

Revisiting Dimensionality Reduction with Optimal Transport: Within and Across Spaces

Hugues Van Assel



Cédric Vincent-Cuaz



Rémi Flamary



Nicolas Courty



Pascal Frossard



Titouan Vayer

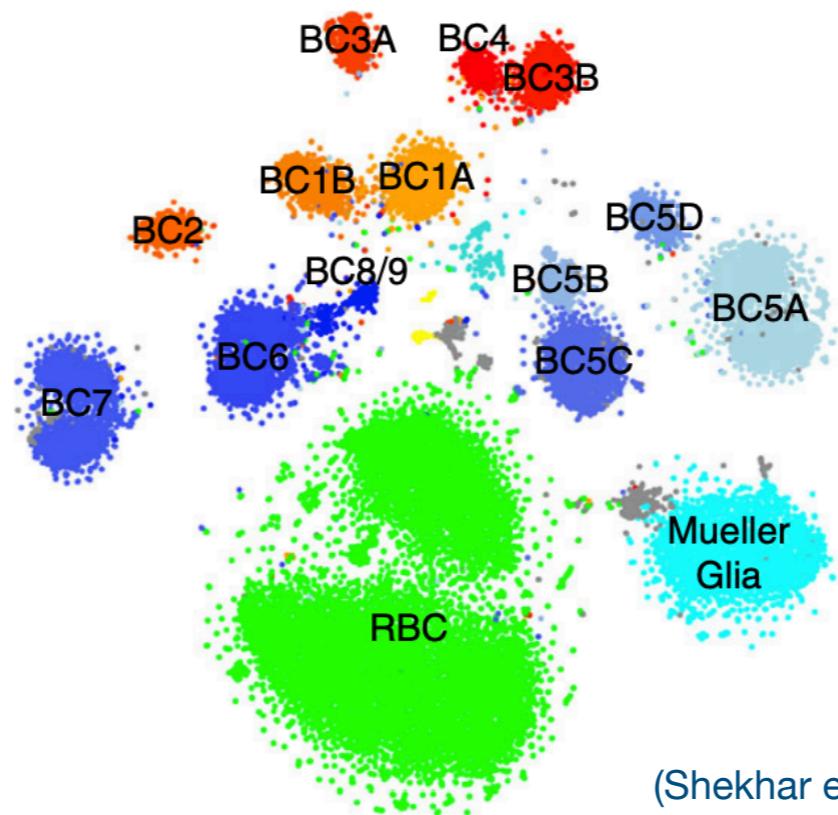
My talk

Overview of Dimensionality Reduction
as **Affinity Matching**

**Constructing Adaptive Symmetric
Affinities**

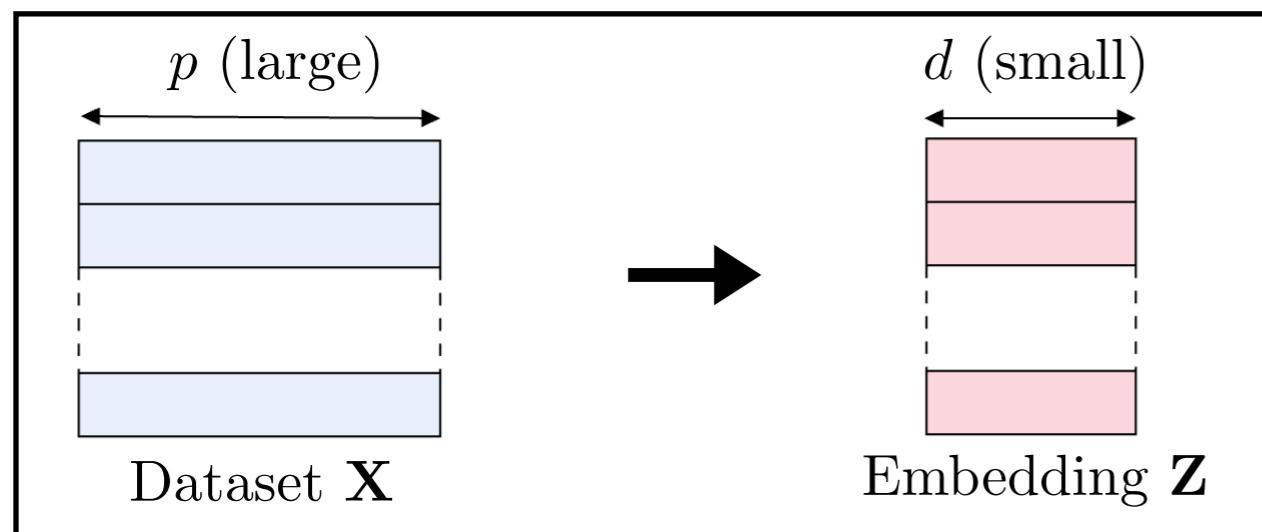
**Distributional Reduction : a Framework
to Embed Distributions**

Dimension Reduction

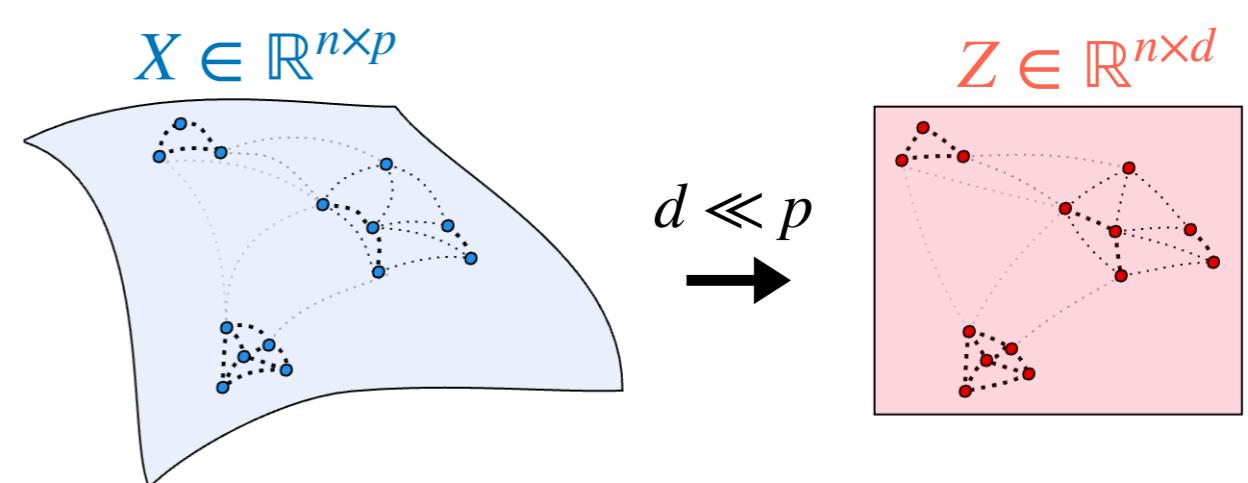


(Shekhar et al., 2016)

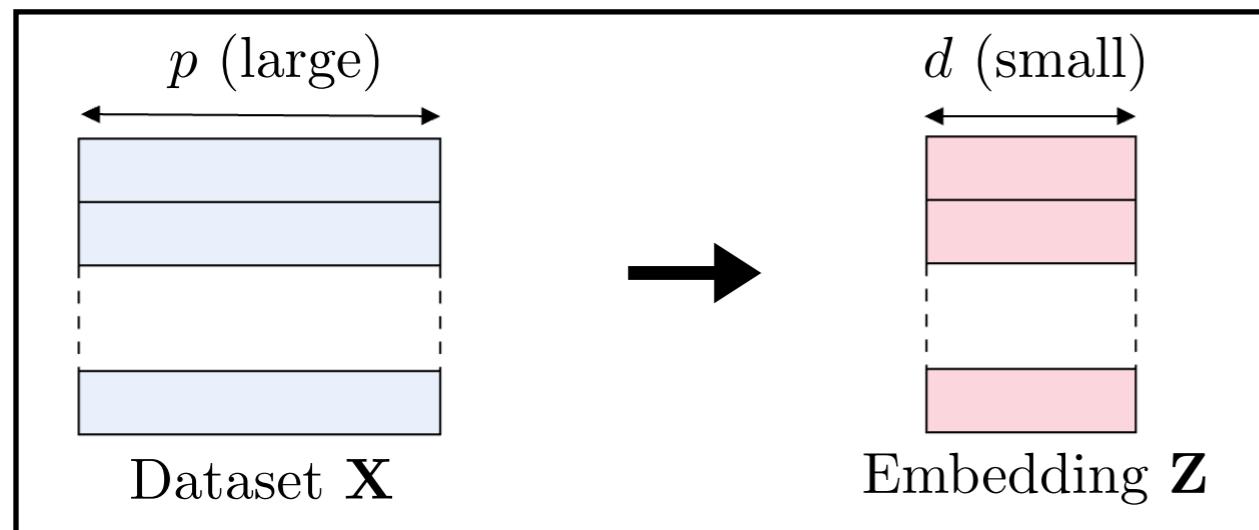
Dimension reduction



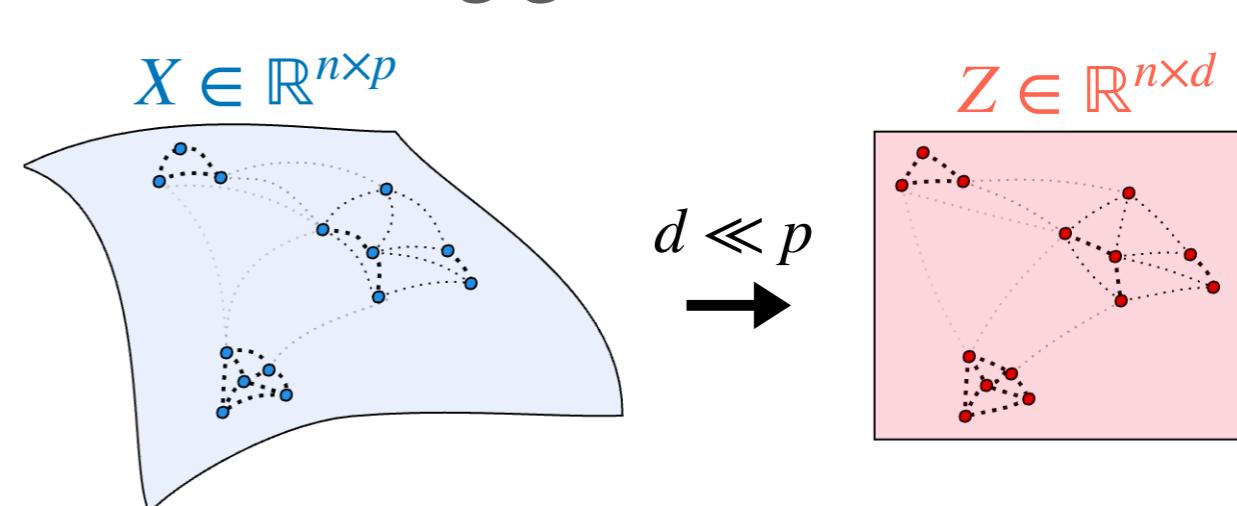
◆ Preserving geometric structure



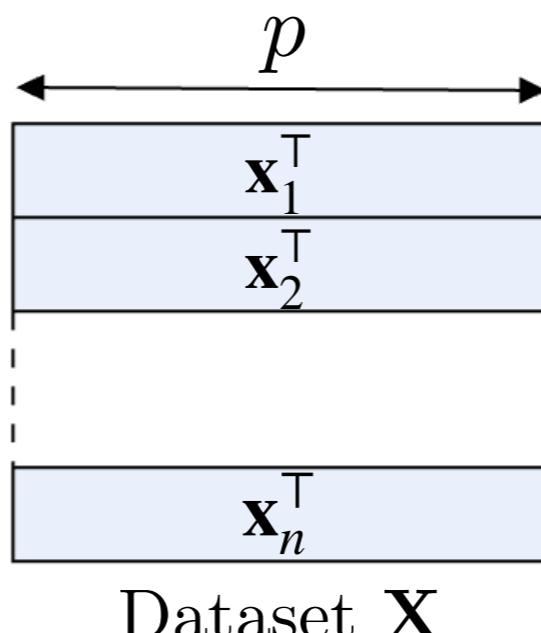
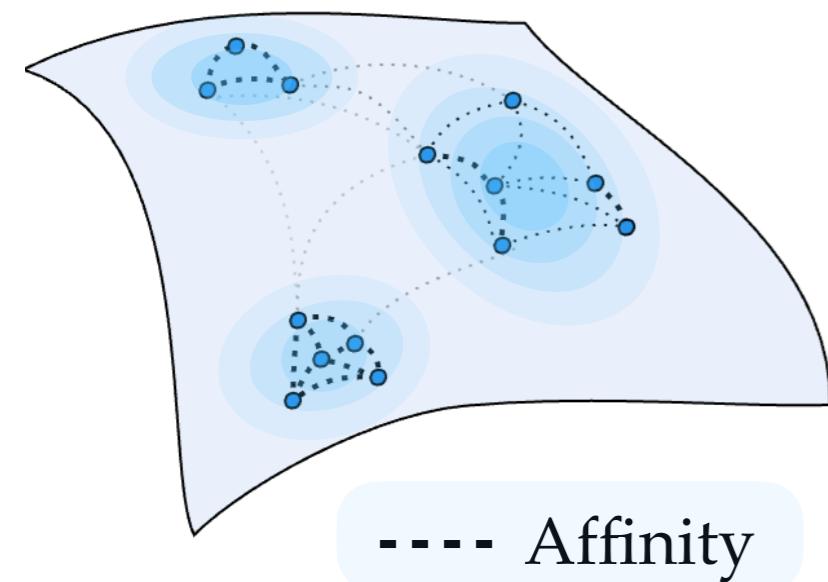
Dimension reduction



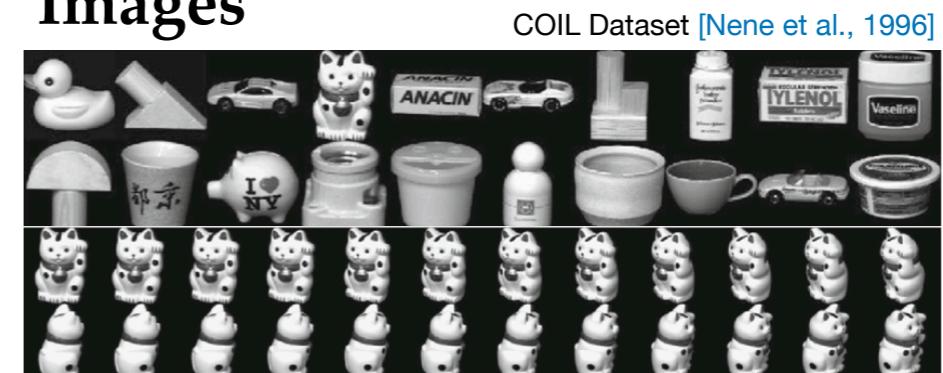
◆ Preserving geometric structure



◆ Affinity Matrices



Images

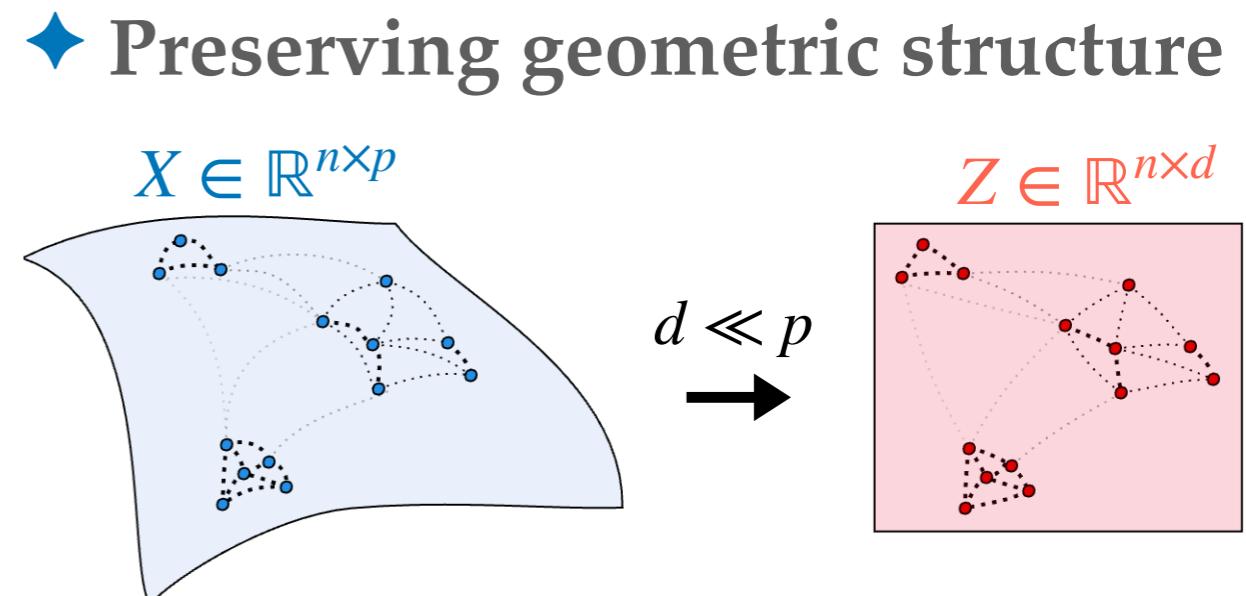
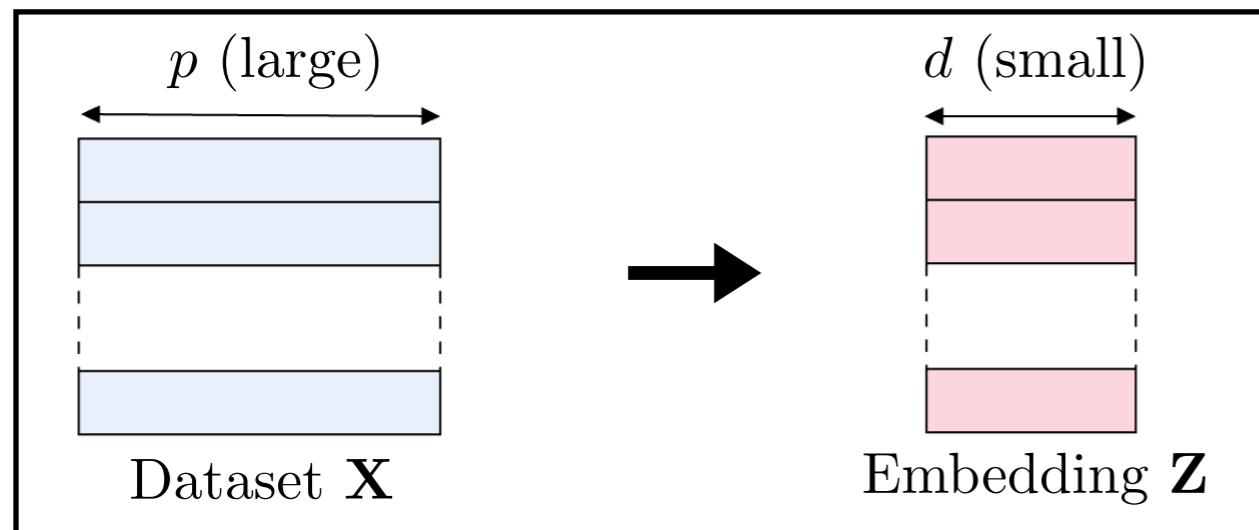


Cells

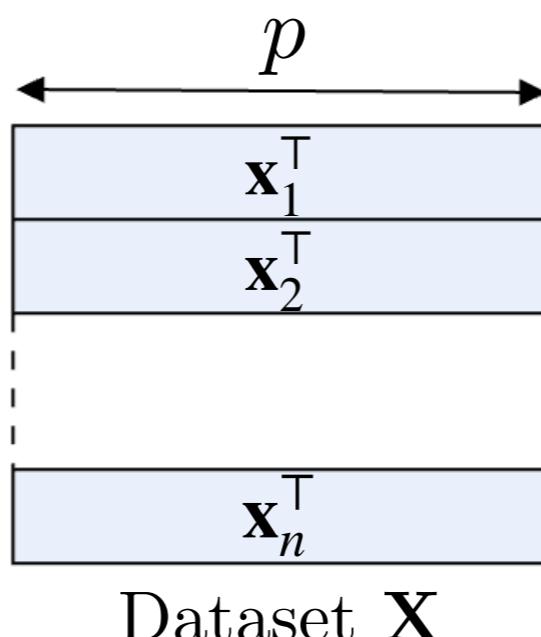
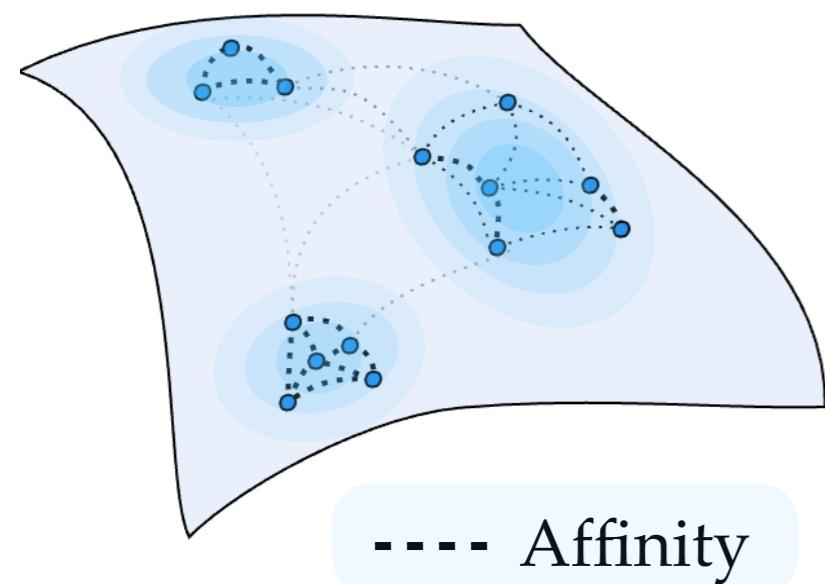


Coefficient (i, j) = similarity between \mathbf{x}_i and \mathbf{x}_j .

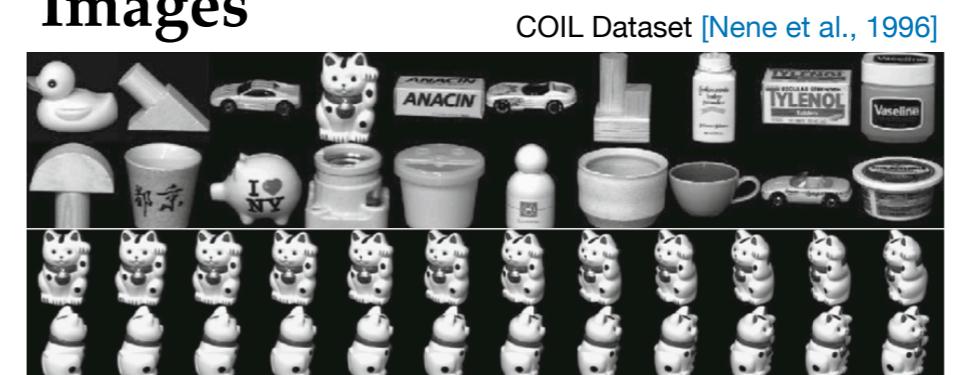
Dimension reduction



◆ Affinity Matrices



Images

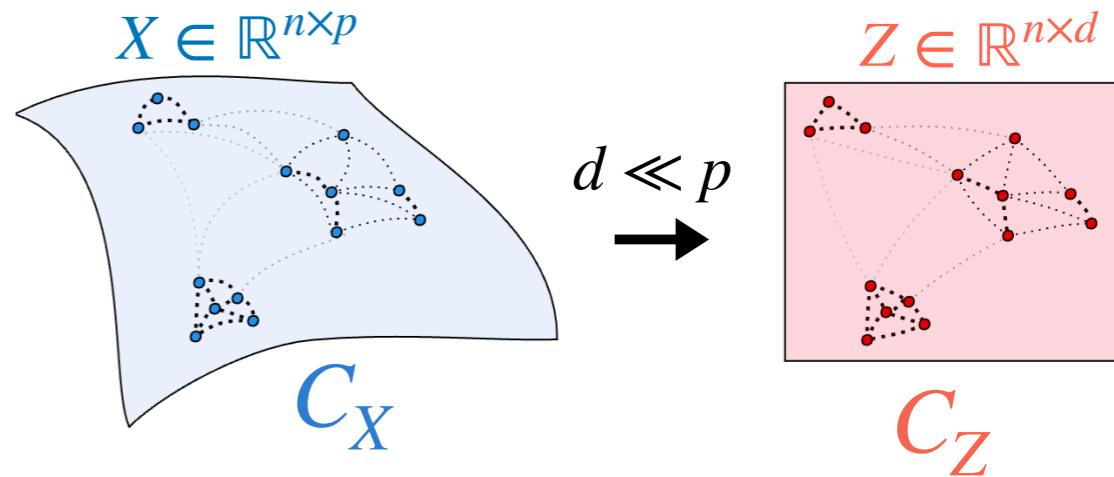


Cells



Requirement : a **meaningful metric** in input space !

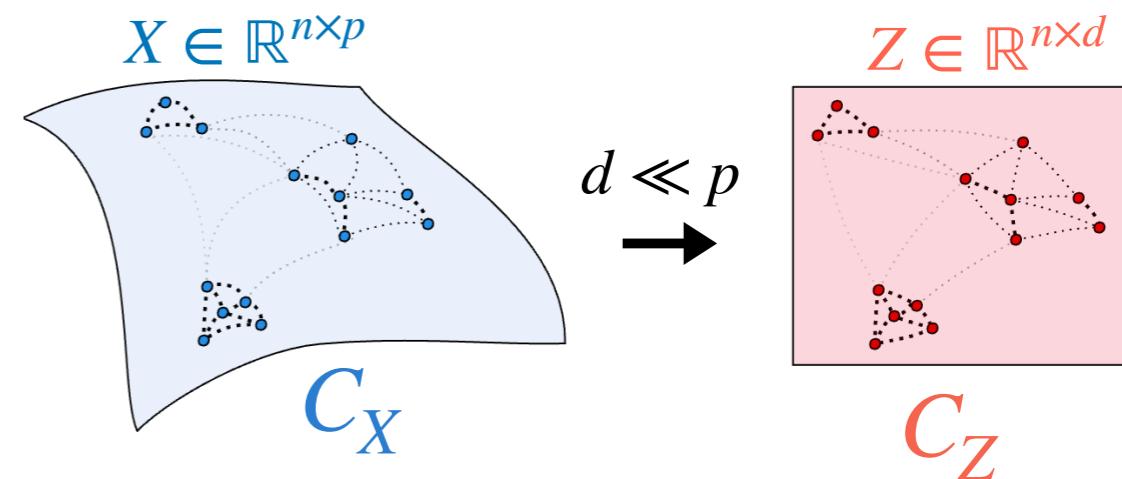
Dimension reduction



◆ A general optimization problem

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right) \text{ for some loss } L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

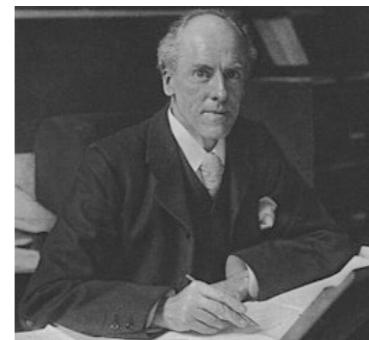
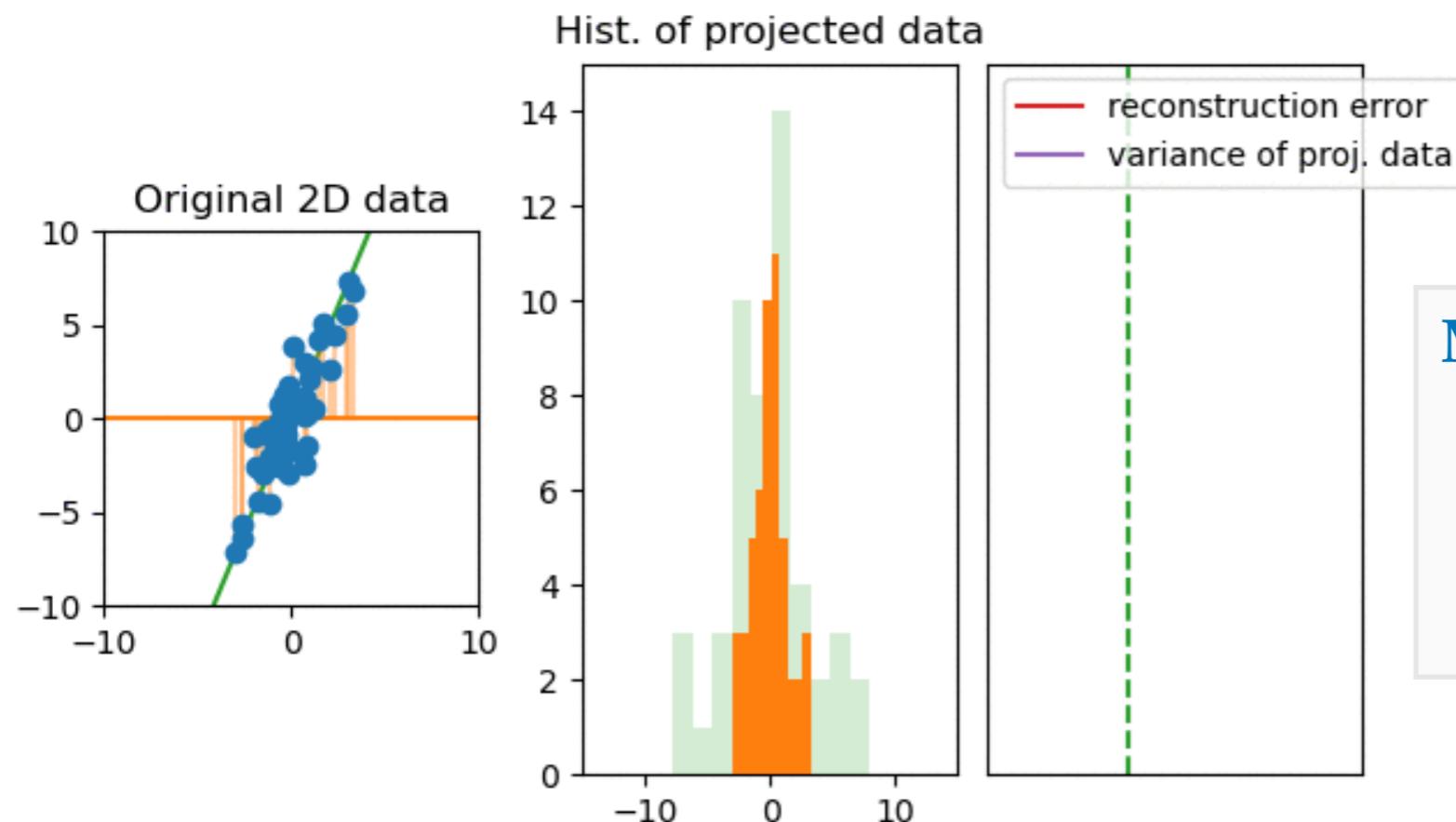
Dimension reduction



♦ A general optimization problem

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right) \text{ for some loss } L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

♦ Principal components analysis

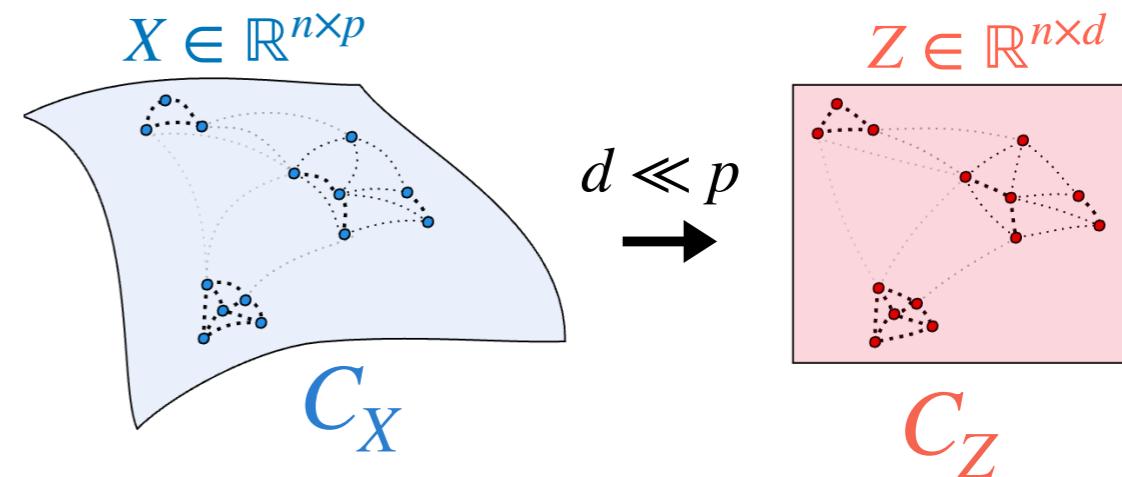


(Pearson, 1901)

Minimizing the reconstruction error

$$\min_{H: \dim(H)=d} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - P_H(\mathbf{x}_i)\|_2^2$$

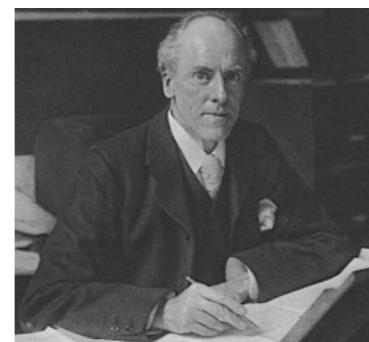
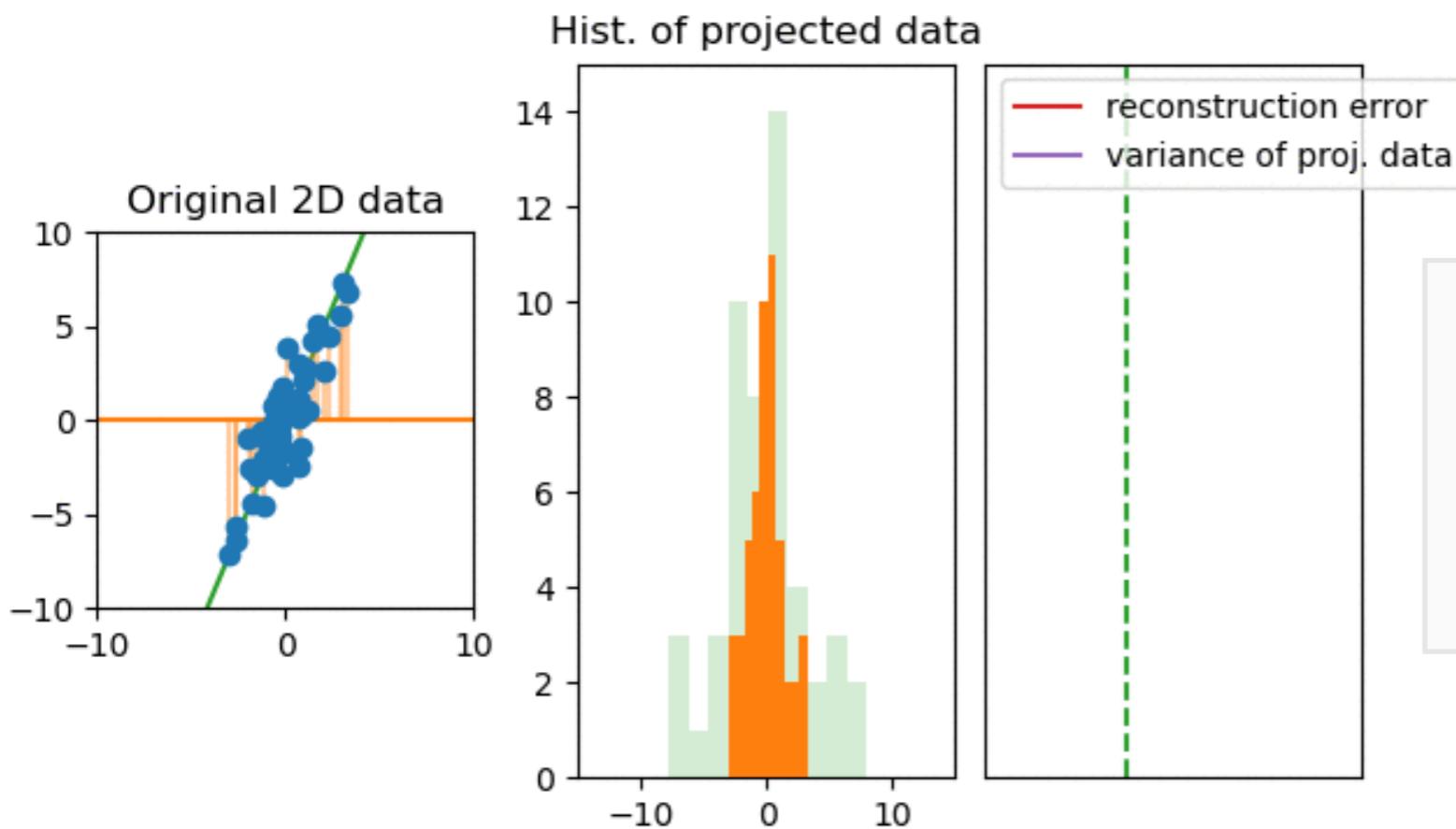
Dimension reduction



◆ A general optimization problem

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right) \text{ for some loss } L : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$$

◆ Principal components analysis



(Pearson, 1901)



(Torgerson, 1958)

Preserving the inner products

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left(\langle \mathbf{x}_i, \mathbf{x}_j \rangle - \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right)^2$$

$$Z \leftarrow \text{EVD}\left(\frac{1}{n} \mathbf{X} \mathbf{X}^\top\right)$$

Dimension reduction

♦ Spectral methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left([C_X]_{ij} - \langle z_i, z_j \rangle \right)^2$$

Dimension reduction

♦ Spectral methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left([C_X]_{ij} - \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right)^2 \xrightarrow[\substack{C_X \succeq 0 \\ \text{solution} \\ (\text{Eckart \& Young, 1936})}]{} Z^\star = (\sqrt{\lambda_1} \mathbf{v}_1, \dots, \sqrt{\lambda_d} \mathbf{v}_d)^\top$$

λ_i i-th largest eigenvalue of C_X
with eigenvector \mathbf{v}_i

$$[C_X]_{ij} = \langle \phi(X_i), \phi(X_j) \rangle_{\mathcal{H}}$$

Dimension reduction

♦ Spectral methods

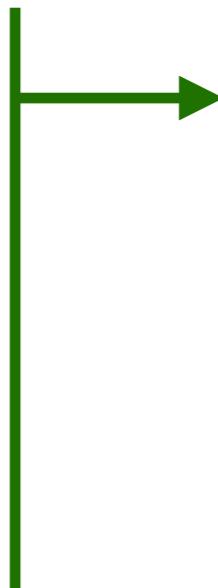
$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left([C_X]_{ij} - \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right)^2 \xrightarrow[\text{(Eckart & Young, 1936)}]{\begin{matrix} C_X \succeq 0 \\ \text{solution} \end{matrix}}$$

$Z^\star = (\sqrt{\lambda_1} \mathbf{v}_1, \dots, \sqrt{\lambda_d} \mathbf{v}_d)^\top$
 λ_i i-th largest eigenvalue of C_X
with eigenvector \mathbf{v}_i

♦ Kernel PCA $C_X \succeq 0$ $[C_X]_{ij} = \langle \phi(X_i), \phi(X_j) \rangle_{\mathcal{H}}$



(Schölkopf, 1997)



$$\text{PCA: } C_X = XX^\top \quad (Z \leftarrow \text{SVD}(X))$$

Dimension reduction

♦ Spectral methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left([C_X]_{ij} - \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right)^2 \quad \begin{matrix} C_X \succeq 0 \\ \text{solution} \\ (\text{Eckart \& Young, 1936}) \end{matrix}$$

$$Z^\star = (\sqrt{\lambda_1} \mathbf{v}_1, \dots, \sqrt{\lambda_d} \mathbf{v}_d)^\top$$

λ_i i-th largest eigenvalue of C_X
with eigenvector \mathbf{v}_i

♦ Kernel PCA $C_X \succeq 0$ $[C_X]_{ij} = \langle \phi(X_i), \phi(X_j) \rangle_{\mathcal{H}}$



(Schölkopf, 1997)

→ PCA: $C_X = XX^\top$ ($Z \leftarrow \text{SVD}(X)$)
→ (classical) Multidimensional scaling: $C_X = -\frac{1}{2}HD_XH$

Dimension reduction

◆ Spectral methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i=1}^n \left([C_X]_{ij} - \langle \mathbf{z}_i, \mathbf{z}_j \rangle \right)^2$$

C_X ≥ 0
solution →
(Eckart & Young, 1936)

$$Z^{\star} = (\sqrt{\lambda_1} \mathbf{v}_1, \dots, \sqrt{\lambda_d} \mathbf{v}_d)^{\top}$$

λ_i i-th largest eigenvalue of C_X
with eigenvector \mathbf{v}_i

◆ Kernel PCA $C_X \succeq 0$ $[C_X]_{ij} = \langle \phi(X_i), \phi(X_j) \rangle_{\mathcal{H}}$



(Schölkopf, 1997)

• PCA: $C_X = XX^\top$ ($Z \leftarrow \text{SVD}(X)$)

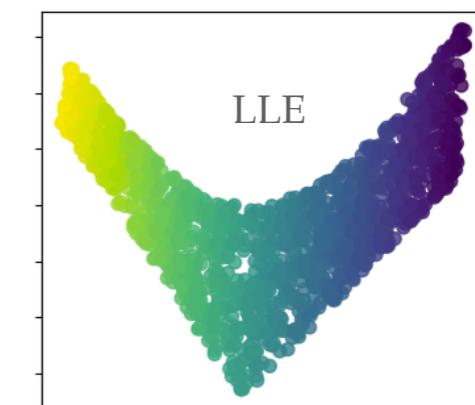
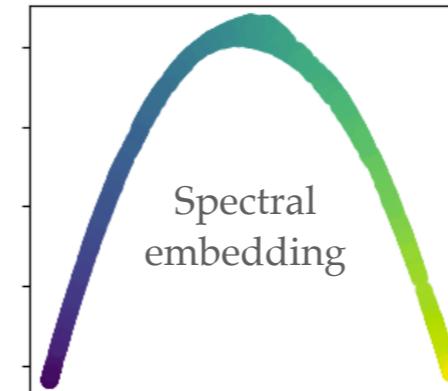
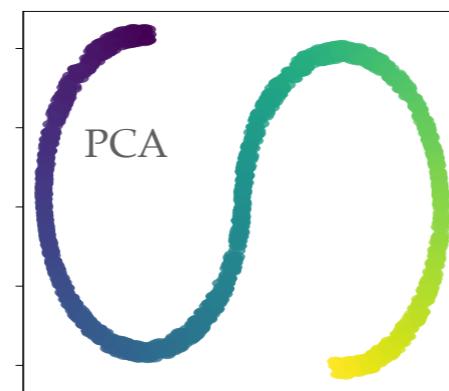
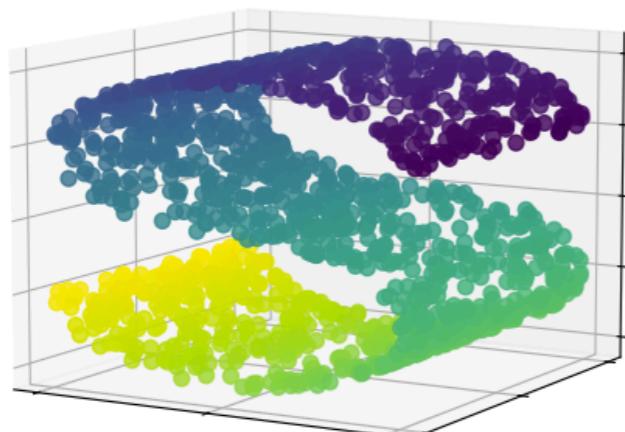
(classical) Multidimensional scaling: $C_X = -\frac{1}{\sigma} HD_X H$

- Laplacian Eigenmap (spectral embedding): $C_X = L_X^\dagger$
(Belkin & Niyogi, 2003)

Locally Linear Embedding, Diffusion Map ...

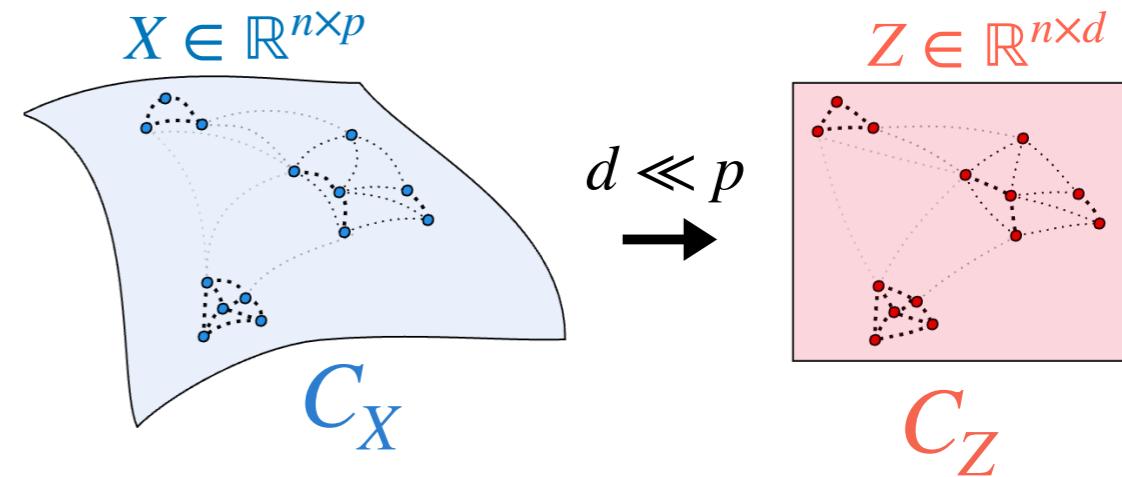
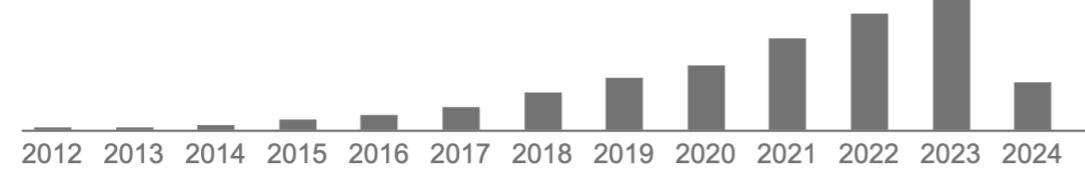
(Roweis & Saul, 2000)

(Coifman & Lafon, 2006)



Dimension reduction

Total citations Cited by 36223



♦ Neighbor embedding methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n \text{KL} \left([C_X]_{ij}, [C_Z]_{ij} \right)$$

TSNE



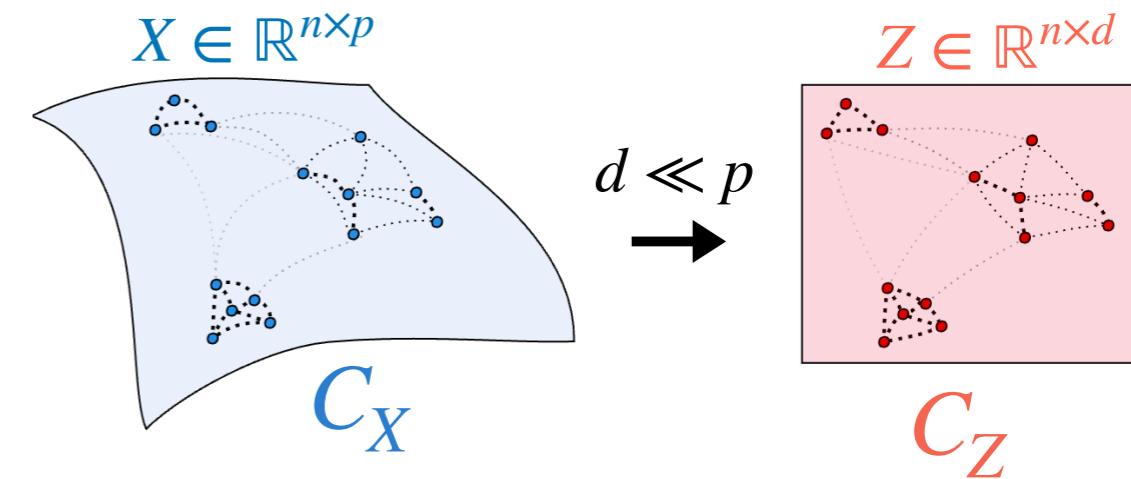
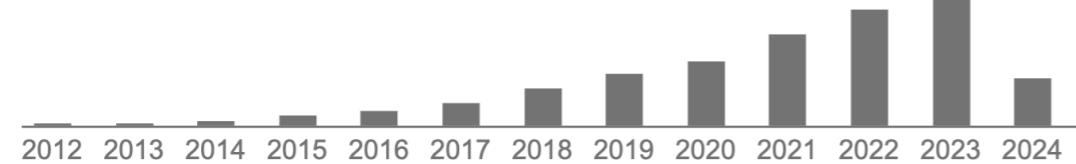
PCA



Kuzushiji - Japanese Letters

Dimension reduction

Total citations Cited by 36223



♦ Neighbor embedding methods

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n \text{KL} \left([C_X]_{ij}, [C_Z]_{ij} \right)$$



♦ Soft neighborhood graphs
(We'll see in a few slides)

TSNE



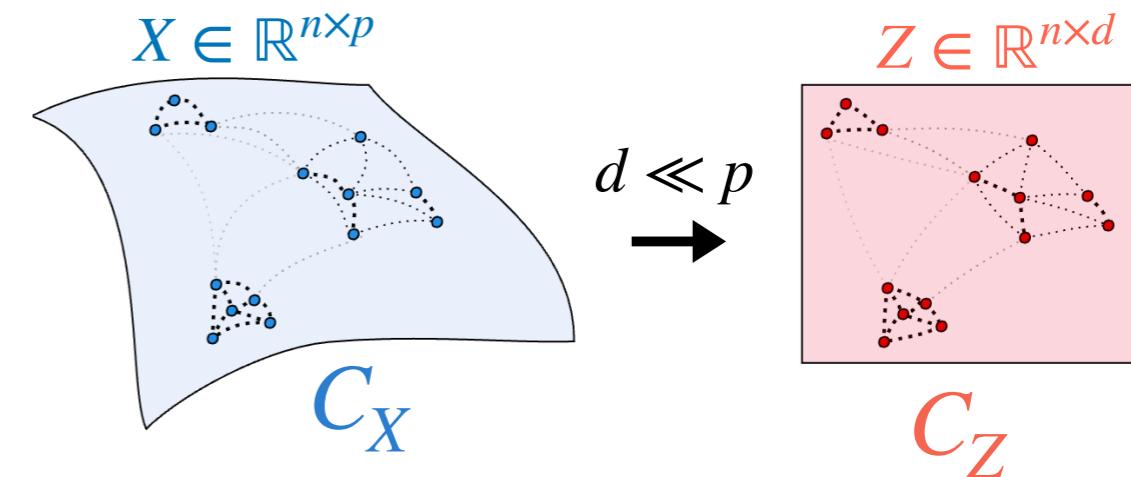
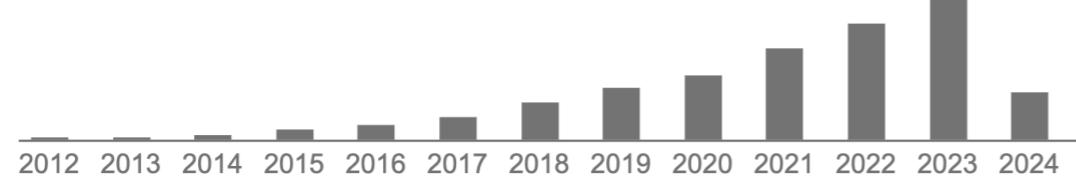
PCA



Kuzushiji - Japanese Letters

Dimension reduction

Total citations Cited by 36223

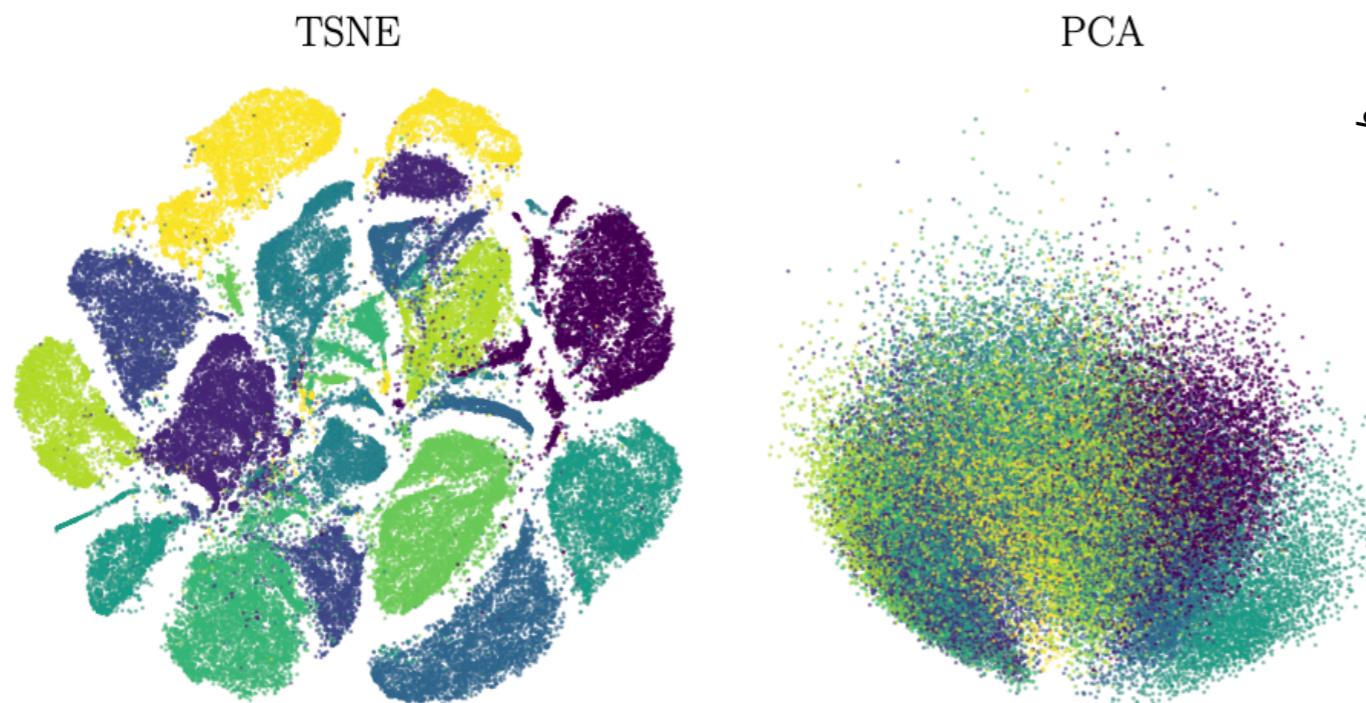


♦ Neighbor embedding methods

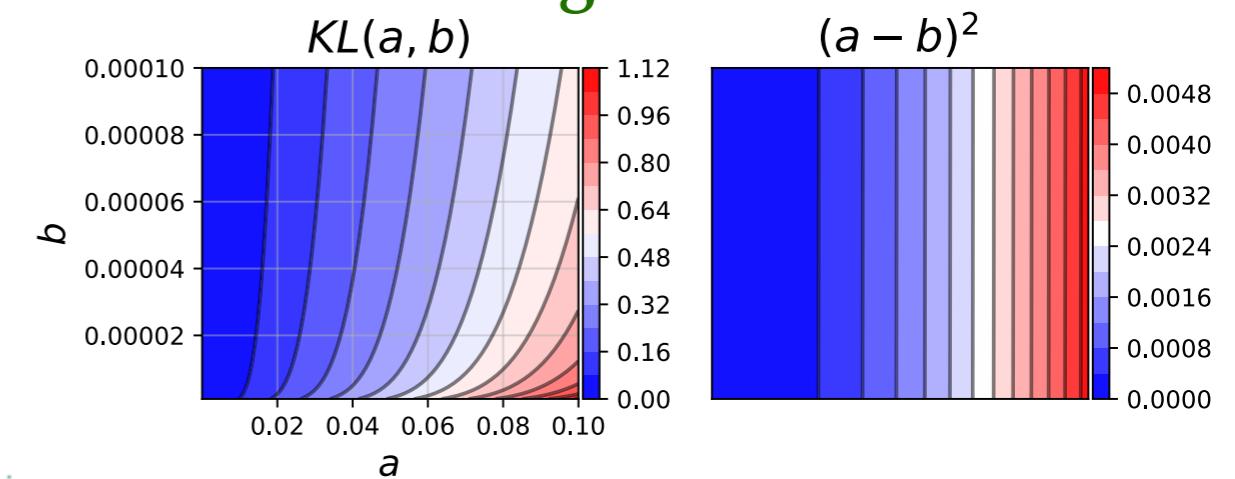
$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n \text{KL}\left([C_X]_{ij}, [C_Z]_{ij}\right)$$

When $\sum_{i,j} [C_X]_{ij} = \sum_{i,j} [C_Z]_{ij}$ (same mass)

$$\sim \frac{\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n [C_X]_{ij} \log\left(\frac{[C_X]_{ij}}{[C_Z]_{ij}}\right)}{\text{KL}}$$

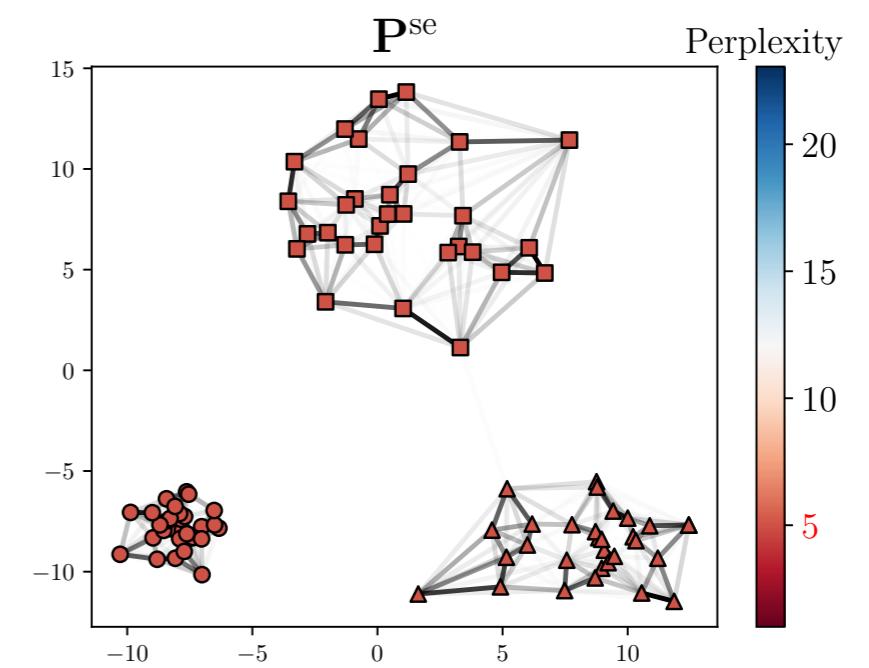
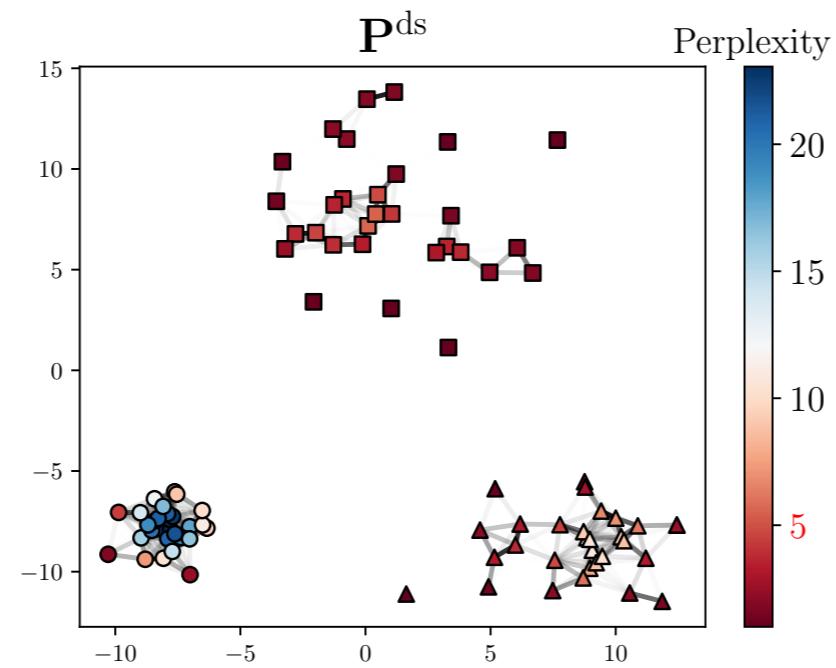


Kuzushiji - Japanese Letters



How to construct good affinities
for neighbor embedding ?

Adaptive Affinities for DR



Gaussian Affinity (or Gibbs kernel)

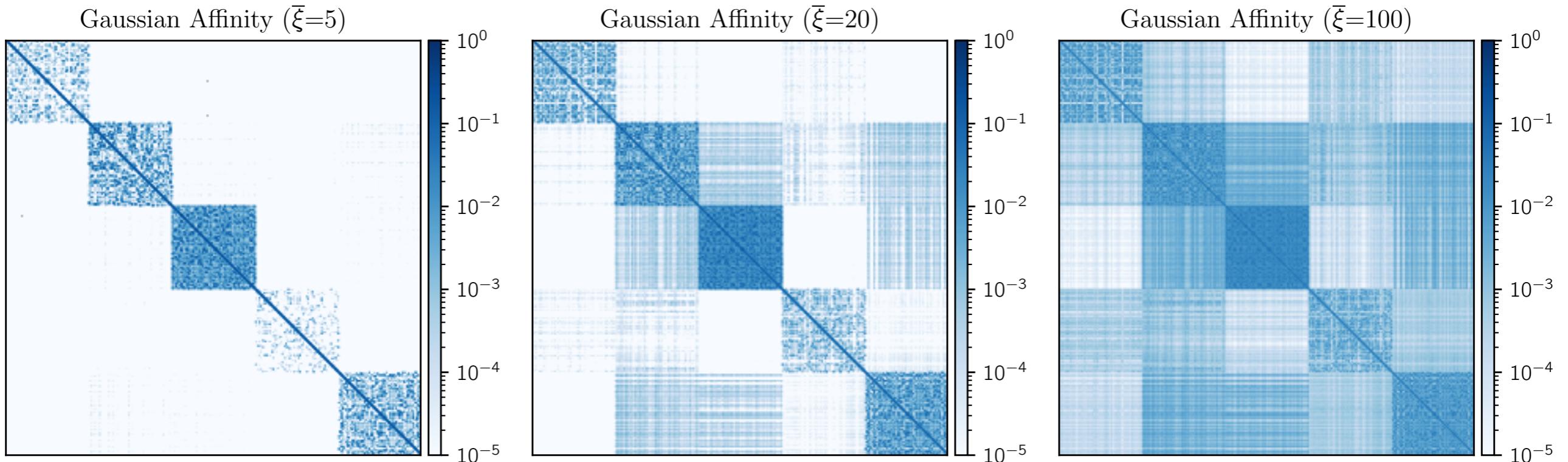


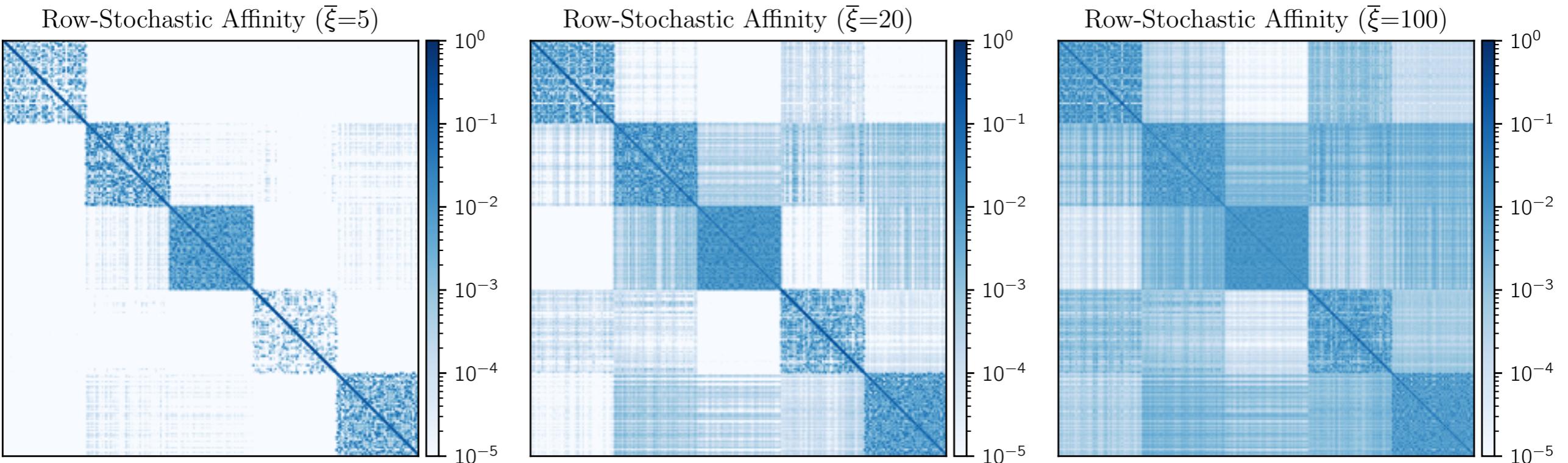
Fig : Affinity on 5 classes of the COIL Dataset [Nene et al., 1996]

Cost matrix: $\mathbf{C} \in \mathbb{R}_+^{n \times n}$ such that $\mathbf{C} = \mathbf{C}^\top$ and $C_{ij} = 0 \iff i = j$.

Example: $C_{ij} = \|\mathbf{x}_i - \mathbf{x}_j\|_2^2$.

Gibbs kernel : $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$.

ℓ_1 Norm - Row Stochastic

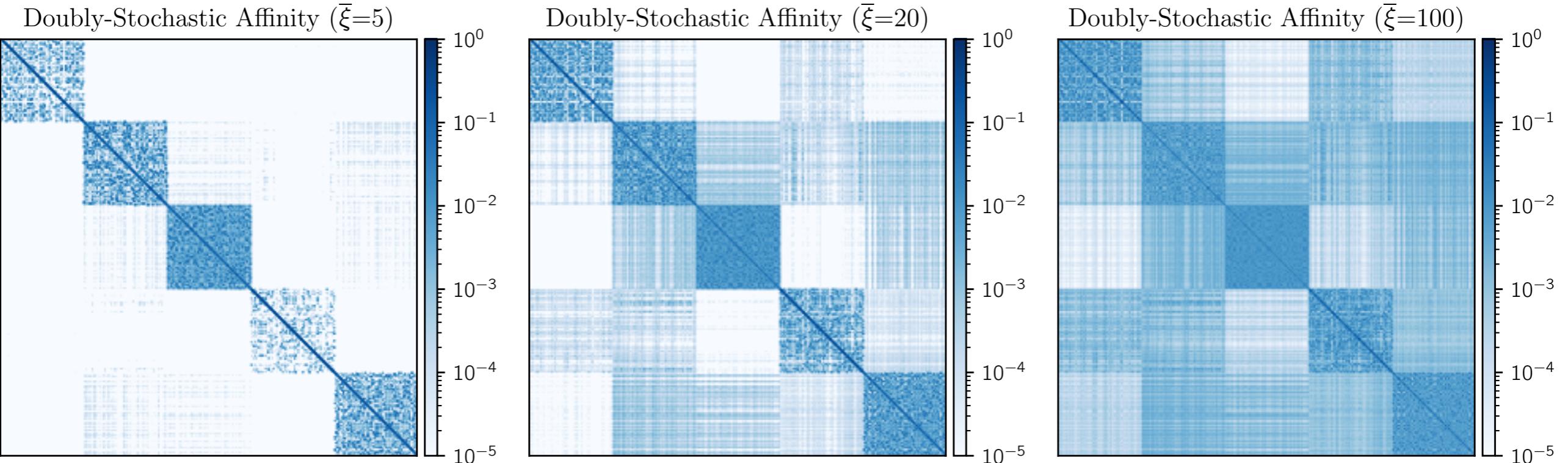


init $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K1})^{-1} \mathbf{K}$$

normalize rows

ℓ_1 Norm - Doubly Stochastic



Sinkhorn Algorithm

init $\mathbf{K} = \exp(-\mathbf{C}/\sigma)$

While not converged:

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K}\mathbf{1})^{-1}\mathbf{K}$$

normalize rows

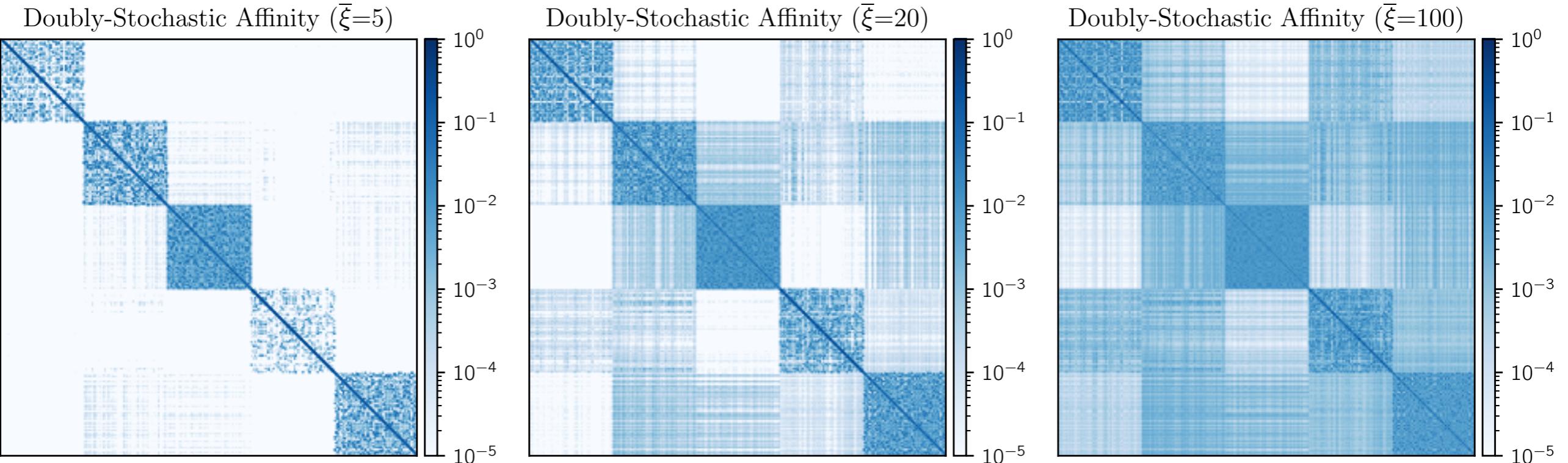
$$\mathbf{K} \leftarrow \mathbf{K} \text{diag}(\mathbf{1}\mathbf{K})^{-1}$$

normalize columns

converges to
→

$$\mathcal{DS} = \{\mathbf{P} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{P}^\top \mathbf{1} = \mathbf{1}\}.$$

ℓ_1 Norm - Doubly Stochastic



Sinkhorn Algorithm

While not converged:

$$\mathbf{K} \leftarrow \text{diag}(\mathbf{K}\mathbf{1})^{-1}\mathbf{K} \quad \# \text{ normalize rows}$$

$$\mathbf{K} \leftarrow \mathbf{K} \text{diag}(\mathbf{1}\mathbf{K})^{-1} \quad \# \text{ normalize columns}$$

converges to $\mathcal{DS} = \{\mathbf{P} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{P}^\top\mathbf{1} = \mathbf{1}\}$.

$$\text{init } \mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

→ Spectral clustering
(Zass & Sashua, 2006)

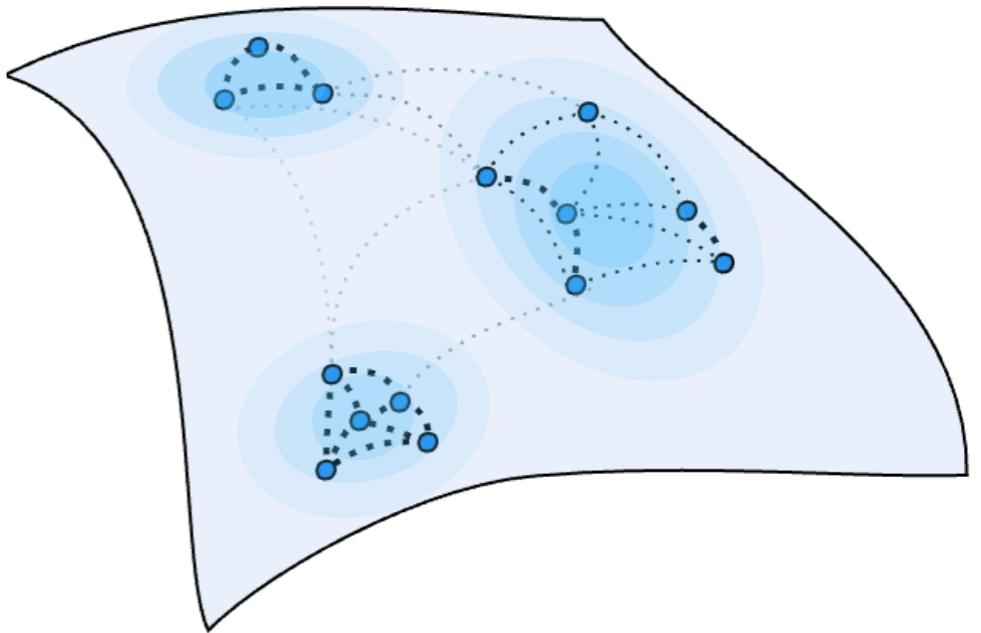
→ Transformers
(Sander & al., 2021)

→ Dim. Reduction
(Yu & al., 2016)

Entropic Affinity

Data has **varying noise levels**.

We can control the entropy in each point with **adaptive bandwidths**.



Definition [Hinton, Roweis 2002]

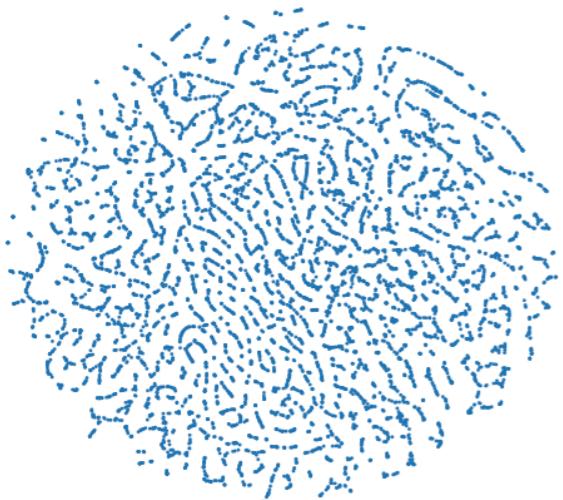
$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_\ell \exp(-C_{i\ell}/\varepsilon_i^*)}$$

$$\text{with } \varepsilon_i^* \in \mathbb{R}_+^* \text{ s.t. } H(\mathbf{P}_{i:}^e) = \log \xi + 1.$$

$H(p) = -\langle p, \log p - 1 \rangle$ is the Shannon entropy.

$\xi \in [\![1, n]\!]$ is the **perplexity** parameter.

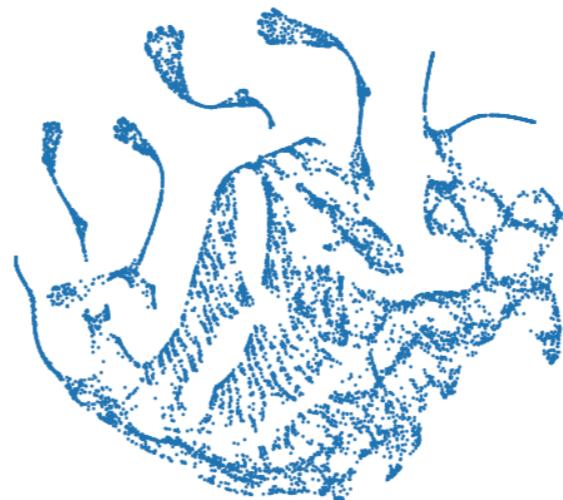
Perplexity = 5



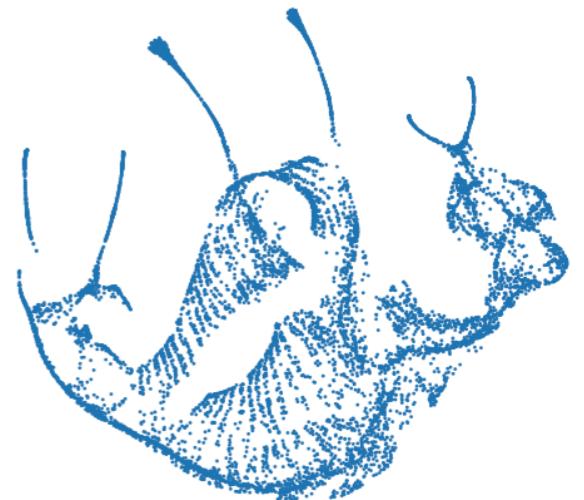
Perplexity = 30



Perplexity = 100



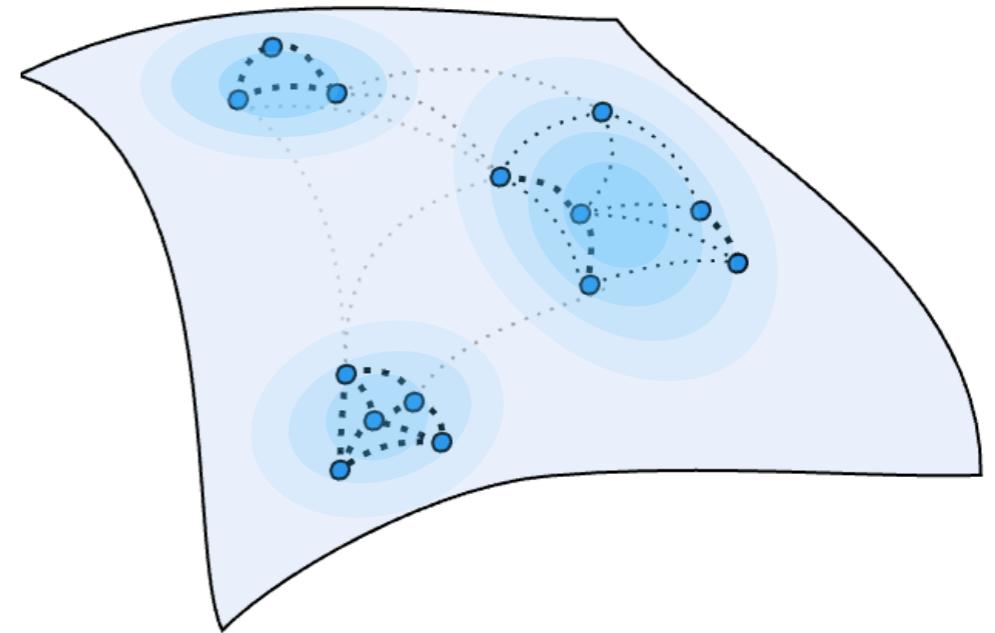
Perplexity = 300



Entropic Affinity

Data has **varying noise levels**.

We can control the entropy in each point with **adaptive bandwidths**.

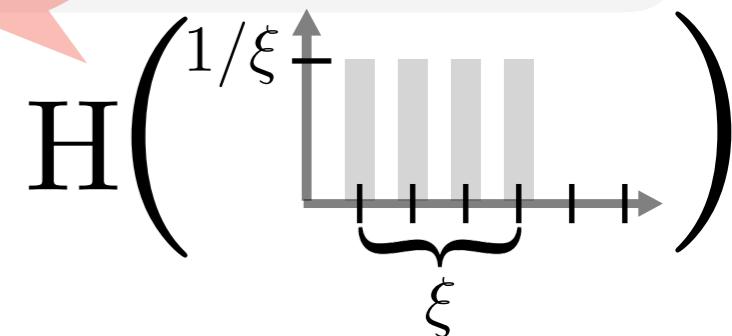


Definition [Hinton, Roweis 2002]

$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_\ell \exp(-C_{i\ell}/\varepsilon_i^*)}$$

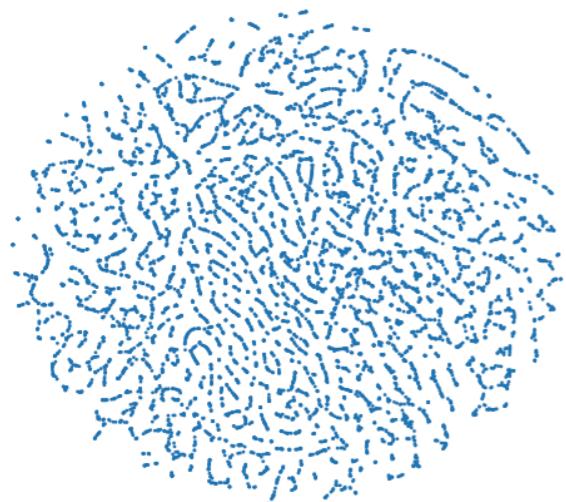
with $\varepsilon_i^* \in \mathbb{R}_+^*$ s.t. $H(\mathbf{P}_{i:}^e) = \log \xi + 1$.

$H(p) = -\langle p, \log p - 1 \rangle$ is the Shannon entropy.



$\xi \in [\![1, n]\!]$ is the **perplexity** parameter.

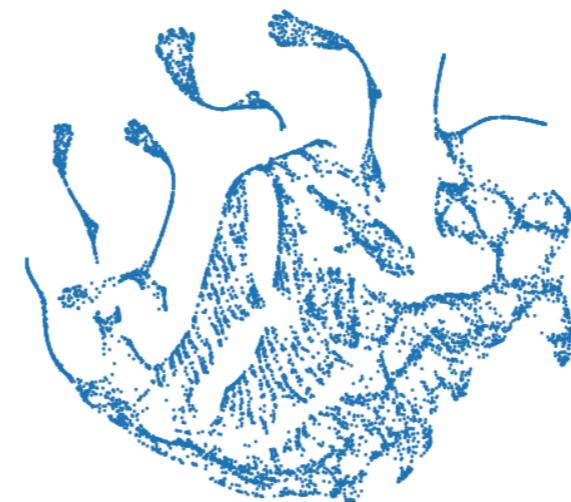
Perplexity = 5



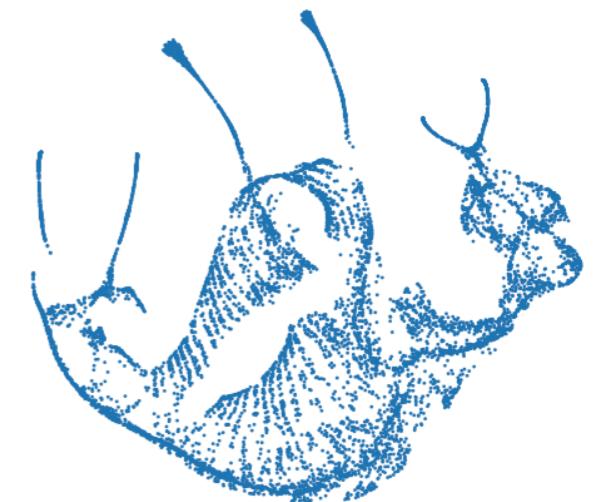
Perplexity = 30



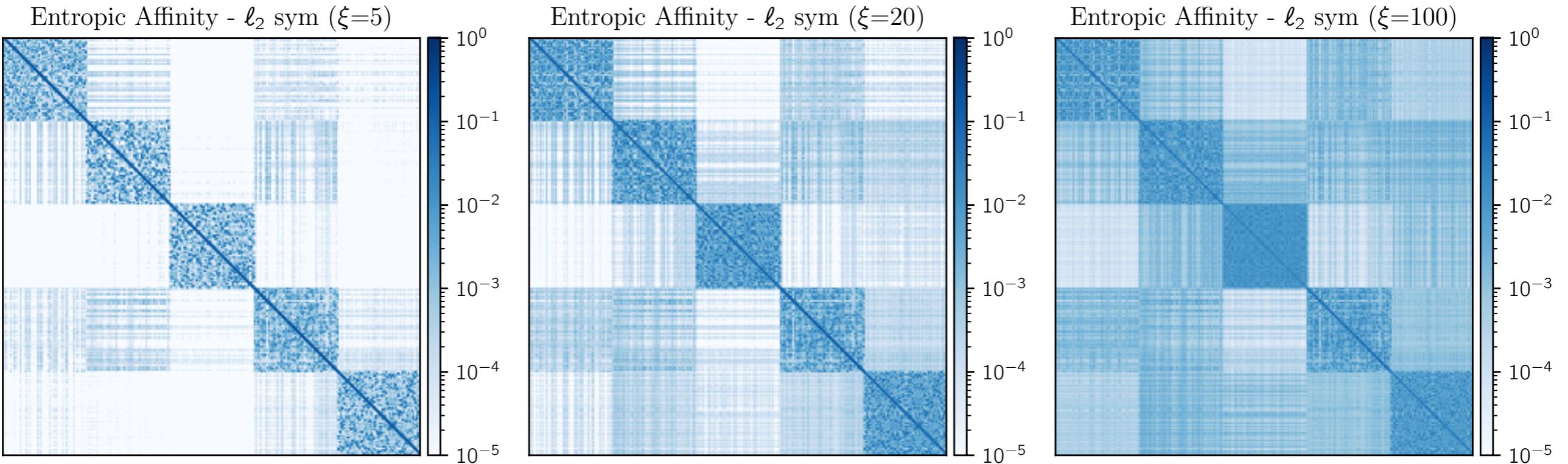
Perplexity = 100



Perplexity = 300



Entropic Affinity



Definition [Hinton, Roweis 2002]

$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_\ell \exp(-C_{i\ell}/\varepsilon_i^*)}$$

$$\text{with } \varepsilon_i^* \in \mathbb{R}_+^* \text{ s.t. } H(\mathbf{P}_{i:}^e) = \log \xi + 1 .$$

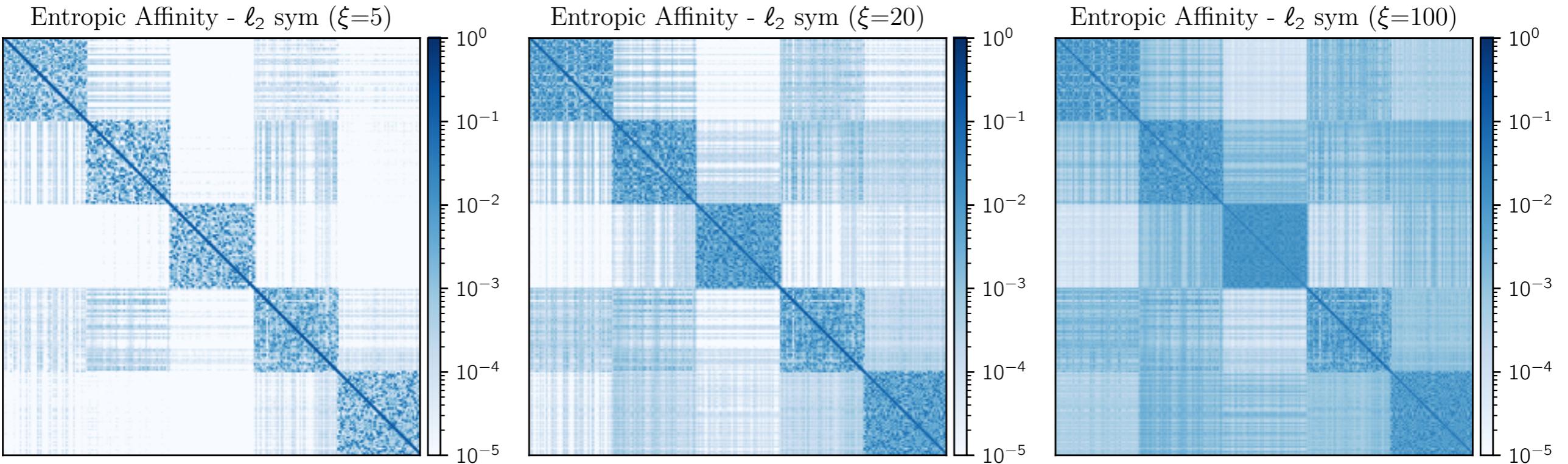
\mathbf{P}^e is not symmetric.

→ $\overline{\mathbf{P}}^e = \frac{1}{2}(\mathbf{P}^e + \mathbf{P}^{e\top})$ is used in practice. [Van der Maaten and Hinton, 2008]

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

t-SNE algorithm

Entropic Affinity



Definition [Hinton, Roweis 2002]

$$\forall i, \forall j, P_{ij}^e = \frac{\exp(-C_{ij}/\varepsilon_i^*)}{\sum_\ell \exp(-C_{i\ell}/\varepsilon_i^*)}$$

$$\text{with } \varepsilon_i^* \in \mathbb{R}_+^* \text{ s.t. } H(\mathbf{P}_{i:}^e) = \log \xi + 1 .$$

\mathbf{P}^e is not symmetric.

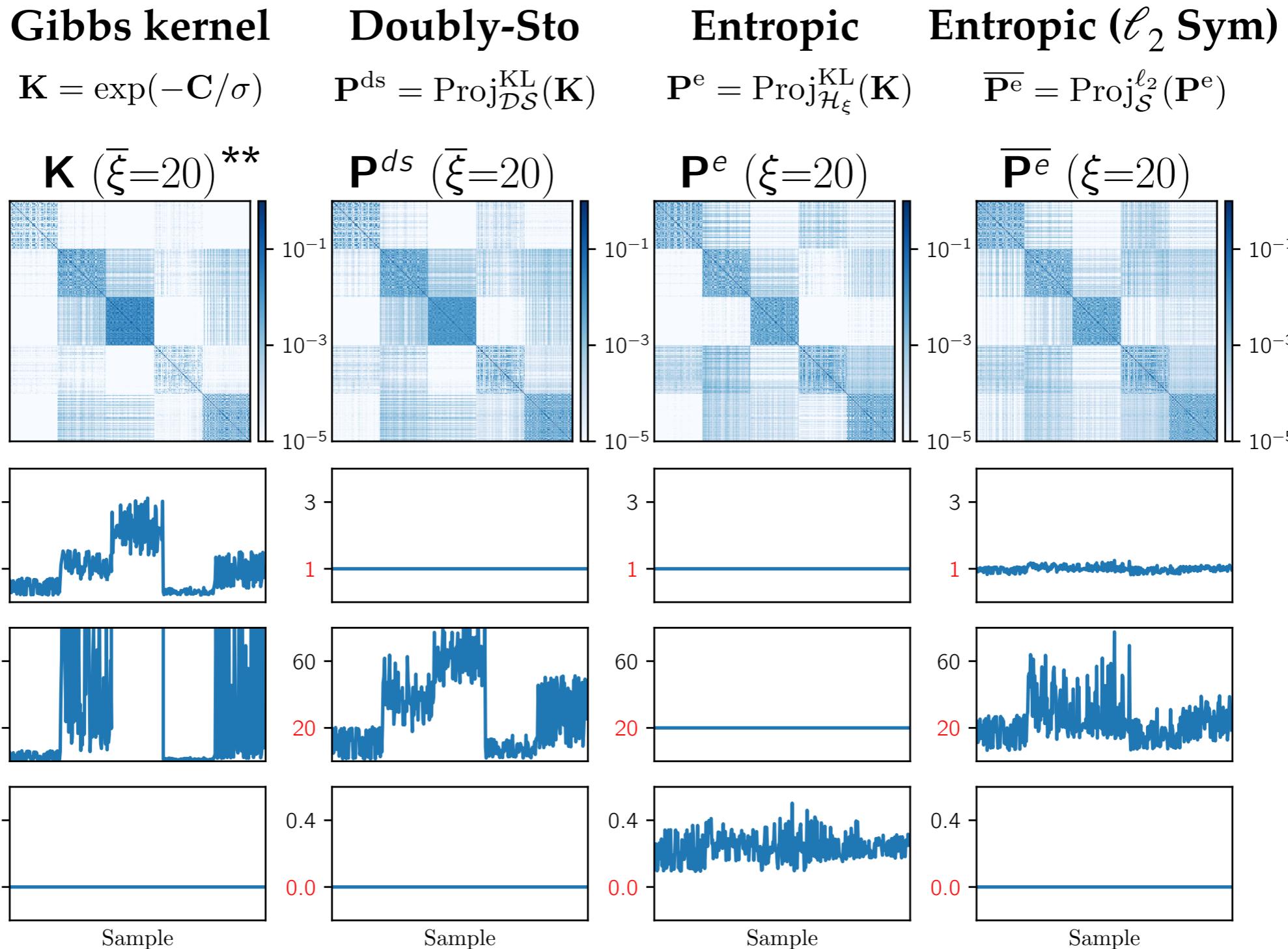
→ $\overline{\mathbf{P}}^e = \frac{1}{2}(\mathbf{P}^e + \mathbf{P}^{e\top})$ is used in practice. [Van der Maaten and Hinton, 2008]

$$\overline{\mathbf{P}}^e = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

t-SNE algorithm

Breaks the construction of entropic affinities.

Affinity Panorama*



* On 5 classes of the COIL Dataset [Nene et al., 1996]

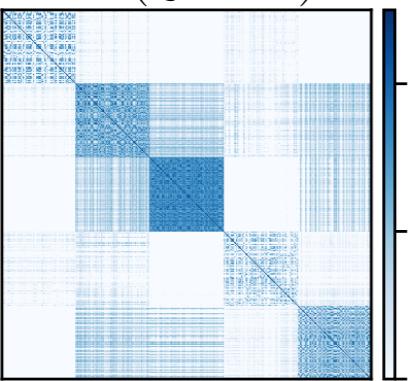
** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Affinity Panorama*

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

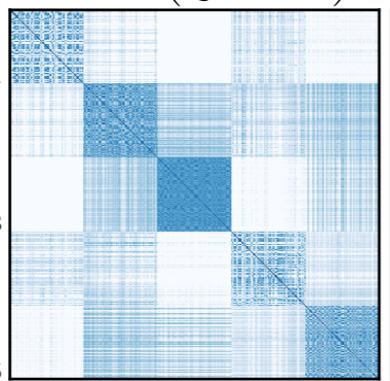
\mathbf{K} ($\bar{\xi}=20$)**



Doubly-Sto

$$\mathbf{P}^{\text{ds}} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

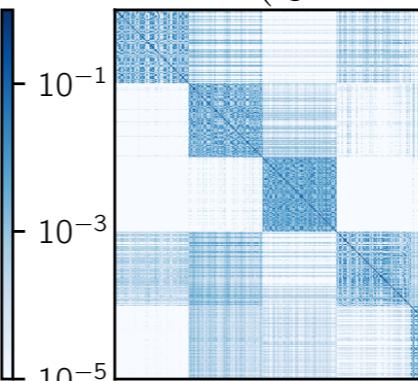
\mathbf{P}^{ds} ($\bar{\xi}=20$)



Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

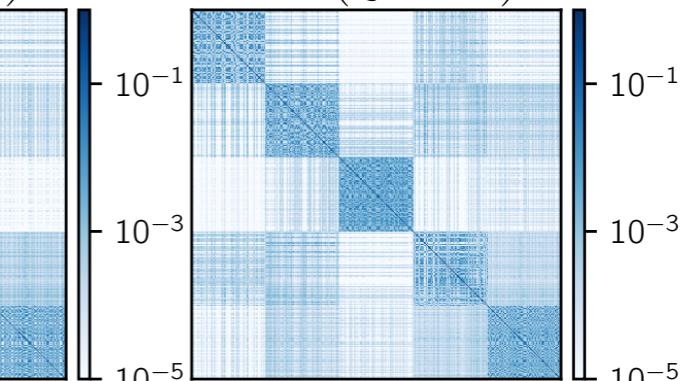
\mathbf{P}^e ($\xi=20$)



Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}^e} = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

$\overline{\mathbf{P}^e}$ ($\xi=20$)



Affinity

ℓ_1 norm

Perplexity

Symmetry

Sample

Sample

Sample

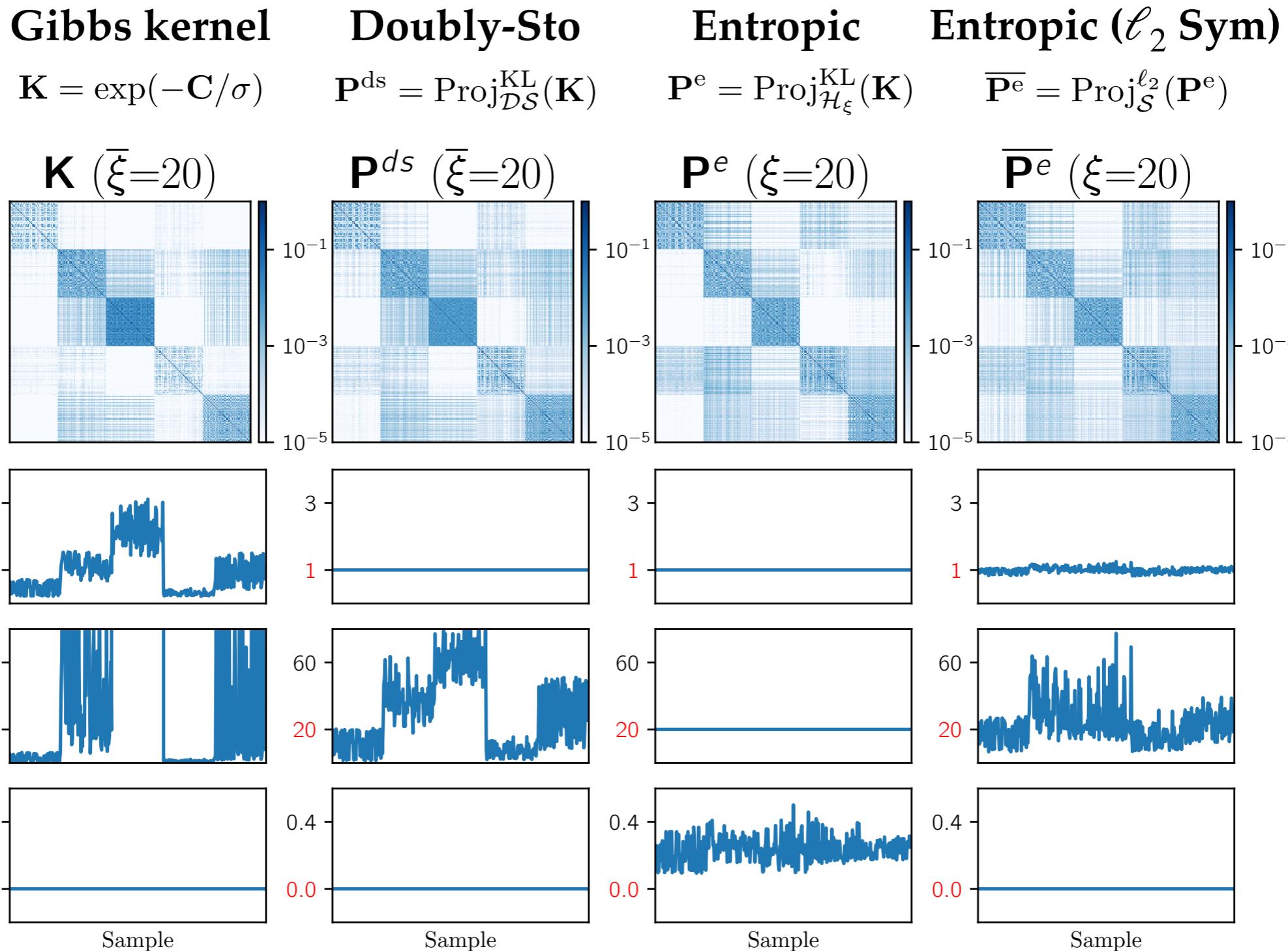
Sample

* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Entropies not controlled.

Affinity Panorama



Can we control ℓ_1 norm, entropy and symmetry ?

| Symmetric Entropic Affinity

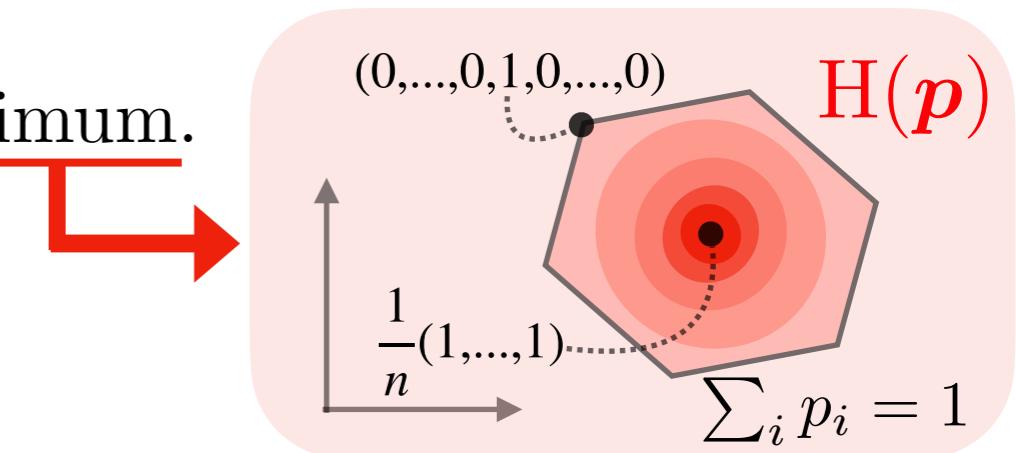
$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Entropic Affinity as OT

$$\mathbf{P}^e = \arg \min_{\mathbf{P} \in \mathcal{H}_\xi} \langle \mathbf{P}, \mathbf{C} \rangle.$$

| The constraints in \mathcal{H}_ξ are saturated at the optimum.

| Symmetric matrices $\mathcal{S} = \{\mathbf{P} \text{ s.t. } \mathbf{P} = \mathbf{P}^\top\}$.



Definition

$$\mathbf{P}^{se} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

OURS

| Symmetric Entropic Affinity

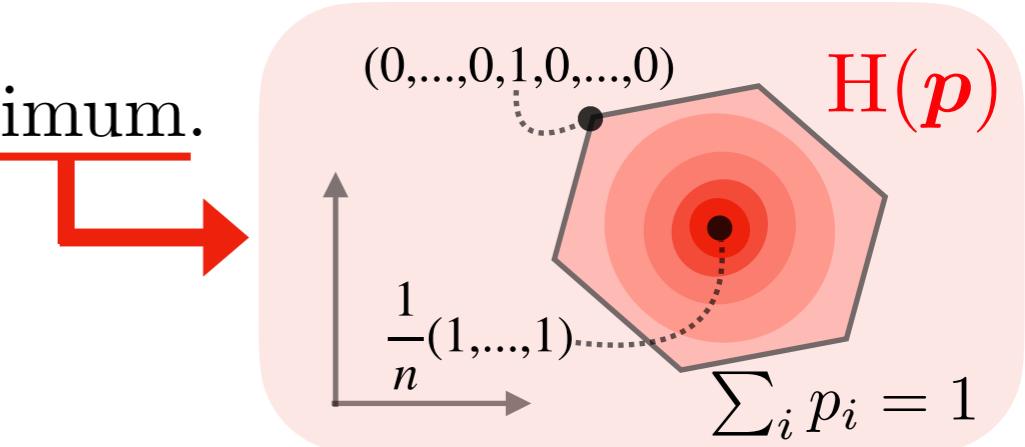
$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Entropic Affinity as OT

$$\mathbf{P}^e = \arg \min_{\mathbf{P} \in \mathcal{H}_\xi} \langle \mathbf{P}, \mathbf{C} \rangle.$$

| The constraints in \mathcal{H}_ξ are saturated at the optimum.

| Symmetric matrices $\mathcal{S} = \{\mathbf{P} \text{ s.t. } \mathbf{P} = \mathbf{P}^\top\}$.



Definition

$$\mathbf{P}^{\text{se}} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

Enforce Symmetry

OURS

| Symmetric Entropic Affinity

$$\mathcal{H}_\xi := \{\mathbf{P} \in \mathbb{R}_+^{n \times n} \text{ s.t. } \mathbf{P}\mathbf{1} = \mathbf{1} \text{ and } \forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1\}$$

Definition

$$\mathbf{P}^{\text{se}} := \arg \min_{\mathbf{P} \in \mathcal{H}_\xi \cap \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle.$$

\mathcal{S}

Enforce Symmetry

OURS

Property

For at least $n - 1$ indices $i \in \llbracket n \rrbracket$, it holds $H(\mathbf{P}_{i:}^{\text{se}}) = \log \xi + 1$.

| In practice, we have n saturated entropies.

Dual Ascent

$$\mathbf{P}^{\text{se}} = \exp((\boldsymbol{\lambda}^* \oplus \boldsymbol{\lambda}^* - 2\mathbf{C}) \oslash (\boldsymbol{\gamma}^* \oplus \boldsymbol{\gamma}^*))$$

where $\boldsymbol{\lambda}^*$ and $\boldsymbol{\gamma}^*$ are computed using dual ascent.

Affinity Panorama *

t-SNE

OURS

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

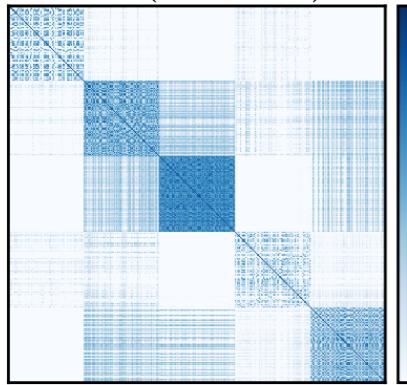
Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}^e} = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

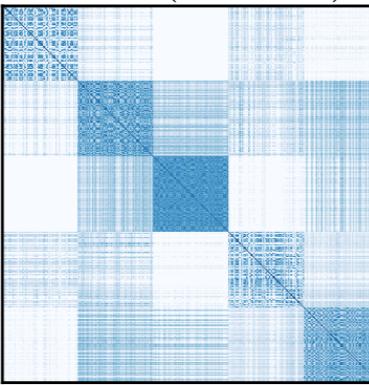
Sym-Entropic

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

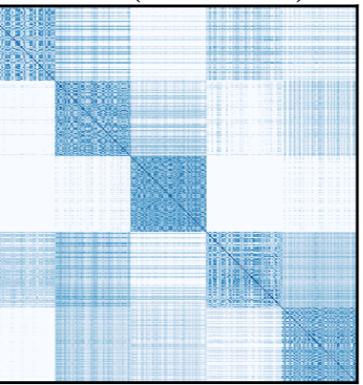
\mathbf{K} ($\bar{\xi}=20$) **



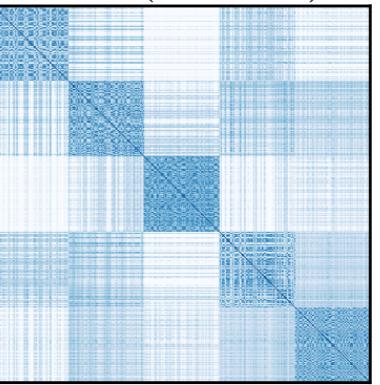
\mathbf{P}^{ds} ($\bar{\xi}=20$)



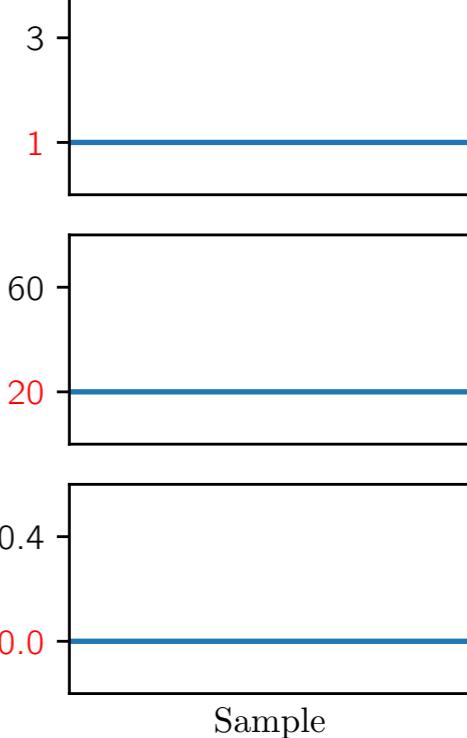
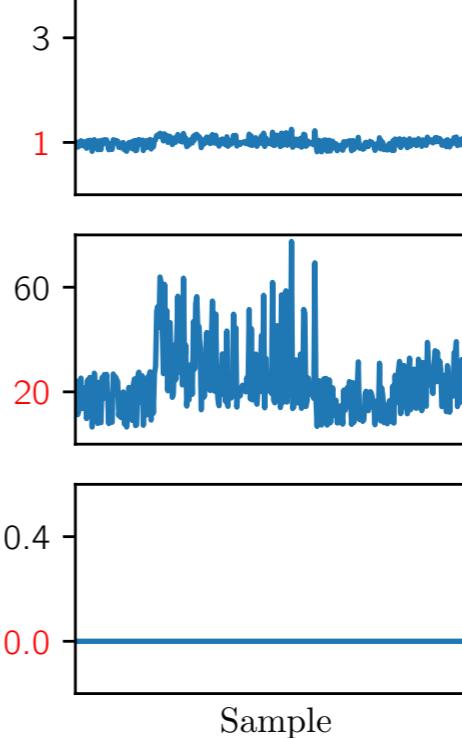
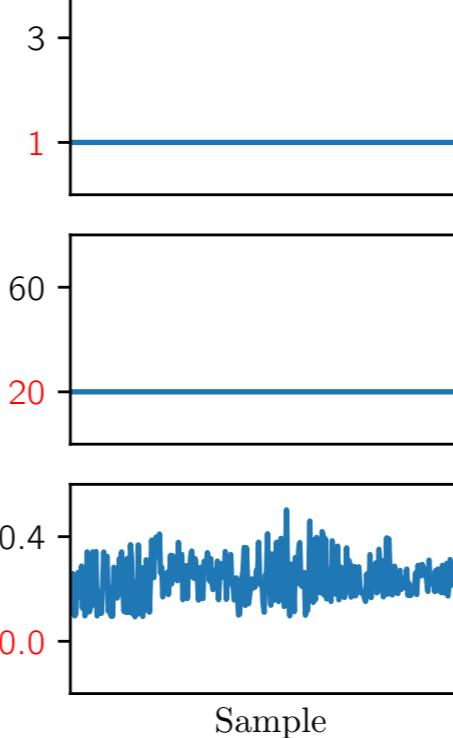
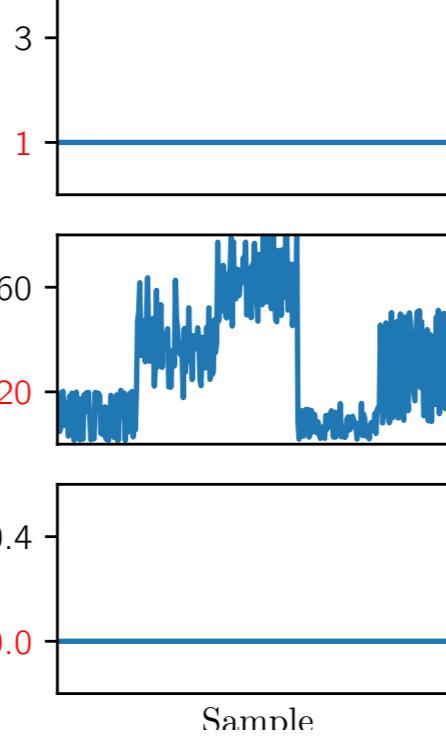
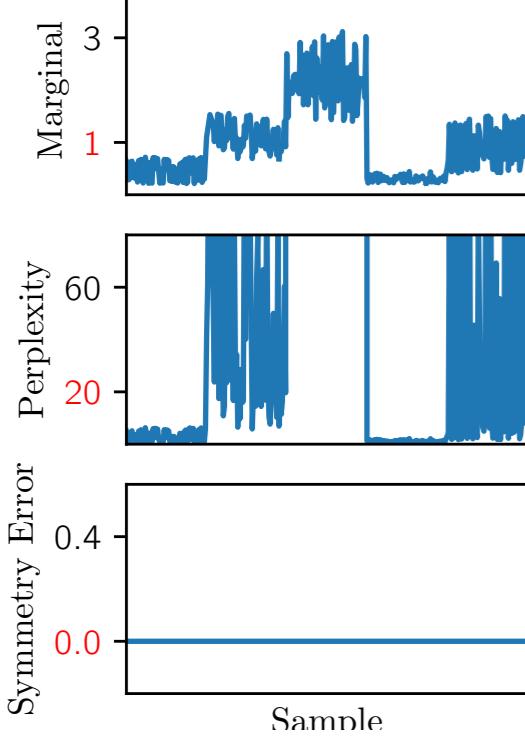
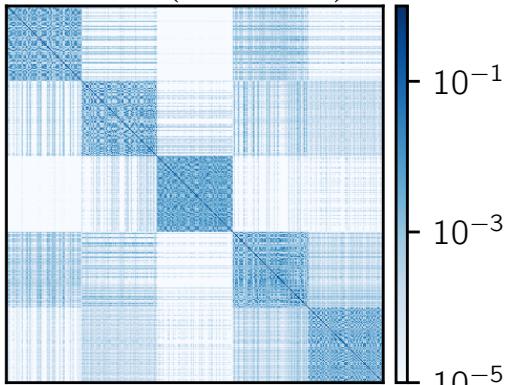
\mathbf{P}^e ($\xi=20$)



$\overline{\mathbf{P}^e}$ ($\xi=20$)



\mathbf{P}^{se} ($\xi=20$)



* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Affinity Panorama *

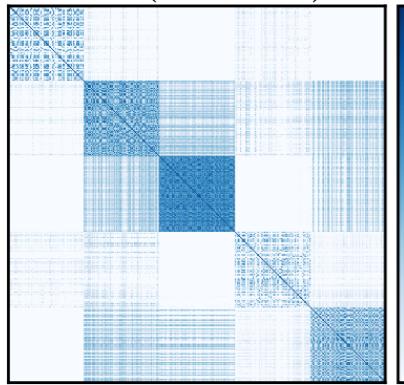
t-SNE

OURS

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

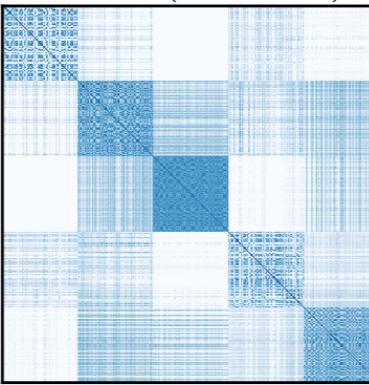
\mathbf{K} ($\bar{\xi}=20$) **



Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

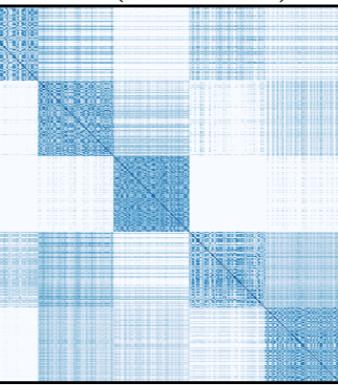
\mathbf{P}^{ds} ($\bar{\xi}=20$)



Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

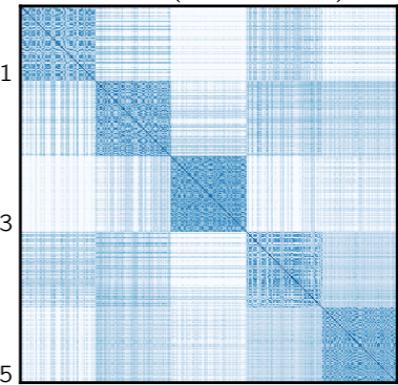
\mathbf{P}^e ($\xi=20$)



Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}^e} = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

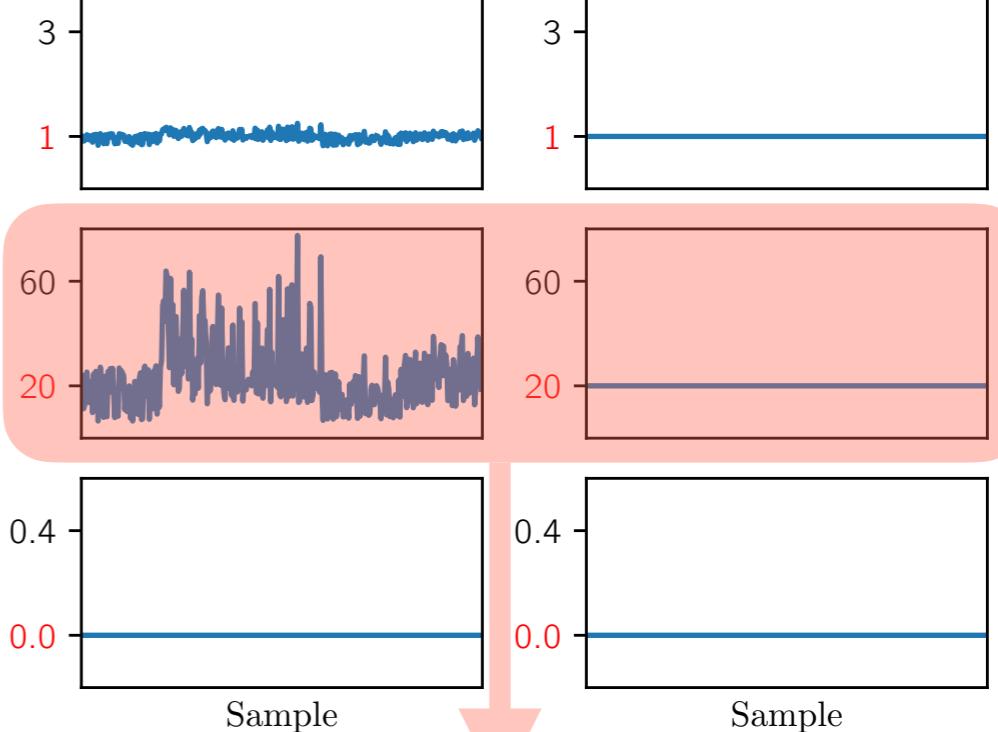
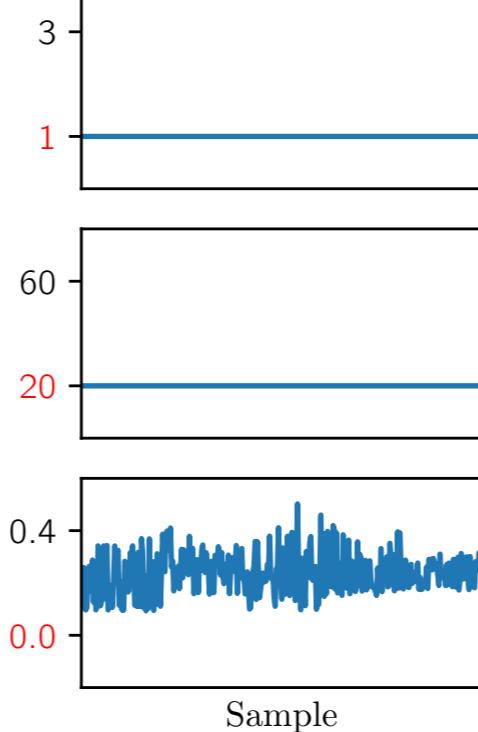
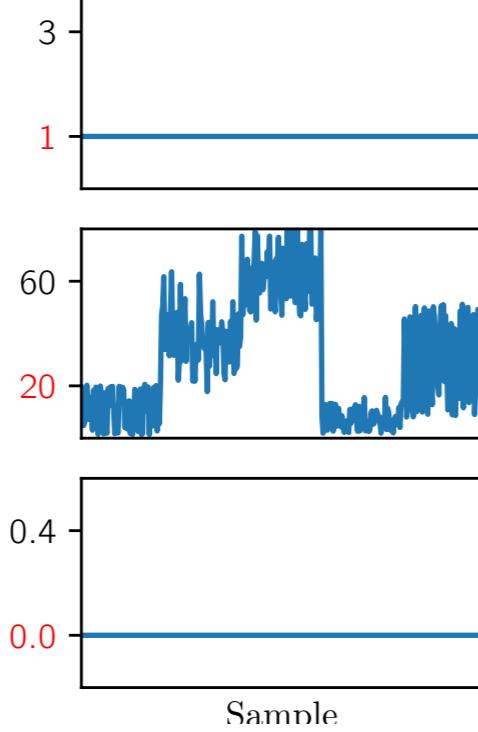
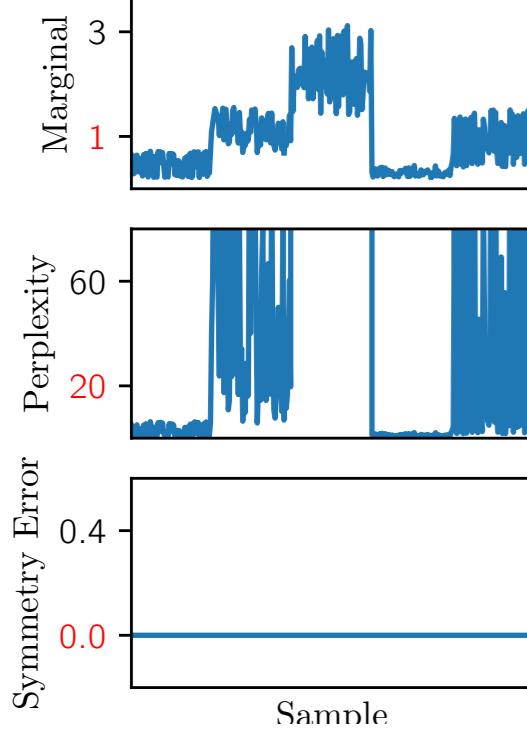
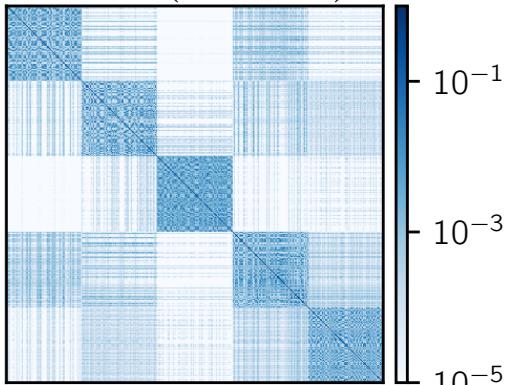
$\overline{\mathbf{P}^e}$ ($\xi=20$)



Sym-Entropic

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

\mathbf{P}^{se} ($\xi=20$)



Effective control over entropies.

* On 5 classes of the COIL Dataset [Nene et al., 1996]

** $\bar{\xi}$ is average perplexity \rightarrow same global entropy as with $\xi = \bar{\xi}$.

Affinity Panorama *

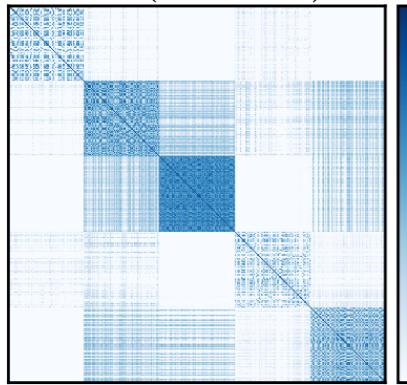
t-SNE

OURS

Gibbs kernel

$$\mathbf{K} = \exp(-\mathbf{C}/\sigma)$$

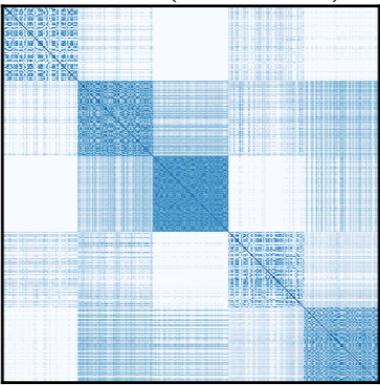
\mathbf{K} ($\xi=20$) **



Doubly-Sto

$$\mathbf{P}^{ds} = \text{Proj}_{\mathcal{DS}}^{\text{KL}}(\mathbf{K})$$

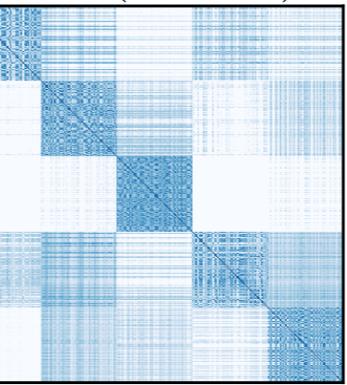
\mathbf{P}^{ds} ($\xi=20$)



Entropic

$$\mathbf{P}^e = \text{Proj}_{\mathcal{H}_\xi}^{\text{KL}}(\mathbf{K})$$

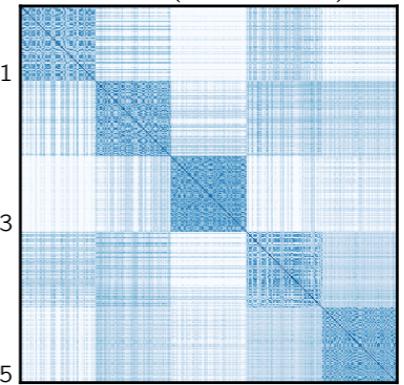
\mathbf{P}^e ($\xi=20$)



Entropic (ℓ_2 Sym)

$$\overline{\mathbf{P}^e} = \text{Proj}_{\mathcal{S}}^{\ell_2}(\mathbf{P}^e)$$

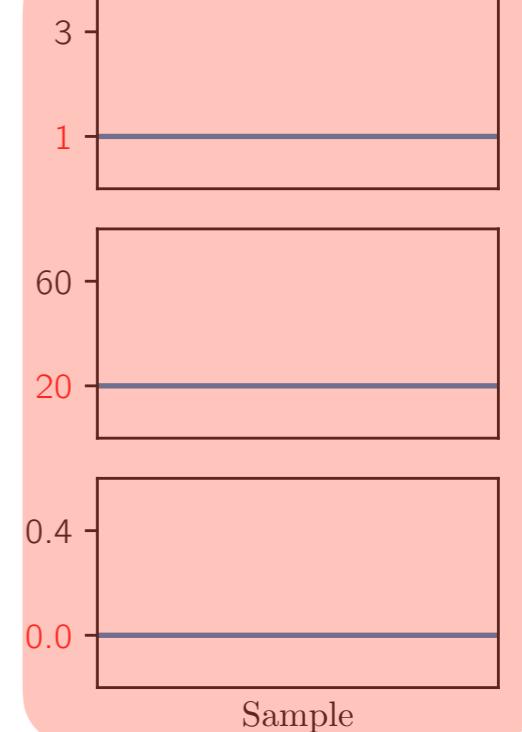
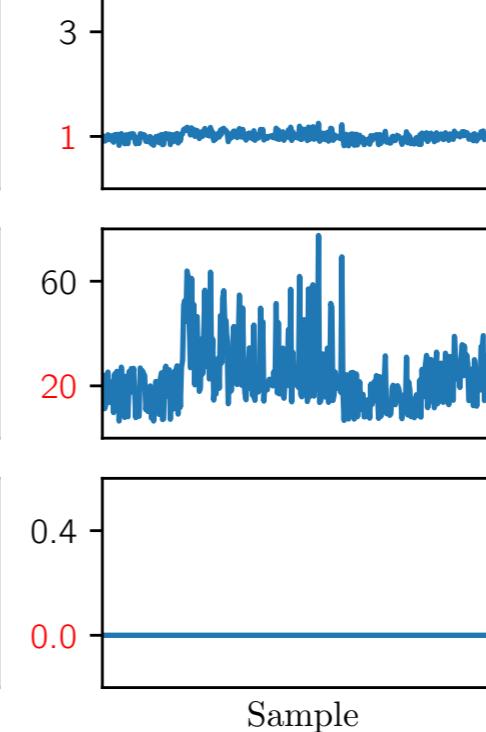
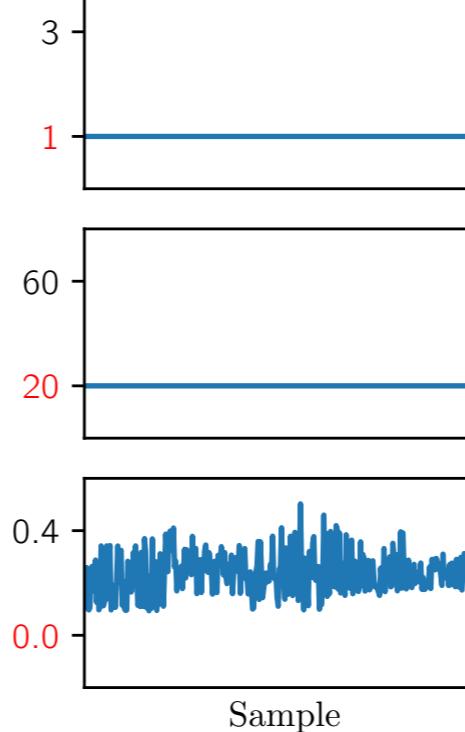
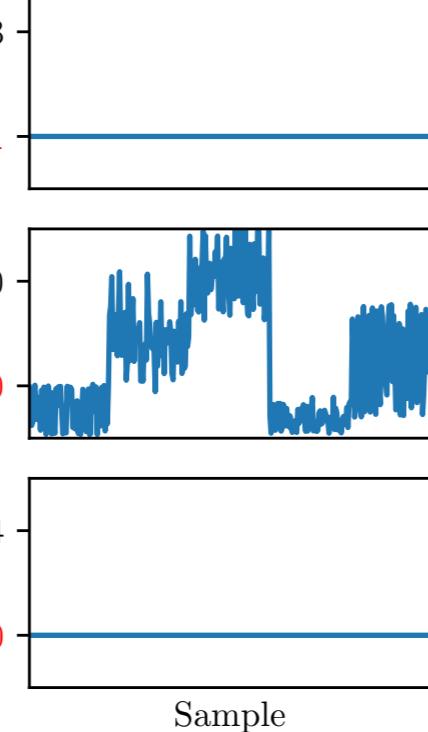
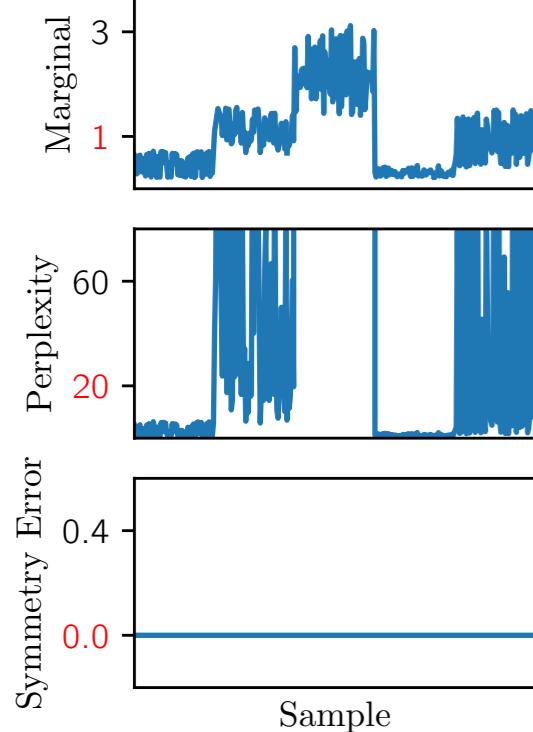
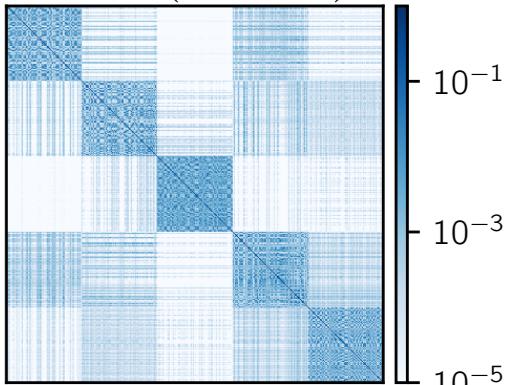
$\overline{\mathbf{P}^e}$ ($\xi=20$)



Sym-Entropic

$$\mathbf{P}^{se} = \text{Proj}_{\mathcal{H}_\xi \cap \mathcal{S}}^{\text{KL}}(\mathbf{K})$$

\mathbf{P}^{se} ($\xi=20$)



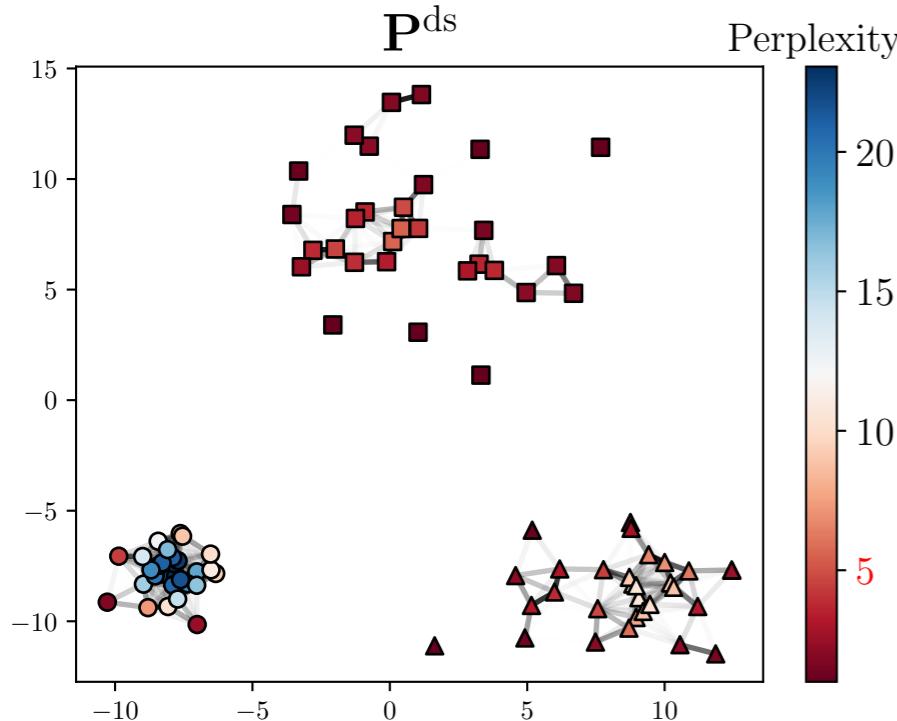
→ Controls ℓ_1 norm, entropy and symmetry at the same time.

Visualization on toy example

Doubly-Stochastic

$$\min_{\mathbf{P} \geq 0, \mathbf{P}\mathbf{1}=\mathbf{1}, \mathbf{P} \in \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle$$

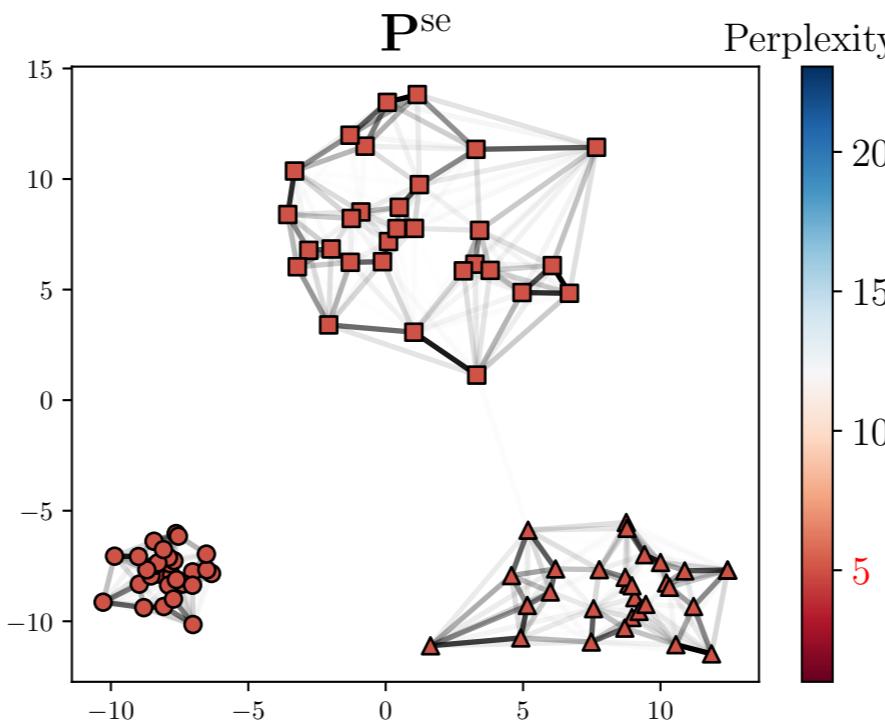
$$\sum_i H(\mathbf{P}_{i:}) \geq n(\log \xi + 1)$$



Symmetric-Entropic (OURS)

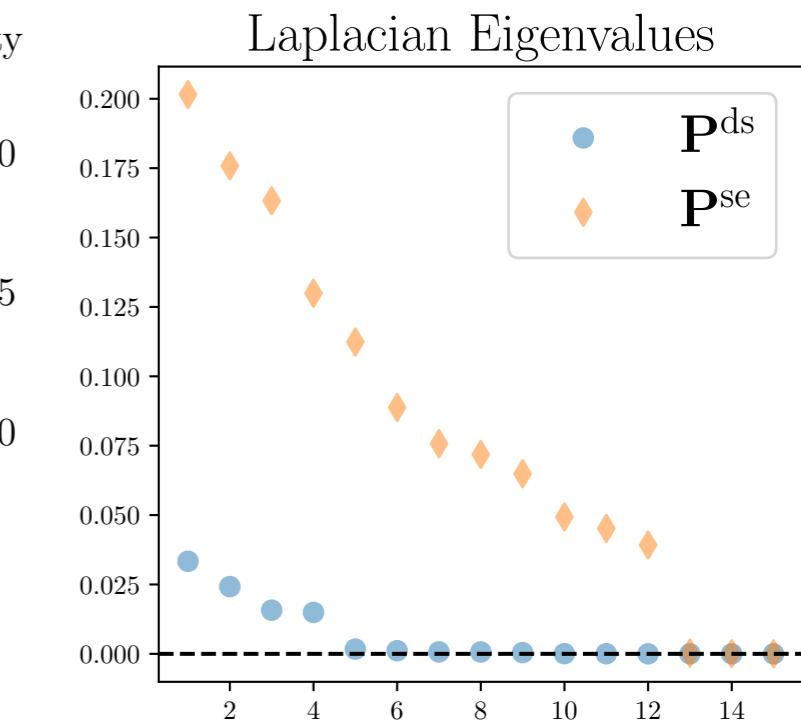
$$\min_{\mathbf{P} \geq 0, \mathbf{P}\mathbf{1}=\mathbf{1}, \mathbf{P} \in \mathcal{S}} \langle \mathbf{P}, \mathbf{C} \rangle$$

$$\forall i, H(\mathbf{P}_{i:}) \geq \log \xi + 1$$



Symmetric OT

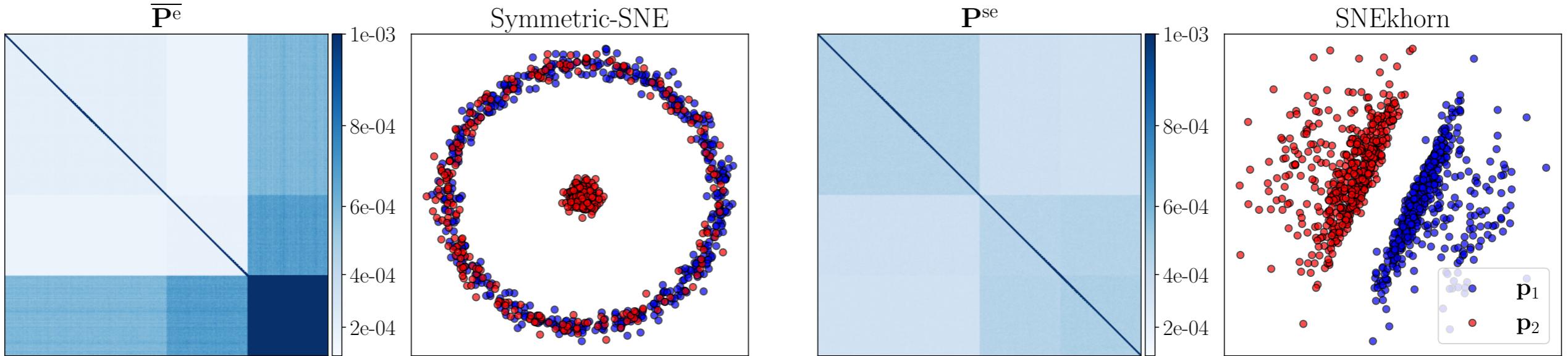
Global vs Pointwise



- Perplexity is set to $\xi = 5$ (average for \mathbf{P}^{ds}).
- \mathbf{P}^{se} can adapt to the varying noise levels.
- \mathbf{P}^{ds} retrieves many unwanted clusters.

null eigenvalues
= # clusters

| Toy example: varying noise levels



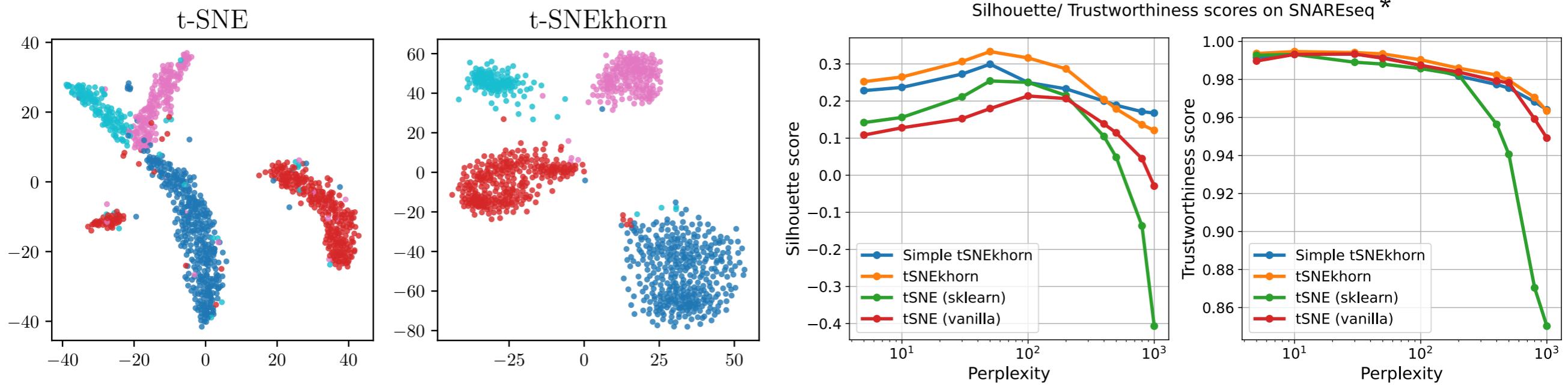
\mathbf{p}_1 and \mathbf{p}_2 taken in the 10⁴-dimensional probability simplex.

$$x_i = \tilde{x}_i / (\sum_j \tilde{x}_{ij}), \quad \tilde{x}_i \sim \begin{cases} \mathcal{M}(1000, \mathbf{p}_1), & 1 \leq i \leq 500 \\ \mathcal{M}(1000, \mathbf{p}_2), & 501 \leq i \leq 750 \\ \mathcal{M}(2000, \mathbf{p}_2), & 751 \leq i \leq 1000 . \end{cases}$$

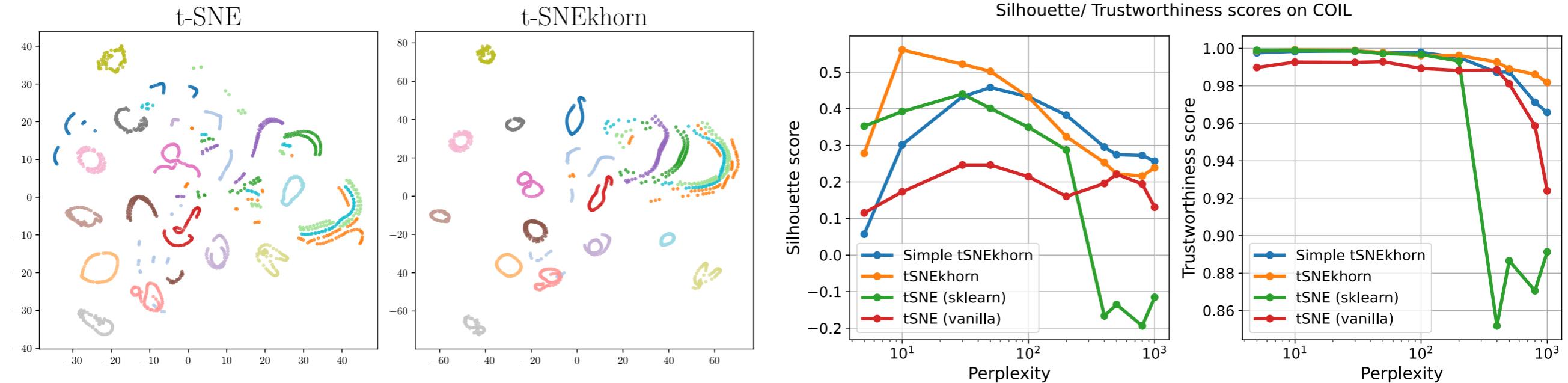
| SNE is misled by the batch effect unlike SNEkhorn.

Dimension Reduction Results

SNAREseq Single Cell data



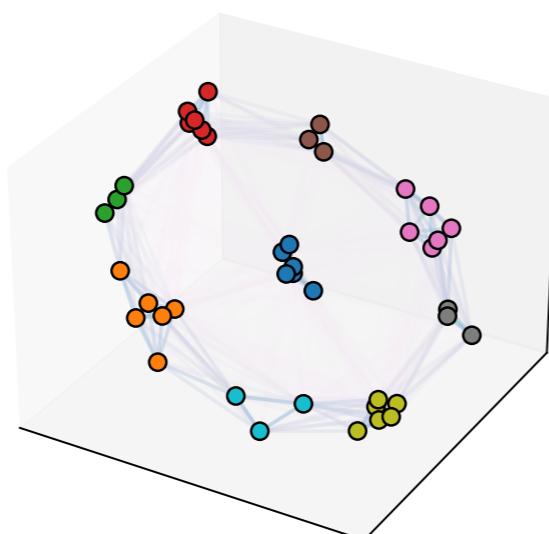
COIL-20 Image data



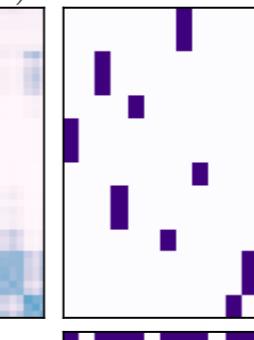
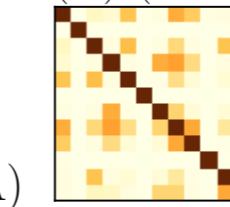
* Simple tSNEkhorn → tSNEkhorn with same affinity as t-SNE for the embeddings \mathbf{Z} (not doubly stochastic).

Distributional reduction

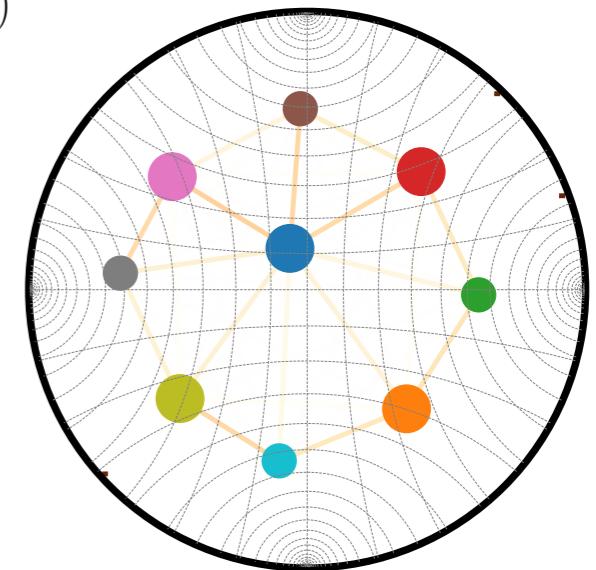
Input \mathbf{X}



$C_Z(\mathbf{Z})$ (Lorentz)



Embedding \mathbf{Z}



Coupling \mathbf{T}

Motivation

◆ Single-cell RNA-seq

Technical noise due to partial sampling of RNA molecules within cells.

METHOD

Open Access

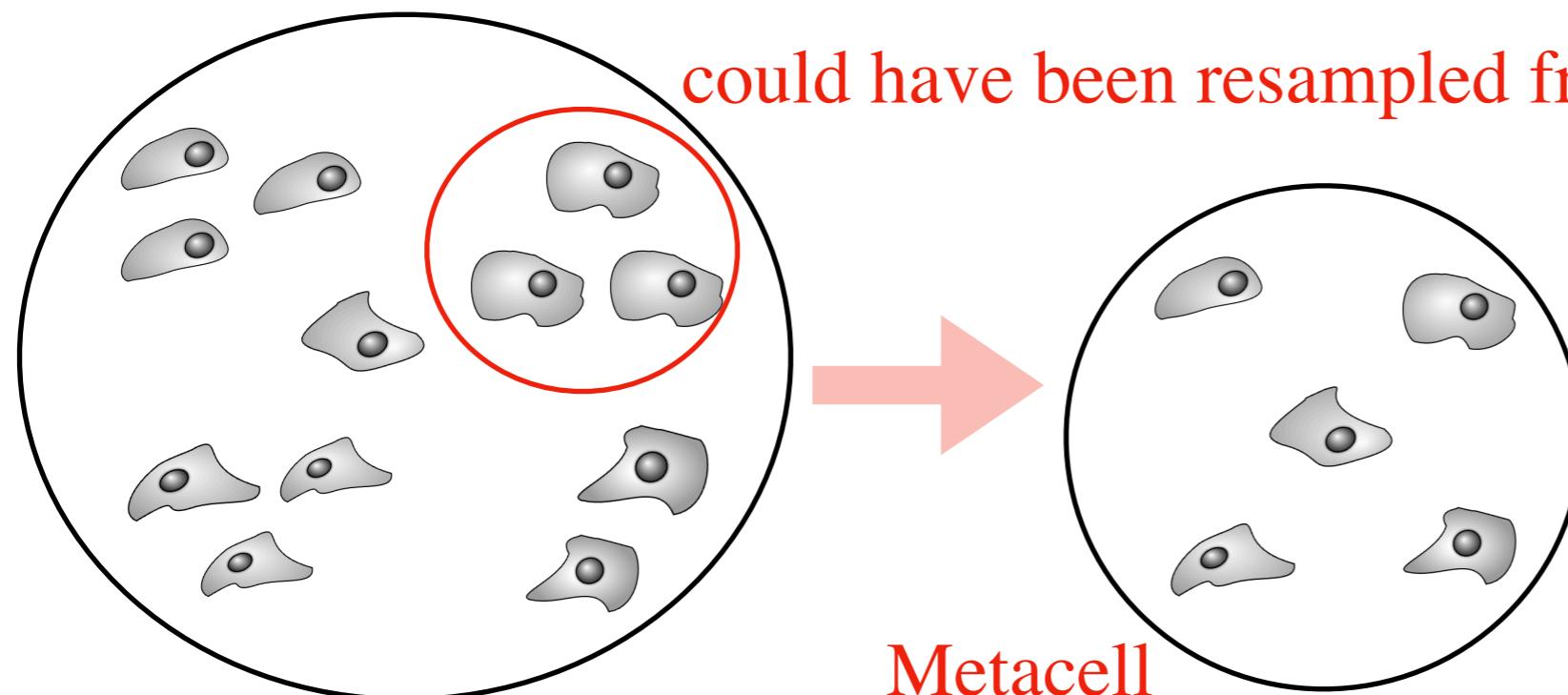


MetaCell: analysis of single-cell RNA-seq data using *K*-nn graph partitions

Yael Baran¹, Akhiad Bercovich¹, Arnau Sebe-Pedros¹, Yaniv Lubling¹, Amir Giladi², Elad Chomsky¹, Zohar Meir¹, Michael Hoichman¹, Aviezer Lifshitz¹ and Amos Tanay^{1*}

Problem : impossible to resample a cell

- integration of data from different cells
- need to **separate the sampling effect from biological variance**

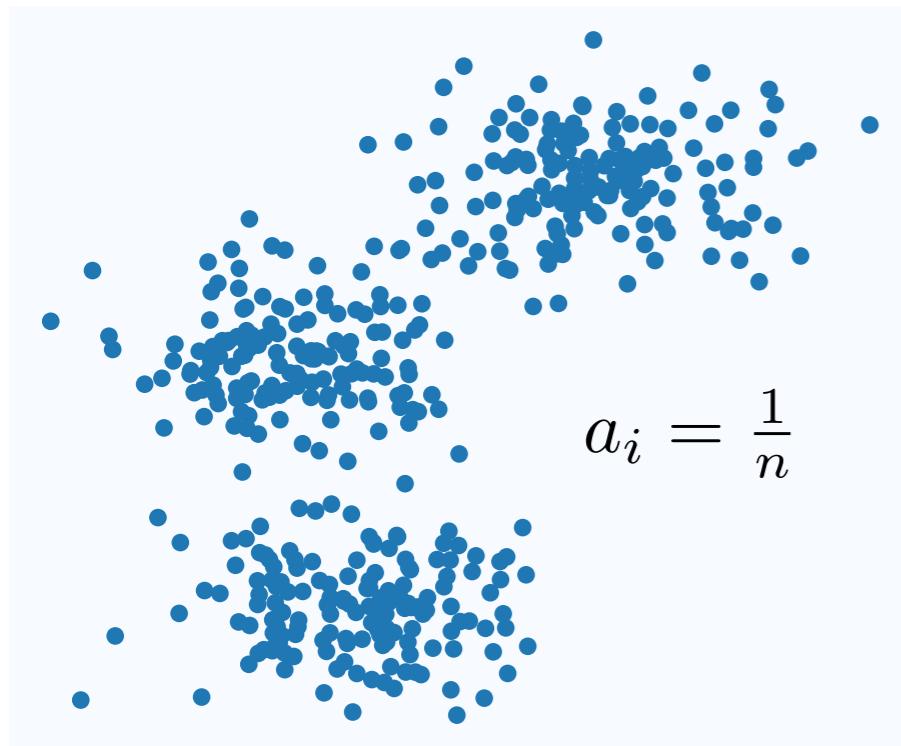


- ◆ We would like to choose the granularity of the output data

Framing DR in terms of distributions

Data: $(\mathbf{x}_i)_{i \in [\![n]\!]} ; \mathbf{x}_i \in \mathbb{R}^d \longrightarrow$ A probability distribution describing the data

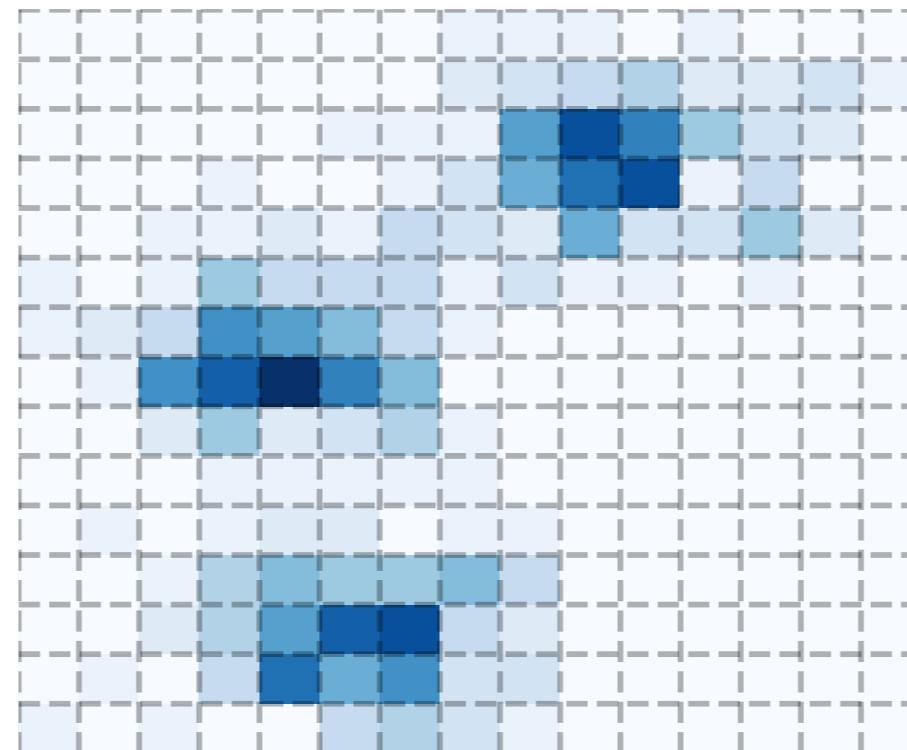
Lagrangian: $\sum_{i=1}^n a_i \delta_{x_i}$



(point clouds)

$$\delta_{\mathbf{x}_i}(\mathbf{x}) = 1 \text{ if } \mathbf{x} = \mathbf{x}_i \text{ else } 0$$

Eulerian: $\sum_{i=1}^N a_i \delta_{\hat{x}_i}$



(histograms)

\hat{x}_i fixed position (grid)

Probability simplex

$$\mathbf{a} = (a_i)_{i \in [\![n]\!]} \in \Sigma_n$$

$$a_i \geq 0, \sum_{i=1}^n a_i = 1$$



From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)

Kantorovitch Formulation

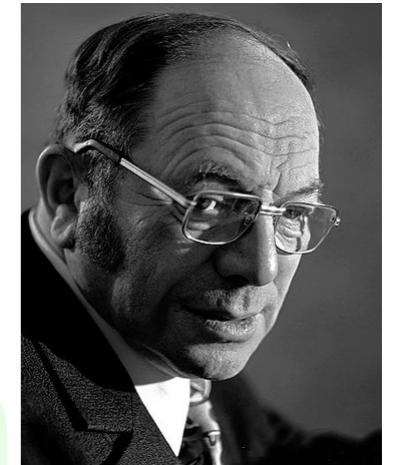
Two probability distributions

$$\mu \in \mathcal{P}(\mathcal{X})$$

$$\nu \in \mathcal{P}(\mathcal{Y})$$

A cost function

$$c(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$$



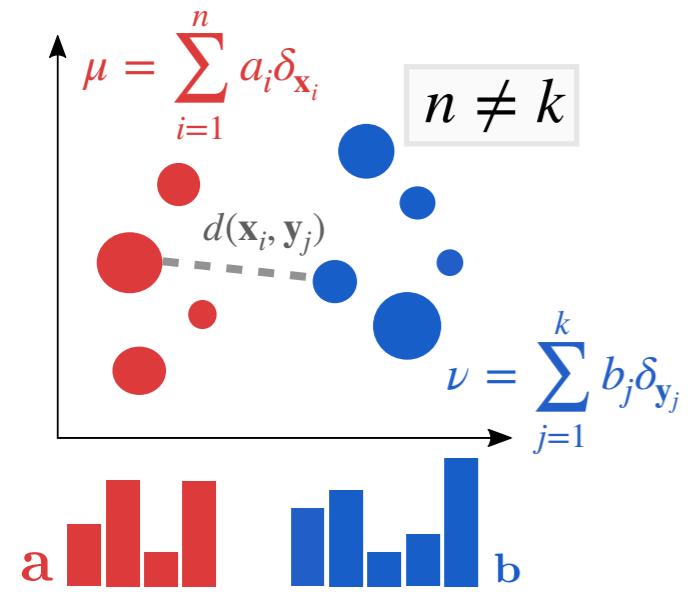
$\mu \in \mathcal{P}(\mathcal{X})$ is transported to $\nu \in \mathcal{P}(\mathcal{Y})$ by a **transport plan** $T \in \mathcal{P}(\mathcal{X} \times \mathcal{Y})$

We want to find the plan that **minimizes the overall cost** of moving all the points.

$$\inf_{T \in \Pi(\mu, \nu)} \int c(x, y) dT(x, y)$$

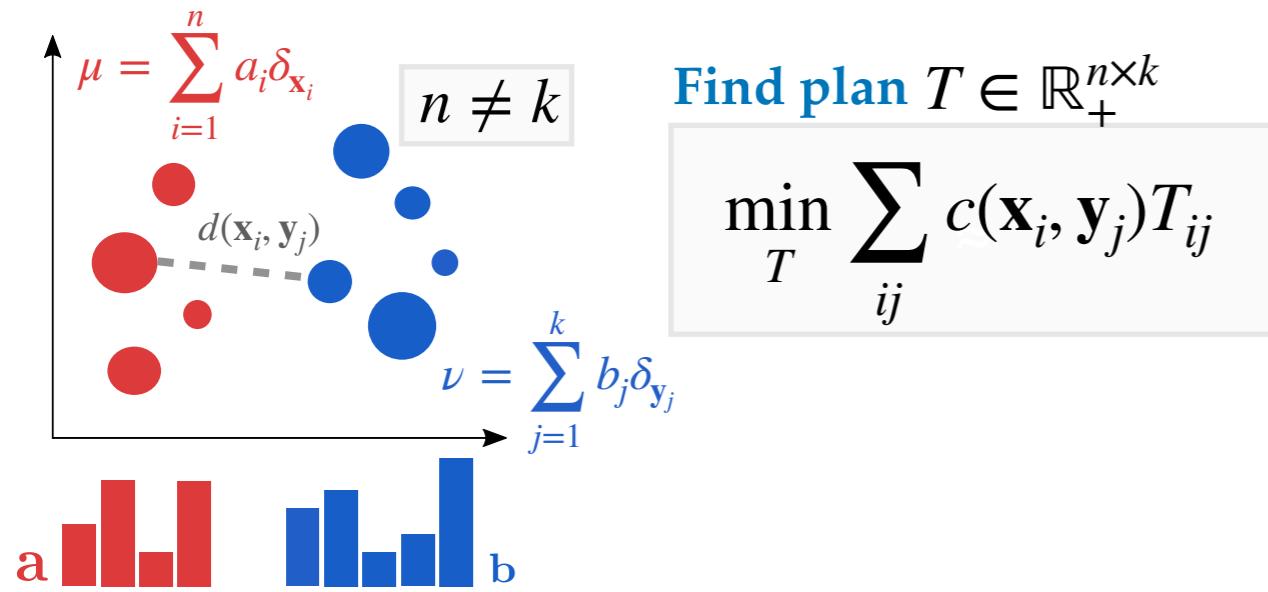
From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



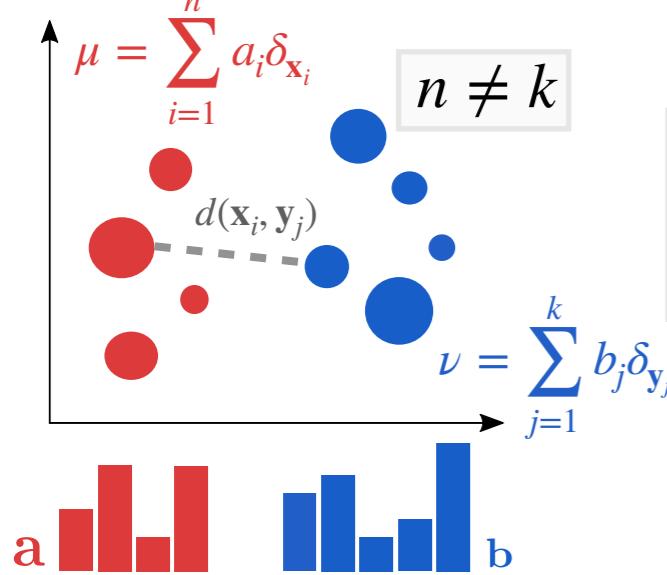
From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

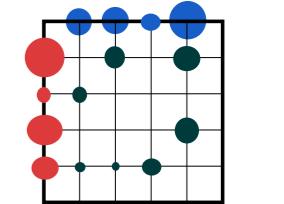
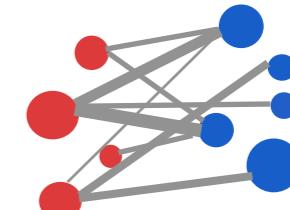
$$\min_T \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

which constraints ?

Coupling
 $\Pi(a, b)$

$$T^\top \mathbf{1}_n = a$$

$$T \mathbf{1}_k = b$$



Bakeries = quantity of breads

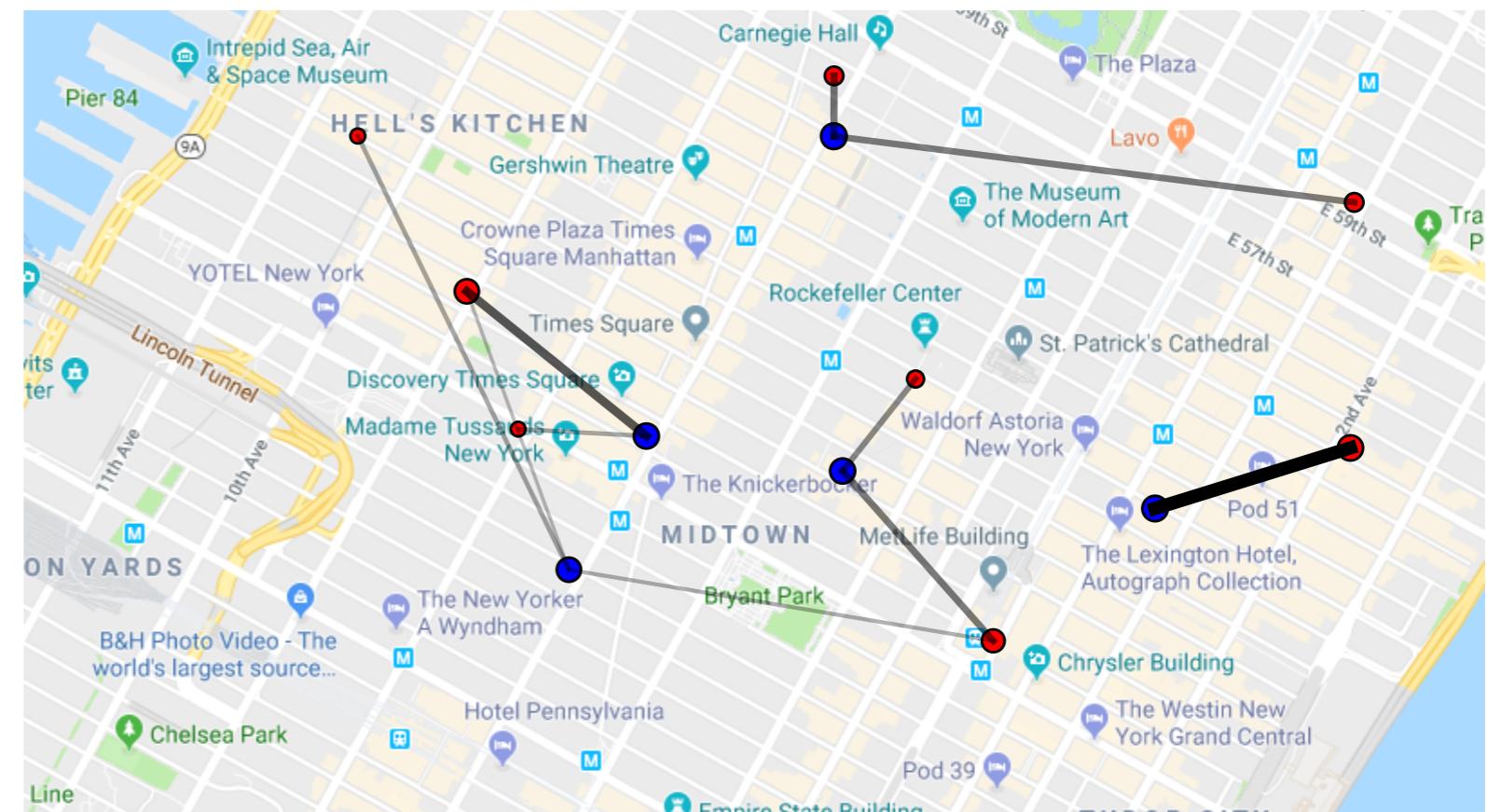
loc: \mathbf{x}_i quantity: a_i

Cafés = demand of breads

loc: \mathbf{y}_j demand: b_j

Distance between bakeries and cafés

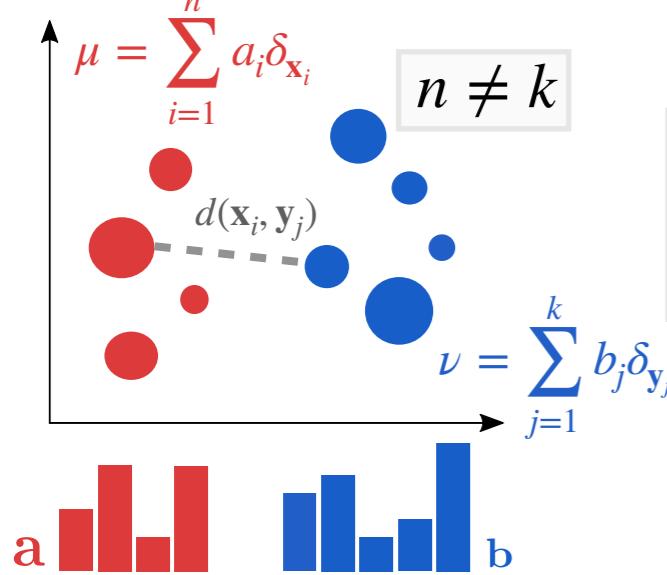
$$c(\mathbf{x}_i, \mathbf{y}_j)$$



We want to route all the breads from bakeries to cafés the cheapest way

From Wasserstein to Gromov-Wasserstein

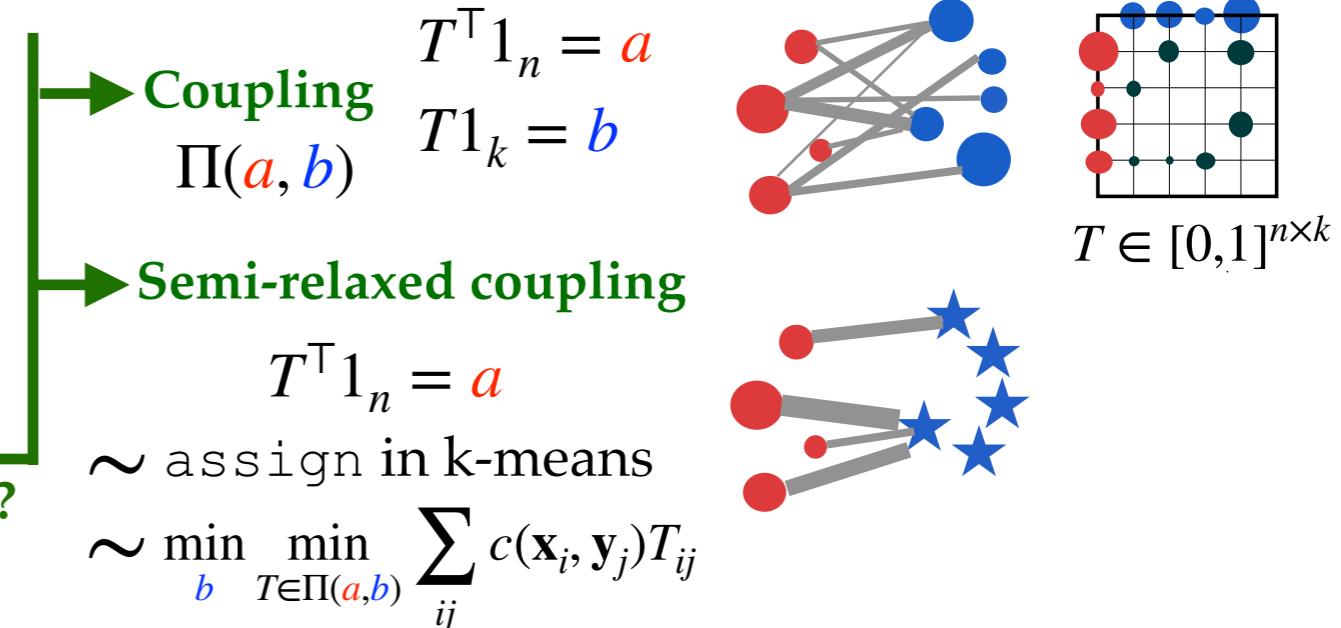
♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

$$\min_T \sum_{ij} c(x_i, y_j) T_{ij}$$

which constraints ?



Bakeries = quantity of breads

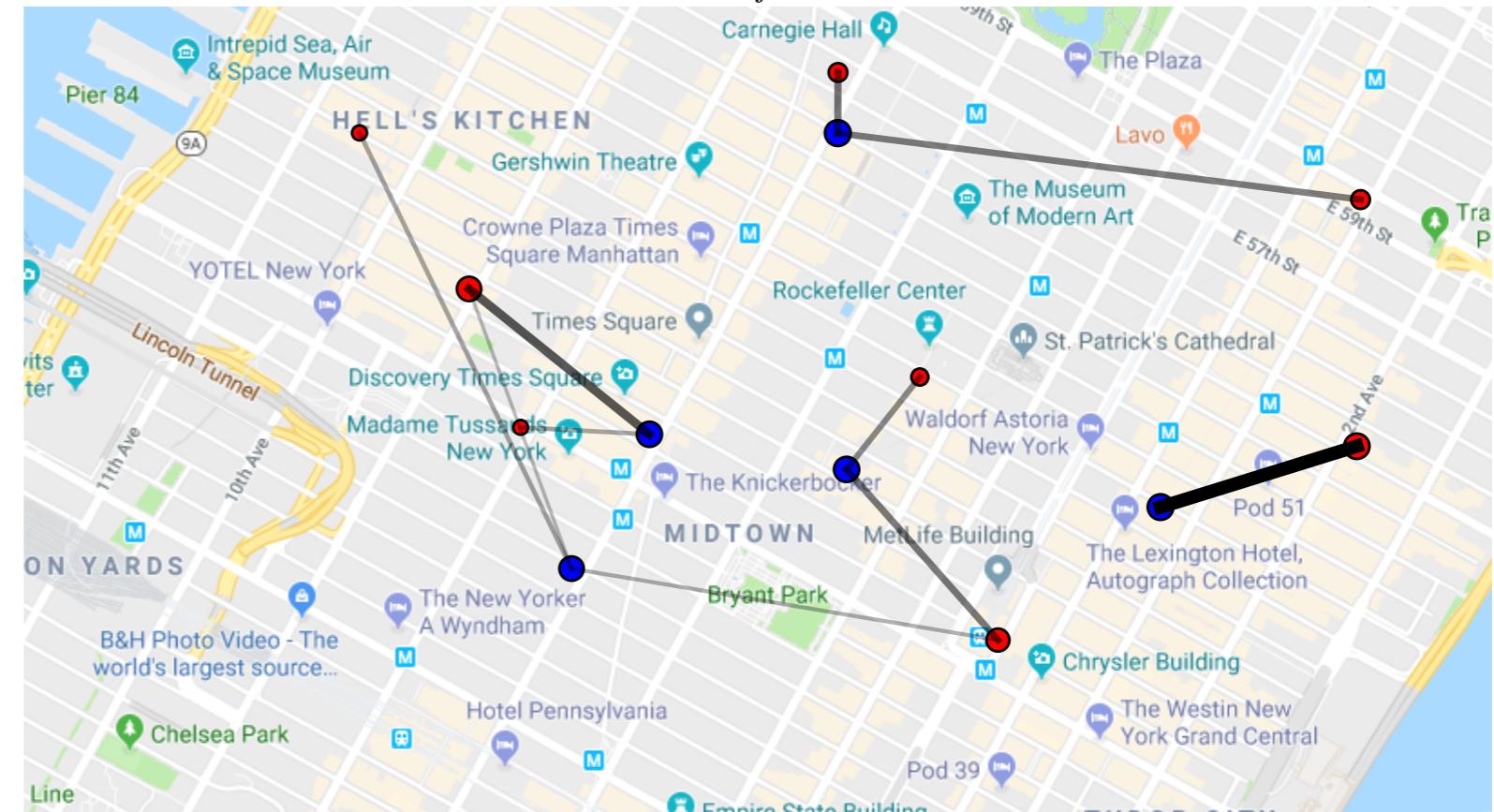
loc: \mathbf{x}_i quantity: a_i

Cafés = demand of breads

loc: \mathbf{y}_j demand: b_j

Distance between bakeries and cafés

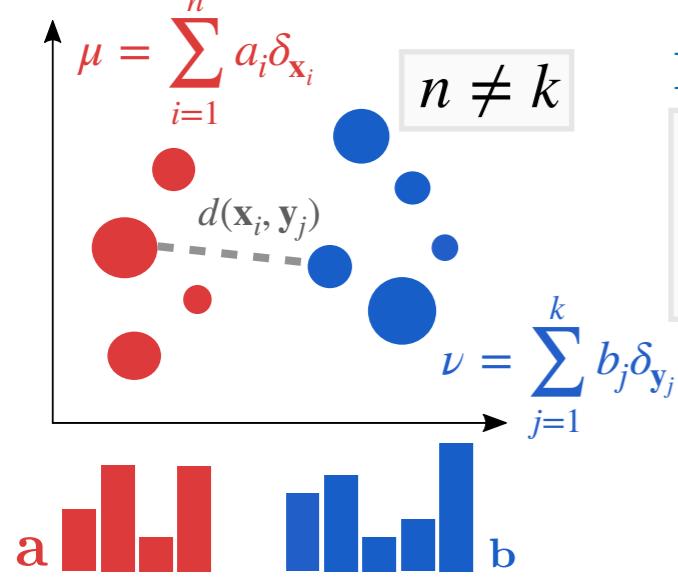
$$c(\mathbf{x}_i, \mathbf{y}_j)$$



We want to route all the breads from bakeries to cafés the cheapest way

From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

$$\min_T \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

which constraints ?

→ **Coupling** $\Pi(a, b)$

$$T^\top \mathbf{1}_n = a$$

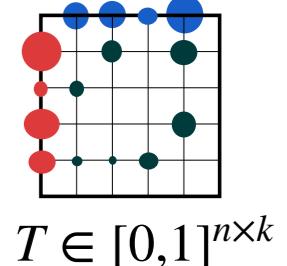
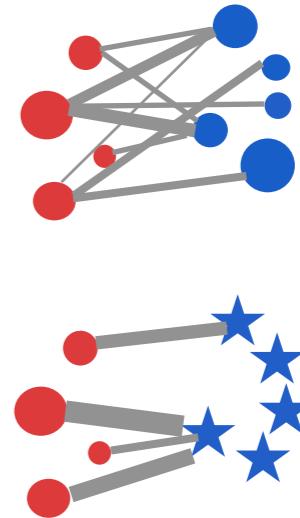
$$T \mathbf{1}_k = b$$

→ **Semi-relaxed coupling**

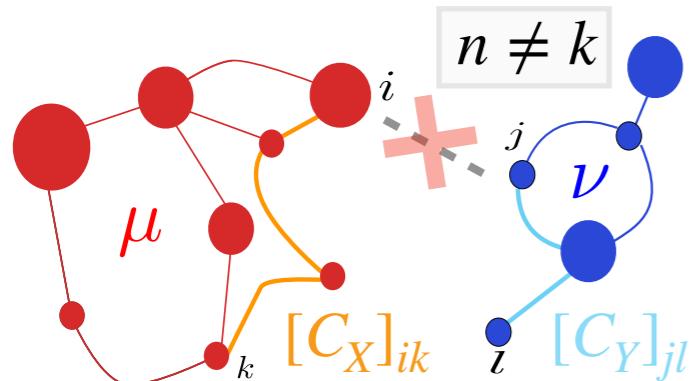
$$T^\top \mathbf{1}_n = a$$

~ assign in k-means

~ $\min_b \min_{T \in \Pi(a, b)} \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$

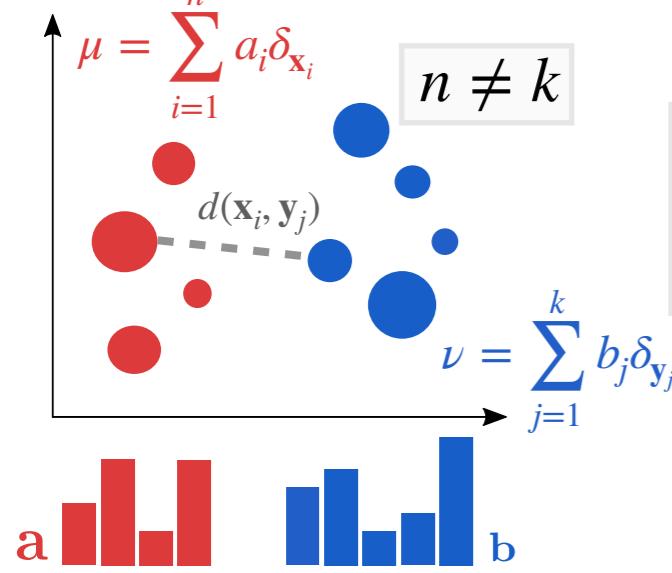


♦ Gromov-Wasserstein



From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

$$\min_T \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

which constraints ?

→ **Coupling**
 $\Pi(a, b)$

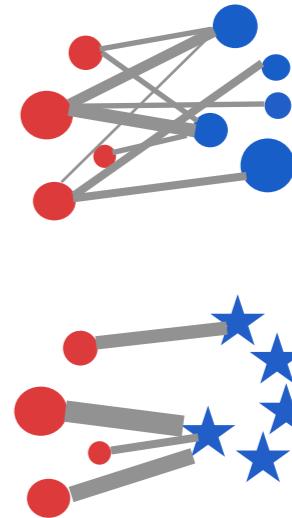
$$T^\top \mathbf{1}_n = a$$

$$T \mathbf{1}_k = b$$

→ **Semi-relaxed coupling**

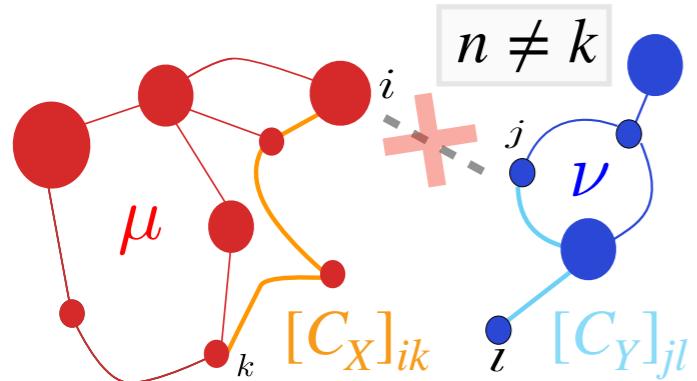
$$T^\top \mathbf{1}_n = a$$

~ assign in k-means

$$\sim \min_b \min_{T \in \Pi(a, b)} \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$


$$T \in [0,1]^{n \times k}$$

♦ Gromov-Wasserstein

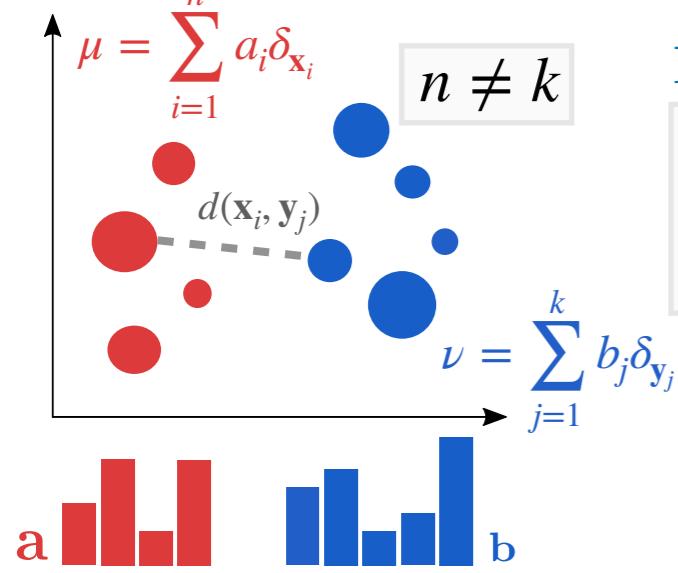


Quadratic OT: find the plan

$$\min_{T \in \Pi(a, b)} \sum_{ijkl} L([C_X]_{ik}, [C_Y]_{jl}) T_{ij} T_{kl}$$

From Wasserstein to Gromov-Wasserstein

♦ Classical optimal transport (in a nutshell)



Find plan $T \in \mathbb{R}_+^{n \times k}$

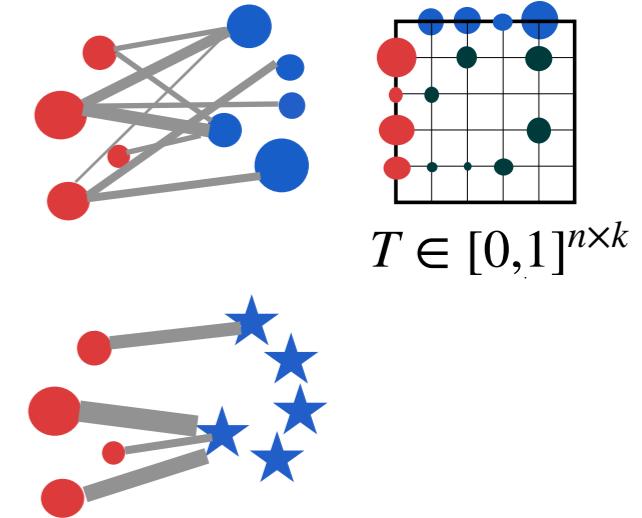
$$\min_T \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$$

which constraints ?

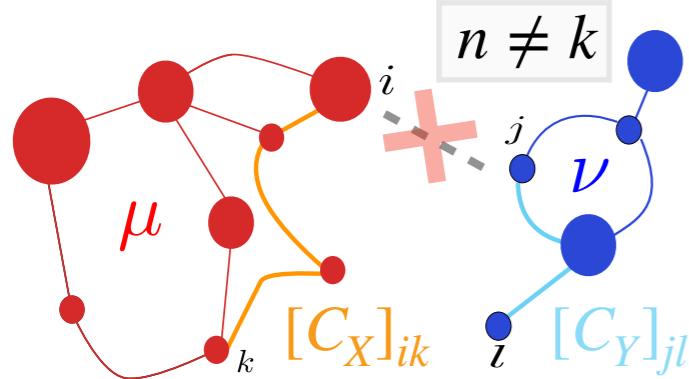
→ **Coupling** $\Pi(a, b)$ $T^\top \mathbf{1}_n = a$
 $T \mathbf{1}_k = b$

→ **Semi-relaxed coupling**

$T^\top \mathbf{1}_n = a$
~ assign in k-means
~ $\min_b \min_{T \in \Pi(a, b)} \sum_{ij} c(\mathbf{x}_i, \mathbf{y}_j) T_{ij}$



♦ Gromov-Wasserstein



♦ L measures distortion

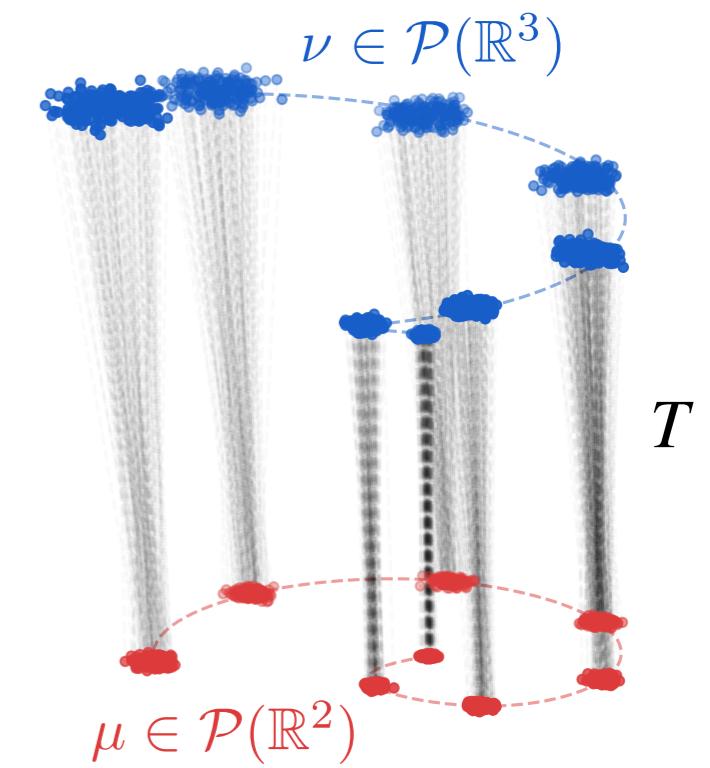
$$\left| [C_X]_{ik} - [C_Y]_{jl} \right|^2$$

♦ Goal : preserving pairwise connectivity

- ♦ Distance w.r.t. isomorphisms
- ♦ Difficult quadratic problem (NP-hard)

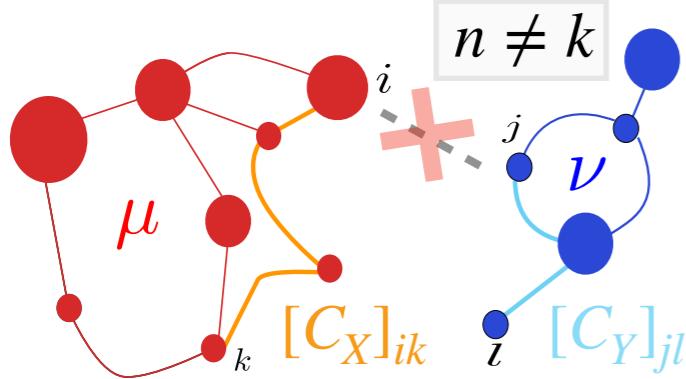
Quadratic OT: find the plan

$$\min_{T \in \Pi(a, b)} \sum_{ijkl} L([C_X]_{ik}, [C_Y]_{jl}) T_{ij} T_{kl}$$



From Wasserstein to Gromov-Wasserstein

♦ Gromov-Wasserstein



(Sturm, 2012) (Mémoli, 2011)

♦ L measures distortion

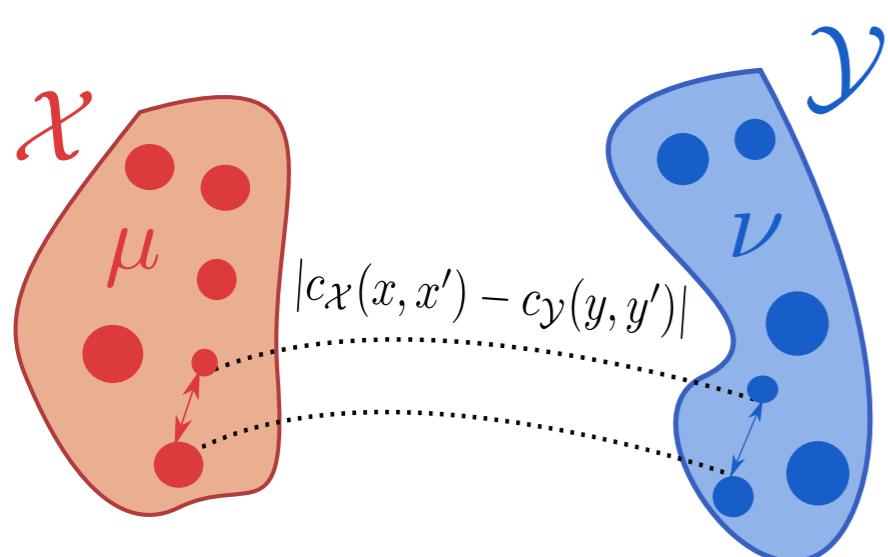
$$\left| [C_X]_{ik} - [C_Y]_{jl} \right|^2$$

♦ Goal : preserving pairwise connectivity

♦ Non-convex quadratic problem (NP-hard)

Quadratic OT: find the plan

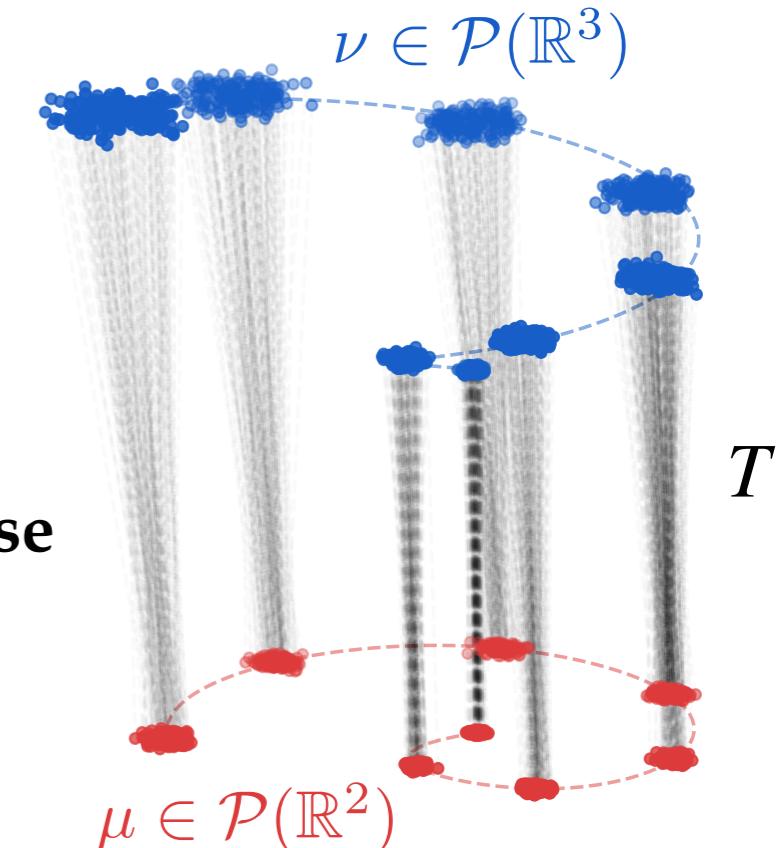
$$\min_{T \in \Pi(\mathbf{a}, \mathbf{b})} \sum_{ijkl} L \left([C_X]_{ik}, [C_Y]_{jl} \right) T_{ij} T_{kl}$$



♦ Distance w.r.t. isomorphisms, on the space of metric measure spaces

$$\mathbb{X} = (\mathcal{X}, c_{\mathcal{X}}, \mu \in \mathcal{P}(\mathcal{X}))$$

$$\mathbb{Y} = (\mathcal{Y}, c_{\mathcal{Y}}, \nu \in \mathcal{P}(\mathcal{Y}))$$



$$GW(\mathbb{X}, \mathbb{Y}) = 0 \text{ iff}$$

$$\exists \phi : \mathcal{X} \rightarrow \mathcal{Y}$$

Isometry

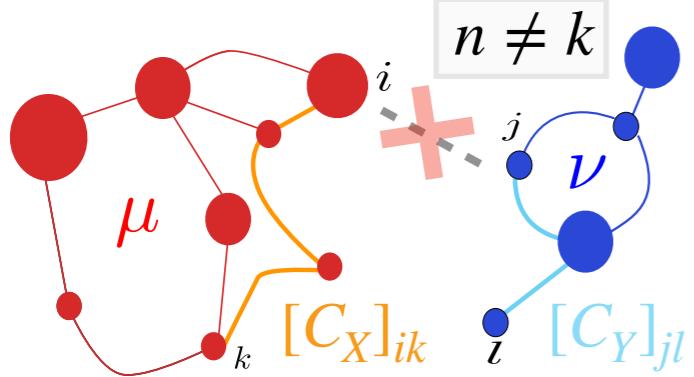
$$c_{\mathcal{X}}(x, x') = c_{\mathcal{Y}}(\phi(x), \phi(x'))$$

Measure preserving

$$\phi \# \mu = \nu$$

From Wasserstein to Gromov-Wasserstein

◆ Gromov-Wasserstein



(Sturm, 2012) (Mémoli, 2011)

◆ L measures distortion

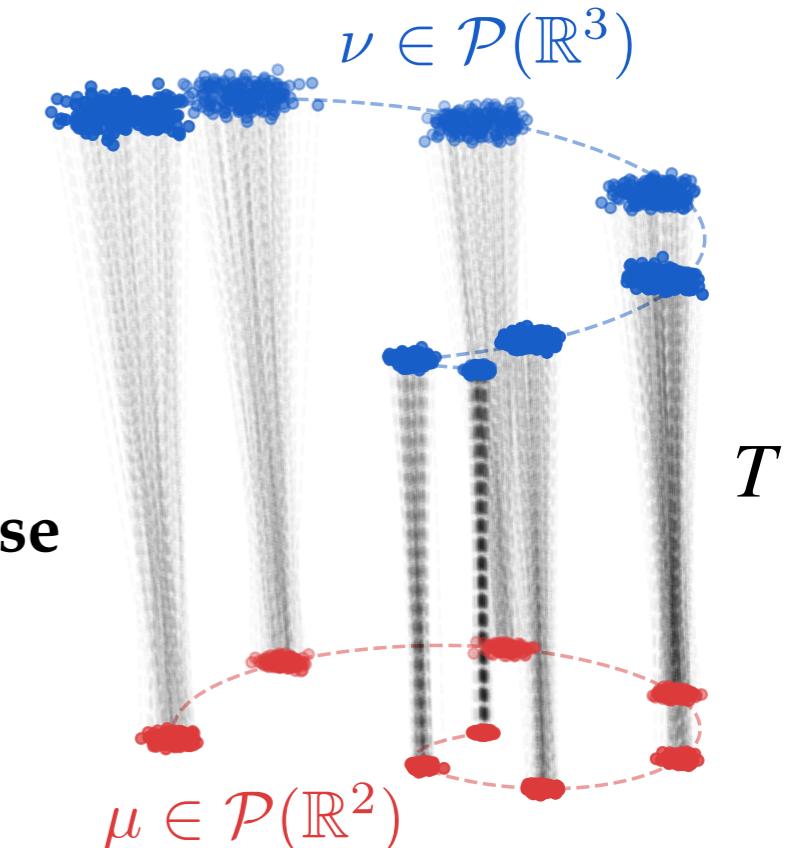
$$\left| [C_X]_{ik} - [C_Y]_{jl} \right|^2$$

◆ Goal : preserving pairwise connectivity

◆ Non-convex quadratic problem (NP-hard)

Quadratic OT: find the plan

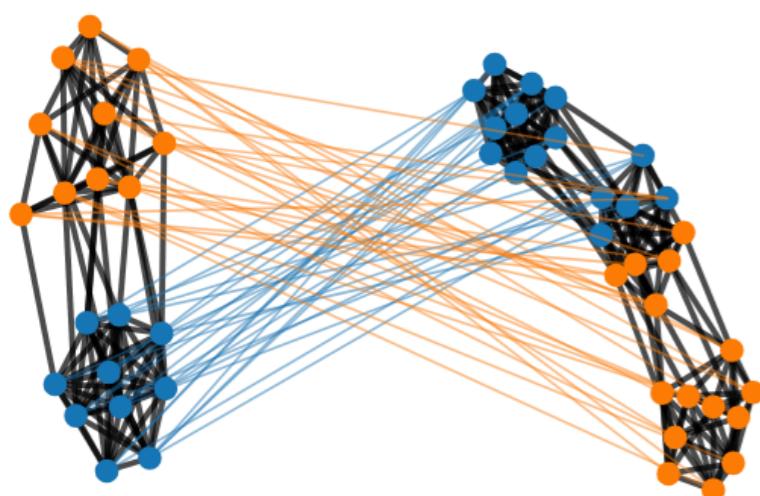
$$\min_{T \in \Pi(\mathbf{a}, \mathbf{b})} \sum_{ijkl} L \left([C_X]_{ik}, [C_Y]_{jl} \right) T_{ij} T_{kl}$$



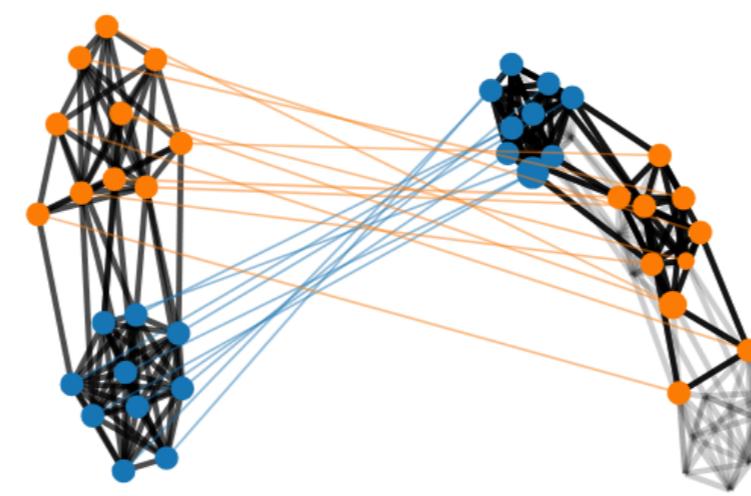
◆ Semi relaxed Gromov-Wasserstein

(Vincent-Cuaz, 2022)

$$GW(\mathbf{C}, \mathbf{h}, \bar{\mathbf{C}}, \bar{\mathbf{h}}) = 0.219$$



$$srGW(\mathbf{C}, \mathbf{h}, \bar{\mathbf{C}}) = 0.05$$

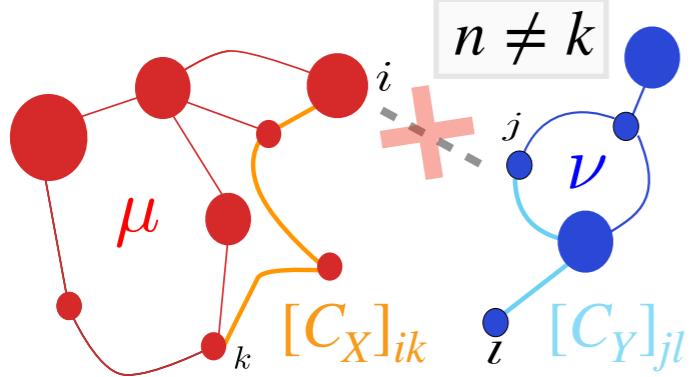


$$\min_T \sum_{ijkl} L \left([C_X]_{ik}, [C_Y]_{jl} \right) T_{ij} T_{kl}$$

| $T^\top \mathbf{1}_n = \mathbf{a}$

From Wasserstein to Gromov-Wasserstein

◆ Gromov-Wasserstein



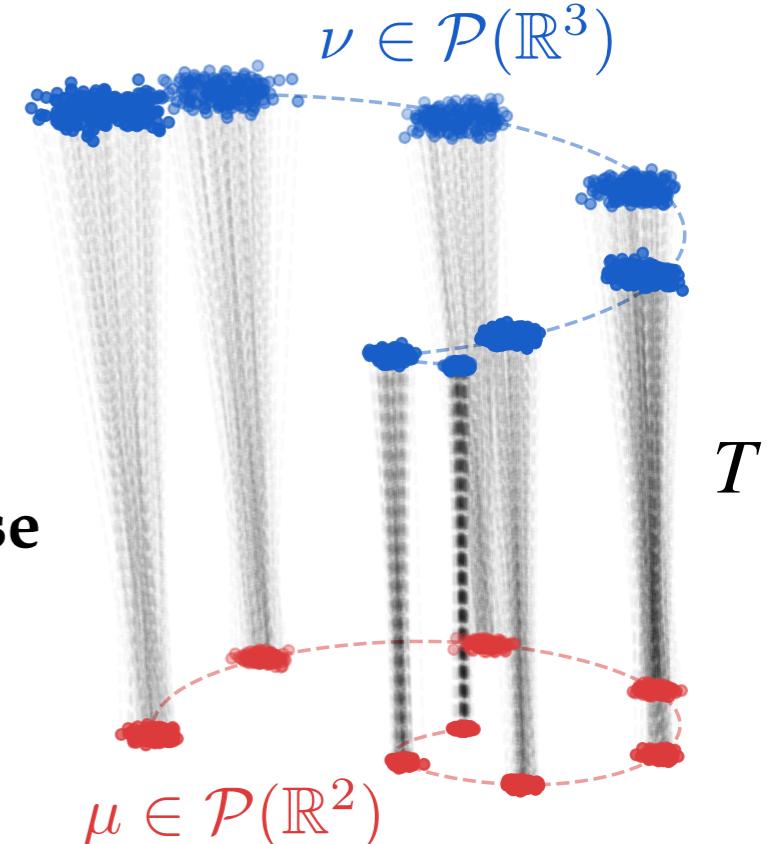
(Sturm, 2012) (Mémoli, 2011)

◆ L measures distortion

$$\left| [C_X]_{ik} - [C_Y]_{jl} \right|^2$$

◆ Goal : preserving pairwise connectivity

◆ Non-convex quadratic problem (NP-hard)



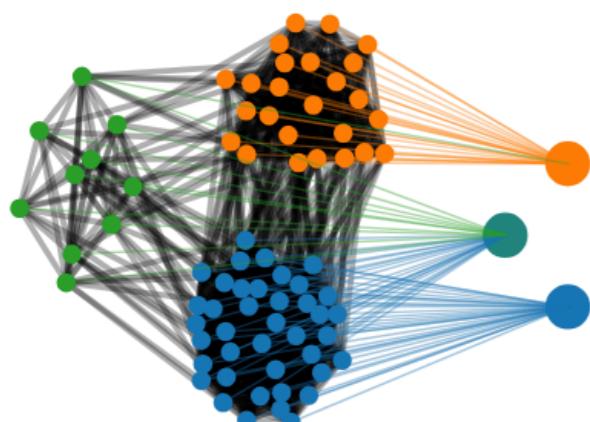
Quadratic OT: find the plan

$$\min_{T \in \Pi(\mathbf{a}, \mathbf{b})} \sum_{ijkl} L \left([C_X]_{ik}, [C_Y]_{jl} \right) T_{ij} T_{kl}$$

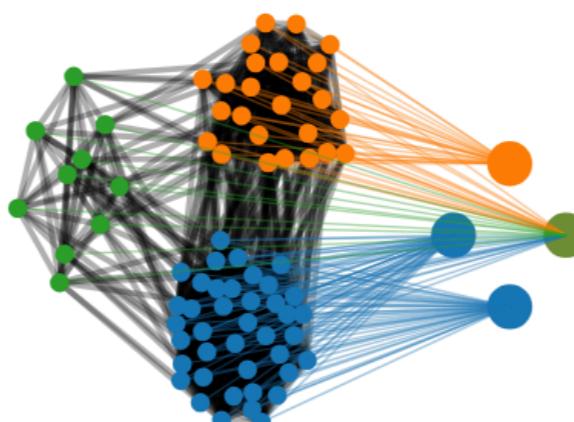
◆ Semi relaxed Gromov-Wasserstein

◆ Clustering of nodes

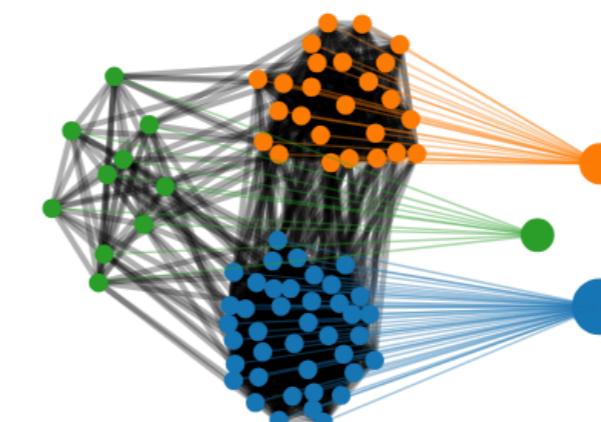
$$GW(\mathbf{C}, \mathbf{h}, \mathbf{I}_3, \bar{\mathbf{h}}) = 0.235 \quad (\text{ami}=0.66)$$



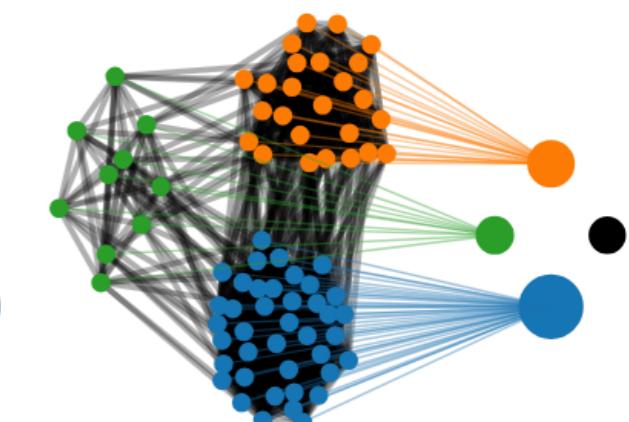
$$GW(\mathbf{C}, \mathbf{h}, \mathbf{I}_4, \bar{\mathbf{h}}) = 0.274 \quad (\text{ami}=0.54)$$



$$srGW(\mathbf{C}, \mathbf{h}, \mathbf{I}_3) = 0.087 \quad (\text{ami}=1.0)$$



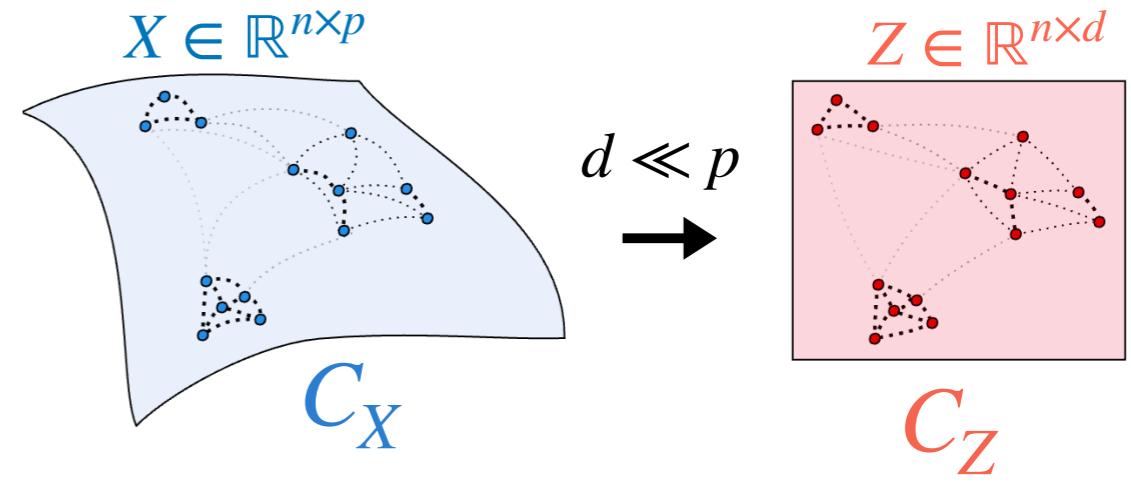
$$srGW(\mathbf{C}, \mathbf{h}, \mathbf{I}_4) = 0.087 \quad (\text{ami}=1.0)$$



DR as OT in disguise

♦ Dimension reduction

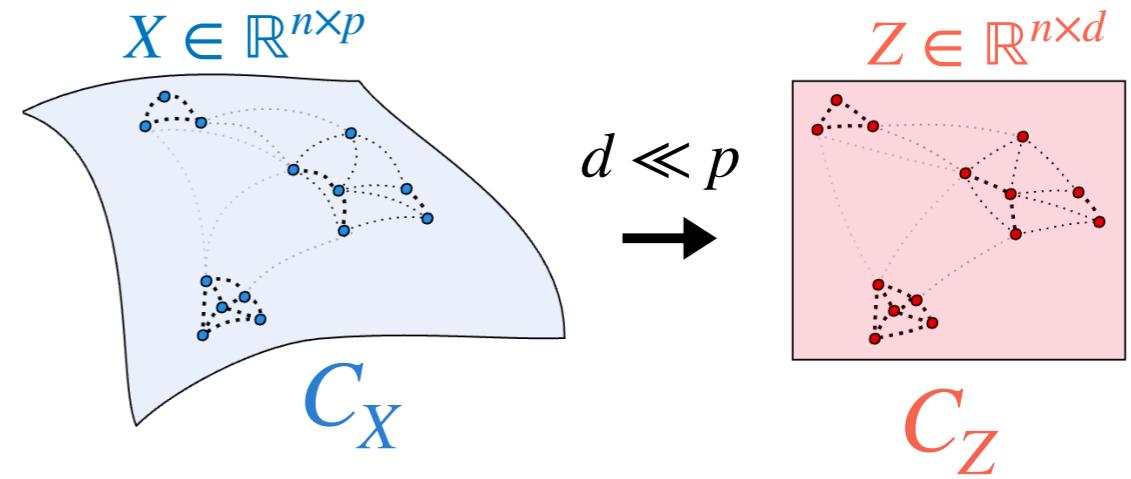
$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$



DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$



↑
equiv
↓

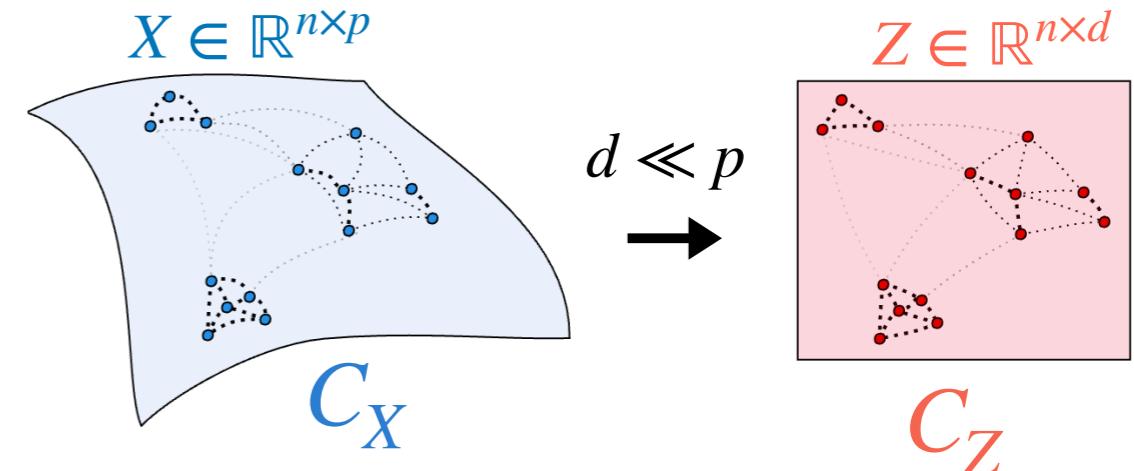
Permutation equivariance
 $\forall P, C_{PZ} = PC_ZP^\top$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$



↑
 equiv
↓

Permutation equivariance
 $\forall P, C_{PZ} = PC_ZP^\top$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

↑
 equiv
↓

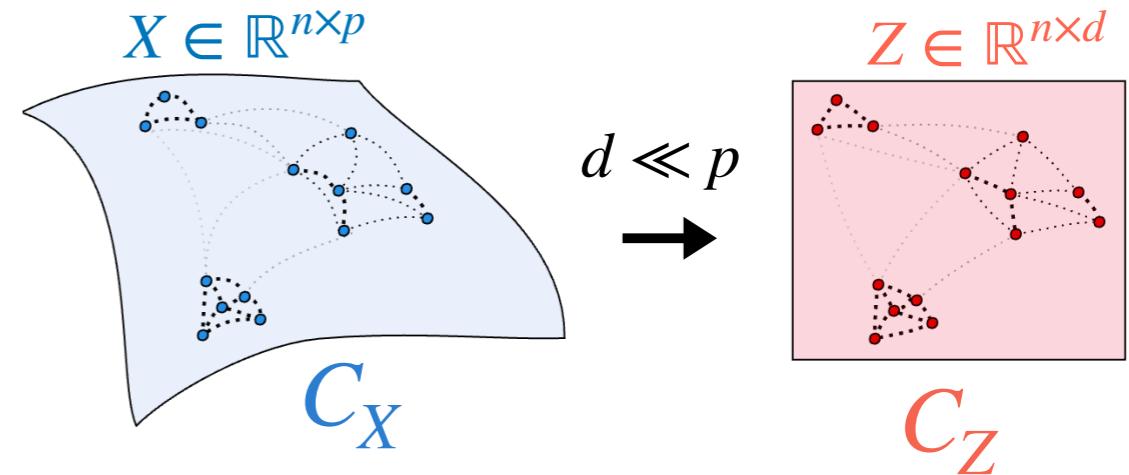
$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij}P_{kl}$$

Gromov-Monge

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$



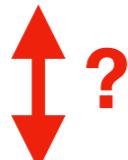
Permutation equivariance

$$\forall P, C_{PZ} = PC_ZP^\top$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij}P_{kl}$$

Gromov-Monge



♦ Gromov-Wasserstein projection

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T \in \Pi(\frac{1_n}{n}, \frac{1_n}{n})} \sum_{ijkl} L\left([C_X]_{ik}, [C_Z]_{jl}\right) T_{ij}T_{kl}$$

Gromov-Wasserstein

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$

equiv
↔

Permutation equivariance

$$\forall P, C_{PZ} = PC_ZP^\top$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

equiv
↔

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij}P_{kl}$$

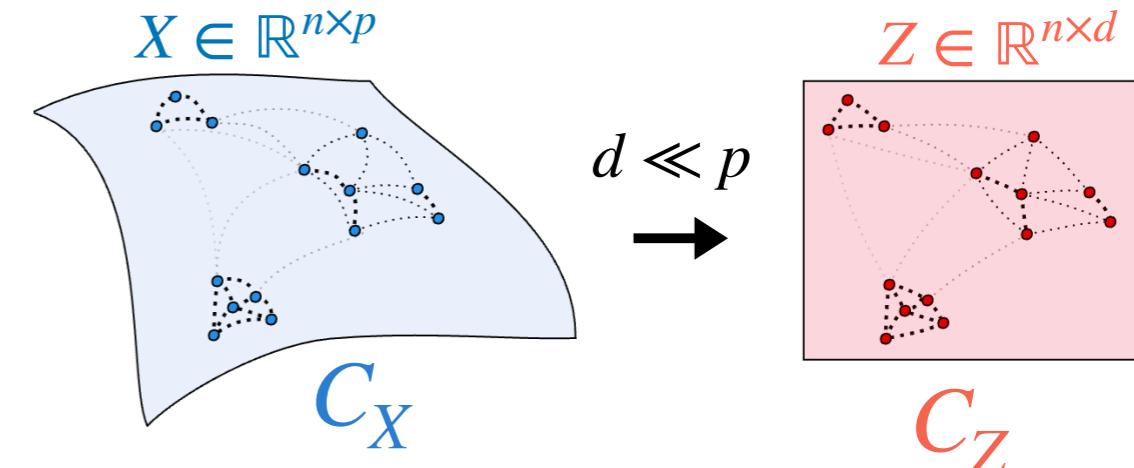
Gromov-Monge

↑?
↓

♦ Gromov-Wasserstein projection

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T \in \Pi(\frac{\mathbf{1}_n}{n}, \frac{\mathbf{1}_n}{n})} \sum_{ijkl} L\left([C_X]_{ik}, [C_Z]_{jl}\right) T_{ij}T_{kl}$$

Gromov-Wasserstein



♦ Equivalence holds for

Spectral methods

♦ C_X any matrix, $L = |\cdot|^2$, $C_Z = ZZ^\top$

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$

equiv
↔

Permutation equivariance

$$\forall P, C_{PZ} = PC_Z P^\top$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

equiv
↔

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij} P_{kl}$$

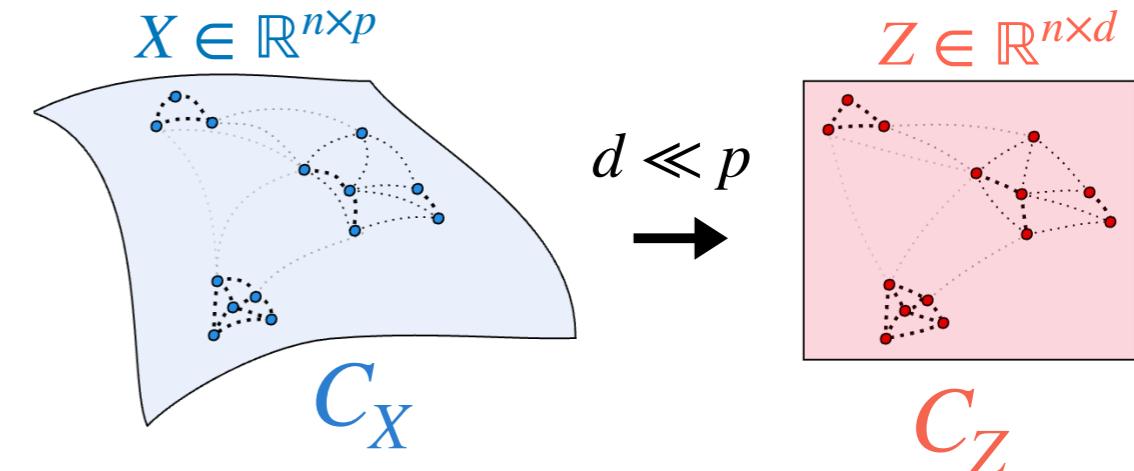
Gromov-Monge

↑?
?

♦ Gromov-Wasserstein projection

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T \in \Pi(\frac{\mathbf{1}_n}{n}, \frac{\mathbf{1}_n}{n})} \sum_{ijkl} L\left([C_X]_{ik}, [C_Z]_{jl}\right) T_{ij} T_{kl}$$

Gromov-Wasserstein



♦ Equivalence holds for

Spectral methods

♦ C_X any matrix, $L = |\cdot|^2$, $C_Z = ZZ^\top$

A is CPD : $\forall x$ **s.t.** $x^\top 1 = 0$, $x^\top Ax \geq 0$

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$

equiv
↔

Permutation equivariance

$$\forall P, C_{PZ} = PC_ZP^\top$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

equiv
↔

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij}P_{kl}$$

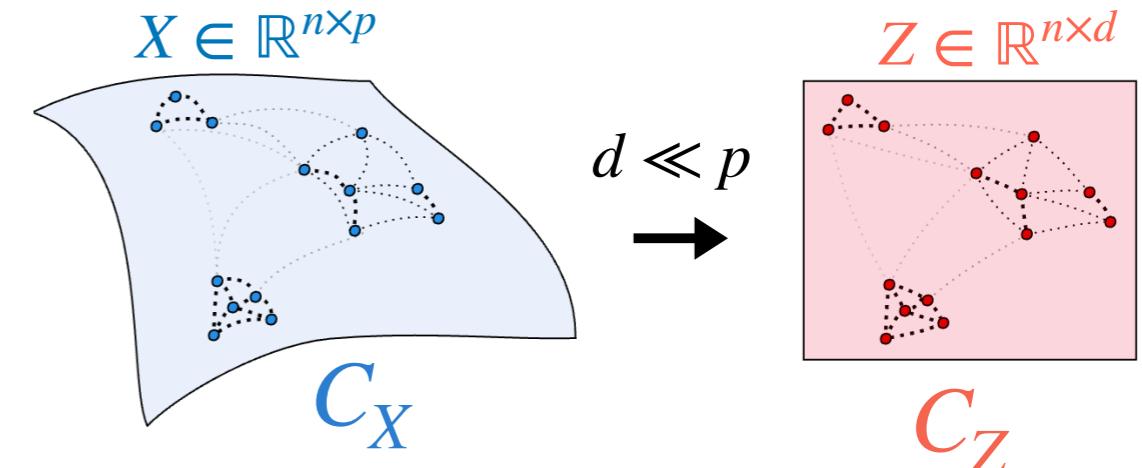
Gromov-Monge

↑?
?

♦ Gromov-Wasserstein projection

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T \in \Pi(\frac{\mathbf{1}_n}{n}, \frac{\mathbf{1}_n}{n})} \sum_{ijkl} L\left([C_X]_{ik}, [C_Z]_{jl}\right) T_{ij}T_{kl}$$

Gromov-Wasserstein



♦ Equivalence holds for

Spectral methods

♦ C_X any matrix, $L = |\cdot|^2$, $C_Z = ZZ^\top$

A is CPD : $\forall x$ **s.t.** $x^\top 1 = 0$, $x^\top Ax \geq 0$

Neighbor embedding methods

♦ C_X is CPD, $L = KL$

$$C_Z = \text{diag}(\alpha_Z) K_Z \text{ diag}(\beta_Z)$$

where $\log(K_Z)$ is CPD

DR as OT in disguise

♦ Dimension reduction

$$\min_{Z \in \mathbb{R}^{n \times d}} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{ij}\right)$$

equiv
↔

Permutation equivariance

$$\forall P, C_{PZ} = PC_ZP^\top$$

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{\sigma \in S_n} \sum_{i,j=1}^n L\left([C_X]_{ij}, [C_Z]_{\sigma(i)\sigma(j)}\right)$$

equiv
↔

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_P \sum_{i,j,k,l=1}^n L\left([C_X]_{ik}, [C_Z]_{jl}\right) P_{ij}P_{kl}$$

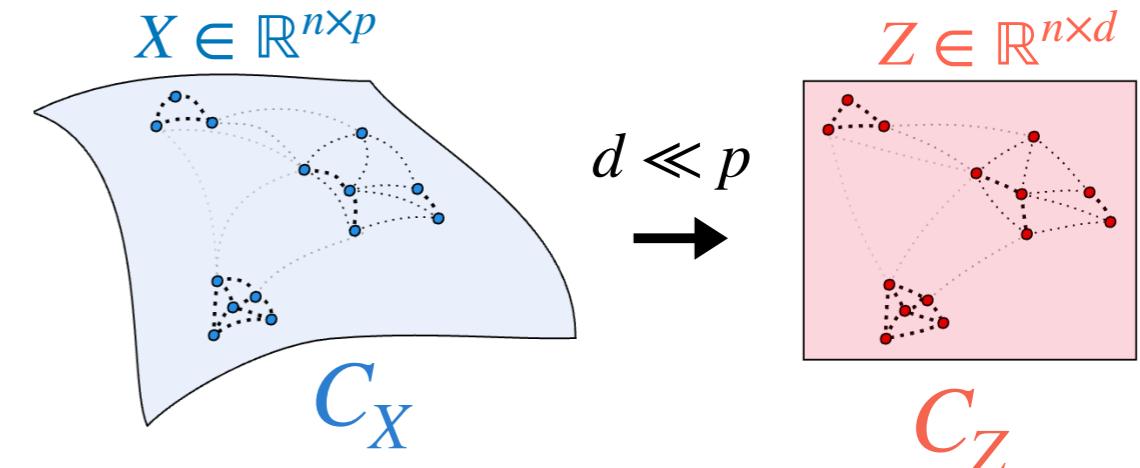
Gromov-Monge

↑?
?

♦ Gromov-Wasserstein projection

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T \in \Pi(\frac{1_n}{n}, \frac{1_n}{n})} \sum_{ijkl} L\left([C_X]_{ik}, [C_Z]_{jl}\right) T_{ij}T_{kl}$$

Gromov-Wasserstein



♦ Equivalence holds for

Spectral methods

♦ C_X any matrix, $L = |\cdot|^2$, $C_Z = ZZ^\top$

A is CPD : $\forall x$ **s.t.** $x^\top 1 = 0$, $x^\top Ax \geq 0$

Neighbor embedding methods

♦ C_X is CPD, $L = KL$

$$C_Z = \text{diag}(\alpha_Z) K_Z \text{ diag}(\beta_Z)$$

where $\log(K_Z)$ is CPD

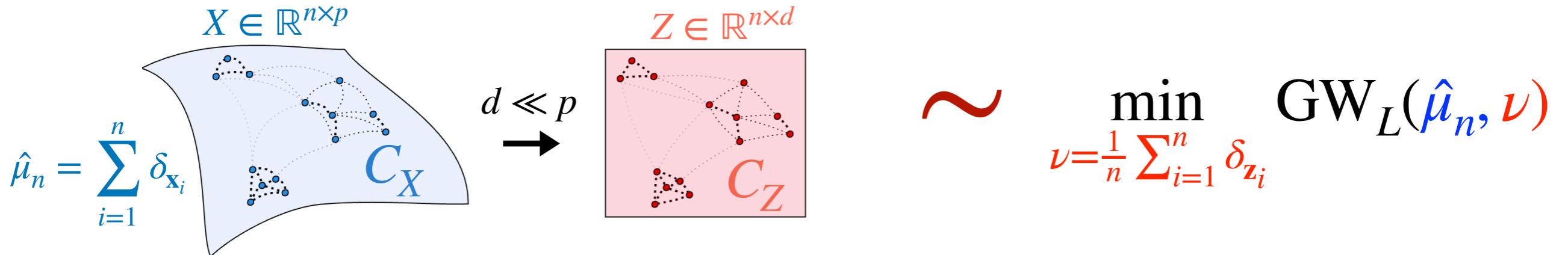
| e.g. $K_Z = \exp(-\|\mathbf{z}_i - \mathbf{z}_j\|_2^2)$
and its usual normalizations

$$1_n^\top K_Z 1_n = 1, K_Z 1_n = 1_n, + K_Z^\top 1_n = 1_n \\ K_Z 1_n = 1_n$$

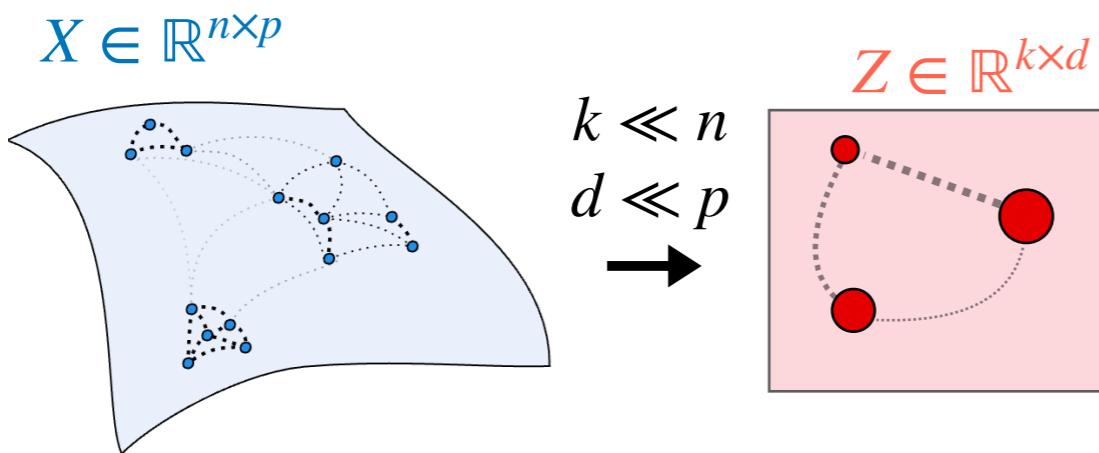
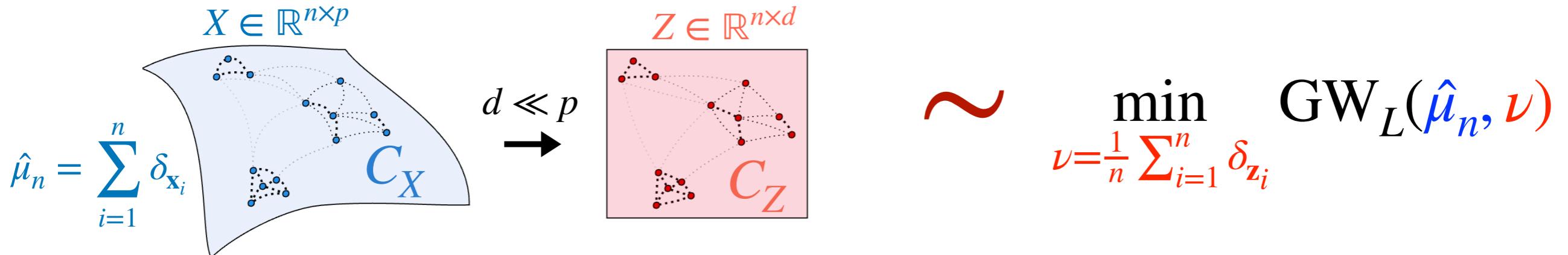
(Sinkhorn & Knopp, 1967)

| Beware that C_X is not always CPD.

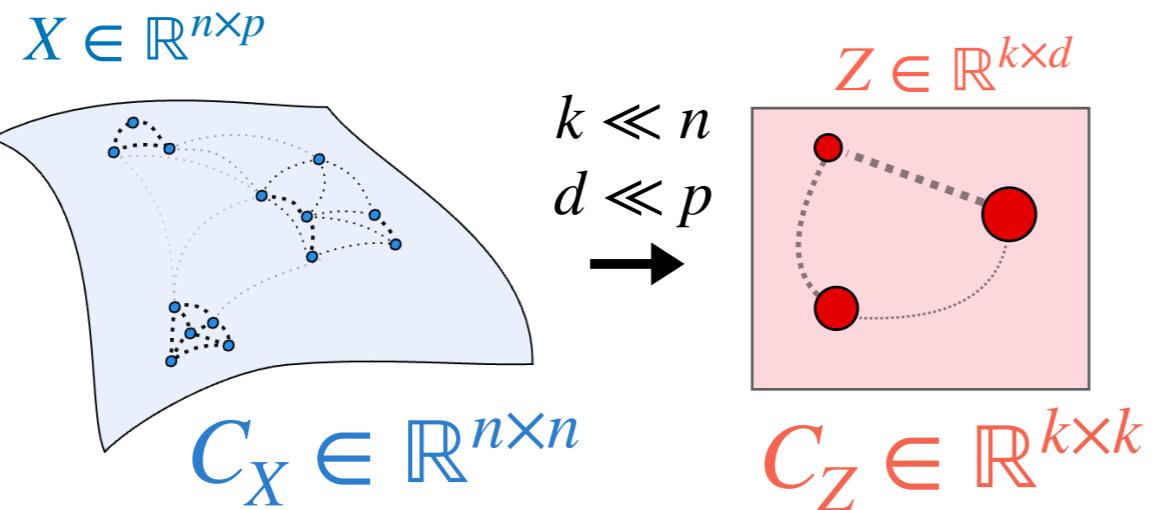
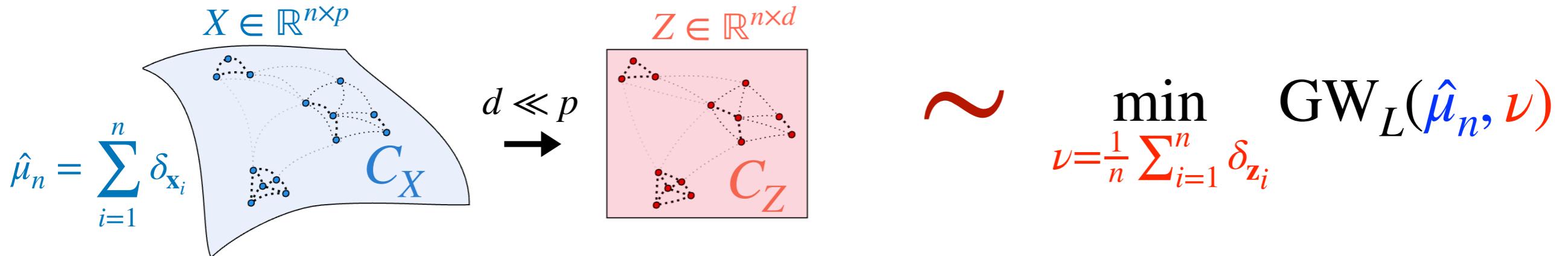
Distributional Reduction



Distributional Reduction



Distributional Reduction

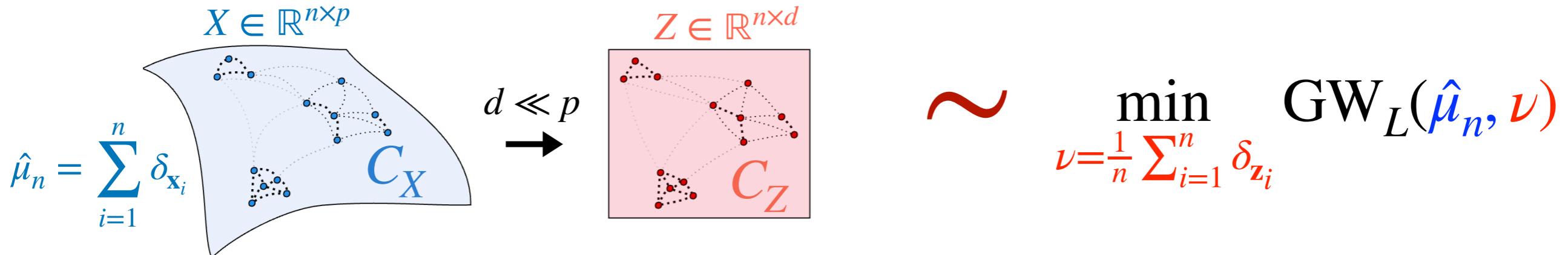


◆ **GW projection**

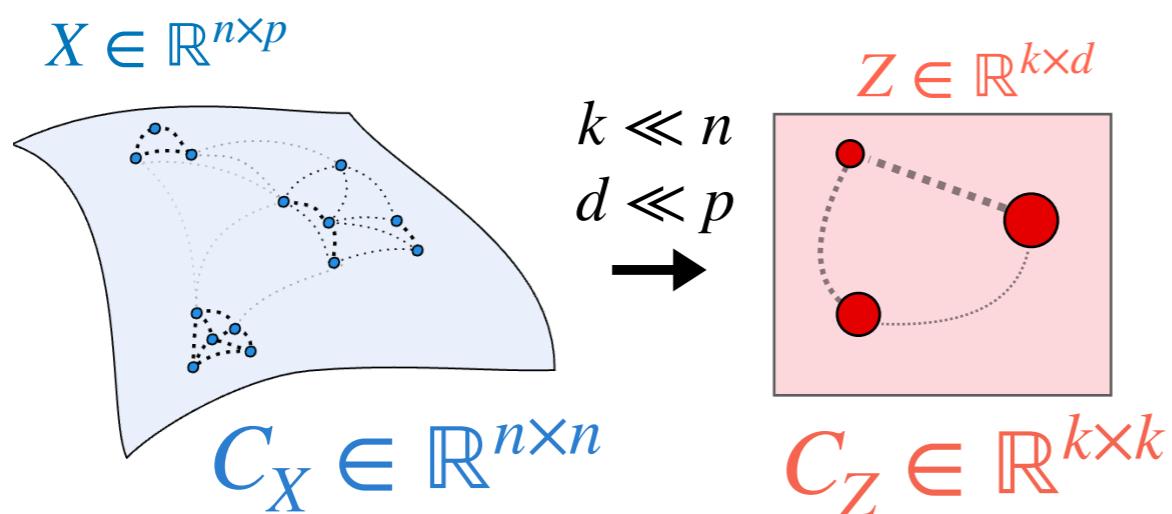
$$\min_{\nu \in \mathcal{P}_k(\mathbb{R}^d)} \text{GW}(\hat{\mu}_n, \nu)$$

$$\nu = \sum_{j=1}^k b_j \delta_{\mathbf{z}_j}$$

Distributional Reduction



$$\sim \min_{\nu = \frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{z}_i}} \text{GW}_L(\hat{\mu}_n, \nu)$$



◆ GW projection

$$\min_{\nu \in \mathcal{P}_k(\mathbb{R}^d)} \text{GW}(\hat{\mu}_n, \nu)$$

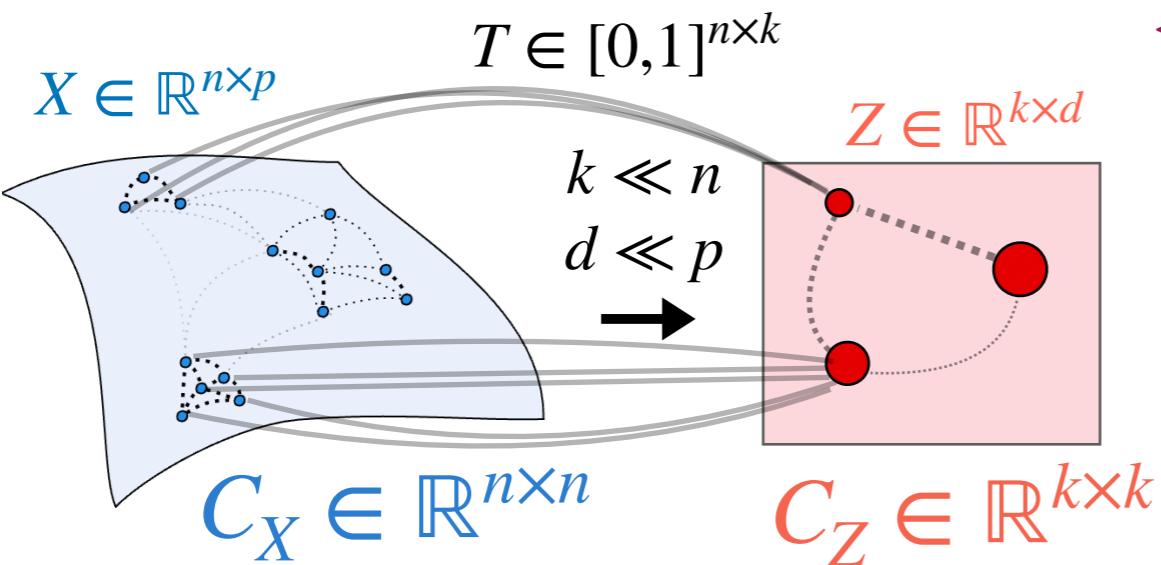
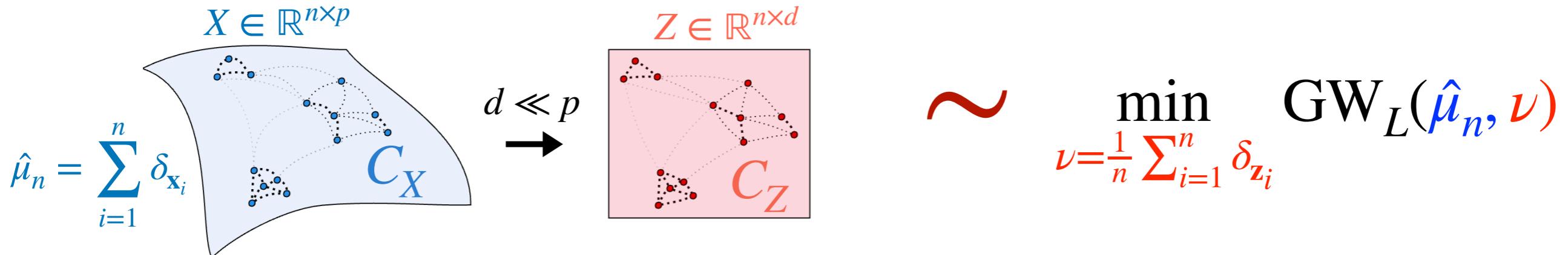
$$\nu = \sum_{j=1}^k b_j \delta_{\mathbf{z}_j}$$

◆ Optimization problem

$$\min_{Z \in \mathbb{R}^{k \times d}} \min_{\mathbf{b} \in \Sigma_k} \text{GW}_L(C_X, C_Z, \frac{1}{n}, \mathbf{b})$$

- ◆ Find few prototypes in low dim.
- ◆ Find the weights/cluster size

Distributional Reduction



◆ **GW projection**

$$\min_{\nu \in \mathcal{P}_k(\mathbb{R}^d)} \text{GW}(\hat{\mu}_n, \nu)$$

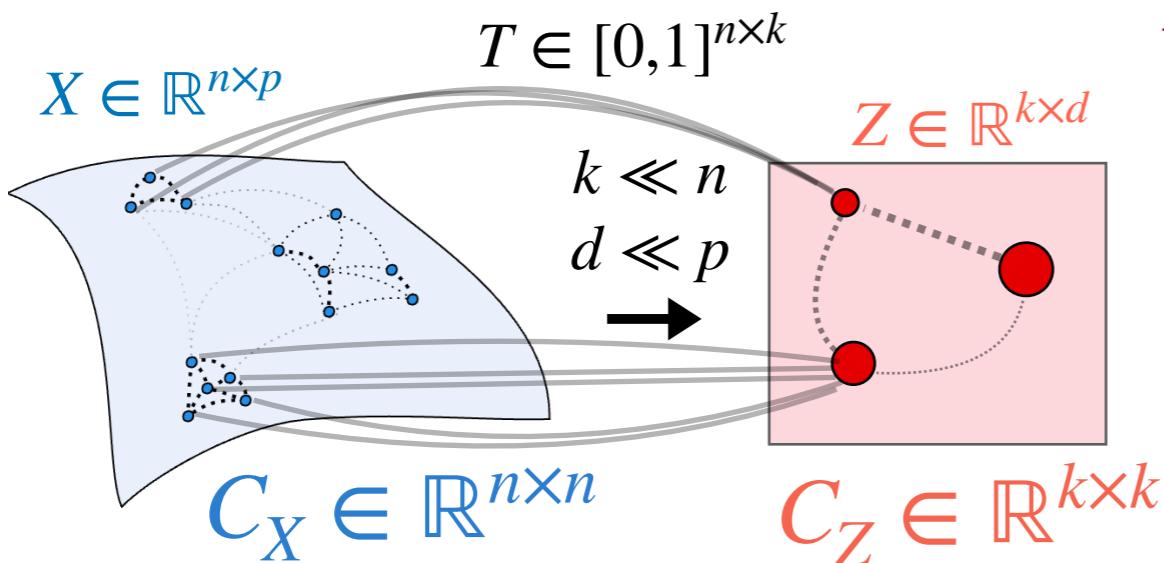
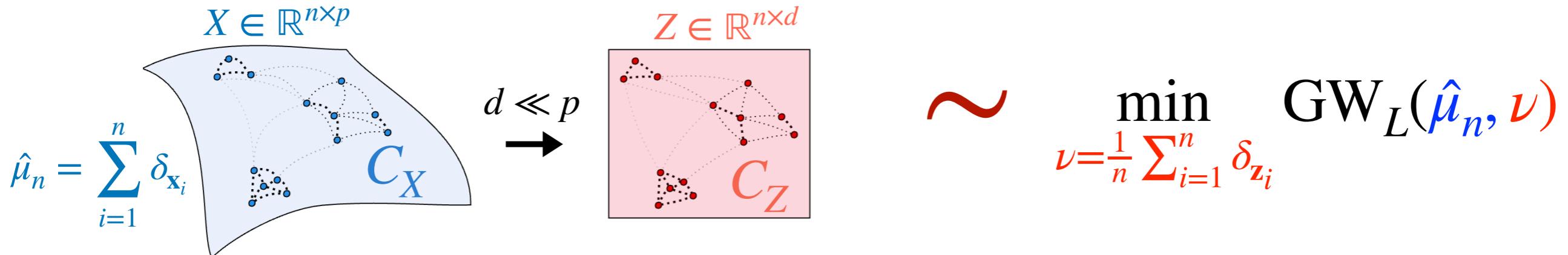
$$\nu = \sum_{j=1}^k b_j \delta_{\mathbf{z}_j}$$

◆ **Optimization problem**

$$\min_{Z \in \mathbb{R}^{k \times d}} \min_{b \in \Sigma_k} \text{GW}_L(C_X, C_Z, \frac{1}{n}, b)$$

- ◆ Find few prototypes in low dim.
- ◆ Find the weights/cluster size
- ◆ Clustering via the coupling T (soft-assignment)
- ◆ Sufficient conditions for hard assignment (see paper)

Distributional Reduction



◆ GW projection

$$\min_{\nu \in \mathcal{P}_k(\mathbb{R}^d)} \text{GW}(\hat{\mu}_n, \nu)$$

$$\nu = \sum_{j=1}^k b_j \delta_{\mathbf{z}_j}$$

◆ Optimization problem

$$\min_{Z \in \mathbb{R}^{k \times d}} \min_{b \in \Sigma_k} \text{GW}_L(C_X, C_Z, \frac{1}{n}, b)$$

- ◆ Find few prototypes in low dim.
- ◆ Find the weights/cluster size
- ◆ Clustering via the coupling T (soft-assignment)
- ◆ Sufficient conditions for hard assignment (see paper)

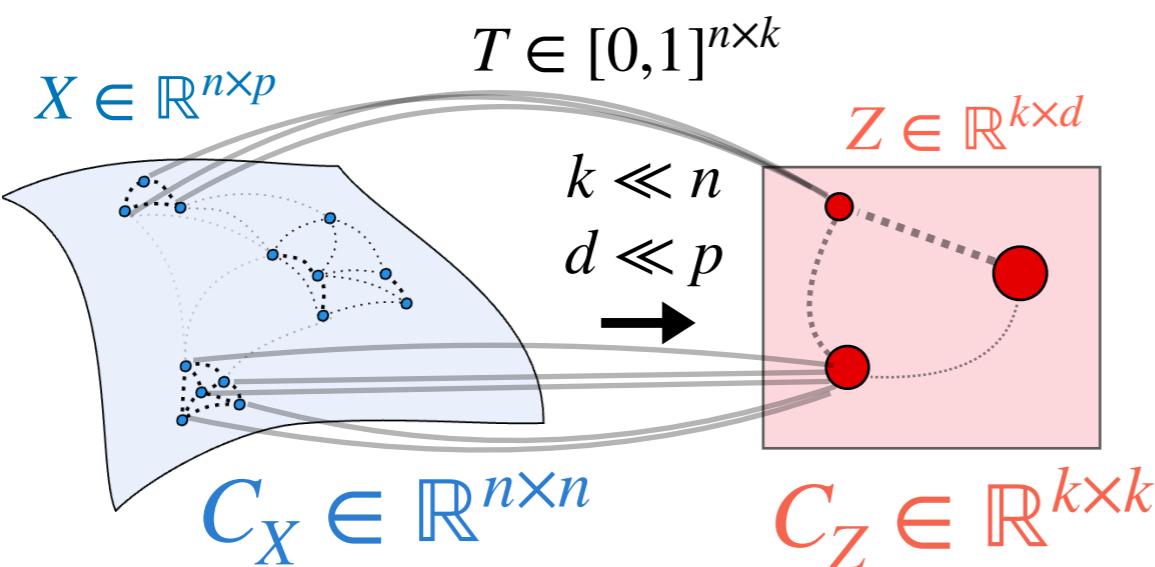
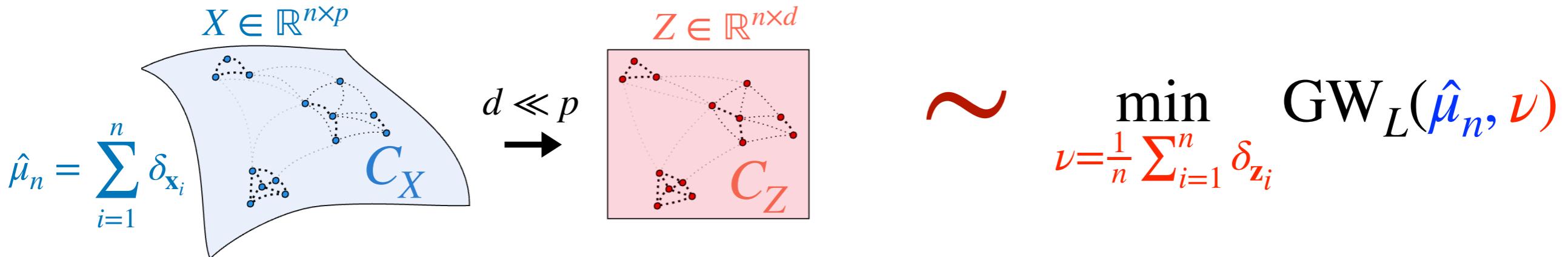
◆ A semi-relaxed objective

(Vincent-Cuaz et al., 2022)

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T: T1_k = \frac{1}{n}} \sum_{ijkl} L([C_X]_{ik}, [C_Z]_{jl}) T_{ij} T_{kl}$$

easier than GW

Distributional Reduction



◆ GW projection

$$\min_{\nu \in \mathcal{P}_k(\mathbb{R}^d)} \text{GW}(\hat{\mu}_n, \nu)$$

$$\nu = \sum_{j=1}^k b_j \delta_{\mathbf{z}_j}$$

◆ Optimization problem

$$\min_{Z \in \mathbb{R}^{k \times d}} \min_{b \in \Sigma_k} \text{GW}_L(C_X, C_Z, \frac{1}{n}, b)$$

- ◆ Find few prototypes in low dim.
- ◆ Find the weights/cluster size
- ◆ Clustering via the coupling T (soft-assignment)
- ◆ Sufficient conditions for hard assignment (see paper)

◆ A semi-relaxed objective

(Vincent-Cuaz et al., 2022)

$$\min_{Z \in \mathbb{R}^{n \times d}} \min_{T: T1_k = \frac{1}{n}} \sum_{ijkl} L([C_X]_{ik}, [C_Z]_{jl}) T_{ij} T_{kl}$$

→ easier than GW

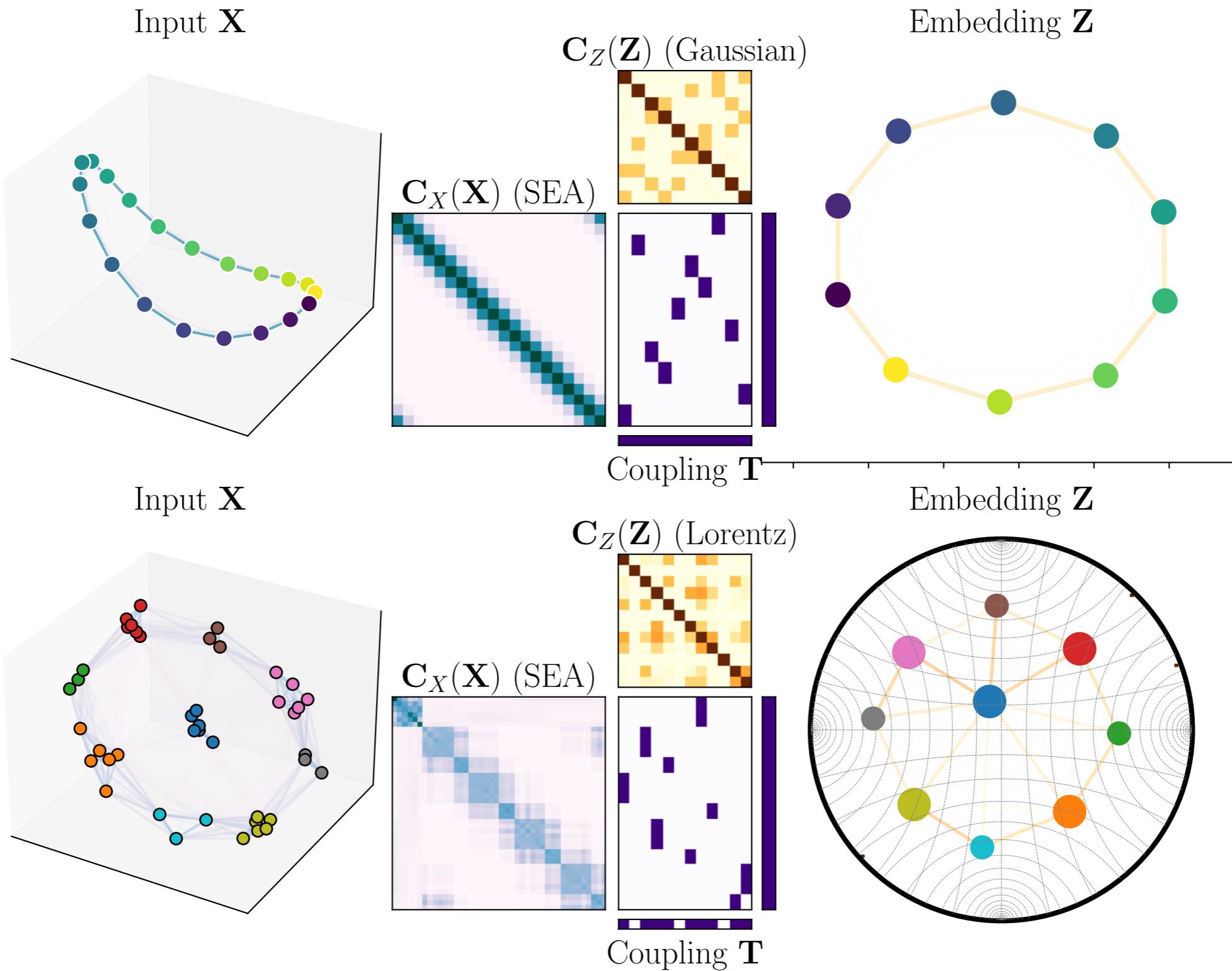
- ◆ Non-convex problem

- ◆ BCD: alternates optim in Z, in T

- ◆ Optim in T: CG solver in $O(n^2 k)$ for $L \in \{\text{KL}, |\cdot|^2\}$

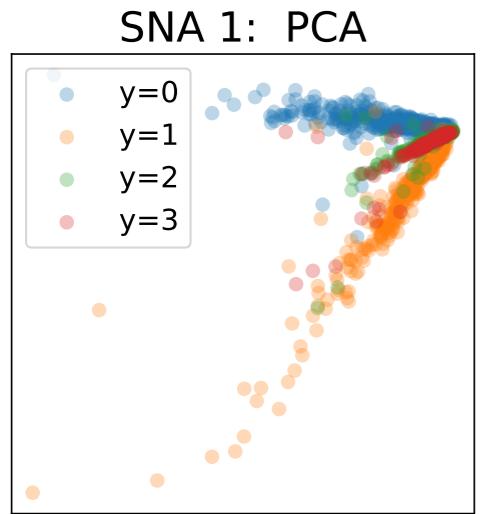
- ◆ With low-rank structures $O(nkr + n^2)$

Distributional Reduction



Distributional Reduction

- ◆ **Single-cell dataset** (Chen et al., 2019) Only solve $\min_{\mathbf{b} \in \Sigma_k} \text{GW}_L(\mathbf{C}_X, \mathbf{C}_Z, \frac{\mathbf{1}_n}{n}, \mathbf{b})$



Distributional Reduction

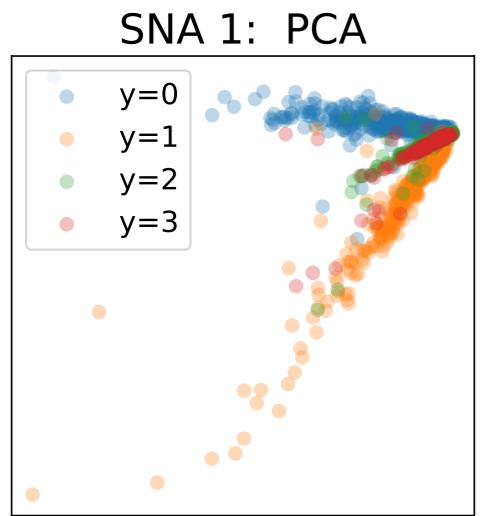
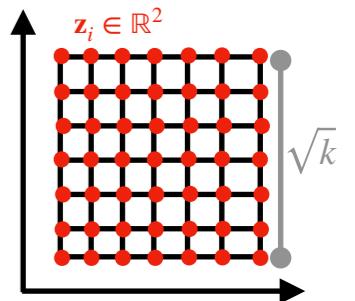
- ◆ **Single-cell dataset** (Chen et al., 2019) Only solve $\min_{\mathbf{b} \in \Sigma_k} \text{GW}_L(\mathbf{C}_X, \mathbf{C}_Z, \frac{\mathbf{1}_n}{n}, \mathbf{b})$ with $\mathbf{C}_X = \mathbf{X}\mathbf{X}^\top$

and

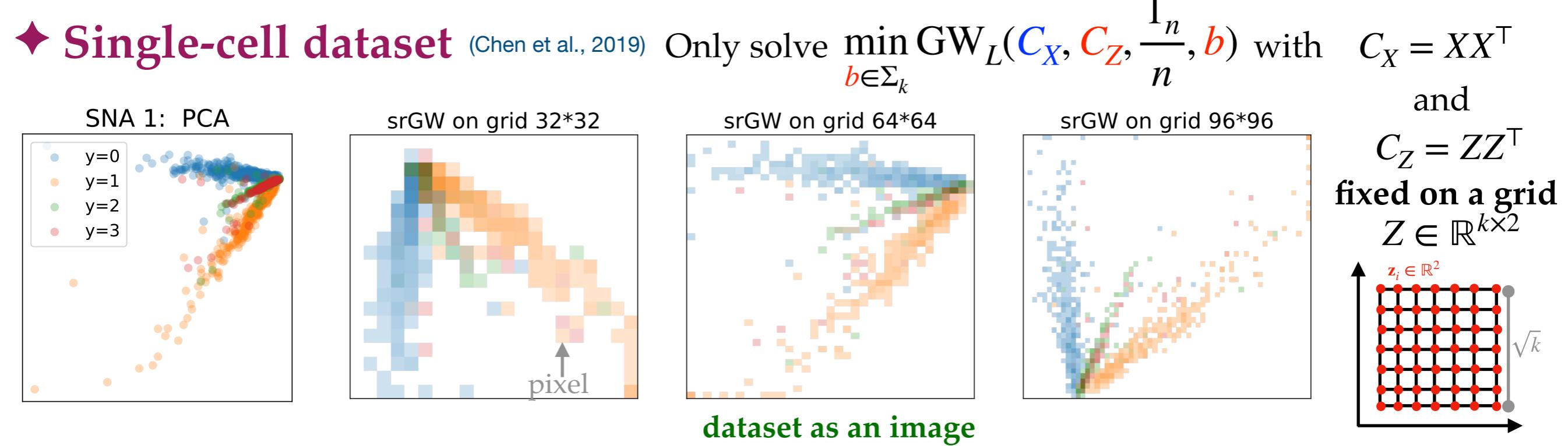
$$\mathbf{C}_Z = \mathbf{Z}\mathbf{Z}^\top$$

fixed on a grid

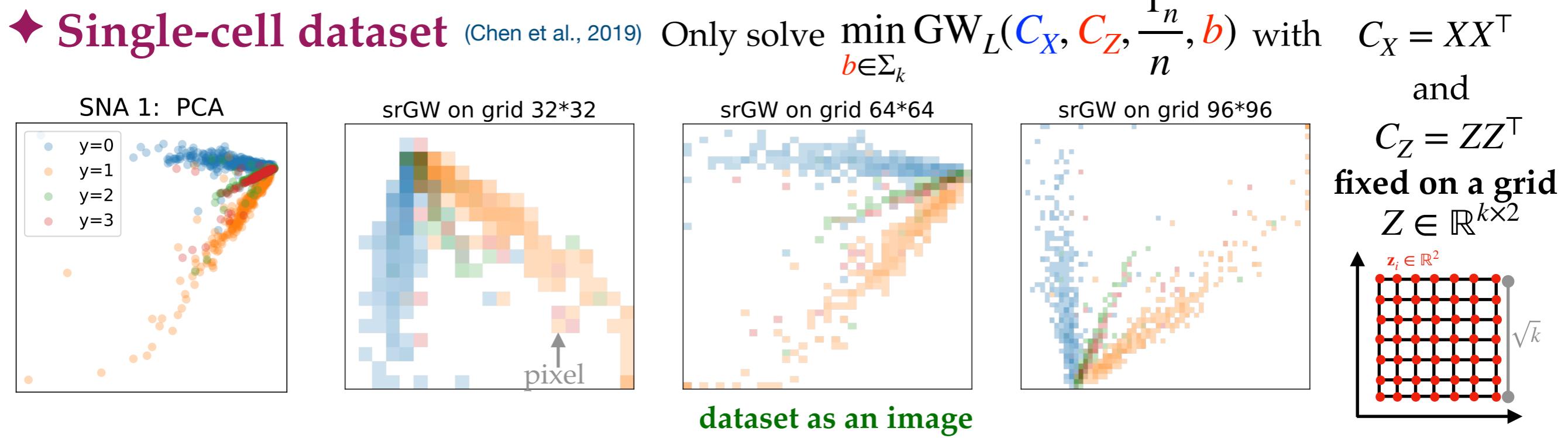
$$\mathbf{Z} \in \mathbb{R}^{k \times 2}$$



Distributional Reduction



Distributional Reduction

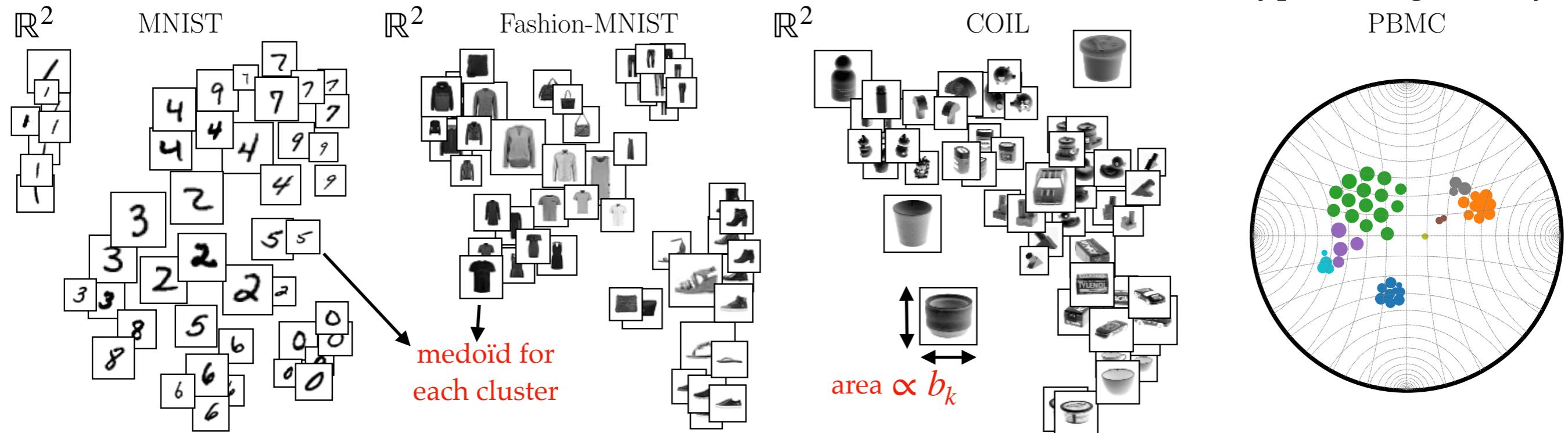


◆ Image datasets

\mathbf{C}_X symmetric entropic aff. (Van Assel et al., 2023) \mathbf{C}_Z Student t-kernel

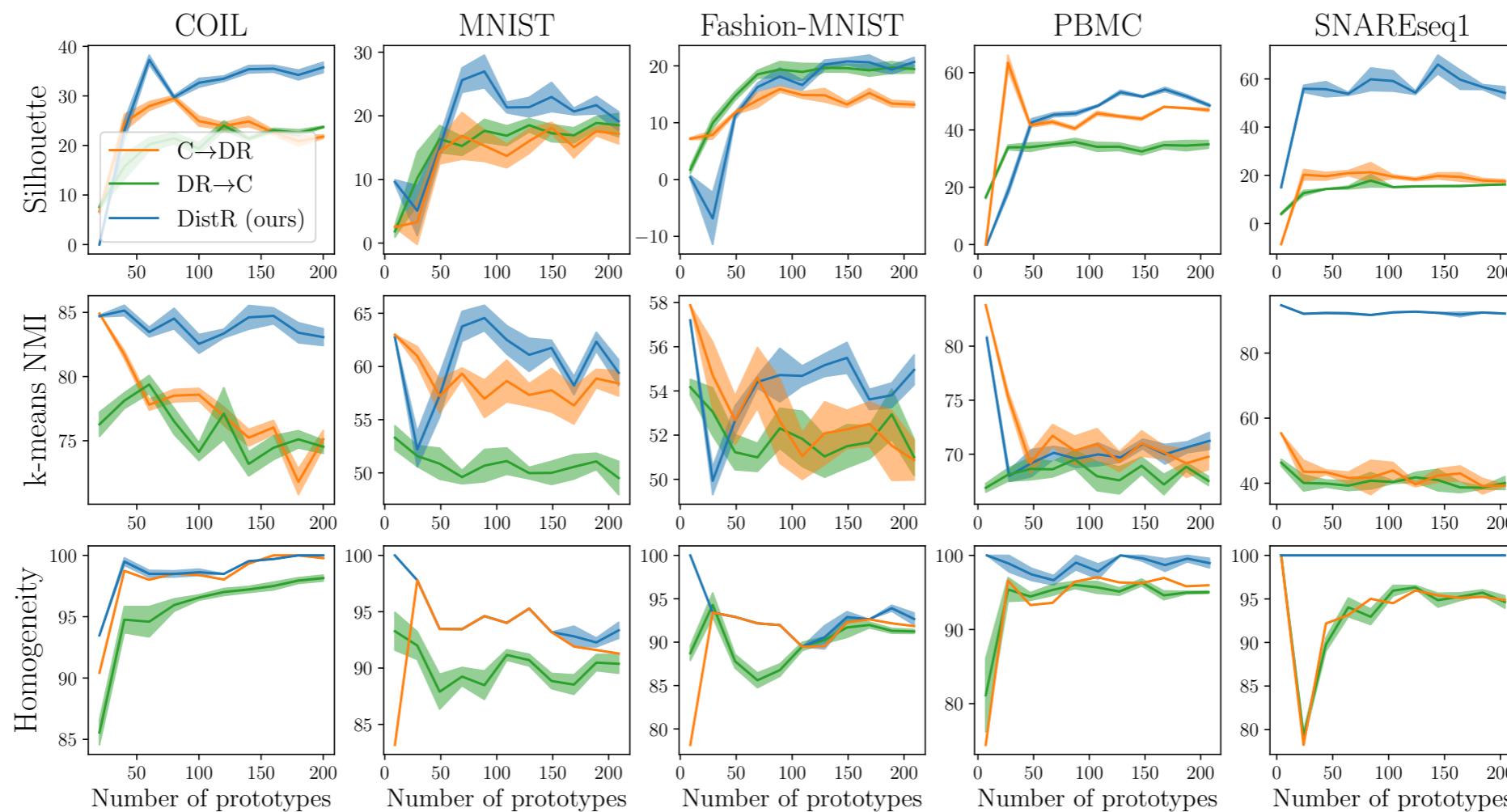
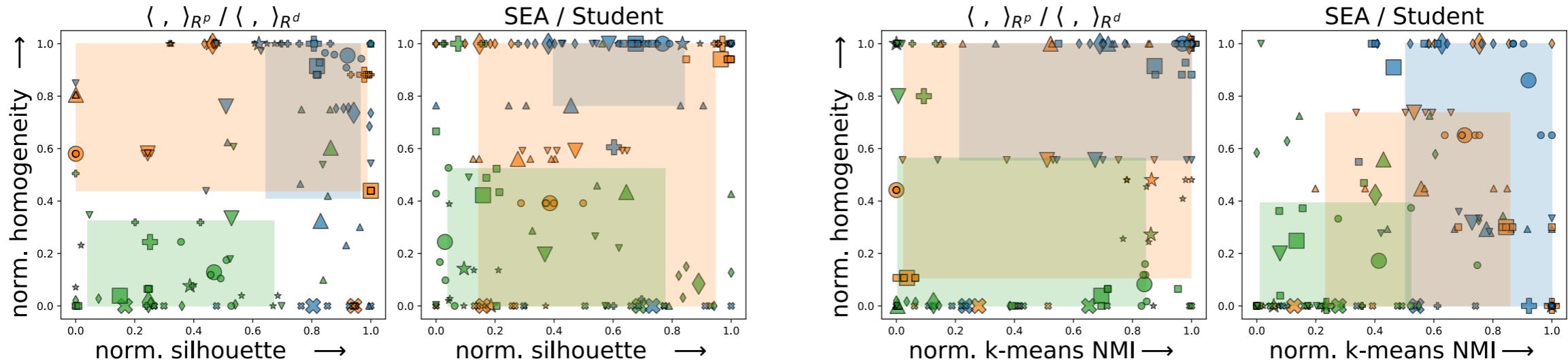
Hyperbolic geometry

PBMC



Distributional Reduction

♦ Comparison with DR then clustering or clustering then DR



Thank you!

Open-source implementations are available here:

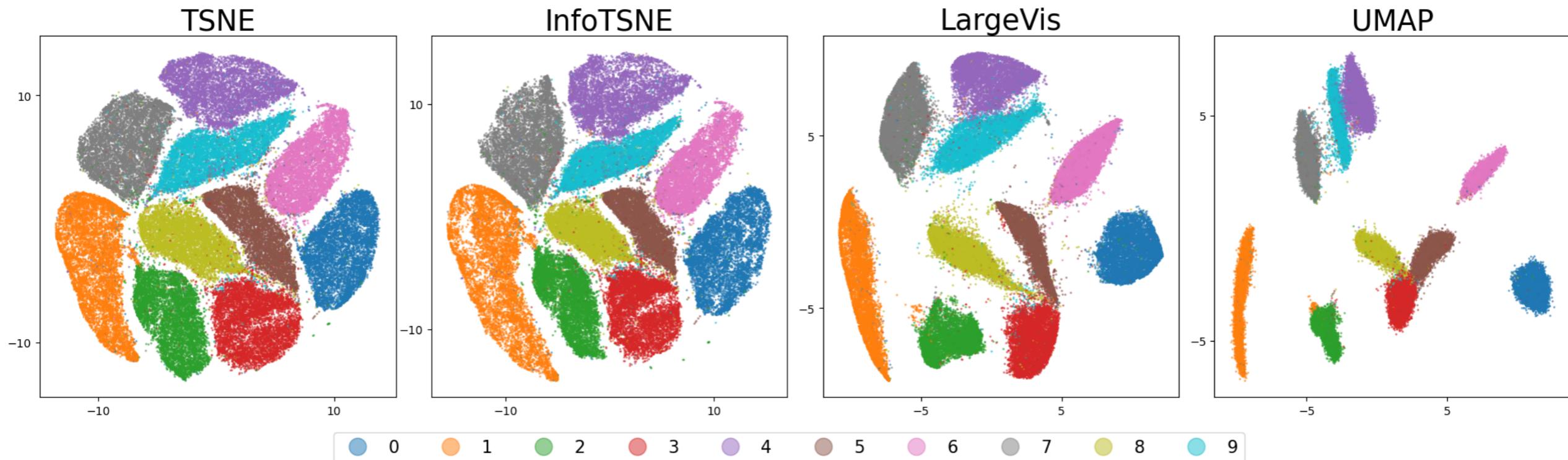
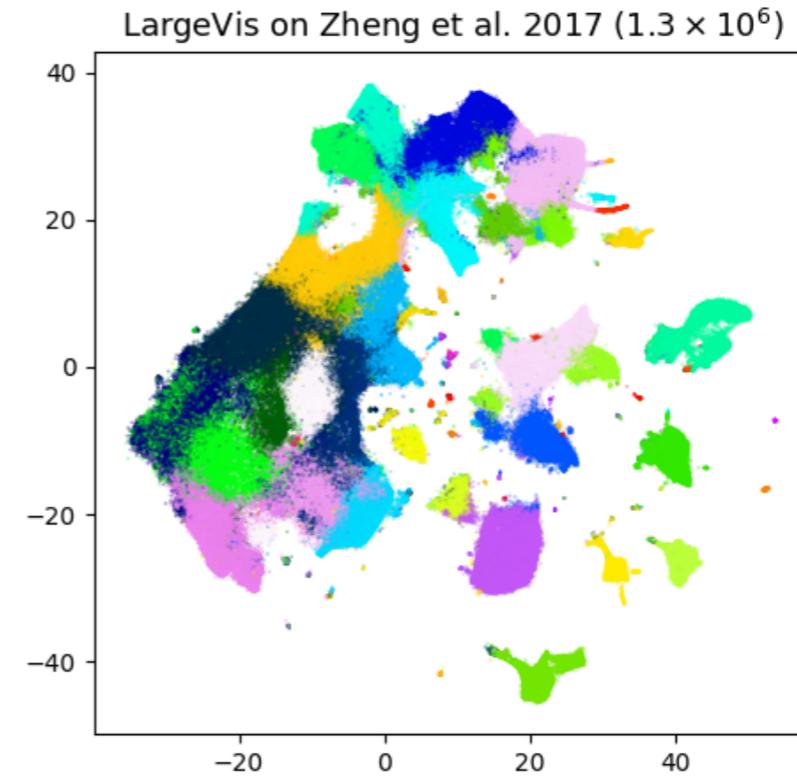
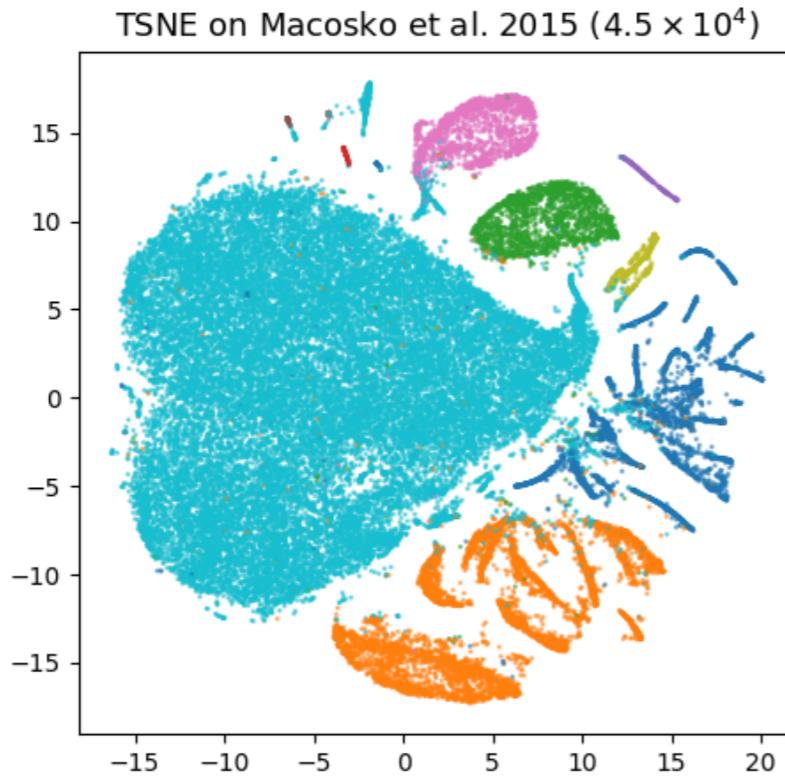


TorchDR

```
from sklearn.datasets import fetch_openml  
from torchdr import PCA, TSNE  
  
x = fetch_openml("mnist_784").data.astype("float32")  
  
x_ = PCA(n_components=50).fit_transform(x)  
z = TSNE(perplexity=30).fit_transform(x_)
```

Modularity	All of it is written in python in a highly modular way, making it easy to create or transform components.
Speed	Supports GPU acceleration , leverages sparsity and batching strategies with contrastive learning techniques.
Memory efficiency	Relies on sparsity and/or pykeops [19] symbolic tensors to avoid memory overflows .
Compatibility	Implemented methods are fully compatible with the sklearn [21] API and torch [20] ecosystem.

Thank you!



Appendix

Solving OT

A linear problem

Discrete probability measures

$$\mu = \sum_{i=1}^n a_i \delta_{x_i} \quad \nu = \sum_{j=1}^m b_j \delta_{y_j}$$

Linear Program:

$$\min_{\pi \in \Pi(\mathbf{a}, \mathbf{b})} \sum_{ij} c_{i,j} \pi_{i,j} = \min_{\pi \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \boldsymbol{\pi} \rangle$$

| Simplex, Network flow, Hungarian algorithms $\sim O(n^3 \log(n))$

Uniform weights

$$\mathbf{a} = \mathbf{b} = \frac{\mathbf{1}_n}{n}$$

Monge Problem

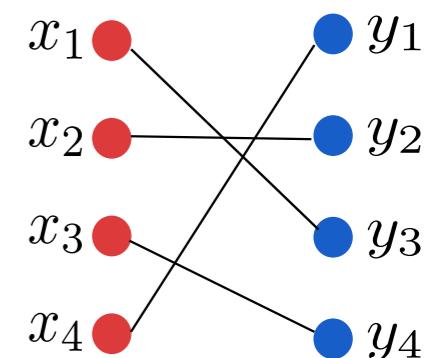
$$\min_{\sigma \in S_n} \sum_{i=1}^n c_{i,\sigma(i)}$$

Fundamental theorem LP:

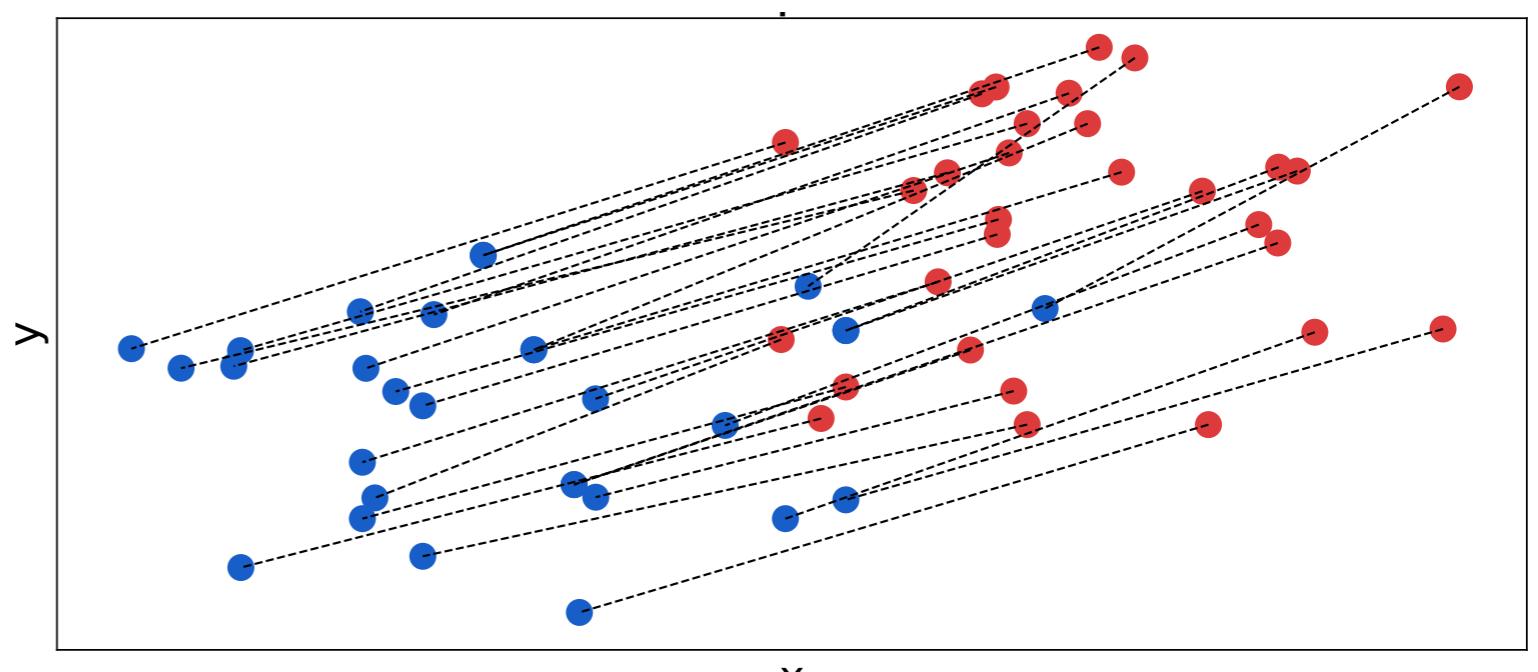
$$\boldsymbol{\pi}^* \leftrightarrow \sigma^* \in S_n$$

Optimal coupling is a permutation

Solves the Monge Problem



One-to-one



Solving OT

Entropic regularization

Discrete probability measures

$$\mu = \sum_{i=1}^n a_i \delta_{x_i} \quad \nu = \sum_{j=1}^m b_j \delta_{y_j}$$

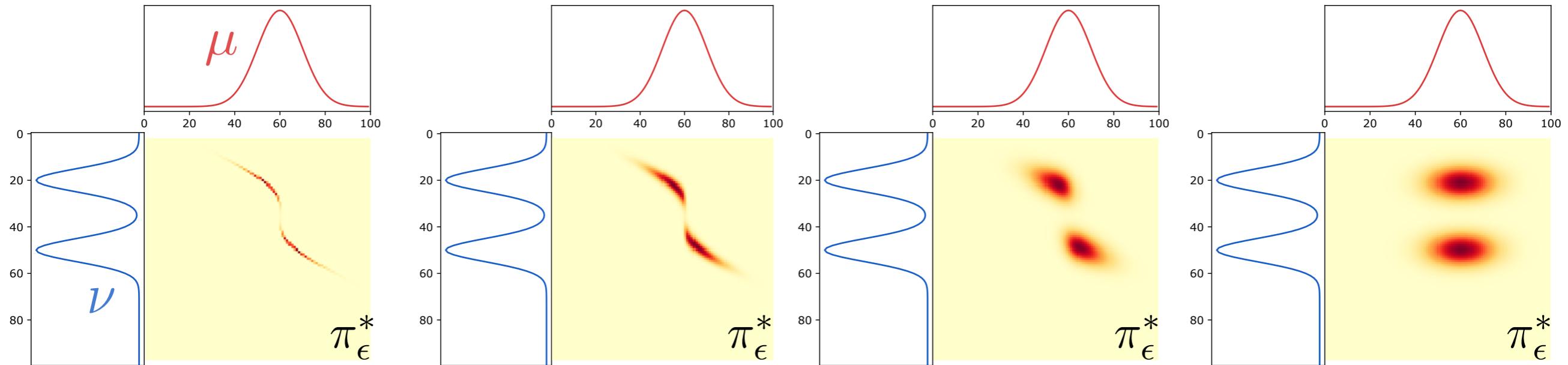
Strongly convex problem:

$$\min_{\pi \in \Pi(\mathbf{a}, \mathbf{b})} \langle \mathbf{C}, \boldsymbol{\pi} \rangle - \varepsilon H(\boldsymbol{\pi})$$

| Entropy term $H(\boldsymbol{\pi}) = - \sum_{ij} (\log(\pi_{ij}) - 1) \pi_{ij}$

| Sinkhorn-Knopp algorithm: 1) fast 2) based on matrix multiplication

| Complexity $O(n^2)$



$0 \leftarrow \epsilon$

$\epsilon \rightarrow +\infty$

Solving OT

Computing GW

Solving FGW: a non convex QP

$$\min_{\pi \in \Pi(\mathbf{a}, \mathbf{b})} \sum_{ijkl} |C_1(i, k) - C_2(j, l)|^p \pi_{ij} \pi_{kl}$$

Quadratic function over polytope -> Conditional Gradient algorithm (a.k.a Frank-Wolfe)

Non convex but converges to a **local optimal solution** [Lacoste-Julien 2016]

Find a **sparse** solution. FW gap = $O(\frac{1}{\sqrt{n_{iter}}})$

Algorithm 1 Conditional Gradient (CG) for FGW

```

1:  $\pi^{(0)} \leftarrow \mathbf{h}\mathbf{g}^\top$ 
2: for  $i = 1, \dots$ , do
3:    $\mathbf{G} \leftarrow$  Gradient from GW loss w.r.t.  $\pi^{(i-1)}$ 
4:    $\tilde{\pi}^{(i)} \leftarrow$  Solve OT with ground loss  $\mathbf{G}$ 
5:    $\tau^{(i)} \leftarrow$  Line-search for GW loss with  $\tau \in (0, 1)$  (closed-form)
6:    $\pi^{(i)} \leftarrow (1 - \tau^{(i)})\pi^{(i-1)} + \tau^{(i)}\tilde{\pi}^{(i)}$ 
7: end for

```

Complexity
 $O(n_{iter} n^3)$

