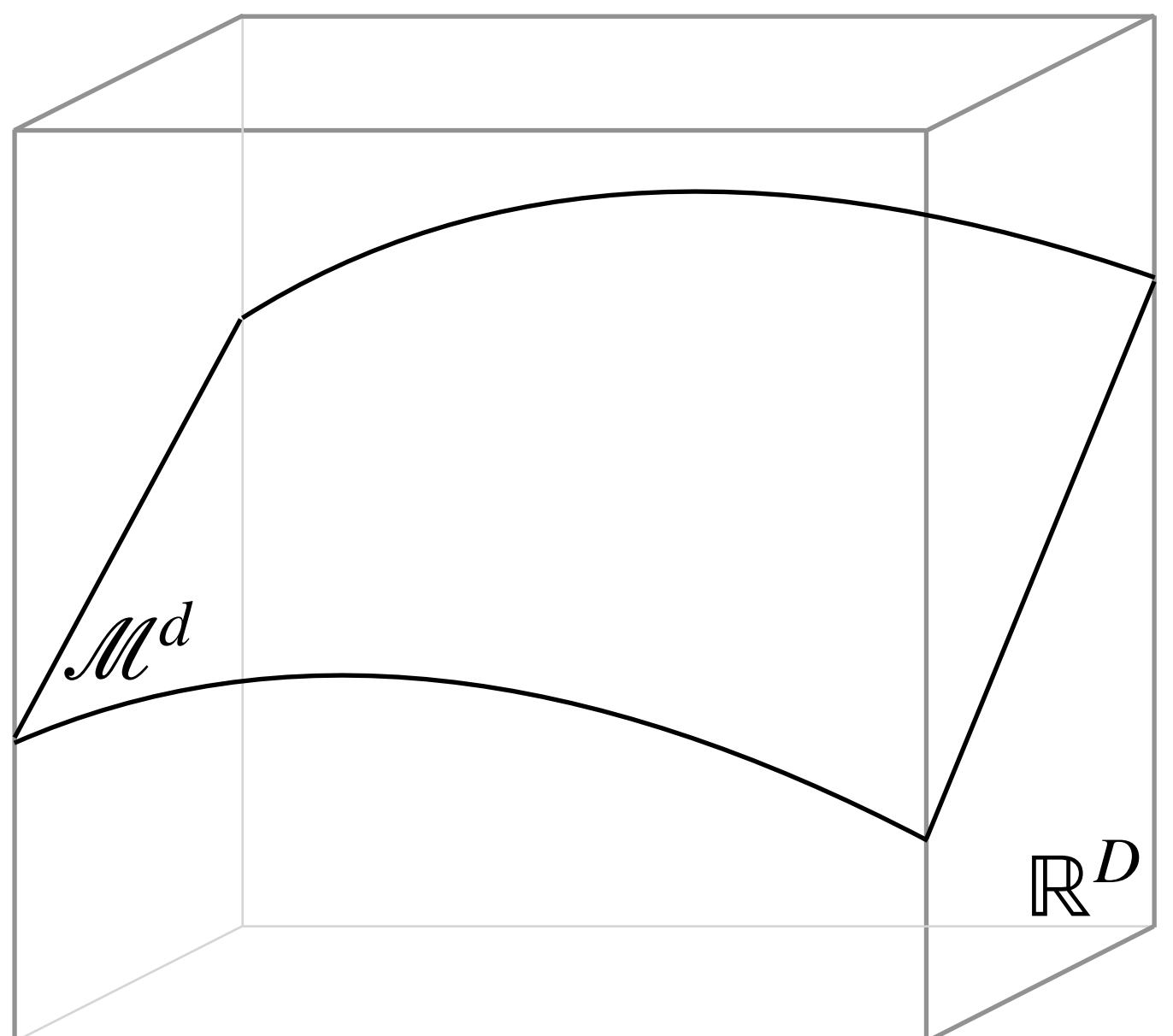


Topological Data Analysis

Lecture 12
Topological Data Analysis for Representation Learning

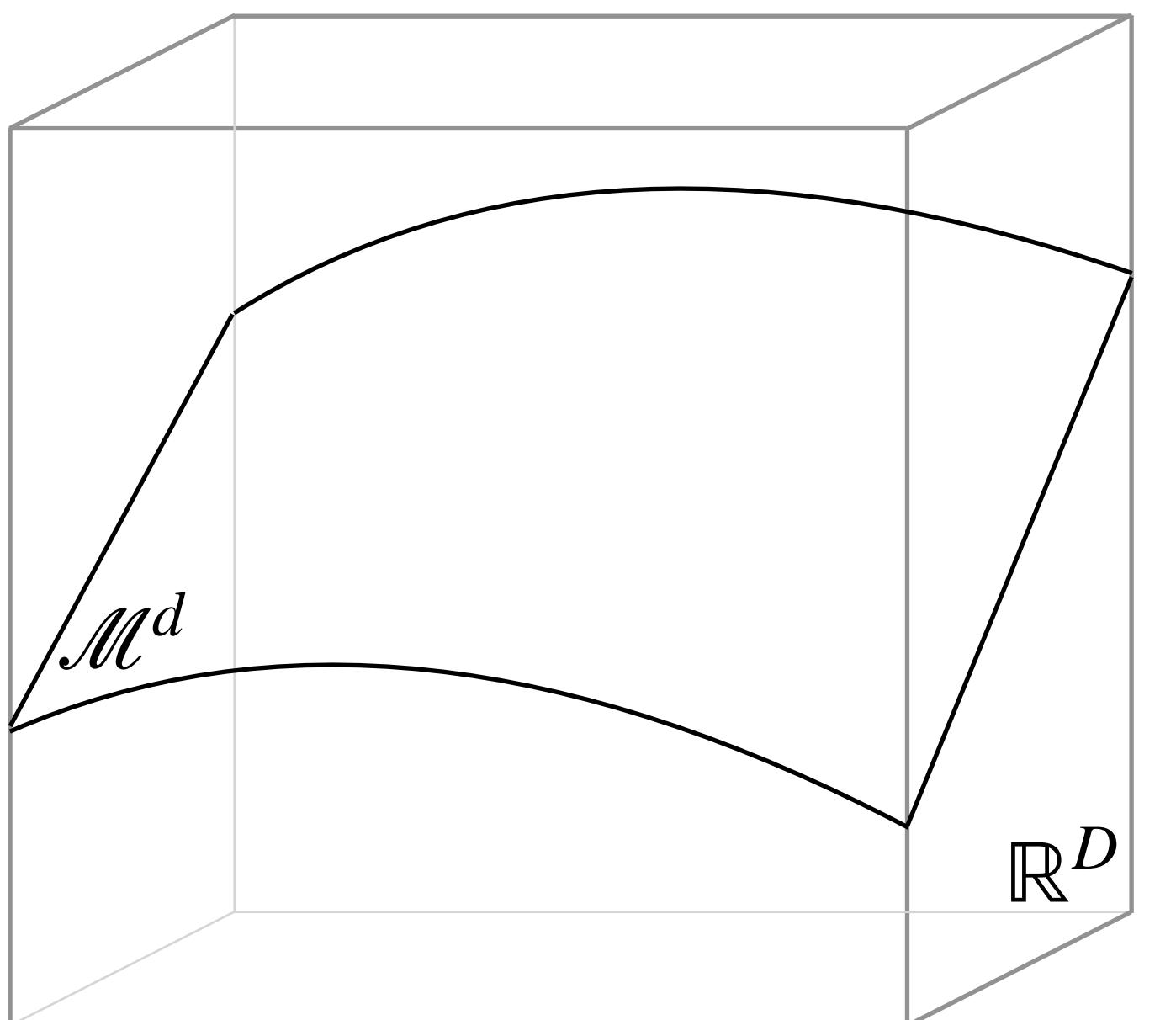
Oleg Kachan

Manifold hypothesis

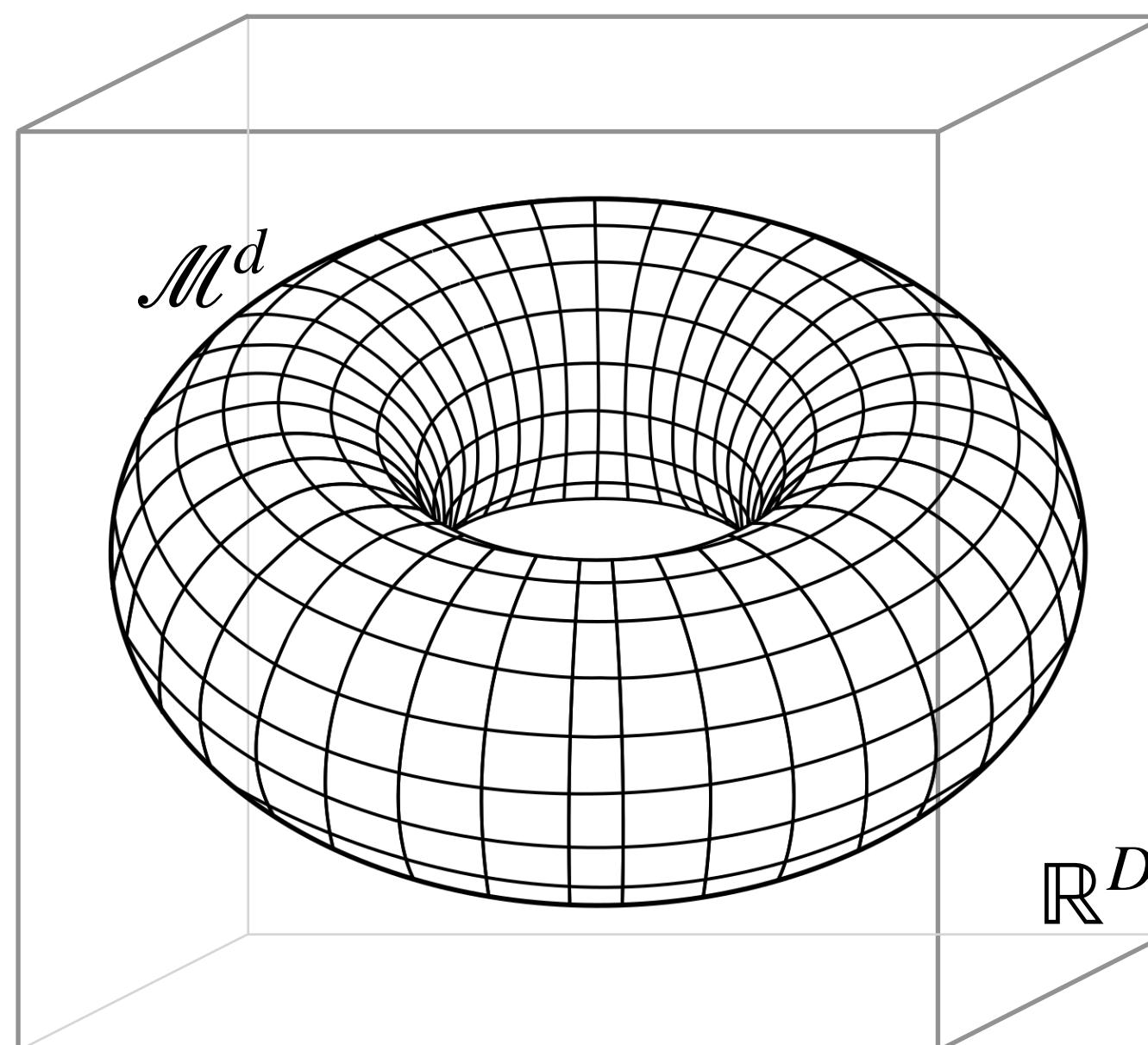


$$d \ll D$$

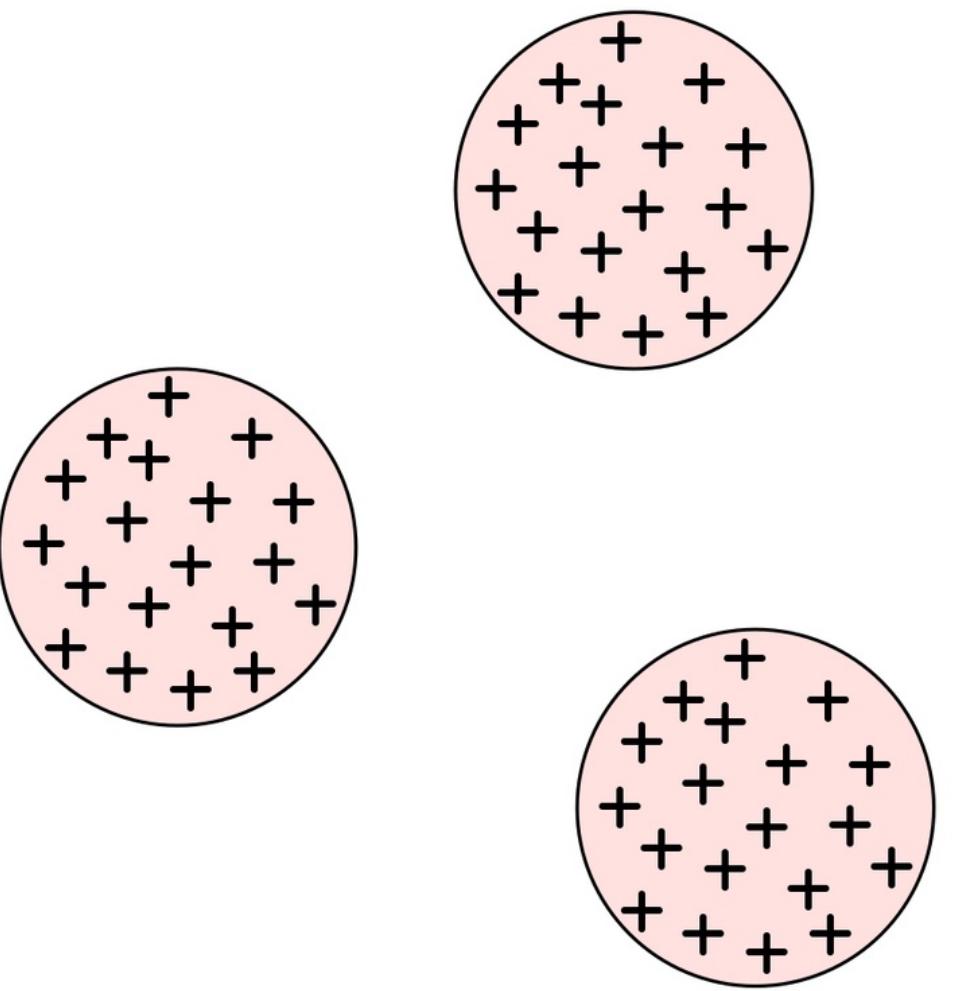
Manifold hypothesis



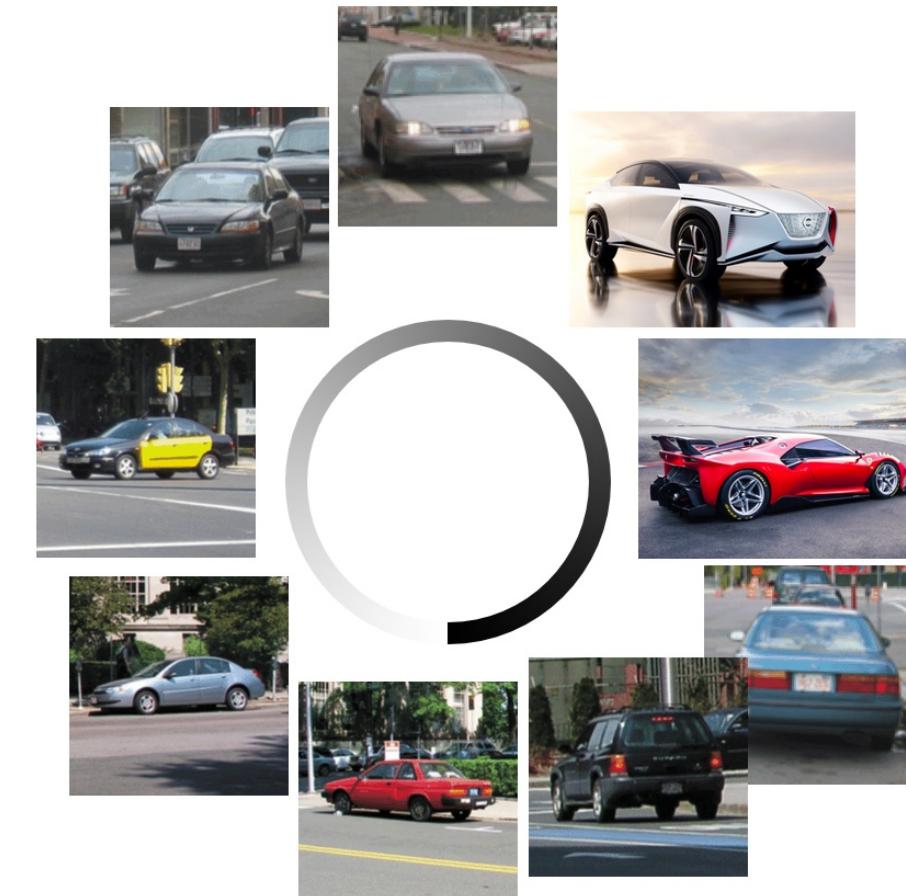
$$d \ll D$$



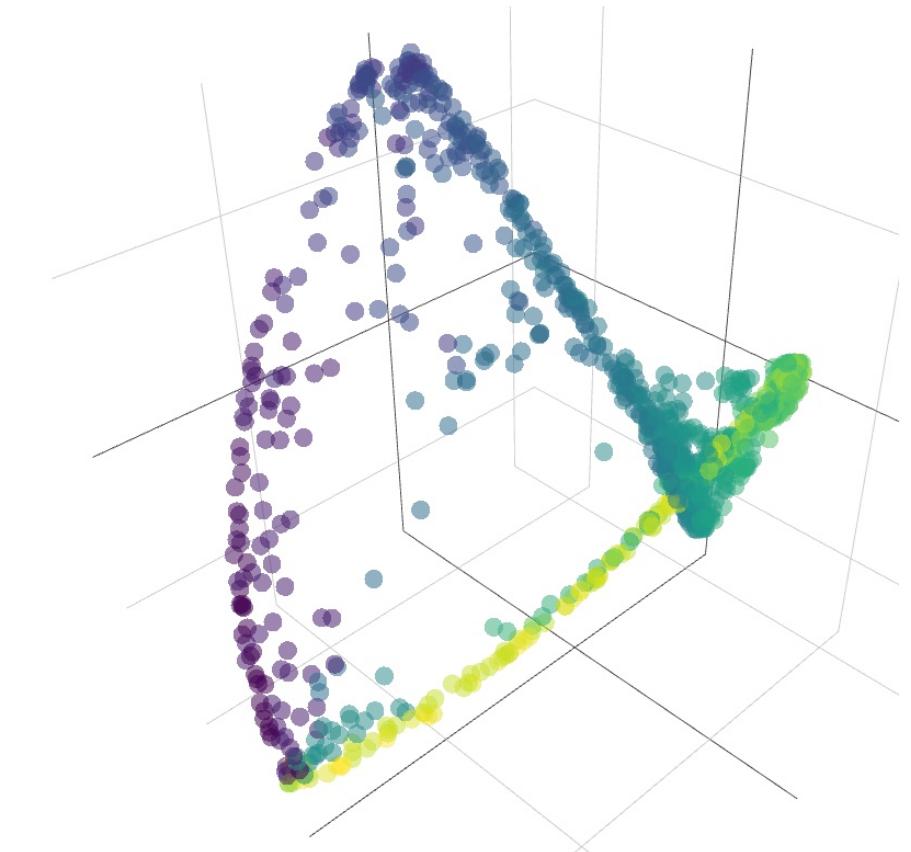
Manifold hypothesis



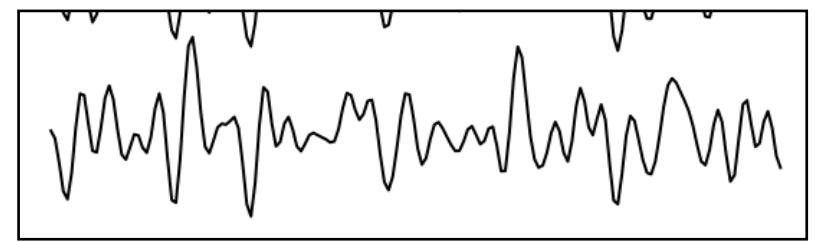
Connected components =
clusters



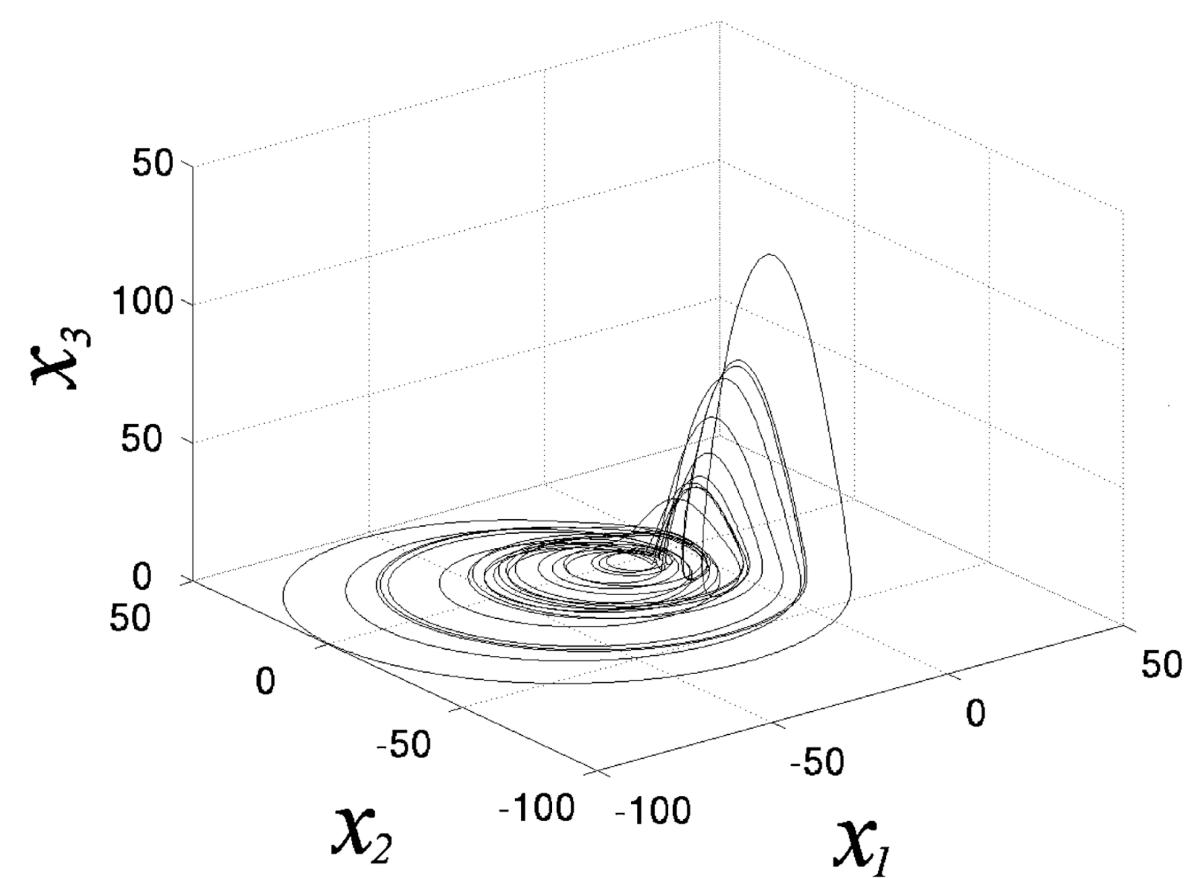
Cycles =
periodicity, regions of lower density



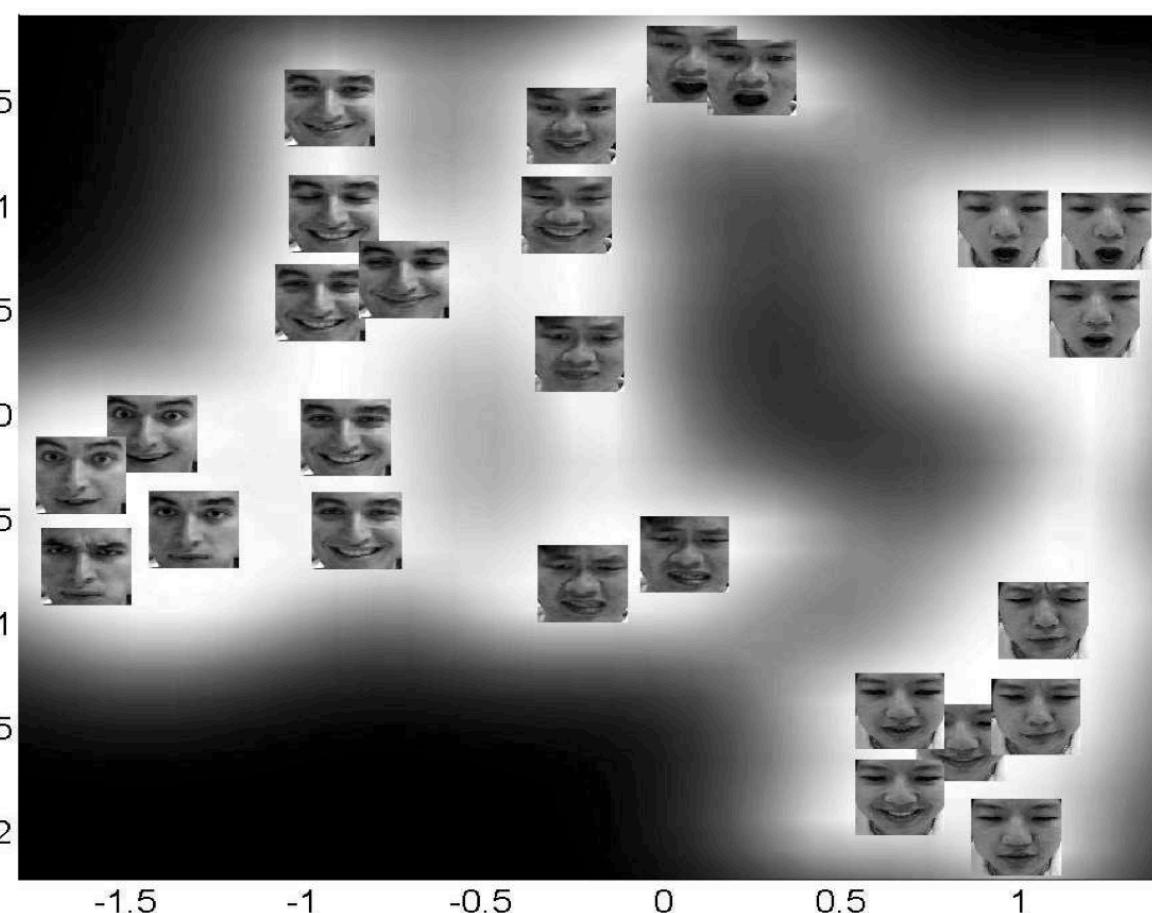
Manifold hypothesis



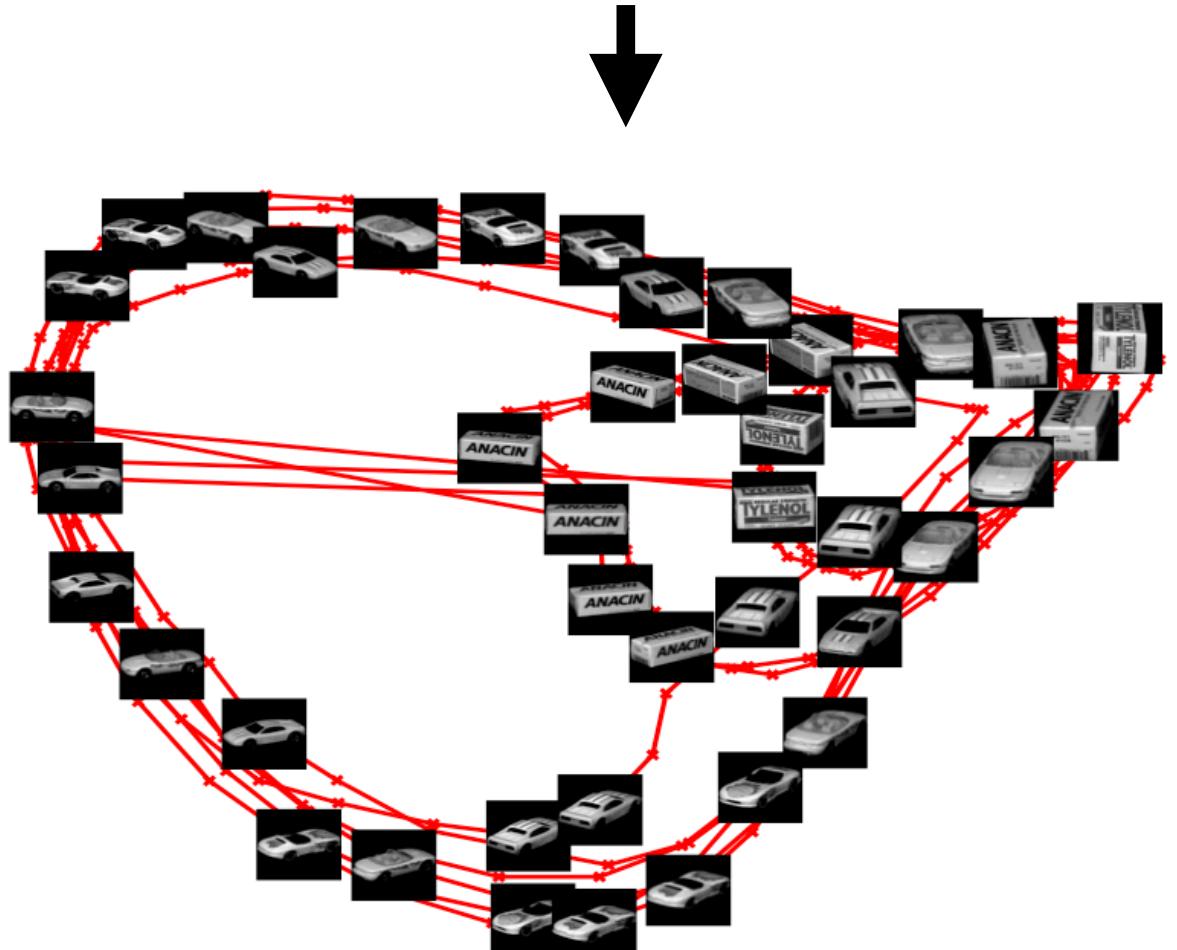
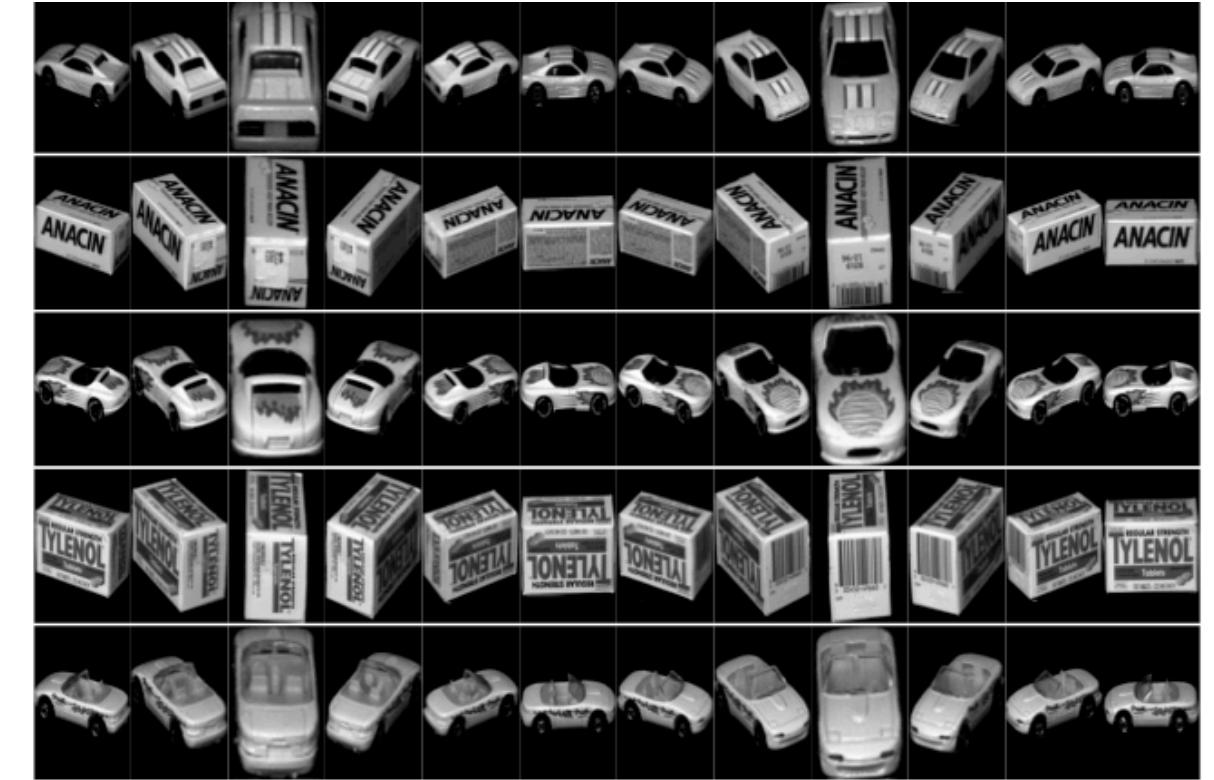
Time series



Delay embedding



COIL dataset



Dimensionality reduction

Given a dataset $X \in \mathbb{R}^m$ find a mapping $f: \mathcal{X} \in \mathbb{R}^D \rightarrow \mathcal{Z} \in \mathbb{R}^d$ where $d \ll D$, while optimizing/preserving some relevant properties of the data.

$$X \xrightarrow{f} Z$$

Properties

- variance/distances
- statistical independence

TDA for dimensionality reduction

Given a dataset $X \in \mathbb{R}^m$ find a mapping $f: \mathcal{X} \in \mathbb{R}^D \rightarrow \mathcal{Z} \in \mathbb{R}^d$ where $d \ll D$, while preserving topology of the input data.

$$X \xrightarrow{f} Z$$

Properties

$$Z \cong X$$

- variance/distances
- statistical independence
- *persistent homology groups*

TDA for representation learning

Given a dataset $X \in \mathbb{R}^m$ and mappings $f: \mathcal{X} \rightarrow \mathcal{Z}$ and $g: \mathcal{Z} \rightarrow \mathcal{Y}$
find a representation Z , close in topology to the input or output data.

$$X \xrightarrow{f} Z \xrightarrow{g} Y$$

$$Z \cong X \qquad Z \cong Y$$

Comparing data point clouds

Persistence diagram

- Wasserstein distance

Persistence diagram statistics

- Geometry Score (GS)

Cross-diagram

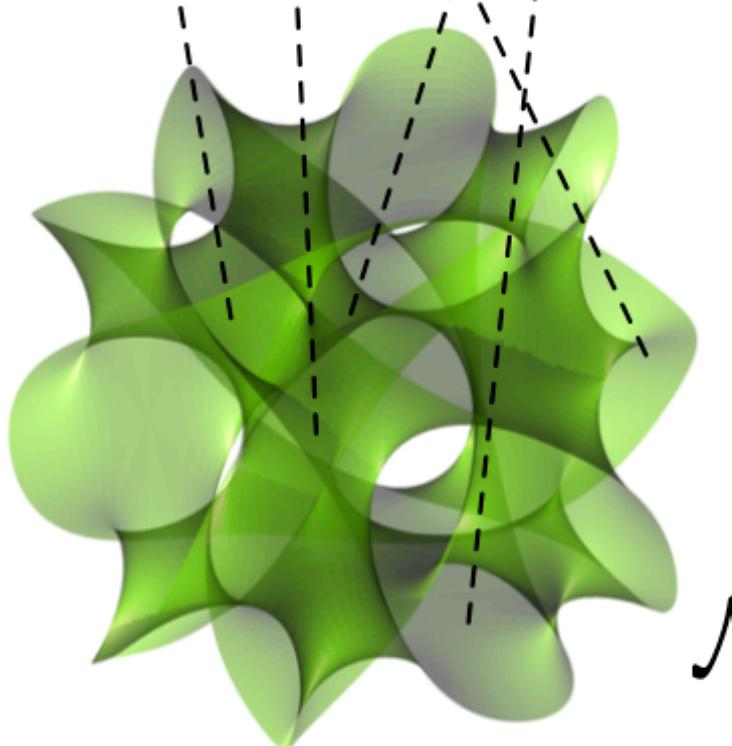
- Manifold topology divergence
- Representation topology divergence

$p_{\text{data}}(\mathbf{x})$

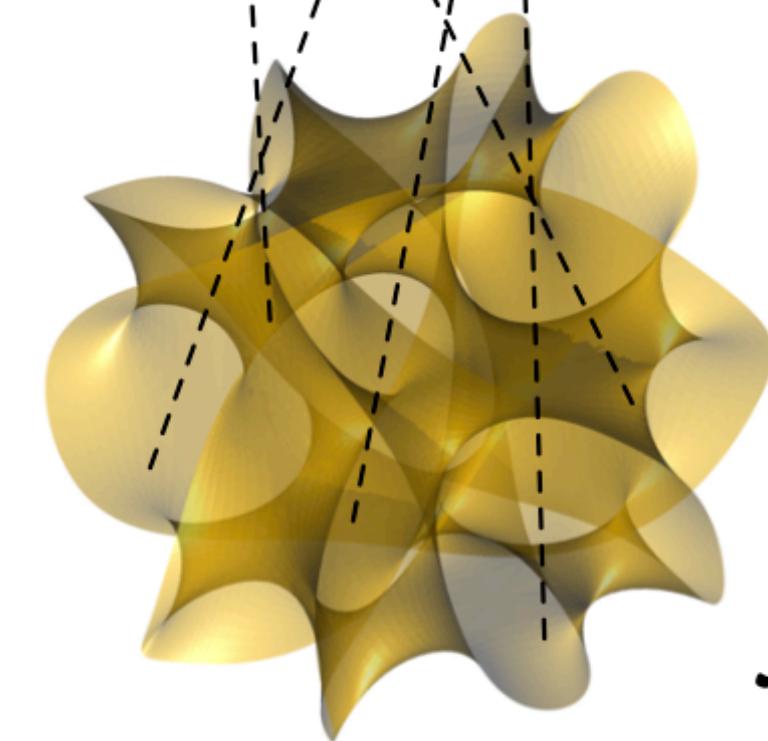
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8

$p_{\text{model}}(\mathbf{x})$

8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8
8	8	8	8	8	8	8	8

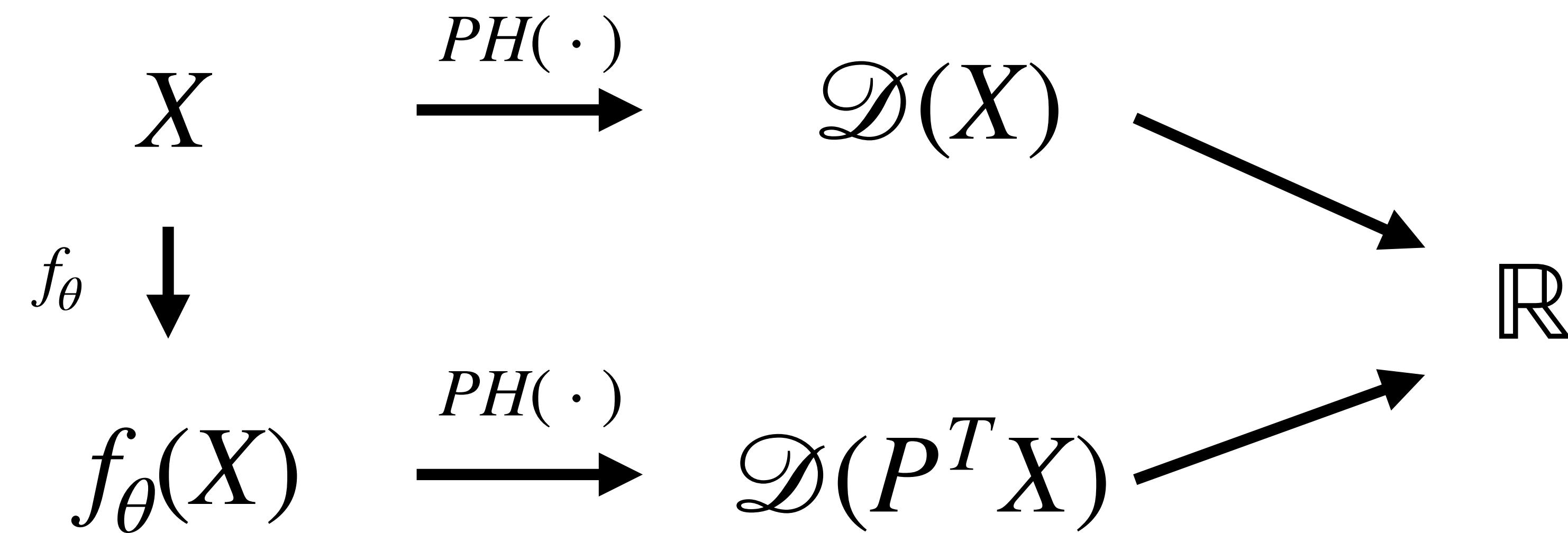


$\mathcal{M}_{\text{data}}$



$\mathcal{M}_{\text{model}}$

Comparing diagrams



Comparing data point clouds

Wasserstein distance

Given a ground metric $d : X \times X \rightarrow \mathbb{R}$ optimal transport equips the space of measures $\mathcal{P}(X)$ with a metric referred to as the p -th Wasserstein distance

$$W_p^p(\mu, \nu) = \min_{\mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{T}, \tilde{\mathbf{M}} \rangle_F$$

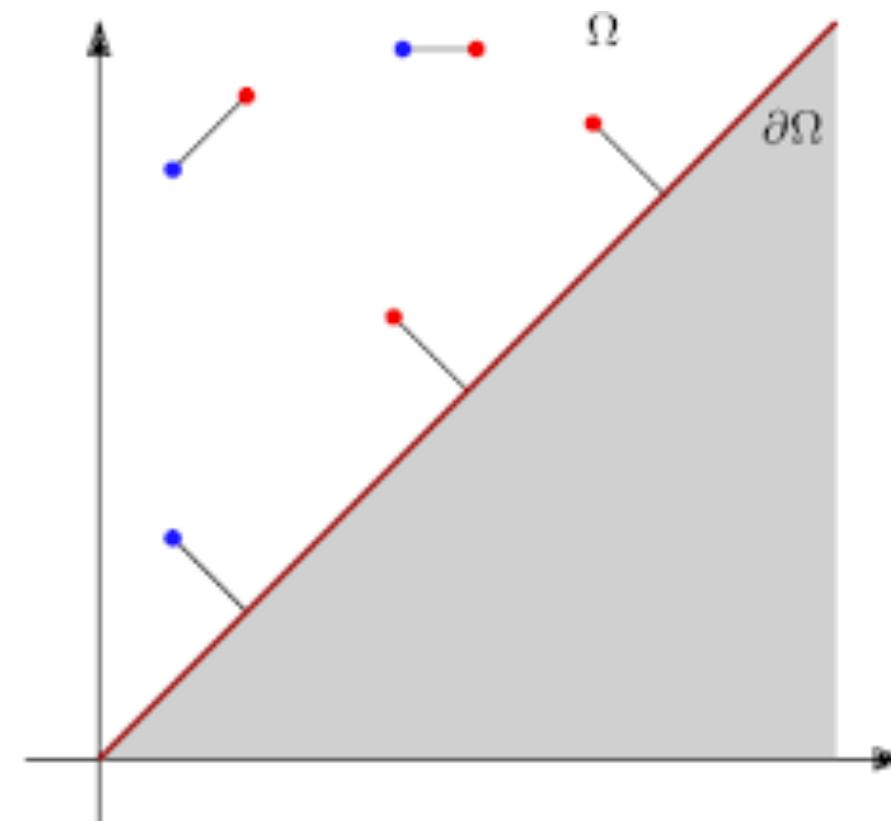
$$\Pi(\mathbf{a}, \mathbf{b}) = \{ \mathbf{T} \in \mathbb{R}_+^{n \times m} : \mathbf{T}\mathbf{1}_m = \mathbf{a}, \mathbf{T}^T\mathbf{1}_n = \mathbf{b} \}$$

Augmented cost matrix

$$\tilde{\mathbf{M}} = \begin{pmatrix} \mathbf{M} & \Delta_{D(X)} \\ \Delta_{D(Y)} & 0 \end{pmatrix}$$

$$(\mathbf{M})_{ij} = d^p(\mathbf{x}_i, \mathbf{y}_j) \quad (\Delta_{D(X)})_i = d^p(\mathbf{x}_i, \partial\Omega)$$

$$(\Delta_{D(Y)})_j = d^p(\mathbf{y}_j, \partial\Omega)$$



Approximations

- Sinkhorn divergence, $O(n^2)$
- Sliced Wasserstein distance, $O(n \log n)$

Comparing data point clouds

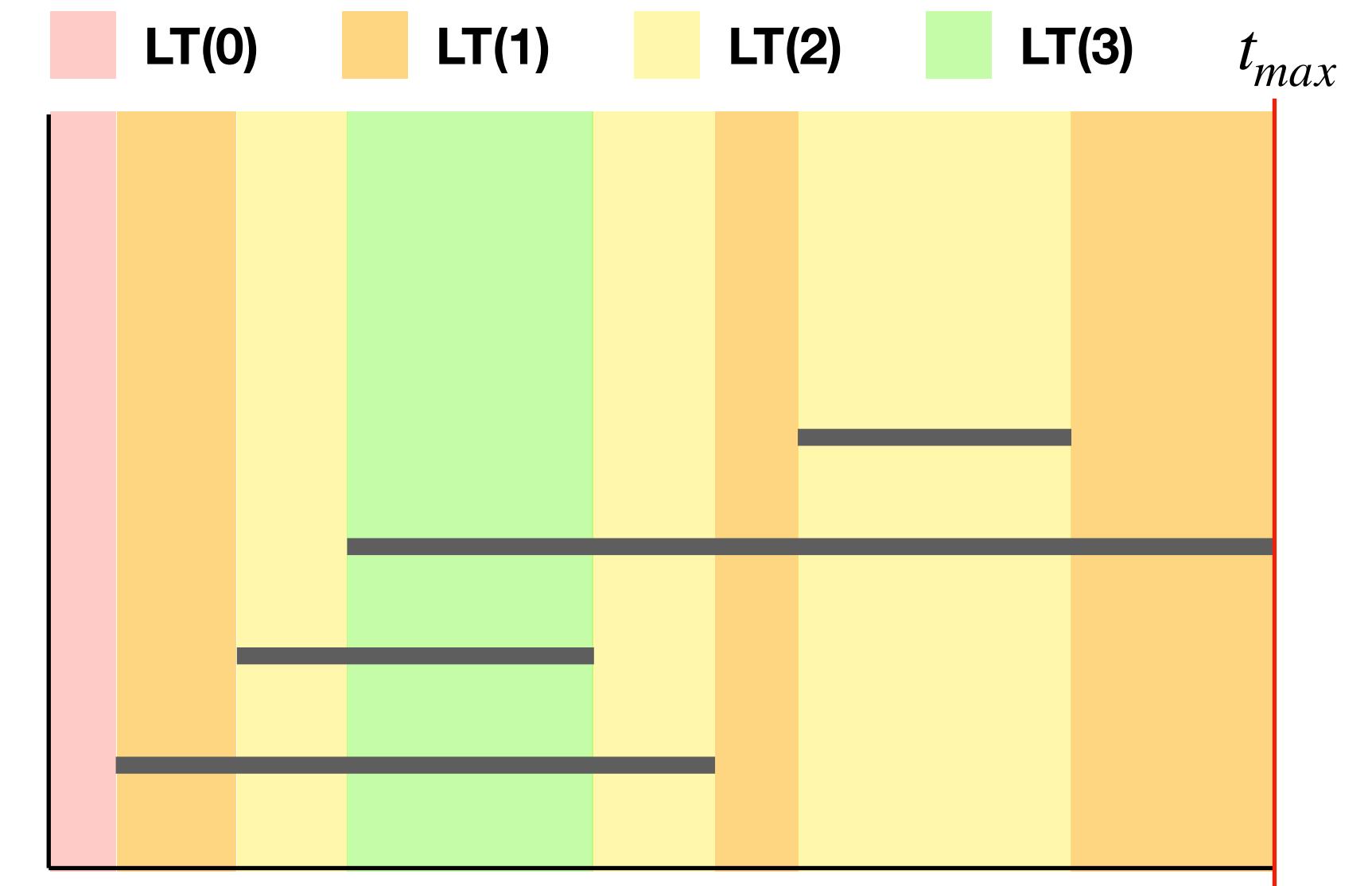
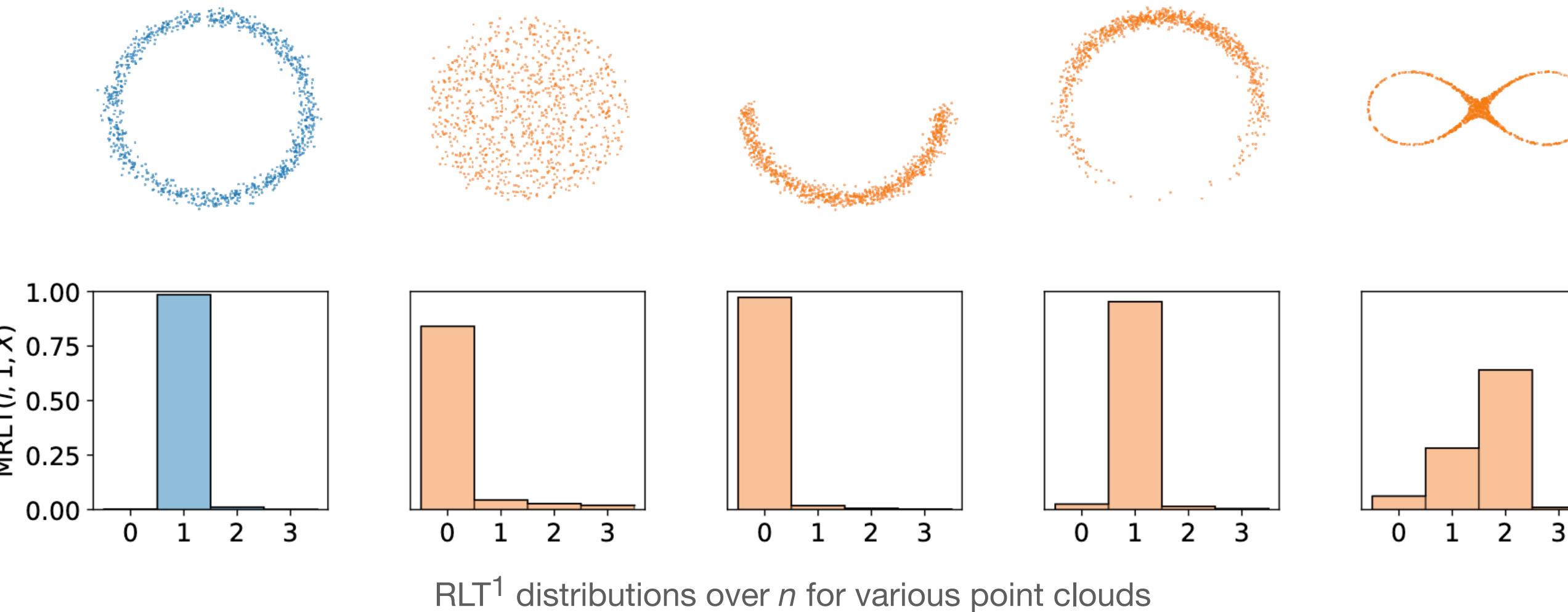
Geometry score

Given a k -th persistence diagram $D^k(X)$ the relative living time (RLT) is defined

$$\text{RLT}_X^k(n) = \frac{\mu\{t \in [0, t_{max}] \mid \beta_k(t) = n\}}{t_{max}}, \quad n \in \mathbb{Z}_{\geq 0}$$

Minimum enclosing radius

$$t_{max} = \min_{x \in X} \max_{x' \in X} d(x, x')$$

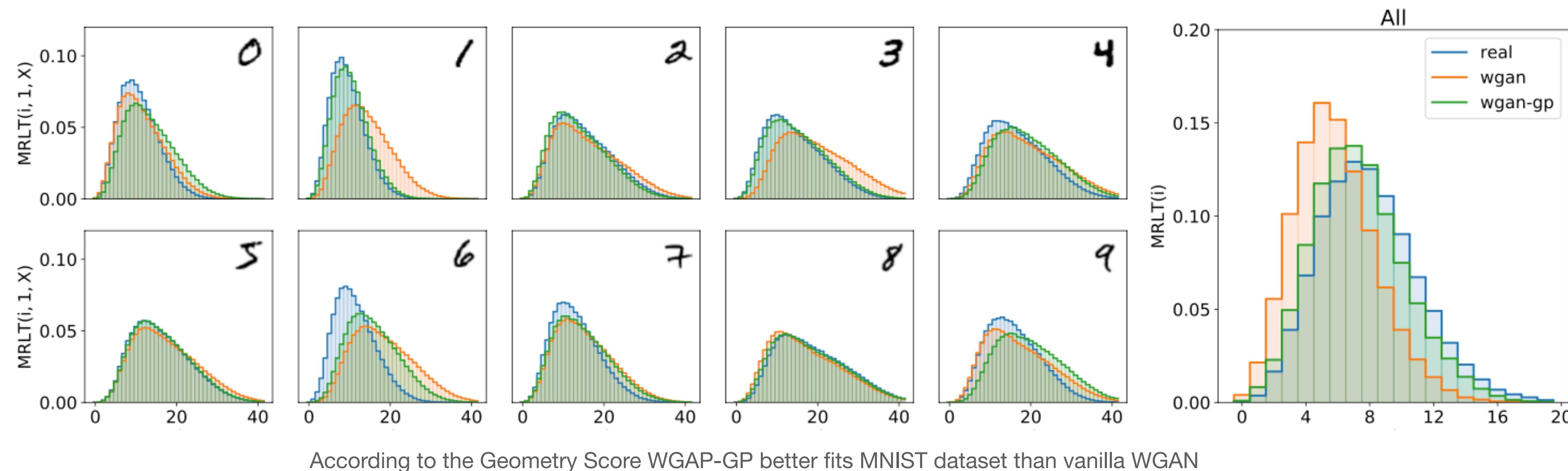


Comparing data point clouds

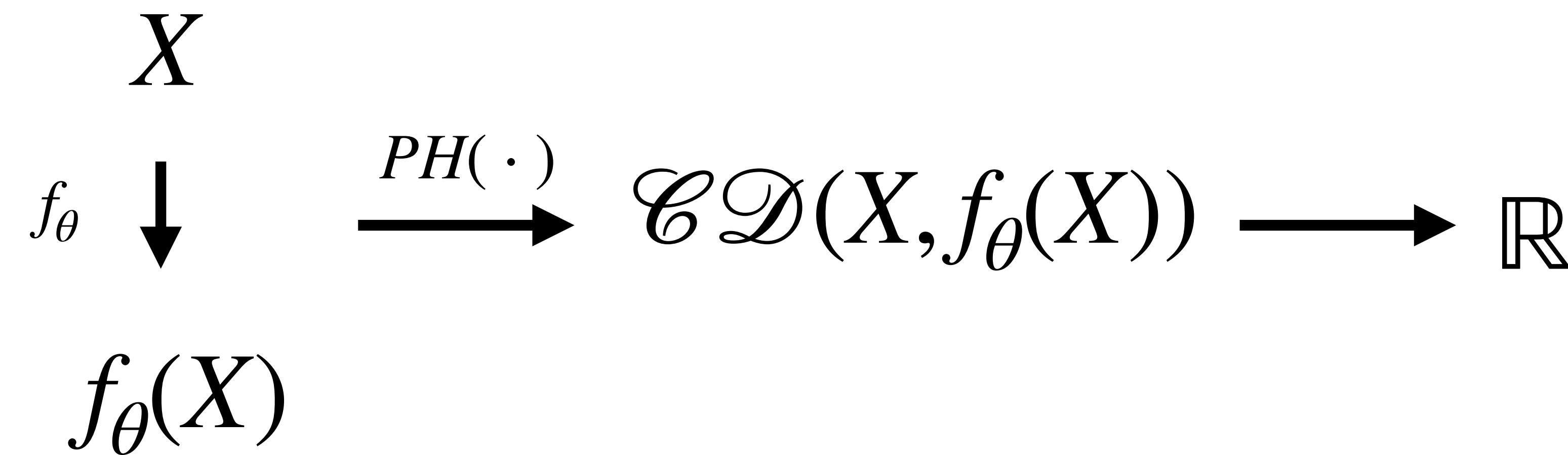
Geometry score

Given two point clouds X and Y the geometry score $GS(\mathcal{X}, \mathcal{Y})$ is given by

$$GS(\mathcal{X}, \mathcal{Y}) = \sum_{n=0}^{n_{max}-1} (RLT_X^1(n) - RLT_Y^1(n))^2$$



Comparing cross-diagrams



Manifold topology divergence

Cross-barcode

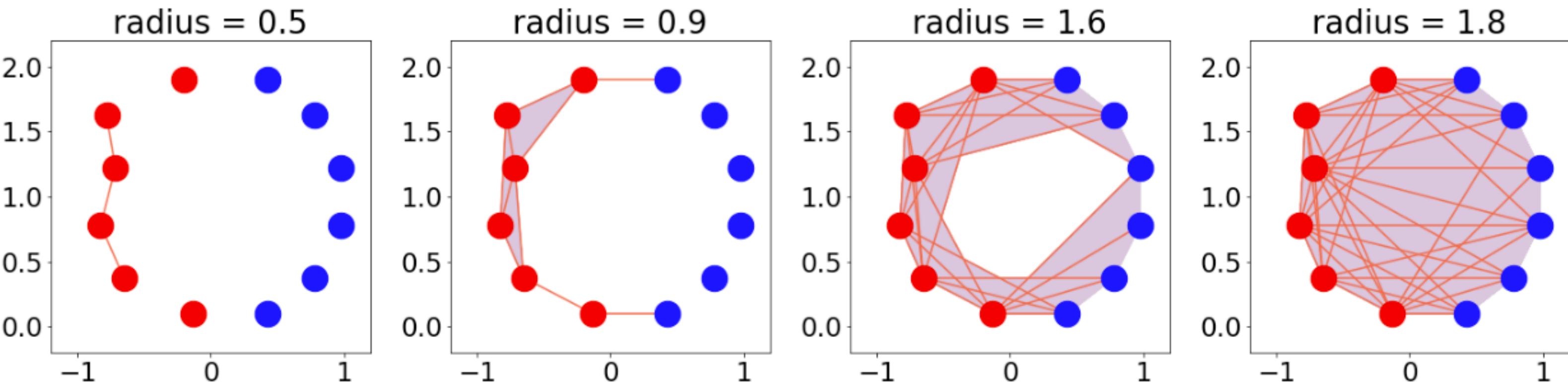
Given two point clouds $X, Y \in \mathcal{X}$ the cross-barcode $CB(X, Y)$ is a map to the persistent diagram of the Vietoris-Rips filtration of the union of X and Y quotiented by Y.

$$CB : (X, Y) \rightarrow D_{VR}((X \cup Y)/Y)$$

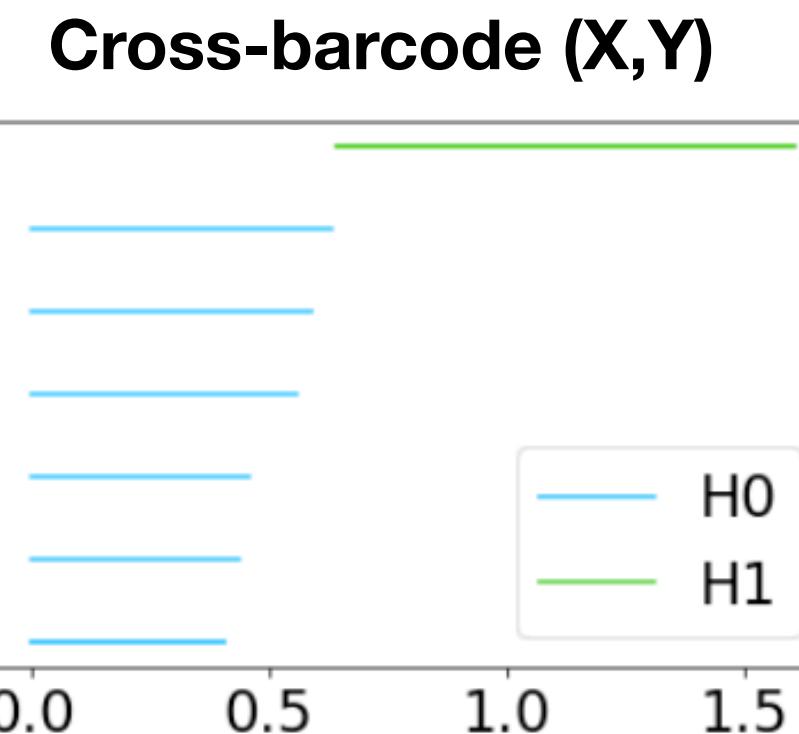
If point clouds X and Y coincide then $CB(X, Y) = \emptyset$.

	x_1	x_2	x_3	y_1	y_2
x_1					
x_2					
x_3					
y_1				0	0
y_2				0	0

$a_{ij} = d(x_i, x_j), d(x_i, y_j), d(x_j, y_i)$
 $a_{ij} = 0, \forall (y_i, y_j)$



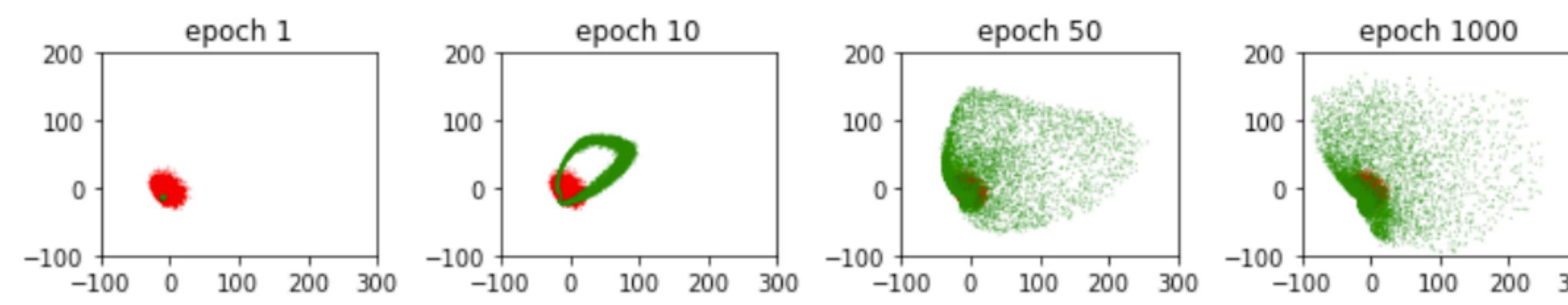
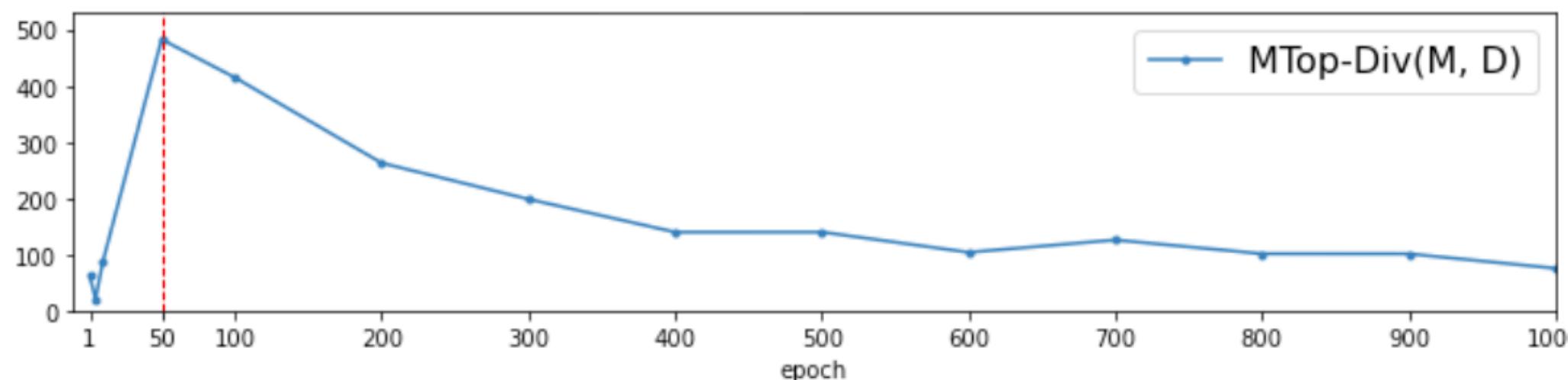
All simplices spanned by $\textcolor{blue}{Y}$ are born at 0 and not shown for perception ease.



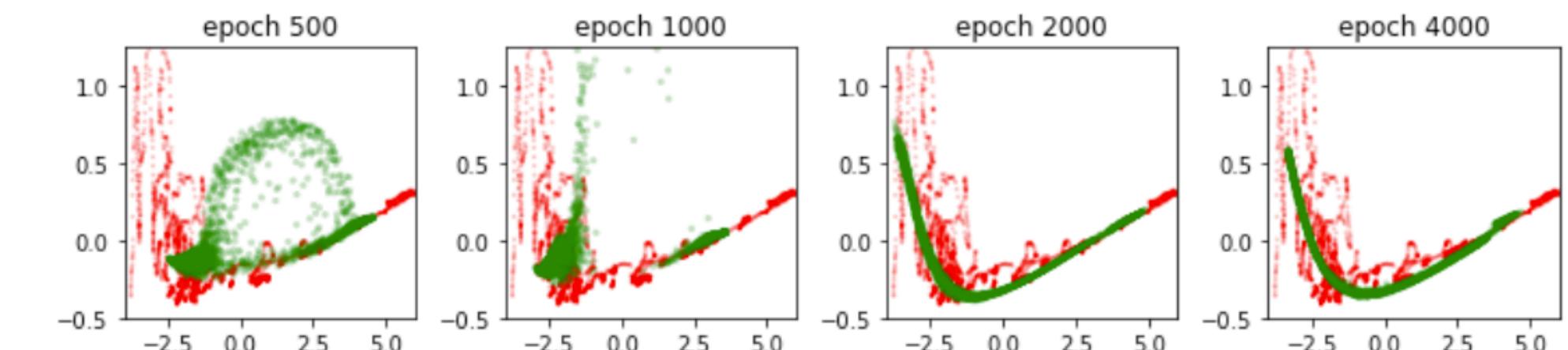
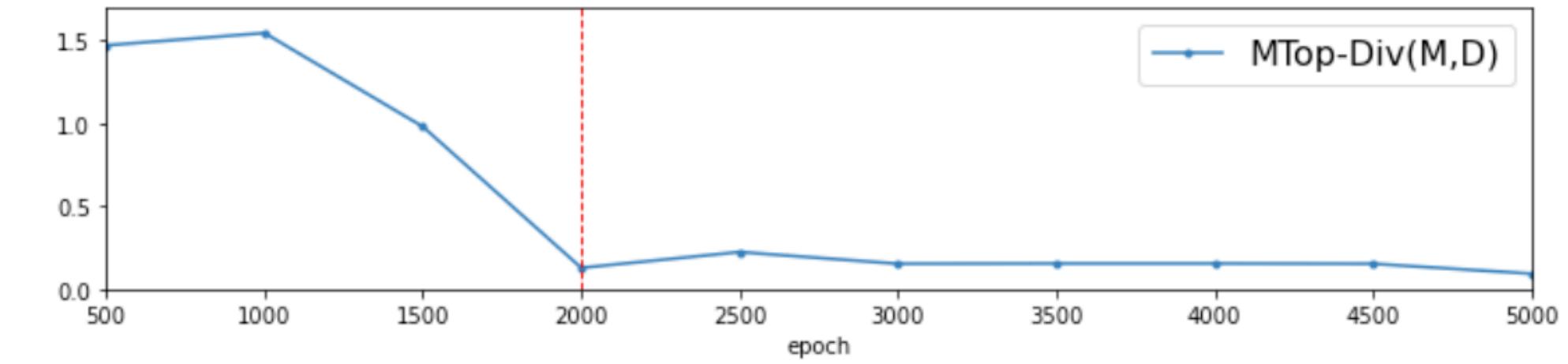
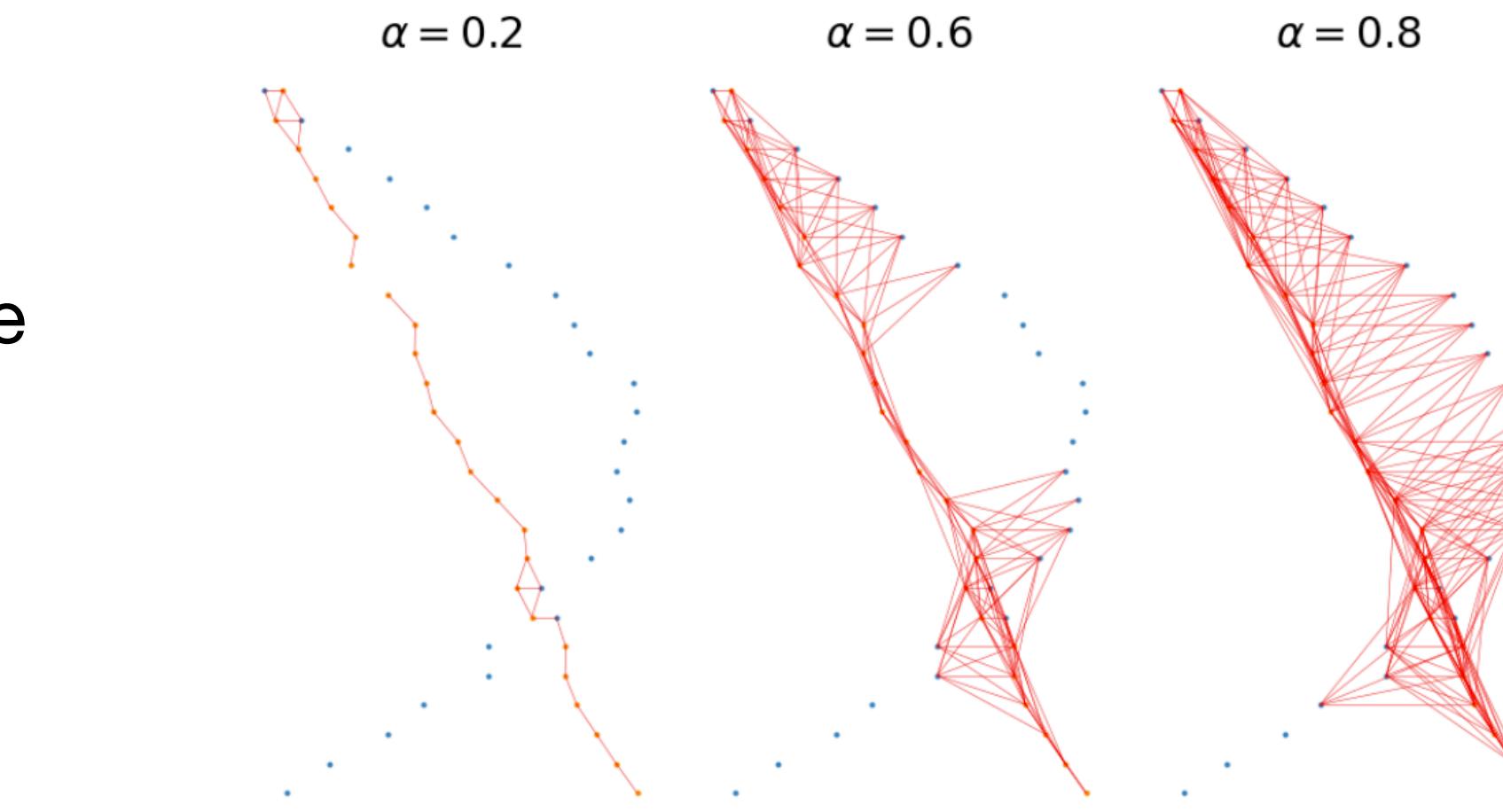
Manifold topology divergence

Given two point clouds $X, Y \in \mathcal{X}$ the manifold topology divergence $\text{MTop-Div}(X, Y)$ is given by the sum of persistences in $CB(X, Y)$

$$\text{MTop-Div}(X, Y) = \sum_i^{|CB(X, Y)|} (d_i - b_i)$$



Training dynamics of GAN applied to 3D shapes



Training dynamics of TimeGAN applied to market stock data

Representation topology divergence

R-cross-barcode

Given two point clouds $X \in \mathcal{X}, Y \in \mathcal{Y}$ with point-to-point correspondence, i.e. two embeddings of the same data, possibly in different ambient spaces, an r-cross-barcode $RCB(X, Y)$

$$RCB : (X, Y) \rightarrow D_{VR}((X, \min(X, Y))$$

$$\text{RTop} - \text{Div}(X, Y) = \sum_i^{|RCB(X, Y)|} (d_i - b_i)$$

	x_1	x_2	x_3
x_1			
x_2			
x_3			

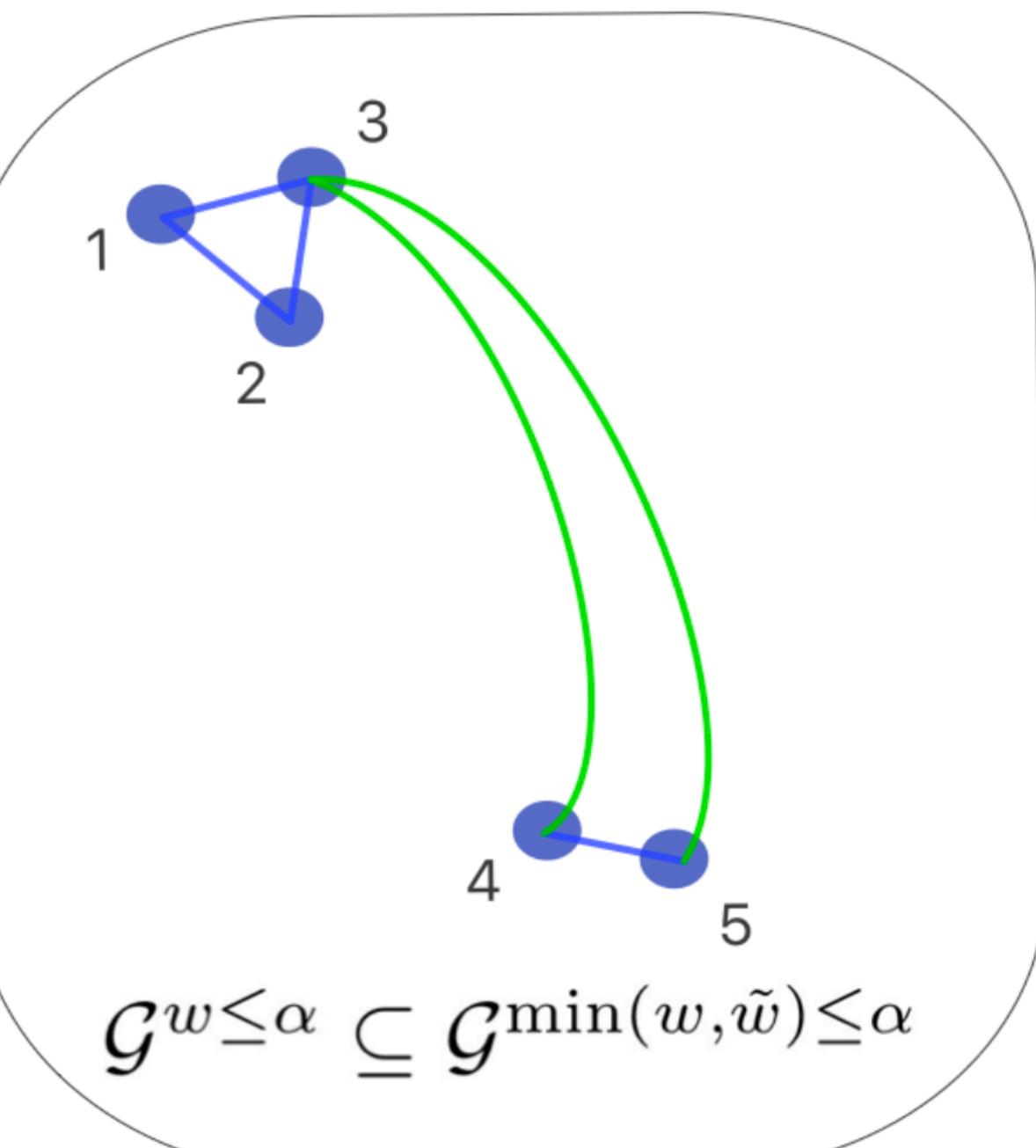
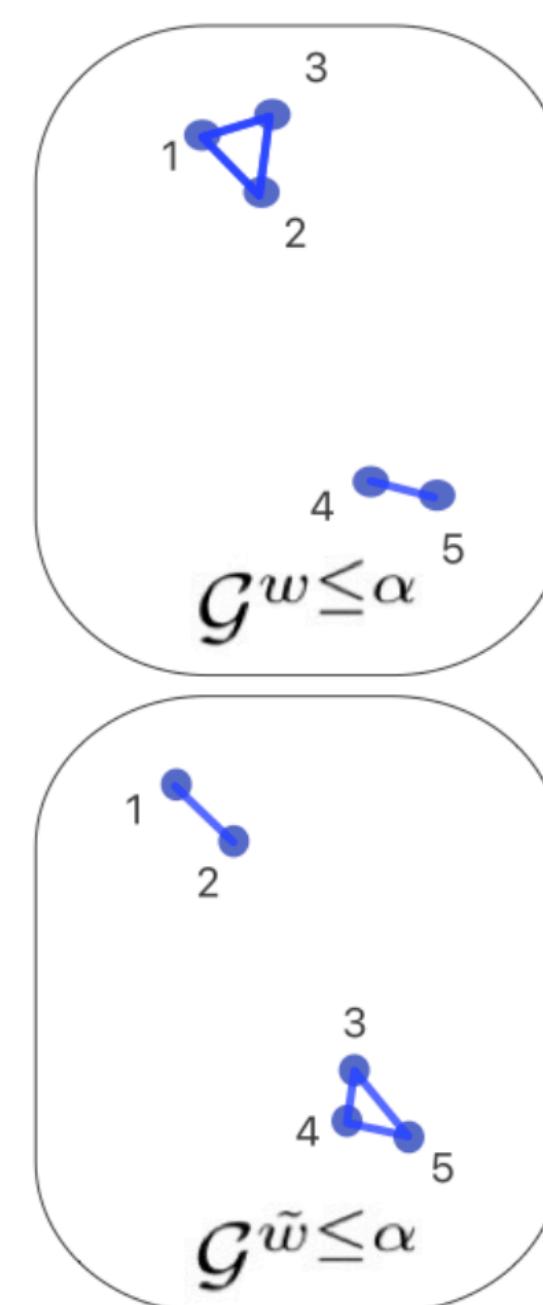
$$a_{ij} = d(x_i, x_j)$$

	y_1	y_2	y_3
y_1			
y_2			
y_3			

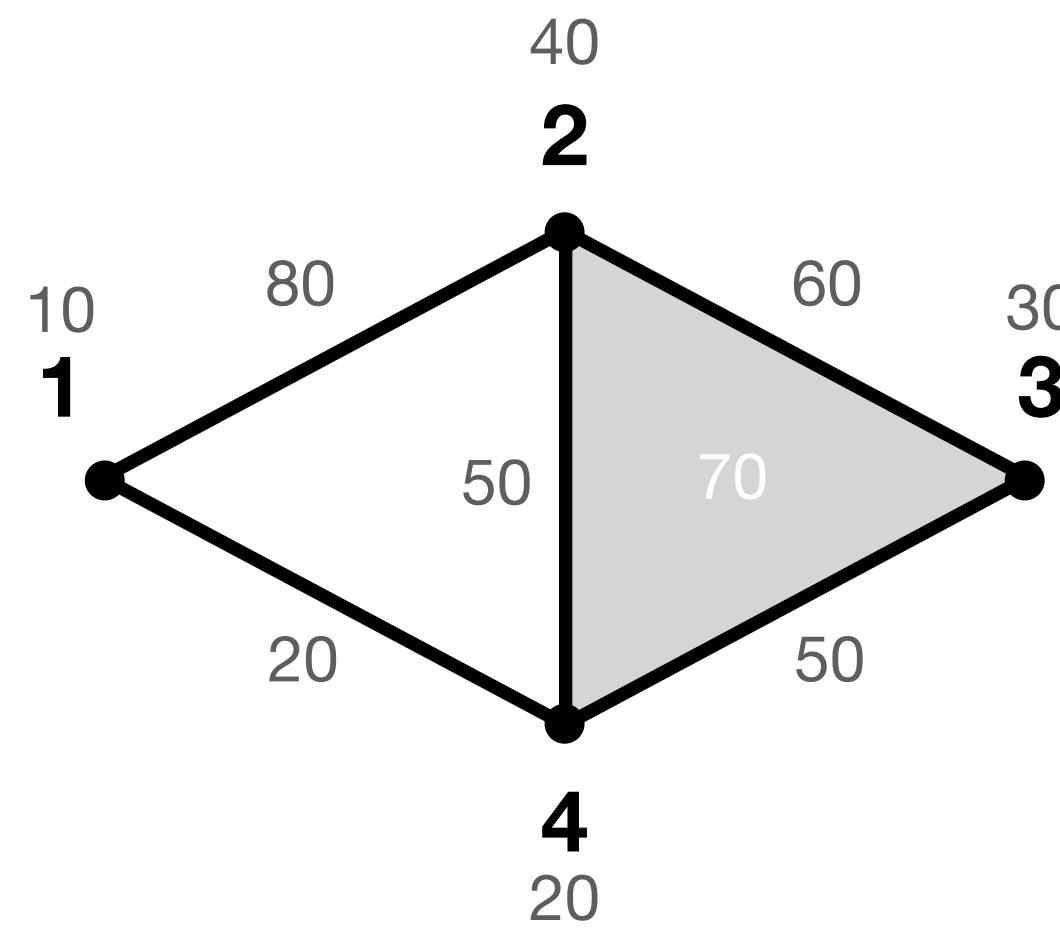
$$a_{ij} = d(y_i, y_j)$$

	x_1	x_2	x_3	x_1	x_2	x_3
x_1						
x_2						
x_3						
x_1				∞		
x_2				∞	∞	
x_3	∞	∞				
x_1						
x_2						
x_3						

$$a_{ij} = \min(d(x_i, x_j), d(y_i, y_j))$$



Differentiating persistent homology



R =

	10	20	20	30	40	50	50	60	70	80
1										
4										
14										
3										
2										
24										
34										
23										
12										

Persistence pairing π_X

- | | |
|-------------|----------------------|
| (4, 14) 0 | (1, \emptyset) 0 |
| (2, 24) 0 | (12, \emptyset) 1 |
| (3, 34) 0 | |
| (23, 234) 1 | |

Persistence diagram D_X

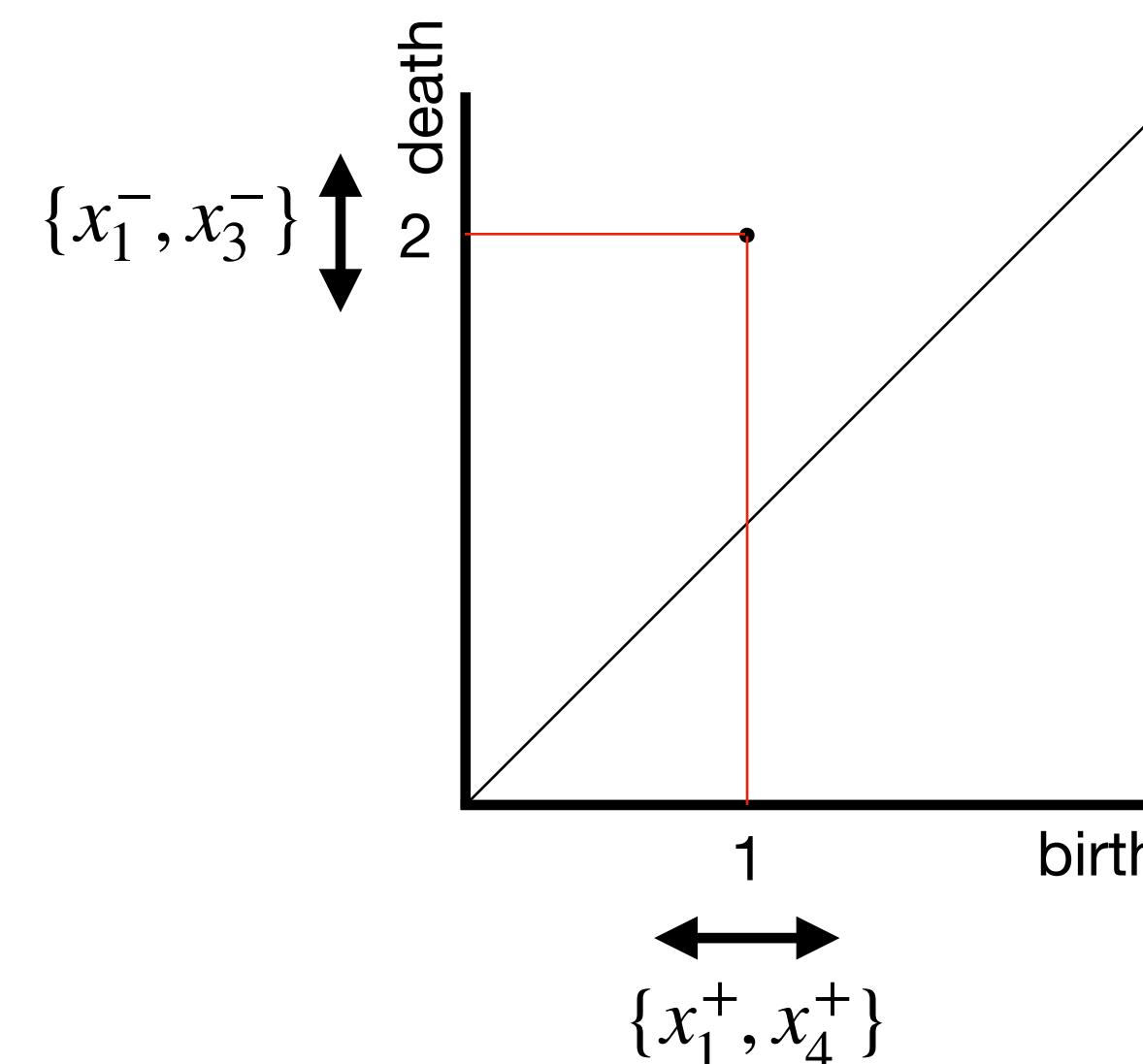
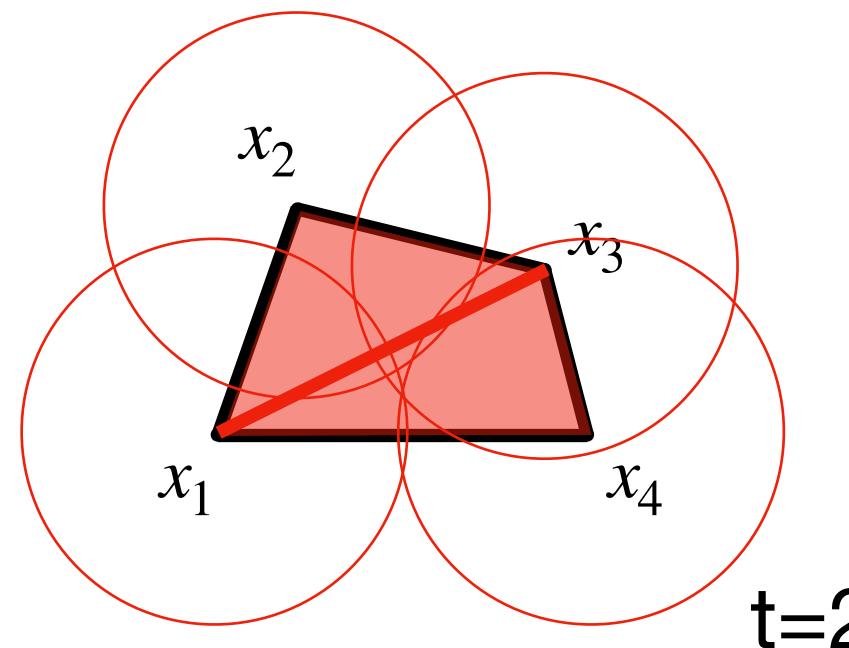
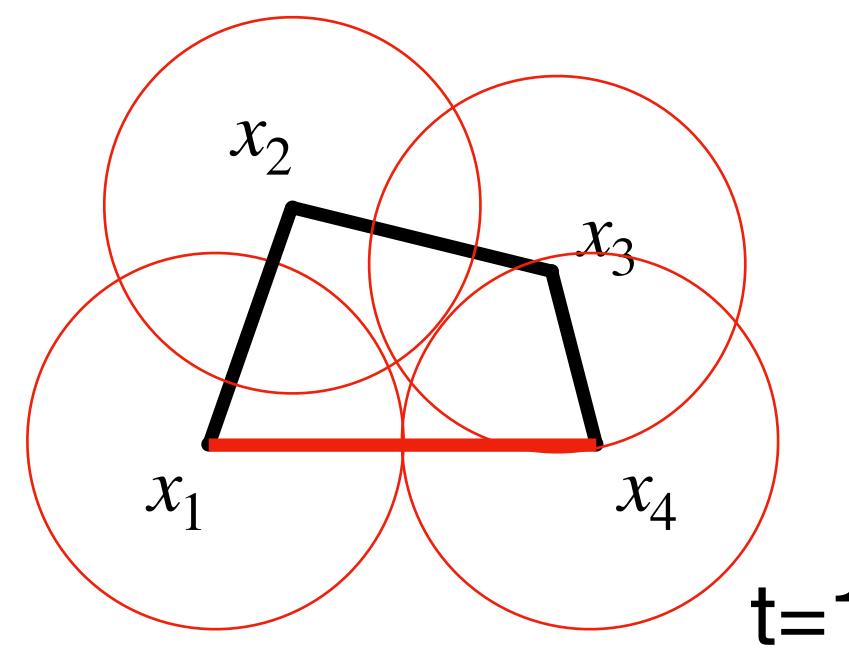
- | | |
|------------|----------------------|
| (20, 20) 0 | (10, \emptyset) 0 |
| (40, 50) 0 | (80, \emptyset) 1 |
| (30, 50) 0 | |
| (60, 70) 1 | |

Differentiating persistent homology

Vietoris-Rips filtration

Persistent homology as a map: $PH : (X, f) \rightarrow D_f(X)$

Given a dataset X , the Vietoris-Rips simplicial complex $K_f(X)$ with the filtration function $f(\sigma) = \max_{(x_i, x_j) \in \sigma} \|x_i - x_j\|_2$ the gradient of point in the persistent diagram



Inverse function

$$\pi(\sigma) = \arg \max_{(x_i, x_j) \in \sigma} \|x_i - x_j\|_2$$

Gradient

$$\frac{\partial D(X)}{\partial X}$$

$$\frac{\partial b_k}{\partial x_i^+} = \partial_{x_i^+} f(\sigma) = + \frac{x_i^+ - x_j^+}{\|x_i^+ - x_j^+\|_2}, \quad \frac{\partial b_k}{\partial x_j^+} = \partial_{x_j^+} f(\sigma) = - \frac{x_i^+ - x_j^+}{\|x_i^+ - x_j^+\|_2}$$

$$\frac{\partial d_k}{\partial x_i^-} = \partial_{x_i^-} f(\sigma) = + \frac{x_i^- - x_j^-}{\|x_i^- - x_j^-\|_2}, \quad \frac{\partial d_k}{\partial x_j^-} = \partial_{x_j^-} f(\sigma) = - \frac{x_i^- - x_j^-}{\|x_i^- - x_j^-\|_2}$$

Persistent homology-based projection pursuit

Given a dataset $X \in \mathbb{R}^m$ find a mapping $f: \mathcal{X} \in \mathbb{R}^m \rightarrow \mathcal{Z} \in \mathbb{R}^n$ where $n << m$, while optimizing/preserving some relevant properties of the data.

Properties

- variance/distances
- statistical independence
- persistent homology groups

Projection pursuit framework

- mapping in linear
- find best *index*, the value of loss function assigned for each projection Z

Persistent homology-based projection pursuit

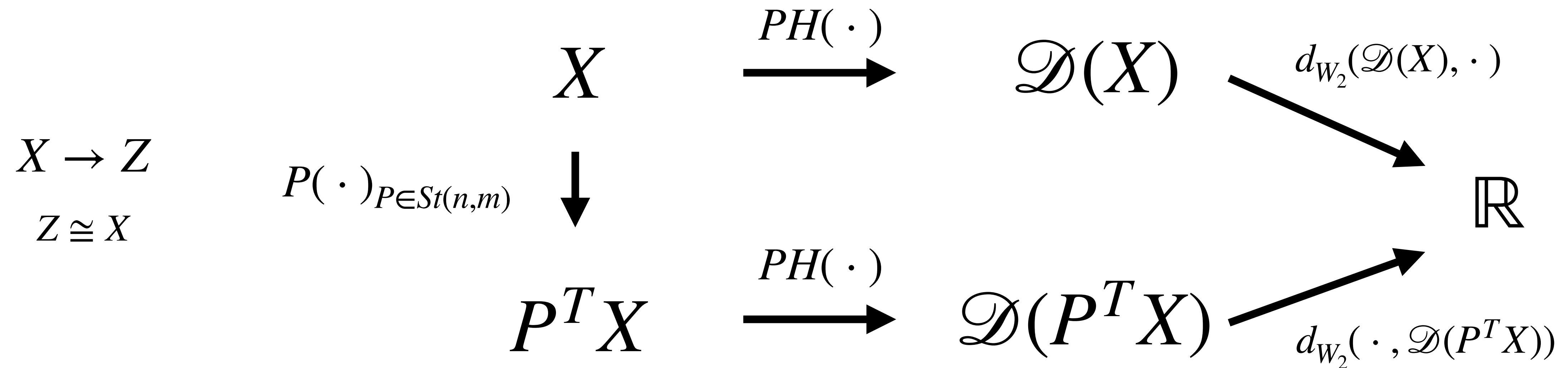
Given a dataset $X \in \mathbb{R}^m$ find a mapping $f: \mathcal{X} \in \mathbb{R}^m \rightarrow \mathcal{Z} \in \mathbb{R}^n$, where $n \ll m$, while preserving topological properties of the data in terms of persistent homology groups

$$\min_{P \in St(n,m)} W_2(\mathcal{D}(X), \mathcal{D}(P^T X))$$

Bottleneck distance bounds Gromov-Hausdorff distance

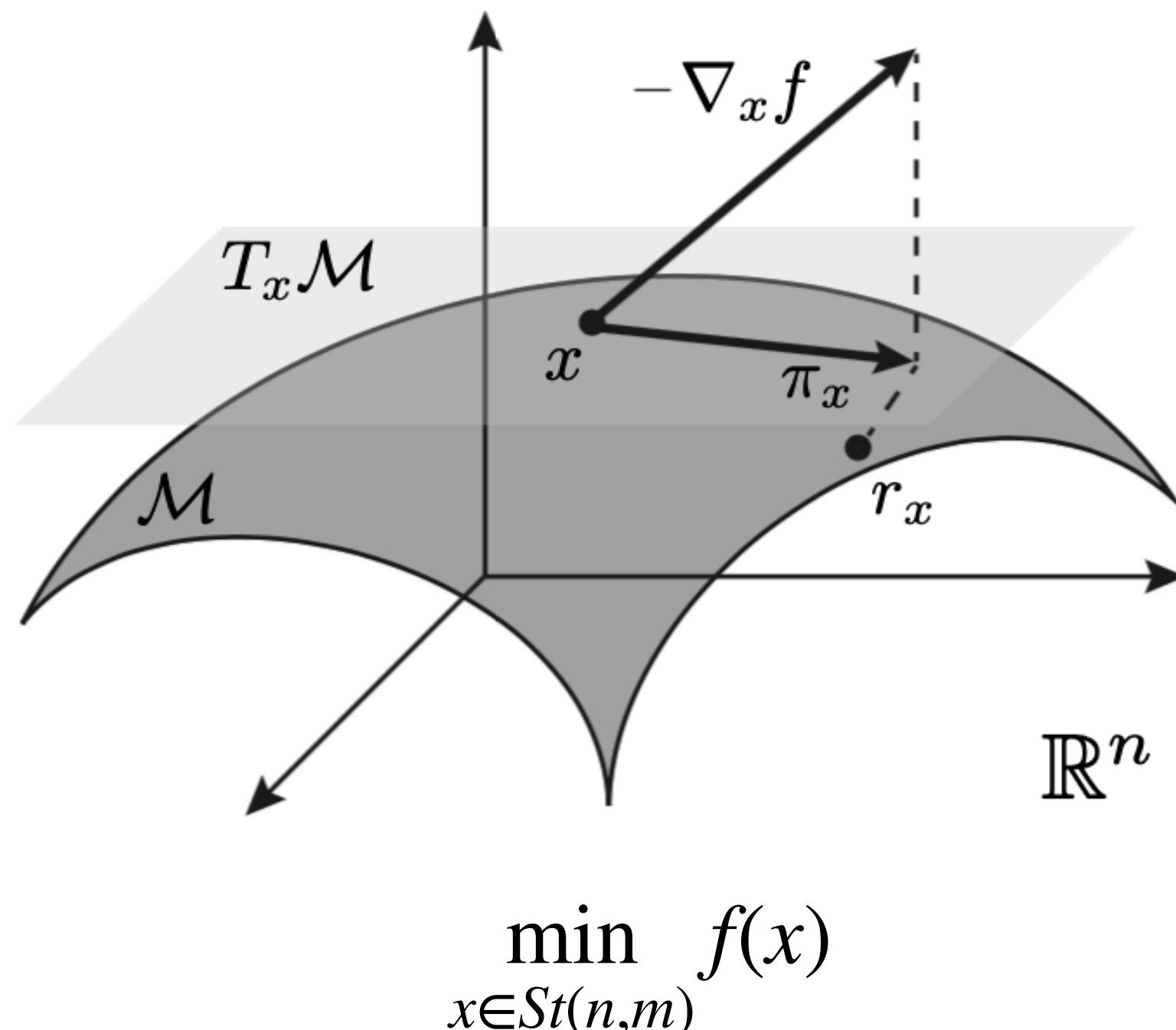
$$d_{W_\infty}(\mathcal{D}(X), \mathcal{D}(Y)) \leq d_{GH}(X, Y)$$

Persistent homology-based projection pursuit



Persistent homology-based projection pursuit

Riemannian optimization generalizes optimization algorithms to the Riemannian manifolds other than \mathbb{R}^n .



Data: A manifold \mathcal{M} , a scalar field f on \mathcal{M} , a projection $\pi_x : \mathbb{R}^n \rightarrow T_x\mathcal{M}$, a retraction $r_x : T_x\mathcal{M} \rightarrow \mathcal{M}$, an initial iterate $x_0 \in \mathcal{M}$, step size $\alpha \in \mathbb{R}_+$

Result: Sequence of iterates $\{x_k\}$

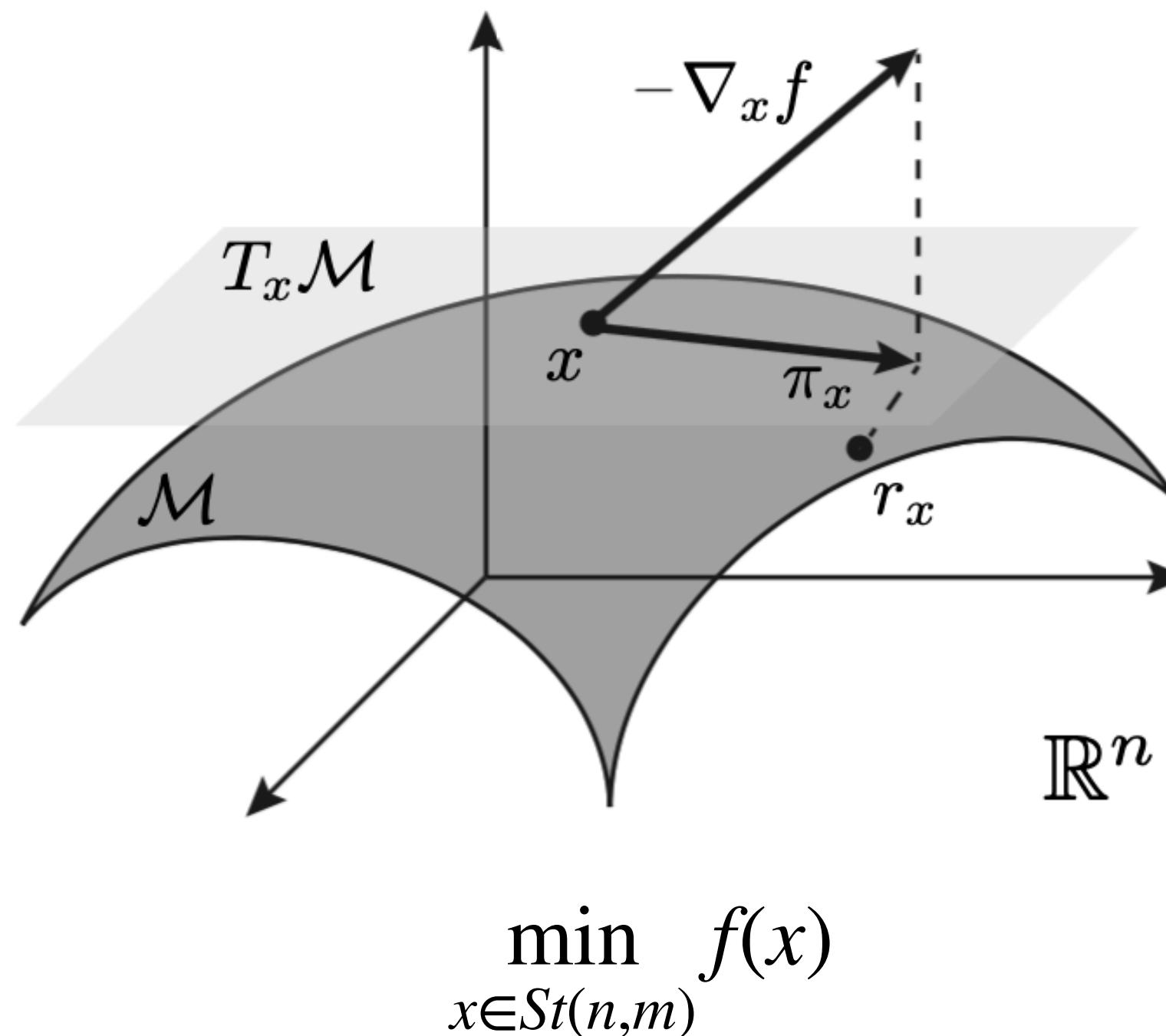
for $k = 1, 2, \dots$ **do**

$x_{k+1} = (r_{x_k} \circ \alpha \pi_{x_k})(-\nabla_{x_k} f(x_k));$

return $\{x_k\}$

Persistent homology-based projection pursuit

Riemannian optimization generalizes optimization algorithms to the Riemannian manifolds other than \mathbb{R}^n .



Stiefel manifold

$$St(n, m) = \{\mathbf{X} \in \mathbb{R}^{m \times n} \mid \mathbf{X}^T \mathbf{X} = \mathbf{I}_n\}$$

Projection map

$$\pi_{\mathbf{X}} : \mathbb{R}^{m \times n} \mapsto T_{\mathbf{X}} St(n, m)$$

Let $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$, then $\pi_{\mathbf{X}}(\mathbf{A}) = \mathbf{U}\mathbf{V}^T$.

Retraction map

Via polar decomp.

Via QR decomp.

$$r_{\mathbf{X}} : T_{\mathbf{X}} St(n, m) \mapsto St(n, m)$$

$$r_{\mathbf{X}}(\mathbf{V}) = (\mathbf{X} + \mathbf{V})(\mathbf{I}_p + \mathbf{V}^T \mathbf{V})^{-1/2}$$

$$r_{\mathbf{X}}(\mathbf{V}) = qf(\mathbf{X} + \mathbf{V})$$

Persistent homology-based projection pursuit

Given a ground metric $d : X \times X \rightarrow \mathbb{R}$ optimal transport equips the space of measures $\mathcal{P}(X)$ with a metric referred to as the Wasserstein distance

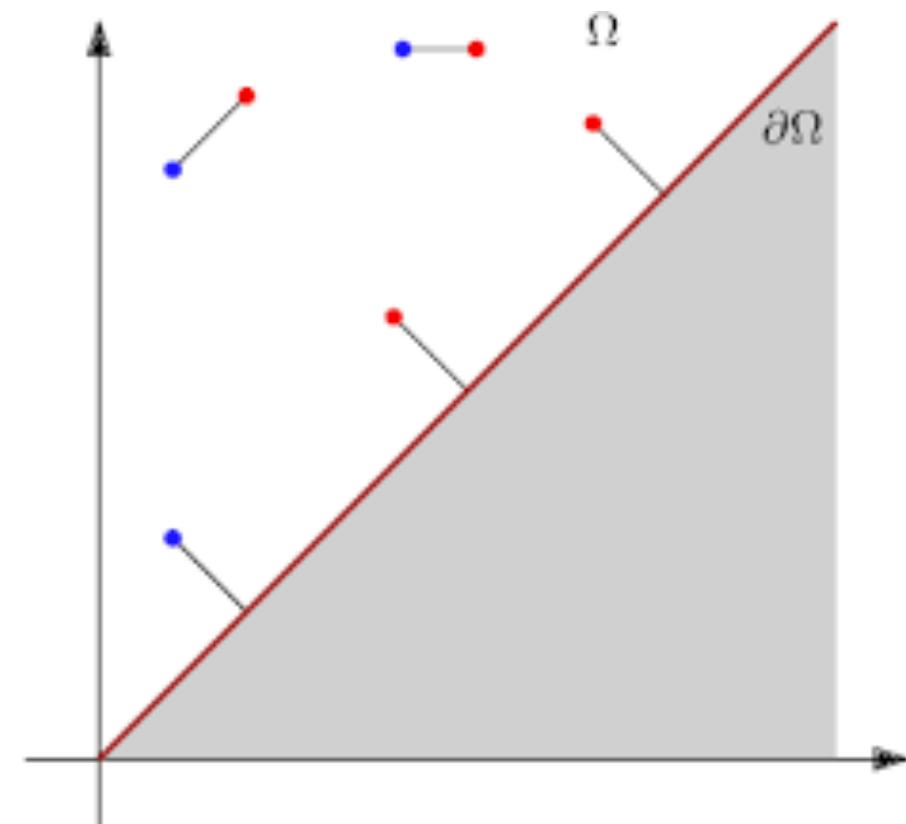
$$W_p^p(\mu, \nu) = \min_{\mathbf{T} \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{T}, \tilde{\mathbf{M}} \rangle_F$$

$$\Pi(\mathbf{a}, \mathbf{b}) = \{ \mathbf{T} \in \mathbb{R}_+^{n \times m} : \mathbf{T}\mathbf{1}_m = \mathbf{a}, \mathbf{T}^T\mathbf{1}_n = \mathbf{b} \}$$

Augmented cost matrix

$$\tilde{\mathbf{M}} = \begin{pmatrix} \mathbf{M} & \Delta_{D(X)} \\ \Delta_{D(Y)} & 0 \end{pmatrix}$$

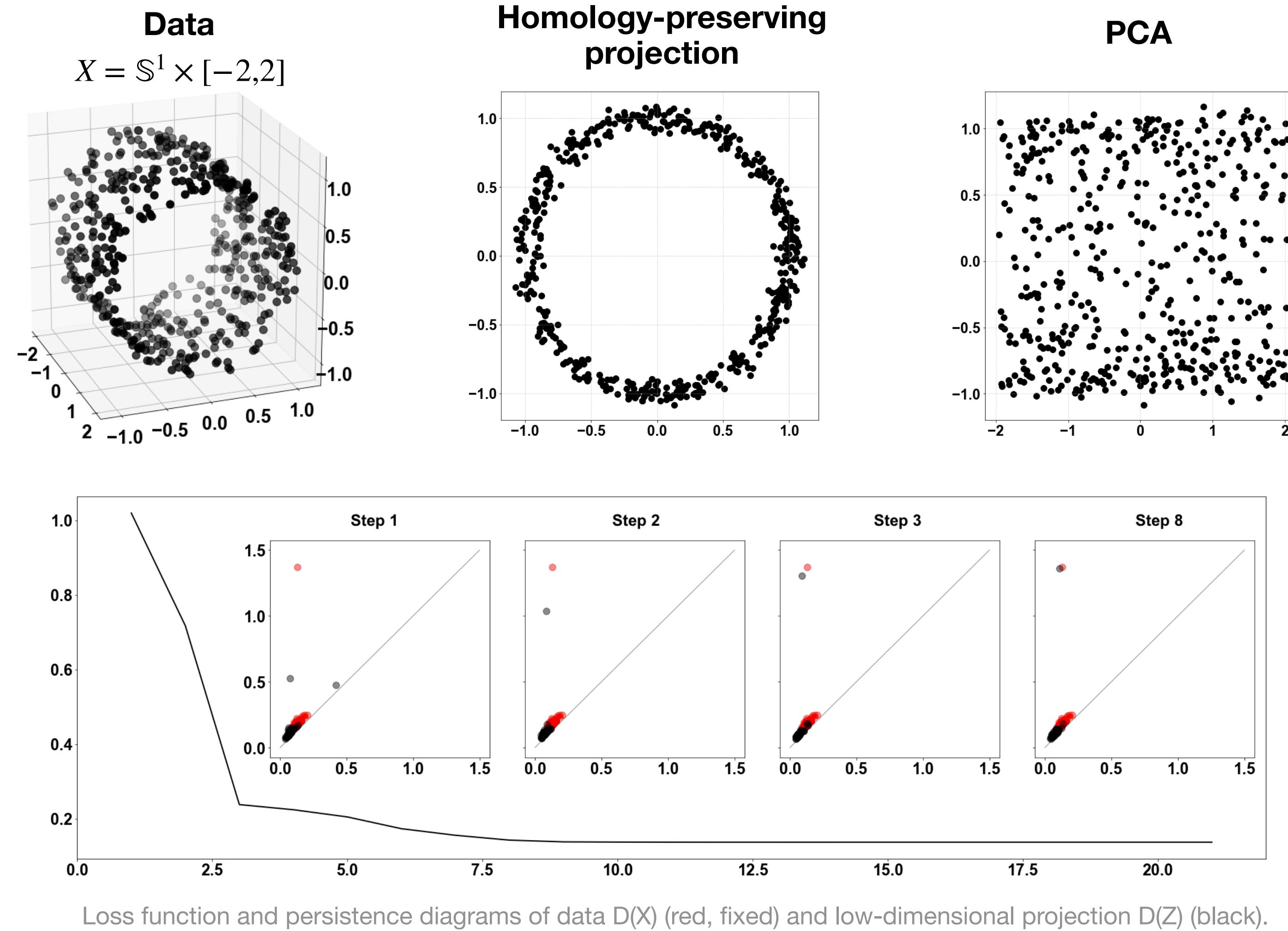
$$(\mathbf{M})_{ij} = d^p(\mathbf{x}_i, \mathbf{y}_j) \quad (\Delta_{D(X)})_i = d^p(\mathbf{x}_i, \partial\Omega) \\ (\Delta_{D(Y)})_j = d^p(\mathbf{y}_j, \partial\Omega)$$



Approximations

- Sinkhorn divergence, $O(n^2)$
- Sliced Wasserstein distance, $O(n \log n)$

Persistent homology-based projection pursuit



Topological autoencoders

$$L = L_{REC}(X, (h \circ g)(X)) + \lambda L_{TOPO}$$

$$L_{TOPO} = L_{X \rightarrow Z} + L_{Z \rightarrow X}$$

$$L_{X \rightarrow Z} = \|A_X[\pi_X] - A_Z[\pi_X]\|^2$$

$$L_{Z \rightarrow X} = \|A_Z[\pi_Z] - A_X[\pi_Z]\|^2$$

$$D_X = \{(10,20), (12,30), (15,50)\}$$

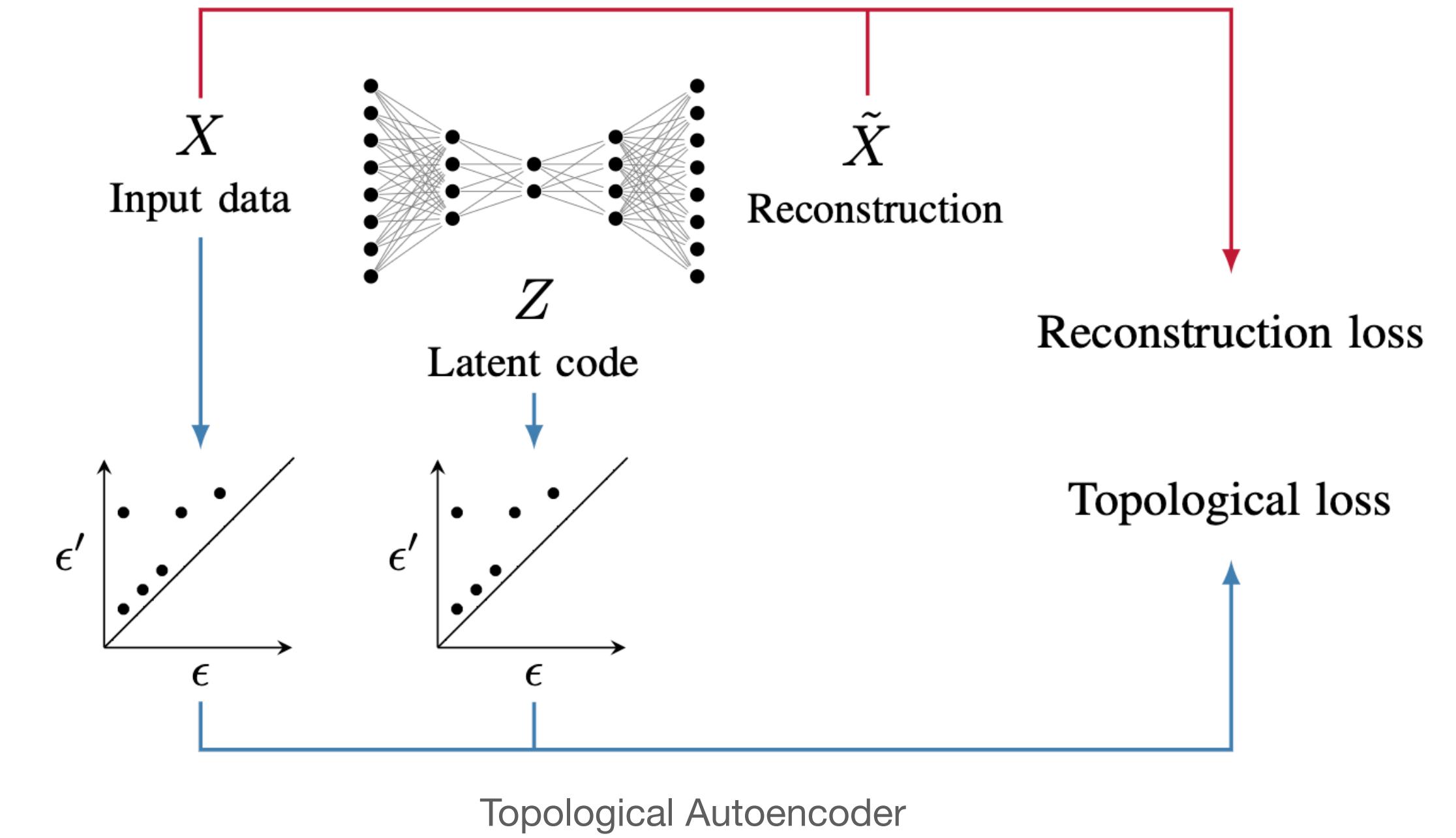
$$\pi_X = \{(1,12), (2,24), (4,34)\}$$

	1	2	3	4
1		20		
2				30
3				50
4				

$$A_X[\pi_X]$$

	1	2	3	4
1		2		
2				3
3				5
4				

$$A_Z[\pi_X]$$

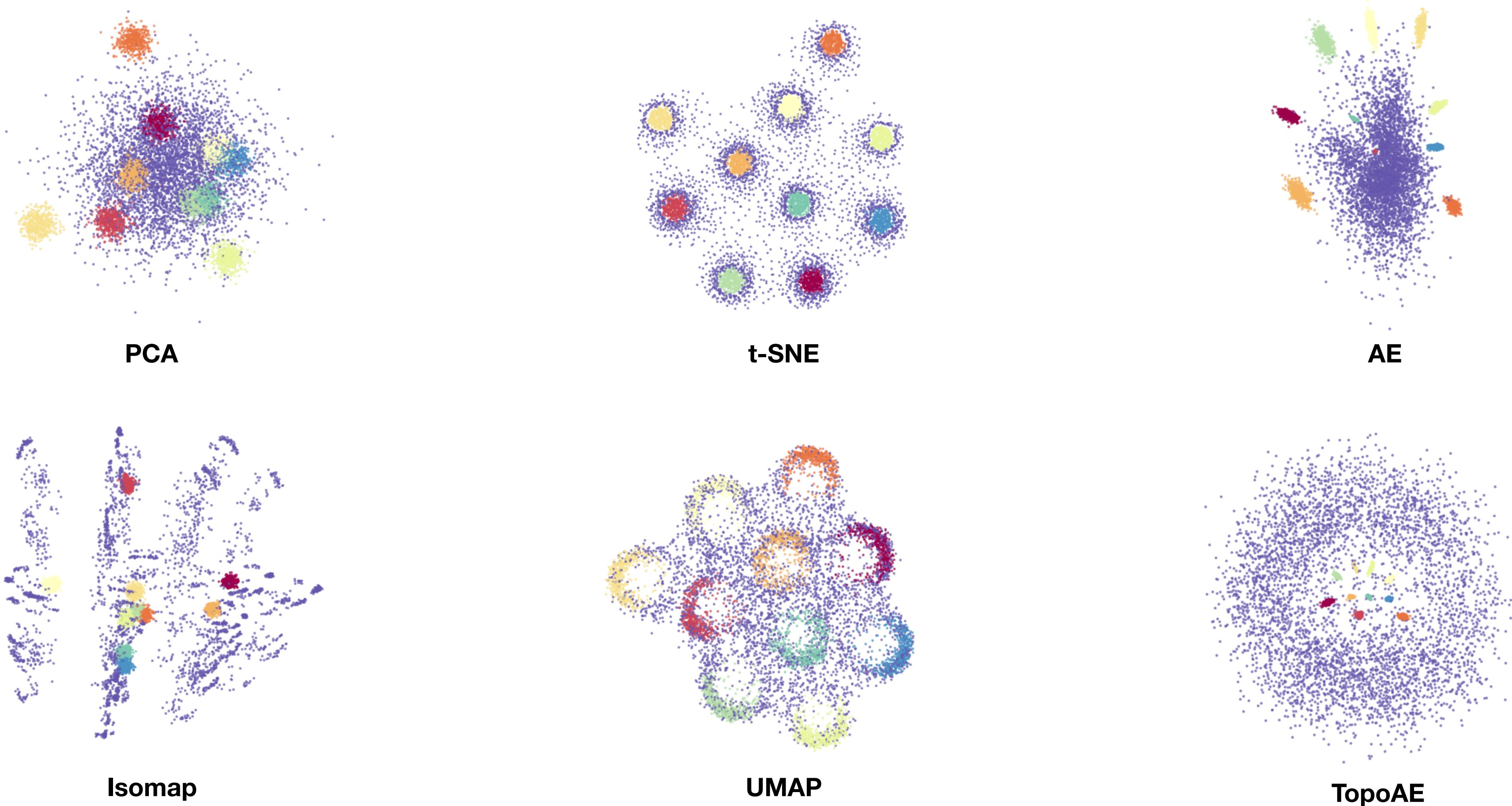


$$X \rightarrow Z \rightarrow \tilde{X}$$

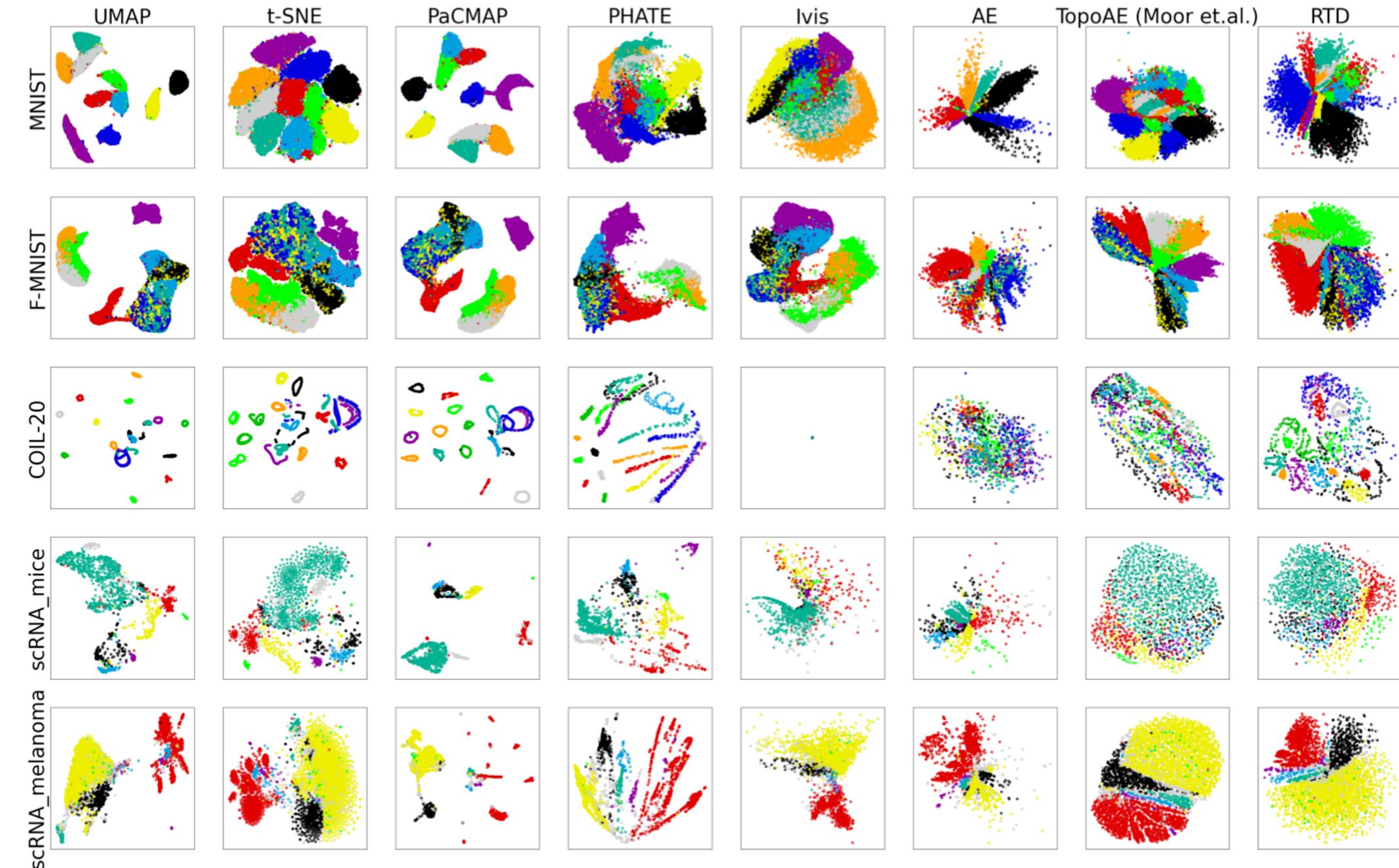
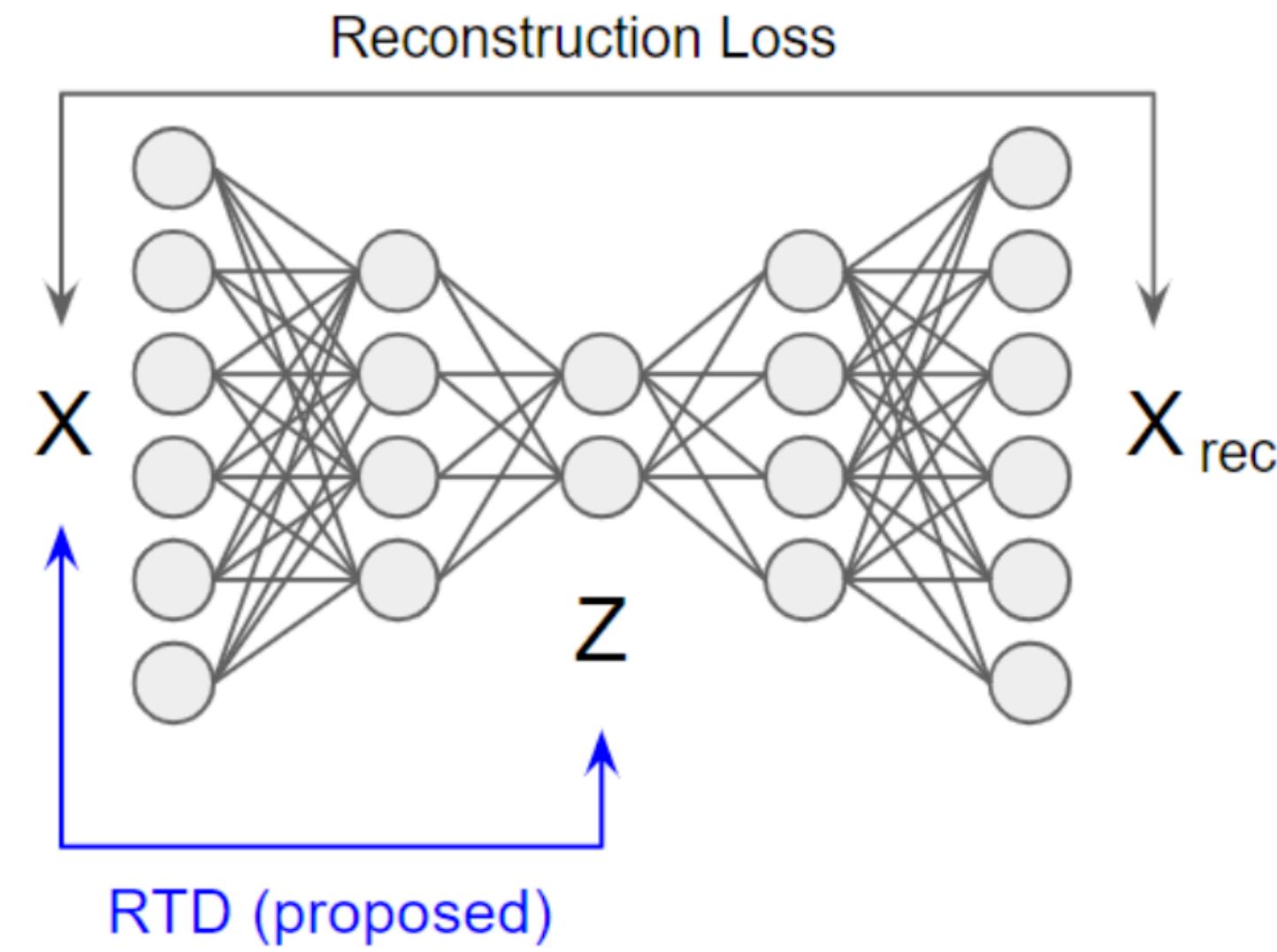
$$Z \cong X$$

Topological autoencoders

SPHERES data set that consists of ten high-dimensional 100-spheres or radius R living in a 101-dimensional space that are enclosed by one larger sphere of radius $5R$ that consists of the same number of points as the total of inner spheres.

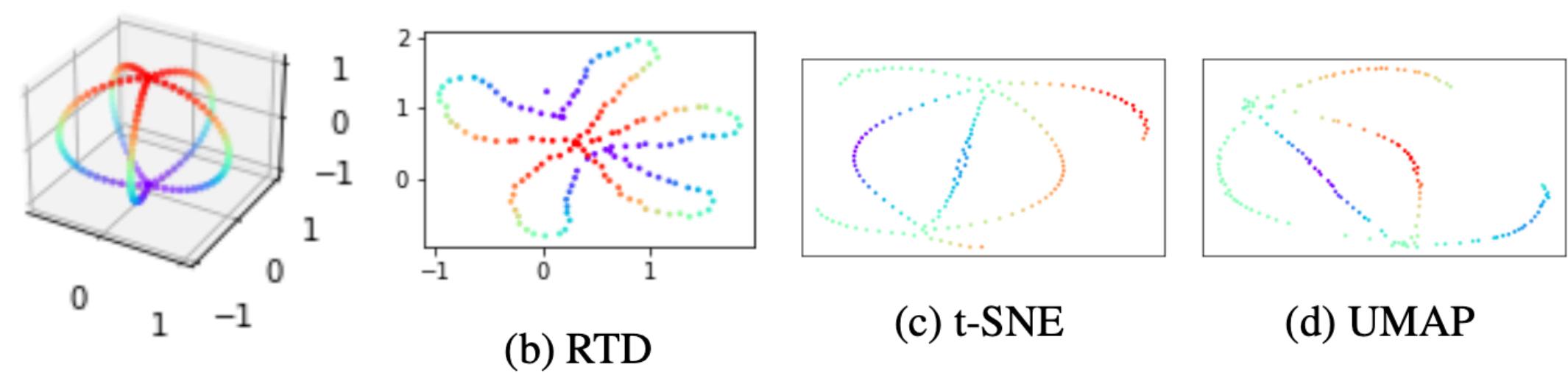


Topology-preserving representation learning

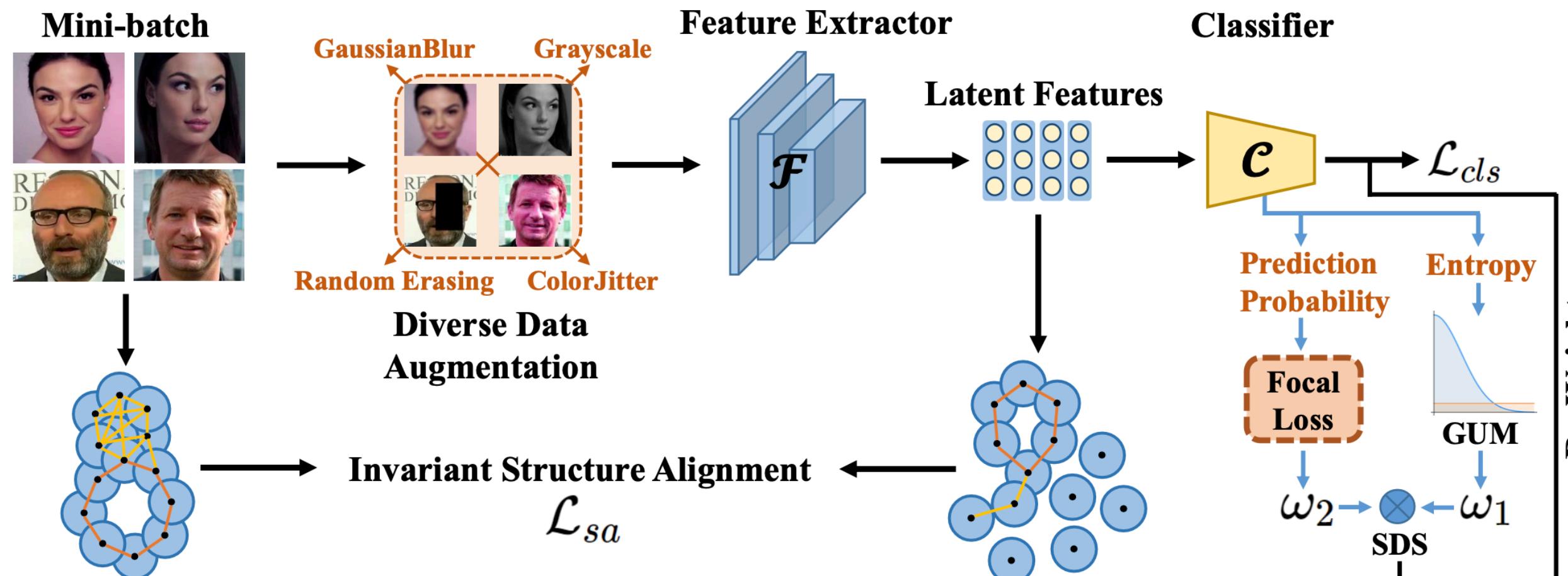


$$X \rightarrow Z \rightarrow \tilde{X}$$

$$Z \cong X$$

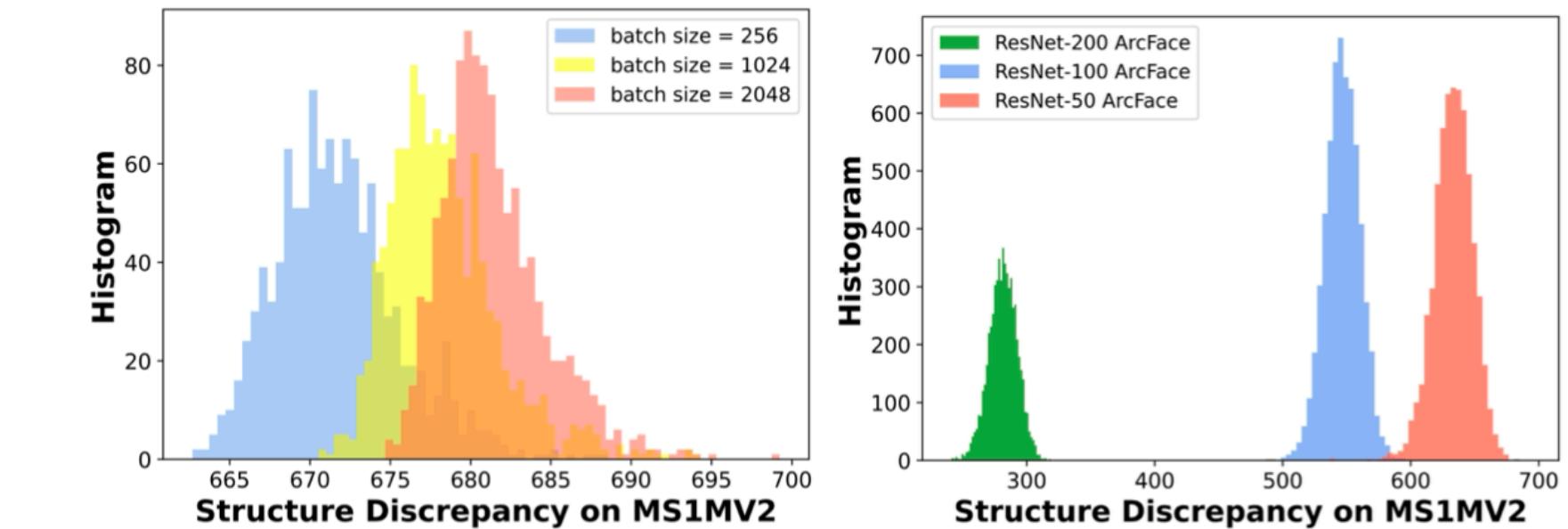


Topology alignment for face recognition



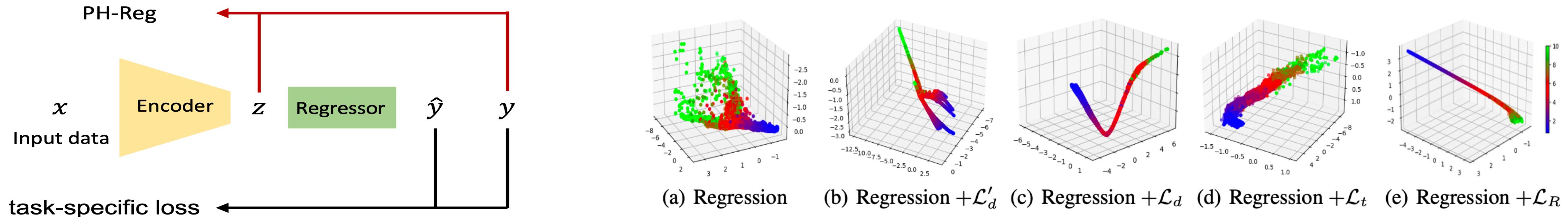
$$X \rightarrow Z \rightarrow Y$$

$$Z \cong X$$



Training Data	Method	IJB-C(1e-6)	IJB-C(1e-5)	IJB-C(1e-4)
MS1MV2	R50, ArcFace Deng et al. (2019)	80.52	88.36	92.52
	R50, MagFace Meng et al. (2021)	81.69	88.95	93.34
	R50, AdaFace Kim et al. (2022)	-	-	96.27
	R50, TopoFR [†]	89.32	94.77	96.40
	R50, TopoFR	90.52	94.71	96.49
	R100, CosFace Wang et al. (2018b)	87.96	92.68	95.56
	R100, ArcFace Deng et al. (2019)	85.65	92.69	95.74
	R100, MV-Softmax Wang et al. (2020)	-	-	95.20
	R100, CircleLoss Sun et al. (2020)	-	89.60	93.95
	R100, URL Shi et al. (2020)	-	95.00	96.60
Glint360K	R100, BroadFace Kim et al. (2020)	85.96	94.59	96.38
	R100, CurricularFace Huang et al. (2020)	-	-	96.10
	R100, MagFace+ Meng et al. (2021)	90.24	94.08	95.97
	R100, SCF Li et al. (2021b)	-	94.78	96.22
	R100, DAM-CurricularFace Liu et al. (2021)	-	-	96.20
	R100, ElasticFace+ Boutros et al. (2022)	-	-	96.65
	R100, AdaFace Kim et al. (2022)	-	-	96.89
	R100, TopoFR [†]	87.90	95.27	96.90
	R100, TopoFR	90.21	95.23	96.95
	R200, ArcFace Deng et al. (2019)	85.75	94.67	96.53

Deep regression representation learning with topology

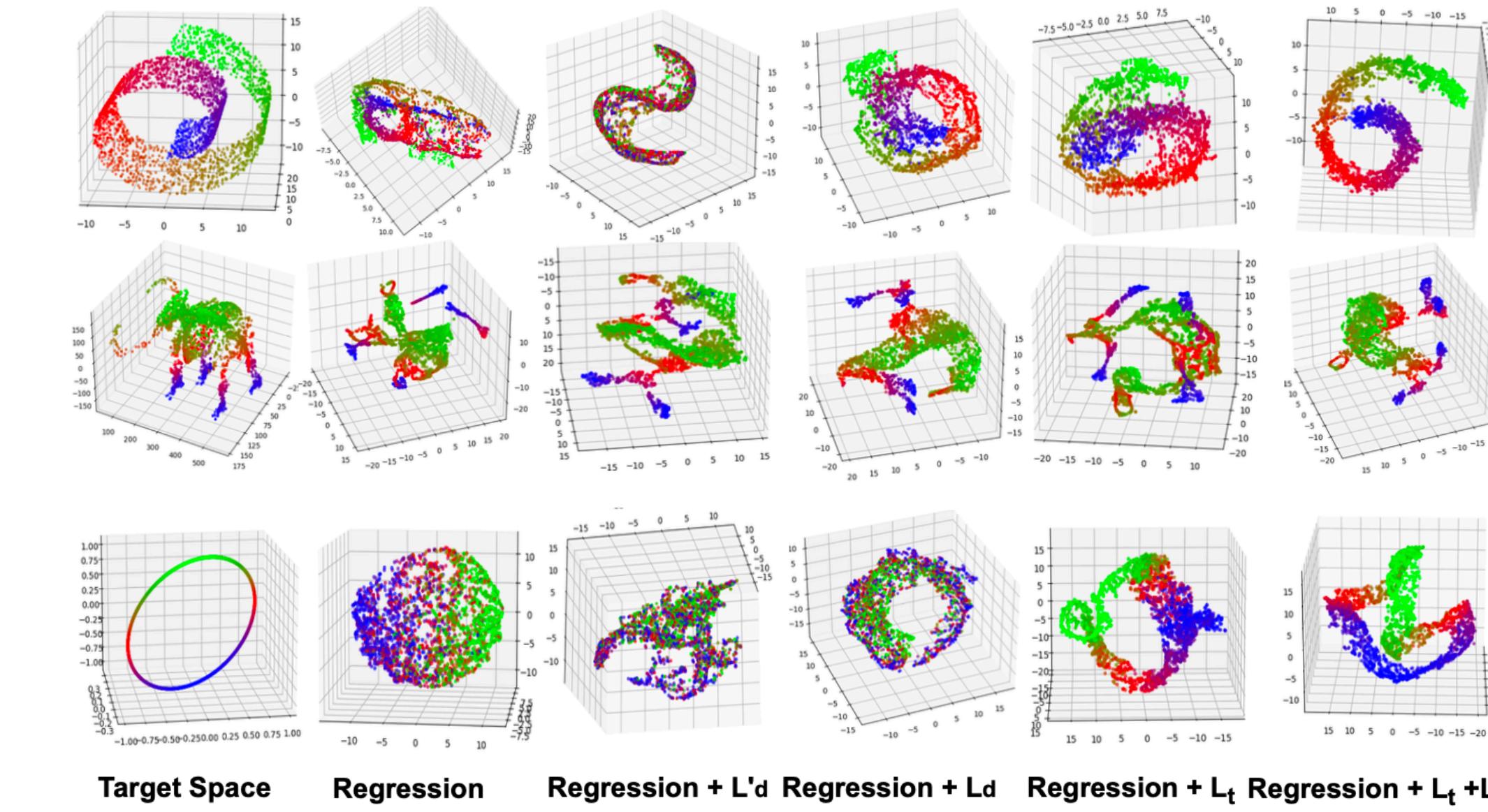


$$X \rightarrow Z \rightarrow Y$$

$$Z \cong Y$$

Table 3: Results on AgeDB

Method	MAE ↓		
	ALL	Many	Med.
Baseline (Yang et al., 2021)	7.77	6.62	9.55
+ \mathcal{L}'_d	7.81	6.96	8.88
+ \mathcal{L}_d	7.55	6.81	8.43
+ \mathcal{L}_t	7.50	6.58	8.79
+ $\mathcal{L}_d + \mathcal{L}_t$	7.48	6.52	8.71



Self-supervised learning

