

# Scaling Optimal Transport for Machine Learning



**Gabriel Peyré**



<https://optimaltransport.github.io>

Home

BOOK

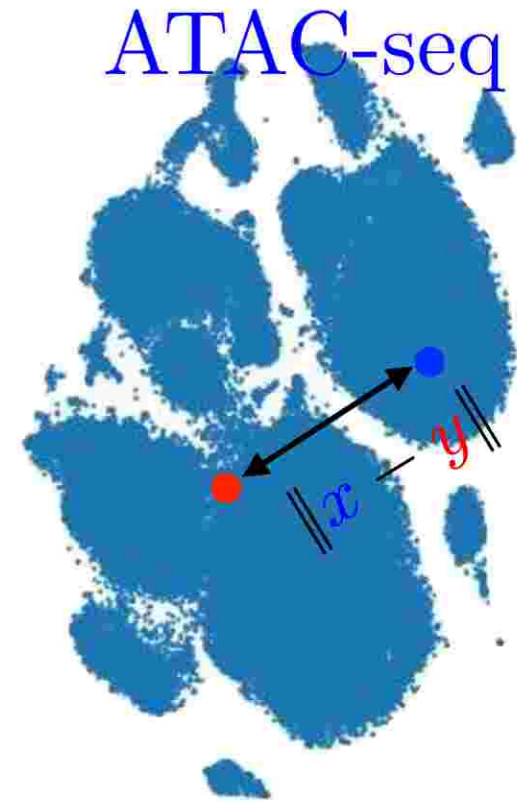
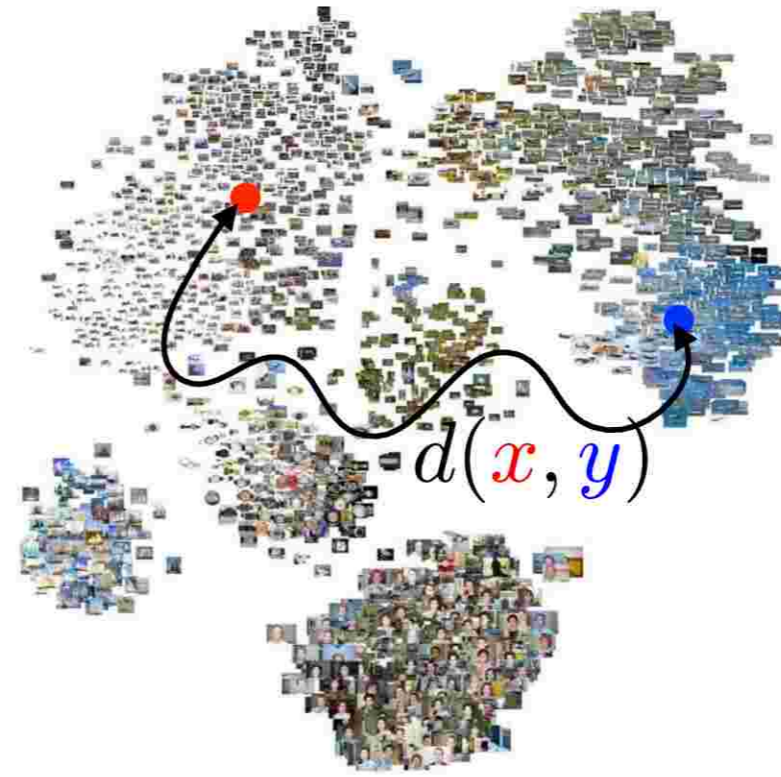
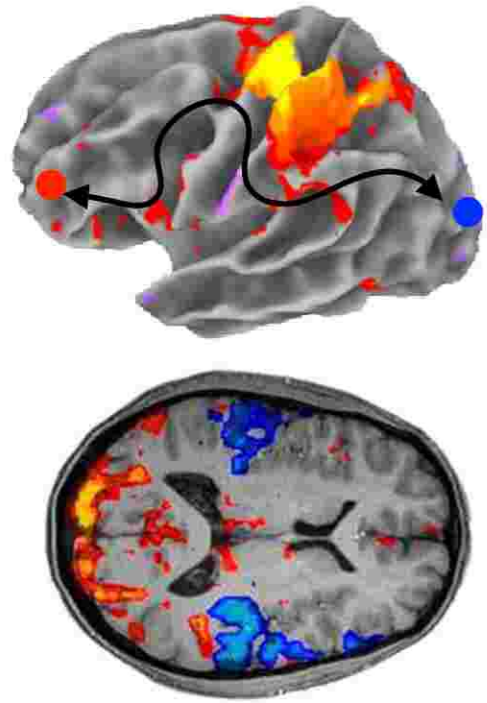
CODE

SLIDES

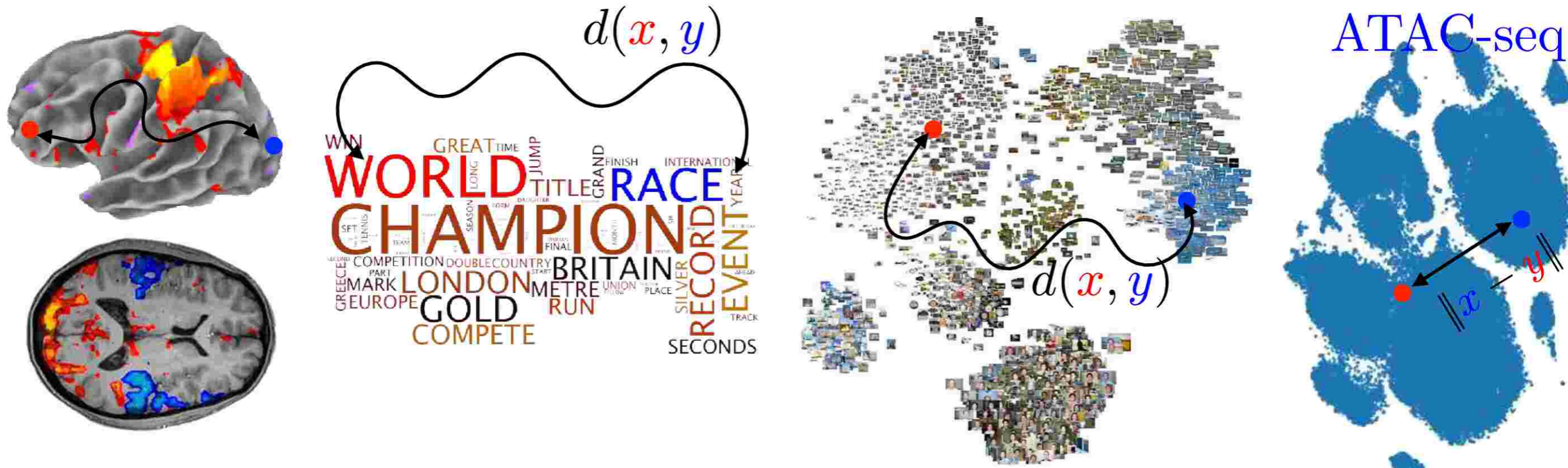
# Computational Optimal Transport

---

# Comparing Distributions for Learning



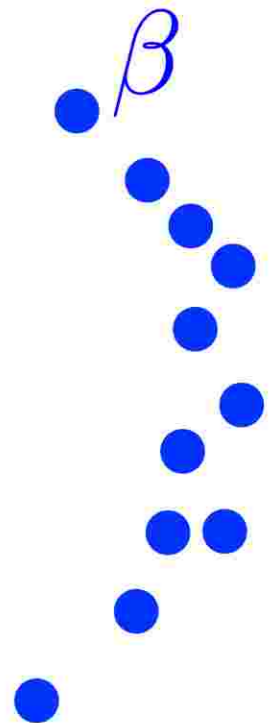
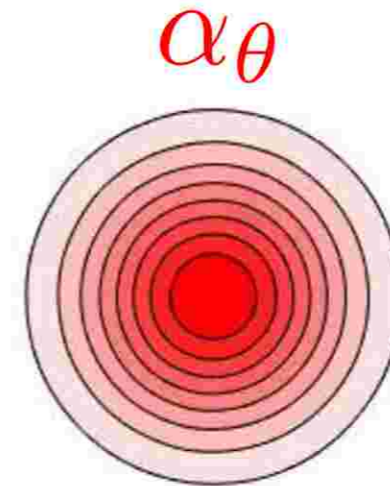
# Comparing Distributions for Learning



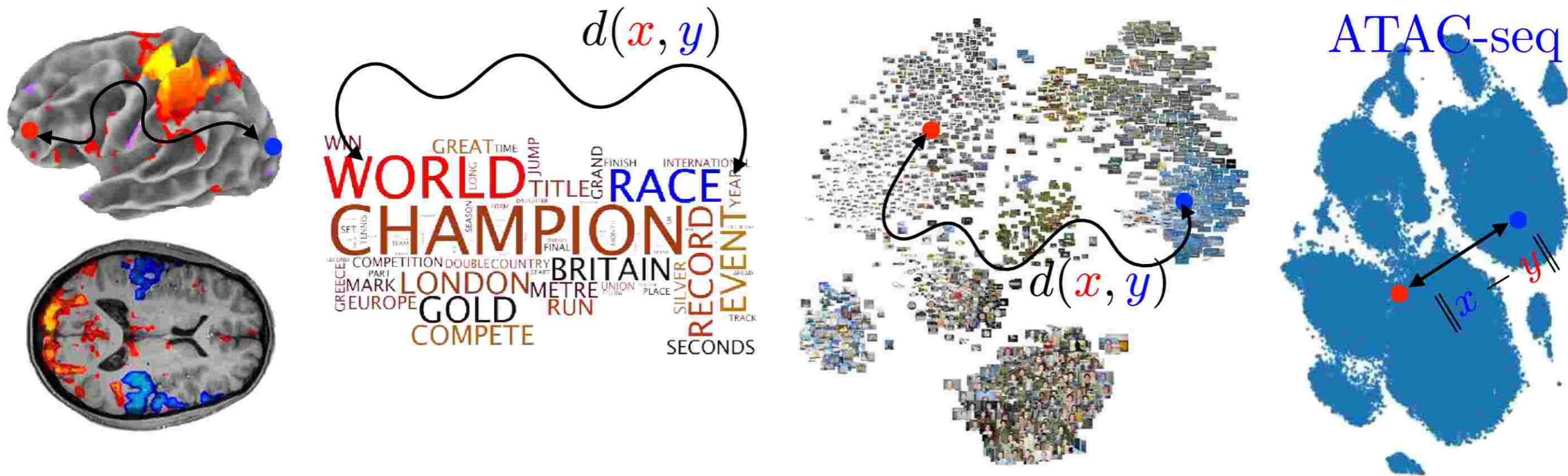
## Unsupervised learning

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

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



# Comparing Distributions for Learning

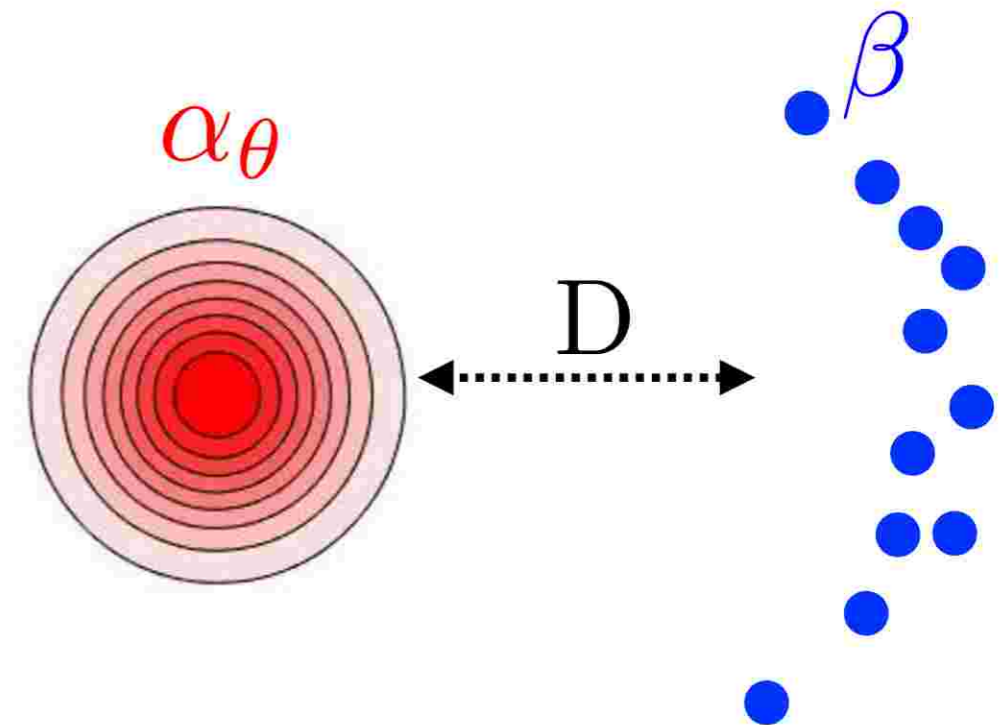


## Unsupervised learning

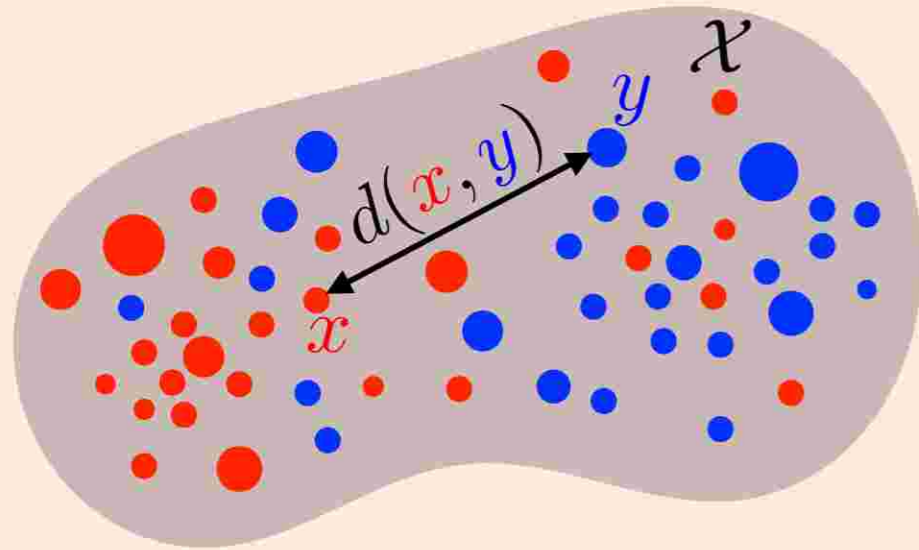
Observations:  $\beta \stackrel{\text{def.}}{=} \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

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

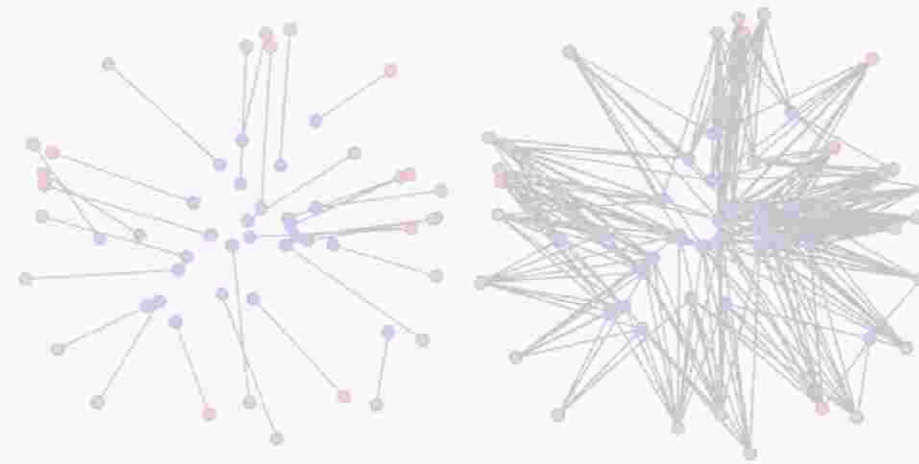
Density fitting:  $\min_{\theta} D(\alpha_\theta, \beta)$   
→ takes into account a metric  $d$ .



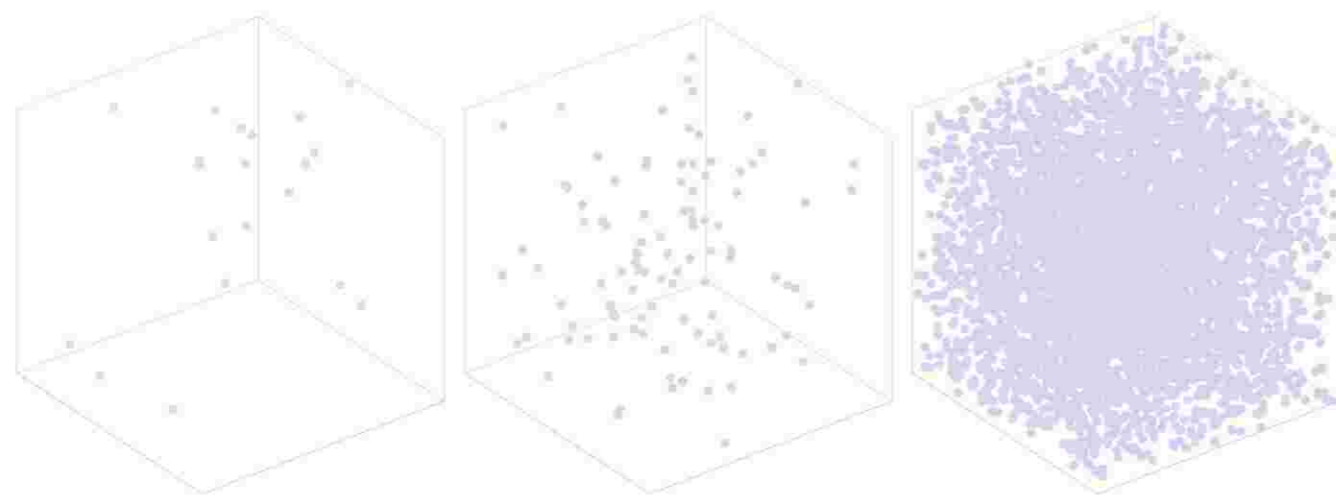
# 1. Optimal Transport



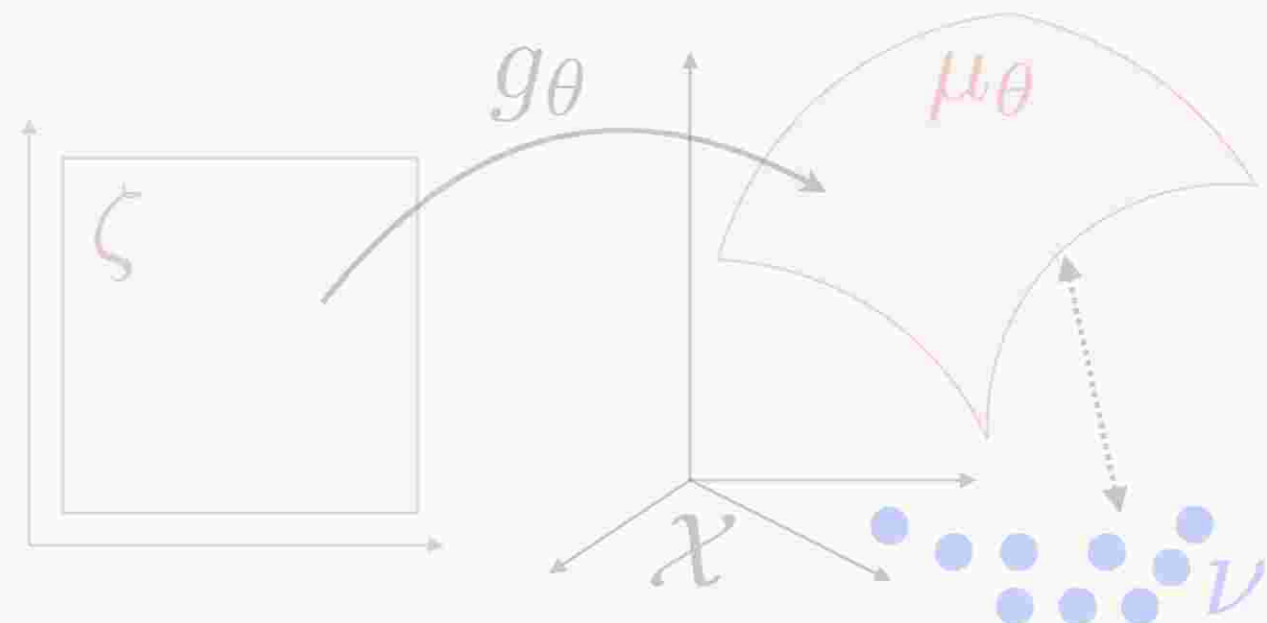
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



# 4. Application to Generative Models

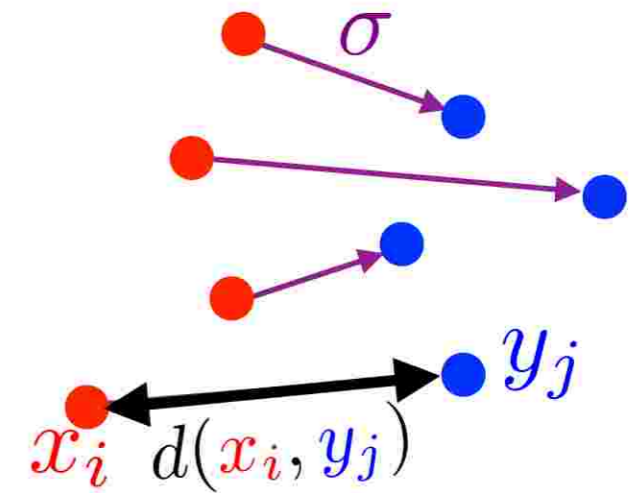


# Monge's Problem

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

Permutation:

$$\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$$



Monge optimal matching:

$$D(X, Y) = \min_{\sigma} \sum_{i=1}^n d(x_i, y_{\sigma(i)})$$



[Monge 1784]

M É M O I R E  
SUR LA  
THÉORIE DES DÉBLAIS  
ET DES REMBLAIS.  
Par M. M O N G E.

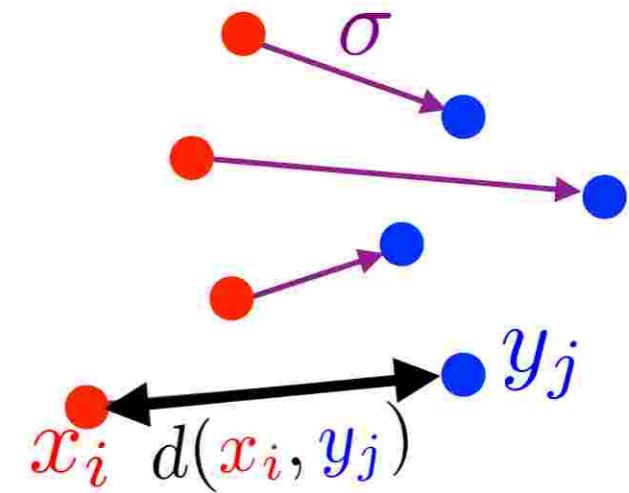
Lorsqu'on doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport. Le prix du transport d'une molécule étant, toutes choses d'ailleurs égales, proportionnel à son poids & à l'espace qu'on lui fait parcourir, & par conséquent le prix du transport total devant être proportionnel à la somme des produits des molécules multipliées chacune par l'espace parcouru, il s'ensuit que le déblai & le remblai étant donnés de figure & de position, il n'est pas indifférent que telle molécule du déblai soit transportée dans tel ou tel autre endroit du remblai, mais qu'il y a une certaine distribution à faire des molécules du premier dans le second, d'après laquelle la somme de ces produits fera la moindre possible, & le prix du transport total fera un *minimum*.

# Monge's Problem

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

Permutation:

$$\sigma : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$$



Monge optimal matching:

$$D(X, Y) = \min_{\sigma} \sum_{i=1}^n d(x_i, y_{\sigma(i)})$$



[Monge 1784]

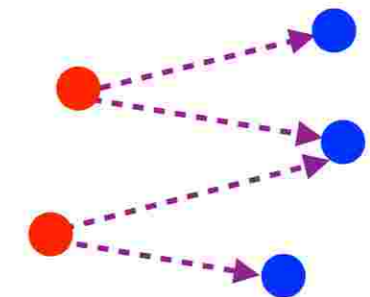
M É M O I R E  
SUR LA  
THÉORIE DES DÉBLAIS  
ET DES REMBLAIS.

Par M. M O N G E.

Lorsqu'on doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport. Le prix du transport d'une molécule étant, toutes choses d'ailleurs égales, proportionnel à son poids & à l'espace qu'on lui fait parcourir, & par conséquent le prix du transport total devant être proportionnel à la somme des produits des molécules multipliées chacune par l'espace parcouru, il s'ensuit que le déblai & le remblai étant donnés de figure & de position, il n'est pas indifférent que telle molécule du déblai soit transportée dans tel ou tel autre endroit du remblai, mais qu'il y a une certaine distribution à faire des molécules du premier dans le second, d'après laquelle la somme de ces produits fera la moindre possible, & le prix du transport total fera un *minimum*.

→ Seems intractable:  $n!$  possibilities.

→ Different number of points?





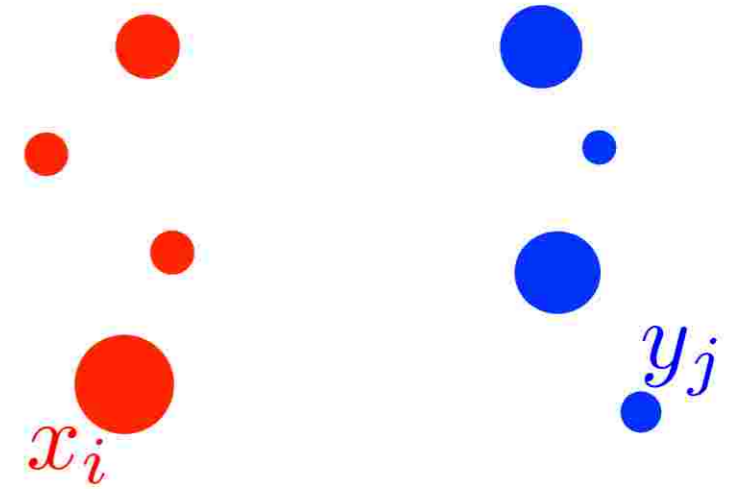
# Kantorovitch's Formulation

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

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

Discrete distributions:

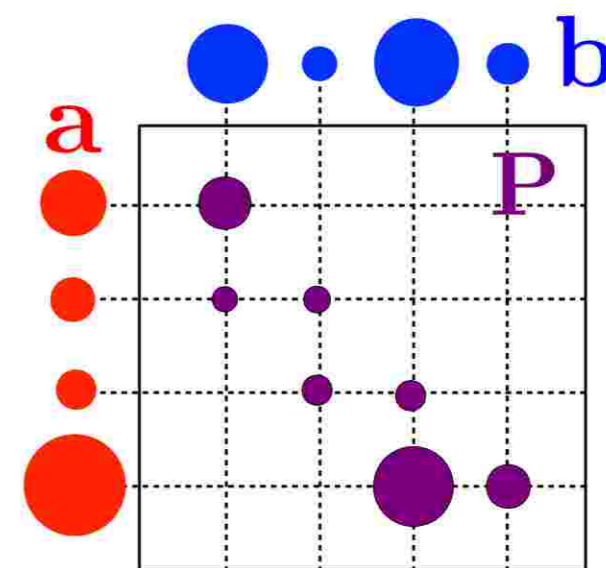
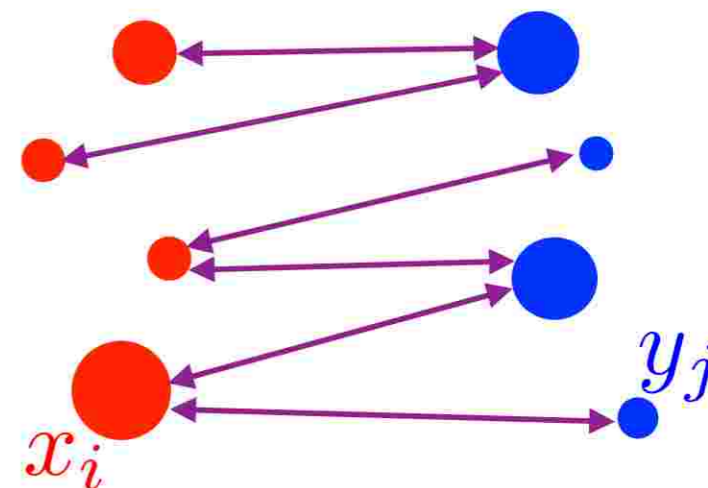
$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$$

$$\beta = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$$

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:

$$\sum_j \mathbf{P}_{i,j} = \mathbf{a}_i$$

$$\sum_i \mathbf{P}_{i,j} = \mathbf{b}_j$$

$$\mathbf{P} \geq 0, \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_n = \mathbf{b}$$

# Kantorovitch's Formulation

Discrete distributions:

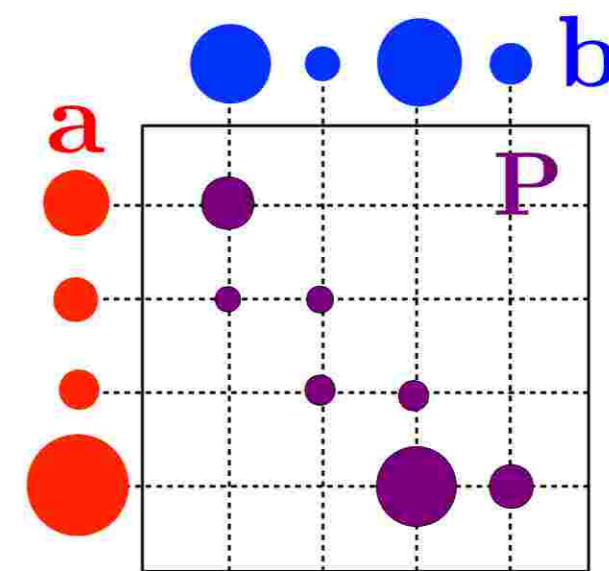
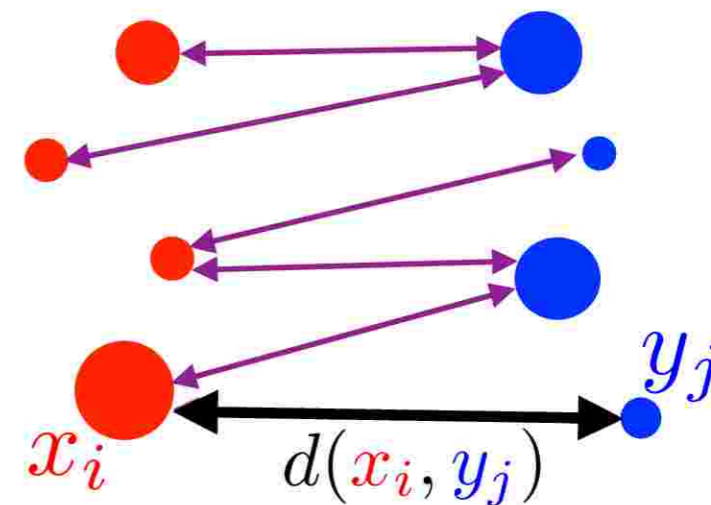
$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$$

$$\beta = \sum_{j=1}^m \mathbf{b}_j \delta_{y_j}$$

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

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

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



Leonid Kantorovitch

George Dantzig

Couplings:

$$\sum_j P_{i,j} = \mathbf{a}_i$$

$$\sum_i P_{i,j} = \mathbf{b}_j$$

[Kantorovich 1942]

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(x_i, y_j)^p P_{i,j} ; \mathbf{P} \geq 0, \mathbf{P} \mathbf{1}_m = \mathbf{a}, \mathbf{P}^\top \mathbf{1}_n = \mathbf{b} \right\}$$



# Optimal Transport Distances

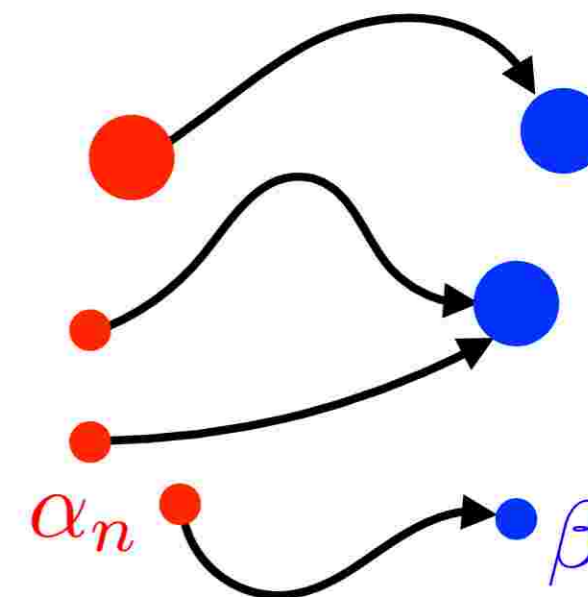
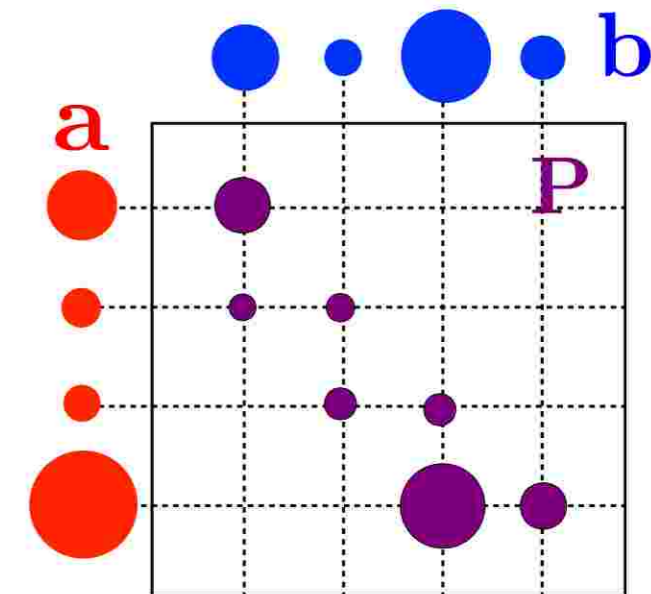
$$W_p(\alpha, \beta) \stackrel{\text{def.}}{=} \left( \min_{\mathbf{P} \mathbf{1}=\mathbf{a}, \mathbf{P}^\top \mathbf{1}=\mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} \right)^{\frac{1}{p}}$$

Convergence in law:  $\alpha_n \rightarrow \beta$

$$\Leftrightarrow \forall f \in \mathcal{C}(\mathcal{X}), \int_{\mathcal{X}} f d\alpha_n \rightarrow \int_{\mathcal{X}} f d\beta$$

*Theorem:*  $W_p$  is a distance and

$$\alpha_n \rightarrow \beta \Leftrightarrow W_p(\alpha_n, \beta) \rightarrow 0$$



# Optimal Transport Distances

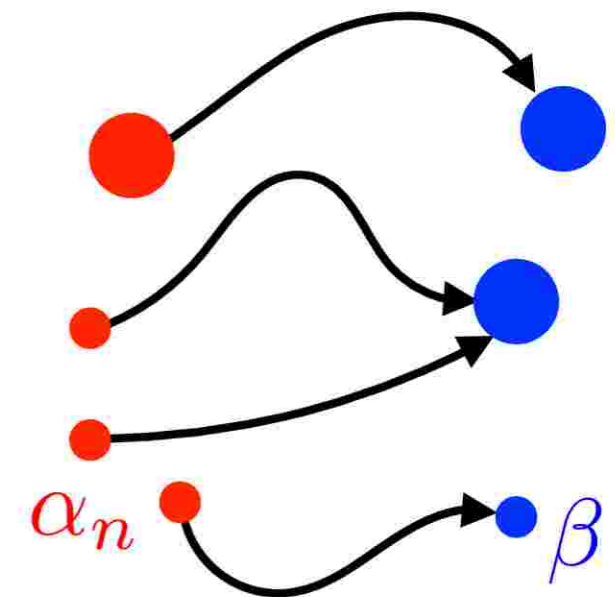
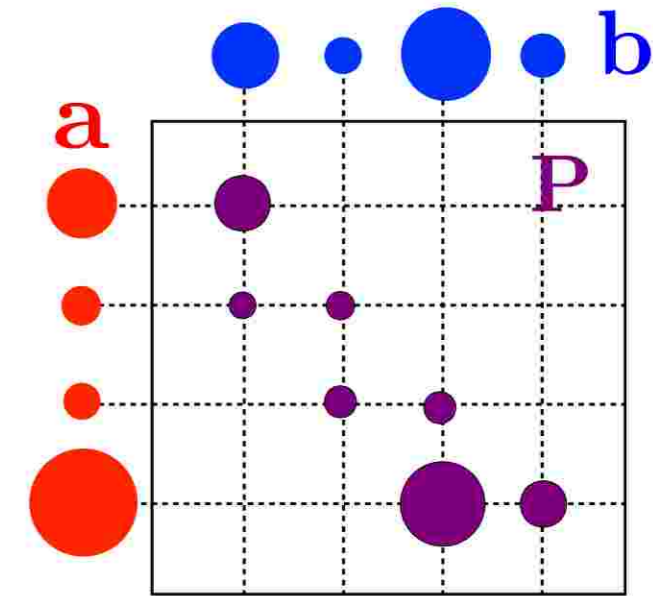
$$W_p(\alpha, \beta) \stackrel{\text{def.}}{=} \left( \min_{\mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} \right)^{\frac{1}{p}}$$

Convergence in law:  $\alpha_n \rightarrow \beta$

$$\Leftrightarrow \forall f \in \mathcal{C}(\mathcal{X}), \int_{\mathcal{X}} f d\alpha_n \rightarrow \int_{\mathcal{X}} f d\beta$$

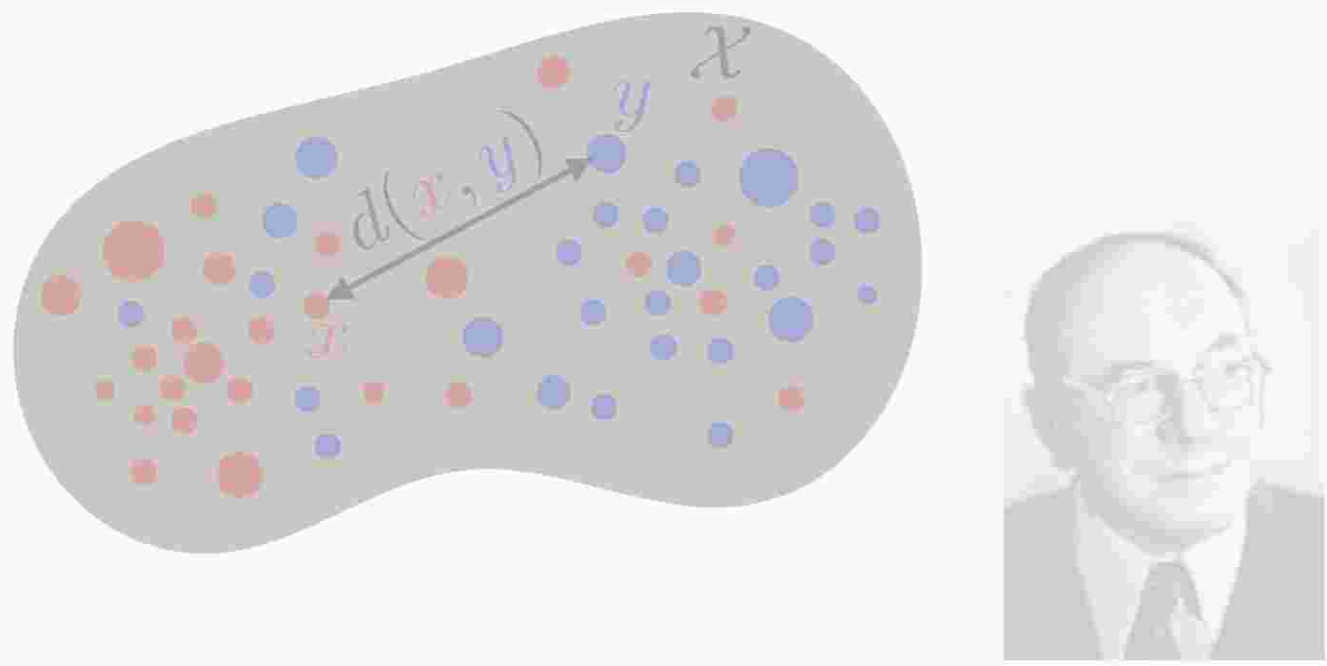
*Theorem:*  $W_p$  is a distance and

$$\alpha_n \rightarrow \beta \Leftrightarrow W_p(\alpha_n, \beta) \rightarrow 0$$

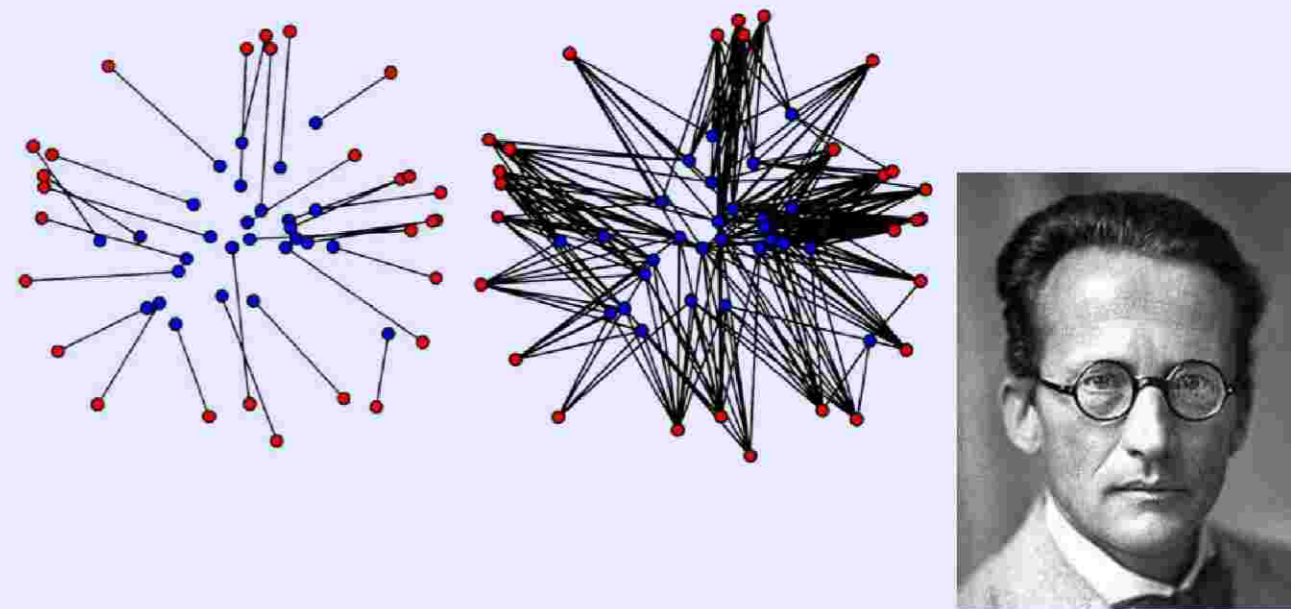


$$\|\delta_{x_n} - \delta_x\|_{\text{TV}} = 2 \quad \text{vs.} \quad W_p(\delta_{x_n}, \delta_x) = d(x_n, x)$$

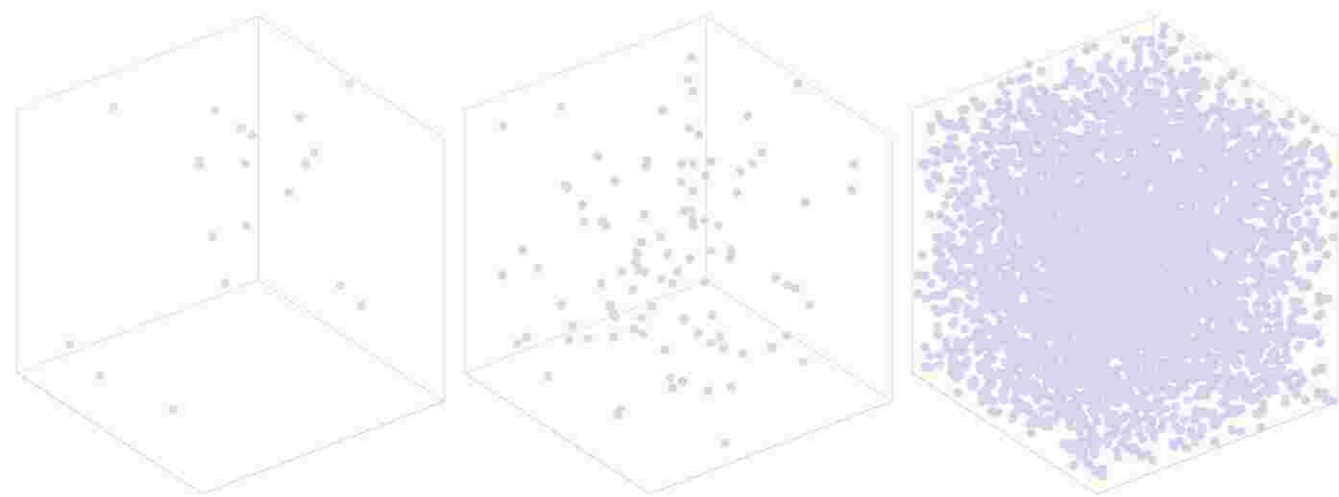
# 1. Optimal Transport



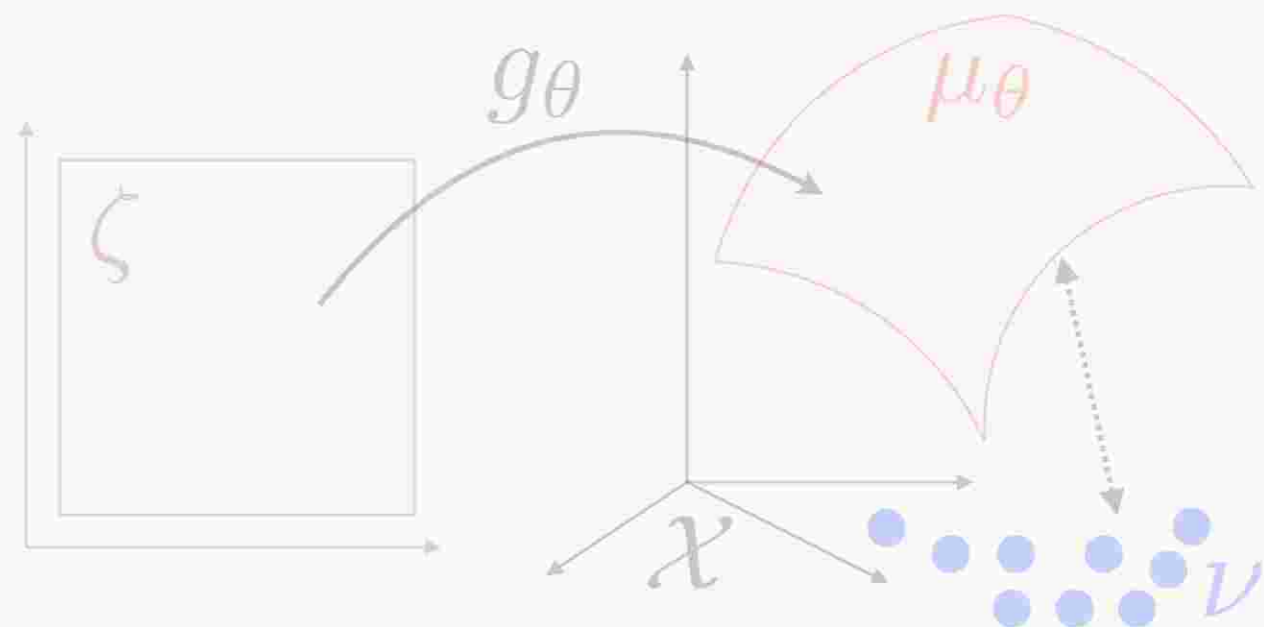
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



# 4. Application to Generative Models

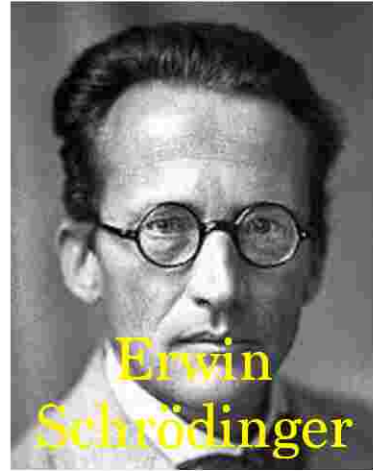


# Entropic Regularization

*Schrödinger's problem:*

[1931]

$$\min_{\mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j})$$



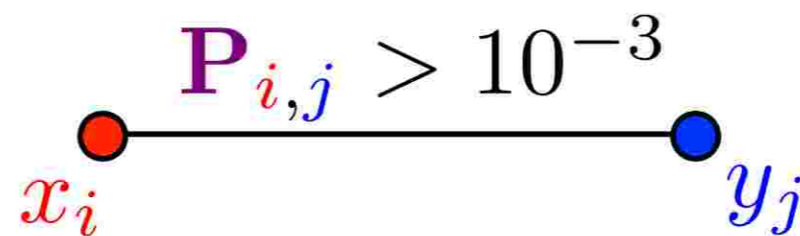
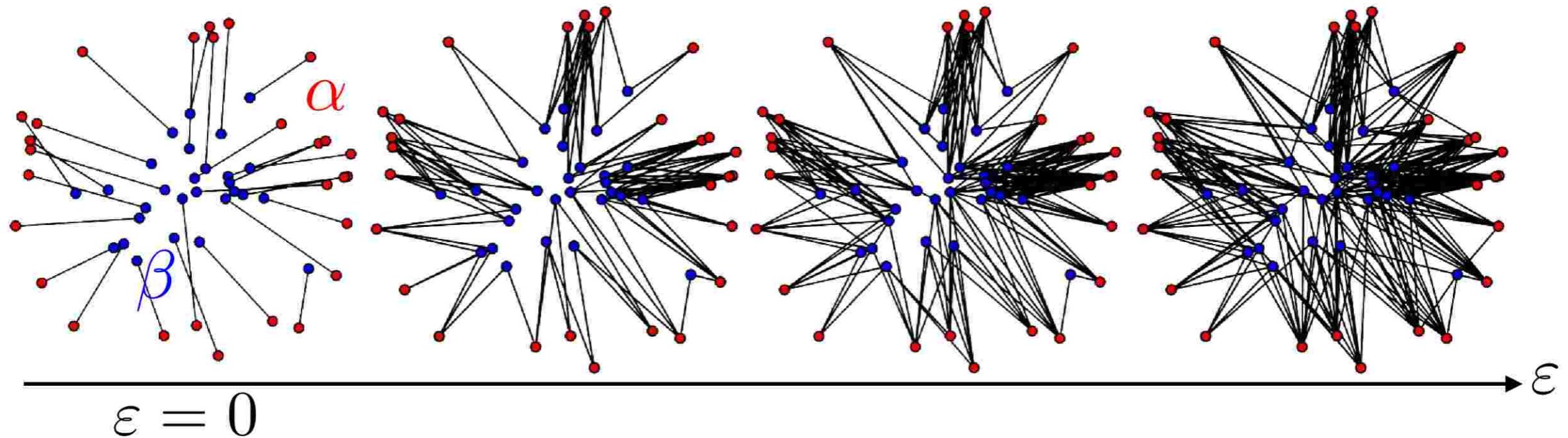
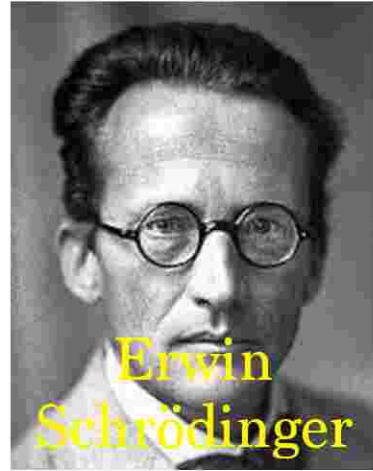


# Entropic Regularization

*Schrödinger's problem:*

[1931]

$$\min_{\mathbf{P} \mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j})$$



# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*

$$\mathbf{P} \text{ solution} \Leftrightarrow \begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases} \quad \mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$$

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*

$$\mathbf{P} \text{ solution} \Leftrightarrow \begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases} \quad \mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

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

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

# Sinkhorn's Algorithm

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\mathbf{P}$  solution  $\Leftrightarrow \begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

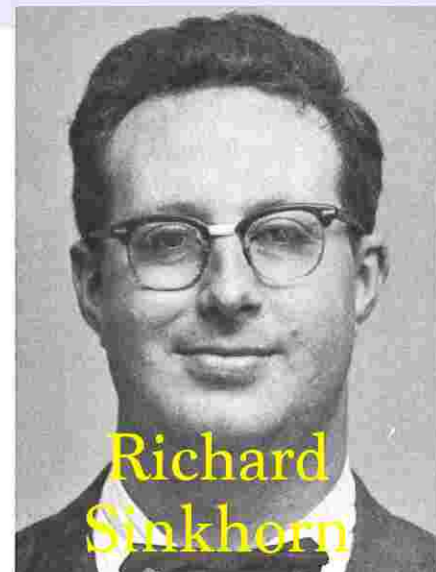
Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$       Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

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

$$\min_{\mathbf{P}} \left\{ \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} + \varepsilon \mathbf{P}_{i,j} \log(\mathbf{P}_{i,j}) ; \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \right\}$$

*Proposition:*  $\mathbf{P}$  solution  $\Leftrightarrow \begin{cases} \mathbf{P}\mathbf{1} = \mathbf{a}, \mathbf{P}^\top \mathbf{1} = \mathbf{b} \text{ and} \\ \exists \mathbf{u}, \mathbf{v}, \mathbf{P}_{i,j} = \mathbf{u}_i \mathbf{K}_{i,j} \mathbf{v}_j \end{cases}$   $\mathbf{K}_{i,j} \stackrel{\text{def.}}{=} e^{-\frac{d(x_i, y_j)^p}{\varepsilon}}$

$$\mathbf{P} = \text{diag}(\mathbf{u}) \mathbf{K} \text{diag}(\mathbf{v}) \implies \mathbf{a} = \mathbf{P}\mathbf{1} = \text{diag}(\mathbf{u})(\mathbf{K}\mathbf{v}) = \mathbf{u} \odot (\mathbf{K}\mathbf{v})$$

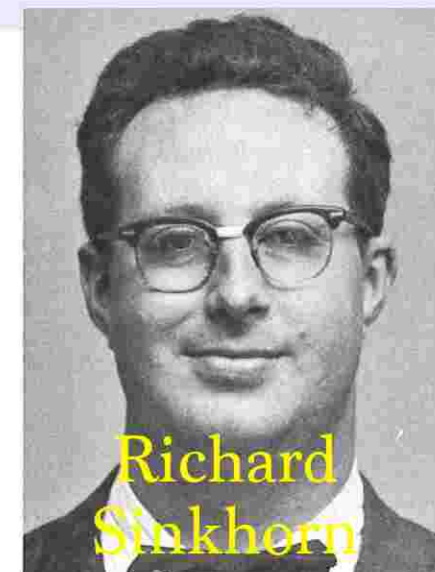
Row constraint:  $\mathbf{u} \odot (\mathbf{K}\mathbf{v}) = \mathbf{a}$       Col. constraint:  $\mathbf{v} \odot (\mathbf{K}^\top \mathbf{u}) = \mathbf{b}$

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.

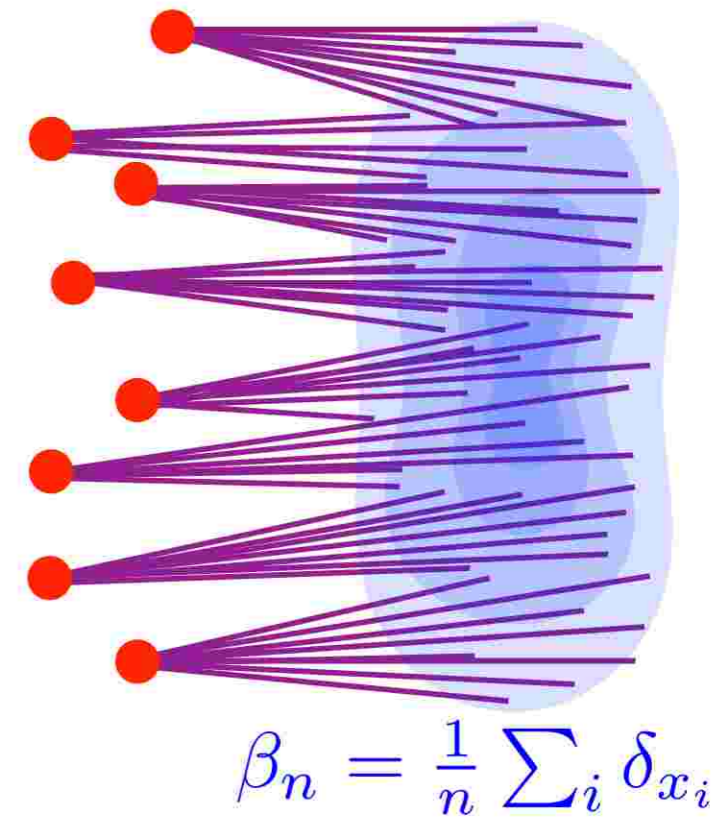
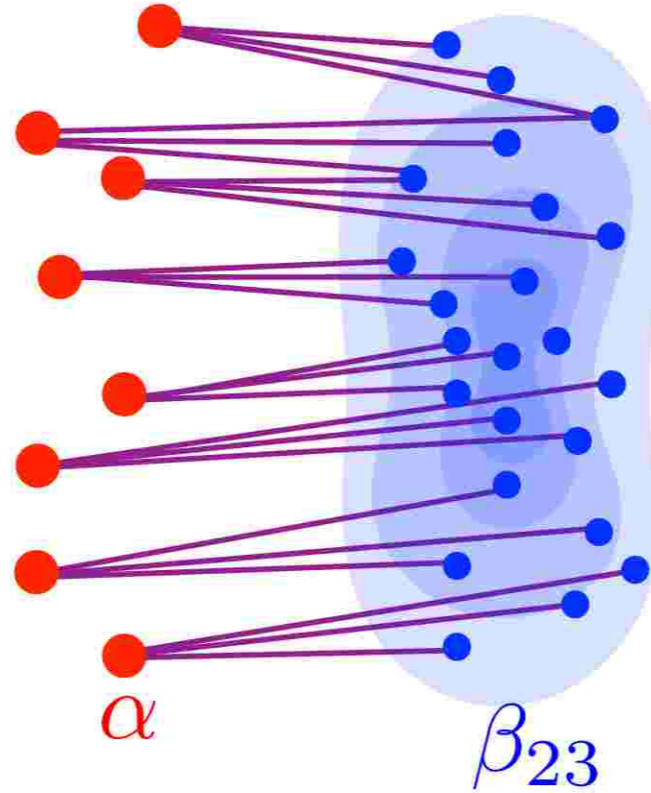
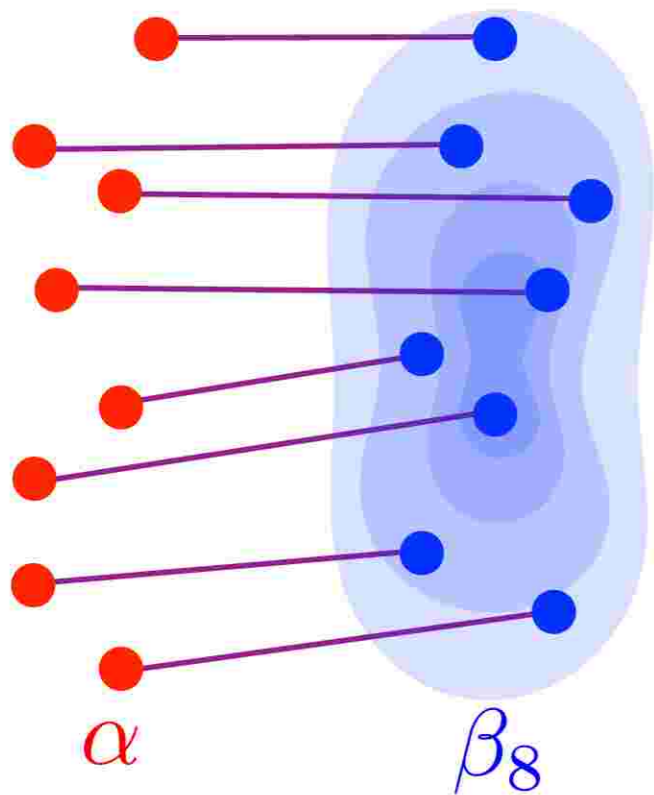


Matrix/vector multiplications:  $\rightarrow O(n^2/\varepsilon^2)$  complexity.

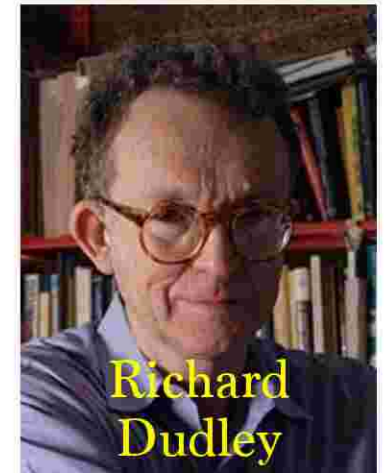
$\rightarrow$  Parallelizable on GPUs.

$\rightarrow$  Convolution on regular grids, separable kernels.

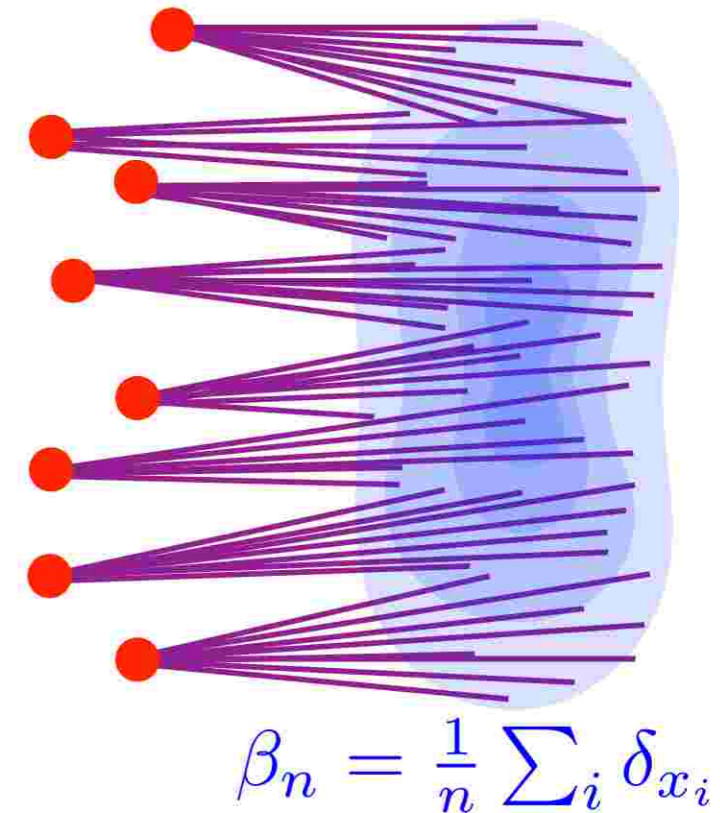
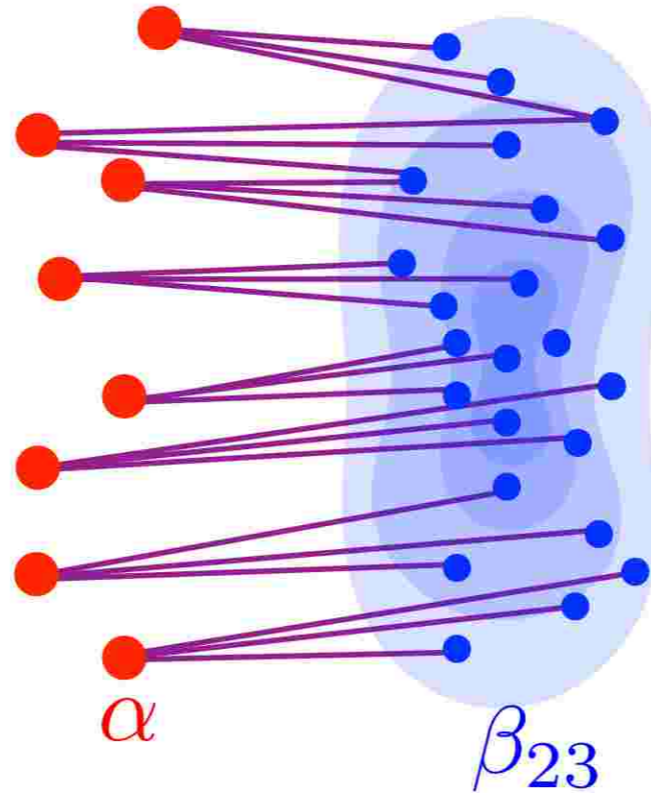
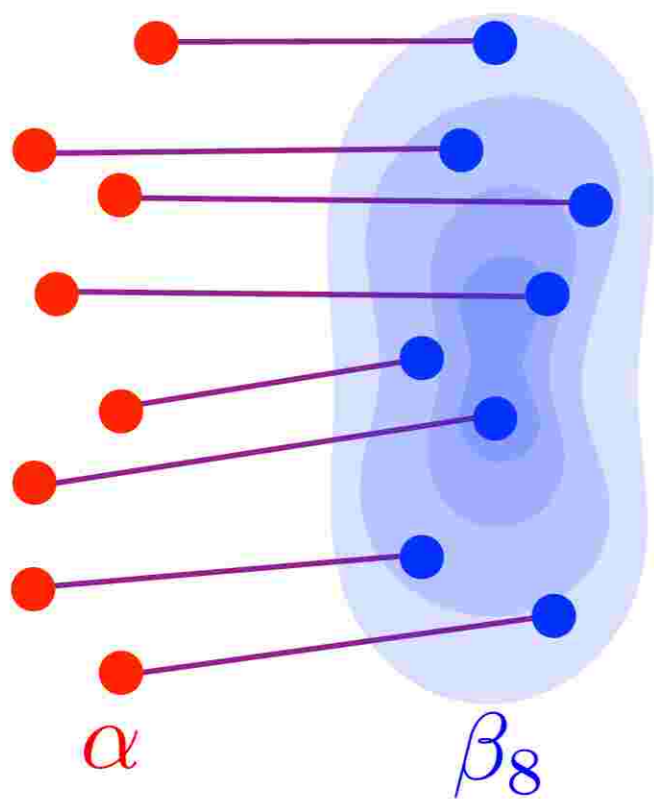
# The Curse of Dimensionality



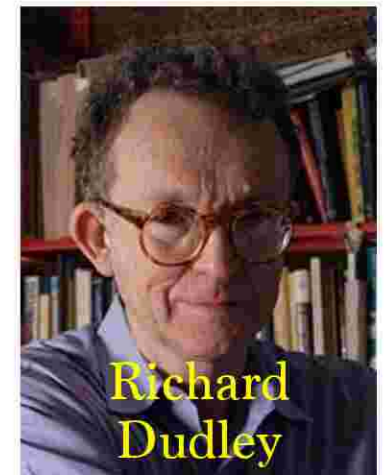
*Theorem:*  $\mathbb{E}|W_p(\alpha, \beta_n) - W_p(\alpha, \beta_\infty)| \leq \delta$   
[Dudley 1968] requires  $n \sim (1/\delta)^{\text{dimension}}$



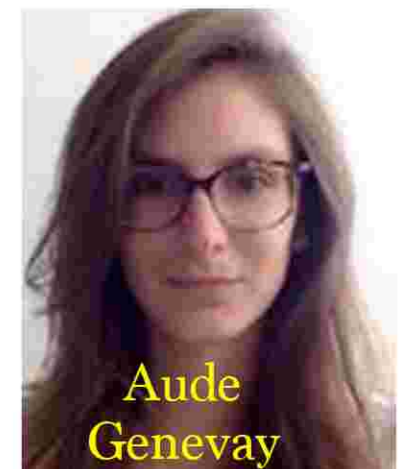
# The Curse of Dimensionality



*Theorem:*  $\mathbb{E}|W_p(\alpha, \beta_n) - W_p(\alpha, \beta_\infty)| \leq \delta$   
[Dudley 1968] requires  $n \sim (1/\delta)^{\text{dimension}}$



*Theorem:*  $\mathbb{E}|W_p^\varepsilon(\alpha, \beta_n) - W_p^\varepsilon(\alpha, \beta_\infty)| \leq \delta$   
[Genevay 2019] requires  $n \sim (1/\varepsilon)^{\text{dimension}} \times (1/\delta)^2$



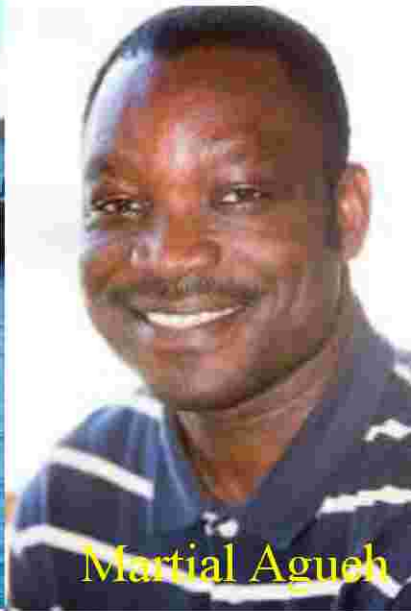
# Wasserstein Barycenters

Barycenters of measures  $(\alpha_s)_s$ :  $\sum_s \lambda_s = 1$

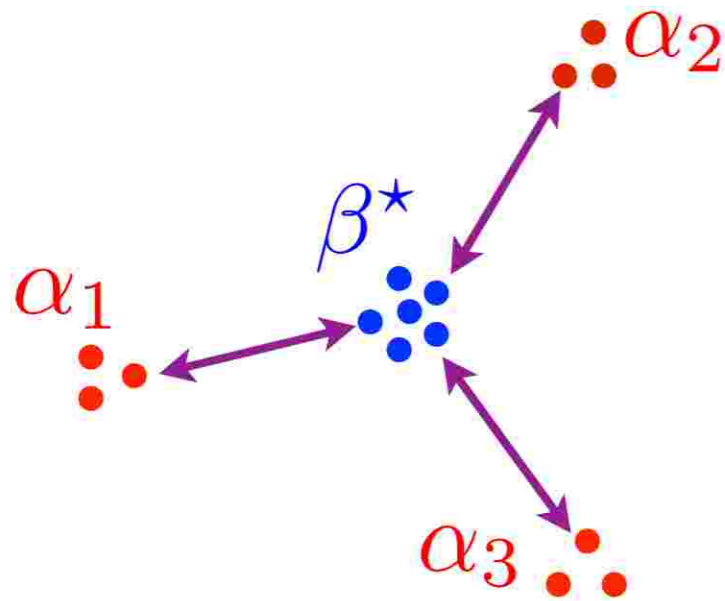
$$\beta^* \in \operatorname{argmin}_{\beta} \sum_s \lambda_s W_p^p(\alpha_s, \beta)$$



Guillaume Carlier



Martial Agueh



[Solomon et al, SIGGRAPH 2015]



# Wasserstein Barycenters

Barycenters of measures  $(\alpha_s)_s$ :  $\sum_s \lambda_s = 1$

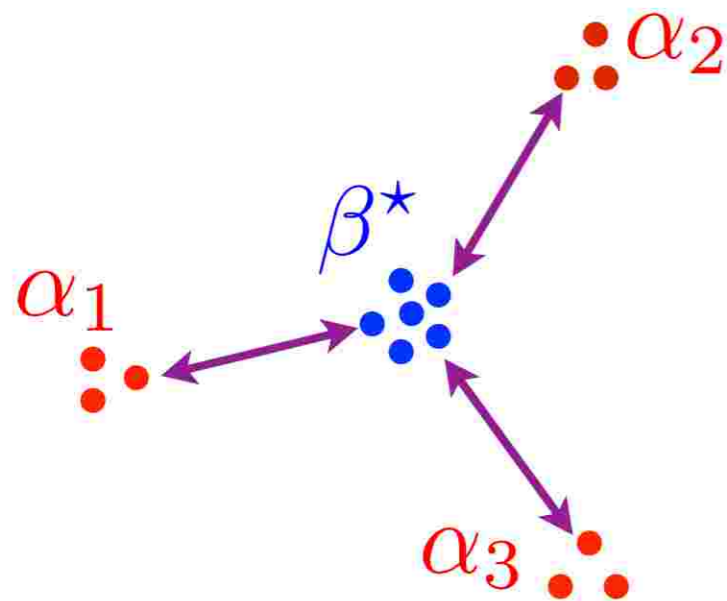
$$\beta^* \in \operatorname{argmin}_{\beta} \sum_s \lambda_s W_p^p(\alpha_s, \beta)$$



Guillaume Carlier



Martial Agueh



[Solomon et al, SIGGRAPH 2015]

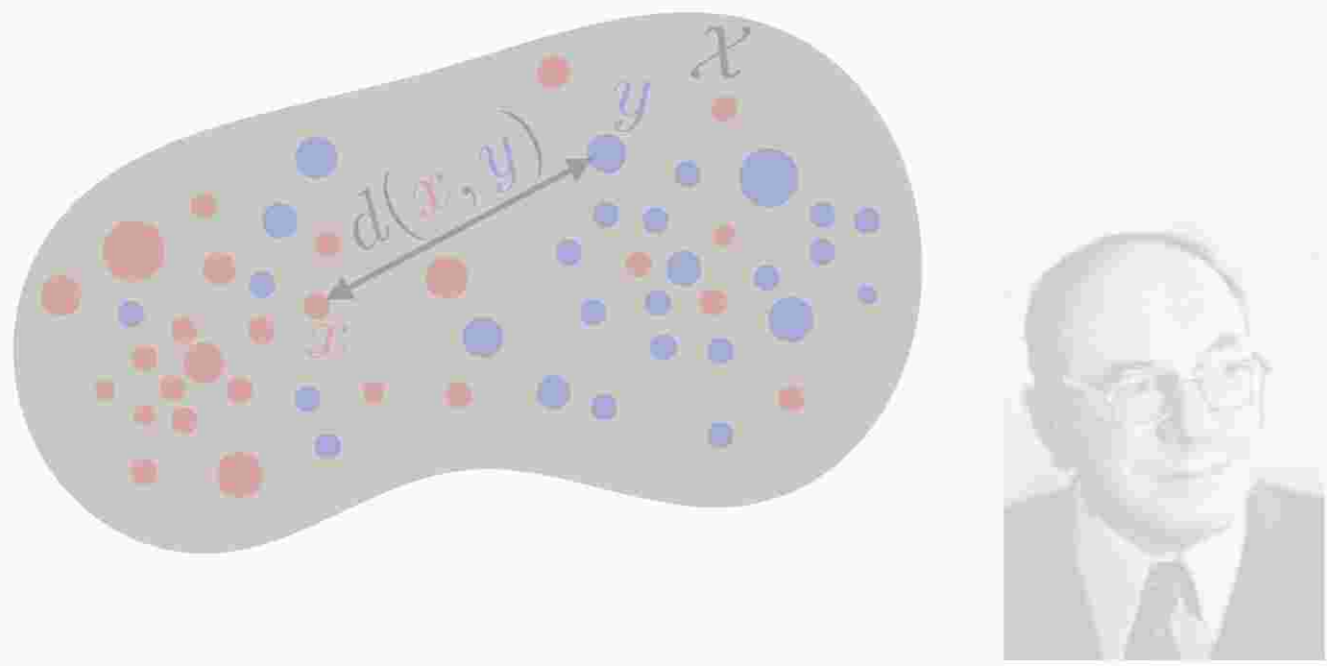
Sinkhorn's algorithm:

$$\left( \mathbf{u}_s \leftarrow \frac{\mathbf{a}_s}{\mathbf{K} \mathbf{v}_s} \right)_s$$

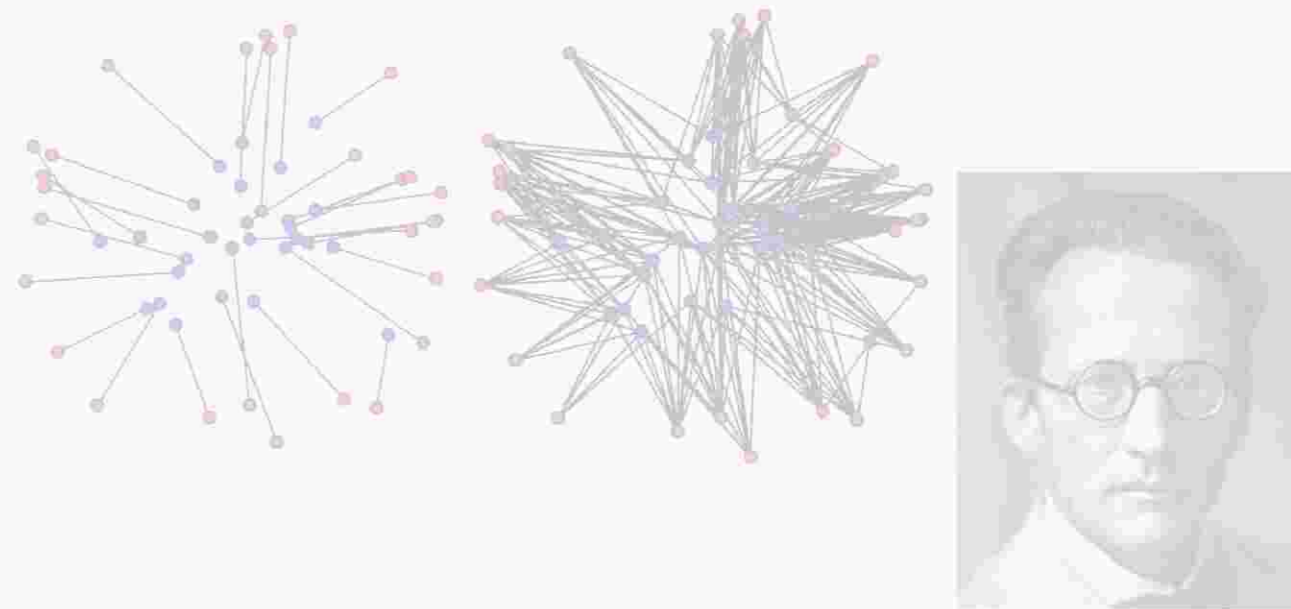
$$\left( \mathbf{v}_s \leftarrow \frac{\mathbf{b}}{\mathbf{K}^\top \mathbf{u}_s} \right)_s$$

$$\mathbf{b} \leftarrow \prod_s (\mathbf{K}^\top \mathbf{u}_s)^{\lambda_s}$$

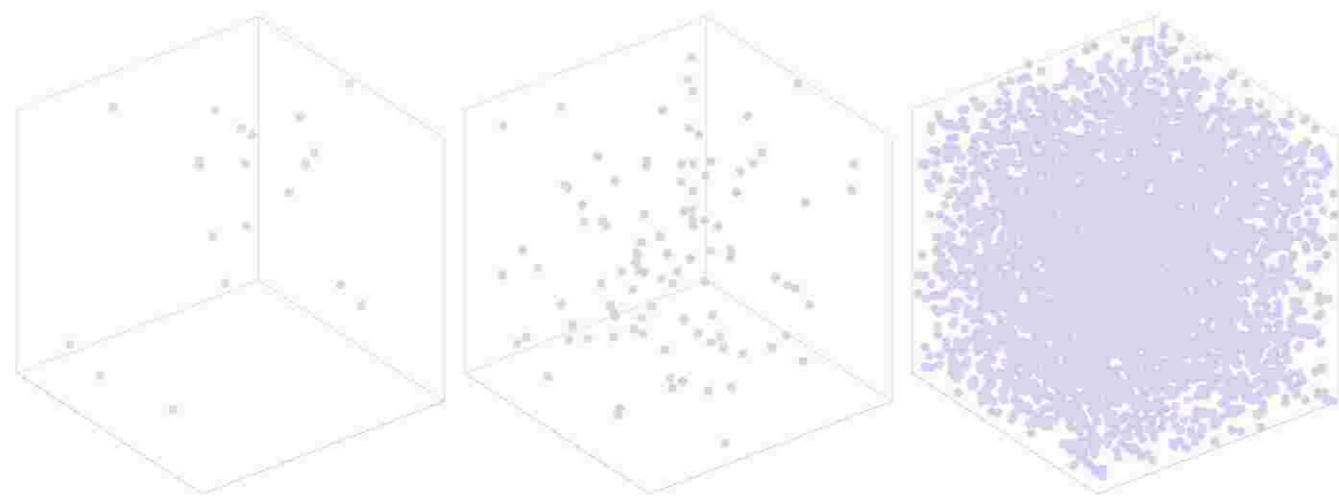
# 1. Optimal Transport



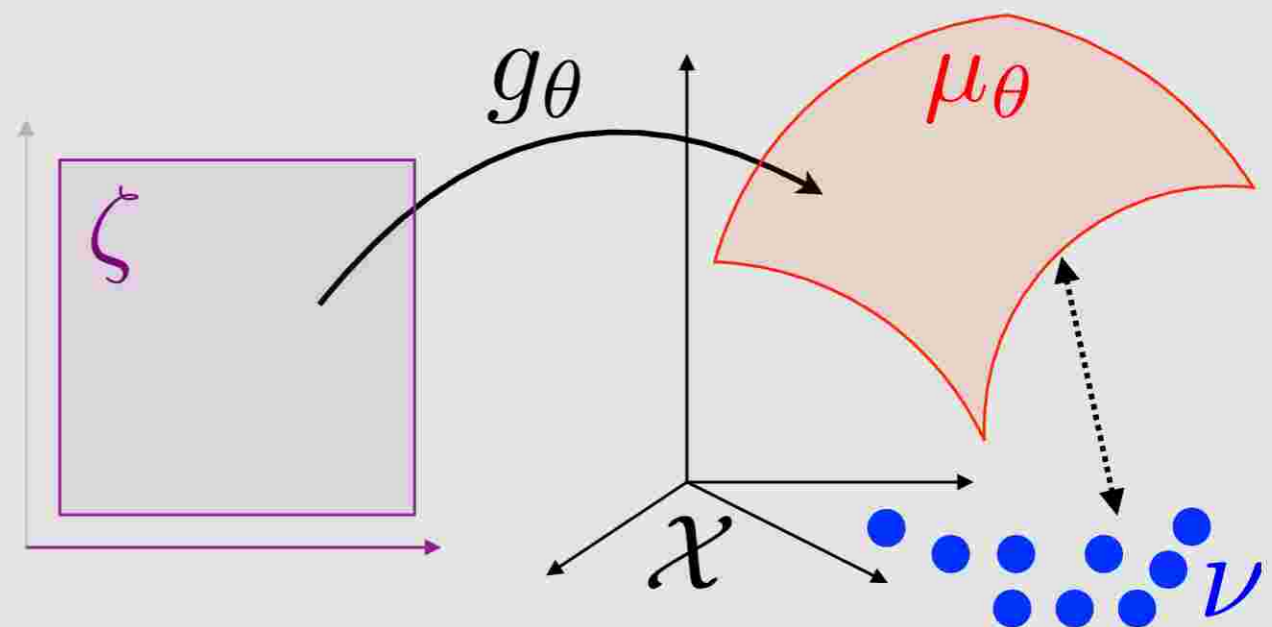
# 2. Entropic Regularization



# 3. Sinkhorn Divergences



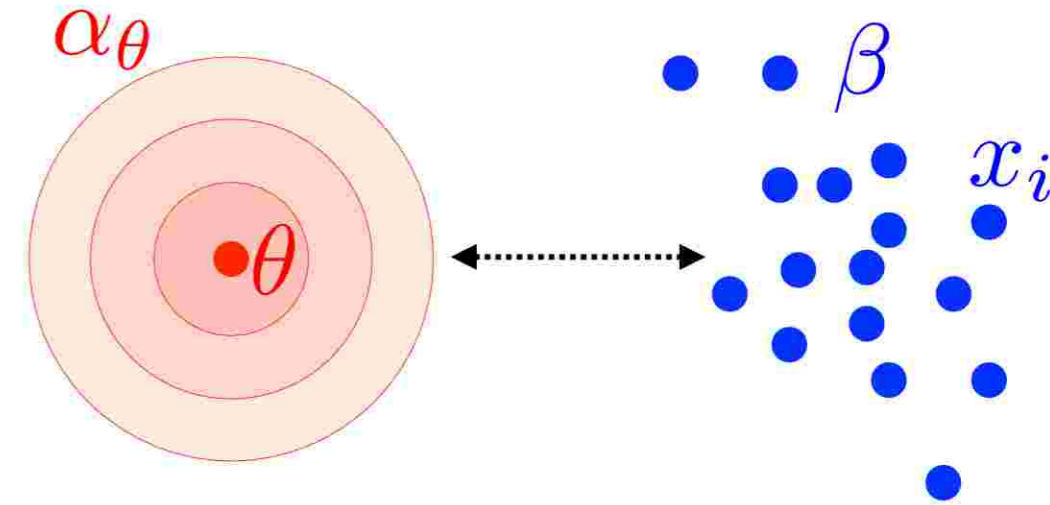
# 4. Application to Generative Models



# Density Fitting and Generative Models

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

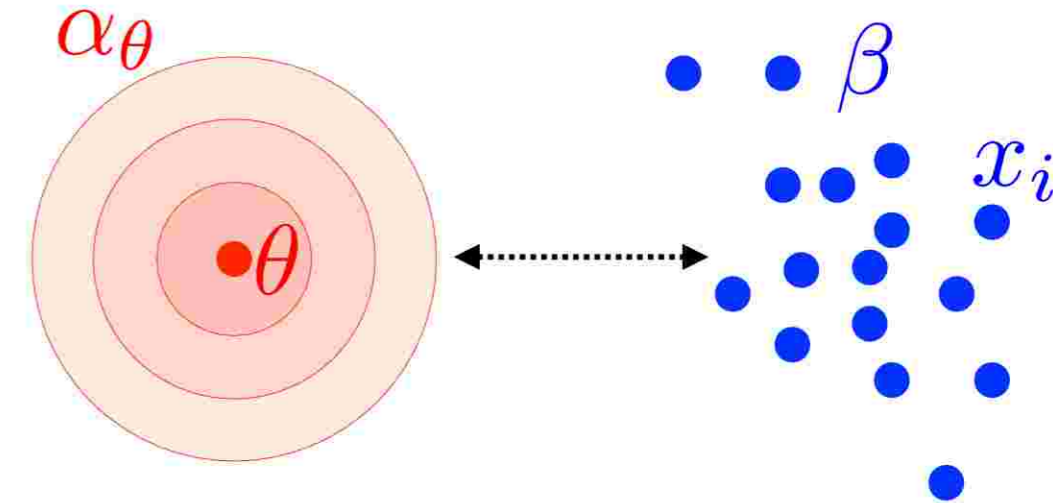
*Parametric model:*  $\theta \mapsto \alpha_\theta$



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

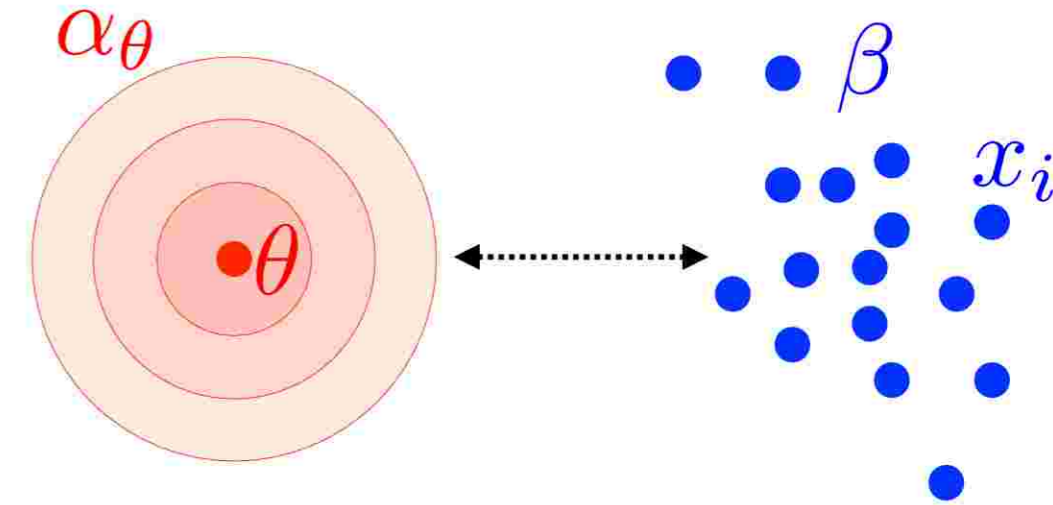
$$\min_{\theta} - \sum_i \log(\rho_\theta(x_i)) \xrightarrow{n \rightarrow +\infty} \text{KL}(\beta | \alpha_\theta)$$

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

$$\min_{\theta} - \sum_i \log(\rho_\theta(x_i)) \xrightarrow{n \rightarrow +\infty} \text{KL}(\beta | \alpha_\theta)$$

Maximum likelihood (MLE)

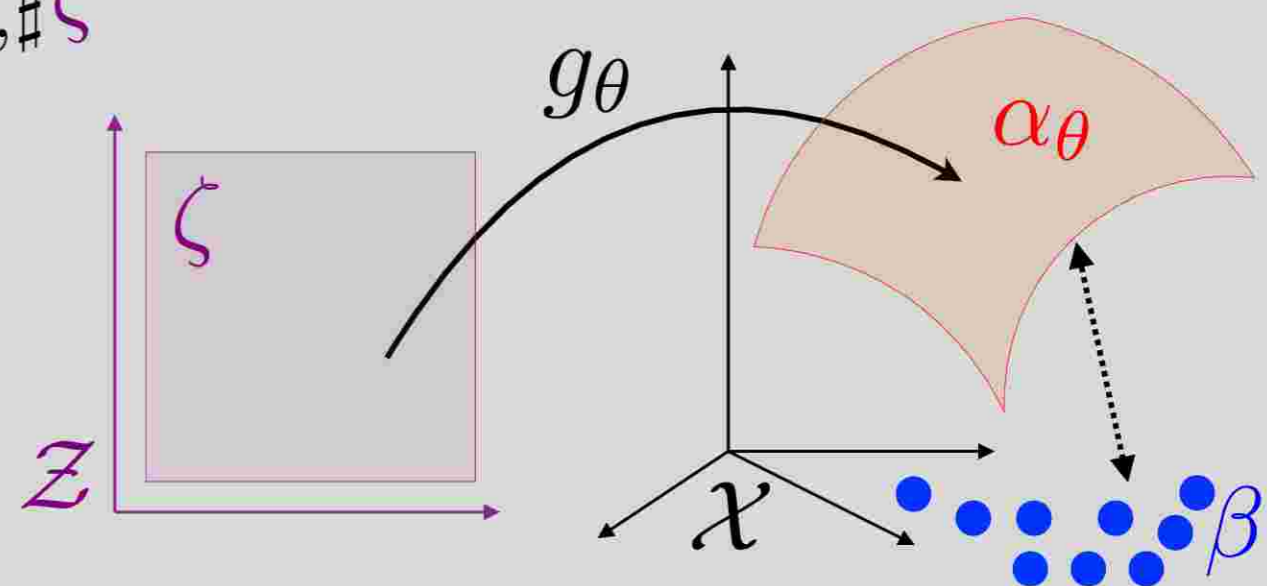
Generative model fit:  $\alpha_\theta = g_{\theta, \#} \zeta$

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

→ MLE undefined.

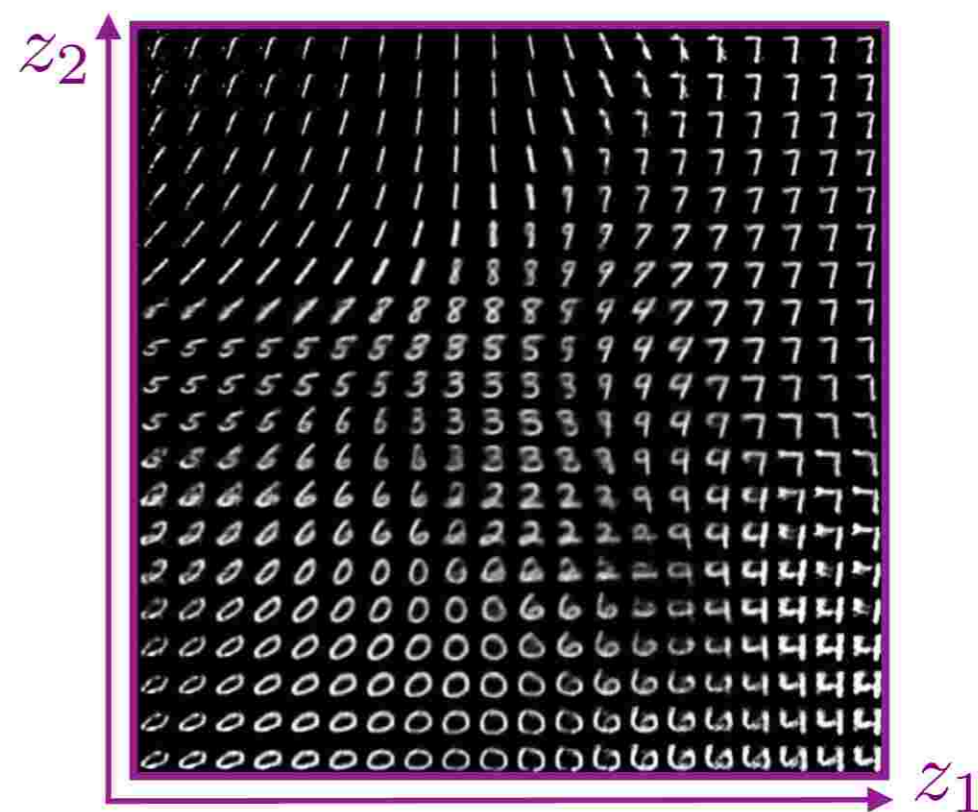
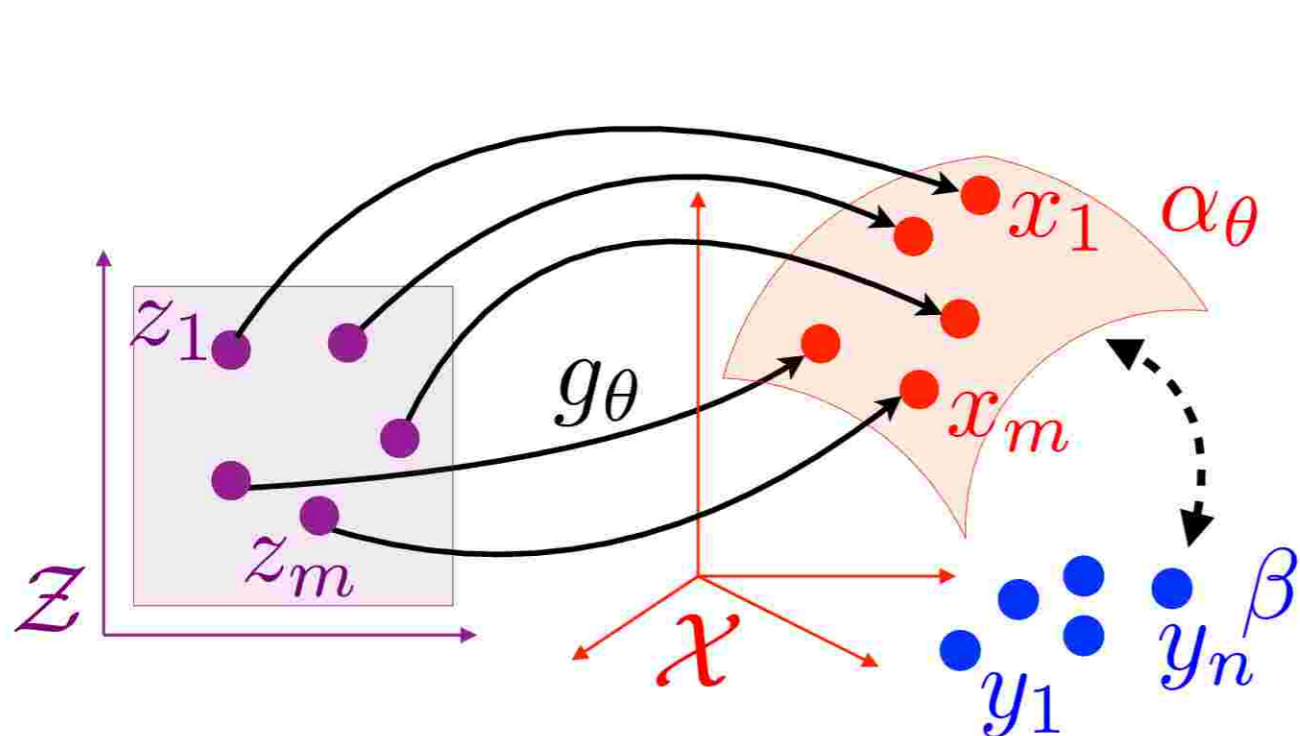
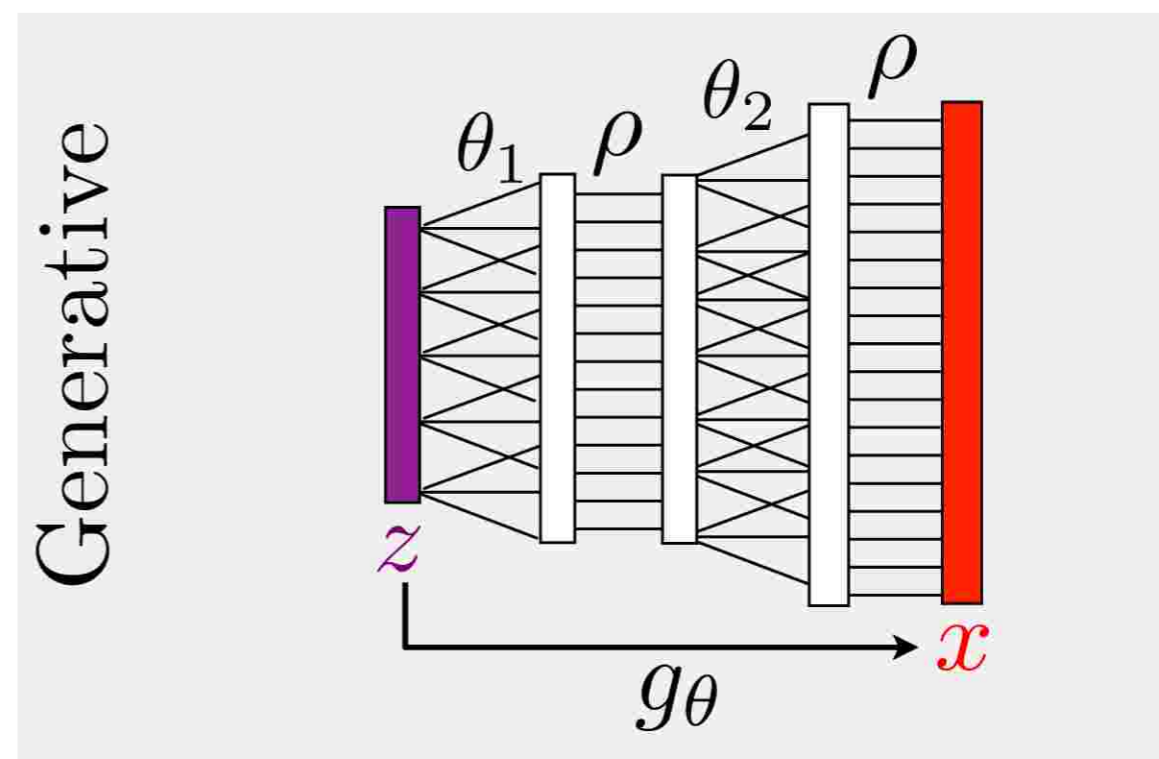
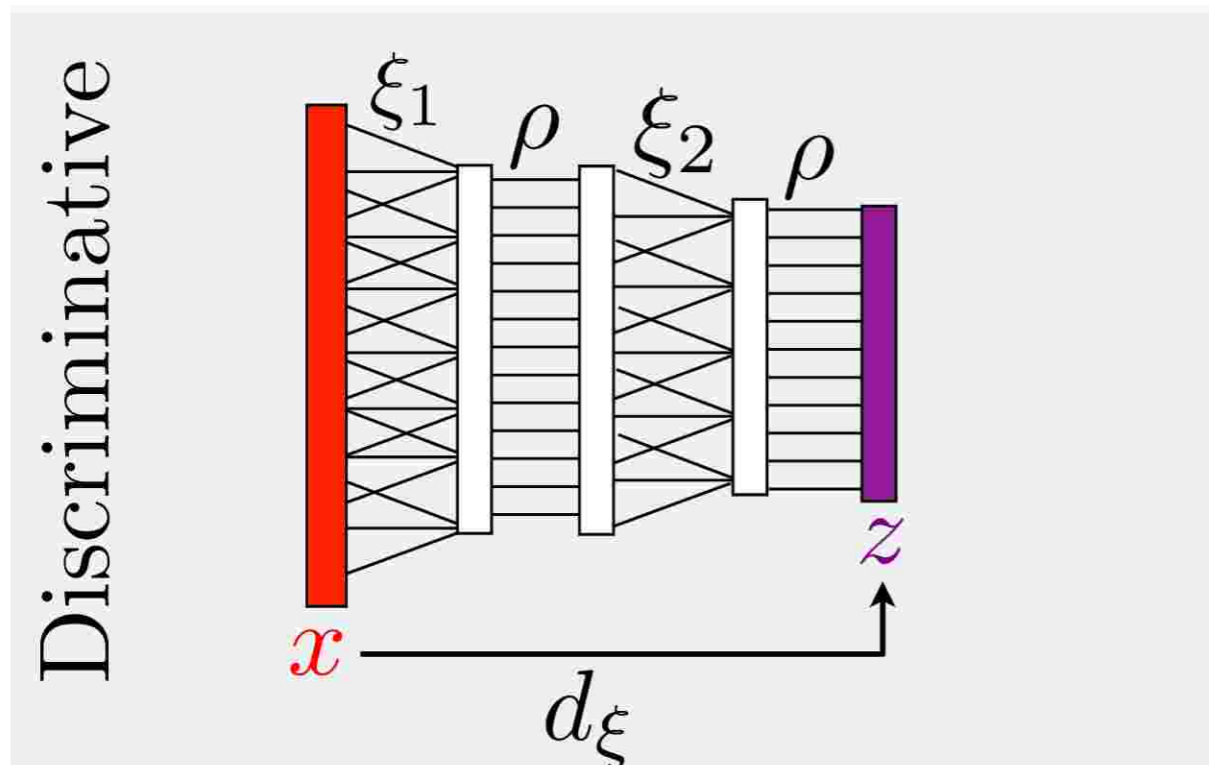
→ Need a weaker metric.

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

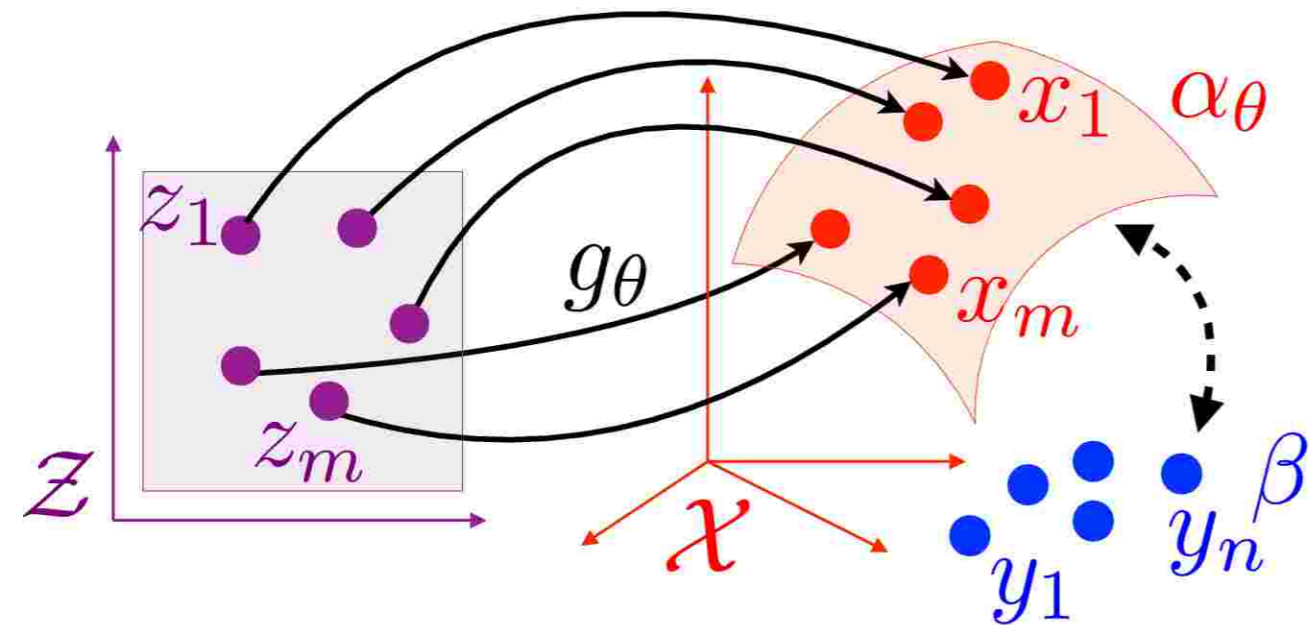


# Deep Discriminative vs Generative Models

Deep networks:  $d_{\xi}(x) = \rho(\xi_K(\dots \rho(\xi_2(\rho(\xi_1(x) \dots)))$   
 $g_{\theta}(z) = \rho(\theta_K(\dots \rho(\theta_2(\rho(\theta_1(z) \dots)))$



# Training Architecture



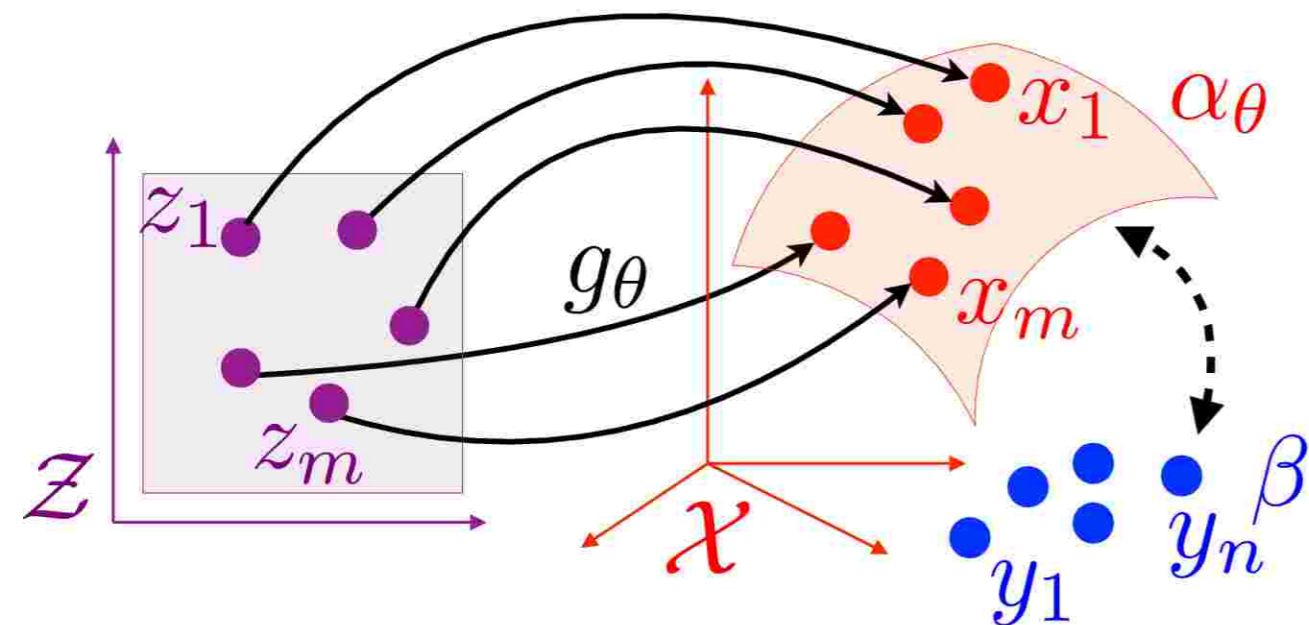
$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$

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

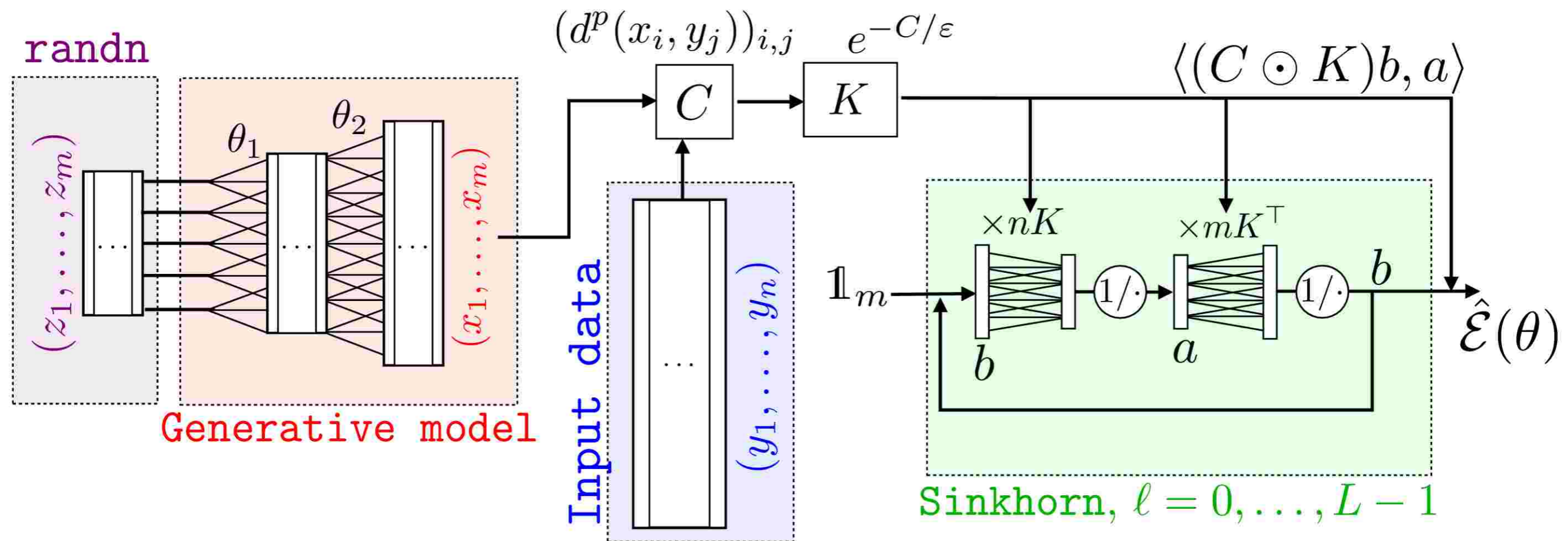


$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

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

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



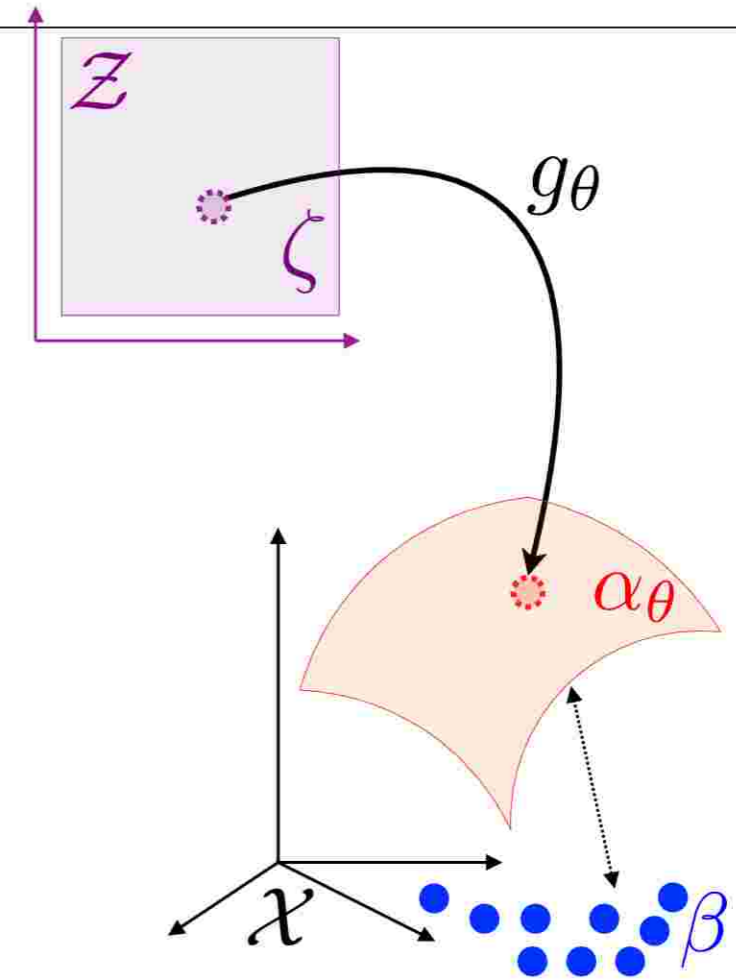


# Examples of Images Generation

Inputs  $\beta$



Generated  $\alpha_\theta$

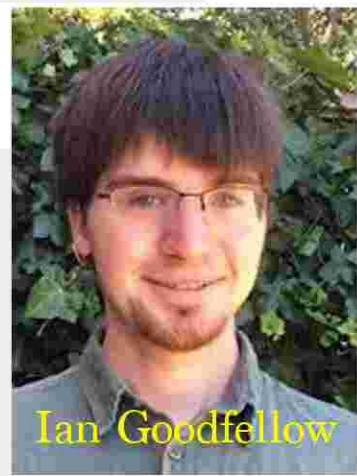
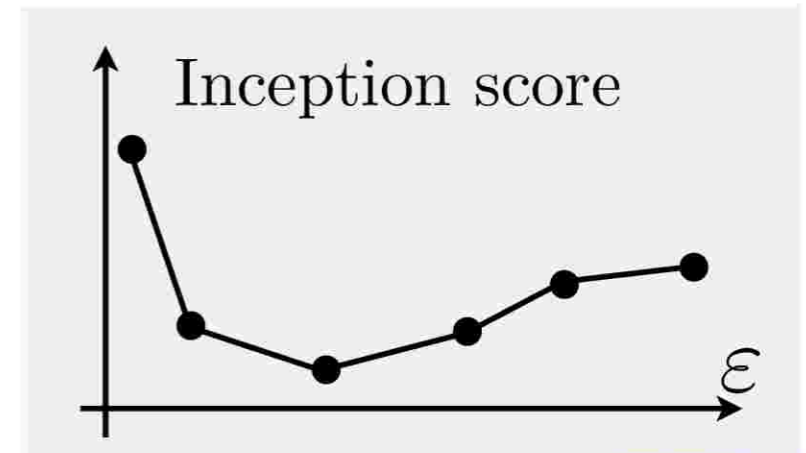
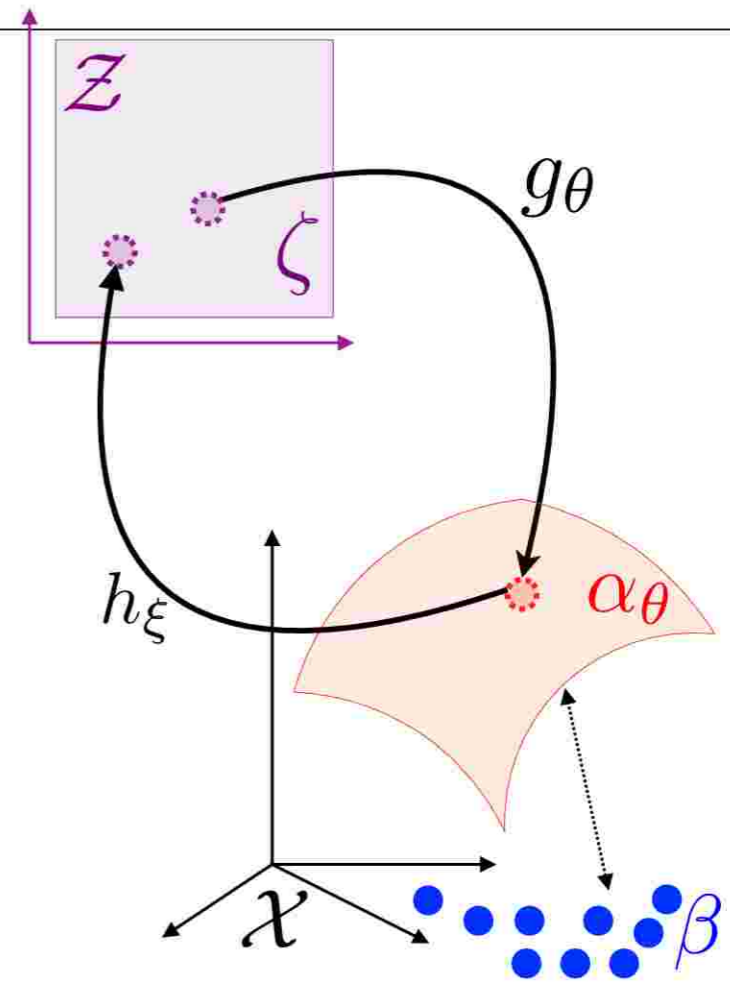


# Examples of Images Generation

Inputs  $\beta$



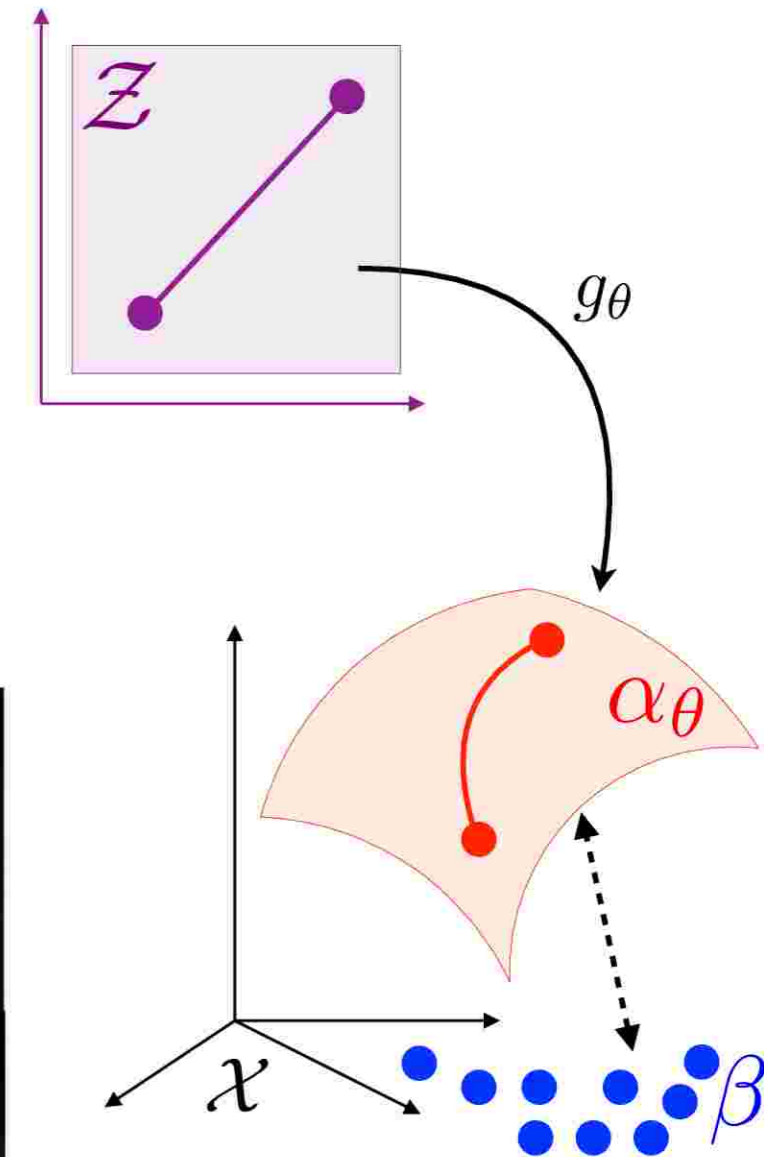
Generated  $\alpha_\theta$



Ian Goodfellow

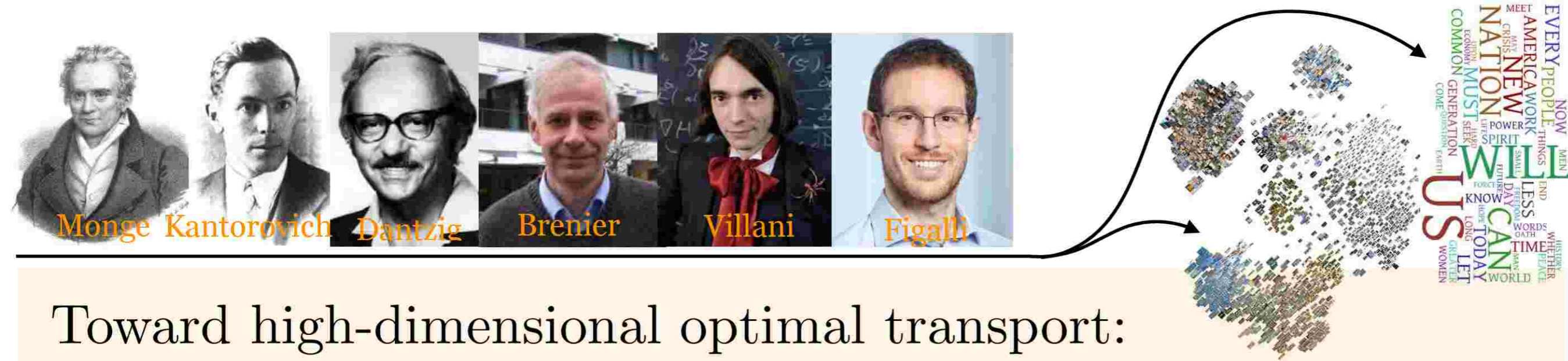
- Need to learn the metric  $d(x, y) = \|h_\xi(x) - h_\xi(y)\|$  (GANs)
- Influence of  $\epsilon$ ?
- Performance evaluation of generative models is an open problem.

# Generative Adversarial Networks



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

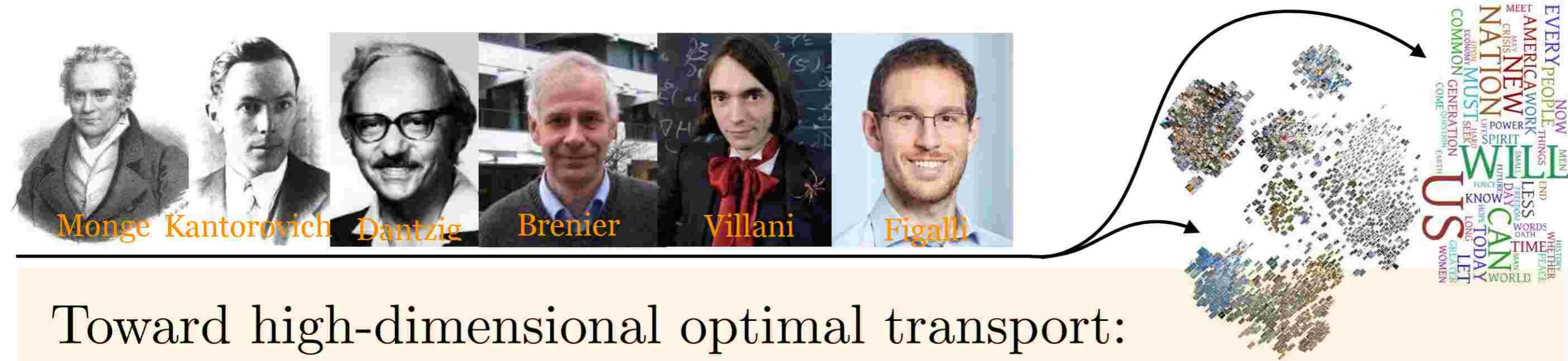
# Conclusion



Toward high-dimensional optimal transport:

→ Scalable geometrical loss functions in high dimension?

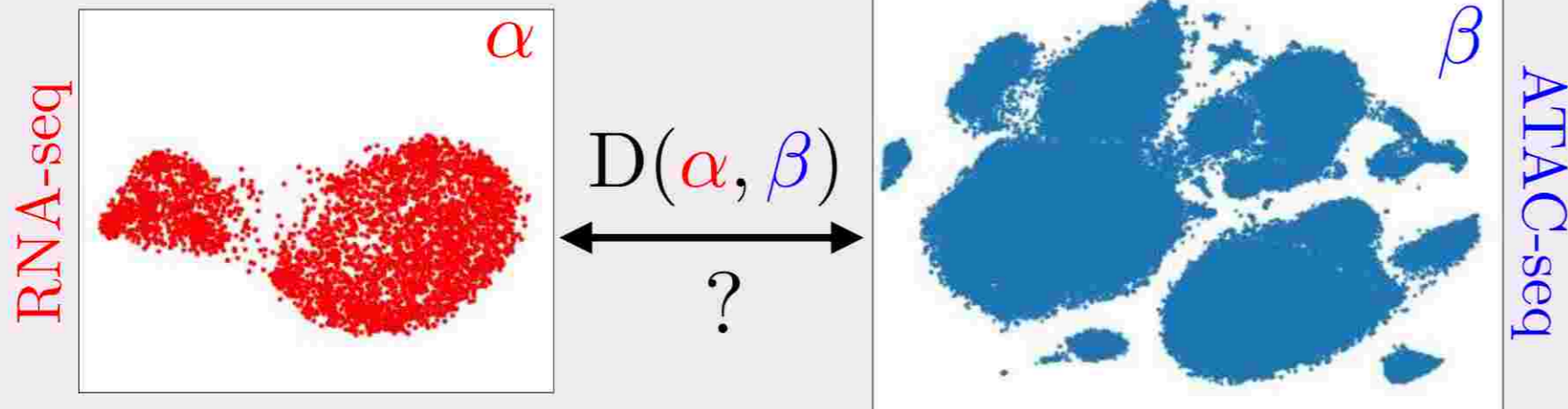
# Conclusion



Toward high-dimensional optimal transport:

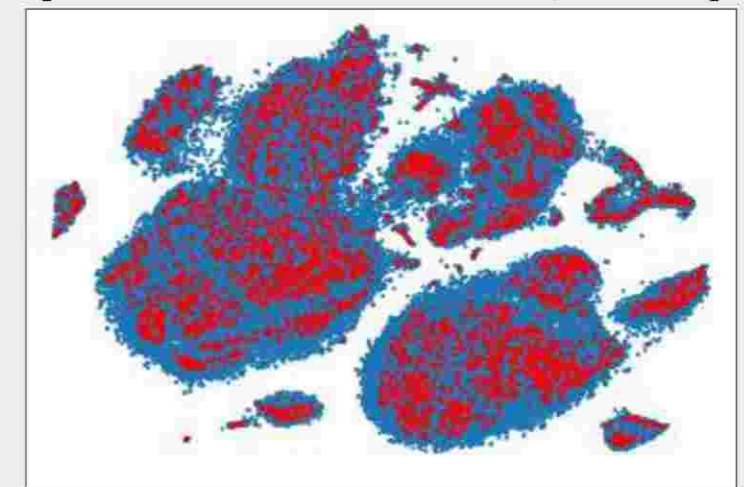
→ Scalable geometrical loss functions in high dimension?

Comparing datasets *across* spaces:



Single-cell multi-omics

[Othmane Sebbouh, 2021]



Gromov-Wasserstein registration