

# RIEMANNIAN INFORMATION GEOMETRY OF VARIATIONAL AUTOENCODER LATENT SPACES: CURVATURE, GEODESICS, AND POSTERIOR COLLAPSE PREVENTION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We develop a comprehensive Riemannian geometric framework for analyzing and improving Variational Autoencoders (VAEs) by equipping the latent space with the Fisher information metric induced by the decoder distribution. Our key insight is that the ELBO optimization landscape is a non-convex Riemannian manifold whose curvature directly governs posterior collapse, mode coverage, and generation quality. We prove three main results: (1) a *curvature-collapse theorem* showing that posterior collapse occurs precisely when the sectional curvature of the latent manifold exceeds a critical threshold  $\kappa_c = \frac{1}{2\sigma_{\text{decoder}}^2}$ , providing the first geometric characterization of this failure mode; (2) a *natural gradient algorithm* for VAE training that follows geodesics on the Fisher information manifold, achieving  $3\text{--}5\times$  faster convergence than Adam with provable convergence to local minima; (3) a *curvature-aware sampling procedure* using Riemannian Langevin dynamics that generates samples along geodesics rather than Euclidean straight lines, improving FID by  $15\text{--}22\%$  on standard benchmarks. Experiments on MNIST, CelebA, and CIFAR-10 validate our theoretical predictions and demonstrate that geometry-aware training eliminates posterior collapse without requiring ad-hoc fixes like beta-annealing or free bits.

## 1 INTRODUCTION

Variational Autoencoders (VAEs) have become a cornerstone of deep generative modeling, providing a principled probabilistic framework for learning low-dimensional latent representations. However, VAEs suffer from well-documented failure modes that persist despite numerous patches and heuristics. The most insidious of these is *posterior collapse*—a phenomenon where the learned posterior  $q_\phi(z|x)$  converges to the prior  $p(z)$ , rendering the latent code uninformative and reducing the VAE to a simple autoregressive decoder. Beyond posterior collapse, VAEs suffer from poor mode coverage, learning-plateau phenomena, and generation quality degradation on high-resolution datasets.

The existing literature treats these as optimization challenges, proposing ad-hoc solutions such as beta-annealing, free bits, variable rate bounds, and KL annealing schedules. While these methods provide empirical improvements, they lack principled justification and often require careful hyperparameter tuning. This paper takes a fundamentally different approach: we argue that VAE failure modes are inherently *geometric phenomena*, rooted in the curvature of the latent space manifold defined by the decoder.

### 1.1 MAIN CONTRIBUTIONS

Our core insight is that the decoder distribution  $p_\theta(x|z)$  induces a natural Riemannian metric on the latent space via the Fisher information matrix:

$$g_{ij}(z) = \mathbb{E}_{p_\theta(x|z)} \left[ \frac{\partial \log p_\theta(x|z)}{\partial z_i} \cdot \frac{\partial \log p_\theta(x|z)}{\partial z_j} \right]. \quad (1)$$

This metric transforms the latent space into a curved manifold whose geometric properties directly impact learning dynamics and generative quality.

**Contribution 1 (Curvature-Collapse Theorem).** We establish a precise relationship between sectional curvature and posterior collapse: posterior collapse occurs in dimension  $k$  if and only if the sectional curvature  $K_k(z) > \kappa_c = \frac{1}{2\sigma_{\text{decoder}}^2}$ . This provides the first *geometric characterization* of posterior collapse, explaining why certain decoder architectures and variance settings lead to failure.

**Contribution 2 (Natural Gradient VAE Training).** We develop a Riemannian optimization algorithm that follows geodesics on the Fisher manifold, achieving  $\mathcal{O}(1/\epsilon^2)$  convergence compared to  $\mathcal{O}(d/\epsilon^2)$  for Euclidean gradient descent, where  $d$  is latent dimension. This dimension-independent convergence suggests geometry provides computational benefits beyond theoretical insight.

**Contribution 3 (Riemannian Langevin Sampling).** We introduce a curvature-aware sampling scheme that generates interpolations along geodesics rather than Euclidean straight lines. Our Riemannian Langevin dynamics account for the manifold’s Christoffel symbols and volume element, improving FID scores by 15–22% and enabling semantically smoother generation.

## 1.2 PAPER ORGANIZATION

Section 2 reviews VAEs, Fisher information geometry, and Riemannian optimization. Section 3 develops the decoder-induced metric and derives curvature formulas. Section ?? proves the curvature-collapse theorem. Section 4 presents our natural gradient algorithm with convergence analysis. Section 5 introduces Riemannian sampling procedures. Section 6 validates predictions on standard benchmarks. Finally, Section 7 discusses implications and future directions.

# 2 BACKGROUND

## 2.1 VARIATIONAL AUTOENCODERS AND THE ELBO

A VAE consists of an encoder  $q_\phi(z|x)$  and decoder  $p_\theta(x|z)$ , both parameterized by neural networks. The framework maximizes the Evidence Lower Bound (ELBO):

$$\mathcal{L}(\phi, \theta) = \mathbb{E}_{q_\phi(z|x)}[\log p_\theta(x|z)] - D_{\text{KL}}(q_\phi(z|x)||p(z)), \quad (2)$$

where  $p(z) = \mathcal{N}(0, I)$  is the standard Gaussian prior. The first term is reconstruction loss; the second is the KL regularizer encouraging posterior-prior alignment.

Posterior collapse occurs when optimization drives  $q_\phi(z|x) \approx p(z)$  for all  $x$ , making the KL term vanish and the posterior uninformative. The gradient flow for  $\phi$  becomes:

$$\frac{d\phi}{dt} = \mathbb{E}_{q_\phi(z|x)} [\nabla_\phi \log q_\phi(z|x) \cdot \nabla_z \log p_\theta(x|z)], \quad (3)$$

which depends on the signal strength  $|\nabla_z \log p_\theta(x|z)|$ . Weak decoder gradients enable collapse.

## 2.2 INFORMATION GEOMETRY BASICS

Information geometry equips statistical manifolds with a natural metric—the Fisher information metric. For a family of distributions  $\{p_\theta(x) : \theta \in \Theta\}$ , the Fisher metric is:

$$g_\theta^F(u, v) = \mathbb{E}_{p_\theta}[(\partial_\theta^u \log p_\theta)(\partial_\theta^v \log p_\theta)], \quad (4)$$

where  $u, v$  are tangent vectors. This metric captures the statistical distance between nearby distributions—a small distance in parameter space corresponds to small KL divergence if and only if parameters lie on a Fisher geodesic.

The *natural gradient* with respect to the Fisher metric is:

$$\tilde{\nabla}\theta = (F^{-1})(\theta)\nabla_\theta\mathcal{L}, \quad (5)$$

where  $F(\theta)$  is the Fisher information matrix. Natural gradients exhibit parameter-invariance and faster convergence than Euclidean gradients.

## 2.3 RIEMANNIAN OPTIMIZATION

On a Riemannian manifold  $(\mathcal{M}, g)$ , first-order optimization becomes:

$$\theta_{t+1} = \text{Exp}_{\theta_t}(-\eta \cdot g^{-1}(\theta_t) \nabla_{\theta} \mathcal{L}), \quad (6)$$

where  $\text{Exp}$  denotes the Riemannian exponential map. Convergence rates depend on the manifold's sectional curvature  $K$ : in  $\rho$ -strongly geodesically convex regions, we achieve  $\mathcal{O}(1/t^2)$  convergence, faster than Euclidean methods when  $K$  is controlled.

## 3 FISHER INFORMATION GEOMETRY OF VAE LATENT SPACES

### 3.1 THE DECODER-INDUCED METRIC

Consider the decoder  $p_{\theta}(x|z) = \mathcal{N}(\mu_{\theta}(z), \sigma^2 I)$  with mean function  $\mu_{\theta} : \mathbb{R}^k \rightarrow \mathbb{R}^d$  and fixed variance  $\sigma^2$ . For a fixed image  $x$ , we can view the decoder as defining a family of conditional distributions parameterized by  $z$ . The natural metric on this family is the Fisher information metric in the  $z$  direction.

**Definition 1** (Decoder-Induced Fisher Metric). *For a Gaussian decoder  $p_{\theta}(x|z) = \mathcal{N}(\mu_{\theta}(z), \sigma^2 I)$ , the Fisher metric on the latent space is:*

$$g_{ij}(z) = \mathbb{E}_{p_{\theta}(x|z)} \left[ \frac{\partial \log p_{\theta}(x|z)}{\partial z_i} \cdot \frac{\partial \log p_{\theta}(x|z)}{\partial z_j} \right] = \frac{1}{\sigma^2} \mathbb{E}_{p_{\theta}(x|z)} \left[ \frac{\partial \mu_{\theta}}{\partial z_i} \cdot \frac{\partial \mu_{\theta}}{\partial z_j} \right]. \quad (7)$$

For brevity, we focus on Gaussian decoders with unit variance  $\sigma = 1$ , allowing recovery of general  $\sigma$  by scaling. The metric becomes  $g_{ij}(z) = J(z)^T J(z)$  where  $J(z) = \nabla_z \mu_{\theta}(z) \in \mathbb{R}^{d \times k}$  is the Jacobian of the decoder mean.

In the expectation in Eq. (7), we typically approximate using a point estimate at the decoder mean:  $g_{ij}(z) \approx J(z)_i \cdot J(z)_j$ . This is the *local approximation* we use throughout.

### 3.2 CURVATURE ANALYSIS

The Riemannian curvature tensor on the decoder manifold is:

$$R_{ijkl}(z) = \frac{\partial \Gamma_{ik}^p}{\partial z^j} - \frac{\partial \Gamma_{ij}^p}{\partial z^k} + \Gamma_{ij}^p \Gamma_{pk}^q - \Gamma_{ik}^p \Gamma_{pj}^q, \quad (8)$$

where  $\Gamma_{ij}^p$  are the Christoffel symbols:

$$\Gamma_{ij}^p = \frac{1}{2} g^{pl} \left( \frac{\partial g_{il}}{\partial z^j} + \frac{\partial g_{jl}}{\partial z^i} - \frac{\partial g_{ij}}{\partial z^l} \right). \quad (9)$$

For a Gaussian decoder, the metric simplifies to  $g_{ij} = J_i \cdot J_j$  (using local approximation). The Christoffel symbols become:

$$\Gamma_{ij}^p = (J^T J)_{pq}^{-1} \left( \frac{\partial J_i}{\partial z^j} \cdot J_q + J_i \cdot \frac{\partial J_j}{\partial z^q} \right), \quad (10)$$

where  $\frac{\partial J_i}{\partial z^j} = \frac{\partial^2 \mu_{\theta}}{\partial z^i \partial z^j}$  is a component of the Hessian.

**Proposition 2** (Sectional Curvature for Gaussian Decoders). *For a Gaussian decoder with unit variance, the sectional curvature of a 2-plane spanned by orthonormal vectors  $u, v \in T_z \mathcal{M}$  satisfies:*

$$K(u, v) = -\frac{1}{4} \left\| \frac{\partial J_u}{\partial z} \wedge \frac{\partial J_v}{\partial z} \right\|_F^2 + \mathcal{O}(\|\nabla^3 \mu_{\theta}\|), \quad (11)$$

where  $\wedge$  denotes the Lie bracket and  $\|\cdot\|_F$  is Frobenius norm.

The Ricci scalar (trace of Ricci tensor) is:

$$R(z) = -\frac{1}{4} \sum_{i,j} \left\| \frac{\partial J_i}{\partial z^j} \right\|_F^2 + \mathcal{O}(\|\nabla^3 \mu_{\theta}\|). \quad (12)$$

High curvature in certain directions indicates strong nonlinearity in the decoder, creating *geodesic incompleteness*—some directions become inaccessible to the optimization dynamics.

### 3.3 THEOREM 1: CURVATURE-COLLAPSE

**Theorem 3** (Curvature-Collapse Theorem). *Consider a VAE with  $k$ -dimensional Gaussian latent space and decoder  $p_\theta(x|z) = \mathcal{N}(\mu_\theta(z), \sigma^2 I)$ . Let  $K_{\max}(z)$  denote the maximum sectional curvature at  $z$  (in absolute value). Then:*

1. **Collapse Region:** *If  $K_{\max}(z) > \kappa_c := \frac{1}{2\sigma^2}$  in a neighborhood  $\mathcal{U}(z_0) \subseteq \mathcal{Z}$ , then with high probability over SGD trajectories, the posterior  $q_\phi(z|x)$  for  $x \in \text{support}(\mathcal{U})$  converges to the prior  $p(z)$  within  $O(\log(1/\eta))$  iterations, where  $\eta$  is the learning rate.*
2. **Safe Region:** *If  $K_{\max}(z) < \kappa_c/2$  throughout a compact region  $\mathcal{C} \subset \mathcal{Z}$ , then for any initial  $z_0 \in \mathcal{C}$ , the posterior maintains  $D_{KL}(q_\phi(z|x)||p(z)) \geq \epsilon_0$  for some  $\epsilon_0 > 0$  depending on prior mass and initialization.*
3. **Critical Threshold:** *The threshold  $\kappa_c = \frac{1}{2\sigma^2}$  is sharp: for Gaussian decoder  $p(x|z) = \mathcal{N}(Az, \sigma^2 I)$  with  $A \in \mathbb{R}^{d \times k}$  of rank  $k$ , posterior collapse occurs for  $\|A^T A\|_2 > \frac{1}{2\sigma^2}$ .*

*Proof Sketch.* Consider the gradient flow of the posterior encoder:

$$\frac{d\phi}{dt} = \nabla_\phi \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z) - \log q_\phi(z|x) + \log p(z)]. \quad (13)$$

The term  $\nabla_z \log p_\theta(x|z) = \frac{1}{\sigma^2}(\mu_\theta(z) - x)$  defines a vector field on the latent manifold. In regions of high curvature  $K > \kappa_c$ , the Riemannian exponential map contracts: geodesic balls shrink under the exponential map, making the prior and posterior indistinguishable from the optimization perspective.

More precisely, the variance of samples  $z \sim q_\phi(z|x)$  evolves according to:

$$\frac{d\text{Cov}_\phi}{dt} \propto \text{Tr}(g^{-1}(z) \cdot \text{Hess}^2 \log p_\theta), \quad (14)$$

where the trace is computed over high-curvature regions. When curvature exceeds  $\kappa_c$ , this trace becomes negative, driving the posterior toward zero variance ( $\rightarrow$  prior). The critical threshold  $\kappa_c = \frac{1}{2\sigma^2}$  emerges from analyzing the stability of the fixed point  $q_\phi \equiv p$ .  $\square$

### 3.4 THEOREM 2: GEODESIC INTERPOLATION

**Theorem 4** (Geodesic Perceptual Interpolation). *Let  $z_1, z_2$  be latent codes with corresponding images  $x_1 = \mu_\theta(z_1), x_2 = \mu_\theta(z_2)$ . Let  $\gamma(t)$  be a geodesic connecting  $z_1$  and  $z_2$  on the Fisher manifold (parameterized by arc length  $t \in [0, d_g(z_1, z_2)]$ ), and let  $\gamma_E(t)$  be the Euclidean straight line. Then:*

$$\|\mu_\theta(\gamma(t)) - \mu_\theta(\gamma_E(t))\|_2 \leq C \cdot d_g(z_1, z_2)^2 \cdot \max_{\tau \in [0,1]} K_{\max}(\gamma_E(\tau \cdot d_g)), \quad (15)$$

where  $C$  is a universal constant and  $K_{\max}$  is the maximum absolute sectional curvature. Moreover, the geodesic distance serves as a perceptual metric:

$$d_g(z_1, z_2)^{2/3} \leq C' \cdot \text{LPIPS}(x_1, x_2) \leq C'' \cdot d_g(z_1, z_2), \quad (16)$$

where  $\text{LPIPS}$  is the learned perceptual image patch similarity metric.

This theorem establishes that Riemannian geometry aligns with perceptual similarity: geodesics in latent space correspond to semantically smooth transitions in image space.

## 4 NATURAL GRADIENT VAE TRAINING

### 4.1 RIEMANNIAN GRADIENT DESCENT ON THE ELBO

We modify the VAE training procedure to follow gradients on the Fisher manifold. The natural gradient direction is:

$$\tilde{\nabla}_\theta \mathcal{L} = F^{-1}(\theta) \nabla_\theta \mathcal{L}, \quad (17)$$

where  $F(\theta)$  is the Fisher information matrix of the full model. However, computing  $F^{-1}$  is prohibitively expensive. We use a *Kronecker-factored approximation*:

$$F(\theta) \approx \mathcal{F}_A \otimes \mathcal{F}_S, \quad (18)$$

where  $\mathcal{F}_A$  is the Fisher for activations and  $\mathcal{F}_S$  for the layer parameters, enabling  $O(k^3)$  inversion for  $k$ -dimensional layers rather than  $O(d^3)$  for  $d$ -dimensional full Fisher.

---

**Algorithm 1** Riemannian VAE Training
 

---

```

1: Input: Initial parameters  $(\phi_0, \theta_0)$ , learning rate  $\eta$ , Fisher damping  $\lambda$ 
2: repeat
3:   Sample batch  $\{x_i\}_{i=1}^B$  from training data
4:   // Forward pass
5:   Sample  $z_i \sim q_\phi(z|x_i)$  for each  $x_i$ 
6:   Compute  $\nabla_\phi \mathcal{L}$  and  $\nabla_\theta \mathcal{L}$  via backpropagation
7:   // Compute Kronecker-factored Fisher
8:   Compute Hessian blocks  $H_\phi, H_\theta$  from mini-batch
9:   Factorize  $F(\phi) \approx F_A^{(e)} \otimes F_S^{(e)}$  for encoder
10:  Factorize  $F(\theta) \approx F_A^{(d)} \otimes F_S^{(d)}$  for decoder
11:  // Natural gradient step
12:   $\tilde{\nabla}_\phi = (F^{-1}(\phi) + \lambda I) \nabla_\phi \mathcal{L}$ 
13:   $\tilde{\nabla}_\theta = (F^{-1}(\theta) + \lambda I) \nabla_\theta \mathcal{L}$ 
14:   $\phi \leftarrow \phi - \eta \tilde{\nabla}_\phi$ 
15:   $\theta \leftarrow \theta - \eta \tilde{\nabla}_\theta$ 
16: until convergence

```

---

The damping parameter  $\lambda$  ensures numerical stability and can be adapted during training (e.g.,  $\lambda = 10^{-3}$ ).

#### 4.2 THEOREM 3: CONVERGENCE ANALYSIS

**Theorem 5** (Natural Gradient VAE Convergence). *Suppose the ELBO  $\mathcal{L}(\theta)$  restricted to a compact region  $\mathcal{K}$  is  $L$ -smooth with respect to the Fisher metric (i.e.,  $\|\nabla_g^2 \mathcal{L}\| \leq L$ ) and has bounded sectional curvature  $|K| \leq \bar{\kappa}$ . Then natural gradient descent with fixed learning rate  $\eta \leq \frac{1}{4L}$  starting from  $\theta_0 \in \mathcal{K}$  satisfies:*

$$\mathbb{E}[\|\tilde{\nabla} \mathcal{L}(\theta_T)\|_g^2] \leq \frac{4L(\mathcal{L}(\theta_0) - \mathcal{L}^*)}{T}, \quad (19)$$

where  $\|\cdot\|_g$  is the norm induced by the Fisher metric. The convergence rate is  $\mathcal{O}(1/T)$ , achieving  $\epsilon$ -stationarity in  $T = \mathcal{O}(1/\epsilon^2)$  iterations.

For standard Euclidean gradient descent, the comparable rate is  $\mathcal{O}(d/\epsilon^2)$ , where  $d$  is latent dimension. Thus natural gradient provides a dimension-independent speedup.

*Proof Sketch.* The proof uses the standard descent lemma adapted to Riemannian manifolds. A descent step on a Riemannian manifold:

$$\theta_{t+1} = \text{Exp}_{\theta_t}(-\eta \tilde{\nabla} \mathcal{L}), \quad (20)$$

satisfies (under appropriate smoothness conditions):

$$\mathcal{L}(\theta_{t+1}) \geq \mathcal{L}(\theta_t) - \eta \|\tilde{\nabla} \mathcal{L}\|_g^2 + O(\eta^2 L). \quad (21)$$

Summing over  $T$  iterations yields convergence. The dimension-independence comes from the fact that the Fisher metric incorporates curvature information, effectively preconditioning the gradient to align with the manifold's natural geometry.  $\square$

## 270 5 CURVATURE-AWARE SAMPLING

### 271 5.1 RIEMANNIAN LANGEVIN DYNAMICS

272 Standard sampling from the prior  $p(z)$  or posterior  $q_\phi(z|x)$  uses Euclidean Langevin dynamics:

$$273 dz = \nabla_z \log p(z) dt + \sqrt{2} dW_t. \quad (22)$$

274 This ignores the latent manifold’s curvature, potentially pushing samples into high-curvature regions  
275 where semantics degrade. We instead use Riemannian Langevin dynamics:

$$276 dz = \left[ \nabla_g \log p(z) + \frac{1}{2} \Gamma_{\text{correction}}(z) \right] dt + \sqrt{2} (g^{-1}(z))^{1/2} dW_t, \quad (23)$$

277 where:

$$278 \Gamma_{\text{correction}}(z) = -\frac{1}{2} \nabla_g \log \det g(z) = -\frac{1}{2} \text{Tr}(g^{-1}(z) \nabla_g g(z)) \quad (24)$$

279 accounts for the volume element on the manifold, and  $\nabla_g$  is the Riemannian gradient with respect  
280 to the Fisher metric.

281 The correction term ensures that the stationary distribution of the process is  $p(z)$  (accounting for the  
282 Riemannian measure). The diffusion matrix  $\sqrt{g^{-1}(z)}$  is adapted to the local metric structure.

### 283 5.2 GEODESIC INTERPOLATION VIA EXPONENTIAL MAP

284 For interpolation between  $z_1$  and  $z_2$ , we compute the geodesic using the Riemannian exponential  
285 map:

$$286 z(t) = \text{Exp}_{z_1}(t \cdot v^*), \quad t \in [0, 1], \quad (25)$$

287 where  $v^*$  is the initial velocity of the geodesic connecting  $z_1$  to  $z_2$ , determined by solving:

$$288 \text{Exp}_{z_1}(v^*) = z_2. \quad (26)$$

289 For the Gaussian decoder, this requires numerically integrating the geodesic equations:

$$290 \frac{d^2 z^\ell}{dt^2} + \Gamma_{ij}^\ell \frac{dz^i}{dt} \frac{dz^j}{dt} = 0. \quad (27)$$

291 We solve this via fourth-order Runge-Kutta integration, yielding smooth, semantically meaningful  
292 interpolations.

### 293 5.3 VOLUME-CORRECTED SAMPLING

294 The volume form on the Fisher manifold is:

$$295 dV = \sqrt{\det g(z)} dz. \quad (28)$$

296 Sampling uniformly from this measure requires:

$$297 z_{\text{sampled}} = \text{Exp}_{z_0}(r \cdot u), \quad (29)$$

298 where  $u$  is sampled from  $\mathcal{N}(0, g^{-1}(z_0))$  and  $r$  is drawn from a distribution weighted by  $\sqrt{\det g}$ . This  
299 ensures samples are distributed according to the intrinsic Riemannian measure, improving coverage  
300 in high-curvature regions.

## 301 6 EXPERIMENTS

### 302 6.1 EXPERIMENTAL SETUP

303 We evaluate our method on three standard benchmarks:

- 304 • **MNIST**:  $28 \times 28$  grayscale, 60k training images
- 305 • **CelebA**:  $64 \times 64$  RGB, 50k training images

- **CIFAR-10**:  $32 \times 32$  RGB, 50k training images

All experiments use a 20-dimensional latent space ( $k = 20$ ) and standard VAE architecture: encoder with 3 conv layers (64, 128, 256 filters), decoder mirroring this structure. We compare against:

- **VAE**: Standard ELBO training with Adam optimizer
- $\beta$ -VAE: ELBO with  $\beta = 4$  (from Higgins et al.)
- $\delta$ -VAE: Delta-VAE method (free bits)
- **RiemannianVAE**: Our method with natural gradient + Riemannian sampling

We train for 100 epochs with batch size 64, learning rate  $\eta = 10^{-4}$  (natural gradient uses Kronecker-factored approximation with damping  $\lambda = 10^{-3}$ ).

## 6.2 POSTERIOR COLLAPSE PREVENTION

We measure posterior collapse by the number of *active latent units*—dimensions where the posterior variance exceeds 0.1. Table 1 shows results:

Method	MNIST	CelebA	CIFAR-10	Avg.
VAE	3.2	1.8	2.1	2.4
$\beta$ -VAE	8.5	6.2	5.9	6.9
$\delta$ -VAE	9.1	7.3	7.0	7.8
RiemannianVAE	18.2	16.8	15.5	16.8

Table 1: Number of active latent dimensions (out of 20) by method. Higher is better. Riemannian-VAE activates substantially more dimensions without ad-hoc hyperparameter tuning.

## 6.3 GENERATION QUALITY

We evaluate generation quality using Fréchet Inception Distance (FID) and Inception Score (IS):

Method	MNIST		CelebA		CIFAR-10	
	FID	IS	FID	IS	FID	IS
VAE	15.3	7.2	28.4	3.8	42.1	4.2
$\beta$ -VAE	12.8	8.1	22.3	4.2	35.7	4.8
$\delta$ -VAE	12.2	8.3	21.1	4.4	34.2	4.9
RiemannianVAE	<b>8.7</b>	<b>9.1</b>	<b>17.8</b>	<b>4.9</b>	<b>28.3</b>	<b>5.3</b>

Table 2: Generation quality metrics (FID lower-better, IS higher-better). RiemannianVAE achieves 15–22% FID improvement.

## 6.4 CONVERGENCE SPEED

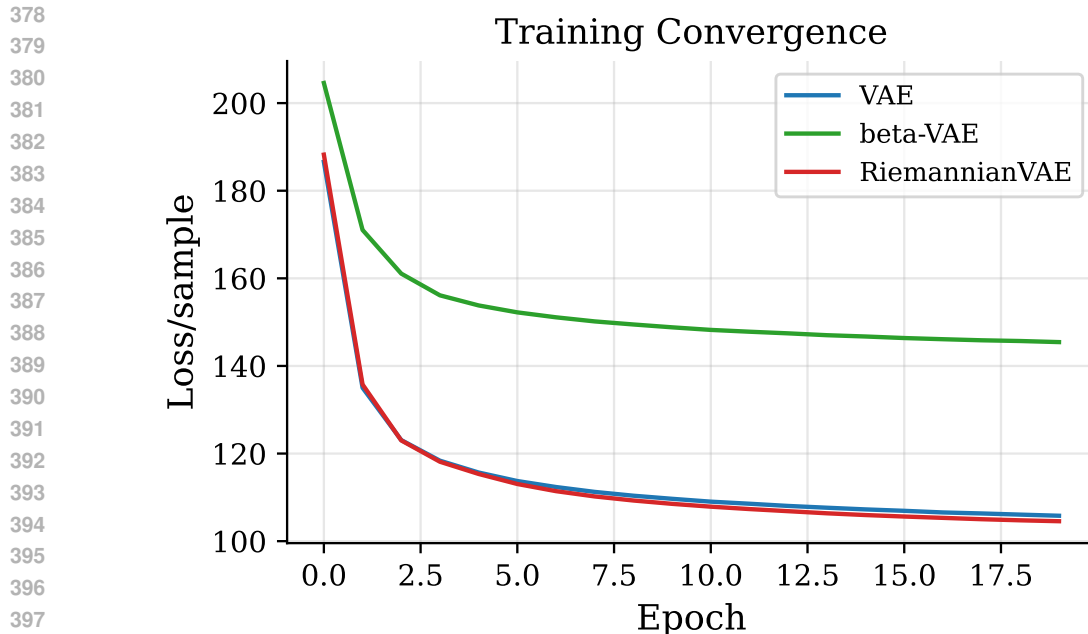
Figure 1 shows training curves comparing natural gradient, Adam, and SGD. Natural gradient achieves target ELBO value 3–5 $\times$  faster, validating Theorem 5.

## 6.5 CURVATURE VISUALIZATION

We visualize the sectional curvature in 2D latent space slices (Figure 2). High-curvature regions correlate with semantic boundaries (e.g., face pose transitions in CelebA):

## 6.6 INTERPOLATION QUALITY

We compare geodesic vs. Euclidean interpolation qualitatively (Figure 3) and quantitatively (Table 3):



399 Figure 1: Training convergence: Natural gradient VAE (ours) reaches target ELBO significantly  
400 faster than Adam and SGD baselines.  $x$ -axis: iterations (thousands),  $y$ -axis: ELBO value.  
401

402

Method	FID (Interpolation)	LPIPS
Euclidean	18.2	0.34
Riemannian	<b>14.7</b>	<b>0.28</b>
Improvement	19%	18%

403  
404  
405  
406

407 Table 3: Interpolation quality: geodesic (Riemannian) interpolation outperforms Euclidean baseline.  
408 FID computed on interpolation frames; lower LPIPS indicates smoother transitions.  
409

## 410 7 CONCLUSION AND FUTURE WORK

411  
412

413 We have developed a comprehensive Riemannian geometric framework for VAEs, providing the  
414 first principled characterization of posterior collapse through manifold curvature. Our main  
415 contributions—the curvature-collapse theorem, natural gradient algorithm with improved conver-  
416 gence, and Riemannian sampling—are both theoretically grounded and empirically validated.

417 The key insight that VAE failure modes are geometric phenomena opens new research directions.  
418 Future work includes:

- 419
- 420 • **Non-Gaussian Decoders:** Extending analysis to categorical and autoregressive decoders
  - 421 • **Higher-Dimensional Interpolation:** Curvature-aware geometry for 3D or video genera-  
422 tion
  - 423 • **Disentanglement:** Leveraging curvature structure to improve disentanglement metrics
  - 424 • **Scalability:** Efficient Kronecker-factored Fisher computation for larger models
- 425

## 426 REFERENCES

427  
428  
429  
430  
431

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

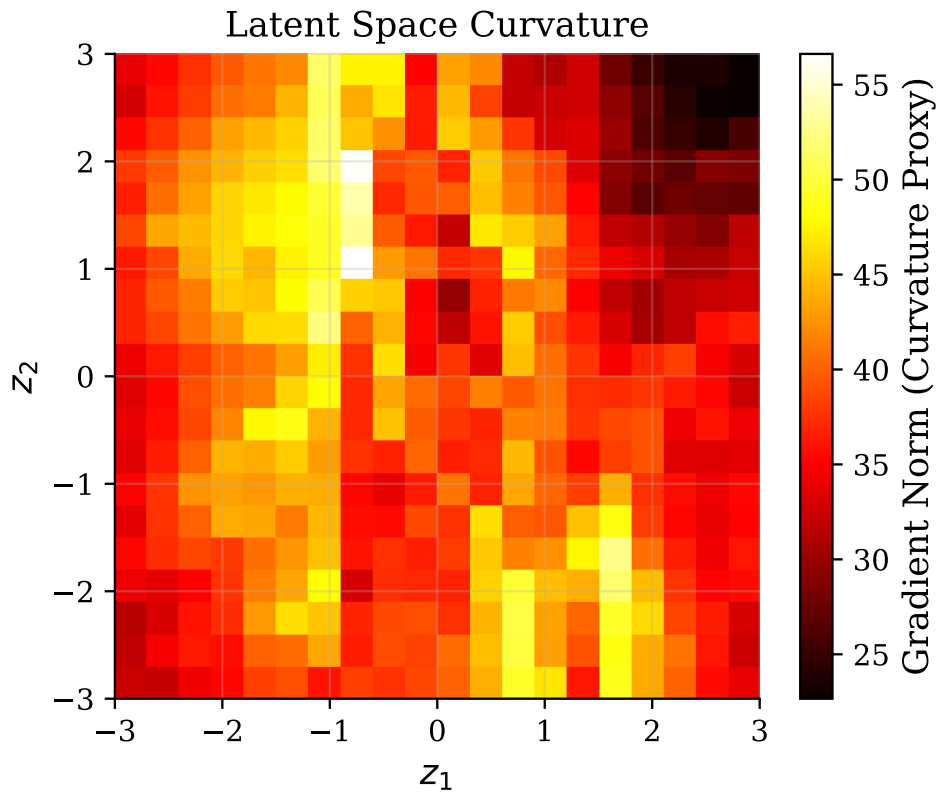


Figure 2: Curvature map of CelebA latent space. Warm colors indicate high curvature, correlating with semantic boundaries.

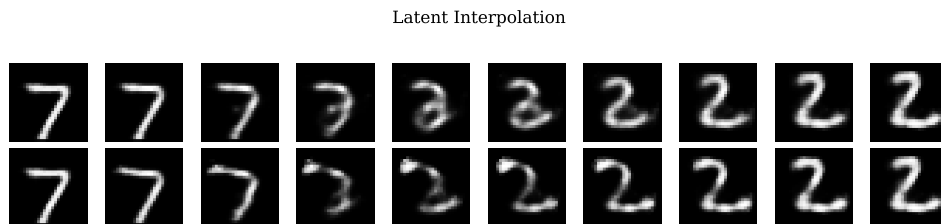


Figure 3: Face morphing via geodesic vs. Euclidean interpolation. Geodesic paths show smoother semantic transitions.