

# Computational Optimal Transport

<http://optimaltransport.github.io>

## *Introduction*

Gabriel Peyré

[www.numerical-tours.com](http://www.numerical-tours.com)



<https://optimaltransport.github.io>

Home

BOOK

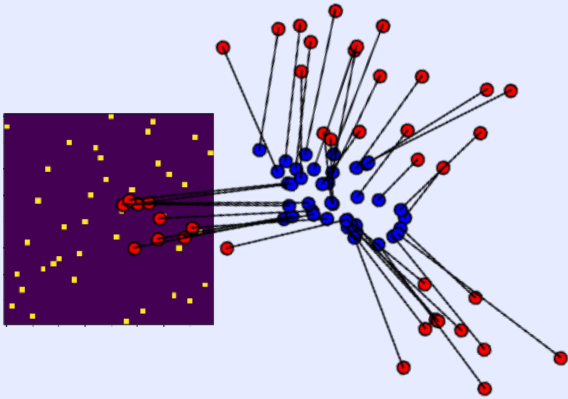
CODE

SLIDES

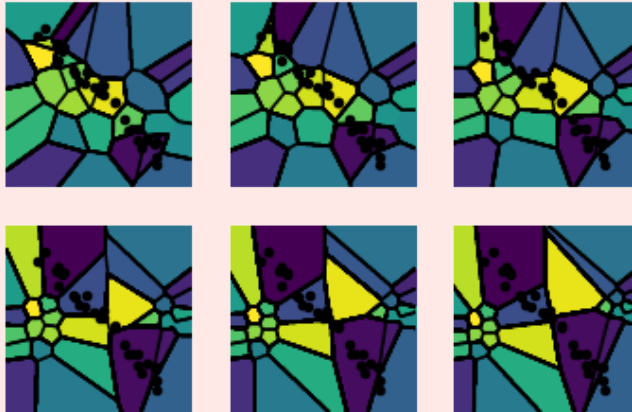
# Computational Optimal Transport

---

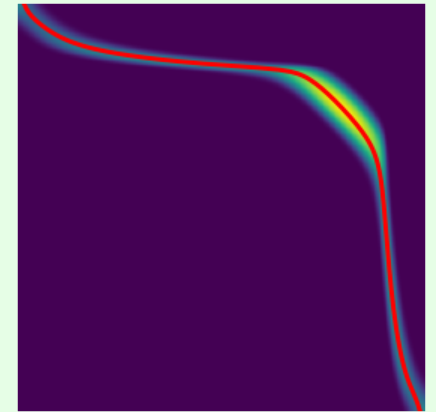
### Optimal Transport with Linear Programming



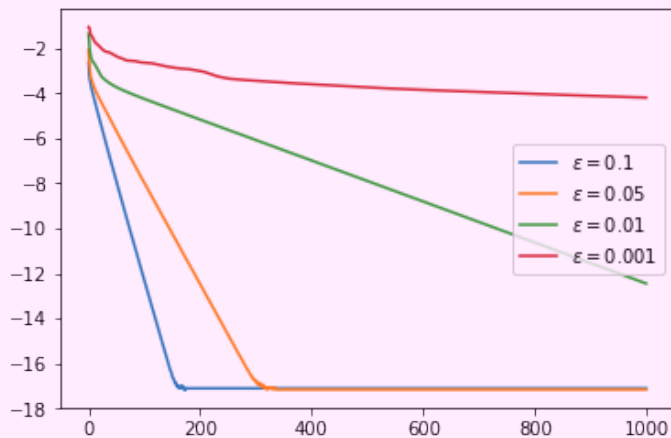
### Semi-discrete Optimal Transport



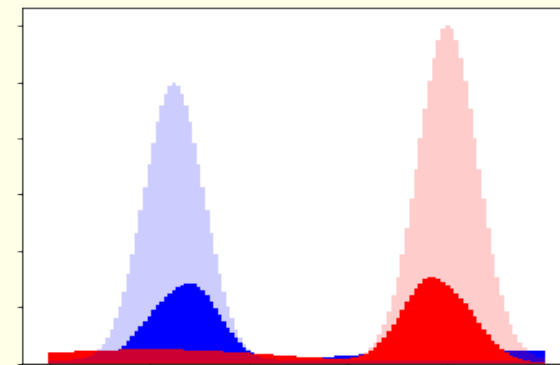
### Entropic Regularization of Optimal Transport



### Advanced Topics on Sinkhorn Algorithm



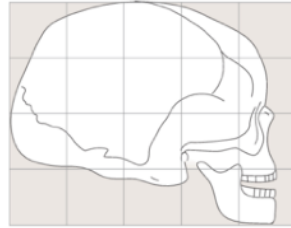
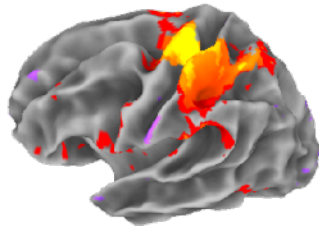
### Unbalanced Optimal Transport



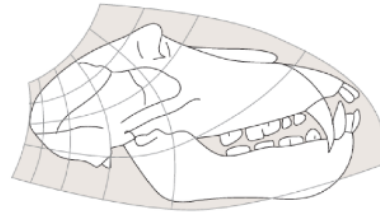
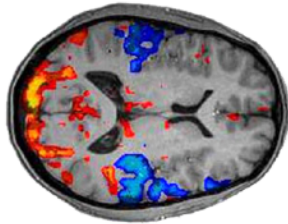
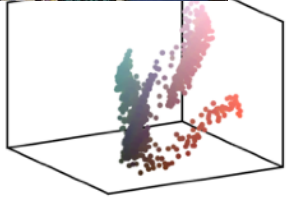
# Probability Distributions in Data Sciences

*Probability distributions and histograms*

→ images, vision, graphics and machine learning, .



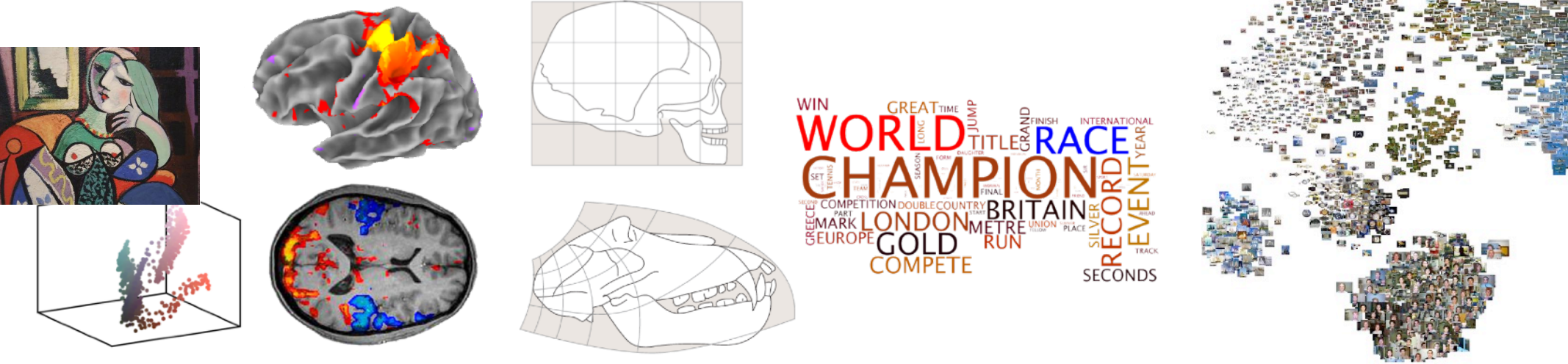
WIN GREAT TIME JUMP GRAND FINISH INTERNATIONAL  
WORLD LONG TITLE RACE  
SEASON OLYMPIC CHAMPION MONTH  
SET TENNIS YEAR  
GREECE COMPETITION DOUBLE COUNTRY BRITAIN PLACE  
PART MARK LONDON METRE UNION  
EUROPE GOLD RUN SILVER RECORD EVENT TRACK  
COMPETE SECONDS



# Probability Distributions in Data Sciences

*Probability distributions and histograms*

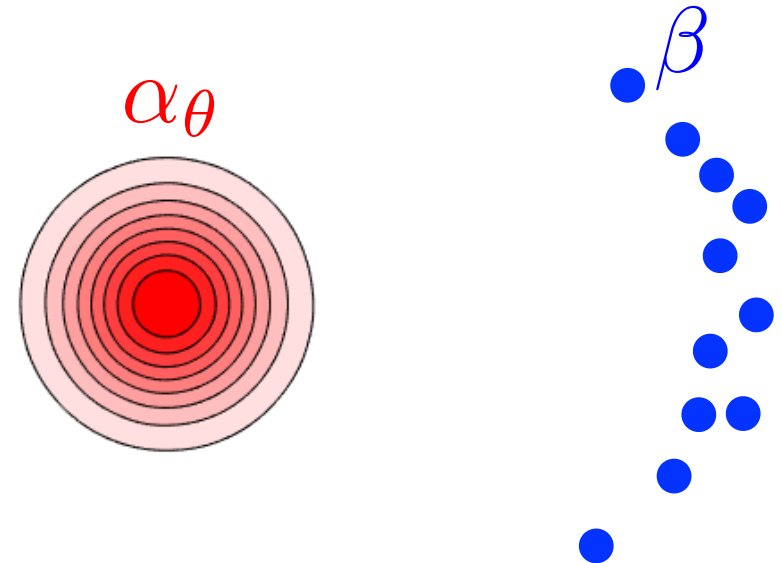
→ images, vision, graphics and machine learning, .



## Unsupervised learning

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

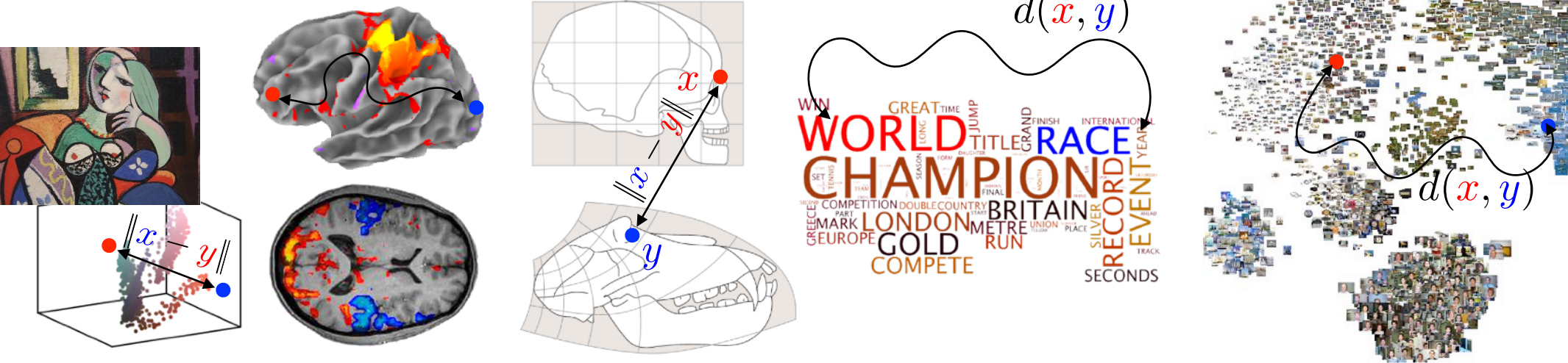
Parametric model:  $\theta \mapsto \alpha_\theta$



# Probability Distributions in Data Sciences

*Probability distributions and histograms*

→ images, vision, graphics and machine learning,



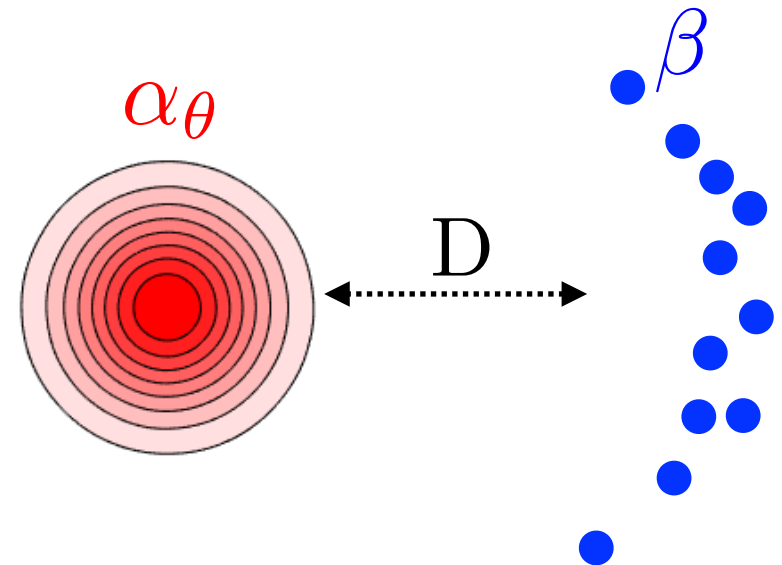
## Unsupervised learning

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  $d$ .



# Overview

---

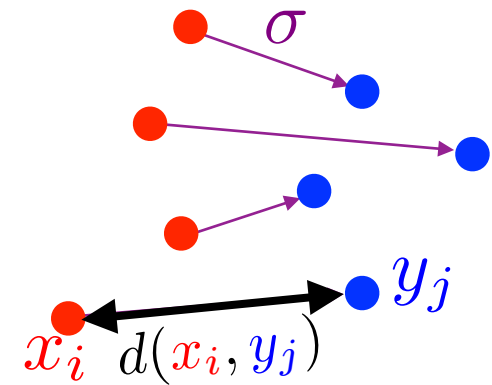
- **Monge Formulation**
- Continuous Optimal Transport
- Kantorovitch Formulation
- Applications

# 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'en suit 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*.

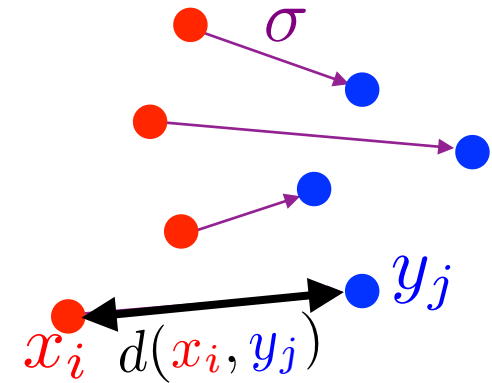


# Monge's Problem

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

Permutation:

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



Monge optimal matching:

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



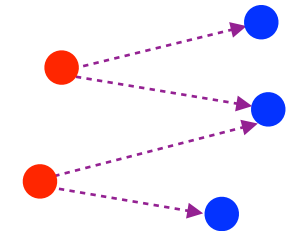
[Monge 1784]

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'en suit 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?



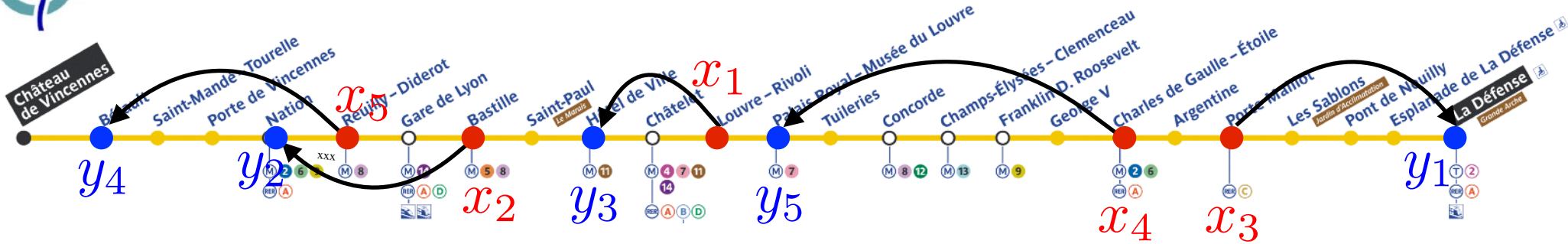
# 1-D Optimal Transport

$$\min_{\sigma \in \Sigma_n} \sum_{i=1}^n |x_i - y_{\sigma(i)}|^p, \quad p \geq 1$$



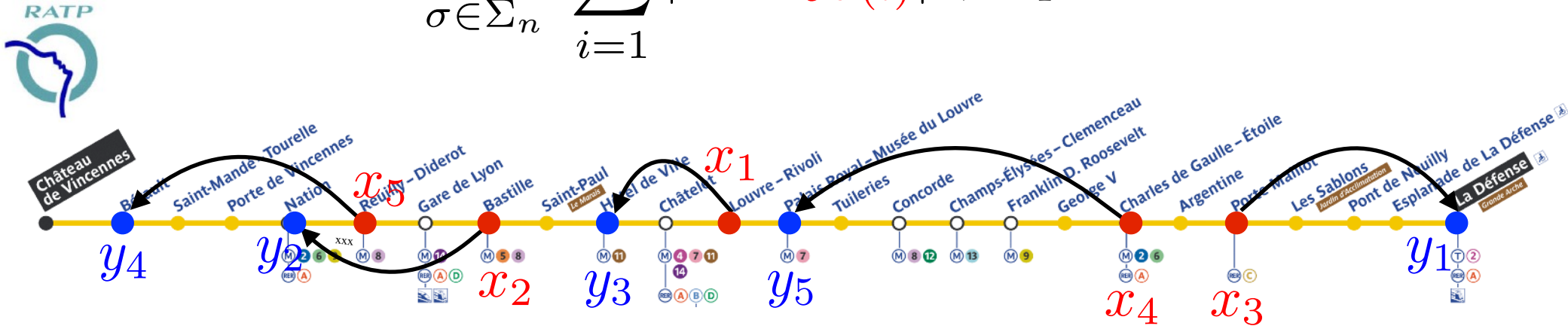
# 1-D Optimal Transport

$$\min_{\sigma \in \Sigma_n} \sum_{i=1}^n |x_i - y_{\sigma(i)}|^p, \quad p \geq 1$$



# 1-D Optimal Transport

$$\min_{\sigma \in \Sigma_n} \sum_{i=1}^n |x_i - y_{\sigma(i)}|^p, \quad p \geq 1$$

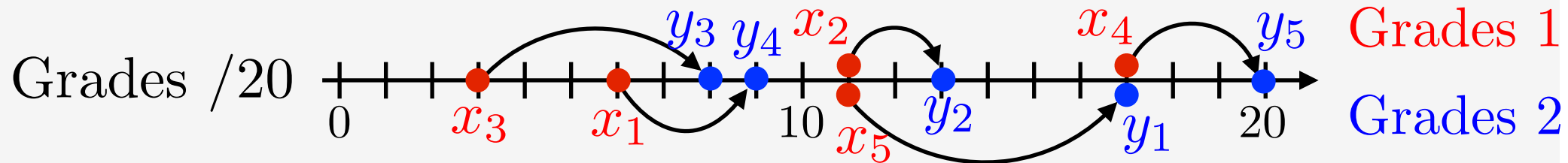


Sorting algorithms: insertion  $n(n - 1)/2$  worst case.

$n$	$n!$	$n(n-1)/2$	$n \log(n)$
10	3628800	45	23
11	39916800	55	26
12	479001600	66	30
25	$1,551 \times 10^{25}$	300	80
70	$1,198 \times 10^{100}$	21415	297

QuickSort:  $O(n \log(n))$ .

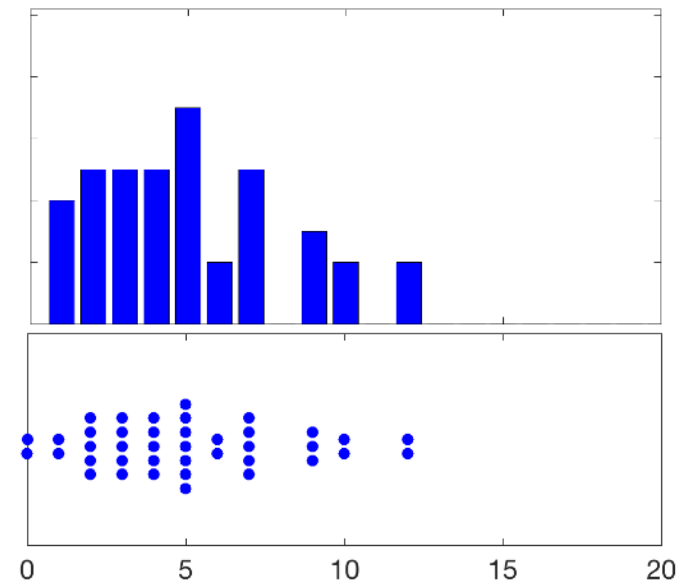
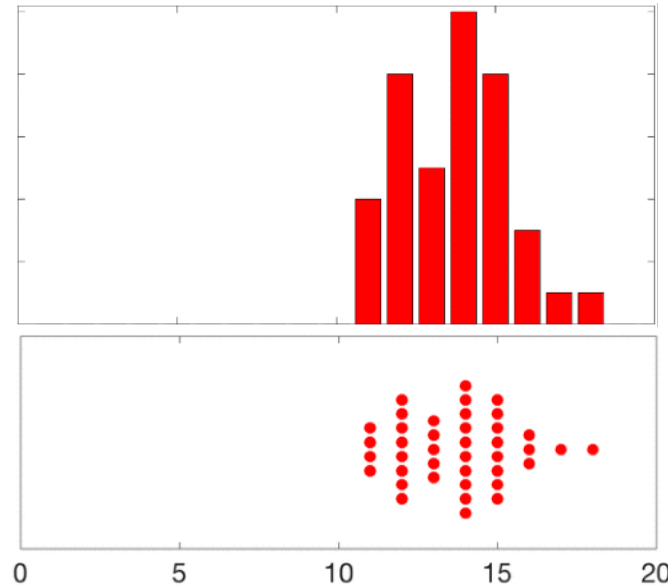
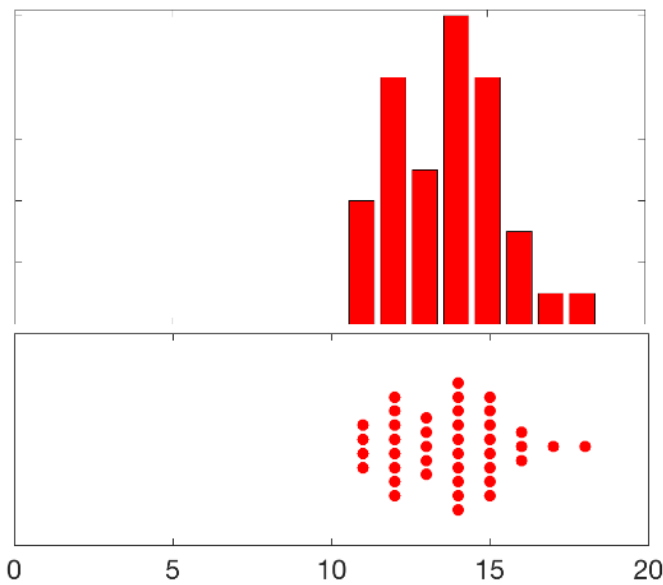
# 1-D OT Interpolation



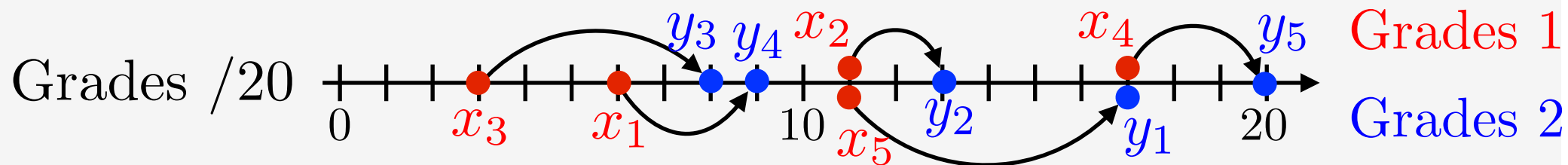
Grades 1

comparison

Grades 2



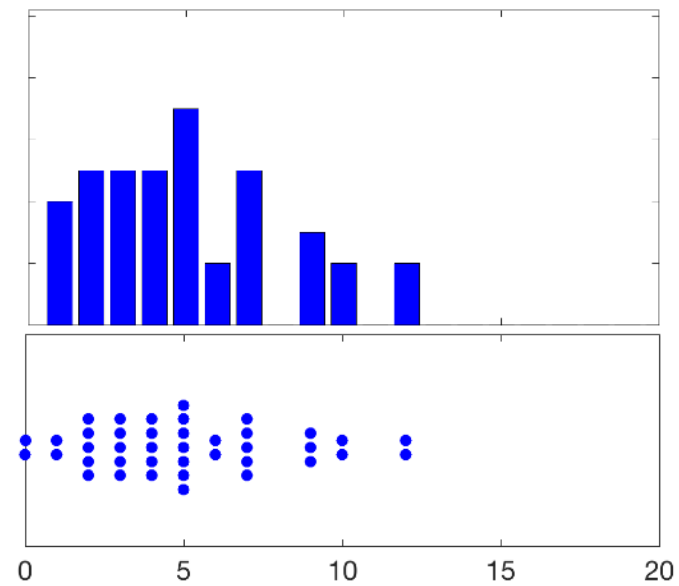
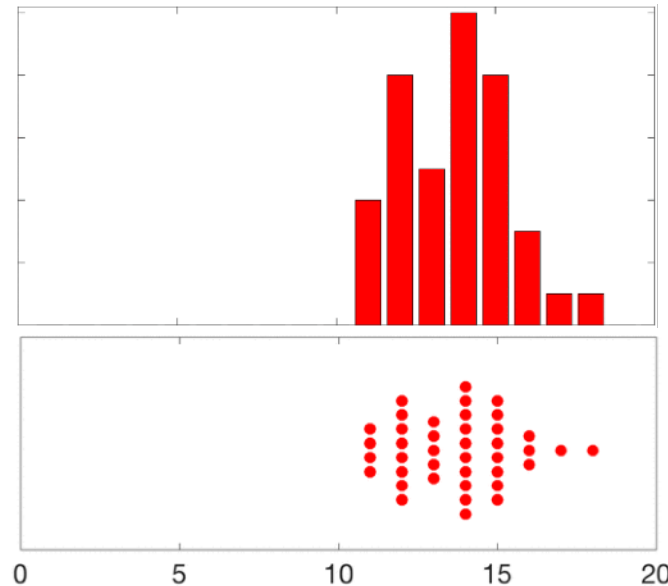
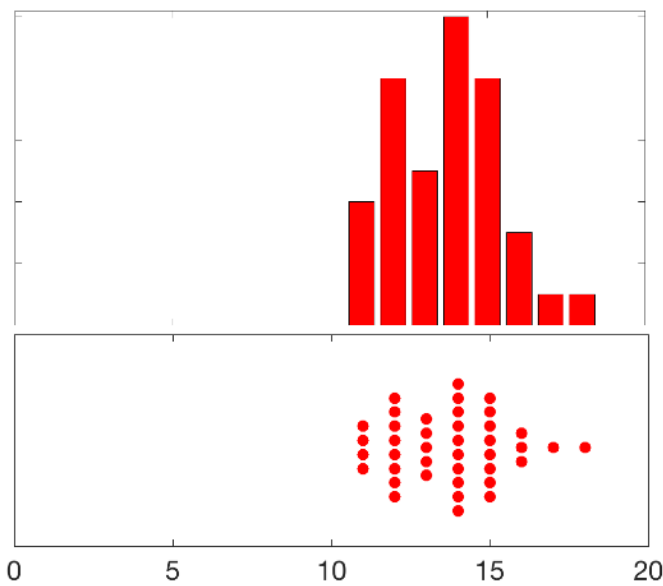
# 1-D OT Interpolation



Grades 1

comparison

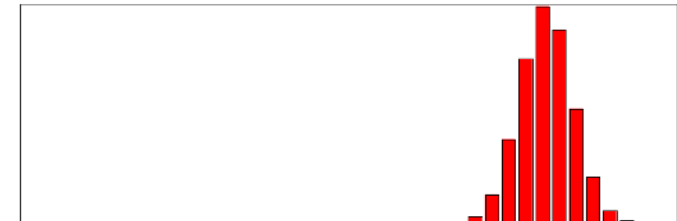
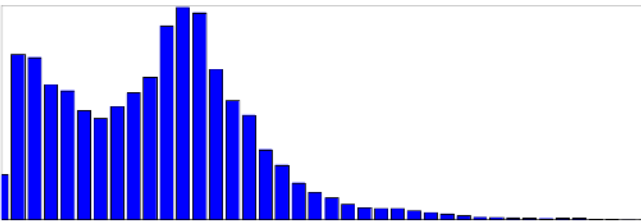
Grades 2



# Grayscale Histogram Equalization



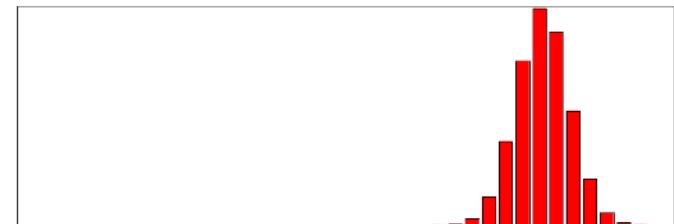
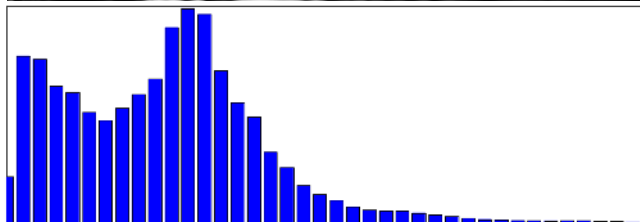
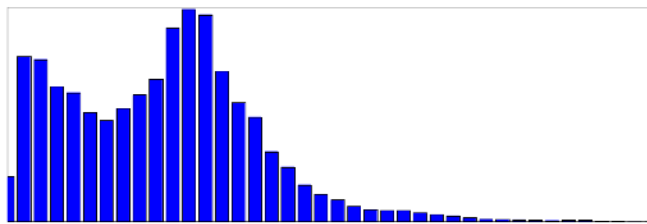
```
f[argsort(f.flatten())] = np.sort(g.flatten())
```



# Grayscale Histogram Equalization



```
f[argsort(f.flatten())] = np.sort(g.flatten())
```

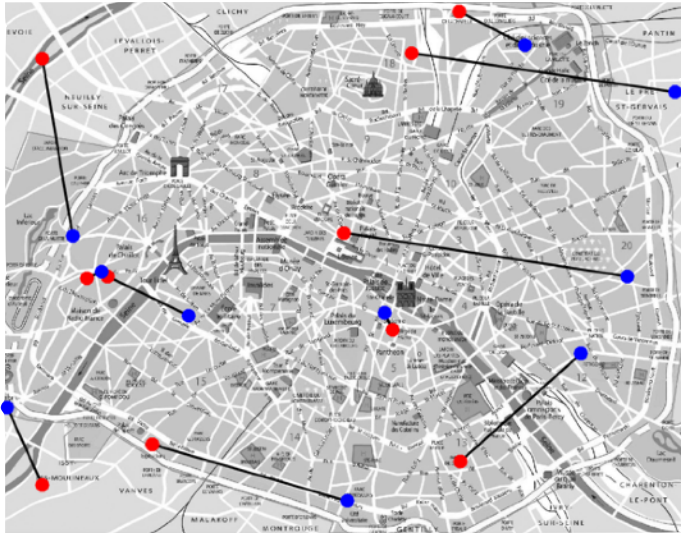
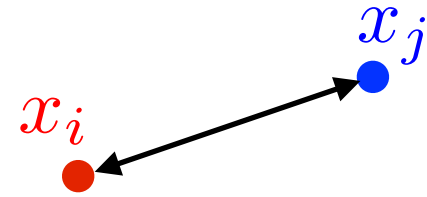




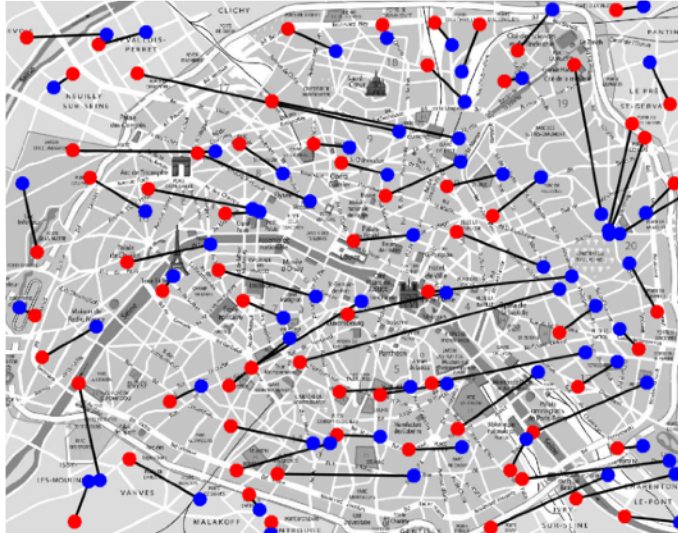
# In 2-D

$$x_i, y_j \in \mathbb{R}^2$$

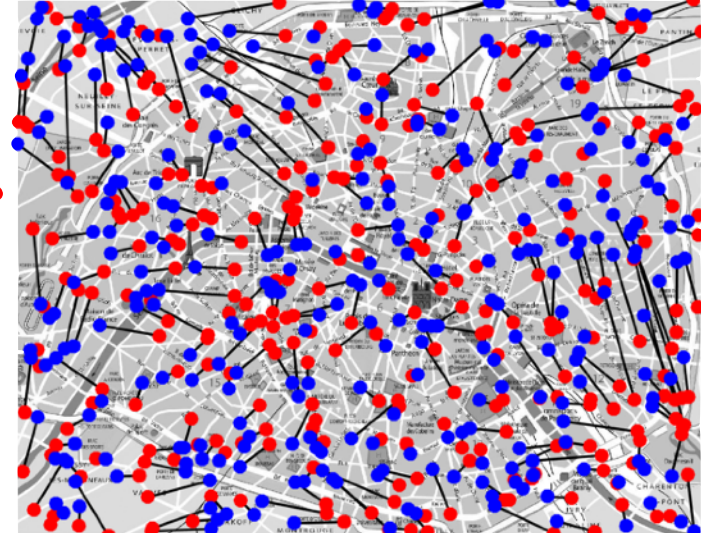
$$c_{i,j} = \|x_i - y_j\| = \sqrt{(x_i^1 - y_j^1)^2 + (x_i^2 - y_j^2)^2}$$



$n = 10$

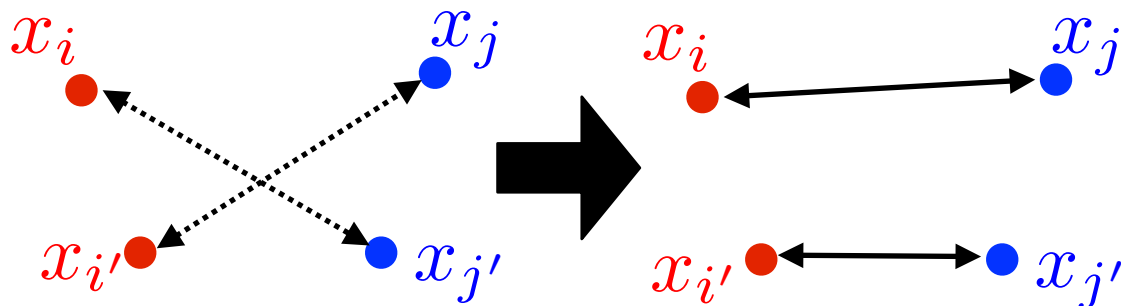


$n = 70$

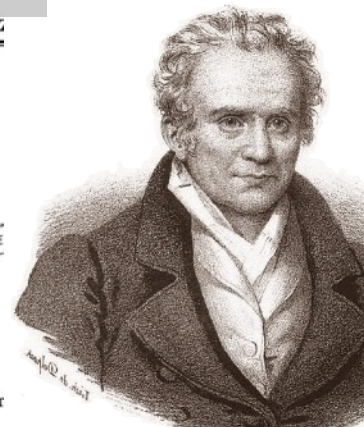
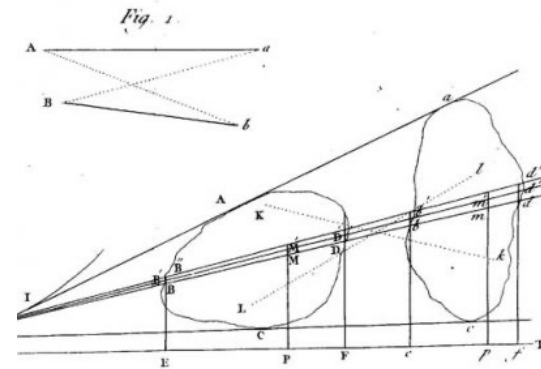


$n = 300$

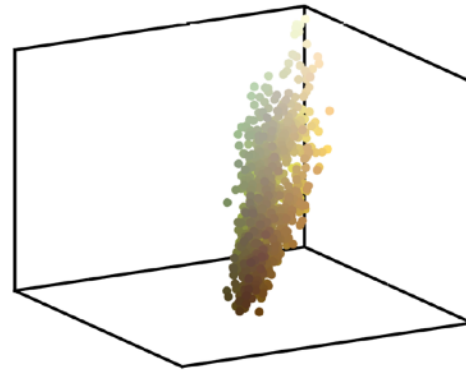
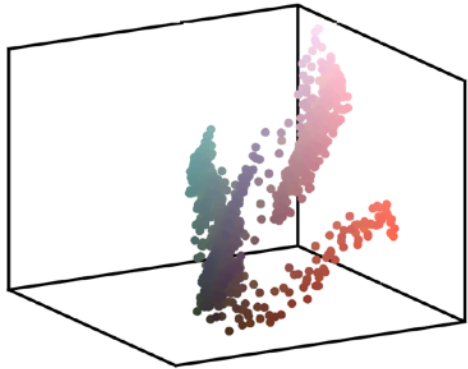
*Proposition: two segments never cross.*



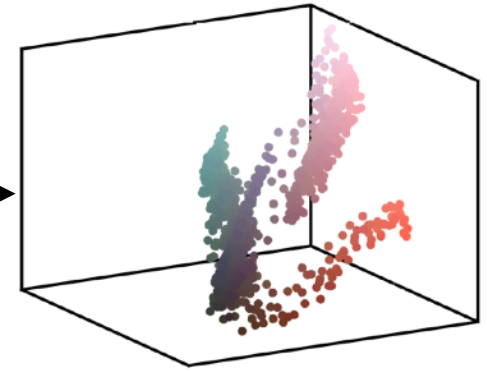
*Mém. de l'Ac. R. des Sc. An. 1781. Page. 704. Pl. XVII.*



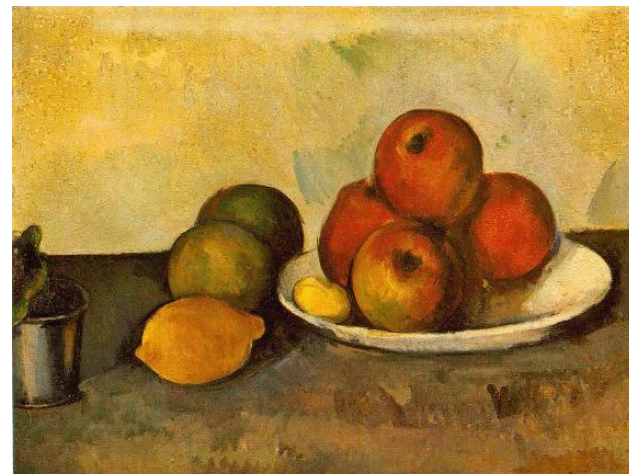
# In 3-D: Color Image Palette Equalization



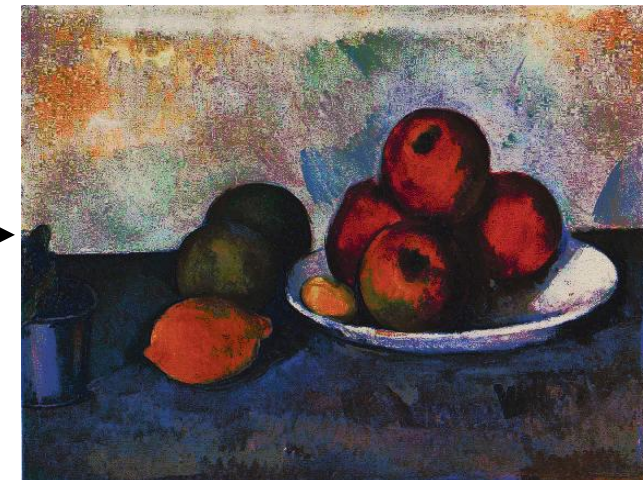
optimal  
transport



Reference

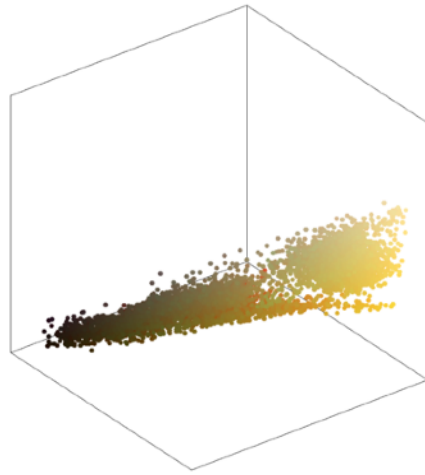


Input

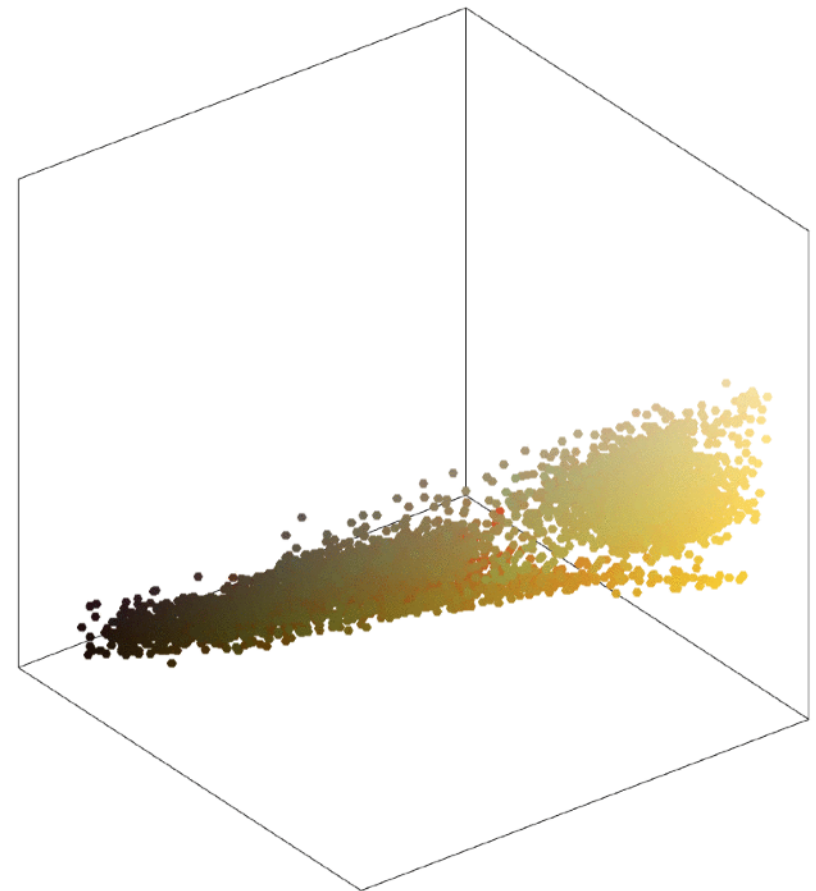
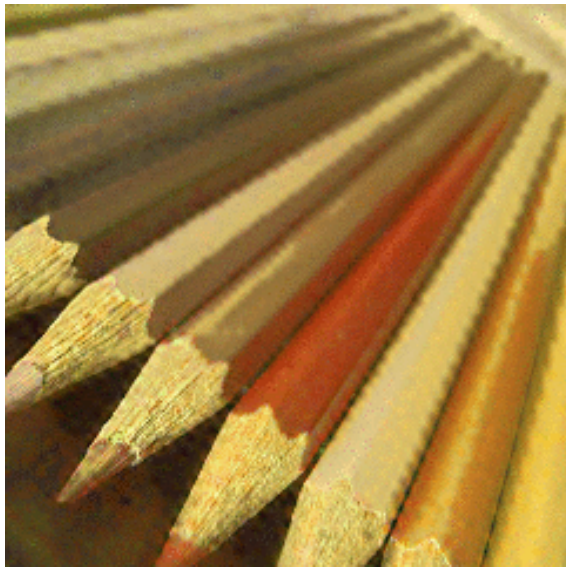
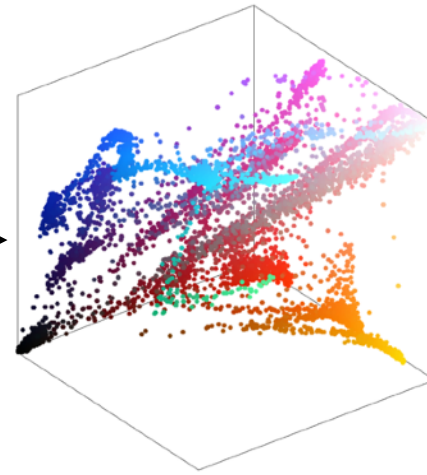


Output

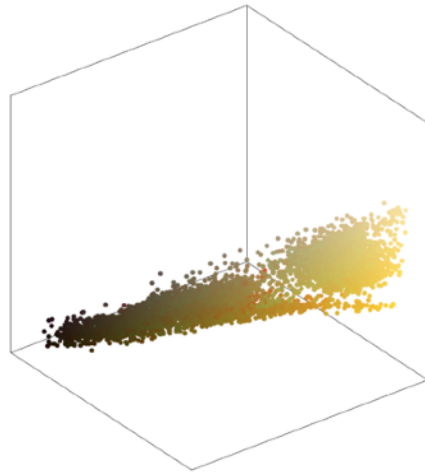
# Color Image Palette Equalization



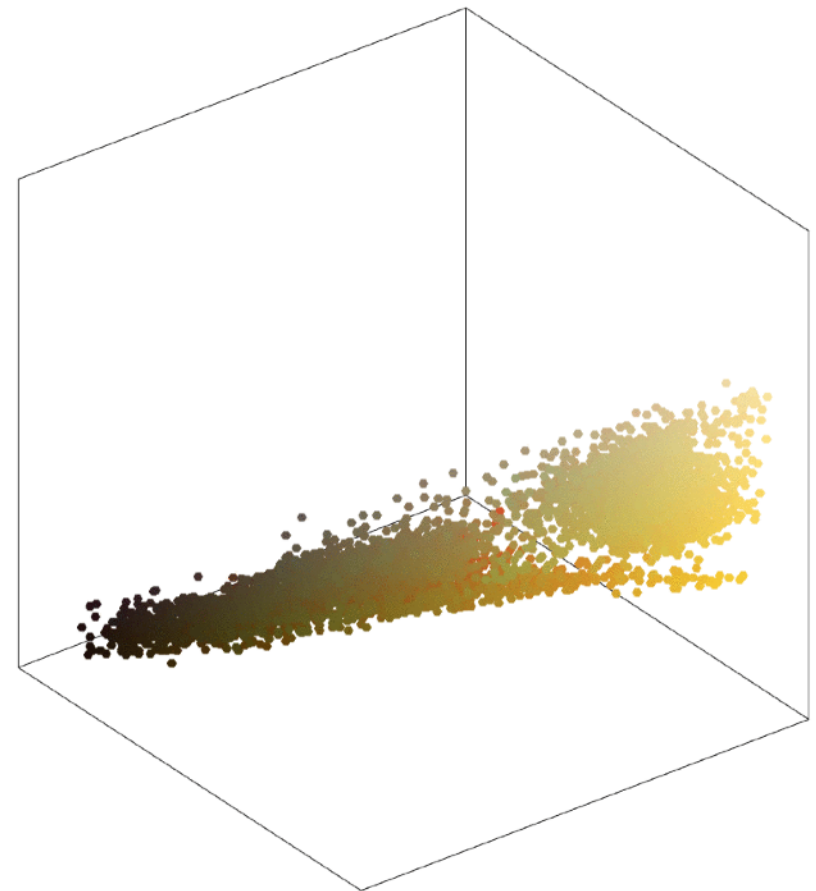
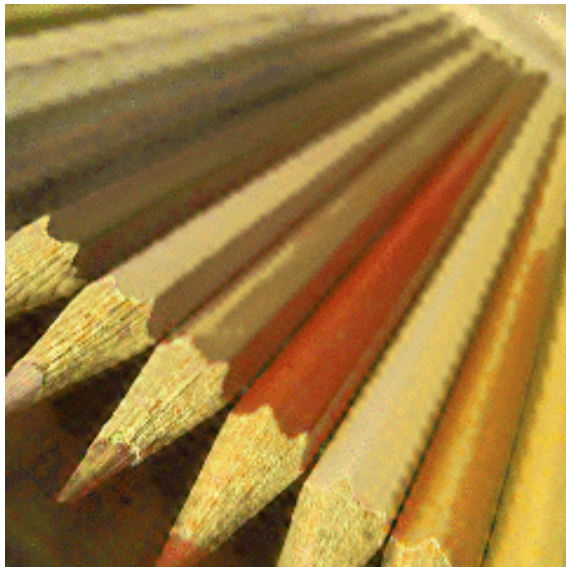
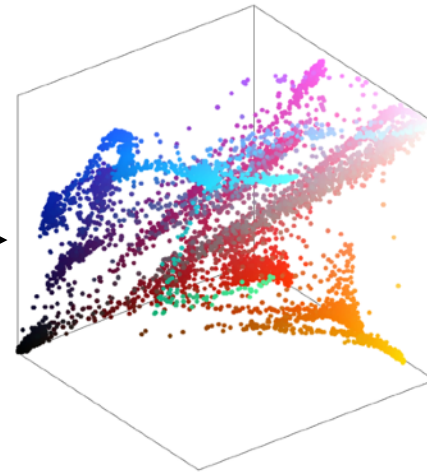
Optimal  
transport



# Color Image Palette Equalization



Optimal  
transport



# Overview

---

- Monge Formulation
- **Continuous Optimal Transport**
- Kantorovitch Formulation
- Applications

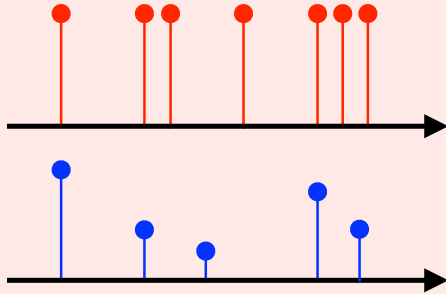
# Probability Measures

Positive Radon measure  $\alpha$  on a metric space  $\mathcal{X}$ .

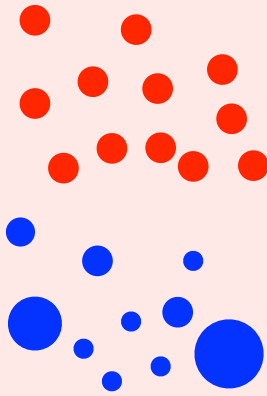
$$d\alpha(x) = \rho_\alpha(x) dx$$

$$\alpha = \sum_i \mathbf{a}_i \delta_{x_i}$$

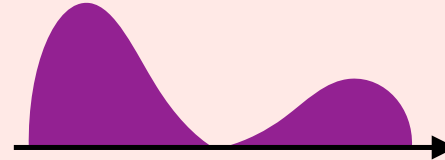
$\mathcal{X} = \mathbb{R}^d$



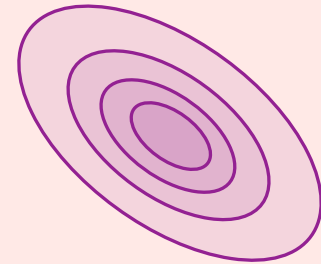
Discrete  $d = 1$



Discrete  $d = 2$



Density  $d = 1$



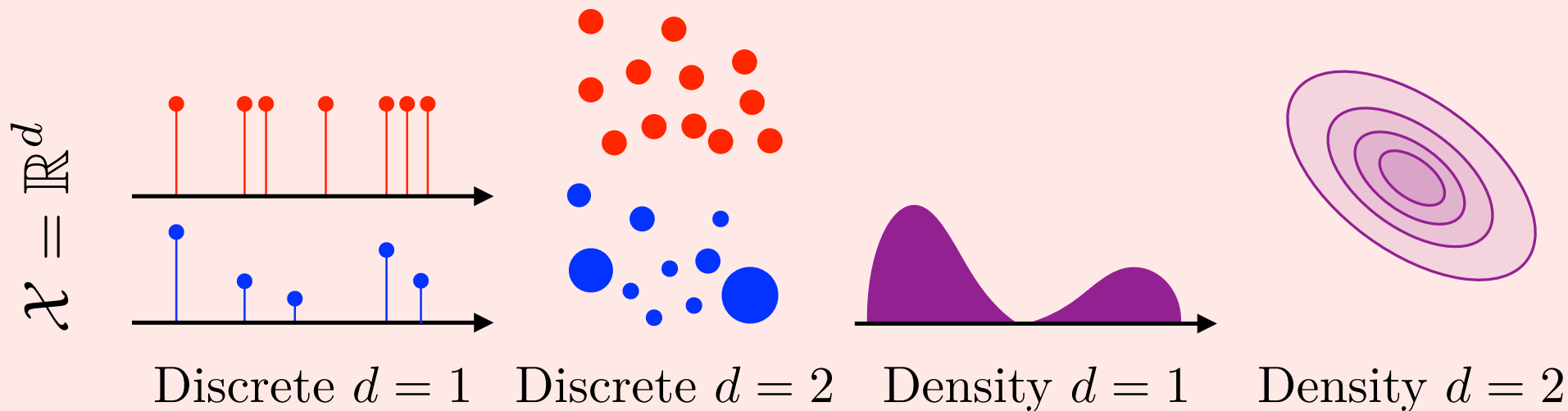
Density  $d = 2$

# Probability Measures

Positive Radon measure  $\alpha$  on a metric space  $\mathcal{X}$ .

$$d\alpha(x) = \rho_\alpha(x) dx$$

$$\alpha = \sum_i \mathbf{a}_i \delta_{x_i}$$



Measure of sets  $A \subset \mathcal{X}$ :  $\alpha(A) = \int_A d\alpha(x) \geq 0$

Integration against continuous functions:  $\int_{\mathcal{X}} g(x) d\alpha(x) \geq 0$

$$d\alpha(x) = \rho_\alpha(x) dx \longrightarrow \int_{\mathcal{X}} g d\alpha = \int_{\mathcal{X}} g(x) \rho_\alpha(x) dx$$

$$\alpha = \sum_i \mathbf{a}_i \delta_{x_i} \longrightarrow \int_{\mathcal{X}} g d\alpha = \sum_i \mathbf{a}_i g(x_i)$$

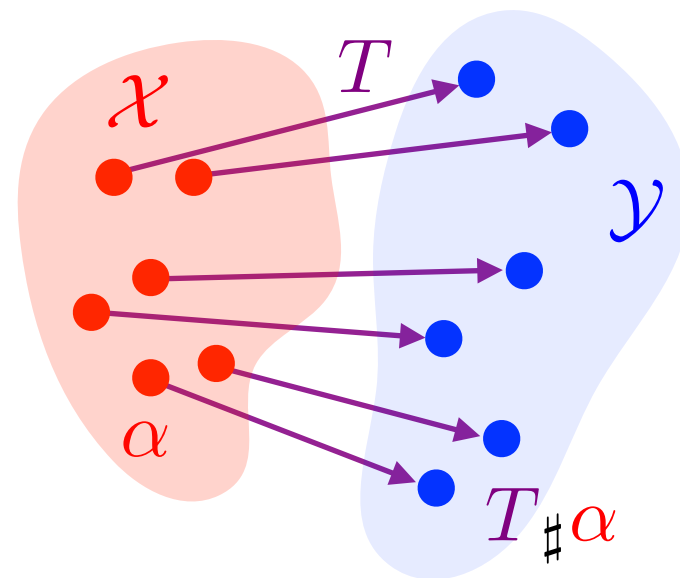
Probability (normalized) measure:  $\alpha(\mathcal{X}) = \int_{\mathcal{X}} d\alpha(x) = 1$

# Push Forward

Map:  $T : \mathcal{X} \rightarrow \mathcal{Y}$

Push-forward:

$$T_{\#} : \begin{cases} \delta_{\mathbf{x}} \longmapsto \delta_{T(\mathbf{x})} \\ \sum_i \delta_{\mathbf{x}_i} \longmapsto \sum_i \delta_{T(\mathbf{x}_i)} \\ \sum_i \mathbf{a}_i \delta_{\mathbf{x}_i} \longmapsto \sum_i \mathbf{a}_i \delta_{T(\mathbf{x}_i)} \end{cases}$$



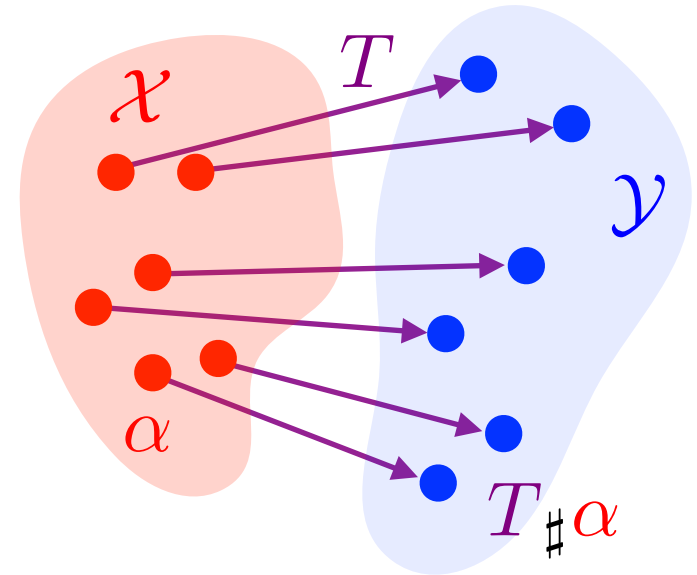


# Push Forward

Map:  $T : \mathcal{X} \rightarrow \mathcal{Y}$

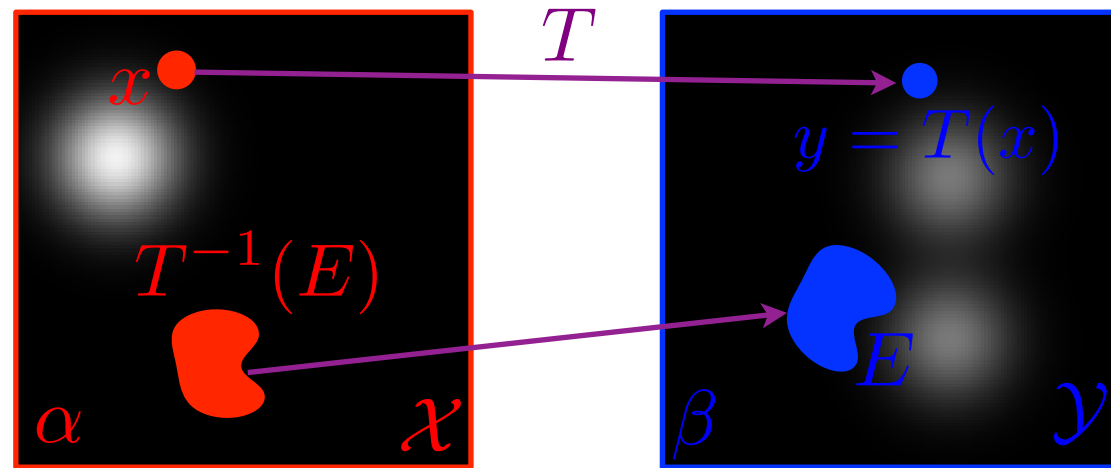
Push-forward:

$$T_{\#} : \begin{cases} \delta_{\mathbf{x}} \mapsto \delta_{T(\mathbf{x})} \\ \sum_i \delta_{\mathbf{x}_i} \mapsto \sum_i \delta_{T(\mathbf{x}_i)} \\ \sum_i \mathbf{a}_i \delta_{\mathbf{x}_i} \mapsto \sum_i \mathbf{a}_i \delta_{T(\mathbf{x}_i)} \end{cases}$$



General case:

$$(T_{\#}\alpha)(E) \stackrel{\text{def.}}{=} \alpha(T^{-1}(E))$$

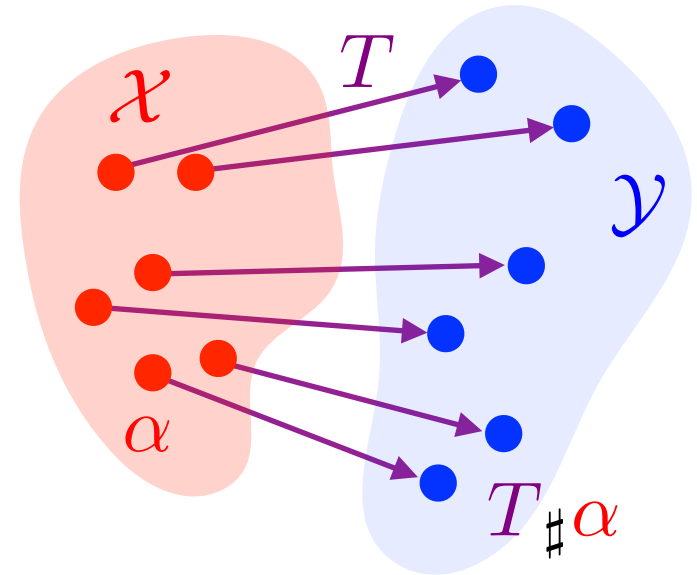


# Push Forward

Map:  $T : \mathcal{X} \rightarrow \mathcal{Y}$

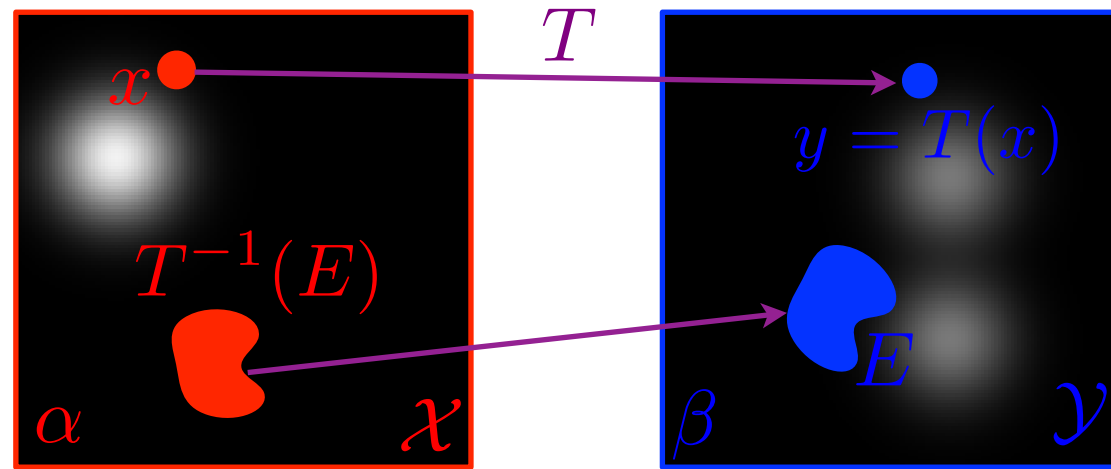
Push-forward:

$$T_{\#} : \begin{cases} \delta_{\mathbf{x}} \longmapsto \delta_{T(\mathbf{x})} \\ \sum_i \delta_{\mathbf{x}_i} \longmapsto \sum_i \delta_{T(\mathbf{x}_i)} \\ \sum_i \mathbf{a}_i \delta_{\mathbf{x}_i} \longmapsto \sum_i \mathbf{a}_i \delta_{T(\mathbf{x}_i)} \end{cases}$$



General case:

$$(T_{\#}\alpha)(E) \stackrel{\text{def.}}{=} \alpha(T^{-1}(E))$$



Change of variables:

$$\beta = T_{\#}\alpha \iff \int_{\mathcal{Y}} g(y) d\beta(y) = \int_{\mathcal{X}} g(T(x)) d\alpha(x)$$

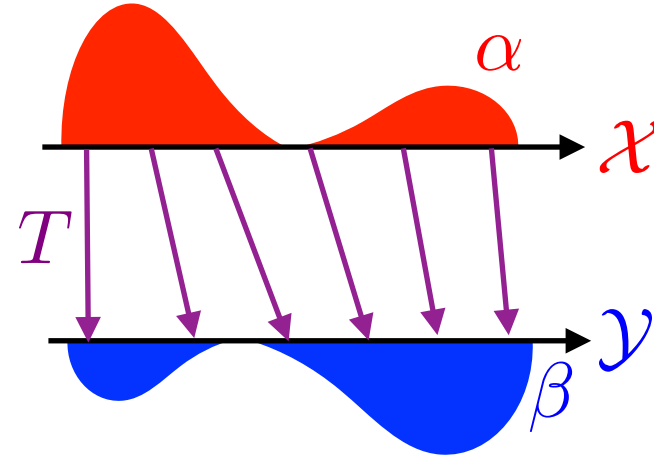
Densities  $\frac{d\alpha}{dx} = \rho_{\alpha}$ :

$$\rho_{\alpha}(x) = |\det(\partial T(x))| \rho_{\beta}(T(x))$$

# Continuous Monge's Problem



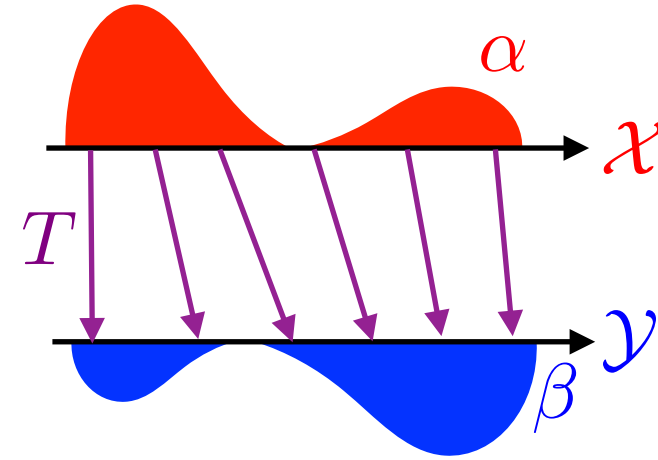
$$\inf_{\beta = T\# \alpha} \int_{\mathcal{X}} c(x, T(x)) d\alpha(x)$$



# Continuous Monge's Problem



$$\inf_{\beta = T\# \alpha} \int_{\mathcal{X}} c(x, T(x)) d\alpha(x)$$



*Discrete case:*

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

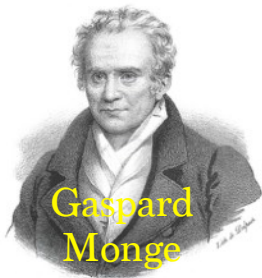
$$\beta = \sum_{j=1}^n \delta_{y_j}$$

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

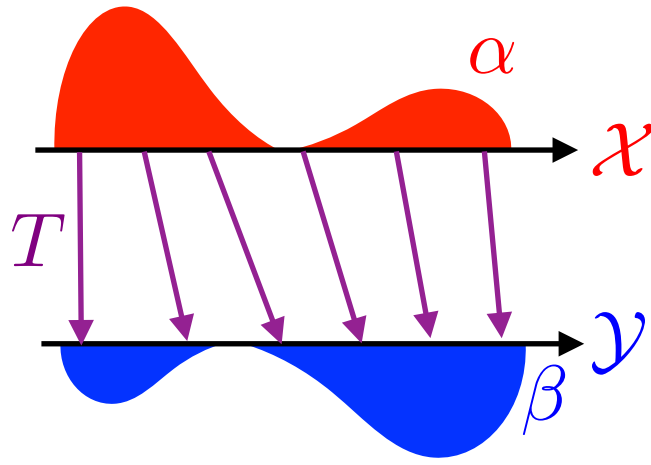
$$T : x_i \mapsto y_{\sigma(i)}$$

$$C_{i,j} = c(x_i, y_j)$$

# Continuous Monge's Problem



$$\inf_{\beta = T\# \alpha} \int_{\mathcal{X}} c(x, T(x)) d\alpha(x)$$



Discrete case:

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

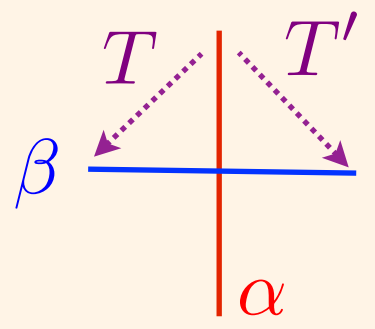
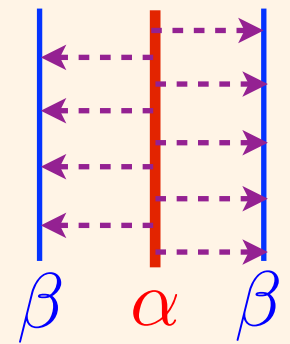
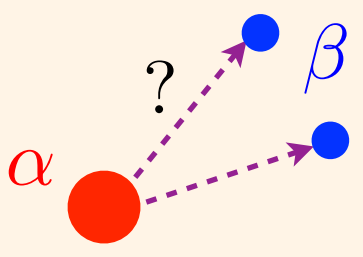
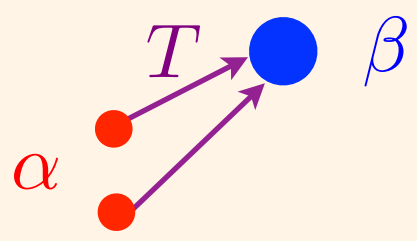
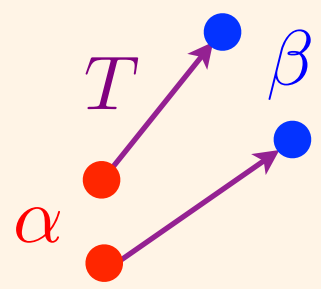
$$\beta = \sum_{j=1}^n \delta_{y_j}$$

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

$$T : x_i \mapsto y_{\sigma(i)}$$

$$C_{i,j} = c(x_i, y_j)$$

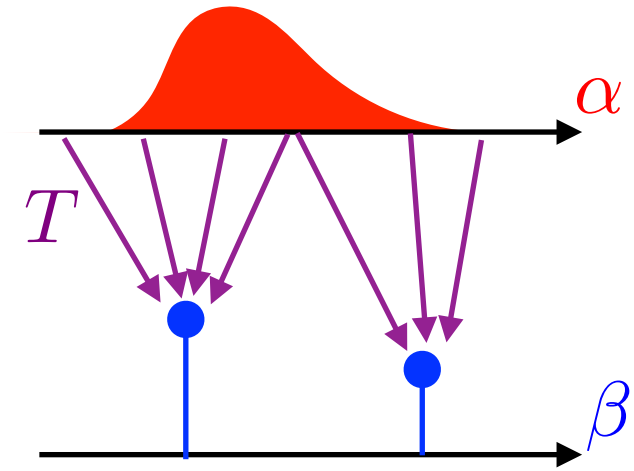
Non-symmetry, non-existence, non-uniqueness:



# Brenier's Theorem

*Hypotheses:*  $c(x, y) = \|x - y\|^2$   
 $\mathcal{X} = \mathbb{R}^d$   $\frac{d\alpha}{dx} = \rho_\alpha$  density.

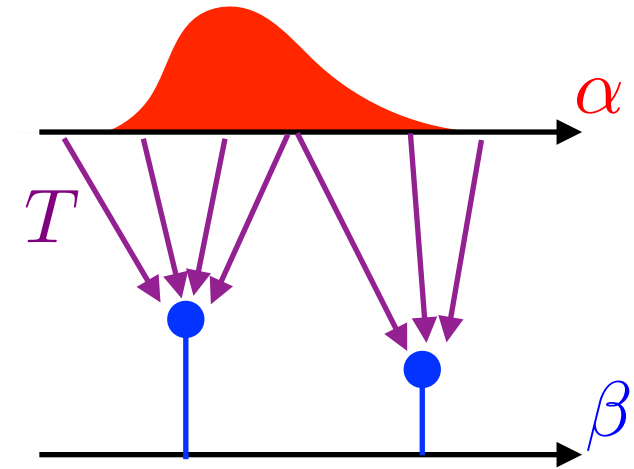
$$W_2^2(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\beta = T\# \alpha} \int_{\mathbb{R}^d} \|x - T(x)\|^2 d\alpha(x)$$



# Brenier's Theorem

*Hypotheses:*  $c(x, y) = \|x - y\|^2$   
 $\mathcal{X} = \mathbb{R}^d$   $\frac{d\alpha}{dx} = \rho_\alpha$  density.

$$W_2^2(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\beta = T\# \alpha} \int_{\mathbb{R}^d} \|x - T(x)\|^2 d\alpha(x)$$



*Theorem:* [Brenier, 1991]

There exists a unique Monge map  $T$ .

It is the unique  $T = \nabla \varphi$  such that

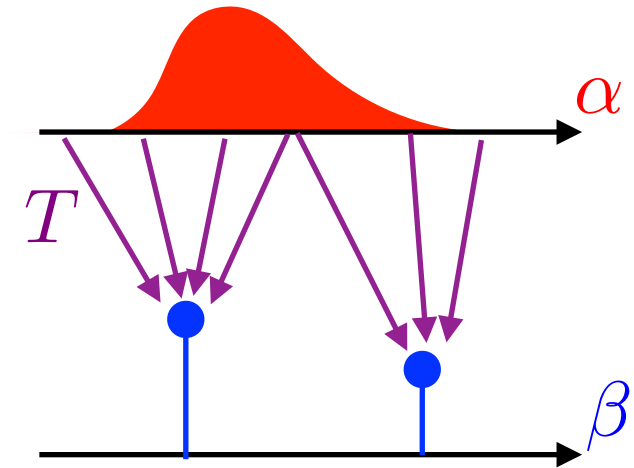
$\varphi$  is convex and  $(\nabla \varphi)\# \alpha = \beta$ .



# Brenier's Theorem

*Hypotheses:*  $c(x, y) = \|x - y\|^2$   
 $\mathcal{X} = \mathbb{R}^d$   $\frac{d\alpha}{dx} = \rho_\alpha$  density.

$$W_2^2(\alpha, \beta) \stackrel{\text{def.}}{=} \inf_{\beta = T\# \alpha} \int_{\mathbb{R}^d} \|x - T(x)\|^2 d\alpha(x)$$



*Theorem:* [Brenier, 1991]

There exists a unique Monge map  $T$ .

It is the unique  $T = \nabla \varphi$  such that

$\varphi$  is convex and  $(\nabla \varphi)\# \alpha = \beta$ .



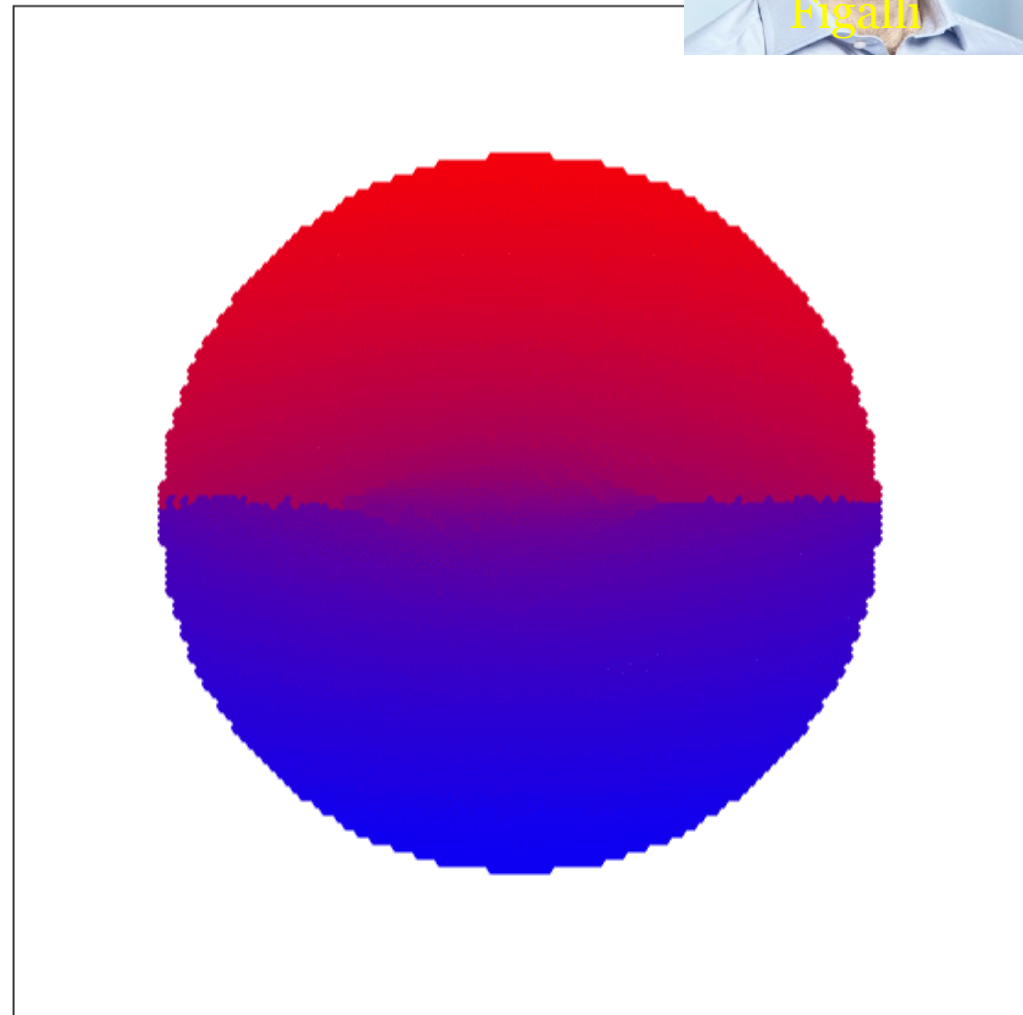
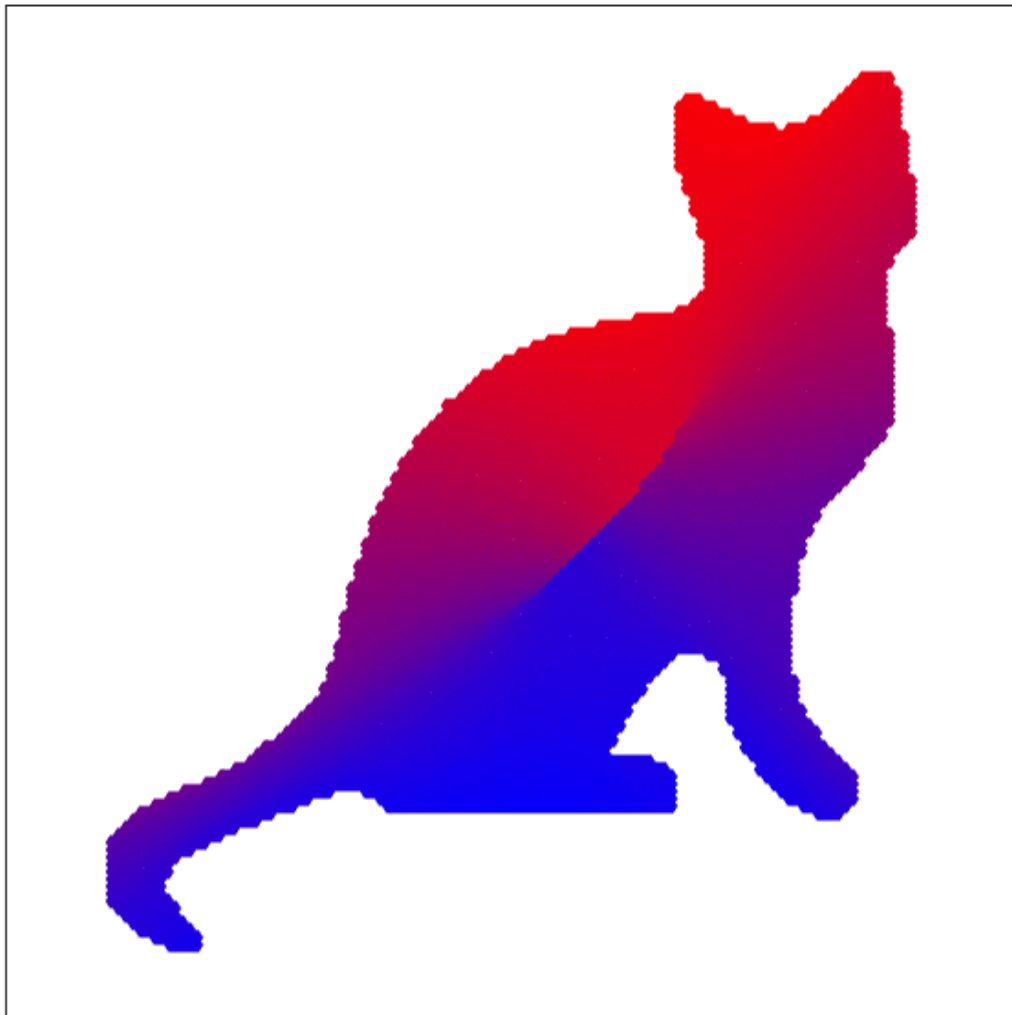
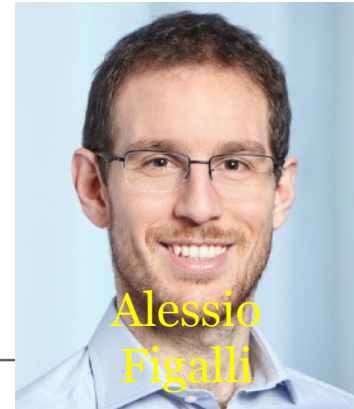
→ Monge-Ampère equation (non-linear, degenerate elliptic).

$$\rho_\alpha(x) = |\det(\partial^2 \varphi(x))| \rho_\beta(T(x)) \quad \text{s.t. } \varphi \text{ convex.}$$



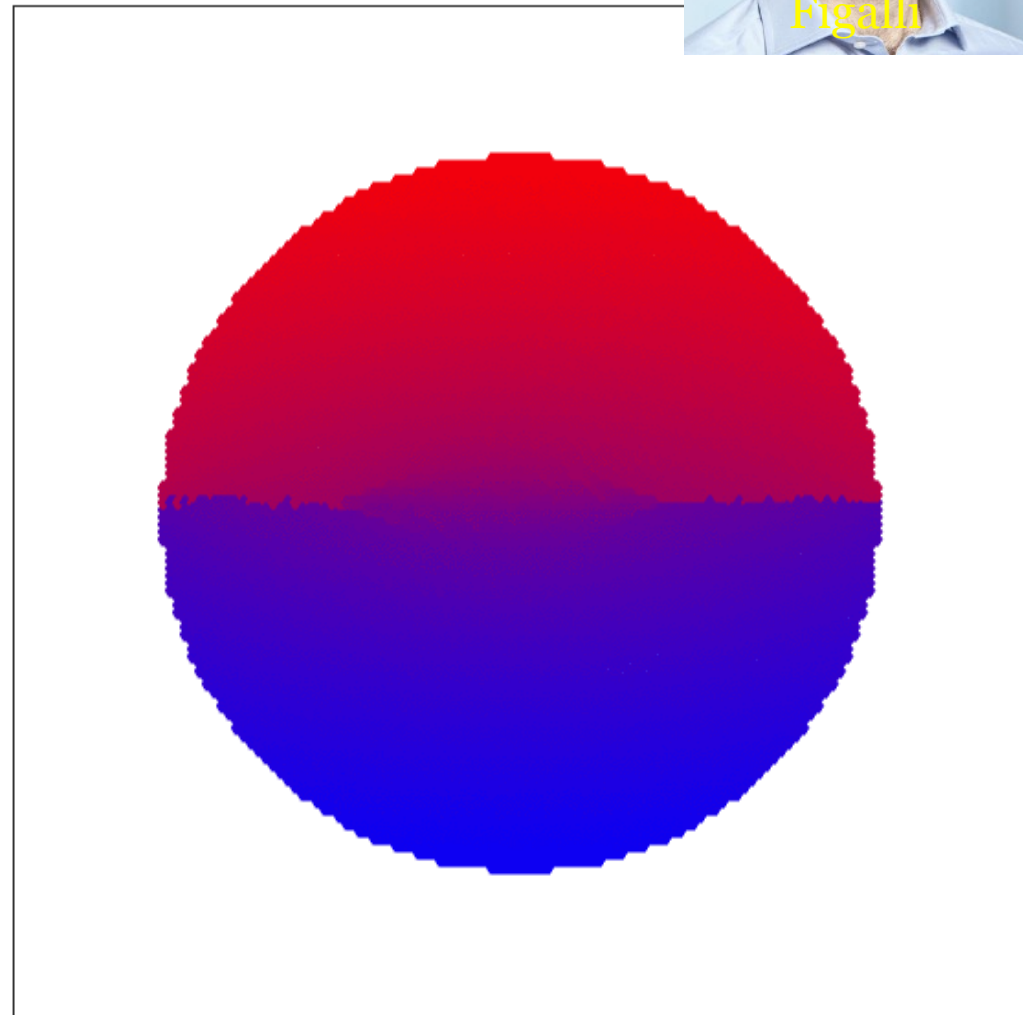
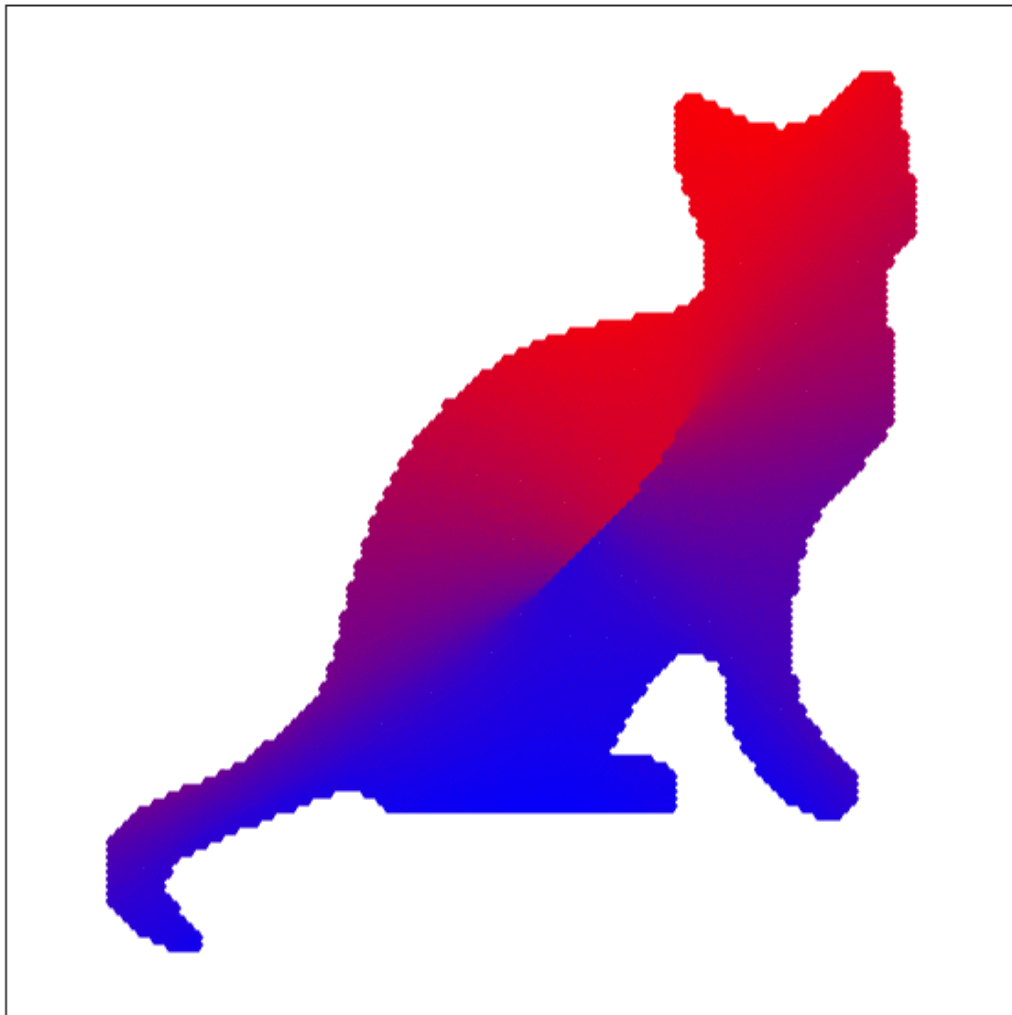
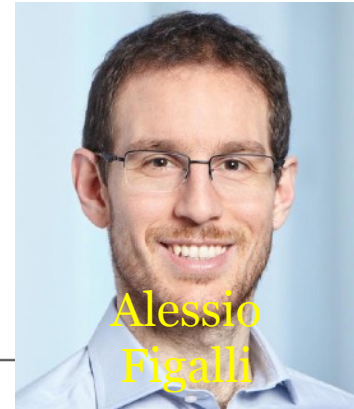
# Regularity Theory

→ Regularity of  $T$  requires convex target.



# Regularity Theory

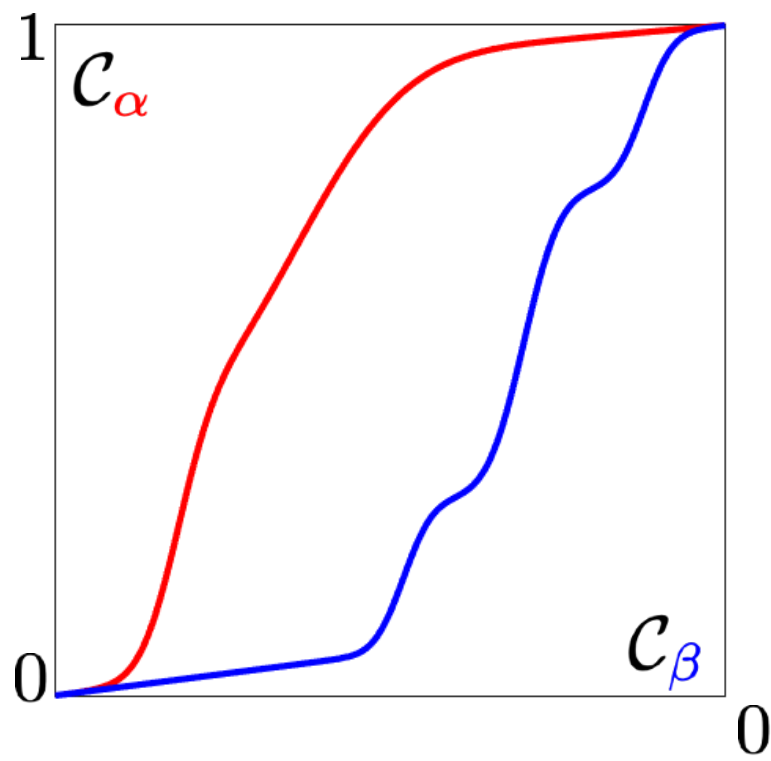
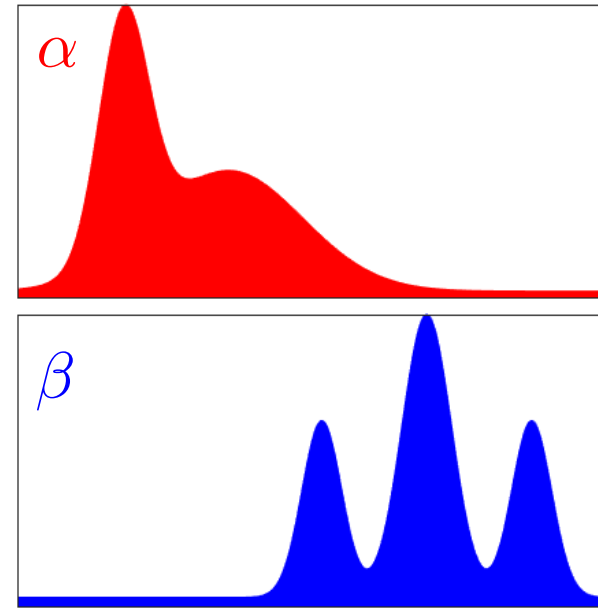
→ Regularity of  $T$  requires convex target.



# 1-D Optimal Transport

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

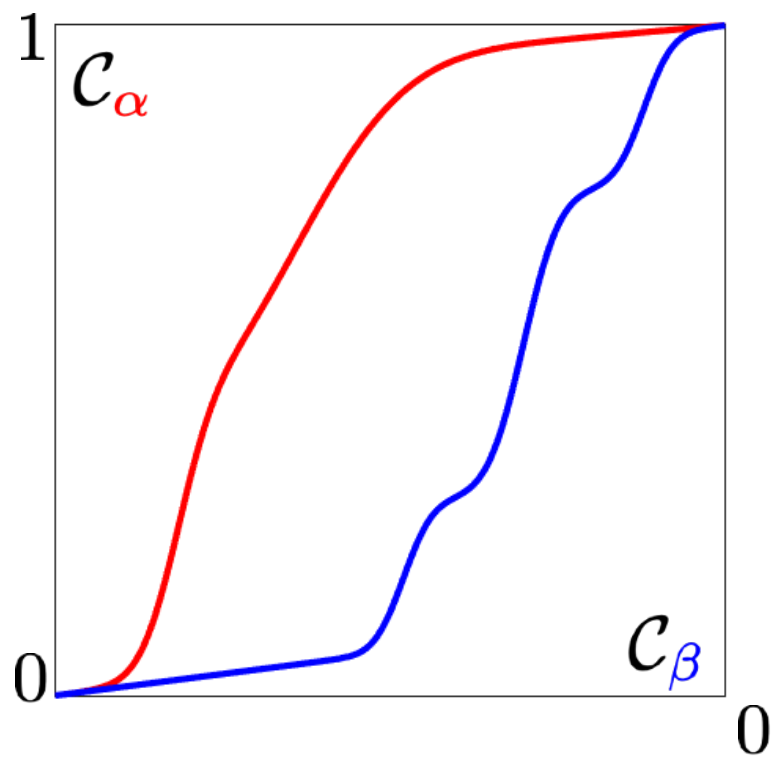
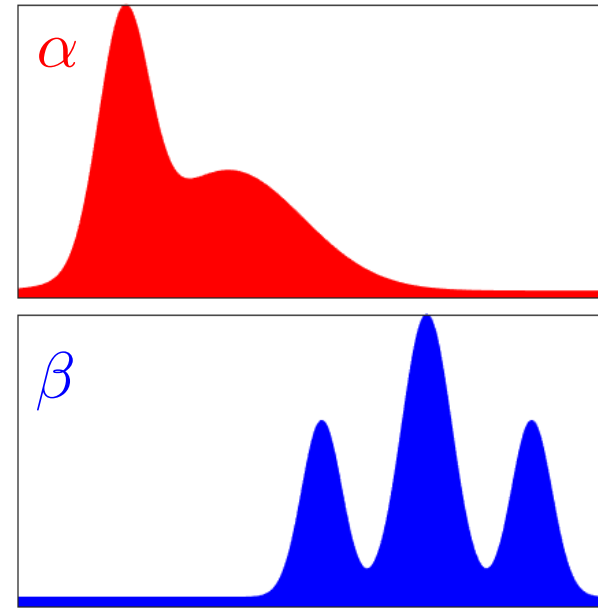
Cumulative function:  $\mathcal{C}_{\alpha\#} : \alpha \mapsto \mathcal{U}_{[0,1]}$



# 1-D Optimal Transport

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

Cumulative function:  $\mathcal{C}_{\alpha\#} : \alpha \mapsto \mathcal{U}_{[0,1]}$

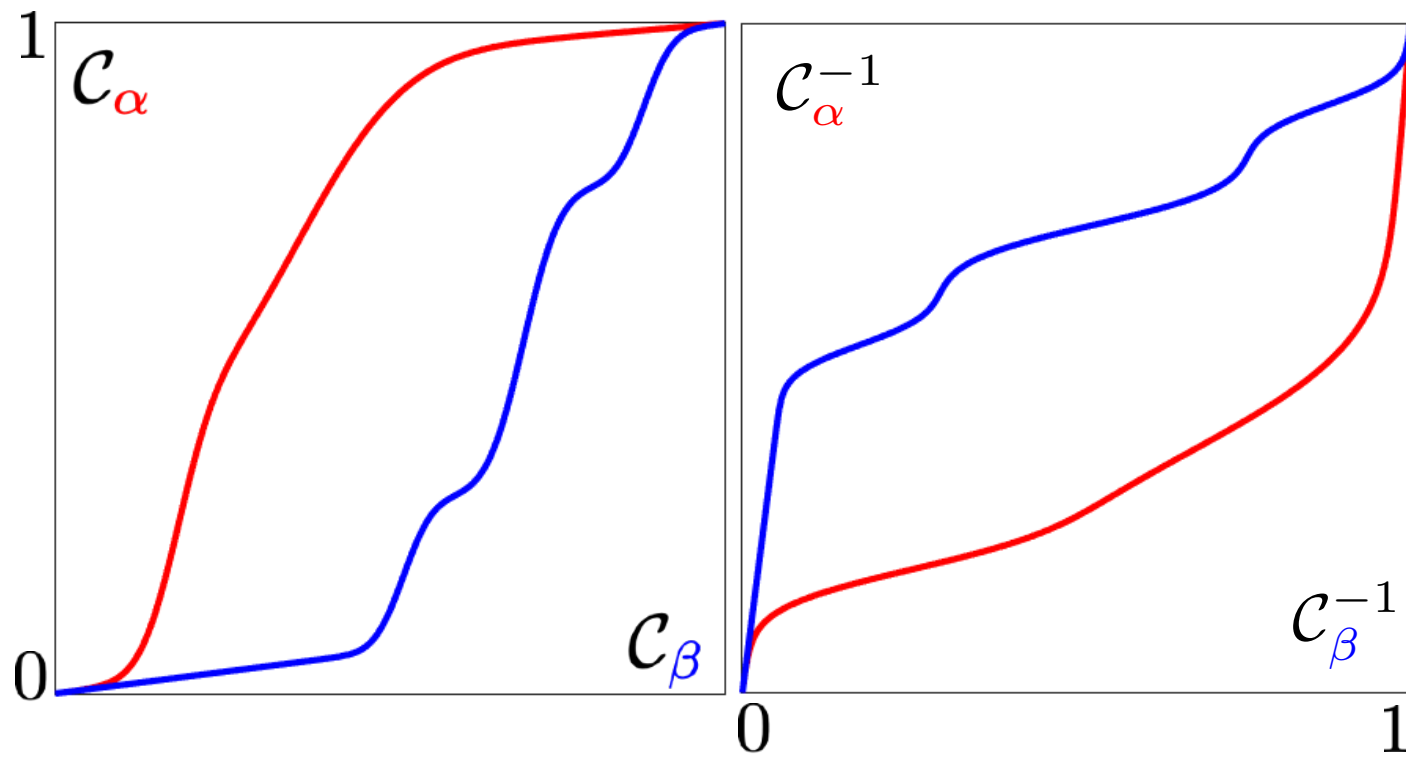
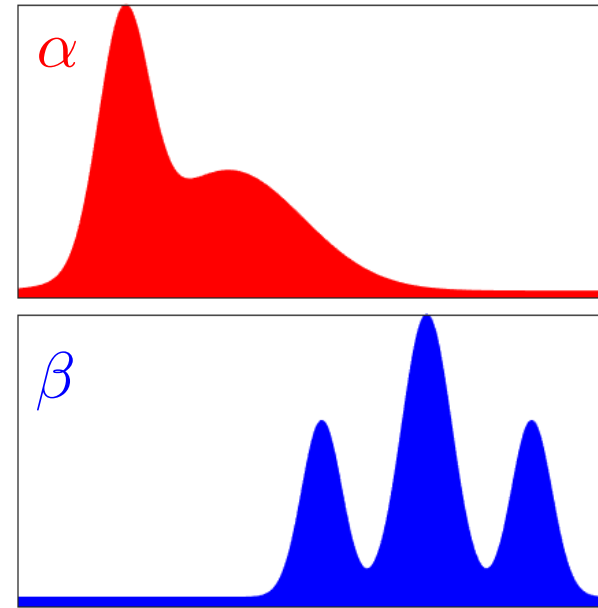


# 1-D Optimal Transport

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

Cumulative function:  $\mathcal{C}_{\alpha\#} : \alpha \mapsto \mathcal{U}_{[0,1]}$

Quantile function:  $\mathcal{C}_\beta^{-1\#} : \mathcal{U}_{[0,1]} \mapsto \beta$



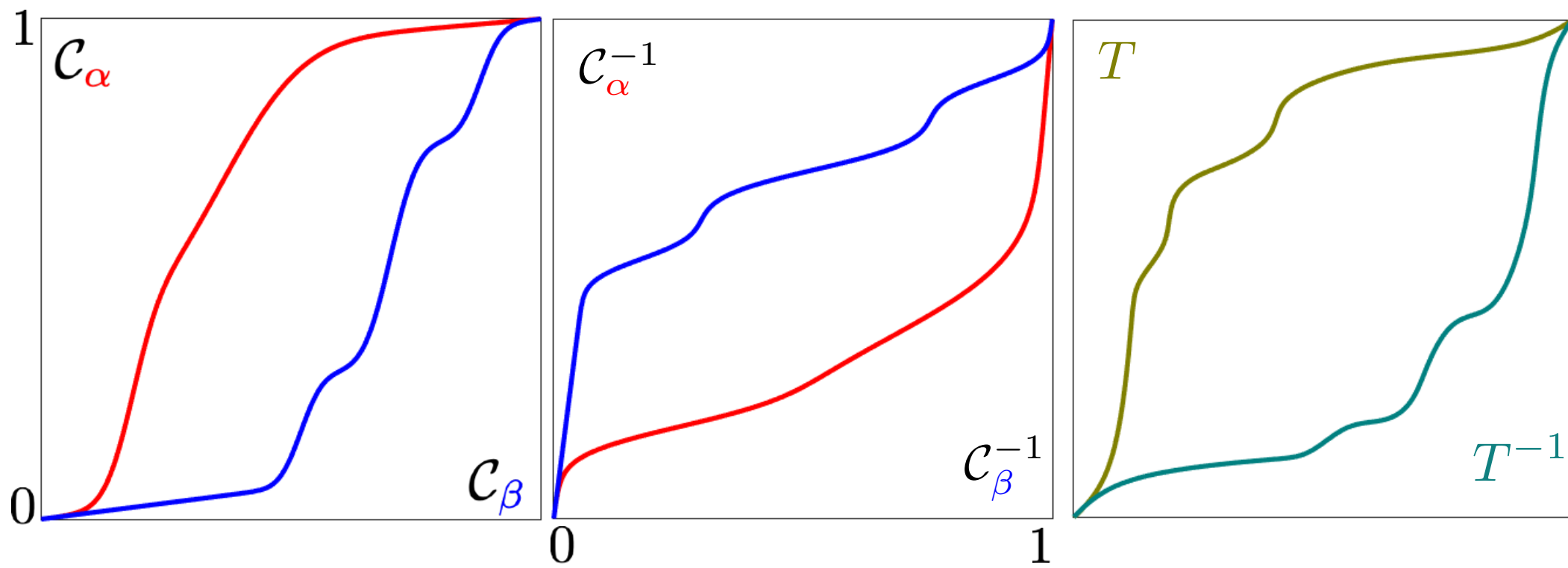
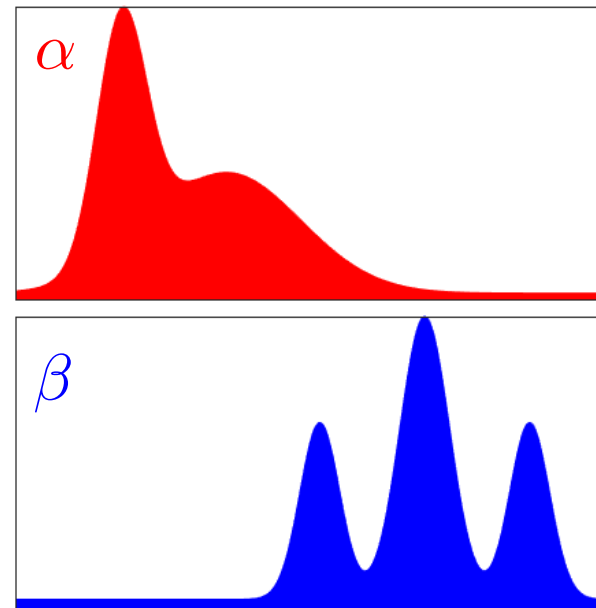
# 1-D Optimal Transport

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

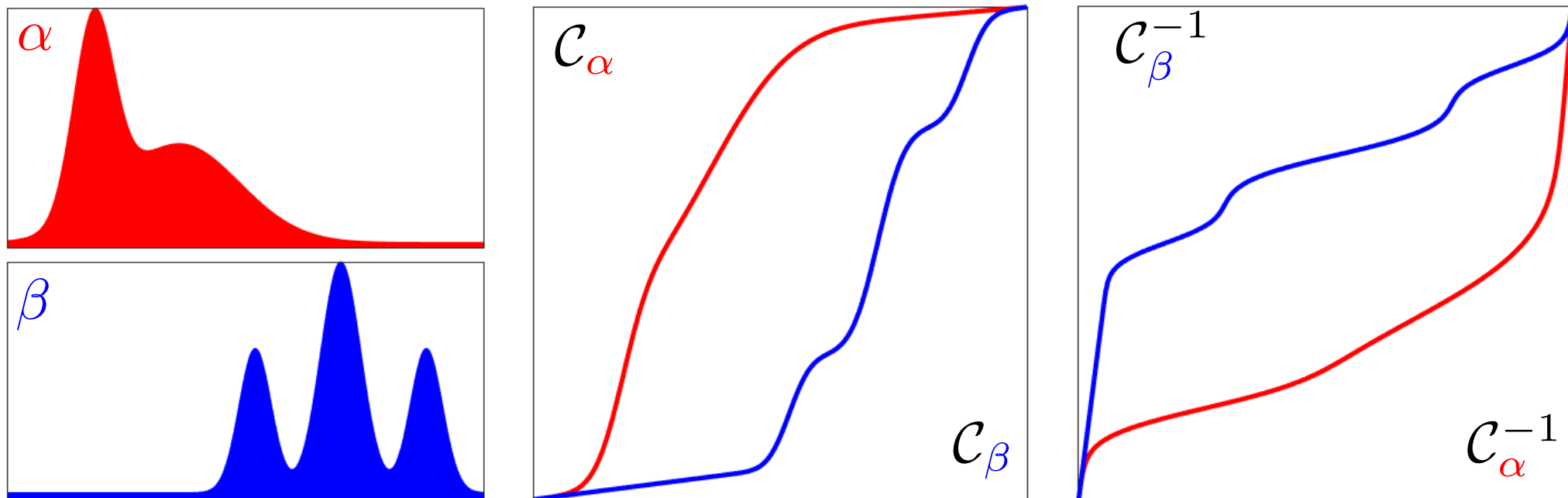
Cumulative function:  $\mathcal{C}_{\alpha\#} : \alpha \mapsto \mathcal{U}_{[0,1]}$

Quantile function:  $\mathcal{C}_\beta^{-1\#} : \mathcal{U}_{[0,1]} \mapsto \beta$

Optimal transport  $\alpha \mapsto \beta$ :  $T = \mathcal{C}_\beta^{-1} \circ \mathcal{C}_\alpha$



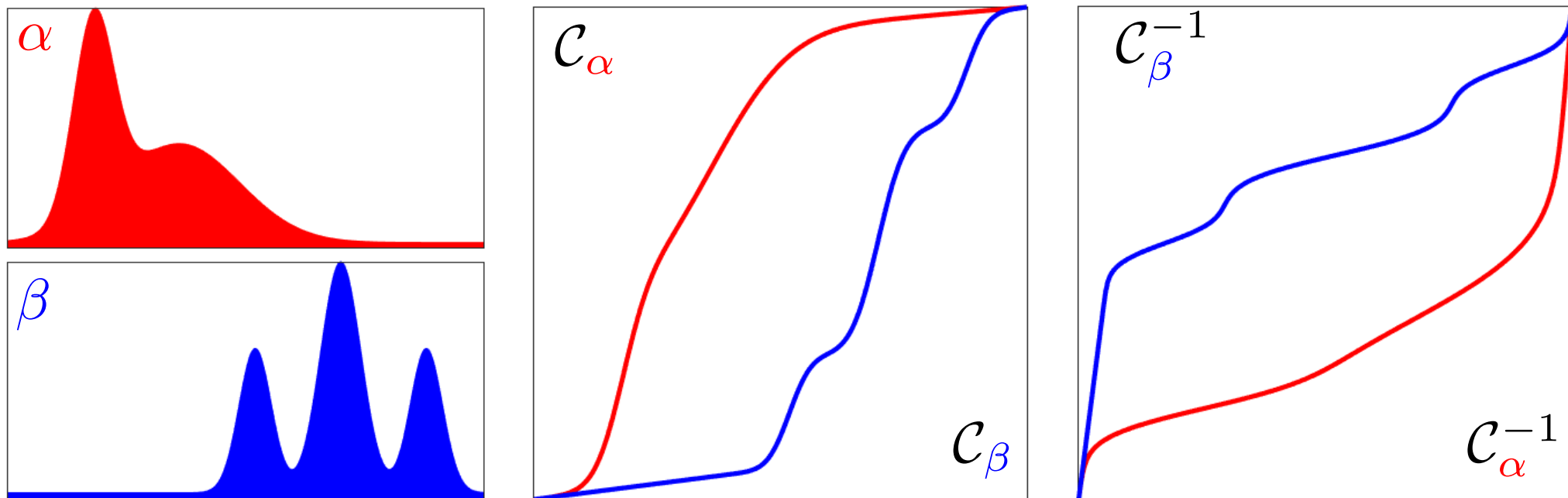
Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$



$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \int \|T(x) - x\|^p d\alpha(x) = \int_0^1 |\mathcal{C}_\alpha^{-1}(t) - \mathcal{C}_\beta^{-1}(t)|^p dt$$

$$W_1(\alpha, \beta) = \|\alpha - \beta\|_{W_1} = \int_{\mathbb{R}} |\mathcal{C}_\alpha(x) - \mathcal{C}_\beta(x)| dx$$

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

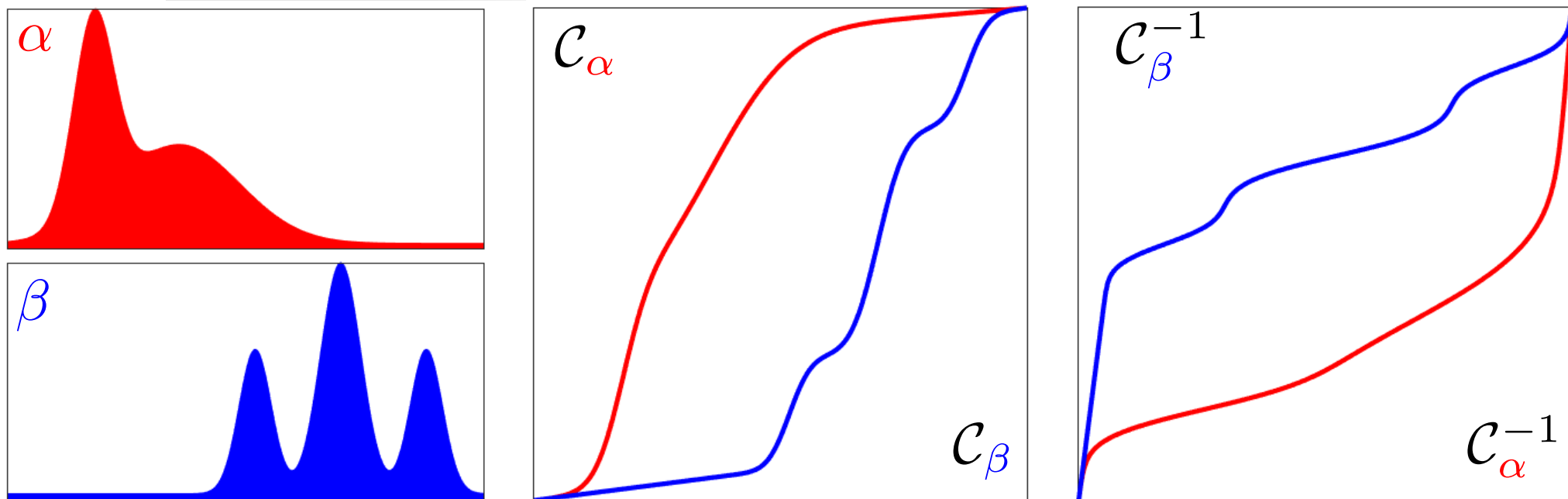


$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \int \|T(x) - x\|^p d\alpha(x) = \int_0^1 |\mathcal{C}_\alpha^{-1}(t) - \mathcal{C}_\beta^{-1}(t)|^p dt$$

$$W_1(\alpha, \beta) = \|\alpha - \beta\|_{W_1} = \int_{\mathbb{R}} |\mathcal{C}_\alpha(x) - \mathcal{C}_\beta(x)| dx$$



Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$



$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \int \|T(x) - x\|^p d\alpha(x) = \int_0^1 |\mathcal{C}_\alpha^{-1}(t) - \mathcal{C}_\beta^{-1}(t)|^p dt$$

$$W_1(\alpha, \beta) = \|\alpha - \beta\|_{W_1} = \int_{\mathbb{R}} |\mathcal{C}_\alpha(x) - \mathcal{C}_\beta(x)| dx$$

$$\text{Kramer (Sobolev) norm: } \|\alpha - \beta\|_K^2 = \int_0^1 |\mathcal{C}_\alpha(t) - \mathcal{C}_\beta(t)|^2 dt$$

$$\text{Kolmogorov-Smirnov norm: } \|\alpha - \beta\|_{KS} = \sup_x |\mathcal{C}_\alpha(x) - \mathcal{C}_\beta(x)|$$

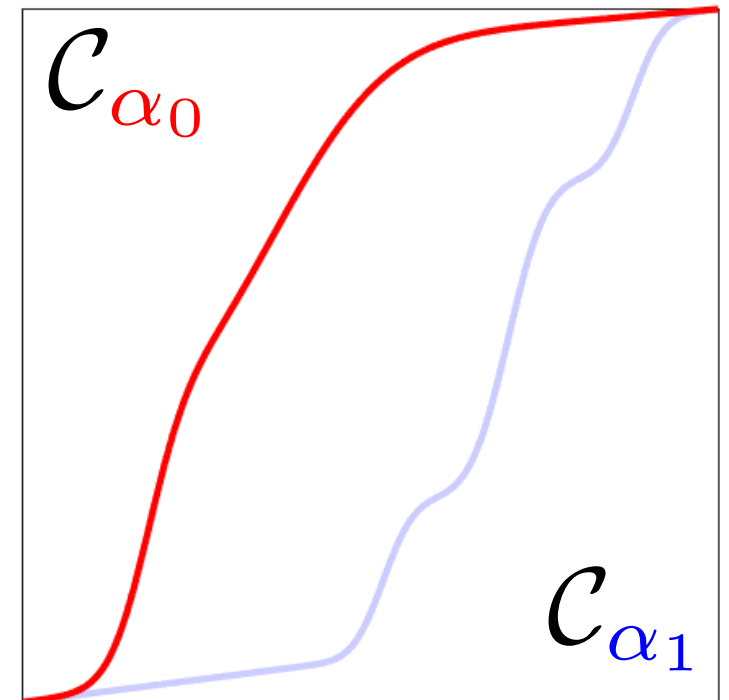
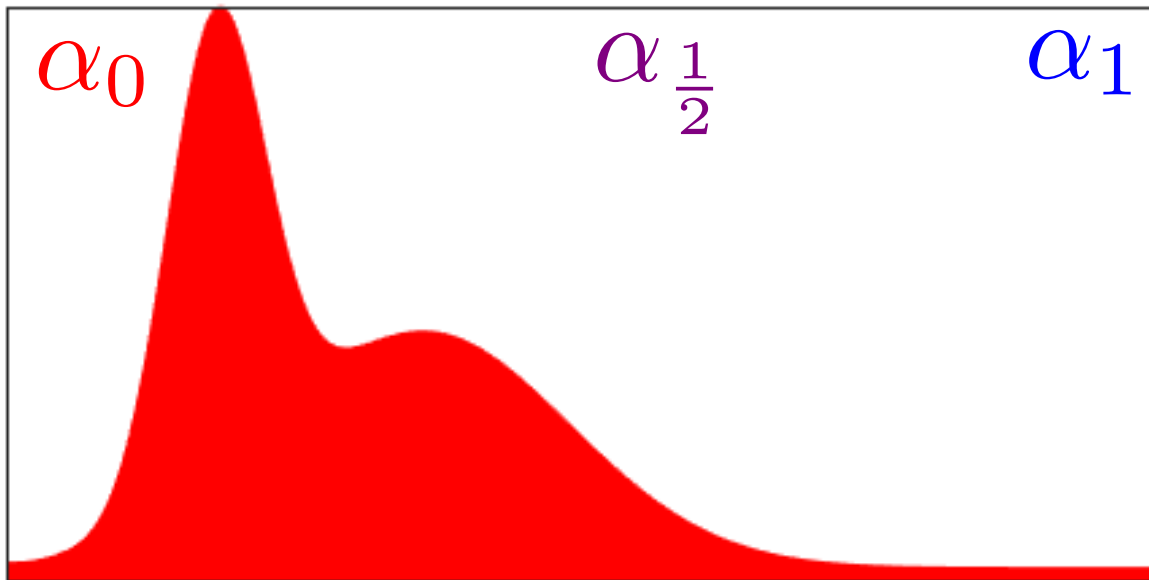
$$\text{Area under the curve: } \text{AUC}(\alpha, \beta) = 1 - \int_0^1 \mathcal{C}_\alpha \circ \mathcal{C}_\beta^{-1}(x)$$

# 1-D Optimal Transport Interpolation

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

Optimal transport interpolation  $\alpha_0 \leftrightarrow \alpha_1$

$$\forall t \in [0, 1], \mathcal{C}_{\alpha_t}^{-1} = (1-t)\mathcal{C}_{\alpha_0}^{-1} + t\mathcal{C}_{\alpha_1}^{-1}$$

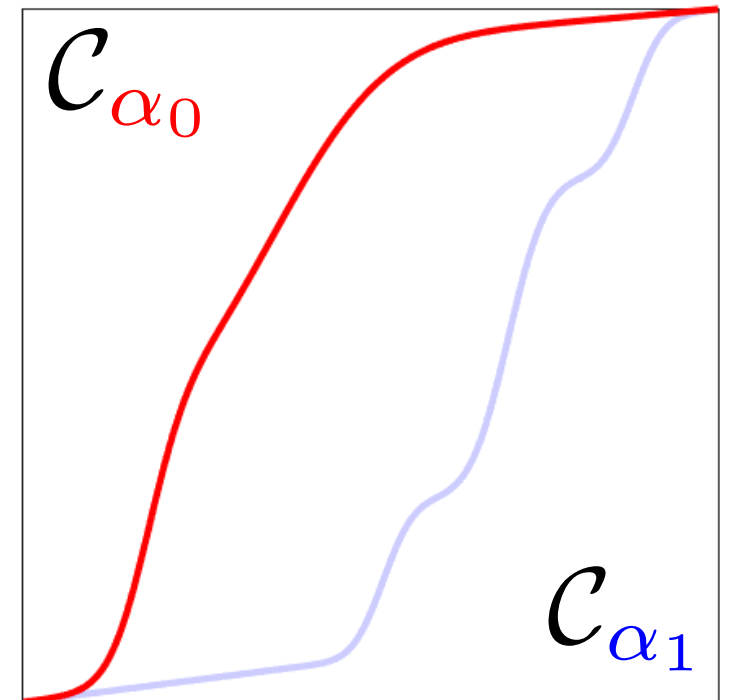
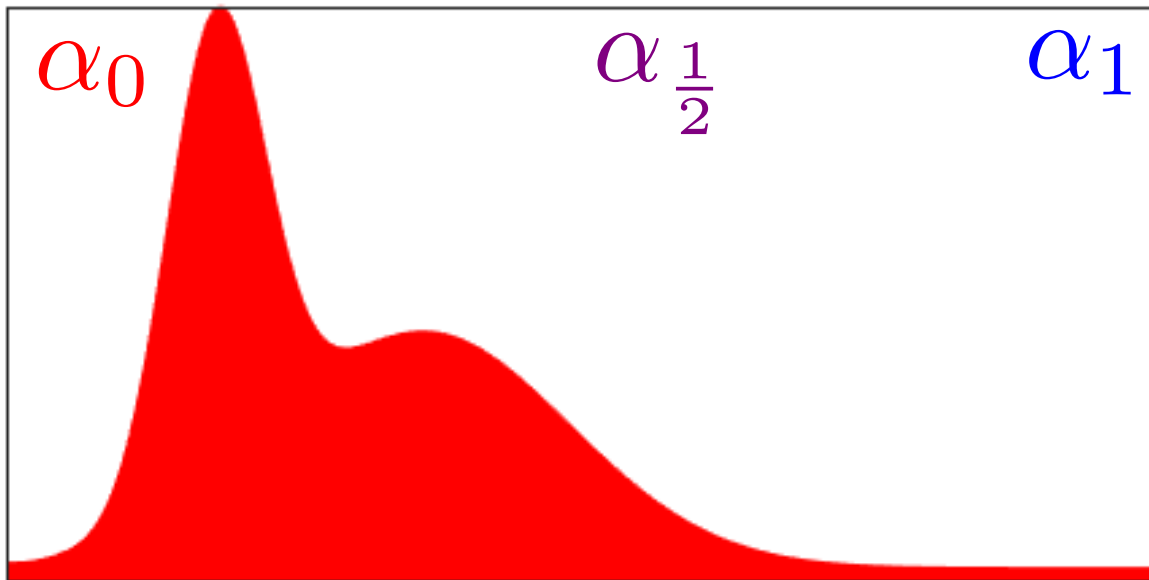


# 1-D Optimal Transport Interpolation

Cumulative function:  $\mathcal{C}_\alpha(x) \stackrel{\text{def.}}{=} \int_{-\infty}^x d\alpha$

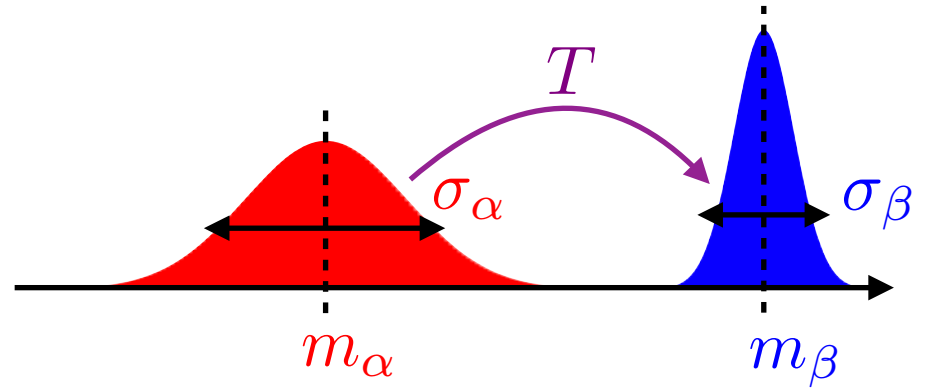
Optimal transport interpolation  $\alpha_0 \leftrightarrow \alpha_1$

$$\forall t \in [0, 1], \mathcal{C}_{\alpha_t}^{-1} = (1-t)\mathcal{C}_{\alpha_0}^{-1} + t\mathcal{C}_{\alpha_1}^{-1}$$



# OT Between 1D Gaussians

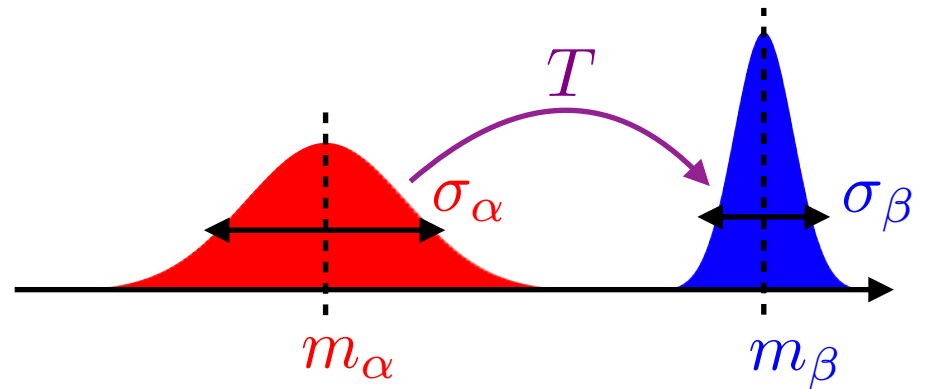
$$\frac{d\alpha}{dx} = \frac{1}{\sigma_\alpha \sqrt{2\pi}} e^{-\frac{(x-m_\alpha)^2}{2\sigma_\alpha^2}}$$



$$T(x) = \frac{\sigma_\beta}{\sigma_\alpha} (x - m_\alpha) + m_\beta$$

# OT Between 1D Gaussians

$$\frac{d\alpha}{dx} = \frac{1}{\sigma_\alpha \sqrt{2\pi}} e^{-\frac{(x-m_\alpha)^2}{2\sigma_\alpha^2}}$$



$$T(x) = \frac{\sigma_\beta}{\sigma_\alpha} (x - m_\alpha) + m_\beta$$

$$T = \nabla \varphi \quad \varphi \text{ is convex.}$$

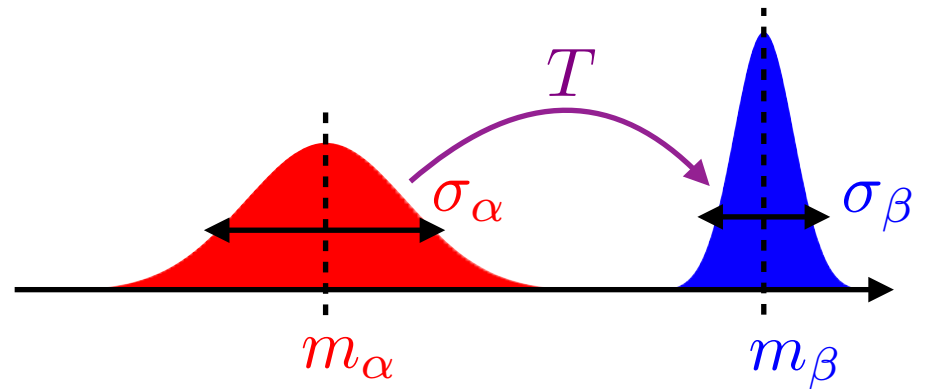
Brenier  
 $\implies$

$$T \equiv \text{OT}$$

$$\varphi(x) = \frac{\sigma_\beta}{2\sigma_\alpha} (x - m_\alpha)^2 + m_\beta x$$

# OT Between 1D Gaussians

$$\frac{d\alpha}{dx} = \frac{1}{\sigma_\alpha \sqrt{2\pi}} e^{-\frac{(x-m_\alpha)^2}{2\sigma_\alpha^2}}$$



$$T(x) = \frac{\sigma_\beta}{\sigma_\alpha} (x - m_\alpha) + m_\beta$$

$$T = \nabla \varphi \quad \varphi \text{ is convex.}$$

Brenier  
 $\implies$

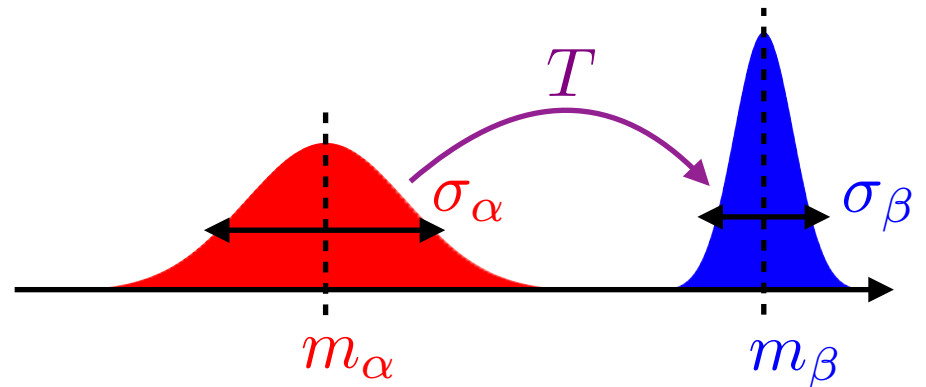
$$T \equiv \text{OT}$$

$$\varphi(x) = \frac{\sigma_\beta}{2\sigma_\alpha} (x - m_\alpha)^2 + m_\beta x$$

$$W_2^2(\alpha, \beta) = (m_\alpha - m_\beta)^2 + (\sigma_\alpha - \sigma_\beta)^2$$

# OT Between 1D Gaussians

$$\frac{d\alpha}{dx} = \frac{1}{\sigma_\alpha \sqrt{2\pi}} e^{-\frac{(x-m_\alpha)^2}{2\sigma_\alpha^2}}$$



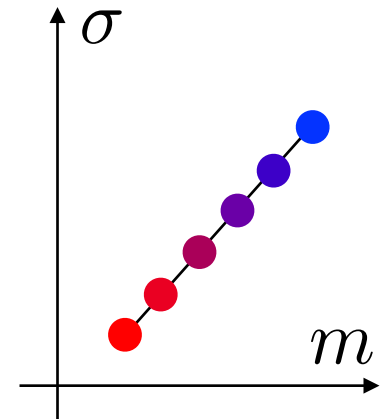
$$T(x) = \frac{\sigma_\beta}{\sigma_\alpha}(x - m_\alpha) + m_\beta$$

$$T = \nabla \varphi \quad \varphi \text{ is convex.}$$

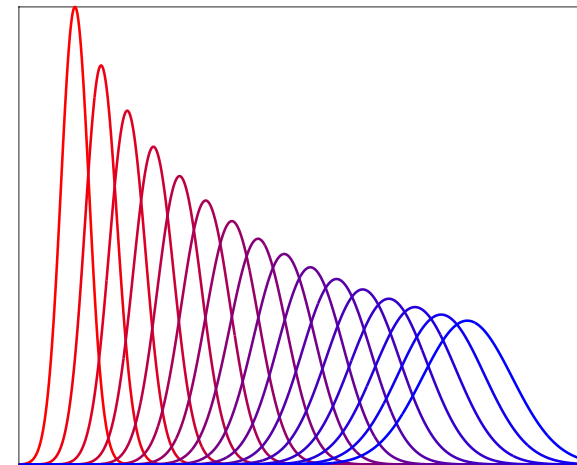
Brenier  
 $\implies$

$$T \equiv \text{OT}$$

$$\varphi(x) = \frac{\sigma_\beta}{2\sigma_\alpha}(x - m_\alpha)^2 + m_\beta x$$



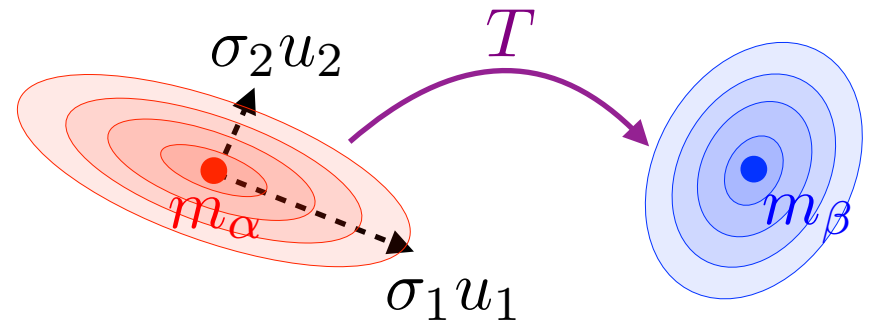
$$W_2^2(\alpha, \beta) = (m_\alpha - m_\beta)^2 + (\sigma_\alpha - \sigma_\beta)^2$$



# OT Between Gaussians

$$\frac{d\alpha}{dx} = \frac{1}{(2\pi)^{d/2} |\Sigma_\alpha|} e^{-\frac{\|x - m_\alpha\|_{\Sigma_\alpha^{-1}}^2}{2}}$$

$$\Sigma_\alpha = U_\alpha \text{diag}(\sigma_\alpha) U_\alpha^\top$$

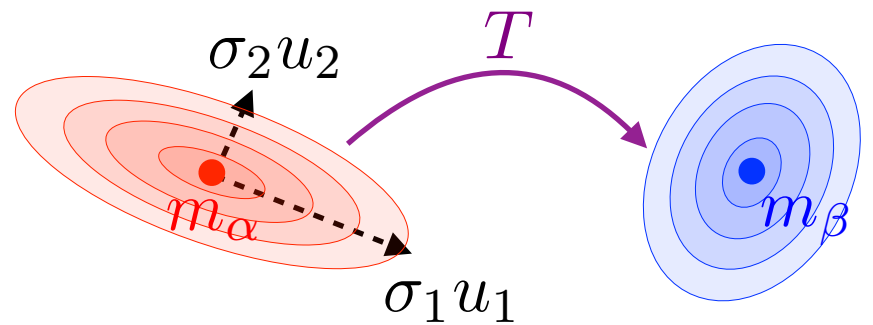




# OT Between Gaussians

$$\frac{d\alpha}{dx} = \frac{1}{(2\pi)^{d/2} |\Sigma_\alpha|} e^{-\frac{\|x - m_\alpha\|_{\Sigma_\alpha^{-1}}^2}{2}}$$

$$\Sigma_\alpha = U_\alpha \text{diag}(\sigma_\alpha) U_\alpha^\top$$



$$\text{Ansatz: } T(x) = A(x - m_\alpha) + m_\beta$$

*Proposition:*  $T$  is the optimal transport if

$$T = \nabla \varphi \Leftrightarrow \varphi(x) = \frac{1}{2} \langle A(x - m_\alpha), x - m_\alpha \rangle + \langle x, m_\beta \rangle + \text{cst}$$

$$\varphi \text{ convex} \Leftrightarrow A \in \mathcal{S}_+^n$$

$$T_\# \alpha = \beta \Leftrightarrow A \Sigma_\alpha A = \Sigma_\beta$$

# Resolution of Algebraic Riccati equations

$$A\Sigma_{\alpha}A = \Sigma_{\beta}$$

$$\left(\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}}\right)\left(\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}}\right) = \Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\beta}\Sigma_{\alpha}^{\frac{1}{2}}$$



Jacopo Riccati

*Proposition:* If  $\Sigma \in \mathcal{S}_+$ ,  $\exists! \sqrt{\Sigma} \in \mathcal{S}_+$  s.t.  $(\sqrt{\Sigma})^2 = \Sigma$ .

*Proof:* eigen-decomposition  $\Sigma = U \text{diag}(\sigma_i)U^{\top}$ , take  $\sqrt{\Sigma} = U \text{diag}(\sqrt{\sigma_i})U^{\top}$ .

Uniqueness: one has  $\sqrt{\Sigma}\Sigma = \sqrt{\Sigma}^3 = \Sigma\sqrt{\Sigma}$ , they co-diagonalize.

# Resolution of Algebraic Riccati equations

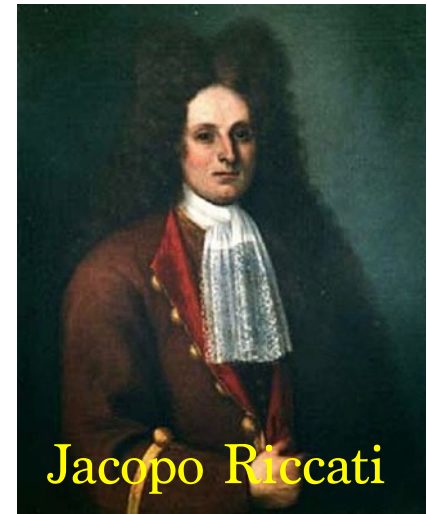
$$A\Sigma_{\alpha}A = \Sigma_{\beta}$$

$$\left(\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}}\right)\left(\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}}\right) = \Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\beta}\Sigma_{\alpha}^{\frac{1}{2}}$$

$$\left(\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}}\right)^2 = \Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\beta}\Sigma_{\alpha}^{\frac{1}{2}}$$

$$\Sigma_{\alpha}^{\frac{1}{2}}A\Sigma_{\alpha}^{\frac{1}{2}} = \sqrt{\Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\beta}\Sigma_{\alpha}^{\frac{1}{2}}}$$

$$A = \Sigma_{\alpha}^{-\frac{1}{2}}\sqrt{\Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\beta}\Sigma_{\alpha}^{\frac{1}{2}}\Sigma_{\alpha}^{-\frac{1}{2}}}$$



Jacopo Riccati

*Proposition:* If  $\Sigma \in \mathcal{S}_+$ ,  $\exists! \sqrt{\Sigma} \in \mathcal{S}_+$  s.t.  $(\sqrt{\Sigma})^2 = \Sigma$ .

*Proof:* eigen-decomposition  $\Sigma = U \text{diag}(\sigma_i)U^{\top}$ , take  $\sqrt{\Sigma} = U \text{diag}(\sqrt{\sigma_i})U^{\top}$ .

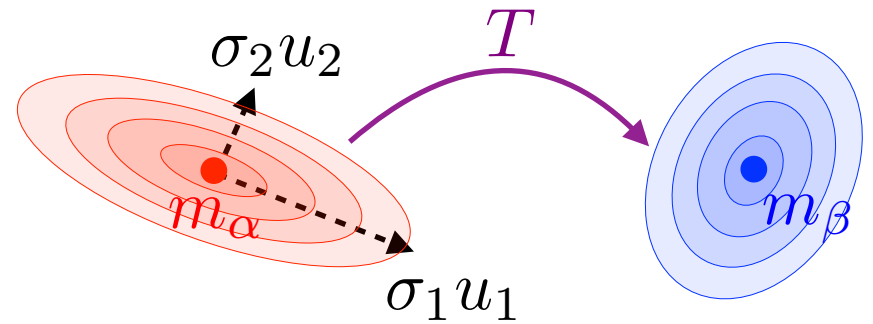
Uniqueness: one has  $\sqrt{\Sigma}\Sigma = \sqrt{\Sigma}^3 = \Sigma\sqrt{\Sigma}$ , they co-diagonalize.

# OT Between Gaussians: Bures Distance

Optimal transport:

$$T(x) = A(x - m_\alpha) + m_\beta$$

$$A = \Sigma_\alpha^{-\frac{1}{2}} \sqrt{\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}}} \Sigma_\alpha^{-\frac{1}{2}}$$



$$W_2^2(\alpha, \beta) = \|m_\alpha - m_\beta\|^2 + \mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2$$

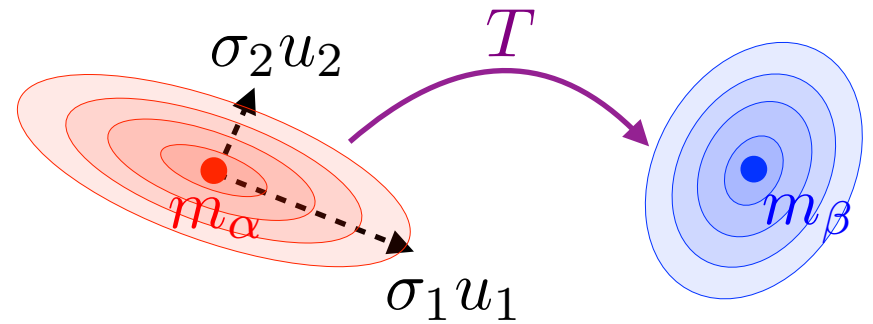
$$\text{Bures distance: } \mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 \stackrel{\text{def.}}{=} \text{tr} \left( \Sigma_\alpha + \Sigma_\beta - 2 \sqrt{\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}}} \right)$$

# OT Between Gaussians: Bures Distance

Optimal transport:

$$T(x) = A(x - m_\alpha) + m_\beta$$

$$A = \Sigma_\alpha^{-\frac{1}{2}} \sqrt{\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}}} \Sigma_\alpha^{-\frac{1}{2}}$$

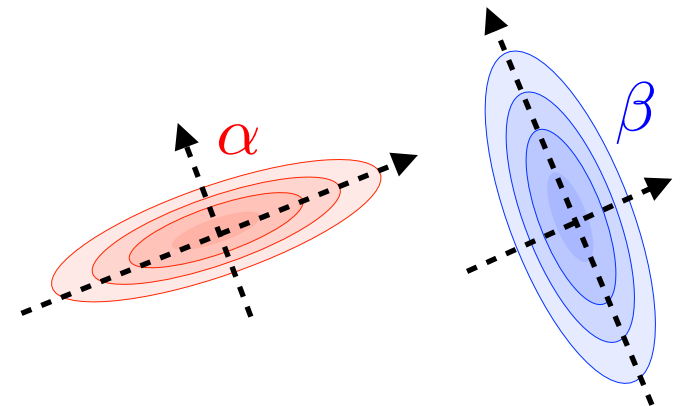


$$W_2^2(\alpha, \beta) = \|m_\alpha - m_\beta\|^2 + \mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2$$

Bures distance:  $\mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 \stackrel{\text{def.}}{=} \text{tr} \left( \Sigma_\alpha + \Sigma_\beta - 2\sqrt{\Sigma_\alpha^{\frac{1}{2}} \Sigma_\beta \Sigma_\alpha^{\frac{1}{2}}} \right)$

If  $\Sigma_\alpha \Sigma_\beta = \Sigma_\beta \Sigma_\alpha$ :

$$\mathcal{B}(\Sigma_\alpha, \Sigma_\beta)^2 = \|\sqrt{\Sigma_\alpha} - \sqrt{\Sigma_\beta}\|^2$$

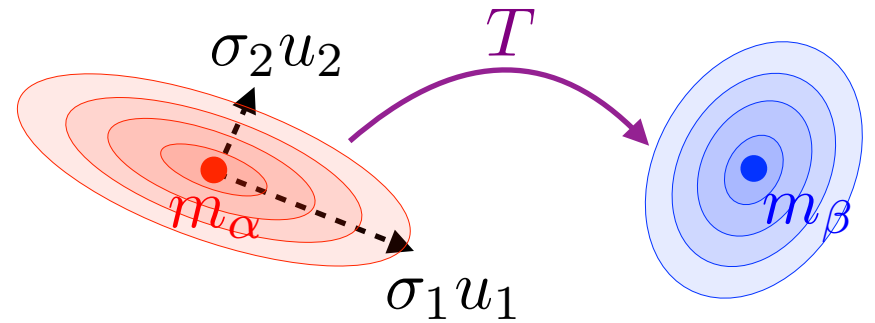


# Interpolation Between Gaussians

Optimal transport map  $T_{\#}\alpha = \beta$ .

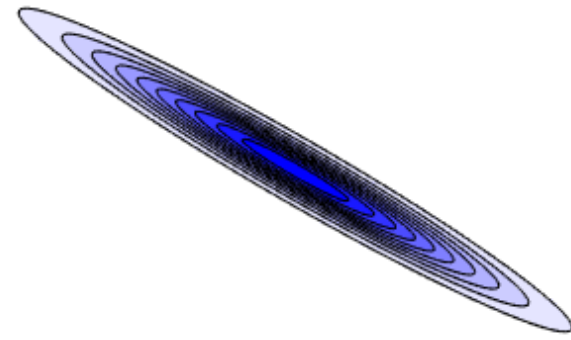
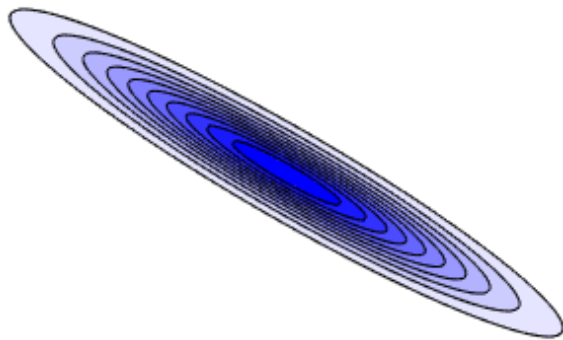
$$T(x) = A(x - m_{\alpha}) + m_{\beta}$$

$$A = \Sigma_{\alpha}^{-\frac{1}{2}} \left( \Sigma_{\alpha}^{\frac{1}{2}} \Sigma_{\beta} \Sigma_{\alpha}^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_{\alpha}^{-\frac{1}{2}}$$



Displacement interpolation:  $\alpha_t \stackrel{\text{def.}}{=} ((1 - t)\text{Id} + tT)_{\#}\alpha = \mathcal{N}(m_t, \Sigma_t)$

$$m_t = (1 - t)m_{\alpha} + tm_{\beta} \quad \Sigma_t = [(1 - t)\text{Id} + tA]\Sigma_{\alpha}[(1 - t)\text{Id} + tA]$$

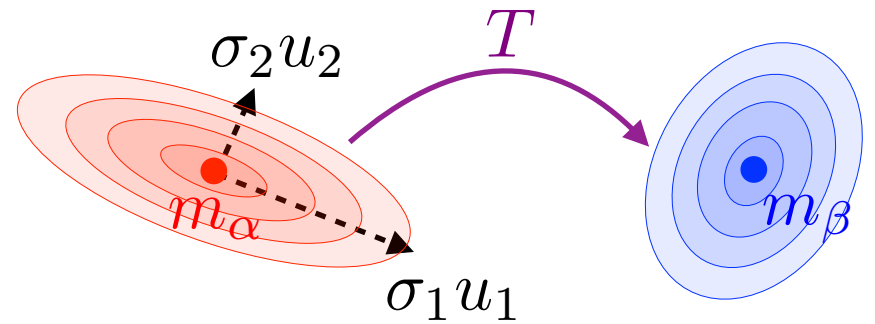


# Interpolation Between Gaussians

Optimal transport map  $T_{\#}\alpha = \beta$ .

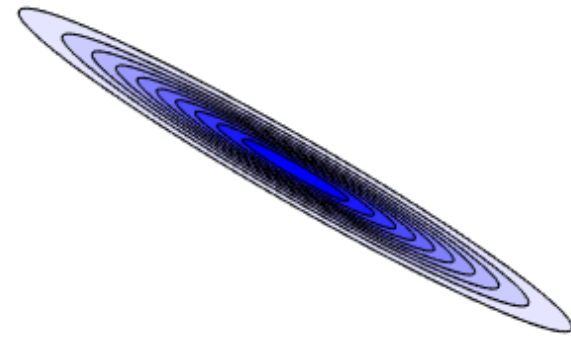
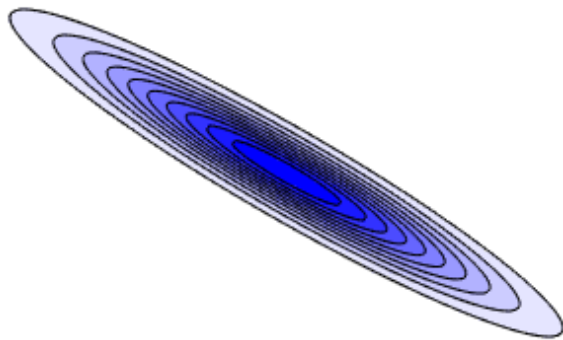
$$T(x) = A(x - m_{\alpha}) + m_{\beta}$$

$$A = \Sigma_{\alpha}^{-\frac{1}{2}} \left( \Sigma_{\alpha}^{\frac{1}{2}} \Sigma_{\beta} \Sigma_{\alpha}^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_{\alpha}^{-\frac{1}{2}}$$



Displacement interpolation:  $\alpha_t \stackrel{\text{def.}}{=} ((1 - t)\text{Id} + tT)_{\#}\alpha = \mathcal{N}(m_t, \Sigma_t)$

$$m_t = (1 - t)m_{\alpha} + tm_{\beta} \quad \Sigma_t = [(1 - t)\text{Id} + tA]\Sigma_{\alpha}[(1 - t)\text{Id} + tA]$$



# Overview

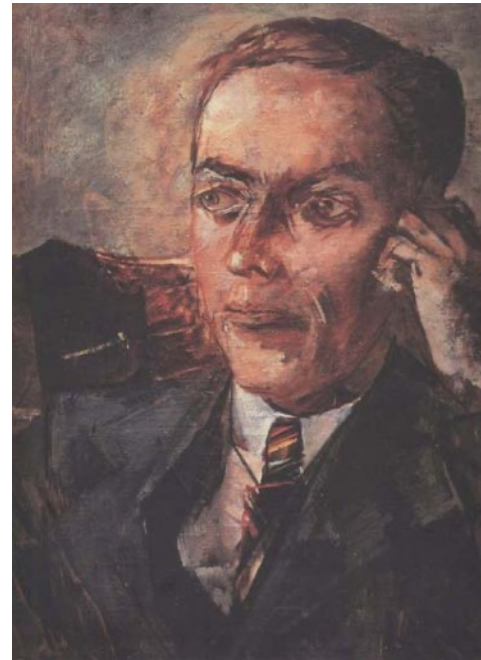
---

- Monge Formulation
- Continuous Optimal Transport
- **Kantorovitch Formulation**
- Applications



# Leonid Kantorovich (1912-1986)

Леонид Витальевич Канторович



[Kantorovich 1942]

*Journal of Mathematical Sciences, Vol. 133, No. 4, 8006*

**ON THE TRANSLLOCATION OF MASSES**

**L. V. Kantorovich\***

*The original paper was published in Dokl. Akad. Nauk SSSR, 37, No. 7-8, 227-229 (1942).*

We assume that  $R$  is a compact metric space, though some of the definitions and results given below can be formulated for more general spaces.

Let  $\Phi(e)$  be a mass distribution, i.e., a set function such that: (1) it is defined for Borel sets, (2) it is nonnegative:  $\Phi(e) \geq 0$ , (3) it is absolutely additive: if  $e = e_1 + e_2 + \dots$ ;  $e_i \cap e_k = 0$  ( $i \neq k$ ), then  $\Phi(e) = \Phi(e_1) + \Phi(e_2) + \dots$ . Let  $\Phi'(e')$  be another mass distribution such that  $\Phi(R) = \Phi'(R)$ . By definition, a translocation of masses is a function  $\Psi(e, e')$  defined for pairs of (B)-sets  $e, e' \in R$  such that: (1) it is nonnegative and absolutely additive with respect to each of its arguments, (2)  $\Psi(e, R) = \Phi(e)$ ,  $\Psi(R, e') = \Phi'(e')$ .

Let  $r(x, y)$  be a known continuous nonnegative function representing the work required to move a unit mass from  $x$  to  $y$ .

We define the work required for the translocation of two given mass distributions as

$$W(\Psi, \Phi, \Phi') = \int \int_{R \times R} r(x, y) \Psi(dy, dx) = \lim_{\lambda \rightarrow 0} \sum_{i,k} r(x_i, x'_k) \Psi(e_i, e'_k),$$

where  $e_i$  are disjoint and  $\sum_1^n e_i = R$ ,  $e'_k$  are disjoint and  $\sum_1^m e'_k = R$ ,  $x_i \in e_i$ ,  $x'_k \in e'_k$ , and  $\lambda$  is the largest of the numbers  $\text{diam } e_i$  ( $i = 1, 2, \dots, n$ ) and  $\text{diam } e'_k$  ( $k = 1, 2, \dots, m$ ).

Clearly, this integral does exist.

We call the quantity

$$W(\Phi, \Phi') = \inf_{\Psi} W(\Psi, \Phi, \Phi')$$

the minimal translocation work. Since the set of all functions  $\{\Psi\}$  is compact, there exists a function  $\Psi_0$  realizing this minimum, so that

$$W(\Phi, \Phi') = W(\Psi_0, \Phi, \Phi'),$$

# Kantorovitch's Formulation

Discrete distributions:

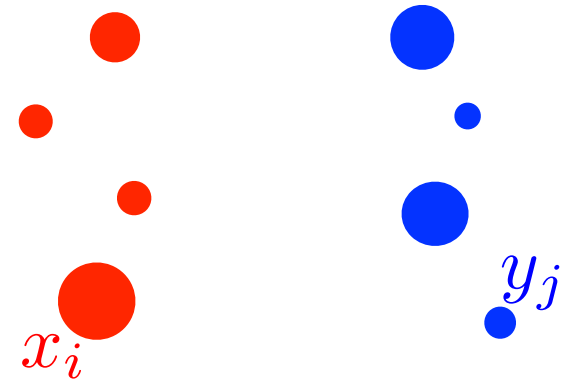
$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$$

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

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:

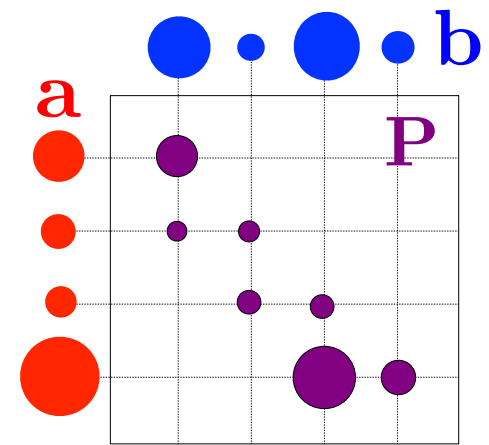
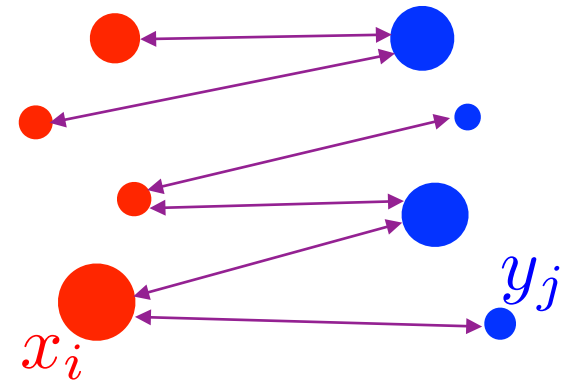
$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$$

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

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:

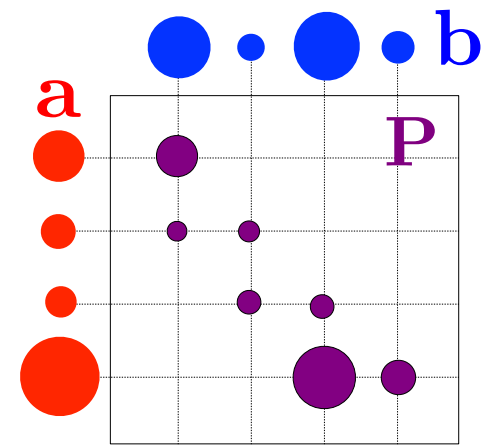
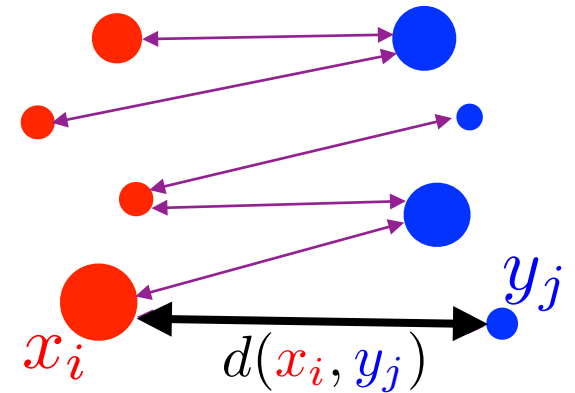
$$\alpha = \sum_{i=1}^n \mathbf{a}_i \delta_{x_i}$$

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

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

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

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



Couplings:

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

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

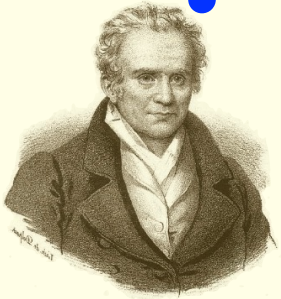
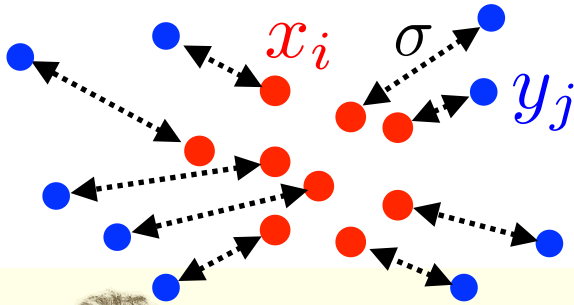
[Kantorovitch 1942]

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

# Kantorovitch's Exact Relaxation

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

$$\beta = \sum_{j=1}^n \delta_{y_j}$$

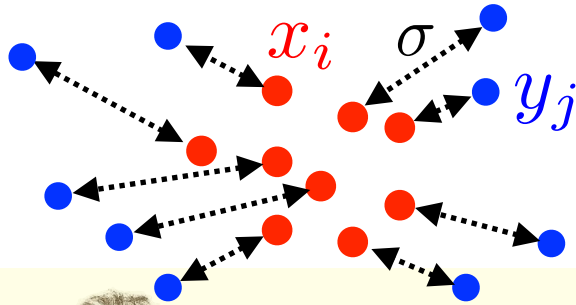


Monge (1784):

$$\min_{\sigma \in \text{Perm}_n} \sum_{i=1}^n C_{i, \sigma(i)}$$

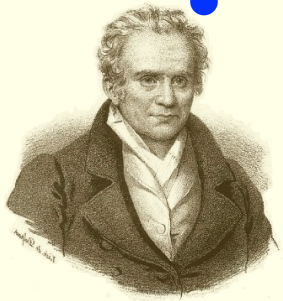
# Kantorovitch's Exact Relaxation

$$\alpha = \sum_{i=1}^n \delta_{x_i} \quad \beta = \sum_{j=1}^n \delta_{y_j}$$



Permutations “C” Bi-stochastic matrices:

$$\text{Bist}_n \stackrel{\text{def.}}{=} \{ \mathbf{P} \in \mathbb{R}_+^{n \times n} ; \mathbf{P}\mathbf{1} = \mathbf{1}, \mathbf{P}^\top \mathbf{1} = \mathbf{1} \}$$



Monge (1784):

$$\min_{\sigma \in \text{Perm}_n} \sum_{i=1}^n C_{i, \sigma(i)}$$

$\cong$  (relaxation)

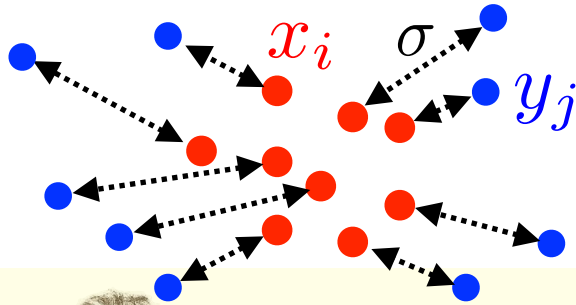


Kantorovitch (1942):

$$\min_{\mathbf{P} \in \text{Bist}_n} \sum_{i=1}^n \sum_{j=1}^n P_{i,j} C_{i,j}$$

# Kantorovitch's Exact Relaxation

$$\alpha = \sum_{i=1}^n \delta_{x_i} \quad \beta = \sum_{j=1}^n \delta_{y_j}$$



Permutations “C” Bi-stochastic matrices:

$$\text{Bist}_n \stackrel{\text{def.}}{=} \{ \mathbf{P} \in \mathbb{R}_+^{n \times n} ; \mathbf{P}\mathbf{1} = \mathbf{1}, \mathbf{P}^\top \mathbf{1} = \mathbf{1} \}$$



Monge (1784):

$$\min_{\sigma \in \text{Perm}_n} \sum_{i=1}^n C_{i, \sigma(i)}$$

$\cong$  (relaxation)



Kantorovitch (1942):

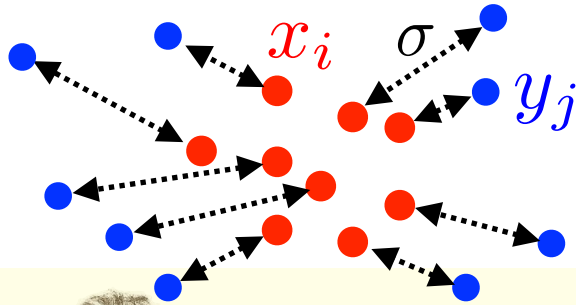
$$\min_{\mathbf{P} \in \text{Bist}_n} \sum_{i=1}^n \sum_{j=1}^n P_{i,j} C_{i,j}$$

$n!$  permutations

$O(n^3)$  algorithm

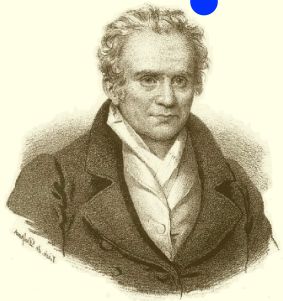
# Kantorovitch's Exact Relaxation

$$\alpha = \sum_{i=1}^n \delta_{x_i} \quad \beta = \sum_{j=1}^n \delta_{y_j}$$



Permutations “C” Bi-stochastic matrices:

$$\text{Bist}_n \stackrel{\text{def.}}{=} \{ \mathbf{P} \in \mathbb{R}_+^{n \times n} ; \mathbf{P}\mathbf{1} = \mathbf{1}, \mathbf{P}^\top \mathbf{1} = \mathbf{1} \}$$



Monge (1784):

$$\min_{\sigma \in \text{Perm}_n} \sum_{i=1}^n C_{i, \sigma(i)}$$

$\cong$  (relaxation)



Kantorovitch (1942):

$$\min_{\mathbf{P} \in \text{Bist}_n} \sum_{i=1}^n \sum_{j=1}^n P_{i,j} C_{i,j}$$

$n!$  permutations

$O(n^3)$  algorithm



George David Birkhoff

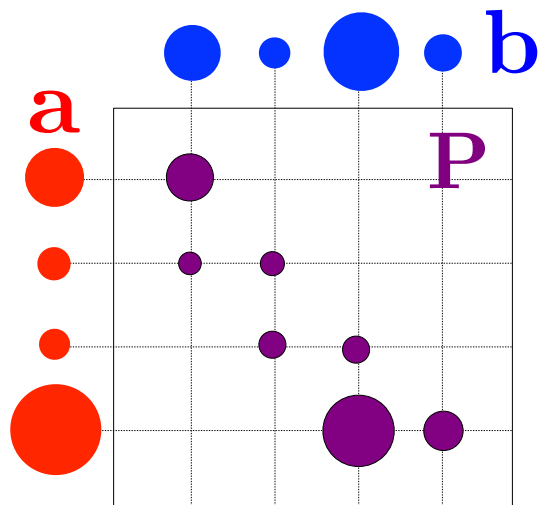
John von Neumann

*Theorem:* [Birkhoff-von Neumann]

“Monge  $\Leftrightarrow$  Kantorovitch”

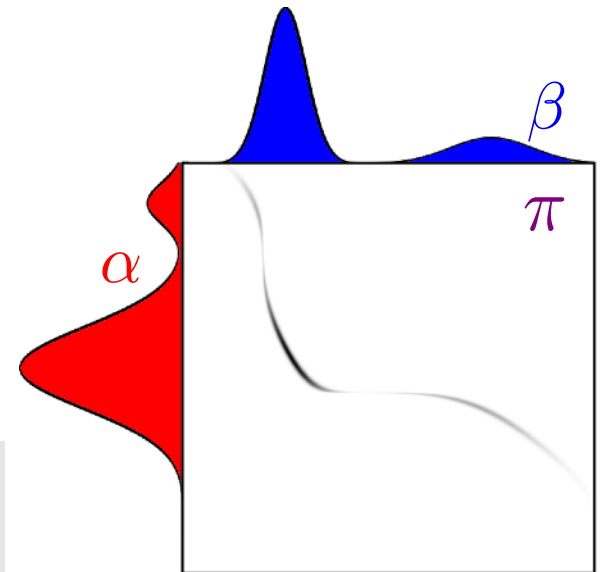


# General Formulation



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

$$c(x, y) = d(x, y)^p$$

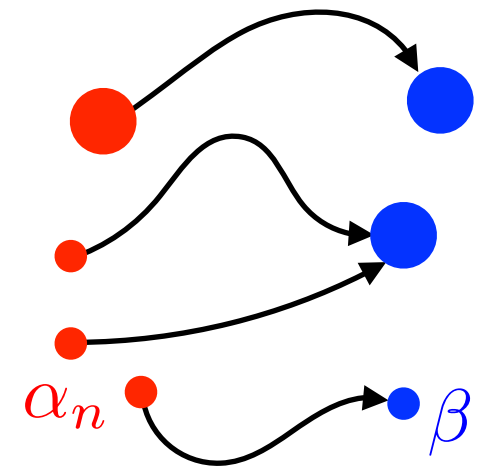
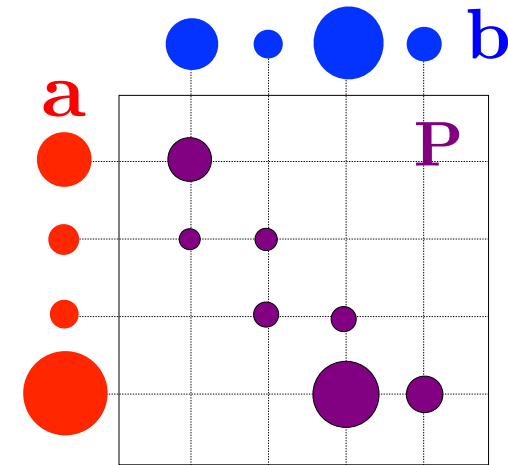


$$W_p(\alpha, \beta)^p \stackrel{\text{def.}}{=} \min_{\pi \in \mathcal{M}_+^1(\mathcal{X}^2)} \left\{ \int_{\mathcal{X}^2} d(x, y)^p d\pi(x, y) ; \pi_1 = \alpha, \pi_2 = \beta \right\}$$

# Optimal Transport Distances

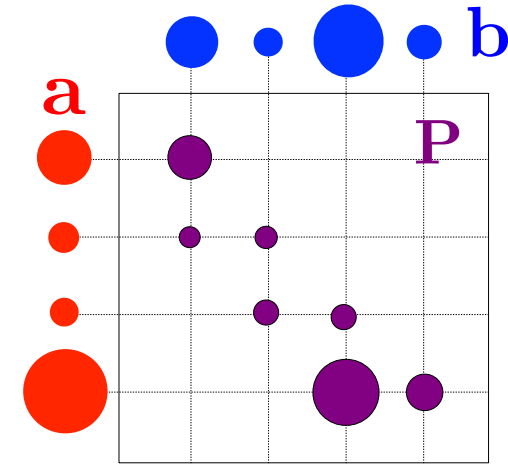
$$W_p(\alpha, \beta) \stackrel{\text{def.}}{=} \left( \min_{\mathbf{P} \mathbf{1}=\mathbf{a}, \mathbf{P}^\top \mathbf{1}=\mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} \right)^{\frac{1}{p}}$$

Convergence in law:  $\alpha_n \rightarrow^* \beta$   
 $\Leftrightarrow \forall f \in \mathcal{C}_c(\mathcal{X}), \int_{\mathcal{X}} f d\alpha_n \rightarrow \int_{\mathcal{X}} f d\beta$

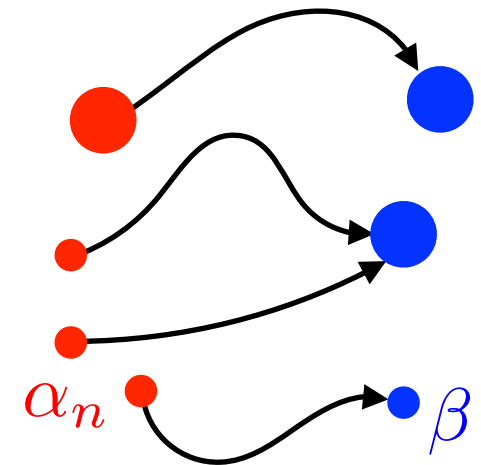


# Optimal Transport Distances

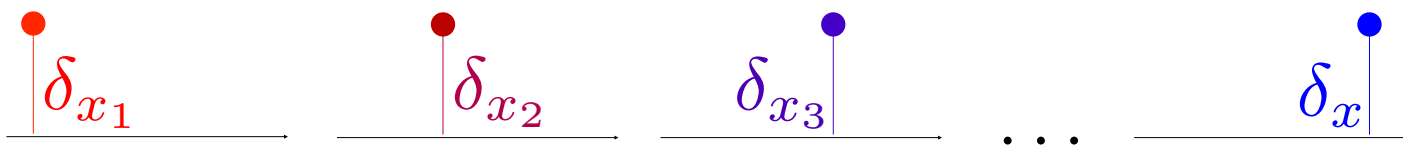
$$W_p(\alpha, \beta) \stackrel{\text{def.}}{=} \left( \min_{\mathbf{P} \mathbf{1}=\mathbf{a}, \mathbf{P}^\top \mathbf{1}=\mathbf{b}} \sum_{i,j} d(x_i, y_j)^p \mathbf{P}_{i,j} \right)^{\frac{1}{p}}$$



Convergence in law:  $\alpha_n \rightharpoonup^* \beta$   
 $\Leftrightarrow \forall f \in \mathcal{C}_c(\mathcal{X}), \int_{\mathcal{X}} f d\alpha_n \rightarrow \int_{\mathcal{X}} f d\beta$



*Theorem:*  $W_p$  is a distance and  
 $\alpha_n \rightharpoonup^* \beta \Leftrightarrow \int \|\cdot\|^p d\alpha_n \rightarrow \int \|\cdot\|^p d\beta \Leftrightarrow W_p(\alpha_n, \beta) \rightarrow 0$



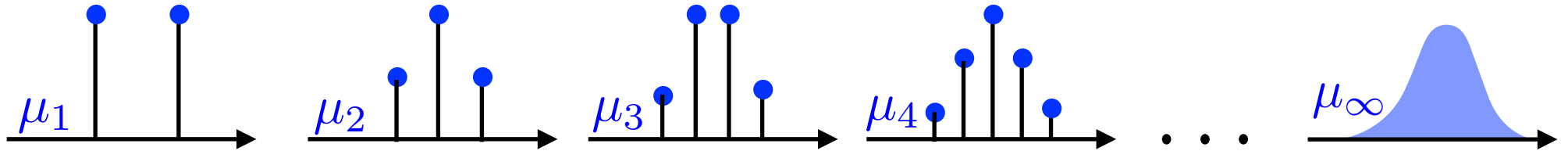
$$\|\delta_{x_n} - \delta_x\|_{\text{TV}} = 2 \quad \text{vs.} \quad W_p(\delta_{x_n}, \delta_x) = d(x_n, x)$$

# Application: Central Limit Theorem

Central limit theorem: If  $\mathbb{E}(X) = 0$ ,  $\mathbb{E}(X^2) = 1$  and  $(X_i)_i \stackrel{\text{i.i.d.}}{\sim} X$

$$Y_n \stackrel{\text{def.}}{=} \frac{X_1 + \dots + X_n}{\sqrt{n}} \xrightarrow{\text{law}} \mathcal{N}(0, 1)$$

$$\mu_n = \text{Law}(Y_n) = (\mu_1 \star \dots \star \mu_1)(\cdot/\sqrt{n}) \xrightarrow{*} \mu_\infty = \text{Law}(\mathcal{N}(0, 1))$$

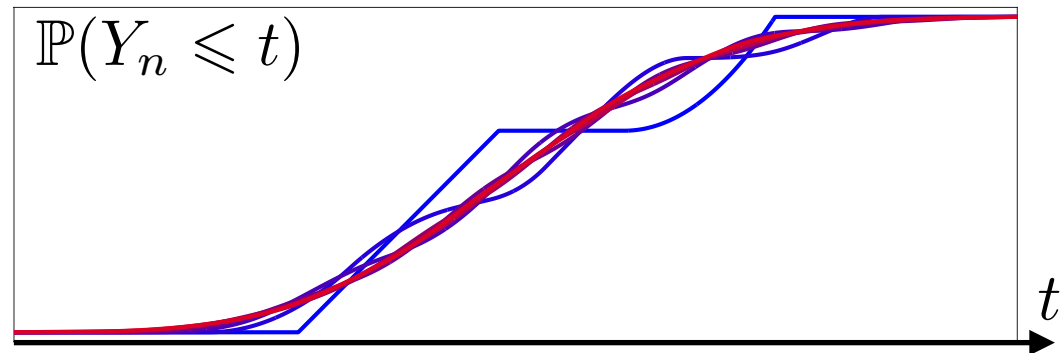
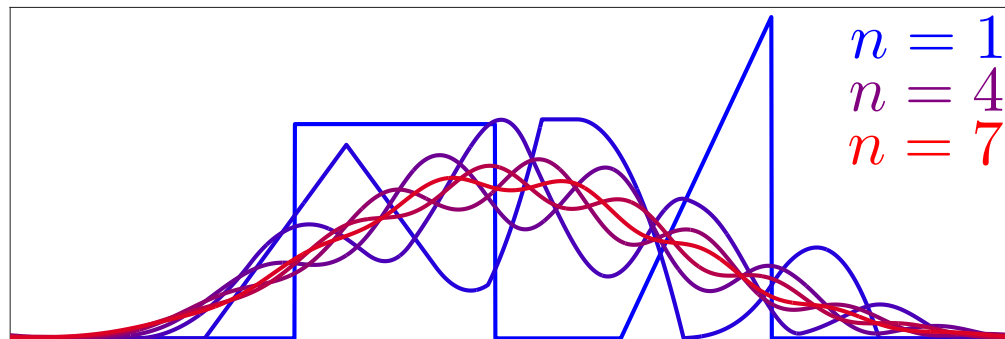
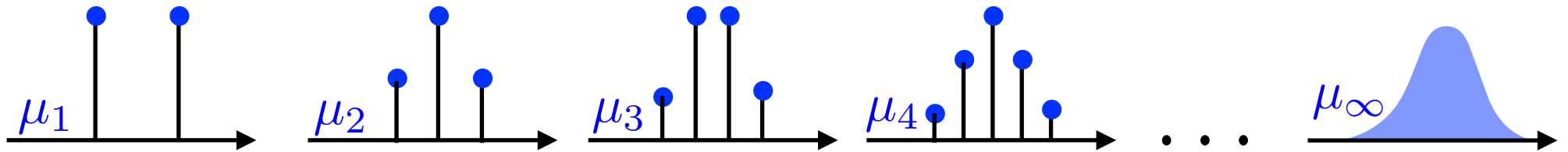


# Application: Central Limit Theorem

Central limit theorem: If  $\mathbb{E}(X) = 0$ ,  $\mathbb{E}(X^2) = 1$  and  $(X_i)_i \stackrel{\text{i.i.d.}}{\sim} X$

$$Y_n \stackrel{\text{def.}}{=} \frac{X_1 + \dots + X_n}{\sqrt{n}} \xrightarrow{\text{law}} \mathcal{N}(0, 1)$$

$$\mu_n = \text{Law}(Y_n) = (\mu_1 \star \dots \star \mu_1)(\cdot/\sqrt{n}) \xrightarrow{*} \mu_\infty = \text{Law}(\mathcal{N}(0, 1))$$

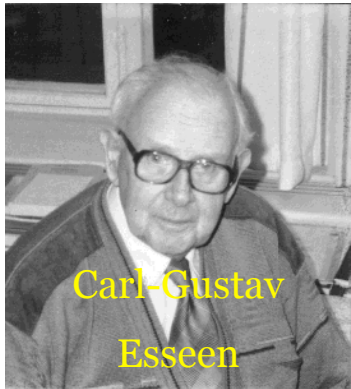
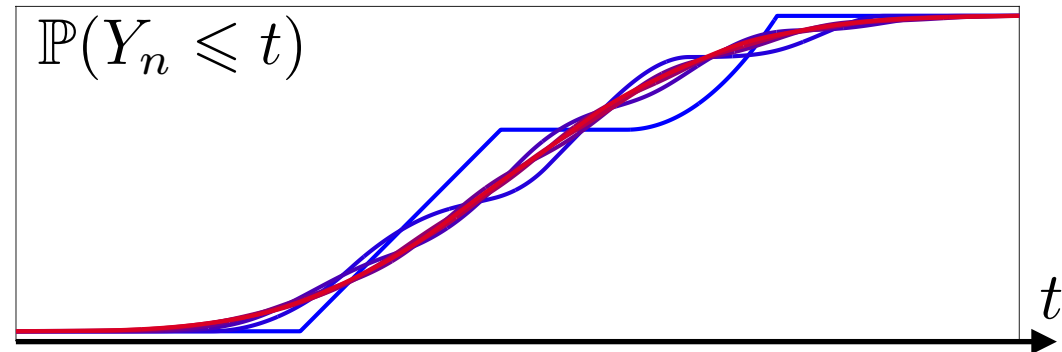
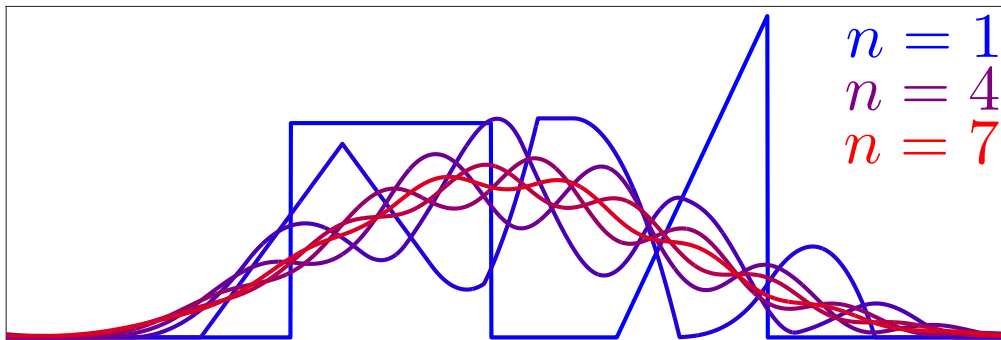
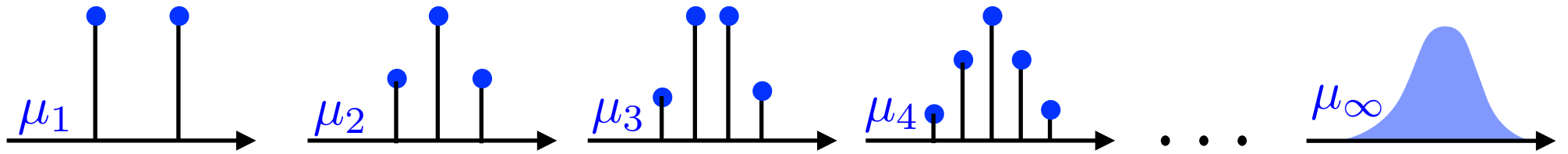


# Application: Central Limit Theorem

Central limit theorem: If  $\mathbb{E}(X) = 0, \mathbb{E}(X^2) = 1$  and  $(X_i)_i \stackrel{\text{i.i.d.}}{\sim} X$

$$Y_n \stackrel{\text{def.}}{=} \frac{X_1 + \dots + X_n}{\sqrt{n}} \xrightarrow{\text{law}} \mathcal{N}(0, 1)$$

$$\mu_n = \text{Law}(Y_n) = (\mu_1 \star \dots \star \mu_1)(\cdot/\sqrt{n}) \xrightarrow{*} \mu_\infty = \text{Law}(\mathcal{N}(0, 1))$$



Carl-Gustav  
Esseen

*Theorem:*

[Berry 1941]

[Esseen, 1942]

$$W_1(\mu_n, \mathcal{N}(0, 1)) \leq \frac{C\mathbb{E}(|X|^3)}{\sqrt{n}} \quad C \leq 1/2$$

→ Generalizes to higher dimensions.

# Overview

---

- Monge Formulation
- Continuous Optimal Transport
- Kantorovitch Formulation
- **Applications**

# Wasserstein Barycenters

Barycenters of measures  $(\alpha_s)_{s=1}^S$ :  $\sum_s \lambda_s = 1$

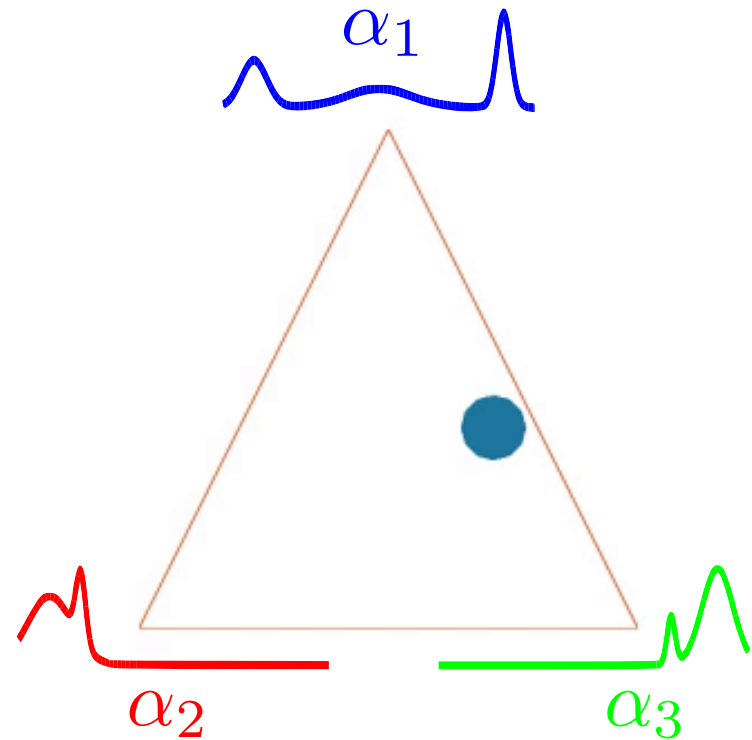
$$\alpha^* \in \operatorname{argmin}_{\alpha} \sum_s \lambda_s W_p^p(\alpha, \alpha_s)$$



Guillaume Carlier



Martial Agueh



$$\lambda \in \Sigma_3$$

$$\min_{\alpha} \sum_s \lambda_s W_p^p(\alpha, \alpha_s)$$

Wasserstein



$$\sum_s \lambda_s \alpha_s$$

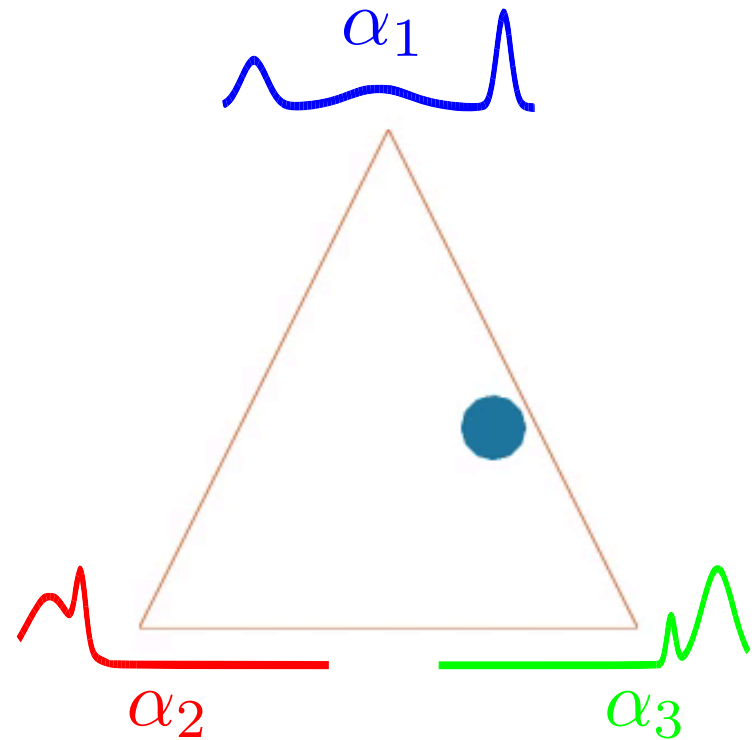
Euclidean



# Wasserstein Barycenters

Barycenters of measures  $(\alpha_s)_{s=1}^S$ :  $\sum_s \lambda_s = 1$

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



$$\lambda \in \Sigma_3$$

$$\min_{\alpha} \sum_s \lambda_s W_p^p(\alpha, \alpha_s)$$

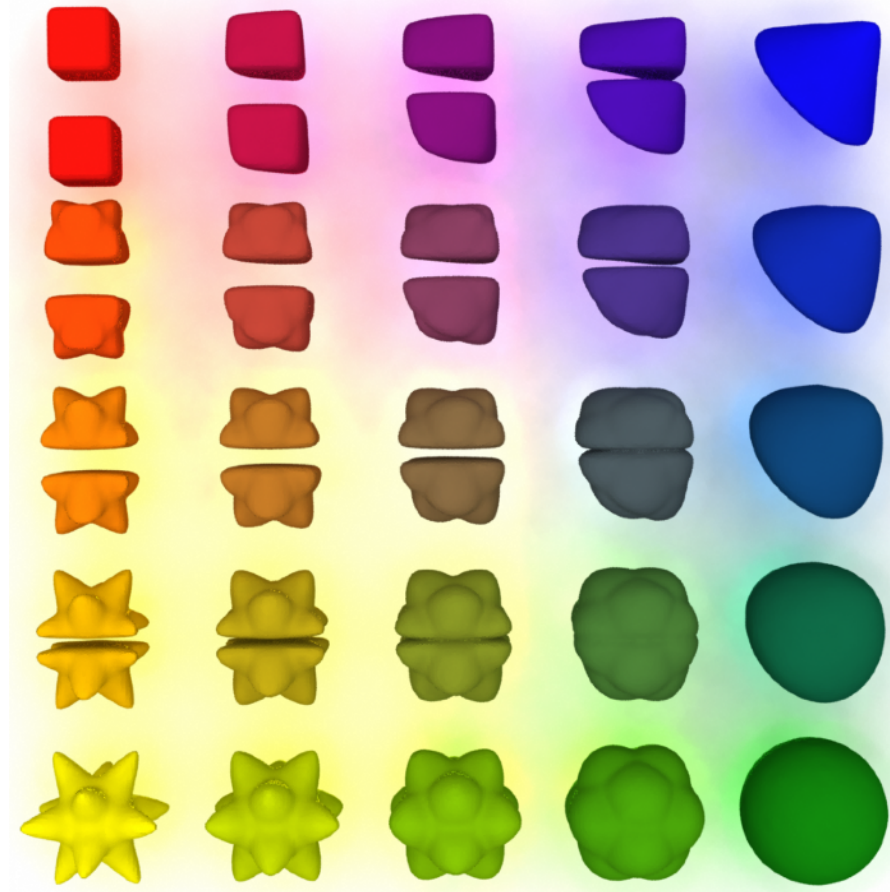
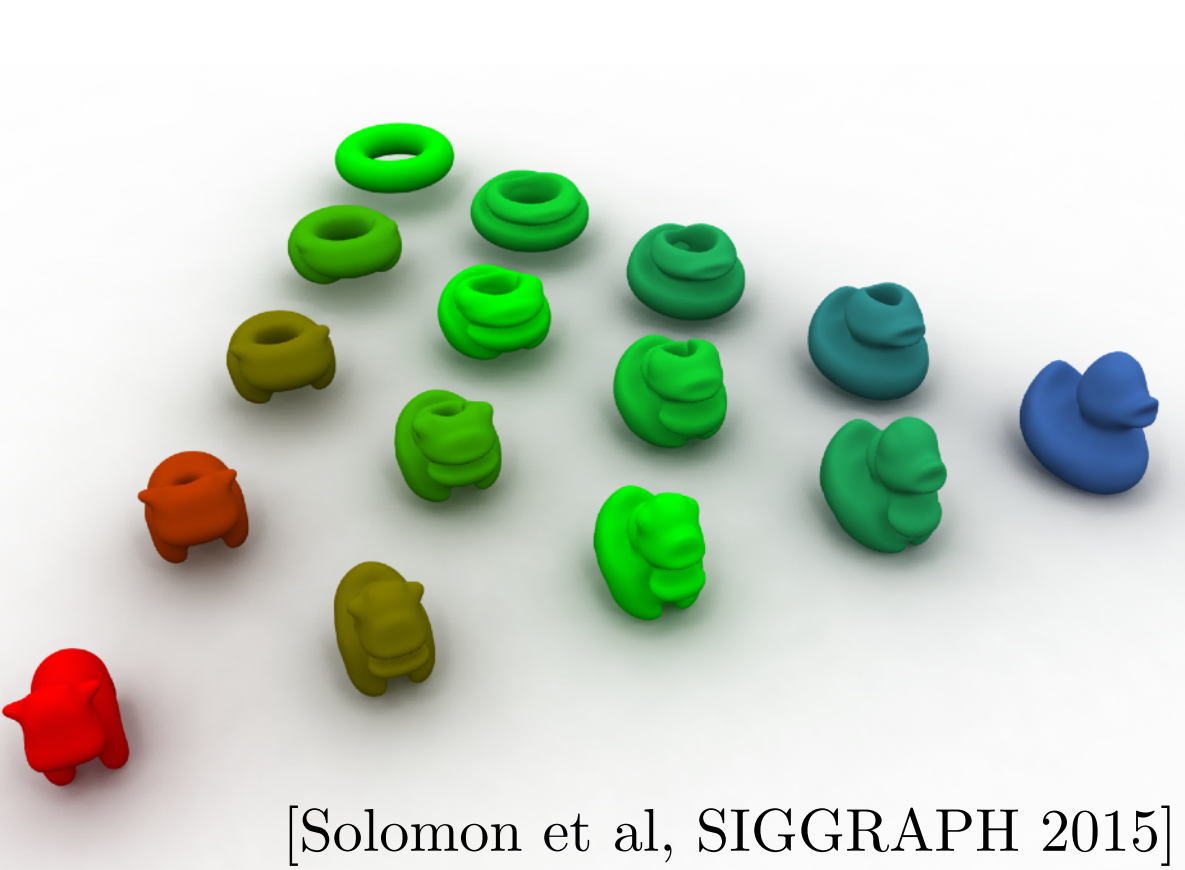
Wasserstein



$$\sum_s \lambda_s \alpha_s$$

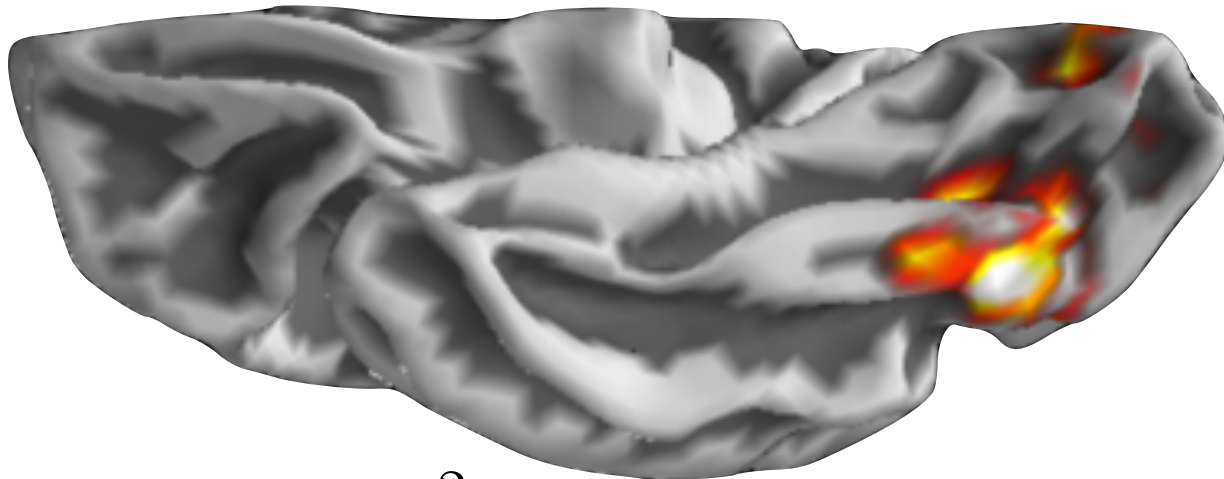
Euclidean

# Examples

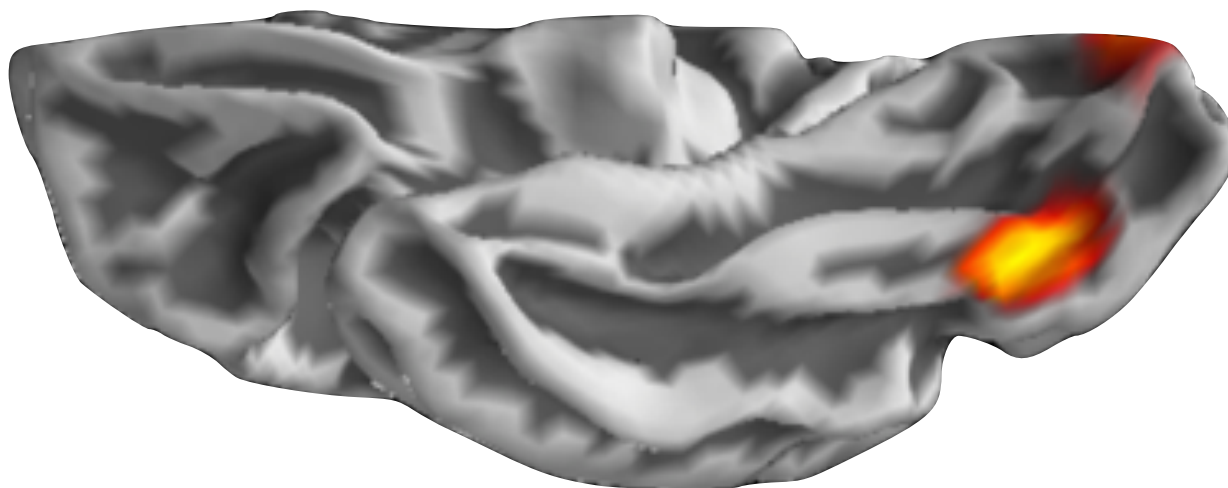


# MRI Data Processing [with A. Gramfort]

Ground cost  $c = d_M$ : geodesic on cortical surface  $M$ .



$L^2$  barycenter

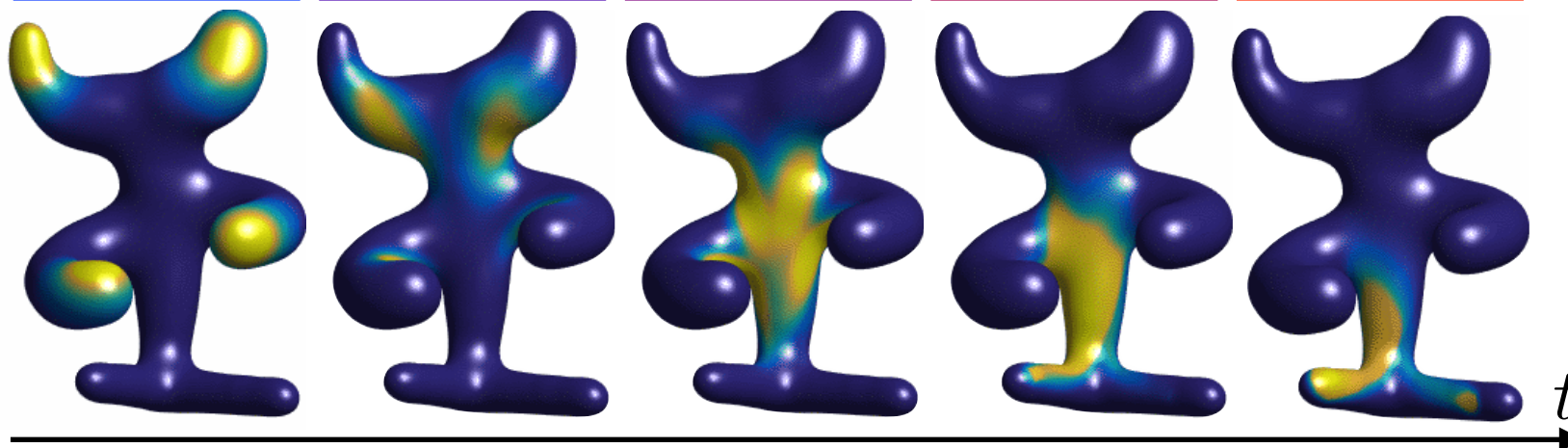
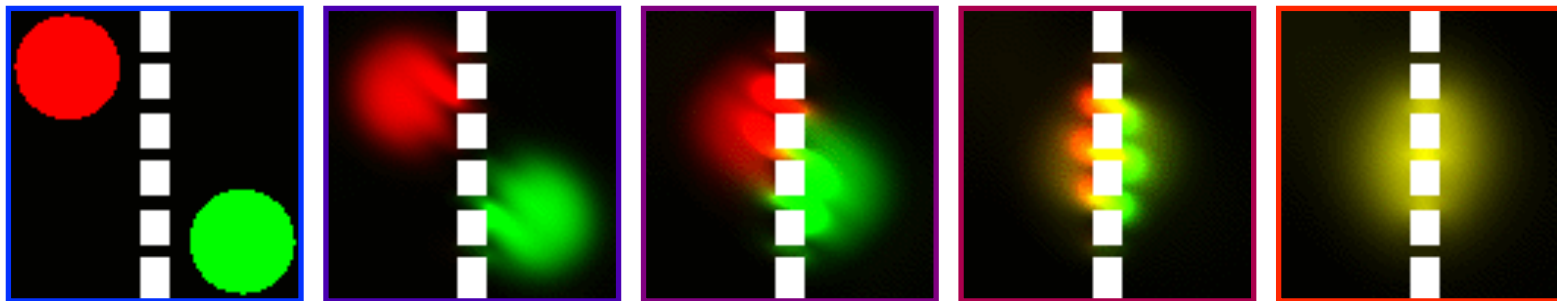


$W_2^2$  barycenter

# Generalizations: Gradient Flows

*Implicit stepping:*  $\alpha_{t+\tau} = \operatorname{argmin}_{\alpha} W_p^p(\alpha_t, \alpha) + \tau f(\alpha)$

Limit  $\tau \rightarrow 0$ :  $\frac{\partial \alpha}{\partial t} = \operatorname{div}(\alpha \nabla(f'(\alpha)))$



*Also:* mean field analysis of 1-hidden layer neural networks.

# Gradient Flows Simulation



<https://www.youtube.com/watch?v=tDQw21ntR64>

Tim Whittaker (New Zealand)



# Gradient Flows Simulation

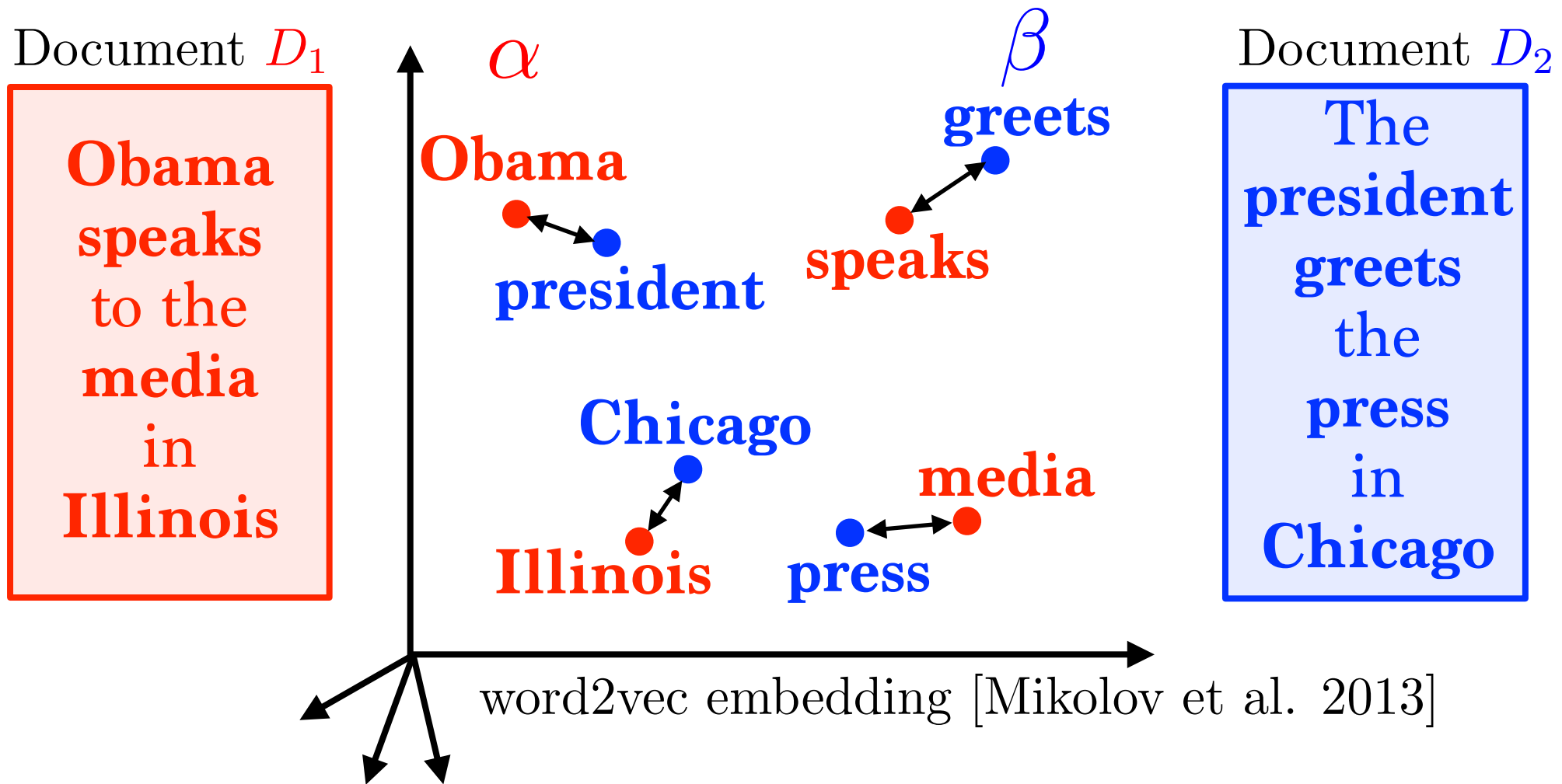


<https://www.youtube.com/watch?v=tDQw21ntR64>

Tim Whittaker (New Zealand)



# Bag of Words



Word mover's distance: [Kusner et al 2015]

$$\text{Dist}(D_1, D_2) = W_2(\alpha, \beta)$$

# Overview

---

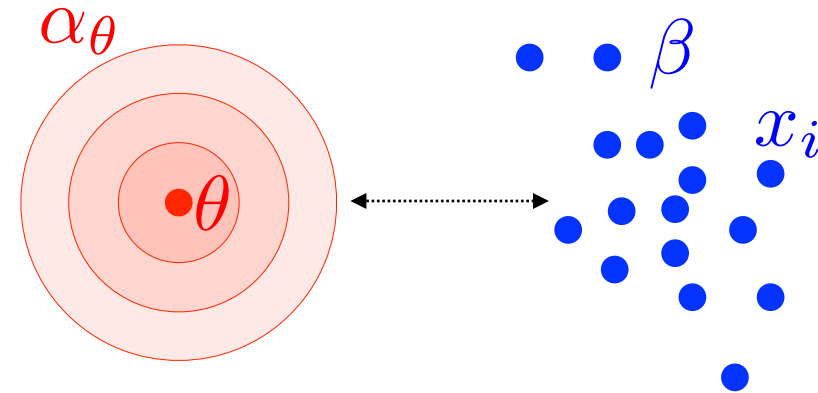
- Monge Formulation
- Continuous Optimal Transport
- Kantorovitch Formulation
- Applications
- **Generative models**



# Density Fitting and Generative Models

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

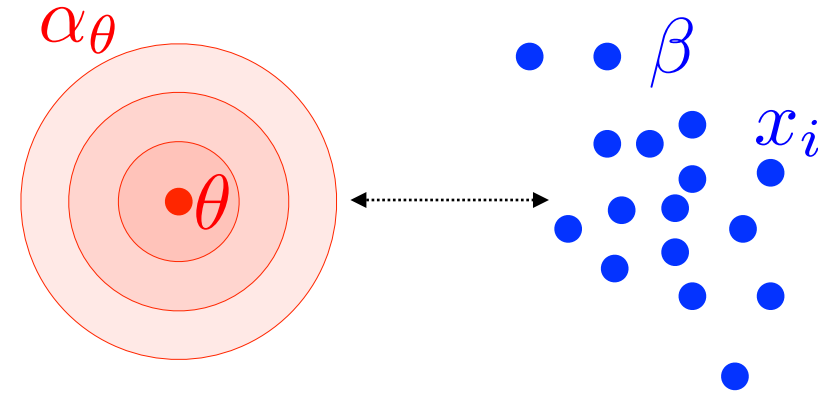
Parametric model:  $\theta \mapsto \alpha_\theta$



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

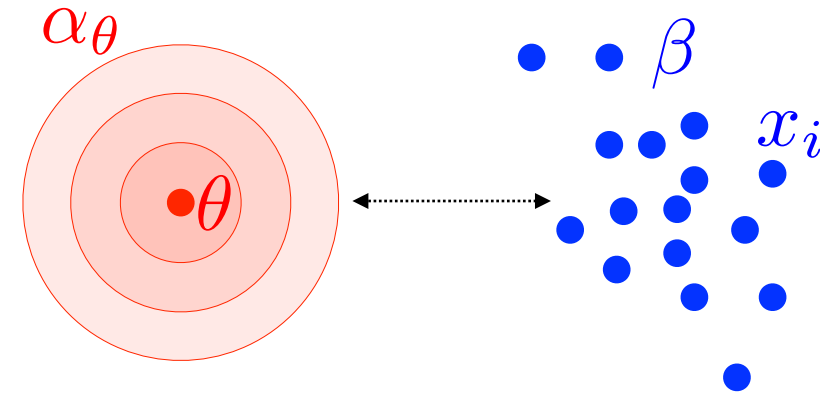
$$\min_{\theta} - \sum_i \log(\rho_\theta(x_i)) \xrightarrow{n \rightarrow +\infty} \text{KL}(\beta | \alpha_\theta)$$

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)

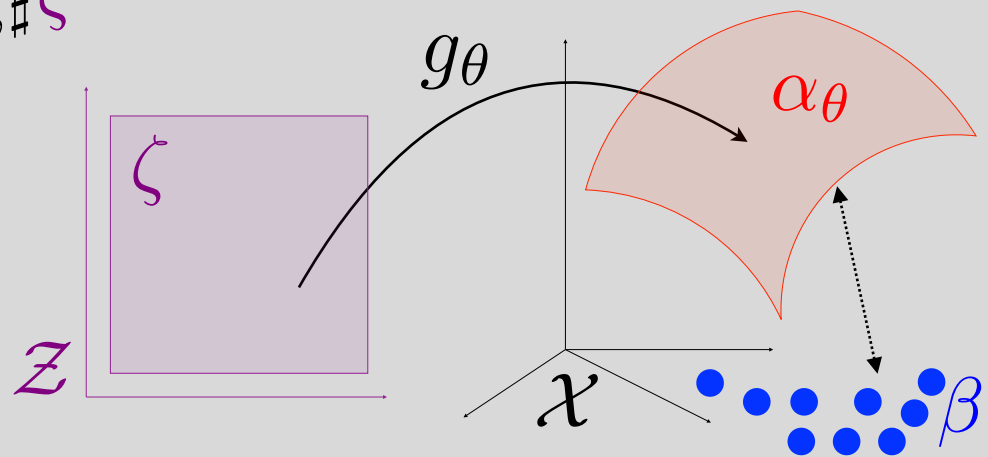
Generative model fit:  $\alpha_\theta = g_{\theta, \#} \zeta$

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

→ MLE undefined.

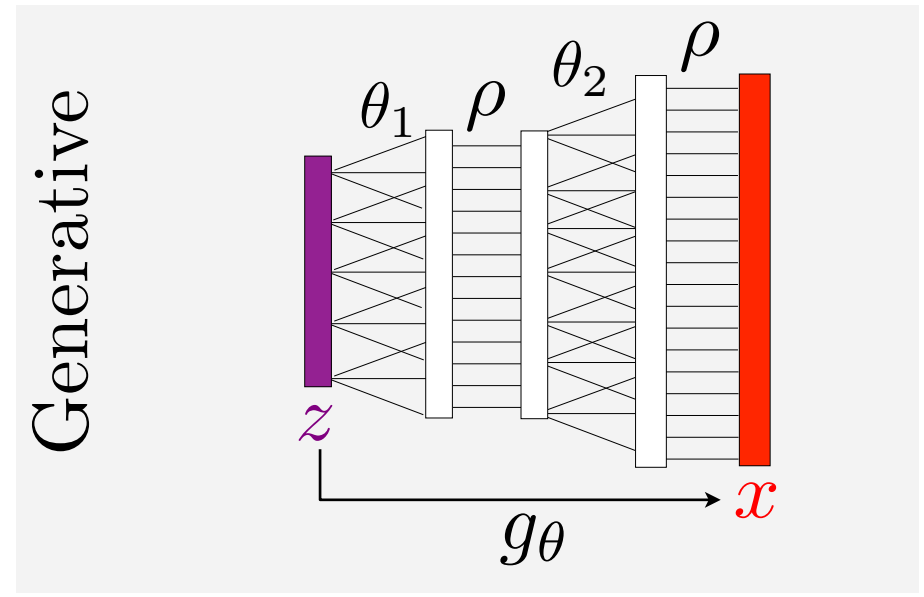
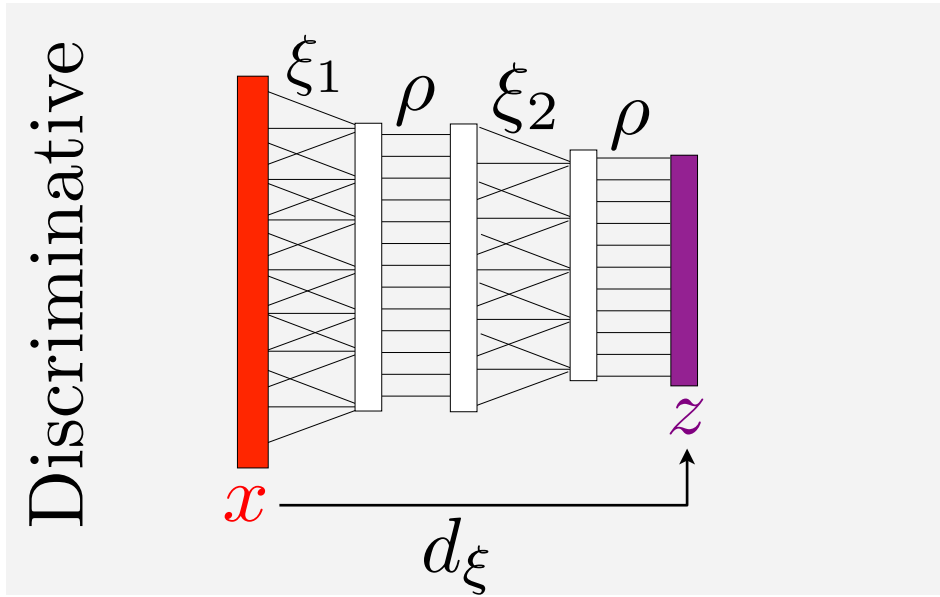
→ Need a weaker metric.

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



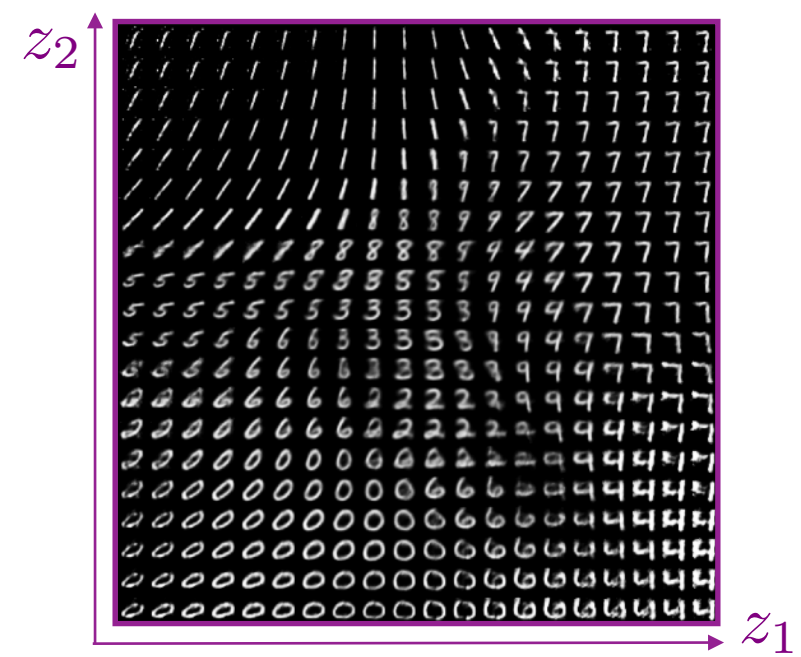
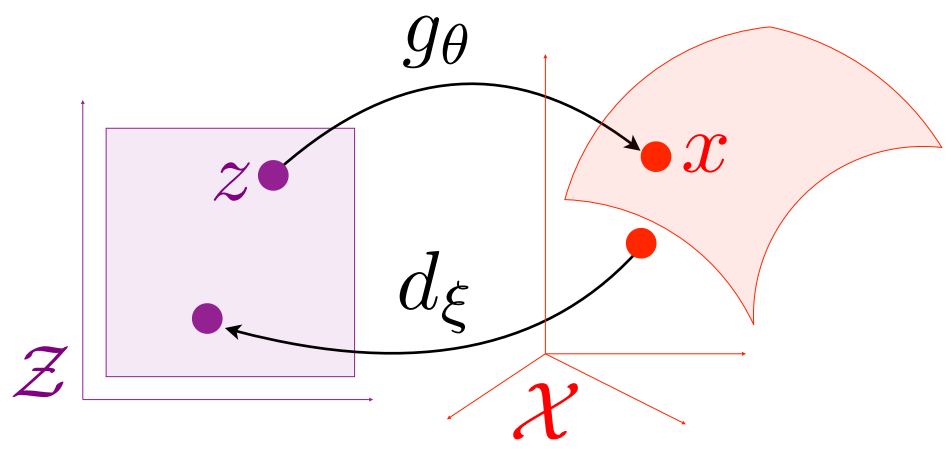
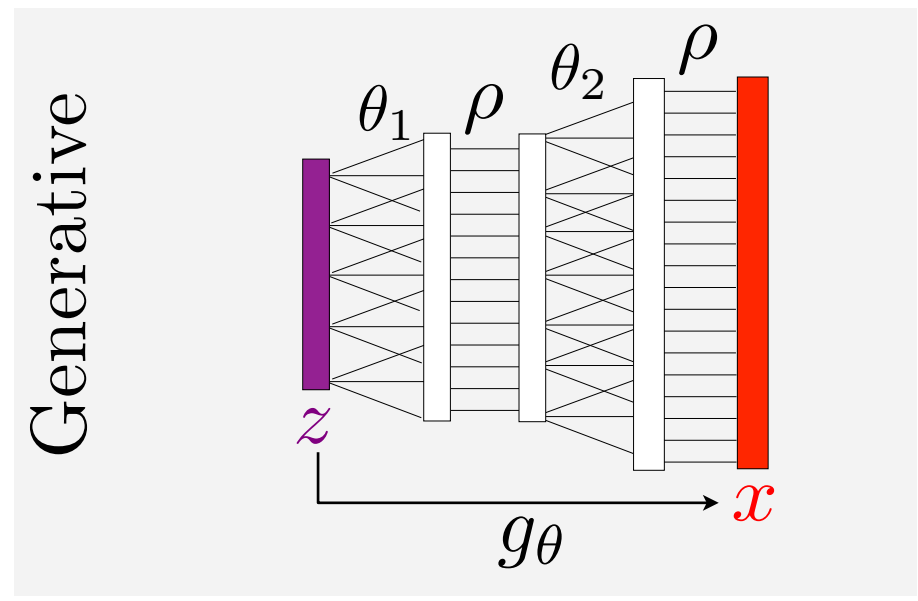
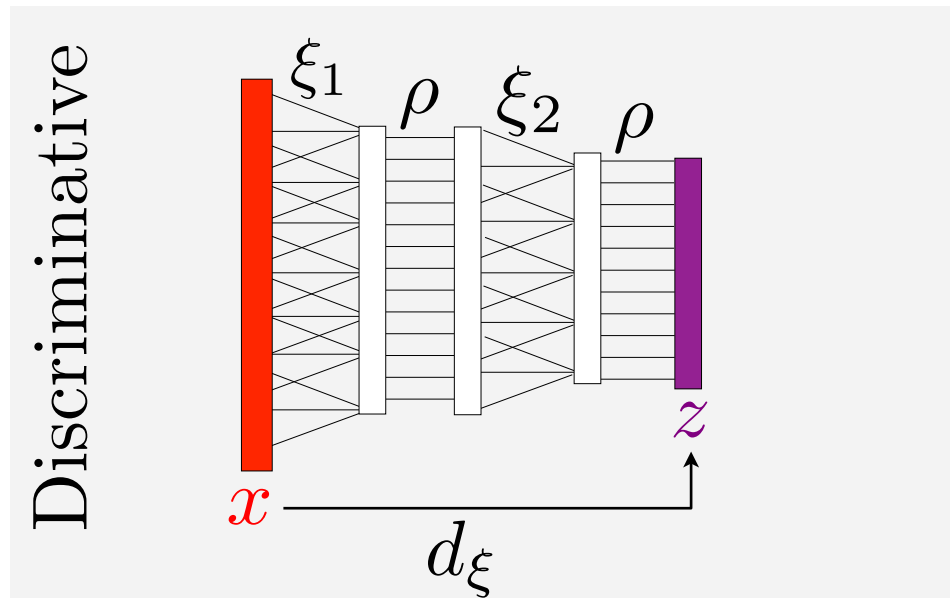
# Deep Discriminative vs Generative Models

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

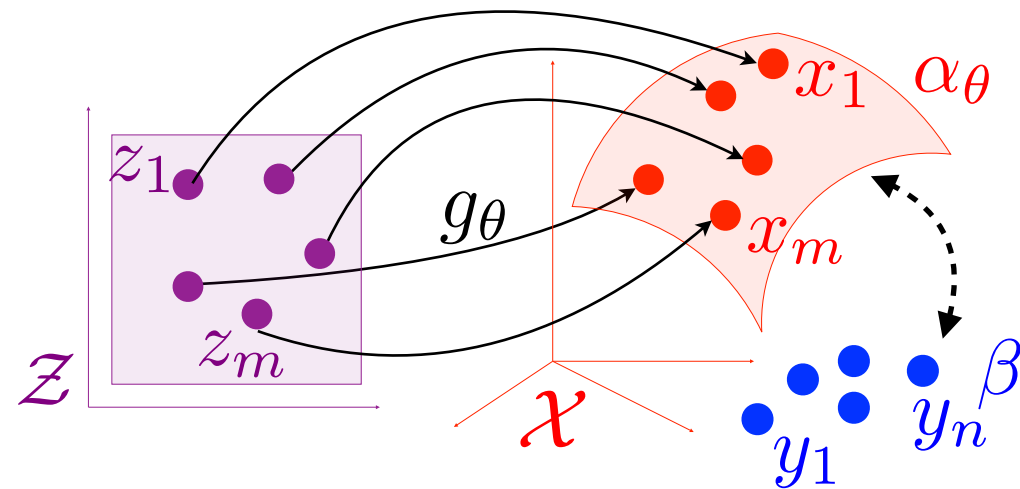


# Deep Discriminative vs Generative Models

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



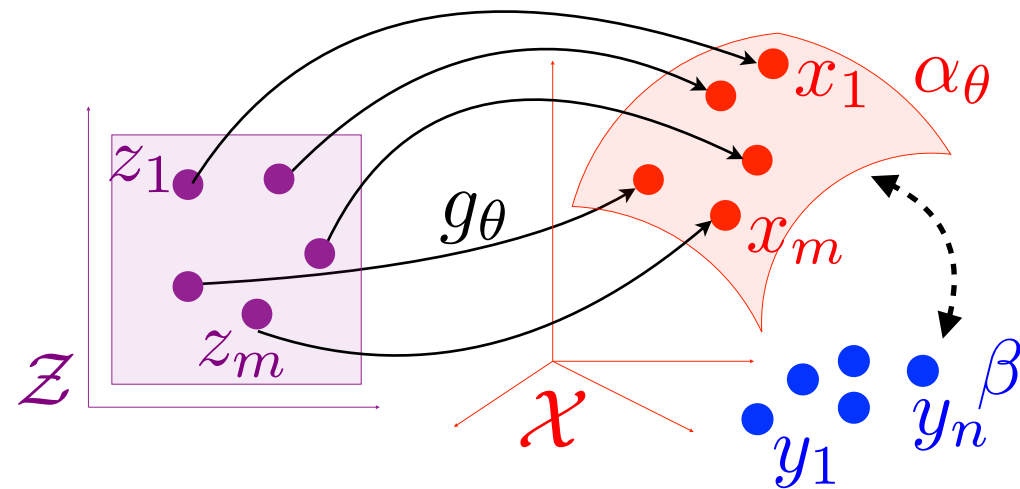
$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$

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

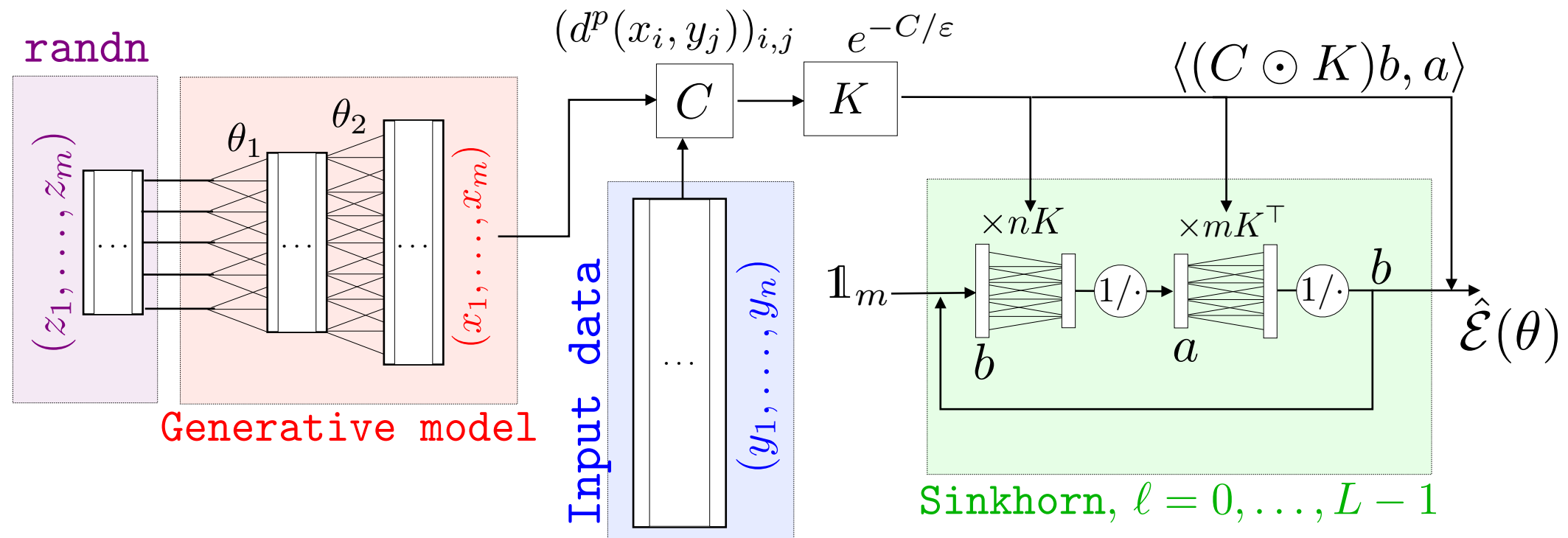


$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} \overline{W}_{\varepsilon, p}^p(\alpha_\theta, \beta)$$

Stochastic gradient descent

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

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

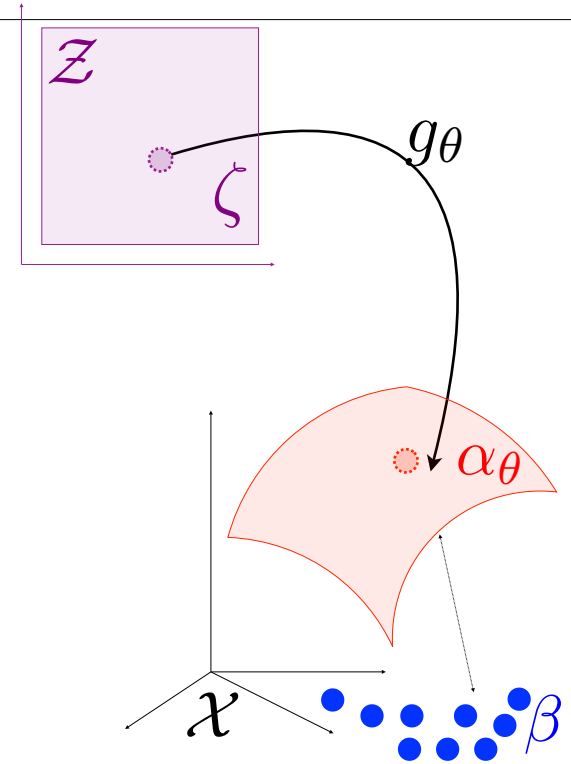


# Examples of Images Generation

Inputs  $\beta$



Generated  $\alpha_\theta$



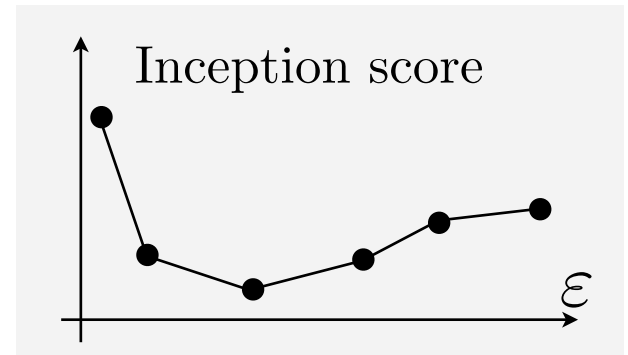
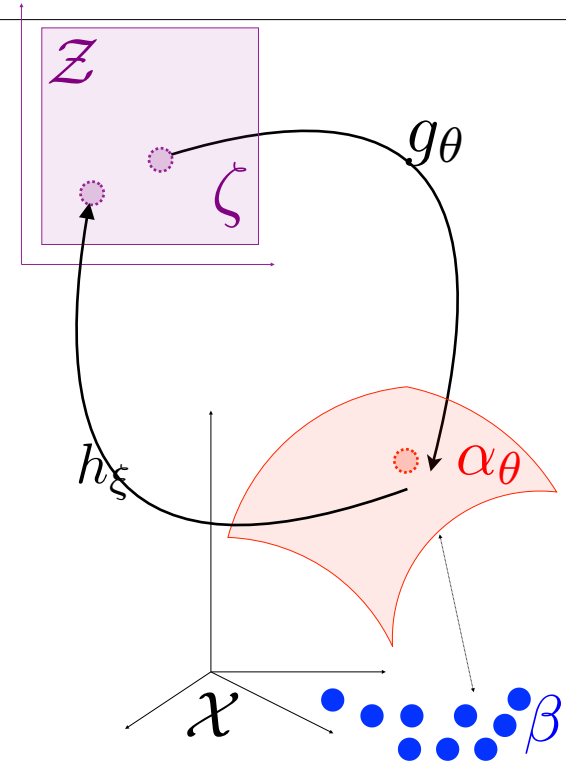


# Examples of Images Generation

Inputs  $\beta$



Generated  $\alpha_\theta$



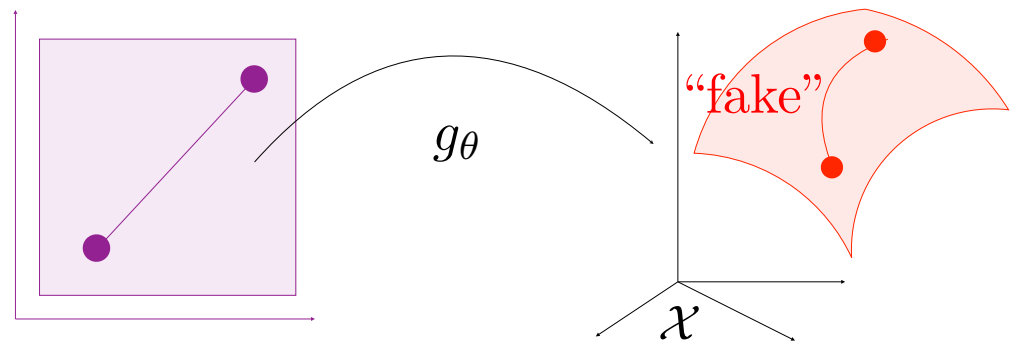
- Need to learn the metric  $d(x, y) = \|h_\xi(x) - h_\xi(y)\|$  (GANs)
- Influence of  $\epsilon$ ?
- Performance evaluation of generative models is an open problem.



Ian Goodfellow

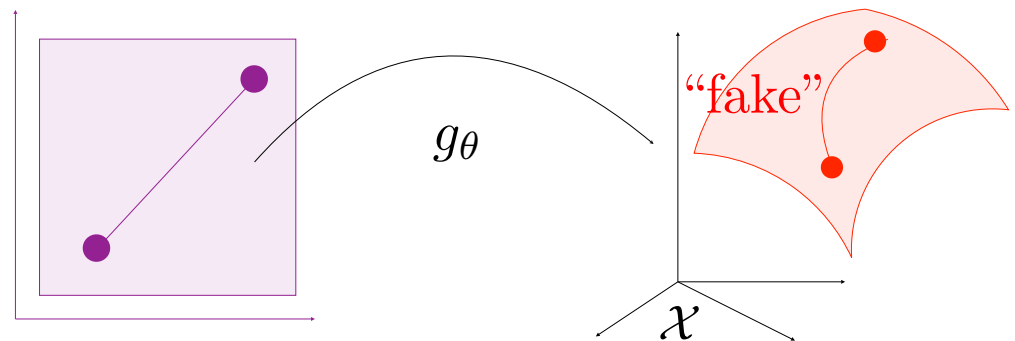


*Progressive Growing of GANs for Improved Quality, Stability, and Variation*  
Tero Karras, Timo Aila, Samuli Laine,  
Jaakko Lehtinen, ICLR 2018



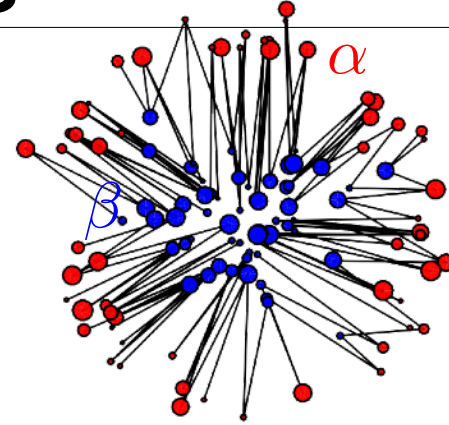


*Progressive Growing of GANs for Improved  
 Quality, Stability, and Variation*  
 Tero Karras, Timo Aila, Samuli Laine,  
 Jaakko Lehtinen, ICLR 2018



# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

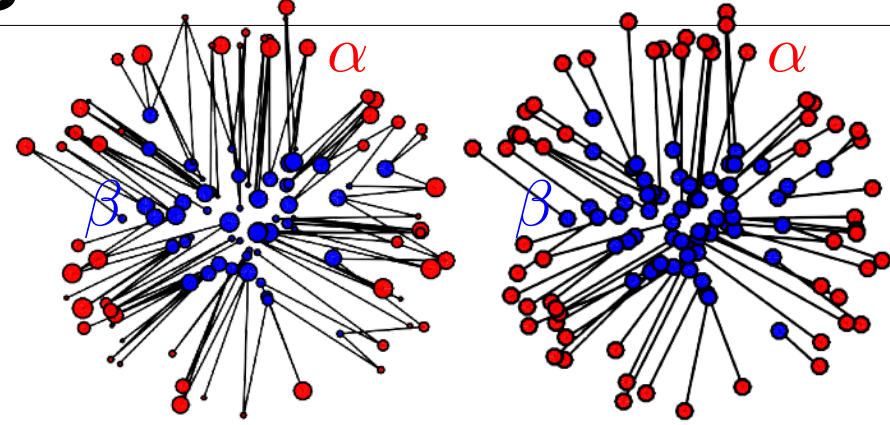


# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

Hungarian/Auction:  $O(n^3)$

$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$



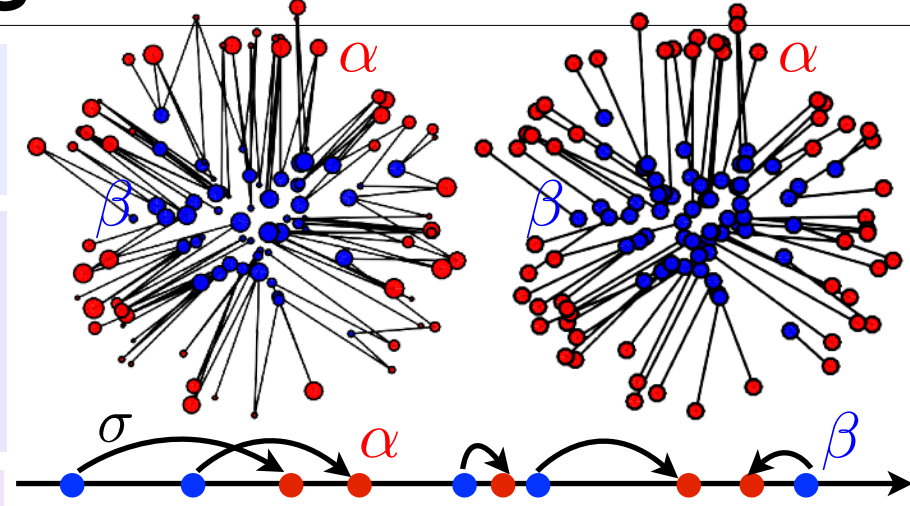
# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

Hungarian/Auction:  $O(n^3)$

$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

1-D case: sorting  $O(n \log(n))$ .



# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

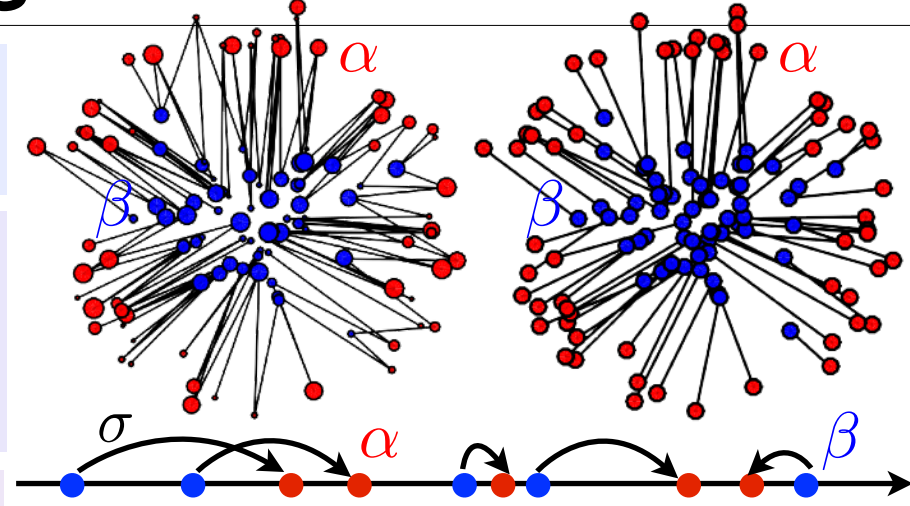
Hungarian/Auction:  $O(n^3)$

$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

1-D case: sorting  $O(n \log(n))$ .

$$p = 1 \quad W_1(\alpha, \beta) = \min_{\text{div}(u) = \alpha - \beta} \int \|u(x)\| dx$$
$$d = \|\cdot\|$$

→ min-cost flow, on graphs  $O(n^2 \log(n))$ .



# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

Hungarian/Auction:  $O(n^3)$

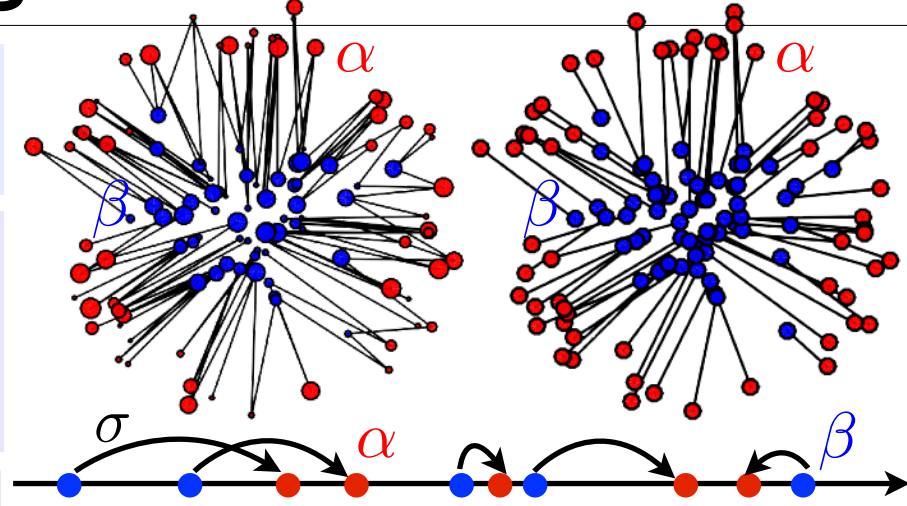
$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

1-D case: sorting  $O(n \log(n))$ .

$$p = 1 \quad W_1(\alpha, \beta) = \min_{\text{div}(u) = \alpha - \beta} \int \|u(x)\| dx$$
$$d = \|\cdot\|$$

→ min-cost flow, on graphs  $O(n^2 \log(n))$ .

Monge-Ampère/Benamou-Brenier,  $d = \|\cdot\|_2$ .





# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

Hungarian/Auction:  $O(n^3)$

$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

1-D case: sorting  $O(n \log(n))$ .

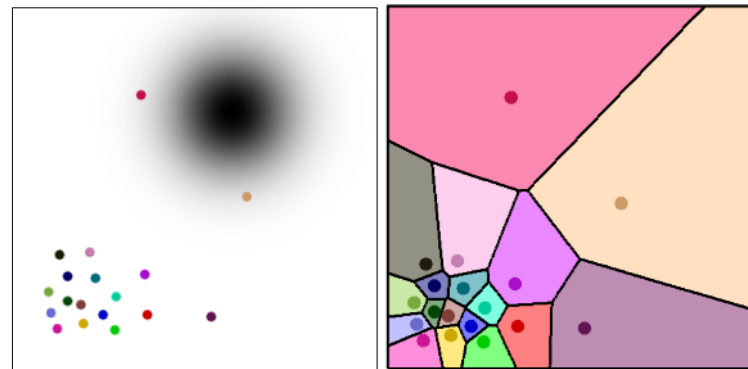
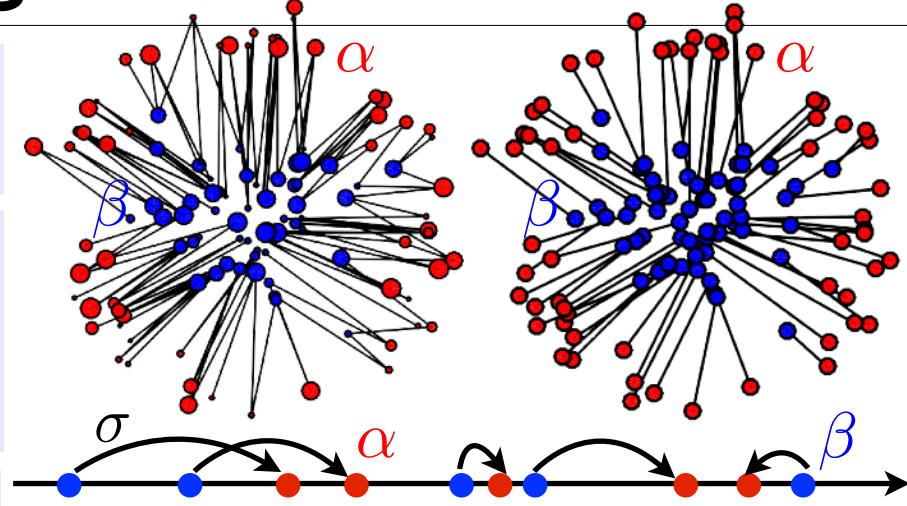
$$p = 1 \quad W_1(\alpha, \beta) = \min_{\text{div}(u) = \alpha - \beta} \int \|u(x)\| dx$$

$$d = \|\cdot\|$$

→ min-cost flow, on graphs  $O(n^2 \log(n))$ .

Monge-Ampère/Benamou-Brenier,  $d = \|\cdot\|_2$ .

Semi-discrete: Laguerre cells,  $d = \|\cdot\|_2$ .  
[Merigot 2013]



# A Glimpse at Algorithms

Linear programming:  $O(n^3 \log(n)^2)$

Hungarian/Auction:  $O(n^3)$

$$\alpha = \frac{1}{n} \sum_{i=1}^n \delta_{x_i} \quad \beta = \frac{1}{n} \sum_{j=1}^n \delta_{y_j}$$

1-D case: sorting  $O(n \log(n))$ .

$$p = 1 \quad W_1(\alpha, \beta) = \min_{\text{div}(u) = \alpha - \beta} \int \|u(x)\| dx$$

$$d = \|\cdot\|$$

→ min-cost flow, on graphs  $O(n^2 \log(n))$ .

Monge-Ampère/Benamou-Brenier,  $d = \|\cdot\|_2$ .

Semi-discrete: Laguerre cells,  $d = \|\cdot\|_2$ .  
[Merigot 2013]

Entropic regularization: generic  $d$ .

