# The Usual Suspects?
# Reassessing Blame for VAE Posterior Collapse

**Bin Dai**                                    DAIB13@MAILS.TSINGHUA.EDU.CN
*Institute for Advanced Study*
*Tsinghua University*
*Beijing, China*

**Ziyu Wang**                                  WZY196@GMAIL.COM>
*Department of Computer Science and Technology*
*Tsinghua University*
*Beijing, China*

**David Wipf**                                 DAVIDWIPF@GMAIL.COM
*Microsoft Research*
*Beijing, China*

## Abstract

In narrow asymptotic settings Gaussian VAE models of continuous data have been shown to possess global optima aligned with ground-truth distributions. Even so, it is well known that poor solutions whereby the latent posterior collapses to an uninformative prior are sometimes obtained in practice. However, contrary to conventional wisdom that largely assigns blame for this phenomena on the undue influence of KL-divergence regularization, we will argue that posterior collapse is, at least in part, a direct consequence of bad local minima inherent to the loss surface of deep autoencoder networks. In particular, we prove that even small nonlinear perturbations of affine VAE decoder models can produce such minima, and in deeper models, analogous minima can force the VAE to behave like an aggressive truncation operator, provably discarding information along all latent dimensions in certain circumstances. Regardless, the underlying message here is not meant to undercut valuable existing explanations of posterior collapse, but rather, to refine the discussion and elucidate alternative risk factors that may have been previously underappreciated.

## 1. Introduction

The variational autoencoder (VAE) (Kingma & Welling, 2014; Rezende et al., 2014) represents a powerful generative model of data points that are assumed to possess some complex yet unknown latent structure. This assumption is instantiated via the marginalized distribution

$$p_\theta(\boldsymbol{x}) = \int p_\theta(\boldsymbol{x}|\boldsymbol{z})p(\boldsymbol{z})d\boldsymbol{z}, \tag{1}$$

which forms the basis of prevailing VAE models. Here $\boldsymbol{z} \in \mathbb{R}^\kappa$ is a collection of unobservable latent factors of variation that, when drawn from the prior $p(\boldsymbol{z})$, are colloquially said to generate an observed data point $\boldsymbol{x} \in \mathbb{R}^d$ through the conditional distribution $p_\theta(\boldsymbol{x}|\boldsymbol{z})$. The latter is controlled by parameters $\theta$ that can, at least conceptually speaking, be optimized by maximum likelihood over $p_\theta(\boldsymbol{x})$ given available training examples.

In particular, assuming $n$ training points $\boldsymbol{X} = [\boldsymbol{x}^{(1)}, \ldots, \boldsymbol{x}^{(n)}]$, maximum likelihood estimation is tantamount to minimizing the negative log-likelihood expression $\frac{1}{n}\sum_i -\log\left[p_\theta\left(\boldsymbol{x}^{(i)}\right)\right]$. Proceeding further, because the marginalization over $\boldsymbol{z}$ in (1) is often intractable, the VAE instead minimizes a convenient variational upper bound given by $\quad \mathcal{L}(\theta,\phi) \triangleq$

$$\frac{1}{n}\sum_{i=1}^{n}\left\{-\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)})}\left[\log p_\theta\left(\boldsymbol{x}^{(i)}|\boldsymbol{z}\right)\right] + \mathbb{KL}\left[q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)}||p(\boldsymbol{z})\right]\right\} \geq \frac{1}{n}\sum_{i=1}^{n}-\log\left[p_\theta\left(\boldsymbol{x}^{(i)}\right)\right], \tag{2}$$

with equality iff $q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)}) = p_\theta(\boldsymbol{z}|\boldsymbol{x}^{(i)})$ for all $i$. The additional parameters $\phi$ govern the shape of the variational distribution $q_\phi(\boldsymbol{z}|\boldsymbol{x})$ that is designed to approximate the true but often intractable latent posterior $p_\theta(\boldsymbol{z}|\boldsymbol{x})$.

The VAE energy from (2) is composed of two terms, a data-fitting loss that borrows the basic structure of an autoencoder (AE), and a KL-divergence-based regularization factor. The former incentivizes assigning high probability to latent codes $\boldsymbol{z}$ that facilitate accurate reconstructions of each $\boldsymbol{x}^{(i)}$. In fact, if $q_\phi(\boldsymbol{z}|\boldsymbol{x})$ is a Dirac delta function, this term is exactly equivalent to a deterministic AE with data reconstruction loss defined by $-\log p_\theta(\boldsymbol{x}|\boldsymbol{z})$. Overall, it is because of this association that $q_\phi(\boldsymbol{z}|\boldsymbol{x})$ is generally referred to as the *encoder* distribution, while $p_\theta(\boldsymbol{x}|\boldsymbol{z})$ denotes the *decoder* distribution. Additionally, the KL regularizer $\mathbb{KL}[q_\phi(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})]$ pushes the encoder distribution towards the prior without violating the variational bound.

For continuous data, which will be our primary focus herein, it is typical to assume that

$$p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{0},\boldsymbol{I}), \quad p_\theta(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}|\boldsymbol{\mu}_x,\gamma\boldsymbol{I}), \text{ and } q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{\mu}_z,\boldsymbol{\Sigma}_z), \tag{3}$$

where $\gamma > 0$ is a scalar variance parameter, while the Gaussian moments $\boldsymbol{\mu}_x \equiv \boldsymbol{\mu}_x(\boldsymbol{z};\theta)$, $\boldsymbol{\mu}_z \equiv \boldsymbol{\mu}_z(\boldsymbol{x};\phi)$, and $\boldsymbol{\Sigma}_z \equiv \text{diag}[\boldsymbol{\sigma}_z(\boldsymbol{x};\phi)]^2$ are computed via feedforward neural network layers. The encoder network parameterized by $\phi$ takes $\boldsymbol{x}$ as an input and outputs $\boldsymbol{\mu}_z$ and $\boldsymbol{\Sigma}_z$. Similarly the decoder network parameterized by $\theta$ converts a latent code $\boldsymbol{z}$ into $\boldsymbol{\mu}_x$. Given these assumptions, the generic VAE objective from (2) can be refined to

$$\mathcal{L}(\theta,\phi) = \frac{1}{n}\sum_{i=1}^{n}\left\{\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)})}\left[\frac{1}{\gamma}\|\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x(\boldsymbol{z};\theta)\|_2^2\right]\right. \tag{4}$$
$$\left. + d\log\gamma + \left\|\boldsymbol{\sigma}_z\left(\boldsymbol{x}^{(i)};\phi\right)\right\|_2^2 - \log\left|\text{diag}\left[\boldsymbol{\sigma}_z\left(\boldsymbol{x}^{(i)};\phi\right)\right]^2\right| + \left\|\boldsymbol{\mu}_z\left(\boldsymbol{x}^{(i)};\phi\right)\right\|_2^2\right\},$$

excluding an inconsequential factor of $1/2$. This expression can be optimized over using SGD and a simple reparameterization strategy (Kingma & Welling, 2014; Rezende et al., 2014) to produce parameter estimates $\{\theta^*,\phi^*\}$. Among other things, new samples approximating the training data can then be generated via the ancestral process $\boldsymbol{z}^{new} \sim \mathcal{N}(\boldsymbol{z}|\boldsymbol{0},\boldsymbol{I})$ and $\boldsymbol{x}^{new} \sim p_{\theta^*}(\boldsymbol{x}|\boldsymbol{z}^{new})$.

Although it has been argued that global minima of (4) may correspond with the optimal recovery of ground truth distributions in certain asymptotic settings (Dai & Wipf, 2019), it is well known that in practice, VAE models are at risk of converging to degenerate solutions where, for example, it may be that $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z})$. This phenomena, commonly referred to as VAE posterior collapse (He et al., 2019; Razavi et al., 2019), has been acknowledged and analyzed from a variety of different perspectives as we detail in Section 2. That being

said, we would argue that there remains lingering ambiguity regarding the different types and respective causes of posterior collapse. Consequently, Section 3 provides a useful taxonomy that will serve to contextualize our main technical contributions. These include the following:

- Building upon existing analysis of affine VAE decoder models, in Section 4 we prove that even arbitrarily small nonlinear activations can introduce suboptimal local minima exhibiting posterior collapse.

- We demonstrate in Section 5 that if the encoder/decoder networks are incapable of sufficiently reducing the VAE reconstruction errors, even in a deterministic setting with no KL-divergence regularizer, there will exist an implicit lower bound on the optimal value of $\gamma$. Moreover, we prove that if this $\gamma$ is sufficiently large, the VAE will behave like an aggressive thresholding operator, enforcing exact posterior collapse, i.e., $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z})$.

- Based on these observations, we present experiments in Section 6 establishing that as network depth/capacity is increased, even for deterministic AE models with no regularization, reconstruction errors become worse. This bounds the effective VAE trade-off parameter $\gamma$ such that posterior collapse is essentially inevitable. Collectively then, we provide convincing evidence that posterior collapse is, at least in certain settings, the fault of deep AE local minima, and need not be exclusively a consequence of usual suspects such as the KL-divergence term.

We conclude in Section 7 with practical take-home messages, and motivate the search for improved AE architectures and training regimes that might be leveraged by analogous VAE models.

## 2. Recent Work and the Usual Suspects for Instigating Collapse

Posterior collapse under various guises is one of the most frequently addressed topics related to VAE performance. Depending on the context, arguably the most common and seemingly transparent suspect for causing collapse is the KL regularization factor that is obviously minimized by $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z})$. This perception has inspired various countermeasures, including heuristic annealing of the KL penalty or KL warm-start (Bowman et al., 2015; Huang et al., 2018; Sønderby et al., 2016), tighter bounds on the log-likelihood (Burda et al., 2015; Rezende & Mohamed, 2015), more complex priors (Bauer & Mnih, 2018; Tomczak & Welling, 2018), modified decoder architectures (Cai et al., 2017; Dieng et al., 2018; Yeung et al., 2017), or efforts to explicitly disallow the prior from ever equaling the variational distribution (Razavi et al., 2019). Thus far though, most published results do not indicate success generating high-resolution images, and in the majority of cases, evaluations are limited to small images and/or relatively shallow networks. This suggests that there may be more nuance involved in pinpointing the causes and potential remedies of posterior collapse. One notable exception though is the BIVA model from (Maaløe et al., 2019), which employs a bidirectional hierarchy of latent variables, in part to combat posterior collapse. While improvements in NLL scores have been demonstrated with BIVA using relatively deep encoder/decoders, this model is significantly more complex and difficult to analyze.

3

On the analysis side, there have been various efforts to explicitly characterize posterior collapse in restricted settings. For example, Lucas et al. (2019) demonstrate that if $\gamma$ is fixed to a sufficiently large value, then a VAE energy function with an affine decoder mean will have minima that overprune latent dimensions. A related linearized approximation to the VAE objective is analyzed in (Rolinek et al., 2019); however, collapsed latent dimensions are excluded and it remains somewhat unclear how the surrogate objective relates to the original. Posterior collapse has also been associated with data-dependent decoder covariance networks $\boldsymbol{\Sigma}_x(\boldsymbol{z}; \theta) \neq \gamma \boldsymbol{I}$ (Mattei & Frellsen, 2018), which allows for degenerate solutions assigning infinite density to a single data point and a diffuse, collapsed density everywhere else. Finally, from the perspective of training dynamics, (He et al., 2019) argue that a lagging inference network can also lead to posterior collapse.

## 3. Taxonomy of Posterior Collapse

Although there is now a vast literature on the various potential causes of posterior collapse, there remains ambiguity as to exactly what this phenomena is referring to. In this regard, we believe that it is critical to differentiate five subtle yet quite distinct scenarios that could reasonably fall under the generic rubric of posterior collapse:

(i) Latent dimensions of $\boldsymbol{z}$ that are not needed for providing good reconstructions of the training data are set to the prior, meaning $q_\phi(z_j|\boldsymbol{x}) \approx p(z_j) = \mathcal{N}(0, 1)$ at any superfluous dimension $j$. Along other dimensions $\boldsymbol{\sigma}_z^2$ will be near zero and $\boldsymbol{\mu}_z$ will provide a usable predictive signal leading to accurate reconstructions of the training data. This case can actually be viewed as a desirable form of *selective* posterior collapse that, as argued in (Dai & Wipf, 2019), is a necessary (albeit not sufficient) condition for generating good samples.

(ii) The decoder variance $\gamma$ is not learned but fixed to a large value[1] such that the KL term from (2) is overly dominant, forcing most or all dimensions of $\boldsymbol{z}$ to follow the prior $\mathcal{N}(0, 1)$. In this scenario, the actual global optimum of the VAE energy (conditioned on $\gamma$ being fixed) will lead to deleterious posterior collapse and the model reconstructions of the training data will be poor. In fact, even the original marginal log-likelihood can potentially default to a trivial/useless solution if $\gamma$ is fixed too large, assigning a small marginal likelihood to the training data, provably so in the affine case (Lucas et al., 2019).

(iii) As mentioned previously, if the Gaussian decoder covariance is learned as a separate network structure (instead of simply $\boldsymbol{\Sigma}_x(\boldsymbol{z}; \theta) = \gamma \boldsymbol{I}$), there can exist degenerate solutions that assign infinite density to a single data point and a diffuse, isotropic Gaussian elsewhere (Mattei & Frellsen, 2018). This implies that (4) can be unbounded from below at what amounts to a posterior collapsed solution and bad reconstructions almost everywhere.

(iv) When powerful non-Gaussian decoders are used, and in particular those that can parameterize complex distributions regardless of the value of $\boldsymbol{z}$ (e.g., PixelCNN-based

1. Or equivalently, a KL scaling parameter such as used by the $\beta$-VAE (Higgins et al., 2017) is set too large.

(Van den Oord et al., 2016)), it is possible for the VAE to assign high-probability to the training data even if $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z})$ (Alemi et al., 2017; Bowman et al., 2015; Chen et al., 2016). This category of posterior collapse is quite distinct from categories (ii) and (iii) above in that, although the reconstructions are similarly poor, the associated NLL scores can still be good.

(v) The previous four categories of posterior collapse can all be directly associated with emergent properties of the VAE *global* minimum under various modeling conditions. In contrast, a fifth type of collapse exists that is the explicit progeny of bad VAE *local* minima. More specifically, as we will argue shortly, when deeper encoder/decoder networks are used, the risk of converging to bad, overregularized solutions increases.

The remainder of this paper will primarily focus on category (v), with brief mention of the other types for comparison purposes where appropriate. Our rationale for this selection bias is that, unlike the others, category (i) collapse is actually advantageous and hence need not be mitigated. In contrast, while category (ii) is undesirable, it be can be avoided by learning $\gamma$. As for category (iii), this represents an unavoidable consequence of models with flexible decoder covariances capable of detecting outliers (Dai et al., 2019). In fact, even simpler inlier/outlier decomposition models such as robust PCA are inevitably at risk for this phenomena (Candès et al., 2011). Regardless, when $\boldsymbol{\Sigma}_z(\boldsymbol{x}; \theta) = \gamma \boldsymbol{I}$ this problem goes away. And finally, we do not address category (iv) in depth simply because it is unrelated to the canonical Gaussian VAE models of continuous data that we have chosen to examine herein. Regardless, it is still worthwhile to explicitly differentiate these five types and bare them in mind when considering attempts to both explain and improve VAE models.

## 4. Insights from Simplified Cases

Because different categories of posterior collapse can be impacted by different global/local minima structures, a useful starting point is a restricted setting whereby we can comprehensively characterize all such minima. For this purpose, we first consider a VAE model with the decoder network set to an affine function. As is often assumed in practice, we choose $\boldsymbol{\Sigma}_x = \gamma \boldsymbol{I}$, where $\gamma > 0$ is a scalar parameter within the parameter set $\theta$. In contrast, for the mean function we choose $\boldsymbol{\mu}_x = \boldsymbol{W}_x \boldsymbol{z} + \boldsymbol{b}_x$ for some weight matrix $\boldsymbol{W}_x$ and bias vector $\boldsymbol{b}_x$. The encoder can be arbitrarily complex (although the optimal structure can be shown to be affine as well).

Given these simplifications, and assuming the training data has $r \geq \kappa$ nonzero singular values, it has been demonstrated that at any global optima, the columns of $\boldsymbol{W}_x$ will correspond with the first $\kappa$ principal components of $\boldsymbol{X}$ provided that we simultaneously learn $\gamma$ or set it to the optimal value (which is available in closed form) (Dai et al., 2019; Lucas et al., 2019; Tipping & Bishop, 1999). Additionally, it has also be shown that no spurious, suboptimal local minima will exist. Note also that if $r < \kappa$ the same basic conclusions still apply; however, $\boldsymbol{W}_x$ will only have $r$ nonzero columns, each corresponding with a different principal component of the data. The unused latent dimensions will satisfy $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, which represents the canonical form of the benign category (i) posterior collapse. Collectively, these results imply that if we converge to any local minima of the VAE energy, we will obtain the best possible linear approximation to the data using a

minimal number of latent dimensions, and malignant posterior collapse is not an issue, i.e., categories (ii)-(v) will not arise.

Even so, if instead of learning $\gamma$, we choose a fixed value that is larger than any of the significant singular values of $\boldsymbol{X}\boldsymbol{X}^\top$, then category (ii) posterior collapse can be inadvertently introduced. More specifically, let $\tilde{r}_\gamma$ denote the number of such singular values that are smaller than some fixed $\gamma$ value. Then along $\kappa - \tilde{r}_\gamma$ latent dimensions $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{0}, \boldsymbol{I})$, and the corresponding columns of $\boldsymbol{W}_x$ will be set to zero at the global optima (conditioned on this fixed $\gamma$), regardless of whether or not these dimensions are necessary for accurately reconstructing the data. And it has been argued that the risk of this type of posterior collapse at a conditionally-optimal global minimum will likely be inherited by deeper models as well (Lucas et al., 2019), although learning $\gamma$ can ameliorate this problem.

Of course when we move to more complex architectures, the risk of bad *local* minima or other suboptimal stationary points becomes a new potential concern, and it is not clear that the affine case described above contributes to reliable, predictive intuitions. To illustrate this point, we will now demonstrate that the introduction of an arbitrarily small nonlinearity can nonetheless produce a pernicious local minimum that exhibits category (v) posterior collapse. For this purpose, we assume the decoder mean function

$$\boldsymbol{\mu}_x = \pi_\alpha\left(\boldsymbol{W}_x\boldsymbol{z}\right) + \boldsymbol{b}_x, \text{ with } \pi_\alpha(u) \triangleq \text{sign}(u)\left(|u| - \alpha\right)_+, \ \alpha \geq 0. \tag{5}$$

The function $\pi_\alpha$ is nothing more than a soft-threshold operator as is commonly used in neural network architectures designed to reflect unfolded iterative algorithms for representation learning (Gregor & LeCun, 2010; Sprechmann et al., 2015). In the present context though, we choose this nonlinearity largely because it allows (5) to reflect arbitrarily small perturbations away from a strictly affine model, and indeed if $\alpha = 0$ the exact affine model is recovered. Collectively, these specifications lead to the parameterization $\theta = \{\boldsymbol{W}_x, \boldsymbol{b}_x, \gamma\}$ and $\phi = \{\boldsymbol{\mu}_z^{(i)}, \boldsymbol{\sigma}_z^{(i)}\}_{i=1}^n$ and energy (excluding irrelevant scale factors and constants) given by

$$\mathcal{L}(\theta, \phi) \ = \ \sum_{i=1}^n \left\{ \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)})} \left[ \frac{1}{\gamma} \left\| \boldsymbol{x}^{(i)} - \pi_\alpha\left(\boldsymbol{W}_x\boldsymbol{z}\right) - \boldsymbol{b}_x \right\|_2^2 \right] \right. \tag{6}$$

$$\left. + d\log\gamma + \left\| \boldsymbol{\sigma}_z^{(i)} \right\|_2^2 - \log\left| \text{diag}\left[ \boldsymbol{\sigma}_z^{(i)} \right]^2 \right| + \left\| \boldsymbol{\mu}_z^{(i)} \right\|_2^2 \right\},$$

where $\boldsymbol{\mu}_z^{(i)}$ and $\boldsymbol{\sigma}_z^{(i)}$ denote arbitrary encoder moments for data point $i$ (this is consistent with the assumption of an arbitrarily complex encoder as used in previous analysis of affine decoder models). Now define $\bar{\gamma} \triangleq \frac{1}{nd}\sum_i \|\boldsymbol{x}^{(i)} - \bar{\boldsymbol{x}}\|_2^2$, with $\bar{\boldsymbol{x}} \triangleq \frac{1}{n}\sum_i \boldsymbol{x}^{(i)}$. We then have the following result (all proofs are deferred to the appendices):

**Proposition 1** *For any $\alpha > 0$, there will always exist data sets $\boldsymbol{X}$ such that (6) has a global minimum that perfectly reconstructs the training data, but also a bad local minimum characterized by*

$$q_\phi(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}|\boldsymbol{0}, \boldsymbol{I}) \quad and \quad p_\theta(\boldsymbol{x}) = \mathcal{N}(\boldsymbol{x}|\bar{\boldsymbol{x}}, \bar{\gamma}\boldsymbol{I}). \tag{7}$$

Hence the moment we allow for nonlinear (or more precisely, non-affine) decoders there can exist a poor local minimum, across all parameters including a learnable $\gamma$, that exhibits

6

category (v) posterior collapse.[2] In other words, no predictive information about $\boldsymbol{x}$ passes through the latent space, and a useless/non-informative distribution $p_\theta(\boldsymbol{x})$ emerges that is incapable of assigning high probability to the data (except obviously in the trivial degenerate case where all the data points are equal to the empirical mean $\bar{\boldsymbol{x}}$). We will next investigate the degree to which such concerns can influence behavior in arbitrarily deep architectures.

## 5. Extrapolating to Practical Deep Architectures

Previously we have demonstrated the possibility of local minima aligned with category (v) posterior collapse the moment we allow for decoders that deviate ever so slightly from an affine model. But nuanced counterexamples designed for proving technical results notwithstanding, it is reasonable to examine what realistic factors are largely responsible for leading optimization trajectories towards such potential bad local solutions. For example, is it merely the strength of the KL regularization term, and if so, why can we not just use KL warm-start to navigate around such points? In this section we will elucidate a deceptively simple, alternative risk factor that will be corroborated empirically in Section 6.

From the outset, we should mention that with deep encoder/decoder architectures commonly used in practice, a stationary point can more-or-less always exist at solutions exhibiting posterior collapse. As a representative and ubiquitous example, please see Appendix D. But of course without further details, this type of stationary point could conceivably manifest as a saddle point (stable or unstable), a local maximum, or a local minimum. For the strictly affine decoder model mentioned in Section 4, there will only be a harmless unstable saddle point at any collapsed solution (the Hessian has negative eigenvalues). In contrast, for the special nonlinear case elucidated via Proposition 1 we can instead have a bad local minima. We will now argue that as the depth of common feedforward architectures increases, the risk of converging to category (v)-like solutions with most or all latent dimensions stuck at bad stationary points can also increase.

Somewhat orthogonal to existing explanations of posterior collapse, our basis for this argument is not directly related to the VAE KL-divergence term. Instead, we consider a deceptively simple yet potentially influential alternative: *Unregularized, deterministic AE models can have bad local solutions with high reconstruction errors when sufficiently deep. This in turn can directly translate to category (v) posterior collapse when training a corresponding VAE model with a matching deep architecture.* Moreover, to the extent that this is true, KL warm-start or related countermeasures will likely be ineffective in avoiding such suboptimal minima. We will next examine these claims in greater depth followed by a discussion of practical implications.

### 5.1 From Deeper Architectures to Inevitable Posterior Collapse

Consider the deterministic AE model formed by composing the encoder mean $\boldsymbol{\mu}_x \equiv \boldsymbol{\mu}_x\left(\cdot\,; \theta\right)$ and decoder mean $\boldsymbol{\mu}_z \equiv \boldsymbol{\mu}_z\left(\cdot\,; \phi\right)$ networks from a VAE model, i.e., reconstructions $\hat{\boldsymbol{x}}$ are computed via $\hat{\boldsymbol{x}} = \boldsymbol{\mu}_x\left[\boldsymbol{\mu}_z\left(\boldsymbol{x}; \phi\right); \theta\right]$. We then train this AE to minimize the squared-error

---

2. This result mirrors related efforts examining linear DNNs, where it has been previously demonstrated that under certain conditions, all local minima are globally optimal (Kawaguchi, 2016), while small nonlinearities can induce bad local optima (Yun et al., 2019). However, the loss surface of these models is completely different from a VAE, and hence we view Proposition 1 as a complementary result.

loss $\frac{1}{nd} \sum_{i=1}^{n} \left\| \boldsymbol{x}^{(i)} - \hat{\boldsymbol{x}}^{(i)} \right\|_2^2$, producing parameters $\{\theta_{ae}, \phi_{ae}\}$. Analogously, the corresponding VAE trained to minimize (4) arrives at a parameter set denoted $\{\theta_{vae}, \phi_{vae}\}$. In this scenario, it will typically follow that

$$\frac{1}{nd} \sum_{i=1}^{n} \left\| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x \left[ \boldsymbol{\mu}_z \left( \boldsymbol{x}^{(i)}; \phi_{ae} \right); \theta_{ae} \right] \right\|_2^2 \leq \frac{1}{nd} \sum_{i=1}^{n} \mathbb{E}_{q_{\phi_{vae}}(\boldsymbol{z}|\boldsymbol{x}^{(i)})} \left[ \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x \left( \boldsymbol{z}; \theta_{vae} \right) \|_2^2 \right],$$
(8)

meaning that the deterministic AE reconstruction error will generally be smaller than the stochastic VAE version. Note that if $\boldsymbol{\sigma}_z^2 \to 0$, the VAE defaults to the same deterministic encoder as the AE and hence will have identical representational capacity; however, the KL regularization prevents this from happening, and any $\boldsymbol{\sigma}_z^2 > 0$ can only make the reconstructions worse.[3] Likewise, the KL penalty factor $\|\boldsymbol{\mu}_z^2\|_2^2$ can further restrict the effective capacity and increase the reconstruction error of the training data. Beyond these intuitive arguments, we have never empirically found a case where (8) does not hold (see Section 6 for examples).

We next define the set

$$\mathcal{S}_\varepsilon \triangleq \left\{ \theta, \phi \; : \; \frac{1}{nd} \sum_{i=1}^{n} \left\| \boldsymbol{x}^{(i)} - \hat{\boldsymbol{x}}^{(i)} \right\|_2^2 \leq \varepsilon \right\}$$
(9)

for any $\epsilon > 0$. Now suppose that the chosen encoder/decoder architecture is such that with high probability, achievable optimization trajectories (e.g., via SGD or related) lead to parameters $\{\theta_{ae}, \phi_{ae}\} \notin \mathcal{S}_\varepsilon$, i.e., $\text{Prob}\left(\{\theta_{ae}, \phi_{ae}\} \in \mathcal{S}_\varepsilon\right) \approx 0$. It then follows that the optimal VAE noise variance denoted $\gamma^*$, when conditioned on practically-achievable values for other network parameters, will satisfy

$$\gamma^* \;=\; \frac{1}{nd} \sum_{i=1}^{n} \mathbb{E}_{q_{\phi_{vae}}(\boldsymbol{z}|\boldsymbol{x}^{(i)})} \left[ \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x \left( \boldsymbol{z}; \theta_{vae} \right) \|_2^2 \right] \;\geq\; \varepsilon.$$
(10)

The equality in (10) can be confirmed by simply differentiating the VAE cost w.r.t. $\gamma$ and equating to zero, while the inequality comes from (8) and the fact that $\{\theta_{ae}, \phi_{ae}\} \notin \mathcal{S}_\varepsilon$.

From inspection of the VAE energy from (4), it is readily apparent that larger values of $\gamma$ will discount the data-fitting term and therefore place greater emphasis on the KL divergence. Since the latter is minimized when the latent posterior equals the prior, we might expect that whenever $\varepsilon$ and therefore $\gamma^*$ is increased per (10), we are at a greater risk of nearing collapsed solutions. But the nature of this approach is not at all transparent, and yet this subtlety has important implications for understanding the VAE loss surface in regions at risk of posterior collapse.

For example, one plausible hypothesis is that only as $\gamma^* \to \infty$ do we risk full category (v) collapse. If this were the case, we might have less cause for alarm since the reconstruction error and by association $\gamma^*$ will typically be bounded from above at any local minimizer. However, we will now demonstrate that even finite values can exactly collapse the posterior. In formally showing this, it is helpful to introduce a slightly narrower but nonetheless representative class of VAE models.

---

3. Except potentially in certain contrived adversarial conditions that do not represent practical regimes.

Specifically, let $f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \theta, \boldsymbol{x}^{(i)}\right) \triangleq \mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)})}\left[\|\boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x\left(\boldsymbol{z};\theta\right)\|_2^2\right]$, i.e., the VAE data term evaluated at a single data point without the $1/\gamma$ scale factor. We then define a *well-behaved VAE* as a model with energy function (4) designed such that $\nabla_{\mu_z} f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \theta, \boldsymbol{x}^{(i)}\right)$ and $\nabla_{\sigma_z} f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \theta, \boldsymbol{x}^{(i)}\right)$ are Lipschitz continuous gradients for all $i$. Furthermore, we specify a *non-degenerate decoder* as any $\boldsymbol{\mu}_x(\boldsymbol{z};\theta = \tilde{\theta})$ with $\theta$ set to a $\tilde{\theta}$ value such that $\nabla_{\sigma_z} f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \tilde{\theta}, \boldsymbol{x}^{(i)}\right) \geq c$ for some constant $c > 0$ that can be arbitrarily small. This ensures that $f$ is an increasing function of $\boldsymbol{\sigma}_z$, a quite natural stipulation given that increasing the encoder variance will generally only serve to corrupt the reconstruction, unless of course the decoder is completely blocking the signal from the encoder. In the latter degenerate situation, it would follow that $\nabla_{\mu_z} f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \theta, \boldsymbol{x}^{(i)}\right) = \nabla_{\sigma_z} f\left(\boldsymbol{\mu}_z, \boldsymbol{\sigma}_z, \theta, \boldsymbol{x}^{(i)}\right) = 0$, which is more-or-less tantamount to category (v) posterior collapse.

Based on these definitions, we can now present the following:

**Proposition 2** *For any well-behaved VAE with arbitrary, non-degenerate decoder $\boldsymbol{\mu}_x(\boldsymbol{z};\theta = \tilde{\theta})$, there will always exist a $\gamma' < \infty$ such that the trivial solution $\boldsymbol{\mu}_x(\boldsymbol{z};\theta \neq \tilde{\theta}) = \bar{\boldsymbol{x}}$ and $q_\phi(\boldsymbol{z}|\boldsymbol{x}) = p(\boldsymbol{z})$ will have lower cost.*

Around any evaluation point, the sufficient condition we applied to demonstrate posterior collapse (see proof details) can also be achieved with some $\gamma'' < \gamma'$ if we allow for partial collapse, i.e., $q_{\phi^*}(z_j|\boldsymbol{x}) = p(z_j)$ along some but not all latent dimensions $j \in \{1, \ldots, \kappa\}$. Overall, the analysis loosely suggests that the number of dimensions vulnerable to exact collapse will increase monotonically with $\gamma$.

Proposition 2 also provides evidence that the VAE behaves like a strict thresholding operator, completely shutting off latent dimensions using a finite value for $\gamma$. This is analogous to the distinction between using the $\ell_1$ versus $\ell_2$ norm for solving regularized regression problems of the standard form $\min_{\boldsymbol{u}} \|\boldsymbol{x} - \boldsymbol{A}\boldsymbol{u}\|_2^2 + \gamma \eta(\boldsymbol{u})$, where $\boldsymbol{A}$ is a design matrix and $\eta$ is a penalty function. When $\eta$ is the $\ell_1$ norm, some or all elements of $\boldsymbol{u}$ can be pruned to exactly zero with a sufficiently large but finite $\gamma$ Zhao & Yu (2006). In contrast, when the $\ell_2$ norm is applied, the coefficients will be shrunk to smaller values but never pushed all the way to zero unless $\gamma \to \infty$.

## 5.2 Practical Implications

In aggregate then, if the AE base model displays unavoidably high reconstruction errors, this implicitly constrains the corresponding VAE model to have a large optimal $\gamma$ value, which can potentially lead to undesirable posterior collapse per Proposition 2. In Section 6 we will demonstrate empirically that training unregularized AE models can become increasingly difficult and prone to bad local minima (or at least bad stable stationary points) as the depth increases; and this difficulty can persist even with counter-measures such as skip connections. Therefore, from this vantage point we would argue that it is *the AE base architecture that is effectively the guilty party when it comes to category (v) posterior collapse.*

The perspective described above also helps to explain why heuristics like KL warm-start are not always useful for improving VAE performance. With the standard Gaussian model (4) considered herein, KL warm-start amounts to adopting a pre-defined schedule for incrementally increasing $\gamma$ starting from a small initial value, the motivation being that a

small $\gamma$ will steer optimization trajectories away from overregularized solutions and posterior collapse.

However, regardless of how arbitrarily small $\gamma$ may be fixed at any point during this process, the VAE reconstructions are not likely to be better than the analogous deterministic AE (which is roughly equivalent to forcing $\gamma = 0$ within the present context). This implies that there can exist an *implicit* $\gamma^*$ as computed by (10) that can be significantly larger such that, even if KL warm-start is used, the optimization trajectory may well lead to a collapsed posterior stationary point that has this $\gamma^*$ as the optimal value in terms of minimizing the VAE cost with other parameters fixed. Note that if full posterior collapse does occur, the gradient from the KL term will equal zero and hence, to be at a stationary point it must be that the data term gradient is also zero. In such situations, varying $\gamma$ manually will not impact the gradient balance anyway.

## 6. Empirical Assessments

In this section we empirically demonstrate the existence of bad AE local minima with high reconstruction errors at increasing depth, as well as the association between these bad minima and imminent VAE posterior collapse. For this purpose, we first train fully connected AE and VAE models with 1, 2, 4, 6, 8 and 10 hidden layers on the Fashion-MNIST dataset (Xiao et al., 2017). Each hidden layer is 512-dimensional and followed by ReLU activations (see Appendix A for further details). The reconstruction error is shown in Figure 1(*left*). As the depth of the network increases, the reconstruction error of the AE model first decreases because of the increased capacity. However, when the network becomes too deep, the error starts to increase, indicating convergence to a bad local minima (or at least stable stationary point/plateau) that is unrelated to KL-divergence regularization. The reconstruction error of a VAE model is always worse than that of the corresponding AE model as expected. Moreover, while KL warm-start/annealing can help to improve the VAE reconstructions to some extent, performance is still worse than the AE as expected.

We next train AE and VAE models using a more complex convolutional network on Cifar100 data (Krizhevsky & Hinton, 2009). At each spatial scale, we use 1 to 5 convolution layers followed by ReLU activations. We also apply $2 \times 2$ max pooling to downsample the feature maps to a smaller spatial scale in the encoder and use a transposed convolution layer to upscale the feature map in the decoder. The reconstruction errors are shown in Figure 1(*middle*). Again, the trend is similar to the fully-connected network results. See Appendix A for an additional ImageNet example.

It has been argued in the past that skip connections can increase the mutual information between observations $\boldsymbol{x}^{(i)}$ and the inferred latent variables $\boldsymbol{z}$ (Dieng et al., 2018), reducing the risk of posterior collapse. And it is well-known that ResNet architectures based on skip connections can improve performance on numerous recognition tasks (He et al., 2016). To this end, we train a number of AE models using ResNet-inspired encoder/decoder architectures on multiple datasets including Cifar10, Cifar100, SVHN and CelebA. Similar to the convolution network structure from above, we use 1, 2, and 4 residual blocks within each spatial scale. Inside each block, we apply 2 to 5 convolution layers. For aggregate comparison purposes, we normalize the reconstruction error obtained on each dataset by dividing it with the corresponding error produced by the most shallow network structure (1 residual
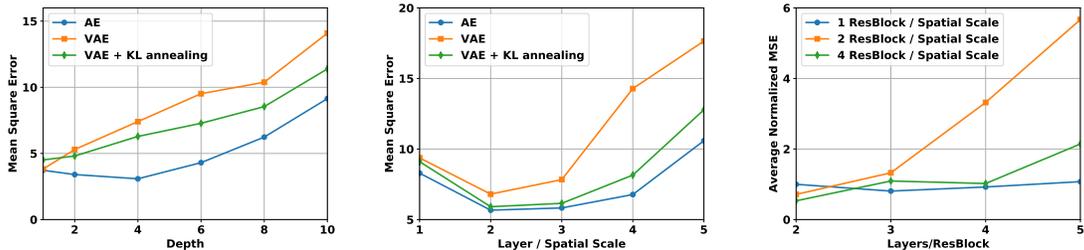
Figure 1: Reconstruction errors for various encoder/decoder models of varying complexity. *Left*: Fully connected networks with different depths trained on Fashion-MNIST. *Middle*: Convolution networks with increasing depth/# of spatial scales trained on Cifar100. *Right*: Averaged AE results from residual networks with varying number of residual blocks and block depth trained on SVHN, Cifar10, Cifar100 and CelebA. In all plots, once the encoder/decoder complexity is sufficiently high, the reconstruction errors begin to increase.

block with 2 convolution layers). We then average the normalized reconstruction errors over all four datasets. The average normalized errors are shown in Figure 1(*right*), where we observe that adding more convolution layers inside each residual block can increase the reconstruction error when the network is too deep. Moreover, adding more residual blocks can also lead to higher reconstruction errors. And empirical results obtained using different datasets and networks architectures, beyond the conditions of Figure 1, also show a general trend of increased reconstruction error once the effective depth is sufficiently deep.

We emphasize that in all these models, as the network complexity/depth increases, the simpler models are always contained within the capacity of the larger ones. Therefore, because the reconstruction error on the training data is becoming worse, it must be the case that the AE is becoming stuck at bad local minima or plateaus. Again since the AE reconstruction error serves as a probable lower bound for that of the VAE model, a deeper VAE model will likely suffer the same problem, only exacerbated by the KL-divergence term in the form of posterior collapse. This implies that there will be more $\boldsymbol{\sigma}_z$ values moving closer to 1 as the VAE model becomes deeper; similarly $\boldsymbol{\mu}_z$ values will push towards 0. The corresponding dimensions will encode no information and become completely useless.

To help corroborate this association between bad AE local minima and VAE posterior collapse, we plot histograms of VAE $\boldsymbol{\sigma}_z$ values as network depth is varied in Figure 2. The models are trained on CelebA and the number of convolution layers in each spatial scale is 2, 4 and 5 from left to right. As the depth increases, the reconstruction error becomes larger and there are more $\boldsymbol{\sigma}_z$ near 1.

## 7. Discussion

In this work we have emphasized the previously-underappreciated role of bad local minima in trapping VAE models at posterior collapsed solutions. Unlike affine decoder models whereby all local minima are provably global, Proposition 1 stipulates that even infinitesimal
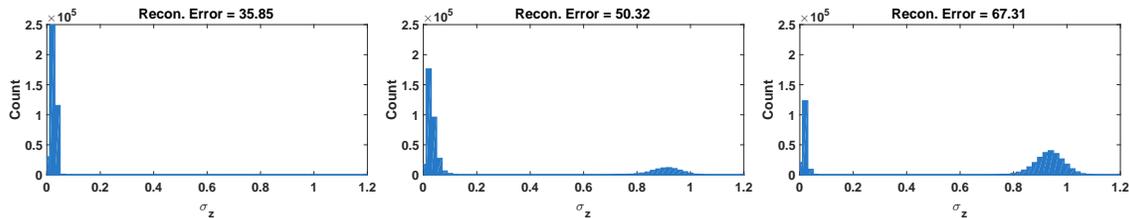
Figure 2: Histogram of $\boldsymbol{\sigma}_z$ values as VAE encoder/decoder network depth is varied. There are 2, 4 and 5 convolution layers in each spatial scale from left to right. As depth increases, the reconstruction error grows and more $\boldsymbol{\sigma}_z$ values are near 1, indicative of impending posterior collapse.

nonlinear perturbations can introduce suboptimal local minima characterized by deleterious posterior collapse. Furthermore, we have demonstrated that the risk of converging to such a suboptimal minima increases with decoder depth. In particular, we outline the following practically-likely pathway to posterior collapse:

1. Deeper AE architectures are essential for modeling high-fidelity images or similar, and yet counter-intuitively, increasing AE depth can actually produce larger reconstruction errors on the training data because of bad local minima (with or without skip connections). An analogous VAE model with the same architecture will likely produce even worse reconstructions because of the additional KL regularization term, which is not designed to steer optimization trajectories away from poor reconstructions.

2. At any such bad local minima, the value of $\gamma$ will necessarily be large, i.e., if it is not large, we cannot be at a local minimum.

3. But because of the thresholding behavior of the VAE as quantified by Proposition 2, as $\gamma$ becomes larger there is an increased risk of exact posterior collapse along excessive latent dimensions. And complete collapse along all dimensions will occur for some finite $\gamma$ sufficiently large. Furthermore, explicitly forcing $\gamma$ to be small does not fix this problem, since in some sense the *implicit* $\gamma^*$ is still large as discussed in Section 5.2.

While we believe that this message is interesting in and of itself, there are nonetheless several practically-relevant implications. For example, complex hierarchical VAEs like BIVA notwithstanding, skip connections and KL warm-start have modest ability to steer optimization trajectories towards good solutions; however, this underappreciated limitation will not generally manifest until networks are sufficiently deep as we have considered. Fortunately, any advances or insights gleaned from developing deeper unregularized AEs, e.g., better AE architectures, training procedures, or initializations (Li & Nguyen, 2019), could likely be adapted to reduce the risk of posterior collapse in corresponding VAE models.

In closing, we should also mention that, although this work has focused on Gaussian VAE models, many of the insights translate into broader non-Gaussian regimes. For example, a variety of recent VAE enhancements involve replacing the fixed Gaussian latent-space prior $p(\boldsymbol{z})$ with a parameterized non-Gaussian alternative (Bauer & Mnih, 2019; Tomczak &

Welling, 2018). This type of modification provides greater flexibility in modeling the aggregated posterior in the latent space, which is useful for generating better samples (Makhzani et al., 2016). However, it does not immunize VAEs against the bad local minima introduced by deep decoders, and good reconstructions are required by models using Gaussian or non-Gaussian priors alike. Therefore, our analysis herein still applies in much the same way.

## Appendix A. Network Structure, Experimental Settings, and Additional ImageNet Results

Three different kinds of network structures are used in the experiments: fully connected networks, convolution networks, and residual networks. For all these structures, we set the dimension of the latent variable $z$ to 64. We then describe the network details accordingly.

**Fully Connected Netowrk:** This experiment is only applied on the simple Fashion-MNIST dataset, which contains 60000 $28 \times 28$ black-and-while images. These images are first flattened to a 784 dimensional vector. Both the encoder and decoder have multiple number of 512-dimensional hidden layers, each followed by ReLU activations.

**Convolution Netowrk:** The original images are either $32 \times 32 \times 3$ (Cifar10, Cifar100 and SVHN) or $64 \times 64 \times 3$ (CelebA and ImageNet). In the encoder, we use a multiple number (denoted as $t$) of $3 \times 3$ convolution layers for each spatial scale. Each convolution layer is followed by a ReLU activation. Then we use a $2 \times 2$ max pooling to downsample the feature map to a smaller spatial scale. The number of channels is doubled when the spatial scale is halved. We use 64 channels when the spatial scale is $32 \times 32$. When the spatial scale reaches $4 \times 4$ (there should be 512 channels in this feature map), we use an average pooling to transform the feature map to a vector, which is then transformed into the latent variable using a fully connected layer. In the decoder, the latent variable is first transformed to a 4096-dimensional vector using a fully connected layer and then reshaped to $2 \times 2 \times 1024$. Again in each spatial scale, we use 1 transpose convolution layer to upscale the feature map and halve the number of channels followed by $t - 1$ convolution layers. Each convolution and transpose convolution layer is followed by a ReLU activation layer. When the spatial scale reaches that of the original image, we use a convolution layer to transofrm the feature map to 3 channels.

**Residual Network:** The network structure of the residual network is similar to that of a convolution network described above. We simply replace the convolution layer with a residual block. Inside the residual block, we use different numbers of convolution numbers. (The typical number of convolution layers inside a residual block is 2 or 3. In our experiments, we try 2, 3, 4 and 5.)

**Training Details:** All the experiments with different network structures and datasets are trained in the same procedure. We use the Adam optimization method and the default optimizer hyper parameters in Tensorflow. The batch size is 64 and we train the model for $250K$ iterations. The initial learning rate is 0.0002 and it is halved every $100K$ iterations.

**Additional Results on ImageNet:** We also show the reconstruction error for convolution networks with increasing depth trained on ImageNet in Figure 3. The trend is the same as that in Figure 1.
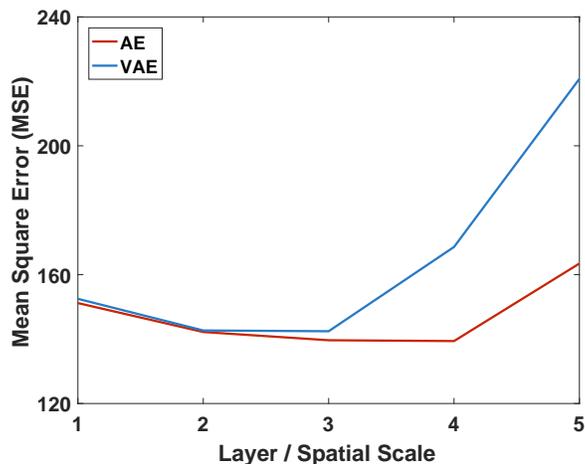
Figure 3: Reconstruction error for Convolution networks with increasing depth/# of spatial scales trained on ImageNet.

## Appendix B. Proof of Proposition 1

While the following analysis could in principle be extended to more complex datasets, for our purposes it is sufficient to consider the following simplified case for ease of exposition. Specifically, we assume that $n > 1, d > \kappa$, set $d = 2, n = 2, \kappa = 1$, and $\boldsymbol{x}^{(1)} = (1,1), \boldsymbol{x}^{(2)} = (-1,-1)$.

Additionally, we will use the following basic facts about the Gaussian tail. Note that (12)-(13) below follow from integration by parts; see Orjebin (2014).

**Lemma 3** *Let* $\epsilon \sim \mathcal{N}(0,1), A > 0$; $\phi(x), \Phi(x)$ *be the pdf and cdf of the standard normal distribution, respectively. Then*

$$1 - \Phi(A) \leq e^{-A^2/2}, \tag{11}$$

$$\mathbb{E}[\epsilon \mathbf{1}_{\{\epsilon > A\}}] = \phi(A), \tag{12}$$

$$\mathbb{E}[\epsilon^2 \mathbf{1}_{\{\epsilon > A\}}] = 1 - \Phi(A) + A\phi(A). \tag{13}$$

### B.1 Suboptimality of (7)

Under the specified conditions, the energy from (7) has a value of $nd$. Thus to show that it is not the global minimum, it suffices to show that the following VAE, parameterized by $\delta$, has energy $\to -\infty$ as $\delta \to 0$:

$$\mu_z^{(1)} = 1, \mu_z^{(2)} = -1,$$
$$\boldsymbol{W}_x = (\alpha + 1, \alpha + 1), \boldsymbol{b}_x = 0,$$
$$\sigma_z^{(1)} = \sigma_z^{(2)} = \delta,$$
$$\gamma = \mathbb{E}_{\mathcal{N}(\varepsilon|0,1)} 2(1 - \pi_\alpha((\alpha + 1)(1 + \delta\varepsilon)))^2.$$

14

This follows because, given the stated parameters, we have that

$$\mathcal{L}(\theta, \phi) = \sum_{i=1}^{2} (1 + 2 \log \mathbb{E}_{\mathcal{N}(\varepsilon|0,1)} 2(1 - \pi_\alpha((\alpha+1)(1+\delta\varepsilon)))^2 - 2\log\delta + \delta^2 + 1)$$

$$= \sum_{i=1}^{2} (\Theta(1) + 2\log \mathbb{E}_{\mathcal{N}(\varepsilon|0,1)}(1 - \pi_\alpha(\alpha + 1 + (\alpha+1)\delta\varepsilon))^2 - 2\log\delta)$$

$$\leq^{(i)} 4\log\delta + \Theta(1).$$

(i) holds when $\delta < \frac{1}{\alpha+1}$; to see this, denote $x := \alpha + 1 + (\alpha+1)(\delta\varepsilon)$. Then

$$\mathbb{E}_{\mathcal{N}(\varepsilon|0,1)}(1 - \pi_\alpha(x))^2$$
$$= \mathbb{E}_\varepsilon[(1 - \pi_\alpha(x))^2 \mathbf{1}_{\{x \geq \alpha\}}] + \mathbb{E}_\varepsilon[(1 - \pi_\alpha(x))^2 \mathbf{1}_{\{|x| < \alpha\}}] + \mathbb{E}_\varepsilon[(1 - \pi_\alpha(x))^2 \mathbf{1}_{\{x < -\alpha\}}]$$
$$\leq \underbrace{\mathbb{E}_\varepsilon[(1 - (x - \alpha))^2]}_{(a)} + \underbrace{\mathbb{P}(|x| < \alpha)}_{(b)} + \underbrace{\mathbb{E}_\varepsilon((1 - x - \alpha)^2 \mathbf{1}_{\{x < -\alpha\}})}_{(c)}.$$

In the RHS above $(a) = [(\alpha+1)\delta]^2$; using (11)-(13) we then have

$$(b) < \mathbb{P}(x < \alpha) = \mathbb{P}\left(\varepsilon < \frac{-1}{(\alpha+1)\delta}\right) \leq \exp\left(-\frac{1}{2[(\alpha+1)\delta]^2}\right).$$
$$(c) < \mathbb{E}_\varepsilon((2\alpha + (\alpha+1)\delta\varepsilon)^2 \mathbf{1}_{\{x < \alpha\}})$$
$$= \int_{-\infty}^{\frac{-1}{(\alpha+1)\delta}} (2\alpha + (\alpha+1)\delta\varepsilon)^2 \frac{1}{\sqrt{2\pi}} e^{-\varepsilon^2/2} d\varepsilon$$
$$< \int_{-\infty}^{\frac{-1}{(\alpha+1)\delta}} (4\alpha^2 + [(\alpha+1)\delta\varepsilon]^2) \frac{1}{\sqrt{2\pi}} e^{-\varepsilon^2/2} d\varepsilon$$
$$< \left\{4\alpha^2 + ((\alpha+1)\delta)^2 \left[1 + \frac{1}{\sqrt{2\pi}}\right]\right\} \exp\left(-\frac{1}{2[(\alpha+1)\delta]^2}\right)$$

when $\delta < \frac{1}{\alpha+1}$. Thus

$$\lim_{\delta \to 0} \frac{\mathbb{E}_{\mathcal{N}(\varepsilon|0,1)}(1 - \pi_\alpha(x))^2}{[(\alpha+1)\delta]^2} = 1,$$

and

$$\lim_{\delta \to 0} \{\log \mathbb{E}_{\mathcal{N}(\varepsilon|0,1)}(1 - \pi_\alpha(x))^2 - 2\log\delta\} = 2\log(\alpha+1),$$

or

$$2\log \mathbb{E}_\epsilon(1 - \pi_\alpha(x))^2 = 4\log\delta + \Theta(1),$$

and we can see (i) holds.

## B.2 Local Optimality of (7)

We will now show that at (7), the Hessian of the energy has structure

$$
\begin{array}{ccccc}
 & (\boldsymbol{W}_x) & (\boldsymbol{b}_x) & (\sigma_z^{(i)}, \mu_z^{(i)}) & (\gamma) \\
(\boldsymbol{W}_x) & 0 & 0 & 0 & 0 \\
(\boldsymbol{b}_x) & 0 & \frac{2}{\gamma}I & 0 & 0 \\
(\sigma_z^{(i)}, \mu_z^{(i)}) & 0 & 0 & (p.d.) & 0 \\
(\gamma) & 0 & 0 & 0 & (p.d.)
\end{array}
$$

where p.d. means the corresponding submatrix is positive definite and independent of other parameters. While the Hessian is 0 in the subspace of $\boldsymbol{W}_x$, we can show that for VAEs that are only different from (7) by $\boldsymbol{W}_x$, the gradient always points back to (7). Thus (7) is a strict local minima.

First we compute the Hessian matrix block-wise. We will identify $\boldsymbol{W}_x \in \mathbb{R}^{2\times 1}$ with the vector $(W_j)_{j=1}^2$, and use the shorthand notations $\boldsymbol{x}^{(i)} = (x_j^{(i)})_{j=1}^2$, $\boldsymbol{b}_x = (b_j)_{j=1}^2$, $z^{(i)} = \mu_z^{(i)} + \sigma_z^{(i)}\varepsilon$, where $\varepsilon \sim \mathcal{N}(0,1)$ (recall that $z^{(i)}$ is a scalar in this proof).

1. The second-order derivatives involving $\boldsymbol{W}_x$ can be expressed as

$$
\frac{\partial \mathcal{L}}{\partial W_j} = \frac{-2}{\gamma} \sum_{i=1}^n \mathbb{E}_\varepsilon [(\pi_\alpha'(W_j z^{(i)}) z^{(i)}) \cdot (x_j^{(i)} - \pi_\alpha(W_j z^{(i)}) - b_j)], \tag{14}
$$

   and therefore all second-order derivatives involving $W_j$ will have the form

$$
\mathbb{E}_\epsilon [\pi_\alpha'(W_j z^{(i)}) F_1 + \pi_\alpha''(W_j z^{(i)}) F_2], \tag{15}
$$

   where $F_1, F_2$ are some arbitrary functions that are finite at (7). Since $\pi_\alpha'(0) = \pi_\alpha''(0) = W_j = 0$, the above always evaluates to 0 at $\boldsymbol{W}_x = 0$.

2. For second-order derivatives involving $\boldsymbol{b}_x$, we have

$$
\frac{\partial \mathcal{L}}{\partial \boldsymbol{b}_x} = \frac{-2}{\gamma} \mathbb{E}_\varepsilon [\boldsymbol{x}^{(i)} - \pi_\alpha(\boldsymbol{W}_x z^{(i)}) - \boldsymbol{b}_x]
$$

   and

$$
\frac{\partial^2 \mathcal{L}}{\partial (\boldsymbol{b}_x)^2} = \frac{2}{\gamma} I,
$$

$$
\frac{\partial^2 \mathcal{L}}{\partial \gamma \partial \boldsymbol{b}_x} = \frac{2}{\gamma^2} \frac{\partial \mathcal{L}}{\partial \boldsymbol{b}_x} = 0, \qquad\qquad (\text{since } \boldsymbol{W}_x = 0);
$$

   and $\frac{\partial^2 \mathcal{L}}{\partial \mu_z^{(i)} \partial \boldsymbol{b}_x}$ and $\frac{\partial^2 \mathcal{L}}{\partial \mu_z^{(i)} \partial \sigma_z^{(i)}}$ will also have the form of (15), thus both equal 0 at $\boldsymbol{W}_x = 0$.

3. Next consider second-order derivatives involving $\mu_z^{(i)}$ or $\sigma_k^{(i)}$. Since the KL part of the energy, $\sum_{i=1}^n \text{KL}(q_\phi(\boldsymbol{z}|\boldsymbol{x}^{(i)})|p(\boldsymbol{z}))$, only depends on $\mu_z^{(i)}$ and $\sigma_k^{(i)}$, and have p.d. Hessian

16

at (7) independent of other parameters, it suffices to calculate the derivatives of the reconstruction error part, denoted as $\mathcal{L}_{\mathrm{recon}}$. Since

$$\frac{\partial \mathcal{L}_{\mathrm{recon}}}{\partial \mu_z^{(i)}} = \frac{-2}{\gamma} \sum_{i,j} \mathbb{E}_\epsilon \left[ (x_j^{(i)} - \pi_\alpha(W_j z^{(i)}) - b_j) W_j \pi_\alpha'(W_j z^{(i)}) \right],$$

$$\frac{\partial \mathcal{L}_{\mathrm{recon}}}{\partial \sigma_z^{(i)}} = \frac{-2}{\gamma} \sum_{i,j} \mathbb{E}_\epsilon \left[ (x_j^{(i)} - \pi_\alpha(W_j z^{(i)}) - b_j) W_j \epsilon \pi_\alpha'(W_j z^{(i)}) \right],$$

all second-order derivatives will have the form of (15), and equal 0 at $\boldsymbol{W}_x = 0$.

4. For $\gamma$, we can calculate that $\partial^2 \mathcal{L} / \partial \gamma^2 = 4/\gamma^2 > 0$ at (7).

Now, consider VAE parameters that are only different from (7) in $\boldsymbol{W}_x$. Plugging $\boldsymbol{b}_x = \bar{\boldsymbol{x}}, \mu_z^{(i)} = 0, \sigma_k^{(i)} = 1$ into (14), we have

$$\frac{\partial \mathcal{L}}{\partial W_j} = \frac{-2}{\gamma} \sum_{i=1}^n \mathbb{E}_\varepsilon [(\pi_\alpha'(W_j \varepsilon)\varepsilon) \cdot (-\pi_\alpha(W_j \varepsilon))].$$

As $(\pi_\alpha'(W_j \varepsilon)\varepsilon) \cdot (-\pi_\alpha(W_j \varepsilon)) \leq 0$ always holds, we can see that the gradient points back to (7). This concludes our proof of (7) being a strict local minima. ∎

## Appendix C. Proof of Proposition 2

We begin by assuming an arbitrarily complex encoder for convenience. This allows us to remove the encoder-sponsored amortized inference and instead optimize independent parameters $\boldsymbol{\mu}_z^{(i)}$ and $\boldsymbol{\sigma}_z^{(i)}$ separately for each data point. Later we will show that this capacity assumption can be dropped and the main result still holds.

We next define

$$\boldsymbol{m}_z \triangleq \left[ \left( \boldsymbol{\mu}_z^{(1)} \right)^\top, \ldots, \left( \boldsymbol{\mu}_z^{(n)} \right)^\top \right]^\top \in \mathbb{R}^{\kappa n} \text{ and } \boldsymbol{s}_z \triangleq \left[ \left( \boldsymbol{\sigma}_z^{(1)} \right)^\top, \ldots, \left( \boldsymbol{\sigma}_z^{(n)} \right)^\top \right]^\top \in \mathbb{R}^{\kappa n}, \quad (16)$$

which are nothing more than the concatenation of all of the decoder means and variances from each data point into the respective column vectors. It is also useful to decompose the assumed non-degenerate decoder parameters via

$$\theta \equiv [\psi, w], \quad \psi \triangleq \theta \backslash w, \quad (17)$$

where $w \in [0, 1]$ is a scalar such that $\mu_x(\boldsymbol{z}; \theta) \equiv \mu_x(w\boldsymbol{z}; \psi)$. Note that we can always reparameterize an existing deep architecture to extract such a latent scaling factor which we can then hypothetically optimize separately while holding the remaining parameters $\psi$ fixed. Finally, with slight abuse of notation, we may then define the function

$$f(w\boldsymbol{m}_z, w\boldsymbol{s}_z) \triangleq \quad (18)$$
$$\sum_{i=1}^n f\left( \boldsymbol{\mu}_z^{(i)}, \boldsymbol{\sigma}_z^{(i)}, [\tilde{\psi}, w], \boldsymbol{x}^{(i)} \right) \equiv \sum_{i=1}^n \mathbb{E}_{\mathcal{N}\left( \boldsymbol{z} | \boldsymbol{\mu}_z^{(i)}, \mathrm{diag}\left[ \boldsymbol{\sigma}_z^{(i)} \right]^2 \right)} \left[ \| \boldsymbol{x}^{(i)} - \boldsymbol{\mu}_x\left( w\boldsymbol{z}; \tilde{\psi} \right) \|_2^2 \right].$$

This is basically just the original function $f$ summed over all training points, with $\psi$ fixed at the corresponding values extracted from $\tilde{\theta}$ while $w$ serves as a free scaling parameter on the decoder.

Based on the assumption of Lipschitz continuous gradients, we can always create the upper bound

$$
\begin{aligned}
f\left(\boldsymbol{u}, \boldsymbol{v}\right) \quad \leq \quad & f\left(\tilde{\boldsymbol{u}}, \tilde{\boldsymbol{v}}\right) \\
& + \left(\boldsymbol{u} - \tilde{\boldsymbol{u}}\right)^{\top} \nabla_u f\left(\boldsymbol{u}, \boldsymbol{v}\right)|_{\boldsymbol{u}=\tilde{\boldsymbol{u}}} \ + \ \tfrac{L}{2} \left\| \boldsymbol{u} - \tilde{\boldsymbol{u}} \right\|_2^2 \ + \ \left(\boldsymbol{v} - \tilde{\boldsymbol{v}}\right)^{\top} \nabla_v f\left(\boldsymbol{u}, \boldsymbol{v}\right)|_{\boldsymbol{v}=\tilde{\boldsymbol{v}}} \ + \ \tfrac{L}{2} \left\| \boldsymbol{v} - \tilde{\boldsymbol{v}} \right\|_2^2,
\end{aligned}
\tag{19}
$$

where $L$ is the Lipschitz constant of the gradients and we have adopted $\boldsymbol{u} \triangleq w\boldsymbol{m}_z$ and $\boldsymbol{v} \triangleq w\boldsymbol{\sigma}_z$ to simplify notation. Equality occurs at the evaluation point $\{\boldsymbol{u}, \boldsymbol{v}\} = \{\tilde{\boldsymbol{u}}, \tilde{\boldsymbol{v}}\}$. However, this bound does not account for the fact that we know $\nabla_v f\left(\boldsymbol{u}, \boldsymbol{v}\right) \geq 0$ (i.e., $f\left(\boldsymbol{u}, \boldsymbol{v}\right)$ is increasing w.r.t. $\boldsymbol{v}$) and that $\boldsymbol{v} \geq 0$. Given these assumptions, we can produce the refined upper bound

$$
f^{ub}\left(\boldsymbol{u}, \boldsymbol{v}\right) \quad \geq \quad f\left(\boldsymbol{u}, \boldsymbol{v}\right),
\tag{20}
$$

where $\quad f^{ub}\left(\boldsymbol{u}, \boldsymbol{v}\right) \quad \triangleq$

$$
f\left(\tilde{\boldsymbol{u}}, \tilde{\boldsymbol{v}}\right) \ + \ \left(\boldsymbol{u} - \tilde{\boldsymbol{u}}\right)^{\top} \nabla_u f\left(\boldsymbol{u}, \boldsymbol{v}\right)|_{\boldsymbol{u}=\tilde{\boldsymbol{u}}} \ + \ \tfrac{L}{2} \left\| \boldsymbol{u} - \tilde{\boldsymbol{u}} \right\|_2^2 + \sum_{j=1}^{nd} g\left( v_j, \tilde{v}_j, \nabla_{v_j} f\left(\boldsymbol{u}, \boldsymbol{v}\right)\big|_{v_j=\tilde{v}_j} \right)
\tag{21}
$$

and the function $g : \mathbb{R}^3 \to \mathbb{R}$ is defined as

$$
g\left(v, \tilde{v}, \delta\right) \triangleq \begin{cases} \left(v - \tilde{v}\right)\delta + \tfrac{L}{2}\left(v - \tilde{v}\right)_2^2 & \text{if } v \geq \tilde{v} - \tfrac{\delta}{L} \ \text{ and } \ \{v, \tilde{v}, \delta\} \geq 0, \\ \frac{-\delta^2}{2L} & \text{if } v < \tilde{v} - \tfrac{\delta}{L} \ \text{ and } \ \{v, \tilde{v}, \delta\} \geq 0, \\ \infty & \text{otherwise.} \end{cases}
\tag{22}
$$

Given that

$$
\tilde{v} - \tfrac{\delta}{L} = \arg\min_v \left[ \left(v - \tilde{v}\right)\delta + \tfrac{L}{2}\left(v - \tilde{v}\right)_2^2 \right] \quad \text{and} \quad \frac{-\delta^2}{2L} = \min_v \left[ \left(v - \tilde{v}\right)\delta + \tfrac{L}{2}\left(v - \tilde{v}\right)_2^2 \right],
\tag{23}
$$

the function $g$ is basically just setting all values of $\left(v - \tilde{v}\right)\delta + \tfrac{L}{2}\left\| v - \tilde{v} \right\|_2^2$ with negative slope to the minimum $\frac{-\delta^2}{2L}$. This change is possible while retaining an upper bound because $f\left(\boldsymbol{u}, \boldsymbol{v}\right)$ is non-decreasing in $\boldsymbol{v}$ by stated assumption. Additionally, $g$ is set to infinity for all $v < 0$ to enforce non-negatively.

While it may be possible to proceed further using $f^{ub}$, we find it useful to consider a final modification. Specifically, we define the approximation

$$
f^{appr}\left(\boldsymbol{u}, \boldsymbol{v}\right) \quad \approx \quad f^{ub}\left(\tilde{\boldsymbol{u}}, \tilde{\boldsymbol{v}}\right),
\tag{24}
$$

where $\quad f^{appr}\left(\boldsymbol{u}, \boldsymbol{v}\right) \quad \triangleq$

$$
f\left(\tilde{\boldsymbol{u}}, \tilde{\boldsymbol{v}}\right) \ + \ \left(\boldsymbol{u} - \tilde{\boldsymbol{u}}\right)^{\top} \nabla_u f\left(\boldsymbol{u}, \boldsymbol{v}\right)|_{\boldsymbol{u}=\tilde{\boldsymbol{u}}} \ + \ \tfrac{L}{2} \left\| \boldsymbol{u} - \tilde{\boldsymbol{u}} \right\|_2^2 + \sum_{j=1}^{nd} g^{appr}\left( v_j, \tilde{v}_j, \nabla_{v_j} f\left(\boldsymbol{u}, \boldsymbol{v}\right)\big|_{v_j=\tilde{v}_j} \right)
\tag{25}
$$

and

$$g^{appr}\left(v,\tilde{v},\delta\right) \triangleq \begin{cases} \frac{-\delta^2}{2L} + \frac{\delta^2}{2L\tilde{v}^2}v^2 & \text{if } \tilde{v} - \frac{\delta}{L} \geq 0 \text{ and } \{v,\tilde{v},\delta\} \geq 0, \\ \left(\frac{L\tilde{v}^2}{2} - \delta\tilde{v}\right) + \left(\frac{\delta}{\tilde{v}} - \frac{L}{2}\right)v^2 & \text{if } \tilde{v} - \frac{\delta}{L} < 0 \text{ and } \{v,\tilde{v},\delta\} \geq 0, \\ \infty & \text{otherwise.} \end{cases} \qquad (26)$$

While slightly cumbersome to write out, $g^{appr}$ has a simple interpretation. By construction, we have that

$$\min_v g^{appr}\left(v,\tilde{v},\delta\right) = g^{appr}\left(0,\tilde{v},\delta\right) = \min_v g\left(v,\tilde{v},\delta\right) = g\left(0,\tilde{v},\delta\right) \qquad (27)$$

and

$$g^{appr}\left(\tilde{v},\tilde{v},\delta\right) = g\left(\tilde{v},\tilde{v},\delta\right) = 0. \qquad (28)$$

At other points, $g^{appr}$ is just a simple quadratic interpolation but without any factor that is linear in $v$. And removal of this linear term, while retaining (27) and (27) will be useful for the analysis that follows below. Note also that although $f^{appr}\left(\boldsymbol{u},\boldsymbol{v}\right)$ is no longer a strict bound on $f\left(\boldsymbol{u},\boldsymbol{v}\right)$, it will nonetheless still be an upper bound whenever $v_j \in \{0,\tilde{v}_j\}$ for all $j$ which will ultimately be sufficient for our purposes.

We now consider optimizing the function

$$h^{appr}(\boldsymbol{m}_z,\boldsymbol{s}_z,w) \triangleq \tfrac{1}{\gamma}f^{appr}\left(w\boldsymbol{m}_z,w\boldsymbol{s}_z\right) + \sum_{i=1}^{n}\left\|\boldsymbol{\mu}_z^{(i)}\right\|_2^2 + \left\|\boldsymbol{\sigma}_z^{(i)}\right\|_2^2 - \log\left|\text{diag}\left[\boldsymbol{\sigma}_z^{(i)}\right]^2\right|. \qquad (29)$$

If we define $\mathcal{L}\left(\boldsymbol{m}_z,\boldsymbol{s}_z,w\right)$ as the VAE cost from (4) under the current parameterization, then by design it follows that

$$h^{appr}(\tilde{\boldsymbol{m}}_z,\tilde{\boldsymbol{s}}_z,\tilde{w}) = \mathcal{L}\left(\tilde{\boldsymbol{m}}_z,\tilde{\boldsymbol{s}}_z,\tilde{w}\right) \qquad (30)$$

and

$$h^{appr}(\boldsymbol{m}_z,\boldsymbol{s}_z,w) \geq \mathcal{L}\left(\boldsymbol{m}_z,\boldsymbol{s}_z,w\right) \qquad (31)$$

whenever $w\sigma_j \in \{0,\tilde{w}\tilde{\sigma}_j\}$ for all $j$. Therefore if we find such a solution $\{\boldsymbol{m}_z',\boldsymbol{s}_z',w'\}$ that satisfies this condition and has $h^{appr}(\boldsymbol{m}_z',\boldsymbol{s}_z',w') < h^{appr}(\tilde{\boldsymbol{m}}_z,\tilde{\boldsymbol{s}}_z,\tilde{w})$, it necessitates that $\mathcal{L}(\boldsymbol{m}_z',\boldsymbol{s}_z',w') < \mathcal{L}(\tilde{\boldsymbol{m}}_z,\tilde{\boldsymbol{s}}_z,\tilde{w})$ as well. This then ensures that $\{\tilde{\boldsymbol{m}}_z,\tilde{\boldsymbol{s}}_z,\tilde{w}\}$ cannot be a local minimum.

We now examine the function $h^{appr}$ more closely. After a few algebraic manipulations and excluding irrelevant constants, we have that

$$h^{appr}(\boldsymbol{m}_z,\boldsymbol{s}_z,w) \equiv$$
$$\sum_{j=1}^{nd}\left\{\tfrac{1}{\gamma}\left[w m_{z,j}\left.\nabla_{u_j}f\left(\boldsymbol{u},\boldsymbol{v}\right)\right|_{u_j=\tilde{w}\tilde{m}_{z,j}} + \tfrac{L}{2}\left(w^2 m_{z,j}^2 - 2w m_{z,j}\tilde{w}\tilde{m}_{z,j}\right) + c_j w^2 s_{z,j}^2\right]\right.$$
$$\left. + m_{z,j}^2 + s_{z,j}^2 - \log s_{z,j}^2\right\}, \qquad (32)$$

where $c_j$ is the coefficient on the $v^2$ term from (26). After rearranging terms, optimizing out $\boldsymbol{m}_z$ and $\boldsymbol{s}_z$, and discarding constants, we can then obtain (with slight abuse of notation) the reduced function

$$h^{appr}(w) \triangleq \sum_{j=1}^{nd}\frac{y_j}{\gamma + \beta w^2} + \log(\gamma + c_j w^2), \qquad (33)$$

19

where $\beta \triangleq \frac{L}{2}$ and $y_j \triangleq \frac{L}{2} \left\| \tilde{w} \tilde{m}_{z,j} - \frac{1}{L} \left. \nabla_{u_j} f \left( \boldsymbol{u}, \boldsymbol{v} \right) \right|_{u_j = \tilde{w} \tilde{m}_{z,j}} \right\|_2^2$. Note that $y_j$ must be bounded since $L \neq 0^4$ and $w \in [0, 1]$, $\left. \nabla_{u_j} f \left( \boldsymbol{u}, \boldsymbol{v} \right) \right|_{u_j = \tilde{w} \tilde{m}_{z,j}} \leq L$, and $\tilde{\boldsymbol{m}}$ are all bounded. The latter is implicitly bounded because the VAE KL term prevents infinite encoder mean functions. Furthermore, $c_j$ must be strictly greater than zero per the definition of a non-degenerate decoder; this guarantees that

$$ g^{appr} \left( \tilde{w} \tilde{s}_j, \tilde{w} \tilde{s}_j, \left. \nabla_{v_j} f \left( \boldsymbol{u}, \boldsymbol{v} \right) \right|_{v_j = \tilde{w} \tilde{s}_j} \right) > g^{appr} \left( 0, \tilde{w} \tilde{s}_j, \left. \nabla_{v_j} f \left( \boldsymbol{u}, \boldsymbol{v} \right) \right|_{v_j = \tilde{w} \tilde{s}_j} \right), \qquad (34) $$

which is only possible with $c_j > 0$. Proceeding further, because

$$ \nabla_{w^2} h^{appr}(w) = \sum_{j=1}^{nd} \left( \frac{-\beta y_j}{(\gamma + \beta w^2)^2} + \frac{c_j}{\gamma + c_j w^2} \right), \qquad (35) $$

we observe that if $\gamma$ is increased sufficiently large, the first term will always be smaller than the second since $\beta$ and all $y_j$ are bounded, and $c_j > 0 \; \forall j$. So there can never be a point whereby $\nabla_{w^2} h^{appr}(w) = 0$ when $\gamma = \gamma'$ sufficiently large. Therefore the minimum in this situation occurs on the boundary where $w^2 = 0$. And finally, if $w^2 = 0$, then the optimal $\boldsymbol{m}_z$ and $\boldsymbol{s}_z$ is determined solely by the KL term, and hence they are set according to the prior. Moreover, the decoder has no signal from the encoder and is therefore optimized by simply setting $\boldsymbol{\mu}_x \left( 0; \tilde{\psi} \right)$ to the mean $\bar{\boldsymbol{x}}$ for all $i$.[5] Additionally, none of this analysis requires and arbitrarily complex encoder; the exact same results hold as long as the encoder can output a 0 for means and 1 for the variances.

Note also that if we proceed through the above analysis using $\boldsymbol{w} \in \mathbb{R}^\kappa$ as parameterizing a separate $w_j$ scaling factor for each latent dimension $j \in \{1, \ldots, \kappa\}$, then a smaller $\gamma$ value would generally force partial collapse. In other words, we could enforce nonzero gradients of $h^{appr}(w)$ along the indices of each latent dimension separately. This loosely criteria would then lead to $q_{\phi^*}(z_j | \boldsymbol{x}) = p(z_j)$ along some but not all latent dimensions as stated in the main text below Proposition 2. ∎

## Appendix D. Representative Stationary Point Exhibiting Posterior Collapse in Deep VAE Models

Here we provide an example of a stationary point that exhibits posterior collapse with an arbitrary deep encoder/decoder architecture. This example is representative of many other possible cases. Assume both encoder and decoder mean functions $\boldsymbol{\mu}_x$ and $\boldsymbol{\mu}_z$, as well as the diagonal encoder covariance function $\boldsymbol{\Sigma}_z = \text{diag}[\boldsymbol{\sigma}_z^2]$, are computed by standard deep neural networks, with layers composed of linear weights followed by element-wise nonlinear activations (the decoder covariance satisfies $\boldsymbol{\Sigma}_x = \gamma \boldsymbol{I}$ as before). We denote the weight matrix from the first layer of the decoder mean network as $\boldsymbol{W}_{\mu_x}^1$, while $\boldsymbol{w}_{\mu_x, \cdot j}^1$ refers to the

---

4. $L = 0$ would violate the stipulated conditions for a non-degenerate decoder since it would imply that no signal from $\boldsymbol{z}$ could pass through the decoder. And of course if $L = 0$, we would already be at a solution exhibiting posterior collapse.

5. We are assuming here that the decoder has sufficient capacity to model any constant value, e.g., the output layer has a bias term.

corresponding $j$-th column. Assuming $\rho$ layers, we denote $\boldsymbol{W}^{\rho}_{\mu_z}$ and $\boldsymbol{W}^{\rho}_{\sigma^2_z}$ as weights from the last layers of the encoder networks producing $\boldsymbol{\mu}_z$ and $\log \boldsymbol{\sigma}^2_z$ respectively, with $j$-th rows defined as $\boldsymbol{w}^{\rho}_{\mu_z,j\cdot}$ and $\boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot}$. We then characterize the following key stationary point:

**Proposition 4** If $\boldsymbol{w}^1_{\mu_x,\cdot j} = \left(\boldsymbol{w}^{\rho}_{\mu_z,j\cdot}\right)^{\top} = \left(\boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot}\right)^{\top} = \boldsymbol{0}$ for any $j \in \{1,2,\ldots,\kappa\}$, then the gradients of (4) with respect to $\boldsymbol{w}^1_{\mu_x,\cdot j}$, $\boldsymbol{w}^{\rho}_{\mu_z,j\cdot}$, and $\boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot}$ are all equal to zero.

If the stated weights are zero along dimension $j$, then obviously it must be that $q_\phi(z_j|\boldsymbol{x}) = p(z_j)$, i.e., a collapsed dimension for better or worse. The proof is straightforward; we provide the details below for completeness.

**Proof:** First we remind that the variational upper bound is defined in (2). We define $\mathcal{L}(\boldsymbol{x};\theta,\phi)$ as the loss at a data point $\boldsymbol{x}$, *i.e.*

$$\mathcal{L}(\boldsymbol{x};\theta,\phi) = -\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})\right] + \mathbb{KL}\left[q_\phi(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})\right]. \tag{36}$$

The total loss is the integration of $\mathcal{L}(\boldsymbol{x};\theta,\phi)$ over $\boldsymbol{x}$. Further more, we denote $\mathcal{L}_{kl}(\boldsymbol{x};\theta)$ and $\mathcal{L}_{gen}(\boldsymbol{x};\theta,\phi)$ as the KL loss and the generation loss at $\boldsymbol{x}$ respectively, *i.e.*

$$\mathcal{L}_{kl}(\boldsymbol{x};\phi) = \mathbb{KL}\left[q_\phi(\boldsymbol{z}|\boldsymbol{x})||p(\boldsymbol{z})\right] = \sum_{i=1}^{\kappa}\mathbb{KL}\left[q_\phi(z_j|\boldsymbol{x})||p(z_j)\right],$$

$$= \frac{1}{2}\sum_{j=1}^{\kappa}\left(\mu^2_{z,j} + \sigma^2_{z,j} - \log\sigma^2_{z,j} - 1\right) \tag{37}$$

$$\mathcal{L}_{gen}(\boldsymbol{x};\phi,\theta) = -\mathbb{E}_{q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\log p_\theta(\boldsymbol{x}|\boldsymbol{z})\right]. \tag{38}$$

The second equality in (37) holds because the covariance of $q_\phi(\boldsymbol{z}|\boldsymbol{x})$ and $p(\boldsymbol{z})$ are both diagonal. The last encoder layer and the first decoder layer are denoted as $\boldsymbol{h}^{\rho}_e$ and $\boldsymbol{h}^1_d$. If $\boldsymbol{w}^{\rho}_{\mu_z,j\cdot} = 0, \boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot} = 0$, then we have

$$\mu_{z,j} = \boldsymbol{w}^{\rho}_{\mu_z,j\cdot}\boldsymbol{h}^{\rho}_e = 0, \quad \sigma^2_{z,j} = \exp\left(\boldsymbol{w}_{\sigma^2_z,j\cdot}\right) = 1, \quad q(z_j|\boldsymbol{x}) = \mathcal{N}(0,1). \tag{39}$$

The gradient of $\mu_{z,j}$ and $\sigma_{z,j}$ from $\mathcal{L}_{kl}(\boldsymbol{x};\phi)$ becomes

$$\frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{\partial\mu_{z,j}} = \mu_{z,j} = 0, \quad \frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{\partial\sigma_{z,j}} = 1 - \sigma^{-1}_{z,j} = 0. \tag{40}$$

So the gradient of $\boldsymbol{w}^{\rho}_{\mu_z,j\cdot}$ and $\boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot}$ from $\mathcal{L}_{kl}$ is

$$\frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{\partial\boldsymbol{w}^{\rho}_{\mu_z,j\cdot}} = \frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{\partial\mu_{z,j}}\boldsymbol{h}^{\rho\top}_e = 0, \tag{41}$$

$$\frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{\partial\boldsymbol{w}^{\rho}_{\sigma^2_z,j\cdot}} = \frac{\partial\mathcal{L}_{kl}(\boldsymbol{x};\phi)}{2\sigma_{z,j}\cdot\partial\sigma_{z,j}}\boldsymbol{h}^{\rho\top}_e = 0. \tag{42}$$

Now we consider the gradient from $\mathcal{L}_{gen}(\boldsymbol{x};\theta,\phi)$. We have

$$\frac{-\partial\log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial z_j} = \frac{-\partial\log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial\boldsymbol{h}^1_d}\frac{\partial\boldsymbol{h}^1_d}{\partial z_j}. \tag{43}$$

Since

$$\boldsymbol{h}_d^1 = \text{act}\left(\sum_{j=1}^{\kappa} \boldsymbol{w}_{\mu_x,\cdot j}^1 z_j\right), \tag{44}$$

where $\text{act}(\cdot)$ is the activation function, we can obtain

$$\frac{\partial \boldsymbol{h}_d^1}{\partial z_j} = \text{act}'\left(\sum_{j=1}^{\kappa} \boldsymbol{w}_{\mu_x,\cdot j}^1 z_j\right) \boldsymbol{w}_{\mu_x,\cdot j}^1 = 0. \tag{45}$$

Plugging this back into (43) gives

$$\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial z_j} = 0. \tag{46}$$

According to the chain rule, we have

$$\frac{\partial \mathcal{L}_{gen}(\boldsymbol{x}; \theta, \phi)}{\partial \boldsymbol{w}_{\mu_z,j\cdot}^\rho} = \mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial z_j} \frac{\partial z_j}{\partial \boldsymbol{w}_{\mu_z,j\cdot}^\rho}\right] = 0, \tag{47}$$

$$\frac{\partial \mathcal{L}_{gen}(\boldsymbol{x}; \theta, \phi)}{\partial \boldsymbol{w}_{\sigma_z^2,j\cdot}^\rho} = \mathbb{E}_{\boldsymbol{z} \sim q_\phi(\boldsymbol{z}|\boldsymbol{x})}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial z_j} \frac{\partial z_j}{\partial \boldsymbol{w}_{\sigma_z^2,j\cdot}^\rho}\right] = 0. \tag{48}$$

After combining these two equations with (41) and (42) and then integrating over $\boldsymbol{x}$, we have

$$\frac{\partial \mathcal{L}(\theta, \phi)}{\partial \boldsymbol{w}_{\mu_z,j\cdot}^\rho} = 0, \tag{49}$$

$$\frac{\partial \mathcal{L}(\theta, \phi)}{\partial \boldsymbol{w}_{\sigma_z^2,j\cdot}^\rho} = 0. \tag{50}$$

Then we consider the gradient with respect to $\boldsymbol{w}_{\mu_x,\cdot j}^1$. Since $\boldsymbol{w}_{\mu_x,\cdot j}$ is part of $\theta$, it only receives gradient from $\mathcal{L}_{gen}(\boldsymbol{x}; \theta, \phi)$. So we do not need to consider the KL loss. If $\boldsymbol{w}_{\mu_x,\cdot j}^1 = 0$, $\boldsymbol{h}_d^1 = \sum_{j=1}^{\kappa} \boldsymbol{w}_{\mu_x,\cdot j}^1 z_j$ is not related to $\boldsymbol{z}_j$. So $p_\theta(\boldsymbol{x}|\boldsymbol{z}) = p_\theta(\boldsymbol{x}|\boldsymbol{z}_{\neg j})$, where $\boldsymbol{z}_{\neg j}$ represents $\boldsymbol{z}$ without the $j$-th dimension. The gradient of $\boldsymbol{w}_{\mu_x,\cdot j}^1$ is

$$\frac{\partial \mathcal{L}_{gen}(\boldsymbol{x}; \theta, \phi)}{\partial \boldsymbol{w}_{\mu_x,\cdot j}^1} = \mathbb{E}_{\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x})}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial \boldsymbol{w}_{\mu_x,\cdot j}^1}\right] = \mathbb{E}_{\boldsymbol{z} \sim q(\boldsymbol{z}|\boldsymbol{x})}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z})}{\partial \boldsymbol{h}_d^1} z_j\right] \tag{51}$$

$$= \mathbb{E}_{\boldsymbol{z}_{\neg j} \sim q(\boldsymbol{z}_{\neg j}|\boldsymbol{x})}\left[\mathbb{E}_{z_j \sim \mathcal{N}(0,1)}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z}_{\neg j})}{\partial \boldsymbol{h}_d^1} z_j\right]\right]$$

$$= \mathbb{E}_{\boldsymbol{z}_{\neg i} \sim q(\boldsymbol{z}_{\neg i}|\boldsymbol{x})}\left[\frac{-\partial \log p_\theta(\boldsymbol{x}|\boldsymbol{z}_{\neg j})}{\partial \boldsymbol{h}_d^1} \mathbb{E}_{z_j \sim \mathcal{N}(0,1)}[z_j]\right] = 0.$$

The integration over $\boldsymbol{x}$ should also be 0. So we obtain

$$\frac{\partial \mathcal{L}(\theta; \phi)}{\partial \boldsymbol{w}_{\mu_x,\cdot j}^1} = 0. \tag{52}$$

$$\blacksquare$$

# References

A. Alemi, B. Poole, I. Fischer, J. Dillon, R. Saurous, and K. Murphy. Fixing a broken ELBO. *arXiv preprint arXiv:1711.00464*, 2017.

M. Bauer and A. Mnih. Resampled priors for variational autoencoders. *arXiv preprint arXiv:1810.11428*, 2018.

M. Bauer and A. Mnih. Resampled priors for variational autoencoders. *International Conference on Artificial Intelligence and Statistics*, 2019.

S. Bowman, L. Vilnis, O. Vinyals, A. Dai, R. Jozefowicz, and S. Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.

Y. Burda, R. Grosse, and R. Salakhutdinov. Importance weighted autoencoders. *arXiv preprint arXiv:1509.00519*, 2015.

L. Cai, H. Gao, and S. Ji. Multi-stage variational auto-encoders for coarse-to-fine image generation. *arXiv preprint arXiv:1705.07202*, 2017.

E. Candès, X. Li, Y. Ma, and J. Wright. Robust principal component analysis? *J. ACM*, 58(2), 2011.

X. Chen, D. Kingma, T. Salimans, Y. Duan, P. Dhariwal, J. Schulman, I. Sutskever, and P. Abbeel. Variational lossy autoencoder. *arXiv preprint arXiv:1611.02731*, 2016.

B. Dai and D. Wipf. Diagnosing and enhancing VAE models. *International Conference on Learning Representations*, 2019.

B. Dai, Y. Wang, J. Aston, G. Hua, and D. Wipf. Hidden talents of the variational autoencoder. *arXiv preprint arXiv:1706.05148*, 2019.

A. Dieng, Y. Kim, A. Rush, and D. Blei. Avoiding latent variable collapse with generative skip models. *arXiv preprint arXiv:1807.04863*, 2018.

K. Gregor and Y. LeCun. Learning fast approximations of sparse coding. *International Conference on Machine Learning*, 2010.

J. He, D. Spokoyny, G. Neubig, and T. Berg-Kirkpatrick. Lagging inference networks and posterior collapse in variational autoencoders. *International Conference on Learning Representations*, 2019.

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

I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, , and A. Lerchner. $\beta$-vae: Learning basic visual concepts with a constrained variational framework. *International Conference on Learning Representations*, 2017.

C. Huang, S. Tan, A. Lacoste, and A. Courville. Improving explorability in variational inference with annealed variational objectives. *Advances in Neural Information Processing Systems*, 2018.

K. Kawaguchi. Deep learning without poor local minima. *Advances in Neural Information Processing Systems*, 2016.

D. Kingma and M. Welling. Auto-encoding variational Bayes. *International Conference on Learning Representations*, 2014.

A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

P. Li and P.M. Nguyen. On random deep weight-tied autoencoders: Exact asymptotic analysis, phase transitions, and implications to training. *International Conference on Learning Representations*, 2019.

J. Lucas, G. Tucker, R. Grosse, and M. Norouzi. Understanding posterior collapse in generative latent variable models. *International Conference on Learning Representations, Workshop Paper*, 2019.

L. Maaløe, M. Fraccaro, V. Liévin, and O. Winther. BIVA: A very deep hierarchy of latent variables for generative modeling. *arXiv preprint arXiv:1902.02102*, 2019.

A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2016.

P.A. Mattei and J. Frellsen. Leveraging the exact likelihood of deep latent variable models. *Advances in Neural Information Processing Systems*, 2018.

E. Orjebin. A recursive formula for the moments of a truncated univariate normal distribution. 2014. URL https://people.smp.uq.edu.au/YoniNazarathy/teaching_projects/studentWork/EricOrjebin_TruncatedNormalMoments.pdf.

A. Razavi, A. Oord, B. Poole, and O. Vinyals. Preventing posterior collapse with $\delta$-VAEs. *International Conference on Learning Representations*, 2019.

D. Rezende and S. Mohamed. Variational inference with normalizing flows. *arXiv preprint arXiv:1505.05770*, 2015.

D. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *International Conference on Machine Learning*, 2014.

M. Rolinek, D. Zietlow, and G. Martius. Variational autoencoders pursue PCA directions (by accident). 2019.

C. Sønderby, T. Raiko, L. Maaløe, S. Sønderby, and O. Winther. How to train deep variational autoencoders and probabilistic ladder networks. *arXiv preprint arXiv:1602.02282*, 2016.

P. Sprechmann, A.M. Bronstein, and G. Sapiro. Learning efficient sparse and low rank models. *IEEE Trans. Pattern Analysis and Machine Intelligence*, 37(9), 2015.

M. Tipping and C. Bishop. Probabilistic principal component analysis. *J. Royal Statistical Society, Series B*, 61(3):611–622, 1999.

J. Tomczak and M. Welling. VAE with a VampPrior. *International Conference on Artificial Intelligence and Statistics*, 2018.

A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, and K. Kavukcuoglu. Conditional image generation with PixelCNN decoders. *Advances in Neural Information Processing Systems*, 2016.

Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

S. Yeung, A. Kannan, Y. Dauphin, and L. Fei-Fei. Tackling over-pruning in variational autoencoders. *arXiv preprint arXiv:1706.03643*, 2017.

C. Yun, S. Sra, and A. Jadbabaie. Small nonlinearities in activation functions create bad local minima in neural networks. *International Conference on Learning Representations*, 2019.

P. Zhao and B. Yu. On model selection consistency of Lasso. *Journal of Machine learning research*, 7:2541–2563, 2006.