

Supplementary for Isometric Manifold Learning using Hierarchical Flow

2

3 1. Proofs and Analysis

4 In this section, we provide detailed proofs and analysis for our theoretical claims for the properties of *distance preserving*
 5 and *rigorous projection* introduced in main text.

6 1.1. The Distance Preserving Property

7 We provide proofs for Prop. 1 (restated as Prop. 1) and Prop. 3 (restated as Prop. 2) in main text for satisfying the *distance*
 8 *preserving* property.

9 **Proposition 1** (*Distance preserving property*). Assume that $\mathcal{U}^1(\mathcal{M})$ is a convex set, and $v^i(x) = 0$ for $\forall i \in [L], x \in \mathcal{M}$.
 10 The length of the geodesic between $g(u_1)$ and $g(u_2)$ on \mathcal{M} for $\forall u_1, u_2 \in \mathcal{U}^1(\mathcal{M})$ equals to $\|u_1 - u_2\|_2$, if $J_g(u)$ is
 11 orthonormal for $\forall u \in \mathcal{U}^1(\mathcal{M})$, where $J_g(u)$ is the Jacobian matrix of the generator g at u .

12 *Proof.* According to Riemannian geometry [12], given a Riemannian manifold \mathcal{M} , a smooth curve $\gamma : [a, b] \rightarrow \mathcal{M}$ on \mathcal{M} is a
 13 geodesic of \mathcal{M} , if $\nabla_{\gamma'(t)}\gamma'(t) \equiv 0, \forall t \in [a, b]$. Since $v^i(x) = 0, \forall i \in [L], x \in \mathcal{M}$, we know that $g(u) \in \mathcal{M}, \forall u \in \mathcal{U}^1(\mathcal{M})$.
 14 Therefore, we can learn that $\gamma(t) \triangleq g(tu_2 + (1-t)u_1)$ is a curve on \mathcal{M} for $t \in [0, 1]$. We firstly prove that such a curve
 15 $\gamma : [0, 1] \rightarrow \mathcal{M}$ is a geodesic connecting $g(u_1)$ and $g(u_2)$ on \mathcal{M} . Since $\gamma' = (u_2 - u_1)^i \frac{\partial}{\partial z^i}$, where $(u_2 - u_1)^i$ is the i -th
 16 element of vector $u_2 - u_1 \in \mathbb{R}^K$, we further have that

$$\nabla_{\gamma'}\gamma' = \nabla_{(u_2 - u_1)^j \frac{\partial}{\partial z^j}} (u_2 - u_1)^i \frac{\partial}{\partial z^i} = \xi^j \xi^i \nabla_{\frac{\partial}{\partial z^j}} \frac{\partial}{\partial z^i}, \quad (1)$$

17 where we use $\xi \triangleq u_2 - u_1$. It is known that the coefficient Γ_{ij}^k of the connection ∇ is given by the *Christoffel symbol*

$$\Gamma_{ij}^k = \frac{1}{2} g^{kl} \left(\frac{\partial g_{il}}{\partial z^j} + \frac{\partial g_{jl}}{\partial z^i} - \frac{\partial g_{ij}}{\partial z^l} \right), \quad (2)$$

18 where g_{ij} and g^{ij} are the covariant and contravariant components of the Riemannian metric \mathcal{G} , respectively, and the matrix
 19 form of \mathcal{G} is induced by $J_g^\top J_g$, i.e., $(J_g^\top J_g)_{ij} = g_{ij}$, where $(J_g^\top J_g)_{ij}$ is the (i, j) -th element of $J_g^\top J_g$. Since J_g is orthonormal,
 20 we have that $g_{ij} \equiv 0$ for $i \neq j$ and $g_{ii} \equiv 1$ for $i = j$. Hence we further have that $\frac{\partial g_{ij}}{\partial z^l} \equiv 0$ and $\Gamma_{ij}^k \equiv 0$. So we
 21 have

$$\nabla_{\frac{\partial}{\partial z^j}} \frac{\partial}{\partial z^i} = \Gamma_{ij}^k \frac{\partial}{\partial z^k} \equiv 0, \quad (3)$$

22 and hence $\nabla_{\gamma'}\gamma' \equiv 0$, i.e., γ is a geodesic connecting $g(u_1)$ and $g(u_2)$ on \mathcal{M} . To further prove that $d_{\mathcal{M}}(g(u_1), g(u_2)) =$
 23 $\|u_1 - u_2\|_2$, according to Riemannian geometry, we have

$$d_{\mathcal{M}}(g(u_1), g(u_2)) \triangleq \int_0^1 \sqrt{\sum_{ij} g_{ij} \frac{dz^i}{dt} \frac{dz^j}{dt}} dt = \int_0^1 \sqrt{\sum_i g_{ii} (\xi^i)^2} dt = \sqrt{\sum_i (\xi^i)^2} \int_0^1 dt = \|u_1 - u_2\|_2, \quad (4)$$

24 where we use $g_{ii} \equiv 1$ due to the orthonormality of J_g . The convex property of $\mathcal{U}^1(\mathcal{M})$ is to ensure that $tu_2 + (1-t)u_1 \in$
 25 $\mathcal{U}^1(\mathcal{M}), \forall t \in [0, 1]$. Summarizing the above analysis completes the proof. \square

26 **Proposition 2.** $J_g(z)$ is orthonormal for $\forall z \in \mathcal{U}^1(\mathcal{M})$, if $J_{f_i}(u)$ is orthonormal for $\forall u \in \mathcal{A}^{i-1}(\mathcal{M}; r), \forall i \in [L]$.

¹The Einstein summation convention is used, namely $X^i \frac{\partial}{\partial z^i} \triangleq \sum_{i=1}^K X^i \frac{\partial}{\partial z^i}$, where z^i is the i -th coordinate.

27 *Proof.* Given $g = f_L \circ p_L \circ f_{L-1} \circ p_{L-1} \circ \cdots \circ f_1 \circ p_1$, we have that $\forall u^1 \in \mathcal{U}^1(\mathcal{M})$,

$$J_g(u^1) = J_{f_L}(u^L) I_{K_L \times K_{L-1}}(u^L) J_{f_{L-1}}(u^{L-1}) I_{K_{L-1} \times K_{L-2}}(u^{L-1}) \cdots J_{f_1}(u^1) I_{K_1 \times K_0}(u^1), \quad (5)$$

where we use $I_{K_i \times K_{i-1}} \in \mathbb{R}^{K_i \times K_{i-1}}$ to denote an identity matrix (*i.e.*, all diagonal elements are 1 meanwhile others are 0), and $u^{i+1} \triangleq f_i(u^i)$. Since $u^1 \in \mathcal{U}^1(\mathcal{M})$, we have $u^{i+1} \in \tilde{\mathcal{A}}^i(\mathcal{M}; r) \subset \mathcal{A}^i(\mathcal{M}; r)$ according to the definition of $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ and $\mathcal{A}^i(\mathcal{M}; r)$ in main text, namely $u^i \in \mathcal{A}^{i-1}(\mathcal{M}; r)$. Hence we know that $J_{f_i}(u^i)$ is orthonormal. Therefore, we have

$$J_g^\top(u^1) J_g(u^1) = I_{K_1 \times K_0}^\top(u^1) J_{f_1}^\top(u^1) \cdots I_{K_{L-1} \times K_{L-2}}^\top(u^{L-1}) J_{f_{L-1}}^\top(u^{L-1}) I_{K_L \times K_{L-1}}^\top(u^L) J_{f_L}^\top(u^L) \quad (6)$$

$$J_{f_L}(u^L) I_{K_L \times K_{L-1}}(u^L) J_{f_{L-1}}(u^{L-1}) I_{K_{L-1} \times K_{L-2}}(u^{L-1}) \cdots J_{f_1}(u^1) I_{K_1 \times K_0}(u^1) \quad (7)$$

$$= I_{K_1 \times K_0}^\top(u^1) J_{f_1}^\top(u^1) \cdots I_{K_{L-1} \times K_{L-2}}^\top(u^{L-1}) J_{f_{L-1}}^\top(u^{L-1}) \quad (8)$$

$$J_{f_{L-1}}(u^{L-1}) I_{K_{L-1} \times K_{L-2}}(u^{L-1}) \cdots J_{f_1}(u^1) I_{K_1 \times K_0}(u^1) \quad (9)$$

$$= I_{K_1 \times K_0}^\top(u^1) J_{f_1}^\top(u^1) \cdots J_{f_1}(u^1) I_{K_1 \times K_0}(u^1) \quad (10)$$

$$= I_{K_0 \times K_0}, \quad (11)$$

28 namely $J_g(u^1)$ is orthonormal, where $I_{K_i \times K_{i-1}}^\top(u^i) I_{K_i \times K_{i-1}}(u^i) = I_{K_{i-1} \times K_{i-1}}$ and $J_{f_i}^\top(u^i) J_{f_i}(u^i) = I_{K_i \times K_i}$. \square

29 1.2. The Rigorous Projection Property

30 We firstly provide proofs for Prop. 2 (restated as Prop. 3) and Prop. 4 (restated as Prop. 4) in main text for satisfying the
31 *rigorous projection* property, and then provide intuitive understandings about the *rigorous projection* property. We finally
32 prove Cor. 1 (restated as Cor. 1) that claims the *rigorous projection* property for an one-layer generator, and provide further
33 analysis on Rem. 1 (restated as Rem. 1) that discusses the *rigorous projection* property for a multi-layer generator.

34 **Proposition 3** (*Rigorous projection* property). *Assume that $v^i(x) = 0$ for $\forall x \in \mathcal{M}$. Given $x \in \mathcal{A}^i(\mathcal{M}; r)$, by denoting
35 $[\alpha; \beta] \triangleq f_i^{-1}(x)$ where $\alpha \in \mathbb{R}^{K_{i-1}}$ and $\beta \in \mathbb{R}^{K_i - K_{i-1}}$, the projection of x onto $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is $\tilde{x} \triangleq f_i(\alpha)$, and the distance
36 from x to $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$, $d_{\mathcal{E}}(x, \tilde{x})$, equals to $\|\beta\|_2$, if $\forall u \in \mathcal{A}^{i-1}(\mathcal{M}; r), v \in \mathcal{V}^i(r)$,*

$$I([u; v]) \triangleq J([u; v])^\top J([u; v]) \in \mathbb{R}^{K_i \times K_i} \quad (12)$$

37 is a diagonal matrix, and $I_{jj}([u; v])$ equals to 1 for $K_{i-1} + 1 \leq j \leq K_i$, where $J([u; v]) \in \mathbb{R}^{K_i \times K_i}$ is the Jacobian of f_i
38 at $[u; v]$, I_{jj} is the (j, j) -th element of I , and $d_{\mathcal{E}}(\cdot, \cdot)$ is the Euclidean distance in \mathbb{R}^{K_i} .

39 *Proof.* Given $\forall u \in \mathcal{A}^{i-1}(\mathcal{M}; r), \forall v \in \mathcal{V}^i(r)$, we firstly prove that the curve $\gamma(t) \triangleq f_i([u; tv])$ is a geodesic on $\mathcal{A}^i(\mathcal{M}; r)$.
40 Because $\gamma' = v^m \frac{\partial}{\partial z^m}$ (note that $m \in \mathcal{I}_v$, where \mathcal{I}_v is the set of indices corresponding to the v coordinates of f_i), we further
41 have that

$$\nabla_{\gamma'} \gamma' = \nabla_{v^n \frac{\partial}{\partial z^n}} v^m \frac{\partial}{\partial z^m} = v^n v^m \nabla_{\frac{\partial}{\partial z^n}} \frac{\partial}{\partial z^m}. \quad (13)$$

42 Let \mathbf{g}_{mn} be the covariant components of the Riemannian metric \mathcal{G} , then from Eq. (12) we have $\mathbf{g}_{mn} \equiv 0$ for $m \neq n \in [K_i]$
43 and $\mathbf{g}_{mm} \equiv 1$ for $m \in \mathcal{I}_v$. Therefore, for $m, n \in \mathcal{I}_v$, we have $\frac{\partial \mathbf{g}_{ml}}{\partial z^n} \equiv 0, \frac{\partial \mathbf{g}_{nl}}{\partial z^m} \equiv 0, \frac{\partial \mathbf{g}_{mn}}{\partial z^l} \equiv 0$. Hence from Eq. (2), we
44 have $\Gamma_{mn}^k \equiv 0$, and $\nabla_{\frac{\partial}{\partial z^n}} \frac{\partial}{\partial z^m} = \Gamma_{mn}^k \frac{\partial}{\partial z^k} \equiv 0$, namely $\nabla_{\gamma'} \gamma' \equiv 0$. Hence γ is a geodesic on $\mathcal{A}^i(\mathcal{M}; r)$. Since $\mathcal{A}^i(\mathcal{M}; r)$
45 is essentially the ambient Euclidean space, we know that γ is a straight line, since the geodesics of an Euclidean space are
46 straight lines. Due to the orthogonality of J , we know that $\gamma'(0) \perp T_{\gamma(0)} \tilde{\mathcal{A}}^i(\mathcal{M}; r)$, where $T_{\gamma(0)} \tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is the tangent
47 space of \mathcal{M} at $\gamma(0)$. Therefore, we can learn that $\gamma \perp T_{\gamma(0)} \tilde{\mathcal{A}}^i(\mathcal{M}; r)$, and the projection of any $x \in \gamma$ onto $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is
48 $\tilde{x} \triangleq f_i(u)$. To further prove that $d_{\mathcal{E}}(x, \tilde{x}) = \|v\|_2$ for $x = f_i([u; v])$, the projection of x onto $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is
49 $\tilde{x} \triangleq f_i(u)$. To further prove that $d_{\mathcal{E}}(x, \tilde{x}) = \|v\|_2$ for $x = f_i([u; v])$, similar to Eq. (4), we have

$$d_{\mathcal{E}}(x, \tilde{x}) = d_{\mathcal{A}^i(\mathcal{M}; r)}(x, \tilde{x}) = \int_0^1 \sqrt{\sum_{m, n \in \mathcal{I}_v} \mathbf{g}_{mn} \frac{dz^m}{dt} \frac{dz^n}{dt}} dt = \sqrt{\sum_m (v^m)^2} \int_0^1 dt = \|v\|_2, \quad (14)$$

50 where we use $\mathbf{g}_{mm} \equiv 1$ for $i \in \mathcal{I}_v$. Summarizing the above analysis completes the proof. \square

51 **Proposition 4.** All singular values of $J_{\phi_i^z}(\mathbb{1})$ equal $a > 0 \Leftrightarrow J_{f_i}(z)^\top J_{f_i}(z) = \text{diag}\left(\frac{a^2}{\omega_1^2}, \dots, \frac{a^2}{\omega_{K_i}^2}\right)$, and $\left\{\frac{a}{\omega_j}\right\}_{j=1}^{K_i}$ are
52 the singular values of $J_{f_i}(z)$.

53 Proof. Note that $J_{\phi_i^z}(\mathbb{1}) = J_{f_i}(z) \text{diag}(\omega_1, \omega_2, \dots, \omega_{K_i})$ from Eq. 13 in main text. Since all singular values of $J_{\phi_i^z}(\mathbb{1})$
54 equal $a > 0$, we have that $J_{\phi_i^z}^\top(\mathbb{1}) J_{\phi_i^z}(\mathbb{1}) = a^2 I$ where I is an identity matrix, and therefore $S^\top J_{f_i}^\top(z) J_{f_i}(z) S = a^2 I$,
55 where $S \triangleq \text{diag}(\omega_1, \dots, \omega_{K_i})$, namely

$$J_{f_i}^\top(z) J_{f_i}(z) = \text{diag}\left(\frac{a^2}{\omega_1^2}, \dots, \frac{a^2}{\omega_{K_i}^2}\right), \quad (15)$$

56 which means that J_{f_i} is orthogonal. Let the SVD of J_{f_i} be $J_{f_i} = U\Sigma V^\top$, where $\Sigma \triangleq \text{diag}(\sigma_1, \dots, \sigma_K)$, then we have that

$$J_{f_i}^\top J_{f_i} = V\Sigma^\top \Sigma V^\top = V \text{diag}(\sigma_1^2, \dots, \sigma_K^2) V^\top = \text{diag}\left(\frac{a^2}{\omega_1^2}, \dots, \frac{a^2}{\omega_{K_i}^2}\right). \quad (16)$$

57 Note that Eq. (16) is a SVD of $\text{diag}(\sigma_1^2, \dots, \sigma_K^2)$, hence up to permutation we have that $\sigma_j^2 = \frac{a^2}{\omega_j^2}$, namely $\left\{\frac{a}{\omega_j}\right\}_{j=1}^{K_i}$ are the
58 singular values of $J_{f_i}(z)$. \square

59 We then reveal that how the following objective

$$\begin{aligned} & \text{constrain all singular values of } J_{\phi_i^z}(\mathbb{1}) \text{ to be equal,} \\ & \text{s.t. } \omega_j = \sqrt{\left(\prod_{l=1}^{K_{i-1}} \omega_l\right) \det J_{f_i}}, \quad K_{i-1} + 1 \leq j \leq K_i \end{aligned} \quad (17)$$

60 mentioned in main text leads to an f_i satisfying the conditions for the *rigorous projection* property in Prop. 3. Given that all
61 singular values of $J_{\phi_i^z}(\mathbb{1})$ are equal, from Prop. 4 we have that $\det J_{f_i} = \prod_{l=1}^{K_i} \frac{a}{\omega_l}$, therefore from Eq. (17) we have that

$$\omega_j = \sqrt{\left(\prod_{l=1}^{K_{i-1}} \omega_l\right) \left(\prod_{l=1}^{K_i} \frac{a}{\omega_l}\right)} = \sqrt{a^{K_i} \prod_{l=K_{i-1}+1}^{K_i} \frac{1}{\omega_l}}, \quad K_{i-1} + 1 \leq j \leq K_i. \quad (18)$$

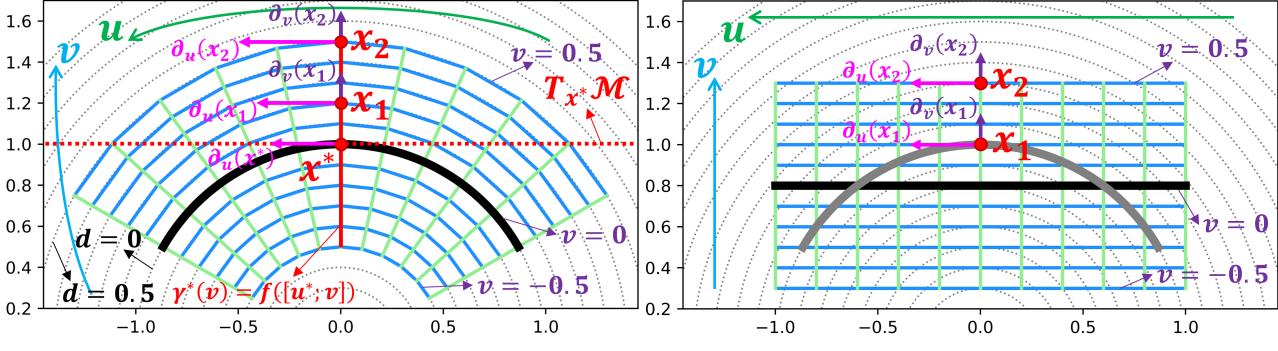
62 Since we know that $\omega_{K_{i-1}+1} = \omega_{K_{i-1}+2} = \dots = \omega_{K_i}$ from Eq. (17), we denote that $\omega_{K_{i-1}+1} = \omega_{K_{i-1}+2} = \dots = \omega_{K_i} = \xi$.
63 From Eq. (18) we know that $\xi^{K_{i-1}} = a^{K_i} \xi^{K_{i-1}-K_i}$, namely $\frac{a^{K_i}}{\xi^{K_i}} = 1$ and hence $\frac{a}{\xi} = 1$. Since we know that $\sigma_j = \frac{a}{\omega_j}$ from
64 Prop. 4, we know that $\sigma_j = \frac{a}{\xi} = 1$ for $K_{i-1} + 1 \leq j \leq K_i$, where σ_j is a singular value of J_{f_i} . Note that J_{f_i} is orthonormal
65 up to a scalar, therefore $J_{f_i}^\top J_{f_i}$ is a diagonal matrix with the diagonal elements being the squared singular values of J_{f_i} . From
66 the above analysis, we know that the conditions for the *rigorous projection* property in Prop. 3 are satisfied.

67 We now provide intuitive understandings about the *rigorous projection* property by using Fig. 1, including how the conditions
68 stated in Prop. 3 leads to a flow module f_i that satisfies the *rigorous projection* property, and why we use Eq. 18 in main
69 text instead of a simpler regularization that trivially constrains J_{f_i} to be orthonormal over $\mathcal{A}^{i-1}(\mathcal{M}; r) \times \mathcal{V}^i(r)$. For intuitive
70 illustration, we use a 1-dimensional manifold $\mathcal{M} \triangleq \{(\cos \theta, \sin \theta) | \theta \in [\frac{\pi}{6}, \frac{5\pi}{6}]\}$ (i.e., a curve) residing in the 2-dimensional
71 Euclidean space as we introduced in main text.

Intuitively Understanding Prop. 3 We consider a flow module $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ that satisfies the conditions in Prop. 3 and successfully recovers \mathcal{M} , namely there exists some $\mathcal{U} \subset \mathbb{R}^1$ such that $\mathcal{M} = \{f([u; 0]) | u \in \mathcal{U}\}$. The coordinate chart of f is plotted in Fig. 1(a). Consider the coordinate curve $\gamma^* \triangleq \{f([u^*; v]) | v \in [-0.5, 0.5]\}$ for $u^* \in \mathcal{U}$ as shown in Fig. 1(a). We know that $x^* \triangleq f([u^*; 0]) \in \gamma^*$ is on \mathcal{M} and hence $\partial_u(x^*)$ is the basis of $T_{x^*}\mathcal{M}$, namely $T_{x^*}\mathcal{M} = \{a\partial_u(x^*) | a \in \mathbb{R}\}$ (see Fig. 1(a)). Since f satisfies the conditions in Prop. 3, we know that for $\forall x \in \gamma^*$,

$$\partial_u(x) \perp \partial_v(x), \quad (19)$$

$$\|\partial_v(x)\|_2 = 1, \quad (20)$$



(a) An intuitive understanding on the conditions stated in Prop. 3. (b) The trivial regularization restricts the flexibility of the flow module.

Figure 1. Intuitive understandings about the *rigorous projection* property by using the synthetic manifold \mathcal{M} introduced in main text. In Fig. 1(a), we visualize the coordinate chart of a flow module $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ that satisfies the conditions in Prop. 3, meanwhile in Fig. 1(b), we show the coordinate chart of a flow module $f' : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ whose Jacobian is trivially constrained to be orthonormal over the ambient space. The ground-truth manifold \mathcal{M} and the manifolds recovered by the flow modules are plotted by using the thick gray solid line and the thick black solid lines, respectively, and the contour lines with respect to distances to \mathcal{M} are plotted by using the thin grey dotted lines, namely samples on the same contour line $d = \alpha$ have equal distances (i.e., α) to \mathcal{M} . For the coordinate charts of the two flow modules, the two embedding dimensions u and v are represented in green and blue, respectively, namely samples on the same green (resp., blue) line corresponds to the same u (resp., v). For both flow modules, we use $\partial_u(x)$ (resp., $\partial_v(x)$) to denote the tangent vector $\frac{\partial}{\partial u}\Big|_x$ (resp., $\frac{\partial}{\partial v}\Big|_x$) of the coordinate curve u (resp., v) at x . In Fig. 1(a), given $x^* \in \mathcal{M}$, we use $T_{x^*} \mathcal{M}$ to denote the tangent space of \mathcal{M} at x^* .

which results in $\partial_u(x)$ being parallel along γ^* , namely $\partial_u(x_1) \parallel \partial_u(x_2)$ for $\forall x_1, x_2 \in \gamma^*$ (see Fig. 1(a)). Combined with Eq. (19), we can learn that γ^* is a straight line perpendicular to $\partial_u(x^*)$, and hence we have $\text{vec}(x^*, x) \perp \partial_u(x^*)$, i.e., $\text{vec}(x^*, x) \perp T_{x^*} \mathcal{M}$. So according to the rigorous definition of projection, we know that the projection of $x = f([u^*; \beta]) \in \gamma^*$ onto \mathcal{M} is $x^* = f([u^*; 0])$. Regarding the distance from x to x^* , given that γ^* is straight, from Eq. (20) we learn that

$$d_{\mathcal{E}}(x, x^*) = \int_0^\beta \|\partial_v dv\|_2 = \int_0^\beta |dv| = |\beta|, \quad (21)$$

which coincides with the results given in Prop. 3. The above analysis provides an intuitive understanding on how the conditions stated in Prop. 3 leads to a flow module f that satisfies the *rigorous projection* property.

Intuitively Understanding the Isometric Regularization on the Ambient Space We consider a flow module $f' : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ that is constrained to be orthonormal over the ambient space $\mathcal{A} \triangleq \{f'([u; v]) \mid u \in \mathcal{U}, v \in [-0.5, 0.5]\}$ as in Fig. 1(b). Since $J_{f'}$ is orthonormal, we know that $\partial_u(x) \perp \partial_v(x)$ and $\|\partial_u(x)\|_2 = \|\partial_v(x)\|_2 = 1, \forall x \in \mathcal{A}$. Hence intuitively, all coordinate curves of u (resp., v) are straight lines and are perpendicular to the coordinate curves of v (resp., u), and all coordinate curves of u (resp., v) are the same length, as we shown in Fig. 1(b). This leads to a trivial case where the generated manifold (plotted in thick black solid line in Fig. 1(b)) cannot be curved, and hence it is impossible to successfully recover the ground-truth manifold (plotted in thick gray solid line in Fig. 1(b)) by using f' . Therefore, trivially constraining $J_{f'}$ to be orthonormal over the ambient space can severely restrict the flexibility of f' , and our proposed isometric regularization (Eq. 18 in main text) allows the flow module to be more flexible meanwhile satisfying the *rigorous projection* property.

1.2.1 Case of One-layer Generator

We prove Cor. 1 that claims the *rigorous projection* property for an one-layer generator.

Corollary 1. Given $g = f_1 \circ p_1$, assume that $v^1(x) = 0$ for $\forall x \in \mathcal{M}$. Given $x \in \mathcal{A}^1(\mathcal{M}; r)$, the projection of x onto \mathcal{M} is $\tilde{x} \triangleq g(u^1(x)) = g \circ g^\dagger(x)$, and the distance from x to \mathcal{M} , $d_{\mathcal{E}}(x, \tilde{x})$, is $\|v^1(x)\|_2$, if f_1 satisfies the conditions for the *rigorous projection* property in Prop. 3.

Proof. Since $x = f_1([u^1(x); v^1(x)])$ and $\tilde{\mathcal{A}}^i(\mathcal{M}; r) = \mathcal{M}$ in the case of one-layer generator, from Prop. 3, the projection of x onto \mathcal{M} is $\tilde{x} = f_1(u^1(x)) = g_1(u^1(x)) = g \circ g^\dagger(x)$, and distance from x to \mathcal{M} , $d_{\mathcal{E}}(x, \tilde{x})$ is $\|v^1(x)\|_2$. \square

94 **1.2.2 Case of Multi-layer Generator**

95 We provide further analysis on Rem. 1 that discusses the *rigorous projection* property for a multi-layer generator.

96 **Remark 1.** Given g comprised of $L \geq 2$ layers, assume that $v^i(x) = 0$ for $\forall x \in \mathcal{M}, i \in [L]$. Given $x \in \mathcal{A}^L(\mathcal{M}; r)$, let
 97 $\tilde{x}^i \triangleq f_i(u^i(x))$, $i \in [L]$, and $\tilde{x} \triangleq g(u^1(x)) = g \circ g^\dagger(x)$. Assume that r is small such that $\mathcal{A}^L(\mathcal{M}; r)$ is an ambient space
 98 close to \mathcal{M} . Given f_i that satisfies the rigorous projection property for $i \in [L]$, we have

99 • The projection of $u^{i+1}(x)$ onto $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is \tilde{x}^i , and the distance from $u^{i+1}(x)$ to $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$, $d_{\mathcal{E}}(u^{i+1}(x), \tilde{x}^i)$,
 100 equals to $\|v^i(x)\|_2$, where $i \in [L]$.

101 • The projection of x onto \mathcal{M} is \tilde{x} , and the distance from x to \mathcal{M} , $d_{\mathcal{E}}(x, \tilde{x})$, is approximately $\sqrt{\sum_{i=1}^L \|v^i(x)\|_2^2}$.

102 Given $x \in \mathcal{A}^L(\mathcal{M}; r)$ and f_i that satisfies the rigorous projection property for $i \in [L]$, let $\tilde{x}^i \triangleq f_i(u^i(x))$, then we can
 103 directly know from Prop. 3 that the projection of $u^{i+1}(x)$ onto $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$ is \tilde{x}^i , and the distance from $u^{i+1}(x)$ to $\tilde{\mathcal{A}}^i(\mathcal{M}; r)$,
 104 $d_{\mathcal{E}}(u^{i+1}(x), \tilde{x}^i)$, equals to $\|v^i(x)\|_2$. Given $\tilde{x} \triangleq g(u^1(x))$, assume that \tilde{x} is in a local region $\mathcal{O}^L \subset \tilde{\mathcal{A}}^L(\mathcal{M}; r)$ around \tilde{x}^L ,
 105 i.e., it can be regarded that $\mathcal{O}^L \approx T_{\tilde{x}^L} \tilde{\mathcal{A}}^L(\mathcal{M}; r)$. Since \tilde{x}^L is the projection of x onto $\tilde{\mathcal{A}}^L(\mathcal{M}; r)$, according to the rigorous
 106 definition of projection, we have $\text{vec}(x, \tilde{x}^L) \perp T_{\tilde{x}^L} \tilde{\mathcal{A}}^L(\mathcal{M}; r)$, and hence $\text{vec}(x, \tilde{x}^L) \perp \mathcal{O}^L$ approximately holds. Because
 107 $\tilde{x} \in \mathcal{O}^L$, we have that $\text{vec}(x, \tilde{x}^L) \perp \text{vec}(\tilde{x}^L, \tilde{x})$ approximately holds, and hence $d_{\mathcal{E}}^2(x, \tilde{x}) \approx d_{\mathcal{E}}^2(x, \tilde{x}^L) + d_{\mathcal{E}}^2(\tilde{x}^L, \tilde{x})$.
 108 Note that $d_{\mathcal{E}}(x, \tilde{x}^L) = \|v^L(x)\|_2^2$. For $d_{\mathcal{E}}(\tilde{x}^L, \tilde{x})$, due to the locality of \mathcal{O}^L , we have that $d_{\mathcal{E}}(\tilde{x}^L, \tilde{x}) \approx d_{\tilde{\mathcal{A}}^L(\mathcal{M}; r)}(\tilde{x}^L, \tilde{x})$. Moreover,
 109 because $\tilde{x} = f_L(u^L(\tilde{x}))$ and f_L is an isometry, we have $d_{\tilde{\mathcal{A}}^L(\mathcal{M}; r)}(\tilde{x}^L, \tilde{x}) = d_{\mathcal{E}}(u^L(x), u^L(\tilde{x}))$. Therefore,
 110 summarizing the above analysis gives us that

$$d_{\mathcal{E}}^2(x, \tilde{x}) \approx \|v^L(x)\|_2^2 + d_{\mathcal{E}}^2(u^L(x), u^L(\tilde{x})). \quad (22)$$

111 The above analysis can be further generalized to the next layer. Since f_{L-1} satisfies the rigorous projection property, we can
 112 learn that the projection of $u^L(x)$ onto $\tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)$ is \tilde{x}^{L-1} , and $d_{\mathcal{E}}(u^L(x), \tilde{x}^{L-1}) = \|v^{L-1}(x)\|_2$. Similarly, assume
 113 that $u^L(\tilde{x})$ is in a local region $\mathcal{O}^{L-1} \subset \tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)$ around \tilde{x}^{L-1} , i.e., it can be regarded that $\mathcal{O}^{L-1} \approx T_{\tilde{x}^{L-1}} \tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)$.
 114 According to the rigorous definition of projection, we have $\text{vec}(u^L(x), \tilde{x}^{L-1}) \perp T_{\tilde{x}^{L-1}} \tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)$, and hence we further
 115 have that $\text{vec}(u^L(x), \tilde{x}^{L-1}) \perp \mathcal{O}^{L-1}$ approximately holds. Because $u^L(\tilde{x}) \in \mathcal{O}^{L-1}$, we have that $\text{vec}(u^L(x), \tilde{x}^{L-1}) \perp$
 116 $\text{vec}(\tilde{x}^{L-1}, u^L(\tilde{x}))$ approximately holds. Hence we have $d_{\mathcal{E}}^2(u^L(x), u^L(\tilde{x})) \approx d_{\mathcal{E}}^2(u^L(x), \tilde{x}^{L-1}) + d_{\mathcal{E}}^2(\tilde{x}^{L-1}, u^L(\tilde{x}))$,
 117 where $d_{\mathcal{E}}(\tilde{x}^{L-1}, u^L(\tilde{x})) \approx d_{\tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)}(\tilde{x}^{L-1}, u^L(\tilde{x}))$ due to the locality of \mathcal{O}^{L-1} . Because $u^L(\tilde{x}) = f_{L-1}(u^{L-1}(\tilde{x}))$
 118 and f_{L-1} is an isometry, we further have $d_{\tilde{\mathcal{A}}^{L-1}(\mathcal{M}; r)}(\tilde{x}^{L-1}, u^L(\tilde{x})) = d_{\mathcal{E}}(u^{L-1}(x), u^{L-1}(\tilde{x}))$. Therefore, summarizing
 119 the above analysis gives us that

$$d_{\mathcal{E}}^2(u^L(x), u^L(\tilde{x})) \approx \|v^{L-1}(x)\|_2^2 + d_{\mathcal{E}}^2(u^{L-1}(x), u^{L-1}(\tilde{x})). \quad (23)$$

120 Therefore, combining Eq. (22) and Eq. (23), we have that

$$d_{\mathcal{E}}^2(x, \tilde{x}) \approx \|v^L(x)\|_2^2 + \|v^{L-1}(x)\|_2^2 + d_{\mathcal{E}}^2(u^{L-1}(x), u^{L-1}(\tilde{x})). \quad (24)$$

121 By further generalizing Eq. (24), we have that

$$d_{\mathcal{E}}^2(x, \tilde{x}) \approx \|v^L(x)\|_2^2 + \|v^{L-1}(x)\|_2^2 + \cdots + \|v^1(x)\|_2^2 + d_{\mathcal{E}}^2(u^1(x), u^1(\tilde{x})). \quad (25)$$

122 From $\tilde{x} = g(u^1(x))$, we know that $u^1(x) = u^1(\tilde{x})$, namely $d_{\mathcal{E}}(u^1(x), u^1(\tilde{x})) = 0$ and hence

$$d_{\mathcal{E}}(x, \tilde{x}) \approx \sqrt{\sum_{j=1}^L \|v^j(x)\|_2^2}. \quad (26)$$

123 2. Supplemental Experiments

124 In this section, we provide supplemental experimental results on the synthetic manifold introduced in main text and natural
 125 image datasets for manifold learning including MNIST, CelebA, FFHQ and AFHQ Cat, etc. We evaluate models in terms of
 126 the *distance preserving* property, the *rigorous projection* property and the *anomaly detection* performance. We also analyze
 127 the time-efficiency for our proposed *two-stage dimensionality reduction* algorithm on natural image datasets. We compare
 128 different manifold learning models including *Generative Adversarial Net* [3] (GAN), *Variational Autoencoder* [8] (VAE),
 129 *Isometric Autoencoder* [4] (IAE), *Pseudo Invertible Encoders* [1] (PIE), \mathcal{M} -flow [2], and our *Hierarchical Flow* (HF).

130 **Architectural Details** Our experiments involve convolutional or fully-connected neural network models including GAN,
 131 VAE, IAE, and NF-based models including PIE, \mathcal{M} -flow and our proposed HF. For an introduction on NF modules such as
 132 *coupling* layers, see [10]. The architectural details of models involved in our experiments are shown in Tbl. 1-2.

133 **Experimental Settings** We introduce implementation details for conducting our experiments. In main text, we provide the
 134 training objectives for the manifold phase and the density phase of our HF. To see the detailed implementation of the LT
 135 technique mentioned in main text for estimating the spectral norm of Jacobians in a sample-efficient manner, see [10]. Based
 136 on the training objectives introduced in main text, we also provide the detailed overall training algorithm for our HF model
 137 as in Alg. 2. In terms of training settings, we train 200,000 iterations with batch size 64, using an Adam optimizer [7] with
 138 fixed learning rate $1e^{-4}$. All models are train from scratch in an end-to-end manner by using PyTorch [11]. For settings of
 139 evaluation metrics corresponding to specific datasets, we introduce along with the respective experiments in the following
 140 Sections 2.1-2.2. All our experiments are done by using a single NVIDIA RTX-2080Ti GPU.

141 2.1. Results on the Synthetic Manifold

142 In main text, we provide demonstrative qualitative results showing that our proposed HF model can learn the manifold in
 143 a manner that satisfies the *distance preserving* property and the *rigorous projection* property while the other models such as
 144 \mathcal{M} -flow and IAE are not guaranteed to satisfy the two properties. In this section, we provide more quantitative results.

145 **Experimental Settings** We introduce the metrics we used to evaluate models. To evaluate the *distance preserving* property,
 146 we propose to compute the *Mean Absolute Error* (MAE) between $\frac{d_{\mathcal{M}}(x_1, x_2)}{\|z_1 - z_2\|_2}$ and 1, because for an isometric g , the manifold
 147 distance $d_{\mathcal{M}}(x_1, x_2)$ equals to the corresponding embedding distance $\|z_1 - z_2\|_2$, where $x_1, x_2 \in \mathcal{M}$ and $x_1 = g(z_1), x_2 =$
 148 $g(z_2)$, and $d_{\mathcal{M}}$ is the manifold distance that can be analytically computed as the arc length between x_1 and x_2 considering
 149 that \mathcal{M} is a part of a circle as we introduce in main text. For VAE, IAE, PIE, \mathcal{M} -flow and HF, the metric is computed as

$$\mathbb{E}_{x_1, x_2 \sim \mathcal{M}} \left| \frac{d_{\mathcal{M}}(x_1, x_2)}{\|g^{-1}(x_1) - g^{-1}(x_2)\|_2} - 1 \right|, \quad (27)$$

150 where for VAE and IAE we choose g^{-1} to be the encoder. Because GAN is incapable of inferring the embedding of a given
 151 sample, the metric for GAN is computed differently as

$$\mathbb{E}_{z_1, z_2 \sim \mathcal{Z}} \left| \frac{d_{\mathcal{M}}(g(z_1), g(z_2))}{\|z_1 - z_2\|_2} - 1 \right|, \quad (28)$$

152 where \mathcal{Z} is a latent space such that $g(z) \in \mathcal{M}, \forall z \in \mathcal{Z}$. In addition to the above metric, we also propose a different metric
 153 for evaluating the *distance preserving* property by measuring how far the Jacobian of g deviates from an orthonormal matrix.
 154 Specifically, the computation of the metric is given as

$$\mathbb{E}_{z \sim \mathcal{Z}} \left| J_g(z)^\top J_g(z) - I \right|, \quad (29)$$

155 where \mathcal{Z} is the latent space for ancestral sampling, $J_g(z)$ is the Jacobian of g at z , and I is an identity matrix. Given a matrix
 156 $M \in \mathbb{R}^{K \times K}$, we denote $|M| \triangleq \frac{1}{K^2} \sum_{i,j}^K |M_{ij}|$, where $|M_{ij}|$ is the absolute value of the (i, j) -th element of M . To evaluate
 157 the *rigorous projection* property of the model, we evaluate how far the projection point \tilde{x}' given by the model g deviates from
 158 the ground-truth \tilde{x} for an off-manifold sample x , the metric is computed as

$$\mathbb{E}_{x \sim \mathcal{A}^1} \|\tilde{x}' - \tilde{x}\|_2, \quad (30)$$

159 where we choose the ambient space $\mathcal{A}^1 = \{(r \cos \theta, r \sin \theta) | r \in [-0.5, 0.5], \theta \in [\frac{\pi}{6}, \frac{5\pi}{6}]\}$, and the ground-truth projection
 160 point \tilde{x} for x can be analytically given as introduce in Fig. 3 in main text. For \mathcal{M} -flow and our proposed HF, in addition to the
 161 above metric, we also propose a different metric for evaluating the *rigorous projection* property for the model by measuring
 162 how far the Jacobian of f_1 deviates from the conditions stated in Prop. 3, which is computed as

$$\underbrace{\mathbb{E}_{u \sim \mathcal{U}^1, v \sim \mathcal{V}^1} \frac{\sum_{i \neq j}^{K_1} |(S_1([u; v]))_{ij}|}{(K_1 - 1) \sum_{i=1}^{K_1} |(S_1([u; v]))_{ii}|}, \quad \textcircled{1}}_{\textcircled{2}} \quad \underbrace{\mathbb{E}_{u \sim \mathcal{U}^1, v \sim \mathcal{V}^1} \frac{1}{K_1 - K_0} \sum_{i=K_0+1}^{K_1} |(S_1([u; v]))_{ii} - 1|}_{\textcircled{2}} \quad (31)$$

163 where f_1 is the *flow* module of the generator g , $S_1([u; v]) \triangleq J_{f_1}([u; v])^\top J_{f_1}([u; v])$, and we use M_{ij} to denote the (i, j) -th
 164 element of a matrix M . The metric $\textcircled{1}$ is used to measure the orthogonality of $S_1([u; v])$, which is computed as the ratio of
 165 the averaged absolute value of non-diagonal elements of $S_1([u; v])$ to the averaged absolute value of the diagonal elements
 166 of $S_1([u; v])$, hence the lower the metric $\textcircled{1}$, the more orthogonal the Jacobian $J_{f_1}([u; v])$. The metric $\textcircled{2}$ is used to measure
 167 how far the last $K_1 - K_0$ diagonal elements of $S_1([u; v])$ deviates from 1. For the synthetic manifold we introduced in main
 168 text, we have $K_0 = K = 1, K_1 = D = 2$, and we use $\mathcal{V}^1 = [-0.5, 0.5]$. We also evaluate how far the projection distance
 169 $d_g(x, \tilde{x}')$ predicted by the model deviates from the ground-truth Euclidean distance $d_{\mathcal{E}}(x, \tilde{x}')$ for \mathcal{M} -flow and HF, which is
 170 computed as

$$\mathbb{E}_{x \sim \mathcal{A}} |d_g(x, \tilde{x}') - d_{\mathcal{E}}(x, \tilde{x}')|, \quad (32)$$

171 where $\tilde{x}' = f_1([u^1(x); 0])$ and $d_g(x, \tilde{x}') = |v^1(x)|$. To evaluate the *anomaly detection* performance, we propose to adopt
 172 the *False Positive Rate at 95% True Positive Rate* following [9, 5]. Specifically, we detect the out-of-distribution samples
 173 based on the sample-to-manifold distance given by the model instead of using a classifier. For the ambient space \mathcal{A}^1 we
 174 used in Eq. (30), we choose the in-distribution samples set $\mathcal{A}_{\text{in}}^1 \subset \mathcal{A}^1$ such that $d(x, \mathcal{M}) < 0.25, \forall x \in \mathcal{A}_{\text{in}}^1$, and choose
 175 the out-distribution samples set $\mathcal{A}_{\text{out}}^1 \subset \mathcal{A}^1$ such that $d(x, \mathcal{M}) \geq 0.25, \forall x \in \mathcal{A}_{\text{out}}^1$, where $d(x, \mathcal{M})$ is the ground-truth
 176 distance from x to the manifold \mathcal{M} . Given in-distribution samples $\mathcal{X}_{\text{in}} \triangleq \{x^{(i)} \in \mathcal{A}_{\text{in}}^1\}_{i=1}^N$ and out-distribution samples
 177 $\mathcal{X}_{\text{out}} \triangleq \{x^{(i)} \in \mathcal{A}_{\text{out}}^1\}_{i=1}^N$, the metric is computed as

$$\underbrace{\frac{|\{|v(x)| < v^* | x \in \mathcal{X}_{\text{out}}\}|}{N}}_{\text{False Positive Rate (FPR)}}, \quad \text{s.t. } \underbrace{\frac{|\{|v(x)| < v^* | x \in \mathcal{X}_{\text{in}}\}|}{N}}_{\text{True Positive Rate (TPR)}} = 0.95, \quad (33)$$

178 where we use $N = 1,000$ and $|\mathcal{X}|$ is the number of elements in the set \mathcal{X} , namely we consider x an in-distribution sample if
 179 the distance to the manifold (*i.e.*, $|v(x)|$) predicted by the model is smaller than v^* .

180 **Quantitative Results** We provide quantitative results in Tbl. 3. As we can see in the table, models including GAN, VAE,
 181 PIE and \mathcal{M} -flow achieve inferior performances in terms of both the global and local metric scores of the *distance preserving*
 182 property, meanwhile IAE and our HF achieve better performances. Because IAE can be regarded as a VAE augmented with
 183 an isometric regularization proposed in [4], the gain of the performance verifies the effectiveness of the regularization in
 184 terms of encouraging the decoder to satisfy the *distance preserving* property. In terms of the *rigorous projection* property, as
 185 we can observe, our HF achieves the best performance across all models. Therefore, we learn that apart from our proposed
 186 HF, none of other models are guaranteed to satisfy the property. We also observe that our HF achieves the best performance in
 187 terms of *anomaly detection*. From the computation method of the *anomaly detection* metric as in Eq. (33), we argue that the
 188 inaccuracy of the projection distance of \mathcal{M} -flow causes the *anomaly detection* performance to drop. To see this, for \mathcal{M} -flow,
 189 the predicted projection distance of an in-distribution sample may larger than that of an out-distribution sample, hence the
 190 distributions of the predicted projection distances of in-distribution and out-distribution samples may not be discriminative,
 191 which leads an inferior score of the *FPR at 95% TPR* metric for \mathcal{M} -flow. However, as our HF can always predict correct
 192 projection distance of an off-manifold sample, the distributions of the predicted projection distances of in-distribution and
 193 out-distribution samples are always discriminative, and hence the *FPR at 95% TPR* metric score of our HF is superior.
 194 Moreover, the inaccuracy of the predicted projection distance $\|v(x)\|_2$ for x also makes it inappropriate for \mathcal{M} -flow to
 195 implement the reconstruction loss by minimizing $\|v(x)\|_2$, as we done in our HF to speed up training.

196 2.2. Results on Natural Image Datasets

197 In this section, we provide qualitative and quantitative results on the aforementioned natural image datasets, which includes
 198 a quantitative demonstration on the time-efficiency for our proposed *two-stage dimensionality reduction* algorithm.

Encoder	FC Neural Network Models Decoder	NF-based Model Generator
Input 2-dimensional sample FC. 1024, ReLU FC. 1024, ReLU FC. 1	Input 1-dimensional embedding FC. 1024, ReLU FC. 1024, ReLU FC. 1024, ReLU FC. 2	Input 1-dimensional embedding Pad. 2, RT. 2, CP. 128, RT. 2 CP. 128, RT. 2

Table 1. The architectures we use for the synthetic manifold \mathcal{M} as we introduced in main text. For VAE and IAE, the encoder $h : \mathbb{R}^2 \rightarrow \mathbb{R}^1$ and the decoder $g : \mathbb{R}^1 \rightarrow \mathbb{R}^2$ are fully-connected neural network models, where we use “FC. n ” to represent a fully-connected layer with n output channels. For GAN, the generator and the discriminator use the architecture of the decoder and the encoder, respectively. For NF-based models, the generator $g = f_1 \circ p_1 : \mathbb{R}^1 \rightarrow \mathbb{R}^2$ consists of a *padding* module $p_1 : \mathbb{R}^1 \rightarrow \mathbb{R}^2$ and a *flow* module $f_1 : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, where we use “Pad. n ” to denote a padding module with the output dimensionality being n , use “RT. n ” to denote a 1×1 convolutional layer with n channels, and use “CP. n ” to denote an *affine coupling* layer with its nonlinear module being a convolutional architecture of “conv3. n , ReLU, conv3. n , ReLU, conv3. l ”, where we use “conv3. n ” to denote a 3×3 convolution layer with output channel n , stride 1 and padding 1, and l is the number of channels of input to the “CP. n ” layer.

Encoder	Conv Neural Network Models Decoder	NF-based Model Generator
Input $1 \times 28 \times 28$ image $4 \times 4@(2, 1)$ conv. 32, ReLU $4 \times 4@(2, 1)$ conv. 64, ReLU $4 \times 4@(1, 0)$ conv. 128, ReLU $4 \times 4@(1, 0)$ conv. 256, ReLU $1 \times 1@(1, 0)$ conv. K	Input $\in \mathbb{R}^K$ $1 \times 1@(1, 0)$ conv. 256, ReLU $7 \times 7@(1, 0)$ upconv. 64, ReLU $4 \times 4@(2, 1)$ upconv. 32, ReLU $4 \times 4@(2, 1)$ upconv. 3	Input $\in \mathbb{R}^K$ Pad. $4 \times 7 \times 7$ RT. 196, Flow@512@6 Pad. $2 \times 14 \times 14$ RT. 8, Flow@256@6 Pad. $1 \times 28 \times 28$ RT. 4, Flow@128@6

Table 2. The architectures we use for MNIST image datasets. For VAE and IAE, the encoder $h : \mathbb{R}^D \rightarrow \mathbb{R}^K$ and the decoder $g : \mathbb{R}^K \rightarrow \mathbb{R}^D$ are convolutional neural network work models, and we use “ $k \times k@(s, p)$ conv. n ” to represent a convolution layer with output channel n , kernel size k , stride s and padding p . The representation of a transposed convolution layer (denoted as upconv.) is similar. For GAN, the generator and the discriminator adopt the architecture of the decoder and encoder, respectively. For NF-based models, we use “Pad. $c \times h \times w$ ” to denote a padding module with output shape being $c \times h \times w$, and use “Flow@ $n@r$ ” to denote a flow layer composed by r 1×1 convolutional layers and r *affine coupling* layers, where for each *affine coupling* layer, the nonlinear module is a convolutional architecture of “conv3. n , ReLU, conv3. n , ReLU, conv3. l ” with l being the number of channels of the input to the *affine coupling* layer. We use “RT. n ” to denote a 1×1 convolutional layer with n channels, and use $K = 10$.

199 **Experimental Settings** We introduce how we evaluate different models on the natural image datasets. In order to evaluate
 200 the *distance preserving* property, since the ground-truth geodesic length $d_{\mathcal{M}}(x_1, x_2)$ between two manifold samples x_1, x_2
 201 is unknown, we only compute the metric score of Eq. (29) for all models. Similarly, since the ground-truth projection point
 202 \tilde{x} of an off-manifold x is unknown, the computation of Eq. (30) is no longer applicable for evaluating the *rigorous projection*
 203 property, hence we only compute the metric score of Eq. (31) and Eq. (32) for \mathcal{M} -flow and our proposed HF. For our HF with
 204 a multi-layer generator, the metric score of Eq. (31) is the average of the computed values of Eq. (31) for all *flow* modules
 205 f_i of the generator. To evaluate the *anomaly detection* performance, the computation of Eq. (33) is fundamentally applicable
 206 but the choices of \mathcal{X}_{in} and \mathcal{X}_{out} are different. Specifically, given a training dataset \mathcal{X} , we use \mathcal{X} itself as the in-distribution
 207 dataset \mathcal{X}_{in} , and use $\mathcal{X}_{\text{out}} \triangleq \{x + \varepsilon | x \in \mathcal{X}\}$, where ε has the same dimensionality as x and each pixel of ε is a noise randomly
 208 sampled from $[-0.1, 0.1]$. We also provide qualitative results for intuitive comparison between different models.

209 **Qualitative Results** We provide qualitative results to visualize the quality and diversity of samples generated by using our
 210 HF model. In addition to train our HF model by using the objectives proposed in main text, we also adapt our HF model to an
 211 adversarial setting by replacing the reconstruction loss (*i.e.*, \mathcal{L}_v in main text) with an adversarial loss, where a StyleGAN2 [6]

Algorithm 1 Manifold Training

Require: Training dataset $\mathcal{X} = \{x^{(i)}\}_{i=1}^N \subset \mathbb{R}^D$, embedding dimensionality K , generator $g : \mathbb{R}^K \rightarrow \mathbb{R}^D$ of L layers, singular values predictor $\{s_i : \mathbb{R}^{K_i} \rightarrow \mathbb{R}^{K_i}\}_{i=1}^L$, ambient radius r , hyper-parameters λ_{dist} , λ_{proj} and λ_v , all modules parameters Θ , optimizer ρ

- 1: Initialize running statistics $\mu_j \leftarrow 0, \sigma_j^2 \leftarrow 1, j \in [K]$
- 2: **while** not converge **do**
- 3: Randomly select a batch of training samples $\mathcal{B} \subset \mathcal{X}$
- 4: Obtain $u(x)$ and $\{v^i(x)\}_{i=1}^L$ by running inference for $x \in \mathcal{B}$ using g^\dagger . Compute \mathcal{L}_v for \mathcal{B}
- 5: Randomly sample a batch $\mathcal{S}^1 = \{u \sim \mathcal{N}(0, \sigma^2)\}$ and a batch $\mathcal{T}^i = \{v \in \mathcal{V}^i(r)\}$, then obtain

$$\mathcal{R}^i = \{f_{i-1}([u; v]) | u \in \mathcal{R}^{i-1}, v \in \mathcal{T}^i\} \quad (34)$$

for $2 \leq i \leq L$ based on $\mathcal{R}^0 \triangleq \mathcal{S}^1$. Obtain $\mathcal{Z}^i \triangleq \mathcal{R}^{i-1} \times \mathcal{T}^i$. Compute $\mathcal{L}_{\text{dist}}$ and $\mathcal{L}_{\text{proj}}$ based on $\mathcal{Z}^i, i \in [L]$

- 6: Compute $\mathcal{L} \leftarrow \lambda_{\text{dist}}\mathcal{L}_{\text{dist}} + \lambda_{\text{proj}}\mathcal{L}_{\text{proj}} + \lambda_v\mathcal{L}_v$
 - 7: Update $\Theta \leftarrow \rho\nabla\mathcal{L}$
 - 8: Update $\mu_j \leftarrow 0.99\mu_j + 0.01\text{avg}\{u_j(x) | x \in \mathcal{B}\}$
 - 9: Update $\sigma_j^2 \leftarrow 0.99\sigma_j^2 + 0.01\text{var}\{u_j(x) | x \in \mathcal{B}\}$
 - 10: **end while**
-

Algorithm 2 The Overall Training Algorithm

Require: Training dataset $\mathcal{X} = \{x^{(i)}\}_{i=1}^N \subset \mathbb{R}^D$, dimensionality K of the first stage, threshold t
/* Manifold Training Phase */

- 1: Perform the *two-stage dimensionality reduction* algorithm as we introduced in Alg. 1 in main text, obtaining the detected manifold dimensionality \tilde{K} and the generator $\tilde{g} : \mathbb{R}^{\tilde{K}} \rightarrow \mathbb{R}^D$. For each stage, the generator is trained by using Alg. 1
 - 2: Obtain the embeddings $\mathcal{U} = \{\tilde{g}^\dagger(x^{(i)})\}_{i=1}^N$ using \tilde{g}
/* Density Training Phase */
 - 3: Initialize a density estimator $h : \mathbb{R}^{\tilde{K}} \rightarrow \mathbb{R}^D$, obtaining the parameters Θ of h . Then, setup an optimizer ρ for Θ
 - 4: **while** not converge **do**
 - 5: Randomly select a batch of training samples $\mathcal{B} \subset \mathcal{U}$
 - 6: Compute \mathcal{L}_d for \mathcal{B}
 - 7: Update $\Theta \leftarrow \rho\nabla\mathcal{L}_d$
 - 8: **end while**
-

discriminator is involved. Our HF model in the adversarial setting is essentially a GAN framework augmented with the two isometric regularizations (*i.e.*, Eq. 16 and Eq. 17) we proposed in main text. As we can observe from Fig. 3, the adversarial setting results in more realistic, sharp and diverse samples, meanwhile the samples generated in the normal setting tends to be blurry. In Fig. 4, we also intuitively compare \mathcal{M} -flow with our HF in terms of the *rigorous projection* property by visualizing how the two models performs regarding generating samples extending to the ambient space. Specifically, for \mathcal{M} -flow, for a constant magnitude of change in v , the generated samples $f_1([u; v])$ should also change constantly if the model satisfies the *rigorous projection* property, as we analyzed in Sec. 2.1. We observe from Fig. 4 that our HF results in more constant change in generated samples in ambient space compared with \mathcal{M} -flow, and also generates samples with higher quality than \mathcal{M} -flow, which intuitively shows the effect of our isometric regularizations. For both models in the adversarial setting, we only show qualitative results, and for all quantitative evaluations below, we use the normal setting for both models.

Quantitative Results We provide quantitative results in Tbl. 3 to compare different models in terms of a series of metrics. We observe that the behaviors of different models are similar to those in the case of the aforementioned synthetic manifold. Specifically, we see that for natural image datasets, our proposed isometric regularizations Eq. 16 and Eq. 17 in main text are still effective in constraining the model to satisfy the properties of *distance preserving* and *rigorous projection*, respectively, and the performance of our HF on *anomaly detection* is superior due to the satisfaction of the *rigorous projection* property.

Model	Synthetic Manifold \mathcal{M}						CelebA					
	GAN	VAE	IAE	PIE	\mathcal{M} -flow	HF	GAN	VAE	IAE	PIE	\mathcal{M} -flow	HF
<i>Distance Pres.</i> ↓	0.76	0.58	0.14	0.34	0.31	0.09	N/A	N/A	N/A	N/A	N/A	N/A
	1.38	0.76	0.17	0.52	0.48	0.14	10.23	9.56	3.28	9.45	4.56	0.96
<i>Rigorous Proj.</i> ↓	N/A	0.45	0.38	0.42	0.46	0.15	N/A	N/A	N/A	N/A	N/A	N/A
	N/A	N/A	N/A	N/A	1.29	0.12	N/A	N/A	N/A	N/A	4.79	1.03
<i>Anomaly Dete.</i> ↓	N/A	N/A	N/A	N/A	0.36	0.10	N/A	N/A	N/A	N/A	0.56	0.11
	N/A	N/A	N/A	N/A	1.35	0.13	N/A	N/A	N/A	N/A	1.58	0.15
<i>Anomaly Dete.</i> ↓	N/A	0.57	0.33	0.35	0.38	0.11	N/A	0.51	0.47	0.42	0.38	0.15
Model	FFHQ						AFHQ Cat					
	GAN	VAE	IAE	PIE	\mathcal{M} -flow	HF	GAN	VAE	IAE	PIE	\mathcal{M} -flow	HF
<i>Distance Pres.</i> ↓	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
	9.85	10.52	4.73	9.86	5.37	1.02	8.42	8.69	4.61	7.52	6.24	0.94
<i>Rigorous Proj.</i> ↓	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
	N/A	N/A	N/A	N/A	5.31	1.02	N/A	N/A	N/A	N/A	4.69	1.04
<i>Anomaly Dete.</i> ↓	N/A	N/A	N/A	N/A	0.52	0.13	N/A	N/A	N/A	N/A	0.49	0.10
	N/A	N/A	N/A	N/A	1.69	0.14	N/A	N/A	N/A	N/A	1.57	0.14
<i>Anomaly Dete.</i> ↓	N/A	0.54	0.48	0.43	0.41	0.14	N/A	0.49	0.51	0.41	0.42	0.12

Table 3. Quantitative comparison between different models on different datasets in terms of various metrics. In terms of evaluating the *distance preserving* property, the first rows and the second rows are the global and local metric scores, respectively, where the global metric scores are computed using Eq. (27) or Eq. (28), and the local metric scores are computed using Eq. (29). In terms of evaluating the *rigorous projection* property, the metric scores of the first and second rows are computed using Eq. (30) and Eq. (32), respectively, and the third and forth rows are computed using ① and ② of Eq. (31), respectively. In terms of *anomaly detection*, the metric scores are computed using Eq. (33). For each metric, ↓ indicates that lower values are better, and the best results across all models are shown in bold.

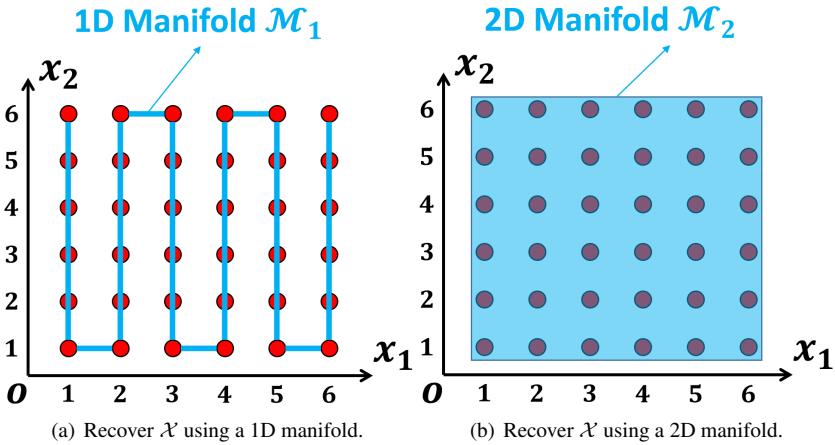


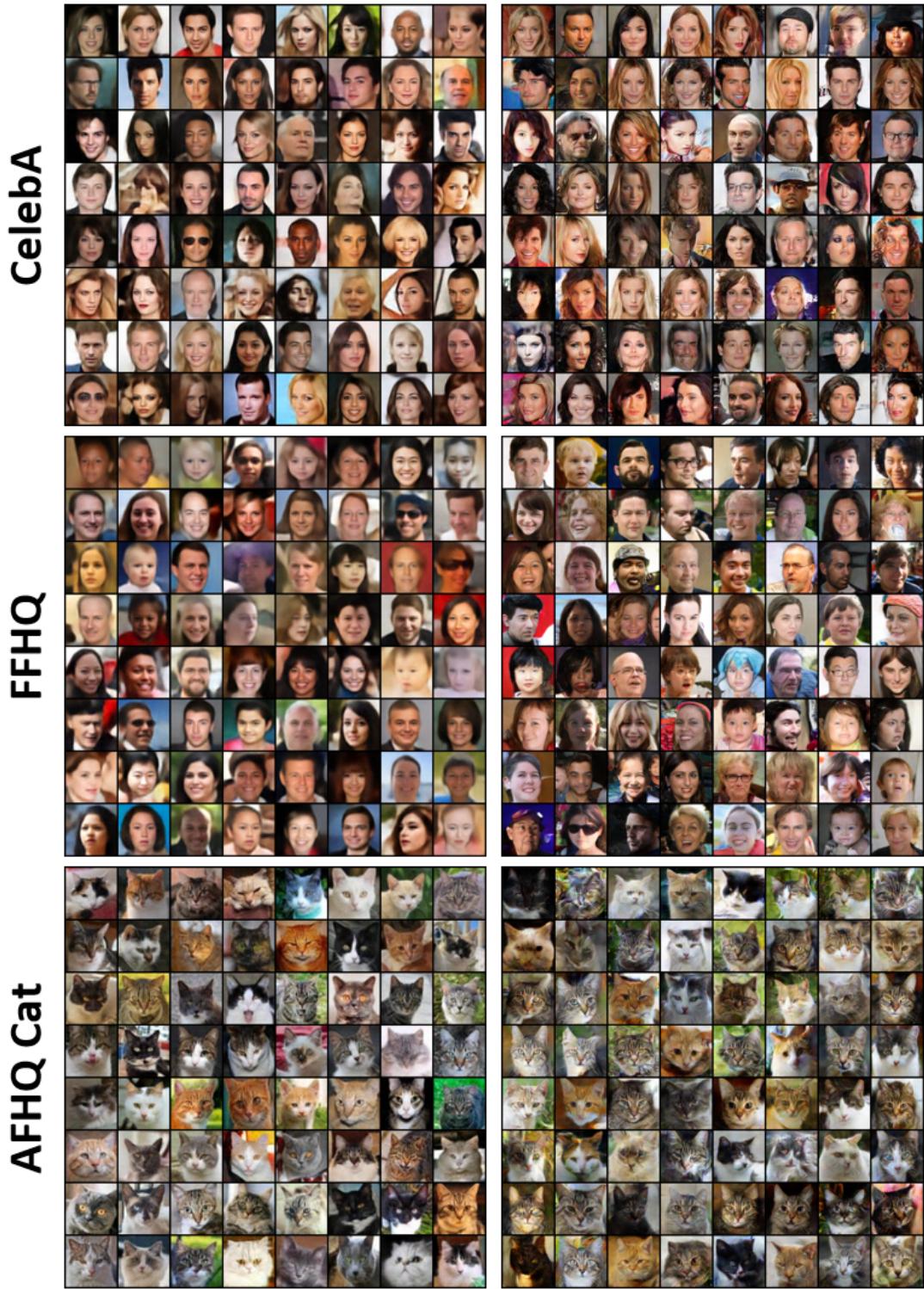
Figure 2. Given samples $\mathcal{X} \triangleq \{(i, j) | i, j \in [6]\}$, we demonstrate that \mathcal{X} can be recovered by manifolds of arbitrary dimensions. We plot samples $\mathcal{X} \triangleq \{(i, j) | i, j \in [6]\}$ by using red points, and the recovered manifolds are shown in blue.

227 **Dimensionality Reduction** We demonstrate that our *two-stage dimensionality reduction* algorithm is an efficient and effective
 228 manner for detecting the ground-truth dimensionality of the data manifold. We firstly provide an intuitive understanding
 229 on how our proposed *distance preserving* property facilitates detecting the manifold dimensionality by using Fig. 2. Given
 230 training samples $\mathcal{X} \triangleq \{(i, j) | i, j \in [6]\}$ (see Fig. 2), we can recover \mathcal{X} by using manifolds of arbitrary dimensions, as long as
 231 the model expressive power is unlimited. For example, we can recover \mathcal{X} by using a 1D manifold \mathcal{M}_1 as shown in Fig. 2(a),
 232 and we can also recover \mathcal{X} by using a 2D manifold \mathcal{M}_2 as we show in Fig. 2(b). Therefore, the criterion for dimensionality
 233 reduction based solely on the reconstruction error could be meaningless if the model expressive power is strong enough. We

Dataset	Search Process of \tilde{K}	Time Cost
CelebA	128 → 127 → ⋯ → 7 → 6 → 5	$124 \times 4.5\text{h} = 558\text{h}$
	128 → 64 → 32 → 16 → 8 → 4 → 6 → 5	$8 \times 4.5\text{h} = 36\text{h}$
FFHQ	128 → 64 → 32 → 16 → 8 → 4 → 6 → 5	$4.5\text{h} + 7 \times 0.5\text{h} = 8\text{h}$
	128 → 127 → ⋯ → 9 → 8 → 7	$122 \times 4.5\text{h} = 549\text{h}$
AFHQ Cat	128 → 64 → 32 → 16 → 8 → 4 → 6 → 7	$8 \times 4.5\text{h} = 36\text{h}$
	128 → 64 → 32 → 16 → 8 → 4 → 6 → 7	$4.5\text{h} + 7 \times 0.5\text{h} = 8\text{h}$
	128 → 127 → ⋯ → 5 → 4 → 3	$126 \times 4.5\text{h} = 567\text{h}$
AFHQ Cat	128 → 64 → 32 → 16 → 8 → 4 → 2 → 3	$8 \times 4.5\text{h} = 36\text{h}$
	128 → 64 → 32 → 16 → 8 → 4 → 2 → 3	$4.5\text{h} + 7 \times 0.5\text{h} = 8\text{h}$

Table 4. Comparison between different strategies for searching the minimum dimensionality \tilde{K} using different models. The criterion used for determining the ground-truth dimensionality is not only based on a given reconstruction threshold t as we introduce in Alg. 1 in main text, but also based on the ratio $r \triangleq \frac{d_u(\tilde{K}_t)}{d_u(\tilde{K}_{t-1})}$ with $d_u(\tilde{K}_t)$ and $d_u(\tilde{K}_{t-1})$ corresponding to the latent dimensionalities of the current step t and the previous step $t - 1$, respectively. Assume that the ground-truth manifold dimensionalities of CelebA, FFHQ and AFHQ Cat are 6, 8 and 4, respectively, and let the first-stage dimensionality K be 128. For each dataset, the first and second rows correspond to the search process using the brute-force algorithm and the binary search algorithm based on a \mathcal{M} -flow model, respectively, and the third row corresponds to the search process using the *two-stage dimensionality reduction* algorithm based on our HF model.

now show intuitively how our proposed *distance preserving* property helps alleviate this problem by using Fig. 2. We know that the ground-truth dimensionality of the manifold corresponding to \mathcal{X} is 2. Let K be the dimensionality of the manifold to recover \mathcal{X} induced by our model. For $K = 2$, assume that the model satisfies the *distance preserving* property, we know that $\max \{\|u(x_1) - u(x_2)\|_2 | x_1, x_2 \in \mathcal{X}\} = \sqrt{6^2 + 6^2} = 6\sqrt{2}$. However, for $K = 1$, given that the model satisfies the *distance preserving* property, we know that $\max \{\|u(x_1) - u(x_2)\|_2 | x_1, x_2 \in \mathcal{X}\} = 6 \times 7 = 42$, which is greatly larger than that of the case $K = 2$. Hence the metric $d_u \triangleq \max \{\|u(x_1) - u(x_2)\|_2 | x_1, x_2 \in \mathcal{X}\}$ increases dramatically if K is smaller than the ground-truth. However, for model such as \mathcal{M} -flow that is not guaranteed to satisfy the *distance preserving* property, the above analysis may not hold. Therefore, our proposed *distance preserving* property facilitates detecting the ground-truth dimensionality of the manifold. In terms of the time-efficiency of our proposed *two-stage dimensionality reduction* algorithm, in Tbl. 4, we compare different dimensionality reduction strategies, including the brute-force and binary search algorithms using \mathcal{M} -flow, and our proposed *two-stage dimensionality reduction* algorithm using our HF. We can learn that the time costs of the brute-force algorithm using \mathcal{M} -flow mentioned in [2] are totally unacceptable compared to the other two algorithms. Although \mathcal{M} -flow can speed up the search process by replacing the brute-force algorithm with the binary search algorithm, the one-layer design of \mathcal{M} -flow leads to training the generator $g : \mathbb{R}^{\tilde{K}} \rightarrow \mathbb{R}^D$ involving a high-dimensionality D repeatedly, which still results in a higher cost of search time compared with our *two-stage dimensionality reduction* algorithm. Thanks to the hierarchical architecture design of our HF, the generator g can be composed as $g = g_1 \circ g_2$, where $g_1 : \mathbb{R}^K \rightarrow \mathbb{R}^D$ and $g_2 : \mathbb{R}^{\tilde{K}} \rightarrow \mathbb{R}^K$ are the first-stage and second-stage generator, respectively. Though g_1 involves a high-dimensionality D , it is trained only once. Though g_2 is repeatedly trained during the search process, the dimensionalities \tilde{K} and K are small and hence the time costs are affordable. Therefore, combining the training of the first and second stage leads to a time-efficient algorithm for dimensionality reduction, as we demonstrated in Tbl. 4.



(a) Normal setting.

(b) Adversarial setting.

Figure 3. Visualization of samples generated by our HF model, where the model generating samples on the right is trained by using the adversarial setting, and the model generating samples on the left is trained by using the normal setting as introduced in main text.

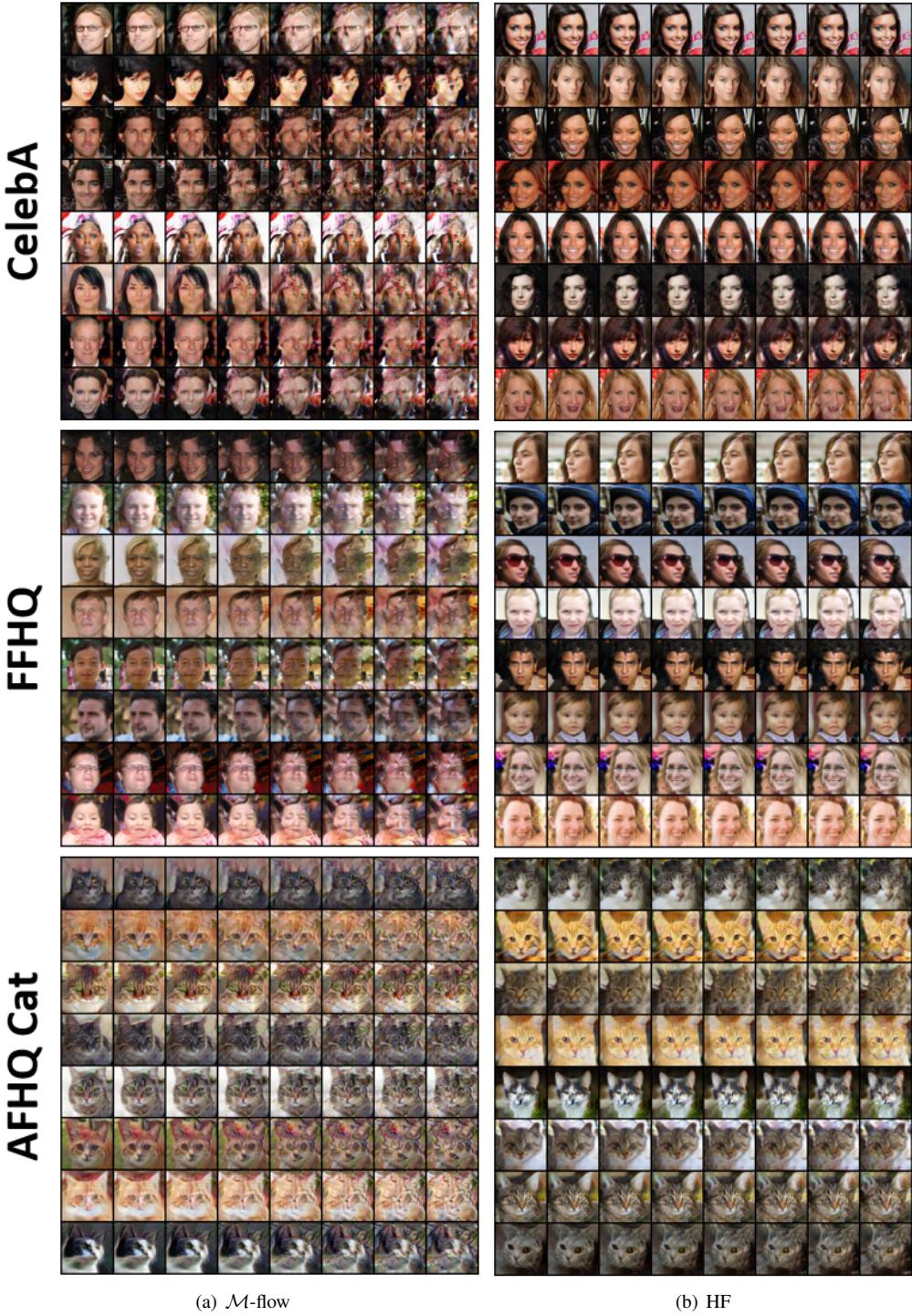


Figure 4. Qualitative comparison between \mathcal{M} -flow and our HF in terms of generating samples extending to the ambient space, where for each dataset, the samples on the left and right are generated using \mathcal{M} -flow and our HF, respectively, and we adopt the adversarial setting for both models. For each row of each model, the samples are generated as $f_1([u; v\varepsilon])$ (resp., u^{L+1} with $u^{i+1} \triangleq f_i([u^i; v\varepsilon])$ and $u^1 \triangleq u$) for \mathcal{M} -flow (resp., our HF with a L -layer generator), where u is shared across the row, ε is a random unit vector, $v \in [0, 0.5]$, and the leftmost and rightmost samples correspond to $v = 0$ and $v = 0.5$, respectively. For intuitive observation, the magnitude of the change in v between any two adjacent samples is constant.

254 **References**

- 255 [1] J. J. Beitler, I. Sosnovik, and A. Smeulders. Pie: Pseudo-invertible encoder. *arXiv preprint arXiv:2111.00619*, 2021.
- 256 [2] J. Brehmer and K. Cranmer. Flows for simultaneous manifold learning and density estimation. *Advances in Neural Information*
257 *Processing Systems*, 33:442–453, 2020.
- 258 [3] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial
259 nets. *Advances in neural information processing systems*, 27, 2014.
- 260 [4] A. Gropp, M. Atzmon, and Y. Lipman. Isometric autoencoders. *arXiv preprint arXiv:2006.09289*, 2020.
- 261 [5] D. Hendrycks and K. Gimpel. A baseline for detecting misclassified and out-of-distribution examples in neural networks. *arXiv*
262 *preprint arXiv:1610.02136*, 2016.
- 263 [6] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila. Analyzing and improving the image quality of stylegan. In
264 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020.
- 265 [7] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- 266 [8] D. P. Kingma and M. Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- 267 [9] S. Liang, Y. Li, and R. Srikant. Enhancing the reliability of out-of-distribution image detection in neural networks. *arXiv preprint*
268 *arXiv:1706.02690*, 2017.
- 269 [10] Z. Pan, L. Niu, and L. Zhang. Unigan: reducing mode collapse using a uniform generator. In *NeurIPS*, 2022.
- 270 [11] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An
271 imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32:8026–8037, 2019.
- 272 [12] P. Petersen. *Riemannian geometry*, volume 171. Springer, 2006.