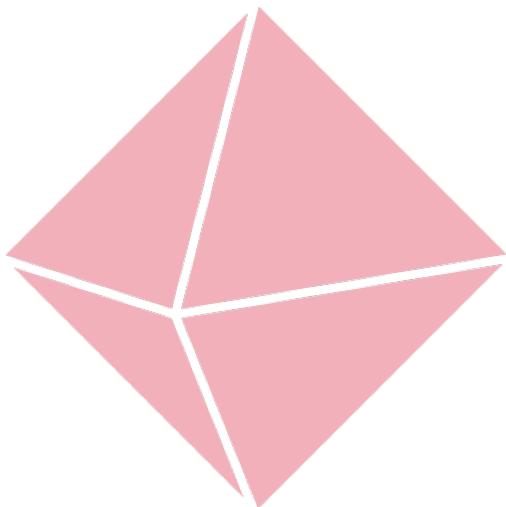


Tutorial: Optimal Transport for Machine Learning

LOGML Summer School 2024

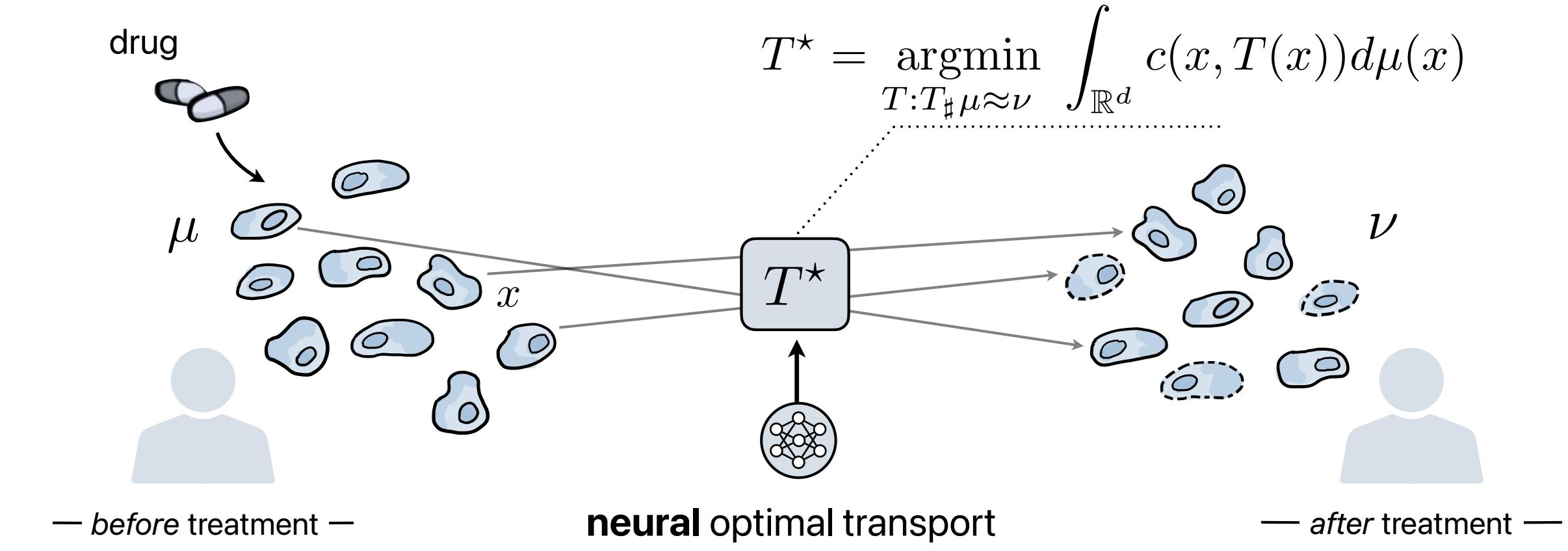
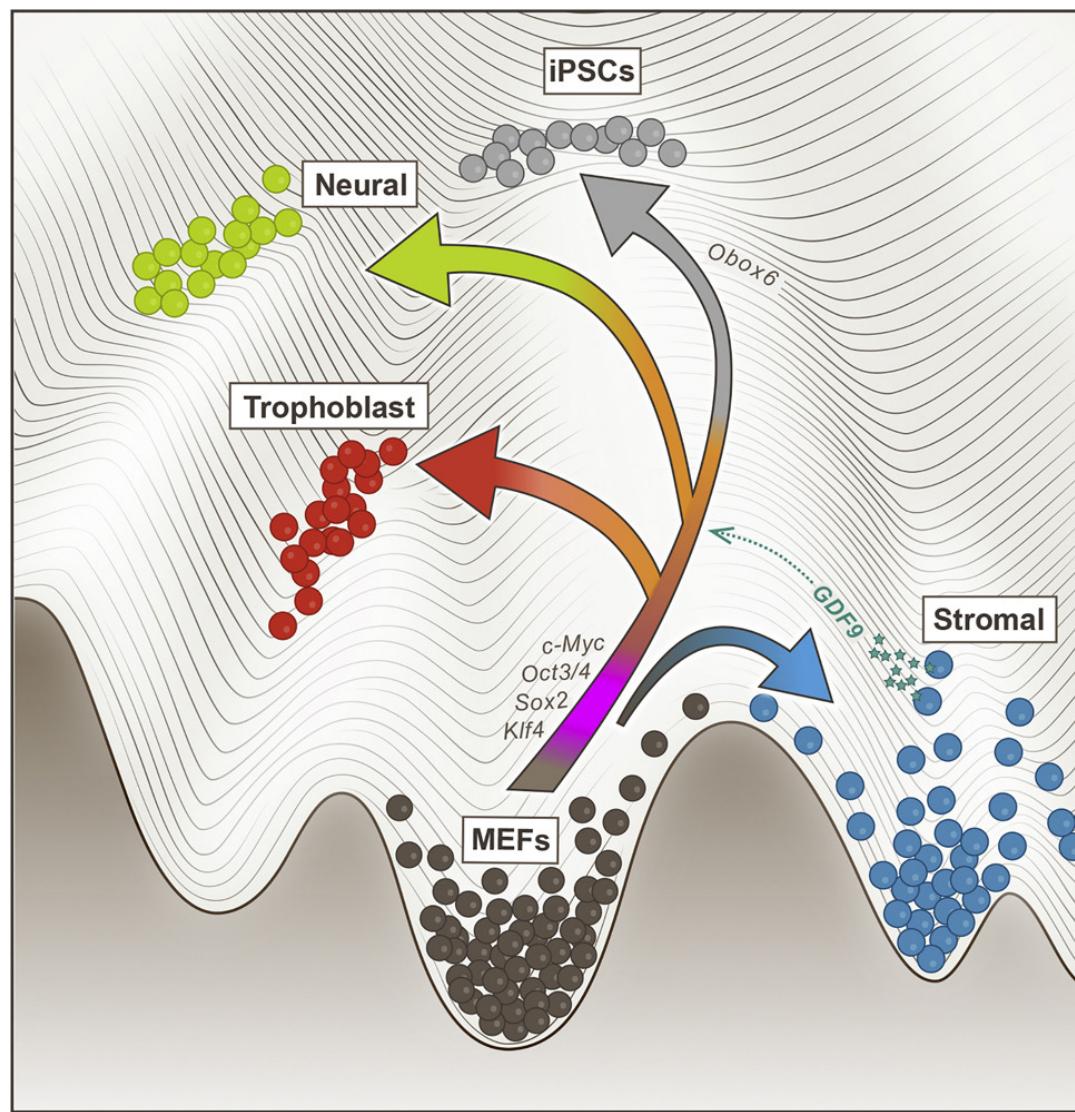
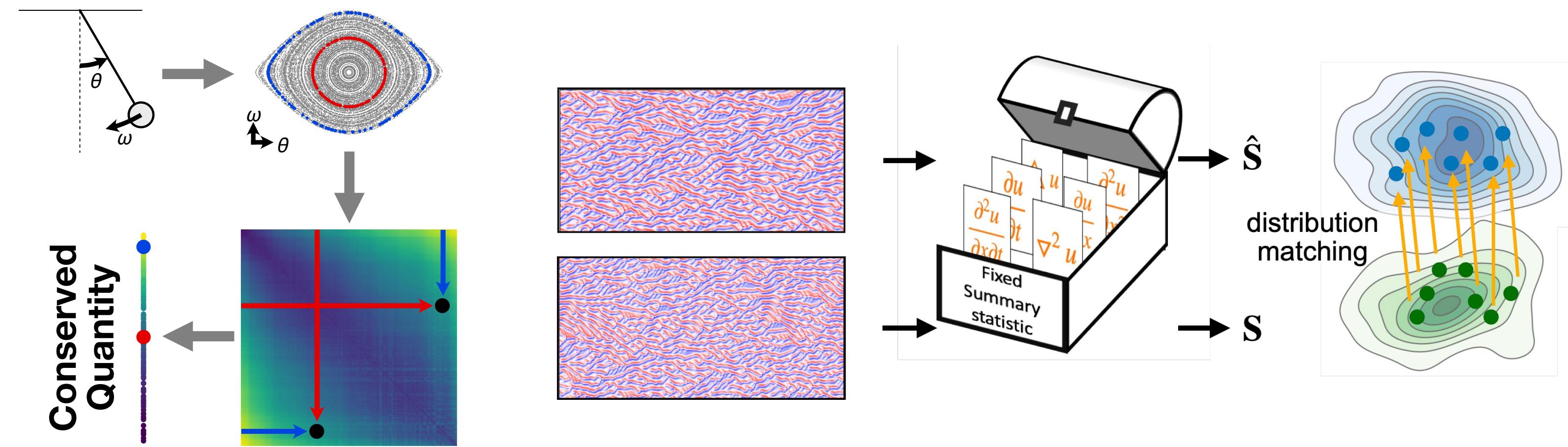
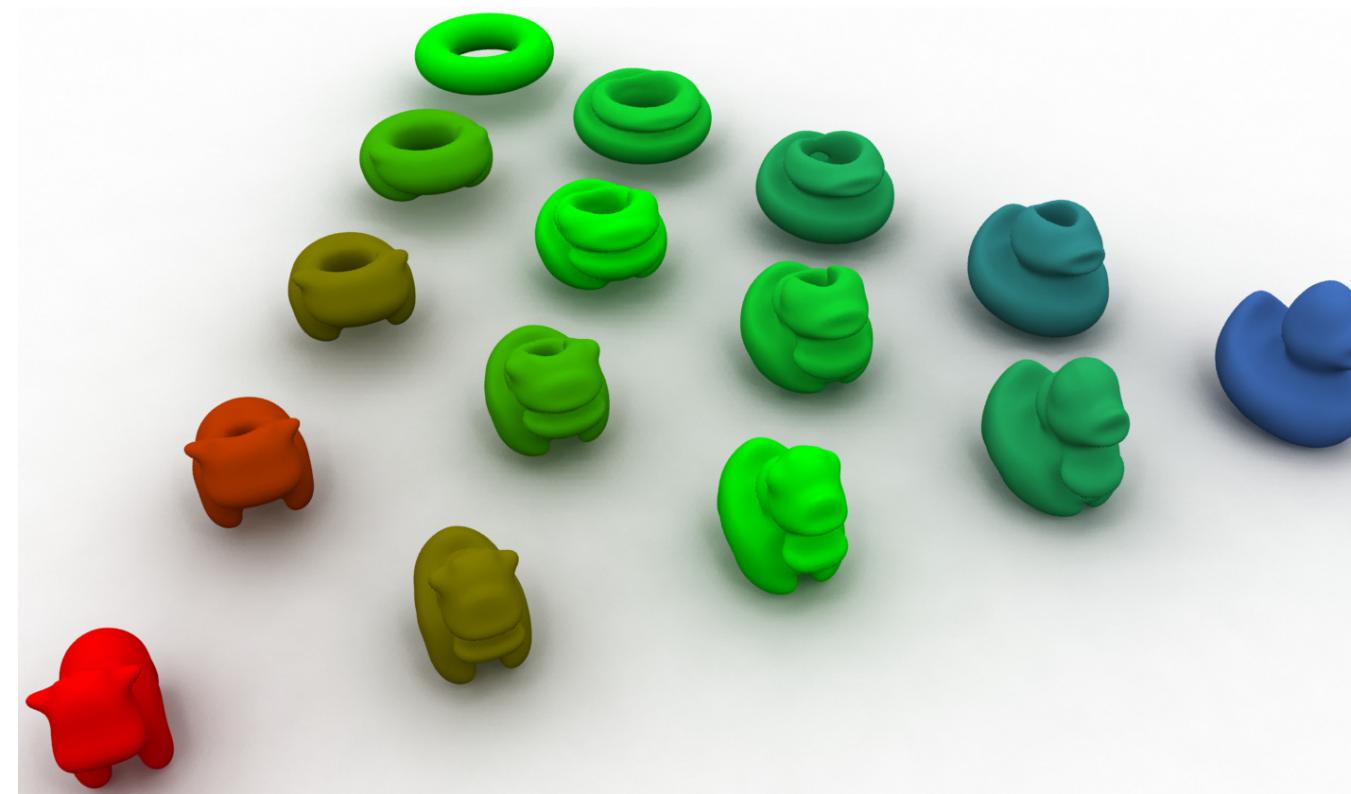
London, UK



Peter Y. Lu



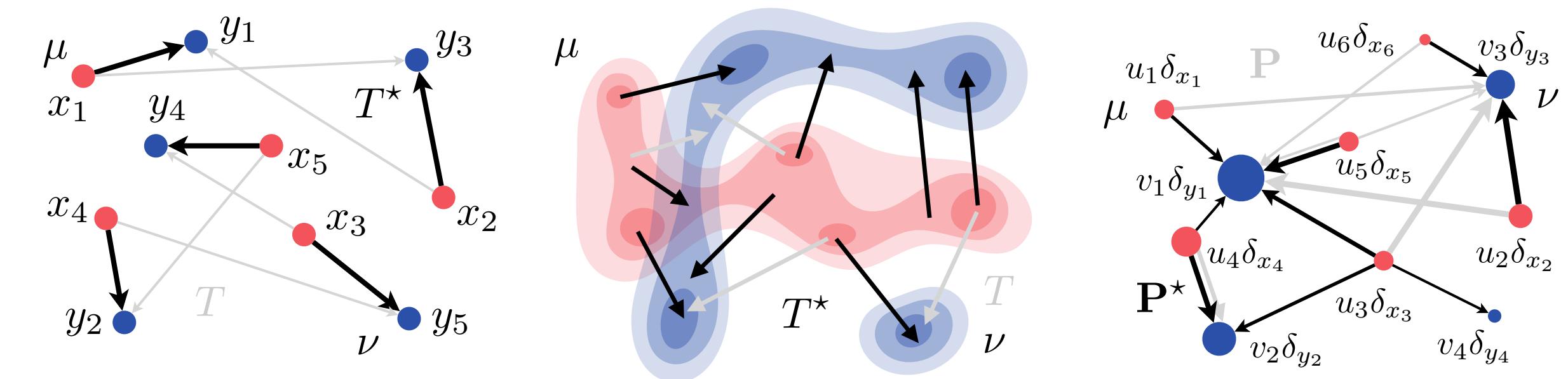
Overview: Optimal Transport for Machine Learning



Overview: Optimal Transport for Machine Learning

Optimal Transport Theory

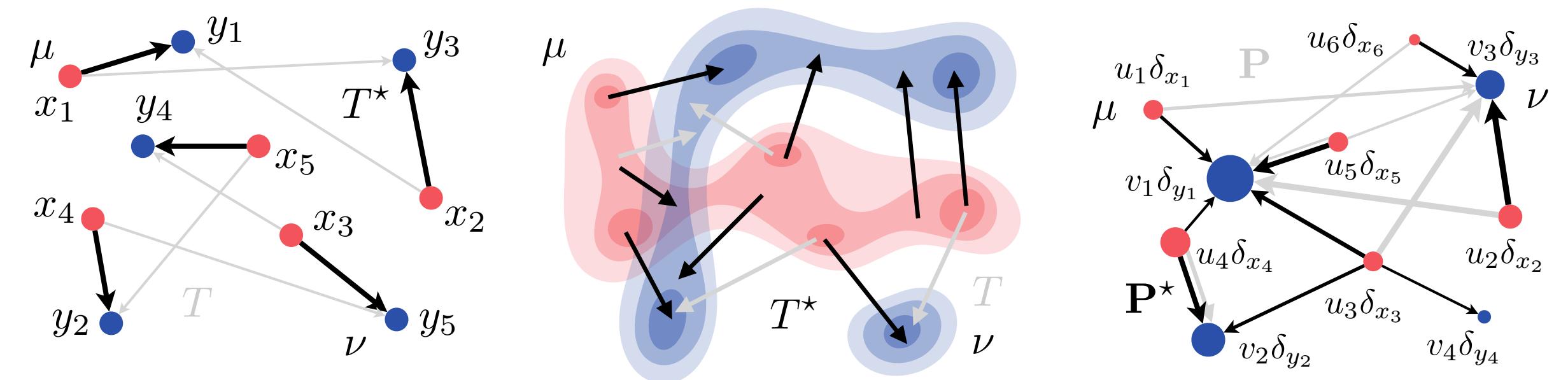
- Monge & Kantorovich Formulations
- Wasserstein Distances



Overview: Optimal Transport for Machine Learning

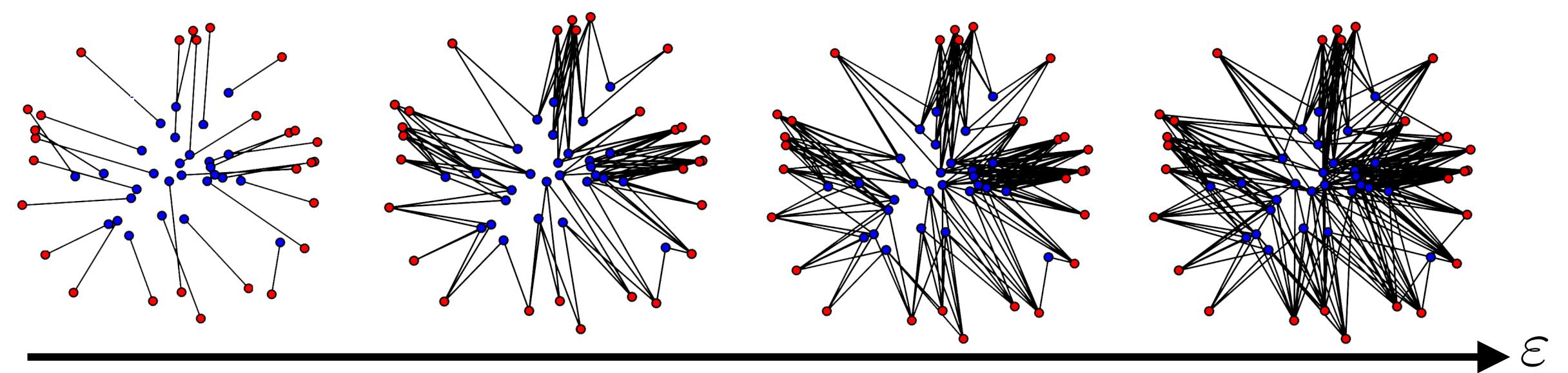
Optimal Transport Theory

- Monge & Kantorovich Formulations
- Wasserstein Distances



Computational Optimal Transport

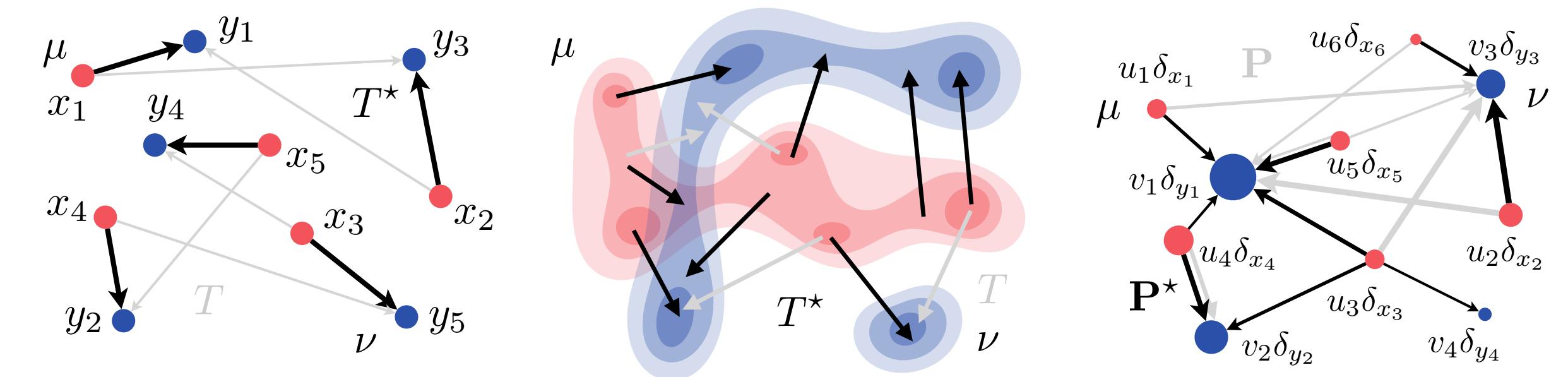
- 1D Optimal Transport
- Entropy-Regularized Optimal Transport



Overview: Optimal Transport for Machine Learning

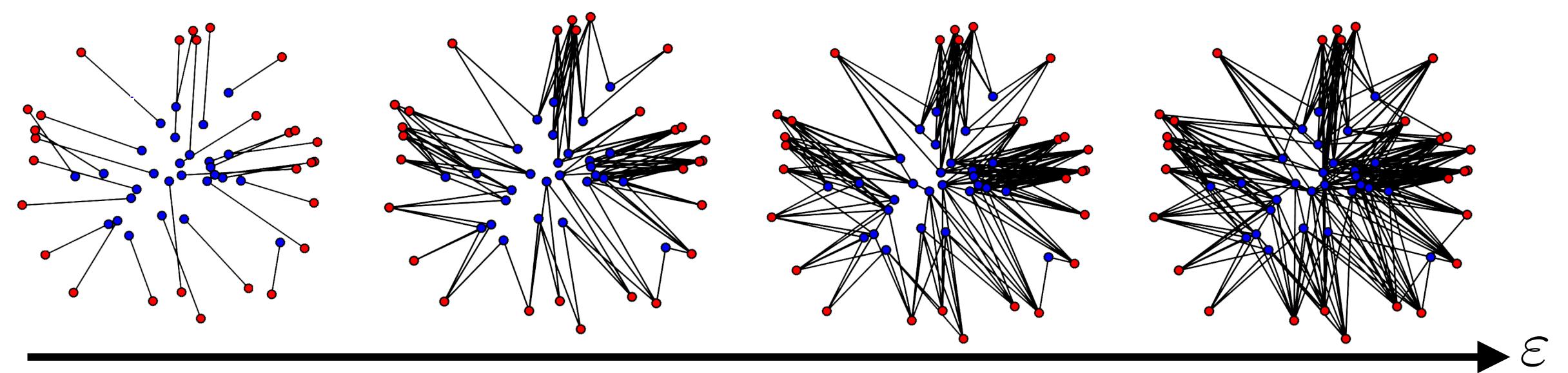
Optimal Transport Theory

- Monge & Kantorovich Formulations
- Wasserstein Distances



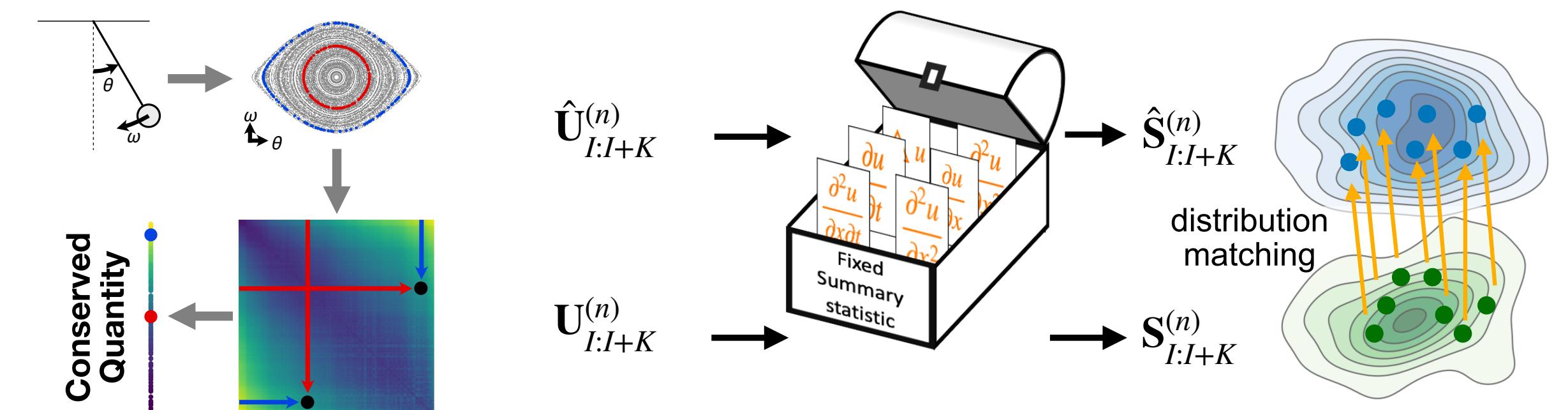
Computational Optimal Transport

- 1D Optimal Transport
- Entropy-Regularized Optimal Transport



Machine Learning Applications

- Shape Analysis
- Deep Learning



Overview: Optimal Transport for Machine Learning

Optimal Transport Theory

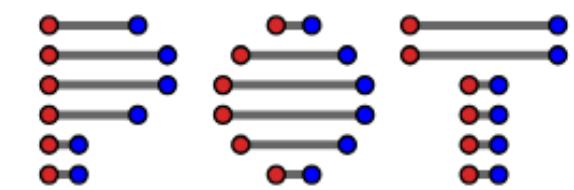
- Monge & Kantorovich Formulations
- Wasserstein Distances

Computational Optimal Transport

- 1D Optimal Transport
- Entropy-Regularized Optimal Transport

Machine Learning Applications

- Shape Analysis
- Deep Learning



POT: Python Optimal Transport (*NumPy*)

pythonot.github.io



OTT: Optimal Transport Tools (*JAX*)

ott-jax.readthedocs.io



GeomLoss: Geometric Loss (*PyTorch*)

kernel-operations.io/geomloss

Gabriel Peyré & Marco Cuturi

Computational Optimal Transport

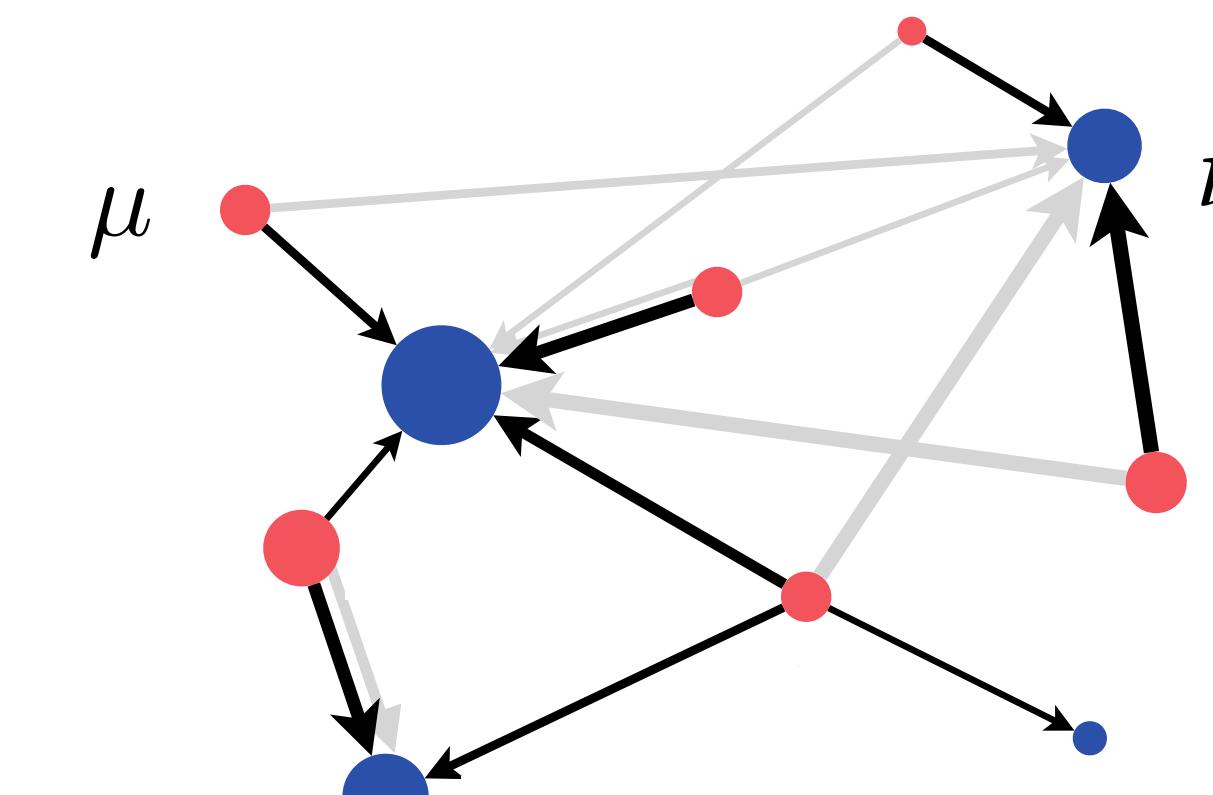
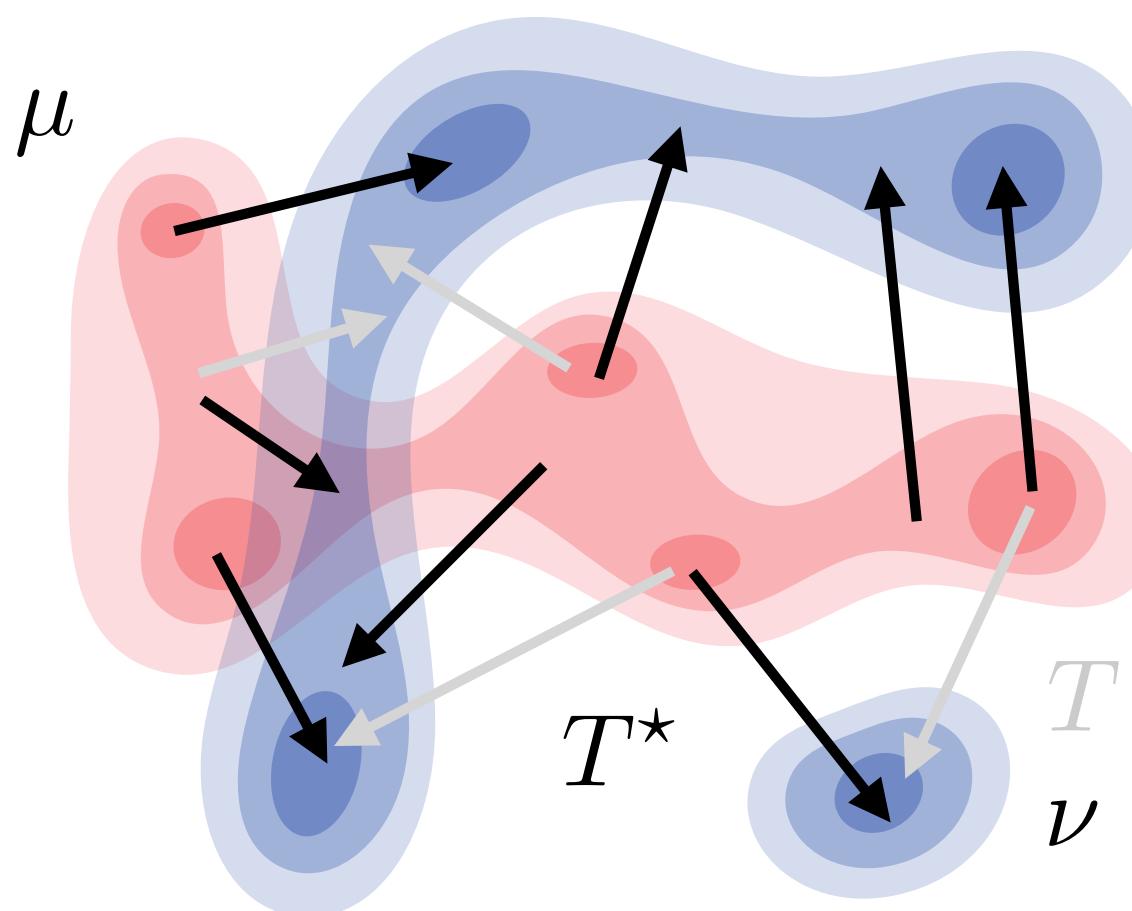
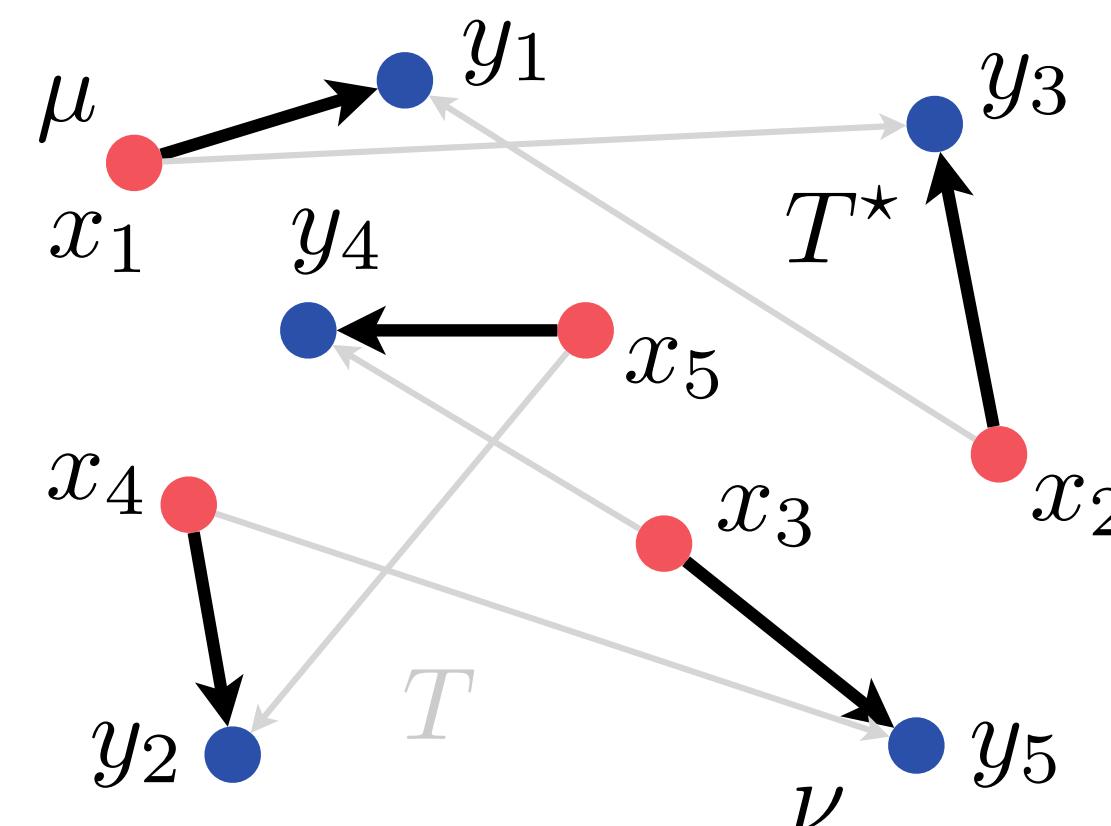
optimaltransport.github.io

Charlotte Bunne & Marco Cuturi

*Optimal Transport in
Learning, Control, and Dynamical Systems*

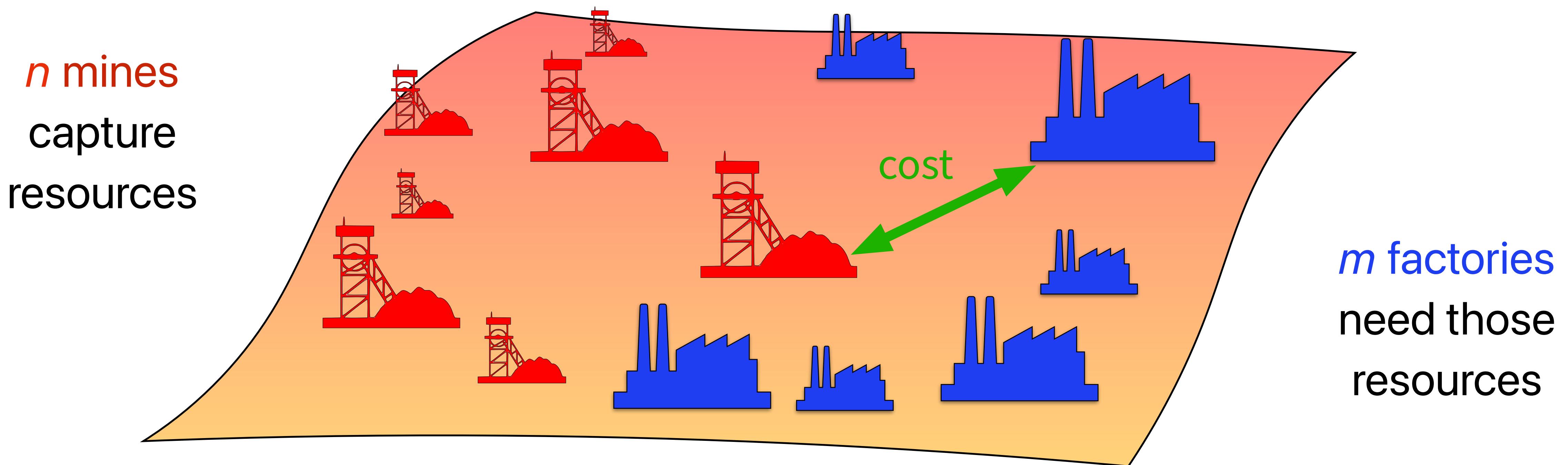
bunne.ch/ot_tutorial

Optimal Transport Theory



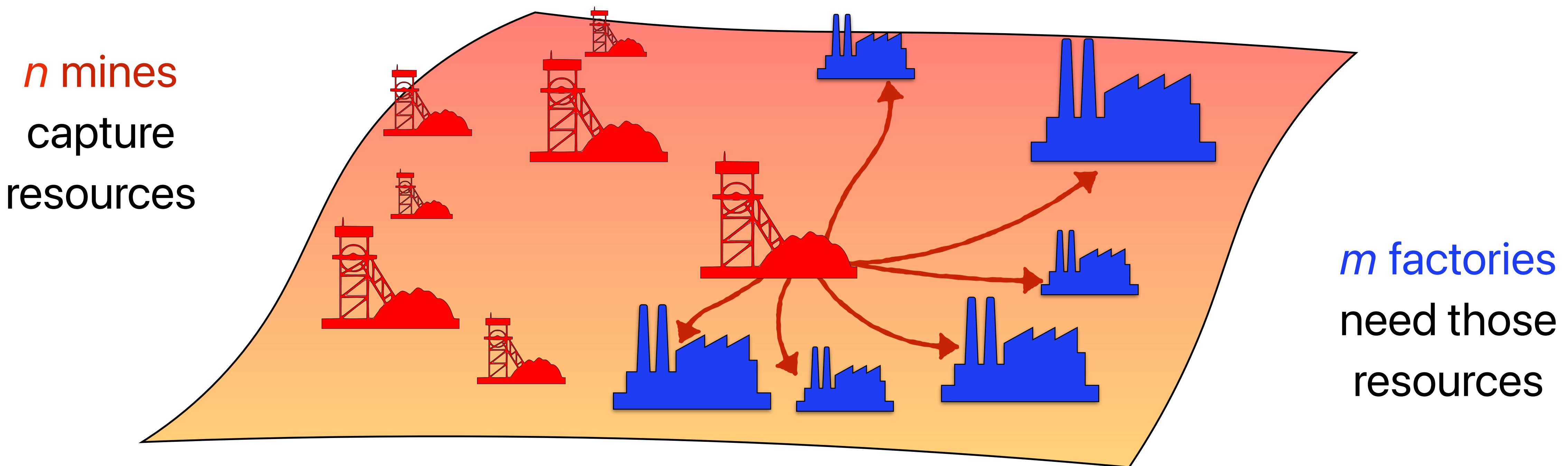
Motivation: Transporting Resources

Goal: Optimize *cost* of moving resources from *sources* to *targets*.

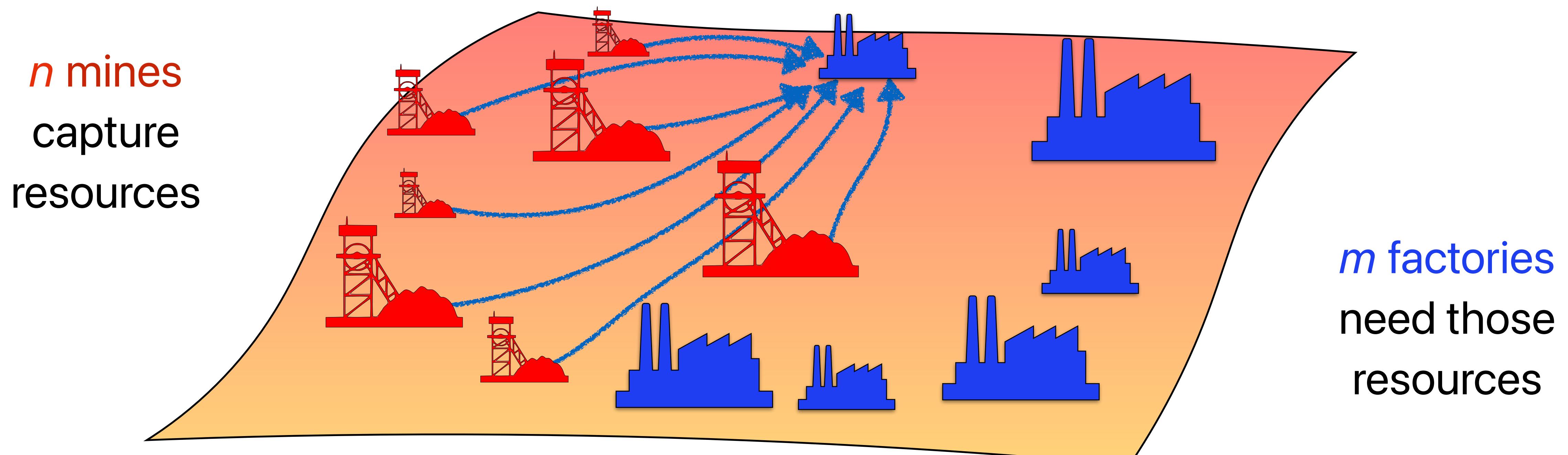


Motivation: Transporting Resources

Each of the n mines must dispatch its production



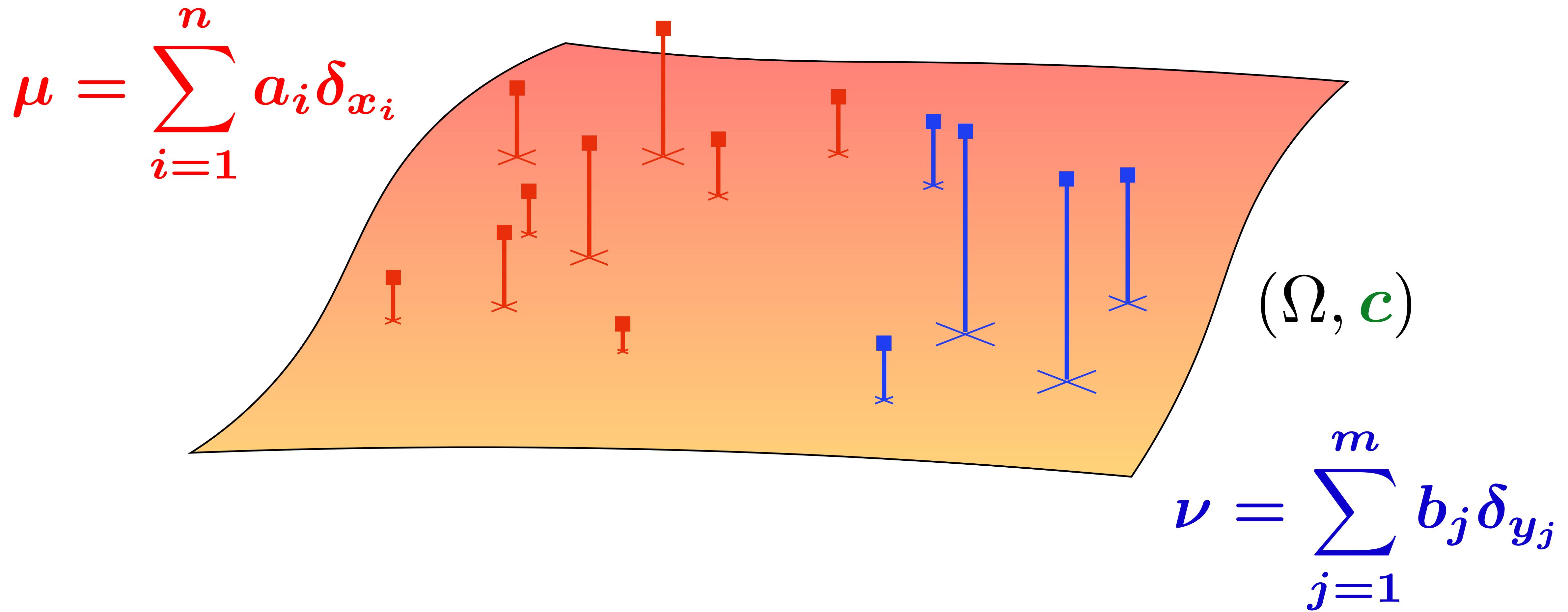
Motivation: Transporting Resources



Each of the m factories must get resource it needs

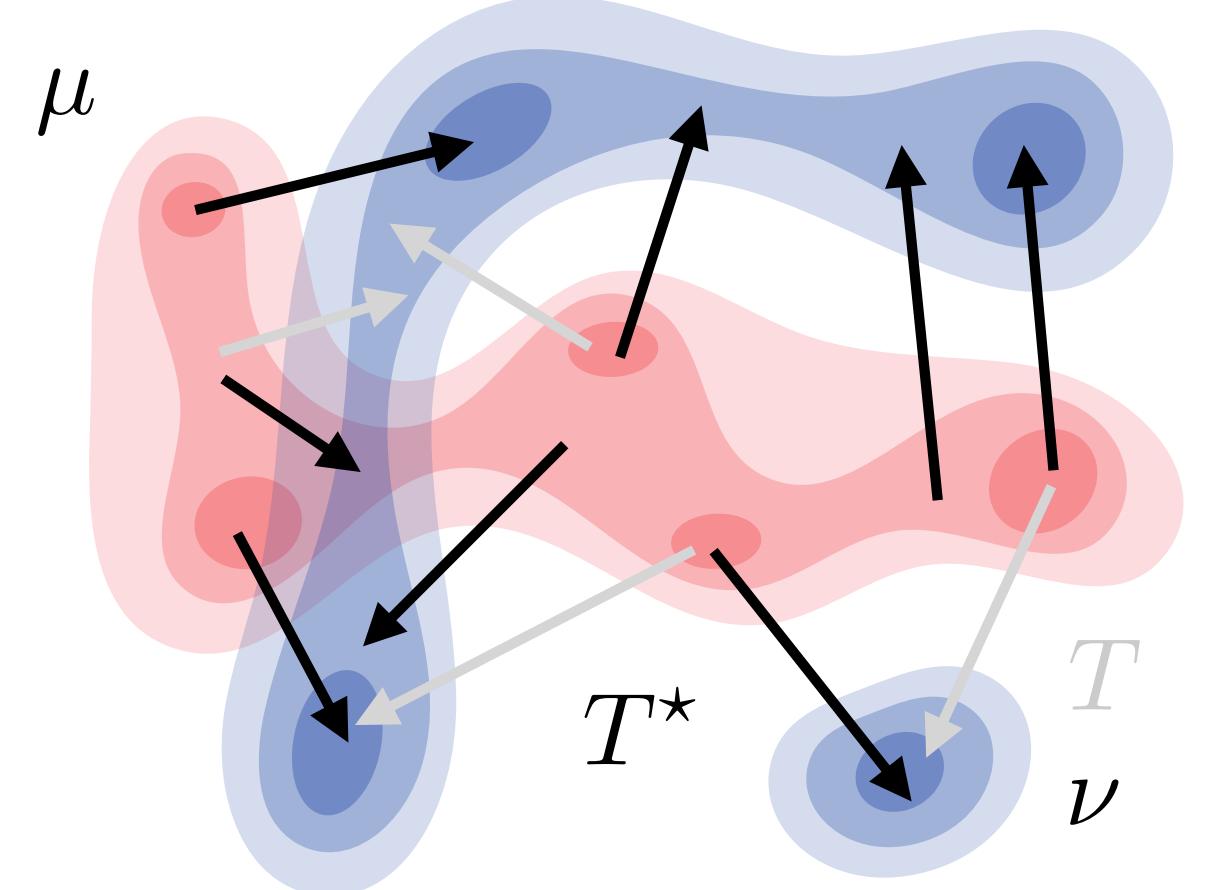
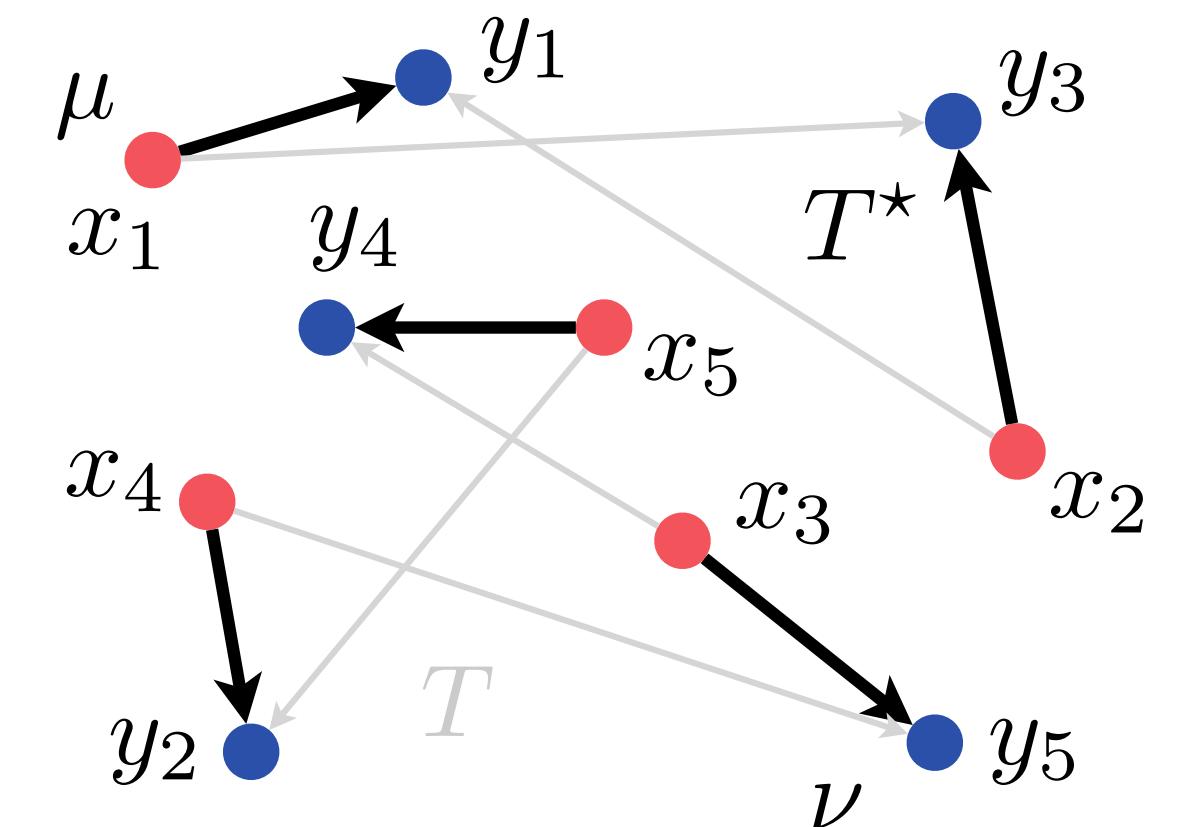
Motivation: Transporting Resources

Goal: Optimize *cost* of moving resources from *sources* to *targets*.



Monge Formulation: Definitions

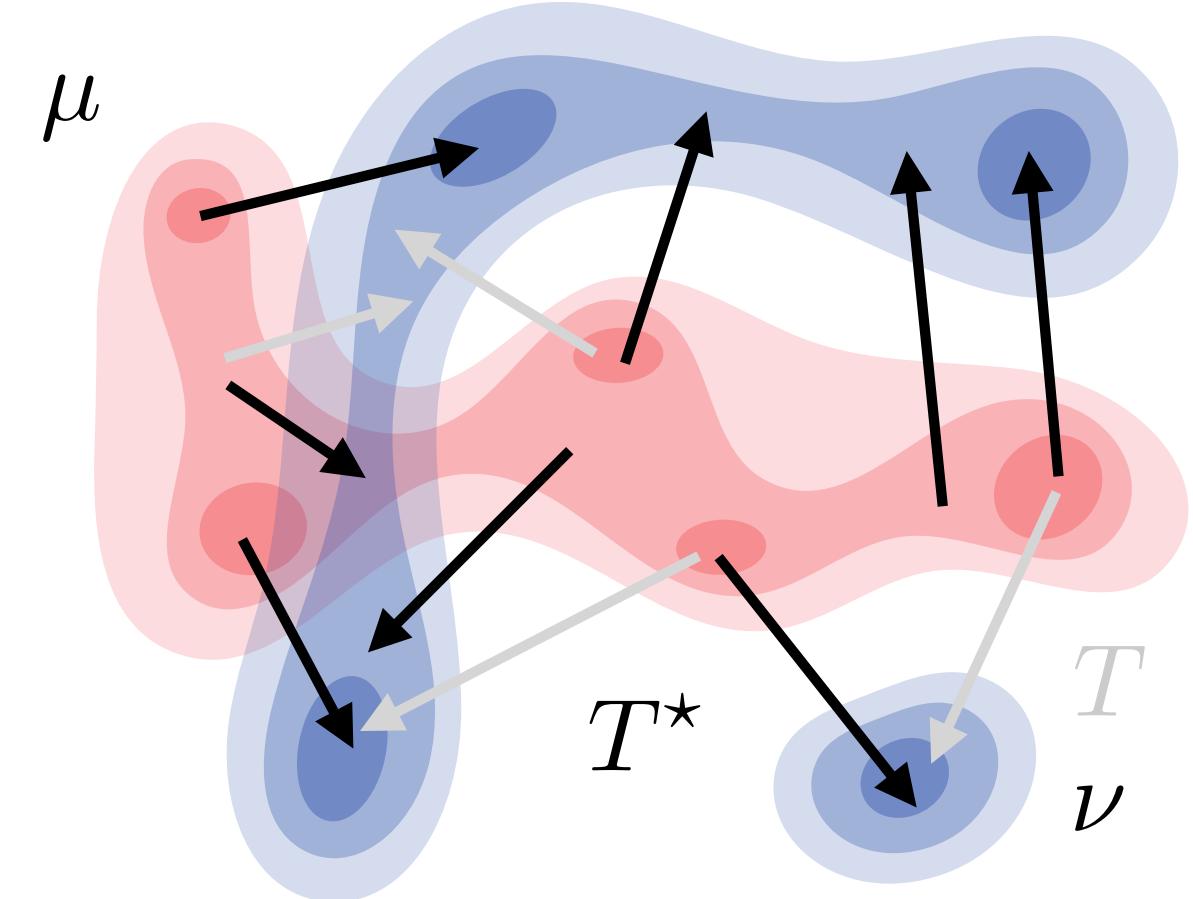
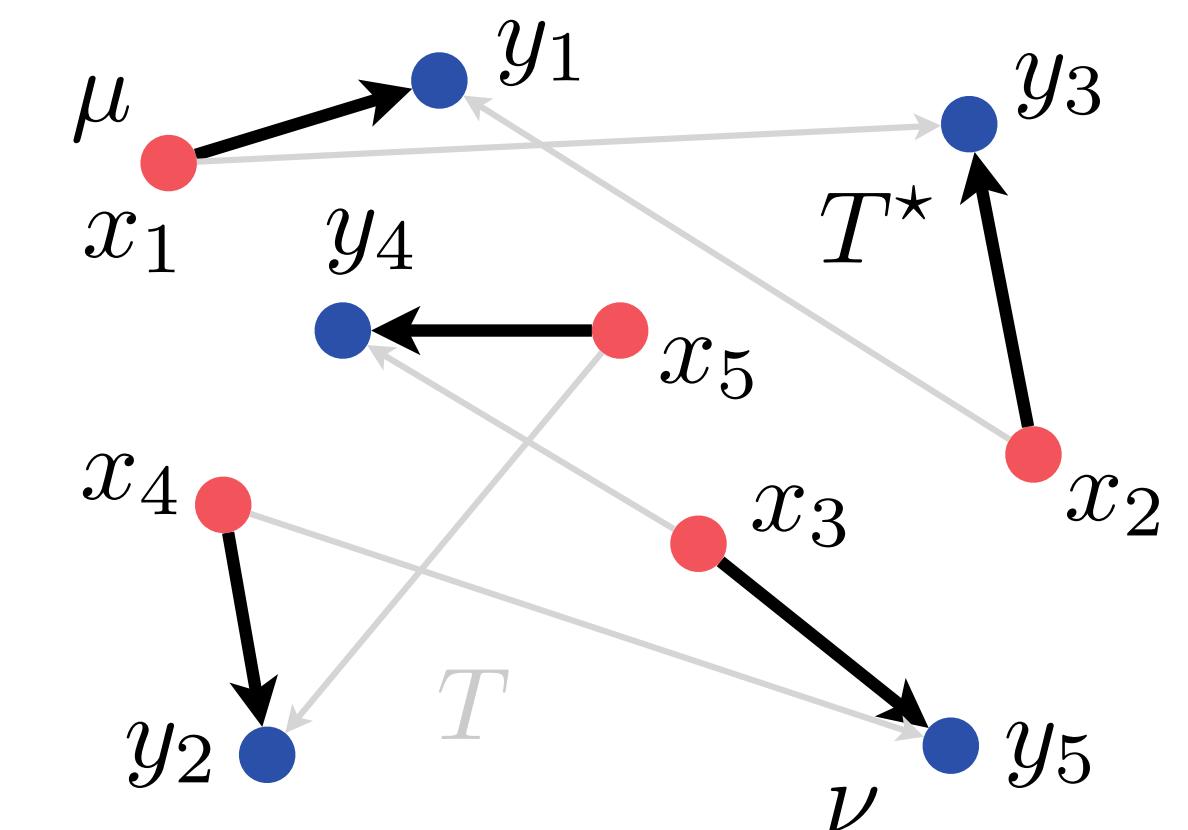
Probability measures μ, ν with support on spaces X, Y define *source* and *target* distributions.



Monge Formulation: Definitions

Probability measures μ, ν with support on spaces X, Y define *source* and *target* distributions.

Cost function $c : X \times Y \rightarrow \mathbb{R}$ defines cost $c(x, y)$ of transporting one unit of resource from $x \in X$ to $y \in Y$.

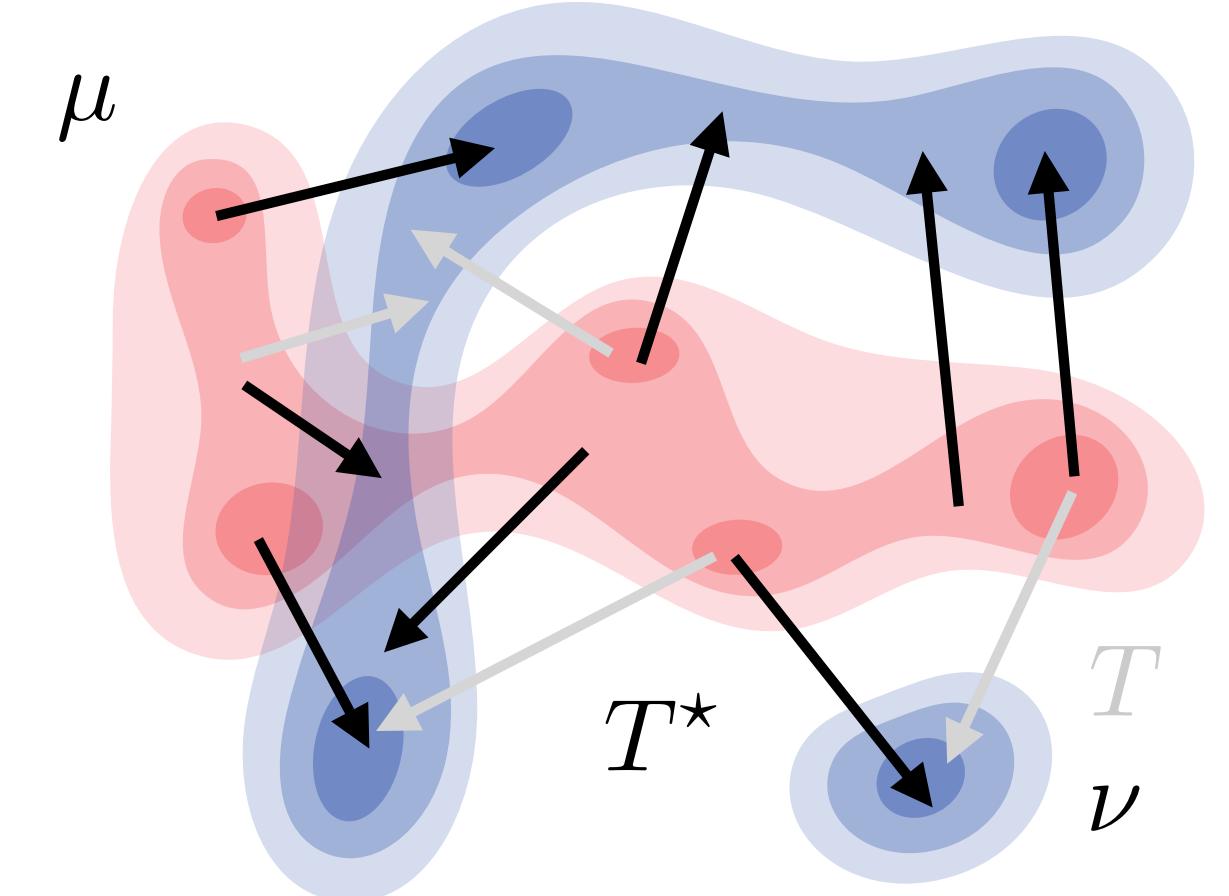
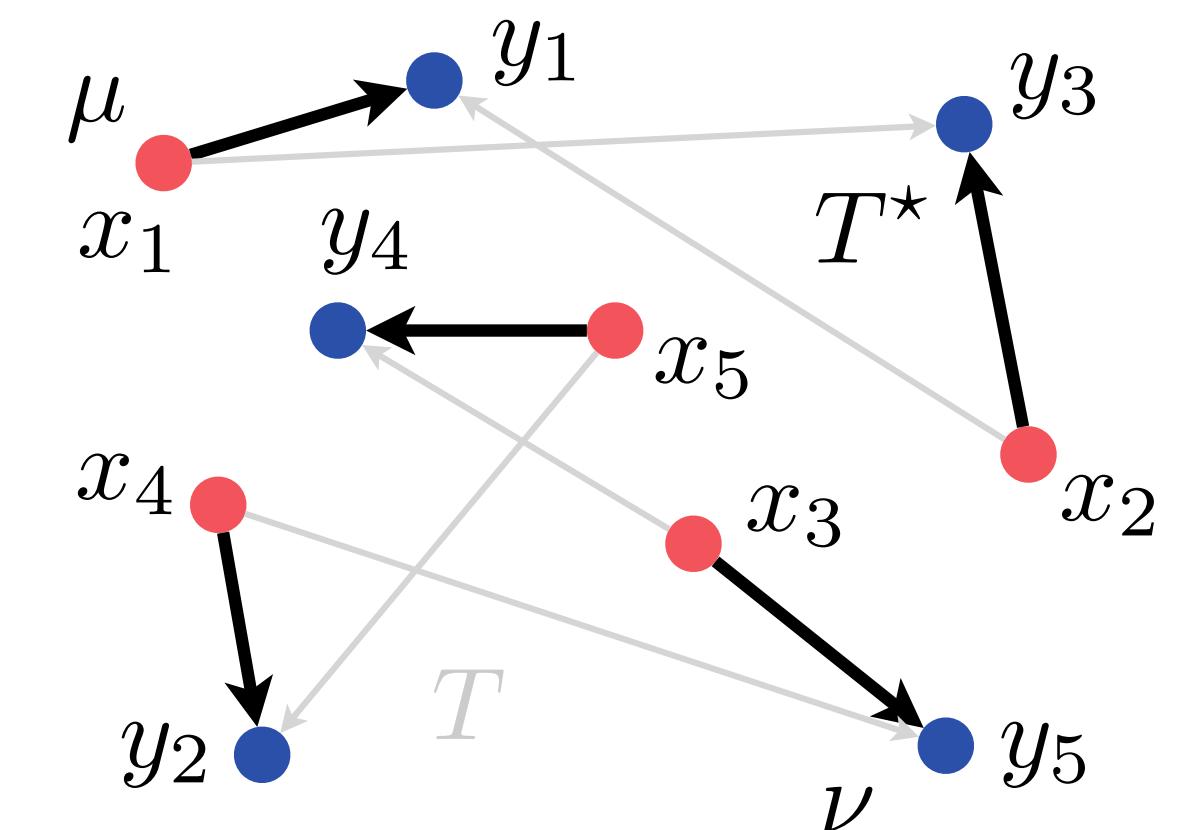


Monge Formulation: Definitions

Probability measures μ, ν with support on spaces X, Y define *source* and *target* distributions.

Cost function $c : X \times Y \rightarrow \mathbb{R}$ defines cost $c(x, y)$ of transporting one unit of resource from $x \in X$ to $y \in Y$.

Most common cost function $c(x, y) = d(x, y)^p$ is written in terms of distance metric d , known as the *ground metric*.



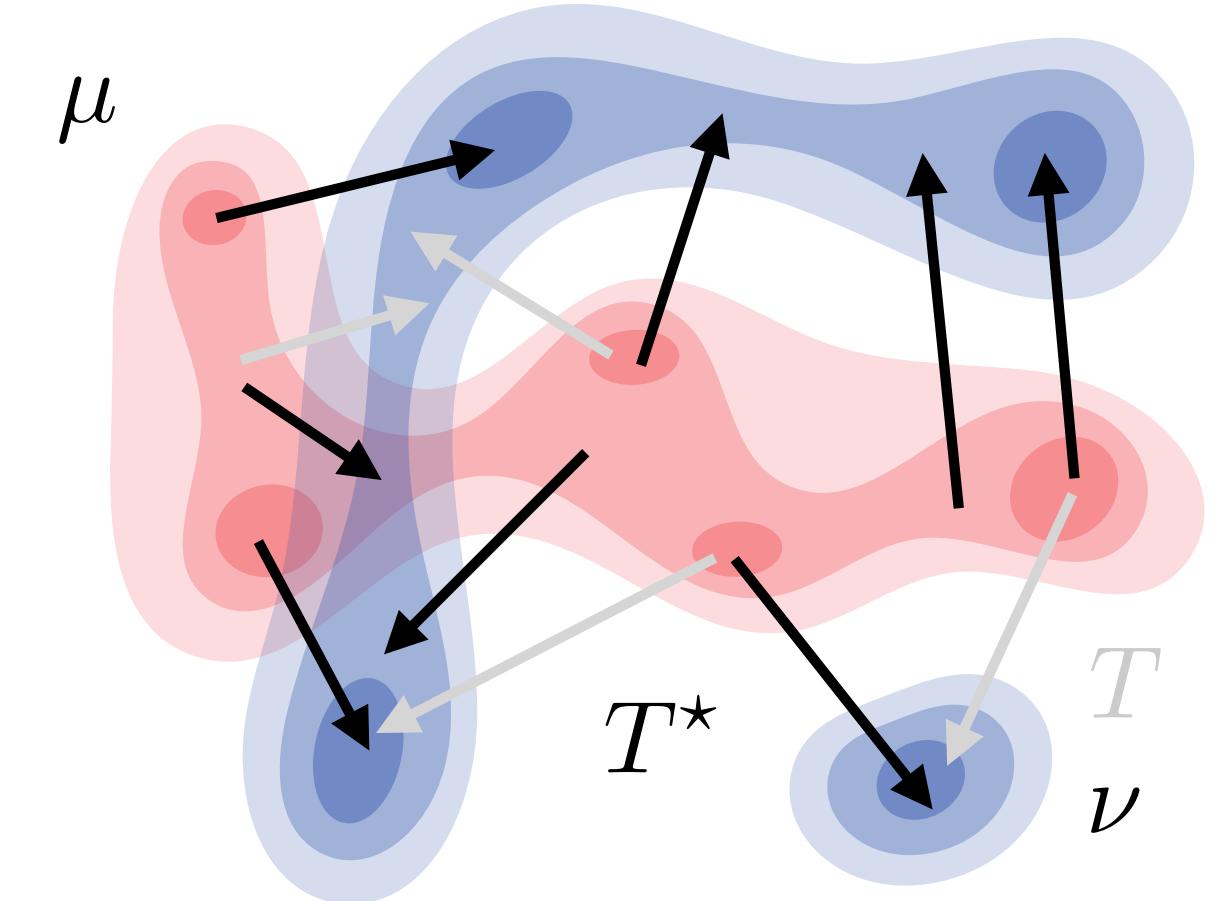
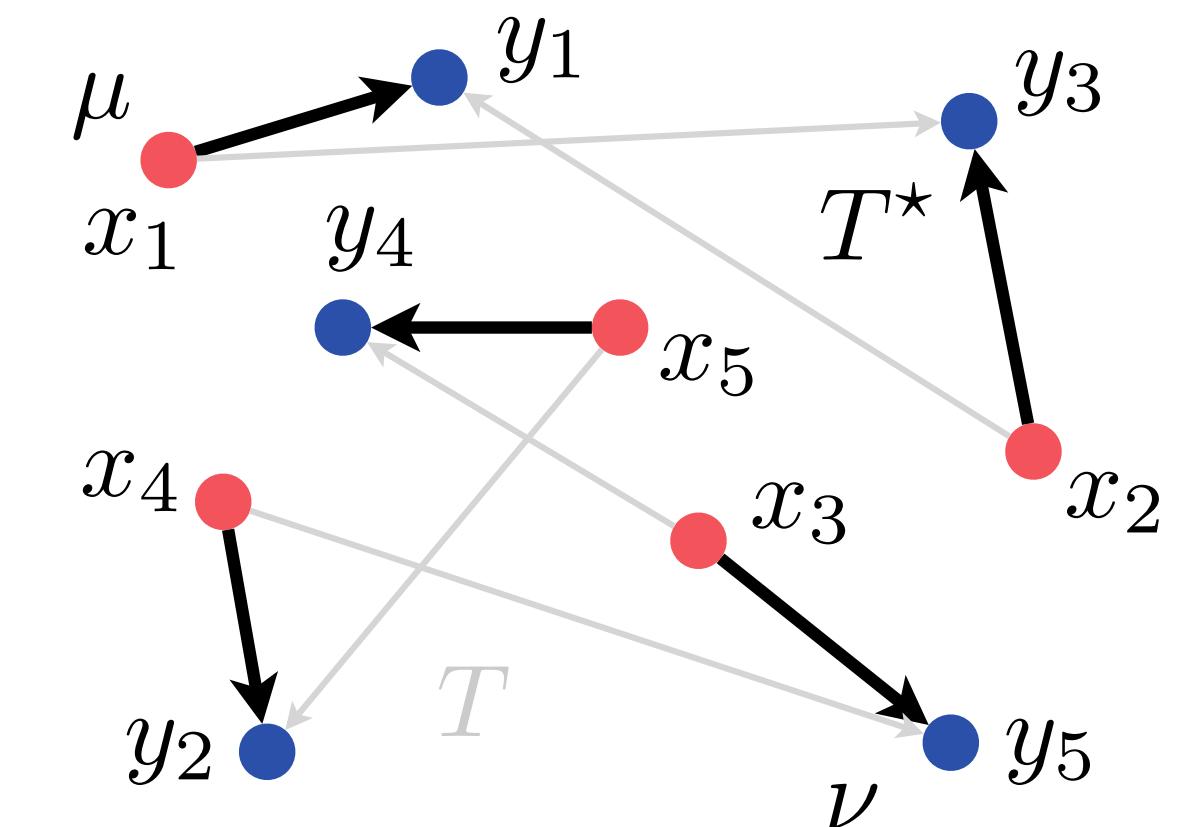
Monge Formulation: Definitions

Probability measures μ, ν with support on spaces X, Y define *source* and *target* distributions.

Cost function $c : X \times Y \rightarrow \mathbb{R}$ defines cost $c(x, y)$ of transporting one unit of resource from $x \in X$ to $y \in Y$.

Most common cost function $c(x, y) = d(x, y)^p$ is written in terms of distance metric d , known as the *ground metric*.

Transport map $T : X \rightarrow Y$ defines plan for transporting resources such that all resources are moved from the *source* to the *target* distribution, that is $T_\# \mu = \nu$.

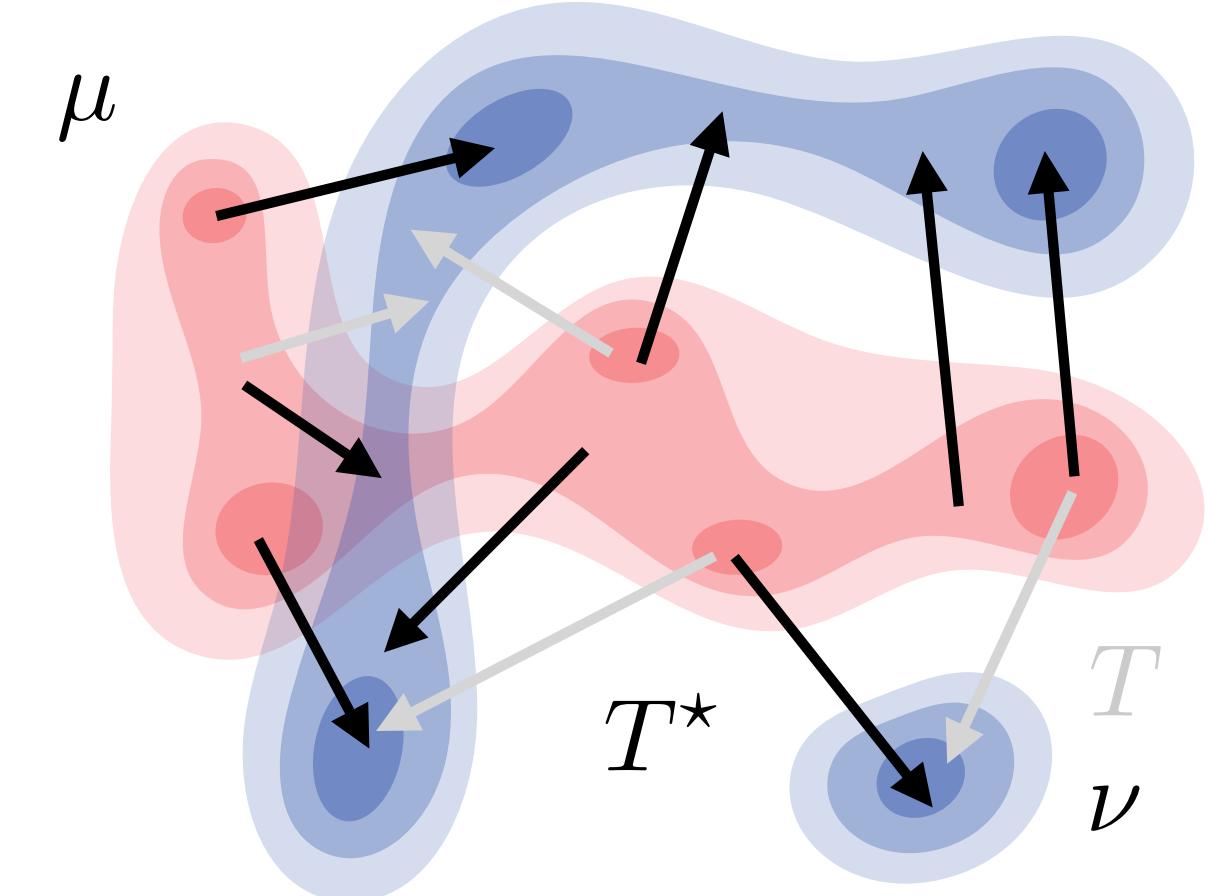
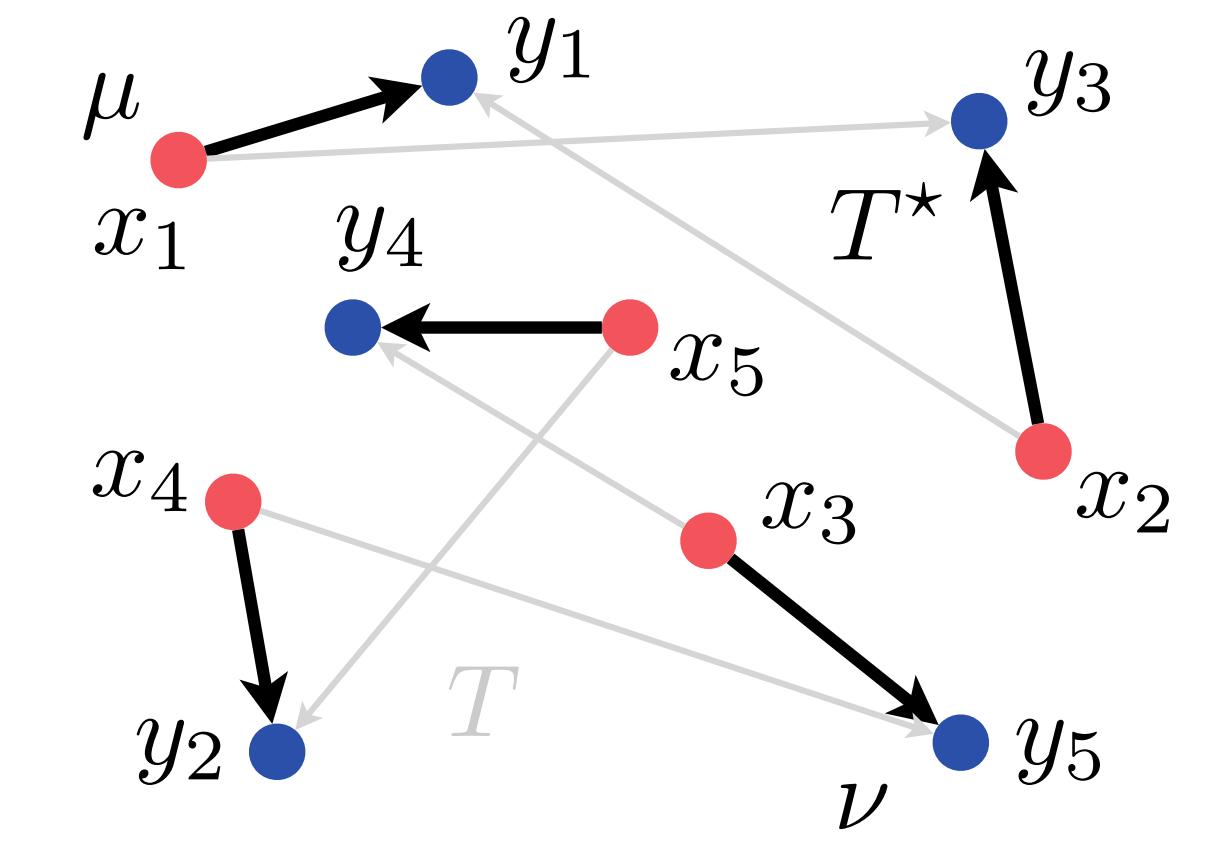


Monge Formulation: Optimal Transport Problem

Optimal transport map T^* has *minimal* total cost

$$\int_X c(x, T(x)) d\mu(x),$$

subject to constraint $T_{\#}\mu = \nu$.



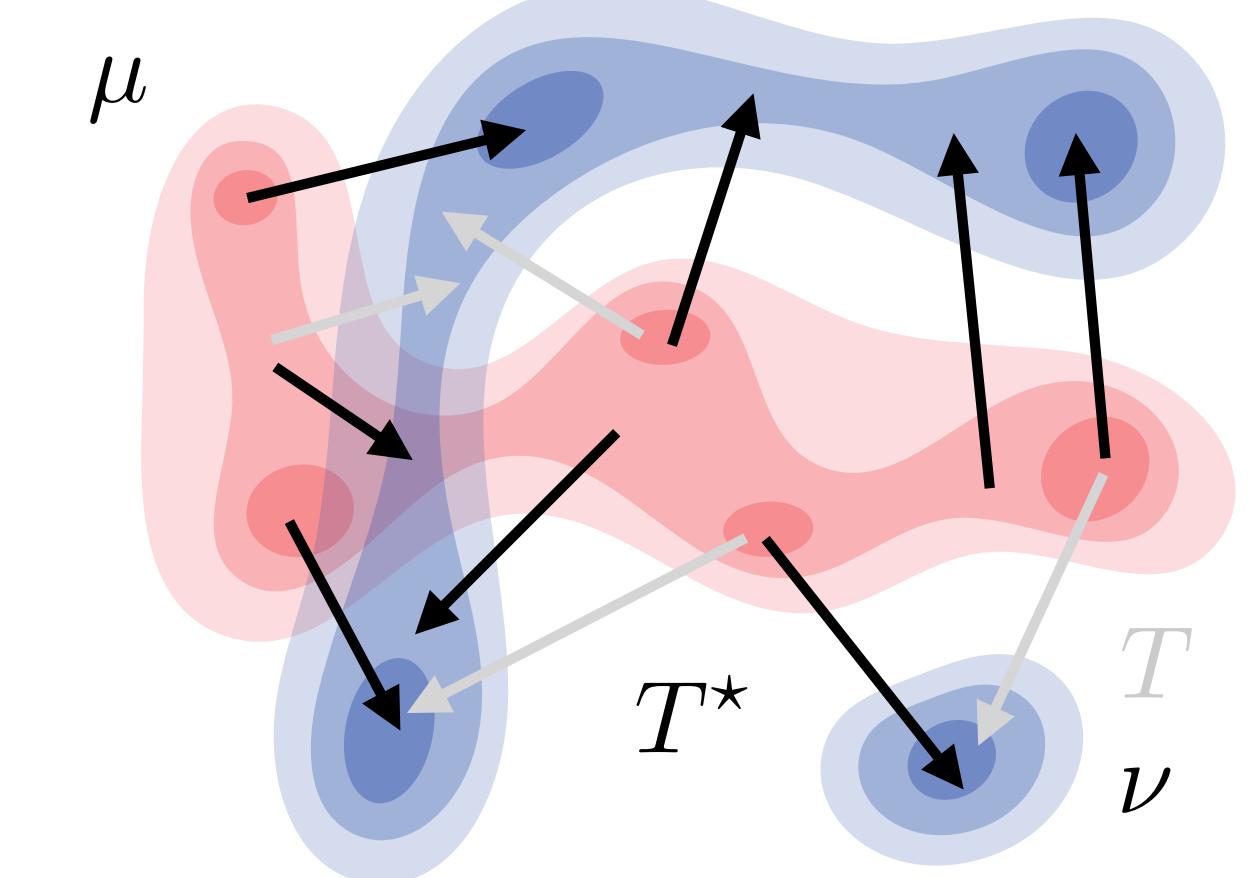
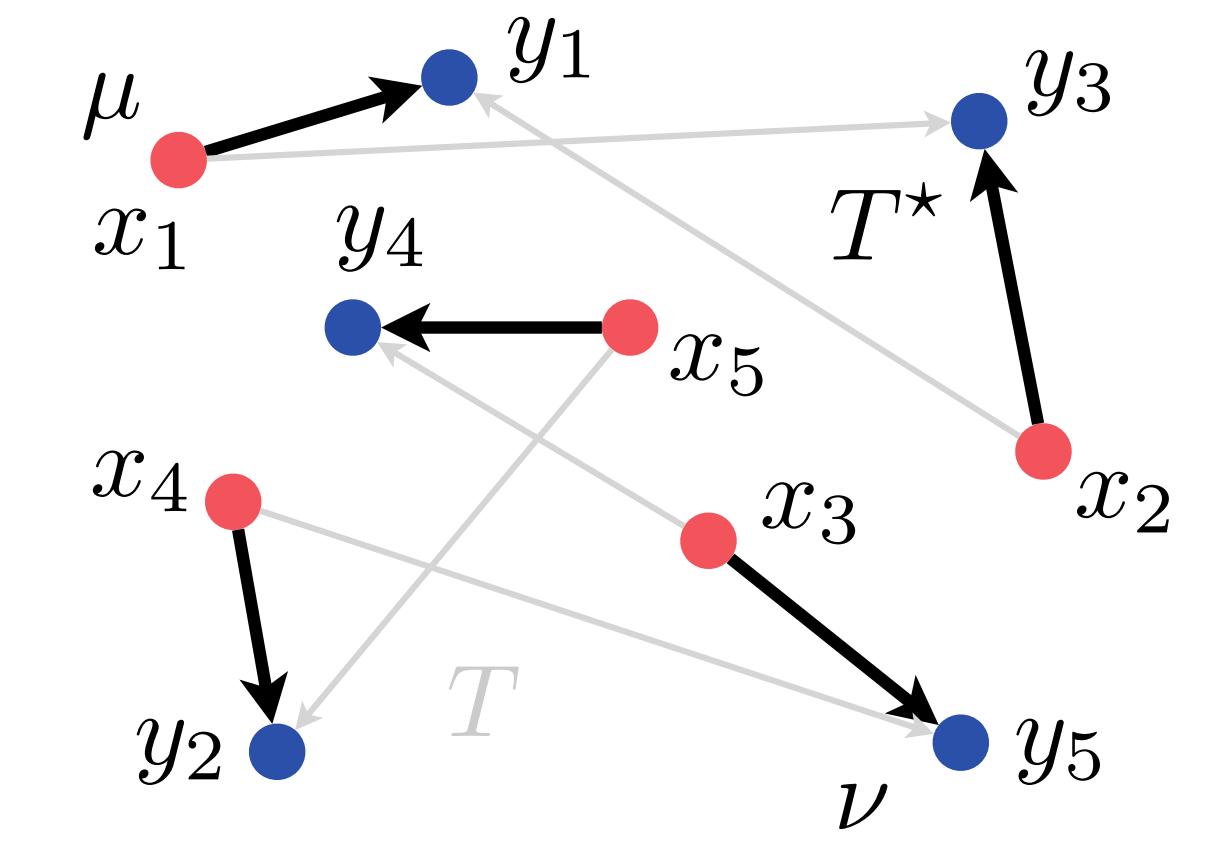
Monge Formulation: Optimal Transport Problem

Optimal transport map T^* has *minimal* total cost

$$\int_X c(x, T(x)) d\mu(x),$$

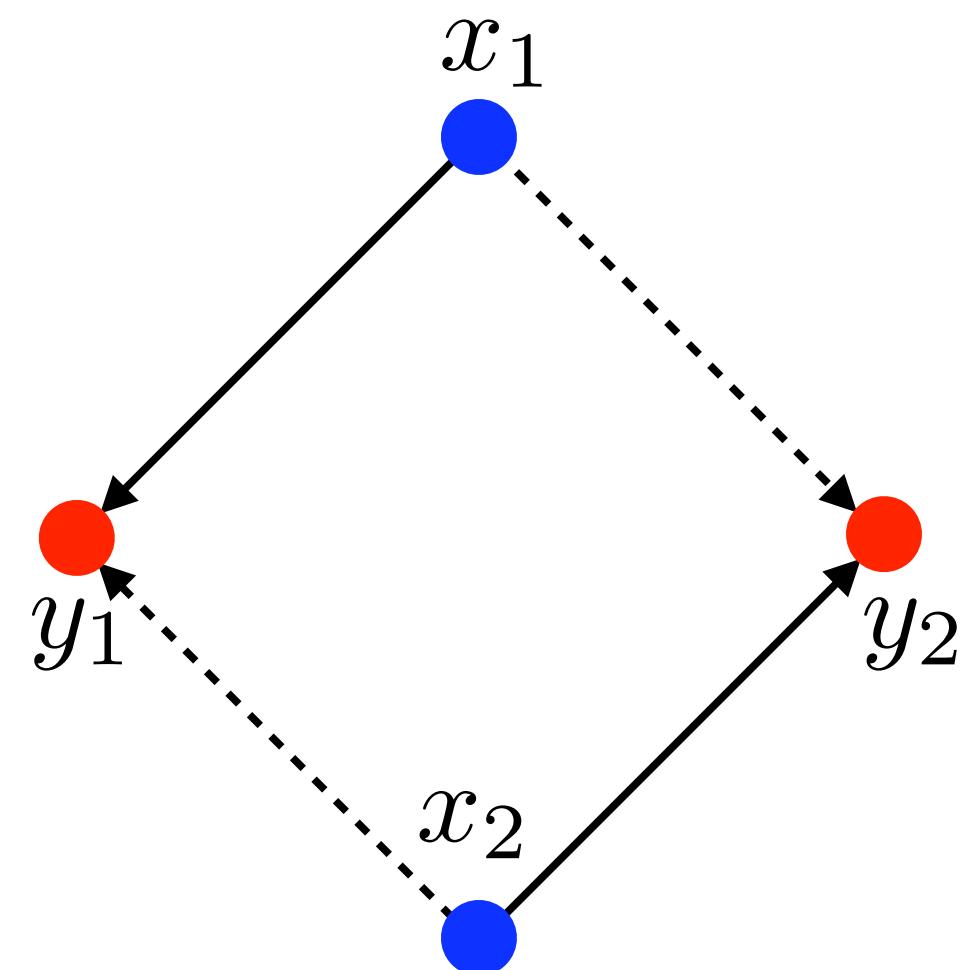
subject to constraint $T_{\#}\mu = \nu$.

$$\begin{aligned} \text{Optimal transport cost } L_c(\mu, \nu) &= \int_X c(x, T^*(x)) d\mu(x) \\ &= \min_T \left\{ \int_X c(x, T(x)) d\mu(x) : T_{\#}\mu = \nu \right\}. \end{aligned}$$



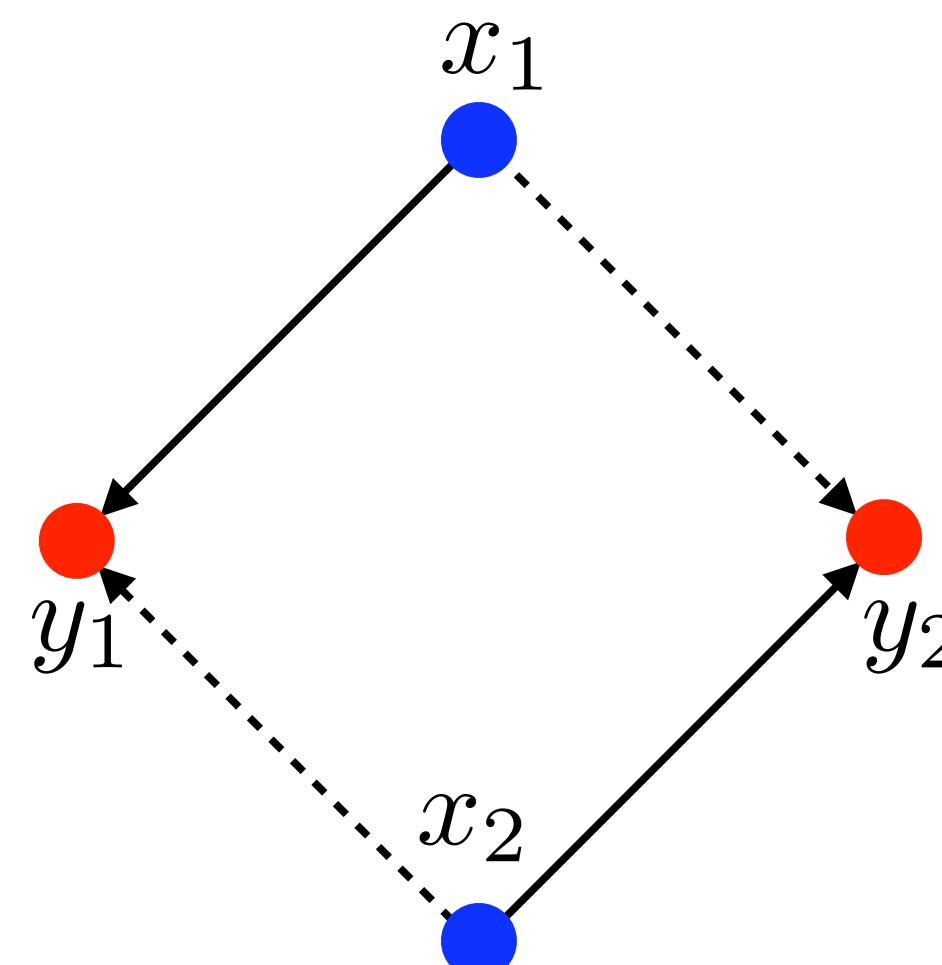
Monge Formulation: Complications

Non-uniqueness

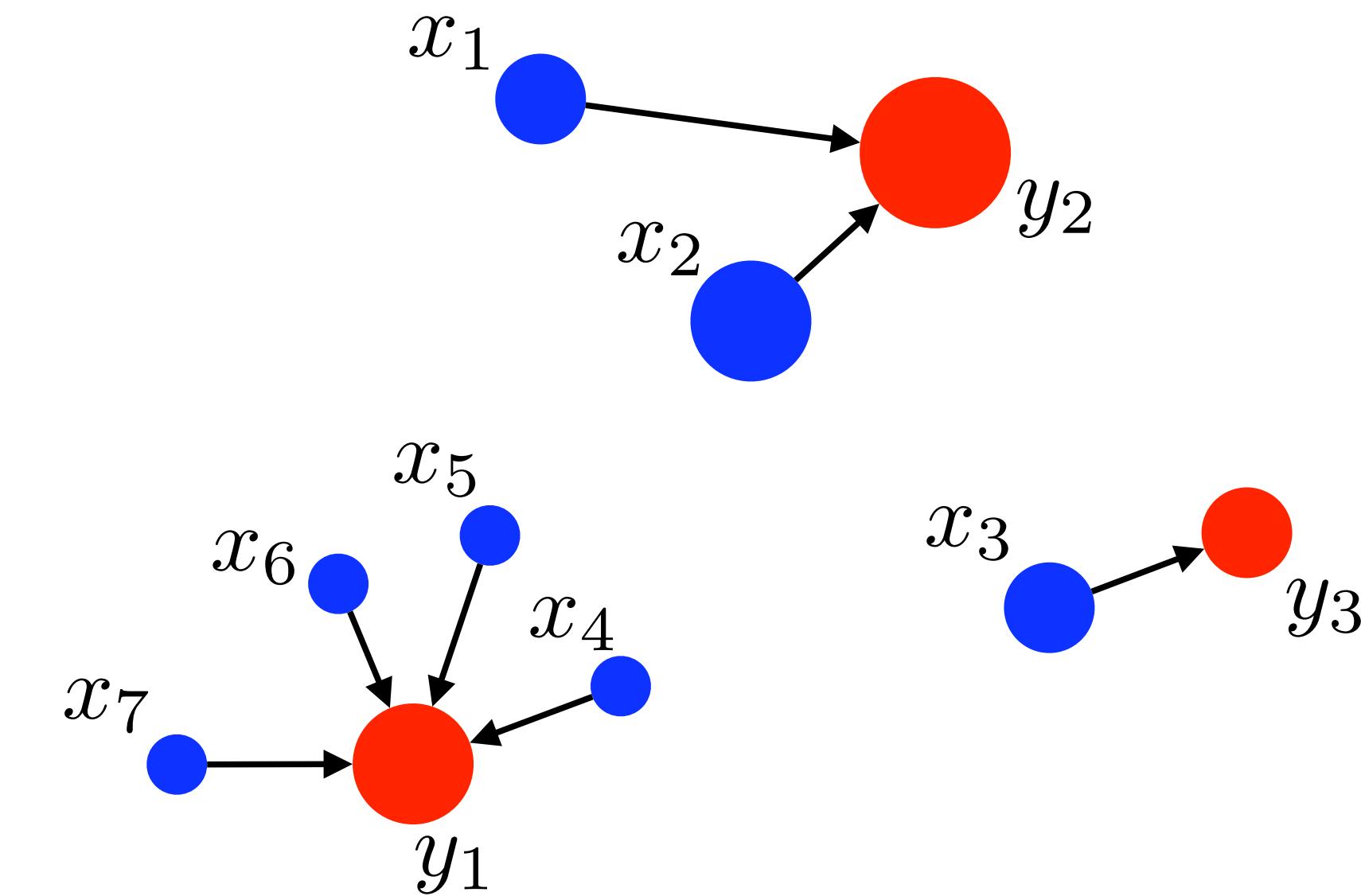


Monge Formulation: Complications

Non-uniqueness

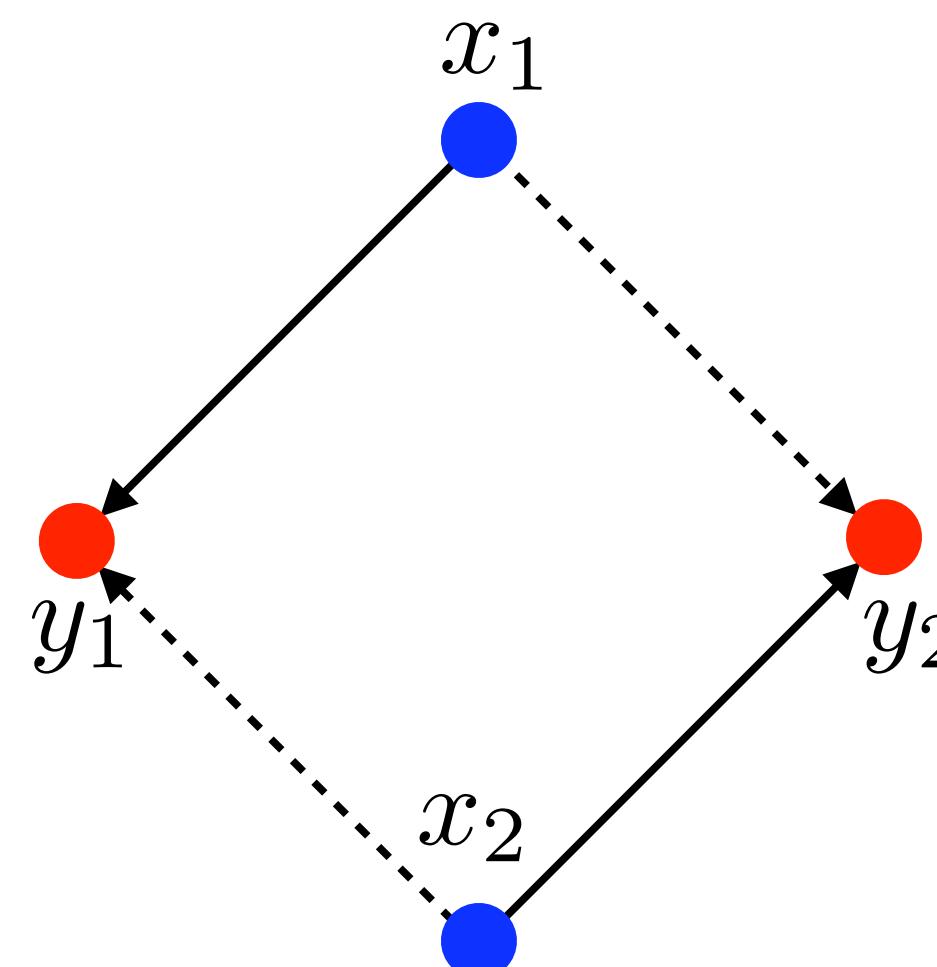


Ill-posed and asymmetric

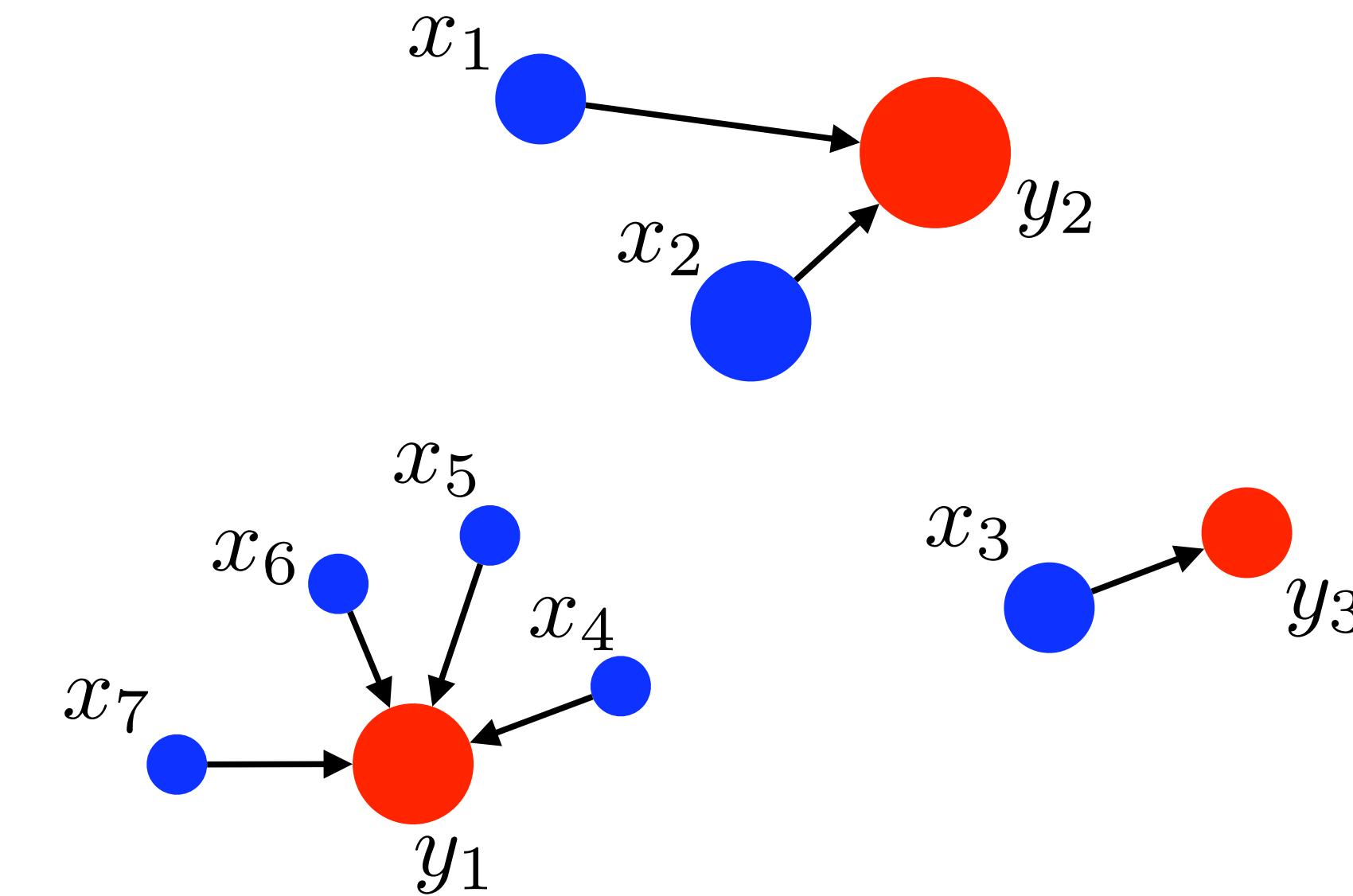


Monge Formulation: Complications

Non-uniqueness



Ill-posed and asymmetric

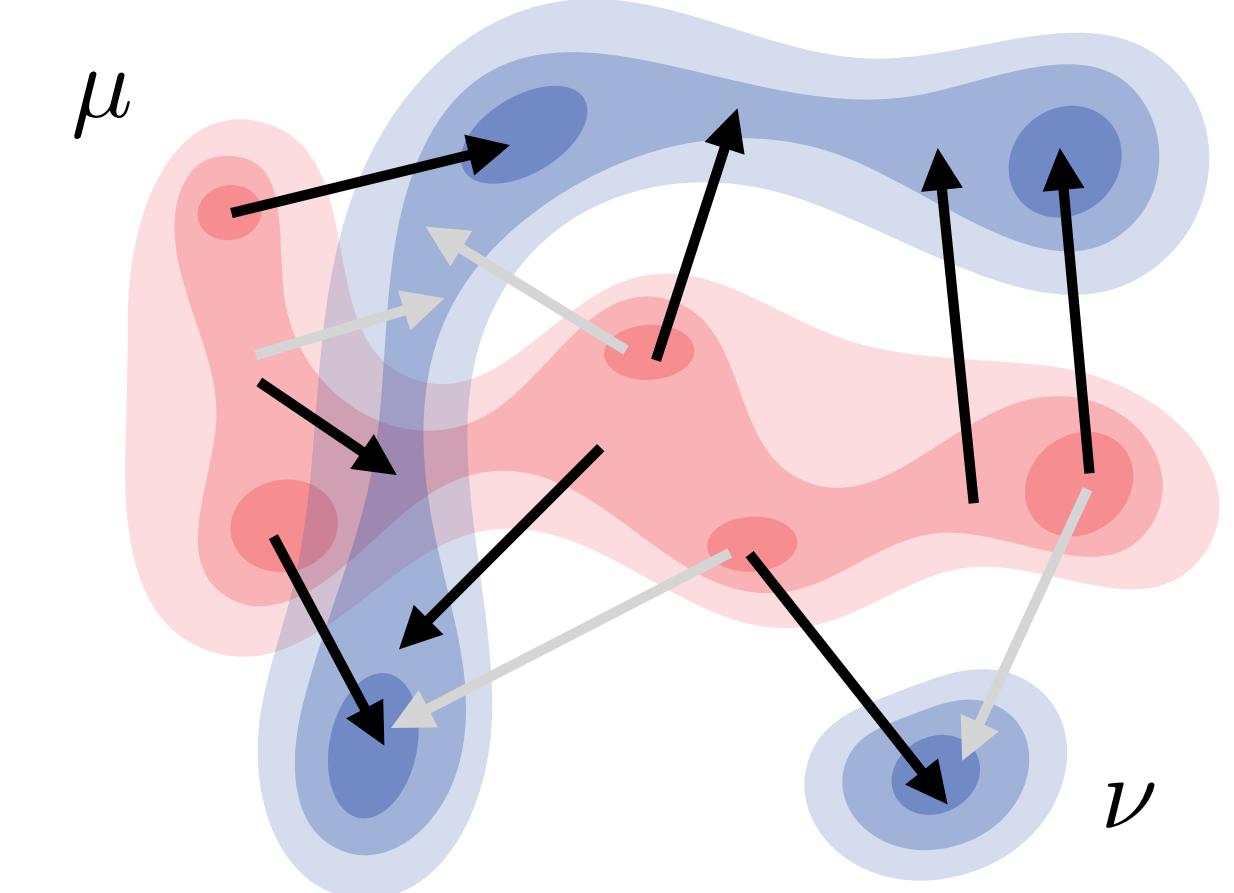
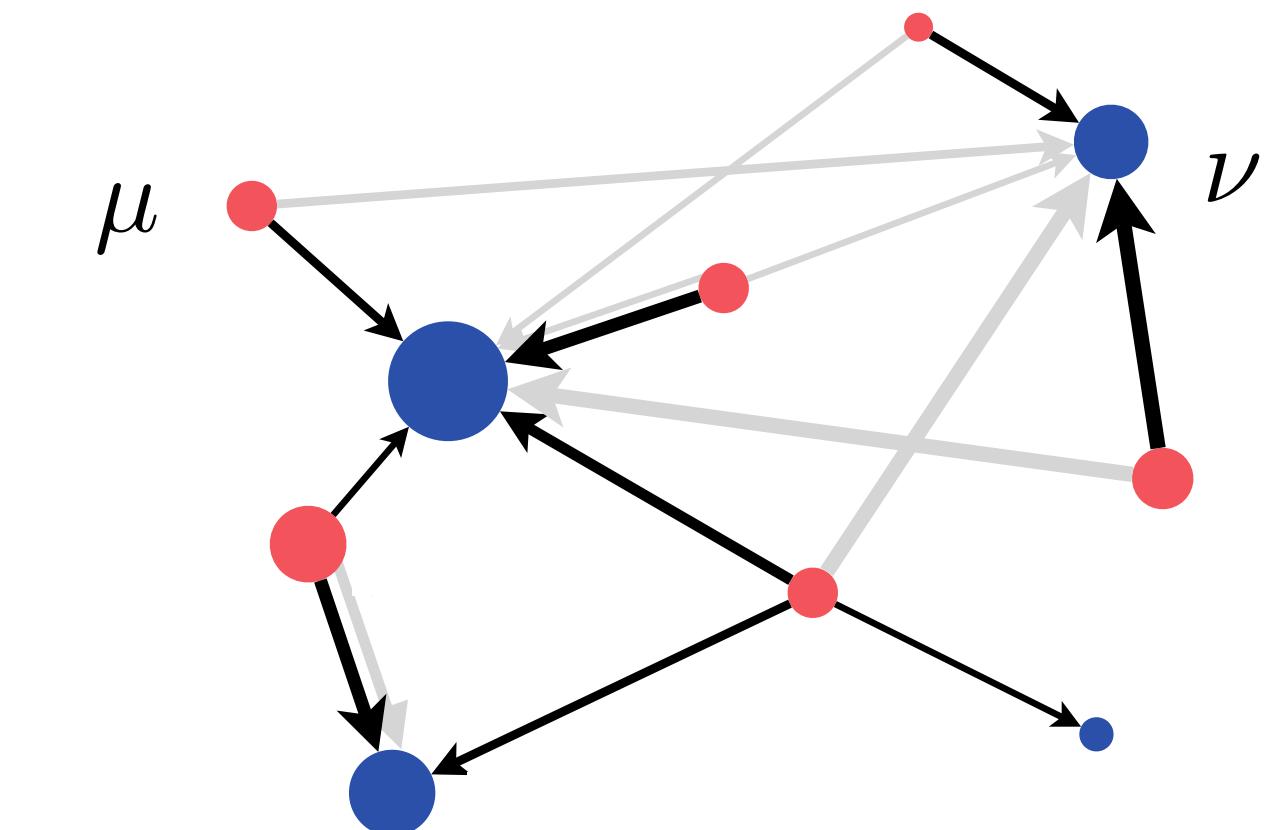


Highly nonlinear constraint $T_{\sharp}\mu = \nu$

Kantorovich Formulation: Convex Relaxation

Transport plan $\pi \in \Pi(\mu, \nu)$ defines probability measure on $X \times Y$ such that marginals of π match μ and ν , that is

$$\Pi(\mu, \nu) = \left\{ \pi : \int_X d\pi(x, y) = d\nu(y), \int_Y d\pi(x, y) = d\mu(x) \right\}.$$



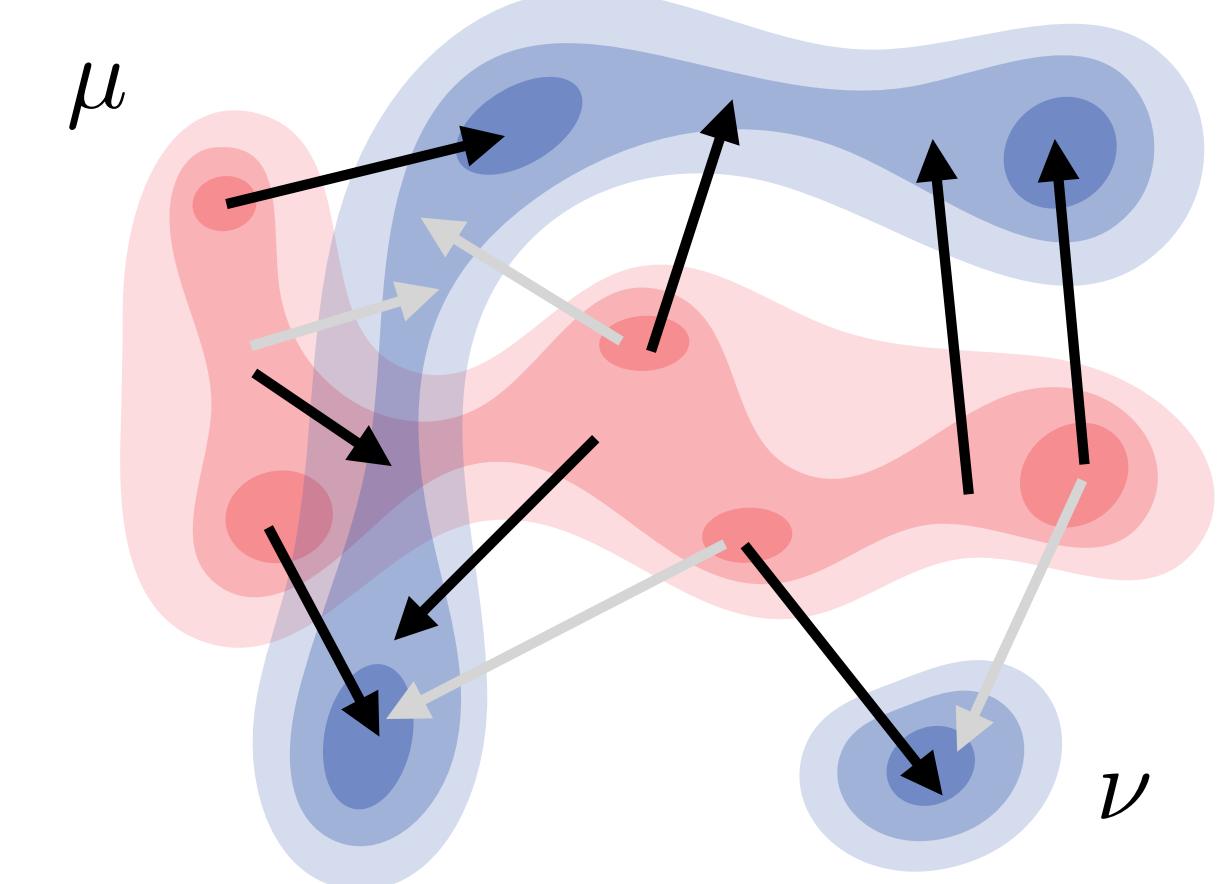
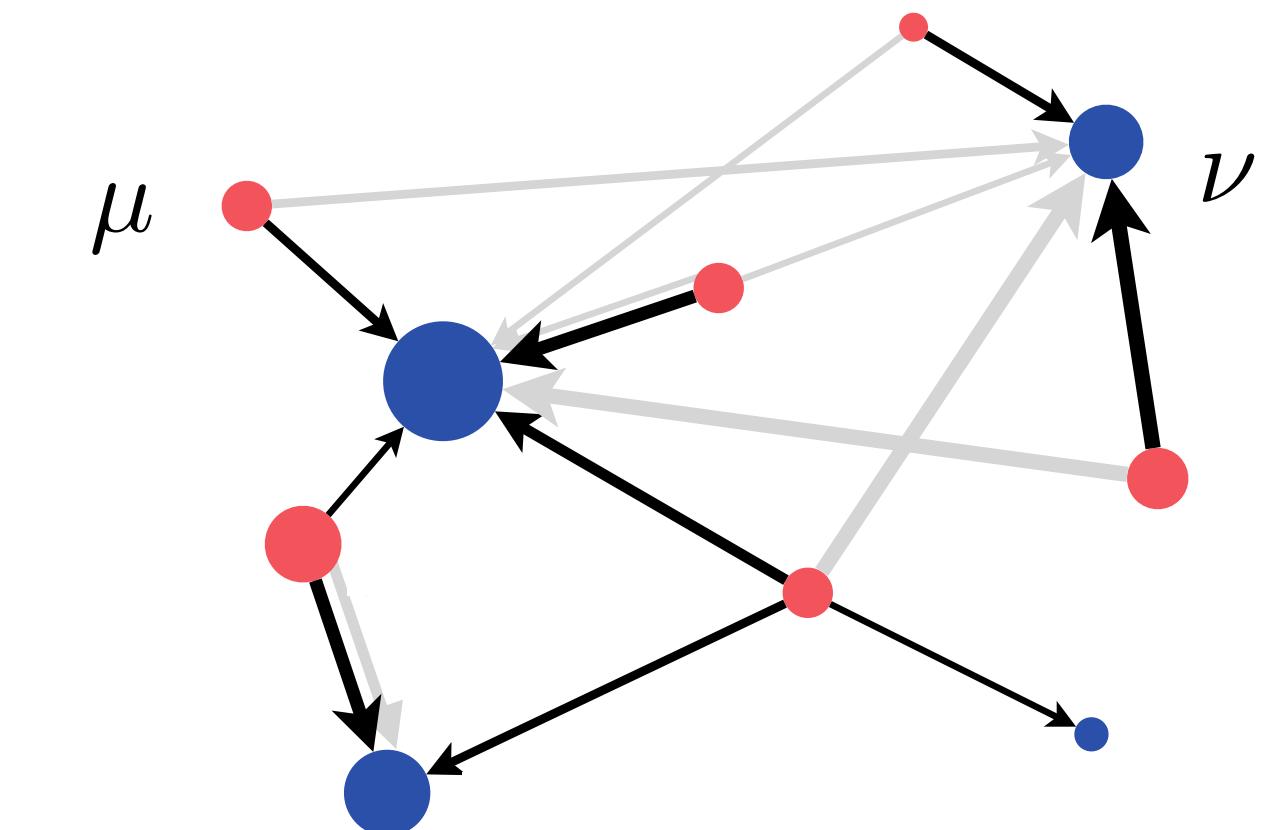
Kantorovich Formulation: Convex Relaxation

Transport plan $\pi \in \Pi(\mu, \nu)$ defines probability measure on $X \times Y$ such that marginals of π match μ and ν , that is

$$\Pi(\mu, \nu) = \left\{ \pi : \int_X d\pi(x, y) = d\nu(y), \int_Y d\pi(x, y) = d\mu(x) \right\}.$$

Optimal transport cost is given by

$$L_c(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) : \pi \in \Pi(\mu, \nu) \right\}.$$



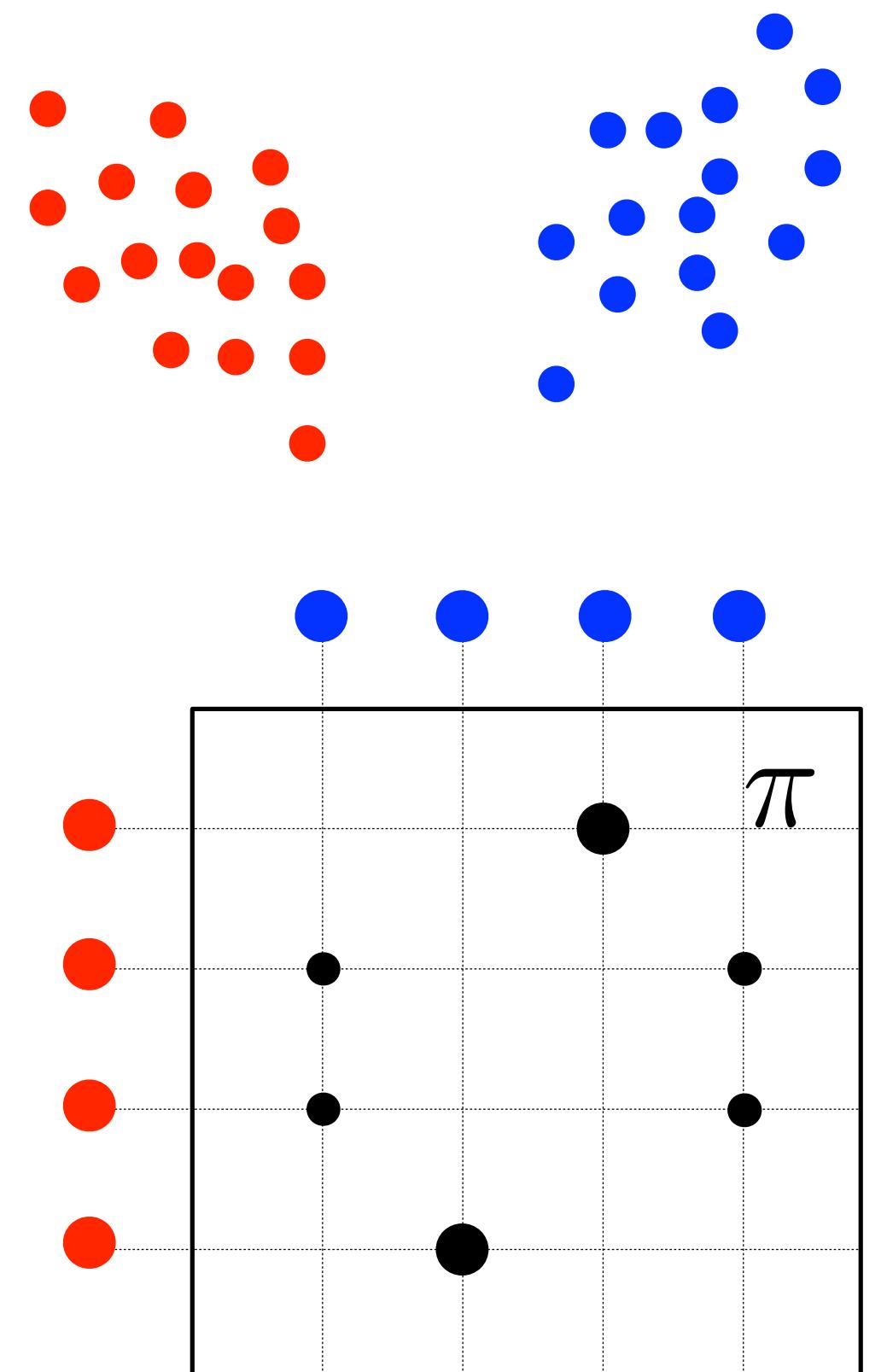
Kantorovich Formulation: Discrete Case

Transport plan matrix $P \in U(\mathbf{a}, \mathbf{b})$, where

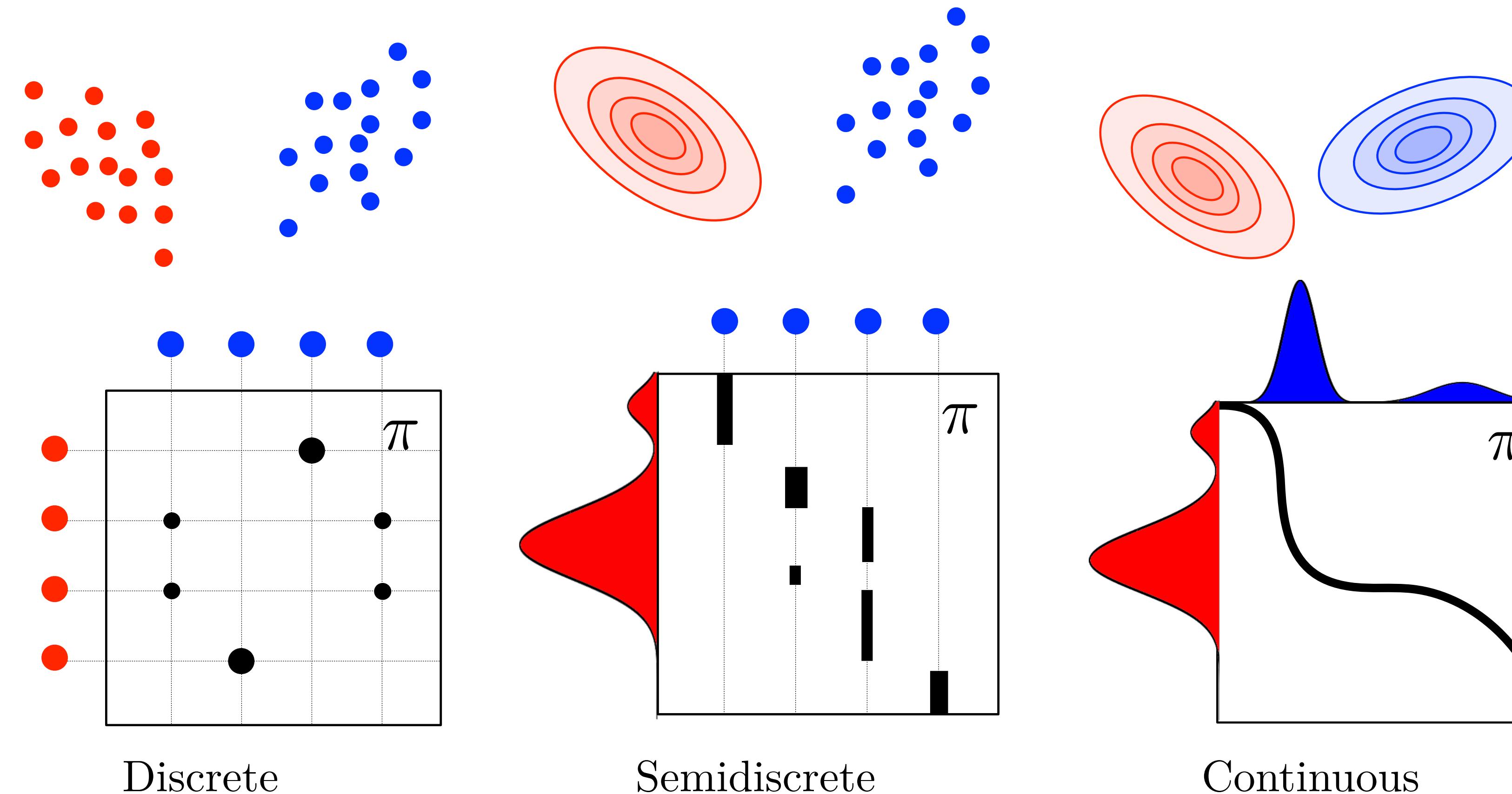
$$U(\mathbf{a}, \mathbf{b}) = \left\{ P : \sum_i P_{ij} = b_j, \sum_j P_{ij} = a_i, P_{ij} \geq 0 \right\}.$$

Optimal transport cost is given by

$$L_C(\mathbf{a}, \mathbf{b}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} : P \in U(\mathbf{a}, \mathbf{b}) \right\}.$$



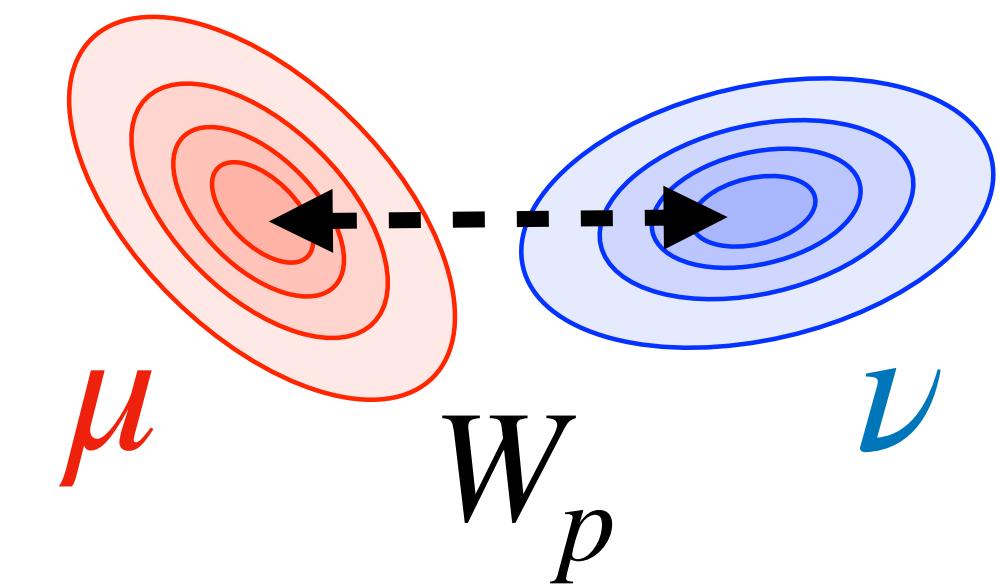
Kantorovich Formulation: Convex Relaxation



Wasserstein Distances: Metric on Space of Distributions

With $c(x, y) = d(x, y)^p$ for ground metric d ,
the p -Wasserstein distance is given by

$$W_p(\mu, \nu) = L_c(\mu, \nu)^{1/p}.$$

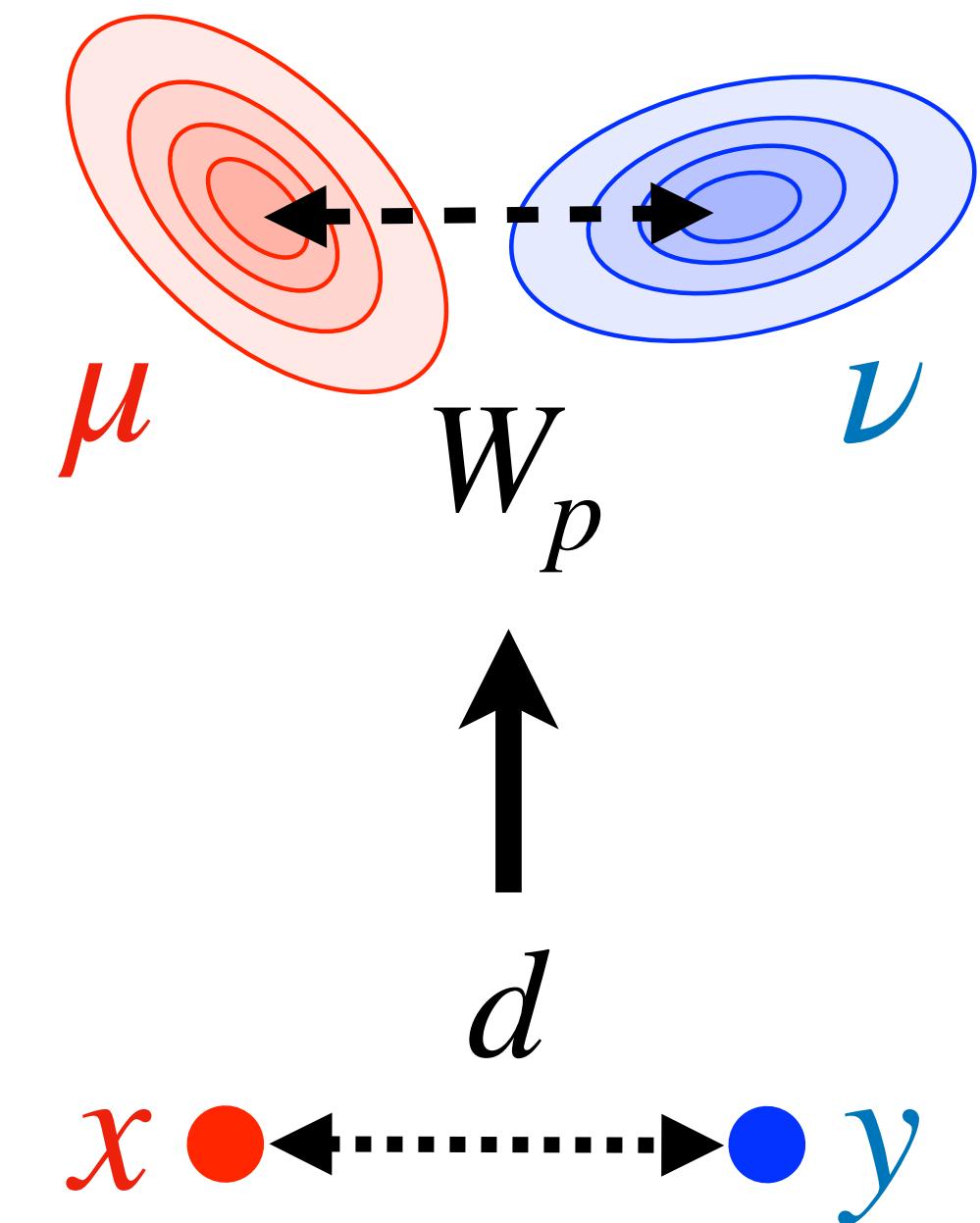


Wasserstein Distances: Metric on Space of Distributions

With $c(x, y) = d(x, y)^p$ for ground metric d ,
the p -Wasserstein distance is given by

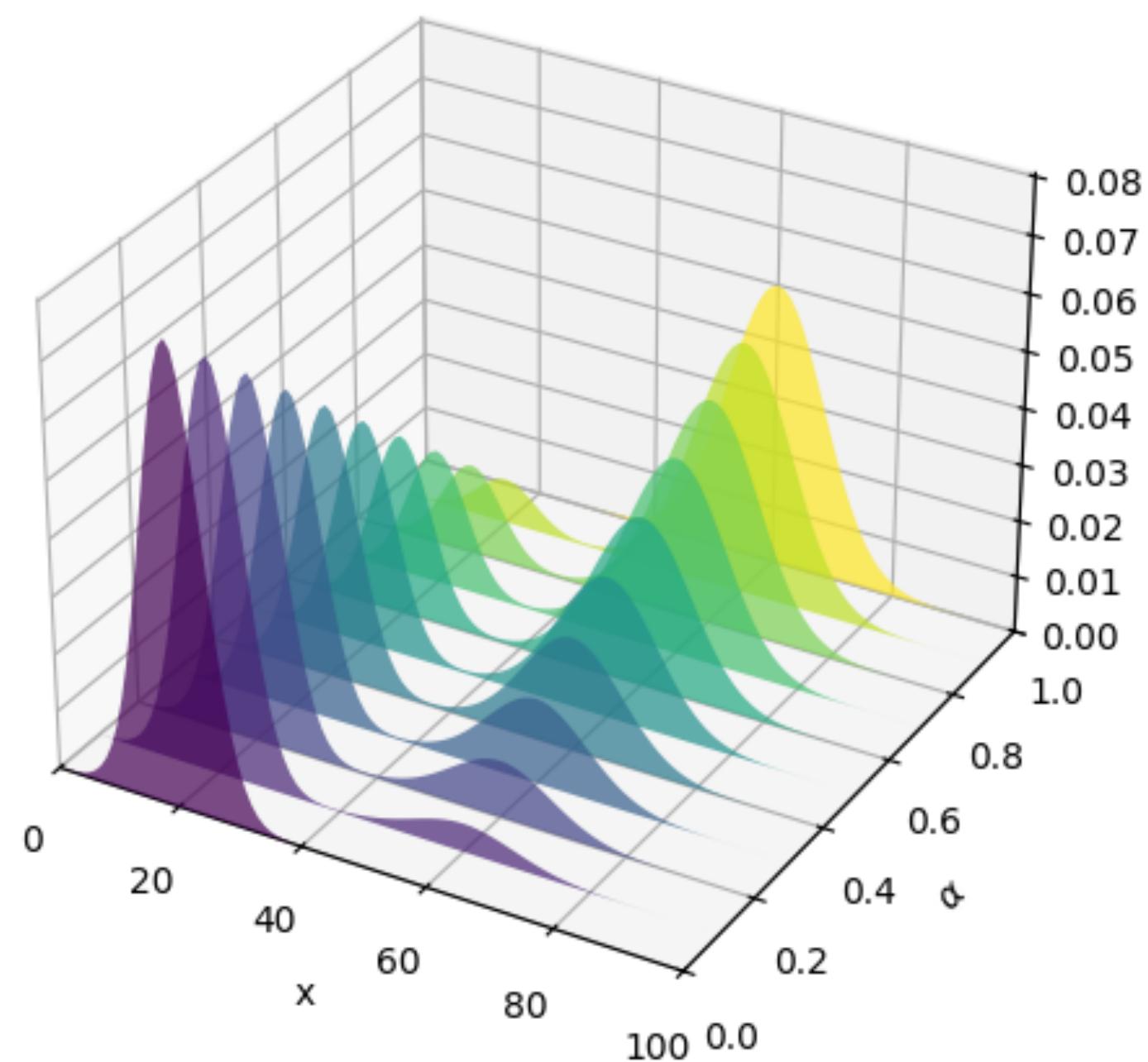
$$W_p(\mu, \nu) = L_c(\mu, \nu)^{1/p}.$$

Optimal transport lifts metric in ambient space to *metric in distribution space!*

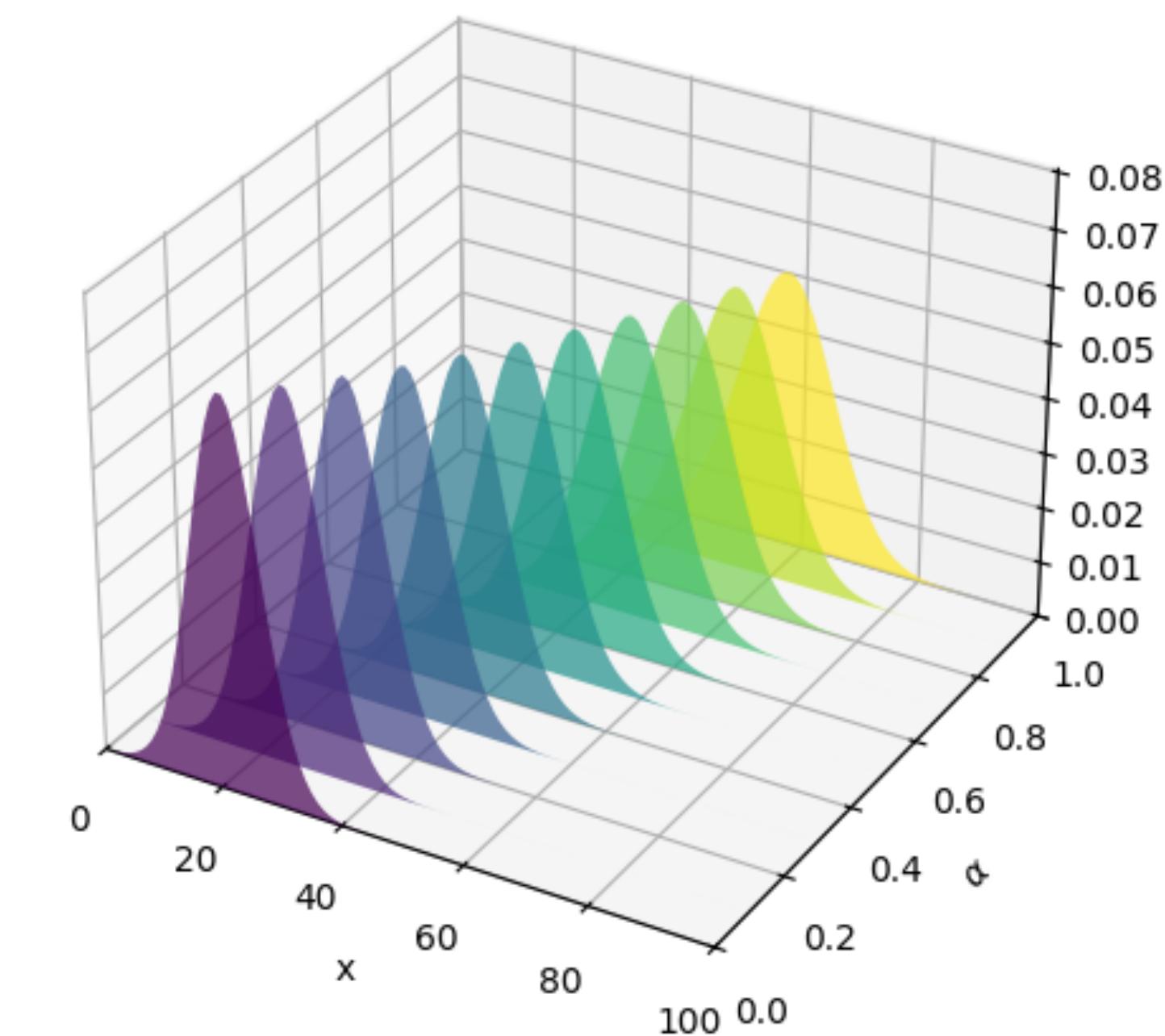


Wasserstein Distances: Barycenter Interpolation

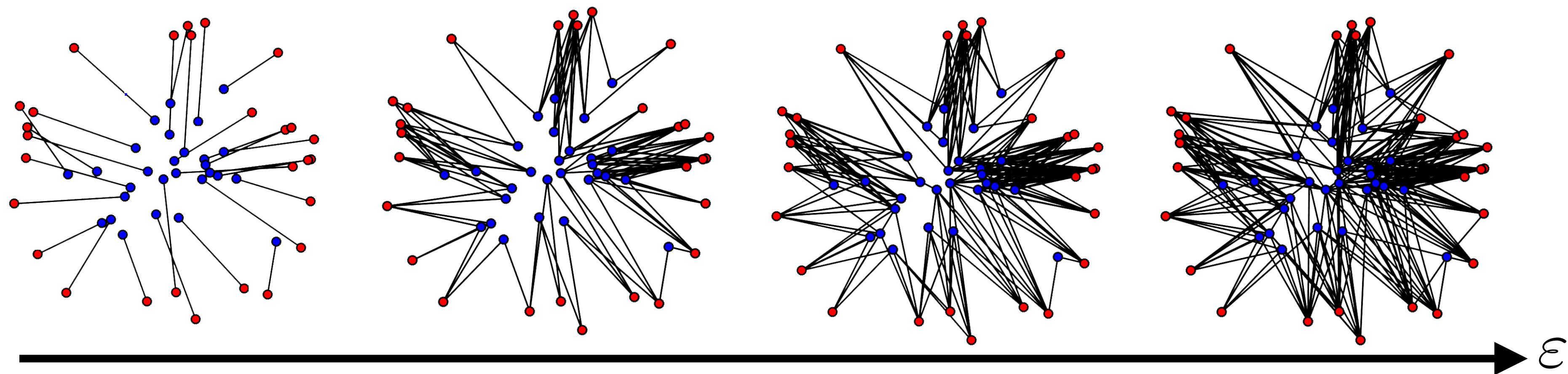
Barycenter interpolation with L^2



Barycenter interpolation with Wasserstein



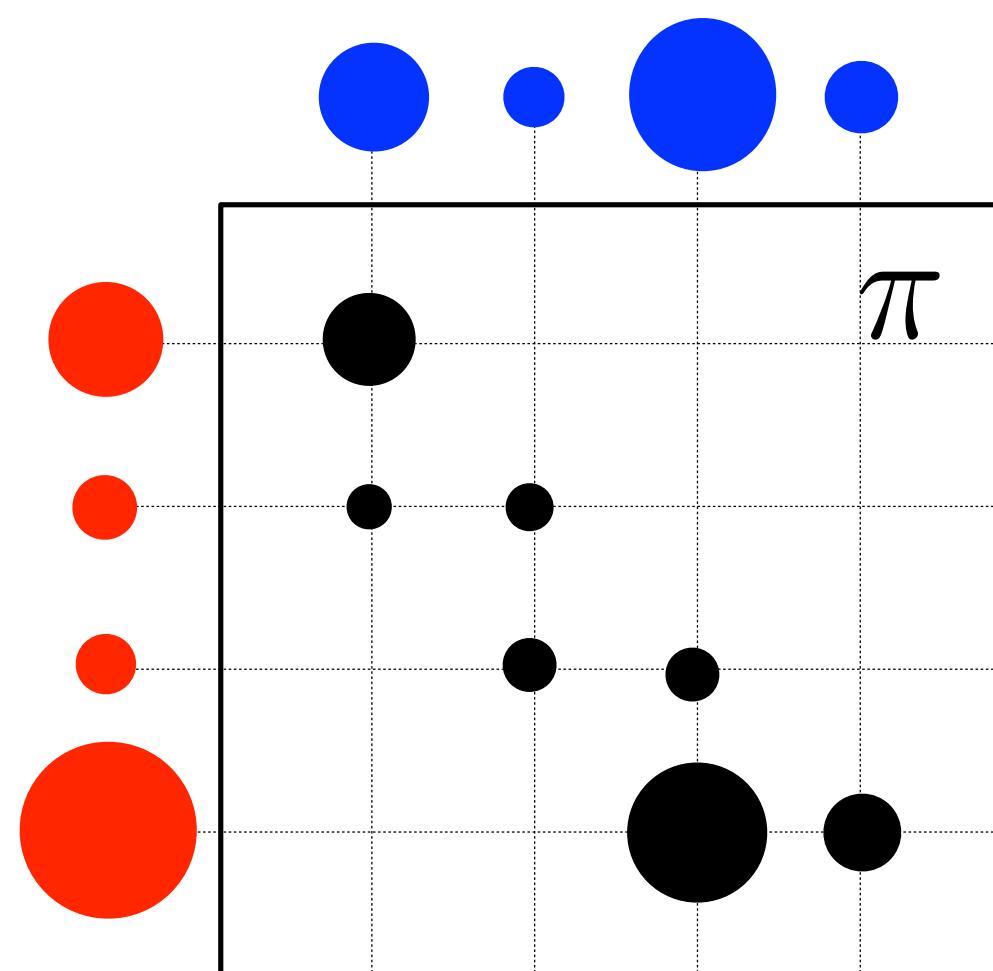
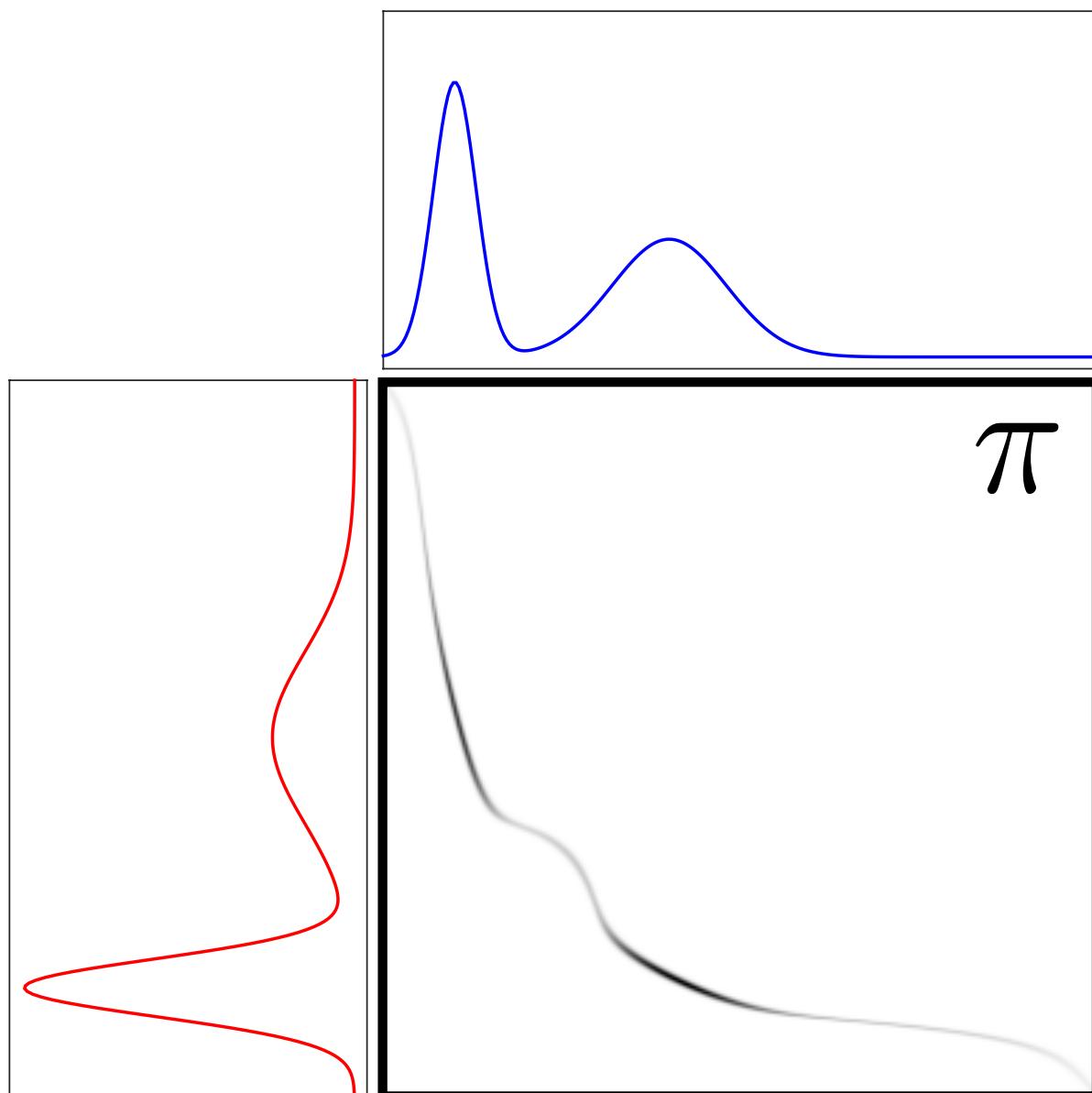
Computational Optimal Transport



1D Optimal Transport: Exact Solution

p -Wasserstein distance for 1D Euclidean ground metric

$$c(x, y) = |x - y|^p$$

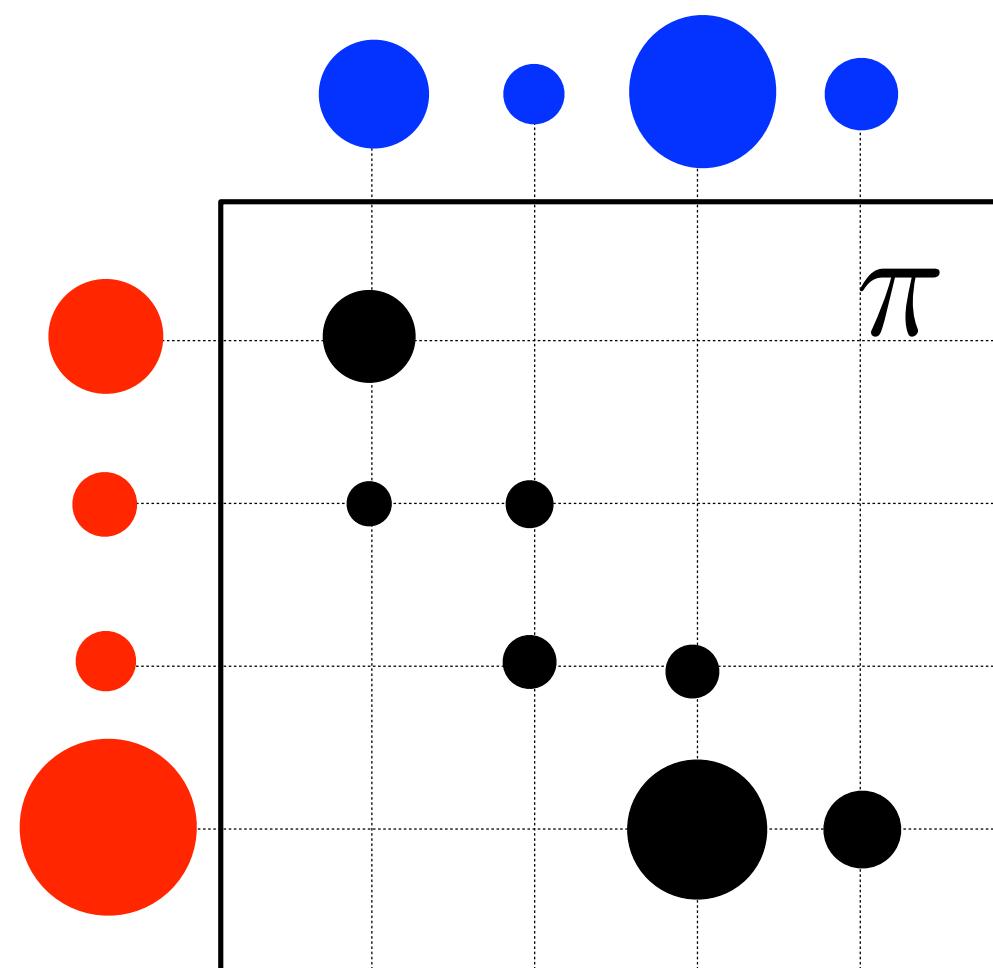
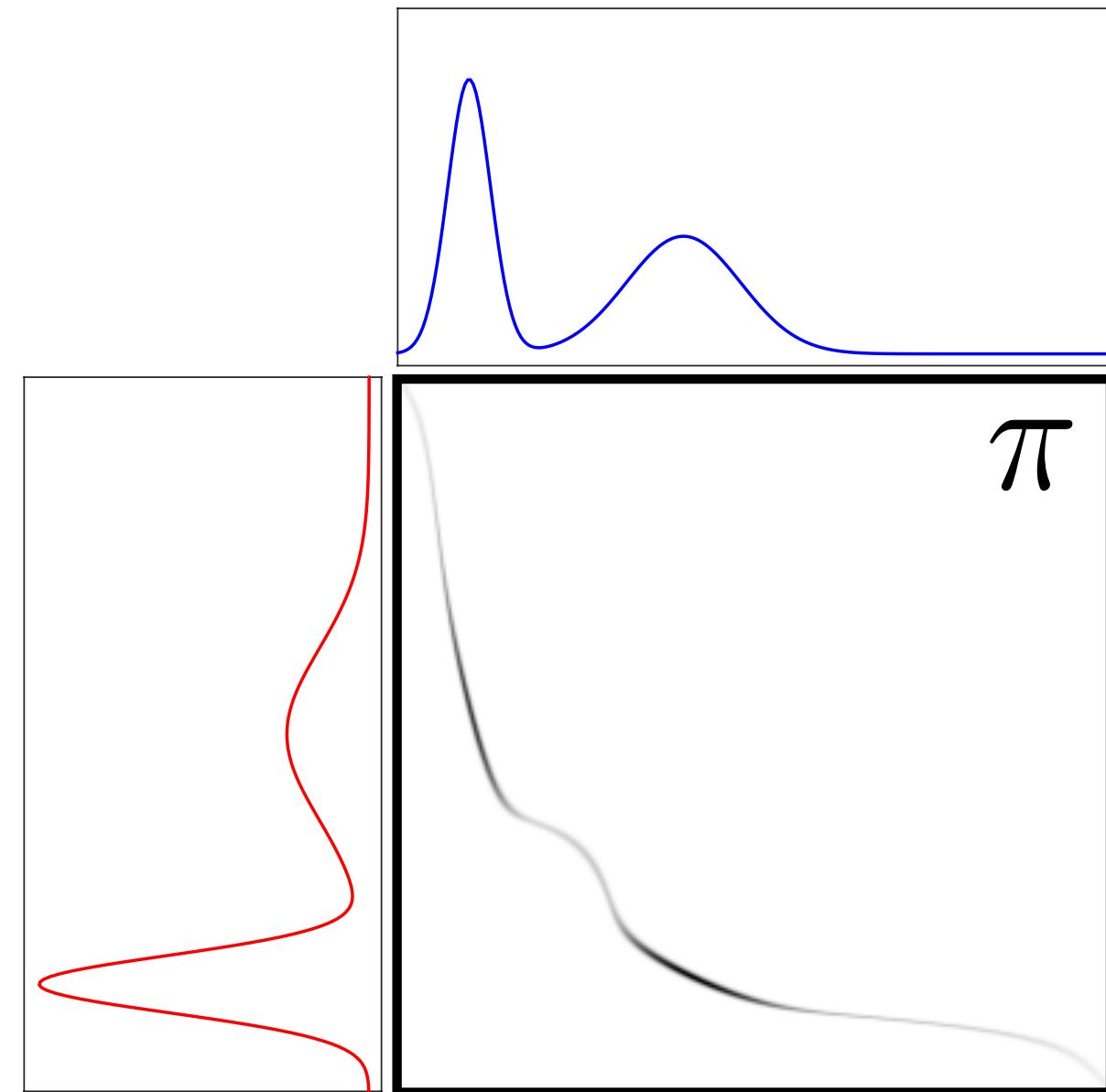


1D Optimal Transport: Exact Solution

p -Wasserstein distance for 1D Euclidean ground metric

$c(x, y) = |x - y|^p$ can be directly computed as

$$W_p(\mu, \nu) = \left(\int_0^1 \left| \text{CDF}_{\mu}^{-1}(r) - \text{CDF}_{\nu}^{-1}(r) \right|^p dr \right)^{1/p}.$$

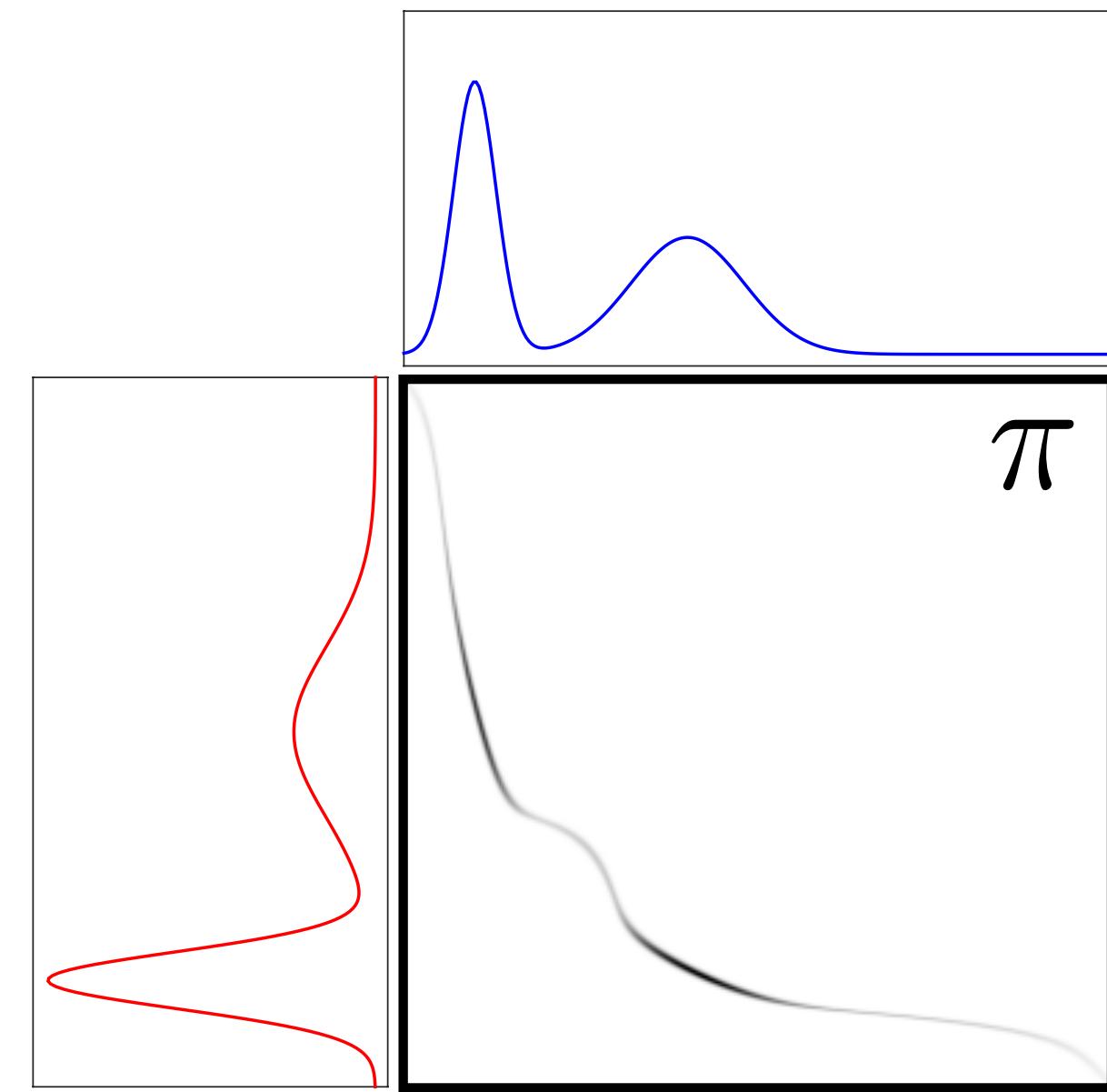


1D Optimal Transport: Exact Solution

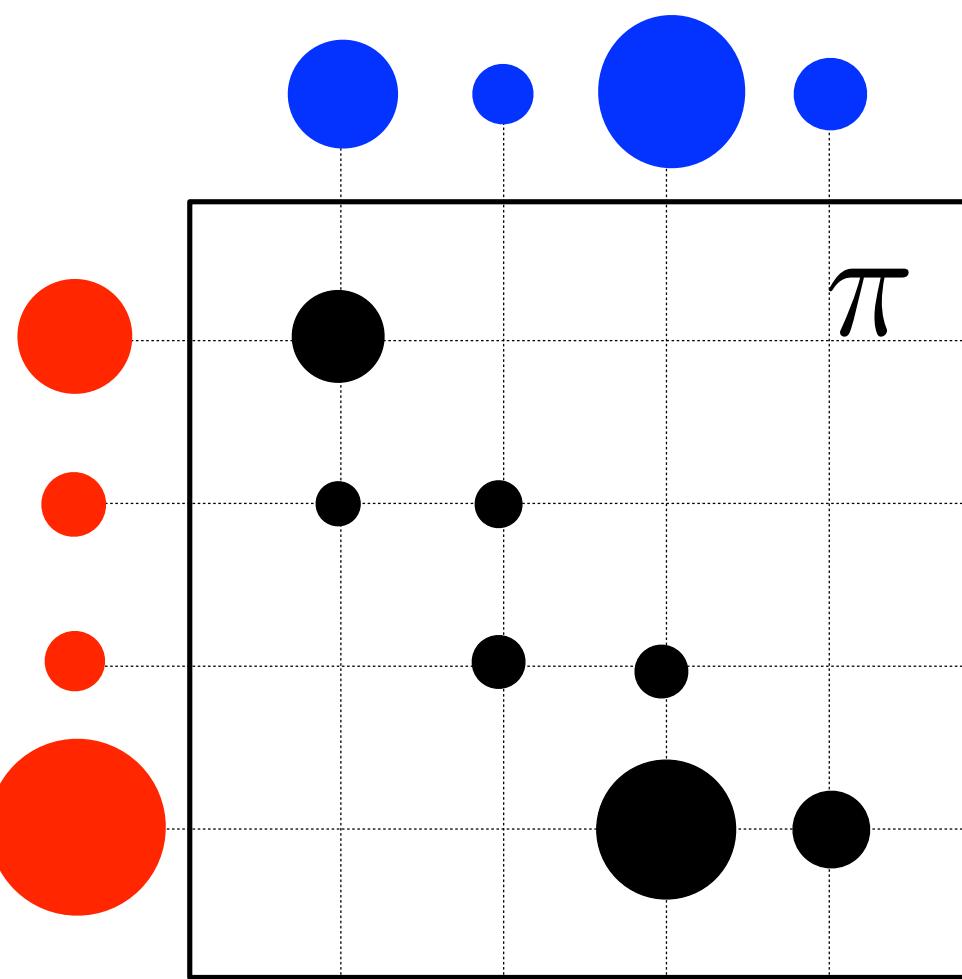
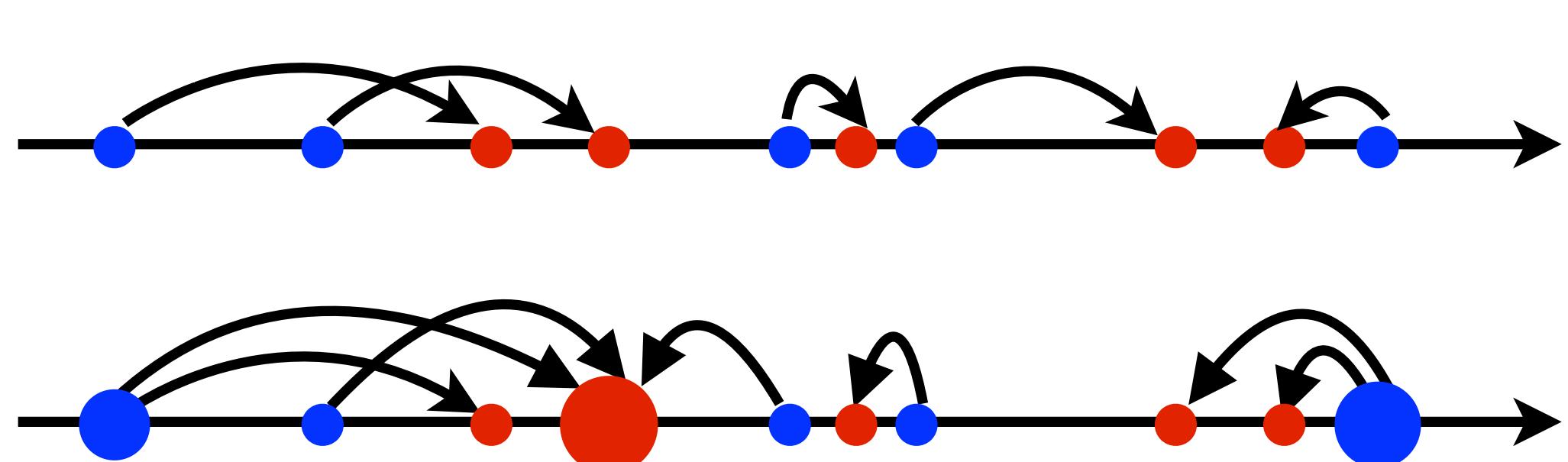
p -Wasserstein distance for 1D Euclidean ground metric

$c(x, y) = |x - y|^p$ can be directly computed as

$$W_p(\mu, \nu) = \left(\int_0^1 \left| \text{CDF}_{\mu}^{-1}(r) - \text{CDF}_{\nu}^{-1}(r) \right|^p dr \right)^{1/p}.$$



Discrete case solved by sorting!



1D Optimal Transport: Sliced Wasserstein Distances

Sliced Wasserstein distance is defined in terms of *1D Wasserstein distances*:

$$\text{SW}_p(\mu, \nu) = \left(\int_{\mathbb{S}^{d-1}} W_p(P_{\theta,\#}\mu, P_{\theta,\#}\nu)^p d\theta \right)^{1/p}$$

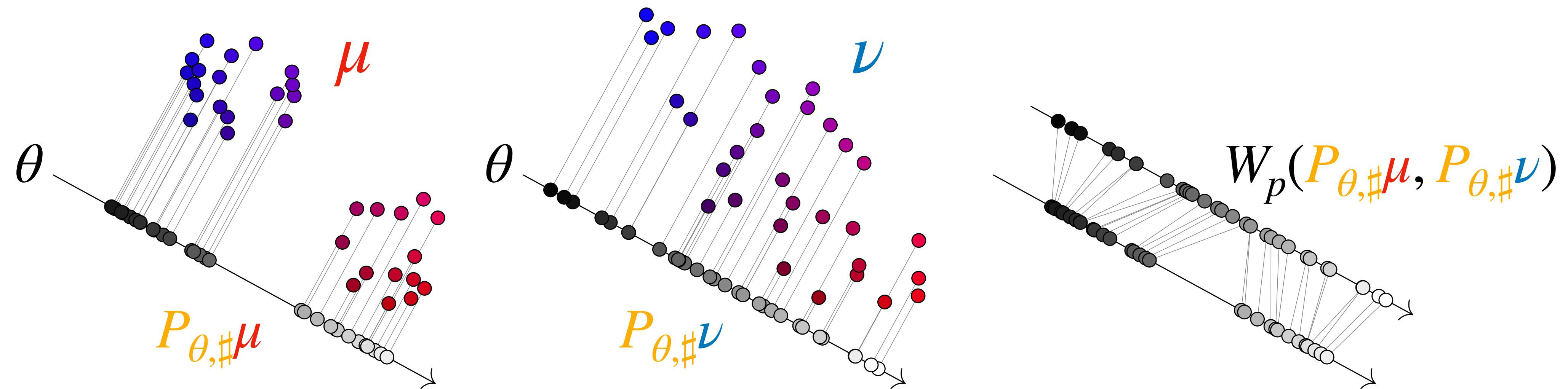
where $P_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$ is a projection onto a 1D subspace $x \mapsto x \cdot \theta$.

1D Optimal Transport: Sliced Wasserstein Distances

Sliced Wasserstein distance is defined in terms of *1D Wasserstein distances*:

$$\text{SW}_p(\mu, \nu) = \left(\mathbb{E}_{\theta \sim \mathbb{S}^{d-1}} [W_p(P_{\theta,\#}\mu, P_{\theta,\#}\nu)^p] \right)^{1/p}$$

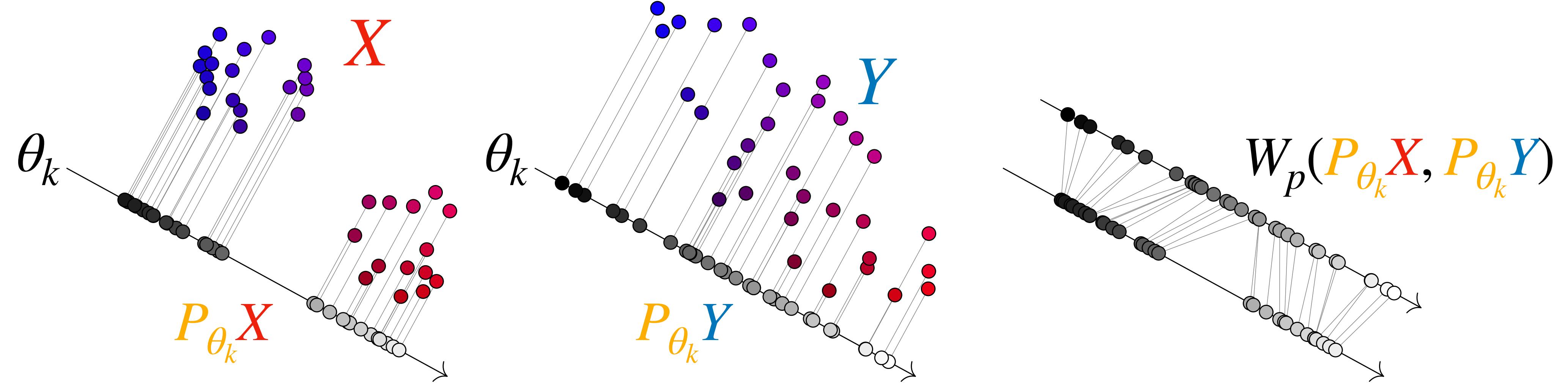
where $P_\theta : \mathbb{R}^d \rightarrow \mathbb{R}$ is a projection onto a 1D subspace $x \mapsto x \cdot \theta$.



1D Optimal Transport: Exercise

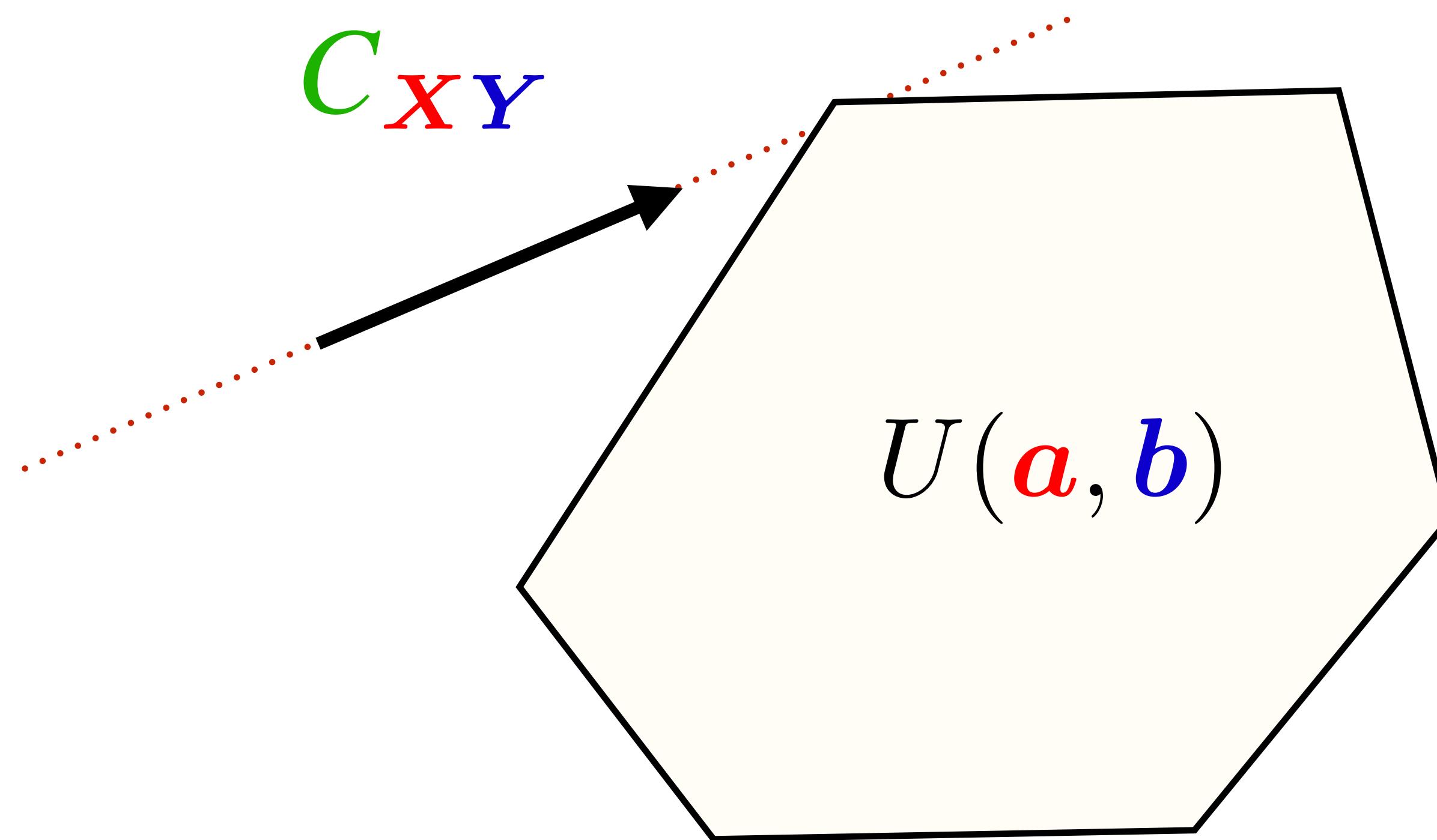
Implement sliced Wasserstein distance for samples $X = \{x_i\}_{i=1}^n$, $Y = \{y_j\}_{j=1}^n$:

$$\text{SW}_p^{(L)}(X, Y) = \left(\frac{1}{L} \sum_{k=1}^L [W_p(\{\theta_k \cdot x_i\}_{i=1}^n, \{\theta_k \cdot y_j\}_{j=1}^n)^p] \right)^{1/p}$$



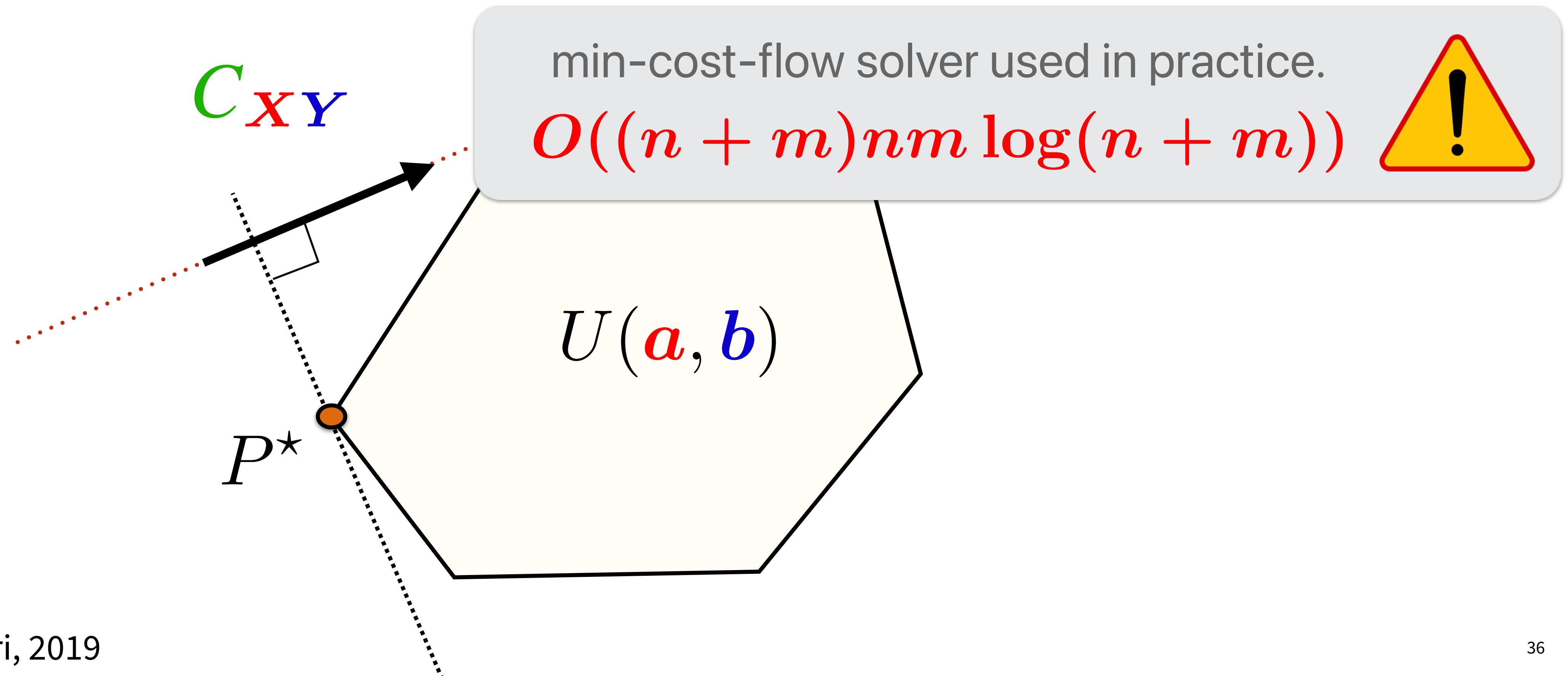
Entropy Regularized OT: Motivation

Efficiency: General linear programming solvers are **computationally expensive**.



Entropy Regularized OT: Motivation

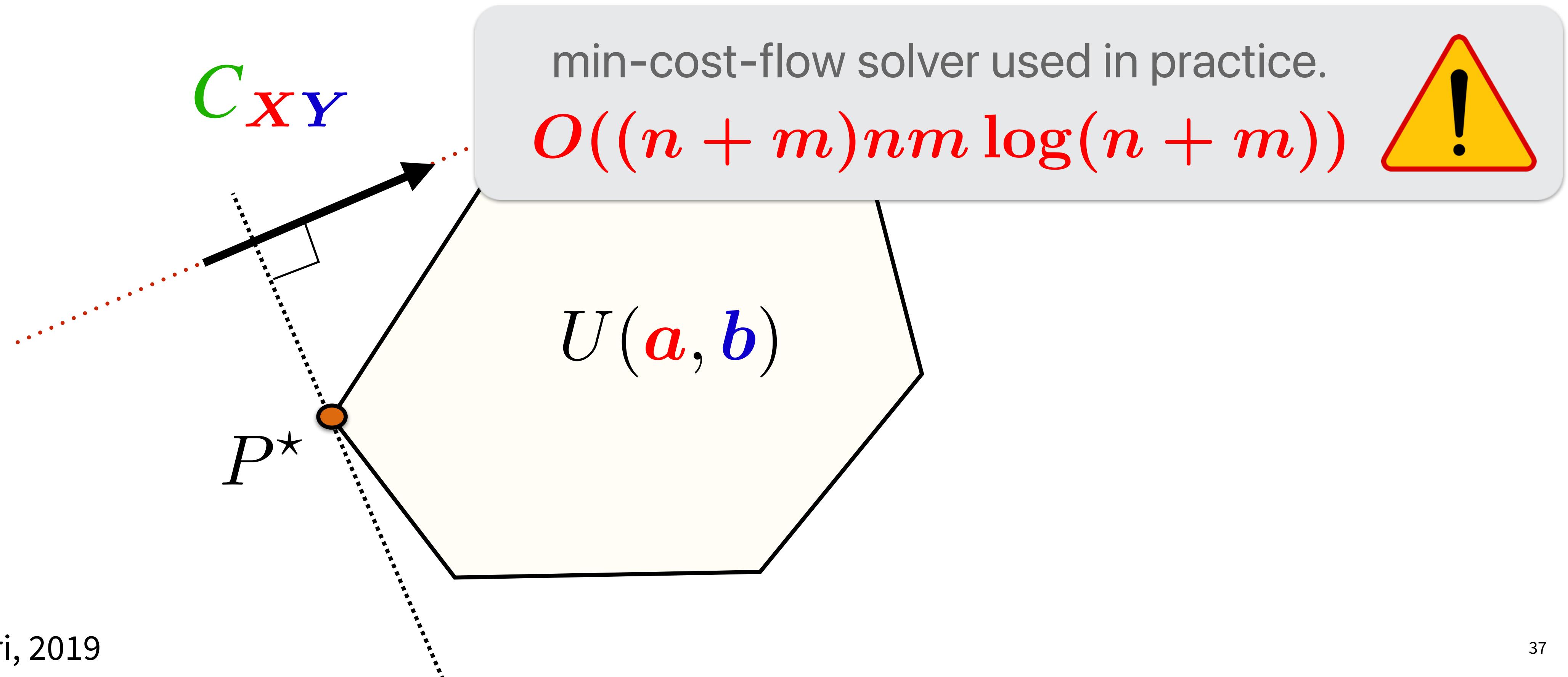
Efficiency: General linear programming solvers are **computationally expensive**.



Entropy Regularized OT: Motivation

Efficiency: General linear programming solvers are **computationally expensive**.

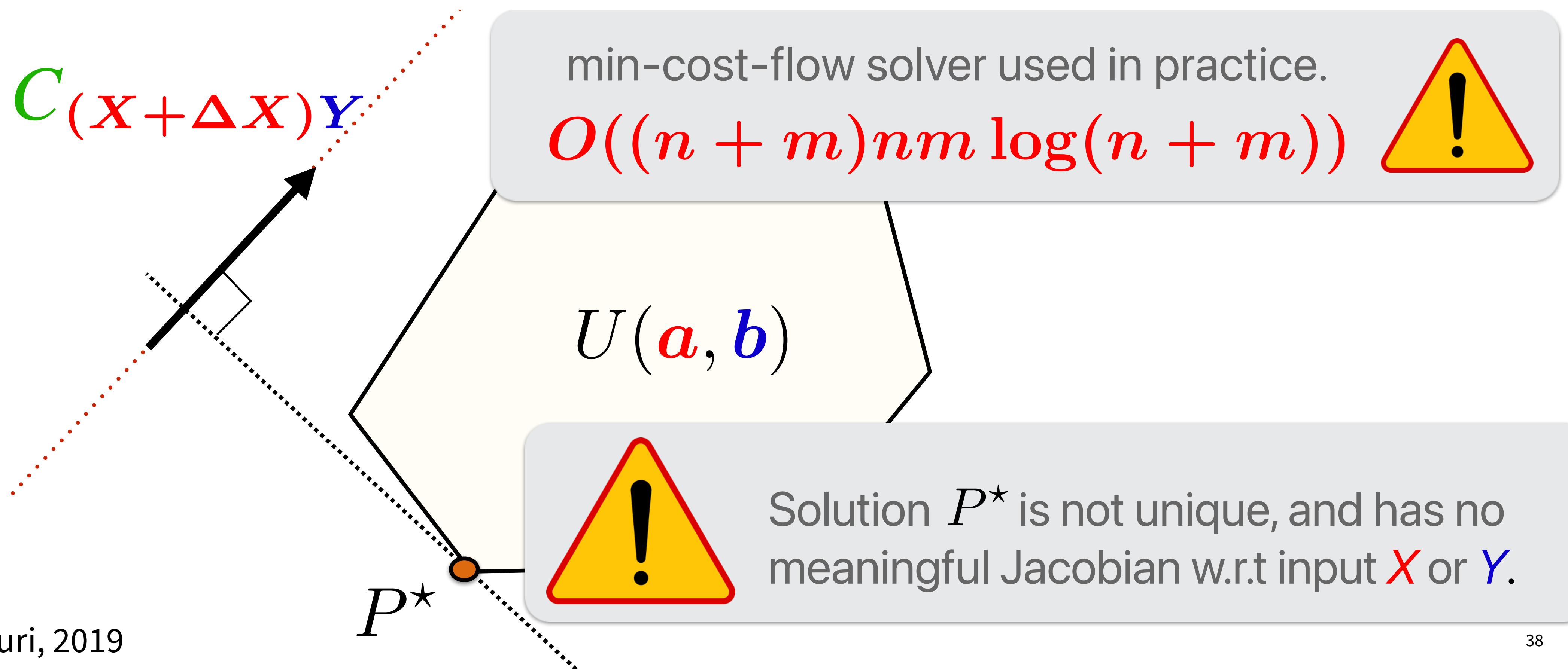
Smoothness: Exact solution may not be **unique** or **differentiable**.



Entropy Regularized OT: Motivation

Efficiency: General linear programming solvers are **computationally expensive**.

Smoothness: Exact solution may not be **unique** or **differentiable**.

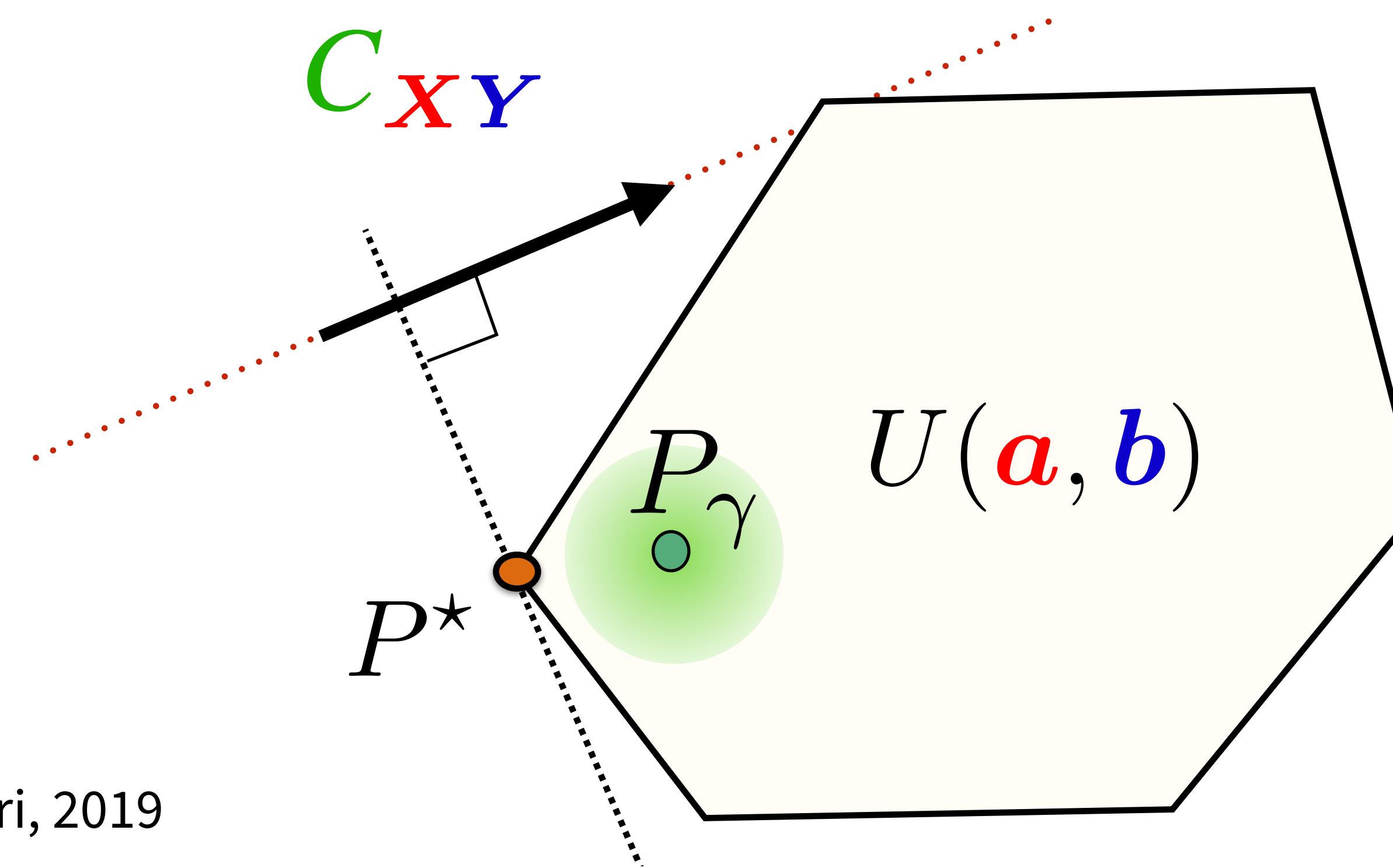


Entropy Regularized OT: Motivation

Efficiency: General linear programming solvers are **computationally expensive**.

Smoothness: Exact solution may not be **unique** or **differentiable**.

Robustness: Do not want to **overfit** on finite samples.

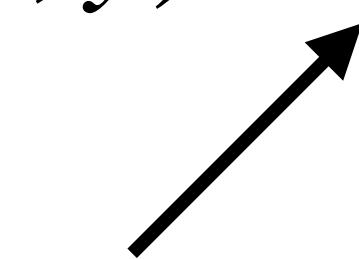


Entropy Regularized OT: Formulation

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \begin{array}{c} \text{total cost} \\ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \end{array} \right\}$$

Entropy Regularized OT: Formulation

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

total cost relative entropy constraint

regularization parameter

Entropy Regularized OT: Formulation

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

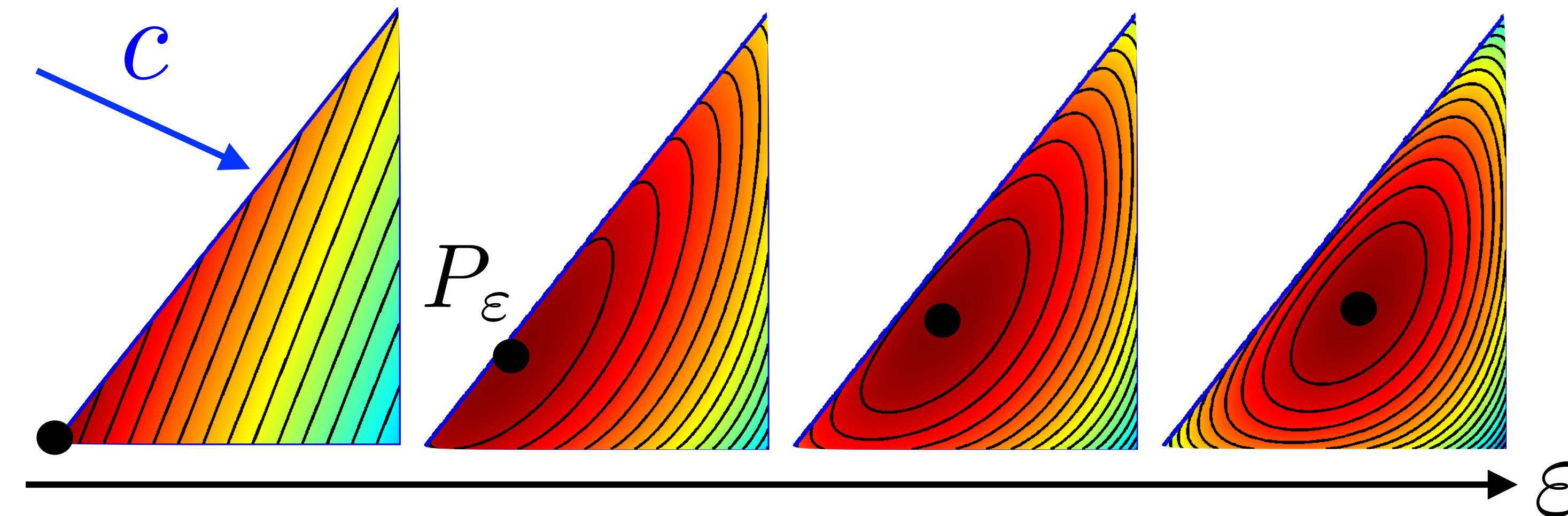
total cost relative entropy constraint
regularization parameter reference distribution



Entropy Regularized OT: Formulation

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

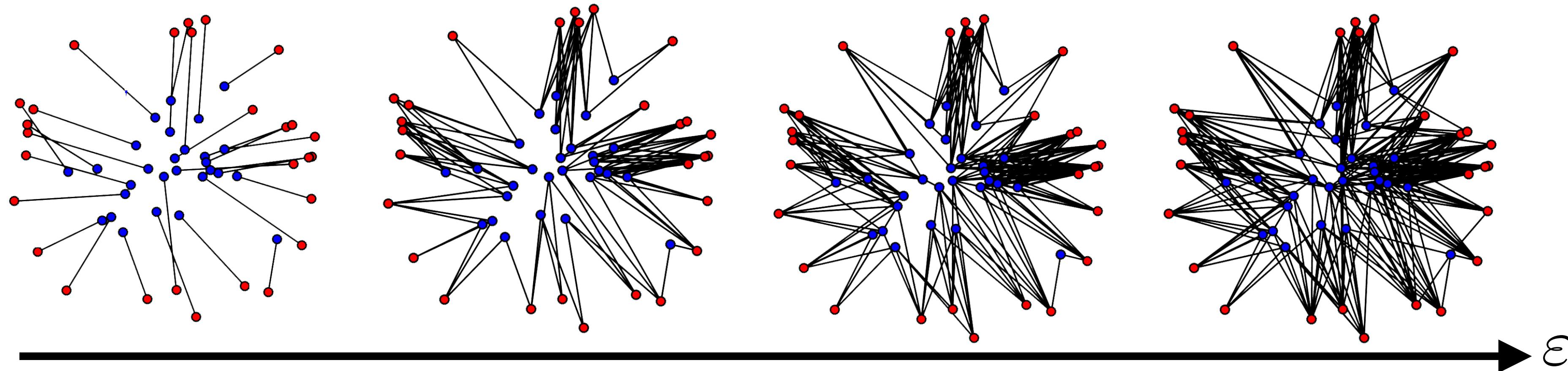
total cost
regularization parameter
relative entropy
reference distribution
constraint



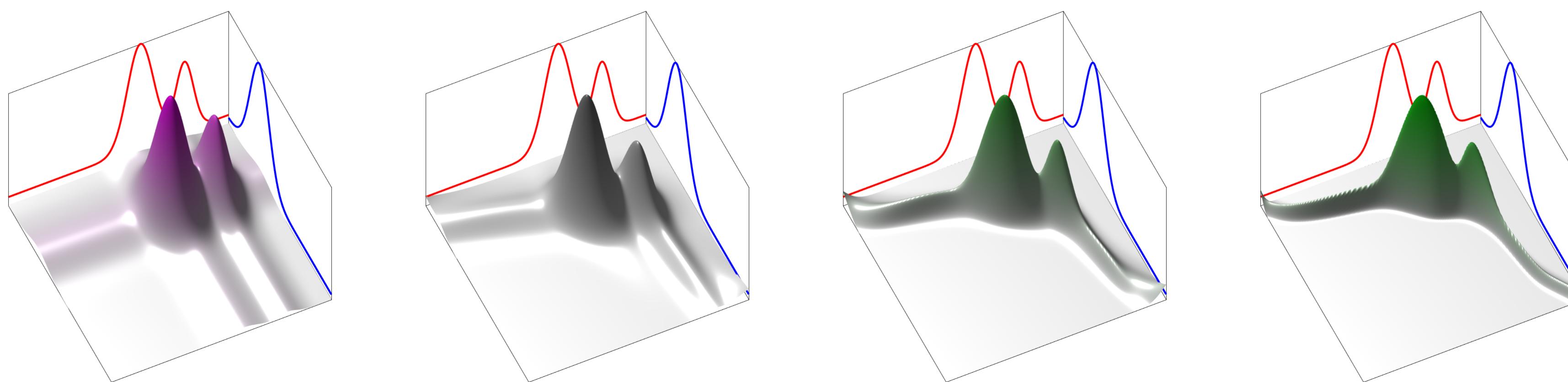
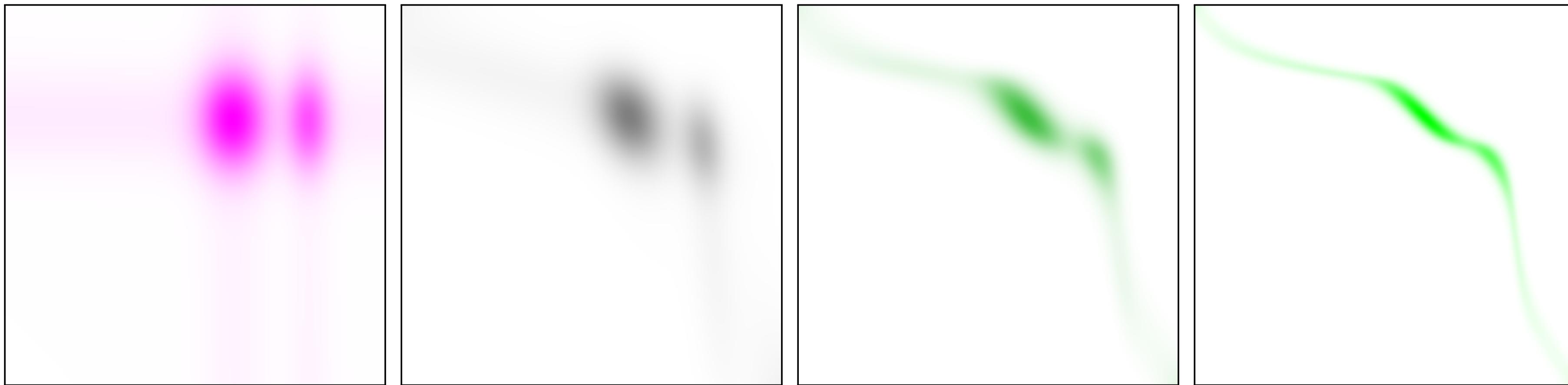
Entropy Regularized OT: Formulation

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

total cost relative entropy constraint
regularization parameter reference distribution



Entropy Regularized OT: Formulation



$\varepsilon = 10$

$\varepsilon = 1$

$\varepsilon = 10^{-1}$

$\varepsilon = 10^{-2}$

Entropy Regularized OT: Discrete Formulation

$$L_C^{(\varepsilon)}(\mathbf{a}, \mathbf{b}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} - \varepsilon H(P) : P \in U(\mathbf{a}, \mathbf{b}) \right\}$$

Entropy Regularized OT: Discrete Formulation

$$L_C^{(\varepsilon)}(\mathbf{a}, \mathbf{b}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} - \varepsilon H(P) : P \in U(\mathbf{a}, \mathbf{b}) \right\}$$

Dual formulation:

$$\tilde{L}_C^{(\varepsilon)}(\mathbf{f}, \mathbf{g}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} - \varepsilon H(P) - \sum_i \mathbf{f}_i \left(\sum_j P_{ij} - a_i \right) - \sum_j \mathbf{g}_j \left(\sum_i P_{ij} - b_j \right) \right\}$$

Entropy Regularized OT: Discrete Formulation

$$L_C^{(\varepsilon)}(\mathbf{a}, \mathbf{b}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} - \varepsilon H(P) : P \in U(\mathbf{a}, \mathbf{b}) \right\}$$

Dual formulation:

$$\tilde{L}_C^{(\varepsilon)}(\mathbf{f}, \mathbf{g}) = \min_P \left\{ \sum_{i,j} C_{ij} P_{ij} - \varepsilon H(P) - \sum_i \mathbf{f}_i \left(\sum_j P_{ij} - a_i \right) - \sum_j \mathbf{g}_j \left(\sum_i P_{ij} - b_j \right) \right\}$$

$$P_{ij} = u_i K_{ij} v_j \quad \text{where} \quad u_i = \exp(f_i/\varepsilon), v_j = \exp(g_j/\varepsilon), K_{ij} = \exp(-C_{ij}/\varepsilon)$$

Entropy Regularized OT: Sinkhorn Algorithm

$$P_{ij} = \textcolor{red}{u}_i K_{ij} \textcolor{blue}{v}_j \quad \text{where} \quad \textcolor{red}{u}_i = \exp(\textcolor{red}{f}_i/\varepsilon), \textcolor{blue}{v}_j = \exp(\textcolor{blue}{g}_j/\varepsilon), K_{ij} = \exp(-C_{ij}/\varepsilon)$$

Impose constraints: $\sum_i P_{ij} = \textcolor{blue}{b}_j, \sum_j P_{ij} = \textcolor{red}{a}_i, P_{ij} \geq 0$

Entropy Regularized OT: Sinkhorn Algorithm

$$P_{ij} = \textcolor{red}{u}_i K_{ij} \textcolor{blue}{v}_j \quad \text{where} \quad \textcolor{red}{u}_i = \exp(\textcolor{red}{f}_i/\varepsilon), \textcolor{blue}{v}_j = \exp(\textcolor{blue}{g}_j/\varepsilon), K_{ij} = \exp(-C_{ij}/\varepsilon)$$

Impose constraints: $\sum_i P_{ij} = \textcolor{blue}{b}_j, \sum_j P_{ij} = \textcolor{red}{a}_i, P_{ij} \geq 0$

Sinkhorn iterations—also known as iterative proportional fitting (IPF):

$$\textcolor{red}{u}_i^{(l+1)} = \textcolor{red}{a}_i / (K \textcolor{blue}{v}^{(l)})_i$$

$$\textcolor{blue}{v}_j^{(l+1)} = \textcolor{blue}{b}_j / (K^\top \textcolor{red}{u}^{(l+1)})_j$$

Entropy Regularized OT: Sinkhorn Algorithm

$$P_{ij} = \textcolor{red}{u}_i K_{ij} \textcolor{blue}{v}_j \quad \text{where} \quad \textcolor{red}{u}_i = \exp(\textcolor{red}{f}_i/\varepsilon), \textcolor{blue}{v}_j = \exp(\textcolor{blue}{g}_j/\varepsilon), K_{ij} = \exp(-C_{ij}/\varepsilon)$$

Impose constraints: $\sum_i P_{ij} = \textcolor{blue}{b}_j, \sum_j P_{ij} = \textcolor{red}{a}_i, P_{ij} \geq 0$

Sinkhorn iterations—also known as iterative proportional fitting (IPF):

$$\mathcal{O}(n^2 \log n)$$

$$\textcolor{red}{u}_i^{(l+1)} = \textcolor{red}{a}_i / (K \textcolor{blue}{v}^{(l)})_i$$

$$\textcolor{blue}{v}_j^{(l+1)} = \textcolor{blue}{b}_j / (K^\top \textcolor{red}{u}^{(l+1)})_j$$

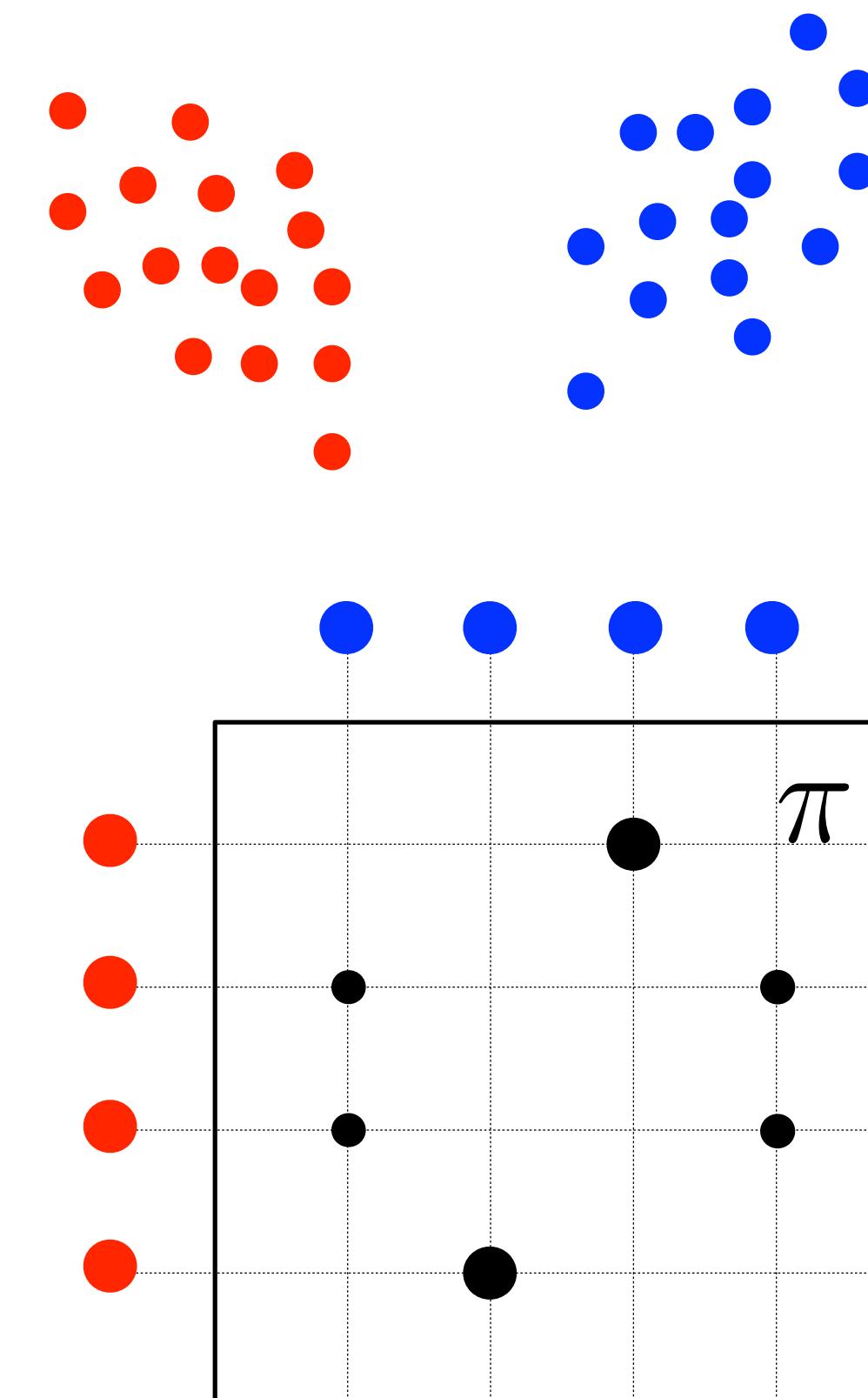
Entropy Regularized OT: Exercise

Implement the **Sinkhorn algorithm** for samples $X = \{x_i\}_{i=1}^n$, $Y = \{y_j\}_{j=1}^n$.

$$P_{ij} = u_i K_{ij} v_j, \quad K_{ij} = \exp(-C_{ij}/\varepsilon)$$

$$u_i^{(l+1)} = a_i / (K v^{(l)})_i$$

$$v_j^{(l+1)} = b_j / (K^\top u^{(l+1)})_j$$



Entropy Regularized OT: Sinkhorn Divergence

Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

Entropy Regularized OT: Sinkhorn Divergence

Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

Entropy Regularized OT: Sinkhorn Divergence

Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

Entropic bias $L_c^{(\varepsilon)}(\mu, \mu) \neq 0$

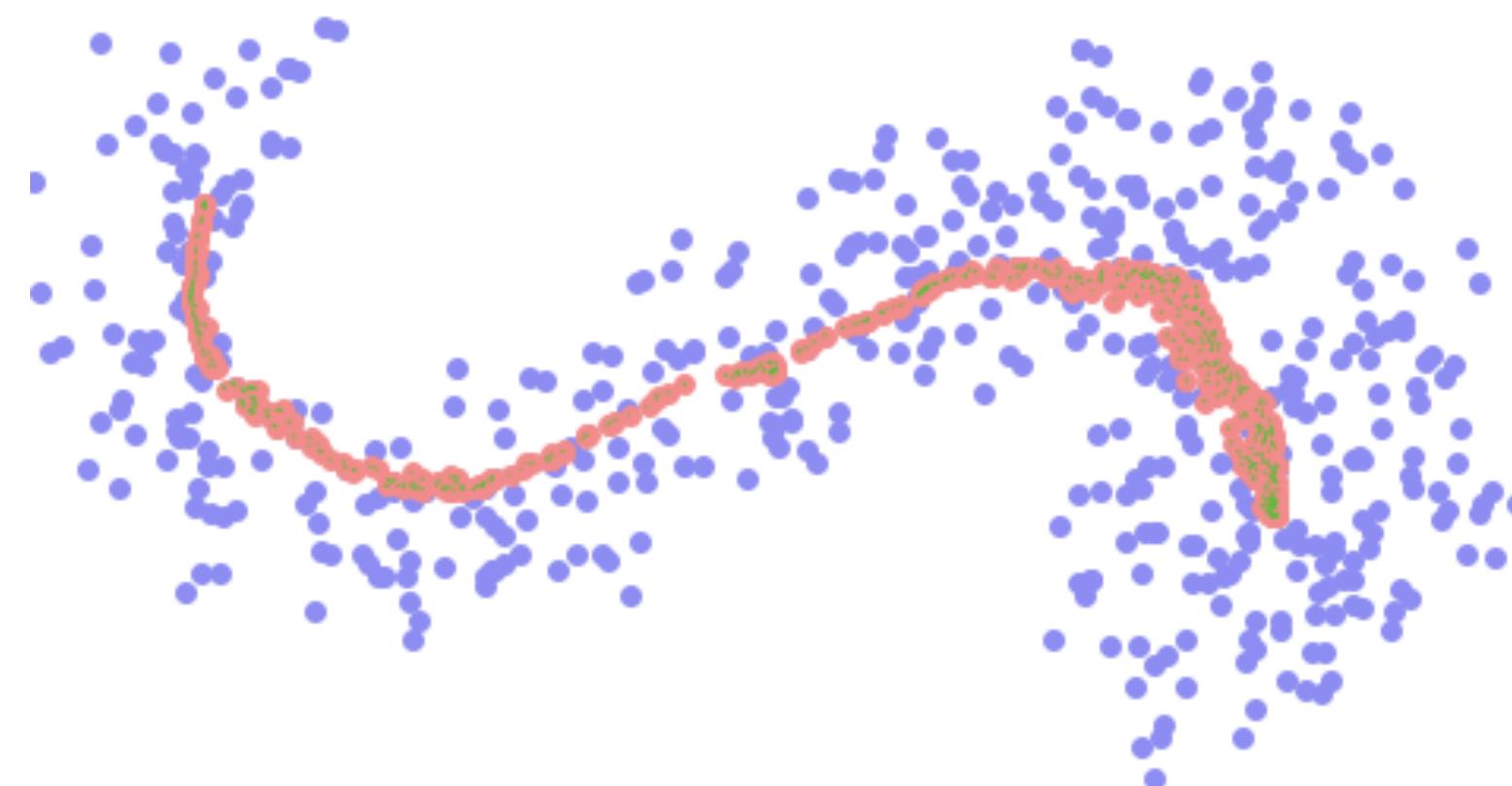
Entropy Regularized OT: Sinkhorn Divergence

Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

$$L_c^{(\varepsilon)}(\mu, \nu) = \min_{\pi} \left\{ \int_{X \times Y} c(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi \mid \mu \otimes \nu) : \pi \in \Pi(\mu, \nu) \right\}$$

Entropic bias $L_c^{(\varepsilon)}(\mu, \mu) \neq 0$

“Shrinkage” $\min_{\mu} L_c^{(\varepsilon)}(\mu, \nu)$ \longrightarrow



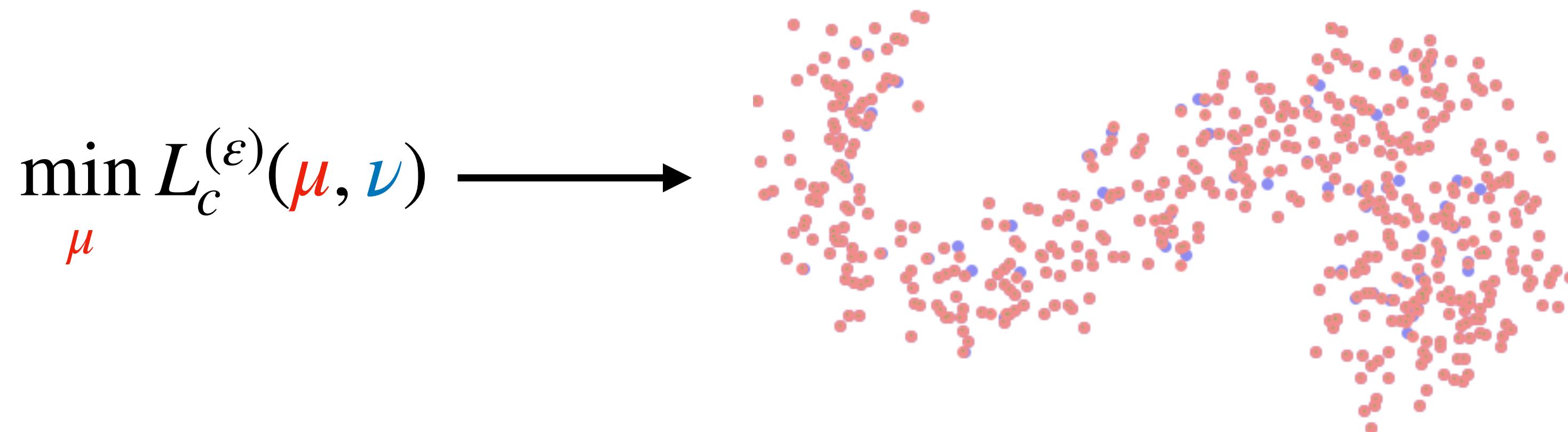
Entropy Regularized OT: Sinkhorn Divergence

Sinkhorn divergence

$$S_c^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu) - \frac{1}{2} \left(L_c^{(\varepsilon)}(\mu, \mu) + L_c^{(\varepsilon)}(\nu, \nu) \right)$$

is an *unbiased, smooth, symmetric, positive definite* divergence with

$$0 = S_c^{(\varepsilon)}(\mu, \mu) \leq S_c^{(\varepsilon)}(\mu, \nu) \text{ and } \mu = \nu \iff S_c^{(\varepsilon)}(\mu, \nu) = 0.$$



Entropy Regularized OT: Sinkhorn Divergence

Sinkhorn divergence

$$S_c^{(\varepsilon)}(\mu, \nu) = L_c^{(\varepsilon)}(\mu, \nu) - \frac{1}{2} \left(L_c^{(\varepsilon)}(\mu, \mu) + L_c^{(\varepsilon)}(\nu, \nu) \right)$$

is an *unbiased, smooth, symmetric, positive definite* divergence with

$$0 = S_c^{(\varepsilon)}(\mu, \mu) \leq S_c^{(\varepsilon)}(\mu, \nu) \text{ and } \mu = \nu \iff S_c^{(\varepsilon)}(\mu, \nu) = 0.$$

ε interpolates between Wasserstein and MMD:

$$W_p(\mu, \nu) \xleftarrow{0 \leftarrow \varepsilon} S_c^{(\varepsilon)}(\mu, \nu) \xrightarrow{\varepsilon \rightarrow \infty} \text{MMD}_{-c}(\mu, \nu)$$

Entropy Regularized OT: Sinkhorn Divergence

Sinkhorn divergence (alternate definition)

$$\overline{S}_c^{(\varepsilon)}(\mu, \nu) = \overline{L}_c^{(\varepsilon)}(\mu, \nu) - \frac{1}{2} (\overline{L}_c^{(\varepsilon)}(\mu, \mu) + \overline{L}_c^{(\varepsilon)}(\nu, \nu)),$$

where $\overline{L}_c^{(\varepsilon)}(\mu, \nu) = \int_{X \times Y} c(x, y) d\pi^{(\varepsilon)\star}(x, y)$ is evaluated *without the entropy term*.

Entropy Regularized OT: Sinkhorn Divergence

Sinkhorn divergence (alternate definition)

$$\overline{S}_c^{(\varepsilon)}(\mu, \nu) = \overline{L}_c^{(\varepsilon)}(\mu, \nu) - \frac{1}{2} (\overline{L}_c^{(\varepsilon)}(\mu, \mu) + \overline{L}_c^{(\varepsilon)}(\nu, \nu)),$$

where $\overline{L}_c^{(\varepsilon)}(\mu, \nu) = \int_{X \times Y} c(x, y) d\pi^{(\varepsilon)\star}(x, y)$ is evaluated *without the entropy term*.

$\overline{S}_c^{(\varepsilon)}(\mu, \nu)$ has similar properties to $S_c^{(\varepsilon)}(\mu, \nu)$.

POT library implements $\overline{S}_c^{(\varepsilon)}(\mu, \nu)$, while OTT and GeomLoss implement $S_c^{(\varepsilon)}(\mu, \nu)$.

Entropy Regularized OT: Sinkhorn Divergence

Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = S_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

Entropy Regularized OT: Sinkhorn Divergence

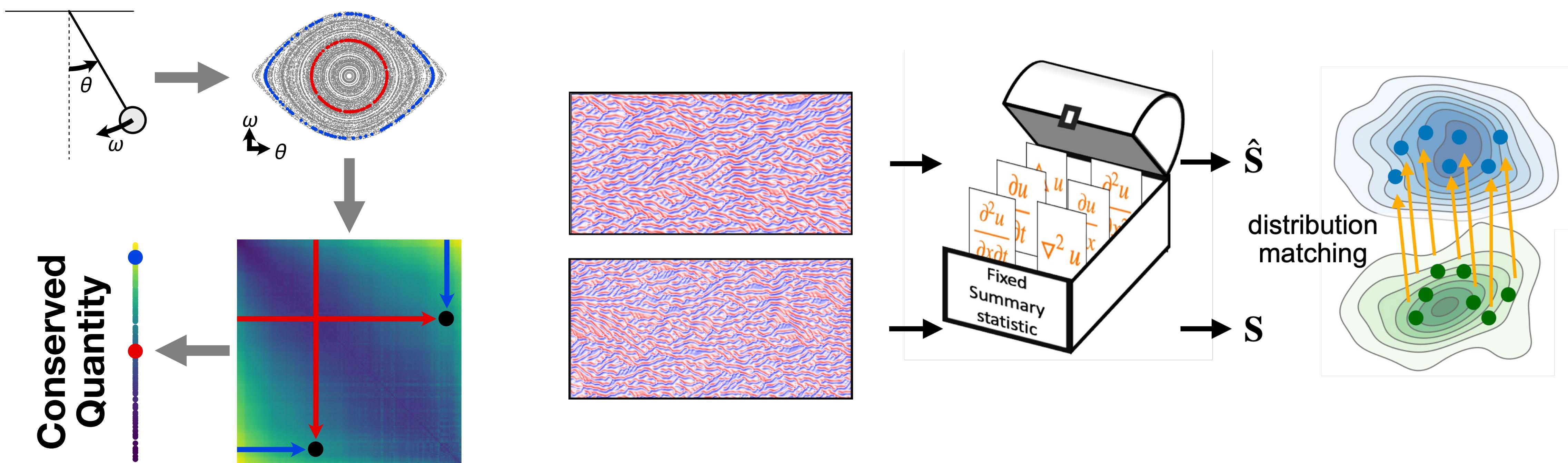
Define a distance $W_p^{(\varepsilon)}(\mu, \nu) = S_c^{(\varepsilon)}(\mu, \nu)^{1/p}$?

Unfortunately, this does not always satisfy triangle inequality

$$W_p^{(\varepsilon)}(\mu, \rho) + W_p^{(\varepsilon)}(\rho, \nu) \not\geq W_p^{(\varepsilon)}(\mu, \nu)$$

but, nevertheless, is often used in practice as a loss function.

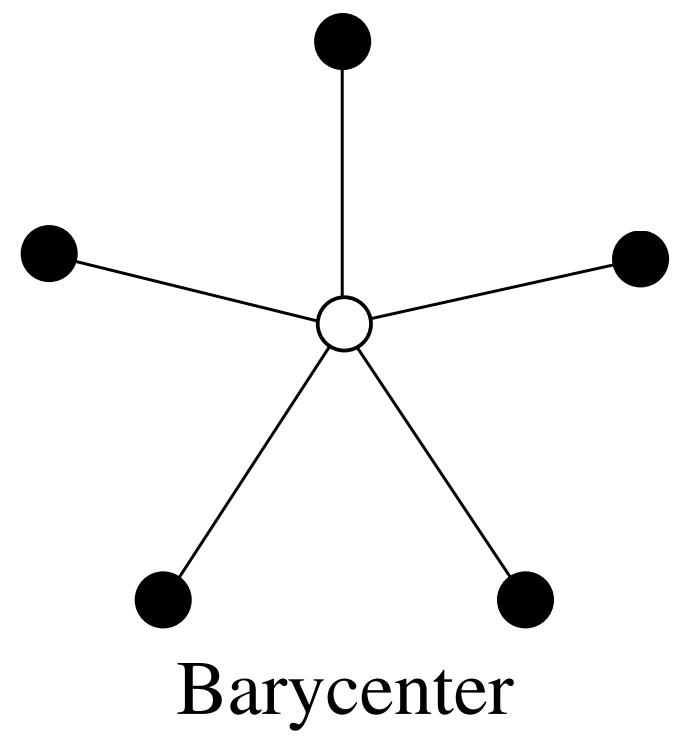
Machine Learning Applications



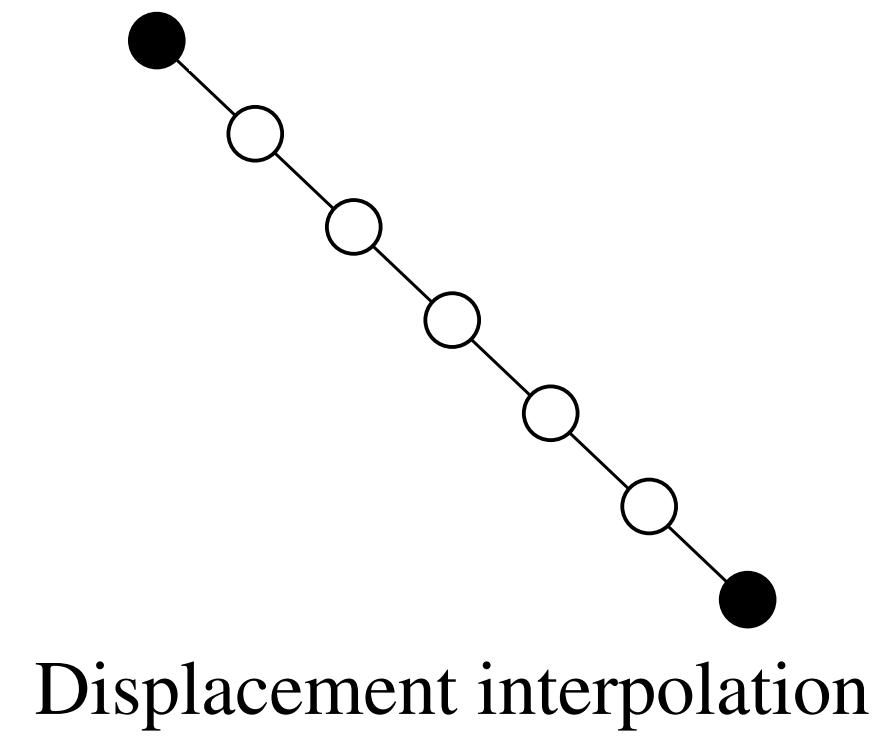
Shape Analysis: Barycenters & Interpolation

Weighted Fréchet mean:

$$\min_{\mu} \sum_i \alpha_i W_2(\mu, \nu_i)^2$$



Barycenter

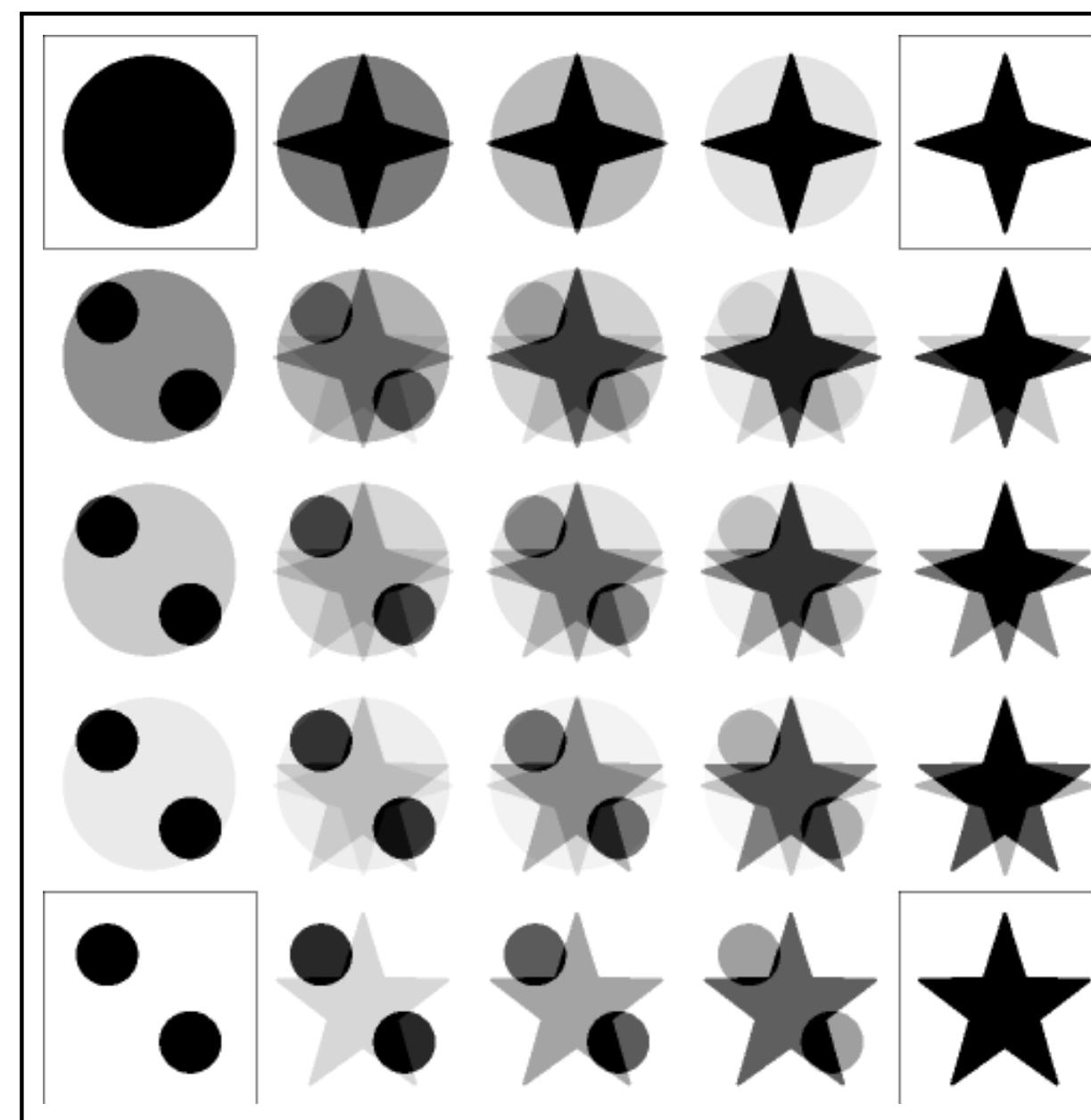
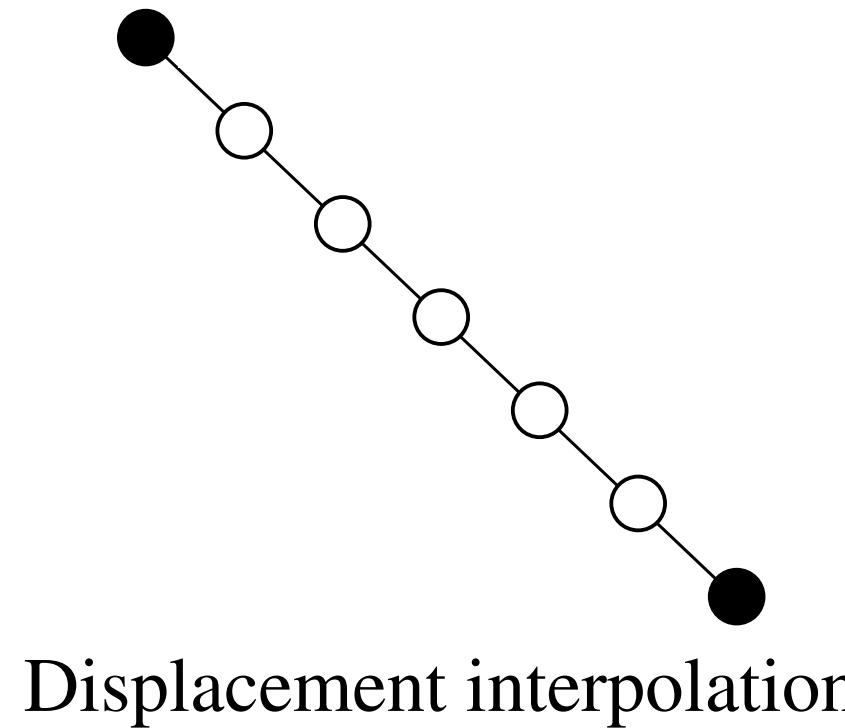
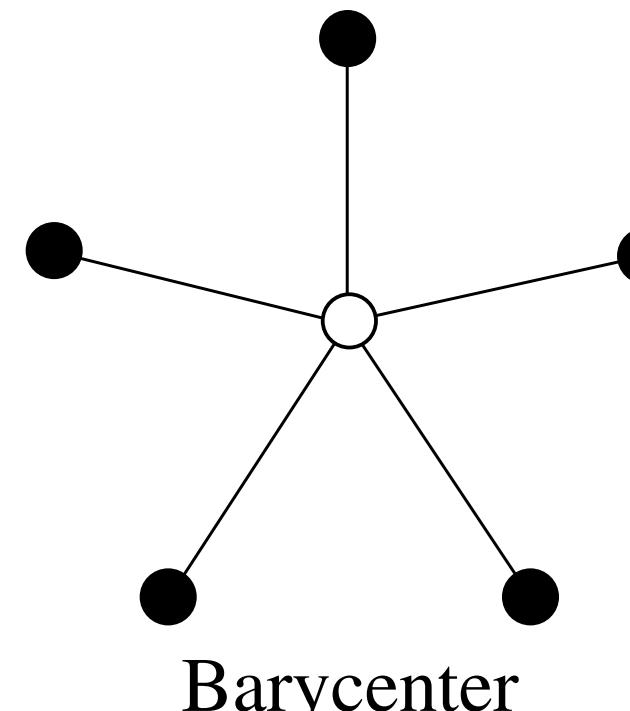


Displacement interpolation

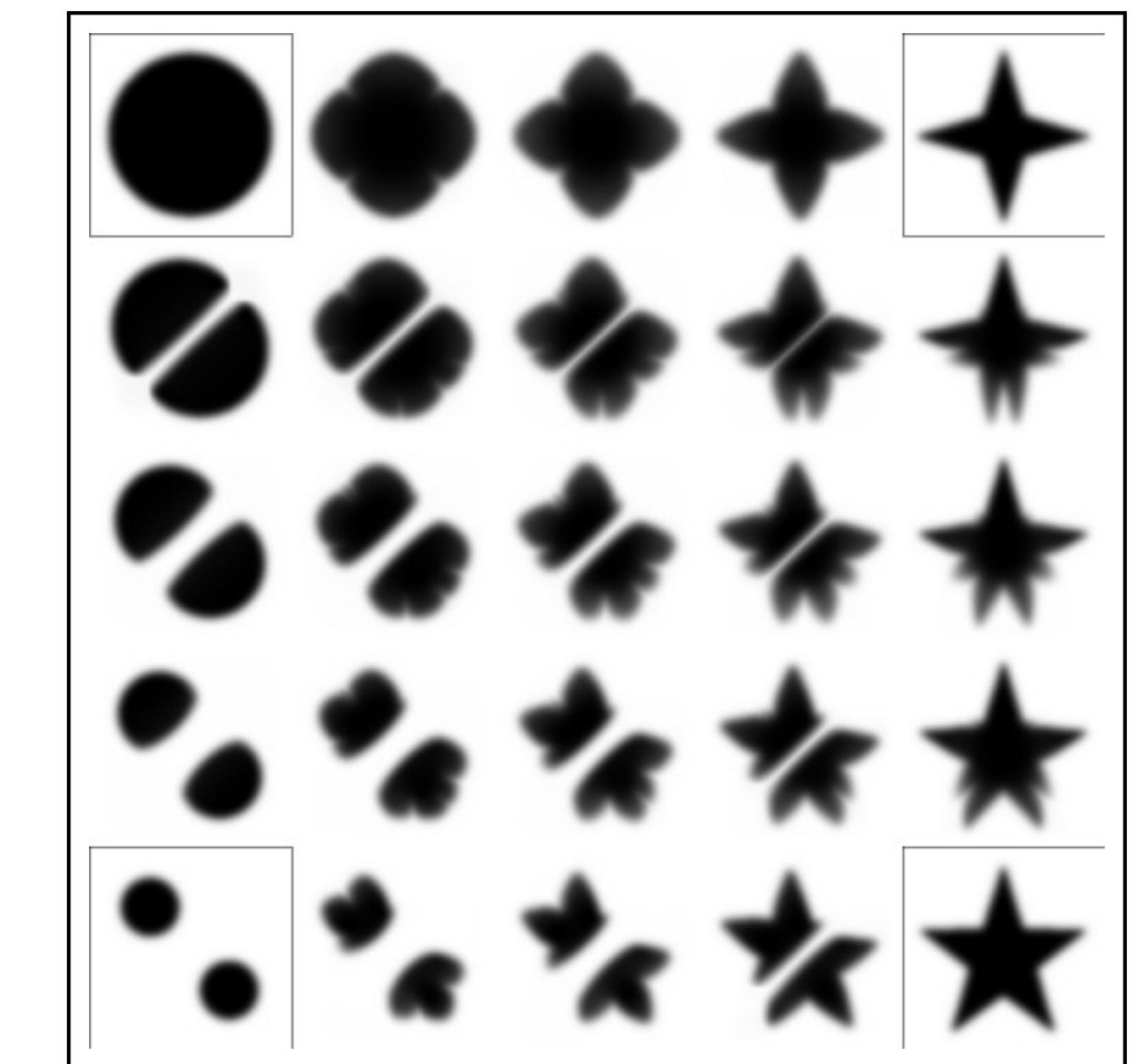
Shape Analysis: Barycenters & Interpolation

Weighted Fréchet mean:

$$\min_{\mu} \sum_i \alpha_i W_2(\mu, \nu_i)^2$$



Euclidean barycenter

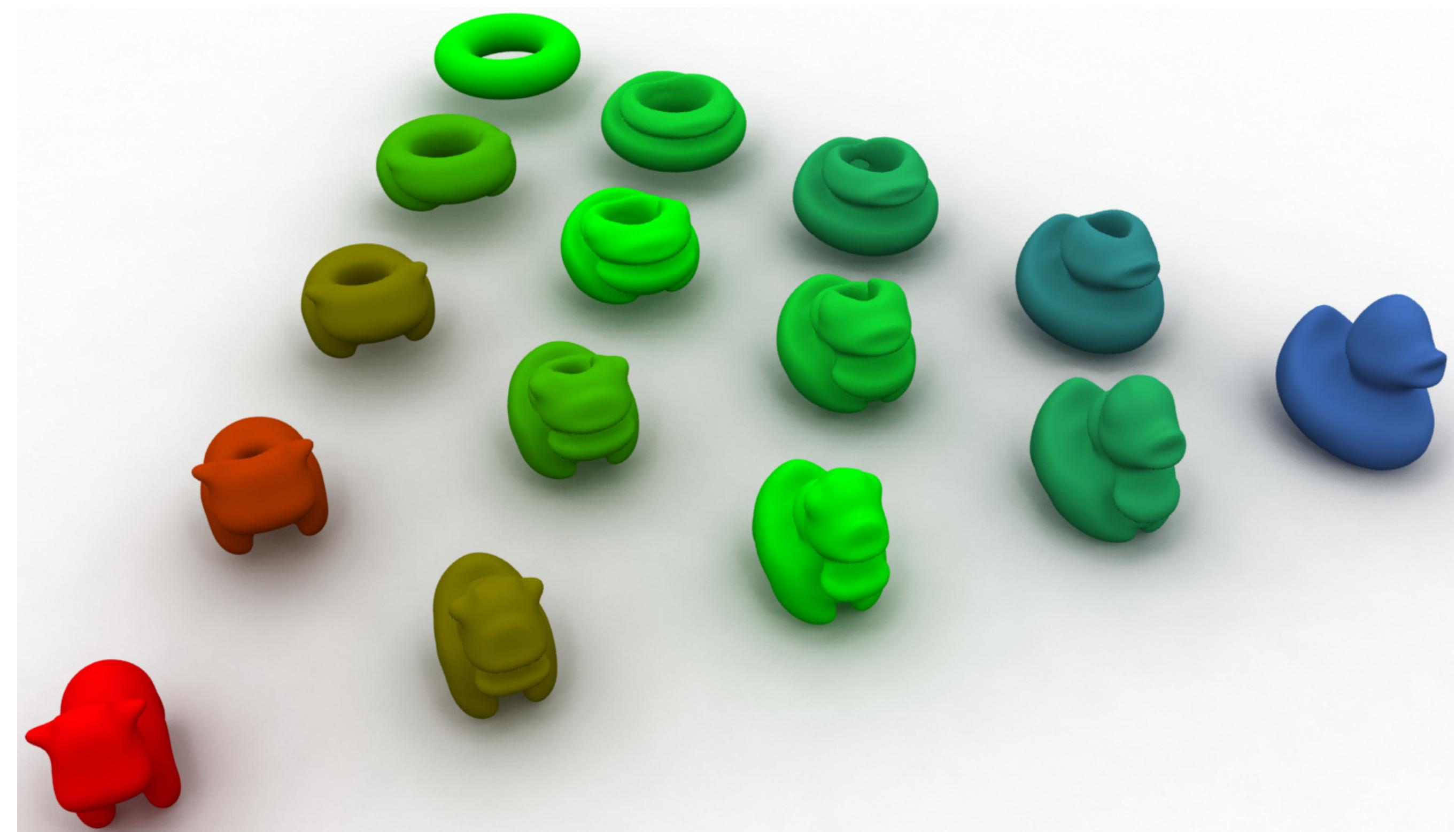
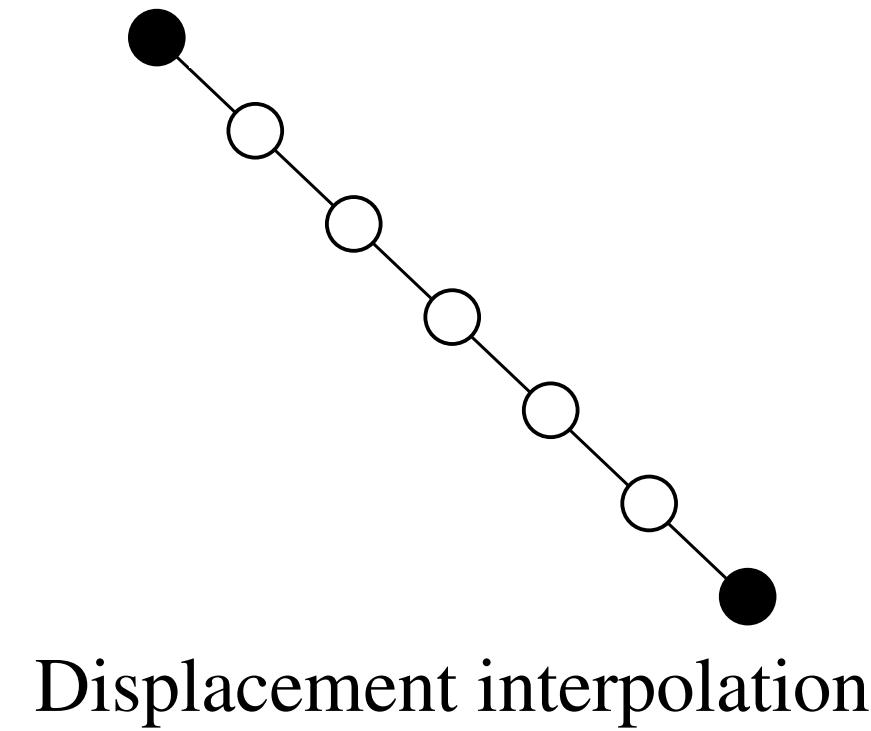
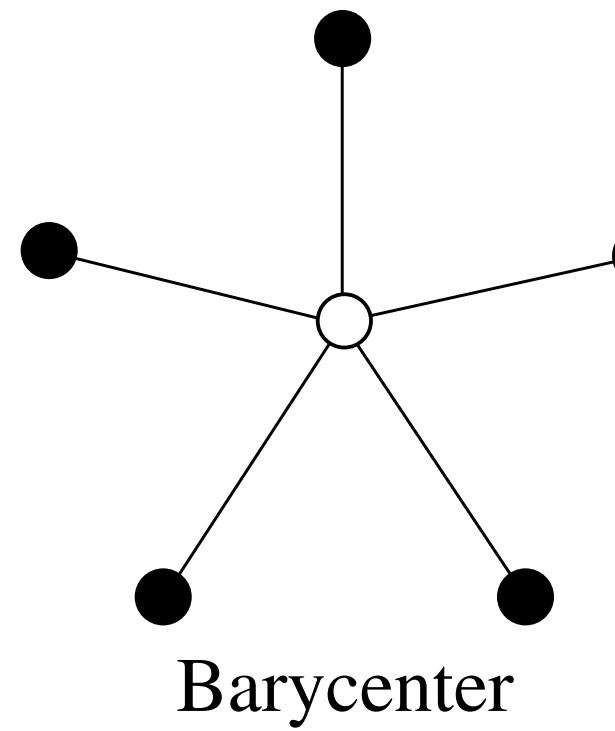


Wasserstein barycenter

Shape Analysis: Barycenters & Interpolation

Weighted Fréchet mean:

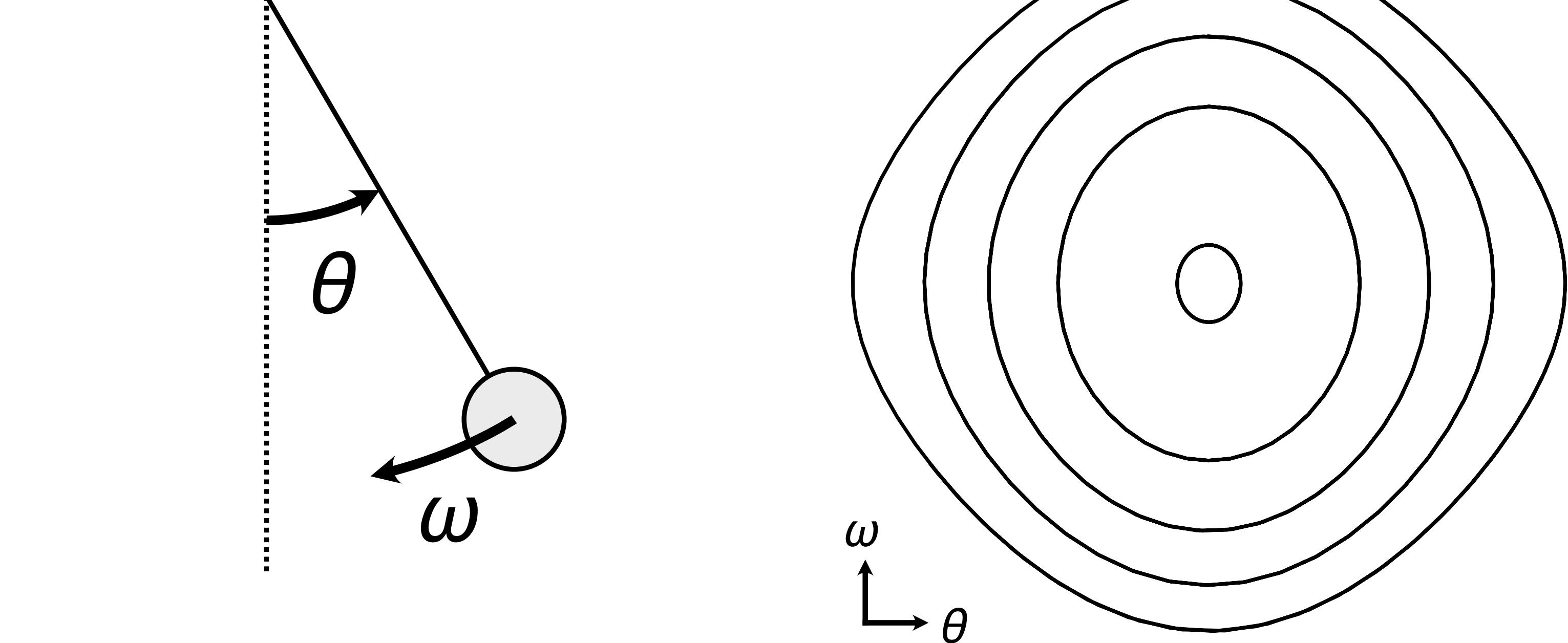
$$\min_{\mu} \sum_i \alpha_i W_2(\mu, \nu_i)^2$$



Shape Analysis: Manifold Learning



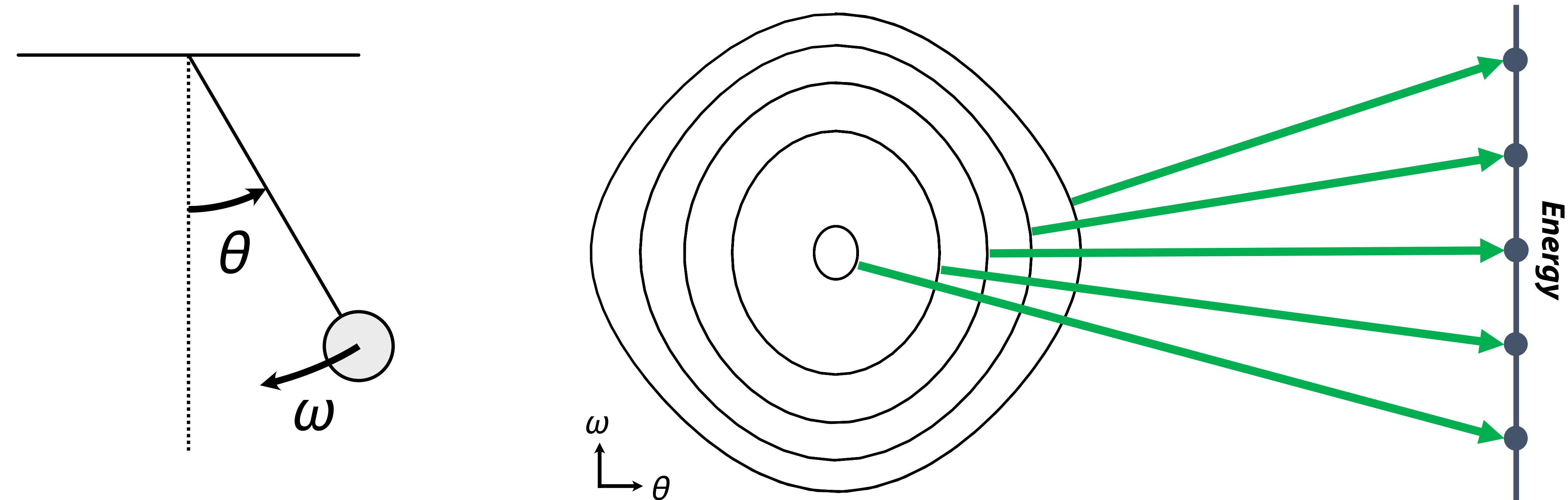
Discovering conservation laws: OT + diffusion maps



Shape Analysis: Manifold Learning



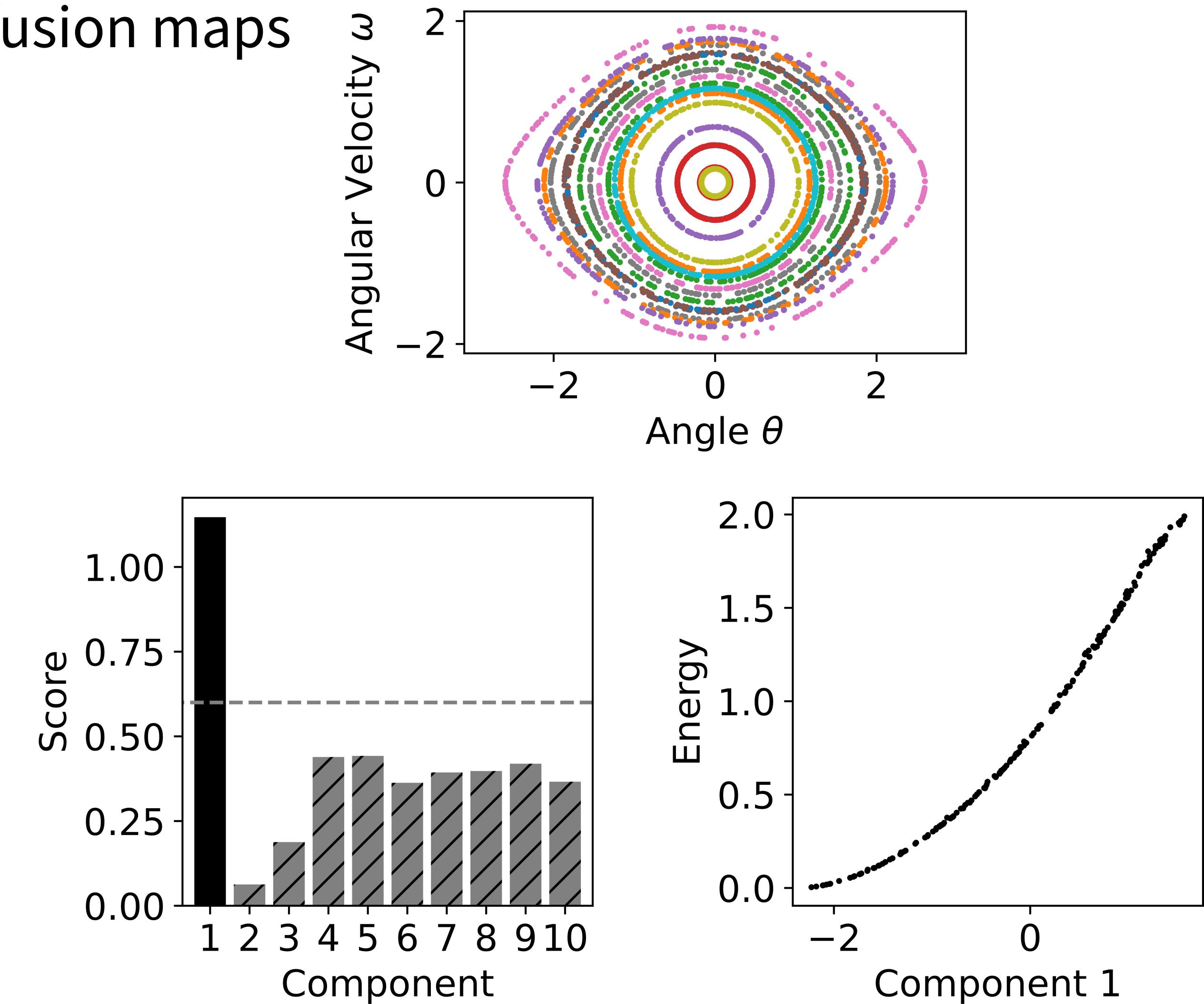
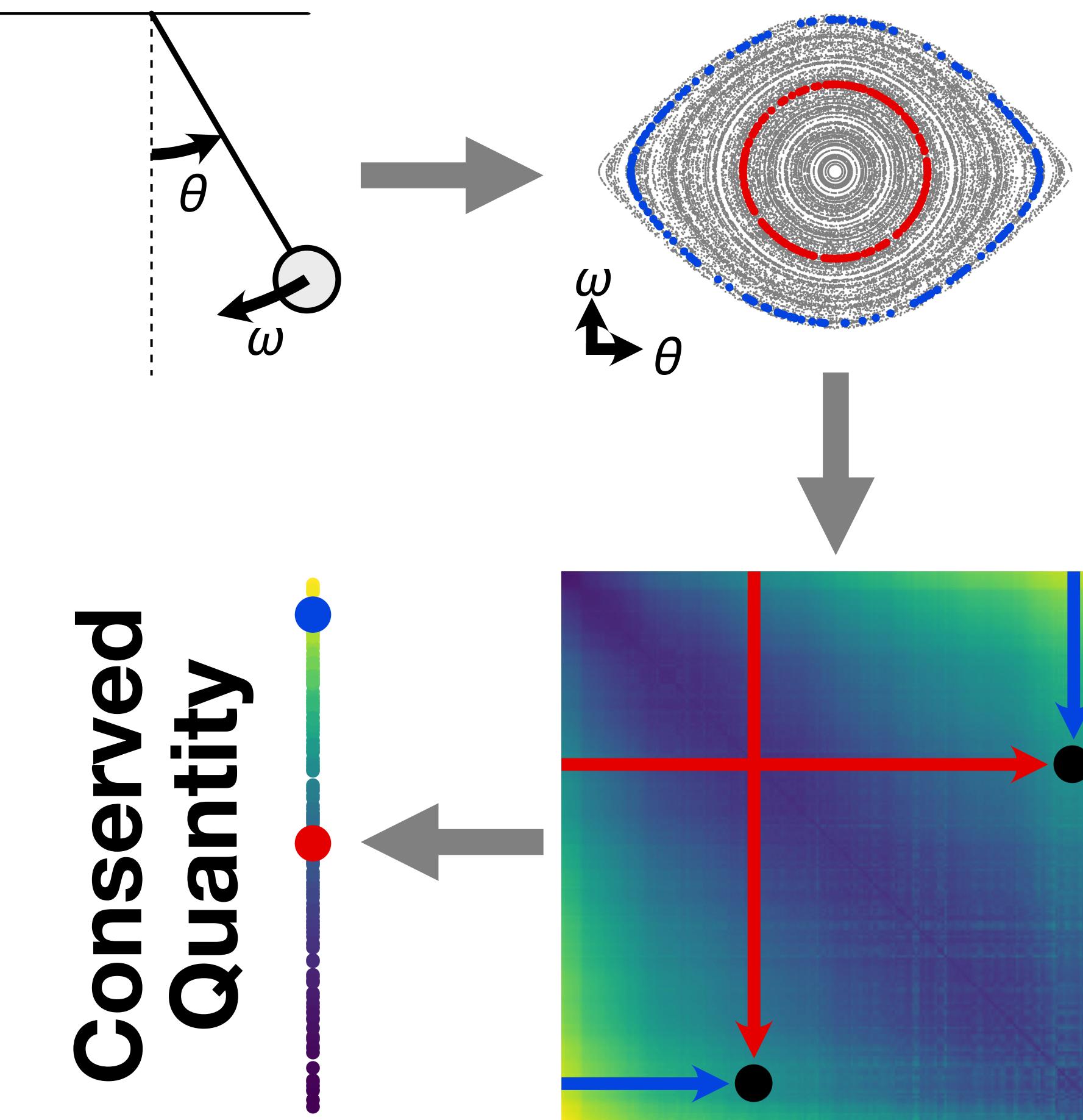
Discovering conservation laws: OT + diffusion maps



Shape Analysis: Manifold Learning

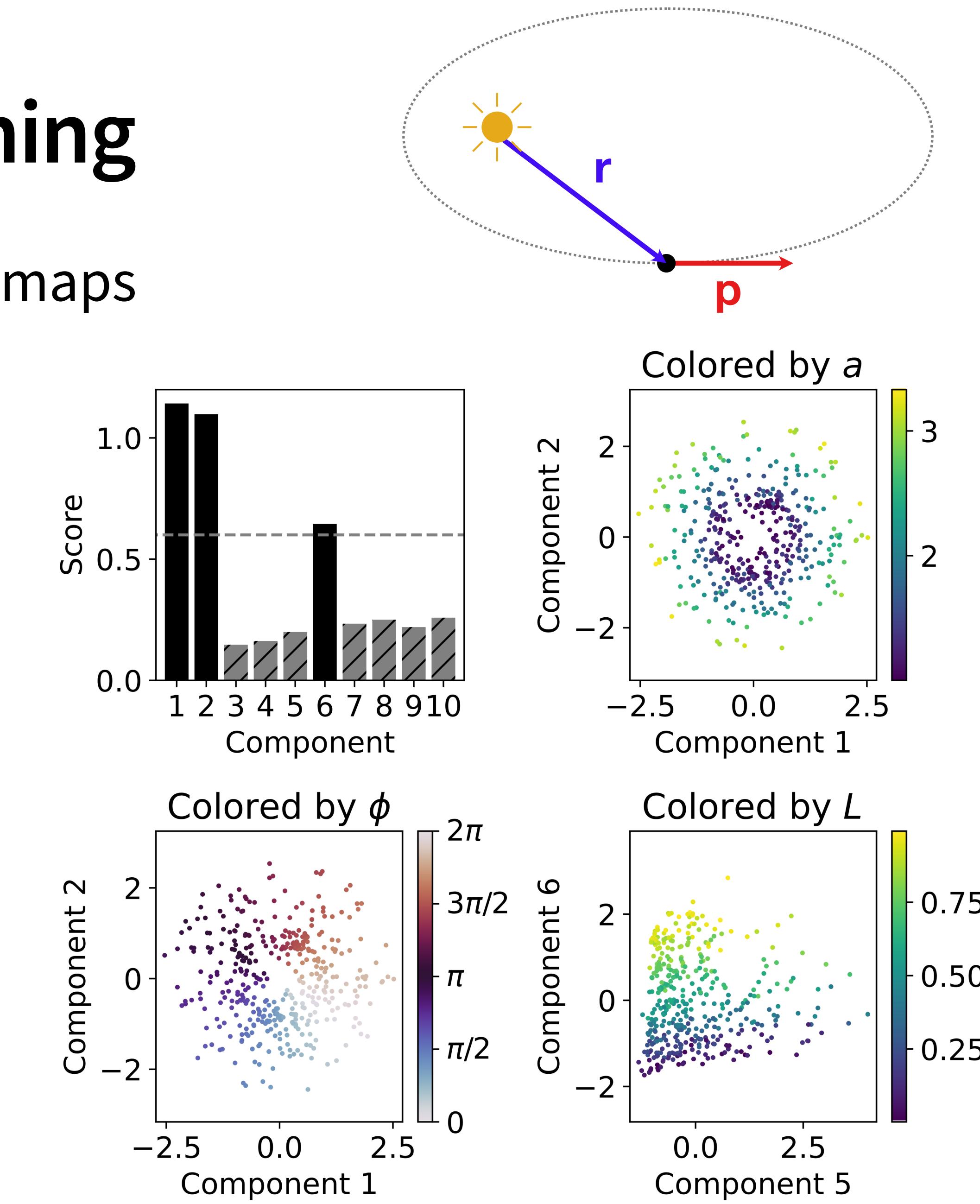
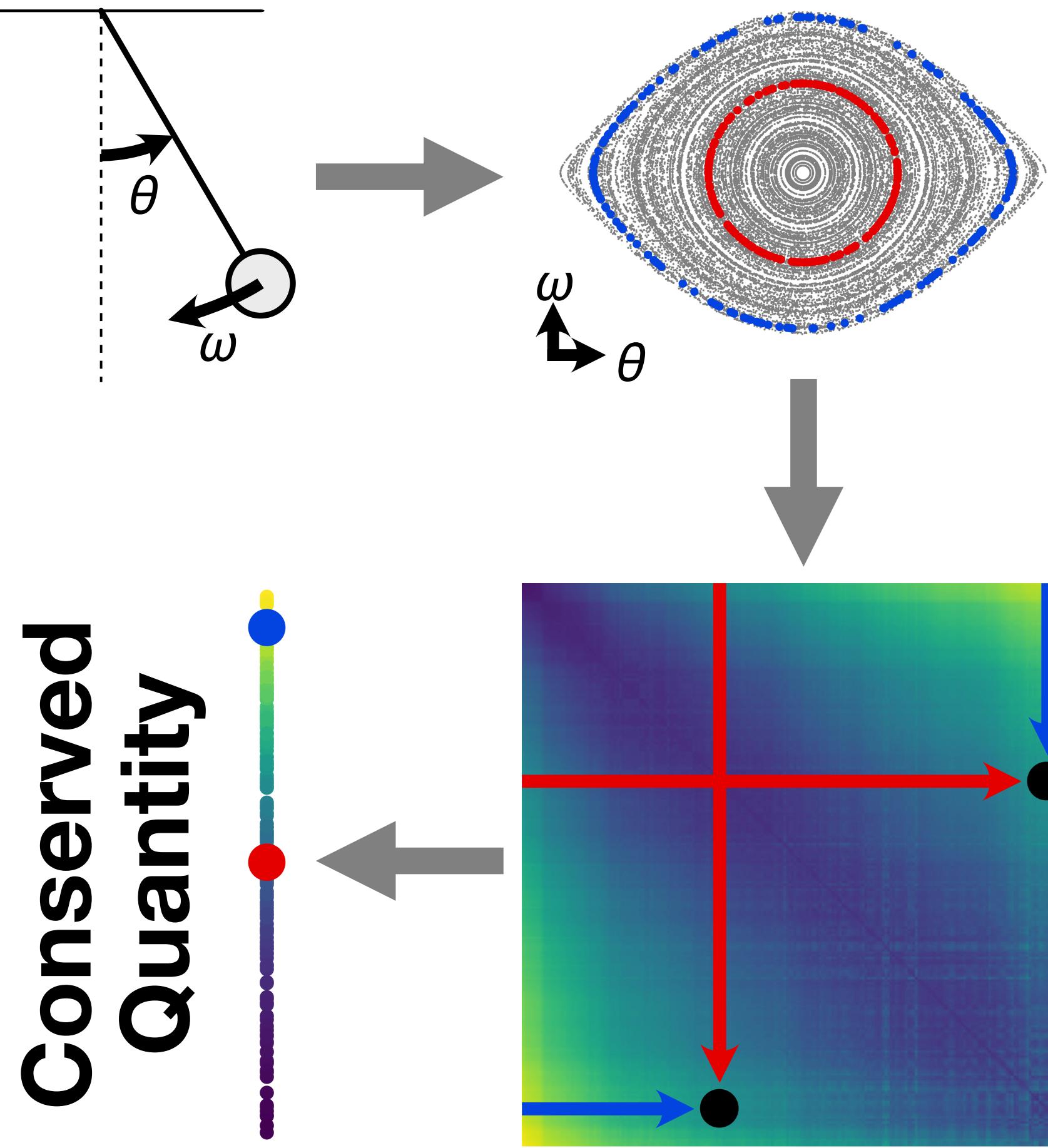


Discovering conservation laws: OT + diffusion maps



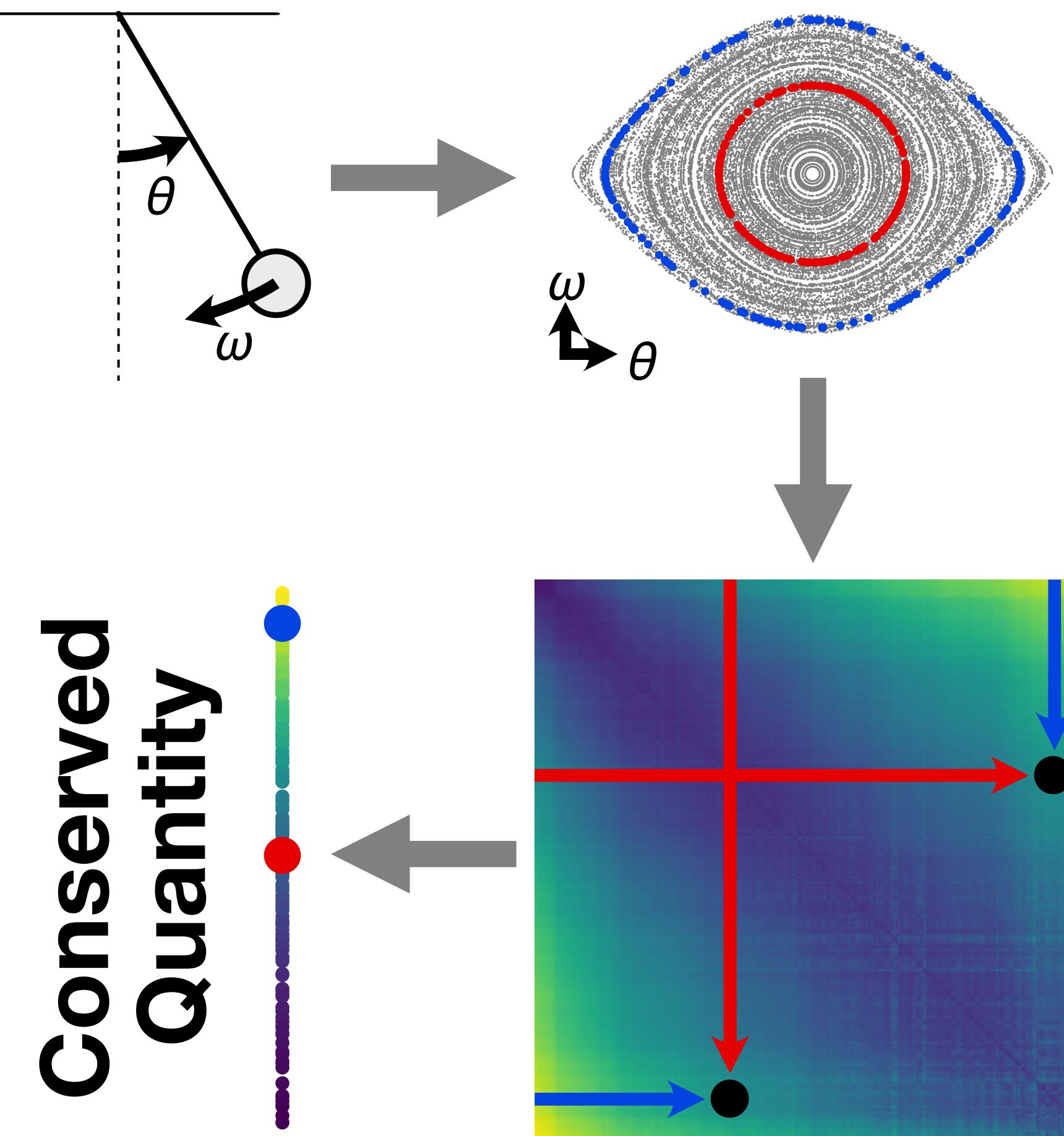
Shape Analysis: Manifold Learning

Discovering conservation laws: OT + diffusion maps

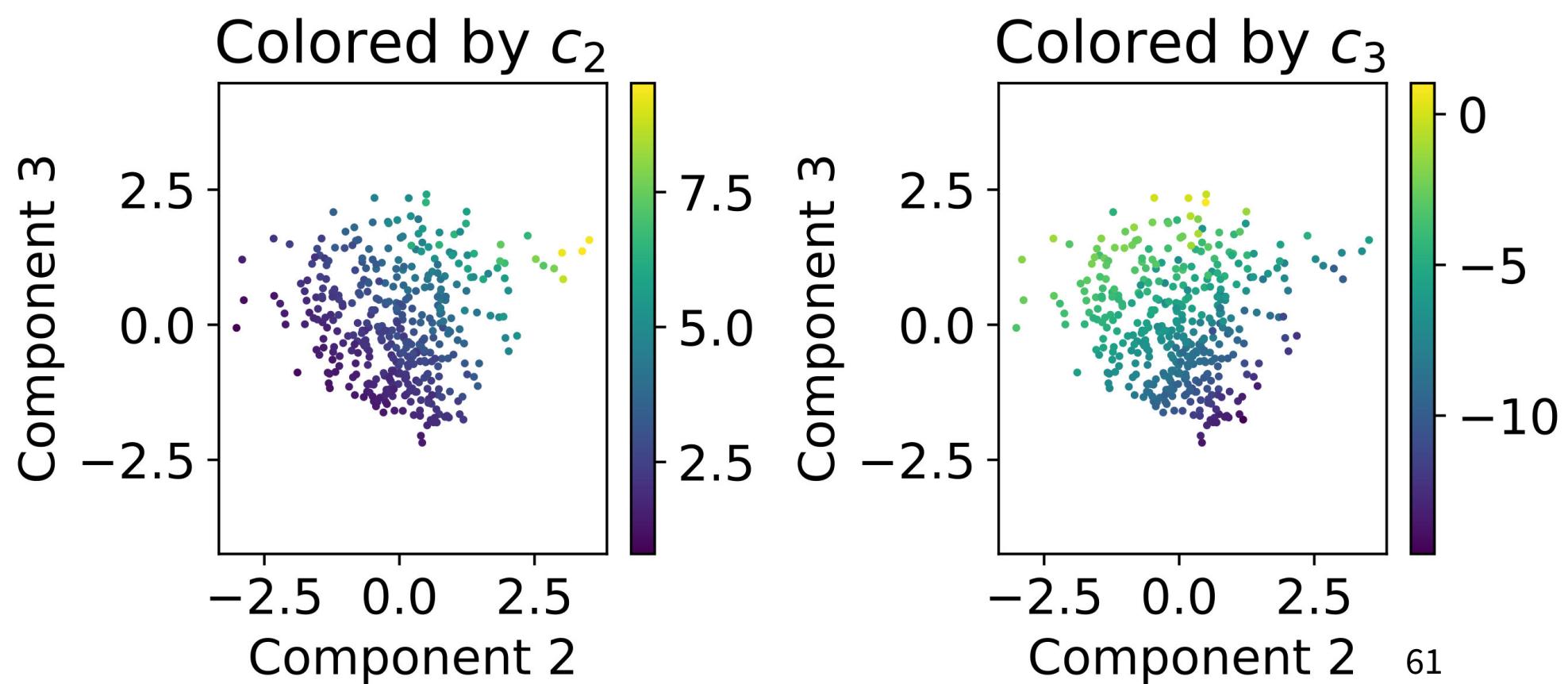
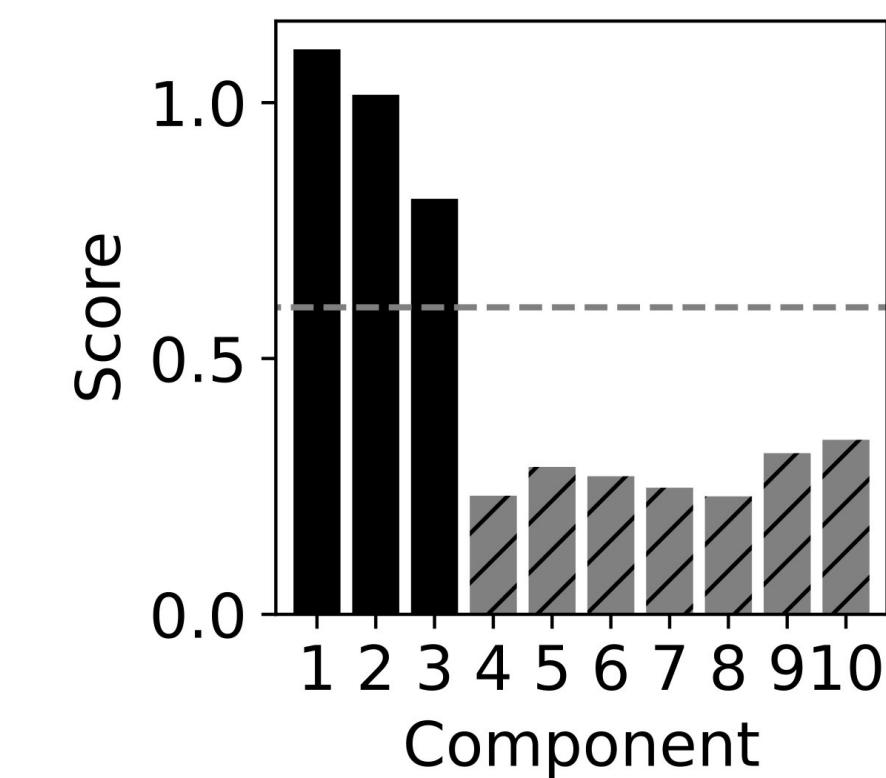
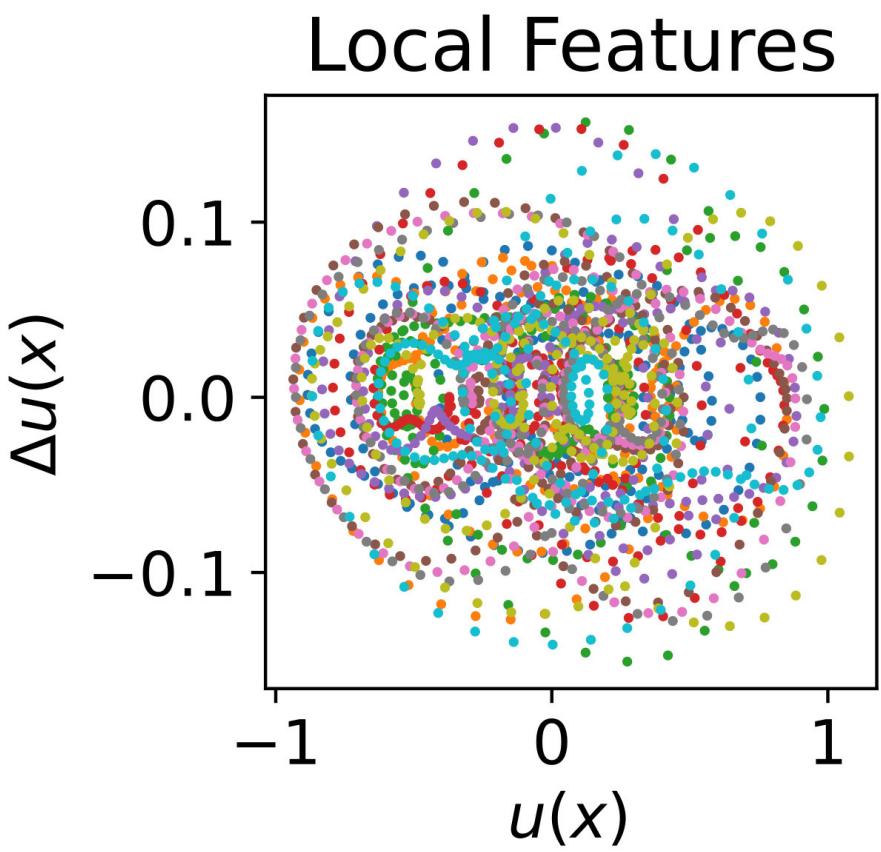
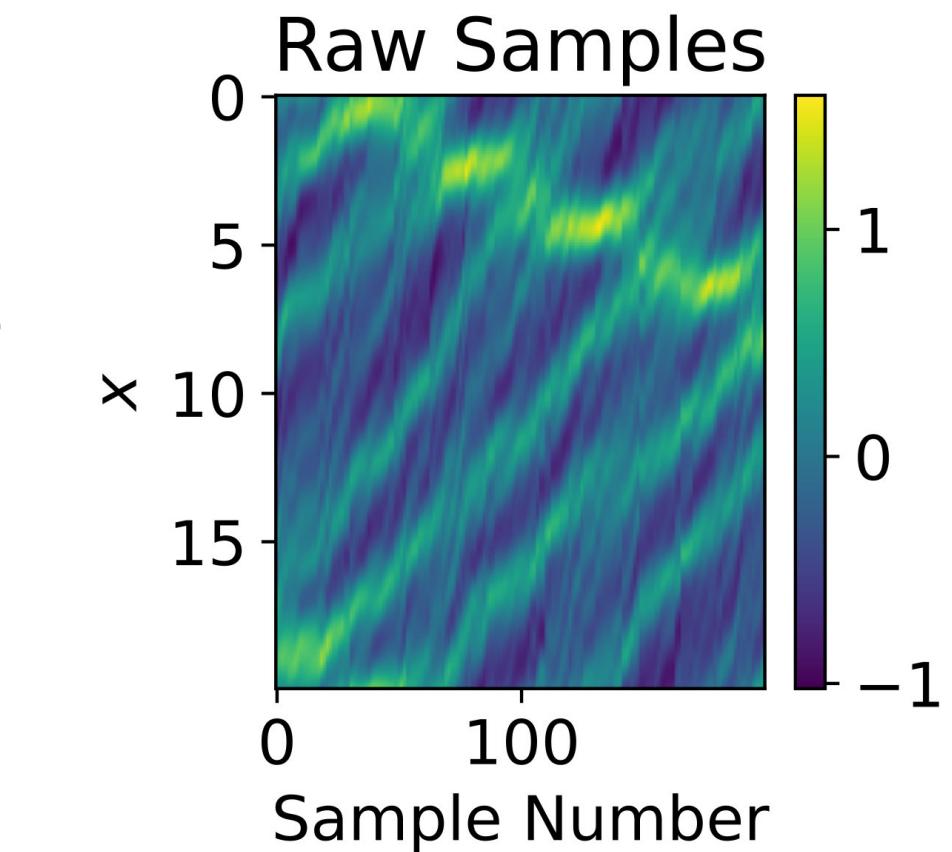


Shape Analysis: Manifold Learning

Discovering conservation laws: OT + diffusion maps

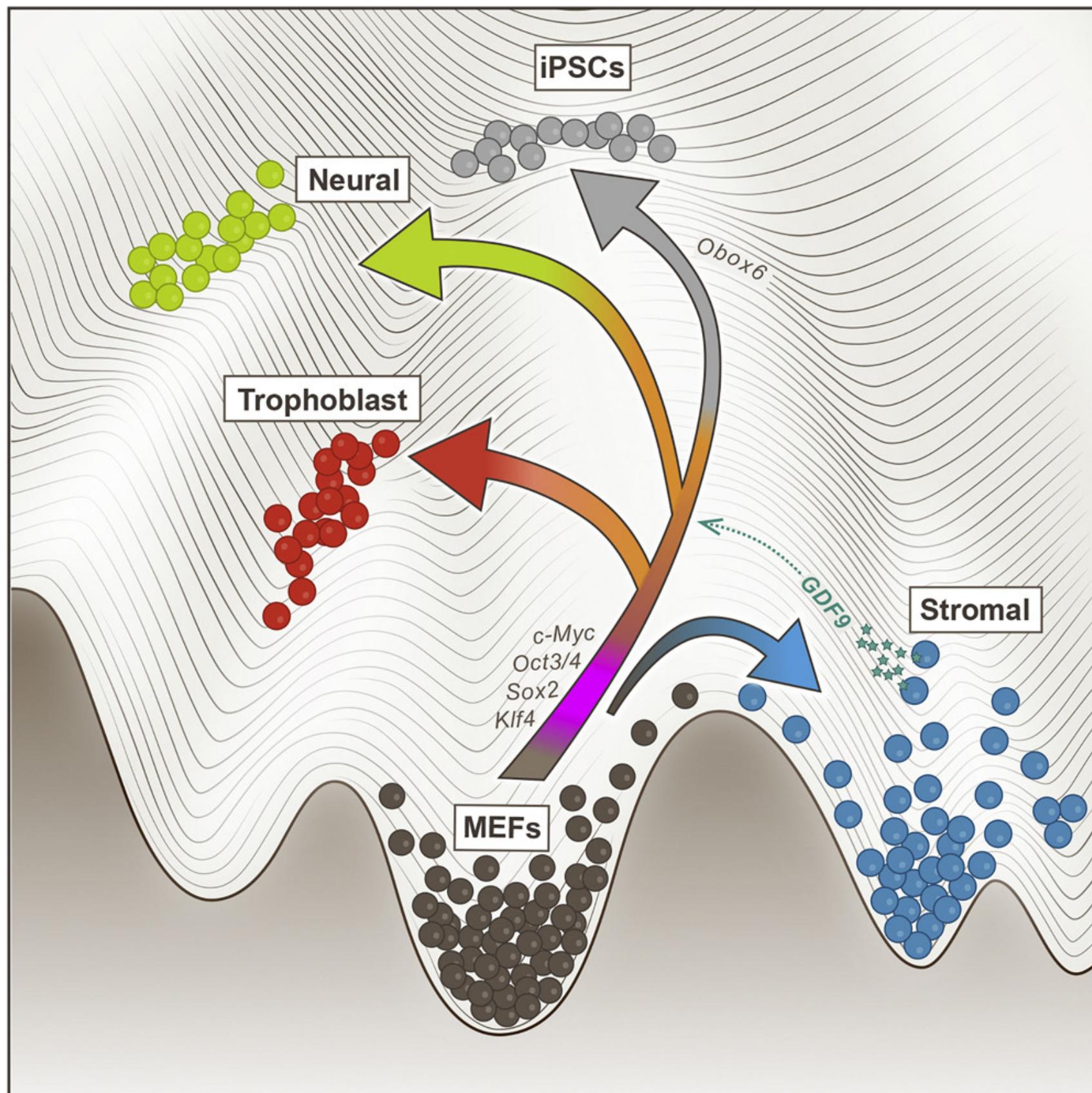


Lu et al., Nat Comm 2023



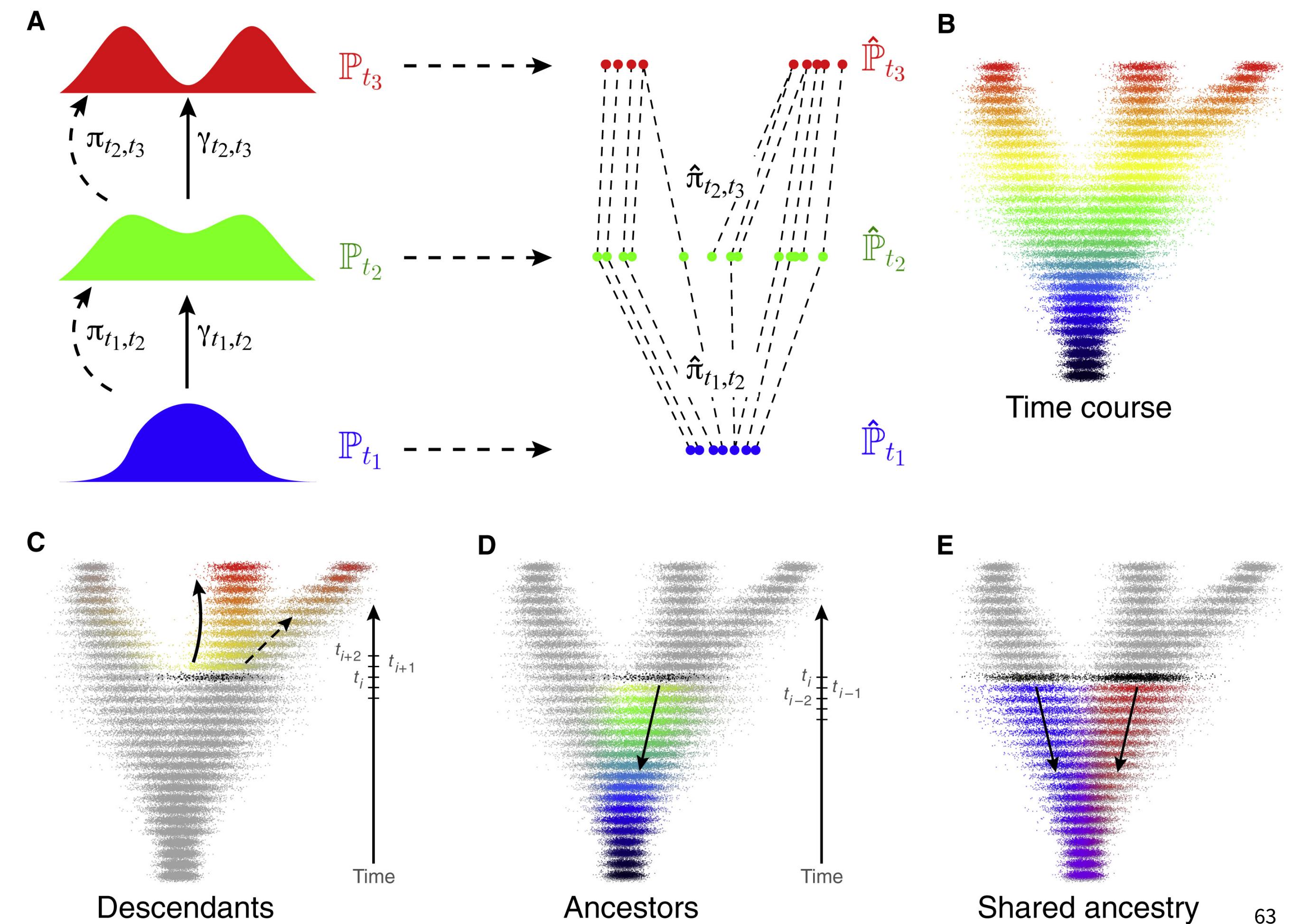
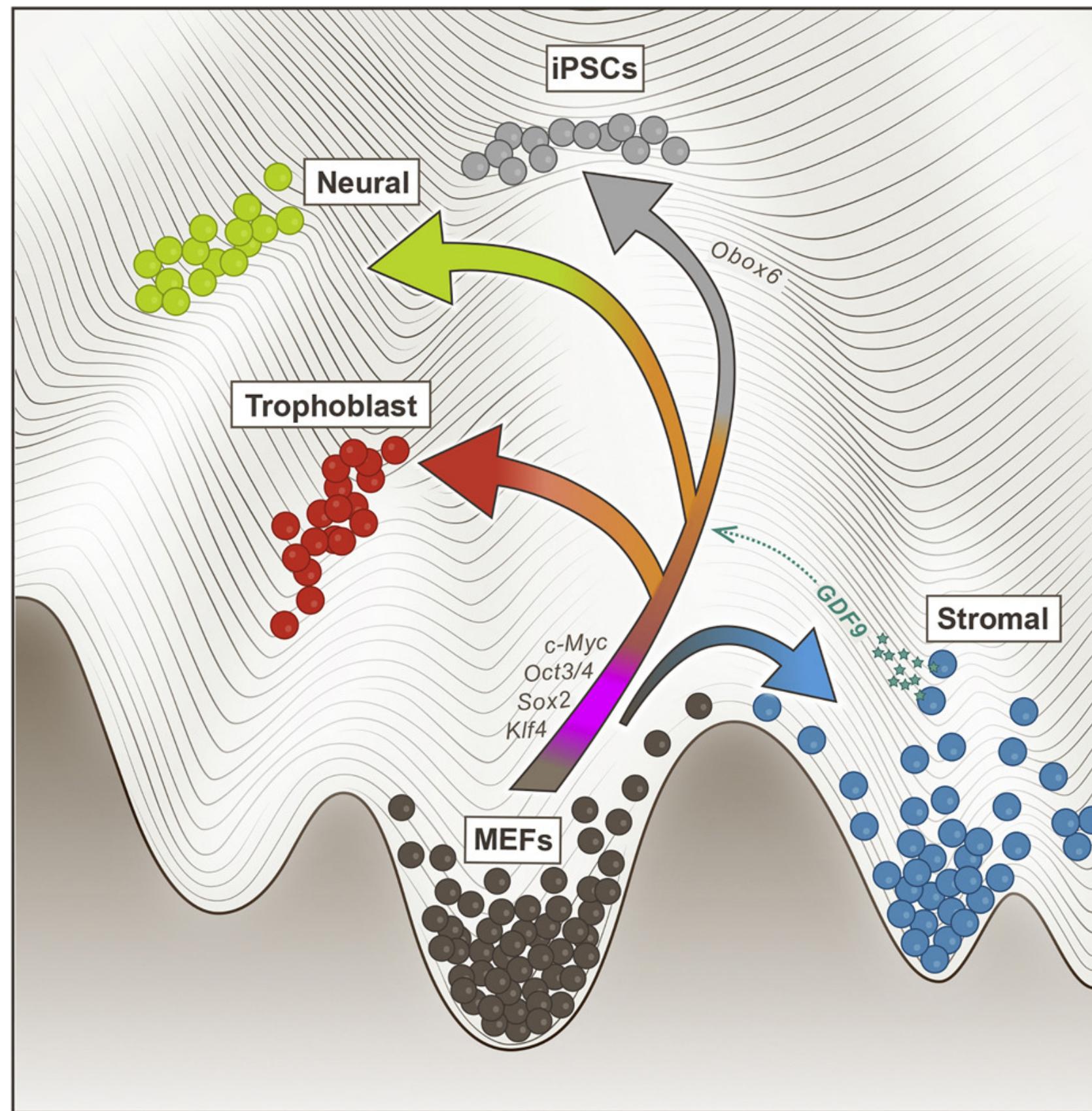
Shape Analysis: Manifold Learning

Developmental biology: inferring population dynamics



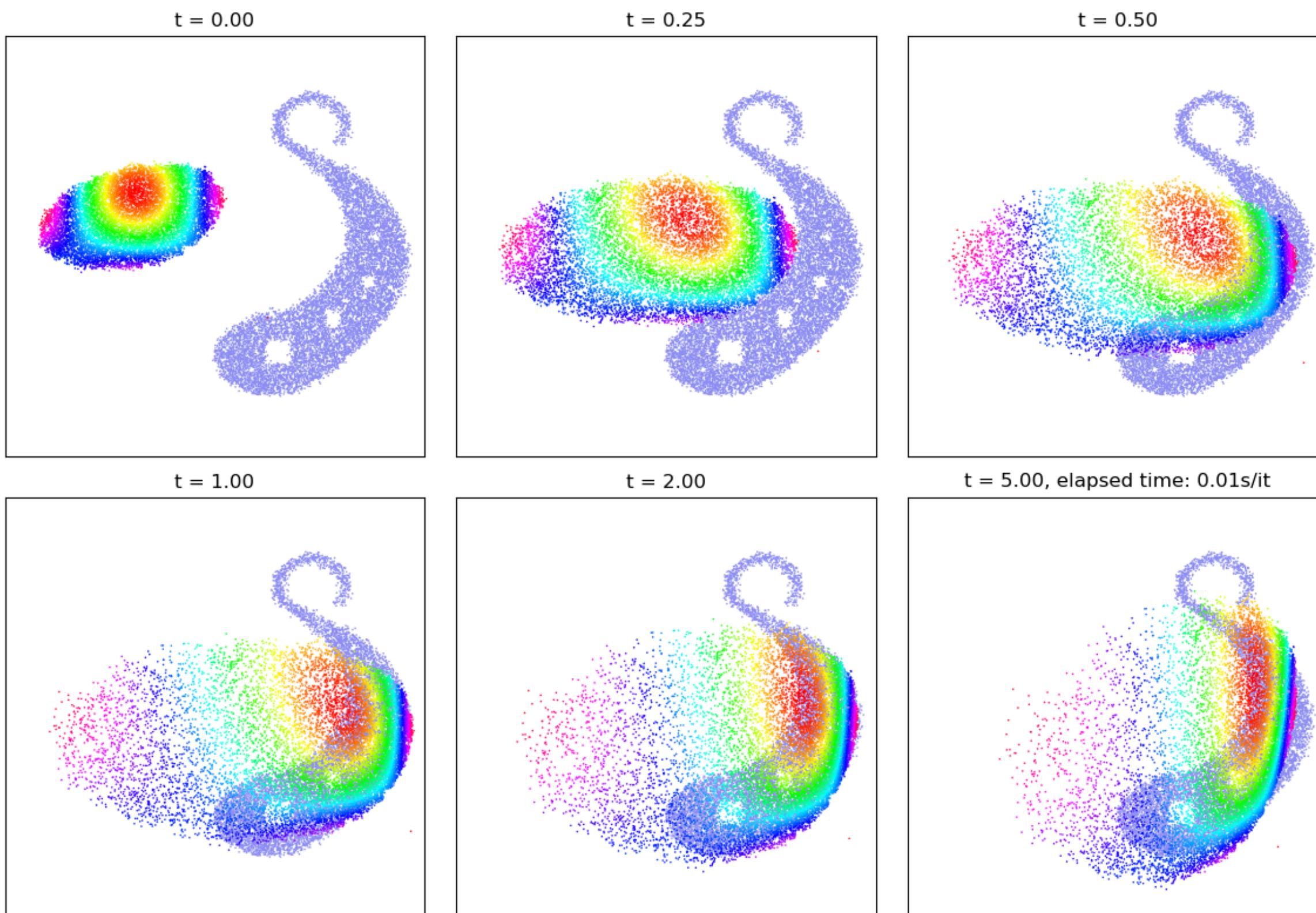
Shape Analysis: Manifold Learning

Developmental biology: inferring population dynamics

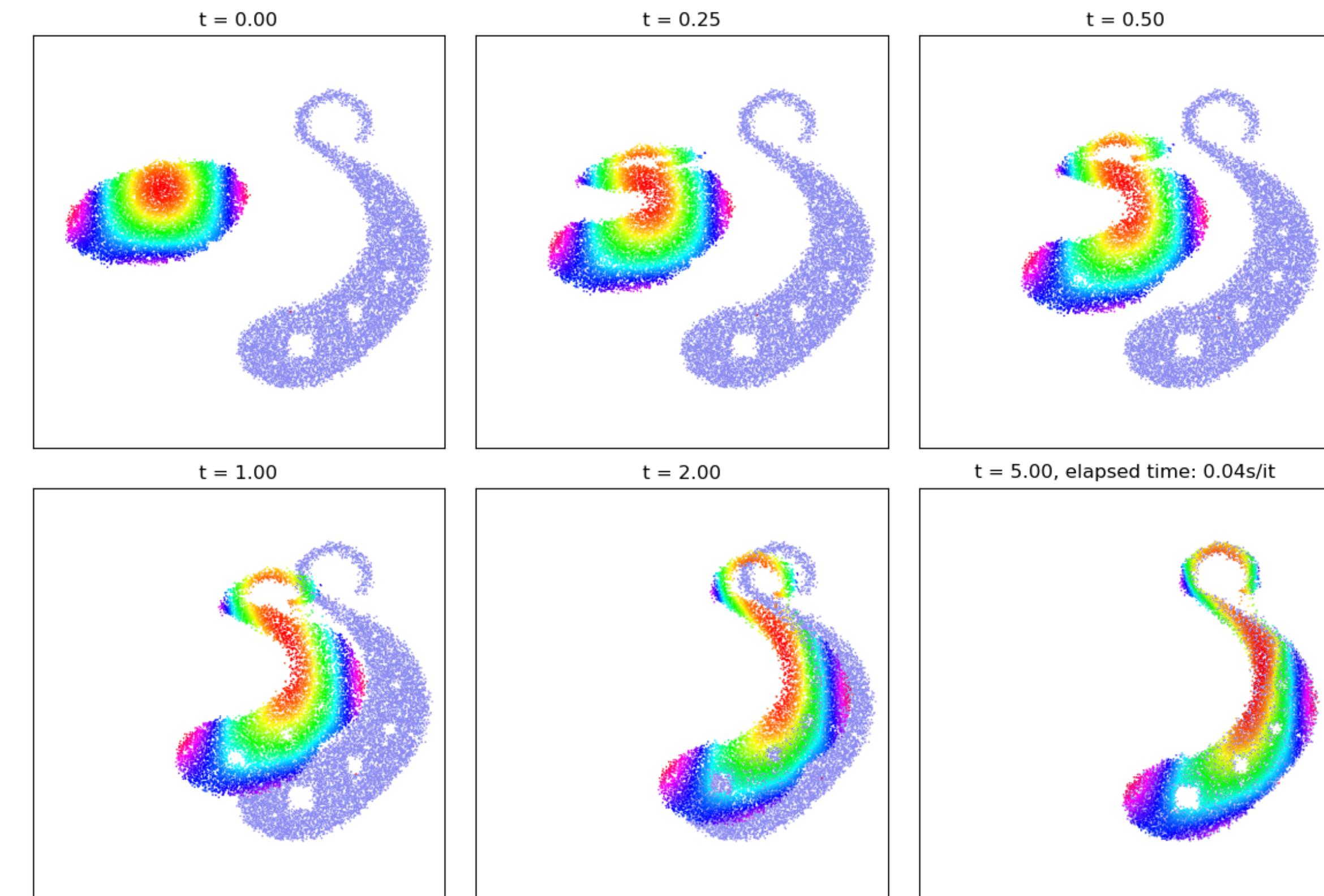


Deep Learning: Sinkhorn Loss

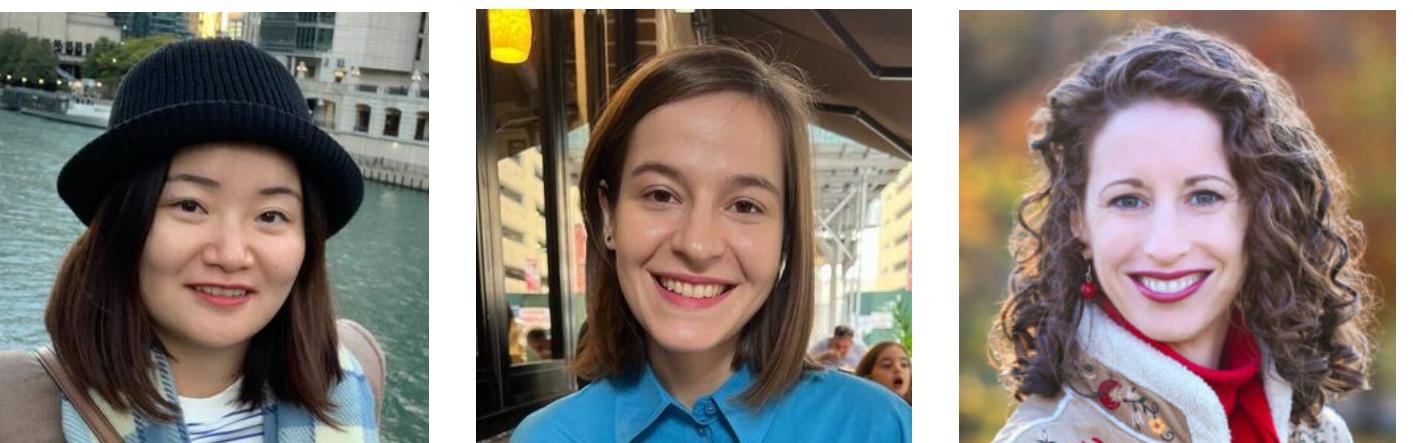
Gaussian MMD



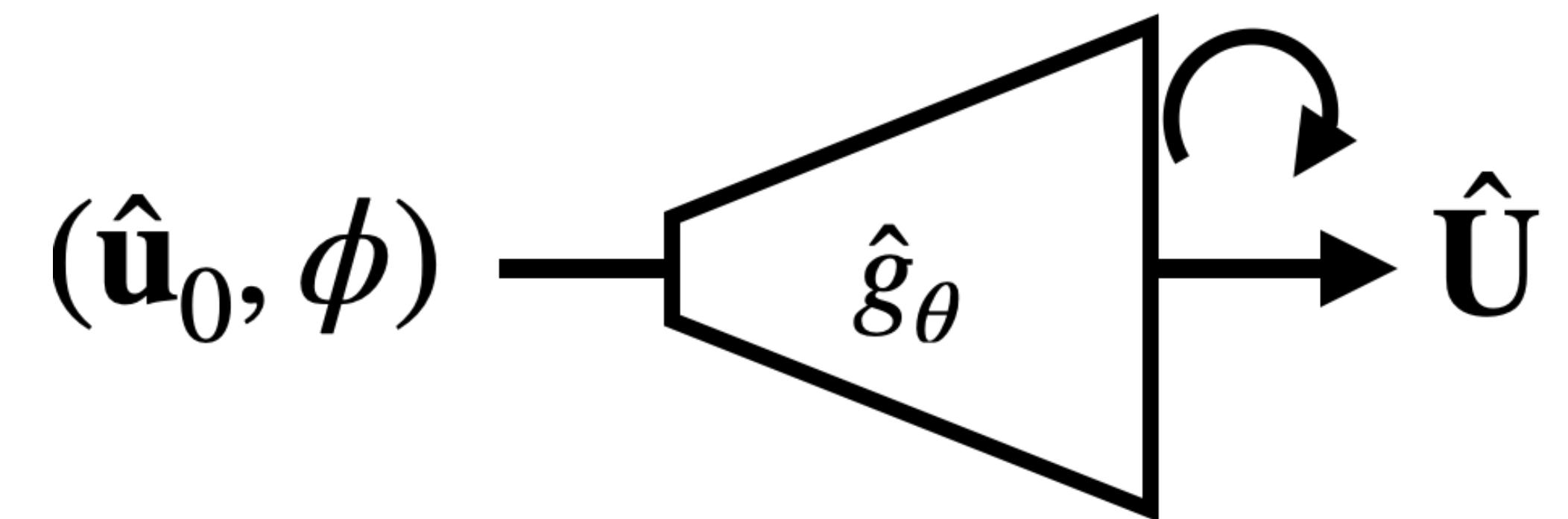
Sinkhorn



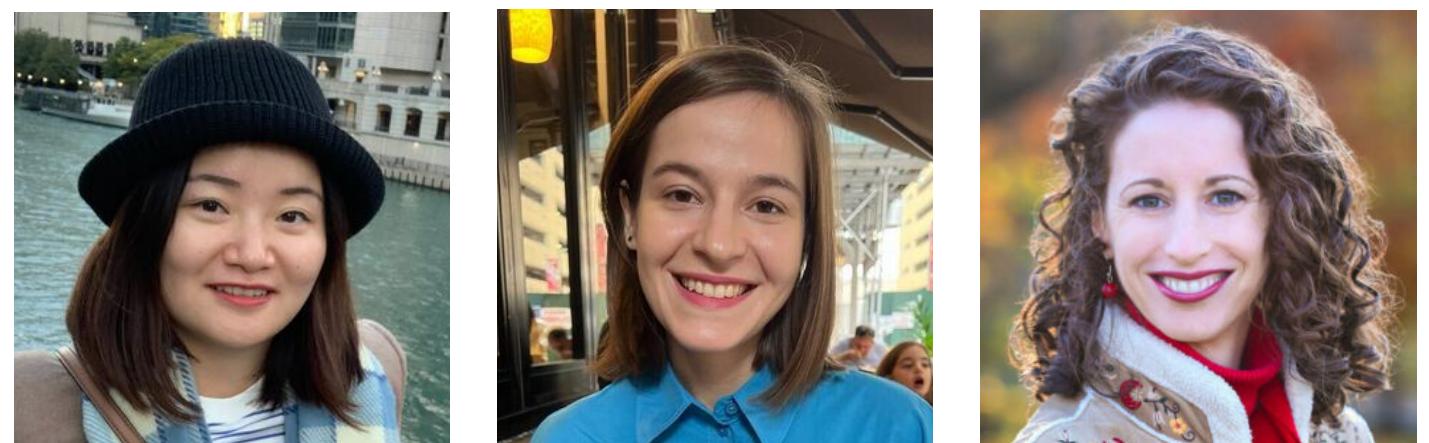
Deep Learning: Sinkhorn Loss



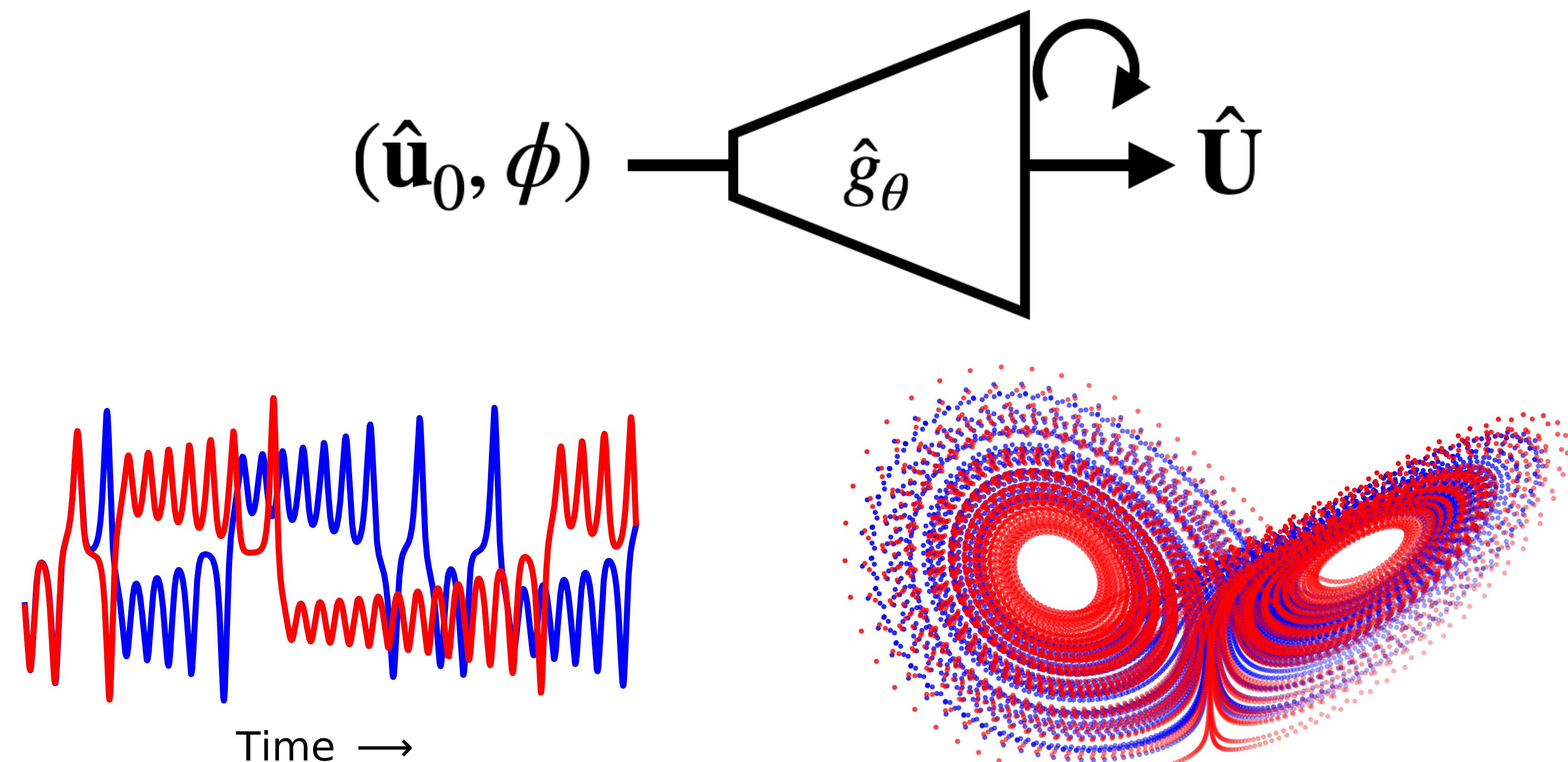
Training emulators for chaos: matching attractor statistics using OT



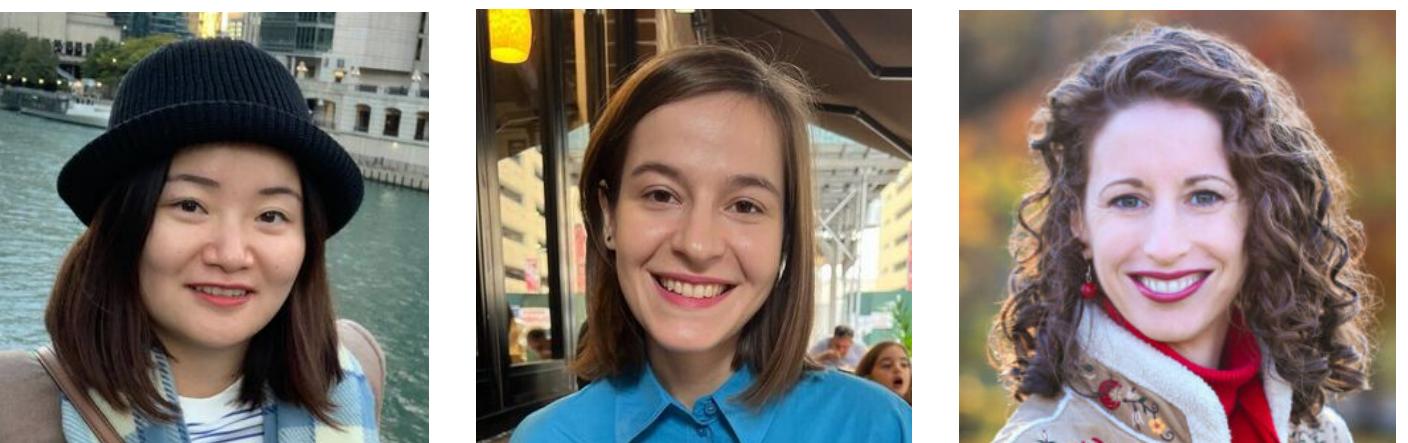
Deep Learning: Sinkhorn Loss



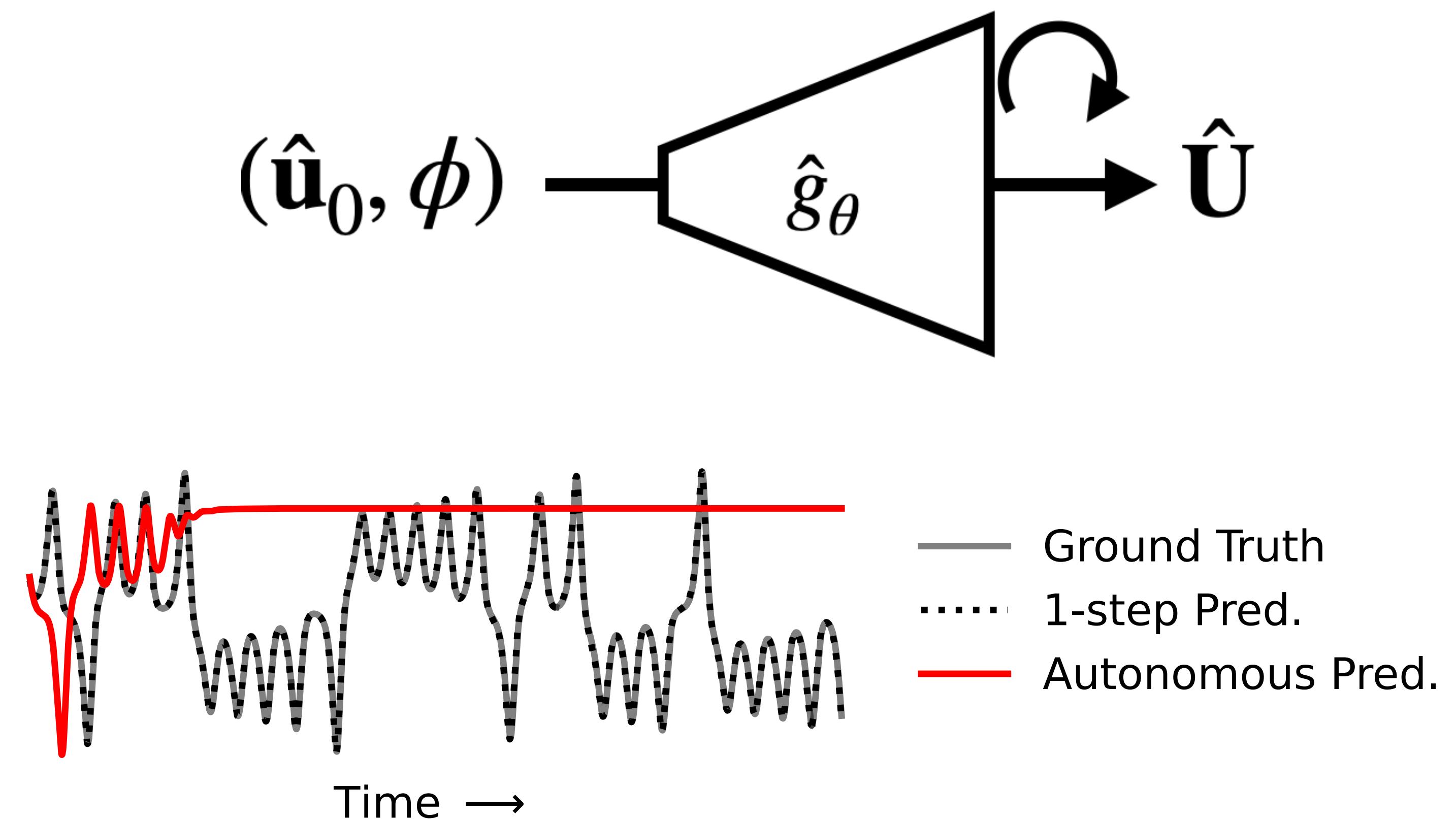
Training emulators for chaos: matching attractor statistics using OT



Deep Learning: Sinkhorn Loss



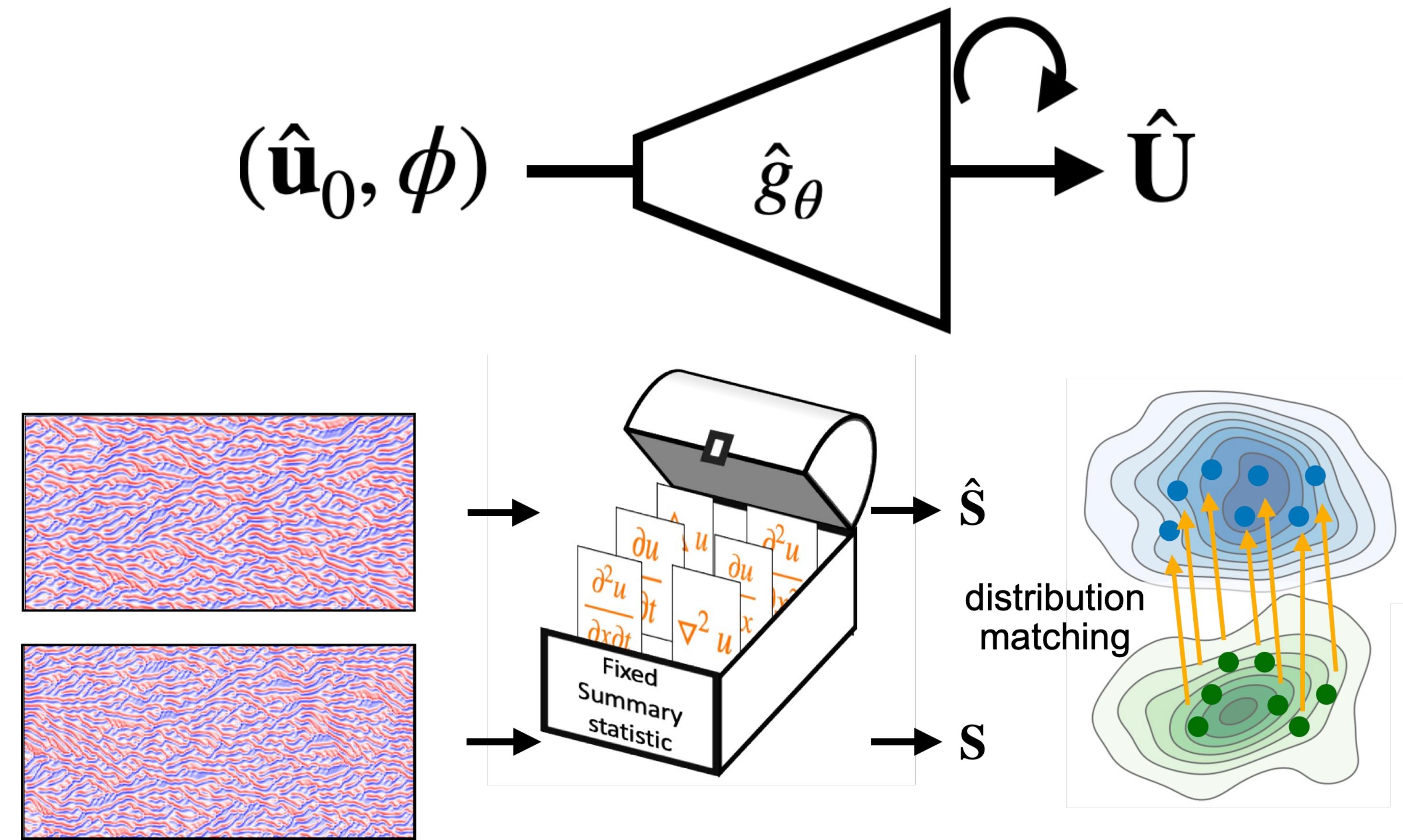
Training emulators for chaos: matching attractor statistics using OT



Deep Learning: Sinkhorn Loss

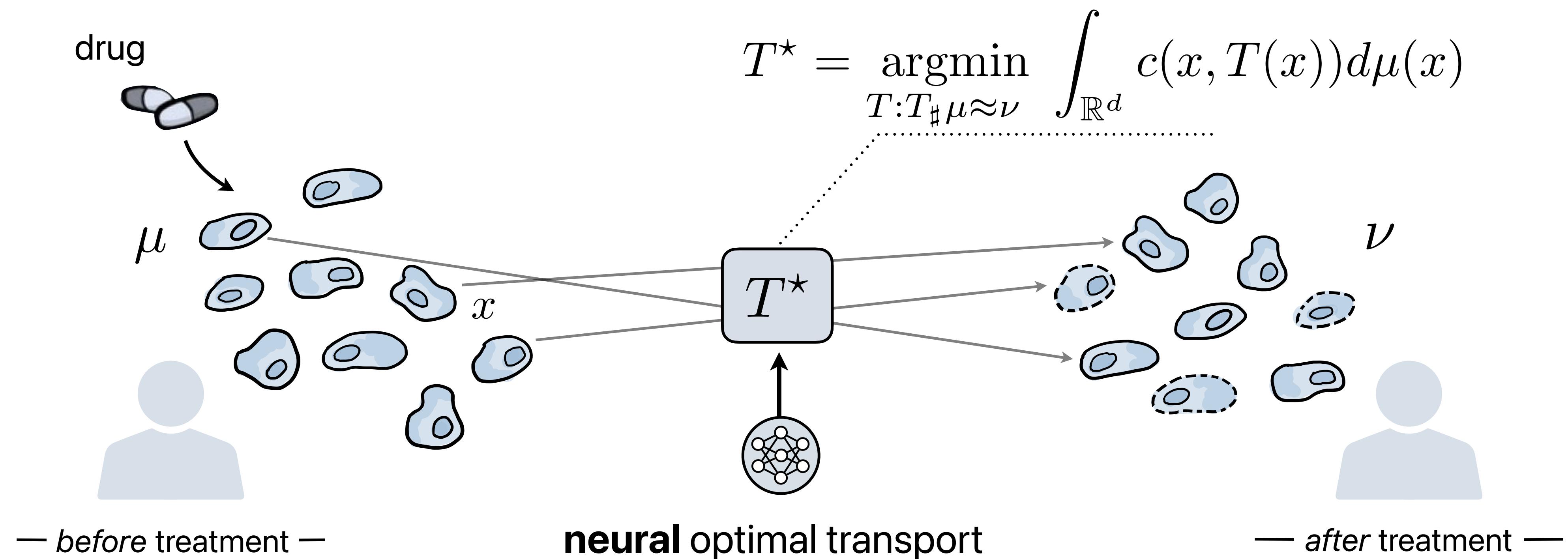


Training emulators for chaos: matching attractor statistics using OT

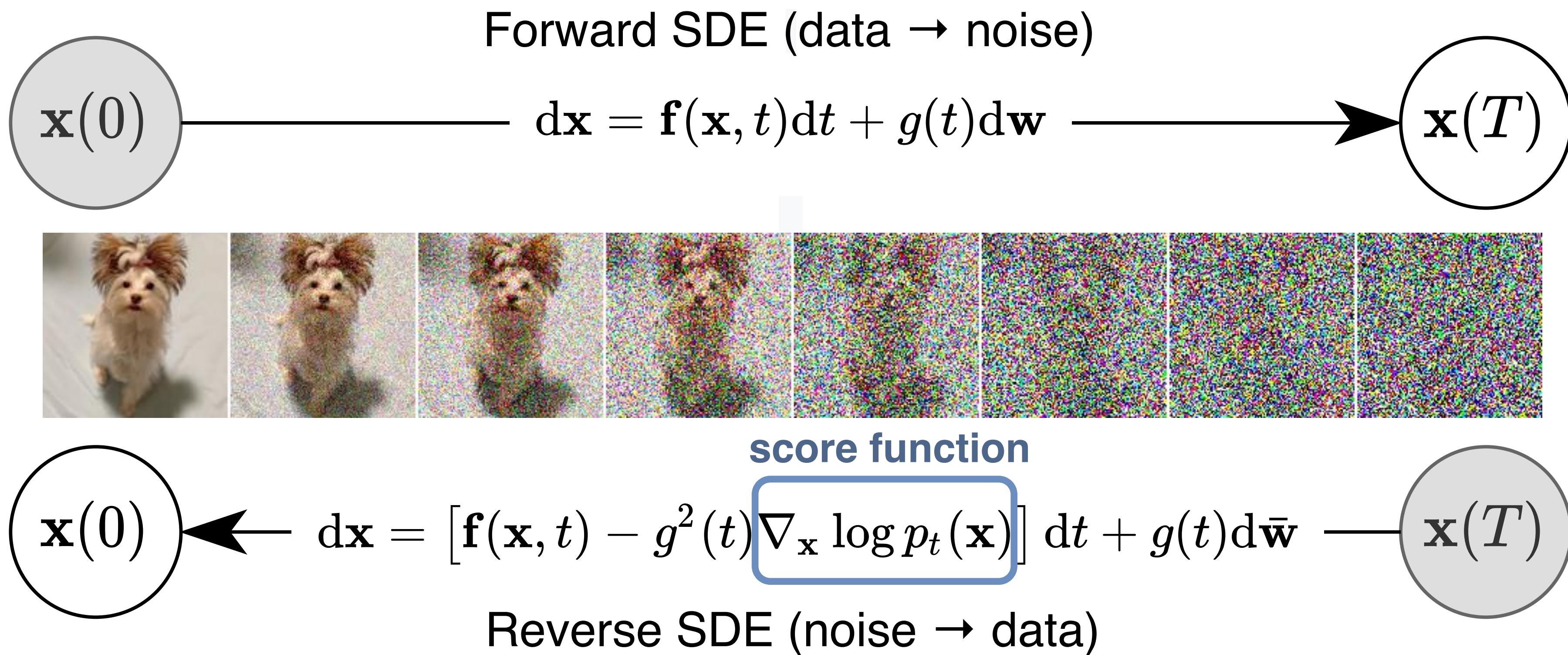


Deep Learning: Neural Optimal Transport

Personalized medicine: predict outcomes of cell populations



Deep Learning: Generative Modeling

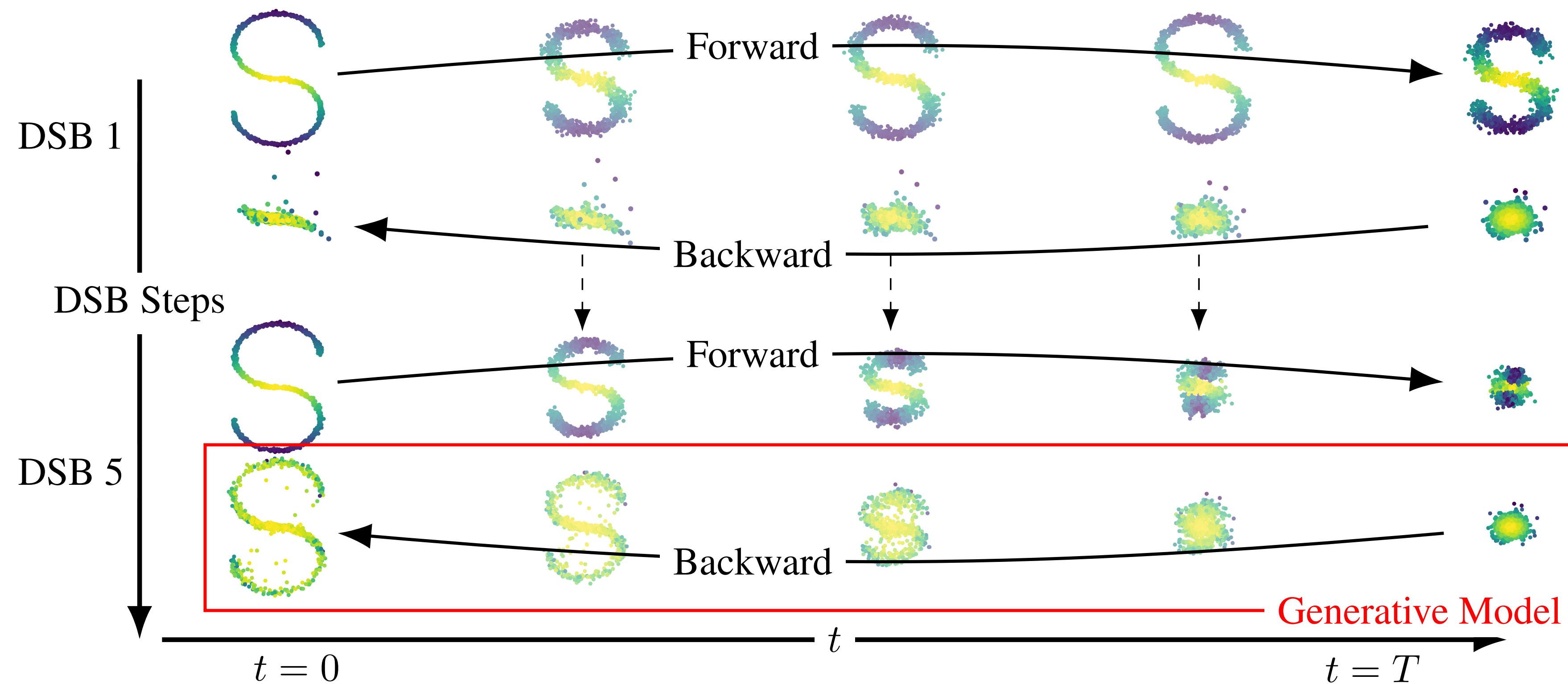


Deep Learning: Generative Modeling



Deep Learning: Generative Modeling

Diffusion Schrödinger bridge: speeding up diffusion models using IPF (Sinkhorn)



Summary: Optimal Transport for Machine Learning

Optimal Transport Theory

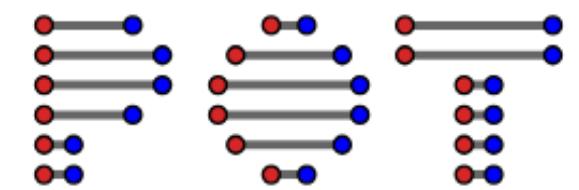
- Deep mathematical foundations
- Wasserstein metric on distributions

Computational Optimal Transport

- 1D OT → slice Wasserstein distance
- Entropy-Regularized OT → Sinkhorn algorithm

Machine Learning Applications

- Shape Analysis: barycenters, manifold learning, ...
- Deep Learning: Sinkhorn loss, generative modeling, ...



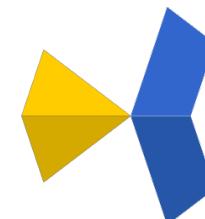
POT: Python Optimal Transport (*NumPy*)

pythonot.github.io



OTT: Optimal Transport Tools (*JAX*)

ott-jax.readthedocs.io



GeomLoss: Geometric Loss (*PyTorch*)

kernel-operations.io/geomloss

Gabriel Peyré & Marco Cuturi

Computational Optimal Transport

optimaltransport.github.io