

# Scaling Optimal Transport to High Dimensional Learning

Gabriel Peyré



Joint works with:



Shun'ichi  
Amari



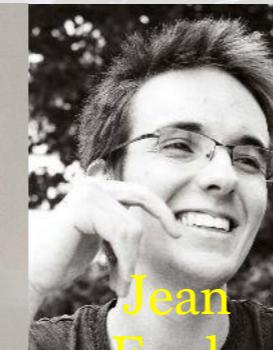
Francis  
Bach



Lénaïc  
Chizat



Marco  
Cuturi



Jean  
Feydy



Aude  
Genevay



Thibault  
Séjourné



Alain  
Trouvé



François-Xavier  
Vialard

<https://optimaltransport.github.io>

Home

# Computational Optimal Transport

BOOK

CODE

SLIDES

# Probability Distributions in Data Sciences

# *Probability distributions and histograms*

→ images, vision, graphics and machine learning,



# Probability Distributions in Data Sciences

*Probability distributions and histograms*

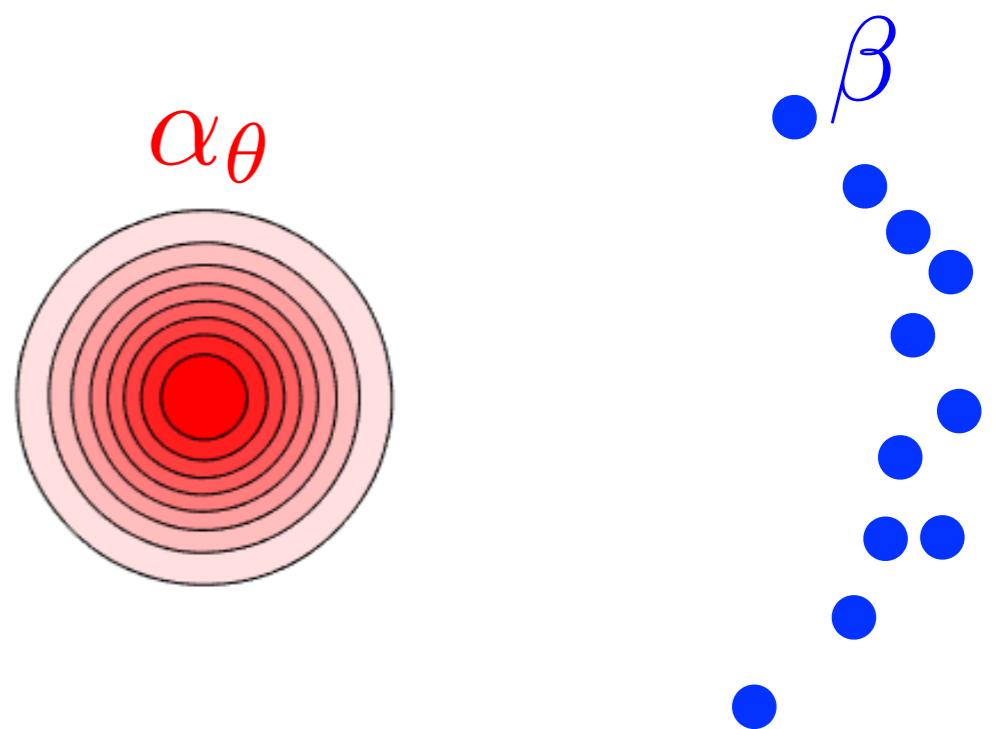
→ images, vision, graphics and machine learning,



## Unsupervised learning

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

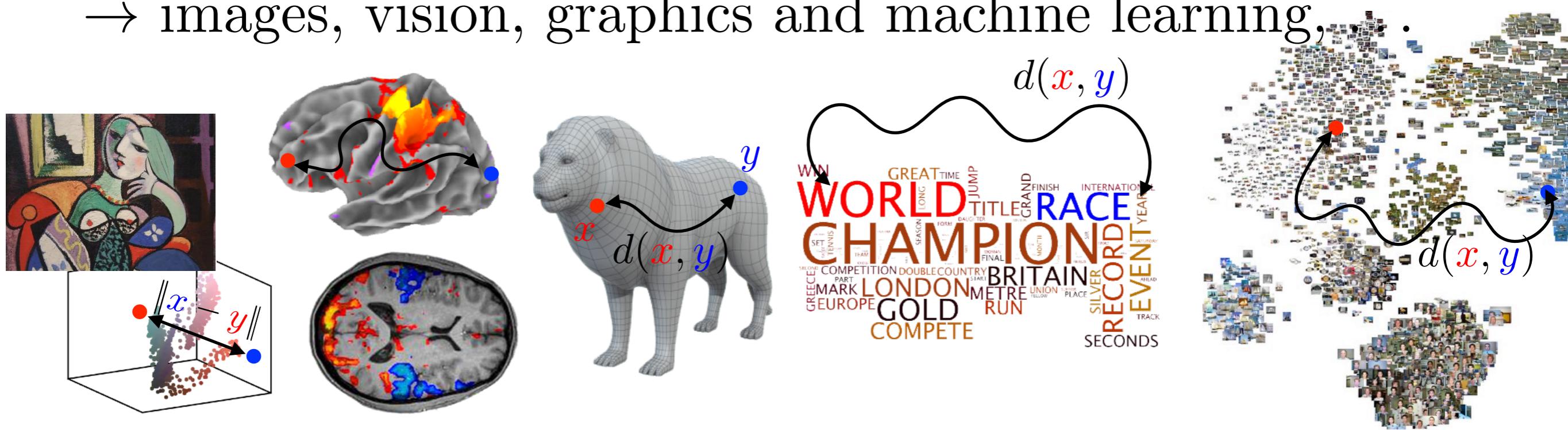
Parametric model:  $\theta \mapsto \alpha_\theta$



# Probability Distributions in Data Sciences

*Probability distributions and histograms*

→ images, vision, graphics and machine learning,

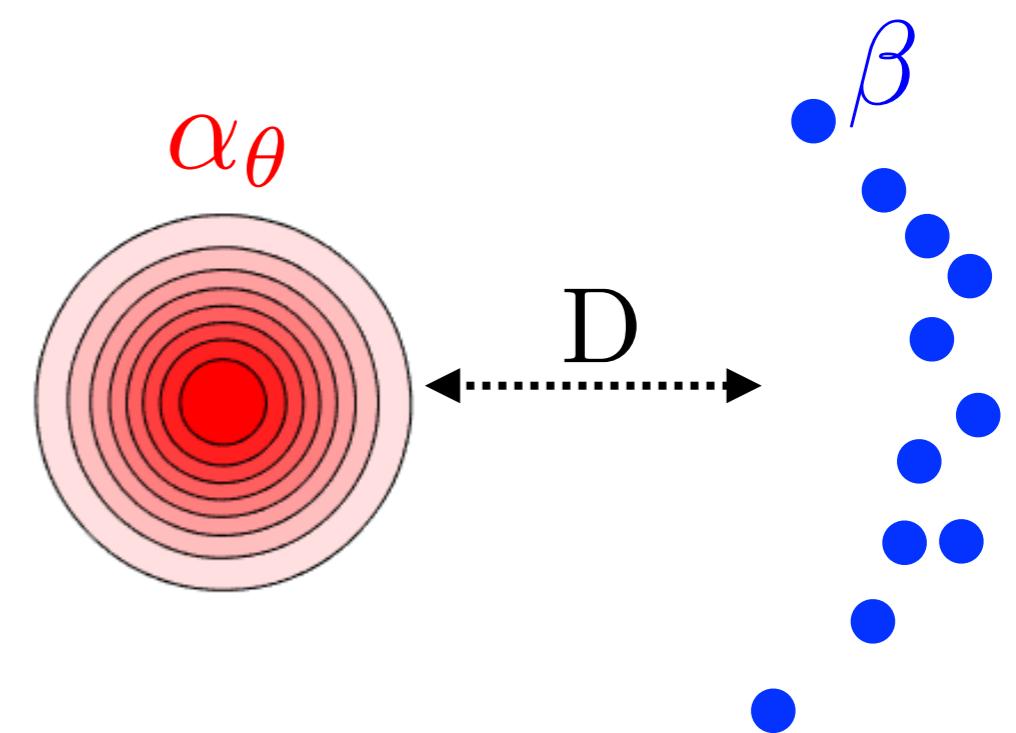


## Unsupervised learning

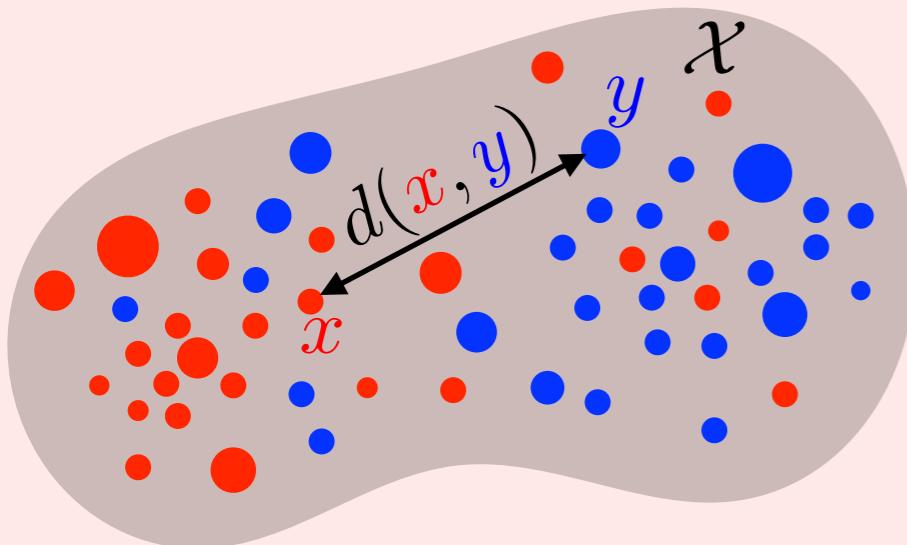
Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$

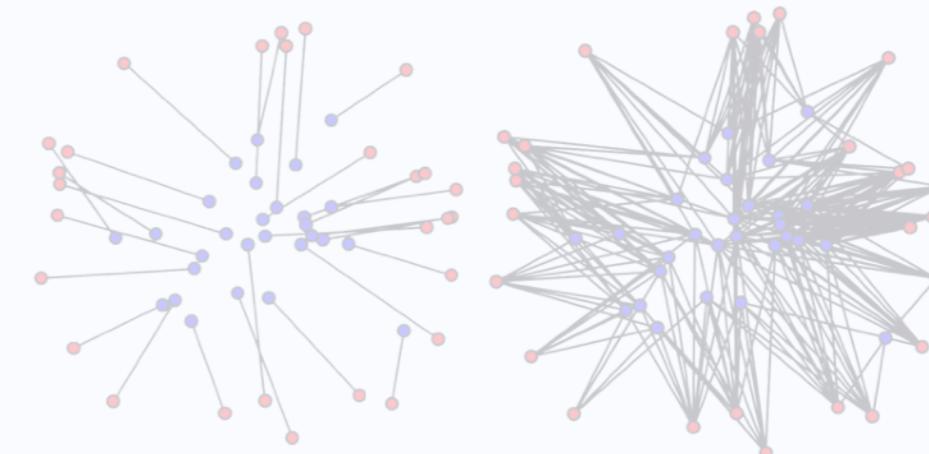
Density fitting:  $\min_{\theta} D(\alpha_\theta, \beta)$   
→ takes into account a metric  $d$ .



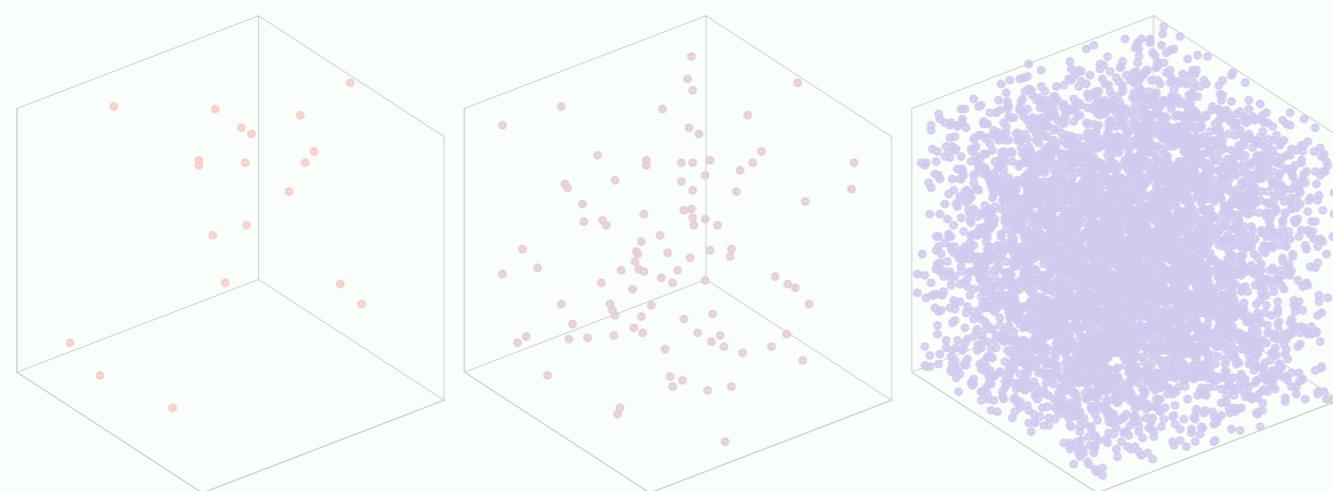
# 1. Optimal Transport



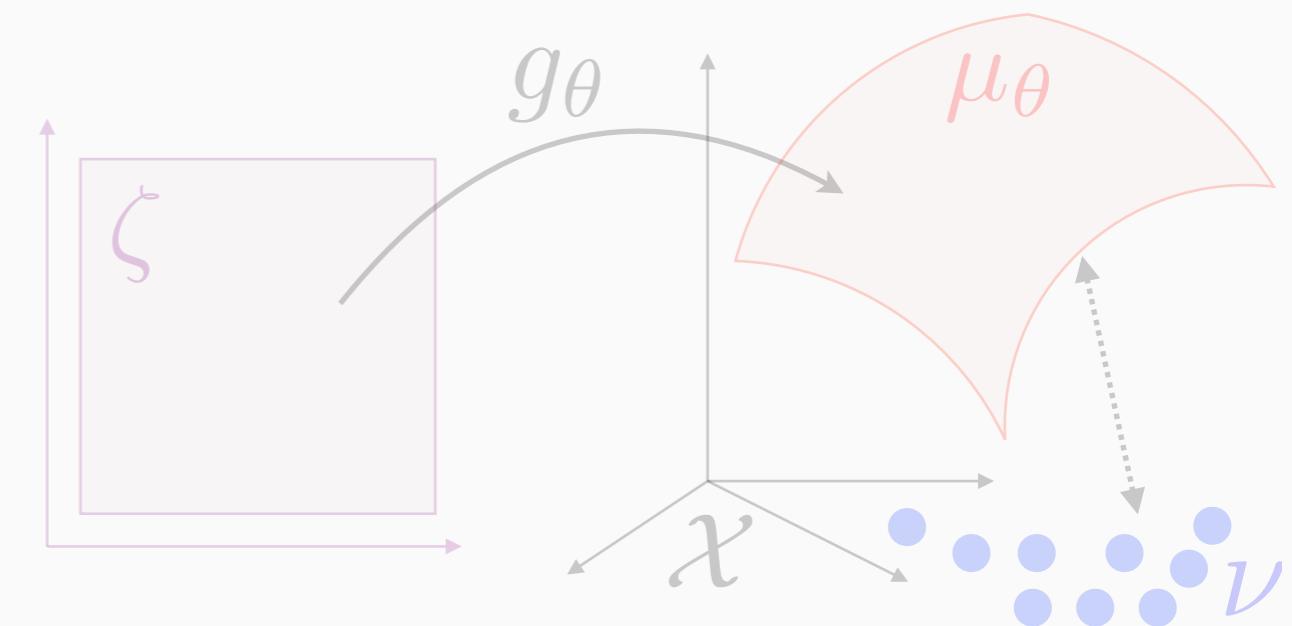
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



# 4. Application to Generative Models



# Kantorovitch's Formulation

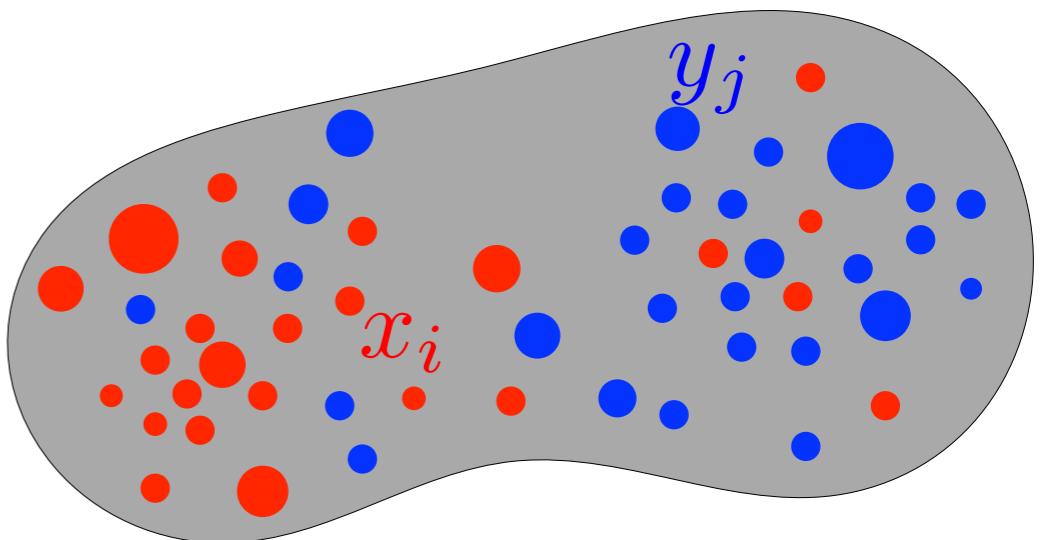
*Input distributions*

$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i} \quad \beta = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$$

Points  $(x_i)_i, (y_j)_j$

Weights  $\mathbf{a}_i \geq 0, \mathbf{b}_j \geq 0.$

$$\sum_{i=1}^n \mathbf{a}_i = \sum_{j=1}^m \mathbf{b}_j = 1$$



# Kantorovitch's Formulation

*Input distributions*

$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i} \quad \beta = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$$

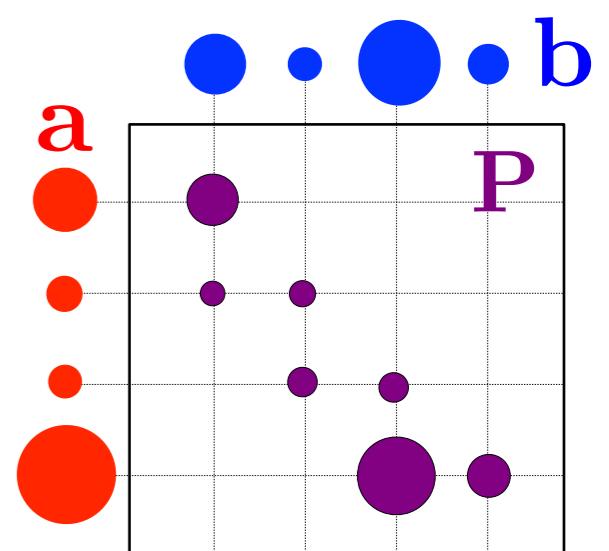
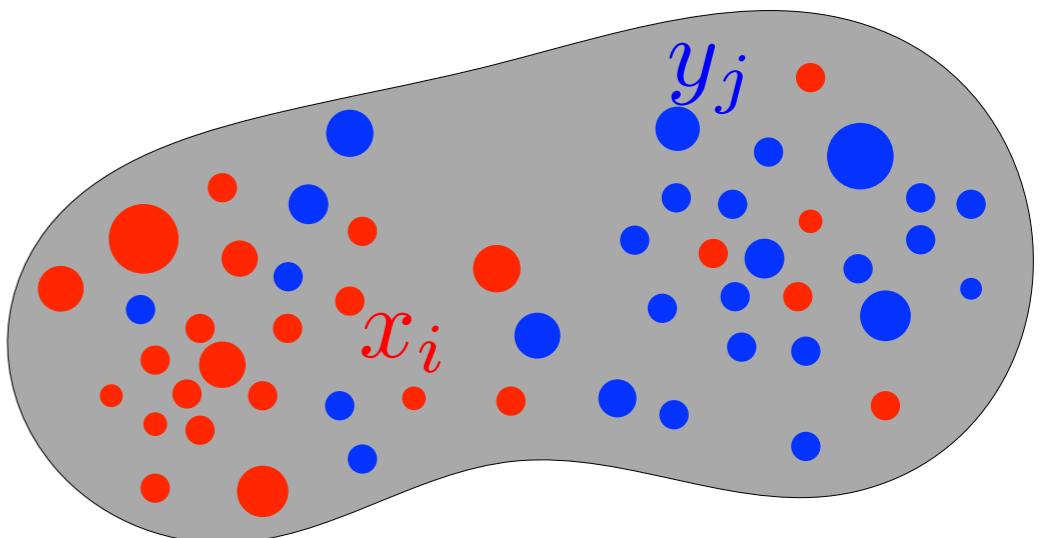
Points  $(x_i)_i, (y_j)_j$

Weights  $\mathbf{a}_i \geq 0, \mathbf{b}_j \geq 0.$

$$\sum_{i=1}^n \mathbf{a}_i = \sum_{j=1}^m \mathbf{b}_j = 1$$

Couplings:

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} ; \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_n = \mathbf{b} \}$$



# Kantorovitch's Formulation

*Input distributions*

$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i} \quad \beta = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$$

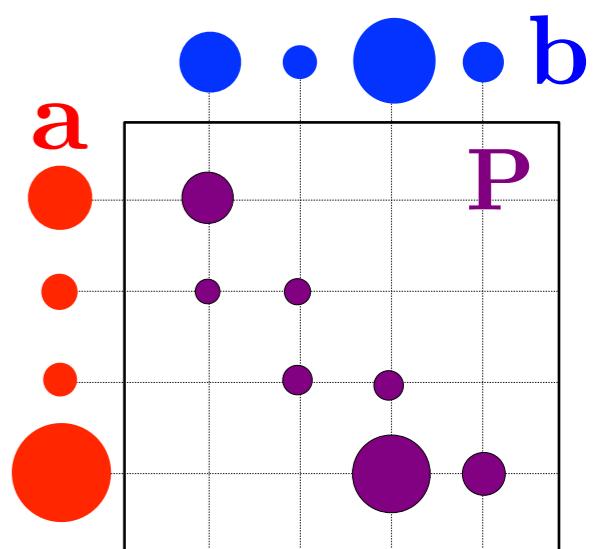
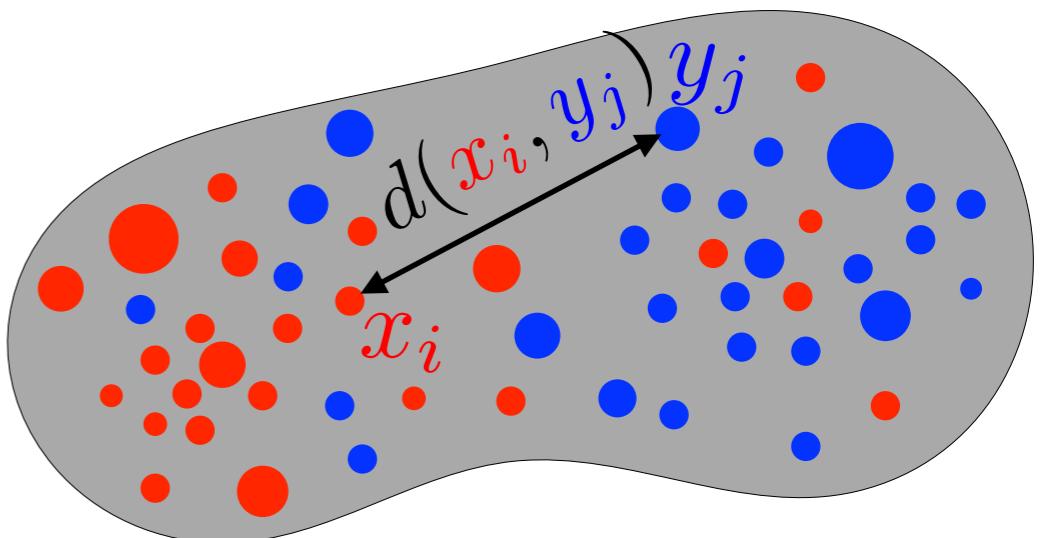
Points  $(x_i)_i, (y_j)_j$

Weights  $\mathbf{a}_i \geq 0, \mathbf{b}_j \geq 0.$

$$\sum_{i=1}^n \mathbf{a}_i = \sum_{j=1}^m \mathbf{b}_j = 1$$

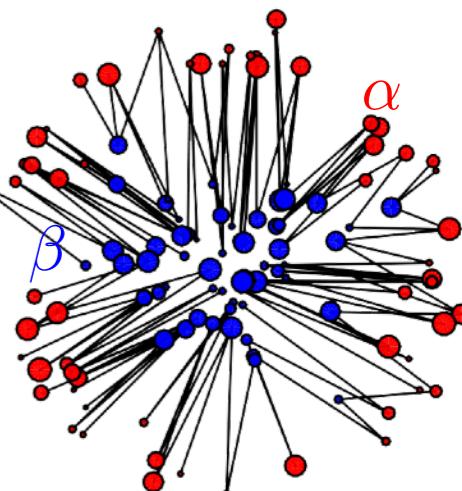
Couplings:

$$\mathbf{U}(\mathbf{a}, \mathbf{b}) \stackrel{\text{def.}}{=} \{ \mathbf{P} \in \mathbb{R}_+^{n \times m} ; \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_n = \mathbf{b} \}$$

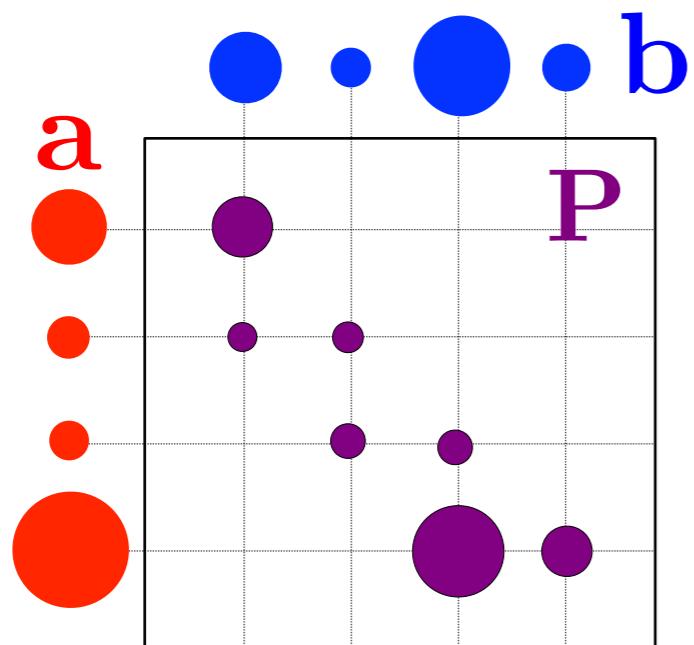


[Kantorovich 1942]

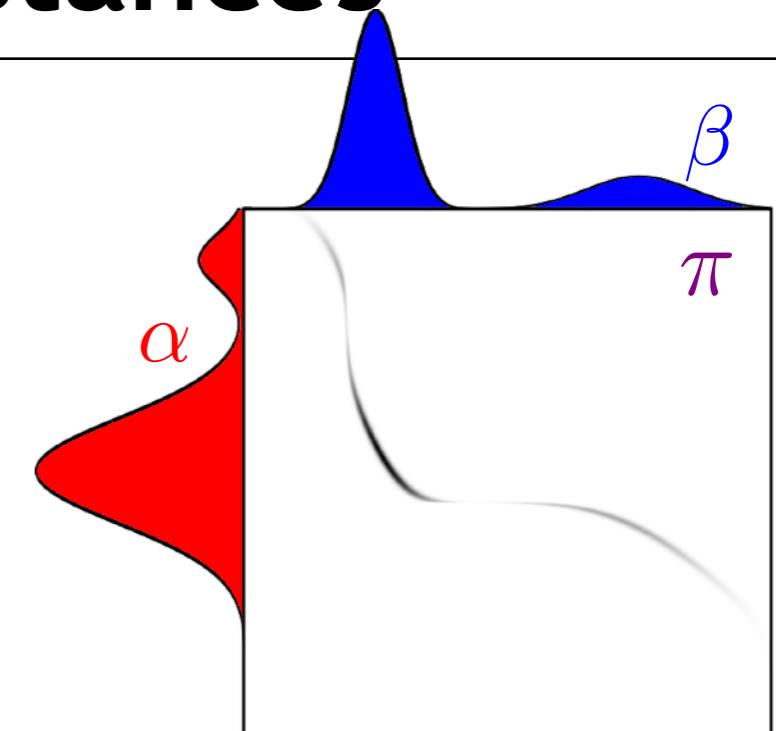
$$\min \left\{ \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} ; \mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b}) \right\}$$



# Optimal Transport Distances

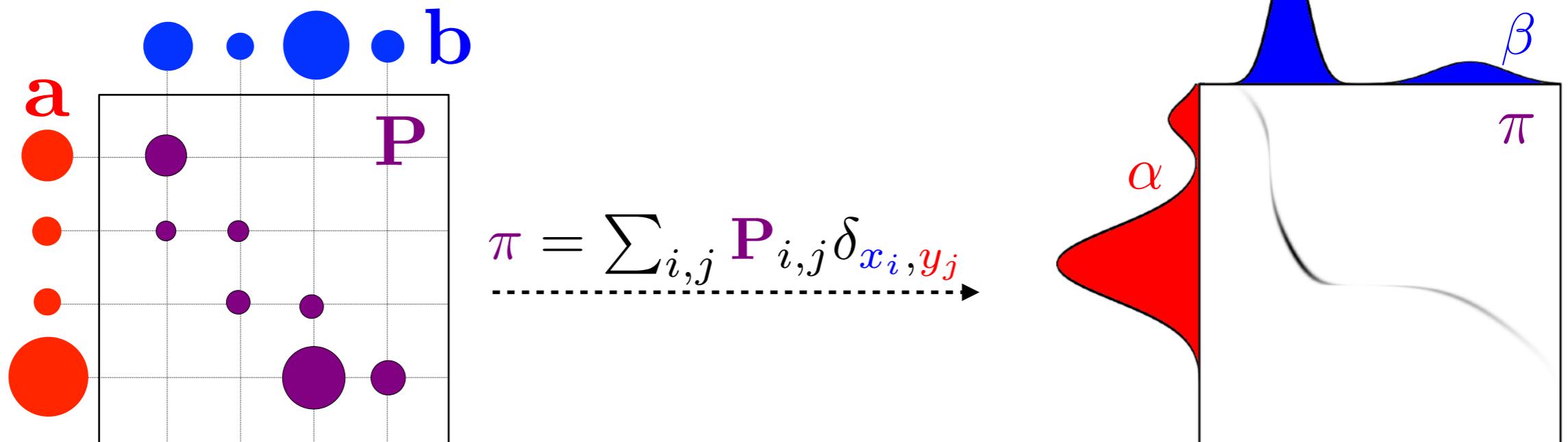


$$\pi = \sum_{i,j} P_{i,j} \delta_{x_i, y_j}$$



$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \min_{\pi \in \mathcal{M}_+^1(\mathcal{X}^2)} \left\{ \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$

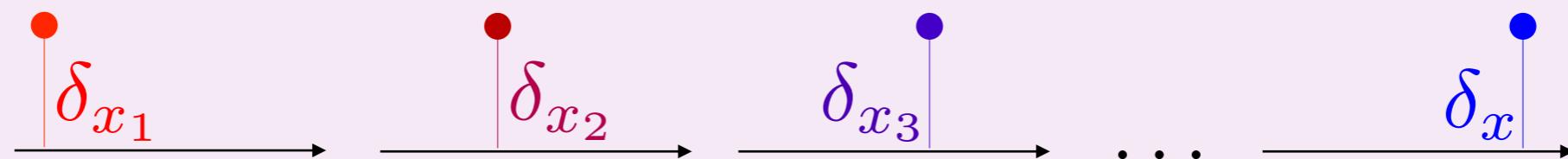
# Optimal Transport Distances



$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \min_{\pi \in \mathcal{M}_+^1(\mathcal{X}^2)} \left\{ \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$

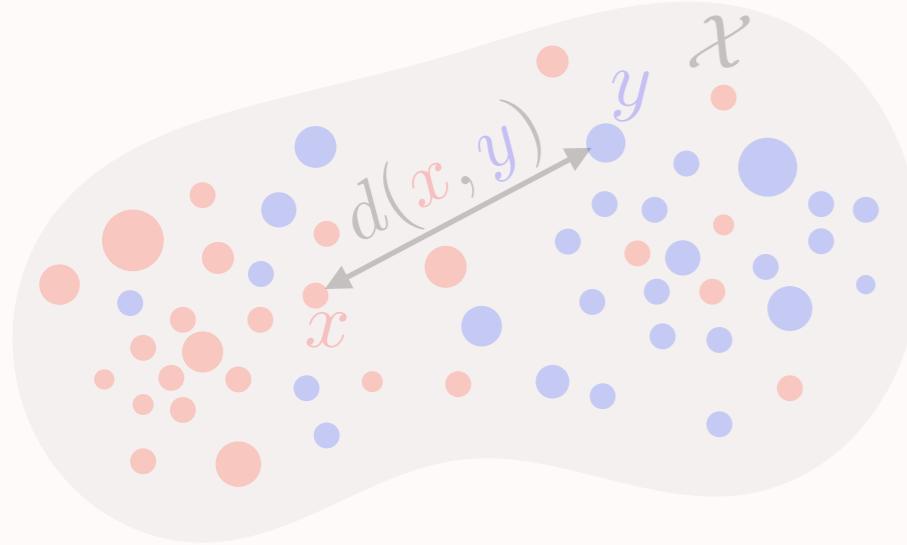
*Theorem:*  $W_p$  is a distance and  $\alpha_n \rightarrow \beta \Leftrightarrow W_p(\alpha_n, \beta) \rightarrow 0$

Weak\* (aka in law) convergence:  $\alpha_n \rightharpoonup \beta \Leftrightarrow \forall f \in \mathcal{C}(\mathcal{X}), \int_{\mathcal{X}} f d\alpha_n \rightarrow \int_{\mathcal{X}} f d\beta$

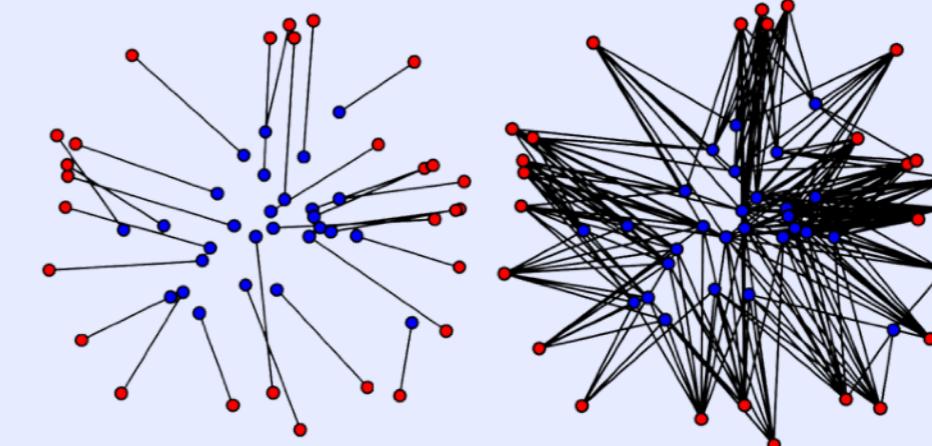


$$\|\delta_{x_n} - \delta_x\|_1 = 2 \quad \text{vs.} \quad W_p(\delta_{x_n}, \delta_x) = d(x_n, x)$$

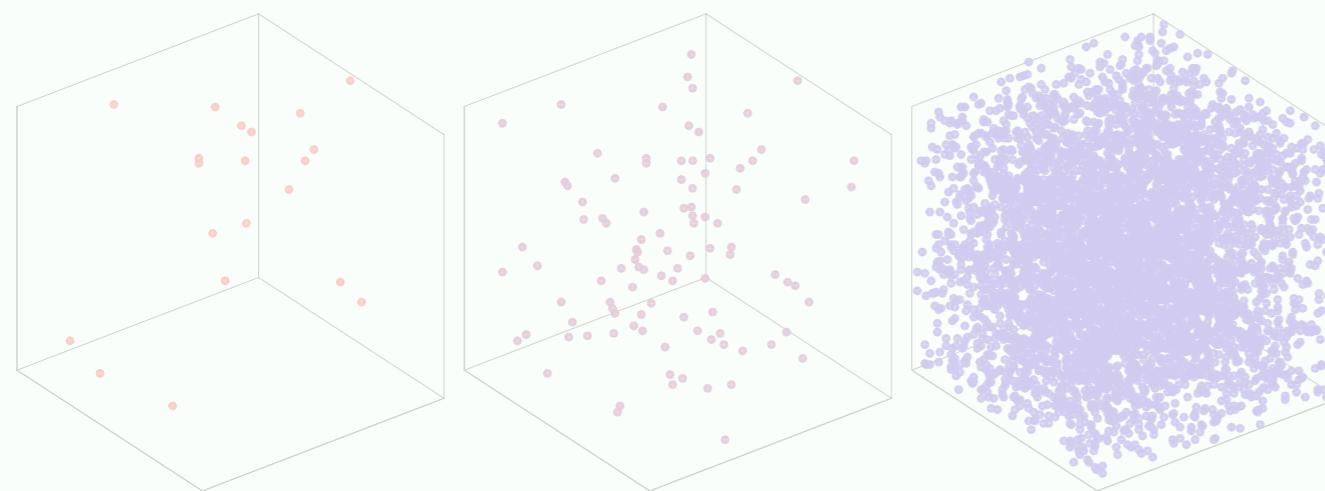
# 1. Optimal Transport



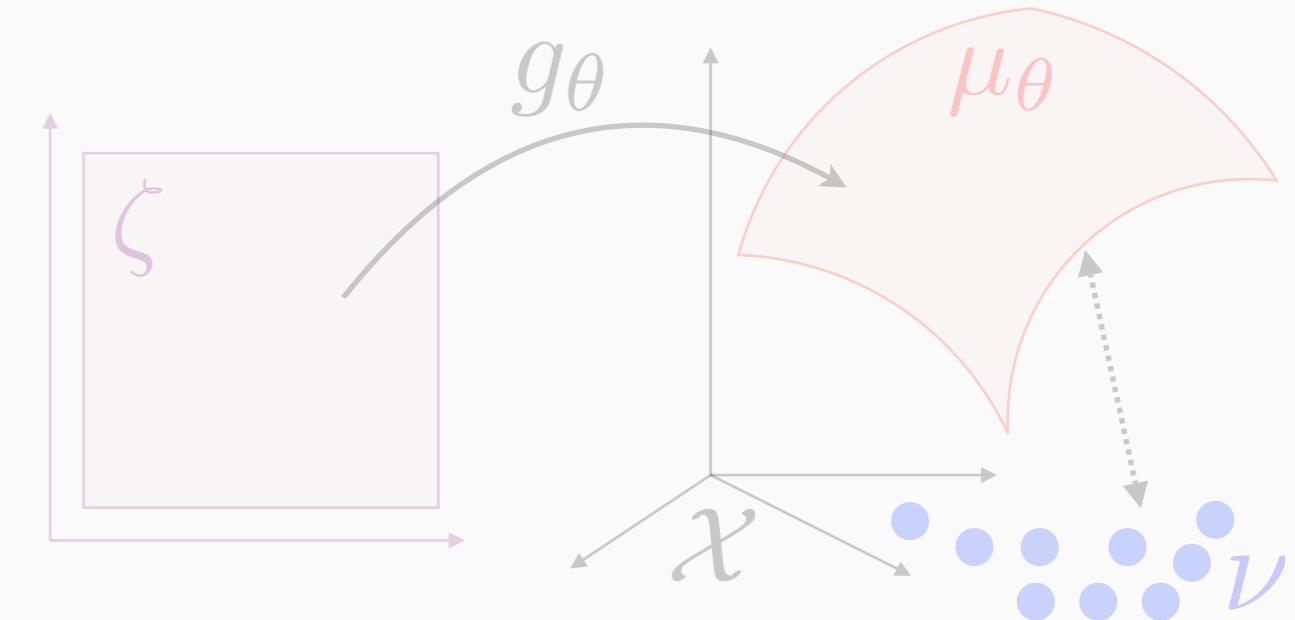
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



# 4. Application to Generative Models

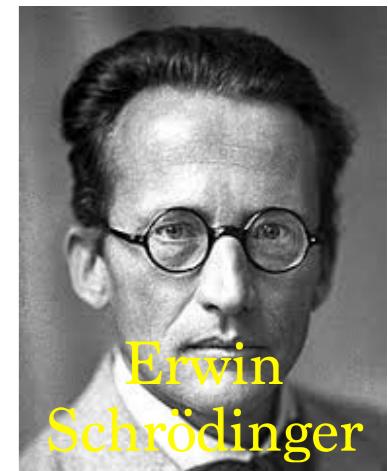


# Entropic Regularization

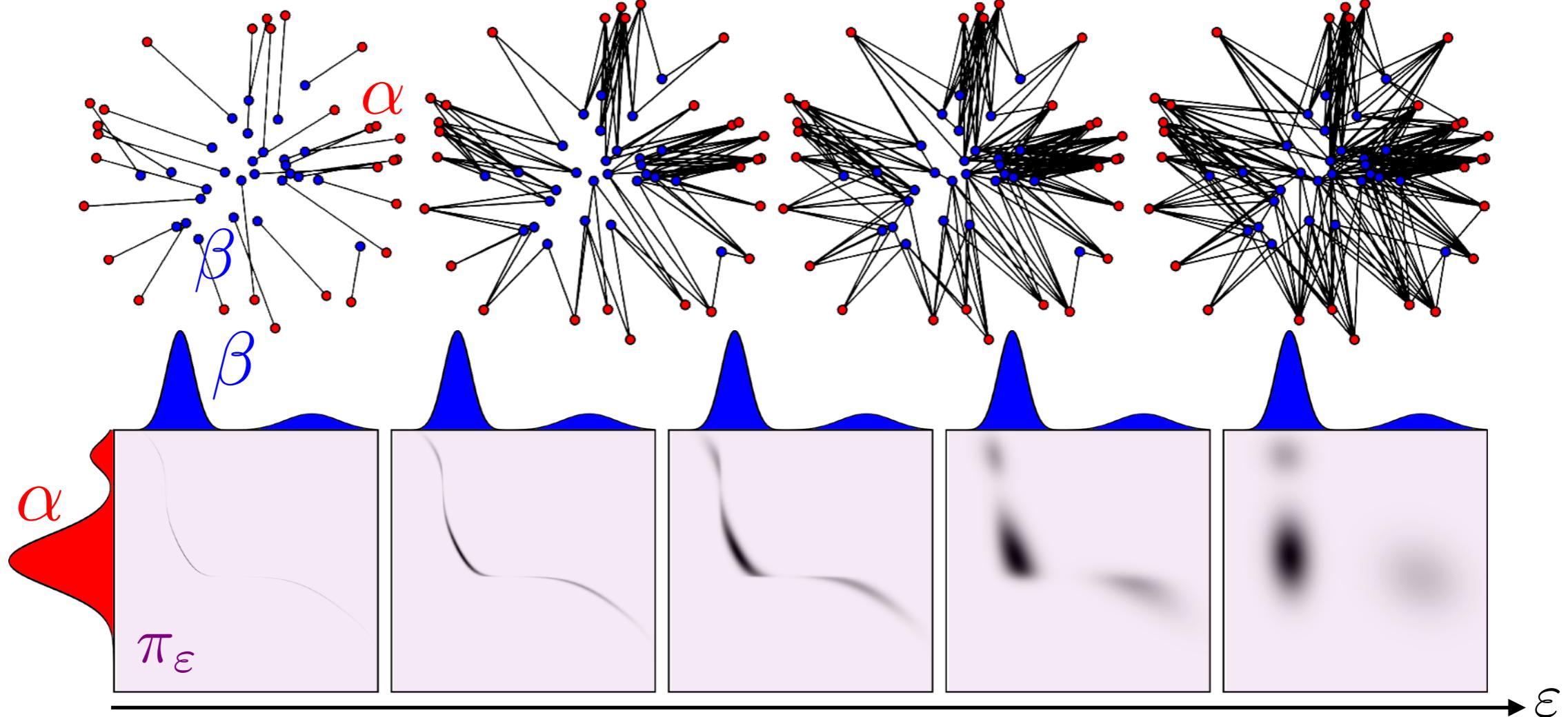
Schrödinger's problem:

[1931]

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$



$$\pi = \sum_{i,j} \mathbf{P}_{i,j} \delta_{x_i, y_j}$$



# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$

$$Proposition: \quad \mathbf{P}_{i,j} = \mathbf{u}_i \ \mathbf{K}_{i,j} \ \mathbf{v}_j \qquad \qquad \mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$

*Proposition:*  $\mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$

*Proposition:*  $\mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

Sinkhorn iterations:

$$\mathbf{u} \leftarrow \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}}$$

$$\mathbf{v} \leftarrow \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$$

*Theorem:* [Sinkhorn 1964]  $(\mathbf{u}, \mathbf{v})$  converges.

# Sinkhorn's Algorithm

$$\min_{\mathbf{P} \in \mathbf{U}(\mathbf{a}, \mathbf{b})} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log \left( \frac{\mathbf{P}_{i,j}}{\mathbf{a}_i \mathbf{b}_j} \right)$$

*Proposition:*  $\mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j$        $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$

Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

Sinkhorn iterations:

$$\mathbf{u} \leftarrow \frac{\mathbf{a}}{\mathbf{K}\mathbf{v}}$$

$$\mathbf{v} \leftarrow \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}}$$

*Theorem:* [Sinkhorn 1964]  $(\mathbf{u}, \mathbf{v})$  converges.

Only matrix/vector multiplications.

Matrix-vectors

$$\mathbf{K} \begin{array}{|c|} \hline \mathbf{v}^1 \\ \hline \vdots \\ \hline \mathbf{v}^q \end{array}, \dots, \mathbf{K} \begin{array}{|c|} \hline \mathbf{v}^1 \\ \hline \vdots \\ \hline \mathbf{v}^q \end{array}$$

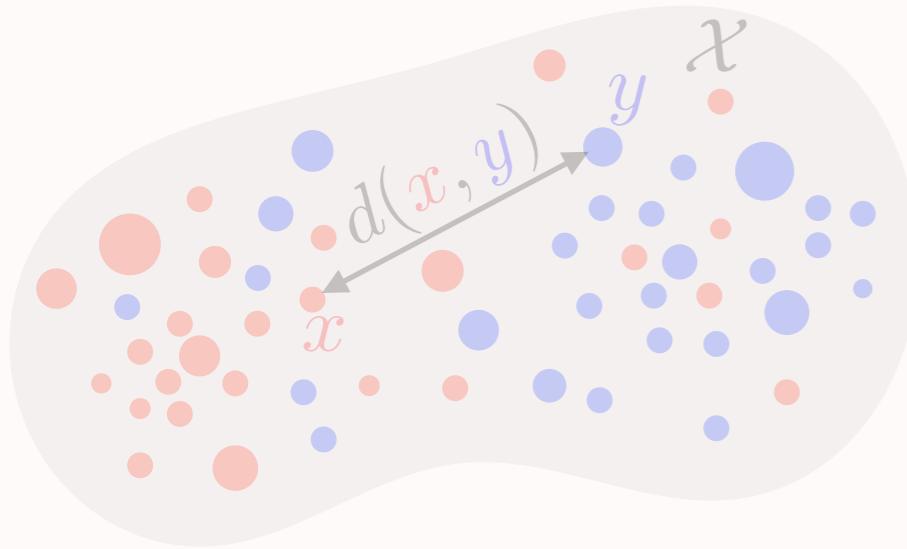
parallelization  
GPU

Matrix-matrix

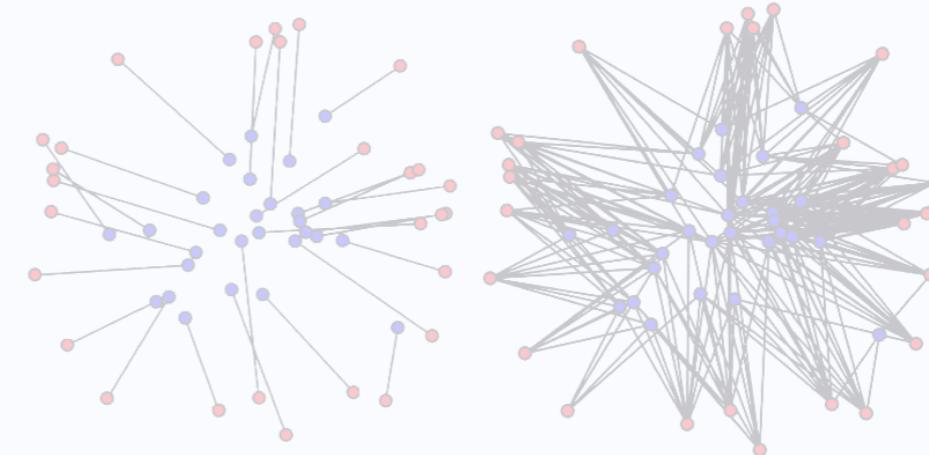
$$\mathbf{K} \begin{array}{|c|} \hline \mathbf{v}^1 \\ \hline \vdots \\ \hline \mathbf{v}^q \end{array}, \dots, \mathbf{V} \begin{array}{|c|} \hline \mathbf{v}^1 \\ \hline \vdots \\ \hline \mathbf{v}^q \end{array}$$

→ Convolution on regular grids, separable kernels.

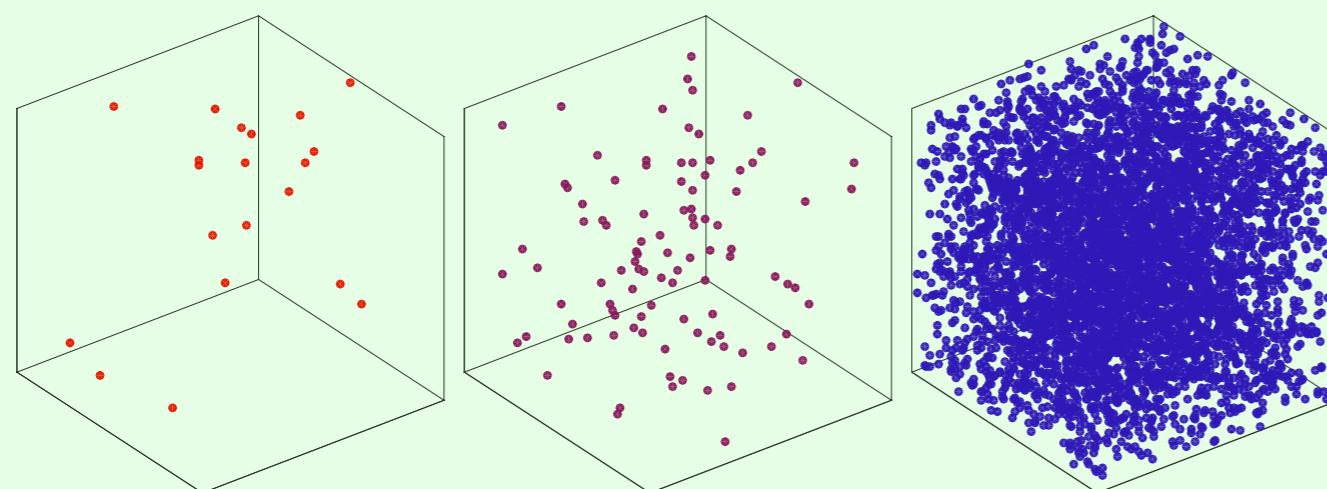
# 1. Optimal Transport



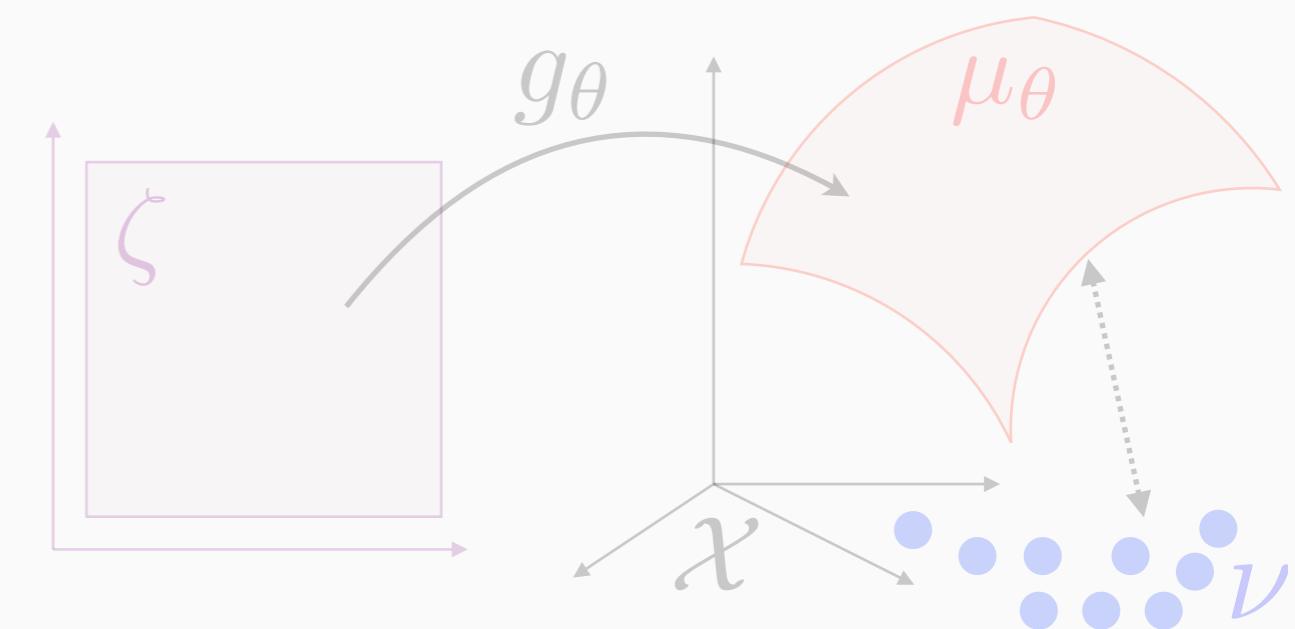
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



# 4. Application to Generative Models

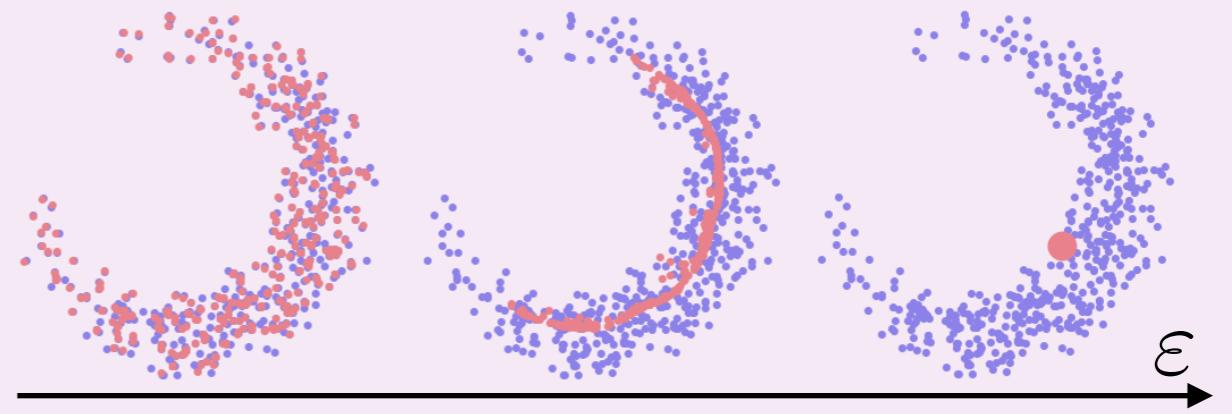


# Sinkhorn Divergences

$$W_{\varepsilon,p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \xi)$$

*Problem:*  $W_\varepsilon(\alpha, \alpha) \neq 0$

$$\min_{\alpha} W_{\varepsilon,p}^p(\alpha, \beta)$$

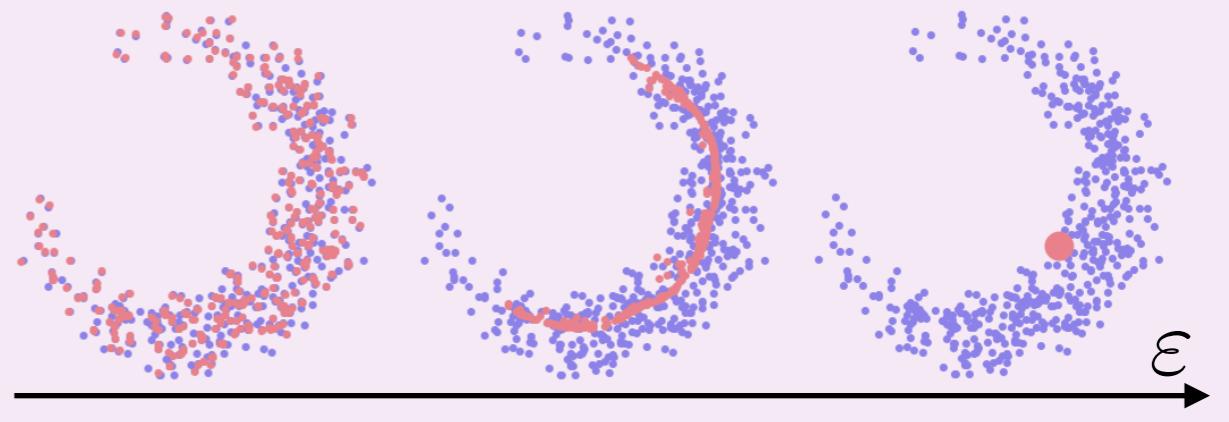


# Sinkhorn Divergences

$$W_{\varepsilon,p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \xi)$$

Problem:  $W_\varepsilon(\alpha, \alpha) \neq 0$

$$\min_{\alpha} W_{\varepsilon,p}^p(\alpha, \beta)$$



$$\overline{W}_{p,\varepsilon}^p(\alpha, \beta) \stackrel{\text{def.}}{=} W_{p,\varepsilon}^p(\alpha, \beta) - \frac{1}{2} W_{p,\varepsilon}^p(\alpha, \alpha) - \frac{1}{2} W_{p,\varepsilon}^p(\beta, \beta)$$

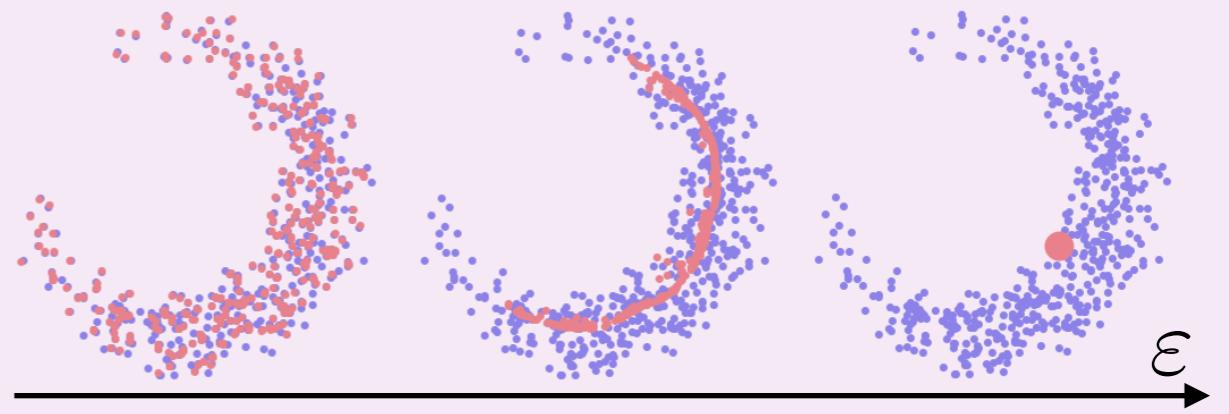
[Ramdas, García Trillos, Cuturi, 2017]

# Sinkhorn Divergences

$$W_{\varepsilon,p}^p(\alpha, \beta) \stackrel{\text{def.}}{=} \min_{\pi_1 = \alpha, \pi_2 = \beta} \int_{\mathcal{X}^2} d^p(x, y) d\pi(x, y) + \varepsilon \text{KL}(\pi | \xi)$$

Problem:  $W_\varepsilon(\alpha, \alpha) \neq 0$

$$\min_{\alpha} W_{\varepsilon,p}^p(\alpha, \beta)$$



$$\overline{W}_{p,\varepsilon}^p(\alpha, \beta) \stackrel{\text{def.}}{=} W_{p,\varepsilon}^p(\alpha, \beta) - \frac{1}{2} W_{p,\varepsilon}^p(\alpha, \alpha) - \frac{1}{2} W_{p,\varepsilon}^p(\beta, \beta)$$

[Ramdas, García Trillos, Cuturi, 2017]

$$\text{Theorem: } W_p^p(\alpha, \beta) \xleftarrow[\substack{[\text{Léonard 2012}] \\ [\text{Carlier et al 2017}]}]{\varepsilon \rightarrow 0} \overline{W}_{\varepsilon,p}^p(\alpha, \beta) \xrightarrow{[\text{Ramdas, García Trillos, Cuturi, 2017}]} \|\alpha - \beta\|_{-d^p}^2$$

Kernel norms (MMD):  $\|\xi\|_{-d^p}^2 \stackrel{\text{def.}}{=} - \int_{\mathcal{X}^2} d(x, y)^p d\xi(x) d\xi(y)$

Proposition:  $\|\cdot\|_{-\|\cdot\|^p}$  is a norm for  $0 < p < 2$ .



# Sinkhorn Divergences

$$\overline{W}_{p,\varepsilon}^p(\alpha, \beta) \stackrel{\text{def.}}{=} W_{p,\varepsilon}^p(\alpha, \beta) - \frac{1}{2} W_{p,\varepsilon}^p(\alpha, \alpha) - \frac{1}{2} W_{p,\varepsilon}^p(\beta, \beta)$$

↓ concave      ↓ concave

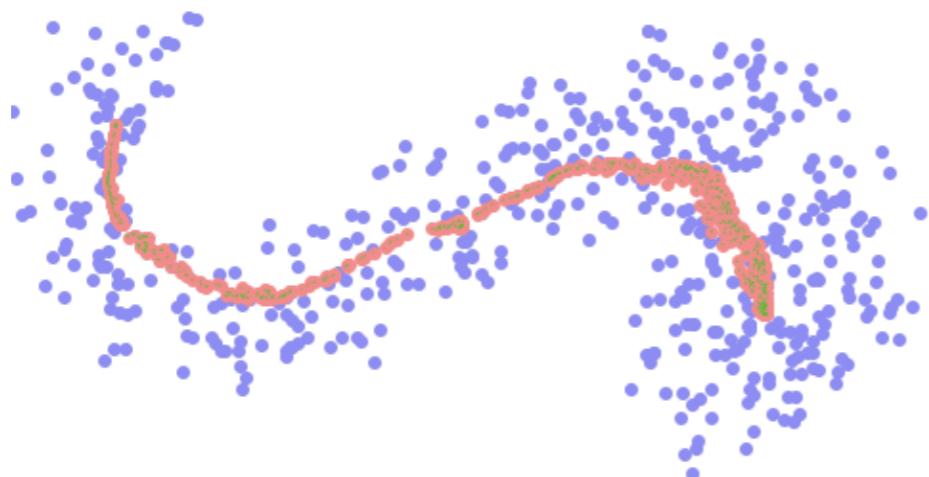
*Theorem:* [Feydy, Séjourné, P, Vialard, Trouvé, Amari 2018]

If  $e^{-\frac{d^p}{\varepsilon}}$  is positive:

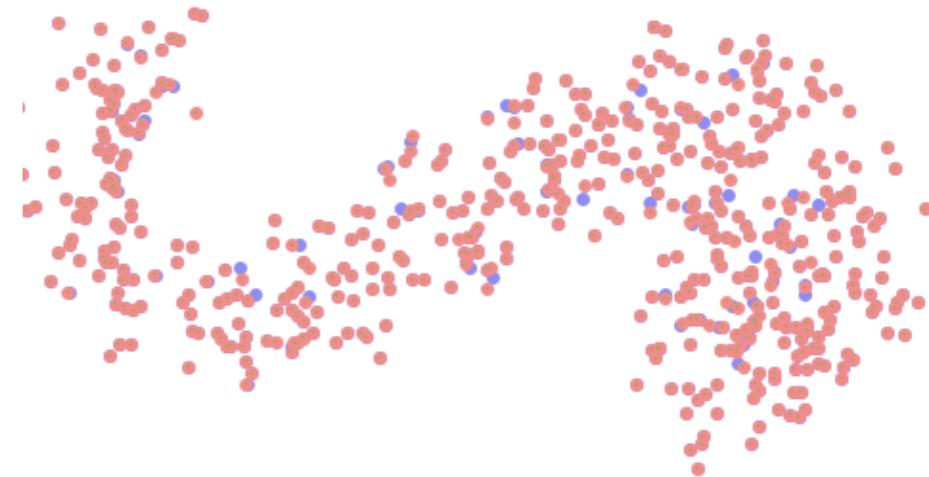
$\overline{W}_{\varepsilon,p} \geqslant 0$  and  $\overline{W}_{\varepsilon,p}^p(\cdot, \beta)$  is convex.

$\overline{W}_{\varepsilon,p}(\alpha_n, \beta) \rightarrow 0 \iff \alpha_n \xrightarrow{\text{weak*}} \beta$

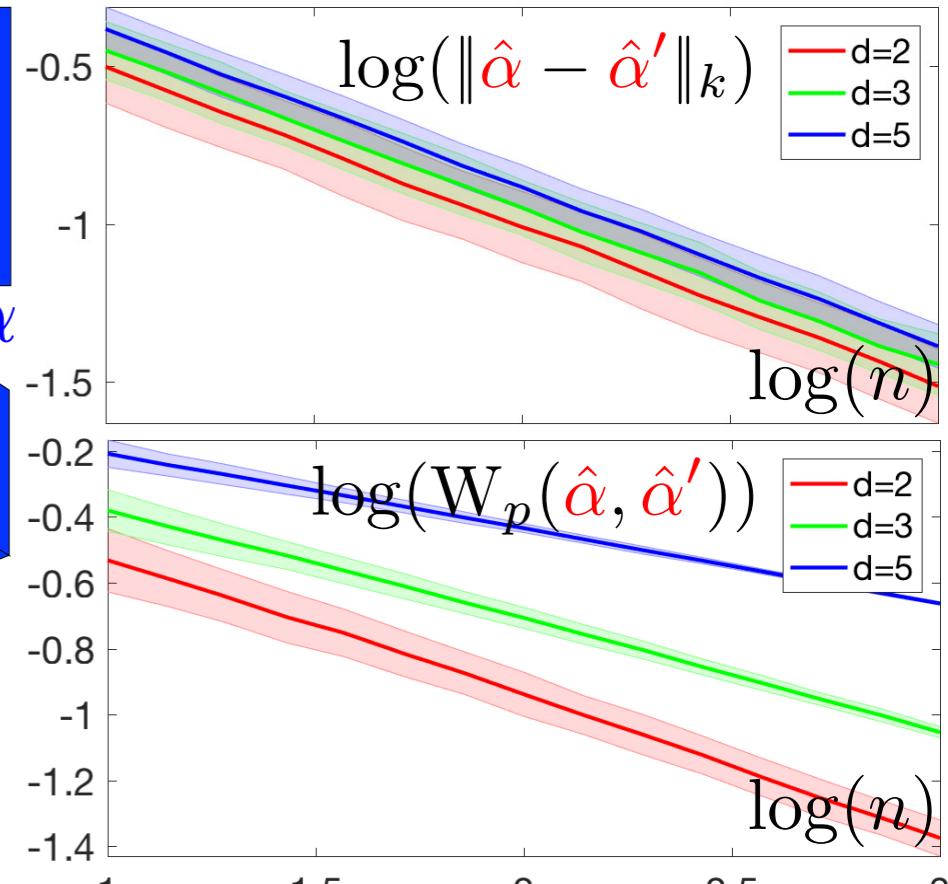
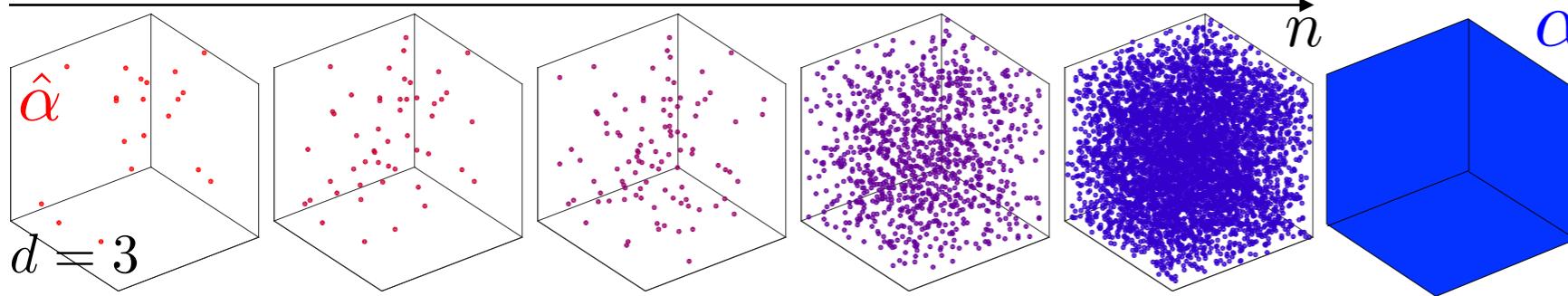
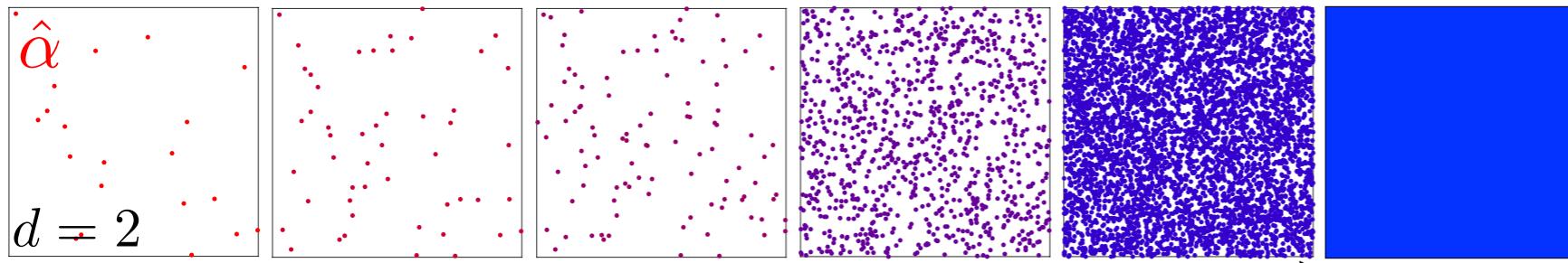
$$\min_{\alpha} W_{\varepsilon,p}^p(\alpha, \beta)$$



$$\min_{\alpha} \overline{W}_{\varepsilon,p}^p(\alpha, \beta)$$



# Sample Complexity



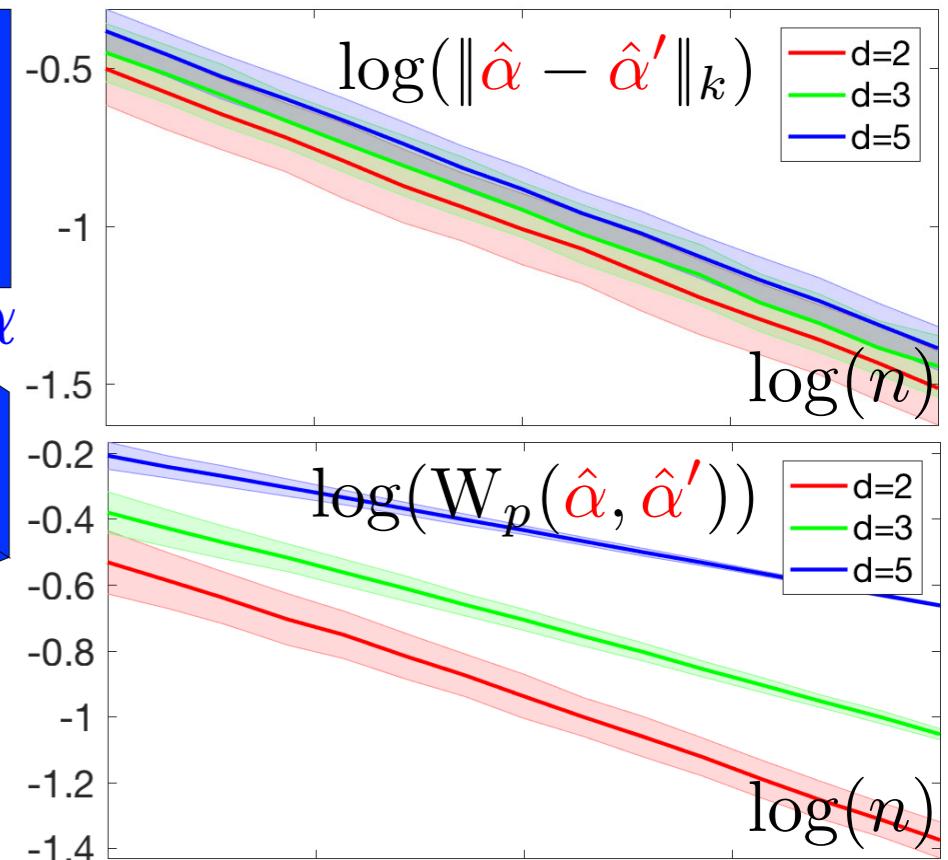
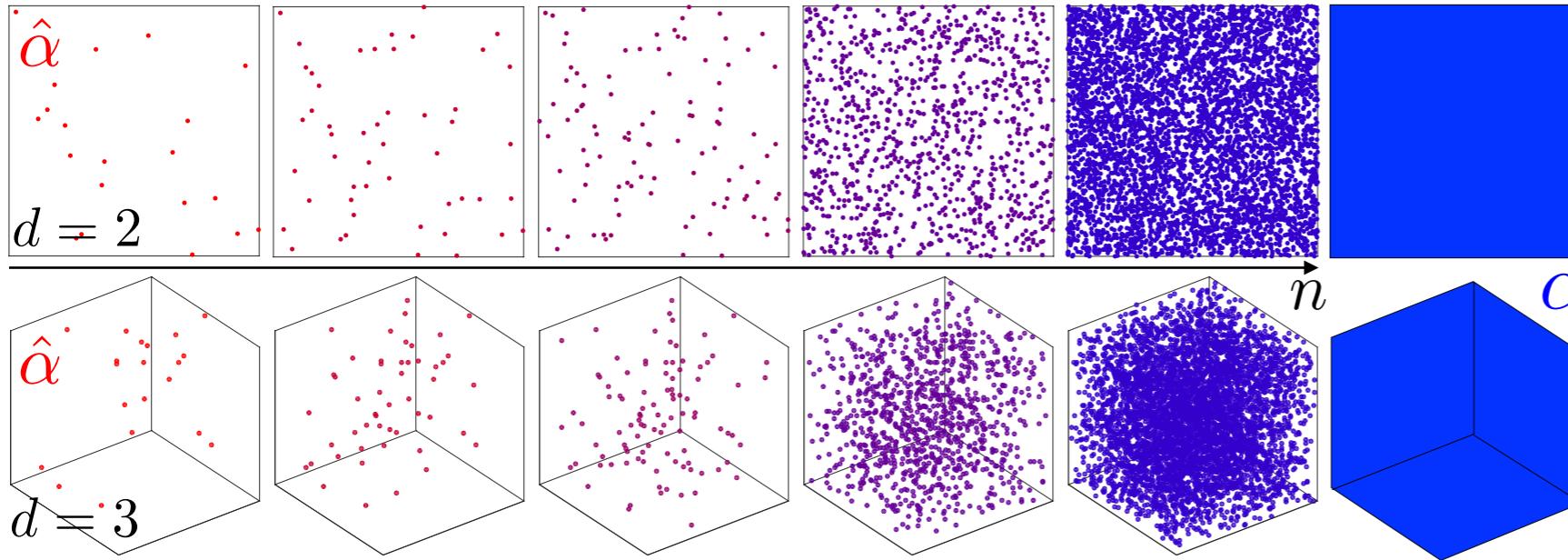
*Theorem:*  $\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}})$

$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

*Optimal transport:* suffers from curse of dimensionality.

→ Adapt to support dimensionality [Weed, Bach 2017]

# Sample Complexity



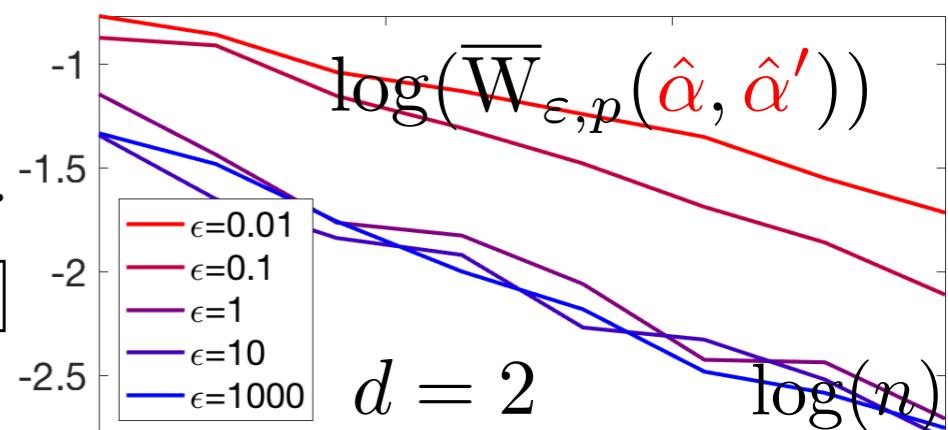
Theorem:

$$\mathbb{E}(|W_p(\hat{\alpha}, \hat{\beta}) - W_p(\alpha, \beta)|) = O(n^{-\frac{1}{d}})$$

$$\mathbb{E}(|\|\hat{\alpha} - \hat{\beta}\|_k - \|\alpha - \beta\|_k|) = O(n^{-\frac{1}{2}})$$

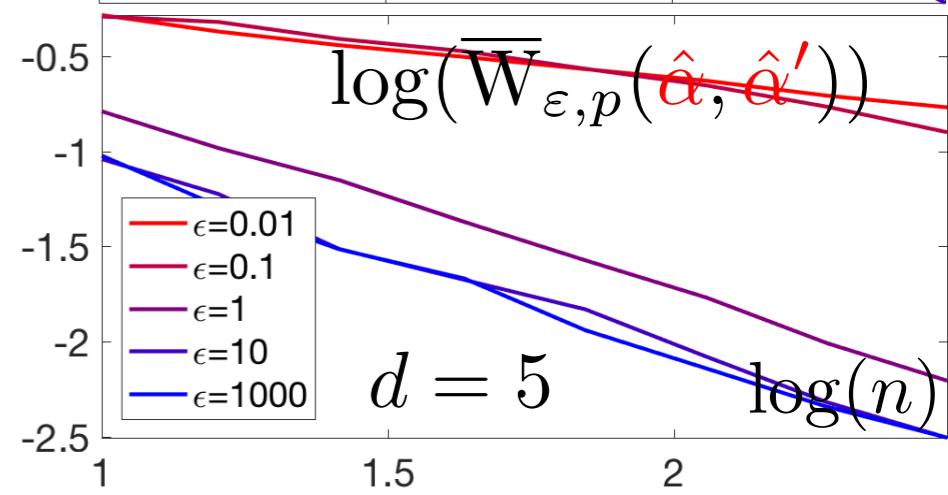
Optimal transport: suffers from curse of dimensionality.

→ Adapt to support dimensionality [Weed, Bach 2017]

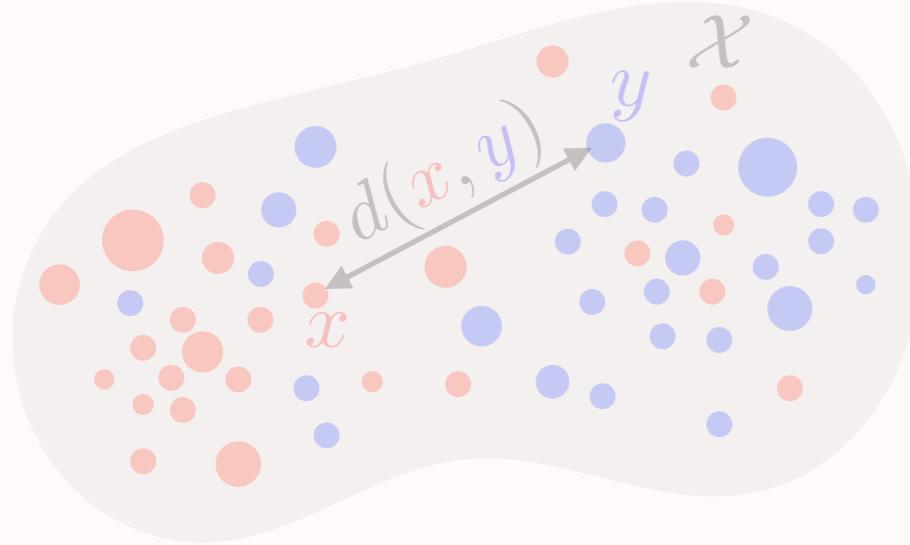


Theorem: [Genevay, Bach, P, Cuturi]

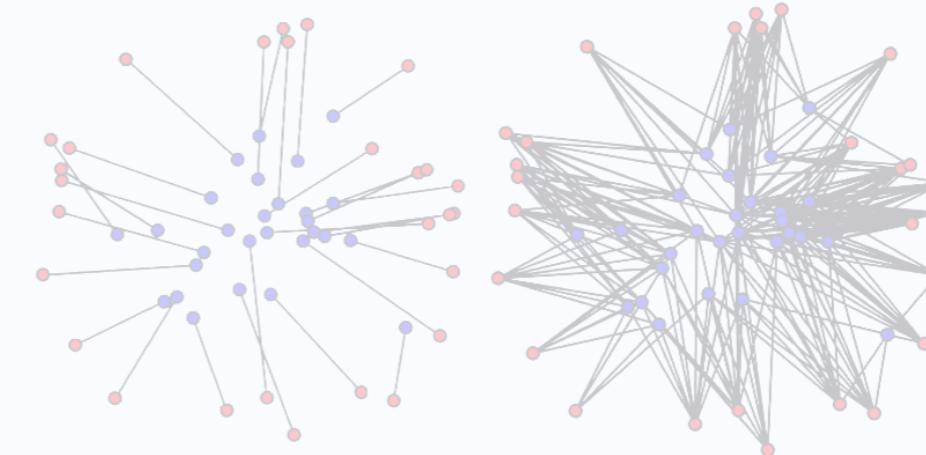
$$\mathbb{E}(|\overline{W}_{\epsilon,p}(\hat{\alpha}, \hat{\beta}) - \overline{W}_{\epsilon,p}(\alpha, \beta)|) = O(\epsilon^{-\frac{d}{2}} n^{-\frac{1}{2}})$$



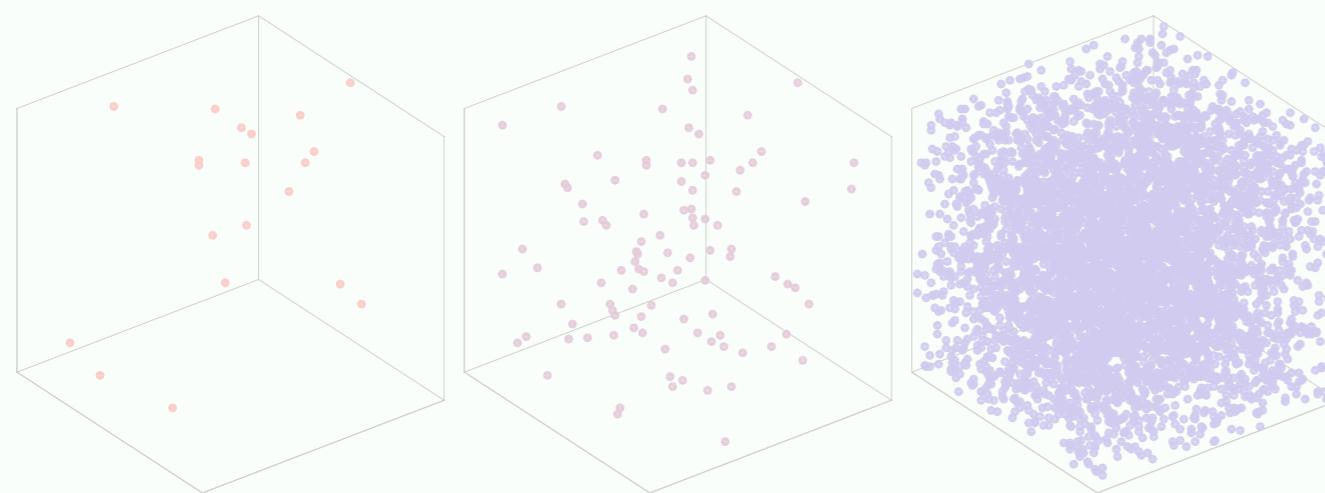
# 1. Optimal Transport



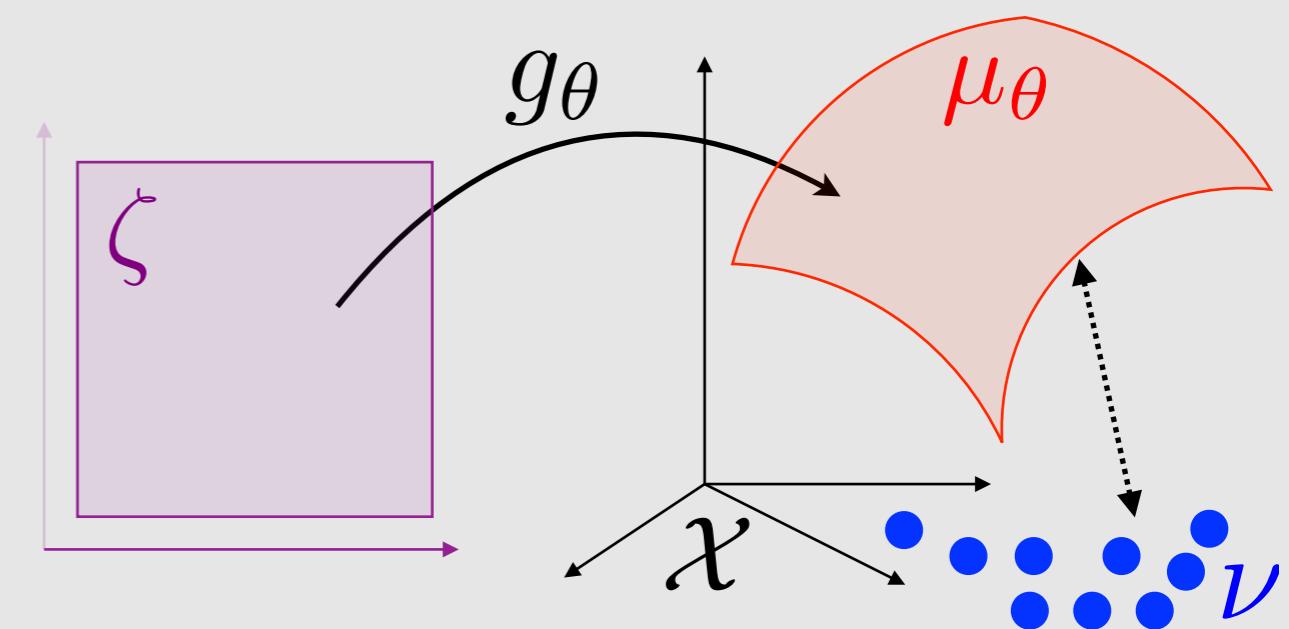
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



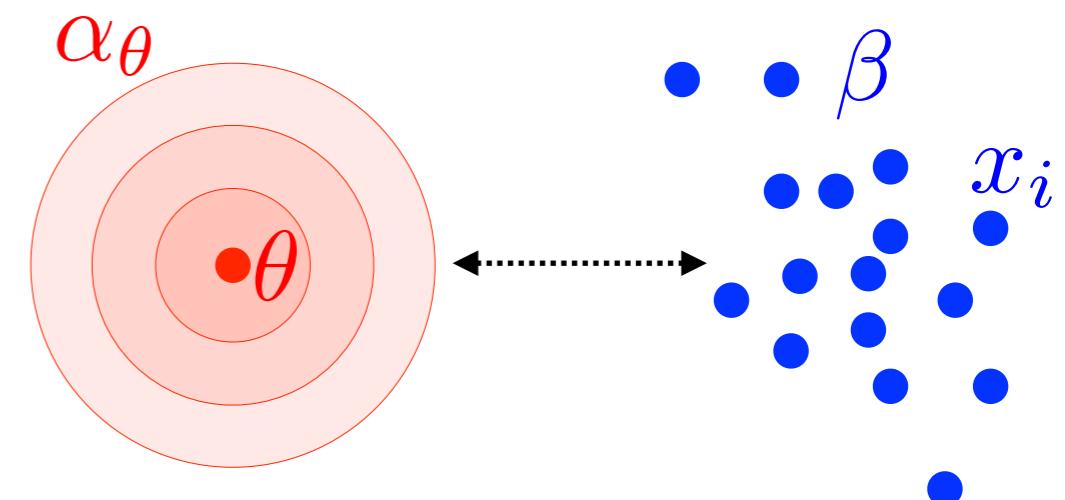
# 4. Application to Generative Models



# Density Fitting and Generative Models

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$



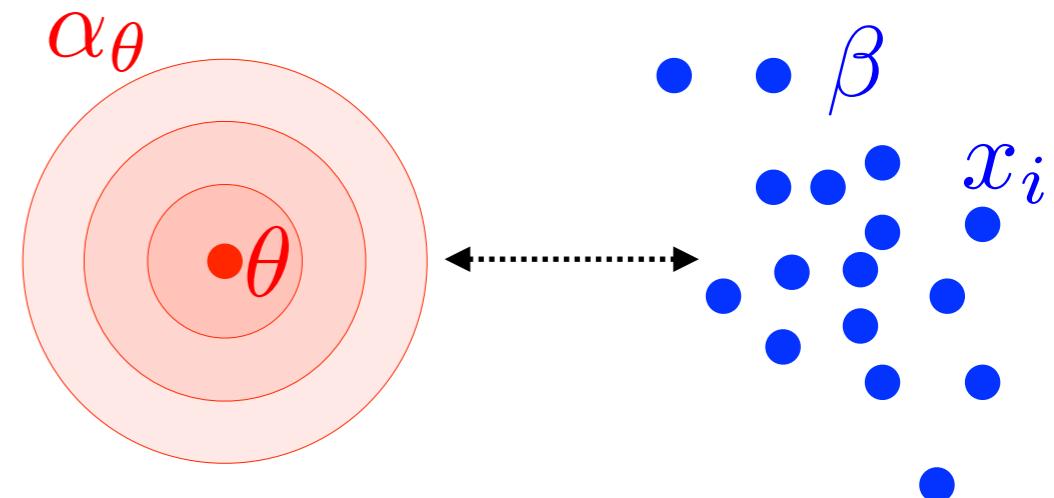
# Density Fitting and Generative Models

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$

Density fitting:  $d\alpha_\theta(x) = \rho_\theta(x)dx$

$$\min_{\theta} \widehat{\text{KL}}(\alpha_\theta | \beta) \stackrel{\text{def.}}{=} - \sum_i \log(\rho_\theta(x_i))$$



Maximum likelihood (MLE)

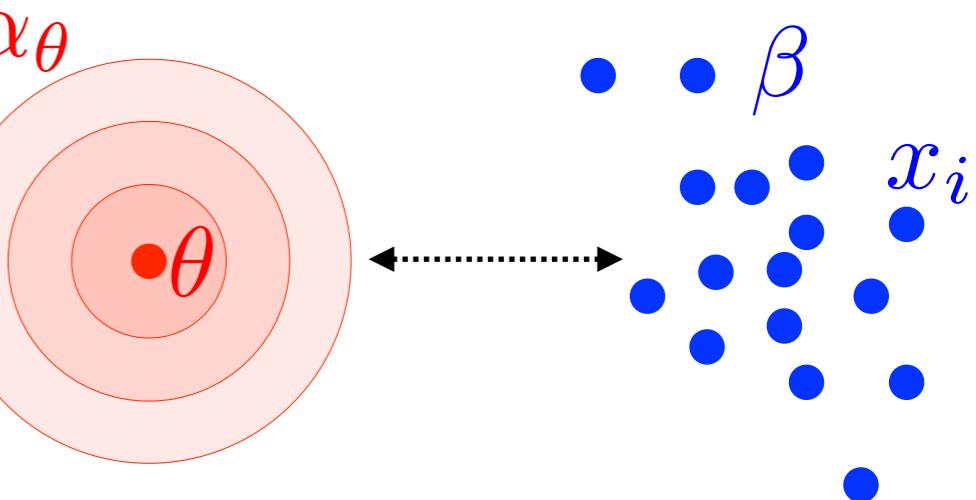
# Density Fitting and Generative Models

Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \alpha_\theta$

Density fitting:  $d\alpha_\theta(x) = \rho_\theta(x)dx$

$$\min_{\theta} \widehat{\text{KL}}(\alpha_\theta | \beta) \stackrel{\text{def.}}{=} - \sum_i \log(\rho_\theta(x_i))$$



Maximum likelihood (MLE)

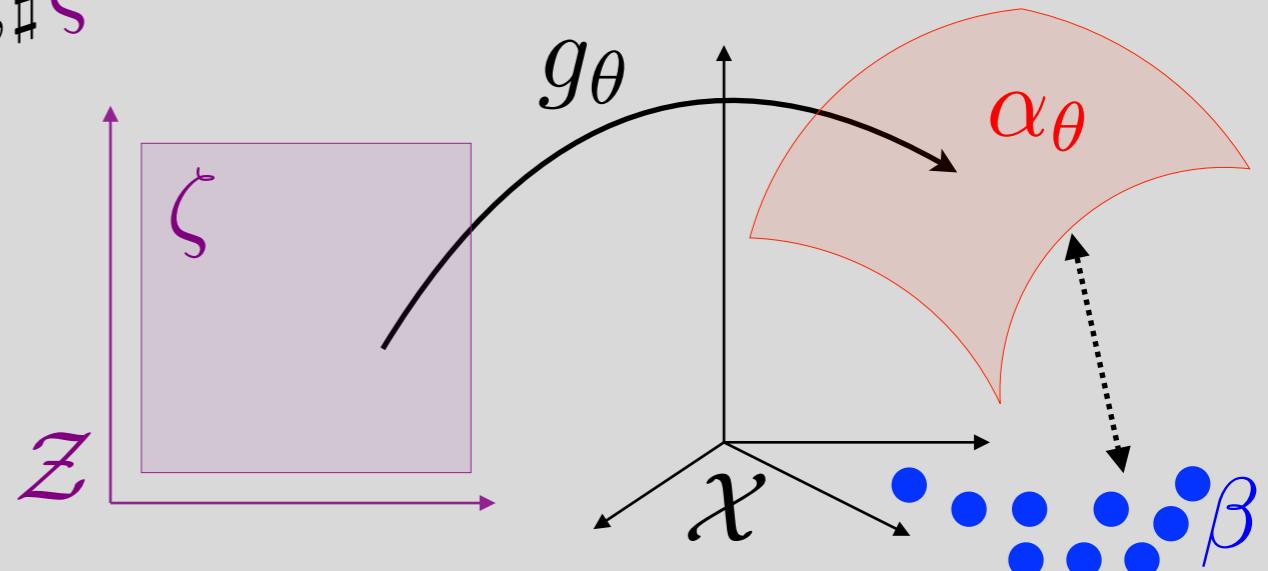
Generative model fit:  $\alpha_\theta = g_{\theta, \sharp} \zeta$

$$\widehat{\text{KL}}(\alpha_\theta | \beta) = +\infty$$

→ MLE undefined.

→ Need a weaker metric.

$$\min_{\theta} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$



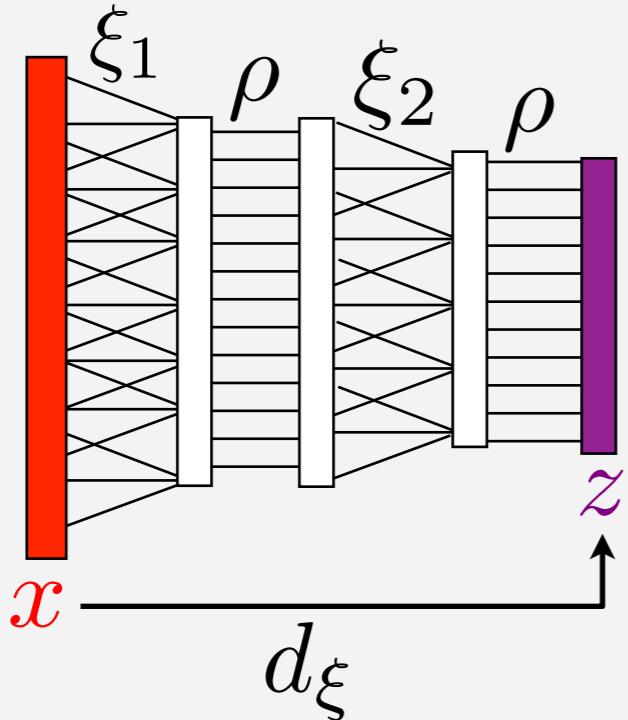
# Deep Discriminative vs Generative Models

Deep networks:

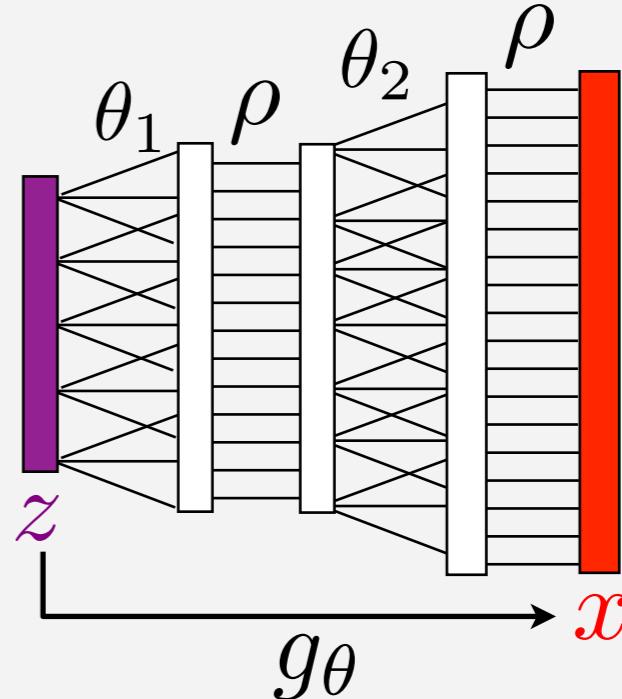
$$d_\xi(\textcolor{red}{x}) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(\textcolor{red}{x}) \dots)$$

$$g_\theta(\textcolor{violet}{z}) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(\textcolor{violet}{z}) \dots)$$

Discriminative



Generative



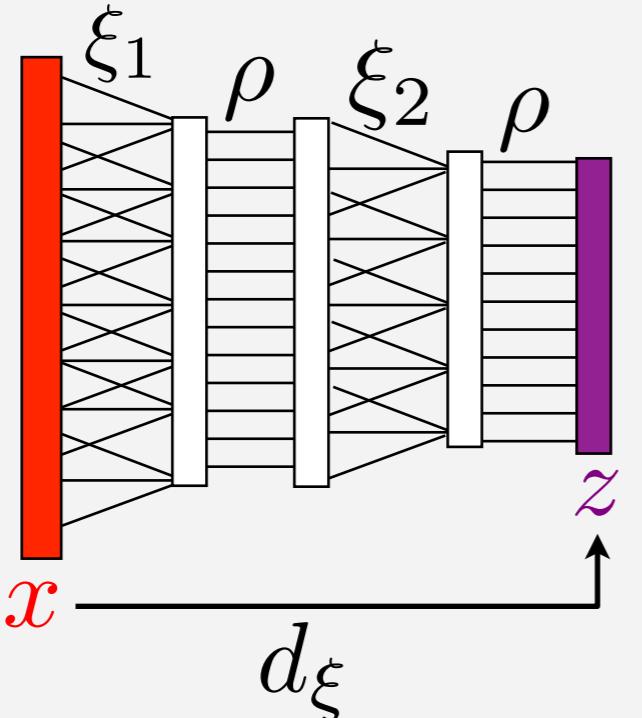
# Deep Discriminative vs Generative Models

Deep networks:

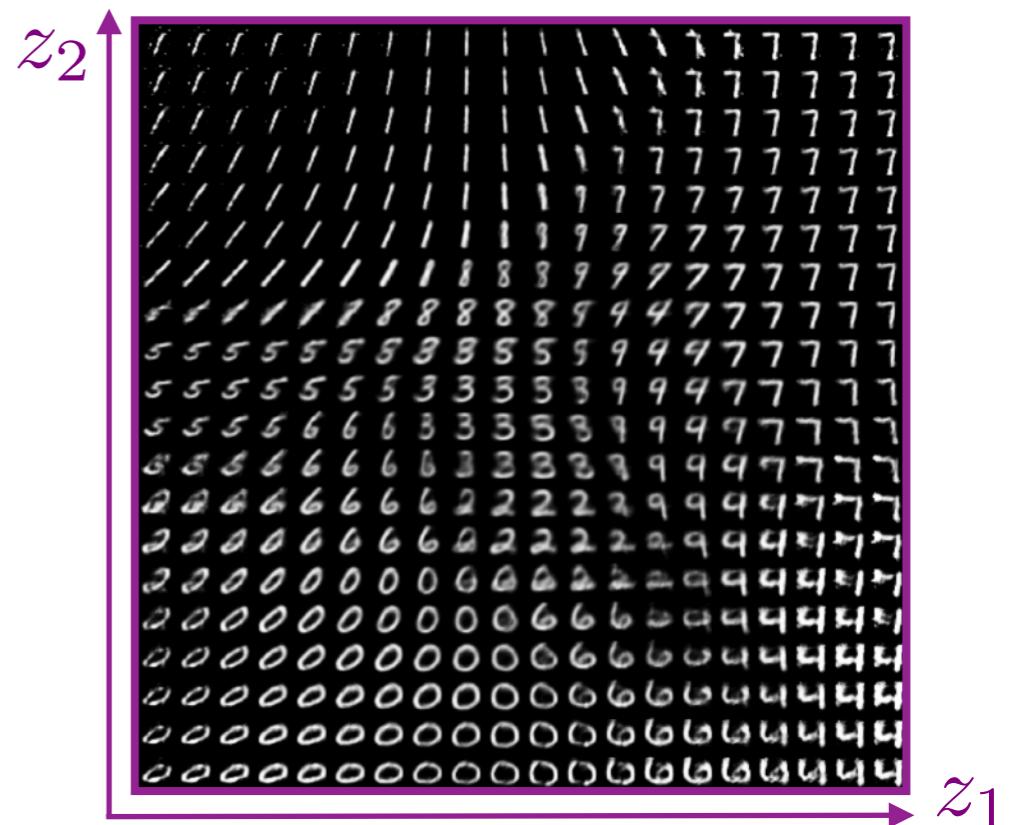
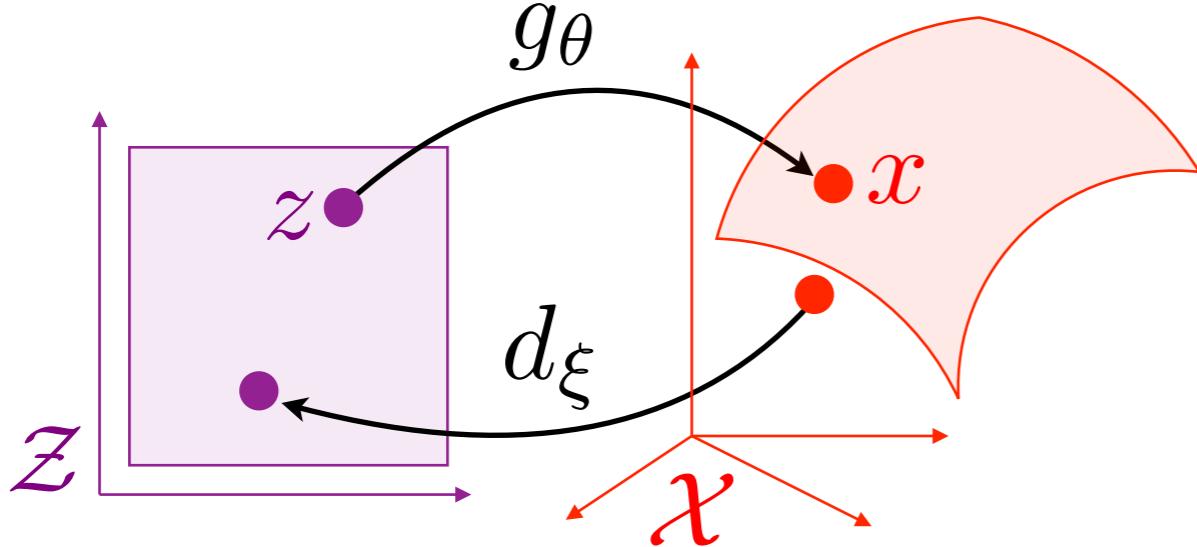
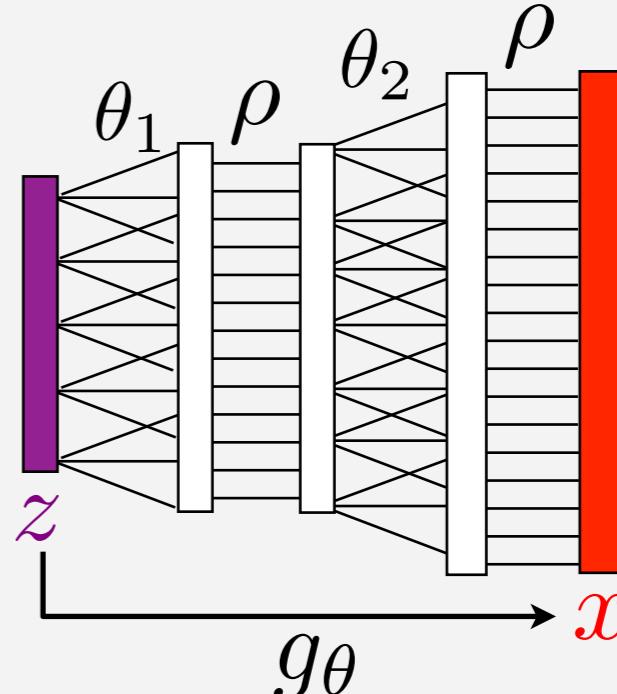
$$d_\xi(\mathbf{x}) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(\mathbf{x}) \dots)$$

$$g_\theta(\mathbf{z}) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(\mathbf{z}) \dots)$$

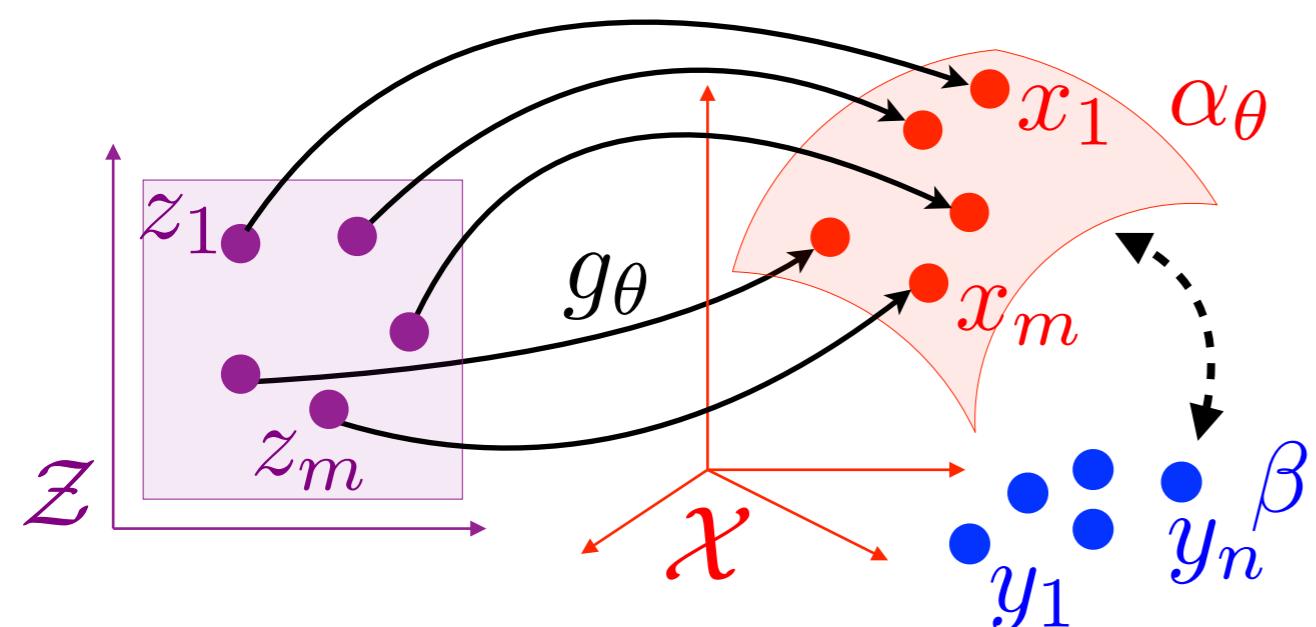
Discriminative



Generative



# Training Architecture



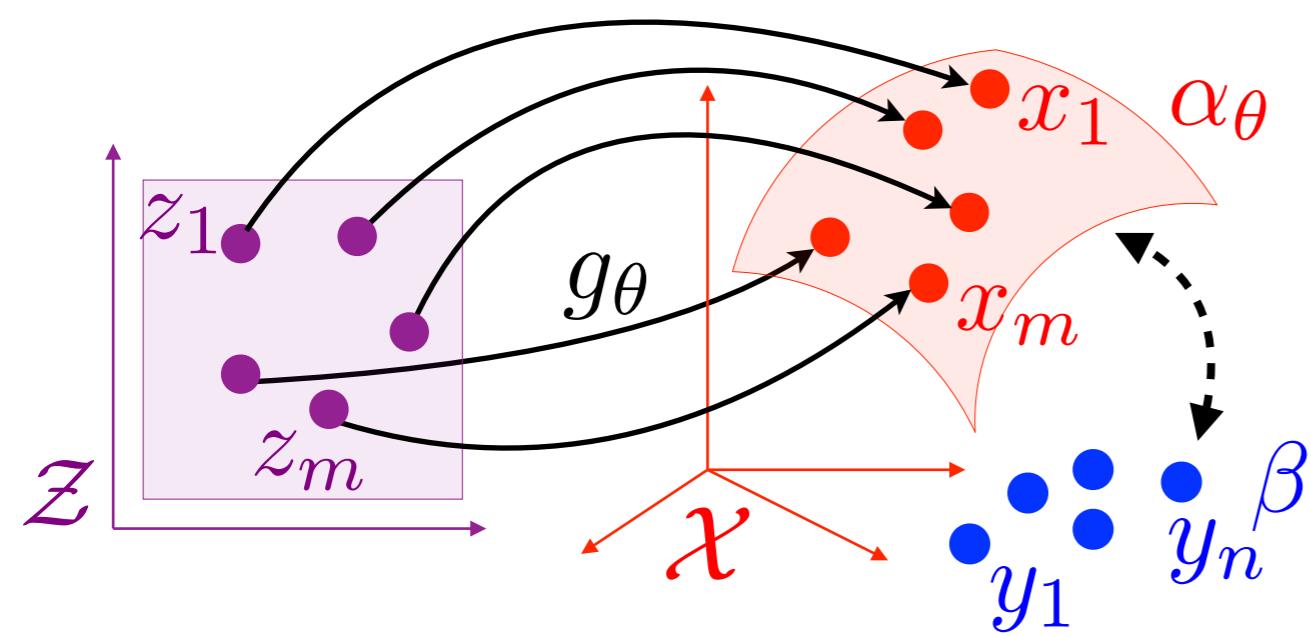
$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon,p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

$$\theta \leftarrow \theta - \tau \nabla \hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon,p}^p\left(\frac{1}{m} \sum_i \delta_{g_\theta(z_i)}, \beta\right)$$

# Training Architecture

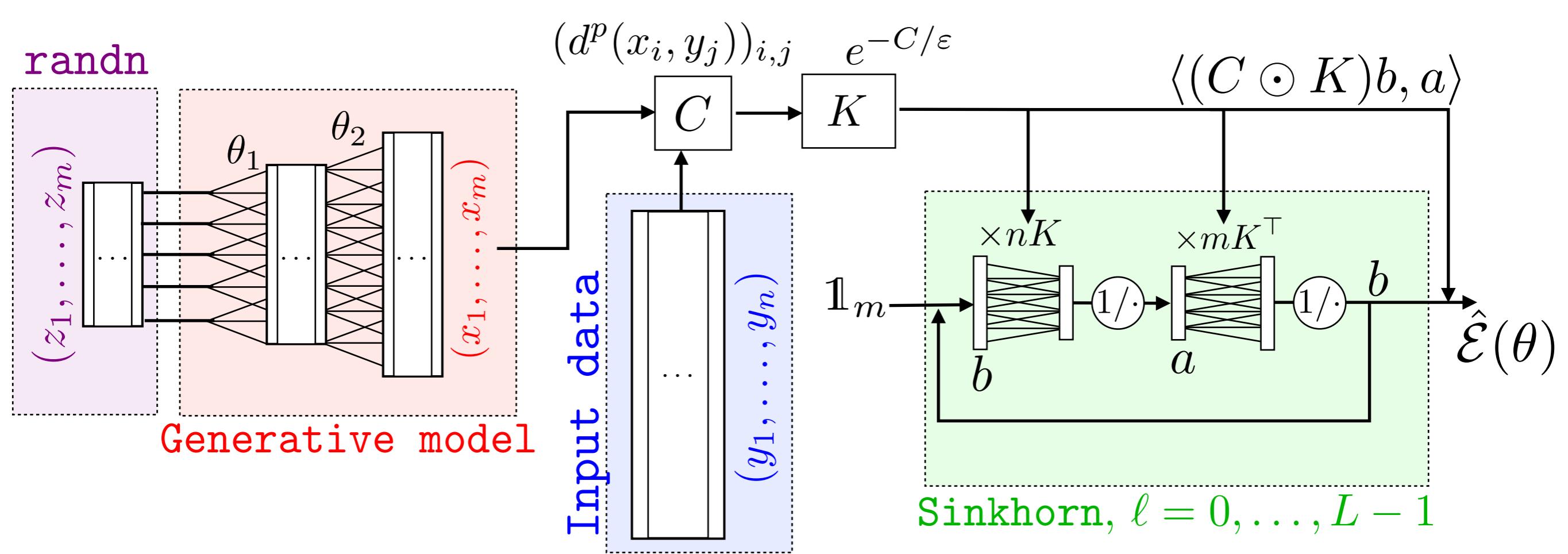


$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon,p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

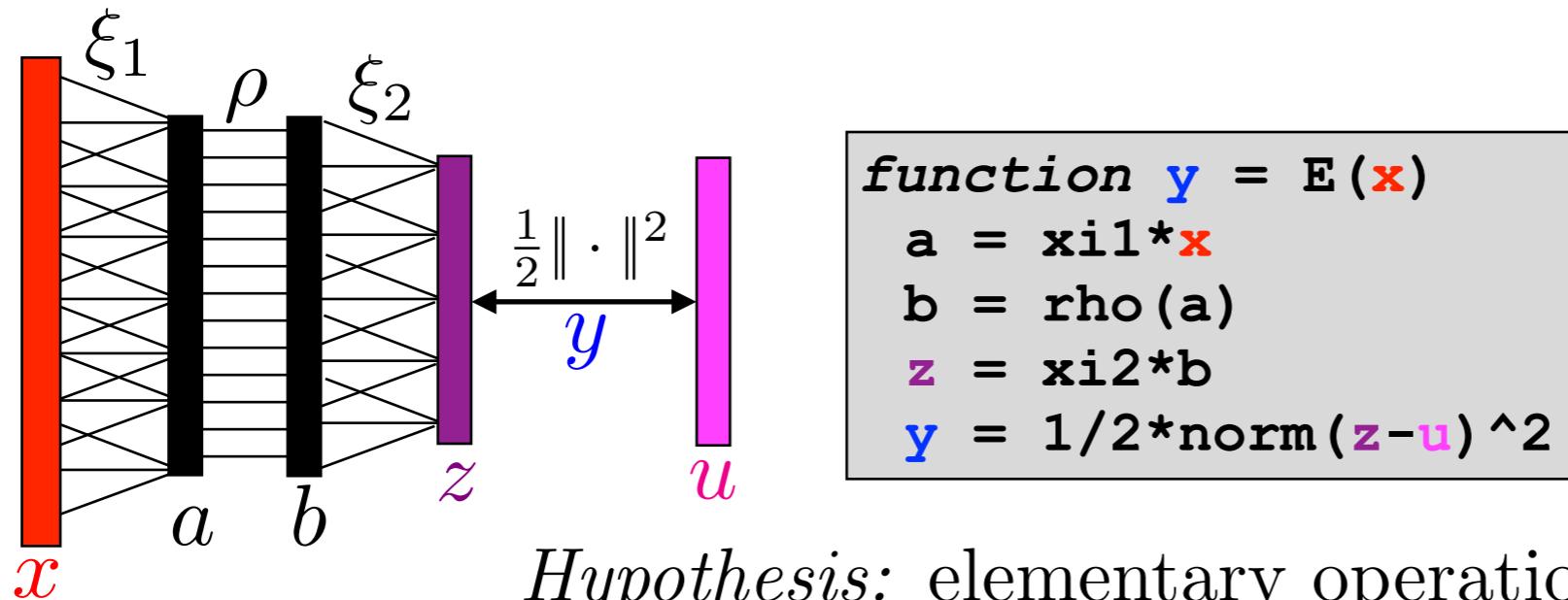
$$\theta \leftarrow \theta - \tau \nabla \hat{\mathcal{E}}(\theta)$$

$$\hat{\mathcal{E}}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon,p}^p\left(\frac{1}{m} \sum_i \delta_{g_\theta(z_i)}, \beta\right)$$



# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.

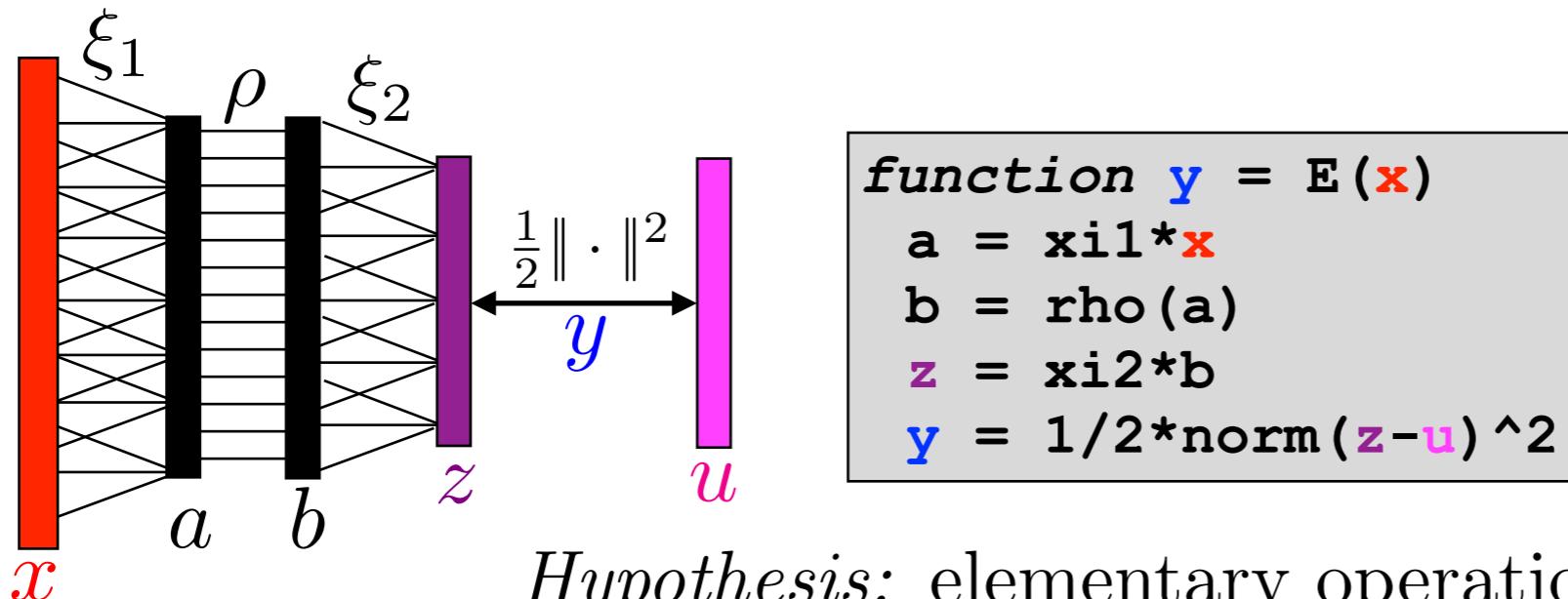


*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.



*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

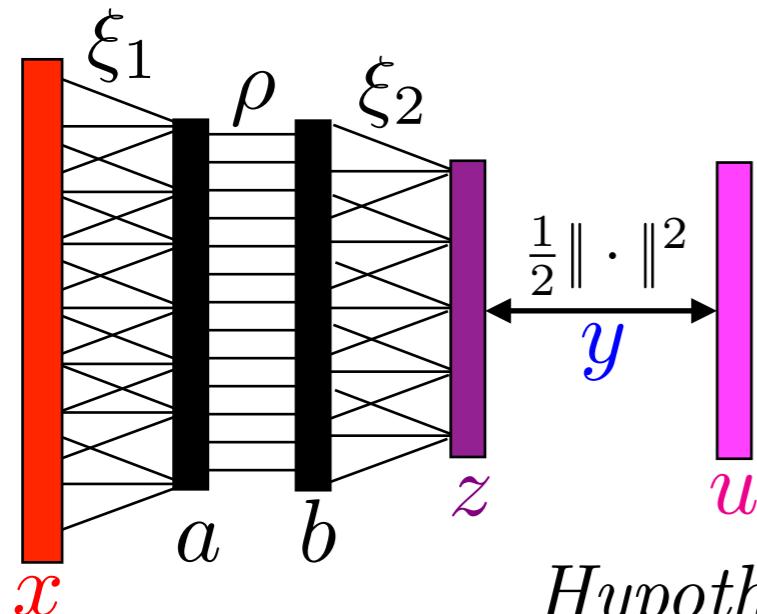
Finite differences:

$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \dots, \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$$

$K(n+1)$  operations, intractable for large  $n$ .

# Automatic Differentiation

**Setup:**  $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$  computable in  $K$  operations.



```
function y = E(x)
    a = xi1*x
    b = rho(a)
    z = xi2*b
    y = 1/2*norm(z-u)^2
```

```
function dx = nablaE(x)
    dz = z-u
    db = xi2'*dz
    da = diag(dphi(a)) * db
    dx = xi1'*da
```

*Hypothesis:* elementary operations ( $a \times b, \log(a), \sqrt{a} \dots$ )  
and their derivatives cost  $O(1)$ .

**Question:** What is the complexity of computing  $\nabla \mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ?

Finite differences:

$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon} (\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \dots, \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$$

$K(n+1)$  operations, intractable for large  $n$ .

*Theorem:* there is an algorithm to compute  $\nabla \mathcal{E}$   
in  $O(K)$  operations. [Seppo Linnainmaa, 1970]

This algorithm is reverse mode automatic differentiation



Seppo  
Linnainmaa

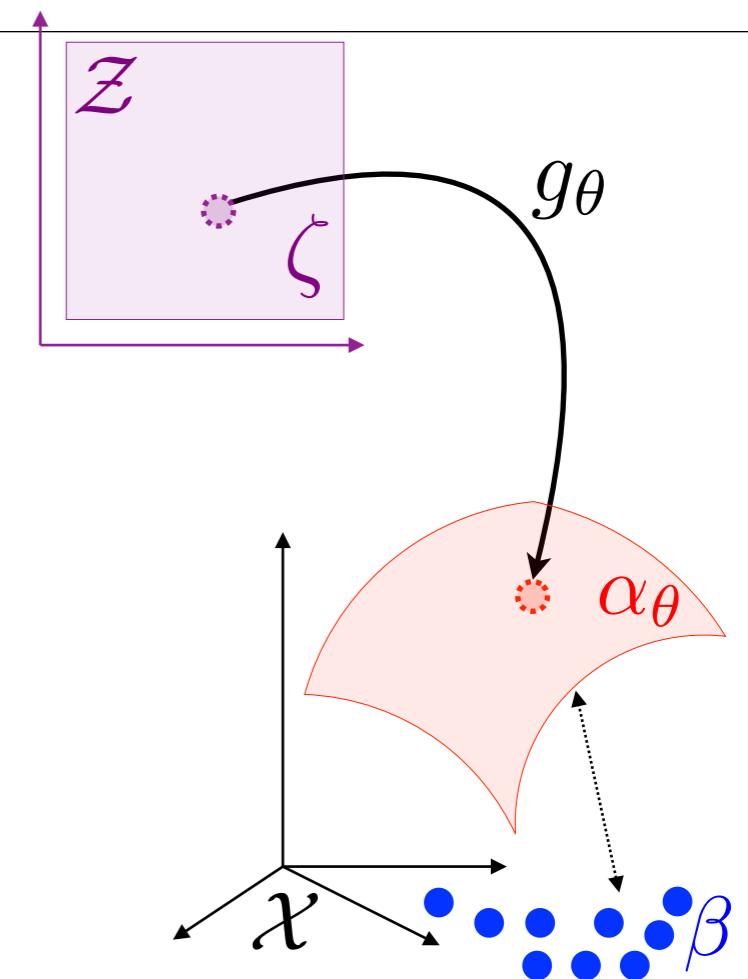
# Examples of Images Generation

Inputs  $\beta$

|   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|
| 3 | 4 | 2 | 1 | 9 | 5 | 6 | 2 | 1 |
| 8 | 9 | 1 | 2 | 5 | 0 | 0 | 6 | 6 |
| 6 | 7 | 0 | 1 | 6 | 3 | 6 | 3 | 7 |
| 3 | 7 | 7 | 9 | 4 | 6 | 6 | 1 | 8 |
| 2 | 9 | 3 | 4 | 3 | 9 | 8 | 7 | 2 |
| 1 | 5 | 9 | 8 | 3 | 6 | 5 | 7 | 2 |
| 9 | 3 | 1 | 9 | 1 | 5 | 8 | 0 | 8 |
| 5 | 6 | 2 | 6 | 8 | 5 | 8 | 8 | 9 |
| 3 | 7 | 7 | 0 | 9 | 4 | 8 | 5 | 4 |

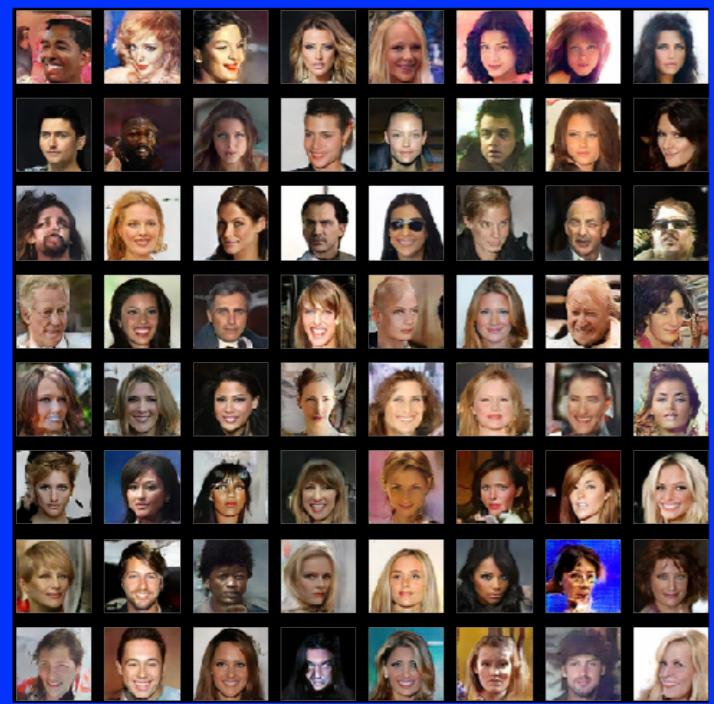
Generated  $\alpha_\theta$

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| 9 | 4 | 7 | 3 | 3 | 9 | 6 | 8 |
| 5 | 5 | 1 | 0 | 8 | 1 | 2 | 0 |
| 5 | 4 | 0 | 8 | 0 | 0 | 5 | 9 |
| 8 | 2 | 6 | 0 | 7 | 2 | 4 | 7 |
| 3 | 9 | 0 | 6 | 1 | 9 | 1 | 8 |
| 4 | 2 | 6 | 7 | 9 | 3 | 6 | 7 |
| 8 | 0 | 0 | 2 | 4 | 8 | 5 | 7 |
| 2 | 6 | 0 | 5 | 3 | 4 | 0 | 3 |

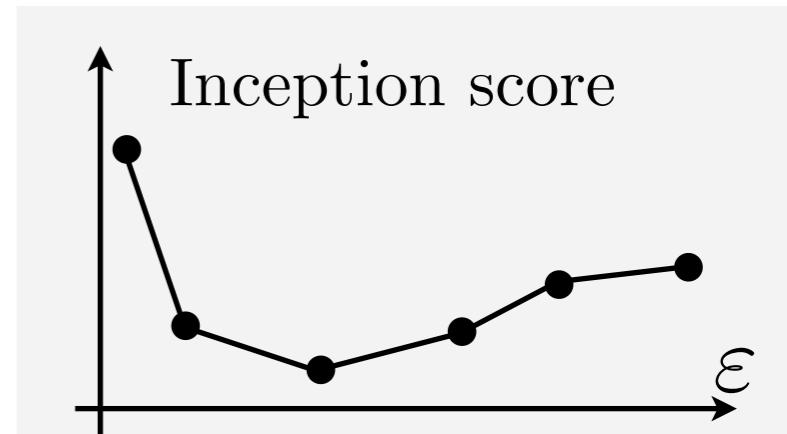
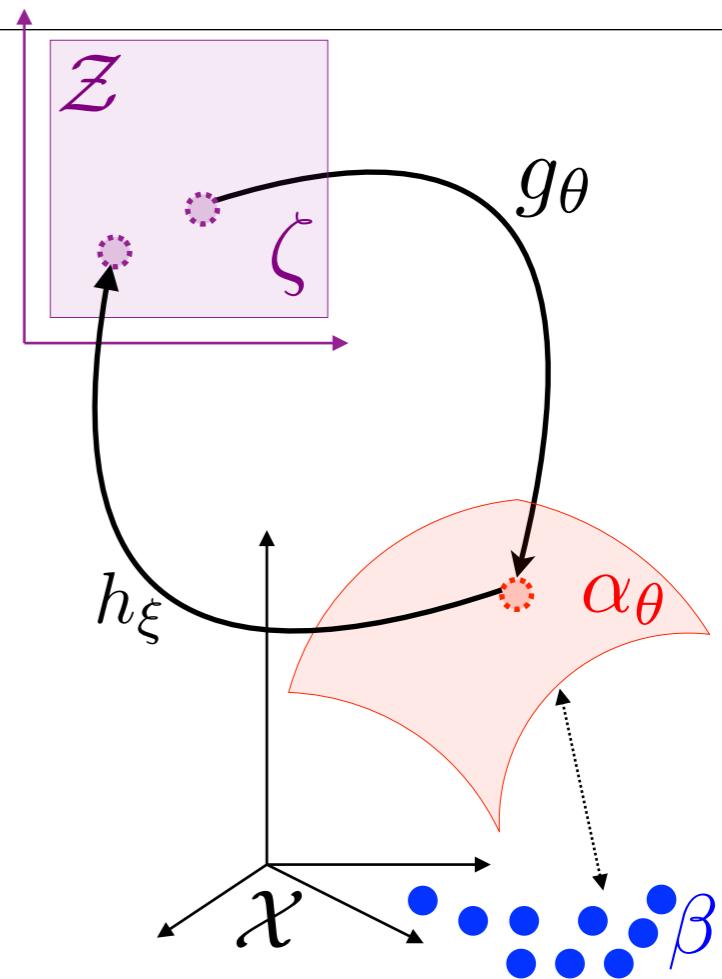
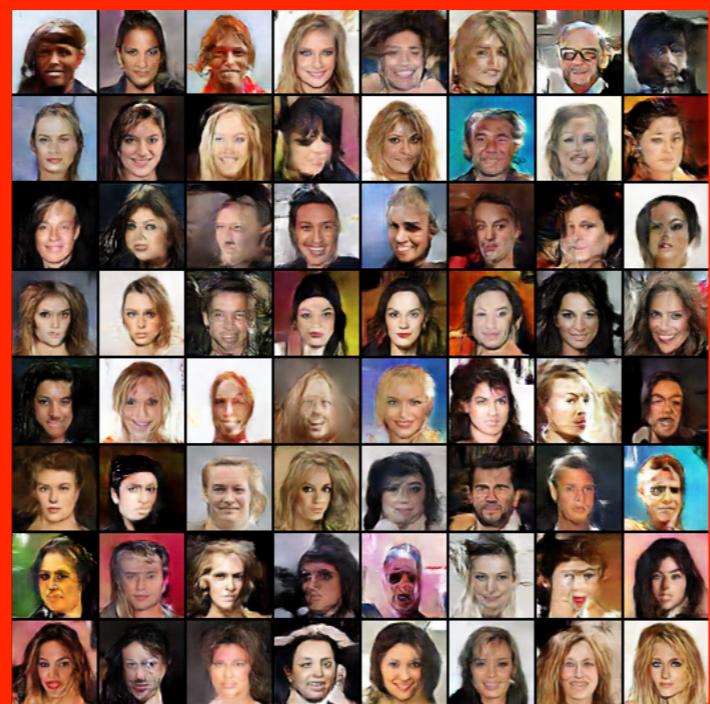


# Examples of Images Generation

Inputs  $\beta$

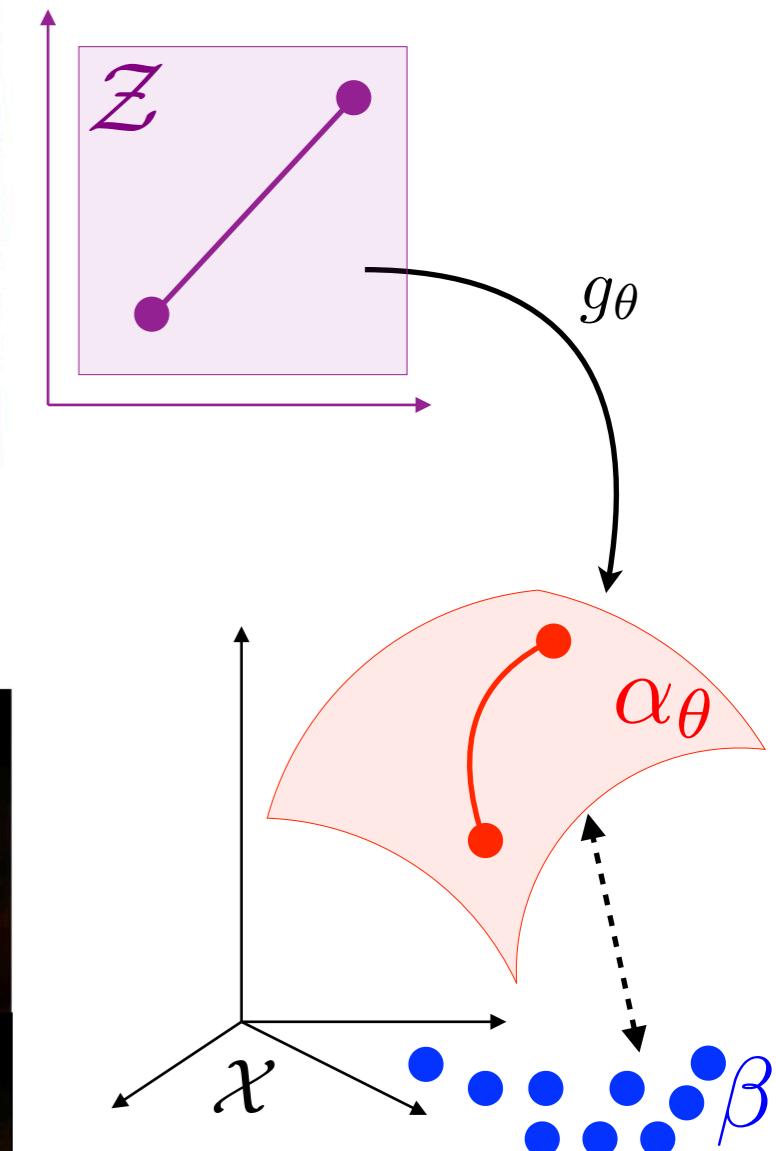


Generated  $\alpha_\theta$



- Need to learn the metric  $d(x, y) = \|h_\xi(x) - h_\xi(y)\|$  (GANs)
- Influence of  $\epsilon$ ?
- Performance evaluation of generative models is an open problem.

# Generative Adversarial Networks



*Progressive Growing of GANs for Improved Quality, Stability, and Variation*  
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018

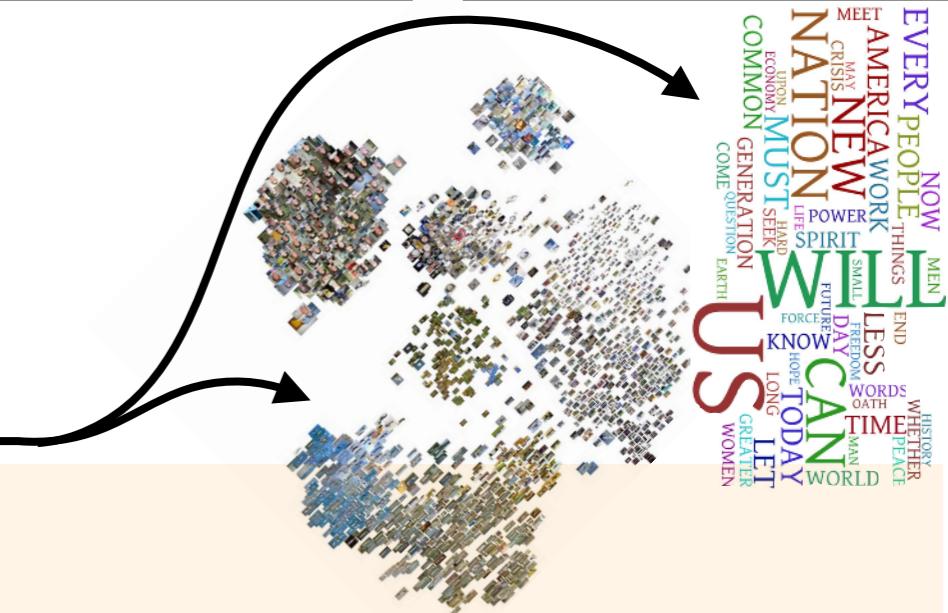


*Progressive Growing of GANs for Improved Quality, Stability, and Variation*  
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018



*Progressive Growing of GANs for Improved Quality, Stability, and Variation*  
Tero Karras, Timo Aila, Samuli Laine, Jaakko Lehtinen, ICLR 2018

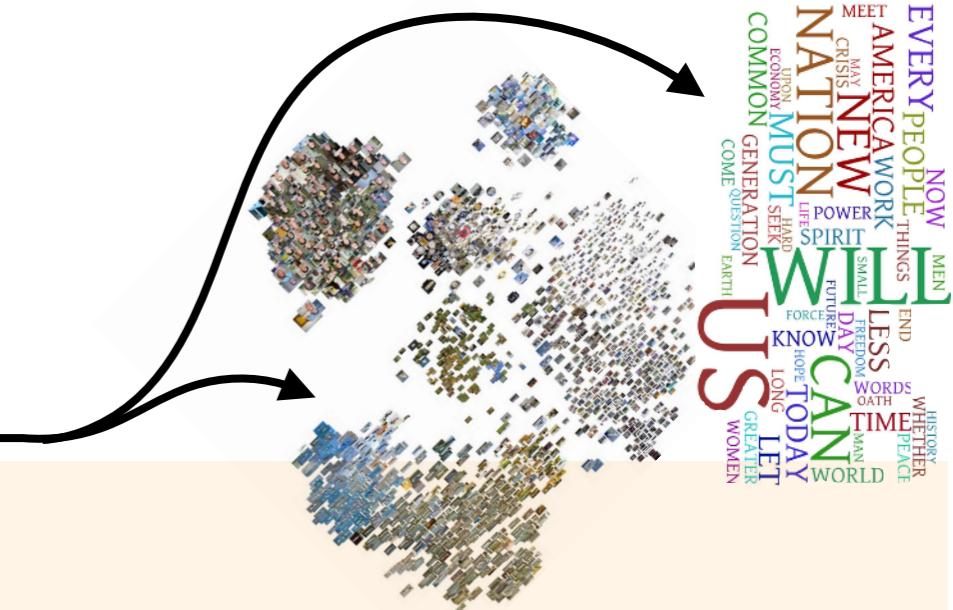
# Open Problems



Toward high-dimensional OT:

- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

# Open Problems

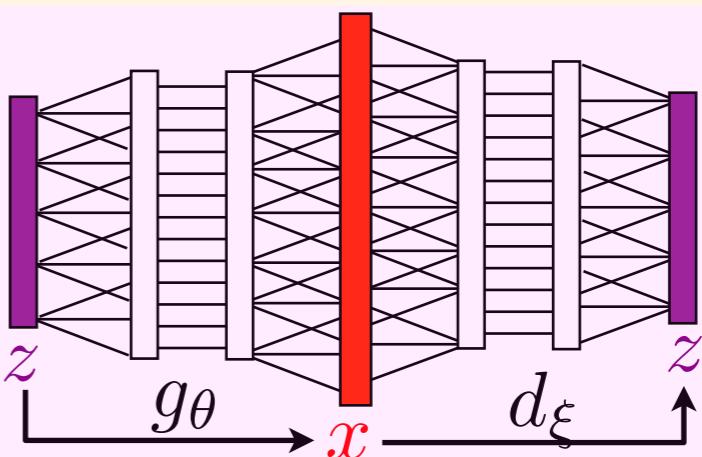


Toward high-dimensional OT:

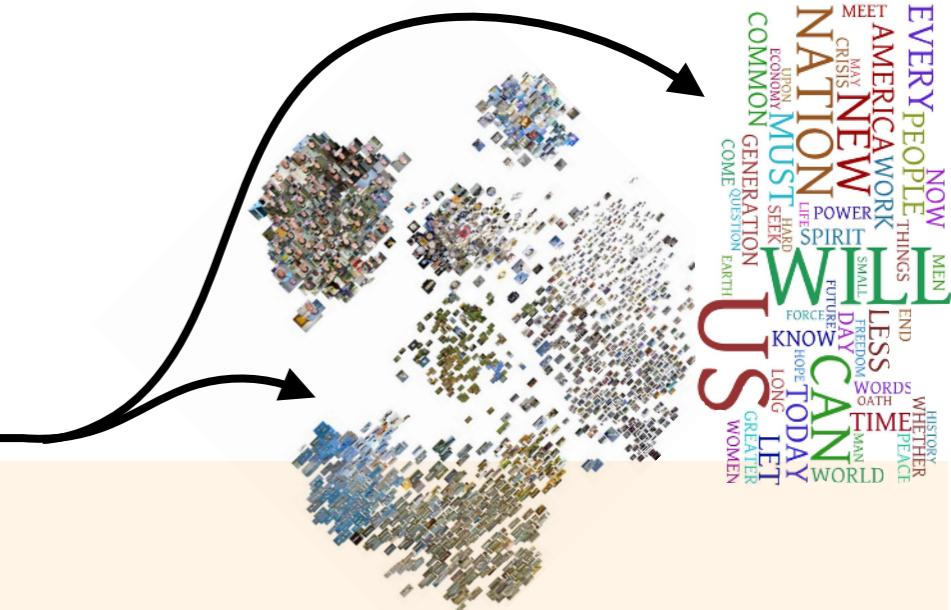
- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

Metric learning for OT:

- Adversarial training to leverage multi scale priors?



# Open Problems

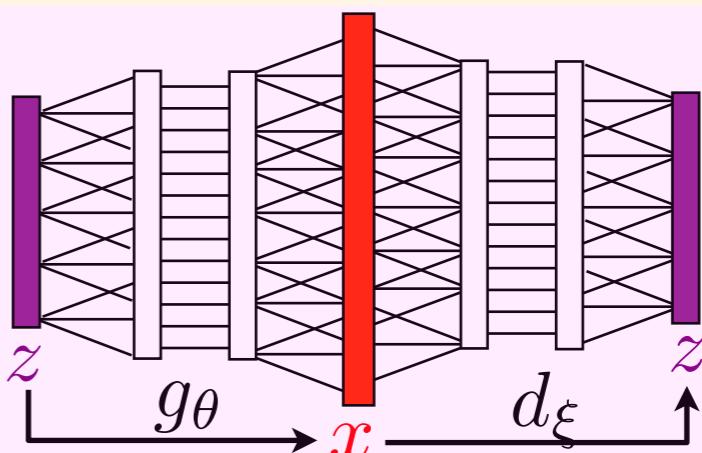


Toward high-dimensional OT:

- Scalable geometrical loss functions in high dimension?
- Performance quality measures for unsupervised learning?

Metric learning for OT:

- Adversarial training to leverage multi scale priors?



Beyond comparing measures:

- Learning for surfaces, graphs, metric spaces?
- Using Gromov-Wasserstein geometry?

