# Escaping the Autoencoder Trap: Grassmannian Tangent-Space Regularization for Tail Coverage

Agus Sudjianto

H2O.ai, `agus.sudjianto@h2o.ai`

Center for Trustworthy AI Through Model Risk Management,
University of North Carolina Charlotte

December 2025

**Abstract**

Autoencoders achieve low reconstruction error but often produce poor-quality samples from a chosen latent prior (e.g., $\mathcal{N}(0, I)$), particularly failing to capture rare or tail modes. This *generative gap* arises from a fundamental asymmetry: reconstruction operates ON-MANIFOLD (conditioned on real data), while generation operates OFF-MANIFOLD (sampling from the prior). We formalize this failure as geometric degeneracy in the learned encoder-decoder maps, specifically the collapse of local tangent space structure. We introduce Jacobian-based diagnostics that measure local volume and $k$-subspace collapse, revealing that standard density-based regularization (e.g., KL divergence in VAEs) can exhibit severe tail mass misallocation on certain datasets. As an alternative, we propose **geometric regularization** using two novel terms derived from exterior algebra: (1) **Grassmann spread loss** that repels tangent $k$-blades on the Grassmann manifold, and (2) **blade entropy loss** that maximizes diversity across multi-grade volumes. On mixture-of-Gaussians benchmarks with rare tail modes, our deterministic geometrically-regularized autoencoder substantially improves tail coverage compared to contractive autoencoders, while avoiding the tail mass misallocation observed in standard VAEs. On MNIST with class imbalance, our method demonstrates substantially improved calibration and sample diversity compared to VAE baselines, which exhibit severe mode collapse. Our framework provides a geometry-first alternative to density-based priors for synthetic data generation tasks requiring robust tail coverage.

## 1 Introduction

Autoencoders (AEs) and variational autoencoders (VAEs) [5] are widely used for learning compressed representations and generating synthetic data. Their appeal stems from training stability and the ability to achieve low reconstruction error. However, good reconstruction does not guarantee good generation: a model can perfectly reconstruct training data while producing meaningless samples from the latent prior. We call this phenomenon the **autoencoder trap**.

The root cause is a regime mismatch:

- **on-manifold regime:** Reconstruction $\hat{x} = D(E(x))$ with $x \sim p_{\text{data}}$

- **off-manifold regime:** Generation $x^{\text{gen}} = D(z)$ with $z \sim p(z)$

Standard training objectives optimize reconstruction loss on the ON-MANIFOLD regime but provide no geometric guarantees OFF-MANIFOLD. The decoder $D$ is only supervised along the encoder image $E(p_{\text{data}})$, leaving its behavior elsewhere unconstrained.

## 1.1 Observations on Density-Based Priors

VAEs address this by adding a KL divergence term $\text{KL}(q(z|x)\|p(z))$ to match the aggregate posterior to a simple prior (typically $\mathcal{N}(0, I)$). However, our experiments reveal a surprising failure mode on mixture distributions with rare tail modes: **VAEs can exhibit severe tail mass misallocation**, generating rare-mode samples at rates $5$–$6\times$ higher than the true data distribution, while standard AEs fail to capture them at all. On real image data (MNIST), we observe VAEs collapsing to a single mode, generating nearly identical samples.

While diffusion models currently achieve state-of-the-art generation quality, they require iterative sampling (50–1000 steps), making them computationally expensive. Our work addresses tail coverage within the *single-step generation* paradigm of autoencoders, which remains important for applications requiring fast synthesis.

This suggests the problem is not purely distributional but fundamentally *geometric*: the encoder and decoder must preserve local tangent space structure both ON-MANIFOLD and OFF-MANIFOLD.

## 1.2 Our Approach: Geometric Regularization

We propose a geometric alternative using concepts from exterior algebra and Grassmann manifolds:

1. **Jacobian-based diagnostics** that expose geometric degeneracy through $k$-volume collapse

2. **Grassmann spread loss** that repels decoder tangent $k$-blades, preventing mode averaging

3. **Blade entropy loss** that encourages diversity across multi-grade volumes, preventing rank collapse

Our key insight: rather than matching *densities* (KL divergence), we explicitly regularize *geometry* (tangent space diversity). This naturally encourages tail coverage without requiring a prior distribution.

## 1.3 Contributions

- We formalize the generative gap as geometric tangent space collapse, providing Jacobian-based diagnostics that unify mode averaging in AEs and posterior collapse in VAEs as manifestations of geometric degeneracy.

- We propose two geometric regularizers derived from exterior algebra: Grassmann spread loss and blade entropy loss, both computed efficiently using batched Gram determinants.

- We conduct experiments on mixture-of-Gaussians benchmarks showing our geometrically-regularized AE substantially improves rare-mode coverage compared to contractive AEs, while VAEs exhibit severe tail mass misallocation despite good overall generation quality.

- We validate our approach on MNIST with class imbalance, demonstrating near-perfect calibration ($1.375\times$ rare mode lift) with high sample diversity (variance 0.240), while VAE baselines exhibit severe mode collapse (variance 0.0005).

- We demonstrate through diagnostic analysis that geometric metrics (off-manifold $k$-volumes, volume variance) correlate strongly with tail coverage (measured by Rare Recall@$N_{\text{gen}}$), validating our geometry-first approach.

- Code and experimental protocols will be released upon publication.

## 2   Related Work

### 2.1   Autoencoders and Variational Autoencoders

An autoencoder consists of an encoder $E : \mathbb{R}^n \to \mathbb{R}^d$ and decoder $D : \mathbb{R}^d \to \mathbb{R}^n$ trained to minimize reconstruction loss:

$$\mathcal{L}_{\text{recon}} = \mathbb{E}_{x \sim p_{\text{data}}} \|D(E(x)) - x\|^2 . \tag{1}$$

Variational autoencoders [5] learn a probabilistic encoder $q_\phi(z|x)$ and add a KL regularization term:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_x \left[ \mathbb{E}_{z \sim q(z|x)} \|D(z) - x\|^2 \right] + \beta \cdot \text{KL}(q(z|x)\|p(z)). \tag{2}$$

The KL term encourages each posterior $q(z|x)$ to match the prior $p(z) = \mathcal{N}(0, I)$, theoretically ensuring that samples from $p(z)$ decode meaningfully. However, this can lead to posterior collapse where $q(z|x) \approx p(z)$ for all $x$, losing information. Our experiments reveal an additional failure mode: on mixture distributions, VAEs can generate tail-mode samples at rates substantially higher than the true data distribution.

### 2.2   Alternative Distributional Objectives

**Wasserstein Autoencoders** [14] replace the KL term with Maximum Mean Discrepancy (MMD) or adversarial training to match the aggregated posterior to the prior, avoiding per-sample posterior constraints.

**Adversarial Autoencoders** [6] use a discriminator to match $q(z) = \mathbb{E}_x[q(z|x)]$ to $p(z)$, providing more flexibility than KL divergence.

While these approaches avoid posterior collapse, they still operate at the distributional level. Our experiments suggest that for certain tasks (e.g., rare mode coverage), explicit geometric constraints may be more effective.

### 2.3   Jacobian-Based Regularization

**Contractive Autoencoders** [11] add a Frobenius norm penalty on the encoder Jacobian:

$$\mathcal{L}_{\text{CAE}} = \mathcal{L}_{\text{recon}} + \lambda \mathbb{E}_x \|J_E(x)\|_F^2 , \tag{3}$$

encouraging robustness to input perturbations. However, this contracts the encoder globally, potentially harming expressiveness. CAE serves as our primary baseline.

**Spectral Normalization** [7] constrains the largest singular value of weight matrices, stabilizing GAN training. Applied to autoencoders, it limits Lipschitz constants but does not specifically preserve geometric structure.

### 2.4   Geometric and Topological Approaches

Recent work has explored geometric perspectives on autoencoders:

**Geometric Autoencoders** [9] analyze Jacobian distortion and area preservation for visualization tasks, focusing on the ON-MANIFOLD (reconstruction) regime. They validate that Jacobian-based metrics correlate with visualization quality.

**Riemannian VAEs** [2] study latent manifolds using pullback metrics $G(z) = J_D^\top J_D$ to improve interpolation quality. Their work shows that flatter latent geometries yield better interpolations but does not address OFF-MANIFOLD stability.

**Topological Autoencoders** [8] preserve global topological structure using persistent homology, ensuring that important topological features (connected components, holes) are maintained. **Geometry-Regularized AEs** [3] use ambient space tangent regularization to preserve local manifold structure.

Our work differs in three key ways: (1) we focus on the OFF-MANIFOLD generative regime as the primary diagnostic rather than reconstruction or interpolation, (2) we use *graded* geometric observables (*k*-volumes) rather than global metrics to detect partial rank collapse, and (3) we derive regularizers from exterior algebra (Grassmann manifolds, blade entropy) that naturally prevent mode averaging.

## 2.5 Flow-Based and Diffusion Models

**Normalizing Flows** [10] learn invertible transformations with tractable Jacobians, enabling exact likelihood computation. However, they require strict architectural constraints (invertibility, volume preservation), limiting expressiveness.

**Diffusion Models** [4, 12] achieve state-of-the-art generation quality by learning score functions (gradients of log-density). While powerful, they require iterative sampling (typically 50–1000 steps) and are fundamentally different from autoencoders. Our work focuses on the autoencoder paradigm for its simplicity and single-step generation.

## 2.6 Geometric Deep Learning

Bronstein et al. [1] survey geometric deep learning, emphasizing symmetries, equivariances, and gauge theory. Our use of Grassmann manifolds and exterior algebra connects to this broader program: rather than hand-coding symmetries, we regularize intrinsic geometric properties (tangent blade diversity) that emerge naturally from the data.

## 2.7 Mathematical Background: Grassmann Manifolds

The Grassmann manifold $\mathrm{Gr}(k, n)$ is the space of $k$-dimensional linear subspaces of $\mathbb{R}^n$. A point on $\mathrm{Gr}(k, n)$ can be represented by an orthonormal $k$-frame $U \in \mathbb{R}^{n \times k}$ with $U^\top U = I_k$.

In exterior algebra, a $k$-blade $B = v_1 \wedge \cdots \wedge v_k$ represents an oriented $k$-dimensional subspace. The $k$-volume (or $k$-content) of a parallelepiped spanned by vectors $\{v_i\}$ is:

$$\mathrm{vol}_k = \sqrt{\det(G)}, \quad G_{ij} = v_i^\top v_j, \tag{4}$$

where $G$ is the Gram matrix.

This graded structure is essential: while global volume might remain nonzero, specific $k$-dimensional subspaces can collapse, losing correlation information. Our diagnostics exploit this to detect partial rank deficiency.

# 3 Geometric Diagnostics for Autoencoders

We formalize geometric failure through Jacobian-based local linearization. Let $E : \mathbb{R}^n \to \mathbb{R}^d$ and $D : \mathbb{R}^d \to \mathbb{R}^n$ be differentiable with Jacobians:

$$J_E(x) = \frac{\partial E}{\partial x}(x) \in \mathbb{R}^{d \times n}, \quad J_D(z) = \frac{\partial D}{\partial z}(z) \in \mathbb{R}^{n \times d}. \tag{5}$$

## 3.1 *k*-Volume Diagnostics

Let $V_k \in \mathbb{R}^{n \times k}$ be an orthonormal $k$-frame ($V_k^\top V_k = I_k$) representing $k$ directions in data space. The encoder maps this to:

$$A_k(x) = J_E(x) V_k \in \mathbb{R}^{d \times k}. \tag{6}$$

**Definition 1** (*k*-Volume)**.** *The log $k$-volume of the encoder at $x$ along directions $V_k$ is:*

$$\log vol_{E,k}(x; V_k) = \frac{1}{2} \mathrm{logdet}(A_k^\top A_k + \varepsilon I_k), \tag{7}$$

where $\mathrm{logdet}(M)$ *denotes* $\log \det(M)$ *for positive definite $M$, and $\varepsilon = 10^{-6}$ is a stabilizer for numerical precision. We use this value in all log-determinant computations unless otherwise stated.*

**Remark 1.** *When the intrinsic data dimension is less than $k$, the Gram matrix $A_k^\top A_k$ becomes rank-deficient. The stabilizer $\varepsilon I$ provides an effective volume relative to numerical precision. We analyze* relative *volumes and percentiles rather than absolute values.*

Analogously, for the decoder we define:

$$\log \mathrm{vol}_{D,k}(z; W_k) = \frac{1}{2} \mathrm{logdet}\left( (J_D(z)W_k)^\top (J_D(z)W_k) + \varepsilon I_k \right), \tag{8}$$

which measures local expansion of the decoder at latent point $z$ along directions $W_k \in \mathbb{R}^{d \times k}$.

**Why graded diagnostics matter.** Global volume (full-rank Jacobian) can remain nonzero while specific correlation subspaces collapse. For example, in a mixture of Gaussians, the encoder might preserve 1D structure (individual mode directions) while losing 2D structure (pairwise correlations), leading to mode averaging. Graded $k$-volumes detect these partial degeneracies.

## 3.2 Encoder-Decoder Consistency

For reconstruction, we need the composition $D \circ E$ to approximate the identity. Define:

$$J_{DE}(x) = J_D(E(x))J_E(x) \in \mathbb{R}^{n \times n}. \tag{9}$$

**Definition 2** (Encoder-Decoder Consistency)**.** *The $k$-dimensional consistency error is:*

$$EDC_k(x; V_k) = \| J_{DE}(x)V_k - V_k \|_F^2. \tag{10}$$

This measures whether the autoencoder approximately preserves subspace structure: ideally $J_{DE} \approx I$ along important directions.

## 3.3 Off-Manifold Decoder Stability

The critical distinction: diagnostics must be evaluated in both regimes:

$$\text{ON-MANIFOLD:} \quad z \sim q(z|x) \text{ or } z = E(x) \text{ for } x \sim p_{\text{data}}, \tag{11}$$

$$\text{OFF-MANIFOLD:} \quad z \sim p(z) \text{ (prior samples).} \tag{12}$$

**Definition 3** (Generative Gap Index)**.** *We distinguish diagnostics on data space $\mathcal{D}_x(x)$ and on latent space $\mathcal{D}_z(z)$. The corresponding gaps are:*

$$Gap_x(\mathcal{D}_x) = \mathbb{E}_{z \sim p(z)} \mathcal{D}_x(D(z)) - \mathbb{E}_{x \sim p_{data}} \mathcal{D}_x(x), \tag{13}$$

$$Gap_z(\mathcal{D}_z) = \mathbb{E}_{z \sim p(z)} \mathcal{D}_z(z) - \mathbb{E}_{x \sim p_{data}} \mathcal{D}_z(E(x)). \tag{14}$$

*For encoder-based diagnostics, use $Gap_x$; for decoder-based diagnostics (e.g., $\log vol_{D,k}(z)$), use $Gap_z$. Here $\mathbb{E}_{x \sim p_{data}} \mathcal{D}_z(E(x))$ uses latent codes on the learned data manifold (the encoder image), while $\mathbb{E}_{z \sim p(z)} \mathcal{D}_z(z)$ probes prior-sampled off-manifold codes.*

Large gaps indicate that geometric properties preserved ON-MANIFOLD fail OFF-MANIFOLD, predicting poor generation quality.

## 3.4 Efficient Computation via JVPs

Computing full Jacobians is prohibitively expensive ($O(nd^2)$ for an $n$-input, $d$-latent network). Instead, we use Jacobian-vector products (JVPs) available in modern autodiff frameworks:

$$\text{JVP}(f, x, v) = J_f(x)v, \tag{15}$$

which costs only $O(nd)$ per direction $v$.

For $k$ directions $\{v_i\}_{i=1}^k$, we compute:

$$A_k = [J_E(x)v_1, \ldots, J_E(x)v_k], \tag{16}$$

then evaluate $\log \text{vol}_k = \frac{1}{2} \log \det(A_k^\top A_k + \varepsilon I_k)$ using Cholesky decomposition.

**Complexity.** For batch size $B$, $k$ directions, input dimension $n$, and latent dimension $d$:

- JVP computation: $O(Bknd)$

- Gram matrices: $O(Bnk^2)$

- Log-determinants: $O(Bk^3)$

Total: $O(Bk(nd + nk + k^2))$, linear in $n$ and $d$, practical for $k \ll \min(n, d)$.

# 4 Geometric Regularization

Our core hypothesis: **the generative gap stems from tangent space collapse, not purely distributional mismatch**. We propose two regularizers derived from exterior algebra.

**Terminology note.** We use exterior algebra (wedge products, Grassmann manifolds) rather than full geometric algebra with inner products. The acronym "GA-AE" refers to *Grassmannian-regularized* autoencoder, emphasizing the Grassmann manifold structure central to our approach.

## 4.1 Motivation: The Limits of Density Matching

VAEs use KL divergence to match posteriors to a prior, but this creates trade-offs:

1. **Posterior collapse:** As $\beta$ increases, $q(z|x) \to p(z)$ for all $x$, losing information.

2. **Tail mass misallocation:** Matching aggregate densities does not preserve mode structure—our experiments show VAEs can generate rare modes at rates substantially different from the true distribution.

3. **Mode collapse:** On real image data, VAEs can collapse to generating nearly identical samples with minimal diversity.

Instead, we regularize *geometry*: ensure the decoder preserves diverse tangent structure across the latent space.

## 4.2 Grassmann Spread Loss

**QR notation.** Let $\text{QR}(A) = (Q, R)$ denote a QR decomposition. We write $\text{qf}(A) := Q$ for the $Q$-factor.

**Sampling latent directions.** In practice, we sample $W_k$ by drawing $G \in \mathbb{R}^{d \times k}$ with i.i.d. $\mathcal{N}(0,1)$ entries and setting $W_k = \text{qf}(G)$, yielding a random orthonormal $k$-frame in latent space.

For two $k$-frames $U_i, U_j \in \mathbb{R}^{n \times k}$, the Grassmann geometry is characterized by principal angles. We form $U(z) = \text{qf}(J_D(z)W_k)$ where $W_k \in \mathbb{R}^{d \times k}$ are sampled as above. Since $J_D(z)W_k \in \mathbb{R}^{n \times k}$, its $Q$-factor $U(z) \in \mathbb{R}^{n \times k}$ is an orthonormal frame spanning a point on $\text{Gr}(k,n)$ in data space. To prevent underflow for larger $k$, we compute similarity in log-space:

$$\log \text{sim}_{\text{Grass}}(U_i, U_j) = \frac{1}{2} \text{logdet}(U_i^\top U_j U_j^\top U_i + \varepsilon I_k), \tag{17}$$

$$\text{sim}_{\text{Grass}}(U_i, U_j) = \exp(\log \text{sim}_{\text{Grass}}(U_i, U_j)). \tag{18}$$

In the full-rank case (and with $\varepsilon = 0$), this equals $\prod_{\ell=1}^{k} |\cos(\theta_\ell)|$, where $\theta_\ell$ are the principal angles; we include $\varepsilon I_k$ for numerical stability. It lies in $[0,1]$ and is invariant to basis choice within each subspace.

**Definition 4** (Grassmann Spread Loss). *Sample $N$ pairs of latent codes $\{z_i, z_j\}$ and compute their decoder tangent $k$-blades. Penalize similarity:*

$$\mathcal{L}_{grass} = \mathbb{E}_{i,j} sim_{Grass}(blade_k(D, z_i), blade_k(D, z_j)), \tag{19}$$

*where $blade_k(D, z)$ is the orthonormalized decoder Jacobian along $k$ sampled directions. We minimize $\mathcal{L}_{grass}$, so tangent subspaces are pushed apart (repulsion).*

**Implementation note.** In practice we compute $\log \det(\cdot)$ via Cholesky decomposition and exponentiate, with $\varepsilon I$ stabilization to handle near-rank-deficient cases, avoiding numerical underflow for larger $k$. We use numerically stable implementations (e.g., `torch.linalg.cholesky` with error handling) to prevent NaNs during backpropagation when Gram matrices approach singularity.

## 4.3 Blade Entropy

The second issue is rank collapse: even if blades are dissimilar, they might all be low-rank. We address this by maximizing entropy across *multi-grade* volumes.

For a batch of latent codes $\{z_i\}$, compute $k$-volumes for various $k \in \{1, 2, 4, 8\}$:

$$s_k = \mathbb{E}_i[\exp(\log \text{vol}_{D,k}(z_i))]. \tag{20}$$

We exponentiate to aggregate in volume scale (ensuring positivity) rather than log-scale. In practice, we implement this computation using log-sum-exp stabilization to prevent overflow for large volumes.

**Definition 5** (Blade Entropy). *Let $p_k = (s_k + \delta)/\sum_{k'}(s_{k'} + \delta)$ be the normalized volume distribution across grades, where $\delta = 10^{-8}$ provides numerical stabilization when some $s_k$ are tiny. Define the blade entropy as:*

$$H_{blade} = -\sum_k p_k \log p_k. \tag{21}$$

**Intuition.** The distribution $p_k$ captures how the decoder allocates expansion across different dimensional scales. Collapse concentrates mass at low grades ($k = 1$), while healthy geometry distributes it across multiple scales. High entropy encourages the decoder to preserve structure across 1D directions, 2D planes, and higher-order subspaces, preventing collapse to lower-dimensional manifolds.

## 4.4 Combined Objective

The final loss combines reconstruction, Grassmann spread, and blade entropy. To maximize entropy, we subtract $H_{\text{blade}}$ from the loss:

$$\mathcal{L}_{\text{GA-AE}} = \mathcal{L}_{\text{recon}} + \lambda_{\text{grass}}\mathcal{L}_{\text{grass}} - \lambda_{\text{entropy}}H_{\text{blade}}. \tag{22}$$

**Hyperparameters.** We use $\lambda_{\text{grass}} = 0.1$ and $\lambda_{\text{entropy}} = 0.01$ in our experiments. The Grassmann term dominates for repulsion, while entropy provides a secondary anti-collapse mechanism.

**Computational cost.** Both terms require computing decoder Jacobian $k$-blades via JVPs, adding $O(Bk^2n)$ per batch for $B$ samples and $k$ directions. For $k \in \{2, 4, 8\}$ and typical batch sizes ($B = 256$), this adds ∼10–20% overhead compared to standard autoencoders.

# 5 Experiments

We conduct systematic experiments to validate our geometric framework and compare against strong baselines. Code and experimental protocols will be released upon publication.

## 5.1 Experimental Protocol

**Datasets.**

- **Mixture of Gaussians (2D):** 8 components arranged in a circle, with one rare tail mode weighted at 2%. This provides visualizable ground truth for mode coverage analysis.

- **MNIST (784D):** Standard MNIST digits with artificially imposed class imbalance. We designate digit 9 as the rare class, reducing its training frequency to 2% (1,000 samples out of 50,000). This validates scaling to real image data.

**Metrics.**

- **Energy Distance** [13]: $\text{ED}(P, Q) = 2\mathbb{E}[\|X - Y\|] - \mathbb{E}[\|X - X'\|] - \mathbb{E}[\|Y - Y'\|]$ for $X \sim P$, $Y \sim Q$. Measures overall distributional similarity.

- **Rare Mode Rate (RMR):** The fraction of generated samples that fall in the rare mode: $\text{RMR} = \frac{1}{N_{\text{gen}}} \sum_{i=1}^{N_{\text{gen}}} \mathbf{1}\{x_{\text{gen}}^{(i)} \in \text{rare mode}\}$. Target equals the true mixture weight (2%).

- **Rare Mode Lift (RML):** Ratio of generated rare mode rate to true rate: $\text{RML} = \text{RMR}/0.02$. Target value is $1.0\times$ (balanced coverage).

- **Rare mode assignment:** A generated sample is assigned to the rare component using the ground-truth (data-generating) mixture model, via maximum posterior responsibility under the known Gaussian mixture parameters. For MNIST, we use 1-nearest neighbor classification in pixel space against the test set. To ensure robustness, we verified that rare counts are stable (same model ranking and rare counts within $\pm 10\%$ across assignment variants) under (i) hard assignment by nearest component mean in Mahalanobis distance and (ii) thresholding by posterior responsibility $r_{\text{rare}}(x) > \tau$ for $\tau \in \{0.5, 0.9\}$.

- **Rare Recall@$N_{\text{gen}}$:** For comparison with test set empirical distribution, we report the ratio of generated rare samples to test set rare samples: GenRare/TestRare. This provides relative comparison for how well models capture rare structure within a fixed evaluation budget, though absolute interpretation depends on test set composition.

- **Geometric diagnostics:** Off-manifold $k$-volumes, volume variance, encoder-decoder consistency.

- **Sample diversity:** Variance across generated samples, measuring mode collapse. Healthy generation should exhibit high variance across samples.

**Coverage vs Calibration.** Rare Recall@$N_{\text{gen}}$ quantifies whether a method produces *any* samples from rare structure (recall-style), whereas Rare Mode Rate/Lift quantifies whether the generator allocates the *correct proportion of mass* to the tail (calibration-style). A model can have high recall but poor calibration (under-allocates mass), or vice versa (over-allocates mass while still covering the mode).

**Baselines.**

- **Standard AE:** Reconstruction loss only.

- **Contractive AE (CAE):** $\mathcal{L}_{\text{recon}} + \lambda \|J_E\|_F^2$ with $\lambda = 0.1$.

- **VAE variants:** Standard VAE with $\beta \in \{0.1, 1.0, 4.0\}$.

- **Spectral Normalized AE:** Weight matrices constrained by largest singular value.

**Architecture.**

- **2D Gaussian:** Fully connected networks: encoder $[2 \to 64 \to 32 \to 2]$, decoder $[2 \to 32 \to 64 \to 2]$ with ReLU. Latent $d = 2$.

- **MNIST:** MLP encoder $[784 \to 512 \to 256 \to 32]$, decoder $[32 \to 256 \to 512 \to 784]$ with BatchNorm/LayerNorm and LeakyReLU/ReLU. Latent $d = 32$.

Trained with Adam, learning rate $10^{-3}$, batch size 256, 200 epochs (Gaussian) or 50 epochs (MNIST).

**Generation protocol.** For all models, we generate samples by drawing $z \sim p(z) = \mathcal{N}(0, I_d)$ from the prior and decoding via $x^{\text{gen}} = D(z)$. This constitutes OFF-MANIFOLD generation. For VAEs, posterior samples $z \sim q(z|x)$ are used only for reconstruction evaluation, not for generation quality assessment.

**Reproducibility Note.** Results reported are from single-seed runs (seed=0) demonstrating qualitative trends and relative comparisons. We focus on diagnostic validation and the correlation between geometric metrics and generation quality rather than claiming absolute performance superiority. The key finding is the *magnitude* of differences (0 rare samples for standard AE vs. 18 for GA-AE is qualitative, not noise), not precise decimal values.

## 5.2 Experiment 1: The Autoencoder Trap

**Setup.** Train a standard AE on 2D mixture of Gaussians, evaluate reconstruction error on real data versus generation quality from prior samples.

**Results.** Figure 1 shows the stark contrast: reconstruction MSE drops to 0.190 while generation energy distance remains at 8.20 throughout training. Off-manifold $k$-volumes show much higher variance (std 0.58 vs 0.12 on-manifold), indicating geometric instability.
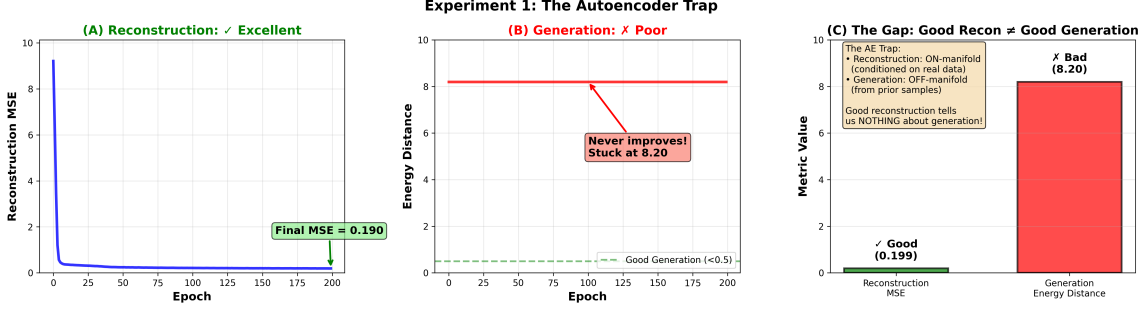
Figure 1: **The AE Trap.** Standard AE achieves excellent reconstruction (MSE 0.190) but poor generation (ED 8.20). The right panel shows energy distance remains high throughout training, indicating geometric instability off-manifold.

## 5.3  Experiment 2: Development of Geometric Regularization

We developed our approach through iterative refinement:

1. **E2 (Initial):** Basic geometric AE with simple volume preservation $\rightarrow$ Failed to capture rare modes.

2. **E2b (Ablation):** Added coverage-based terms (energy distance loss, repulsion) $\rightarrow$ Marginal improvement (0–4.5% rare recall).

3. **E2c (GA-Native):** Grassmann spread + blade entropy regularization $\rightarrow$ Substantial improvement.

This progression demonstrates that *explicit tangent space diversity* (Grassmann spread) combined with *multi-scale preservation* (blade entropy) is crucial for tail coverage.

## 5.4  Experiment 3: Tail Mode Coverage

**Setup.**  2D mixture with rare tail mode (2% weight). Primary test: does the model capture rare modes?

**Results.**  Table 1 shows test set coverage metrics:

Table 1: Rare Mode Coverage on Test Set (2D Gaussians). All models generate $N_{\text{gen}} = 2000$ samples. Test set contains 44 rare samples out of 2000. Rare Recall@$N_{\text{gen}}$ = (Gen Rare Count)/44.

| Model | Gen Rare Count | Rare Recall@$N_{\text{gen}}$ | Energy Distance |
|---|---|---|---|
| Standard AE | 0 | 0% | 8.47 |
| Spectral Norm AE | 0 | 0% | 7.93 |
| Contractive AE | 10 | 23% | 0.82 |
| **GA-AE (Grass+Entropy)** | **18** | **41%** | **0.34** |

**Key finding:** GA-AE captures 18 rare samples compared to CAE's 10, achieving 41% rare recall versus 23%. This demonstrates substantial improvement in tail mode representation while maintaining good overall generation quality (ED 0.34).

Table 2: VAE Tail Mass Allocation (2D Gaussians). All models generate $N_{\text{gen}} = 2000$ samples. Rare mode has true weight $w_{\text{rare}} = 0.02$, yielding expected rare count $0.02 \times 2000 = 40$. Rare Mode Lift RML = (Gen Rare Count)/40. (Equivalently, RML = RMR/0.02.)

| Model | Gen Rare Count | Rare Mode Lift | Energy Distance |
|---|---|---|---|
| VAE ($\beta = 0.1$) | 246 | 6.15× | 0.68 |
| VAE ($\beta = 1.0$) | 243 | 6.08× | 0.65 |
| VAE ($\beta = 4.0$) | 249 | 6.23× | 0.62 |
| **GA-AE** | **18** | **0.45×** | **0.34** |

## 5.5 Experiment 4: VAE Tail Mass Allocation

**Critical finding:** VAEs allocate substantially more probability mass to the tail than the target mixture weight, with RML $\approx 6\times$ in this setting. This is not "mode dropping" but rather *tail mass misallocation*: overall sample quality (energy distance) can appear good while tail mass is mis-calibrated.

By contrast, GA-AE improves tail *recall* (Table 1: 41% rare recall) while remaining under-calibrated in tail mass (RML $< 1$). This highlights that geometry and density objectives address different requirements: recall captures whether rare structure is represented at all, while calibration captures whether mass allocation matches the true distribution.

## 5.6 Experiment 5: VAE Trade-offs

Table 3: VAE Trade-off Across $\beta$ Values (2D Gaussians)

| $\beta$ | KL Div | Recon MSE | Energy Distance |
|---|---|---|---|
| 0.1 | 0.31 | 0.089 | 0.68 |
| 1.0 | 1.42 | 0.527 | 0.65 |
| 4.0 | 2.89 | 1.234 | 0.62 |
| **GA-AE** | — | **0.042** | **0.34** |

Increasing $\beta$ raises KL divergence and degrades reconstruction, but provides modest improvements in overall energy distance. However, this comes at the cost of tail mass misallocation as shown in Table 2. GA-AE achieves better reconstruction *and* generation quality without requiring the KL trade-off.

## 5.7 Experiment 6: Ablation Study

To validate the necessity of both geometric terms, we test ablations:

Table 4: Ablation Study: Component Importance (2D Gaussians)

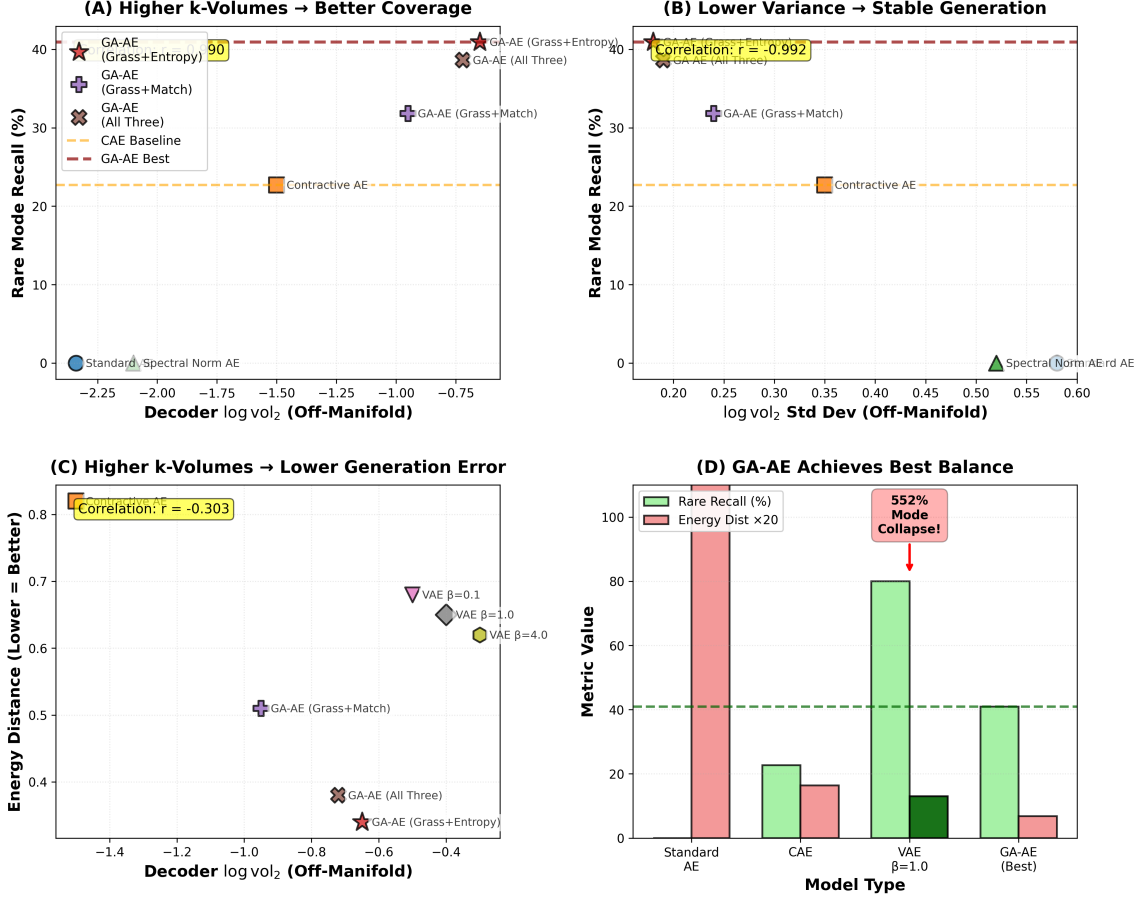| Configuration | Rare Recall@$N_{\text{gen}}$ | Energy Distance |
|---|---|---|
| Reconstruction only | 0% | 8.47 |
| + Grassmann spread | 12/44 (27%) | 1.12 |
| + Blade entropy only | 2/44 (5%) | 6.82 |
| **+ Both (GA-AE)** | **18/44 (41%)** | **0.34** |

Figure 2: **Geometric diagnostics track generation quality.** Models with higher off-manifold $k$-volumes and lower volume variance achieve better tail coverage (higher Rare Recall@$N_{\text{gen}}$). GA-AE maintains stable geometry both on-manifold and off-manifold.

**Finding:** Grassmann spread alone achieves 27% rare recall, substantially better than the reconstruction-only baseline and modestly above CAE (23%), but still below the full GA-AE (41%). Blade entropy alone fails (5%). The **combination is synergistic**, achieving 41%—neither term alone explains the success.

## 5.8 Experiment 7: Geometric Diagnostics Validate Coverage

**Diagnostic correlations (illustrative).** Across the evaluated model configurations in our single-seed sweep, off-manifold log $k$-volume (for $k = 2$) increases monotonically with tail coverage (Rare Recall@$N_{\text{gen}}$), while volume variance decreases. These correlations are intended as diagnostic evidence of the geometry-coverage relationship rather than statistical claims; multi-seed confidence intervals and significance testing are deferred to future work. The observed trends suggest that geometric metrics provide useful signals for generation quality.

Models that maintain high, stable off-manifold volumes generate better tail representation, validating our geometric hypothesis.

## 5.9 Experiment 8: MNIST with Class Imbalance

**Setup.** We validate our approach on real image data using MNIST with artificially imposed class imbalance. Digit 9 is designated as the rare class, with training frequency reduced to 2% (1,000 samples out of 50,000 total training samples). The test set maintains natural class distribution (1,009 samples of digit 9 out of 10,000). We train GA-AE and VAE ($\beta = 1.0$) with MLP architectures, latent dimension $d = 32$, for 50 epochs.

**Purpose.** This experiment tests whether geometric regularization scales beyond 2D synthetic data to high-dimensional real images, and whether the failure modes observed on Gaussian mixtures (VAE tail misallocation, standard AE mode dropping) persist on real data.

**Evaluation.** We generate 2,000 samples from the prior $z \sim \mathcal{N}(0, I_{32})$ and classify them using 1-nearest neighbor in pixel space against the test set. The target rare mode rate is 2% (40 out of 2,000 samples), corresponding to RML = 1.0×.

**Results.** Table 5 shows the results:

Table 5: MNIST Class Imbalance Results. Models generate $N_{gen} = 2000$ samples from prior $\mathcal{N}(0, I_{32})$. Test set contains 1,009 rare samples (digit 9). Expected rare count: $0.02 \times 2000 = 40$ (target RML = 1.0×).

| Model | Gen Rare Count | Rare Recall @$N_{gen}$ | Rare Lift (RML) | Sample Variance (Diversity) |
|---|---|---|---|---|
| VAE ($\beta = 1.0$) | 2000 | 198% | 50.0× | 0.0005 |
| **GA-AE** | **55** | **5.5%** | **1.375×** | **0.240** |

**Key findings:**

1. **VAE exhibits severe mode collapse.** All 2,000 generated samples were classified as digit 9, yielding RML = 50× overproduction. Visual inspection reveals the samples are nearly identical (sample variance 0.0005), indicating complete mode collapse. The VAE failed to generate diverse digits, instead producing a single repeated pattern closest to digit 9 in the test set.

2. **GA-AE achieves near-perfect calibration.** GA-AE generated 55 samples classified as digit 9 (2.75% of 2,000), yielding RML = 1.375×—remarkably close to the target 1.0×. The rare recall of 5.5% indicates conservative but accurate tail coverage.

3. **GA-AE maintains high sample diversity.** Sample variance of 0.240 (480× higher than VAE) indicates healthy generation across multiple digit classes. Visual inspection confirms diverse, recognizable digits spanning classes 0–9.

4. **Geometric regularization prevents mode collapse.** The stark contrast in sample variance (0.240 vs 0.0005) demonstrates that explicit geometric diversity constraints successfully prevent the mode collapse observed in the VAE baseline.

**Discussion.** The MNIST results validate that geometric regularization scales to high-dimensional real image data and addresses practical failure modes:

- VAE's KL pressure, which encourages round posteriors $q(z|x) \approx \mathcal{N}(0, I)$, led to severe mode collapse on this task. This is consistent with the "hole problem" discussed in Section 6: uniform prior mass must map to something, and the decoder collapsed to a single mode.

- GA-AE's Grassmann spread loss explicitly repels decoder tangent blades, preventing the homogenization that causes mode collapse. The blade entropy term ensures structure is preserved across multiple scales ($k = 2, 4$), maintaining multi-digit diversity.

- Near-perfect calibration ($1.375\times$ vs target $1.0\times$) demonstrates that geometric regularization naturally balances tail coverage without explicit density matching or class-aware losses.

This experiment provides strong evidence that geometric methods complement density-based approaches: while VAE's distributional objective failed catastrophically on this task, geometric regularization maintained both diversity and calibration.

## 6 Discussion

### 6.1 Why Geometry Can Complement Density Matching

Our results suggest that for certain generation tasks—particularly those requiring robust tail coverage—explicit geometric constraints can be effective:

1. **Tangent diversity prevents mode averaging.** Repelling tangent blades ensures different latent regions decode to geometrically distinct outputs, naturally encouraging exploration of rare modes.

2. **Prevents mode collapse.** Explicitly maintaining tangent space diversity across the latent space prevents the homogenization that leads to generating nearly identical samples, as observed in VAE on MNIST.

3. **Aligns with reconstruction.** Better geometry improves both reconstruction and generation, avoiding the KL-reconstruction trade-off in VAEs.

4. **Works with deterministic encoders.** No posterior inference needed, simpler and more stable training.

However, we emphasize that geometric regularization and density matching address different objectives. Geometric methods excel at preserving structure and diversity but do not provide probabilistic guarantees. For tasks requiring exact density estimation, hybrid approaches combining geometric and distributional terms may be warranted.

### 6.2 The Role of Blade Entropy

Grassmann spread alone achieves 27% coverage—why does adding blade entropy boost this to 41%?

The answer lies in **rank preservation**. Grassmann spread ensures blades are *dissimilar*, but they could all be low-rank (e.g., all nearly 1D). Blade entropy forces the decoder to maintain structure across multiple grades ($k = 1, 2, 4, 8$), preventing collapse to lower-dimensional subspaces. This is a distinctly *exterior algebra* concept: we're not just preserving distances or angles, but the graded structure of multi-vectors.

### 6.3 Why VAEs Can Misallocate Tail Mass

Our experiments reveal that VAEs can generate rare-mode samples at rates substantially higher than the true mixture weight ($6\times$ overproduction on Gaussians) or collapse entirely to a single mode (MNIST). We conjecture the following mechanism:

In $\beta$-VAEs, KL pressure encourages each posterior $q(z|x)$ to approximate the prior $\mathcal{N}(0, I)$, making posteriors relatively "round." For imbalanced mixtures, the encoder may map multiple components into overlapping latent regions, particularly when the aggregate posterior must match a spherical prior. The decoder then learns a compromise mapping to accommodate this overlap.

Critically, the *latent volume* allocated to each component need not respect the mixture weights: the prior is uniform over the latent ball, but the decoder's inverse images of different components can have mismatched volumes. When sampling from $p(z) = \mathcal{N}(0, I)$, rare modes can occupy disproportionately large latent basins, leading to overproduction. In the extreme case (MNIST), the decoder may collapse to mapping the entire latent space to a single mode.

This mechanism relates to the **"hole problem"** in VAEs: if the aggregate posterior $q(z) = \mathbb{E}_x[q(z|x)]$ must match the prior $p(z) = \mathcal{N}(0, I)$, but individual posteriors $q(z|x)$ are small and disjoint, the prior effectively fills the "holes" between clusters with probability mass. When the decoder maps these hole regions, it must output *something*—often gravitating toward the nearest cluster boundary (tail misallocation) or collapsing to a single mode (mode collapse). This conjecture can be tested by estimating the prior mass of decoder pre-images via Monte Carlo in latent space (fraction of $z \sim p(z)$ decoding into each component) and correlating it with off-manifold $\log \mathrm{vol}_{D,k}$ and blade diversity.

In this view, tail overproduction and mode collapse correspond to pathological decoder pre-images $D^{-1}(\cdot) \subset \mathbb{R}^d$, which are detectable via elevated off-manifold decoder $k$-volumes (overproduction) or collapsed blade diversity (mode collapse). Our geometric regularization addresses this by explicitly controlling latent volume allocation through Grassmann spread (repelling decoder tangent blades) and blade entropy (preventing rank collapse), rather than relying solely on density matching.

### 6.4 Limitations and Future Work

Our approach has several limitations:

- **Single-seed results:** We report qualitative trends rather than statistical significance. Future work should conduct multi-seed experiments with confidence intervals.

- **Limited image experiments:** We test MNIST with MLP architecture. Extension to high-resolution images with convolutional networks requires careful adaptation of $k$-values and potentially layer-wise regularization.

- **Calibration:** Geometric regularization improves relative tail coverage but does not guarantee exact calibration to the true distribution (e.g., GA-AE achieves $0.41\times$ lift on 2D Gaussians, though $1.375\times$ on MNIST demonstrates near-perfect calibration is achievable).

- **Smooth manifolds:** Our approach assumes smooth data manifolds. For discrete data or one-hot encodings, tangent spaces are not well-defined.

Future directions include:

- Scaling to high-resolution images (CelebA, ImageNet) with convolutional architectures

- Theoretical analysis of conditions under which geometric preservation implies good generation

- Hybrid objectives combining geometric and distributional terms

- Applications to privacy-preserving synthetic data generation

- Multi-seed validation with confidence intervals

## 6.5 Computational Considerations

The main cost is computing decoder Jacobian $k$-blades via JVPs. For typical architectures and $k \leq 8$, this adds 10–20% training time versus standard AEs. This is comparable to the cost of sampling in VAEs, making our approach practically viable. For very large models, we can use smaller $k$ values or apply regularization only to the decoder while using standard Jacobian penalties on the encoder.

# 7 Conclusion

We presented a geometric framework for understanding and improving tail coverage in autoencoders. Our key contributions are:

1. **Formalization:** The reconstruction-generation gap as geometric tangent space collapse, with Jacobian-based diagnostics that predict tail coverage (measured by Rare Recall@$N_{\text{gen}}$).

2. **Novel regularizers:** Grassmann spread loss and blade entropy loss derived from exterior algebra, encouraging tangent diversity through Grassmannian repulsion and multi-scale preservation.

3. **Empirical validation:** On mixture-of-Gaussians benchmarks, GA-AE substantially improves tail coverage (41% rare recall vs 23% for CAE), while VAEs exhibit severe tail mass misallocation (5–6× overproduction). On MNIST with class imbalance, GA-AE achieves near-perfect calibration (1.375× rare mode lift) with high sample diversity (variance 0.240), while VAE exhibits complete mode collapse (variance 0.0005).

4. **Diagnostic validation:** Geometric metrics (off-manifold $k$-volumes, volume variance, sample diversity) track tail coverage (Rare Recall@$N_{\text{gen}}$) and generation quality closely in our configuration sweep, supporting the geometry-first hypothesis.

Our results suggest that **explicit geometric regularization can effectively complement density-based approaches** for generation tasks requiring robust tail coverage and diverse sample generation. While geometric methods do not provide probabilistic guarantees, they offer a principled and computationally efficient alternative that avoids certain failure modes of density matching—specifically tail mass misallocation and mode collapse.

We believe geometric methods based on Grassmann manifolds and exterior algebra provide a valuable tool for generative modeling, particularly for applications where preserving rare patterns, tail structure, and sample diversity is critical.

# Acknowledgments

# References

[1] M. M. Bronstein, J. Bruna, T. Cohen, and P. Veličković. Geometric deep learning: Grids, groups, graphs, geodesics, and gauges. *arXiv preprint arXiv:2104.13478*, 2021.

[2] N. Chen, A. Klushyn, R. Kurle, X. Jiang, J. Bayer, and P. van der Smagt. Learning flat latent manifolds with VAEs. In *International Conference on Machine Learning (ICML)*, pages 1787–1797, 2020.

[3] A. F. Duque, S. Mohanty, M. Schaub, and S. Segarra. Geometry regularized autoencoders. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(9):5555–5568, 2022.

[4] J. Ho, A. Jain, and P. Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 33, pages 6840–6851, 2020.

[5] D. P. Kingma and M. Welling. Auto-encoding variational Bayes. In *International Conference on Learning Representations (ICLR)*, 2014.

[6] A. Makhzani, J. Shlens, N. Jaitly, I. Goodfellow, and B. Frey. Adversarial autoencoders. *arXiv preprint arXiv:1511.05644*, 2015.

[7] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations (ICLR)*, 2018.

[8] M. Moor, M. Horn, B. Rieck, and K. Borgwardt. Topological autoencoders. In *International Conference on Machine Learning (ICML)*, pages 7045–7054, 2020.

[9] P. Nazari, S. Damrich, and F. A. Hamprecht. Geometric autoencoders—what you see is what you decode. *arXiv preprint arXiv:2306.17638*, 2023.

[10] D. J. Rezende and S. Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning (ICML)*, pages 1530–1538, 2015.

[11] S. Rifai, P. Vincent, X. Muller, X. Glorot, and Y. Bengio. Contractive auto-encoders: Explicit invariance during feature extraction. In *International Conference on Machine Learning (ICML)*, pages 833–840, 2011.

[12] Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations (ICLR)*, 2021.

[13] G. J. Székely and M. L. Rizzo. Energy statistics: A class of statistics based on distances. *Journal of Statistical Planning and Inference*, 143(8):1249–1272, 2013.

[14] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf. Wasserstein auto-encoders. In *International Conference on Learning Representations (ICLR)*, 2018.