian VAE Models

Bin Dai

Institute for Advanced Study, Tsinghua University
daib13@mails.tsinghua.edu.cn

David Wipf

Microsoft Research
davidwipf@gmail.com

Abstract

Although variational autoencoders (VAEs) represent a widely influential deep generative model, many aspects of the underlying energy function remain poorly understood. In particular, it is commonly believed that Gaussian encoder/decoder assumptions reduce the effectiveness of VAEs in generating realistic samples. In this regard, we rigorously analyze the VAE objective, differentiating situations where this belief is and is not actually true. We then leverage the corresponding insights to develop a simple VAE enhancement that requires no additional hyperparameters or sensitive tuning. Quantitatively, this proposal produces crisp samples and stable FID scores that are actually competitive with a variety of GAN models, all while retaining desirable attributes of the original VAE architecture.

1 Introduction

Our starting point is the desire to learn a probabilistic generative model of observable variables $\mathbf{x} \in \chi$, where χ is a r -dimensional manifold embedded in \mathbb{R}^d . Note that if $r = d$, then this assumption places no restriction on the distribution of $\mathbf{x} \in \mathbb{R}^d$ whatsoever; however, the added formalism is introduced to handle the frequently encountered case where \mathbf{x} possesses low-dimensional structure relative to a high-dimensional ambient space, i.e., $r \ll d$. In fact, the very utility of generative models of continuous data, and their attendant low-dimensional representations, often hinges on this assumption (Bengio et al., 2013). It therefore behooves us to explicitly account for this situation.

Beyond this, we assume that χ is a simple Riemannian manifold, which means there exists a diffeomorphism φ between χ and \mathbb{R}^r , or more explicitly, the mapping $\varphi : \chi \mapsto \mathbb{R}^r$ is invertible and differentiable. Denote a ground-truth probability measure on χ as μ_{gt} such that the probability mass of an infinitesimal $d\mathbf{x}$ on the manifold is $\mu_{gt}(d\mathbf{x})$ and $\int_{\chi} \mu_{gt}(d\mathbf{x}) = 1$.

The variational autoencoder (VAE) (Kingma & Welling, 2014; Rezende et al., 2014) attempts to approximate this ground-truth measure using a parameterized density $p_{\theta}(\mathbf{x})$ defined across all of \mathbb{R}^d since any underlying generative manifold is unknown in advance. This density is further assumed to admit the latent decomposition $p_{\theta}(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$, where $\mathbf{z} \in \mathbb{R}^{\kappa}$ serves as a low-dimensional representation, with $\kappa \approx r$ and prior $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$.

Ideally we might like to minimize the negative log-likelihood $-\log p_{\theta}(\mathbf{x})$ averaged across the ground-truth measure μ_{gt} , i.e., solve $\min_{\theta} \int_{\chi} -\log p_{\theta}(\mathbf{x})\mu_{gt}(d\mathbf{x})$. Unfortunately though, the required marginalization over \mathbf{z} is generally infeasible. Instead the VAE model relies on tractable *encoder* $q_{\phi}(\mathbf{z}|\mathbf{x})$ and *decoder* $p_{\theta}(\mathbf{x}|\mathbf{z})$ distributions, where ϕ represents additional trainable parameters. The canonical VAE cost is a bound on the average negative log-likelihood given by

$$\mathcal{L}(\theta, \phi) \triangleq \int_{\chi} \{-\log p_{\theta}(\mathbf{x}) + \mathbb{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x})]\} \mu_{gt}(d\mathbf{x}) \geq \int_{\chi} -\log p_{\theta}(\mathbf{x})\mu_{gt}(d\mathbf{x}), \quad (1)$$

where the inequality follows directly from the non-negativity of the KL-divergence. Here ϕ can be viewed as tuning the tightness of bound, while θ dictates the actual estimation of μ_{gt} . Using a few standard manipulations, this bound can also be expressed as

$$\mathcal{L}(\theta, \phi) = \int_{\chi} \{-\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] + \mathbb{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]\} \mu_{gt}(d\mathbf{x}), \quad (2)$$

which explicitly involves the encoder/decoder distributions and is conveniently amenable to SGD optimization of $\{\theta, \phi\}$ via a reparameterization trick (Kingma & Welling, 2014; Rezende et al., 2014). The first term in (2) can be viewed as a reconstruction cost (or a stochastic analog of a traditional autoencoder), while the second penalizes posterior deviations from the prior $p(z)$. Additionally, for any realizable implementation via SGD, the integration over χ must be approximated via a finite sum across training samples $\{\mathbf{x}^{(i)}\}_{i=1}^n$ drawn from μ_{gt} . Nonetheless, examining the true objective $\mathcal{L}(\theta, \phi)$ can lead to important, practically-relevant insights.

At least in principle, $q_\phi(z|\mathbf{x})$ and $p_\theta(\mathbf{x}|z)$ can be arbitrary distributions, in which case we could simply enforce $q_\phi(z|\mathbf{x}) = p_\theta(z|\mathbf{x}) \propto p_\theta(\mathbf{x}|z)p(z)$ such that the bound from (1) is tight. Unfortunately though, this is essentially always an intractable undertaking. Consequently, largely to facilitate practical implementation, a commonly adopted distributional assumption for continuous data is that both $q_\phi(z|\mathbf{x})$ and $p_\theta(\mathbf{x}|z)$ are Gaussian. This design choice has previously been cited as a key limitation of VAEs (Burda et al., 2015; Kingma et al., 2016), and existing quantitative tests of generative modeling quality thus far dramatically favor contemporary alternatives such as generative adversarial networks (GAN) (Goodfellow et al., 2014). Regardless, because the VAE possesses certain desirable properties relative to GAN models (e.g., stable training (Tolstikhin et al., 2018), interpretable encoder/inference network (Brock et al., 2016), outlier-robustness (Dai et al., 2018), etc.), it remains a highly influential paradigm worthy of examination and enhancement.

In Section 2 we closely investigate the implications of VAE Gaussian assumptions leading to a number of interesting diagnostic conclusions. In particular, we differentiate the situation where $r = d$, in which case we prove that recovering the ground-truth distribution is actually possible iff the VAE global optimum is reached, and $r < d$, in which case the VAE global optimum can be reached by solutions that reflect the ground-truth distribution almost everywhere, but not necessarily uniquely so. In other words, there could exist alternative solutions that both reach the global optimum and yet do not assign the same probability measure as μ_{gt} .

Section 3 then further probes this non-uniqueness issue by inspecting necessary conditions of global optima when $r < d$. This analysis reveals that an optimal VAE parameterization will provide an encoder/decoder pair capable of perfectly reconstructing all $\mathbf{x} \in \chi$ using *any* z drawn from $q_\phi(z|\mathbf{x})$. Moreover, we demonstrate that the VAE accomplishes this using a degenerate latent code whereby only r dimensions are effectively active. Collectively, these results indicate that the VAE global optimum can in fact uniquely learn a mapping to the correct ground-truth manifold when $r < d$, but not necessarily the correct probability measure *within* this manifold, a critical distinction.

Next we leverage these analytical results in Section 4 to motivate an almost trivially-simple, two-stage VAE enhancement for addressing typical regimes when $r < d$. In brief, the first stage just learns the manifold per the allowances from Section 3, and in doing so, provides a mapping to a lower dimensional intermediate representation with no degenerate dimensions that mirrors the $r = d$ regime. The second (much smaller) stage then only needs to learn the correct probability measure on this intermediate representation, which is possible per the analysis from Section 2. Experiments from Section 5 reveal that this procedure can generate high-quality crisp samples, avoiding the blurriness often attributed to VAE models in the past (Dosovitskiy & Brox, 2016; Larsen et al., 2015). And to the best of our knowledge, this is the first demonstration of a VAE pipeline that can produce stable FID scores, an influential recent metric for evaluating generated sample quality (Heusel et al., 2017), that are comparable to GAN models under neutral testing conditions. Moreover, this is accomplished without additional penalty functions, cost function modifications, or sensitive tuning parameters. Finally, Section 6 provides concluding thoughts and a discussion of broader VAE modeling paradigms.

2 High-Level Impact of VAE Gaussian Assumptions

Conventional wisdom suggests that VAE Gaussian assumptions will introduce a gap between $\mathcal{L}(\theta, \phi)$ and the ideal negative log-likelihood $\int_{\chi} -\log p_\theta(\mathbf{x})\mu_{gt}(d\mathbf{x})$, compromising efforts to learn the ground-truth measure. However, we will now argue that this pessimism is in some sense premature. In fact, we will demonstrate that, even with the stated Gaussian distributions, there exist parameters ϕ and θ that can simultaneously: (i) Globally optimize the VAE objective and, (ii) Recover the ground-truth probability measure in a certain sense described below. This is possible because, at least for some coordinated values of ϕ and θ , $q_\phi(z|\mathbf{x})$ and $p_\theta(\mathbf{x}|z)$ can indeed become arbitrarily close.

Before presenting the details, we first formalize a κ -simple VAE, which is merely a VAE model with explicit Gaussian assumptions and parameterizations:

Definition 1 A κ -simple VAE is defined as a VAE model with $\dim[\mathbf{z}] = \kappa$ latent dimensions, the Gaussian encoder $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)$, and the Gaussian decoder $p_\theta(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x)$. Moreover, the encoder moments are defined as $\boldsymbol{\mu}_z = f_{\mu_z}(\mathbf{x}; \phi)$ and $\boldsymbol{\Sigma}_z = \mathbf{S}_z \mathbf{S}_z^\top$ with $\mathbf{S}_z = f_{\mathbf{S}_z}(\mathbf{x}; \phi)$. Likewise, the decoder moments are $\boldsymbol{\mu}_x = f_{\mu_x}(\mathbf{z}; \theta)$ and $\boldsymbol{\Sigma}_x = \gamma \mathbf{I}$. Here $\gamma > 0$ is a tunable scalar, while f_{μ_z} , $f_{\mathbf{S}_z}$ and f_{μ_x} specify parameterized differentiable functional forms that can be arbitrarily complex, e.g., a deep neural network.

Equipped with these definitions, we will now demonstrate that a κ -simple VAE, with $\kappa \geq r$, can achieve the optimality criteria (i) and (ii) from above. In doing so, we first consider the simpler case where $r = d$, followed by the extended scenario with $r < d$. The distinction between these two cases turns out to be significant, with practical implications to be explored in Section 4.

2.1 Manifold Dimension Equal to Ambient Space Dimension ($r = d$)

We first analyze the specialized situation where $r = d$. Assuming $p_{gt}(\mathbf{x}) \triangleq \mu_{gt}(d\mathbf{x})/d\mathbf{x}$ exists everywhere in \mathbb{R}^d , then $p_{gt}(\mathbf{x})$ represents the ground-truth probability density with respect to the standard Lebesgue measure in Euclidean space. Given these considerations, the minimal possible value of (1) will necessarily occur if

$$\mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}|\mathbf{x})] = 0 \quad \text{and} \quad p_\theta(\mathbf{x}) = p_{gt}(\mathbf{x}) \text{ almost everywhere.} \quad (3)$$

This follows because by VAE design it must be that $\mathcal{L}(\theta, \phi) \geq -\int p_{gt}(\mathbf{x}) \log p_{gt}(\mathbf{x}) d\mathbf{x}$, and in the present context, this lower bound is achievable iff the conditions from (3) hold. Collectively, this implies that the approximate posterior produced by the encoder $q_\phi(\mathbf{z}|\mathbf{x})$ is in fact perfectly matched to the actual posterior $p_\theta(\mathbf{z}|\mathbf{x})$, while the corresponding marginalized data distribution $p_\theta(\mathbf{x})$ is perfectly matched the ground-truth density $p_{gt}(\mathbf{x})$ as desired. Perhaps surprisingly, a κ -simple VAE can actually achieve such a solution:

Theorem 1 Suppose that $r = d$ and there exists a density $p_{gt}(\mathbf{x})$ associated with the ground-truth measure μ_{gt} that is nonzero everywhere on \mathbb{R}^d .¹ Then for any $\kappa \geq r$, there is a sequence of κ -simple VAE model parameters $\{\theta_t^*, \phi_t^*\}$ such that

$$\lim_{t \rightarrow \infty} \mathbb{KL}[q_{\phi_t^*}(\mathbf{z}|\mathbf{x})||p_{\theta_t^*}(\mathbf{z}|\mathbf{x})] = 0 \quad \text{and} \quad \lim_{t \rightarrow \infty} p_{\theta_t^*}(\mathbf{x}) = p_{gt}(\mathbf{x}) \text{ almost everywhere.} \quad (4)$$

All the proofs can be found in the appendix. So at least when $r = d$, the VAE Gaussian assumptions need not actually prevent the optimal ground-truth probability measure from being recovered, as long as the latent dimension is sufficiently large (i.e., $\kappa \geq r$). And contrary to popular notions, a richer class of distributions is not required to achieve this. Of course Theorem 1 only applies to a restricted case that excludes $d > r$; however, later we will demonstrate that a key consequence of this result can nonetheless be leveraged to dramatically enhance VAE performance.

2.2 Manifold Dimension Less Than Ambient Space Dimension ($r < d$)

When $r < d$, additional subtleties are introduced that will be unpacked both here and in the sequel. To begin, if both $q_\phi(\mathbf{z}|\mathbf{x})$ and $p_\theta(\mathbf{x}|\mathbf{z})$ are arbitrary/unconstrained (i.e., not necessarily Gaussian), then $\inf_{\phi, \theta} \mathcal{L}(\theta, \phi) = -\infty$. To achieve this global optimum, we need only choose ϕ such that $q_\phi(\mathbf{z}|\mathbf{x}) = p_\theta(\mathbf{z}|\mathbf{x})$ (minimizing the KL term from (1)) while selecting θ such that all probability mass collapses to the correct manifold χ . In this scenario the density $p_\theta(\mathbf{x})$ will become unbounded on χ and zero elsewhere, such that $\int \chi - \log p_\theta(\mathbf{x}) \mu_{gt}(d\mathbf{x})$ will approach negative infinity.

But of course the stated Gaussian assumptions from the κ -simple VAE model could ostensibly prevent this from occurring by causing the KL term to blow up, counteracting the negative log-likelihood factor. We will now analyze this case to demonstrate that this need not happen. Before proceeding to

¹This nonzero assumption can be replaced with a much looser condition. Specifically, if there exists a diffeomorphism between the set $\{\mathbf{x} | p_{gt}(\mathbf{x}) \neq 0\}$ and \mathbb{R}^d , then it can be shown that Theorem 1 still holds even if $p_{gt}(\mathbf{x}) = 0$ for some $\mathbf{x} \in \mathbb{R}^d$.

this result, we first define a manifold density $\tilde{p}_{gt}(\mathbf{x})$ as the probability density (assuming it exists) of μ_{gt} with respect to the volume measure of the manifold χ . If $d = r$ then this volume measure reduces to the standard Lebesgue measure in \mathbb{R}^d and $\tilde{p}_{gt}(\mathbf{x}) = p_{gt}(\mathbf{x})$; however, when $d > r$ a density $p_{gt}(\mathbf{x})$ defined in \mathbb{R}^d will not technically exist, while $\tilde{p}_{gt}(\mathbf{x})$ is still perfectly well-defined. We then have the following:

Theorem 2 Assume $r < d$ and that there exists a manifold density $\tilde{p}_{gt}(\mathbf{x})$ associated with the ground-truth measure μ_{gt} that is nonzero everywhere on χ . Then for any $\kappa \geq r$, there is a sequence of κ -simple VAE model parameters $\{\theta_t^*, \phi_t^*\}$ such that

$$(i) \lim_{t \rightarrow \infty} \mathbb{KL}[q_{\phi_t^*}(\mathbf{z}|\mathbf{x})||p_{\theta_t^*}(\mathbf{z}|\mathbf{x})] = 0 \quad \text{and} \quad \lim_{t \rightarrow \infty} \int_{\chi} -\log p_{\theta_t^*}(\mathbf{x}) \mu_{gt}(d\mathbf{x}) = -\infty, \quad (5)$$

$$(ii) \lim_{t \rightarrow \infty} \int_{\mathbf{x} \in A} p_{\theta_t^*}(\mathbf{x}) d\mathbf{x} = \mu_{gt}(A \cap \chi) \quad (6)$$

for all measurable sets $A \subseteq \mathbb{R}^d$ with $\mu_{gt}(\partial A \cap \chi) = 0$, where ∂A is the boundary of A .

Technical details notwithstanding, Theorem 2 admits a very intuitive interpretation. First, (5) directly implies that the VAE Gaussian assumptions do not prevent minimization of $\mathcal{L}(\theta, \phi)$ from converging to minus infinity, which can be trivially viewed as a globally optimum solution. Furthermore, based on (6), this solution can be achieved with a limiting density estimate that will assign a probability mass to most all measurable subsets of \mathbb{R}^d that is indistinguishable from the ground-truth measure (which confines all mass to χ). Hence this solution is more-or-less an arbitrarily-good approximation to μ_{gt} for all practical purposes.²

Regardless, there is an absolutely crucial distinction between Theorem 2 and the simpler case quantified by Theorem 1. Although both describe conditions whereby the κ -simple VAE can achieve the minimal possible objective, in the $r = d$ case achieving the lower bound (whether the specific parameterization for doing so is unique or not) necessitates that the ground-truth probability measure has been recovered almost everywhere. But the $r < d$ situation is quite different because we have not ruled out the possibility that a different set of parameters $\{\theta, \phi\}$ could push $\mathcal{L}(\theta, \phi)$ to $-\infty$ and yet not achieve (6). In other words, the VAE could reach the lower bound but fail to closely approximate μ_{gt} . And we stress that this uniqueness issue is not a consequence of the VAE Gaussian assumptions per se; even if $q_{\phi}(\mathbf{z}|\mathbf{x})$ were unconstrained the same lack of uniqueness can persist.

Rather, the intrinsic difficulty is that, because the VAE model does not have access to the ground-truth low-dimensional manifold, it must implicitly rely on a density $p_{\theta}(\mathbf{x})$ defined across *all* of \mathbb{R}^d as mentioned previously. Moreover, if this density converges towards infinity on the manifold during training without increasing the KL term at the same rate, the VAE cost can be unbounded from below, even in cases where (6) is not satisfied, meaning incorrect assignment of probability mass.

To conclude, the key take-home message from this section is that, at least in principle, VAE Gaussian assumptions need not actually be the root cause of any failure to recover ground-truth distributions. Instead we expose a structural deficiency that lies elsewhere, namely, the non-uniqueness of solutions that can optimize the VAE objective without necessarily learning a close approximation to μ_{gt} . But to probe this issue further and motivate possible workarounds, it is critical to further disambiguate these optimal solutions and their relationship with ground-truth manifolds. This will be the task of Section 3, where we will explicitly differentiate the problem of locating the correct ground-truth manifold, from the task of learning the correct probability measure *within* the manifold.

Note that the only comparable prior work we are aware of related to the results in this section comes from Doersch (2016), where the implications of adopting Gaussian encoder/decoder pairs in the specialized case of $r = d = 1$ are briefly considered. Moreover, the analysis there requires additional much stronger assumptions than ours, namely, that $p_{gt}(\mathbf{x})$ should be nonzero and infinitely differentiable everywhere in the requisite 1D ambient space. These requirements of course exclude essentially all practical usage regimes where $d = r > 1$ or $d > r$, or when ground-truth densities are not sufficiently smooth.

²Note that (6) is only framed in this technical way to accommodate the difficulty of comparing a measure μ_{gt} restricted to χ with the VAE density $p_{\theta}(\mathbf{x})$ defined everywhere in \mathbb{R}^d . See the appendix for details.

3 Optimal Solutions and the Ground Truth Manifold

We will now more closely examine the properties of optimal κ -simple VAE solutions, and in particular, the degree to which we might expect them to at least reflect the true χ , even if perhaps not the correct probability measure μ_{gt} defined within χ . To do so, we must first consider some *necessary* conditions for VAE optima:

Theorem 3 *Let $\{\theta_\gamma^*, \phi_\gamma^*\}$ denote an optimal κ -simple VAE solution (with $\kappa \geq r$) where the decoder variance γ is fixed (i.e., it is the sole unoptimized parameter). Moreover, we assume that μ_{gt} is not a Gaussian distribution when $d = r$.³ Then for any $\gamma > 0$, there exists a $\gamma' < \gamma$ such that $\mathcal{L}(\theta_{\gamma'}^*, \phi_{\gamma'}^*) < \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*)$.*

This result implies that we can always reduce the VAE cost by choosing a smaller value of γ , and hence, if γ is not constrained, it must be that $\gamma \rightarrow 0$ if we wish to minimize (2). Despite this necessary optimality condition, in existing practical VAE applications, it is standard to fix $\gamma \approx 1$ during training. This is equivalent to simply adopting a non-adaptive squared-error loss for the decoder and, at least in part, likely contributes to unrealistic/blurred VAE-generated samples. Regardless, there are more significant consequences of this intrinsic favoritism for $\gamma \rightarrow 0$, in particular as related to reconstructing data drawn from the ground-truth manifold χ :

Theorem 4 *Applying the same conditions and definitions as in Theorem 3, then for all \mathbf{x} drawn from μ_{gt} , we also have that*

$$\lim_{\gamma \rightarrow 0} f_{\mu_x} [f_{\mu_z}(\mathbf{x}; \phi_\gamma^*) + f_{S_z}(\mathbf{x}; \phi_\gamma^*)\epsilon; \theta_\gamma^*] = \lim_{\gamma \rightarrow 0} f_{\mu_x} [f_{\mu_z}(\mathbf{x}; \phi_\gamma^*); \theta_\gamma^*] = \mathbf{x}, \quad \forall \epsilon \in \mathbb{R}^\kappa. \quad (7)$$

By design any random draw $\mathbf{z} \sim q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})$ can be expressed as $f_{\mu_z}(\mathbf{x}; \phi_\gamma^*) + f_{S_z}(\mathbf{x}; \phi_\gamma^*)\epsilon$ for some $\epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \mathbf{I})$. From this vantage point then, (7) effectively indicates that any $\mathbf{x} \in \chi$ will be perfectly reconstructed by the VAE encoder/decoder pair at globally optimal solutions, achieving this necessary condition despite any possible stochastic corrupting factor $f_{S_z}(\mathbf{x}; \phi_\gamma^*)\epsilon$.

But still further insights can be obtained when we more closely inspect the VAE objective function behavior at arbitrarily small but explicitly nonzero values of γ . In particular, when $\kappa = r$ (meaning \mathbf{z} has no superfluous capacity), Theorem 4 and attendant analyses in the appendix ultimately imply that the squared eigenvalues of $f_{S_z}(\mathbf{x}; \phi_\gamma^*)$ will become arbitrarily small at a rate proportional to γ , meaning $\frac{1}{\sqrt{\gamma}} f_{S_z}(\mathbf{x}; \phi_\gamma^*) \approx O(1)$ under mild conditions. It then follows that the VAE data term integrand from (2), in the neighborhood around optimal solutions, behaves as $-2\mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} [\log p_{\theta_\gamma^*}(\mathbf{x}|\mathbf{z})] =$

$$2\mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} \left[\frac{1}{\gamma} \|\mathbf{x} - f_{\mu_x}[\mathbf{z}; \theta_\gamma^*]\|_2^2 \right] + d \log 2\pi\gamma \approx \mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} [O(1)] + d \log 2\pi\gamma = d \log \gamma + O(1). \quad (8)$$

This expression can be derived by excluding the higher-order terms of a Taylor series approximation of $f_{\mu_x} [f_{\mu_z}(\mathbf{x}; \phi_\gamma^*) + f_{S_z}(\mathbf{x}; \phi_\gamma^*)\epsilon; \theta_\gamma^*]$ around the point $f_{\mu_z}(\mathbf{x}; \phi_\gamma^*)$, which will be relatively tight under the stated conditions. But because $2\mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} \left[\frac{1}{\gamma} \|\mathbf{x} - f_{\mu_x}[\mathbf{z}; \theta_\gamma^*]\|_2^2 \right] \geq 0$, a theoretical lower bound on (8) is given by $d \log 2\pi\gamma \equiv d \log \gamma + O(1)$. So in this sense (8) cannot be significantly lowered further.

This observation is significant when we consider the inclusion of additional latent dimensions by allowing $\kappa > r$. Clearly based on the analysis above, adding dimensions to \mathbf{z} *cannot* improve the value of the VAE data term in any meaningful way. However, it can have a detrimental impact on the the KL regularization factor in the $\gamma \rightarrow 0$ regime, where

$$2\mathbb{KL}[q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \equiv \text{trace}[\Sigma_z] + \|\mu_z\|_2^2 - \log |\Sigma_z| \approx -\hat{r} \log \gamma + O(1). \quad (9)$$

Here \hat{r} denotes the number of eigenvalues $\{\lambda_j(\gamma)\}_{j=1}^{\kappa}$ of $f_{S_z}(\mathbf{x}; \phi_\gamma^*)$ (or equivalently Σ_z) that satisfy $\lambda_j(\gamma) \rightarrow 0$ if $\gamma \rightarrow 0$. \hat{r} can be viewed as an estimate of how many low-noise latent dimensions

³This requirement is only included to avoid a practically irrelevant form of non-uniqueness that exists with full, non-degenerate Gaussian distributions.

the VAE model is preserving to reconstruct \mathbf{x} . Based on (9), there is obvious pressure to make \hat{r} as small as possible, at least without disrupting the data fit. The smallest possible value is $\hat{r} = r$, since it is not difficult to show that any value below this will contribute consequential reconstruction errors, causing $2\mathbb{E}_{q_{\phi}^*}(\mathbf{z}|\mathbf{x}) \left[\frac{1}{\gamma} \|\mathbf{x} - f_{\mu_x}[\mathbf{z}; \theta_{\gamma}^*]\|_2^2 \right]$ to grow at a rate of $\Omega\left(\frac{1}{\gamma}\right)$, pushing the entire cost function towards infinity.⁴

Therefore, *in the neighborhood of optimal solutions the VAE will naturally seek to produce perfect reconstructions using the fewest number of clean, low-noise latent dimensions*, meaning dimensions whereby $q_{\phi}(\mathbf{z}|\mathbf{x})$ has negligible variance. For superfluous dimensions that are unnecessary for representing \mathbf{x} , the associated encoder variance in these directions can be pushed to one. This will optimize $\mathbb{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})]$ along these directions, and the decoder can selectively block the residual randomness to avoid influencing the reconstructions per Theorem 4. So in this sense the VAE is capable of learning a minimal representation of the ground-truth manifold χ when $r < \kappa$.

But we must emphasize that the VAE can learn χ independently of the actual distribution μ_{gt} within χ . Addressing the latter is a completely separate issue from achieving the perfect reconstruction error defined by Theorem 4. This fact can be understood within the context of a traditional PCA-like model, which is perfectly capable of learning a low-dimensional subspace containing some training data without actually learning the distribution of the data within this subspace. The central issue is that there exists an intrinsic bias associated with the VAE objective such that *fitting the distribution within the manifold will be completely neglected whenever there exists the chance for even an infinitesimally better approximation of the manifold itself*.

Stated differently, if VAE model parameters have learned a near optimal, parsimonious latent mapping onto χ using $\gamma \approx 0$, then the VAE cost will scale as $(d - r) \log \gamma$ regardless of μ_{gt} . Hence there remains a *huge* incentive to reduce the reconstruction error still further, allowing γ to push even closer to zero and the cost closer to $-\infty$. And if we constrain γ to be sufficiently large so as to prevent this from happening, then we risk degrading/blurring the reconstructions and widening the gap between $q_{\phi}(\mathbf{z}|\mathbf{x})$ and $p_{\theta}(\mathbf{z}|\mathbf{x})$, which can also compromise estimation of μ_{gt} . Fortunately though, as will be discussed next there is a convenient way around this dilemma by exploiting the fact that this dominating $(d - r) \log \gamma$ factor goes away when $d = r$.

4 From Theory to Practical VAE Enhancements

Sections 2 and 3 have exposed a collection of VAE properties with useful diagnostic value in and of themselves. But the practical utility of these results, beyond the underappreciated benefit of learning γ , warrant further exploration. In this regard, suppose we wish to develop a generative model of high-dimensional data $\mathbf{x} \in \chi$ where unknown low-dimensional structure is significant (i.e., the $r < d$ case with r unknown). The results from Section 3 indicate that the VAE can partially handle this situation by learning a parsimonious representation of low-dimensional manifolds, but not necessarily the correct probability measure μ_{gt} within such a manifold. In quantitative terms, this means that a decoder $p_{\theta}(\mathbf{x}|\mathbf{z})$ will map all samples from an encoder $q_{\phi}(\mathbf{z}|\mathbf{x})$ to the correct manifold such that the reconstruction error is negligible for any $\mathbf{x} \in \chi$. But if the measure μ_{gt} on χ has not been accurately estimated, then

$$q_{\phi}(\mathbf{z}) \triangleq \int_{\chi} q_{\phi}(\mathbf{z}|\mathbf{x}) \mu_{gt}(d\mathbf{x}) \not\approx \int_{\mathbb{R}^d} p_{\theta}(\mathbf{z}|\mathbf{x}) p_{\theta}(\mathbf{x}) d\mathbf{x} = \int_{\mathbb{R}^d} p_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{x} = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}), \quad (10)$$

where $q_{\phi}(\mathbf{z})$ is sometimes referred to as the aggregated posterior (Makhzani et al., 2016). In other words, the distribution of the latent samples drawn from the encoder distribution, when averaged across the training data, will have lingering latent structure that is errantly incongruous with the original isotropic Gaussian prior. This then disrupts the pivotal ancestral sampling capability of the VAE, implying that samples drawn from $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$ and then passed through the decoder $p_{\theta}(\mathbf{x}|\mathbf{z})$ will *not* closely approximate μ_{gt} . Fortunately, our analysis suggests the following two-stage remedy:

1. Given n observed samples $\{\mathbf{x}^{(i)}\}_{i=1}^n$, train a κ -simple VAE, with $\kappa \geq r$, to estimate the unknown r -dimensional ground-truth manifold χ embedded in \mathbb{R}^d using a minimal number of active latent dimensions. Generate latent samples $\{\mathbf{z}^{(i)}\}_{i=1}^n$ via $\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$. By design, these samples will be distributed as $q_{\phi}(\mathbf{z})$, but likely not $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$.

⁴Note that $\inf_{\gamma>0} \frac{C}{\gamma} + \log \gamma = \infty$ for any $C > 0$.

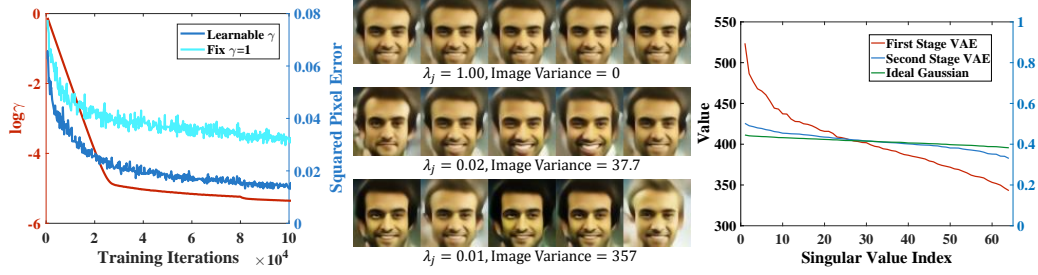


Figure 1: *Demonstrating VAE properties.* (Left) Validation of Theorem 3 and the influence on image reconstructions. (Center) Validation of Theorem 4. (Right) Motivation for two separate VAE stages by comparing the aggregated posteriors $q_\phi(\mathbf{z})$ (1^{st} stage) vs. $q_{\phi'}(\mathbf{u})$ (enhanced 2^{nd} stage).

2. Train a second κ -simple VAE, with independent parameters $\{\theta', \phi'\}$ and latent representation \mathbf{u} , to learn the unknown distribution $q_\phi(\mathbf{z})$, i.e., treat $q_\phi(\mathbf{z})$ as a new ground-truth distribution and use samples $\{\mathbf{z}^{(i)}\}_{i=1}^n$ to learn it.
3. Samples approximating the *original* ground-truth μ_{gt} can then be formed via the extended ancestral process $\mathbf{u} \sim \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{I})$, $\mathbf{z} \sim p_{\phi'}(\mathbf{z}|\mathbf{u})$, and finally $\mathbf{x} \sim p_\theta(\mathbf{x}|\mathbf{z})$.

The efficacy of the second-stage VAE from above is based on the following. If the first stage was successful, then even though they will not generally resemble $\mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$, samples from $q_\phi(\mathbf{z})$ will nonetheless have nonzero measure across the full ambient space \mathbb{R}^κ . If $\kappa = r$, this occurs because the entire latent space is needed to represent an r -dimensional manifold, and if $\kappa > r$, then the extra latent dimensions will be naturally filled in via randomness introduced along dimensions associated with nonzero eigenvalues of the decoder covariance Σ_z per the analysis in Section 3.

Consequently, as long as we set $\kappa \geq r$, the *operational regime of the second-stage VAE is effectively equivalent to the situation described in Section 2.1 where the manifold dimension is equal to the ambient dimension*.⁵ And as we have already shown there via Theorem 1, the VAE can readily handle this situation, since in the narrow context of the second-stage VAE, $d = r = \kappa$, the troublesome $(d - r) \log \gamma$ factor becomes zero, and any globally minimizing solution is uniquely matched to the new ground-truth distribution $q_\phi(\mathbf{z})$. Consequently, the revised aggregated posterior $q_{\phi'}(\mathbf{u})$ produced by the second-stage VAE should now closely resemble $\mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{I})$. And finally, because we generally assume that $d \gg \kappa \geq r$, we have found that the second-stage VAE can be quite small.

5 Empirical Evaluation of VAE Two-Stage Enhancement

We initially describe experiments explicitly designed to corroborate some of our previous analytical results using VAE models trained on CelebA (Liu et al., 2015) data; please see the appendix for training details and more related experiments. First, the leftmost plot of Figure 1 presents support for Theorem 3, where indeed the decoder variance γ does tend towards zero during training. This then allows for tighter image reconstructions with lower average squared error, i.e., a better manifold fit as expected. The center plot bolsters Theorem 4 and the analysis that follows by showcasing the dissimilar impact of noise factors applied to different directions in the latent space before passage through the decoder mean network f_{μ_x} . In a direction where an eigenvalue λ_j of Σ_z is large (i.e., a superfluous dimension), a random perturbation is completely muted by the decoder as predicted. In contrast, in directions where such eigenvalues are small (i.e., needed for representing the manifold), varying the input causes large changes in the image space reflecting reasonable movement along the correct manifold. Finally, the rightmost plot of Figure 1 displays the singular value spectrum of latent sample matrices drawn from the first- and second-stage VAE models. As expected, the latter is much closer to the spectrum from an analogous i.i.d. $\mathcal{N}(\mathbf{0}, \mathbf{I})$ matrix. This indicates a superior latent representation, providing high-level support for our two-stage VAE proposal.

Next we present quantitative evaluation of novel generated samples using the large-scale testing protocol of GAN models from (Lucic et al., 2018). In this regard, GANs are well-known to dramatically

⁵Note that if a regular autoencoder were used to replace the first-stage VAE, then this would no longer be the case, so indeed a VAE is required for both stages.

		MNIST	Fashion	CIFAR-10	CelebA
optimized, data-dependent settings	MM GAN	9.8 ± 0.9	29.6 ± 1.6	72.7 ± 3.6	65.6 ± 4.2
	NS GAN	6.8 ± 0.5	26.5 ± 1.6	58.5 ± 1.9	55.0 ± 3.3
	LSGAN	7.8 ± 0.6	30.7 ± 2.2	87.1 ± 47.5	53.9 ± 2.8
	WGAN	6.7 ± 0.4	21.5 ± 1.6	55.2 ± 2.3	41.3 ± 2.0
	WGAN GP	20.3 ± 5.0	24.5 ± 2.1	55.8 ± 0.9	30.3 ± 1.0
	DRAGAN	7.6 ± 0.4	27.7 ± 1.2	69.8 ± 2.0	42.3 ± 3.0
	BEGAN	13.1 ± 1.0	22.9 ± 0.9	71.4 ± 1.6	38.9 ± 0.9
default settings	Best GAN	~ 10	~ 32	~ 70	~ 49
	VAE (fixed γ)	52.0 ± 0.6	84.6 ± 0.9	160.5 ± 1.1	55.9 ± 0.6
	VAE (learned γ)	54.5 ± 1.0	60.0 ± 1.1	76.7 ± 0.8	60.5 ± 0.6
	2-Stage VAE (ours)	12.6 ± 1.5	29.3 ± 1.0	72.9 ± 0.9	44.4 ± 0.7

Table 1: *FID score comparisons*. For all GAN-based models listed in the top section of the table, reported values represent the optimal FID obtained across a large-scale hyperparameter search conducted separately for each dataset (Lucic et al., 2018). Outlier cases (e.g., severe mode collapse) were omitted, which would have otherwise increased these GAN FID scores. In the lower section of the table, the label *Best GAN* indicates the lowest FID produced across all GAN approaches when trained using settings suggested by original authors; these approximate values were extracted from (Lucic et al., 2018, Figure 4). For the VAE results, only a single default setting was adopted across all datasets and models (no tuning whatsoever), and no cases of mode collapse were removed. Note that specialized architectures and/or random seed optimization can potentially improve the FID score for all models reported here.

outperform existing VAE approaches in terms of the Fréchet Inception Distance (FID) score (Heusel et al., 2017) and related quantitative metrics. For fair comparison, (Lucic et al., 2018) adopted a common neutral architecture for all models, with generator and discriminator networks based on InfoGAN (Chen et al., 2016); the point here is standardized comparisons, not tuning arbitrarily-large networks to achieve the lowest possible absolute FID values. We applied the same architecture to our first-stage VAE decoder and encoder networks respectively for direct comparison. For the low-dimensional second-stage VAE we used small, 3-layer networks contributing negligible additional parameters beyond the first stage (see the appendix for further design details).⁶

We evaluated our proposed VAE pipeline, denoted *2-Stage VAE*, against baseline VAE models differing only in the decoder output layer: a Gaussian layer with fixed γ , and a Gaussian layer with a learned γ (the latter is also used by the two-stage VAE). We also present results from (Lucic et al., 2018) involving numerous competing GAN models, including MM GAN (Goodfellow et al., 2014), WGAN (Arjovsky et al., 2017), WGAN-GP (Gulrajani et al., 2017), NS GAN (Fedus et al., 2017), DRAGAN (Kodali et al., 2017), LS GAN (Mao et al., 2017) and BEGAN (Berthelot et al., 2017). Testing is conducted across four significantly different datasets: MNIST (LeCun et al., 1998), Fashion MNIST (Xiao et al., 2017), CIFAR-10 (Krizhevsky & Hinton, 2009) and CelebA (Liu et al., 2015).

For each dataset we executed 10 independent trials and report the mean and standard deviation of the FID scores in Table 1.⁷ No effort was made to tune VAE training hyperparameters (e.g., learning rates, etc.); rather a single generic setting was first selected and then applied to all VAE models. As an analogous baseline, we also report the value of the best GAN model for each dataset when trained using suggested settings from the authors; no single model was optimal across all datasets, so these values represent performance from different GANs. Even so, our single 2-Stage VAE is still better on two of four datasets, and in aggregate, better than any individual GAN model. For example, when

⁶It should also be emphasized that concatenating the two stages and jointly training does not improve the performance. If trained jointly the few extra second-stage parameters are simply hijacked by the dominant objective from the first stage and forced to work on an incrementally better fit of the manifold. As expected then, on empirical tests (not shown) we have found that this does not improve upon standard VAE baselines.

⁷All reported FID scores for VAE and GAN models were computed using TensorFlow (<https://github.com/bioinf-jku/TTUR>). We have found that alternative PyTorch implementations (<https://github.com/mseitzer/pytorch-fid>) can produce different values in some circumstances. This seems to be due, at least in part, to subtle differences in the underlying Inception models being used for computing the scores. Either way, a consistent implementation is essential for calibrating results across different scenarios.

averaged across datasets, the mean FID score for any individual GAN trained with suggested settings was always approximately 45 or higher (see (Lucic et al., 2018, Figure 4)), while our analogous 2-Stage VAE maintained a mean below 40. The other VAE baselines were not competitive.

Table 1 also displays FID scores from GAN models evaluated using hyperparameters obtained from a large-scale search executed independently across each dataset to achieve the best results; 100 settings per model per dataset, plus an optimal, data-dependent stopping criteria as described in (Lucic et al., 2018). Within this broader paradigm, cases of severe mode collapse were omitted when computing final GAN FID averages. Despite these considerable advantages, the FID performance of the default 2-Stage VAE is well within the range of the heavily-optimized GAN models for each dataset unlike the other VAE baselines. Overall then, these results represent the first demonstration of a VAE pipeline capable of competing with GANs in the arena of generated sample quality. Representative samples generated using our two-stage VAE approach are in the appendix.

6 Discussion

It is often assumed that there exists an unavoidable trade-off between the stable training, valuable attendant encoder network, and resistance to mode collapse of VAEs, versus the impressive visual quality of images produced by GANs. While we certainly are not claiming that our two-stage VAE model is superior to the latest and greatest GAN-based architecture in terms of the realism of generated samples, we do strongly believe that this work at least narrows that gap substantially such that VAEs are worth considering in a broader range of applications.

It is also important to recognize that a variety of alternative VAE enhancements have recently been proposed as well; however, nearly all of these have focused on improving the log-likelihood scores assigned by the model to test data. In particular, multiple elegant VAE modifications involve replacing the Gaussian encoder network with a richer class of distributions instantiated through normalizing flows or related (Burda et al., 2015; Kingma et al., 2016; Rezende & Mohamed, 2015; van den Berg et al., 2018). While impressive log-likelihood gains have been demonstrated, this achievement is largely orthogonal to the goal of improving quantitative measures of visual quality (Theis et al., 2016), which has been our focus herein. Additionally, improving the VAE encoder does not address the uniqueness issue raised in Section 2, and therefore, a second stage could potentially benefit these models too in the right circumstances.

Broadly speaking, if the overriding objective is generating realistic samples using an encoder-decoder-based architecture (VAE or otherwise), two important, well-known criteria must be satisfied:

- (i) Small reconstruction error when passing through the encoder-decoder networks, and
- (ii) An aggregate posterior $q_\phi(z)$ that is close to some known distribution like $p(z) = \mathcal{N}(z|0, I)$ that is easy to sample from.

Without the latter criteria, we have no tractable way of generating random inputs that, when passed through the learned decoder, produce realistic output samples resembling the training data distribution.

Criteria (i) and (ii) can be addressed multiple different ways. For example, (Tomczak & Welling, 2018) replace $\mathcal{N}(z|0, I)$ with a richer parameterized class of prior distributions $p(z)$ such that there exist more flexible pathways for pushing $p(z)$ and $q_\phi(z)$ closer together. Consequently, even if $q_\phi(z)$ is not Gaussian, we can nonetheless sample from a known non-Gaussian alternative. This is certainly an interesting idea, but it has not as of yet been applied to improving FID scores and only log-likelihood values on relatively small black-and-white images are reported.

In fact, the only competing encoder-decoder-based architecture we are aware of that explicitly attempts to improve FID scores comes from (Tolstikhin et al., 2018), which presents what can be viewed as a generalization of the adversarial autoencoder (Makhzani et al., 2016). The basic idea is to minimize an objective function composed of a reconstruction penalty for handling criteria (i), and a Wassenstein distance measure between $p(z)$ and $q_\phi(z)$ for addressing criteria (ii). Two variants of this approach are referred to as WAE-MMD and WAE-GAN because different MMD and GAN regularization factors are involved. Both are evaluated using hyperparameters and encoder-decoder networks specifically adapted for use with the CelebA dataset. Therefore, although FID scores are reported, they are not really comparable with the Table 1 values because of the different architecture and testing conditions. That being said, the WAE-GAN version involves GAN-like adversarial

training, and under the reported testing conditions more-or-less defaults to an adversarial autoencoder with an FID score of 42. This is similar to the other optimized GAN-based models from Table 1.

In contrast, WAE-MMD does not require potentially-difficult adversarial training, just like a VAE as desired, but the corresponding FID score increases to 55. Again, although not directly comparable since a specific network structure has been selected for CelebA, this is nonetheless still significantly higher than our 2-Stage VAE trained using a neutral architecture borrowed from (Lucic et al., 2018) with default settings. Additionally, both WAE-MMD and WAE-GAN models are dependent on having a reasonable estimate for $\kappa \approx r$ (at least for the deterministic encoder-decoder models that were empirically tested), otherwise matching $p(\mathbf{z})$ and $q_\phi(\mathbf{z})$ is not possible (Tolstikhin et al., 2018). For the 2-Stage VAE, we only need choose $\kappa \geq r$ in principle.

As with the approaches mentioned above, the two VAE stages we have proposed can also be motivated in one-to-one correspondence with criteria (i) and (ii). In brief, the first VAE stage addresses criteria (i) by pushing both the encoder variance, and the decoder variances selectively, towards zero such that accurate reconstruction is possible using a minimal number of active latent dimensions. However, our detailed analysis suggests that, although the resulting aggregate posterior $q_\phi(\mathbf{z})$ will occupy nonzero measure in κ -dimensional space (selectively filling out superfluous dimensions with random noise), it need not be close to $\mathcal{N}(\mathbf{z}|0, \mathbf{I})$. This then implies that if we take samples from $\mathcal{N}(\mathbf{z}|0, \mathbf{I})$ and pass them through the learned decoder, the result may not closely resemble real data.

Of course if we could somehow directly sample from $q_\phi(\mathbf{z})$, then we would not need to use $\mathcal{N}(\mathbf{z}|0, \mathbf{I})$. And fortunately, because the first-stage VAE ensures that $q_\phi(\mathbf{z})$ will satisfy the conditions of Theorem 1, we know that a second VAE can in fact be learned to accurately sample from this distribution, which in turn addresses criteria (ii). Specifically, per the arguments from Section 4, sampling $\mathbf{u} \sim \mathcal{N}(\mathbf{u}|\mathbf{0}, \mathbf{I})$ and then $\mathbf{z} \sim p_{\theta'}(\mathbf{z}|\mathbf{u})$ is akin to sampling $\mathbf{z} \sim q_\phi(\mathbf{z})$. Such samples can then be passed through the first-stage VAE decoder to obtain samples of \mathbf{x} . Hence our framework provides a principled alternative to existing encoder-decoder structures designed to handle criteria (i) and (ii), leading to state-of-the-art results for this class of model in terms of FID scores. In any event, we intend to further explore these issues in an extended journal version, including broader empirical testing with alternative VAE baselines.

Appendix

A Comparison of Novel Samples Generated from our Model

Generation results for CelebA, MNIST, Fashion-MNIST and CIFAR-10 datasets are shown in Figures 2–5 respectively. When γ is fixed to be one, the generated samples are very blurry. If a learnable γ is used, the samples becomes sharper; however, there are many lingering artifacts as expected. In contrast, the proposed 2-Stage VAE can remove these artifacts and generate more realistic samples. For comparison purposes, we also show the results from WAE-MMD, WAE-GAN (Tolstikhin et al., 2018) and WGAN-GP (Gulrajani et al., 2017) for the CelebA dataset.



Figure 2: Randomly generated samples on the CelebA dataset (i.e., no cherry-picking).

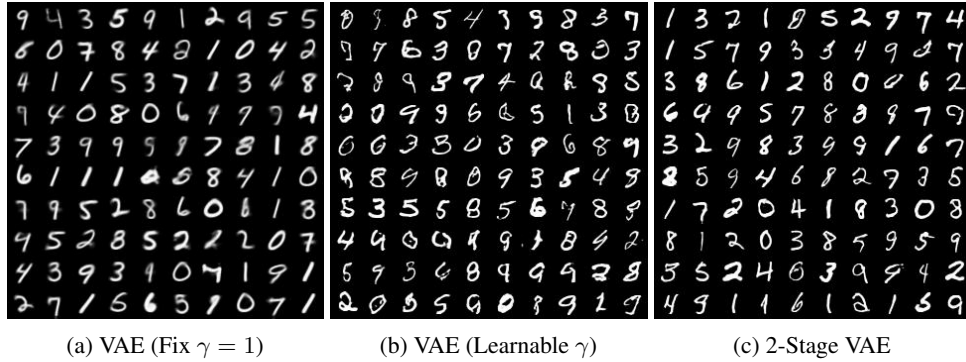


Figure 3: Randomly generated samples on the MNIST dataset (i.e., no cherry-picking).

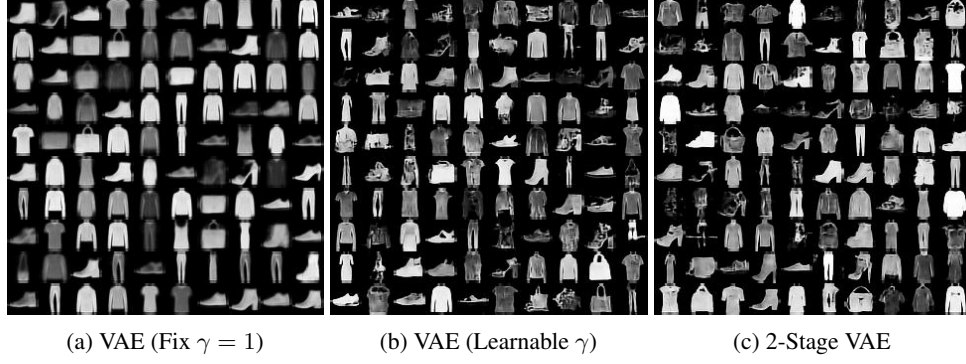


Figure 4: Randomly Generated Samples on Fashion-MNIST Dataset (i.e., no cherry-picking).



Figure 5: Randomly Generated Samples on CIFAR-10 Dataset (i.e., no cherry-picking).

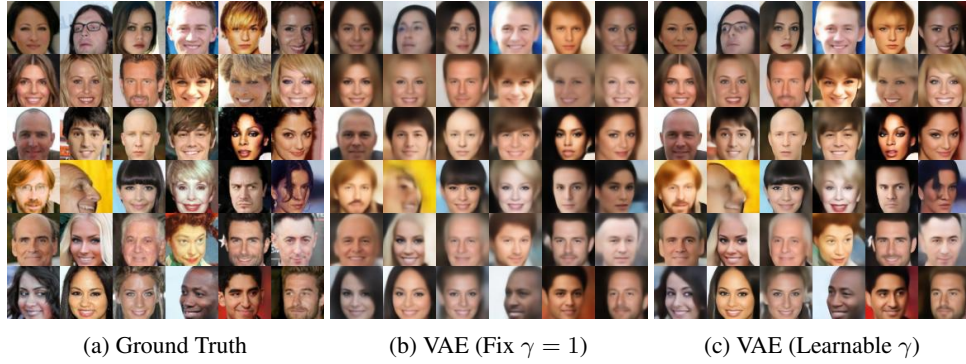


Figure 6: Reconstructions on CelebA Dataset.

B Example Reconstructions of Training Data

Reconstruction results for MNIST, Fashion-MNIST, CIFAR-10 and CelebA datasets are shown in Figures 6–9 respectively. On relatively simple datasets like MNIST and Fashion-MNIST, the VAE with learnable γ achieves almost exact reconstruction because of a better estimate of the underlying manifold consistent with theory. However, the VAE with fixed $\gamma = 1$ produces blurry reconstructions as expected. Note that the reconstruction of a 2-Stage VAE is the same as that of a VAE with learnable γ because the second-stage VAE has nothing to do with facilitating the reconstruction task.

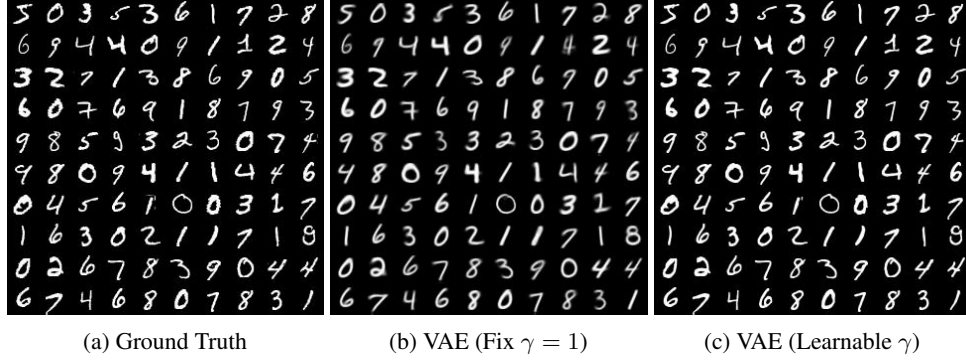


Figure 7: Reconstructions on MNIST Dataset.

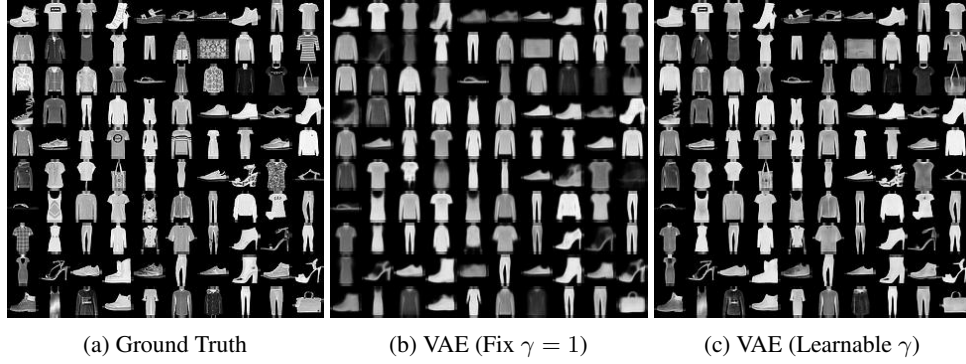


Figure 8: Reconstructions on Fashion-MNIST Dataset.



Figure 9: Reconstructions on CIFAR-10 Dataset.

C Additional Experimental Results Validating Theoretical Predictions

We first present more examples similar to Figure 1(*center*) from the main paper. Random noise is added to μ_z along different directions and the result is passed through the decoder network. Each row corresponds to a certain direction in the latent space and 15 samples are shown for each direction. These dimensions/rows are ordered by the eigenvalues λ_j of Σ_z . The larger λ_j is, the less impact a random perturbation along this direction will have as quantified by the reported image variance values. In the first two or three rows, the noise generates some images from different classes/objects/identities, indicating a significant visual difference. For a slightly larger λ_j , the corresponding dimensions encode relatively less significant attributes as predicted. For example, the fifth row of both MNIST and Fashion-MNIST contains images from the same class but with a slightly different style. The images in the fourth row of the CelebA dataset have very subtle differences. When $\lambda_j = 1$, the

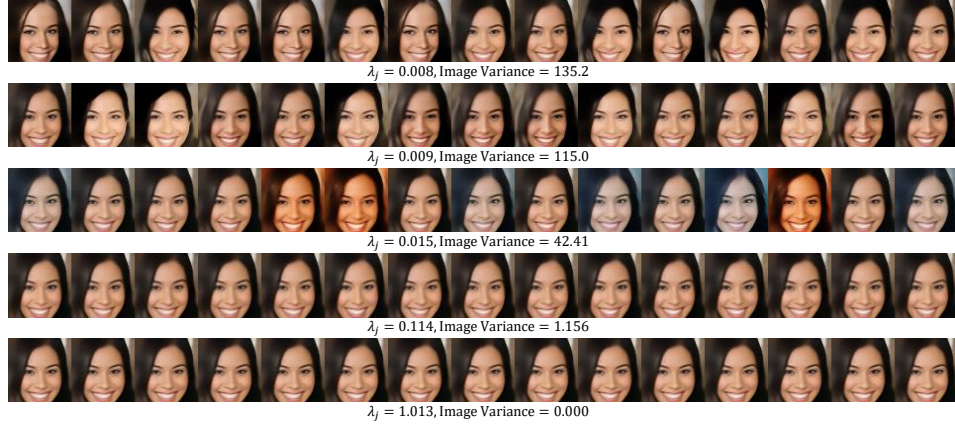
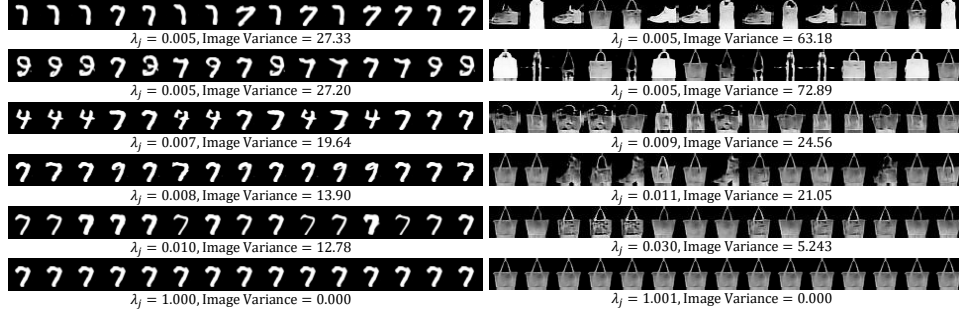


Figure 10: More examples similar to Figure 1(*center*).

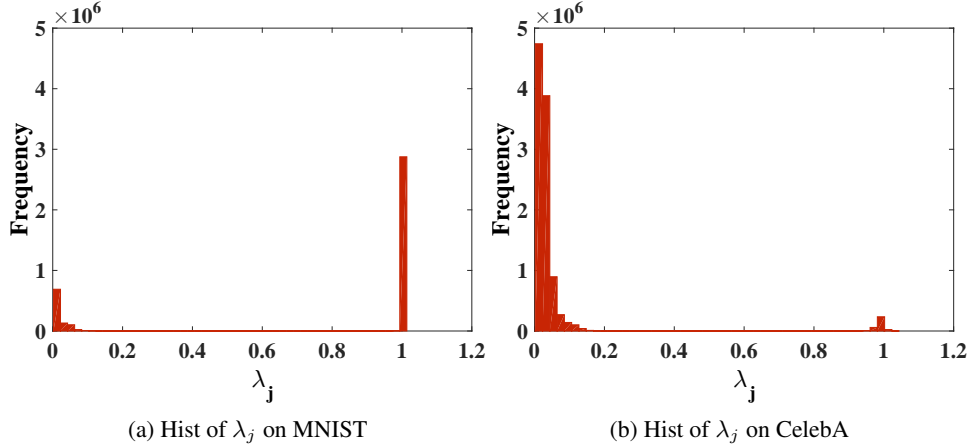


Figure 11: Histogram of λ_j values. There are more values around 0 for CelebA because it is more complicated than MNIST and therefore requires more active dimensions to model the underlying manifold.

corresponding dimensions become completely inactive and all the output images are exactly the same, as shown in the last rows for all the three datasets.

Additionally, as discussed in the main text and below in Section I, there are likely to be r eigenvalues of Σ_z converging to zero and $\kappa - r$ eigenvalues converging to one. We plot the histogram of λ_j values for both MNIST and CelebA datasets in Figure 11. For both datasets, λ_j approximately converges to either 0 or one. However, since CelebA is a more complicated dataset than MNIST, the ground-truth manifold dimension of CelebA is likely to be much larger than that of MNIST.

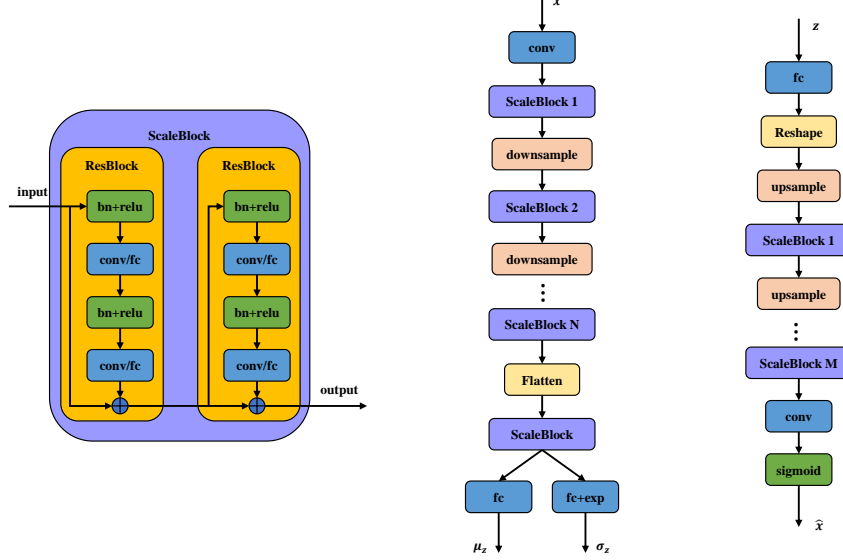


Figure 12: Network structure of the first-stage VAE used in producing Figure 1, and for generating samples and reconstructions. *(Left)* The basic building block of the network called a Scale Block, which consists of two Residual Blocks. *(Center)* The encoder network. For an input image x , we use a convolutional layer to transform it into 32 channels. We then pass it to a Scale Block. After each Scale Block, we downsample using a convolutional layer with stride 2 and double the channels. After N Scale Blocks, the feature map is flattened to a vector. In our experiments, we used $N = 4$ for CelebA dataset and 3 for other datasets. The vector is then passed through another Scale Block, the convolutional layers of which are replaced with fully connected layers of 512 dimensions. The output of this Scale Block is used to produce the κ -dimensional latent code, with $\kappa = 64$. *(Right)* The decoder network. A latent code z is first passed through a fully connected layer. The dimension is 4096 for CelebA dataset and 2048 for other datasets. Then it is reshaped to 2×2 resolution. We upsample the feature map using a deconvolution layer and half the number of channels at the same time. It then goes through some Scale Blocks and upsampling layers until the feature map size becomes the desired value. Then we use a convolutional layer to transform the feature map, which should have 32 channels, to 3 channels for RGB datasets and 1 channel for gray scale datasets.

So more eigenvalues are expected to be near zero for the CelebA dataset. This is indeed the case, demonstrating that VAE has the ability to detect the manifold dimension and select the proper number of latent dimensions in practical environments.

D Network Structure and Experimental Settings

We first describe the network and training details used in producing Figure 1 from the main file, and for generating samples and reconstructions in the supplementary. The first-stage VAE network is shown in Figure 12. Basically we use two Residual Blocks for each resolution scale, and we double the number of channels when downsampling and halve it when upsampling. The specific settings such as the number of channels and the number of scales are specified in the caption. The second VAE is much simpler. Both the encoder and decoder have three 2048-dimensional hidden layers. Finally, the training details are presented below. Note that these settings were not tuned, we simply chose more epochs for more complex data sets and fewer for datasets with larger training samples. For each dataset just a single setting was tested as follows:

- **MNIST and Fashion-MNIST:** The batch size is specified to be 100. We use the ADAM optimizer with the default hyperparameters in TensorFlow. Standard weight decay is set as 5×10^{-4} . The first VAE is trained for 400 epochs. The initial learning rate is 0.0001 and we halve it every 150 epochs. The second VAE is trained for 800 epochs with the same initial learning rate, halved every 300 epochs.

- **CIFAR-10:** Since CIFAR-10 is more complicated than MNIST and Fashion-MNIST, we use more epochs for training. Specifically, we use 1000 and 2000 epochs for the two VAEs respectively and half the learning rate every 300 and 600 epochs for the two stages. The other settings are the same as that for MNIST.
- **CelebA:** Because CelebA has many more examples, in the first stage we train 120 epochs and half the learning rate every 48 epochs. In the second stage, we train 300 epochs and half the learning rate every 120 epochs. The other settings are the same as that for MNIST, etc.

Finally, to fairly compare against various GAN models and VAE baselines using FID scores on a neutral architecture (i.e., the results from Table 1), we simply adopt the InfoGAN network structure consistent with the neutral setup from (Lucic et al., 2018) for the first-stage VAE. For the second-stage VAE we just use three 1024-dimensional hidden layers, which contribute less than 5% to the total number of parameters. Note that the small number of additional parameters contributing to the second stage do not improve the other VAE baselines when aggregated and trained jointly.

E Proof of Theorem 1

We first consider the case where the latent dimension κ equals the manifold dimension r and then extend the proof to allow for $\kappa > r$. The intuition is to build a bijection between χ and \mathbb{R}^r that transforms the ground-truth distribution $p_{gt}(\mathbf{x})$ to a normal Gaussian distribution. The way to build such a bijection is shown in Figure 13. We now fill in the details.

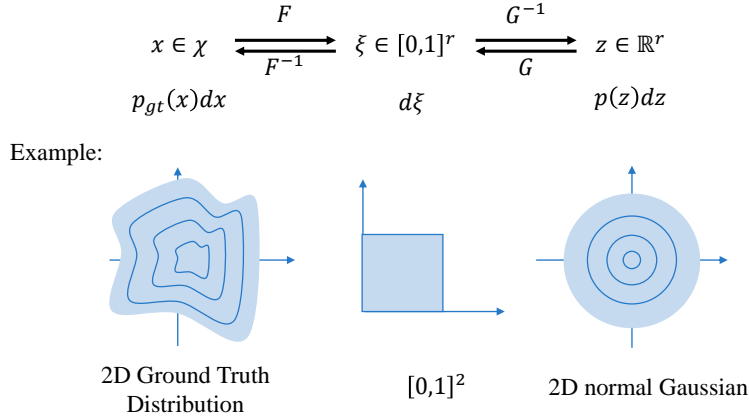


Figure 13: The relationship between different variables.

E.1 Finding a Sequence of Decoders such that $p_{\theta_i^*}(\mathbf{x})$ Converges to $p_{gt}(\mathbf{x})$

Define the function $F : \mathbb{R}^r \mapsto [0, 1]^r$ as

$$F(\mathbf{x}) = [F_1(\mathbf{x}_1), F_2(\mathbf{x}_2; \mathbf{x}_1), \dots, F_r(\mathbf{x}_r; \mathbf{x}_{1:r-1})]^\top, \quad (11)$$

$$F_i(\mathbf{x}_i; \mathbf{x}_{1:i-1}) = \int_{\mathbf{x}'_i=-\infty}^{\mathbf{x}_i} p_{gt}(\mathbf{x}'_i | \mathbf{x}_{1:i-1}) d\mathbf{x}'_i. \quad (12)$$

Per this definition, we have that

$$dF(\mathbf{x}) = p_{gt}(\mathbf{x})d\mathbf{x}. \quad (13)$$

Also, since $p_{gt}(\mathbf{x})$ is nonzero everywhere, $F(\cdot)$ is invertible. Similarly, we define another differentiable and invertible function $G : \mathbb{R}^r \mapsto [0, 1]^r$ as

$$G(\mathbf{z}) = [G_1(\mathbf{z}_1), G_2(\mathbf{z}_2), \dots, G_r(\mathbf{z}_r)]^\top, \quad (14)$$

$$G_i(\mathbf{z}_i) = \int_{\mathbf{z}'_i=-\infty}^{\mathbf{z}_i} \mathcal{N}(\mathbf{z}_i | 0, \mathbf{I}) d\mathbf{z}'_i. \quad (15)$$

Then

$$dG(\mathbf{z}) = p(\mathbf{z})d\mathbf{z} = \mathcal{N}(\mathbf{z} | 0, \mathbf{I})d\mathbf{z}. \quad (16)$$

Now let the decoder be

$$f_{\mu_x}(\mathbf{z}; \theta_t^*) = F^{-1} \circ G(\mathbf{z}), \quad (17)$$

$$\gamma_t^* = \frac{1}{t}. \quad (18)$$

Then we have

$$p_{\theta_t^*}(\mathbf{x}) = \int_{\mathbb{R}^r} p_{\theta_t^*}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} = \int_{\mathbb{R}^r} \mathcal{N}(\mathbf{x}|F^{-1} \circ G(\mathbf{z}), \gamma_t^* \mathbf{I}) dG(\mathbf{z}). \quad (19)$$

Additionally, let $\xi = G(\mathbf{z})$ such that

$$p_{\theta_t^*}(\mathbf{x}) = \int_{[0,1]^r} \mathcal{N}(\mathbf{x}|F^{-1}(\xi), \gamma_t^* \mathbf{I}) d\xi, \quad (20)$$

and let $\mathbf{x}' = F^{-1}(\xi)$ such that $d\xi = dF(\mathbf{x}') = p_{gt}(\mathbf{x}')d\mathbf{x}'$. Plugging this expression into the previous $p_{\theta_t^*}(\mathbf{x})$ we obtain

$$p_{\theta_t^*}(\mathbf{x}) = \int_{\mathbb{R}^r} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) p_{gt}(\mathbf{x}')d\mathbf{x}'. \quad (21)$$

As $t \rightarrow \infty$, γ_t^* becomes infinitely small and $\mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I})$ becomes a Dirac-delta function, resulting in

$$\lim_{t \rightarrow \infty} p_{\theta_t^*}(\mathbf{x}) = \int_{\mathcal{X}} \delta(\mathbf{x}' - \mathbf{x}) p_{gt}(\mathbf{x}')d\mathbf{x}' = p_{gt}(\mathbf{x}). \quad (22)$$

E.2 Finding a Sequence of Encoders such that $\mathbb{KL}[q_{\phi_t^*}(\mathbf{z}|\mathbf{x})||p_{\theta_t^*}(\mathbf{z}|\mathbf{x})]$ Converges to 0

Assume the encoder networks satisfy

$$f_{\mu_x}(\mathbf{x}; \phi_t^*) = G^{-1} \circ F(\mathbf{x}) = f_{\mu_x}^{-1}(\mathbf{x}; \theta_t^*), \quad (23)$$

$$f_{S_z}(\mathbf{x}; \phi_t^*) = \sqrt{\gamma_t^* \left(f'_{\mu_x}(f_{\mu_z}(\mathbf{x}; \phi_t^*); \theta_t^*)^\top f'_{\mu_x}(f_{\mu_z}(\mathbf{x}; \phi_t^*); \theta_t^*) \right)^{-1}}, \quad (24)$$

where $f'_{\mu_x}(\cdot)$ is a $d \times r$ Jacobian matrix. We omit the arguments θ_t^* and ϕ_t^* in $f_{\mu_x}(\cdot)$, $f_{S_z}(\cdot)$ and $f_{\mu_x}(\cdot)$ hereafter to avoid unnecessary clutter. We first explain why $f_{\mu_x}(\cdot)$ is differentiable. Since $f_{\mu_x}(\cdot)$ is a composition of $F^{-1}(\cdot)$ and $G(\cdot)$ according to (17), we only need to explain that both functions are differentiable. For $F^{-1}(\cdot)$, it is the inverse of a differentiable function $F(\cdot)$. Moreover, the derivative of $F(\mathbf{x})$ is $p_{gt}(\mathbf{x})$, which is nonzero everywhere. So $F^{-1}(\cdot)$ and therefore $f_{\mu_x}(\cdot)$ are both differentiable.

The true posterior $p_{\theta_t^*}(\mathbf{z}|\mathbf{x})$ and the approximate posterior are

$$p_{\theta_t^*}(\mathbf{z}|\mathbf{x}) = \frac{\mathcal{N}(\mathbf{z}|0, \mathbf{I})\mathcal{N}(\mathbf{x}|f_{\mu_x}(\mathbf{z}), \gamma_t^* \mathbf{I})}{p_{\theta_t^*}(\mathbf{x})}, \quad (25)$$

$$q_{\phi_t^*}(\mathbf{z}|\mathbf{x}) = \mathcal{N}\left(\mathbf{z}|f_{\mu_z}(\mathbf{x}), \gamma_t^* \left(f'_{\mu_x}(f_{\mu_z}(\mathbf{x}))^\top f'_{\mu_x}(f_{\mu_z}(\mathbf{x})) \right)^{-1}\right) \quad (26)$$

respectively. We now prove that $q_{\phi_t^*}(\mathbf{z}|\mathbf{x})/p_{\theta_t^*}(\mathbf{z}|\mathbf{x})$ converges to a constant not related to \mathbf{z} as t goes to ∞ . If this is true, the constant must be 1 since both $q_{\phi_t^*}(\mathbf{z}|\mathbf{x})$ and $p_{\theta_t^*}(\mathbf{z}|\mathbf{x})$ are probability distributions. Then the KL divergence between them converges to 0 as $t \rightarrow \infty$.

We denote $\left(f'_{\mu_x}(f_{\mu_z}(\mathbf{x}))^\top f'_{\mu_x}(f_{\mu_z}(\mathbf{x})) \right)^{-1}$ as $\tilde{\Sigma}_z(\mathbf{x})$ for short. In addition, we define $\mathbf{z}^* = f_{\mu_z}(\mathbf{x})$. Given these definitions, it follows that

$$\begin{aligned} \frac{q_{\phi_t^*}(\mathbf{z}|\mathbf{x})}{p_{\theta_t^*}(\mathbf{z}|\mathbf{x})} &= \frac{\mathcal{N}(\mathbf{z}|\mathbf{z}^*, \gamma_t^* \tilde{\Sigma}_z) p_{\theta_t^*}(\mathbf{x})}{\mathcal{N}(\mathbf{z}|0, \mathbf{I})\mathcal{N}(\mathbf{x}|f_{\mu_x}(\mathbf{z}), \gamma_t^* \mathbf{I})} \\ &= (2\pi)^{d/2} \gamma_t^{*(d-r)/2} \left| \tilde{\Sigma}_z \right|^{-1/2} \exp \left\{ -\frac{(\mathbf{z} - \mathbf{z}^*)^\top \tilde{\Sigma}_z^{-1} (\mathbf{z} - \mathbf{z}^*)}{2\gamma_t^*} \right. \\ &\quad \left. + \frac{\|\mathbf{z}\|_2^2}{2} + \frac{\|\mathbf{x} - f_{\mu_x}(\mathbf{z})\|_2^2}{2\gamma_t^*} \right\} p_{\theta_t^*}(\mathbf{x}). \end{aligned} \quad (27)$$

At this point, let

$$\mathbf{z} = \mathbf{z}^* + \sqrt{\gamma_t^*} \tilde{\mathbf{z}} \quad (28)$$

According to Lagrangian's mean value theorem, there exists a \mathbf{z}' between \mathbf{z} and \mathbf{z}^* such that

$$f_{\mu_x}(\mathbf{z}) = f_{\mu_x}(\mathbf{z}^*) + f'_{\mu_x}(\mathbf{z}')(\mathbf{z} - \mathbf{z}^*) = \mathbf{x} + f'_{\mu_x}(\mathbf{z}')\sqrt{\gamma_t^*} \tilde{\mathbf{z}}, \quad (29)$$

where $\mathbf{z}' = \mathbf{z}^* + \eta\sqrt{\gamma_t^*} \tilde{\mathbf{z}}$ is between \mathbf{z} and \mathbf{z}^* and η is a value between 0 and 1 ($\mathbf{z}' = \mathbf{z}^*$ if $\eta = 0$ and $\mathbf{z}' = \mathbf{z}$ if $\eta = 1$). Use $C(\mathbf{x})$ to represent the terms not related to \mathbf{z} , i.e., $(2\pi)^{d/2} \gamma_t^{*(d-r)/2} |\tilde{\Sigma}_z|^{-1/2} p_{\theta_t^*}(\mathbf{x})$. Plug (28) and (29) into (27) and consider the limit given by

$$\begin{aligned} \lim_{t \rightarrow \infty} \frac{q_{\phi_t^*}(\mathbf{z}|\mathbf{x})}{p_{\theta_t^*}(\mathbf{z}|\mathbf{x})} &= \lim_{t \rightarrow \infty} C(\mathbf{x}) \exp \left\{ -\frac{\tilde{\mathbf{z}}^\top \tilde{\Sigma}_z^{-1} \tilde{\mathbf{z}}}{2} + \frac{\|\mathbf{z}^* + \sqrt{\gamma_t^*} \tilde{\mathbf{z}}\|_2^2}{2} + \frac{\|f'_{\mu_x}(\mathbf{z}^* + \eta\sqrt{\gamma_t^*} \tilde{\mathbf{z}}) \tilde{\mathbf{z}}\|_2^2}{2} \right\} \\ &= C(\mathbf{x}) \exp \left\{ -\frac{\tilde{\mathbf{z}}^\top \tilde{\Sigma}_z^{-1} \tilde{\mathbf{z}}}{2} + \frac{\|\mathbf{z}^*\|_2^2}{2} + \frac{\|f'_{\mu_x}(\mathbf{z}^*) \tilde{\mathbf{z}}\|_2^2}{2} \right\} \\ &= C(\mathbf{x}) \exp \left\{ -\frac{\tilde{\mathbf{z}}^\top \tilde{\Sigma}_z^{-1} \tilde{\mathbf{z}}}{2} + \frac{\|\mathbf{z}^*\|_2^2}{2} + \frac{\tilde{\mathbf{z}}^\top f'_{\mu_x}(\mathbf{z}^*)^\top f'_{\mu_x}(\mathbf{z}^*) \tilde{\mathbf{z}}}{2} \right\} \\ &= C(\mathbf{x}) \exp \left\{ \frac{\|\mathbf{z}^*\|_2^2}{2} \right\} \end{aligned} \quad (30)$$

The fourth equality comes from the fact that $f'_{\mu_x}(\mathbf{z}^*)^\top f'_{\mu_x}(\mathbf{z}^*) = f'_{\mu_x}(f_{\mu_z}(\mathbf{x}))^\top f'_{\mu_x}(f_{\mu_z}(\mathbf{x})) = \tilde{\Sigma}_z(\mathbf{x})^{-1}$. This expression is not related to \mathbf{z} . Considering both $q_{\phi_t^*}(\mathbf{z}|\mathbf{x})$ and $p_{\theta_t^*}(\mathbf{z}|\mathbf{x})$ are probability distributions, the ratio should be equal to 1. The KL divergence between them thus converges to 0 as $t \rightarrow \infty$.

E.3 Generalization to the Case with $\kappa > r$

When $\kappa > r$, we use the first r latent dimensions to build a projection between \mathbf{z} and \mathbf{x} and leave the remaining $\kappa - r$ latent dimensions unused. Specifically, let $f_{\mu_x}(\mathbf{z}) = \tilde{f}_{\mu_x}(\mathbf{z}_{1:r})$, where $\tilde{f}_{\mu_x}(\mathbf{z}_{1:r})$ is defined as in (17) and $\gamma_t^* = 1/t$. Again consider the case that $t \rightarrow \infty$. Then this decoder can also satisfy $\lim_{t \rightarrow \infty} p_{\theta_t^*}(\mathbf{x}) = p_{gt}(\mathbf{x})$ because it produces exactly the same distribution as the decoder defined by (17) and (18). The last $\kappa - r$ dimensions contribute nothing to the generation process.

Now define the encoder as

$$f_{\mu_z}(\mathbf{x})_{1:r} = \tilde{f}_{\mu_x}^{-1}(\mathbf{x}) \quad (31)$$

$$f_{\mu_z}(\mathbf{x})_{r+1:\kappa} = 0 \quad (32)$$

$$f_{S_z}(\mathbf{x}) = \begin{bmatrix} \tilde{f}_{S_z}(\mathbf{x}) \\ \mathbf{n}_{r+1}^\top \\ \vdots \\ \mathbf{n}_\kappa^\top \end{bmatrix} \quad (33)$$

where $\tilde{f}_{S_z}(\mathbf{x})$ is defined as (24). Denote $\{\mathbf{n}_i\}_{i=r+1}^\kappa$ as a set of κ -dimensional column vectors satisfying

$$\tilde{f}_{S_z}(\mathbf{x}) \mathbf{n}_i = 0 \quad (34)$$

$$\mathbf{n}_i^\top \mathbf{n}_j = \mathbf{1}_{i=j} \quad (35)$$

Such a set always exists because $\tilde{f}_{S_z}(\mathbf{x})$ is a $r \times \kappa$ matrix. So the dimension of the null space of $\tilde{f}_{S_z}(\mathbf{x})$ is at least $\kappa - r$. Assuming that $\{\mathbf{n}_i\}_{i=r+1}^\kappa$ are $\kappa - r$ basis vectors of $\text{null}(\tilde{f}_{S_z})$, then the conditions (34) and (35) will be satisfied. The variance of the approximate posterior then becomes

$$\Sigma_z = f_{S_z}(\mathbf{x}) f_{S_z}(\mathbf{x})^\top = \begin{bmatrix} \tilde{f}_{S_z}(\mathbf{x}) \tilde{f}_{S_z}(\mathbf{x})^\top & 0 \\ 0 & \mathbf{I}_{\kappa-r} \end{bmatrix} \quad (36)$$

The first r dimensions can exactly match the true posterior as we have already shown. The remaining $\kappa - r$ dimensions follow a standardized Gaussian distribution. Since these dimensions contribute nothing to generating \mathbf{x} , the true posterior should be the same as the prior, *i.e.* a standardized Gaussian distribution. Moreover, any of these dimensions is independent of all the other dimensions, so the corresponding off-diagonal elements of the covariance of the true posterior should equal 0. Thus the approximate posterior also matches the true posterior for the last $\kappa - r$ dimensions. As a result, we again have $\lim_{t \rightarrow \infty} \mathbb{KL} [q_{\phi_t^*}(\mathbf{z}|\mathbf{x}) || p_{\theta_t^*}(\mathbf{z}|\mathbf{x})] = 0$.

F Proof of Theorem 2

Similar to Section E, we also construct a bijection between \mathcal{X} and \mathbb{R}^r which transforms the ground-truth measure μ_{gt} to a normal Gaussian distribution. But in this construction, we need one more step that bijects between \mathcal{X} and \mathbb{R}^r using the diffeomorphism $\varphi(\cdot)$, as shown in Figure 14. We will now go into the details.

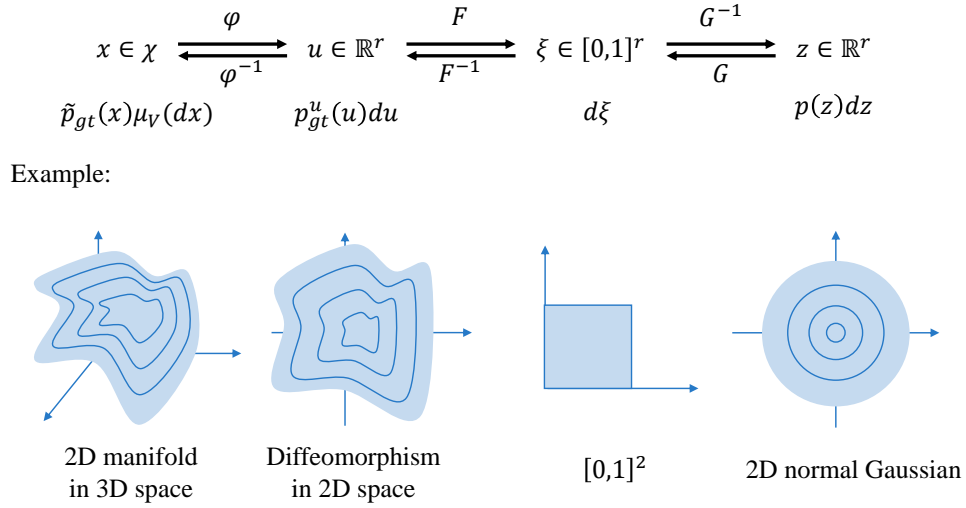


Figure 14: The relationship between different variables.

F.1 Finding a Sequence of Decoders such that $-\log p_{\theta_t^*}(\mathbf{x})$ Converges to $-\infty$

$\varphi(\cdot)$ is a diffeomorphism between \mathcal{X} and \mathbb{R}^r . So it transforms the ground-truth probability distribution $\tilde{p}_{gt}(\mathbf{x})$ to another distribution $p_{gt}^u(\mathbf{u})$, where $\mathbf{u} \in \mathbb{R}^r$. The relationship between the two distributions is

$$p_{gt}^u(\mathbf{u})d\mathbf{u} = \tilde{p}_{gt}(\mathbf{x})\mu_V(d\mathbf{x})|_{\mathbf{x}=\varphi^{-1}(\mathbf{u})} = \mu_{gt}(d\mathbf{x}), \quad (37)$$

where $\mu_V(d\mathbf{x})$ is the volume measure with respect to \mathcal{X} . Because $\varphi(\cdot)$ is a diffeomorphism, both $\varphi(\cdot)$ and $\varphi^{-1}(\cdot)$ are differentiable. Thus $d\mathbf{x}/d\mathbf{u}$ is nonzero everywhere on the manifold. Considering $\tilde{p}_{gt}(\mathbf{x})$ is also nonzero everywhere, $p_{gt}^u(\mathbf{u})$ is nonzero everywhere.

Analogous to the previous proof, define a function $F : \mathbb{R}^r \mapsto [0, 1]^r$ as

$$F(\mathbf{u}) = [F_1(\mathbf{u}_1), F_2(\mathbf{u}_2; \mathbf{u}_1), \dots, F_r(\mathbf{u}_r; \mathbf{u}_{1:r-1})]^\top, \quad (38)$$

$$F_i(\mathbf{u}_i; \mathbf{u}_{1:i-1}) = \int_{\mathbf{u}'_i=-\infty}^{\mathbf{u}_i} p_{gt}^u(\mathbf{u}'_i | \mathbf{u}_{1:i-1}) d\mathbf{u}'_i. \quad (39)$$

According to this definition, we have

$$dF(\mathbf{u}) = p_{gt}^u(\mathbf{u})d\mathbf{u}. \quad (40)$$

Since $p_{gt}^u(\mathbf{u})$ is nonzero everywhere, $F(\cdot)$ is invertible. We also define another differentiable and invertible function $G : \mathbb{R}^r \mapsto [0, 1]^r$ as (15).

Now let the decoder mean function be given by

$$f_{\mu_x}(\mathbf{z}; \theta^*) = \varphi^{-1} \circ F^{-1} \circ G(\mathbf{z}), \quad (41)$$

$$\gamma_t^* = \frac{1}{t}. \quad (42)$$

Then we have

$$\begin{aligned} p_{\theta_t^*}(\mathbf{x}) &= \int_{\mathbb{R}^r} p_{\theta_t^*}(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \\ &= \int_{\mathbb{R}^r} \mathcal{N}(\mathbf{x}|\varphi^{-1}(\mathbf{u}), \gamma_t^* \mathbf{I}) p_{gt}^u(\mathbf{u})d\mathbf{u}. \end{aligned} \quad (43)$$

We next show that $p_{\theta_t^*}(\mathbf{x})$ diverges to infinite as $t \rightarrow \infty$ for any \mathbf{x} . For a given \mathbf{x} , let $\mathbf{u}^* = \varphi(\mathbf{x})$ and $B(\mathbf{u}^*, \sqrt{\gamma_t^*})$ be the closed ball centered at \mathbf{u}^* with radius $\sqrt{\gamma_t^*}$. Then

$$\begin{aligned} p_{\theta_t^*}(\mathbf{x}) &\geq \int_{B(\mathbf{u}^*, \sqrt{\gamma_t^*})} \mathcal{N}(\mathbf{x}|\varphi^{-1}(\mathbf{u}), \gamma_t^* \mathbf{I}) p_{gt}^u(\mathbf{u})d\mathbf{u} \\ &= \int_{B(\mathbf{u}^*, \gamma_t^*)} (2\pi\gamma_t^*)^{-d/2} \exp\left\{-\frac{\|\mathbf{x} - \varphi^{-1}(\mathbf{u})\|_2^2}{2\gamma_t^*}\right\} p_{gt}^u(\mathbf{u})d\mathbf{u}. \end{aligned} \quad (44)$$

According to the Lagrangian's mean value theorem, there exists a \mathbf{u}' between \mathbf{u} and \mathbf{u}^* such that

$$\varphi^{-1}(\mathbf{u}) = \varphi^{-1}(\mathbf{u}^*) + \frac{d\varphi^{-1}(\mathbf{u})}{d\mathbf{u}}|_{\mathbf{u}=\mathbf{u}'}(\mathbf{u} - \mathbf{u}^*) = \mathbf{x} + \frac{d\varphi^{-1}(\mathbf{u})}{d\mathbf{u}}|_{\mathbf{u}=\mathbf{u}'}(\mathbf{u} - \mathbf{u}^*). \quad (45)$$

If we denote $\Lambda(\mathbf{u}') = \left(\frac{d\varphi^{-1}(\mathbf{u})}{d\mathbf{u}}|_{\mathbf{u}=\mathbf{u}'}\right)^\top \left(\frac{d\varphi^{-1}(\mathbf{u})}{d\mathbf{u}}|_{\mathbf{u}=\mathbf{u}'}\right)$, we then have that

$$\begin{aligned} \|\mathbf{x} - \varphi^{-1}(\mathbf{u})\|_2^2 &= (\mathbf{u} - \mathbf{u}^*)^\top \Lambda(\mathbf{u}')(\mathbf{u} - \mathbf{u}^*) = \sum_{i,j} \Lambda(\mathbf{u}')_{i,j} (\mathbf{u}_i - \mathbf{u}_i^*)^\top (\mathbf{u}_j - \mathbf{u}_j^*) \\ &\leq \sum_i \left(\sum_j \Lambda(\mathbf{u}')_{i,j} \right) (\mathbf{u}_i - \mathbf{u}_i^*)^2 \leq \|\Lambda(\mathbf{u}')\|_1 \cdot \|\mathbf{u} - \mathbf{u}^*\|_2^2. \end{aligned} \quad (46)$$

And after defining

$$D(\mathbf{u}^*) = \max_{\mathbf{u} \in B(\mathbf{u}^*, 1)} \|\Lambda(\mathbf{u})\|_1 \geq \max_{\mathbf{u} \in B(\mathbf{u}^*, \sqrt{\gamma_t^*})} \|\Lambda(\mathbf{u})\|, \quad (47)$$

it also follows that

$$\|\mathbf{x} - \varphi^{-1}(\mathbf{u})\|_2^2 \leq \|\Lambda(\mathbf{u}')\|_1 \cdot \|\mathbf{u} - \mathbf{u}^*\|_2^2 \leq D(\mathbf{u}^*)\gamma_t^* \quad \forall \mathbf{u} \in B(\mathbf{u}^*, \sqrt{\gamma_t^*}). \quad (48)$$

Plugging this inequality into (44) gives

$$\begin{aligned} p_{\theta_t^*}(\mathbf{x}) &\geq \int_{B(\mathbf{u}^*, \sqrt{\gamma_t^*})} (2\pi\gamma_t^*)^{-d/2} \exp\left\{-\frac{D(\mathbf{u}^*)\gamma_t^*}{2\gamma_t^*}\right\} p_{gt}^u(\mathbf{u})d\mathbf{u} \\ &\geq (2\pi\gamma_t^*)^{-d/2} \exp\left\{-\frac{D(\mathbf{u}^*)}{2}\right\} \left(\min_{\mathbf{u} \in B(\mathbf{u}^*, 1)} p_{gt}^u(\mathbf{u}) \right) \int_{B(\mathbf{u}^*, \sqrt{\gamma_t^*})} d\mathbf{u} \\ &= (2\pi\gamma_t^*)^{-d/2} \exp\left\{-\frac{D(\mathbf{u}^*)}{2}\right\} \left(\min_{\mathbf{u} \in B(\mathbf{u}^*, 1)} p_{gt}^u(\mathbf{u}) \right) V(B(\mathbf{u}^*, \sqrt{\gamma_t^*})), \end{aligned} \quad (49)$$

where $V(B(\mathbf{u}^*, \sqrt{\gamma_t^*}))$ is the volume of the r -dimensional ball $B(\mathbf{u}^*, \sqrt{\gamma_t^*})$. The volume should be $a_r \gamma_t^{*r/2}$ where a_r is a constant related to the dimension r . So

$$p_{\theta_t^*}(\mathbf{x}) \geq (2\pi)^{-d/2} \gamma_t^{*-(d-r)/2} a_r \exp\left\{-\frac{D(\mathbf{u}^*)}{2}\right\} \left(\min_{\mathbf{u} \in B(\mathbf{u}^*, 1)} p_{gt}^u(\mathbf{u}) \right). \quad (50)$$

Since $\varphi(\cdot)$ defines a diffeomorphism, $D(\mathbf{u}^*) < \infty$. Moreover, $(\min_{\mathbf{u} \in B(\mathbf{u}^*, 1)} p_{gt}^u(\mathbf{u})) > 0$ because $p_{gt}^u(\mathbf{u})$ is nonzero and continuous everywhere. We may then conclude that

$$\lim_{t \rightarrow \infty} -\log p_{\theta_t^*}(\mathbf{x}) = -\infty. \quad (51)$$

for $\mathbf{x} \in \mathcal{X}$. This then implies that the stated average across \mathcal{X} with respect to μ_{gt} will also be $-\infty$.

F.2 Finding a Sequence of Encoders such that $\mathbb{KL} [q_{\phi_t^*}(z|x)||p_{\theta_t^*}(z|x)]$ Converges to 0

Similar to (23) and (24), let the encoder be

$$f_{\mu_z}(\mathbf{x}; \phi_t^*) = G^{-1} \circ F \circ \varphi(\mathbf{x}) = f_{\mu_x}^{-1}(\mathbf{x}; \theta_t^*), \quad (52)$$

$$f_{S_z}(\mathbf{x}; \phi_t^*) = \sqrt{\gamma_t^* \left(f'_{\mu_x}(f_{\mu_z}(\mathbf{x}; \phi_t^*); \theta_t^*)^\top f'_{\mu_x}(f_{\mu_z}(\mathbf{x}; \phi_t^*); \theta_t^*) \right)^{-1}}. \quad (53)$$

Following the proofs in Section E.2, we can prove the KL divergence between $q_{\phi_t^*}(z|x)$ and $p_{\theta_t^*}(z|x)$ converges to 0.

F.3 The Relationship between $\lim_{t \rightarrow \infty} p_{\theta_t^*}(\mathbf{x})$ and $\mu_{gt}(\mathbf{x})$

We then prove our construction from (41) and (42) satisfies (6). Unlike the case $d = r$ where we can compare $p_{\theta_t^*}(\mathbf{x})$ and $p_{gt}(\mathbf{x})$ directly, here $p_{\theta_t^*}(\mathbf{x})$ is a density defined everywhere in \mathbb{R}^d while μ_{gt} is a probability measure defined only on the r -dimensional manifold χ . Consequently, to assess $p_{\theta_t^*}(\mathbf{x})$ relative to μ_{gt} , we evaluate the respective probability mass assigned to any measurable subset of \mathbb{R}^d denoted as A . For $p_{\theta_t^*}(\mathbf{x})$, we integrate the density over A while for μ_{gt} we compute the measure of the intersection of A with χ , i.e., μ_{gt} confines all mass to the manifold.

We begin with the probability distribution given by $p_{\theta_t^*}(\mathbf{x})$:

$$\begin{aligned} p_{\theta_t^*}(\mathbf{x}) &= \int_{\mathbb{R}^r} p_{\theta_t^*}(\mathbf{x}|z)p(z)dz = \int_{\mathbb{R}^r} \mathcal{N}(\mathbf{x}|\varphi^{-1} \circ F^{-1} \circ G(z), \gamma_t^* \mathbf{I}) dG(z) \\ &= \int_{[0,1]^r} \mathcal{N}(\mathbf{x}|\varphi^{-1} \circ F^{-1}(\xi), \gamma_t^* \mathbf{I}) d\xi \\ &= \int_{\mathbb{R}^r} \mathcal{N}(\mathbf{x}|\varphi^{-1}(u), \gamma_t^* \mathbf{I}) p_{gt}^u(u) du \\ &= \int_{\mathbf{x}' \in \chi} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*) \mu_{gt}(d\mathbf{x}'). \end{aligned} \quad (54)$$

Consider a measurable set $A \in \mathbb{R}^d$,

$$\begin{aligned} \lim_{t \rightarrow \infty} \int_{\mathbf{x} \in A} p_{\theta_t^*}(\mathbf{x}) d\mathbf{x} &= \lim_{t \rightarrow \infty} \int_{\mathbf{x} \in A} \left[\int_{\mathbf{x}' \in \chi} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*) \mu_{gt}(d\mathbf{x}') \right] d\mathbf{x} \\ &= \lim_{t \rightarrow \infty} \int_{\mathbf{x}' \in \chi} \left[\int_{\mathbf{x} \in A} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') \\ &= \int_{\mathbf{x}' \in \chi} \lim_{t \rightarrow \infty} \left[\int_{\mathbf{x} \in A} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}'). \end{aligned} \quad (55)$$

The second equation that interchanges the order of the integrations admitted by Fubini's theorem. The third equation that interchanges the order of the integration and the limit is justified by the bounded convergence theorem. We now note that the term inside the first integration, $\mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*)$, converges to a Dirac-delta function as $\gamma_t^* \rightarrow 0$. So the integration over A depends on whether \mathbf{x}' is inside A or not, i.e.,

$$\lim_{t \rightarrow \infty} \left[\int_{\mathbf{x} \in A} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) d\mathbf{x} \right] = \begin{cases} 1 & \text{if } \mathbf{x}' \in A - \partial A, \\ 0 & \text{if } \mathbf{x}' \in A^c - \partial A. \end{cases} \quad (56)$$

We separate the manifold χ into three parts: $\chi \cap (A - \partial A)$, $\chi \cap (A^c - \partial A)$ and $\chi \cap \partial A$. Then (55) can be separated into three parts accordingly. The first two parts can be derived as

$$\int_{\chi \cap (A - \partial A)} \lim_{t \rightarrow \infty} \left[\int_A \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') = \int_{\chi \cap (A - \partial A)} 1 \mu_{gt}(d\mathbf{x}') = \mu_{gt}(\chi \cap (A - \partial A)), \quad (57)$$

$$\int_{\chi \cap (A^c - \partial A)} \lim_{t \rightarrow \infty} \left[\int_A \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') = \int_{\chi \cap (A^c - \partial A)} 0 \mu_{gt}(d\mathbf{x}') = 0. \quad (58)$$

For the third part, given the assumption that $\mu_{gt}(\partial A) = 0$, we have

$$0 \leq \int_{\chi \cap \partial A} \lim_{t \rightarrow \infty} \left[\int_A \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') \leq \int_{\chi \cap \partial A} 1 \mu_{gt}(d\mathbf{x}') = \mu_{gt}(\chi \cap \partial A) = 0. \quad (59)$$

Therefore we have

$$\int_{\chi \cap \partial A} \lim_{t \rightarrow \infty} \left[\int_A \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^* \mathbf{I}) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') = 0 \quad (60)$$

and thus

$$\begin{aligned} \lim_{t \rightarrow \infty} \int_A p_{gt}(\mathbf{x}; \gamma_t^* \mathbf{I}) d\mathbf{x} &= \int_{\mathbf{x}' \in \chi} \lim_{t \rightarrow \infty} \left[\int_{\mathbf{x} \in A} \mathcal{N}(\mathbf{x}|\mathbf{x}', \gamma_t^*) d\mathbf{x} \right] \mu_{gt}(d\mathbf{x}') \\ &= \mu_{gt}(\chi \cap (A - \partial A)) + 0 + 0 \\ &= \mu_{gt}(\chi \cap A), \end{aligned} \quad (61)$$

leading to (6).

Note that this result involves a subtle requirement involving the boundary ∂A . This condition is only included to handle a minor, practically-inconsequential technicality. In brief, as a density $p_{\theta_t^*}(\mathbf{x})$ will apply zero mass exactly on any low-dimensional manifold, although it can apply all of its mass to any region in the neighborhood of χ . But suppose we choose some A is a subset of χ , i.e, it is exclusively confined to the ground-truth manifold. Then the probability mass within A assigned by μ_{gt} will be nonzero while that given by $p_{\theta_t^*}(\mathbf{x})$ can still be zero. Of course this does not mean that $p_{\theta_t^*}(\mathbf{x})$ and μ_{gt} do not match each other in any practical sense. This is because if we expand this specialized A by an arbitrary small d -dimensional volume, then $p_{\theta_t^*}(\mathbf{x})$ and μ_{gt} will now supply essentially the same probability mass on this infinitesimally expanded set (which is arbitrary close to A).

G Proof of Theorem 3

From the main text, $\{\theta_\gamma^*, \phi_\gamma^*\}$ is the optimal solution with a fixed γ . The true posterior and the approximate posterior are

$$p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{z})p_{\theta_\gamma^*}(\mathbf{x}|\mathbf{z})}{p_{\theta_\gamma^*}(\mathbf{x})}, \quad (62)$$

$$q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}_z(\mathbf{x}; \phi_\gamma^*), \boldsymbol{\Sigma}_z(\mathbf{x}; \phi_\gamma^*)). \quad (63)$$

G.1 Case 1: $r = d$

We first argue that the KL divergence between $p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x})$ and $q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})$ is always strictly greater than zero. This can be proved by contradiction. Suppose the KL divergence exactly equals zero. Then $p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x})$ must also be a Gaussian distribution, meaning that the logarithm of $p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x})$ is a quadratic form in \mathbf{z} . In particular, we have

$$\begin{aligned} \log p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x}) &= \log \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) + \log \mathcal{N}(\mathbf{x}|f_{\mu_x}(\mathbf{z}), \gamma \mathbf{I}) - \log p(\mathbf{x}) \\ &= -\frac{1}{2} \|\mathbf{z}\|_2^2 - \frac{1}{2\gamma} \|\mathbf{x} - f_{\mu_x}(\mathbf{z})\|_2^2 + \text{constant}, \end{aligned} \quad (64)$$

where we have absorbed all the terms not related to \mathbf{z} into a constant, and it must be that

$$f_{\mu_x}(\mathbf{z}) = \mathbf{W}\mathbf{z} + \mathbf{b}, \quad (65)$$

for some matrix \mathbf{W} and vector \mathbf{b} . Then we have

$$\begin{aligned} p_{\theta_\gamma^*}(\mathbf{x}) &= \int_{\mathbb{R}^K} p_{\theta_\gamma^*}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \\ &= \int_{\mathbb{R}^K} \mathcal{N}(\mathbf{x}|\mathbf{W}\mathbf{z} + \mathbf{b}, \gamma \mathbf{I}) \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I}) d\mathbf{z}. \end{aligned} \quad (66)$$

This is a Gaussian distribution in \mathbb{R}^d which contradicts our assumption that $p_{gt}(\mathbf{x})$ is not Gaussian. So the KL divergence between $p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x})$ and $q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})$ is always greater than 0. As a result, $\mathcal{L}(\theta_\gamma^*, \phi_\gamma^*)$ cannot reach the theoretical optimal solution, i.e., $\int_{\mathcal{X}} -p_{gt}(\mathbf{x}) \log p_{gt}(\mathbf{x}) d\mathbf{x}$. Denote the gap between $\mathcal{L}(\theta_\gamma^*, \phi_\gamma^*)$ and $\int_{\mathcal{X}} -p_{gt}(\mathbf{x}) \log p_{gt}(\mathbf{x}) d\mathbf{x}$ as ϵ . According to the proof in Section E, there exists a t_0 such that for any $t > t_0$, the gap between the proposed solution in Section E and the theoretical optimal solution is smaller than ϵ . Pick some $t > t_0$ such that $1/t < \gamma$ and let $\gamma' = 1/t$. Then

$$\mathcal{L}(\theta_{\gamma'}^*, \phi_{\gamma'}^*) \leq \mathcal{L}(\theta_t^*, \phi_t^*) < \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*). \quad (67)$$

The first inequality comes from the fact that $\{\theta_{\gamma'}^*, \phi_{\gamma'}^*\}$ is the optimal solution when γ is fixed at γ' while $\{\theta_t^*, \phi_t^*\}$ is just one solution with $\gamma = 1/t = \gamma'$. The second inequality holds because we chose $\{\theta_t^*, \phi_t^*\}$ to be a better solution than $\{\theta_\gamma^*, \phi_\gamma^*\}$.

G.2 Case 2: $r < d$

In this case, $\mathbb{KL}[q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})||p_{\theta_\gamma^*}(\mathbf{z}|\mathbf{x})]$ does not need to be zero because it is possible that $-\log p_{\theta_\gamma^*}(\mathbf{x})$ diverges to negative infinity and absorbs the positive cost caused by the KL divergence. Consider the objective function expression from (2). It can be bounded by

$$\begin{aligned} \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*) &= \int_{\mathcal{X}} \left\{ -\mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} [\log p_{\theta_\gamma^*}(\mathbf{x}|\mathbf{z})] + \mathbb{KL}[q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \right\} \mu_{gt}(d\mathbf{x}) \\ &\geq \int_{\mathcal{X}} \left\{ -\mathbb{E}_{q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} \left[\frac{\|\mathbf{x} - f_{\mu_x}(\mathbf{z})\|_2^2}{2\gamma} + \frac{d}{2} \log(2\pi\gamma) \right] \right\} \mu_{gt}(d\mathbf{x}) \\ &\geq \frac{d}{2} \log \gamma > -\infty. \end{aligned} \quad (68)$$

The first inequality holds discarding the KL term, which is non-negative. The second inequality holds because a quadratic term is removed. Furthermore, according to the proof in Section F, there exists a t_0 such that for any $t > t_0$,

$$\mathcal{L}(\theta_t^*, \phi_t^*) < \frac{d}{2} \log \gamma. \quad (69)$$

Again, we select a $t > t_0$ such that $1/t < \gamma$ and let $\gamma' = 1/t$. Then

$$\mathcal{L}(\theta_{\gamma'}^*, \phi_{\gamma'}^*) \leq \mathcal{L}(\theta_t^*, \phi_t^*) < \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*). \quad (70)$$

H Proof of Theorem 4

Recall that

$$q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z} | f_{\mu_z}(\mathbf{x}; \phi_\gamma^*), f_{S_z}(\mathbf{x}; \phi_\gamma^*) f_{S_z}(\mathbf{x}; \phi_\gamma^*)^\top), \quad (71)$$

$$p_{\theta_\gamma^*}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x} | f_{\mu_x}(\mathbf{z}; \theta_\gamma^*), \gamma \mathbf{I}). \quad (72)$$

Plugging these expressions into (2) we obtain

$$\mathcal{L}(\theta_\gamma^*, \phi_\gamma^*) = \mathbb{E}_{\mathbf{z} \sim q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} \left[\frac{1}{2\gamma} \|\mathbf{x} - f_{\mu_x}(\mathbf{z})\|_2^2 + \frac{d}{2} \log(2\pi\gamma) \right] + \mathbb{KL}[q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})] \quad (73)$$

$$\geq \frac{1}{2\gamma} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|\mathbf{x} - f_{\mu_x}(\mathbf{x}) - f_{S_z}(\mathbf{x})\epsilon\|_2^2 \right] + \frac{d}{2} \log(2\pi\gamma), \quad (74)$$

where we have omitted explicit inclusion of the parameters ϕ_γ^* and θ_γ^* in the functions $f_{\mu_z}(\cdot)$, $f_{S_z}(\cdot)$ and $f_{\mu_x}(\cdot)$ to avoid undue clutter. Now suppose

$$\lim_{\gamma \rightarrow 0} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|\mathbf{x} - f_{\mu_x}(\mathbf{x}) - f_{S_z}(\mathbf{x})\epsilon\|_2^2 \right] = \Delta \neq 0. \quad (75)$$

It then follows that

$$\lim_{\gamma \rightarrow 0} \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*) \geq \lim_{\gamma \rightarrow 0} \frac{\Delta}{2\gamma} + \frac{d}{2} \log(2\pi\gamma) = +\infty, \quad (76)$$

which contradicts the fact that $\mathcal{L}(\theta_\gamma^*, \phi_\gamma^*)$ converges to $-\infty$. So we must have that

$$\lim_{\gamma \rightarrow 0} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \left[\|f_{\mu_x} [f_{\mu_z}(\mathbf{x}) + f_{S_z}(\mathbf{x})\epsilon] - \mathbf{x}\|_2^2 \right] = 0. \quad (77)$$

Because the term inside the expectation, i.e., $\|f_{\mu_x} [f_{\mu_z}(\mathbf{x}) + f_{S_z}(\mathbf{x})\epsilon] - \mathbf{x}\|_2^2$, is always non-negative, we can conclude that

$$\lim_{\gamma \rightarrow 0} f_{\mu_x} [f_{\mu_z}(\mathbf{x}) + f_{S_z}(\mathbf{x})\epsilon] = \mathbf{x}. \quad (78)$$

And if we let $\epsilon = 0$, this equation then becomes

$$\lim_{\gamma \rightarrow 0} f_{\mu_x} [f_{\mu_z}(\mathbf{x})] = \mathbf{x}. \quad (79)$$

I Further Analysis of the VAE Cost as γ becomes small

In the main paper, we mentioned that the squared eigenvalues of $f_{S_z}(\mathbf{x}; \phi_\gamma^*)$ will become arbitrary small at a rate proportional to γ . To justify this, we borrow the simplified notation from the proof of Theorem 4 and expand $f_{\mu_x}(\mathbf{z})$ at $\mathbf{z} = f_{\mu_z}(\mathbf{x})$ using a Taylor series. Omitting the high order terms (in the present narrow context around the neighborhood of VAE global optima these will be small), this gives

$$f_{\mu_x}(\mathbf{z}) \approx f_{\mu_x} [f_{\mu_z}(\mathbf{x})] + f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] (\mathbf{z} - f_{\mu_z}(\mathbf{x})) \approx \mathbf{x} + f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] (\mathbf{z} - f_{\mu_z}(\mathbf{x})). \quad (80)$$

Plug this expression and (71) into (73), we obtain

$$\begin{aligned} \mathcal{L}(\theta_\gamma^*, \phi_\gamma^*) &\approx \mathbb{E}_{\mathbf{z} \sim q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} \left[\frac{1}{2\gamma} \|f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] (\mathbf{z} - f_{\mu_z}(\mathbf{x}))\|_2^2 + \frac{d}{2} \log(2\pi\gamma) \right] \\ &\quad + \frac{1}{2} \{ \|f_{\mu_z}(\mathbf{x})\|_2^2 + \text{tr}(f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top) - \log |f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top| - \kappa \} \\ &= \frac{1}{2\gamma} \text{tr} \left(\mathbb{E}_{\mathbf{z} \sim q_{\phi_\gamma^*}(\mathbf{z}|\mathbf{x})} [(\mathbf{z} - f_{\mu_z}(\mathbf{x}))^\top (\mathbf{z} - f_{\mu_z}(\mathbf{x}))] f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]^\top f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] \right) \\ &\quad + \frac{d}{2} \log(2\pi\gamma) + \frac{1}{2} \{ \|f_{\mu_z}(\mathbf{x})\|_2^2 + \text{tr}(f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top) - \log |f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top| - \kappa \} \\ &= \text{tr} \left(f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top \left[\frac{1}{2} \mathbf{I} + \frac{1}{2\gamma} f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]^\top f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] \right] \right) \\ &\quad + \frac{d}{2} \log(2\pi\gamma) + \frac{1}{2} \{ \|f_{\mu_z}(\mathbf{x})\|_2^2 - \log |f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top| - \kappa \}. \end{aligned} \quad (81)$$

From these manipulations we may conclude that the optimal value of $f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top$ must satisfy

$$\left[\frac{1}{2} \mathbf{I} + \frac{1}{2\gamma} f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]^\top f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] \right] - \frac{1}{2} (f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top)^{-1} = 0. \quad (82)$$

So

$$f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top = \left[\mathbf{I} + \frac{1}{\gamma} f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]^\top f'_{\mu_x} [f_{\mu_z}(\mathbf{x})] \right]^{-1}. \quad (83)$$

Note that $f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]$ is the tangent space of the manifold χ at $f_{\mu_x} [f_{\mu_z}(\mathbf{x})]$, so the rank must be r . $f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]^\top f'_{\mu_x} [f_{\mu_z}(\mathbf{x})]$ can be decomposed as $\mathbf{U}^\top \mathbf{S} \mathbf{U}$, where \mathbf{U} is a κ -dimensional orthogonal matrix and \mathbf{S} is a κ -dimensional diagonal matrix with r nonzero elements. Denote $\text{diag}[\mathbf{S}] = [S_1, S_2, \dots, S_r, 0, \dots, 0]$. Then

$$f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top = \left[\mathbf{U}^\top \text{diag} \left[1 + \frac{S_1}{\gamma}, \dots, 1 + \frac{S_r}{\gamma}, 1, \dots, 1 \right] \mathbf{U} \right]^{-1}. \quad (84)$$

Case 1: $r = \kappa$. In this case, \mathbf{S} has no nonzero diagonal elements, and therefore

$$f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top = \left[\mathbf{U}^\top \text{diag} \left[\frac{1}{1 + \frac{S_1}{\gamma}}, \dots, \frac{1}{1 + \frac{S_r}{\gamma}} \right] \mathbf{U} \right]. \quad (85)$$

As $\gamma \rightarrow 0$, the eigenvalues of $f_{S_z}(\mathbf{x})f_{S_z}(\mathbf{x})^\top$, which are given by $\frac{1}{1+\frac{s_z}{\gamma}}$, converge to 0 at a rate of $O(\gamma)$.

Case 2: $r < \kappa$. In this case, the first r eigenvalues also converge to 0 at a rate of $O(\gamma)$, but the remaining $\kappa - r$ eigenvalues will be 1, meaning the redundant dimensions are simply filled with noise matching the prior $p(\mathbf{z})$ as desired.

References

- Arjovsky, Martin, Chintala, Soumith, and Bottou, Léon. Wasserstein generative adversarial networks. In *International Conference on Machine Learning*, pp. 214–223, 2017.
- Bengio, Yoshua, Courville, Aaron, and Vincent, Pascal. Representation learning: A review and new perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2013.
- Berthelot, David, Schumm, Thomas, and Metz, Luke. BEGAN: Boundary equilibrium generative adversarial networks. *arXiv:1703.10717*, 2017.
- Brock, Andrew, Lim, Theodore, Ritchie, James M, and Weston, Nick. Neural photo editing with introspective adversarial networks. *arXiv:1609.07093*, 2016.
- Burda, Yuri, Grosse, Roger, and Salakhutdinov, Ruslan. Importance weighted autoencoders. *arXiv:1509.00519*, 2015.
- Chen, Xi, Duan, Yan, Houthoofd, Rein, Schulman, John, Sutskever, Ilya, and Abbeel, Pieter. InfoGAN: Interpretable representation learning by information maximizing generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2172–2180, 2016.
- Dai, Bin, Wang, Yu, Aston, John, Hua, Gang, and Wipf, David. Connections with robust PCA and the role of emergent sparsity in variational autoencoder models. *Journal of Machine Learning Research*, 2018.
- Doersch, Carl. Tutorial on variational autoencoders. *arXiv:1606.05908*, 2016.
- Dosovitskiy, Alexey and Brox, Thomas. Generating images with perceptual similarity metrics based on deep networks. In *Advances in Neural Information Processing Systems*, pp. 658–666, 2016.
- Fedus, William, Rosca, Mihaela, Lakshminarayanan, Balaji, Dai, Andrew M, Mohamed, Shakir, and Goodfellow, Ian. Many paths to equilibrium: GANs do not need to decrease a divergence at every step. *arXiv:1710.08446*, 2017.
- Goodfellow, Ian, Pouget-Abadie, Jean, Mirza, Mehdi, Xu, Bing, Warde-Farley, David, Ozair, Sherjil, Courville, Aaron, and Bengio, Yoshua. Generative adversarial nets. In *Advances in Neural Information Processing Systems*, pp. 2672–2680, 2014.
- Gulrajani, Ishaan, Ahmed, Faruk, Arjovsky, Martin, Dumoulin, Vincent, and Courville, Aaron. Improved training of Wasserstein GANs. In *Advances in Neural Information Processing Systems*, pp. 5767–5777, 2017.
- Heusel, Martin, Ramsauer, Hubert, Unterthiner, Thomas, Nessler, Bernhard, and Hochreiter, Sepp. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Advances in Neural Information Processing Systems*, pp. 6626–6637, 2017.
- Kingma, Diederik and Welling, Max. Auto-encoding variational Bayes. In *International Conference on Learning Representations*, 2014.
- Kingma, Diederik, Salimans, Tim, Jozefowicz, Rafal, Chen, Xi, Sutskever, Ilya, and Welling, Max. Improved variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, pp. 4743–4751, 2016.
- Kodali, Naveen, Abernethy, Jacob, Hays, James, and Kira, Zsolt. On convergence and stability of GANs. *arXiv:1705.07215*, 2017.

- Krizhevsky, Alex and Hinton, Geoffrey. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Larsen, Anders Boesen Lindbo, Sønderby, Søren Kaae, Larochelle, Hugo, and Winther, Ole. Autoencoding beyond pixels using a learned similarity metric. *arXiv:1512.09300*, 2015.
- LeCun, Yann, Bottou, Léon, Bengio, Yoshua, and Haffner, Patrick. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Liu, Ziwei, Luo, Ping, Wang, Xiaogang, and Tang, Xiaoou. Deep learning face attributes in the wild. In *IEEE International Conference on Computer Vision*, pp. 3730–3738, 2015.
- Lucic, Mario, Kurach, Karol, Michalski, Marcin, Gelly, Sylvain, and Bousquet, Olivier. Are GANs created equal? A large-scale study. *Advances in Neural Information Processing Systems*, 2018.
- Makhzani, Alireza, Shlens, Jonathon, Jaitly, Navdeep, Goodfellow, Ian, and Frey, Brendan. Adversarial autoencoders. *arXiv:1511.05644*, 2016.
- Mao, Xudong, Li, Qing, Xie, Haoran, Lau, Raymond YK, Wang, Zhen, and Smolley, Stephen Paul. Least squares generative adversarial networks. In *IEEE International Conference on Computer Vision*, pp. 2813–2821, 2017.
- Rezende, Danilo Jimenez and Mohamed, Shakir. Variational inference with normalizing flows. *arXiv:1505.05770*, 2015.
- Rezende, D.J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning*, 2014.
- Theis, L, van den Oord, A, and Bethge, M. A note on the evaluation of generative models. In *International Conference on Learning Representations*, pp. 1–10, 2016.
- Tolstikhin, Ilya, Bousquet, Olivier, Gelly, Sylvain, and Schoelkopf, Bernhard. Wasserstein autoencoders. *International Conference on Learning Representations*, 2018.
- Tomczak, Jakub and Welling, Max. VAE with a VampPrior. In *International Conference on Artificial Intelligence and Statistics*, pp. 1214–1223, 2018.
- van den Berg, Rianne, Hasenclever, Leonard, Tomczak, Jakub M, and Welling, Max. Sylvester normalizing flows for variational inference. In *Uncertainty in Artificial Intelligence*, 2018.
- Xiao, Han, Rasul, Kashif, and Vollgraf, Roland. Fashion-MNIST: A novel image dataset for benchmarking machine learning algorithms. *arXiv:1708.07747*, 2017.