

# Optimal Transport for Imaging and Learning

Gabriel Peyré



Marco Cuturi



Aude Genevay

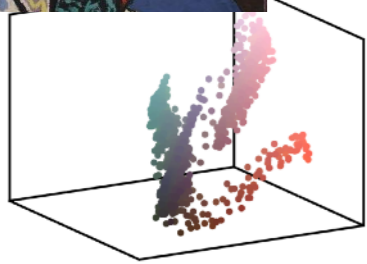
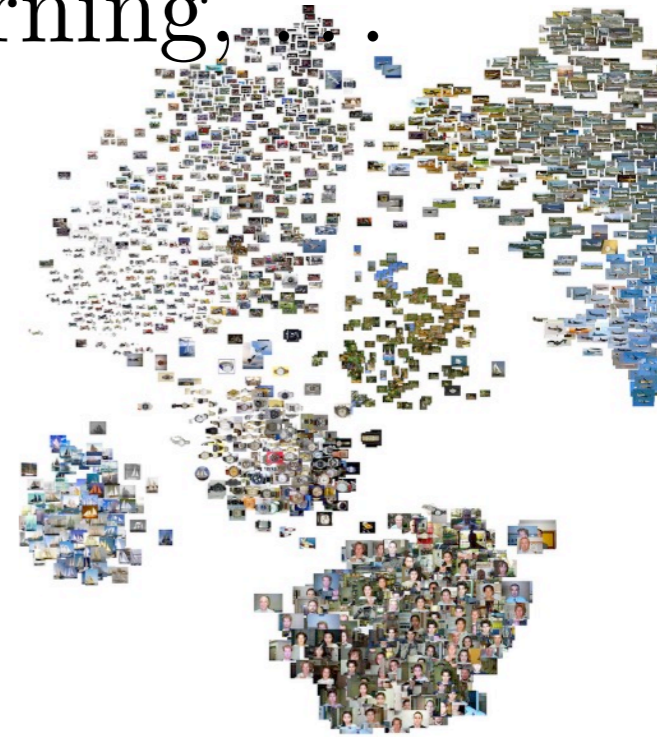
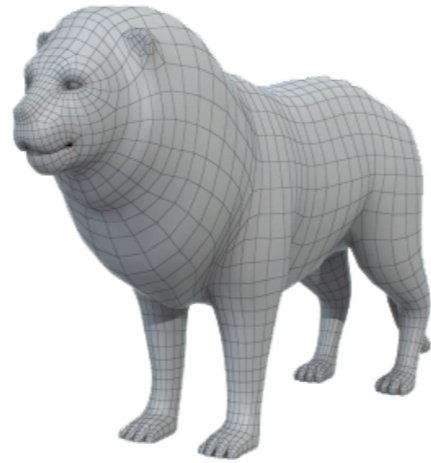
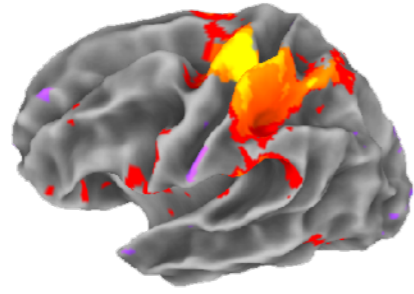


[www.numerical-tours.com](http://www.numerical-tours.com)



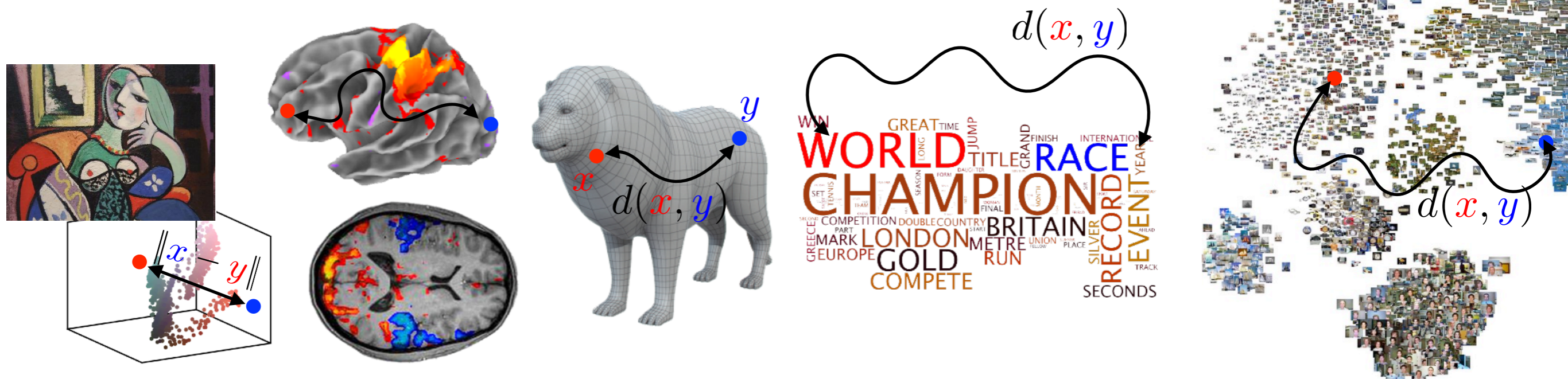
# Comparing Measures and Spaces

- *Probability distributions and histograms*  
→ images, vision, graphics and machine learning, .

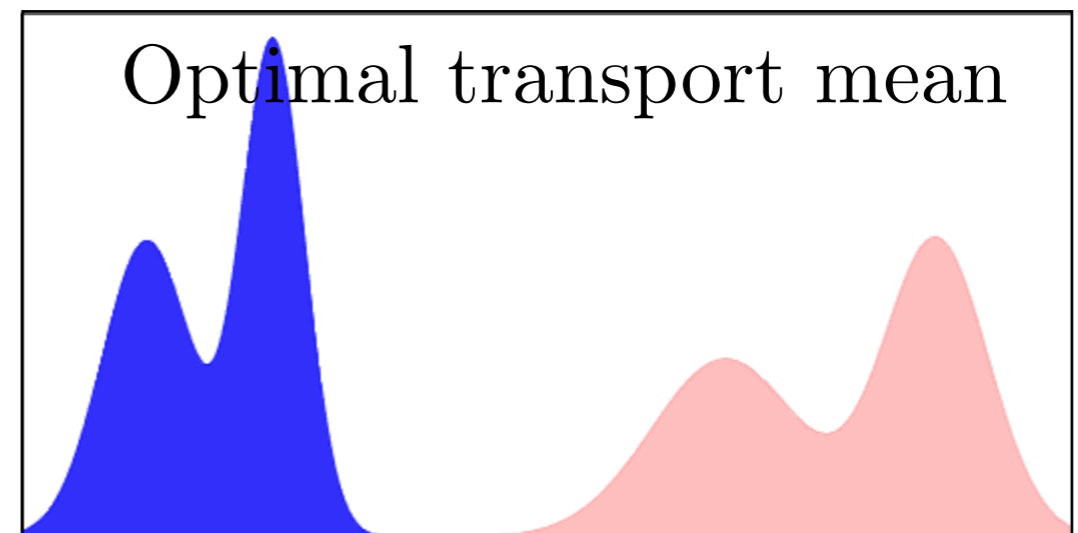
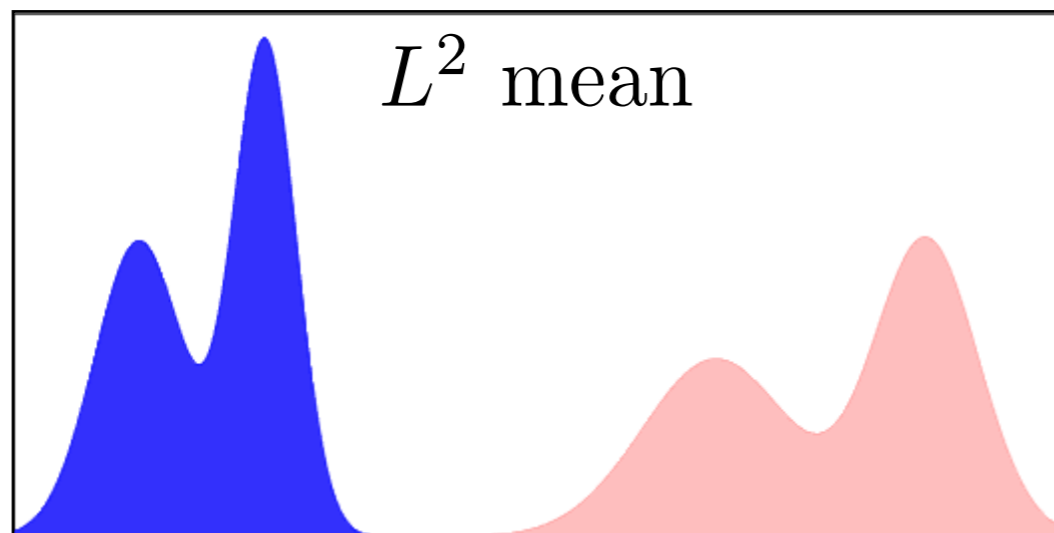


# Comparing Measures and Spaces

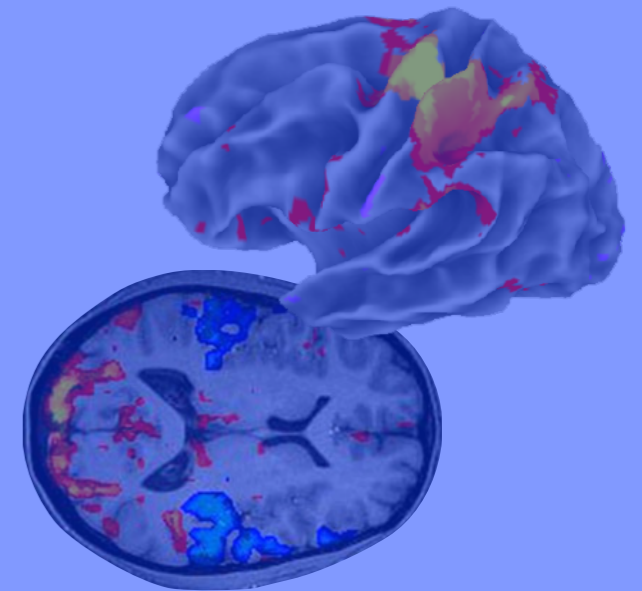
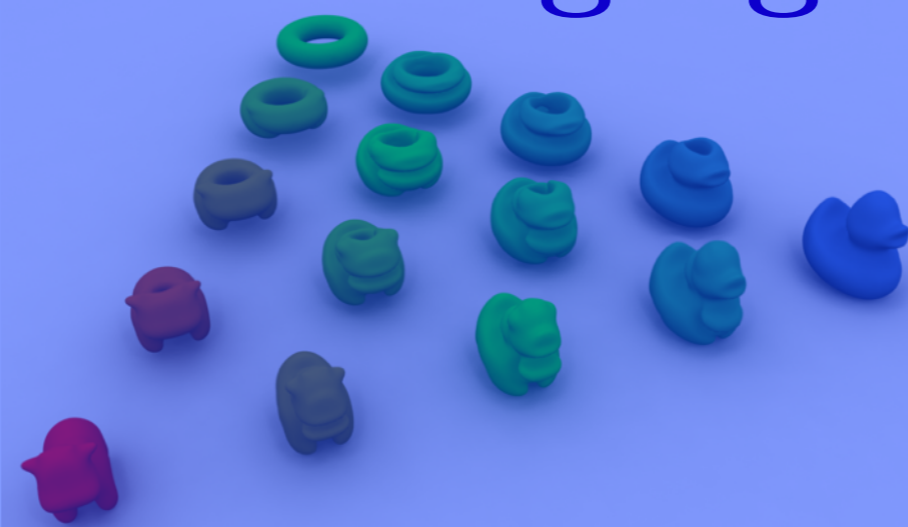
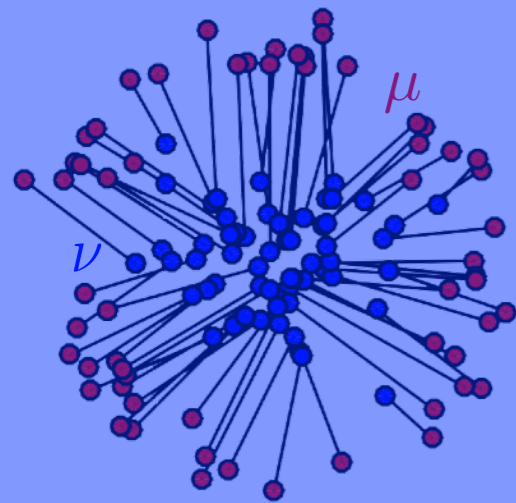
- *Probability distributions and histograms*  
→ images, vision, graphics and machine learning, .



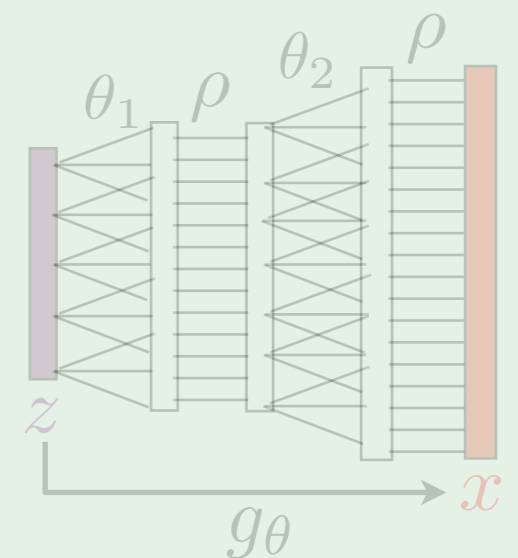
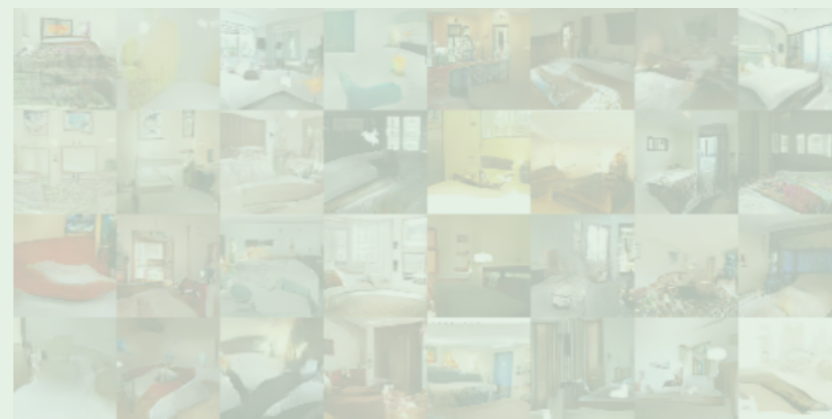
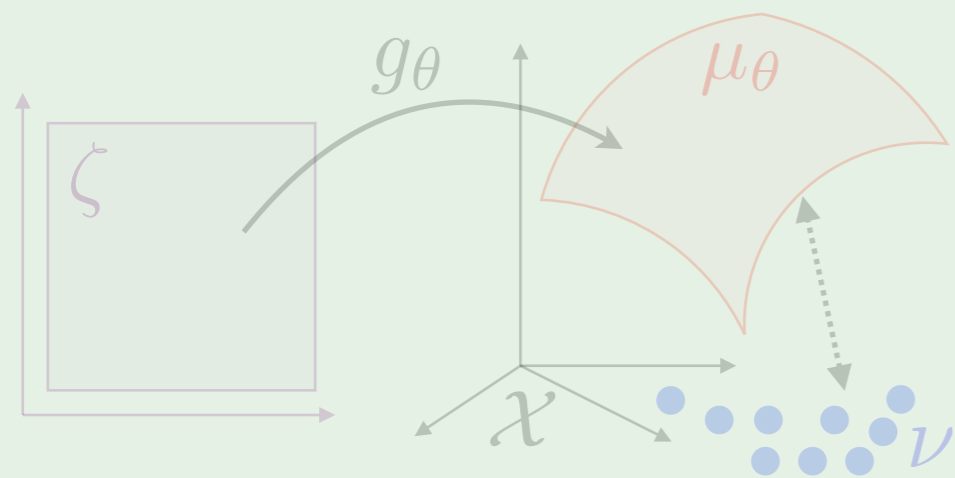
- *Optimal transport*  
→ takes into account a metric  $d$ .



# 1. OT for Imaging Sciences



# 2. OT for Machine Learning



# Couplings and Optimal Transport

Input distributions

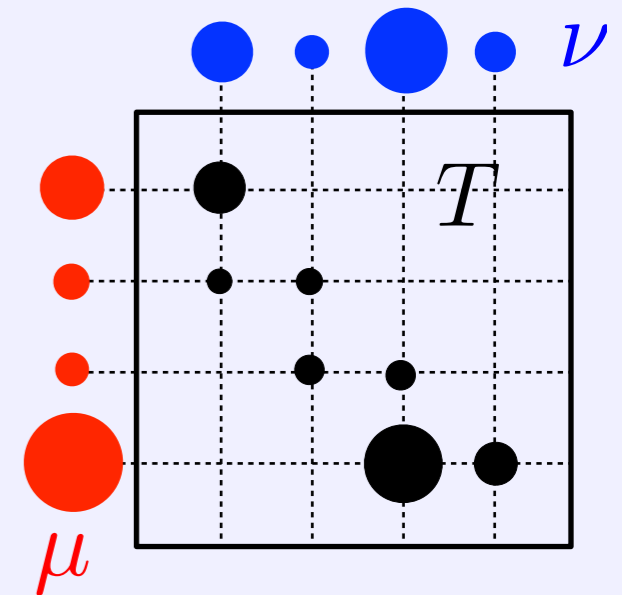
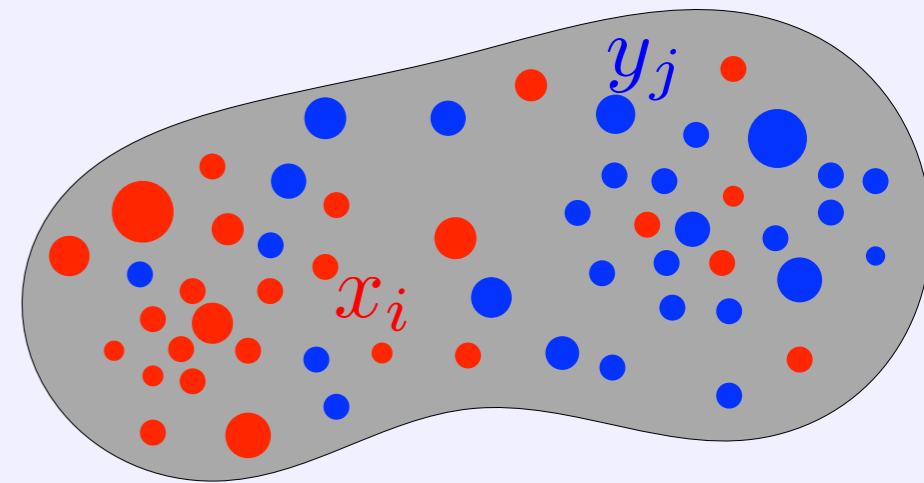
$$\mu = \sum_i \mu_i \delta_{x_i}$$

$$\nu = \sum_j \nu_j \delta_{y_j}$$

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

Weights  $\mu_i \geq 0, \nu_j \geq 0$ .

$$\sum_{i=1}^{N_1} \mu_i = \sum_{j=1}^{N_2} \nu_j = 1$$



**Def.** *Couplings*

$$\mathcal{C}_{\mu, \nu} \stackrel{\text{def.}}{=} \left\{ T \in \mathbb{R}_+^{N_1 \times N_2} ; T \mathbf{1}_{N_1} = \mu, T^\top \mathbf{1}_{N_2} = \nu \right\}$$

# Couplings and Optimal Transport

Input distributions

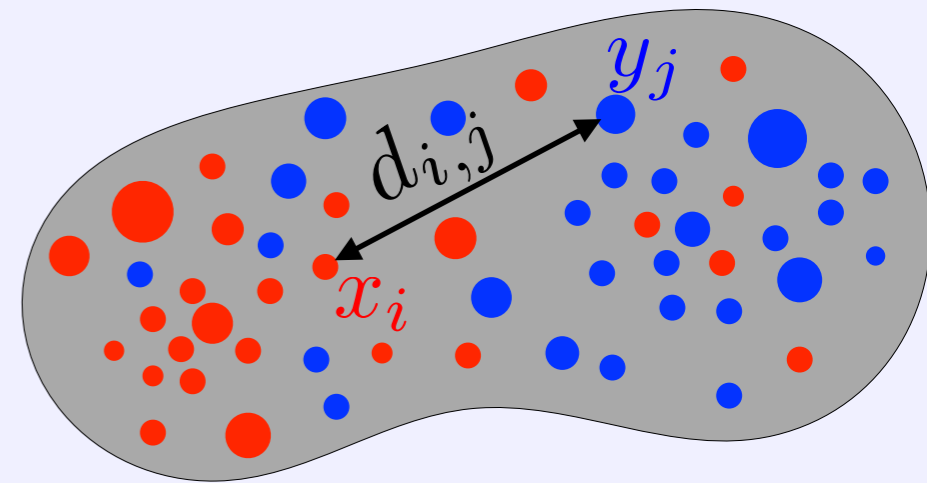
$$\mu = \sum_i \mu_i \delta_{x_i}$$

$$\nu = \sum_j \nu_j \delta_{y_j}$$

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

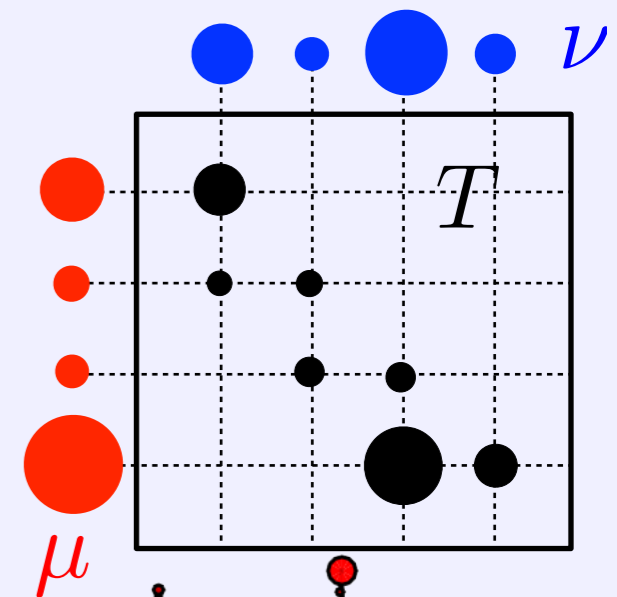
Weights  $\mu_i \geq 0, \nu_j \geq 0$ .

$$\sum_{i=1}^{N_1} \mu_i = \sum_{j=1}^{N_2} \nu_j = 1 \quad d_{i,j} = d(x_i, y_j)$$



**Def.** *Couplings*

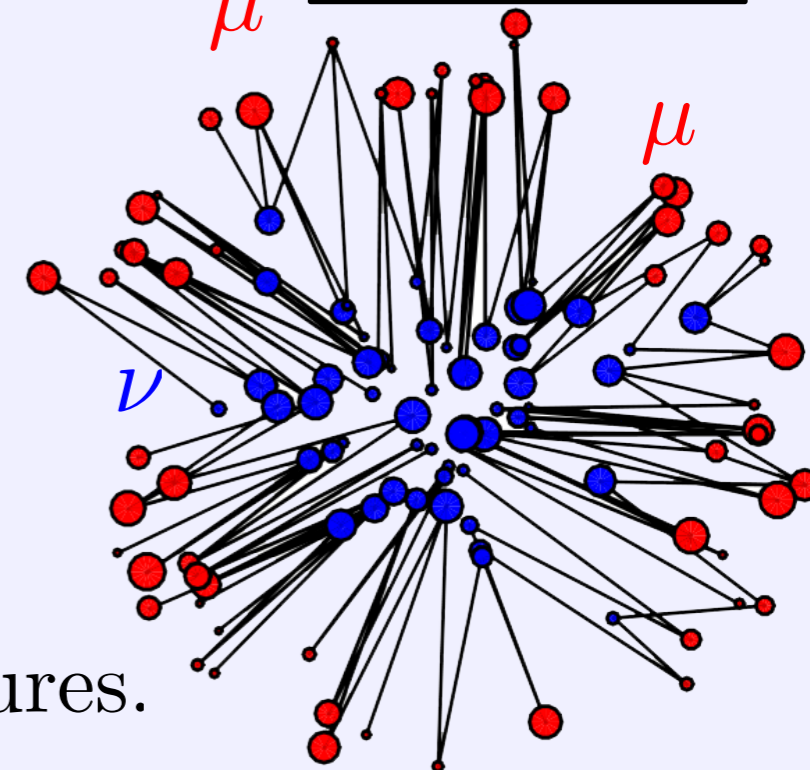
$$\mathcal{C}_{\mu, \nu} \stackrel{\text{def.}}{=} \left\{ T \in \mathbb{R}_+^{N_1 \times N_2} ; T \mathbf{1}_{N_1} = \mu, T^\top \mathbf{1}_{N_2} = \nu \right\}$$



**Def.** *Wasserstein Distance / EMD*

$$W_p^p(\mu, \nu) \stackrel{\text{def.}}{=} \min \left\{ \sum_{i,j} T_{i,j} d_{i,j}^p ; T \in \mathcal{C}_{\mu, \nu} \right\}$$

[Kantorovich 1942]

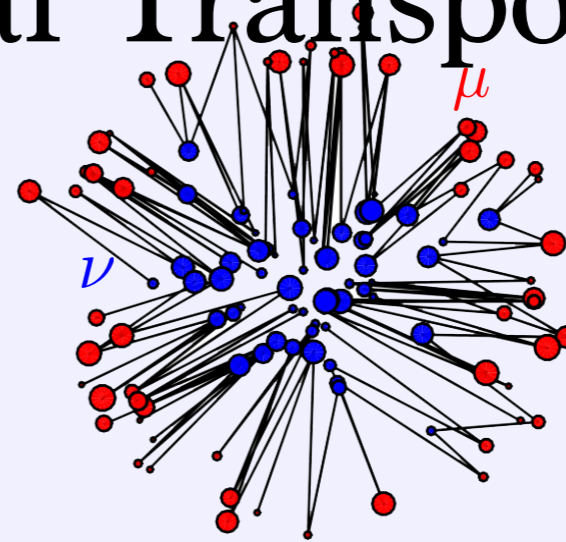


$\rightarrow W_p$  is a distance over Radon probability measures.

# Numerical Optimal Transport

Linear programming:

$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$



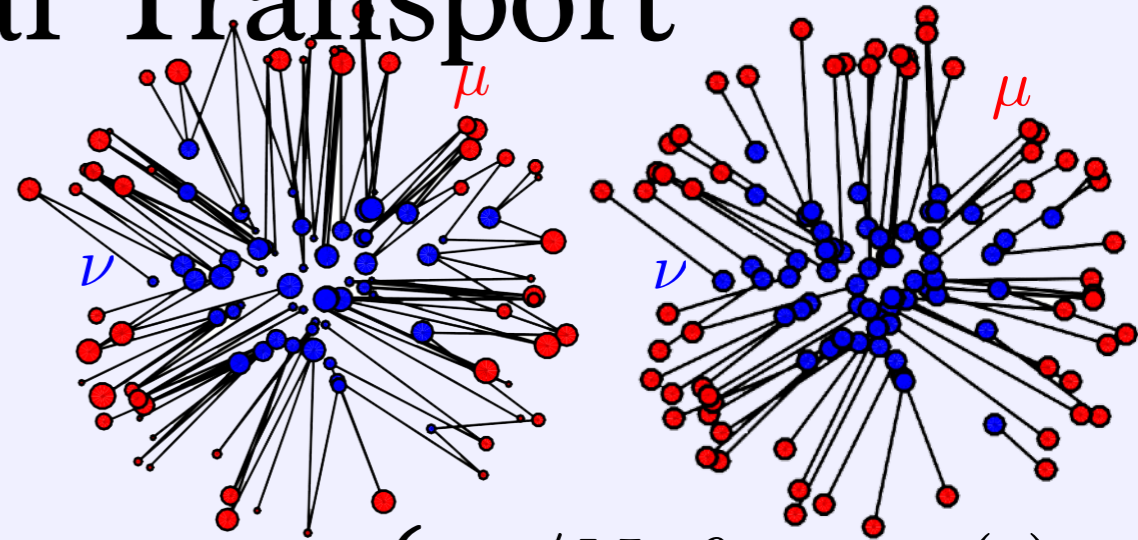
# Numerical Optimal Transport

Linear programming:

$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

Hungarian/Auction:  $\sim O(N^3)$

$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$



# Numerical Optimal Transport

Linear programming:

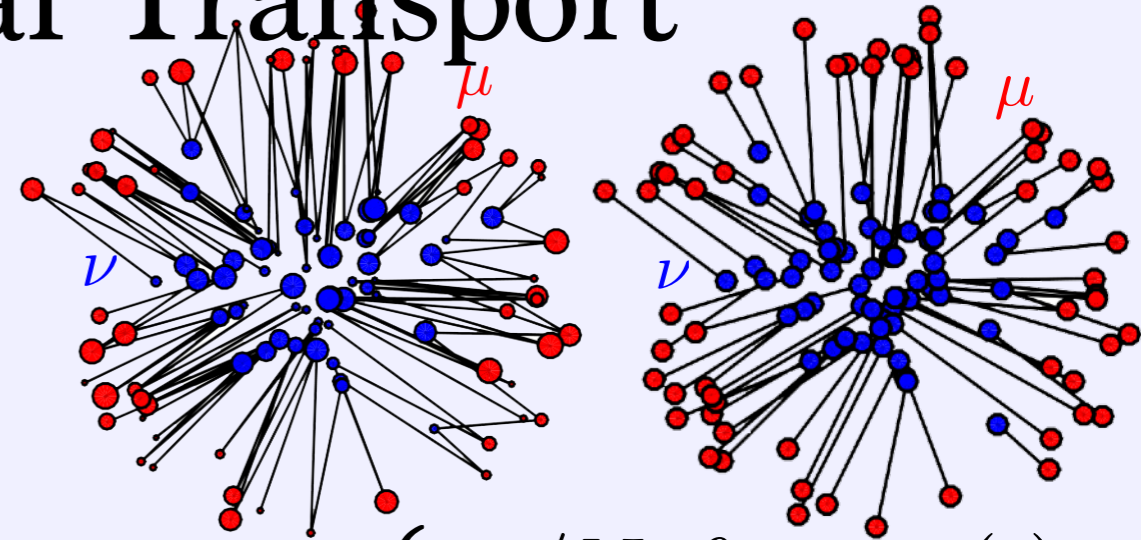
$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

Hungarian/Auction:  $\sim O(N^3)$

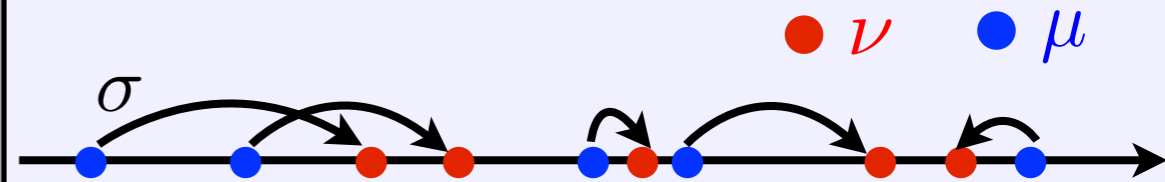
$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$

1-D case,  $d = |\cdot|^p, p \geq 1$ .

→ sorting,  $O(N \log(N))$  operations.



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$



# Numerical Optimal Transport

Linear programming:

$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

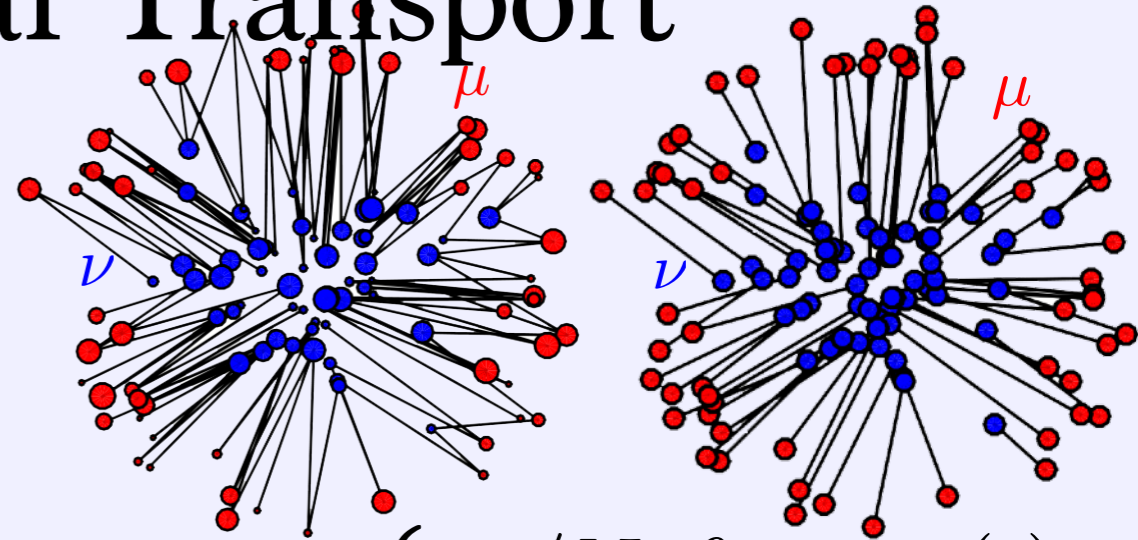
Hungarian/Auction:  $\sim O(N^3)$

$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$

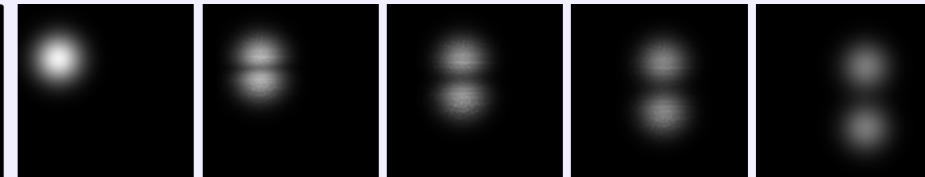
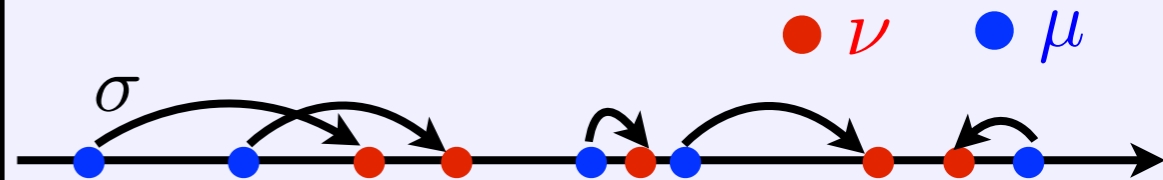
1-D case,  $d = |\cdot|^p, p \geq 1$ .

→ sorting,  $O(N \log(N))$  operations.

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



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$



# Numerical Optimal Transport

Linear programming:

$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

Hungarian/Auction:  $\sim O(N^3)$

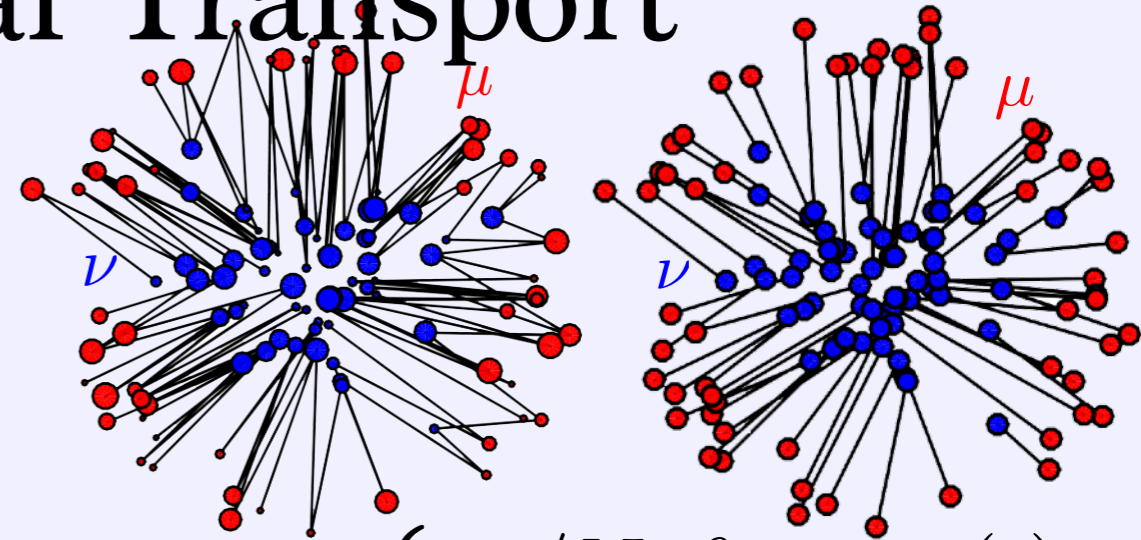
$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$

1-D case,  $d = |\cdot|^p, p \geq 1$ .

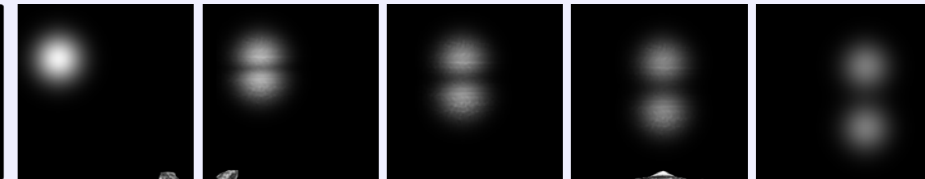
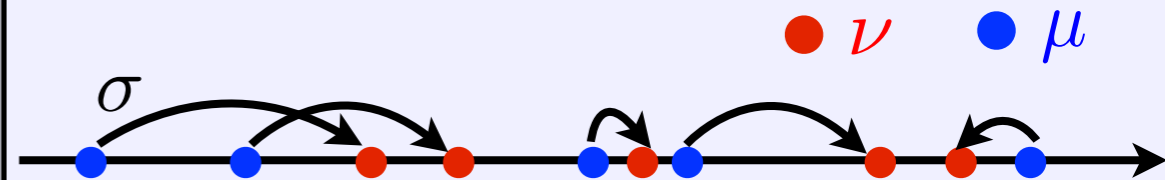
→ sorting,  $O(N \log(N))$  operations.

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

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



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$



[Levy, '15]

# Numerical Optimal Transport

Linear programming:

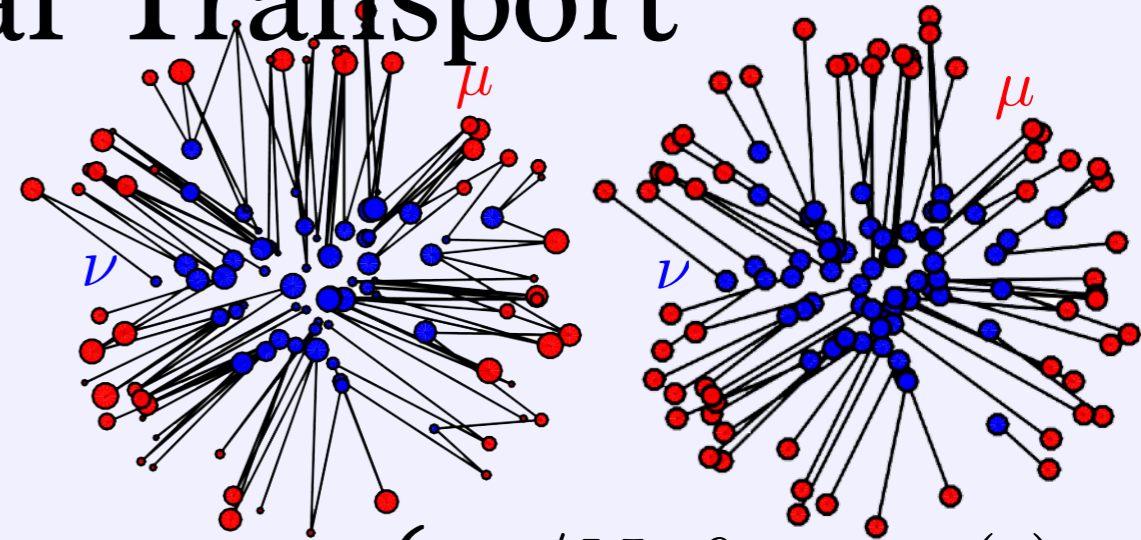
$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

Hungarian/Auction:  $\sim O(N^3)$

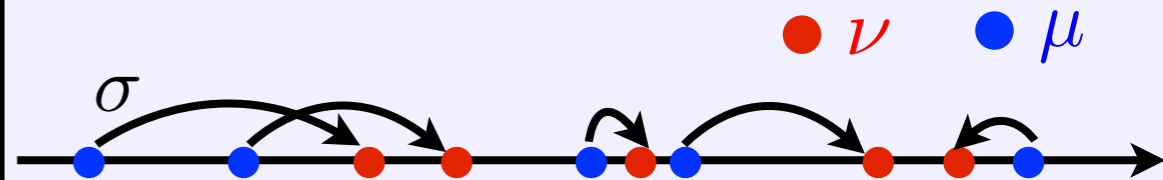
$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$

1-D case,  $d = |\cdot|^p, p \geq 1$ .

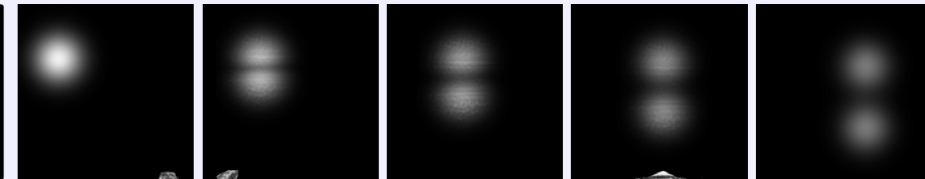
→ sorting,  $O(N \log(N))$  operations.



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$

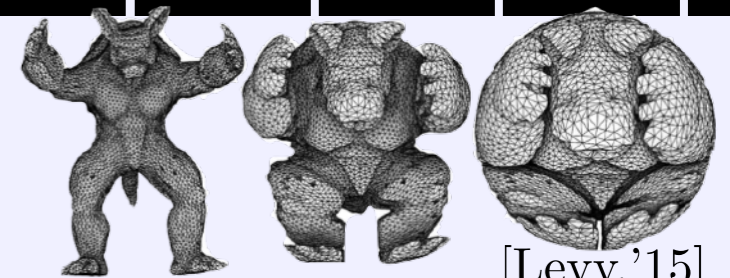


Monge-Ampère/Benamou-Brenier,  $d = \|\cdot\|_2$ .



Semi-discrete: Laguerre cells,  $d = \|\cdot\|_2$ .

[Merigot 2013]



[Levy, '15]

$$d = \|\cdot\|, p = 1 : W_1(\mu, \nu) = \min_{\text{div}(v) = \mu - \nu} \int \|u(x)\| dx \rightarrow \text{max-flow.}$$

# Numerical Optimal Transport

Linear programming:

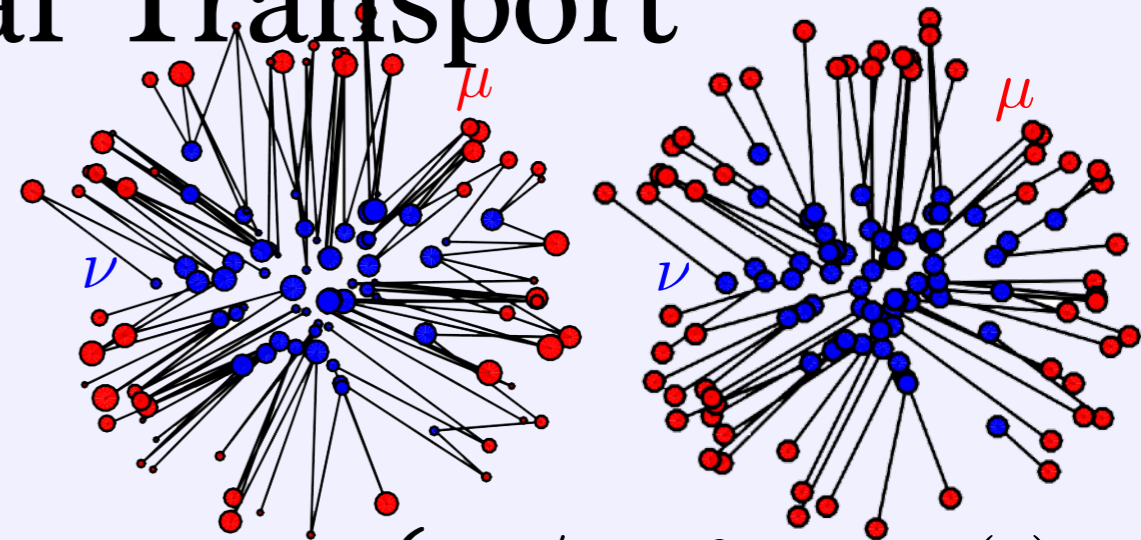
$$\mu = \sum_{i=1}^{N_1} p_i \delta_{x_i}, \nu = \sum_{j=1}^{N_2} p_j \delta_{y_j}$$

Hungarian/Auction:  $\sim O(N^3)$

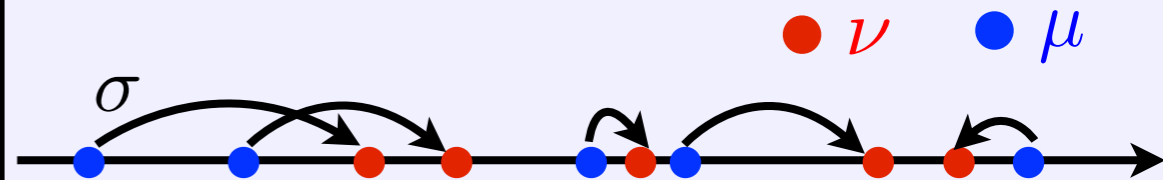
$$\mu = \frac{1}{N} \sum_{i=1}^N \delta_{x_i}, \nu = \frac{1}{N} \sum_{j=1}^N \delta_{y_j}$$

1-D case,  $d = |\cdot|^p, p \geq 1$ .

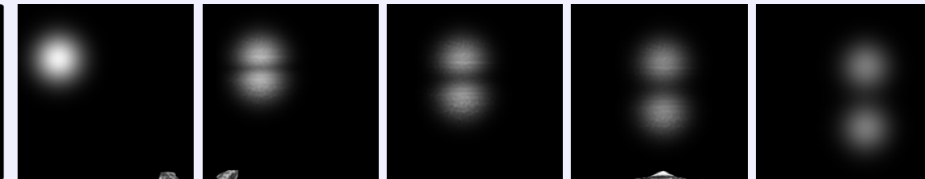
→ sorting,  $O(N \log(N))$  operations.



$$T_{i,j} = \begin{cases} 1/N & \text{if } j = \sigma(i), \\ 0 & \text{otherwise.} \end{cases}$$

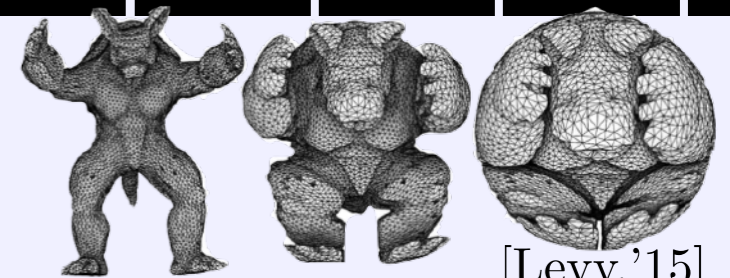


Monge-Ampère/Benamou-Brenier,  $d = \|\cdot\|_2$ .



Semi-discrete: Laguerre cells,  $d = \|\cdot\|_2$ .

[Merigot 2013]



[Levy, '15]

$$d = \|\cdot\|, p = 1 : W_1(\mu, \nu) = \min_{\text{div}(v) = \mu - \nu} \int \|u(x)\| dx \rightarrow \text{max-flow.}$$

Need for fast approximate algorithms for generic  $c$ .

# Entropic Regularization

*Entropy:*  $H(T) \stackrel{\text{def.}}{=} - \sum_{i,j=1}^N T_{i,j} (\log(T_{i,j}) - 1)$

**Def.** *Regularized OT:* [Cuturi NIPS'13]

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} - \varepsilon H(T) ; T \in \mathcal{C}_{\mu,\nu} \right\}$$

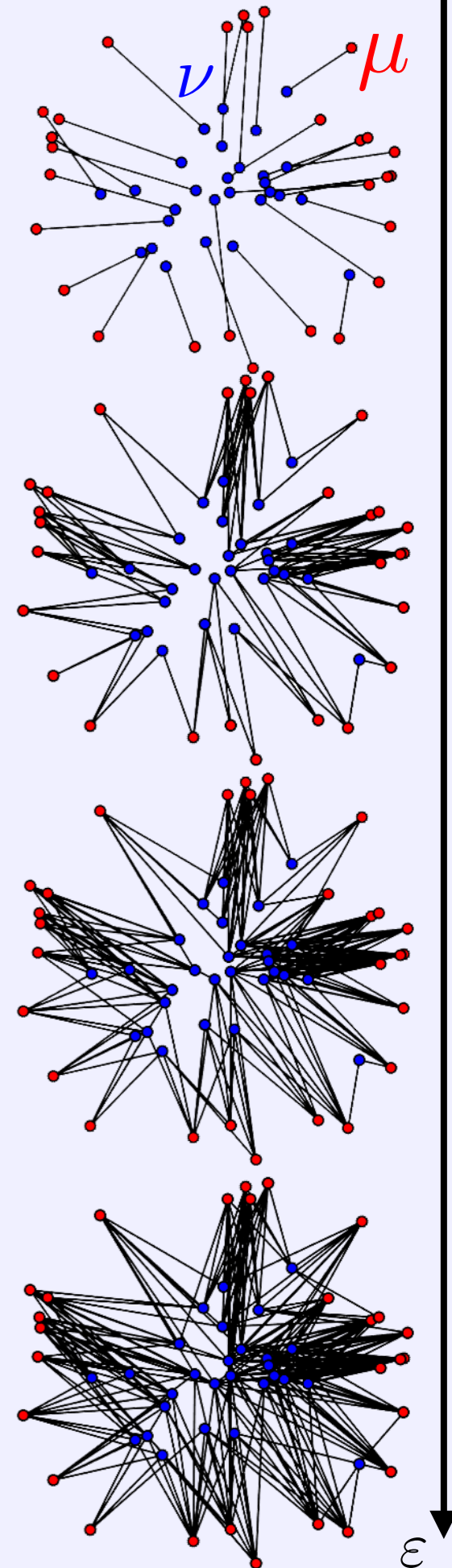
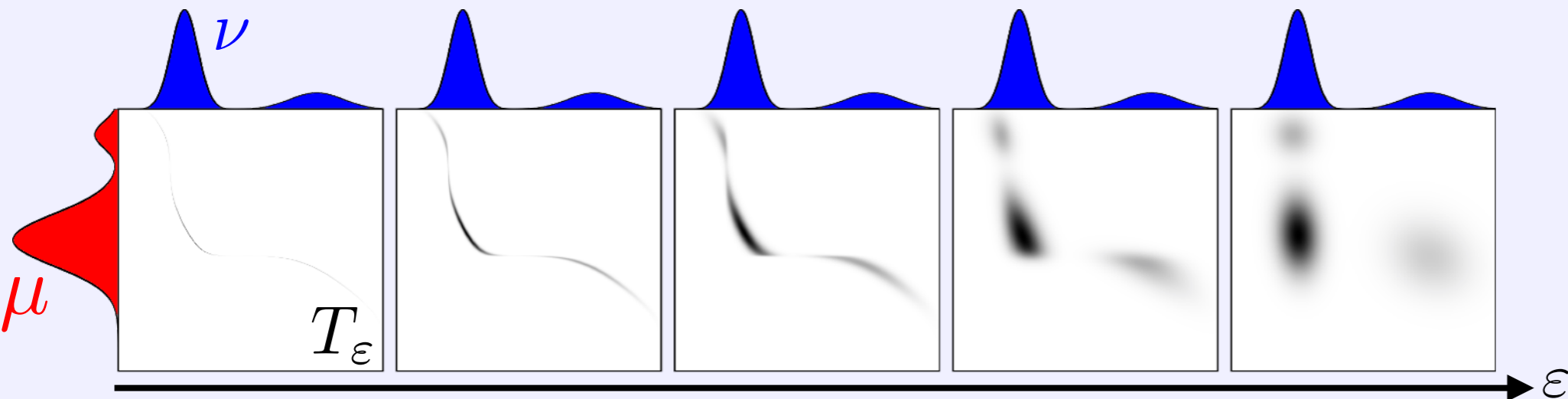
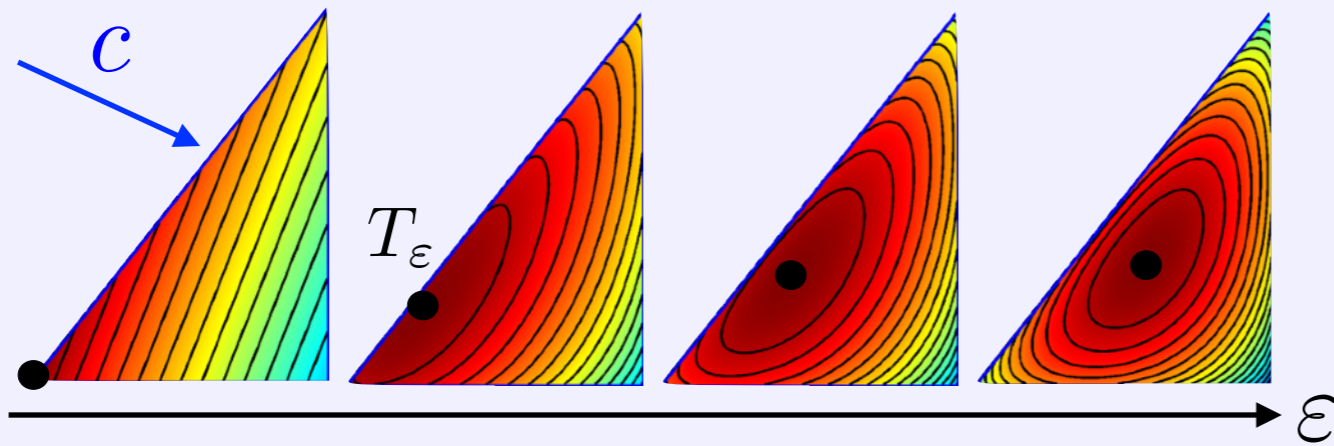
# Entropic Regularization

*Entropy:*  $H(T) \stackrel{\text{def.}}{=} - \sum_{i,j=1}^N T_{i,j} (\log(T_{i,j}) - 1)$

**Def. Regularized OT:** [Cuturi NIPS'13]

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} - \varepsilon H(T) ; T \in \mathcal{C}_{\mu,\nu} \right\}$$

*Regularization impact on solution:*



# Sinkhorn's Algorithm

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} + \varepsilon T_{i,j} \log(T_{i,j}) ; T \in \mathcal{C}_{\mu,\nu} \right\} \quad (\star)$$

**Prop.** One has  $T = \text{diag}(a)K \text{diag}(b)$ , where  $K = e^{-\frac{d^p}{\varepsilon}}$ .



# Sinkhorn's Algorithm

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} + \varepsilon T_{i,j} \log(T_{i,j}) ; T \in \mathcal{C}_{\mu,\nu} \right\} \quad (\star)$$

**Prop.** One has  $T = \text{diag}(a)K \text{diag}(b)$ , where  $K = e^{-\frac{d^p}{\varepsilon}}$ .

Row constraint:  $T \mathbf{1}_{N_2} = \mu \iff a \odot (Kb) = \mu$

Col. constraint:  $T^\top \mathbf{1}_{N_2} = \nu \iff b \odot (K^\top a) = \nu$

Sinkhorn iterations:  $a \leftarrow \frac{\mu}{Kb}$  and  $b \leftarrow \frac{\nu}{K^\top a}$

# Sinkhorn's Algorithm

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} + \varepsilon T_{i,j} \log(T_{i,j}) ; T \in \mathcal{C}_{\mu,\nu} \right\} \quad (\star)$$

**Prop.** One has  $T = \text{diag}(a)K \text{diag}(b)$ , where  $K = e^{-\frac{d^p}{\varepsilon}}$ .

Row constraint:  $T \mathbf{1}_{N_2} = \mu \iff a \odot (Kb) = \mu$

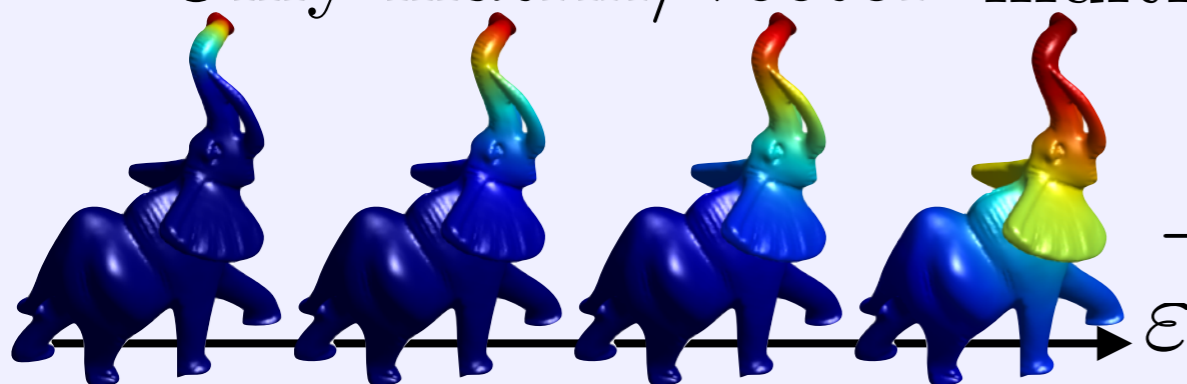
Col. constraint:  $T^\top \mathbf{1}_{N_2} = \nu \iff b \odot (K^\top a) = \nu$

Sinkhorn iterations:  $a \leftarrow \frac{\mu}{Kb}$  and  $b \leftarrow \frac{\nu}{K^\top a}$

Only matrix/vector multiplications.  $\rightarrow$  Parallelizable.

$\rightarrow$  Streams well on GPU.

$\rightarrow$  convolutive/heat structure for  $K$



# Sinkhorn's Algorithm

$$\min_T \left\{ \sum_{i,j} d_{i,j}^p T_{i,j} + \varepsilon T_{i,j} \log(T_{i,j}) ; T \in \mathcal{C}_{\mu,\nu} \right\} \quad (\star)$$

**Prop.** One has  $T = \text{diag}(a)K \text{diag}(b)$ , where  $K = e^{-\frac{d^p}{\varepsilon}}$ .

Row constraint:  $T \mathbf{1}_{N_2} = \mu \iff a \odot (Kb) = \mu$

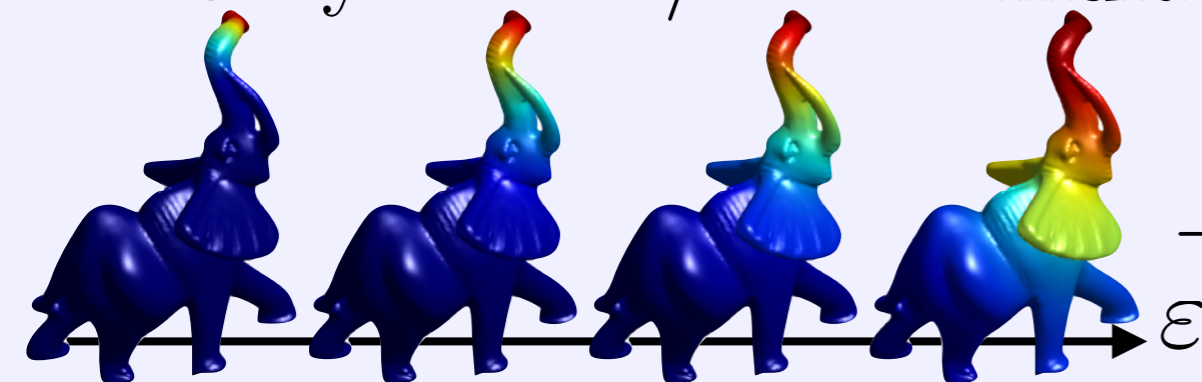
Col. constraint:  $T^\top \mathbf{1}_{N_2} = \nu \iff b \odot (K^\top a) = \nu$

Sinkhorn iterations:  $a \leftarrow \frac{\mu}{Kb}$  and  $b \leftarrow \frac{\nu}{K^\top a}$

Only matrix/vector multiplications.  $\rightarrow$  Parallelizable.

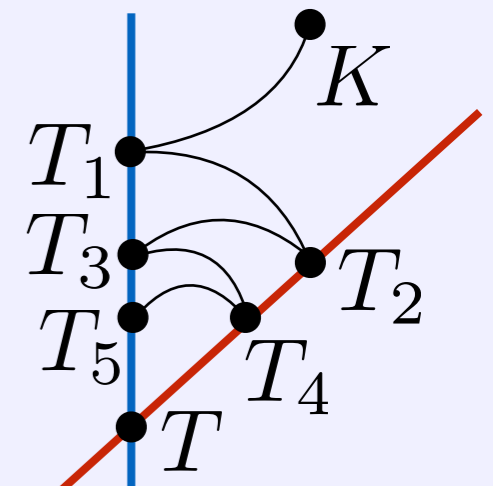
$\rightarrow$  Streams well on GPU.

$\rightarrow$  convolutive/heat structure for  $K$



**Prop.**  $(\star) \iff \min_T \{ \text{KL}(T|K) ; T \in \mathcal{C}_{\mu,\nu} \}$

Sinkhorn  $\iff$  iterative projections.

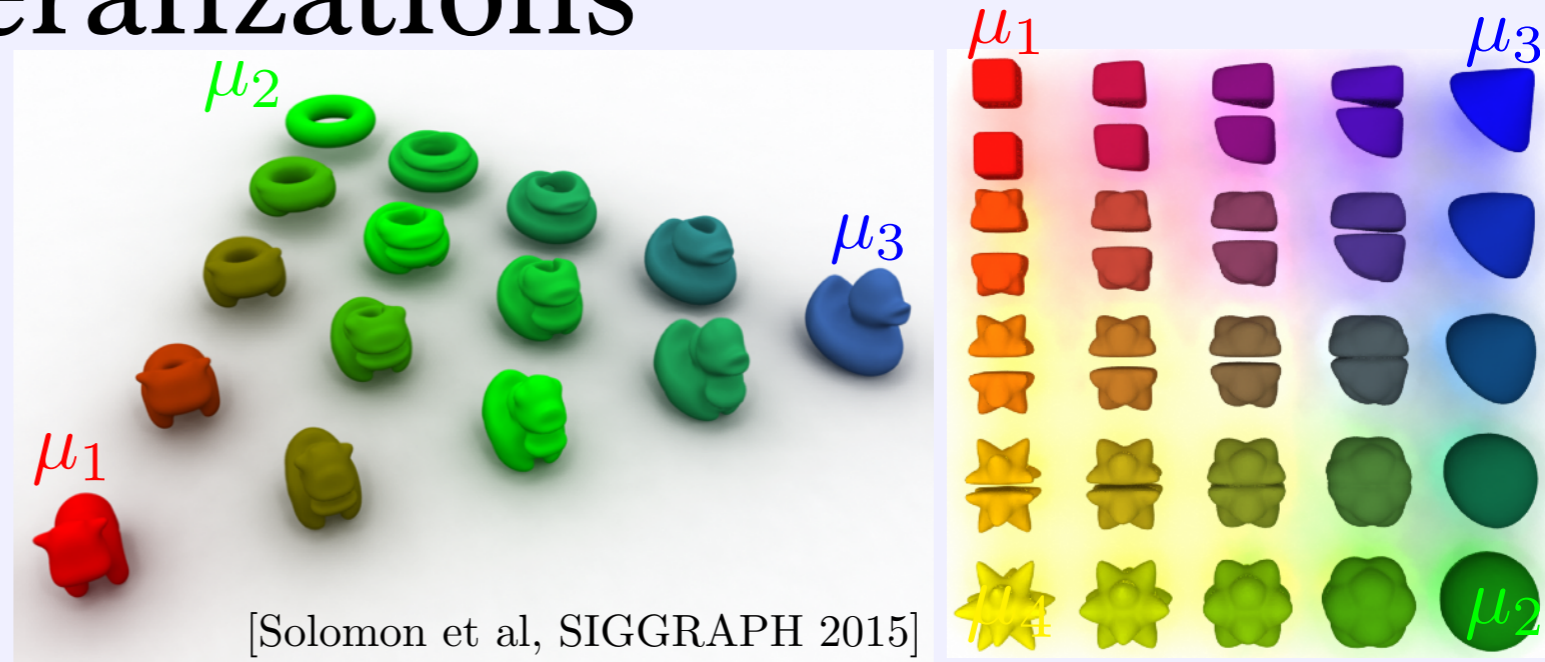


# Generalizations

*OT barycenters:*

$$\min_{\nu} \sum_k \lambda_k W_2^2(\mu_k, \nu)$$

[Agueh, Carlier 2010]

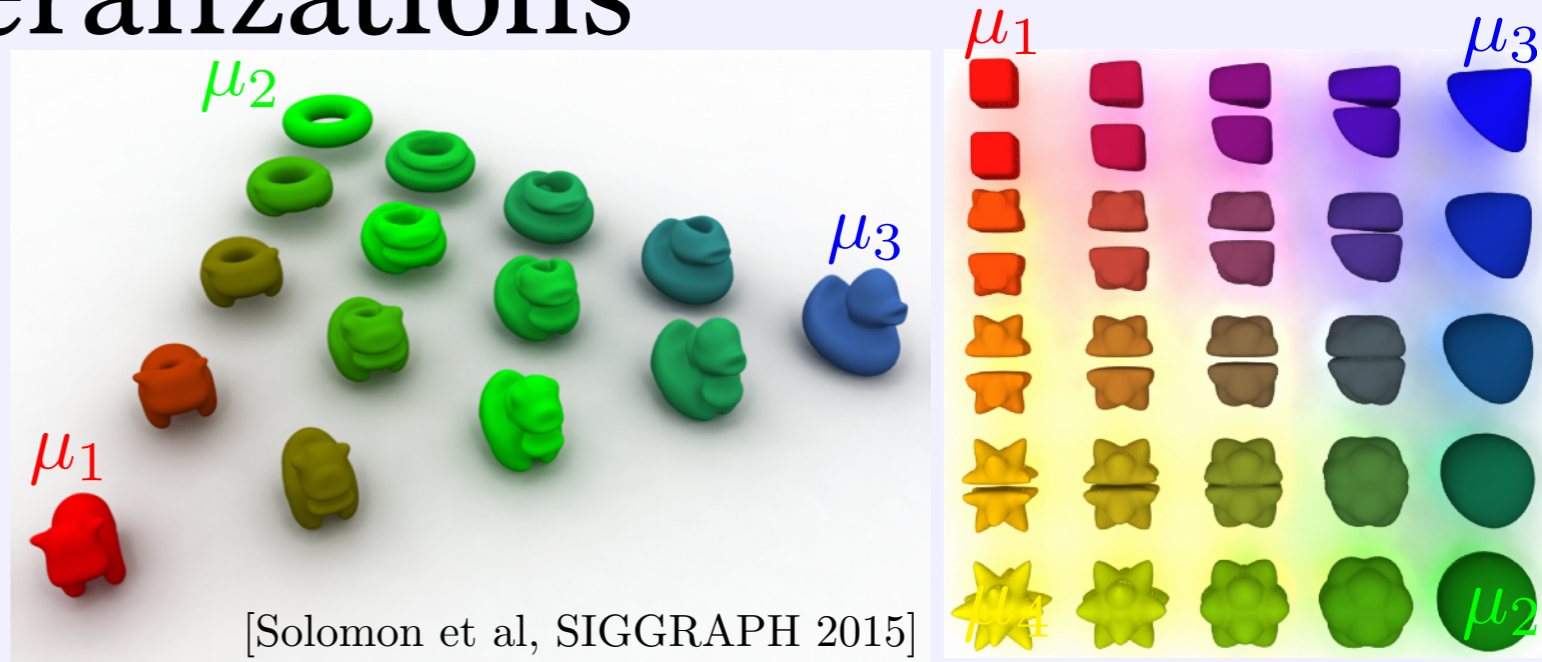


# Generalizations

*OT barycenters:*

$$\min_{\nu} \sum_k \lambda_k W_2^2(\mu_k, \nu)$$

[Agueh, Carlier 2010]



[Solomon et al, SIGGRAPH 2015]

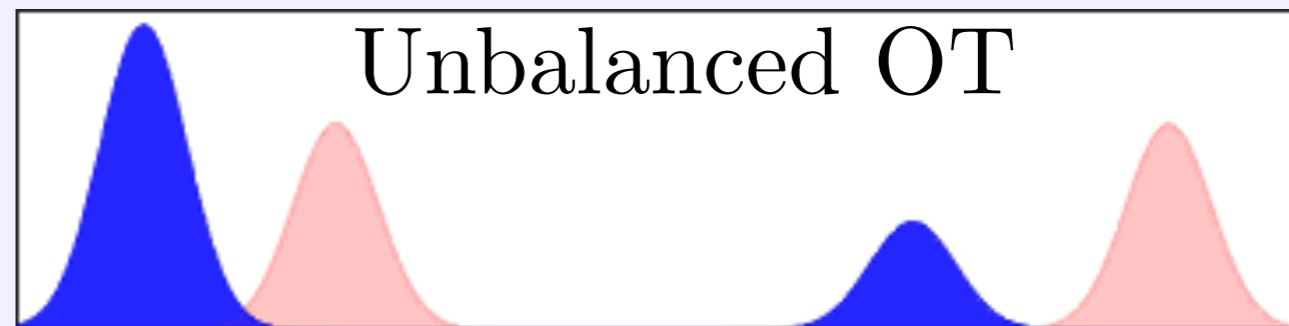
*Unbalanced transport:*

$$\min_T \sum_{i,j} d_{i,j}^p T_{i,j} + \rho \text{KL}(T \mathbf{1}_{N_1} | \mu) + \rho \text{KL}(T^\top \mathbf{1}_{N_2} | \nu)$$

[Liereo, Mielke, Savaré 2015]

[Chizat, Schmitzer, Peyré, Vialard 2015]

[Kondratyev, Monsaingeon, Vorotnikov, 2015]

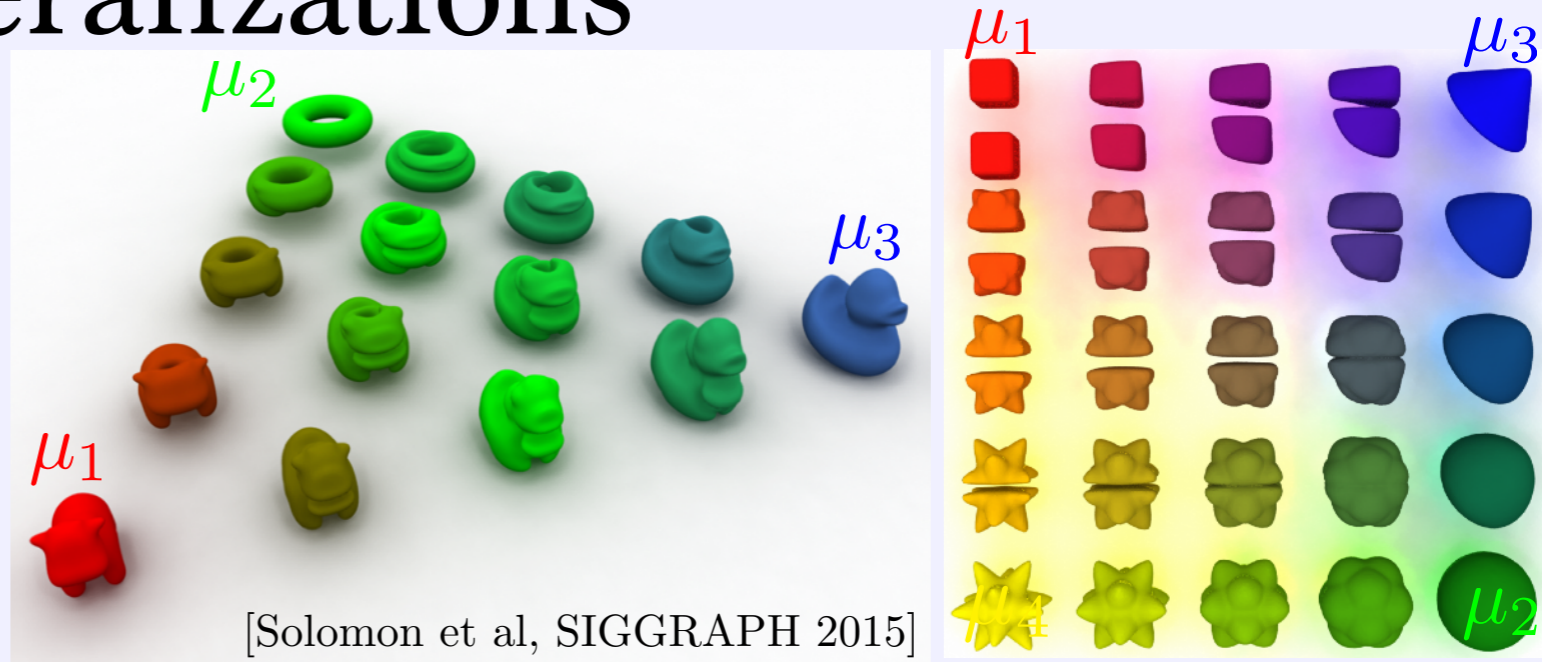


# Generalizations

*OT barycenters:*

$$\min_{\nu} \sum_k \lambda_k W_2^2(\mu_k, \nu)$$

[Agueh, Carlier 2010]



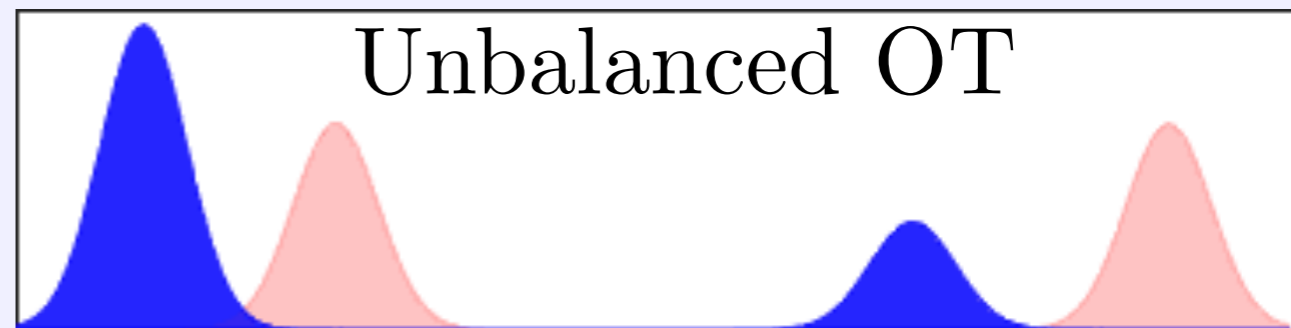
*Unbalanced transport:*

$$\min_T \sum_{i,j} d_{i,j}^p T_{i,j} + \rho \text{KL}(T \mathbf{1}_{N_1} | \mu) + \rho \text{KL}(T^\top \mathbf{1}_{N_2} | \nu)$$

[Liereo, Mielke, Savaré 2015]

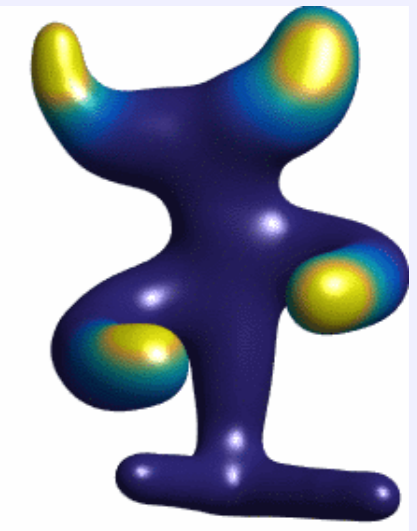
[Chizat, Schmitzer, Peyré, Vialard 2015]

[Kondratyev, Monsaingeon, Vorotnikov, 2015]

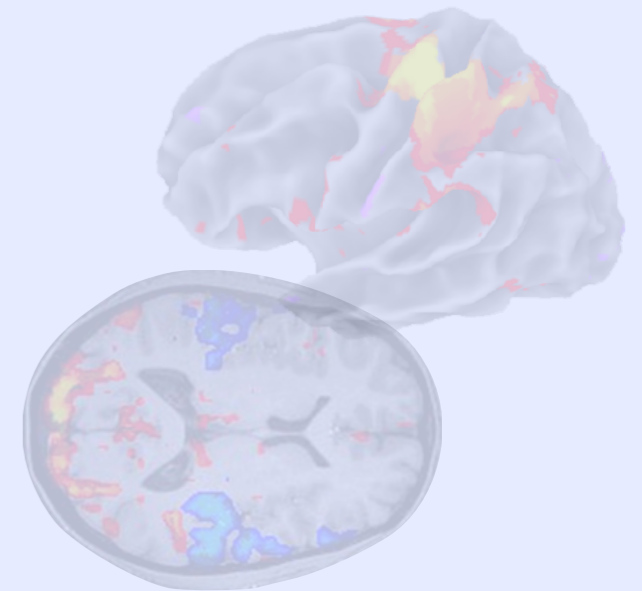
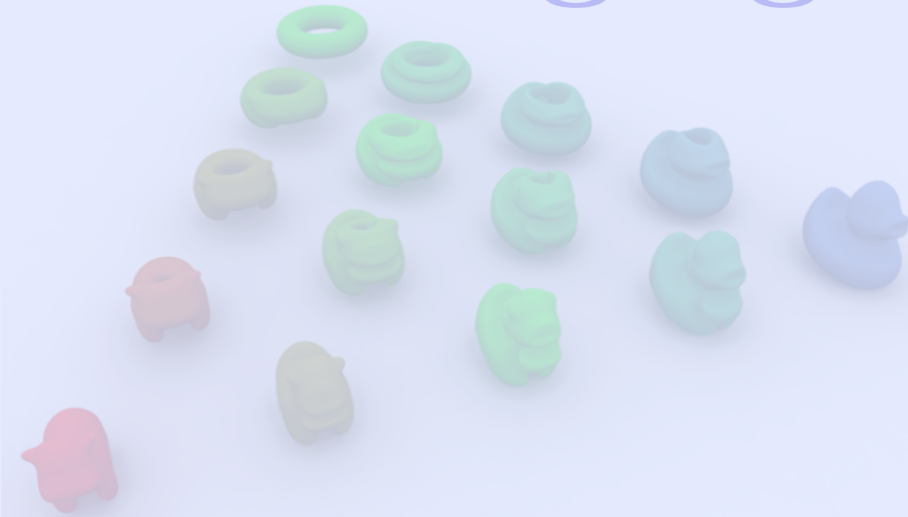


*Gradient flows:*

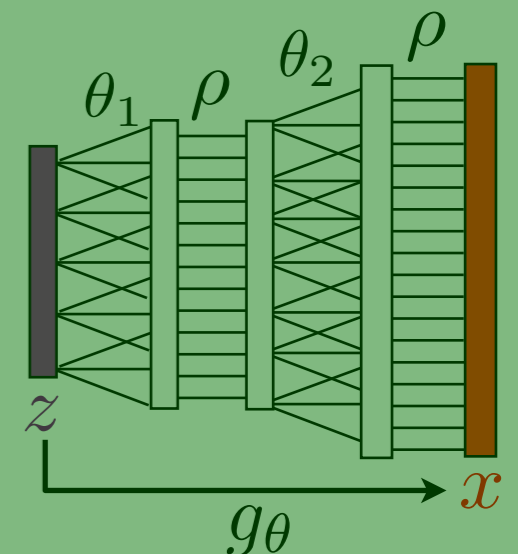
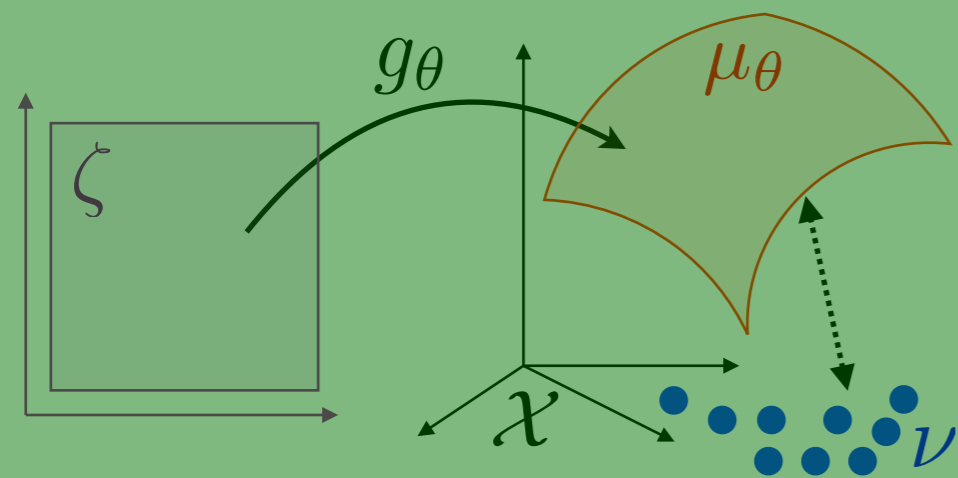
$$\mu_{t+1} = \min_{\mu} \frac{1}{2\tau} W(\mu_t, \mu) + f(\mu)$$



# 1. OT for Imaging Sciences



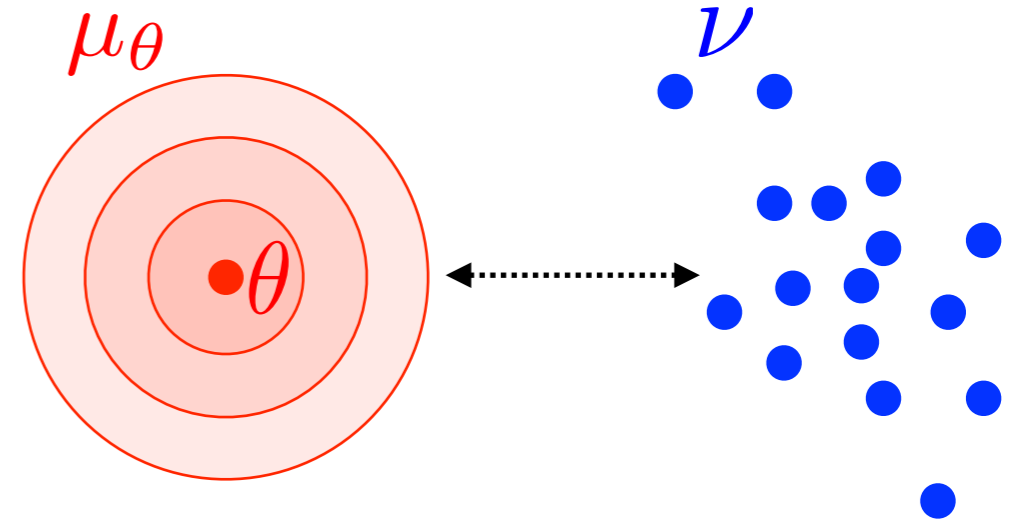
# 2. OT for Machine Learning



# Density Fitting and Generative Models

*Observations:*  $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

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

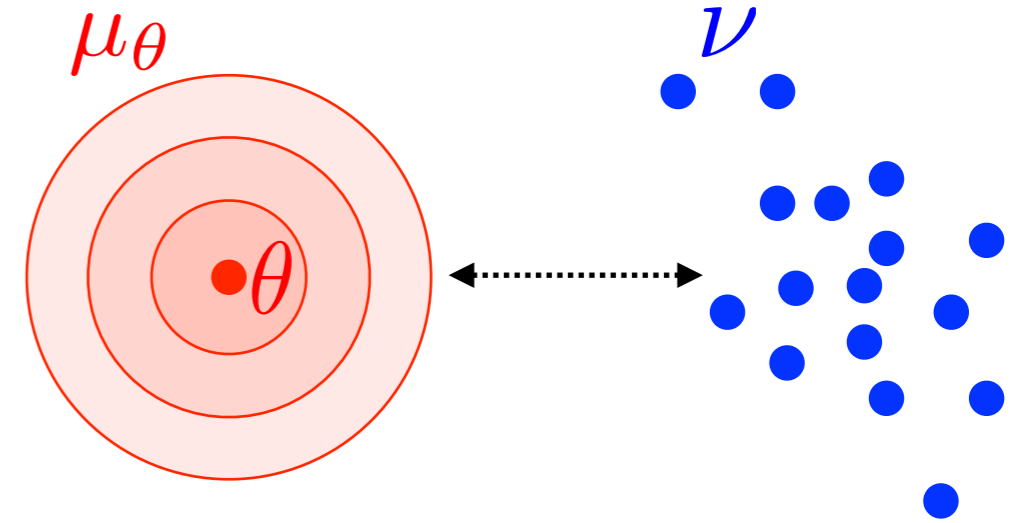




# Density Fitting and Generative Models

Observations:  $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \mu_\theta$



Density fitting:  $d\mu_\theta(y) = f_\theta(y)dy$

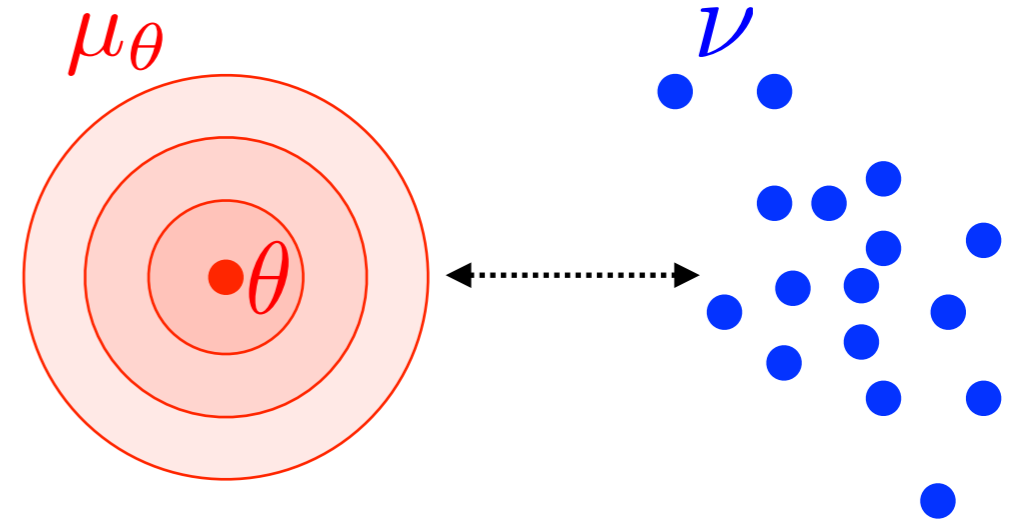
$$\min_{\theta} \widehat{\text{KL}}(\mu_\theta | \nu) \stackrel{\text{def.}}{=} - \sum_j \log(f_\theta(y_j))$$

Maximum likelihood (MLE)

# Density Fitting and Generative Models

Observations:  $\nu = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$

Parametric model:  $\theta \mapsto \mu_\theta$



Density fitting:  $d\mu_\theta(y) = f_\theta(y)dy$

$$\min_{\theta} \widehat{\text{KL}}(\mu_\theta | \nu) \stackrel{\text{def.}}{=} - \sum_j \log(f_\theta(y_j))$$

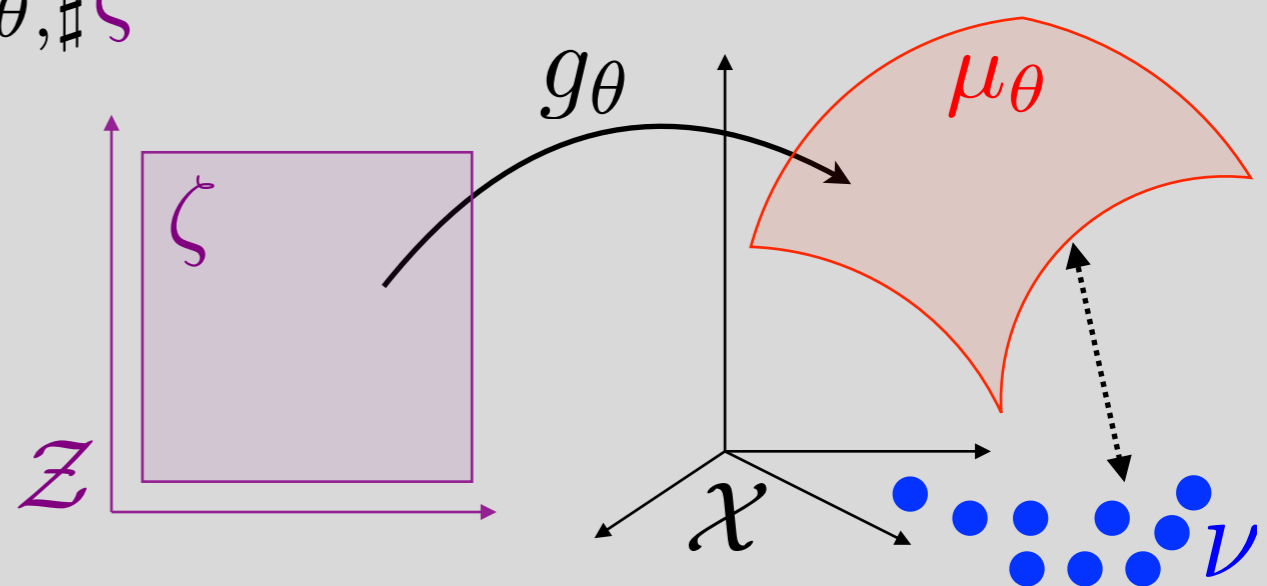
Maximum likelihood (MLE)

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

$$\widehat{\text{KL}}(\mu_\theta | \nu) = +\infty$$

→ MLE undefined.

→ Need a weaker metric.



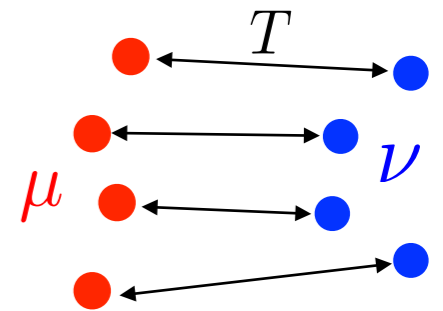
# Loss Functions for Measures

Density fitting:  $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

## Optimal Transport Distances

$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



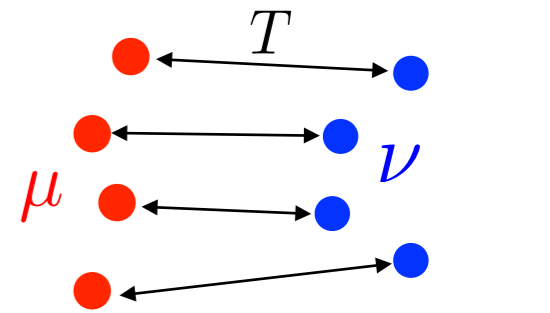
# Loss Functions for Measures

Density fitting:  $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

## Optimal Transport Distances

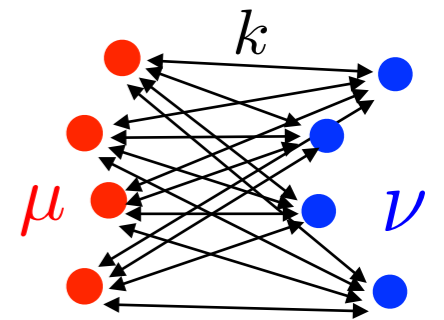
$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



## Maximum Mean Discrepancy (MMD)

$$\|\mu - \nu\|_k^2 \stackrel{\text{def.}}{=} \frac{1}{N^2} \sum_{i,i'} k(x_i, x_{i'}) + \frac{1}{P^2} \sum_{j,j'} k(y_j, y_{j'}) - \frac{2}{NP} \sum_{i,j} k(x_i, y_j)$$

Gaussian:  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$ . Energy distance:  $k(x, y) = -\|x - y\|^2$ .



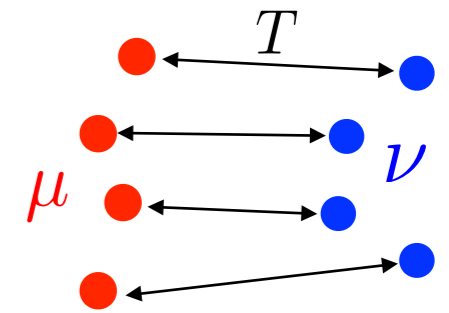
# Loss Functions for Measures

Density fitting:  $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

## Optimal Transport Distances

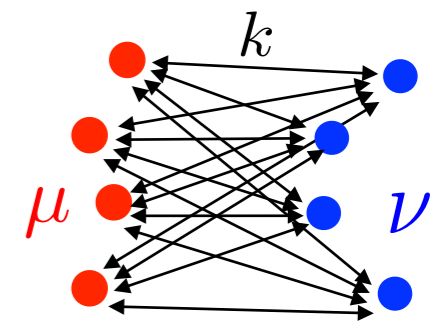
$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



## Maximum Mean Discrepancy (MMD)

$$\|\mu - \nu\|_k^2 \stackrel{\text{def.}}{=} \frac{1}{N^2} \sum_{i,i'} k(x_i, x_{i'}) + \frac{1}{P^2} \sum_{j,j'} k(y_j, y_{j'}) - \frac{2}{NP} \sum_{i,j} k(x_i, y_j)$$

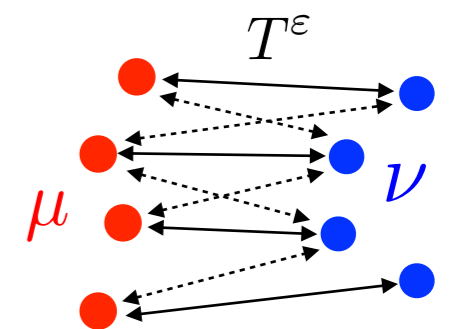
Gaussian:  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$ . Energy distance:  $k(x, y) = -\|x - y\|^2$ .



## Sinkhorn divergences [Cuturi 13]

$$W_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} \sum_{i,j} T_{i,j}^{\varepsilon} \|x_i - y_j\|^p$$

$$\bar{W}_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} W_{\varepsilon}(\mu, \nu)^p - \frac{1}{2} W_{\varepsilon}(\mu, \mu)^p - \frac{1}{2} W_{\varepsilon}(\nu, \nu)^p$$



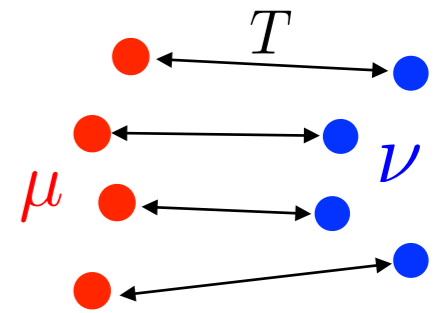
# Loss Functions for Measures

Density fitting:  $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

## Optimal Transport Distances

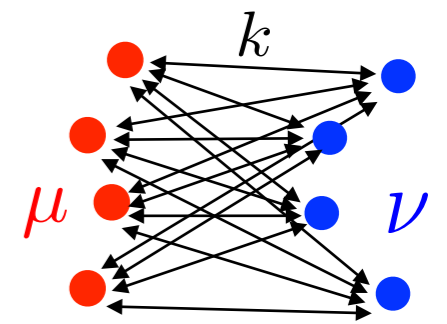
$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



## Maximum Mean Discrepancy (MMD)

$$\|\mu - \nu\|_k^2 \stackrel{\text{def.}}{=} \frac{1}{N^2} \sum_{i,i'} k(x_i, x_{i'}) + \frac{1}{P^2} \sum_{j,j'} k(y_j, y_{j'}) - \frac{2}{NP} \sum_{i,j} k(x_i, y_j)$$

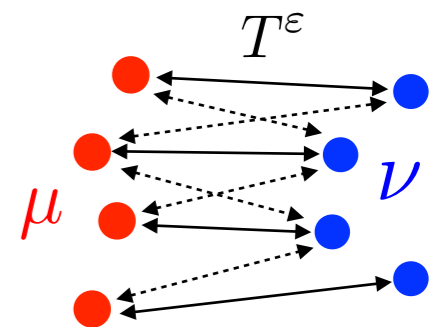
Gaussian:  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$ . Energy distance:  $k(x, y) = -\|x - y\|^2$ .



## Sinkhorn divergences [Cuturi 13]

$$W_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} \sum_{i,j} T_{i,j}^{\varepsilon} \|x_i - y_j\|^p$$

$$\bar{W}_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} W_{\varepsilon}(\mu, \nu)^p - \frac{1}{2} W_{\varepsilon}(\mu, \mu)^p - \frac{1}{2} W_{\varepsilon}(\nu, \nu)^p$$



*Theorem:* [Ramdas, G.Trillos, Cuturi 17]  $\bar{W}_{\varepsilon}(\mu, \nu)^p \xrightarrow[\varepsilon \rightarrow +\infty]{\varepsilon \rightarrow 0} W(\mu, \nu)^p$   
 $\xrightarrow[\varepsilon \rightarrow +\infty]{\varepsilon \rightarrow 0} \|\mu - \nu\|_k^2$

for  $k(x, y) = -\|x - y\|^p$

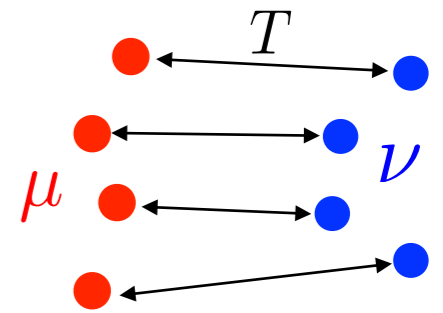
# Loss Functions for Measures

Density fitting:  $\min_{\theta} D(\mu_{\theta}, \nu)$

$$\nu = \frac{1}{P} \sum_j \delta_{y_j} \quad \mu = \frac{1}{N} \sum_i \delta_{x_i}$$

## Optimal Transport Distances

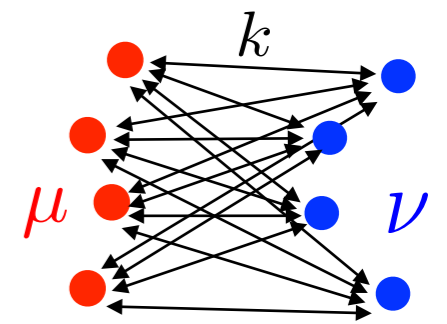
$$W(\mu, \nu)^p \stackrel{\text{def.}}{=} \min_{T \in \mathcal{C}_{\mu, \nu}} \sum_{i,j} T_{i,j} \|x_i - y_j\|^p$$



## Maximum Mean Discrepancy (MMD)

$$\|\mu - \nu\|_k^2 \stackrel{\text{def.}}{=} \frac{1}{N^2} \sum_{i,i'} k(x_i, x_{i'}) + \frac{1}{P^2} \sum_{j,j'} k(y_j, y_{j'}) - \frac{2}{NP} \sum_{i,j} k(x_i, y_j)$$

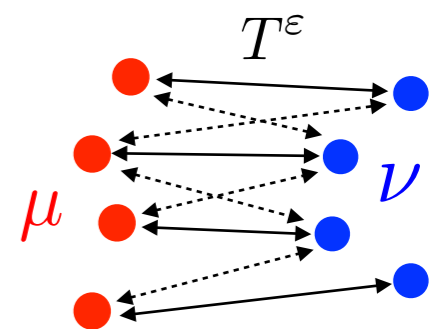
Gaussian:  $k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$ . Energy distance:  $k(x, y) = -\|x - y\|^2$ .



## Sinkhorn divergences [Cuturi 13]

$$W_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} \sum_{i,j} T_{i,j}^{\varepsilon} \|x_i - y_j\|^p$$

$$\bar{W}_{\varepsilon}(\mu, \nu)^p \stackrel{\text{def.}}{=} W_{\varepsilon}(\mu, \nu)^p - \frac{1}{2} W_{\varepsilon}(\mu, \mu)^p - \frac{1}{2} W_{\varepsilon}(\nu, \nu)^p$$



*Theorem:* [Ramdas, G.Trillos, Cuturi 17]  $\bar{W}_{\varepsilon}(\mu, \nu)^p \underset{\varepsilon \rightarrow +\infty}{\xrightarrow{\varepsilon \rightarrow 0}} W(\mu, \nu)^p \underset{\varepsilon \rightarrow +\infty}{\xrightarrow{\varepsilon \rightarrow 0}} \|\mu - \nu\|_k^2$

for  $k(x, y) = -\|x - y\|^p$

## Best of both worlds:

→ cross-validate  $\varepsilon$

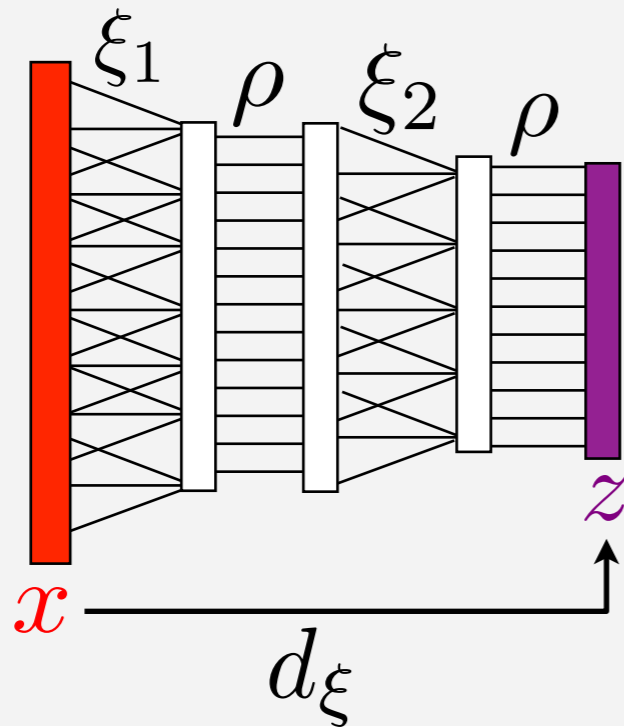
- Scale free (no  $\sigma$ , no heavy tail kernel).
- Non-Euclidean, arbitrary ground distance.
- Less biased gradient.
- No curse of dimension (low sample complexity).

# Deep Discriminative vs Generative Models

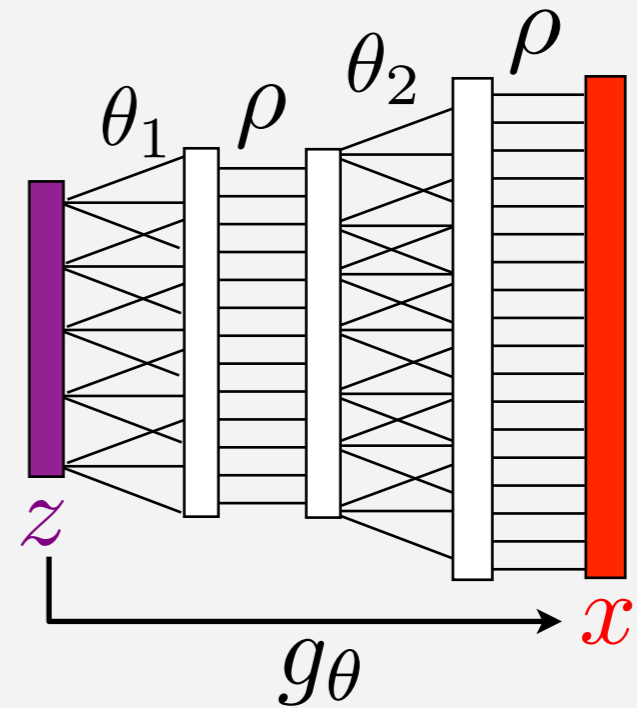
Deep networks:

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

Discriminative



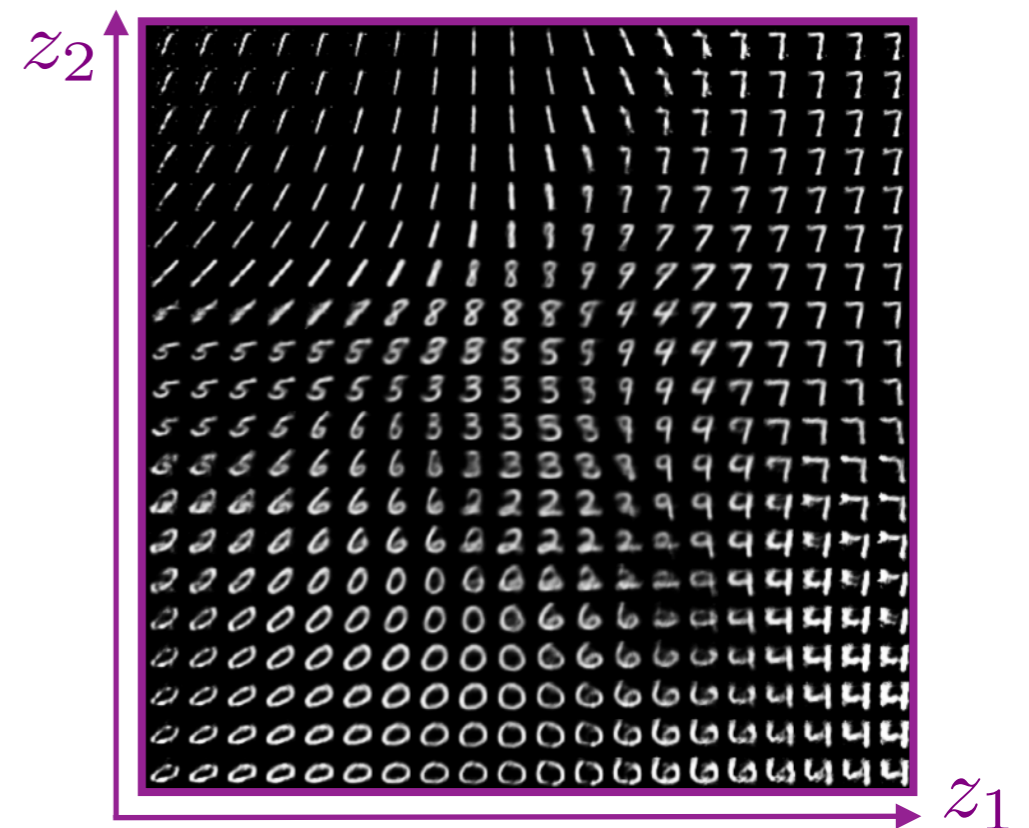
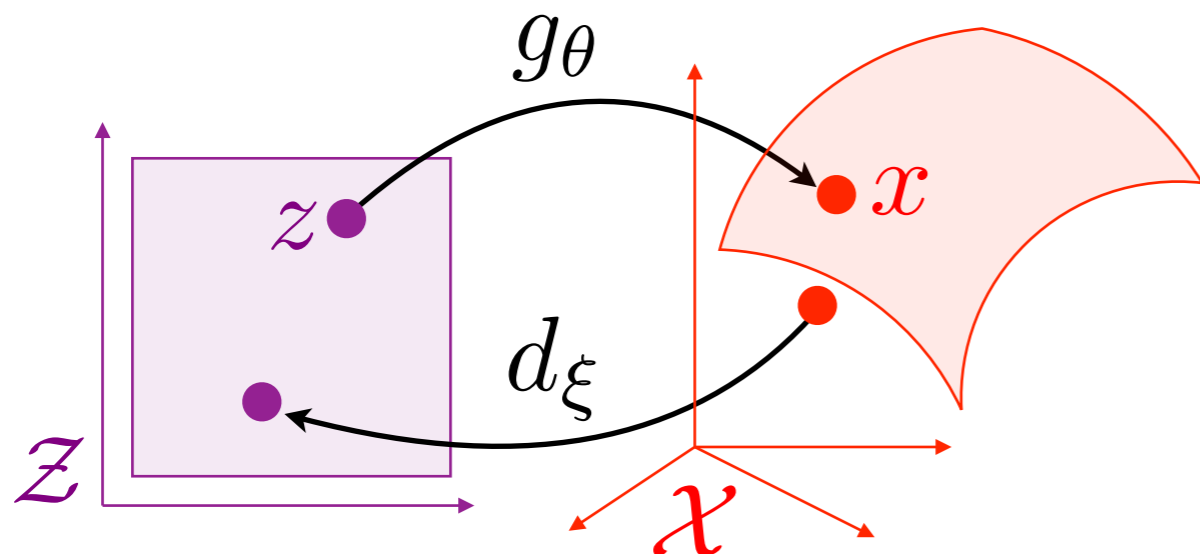
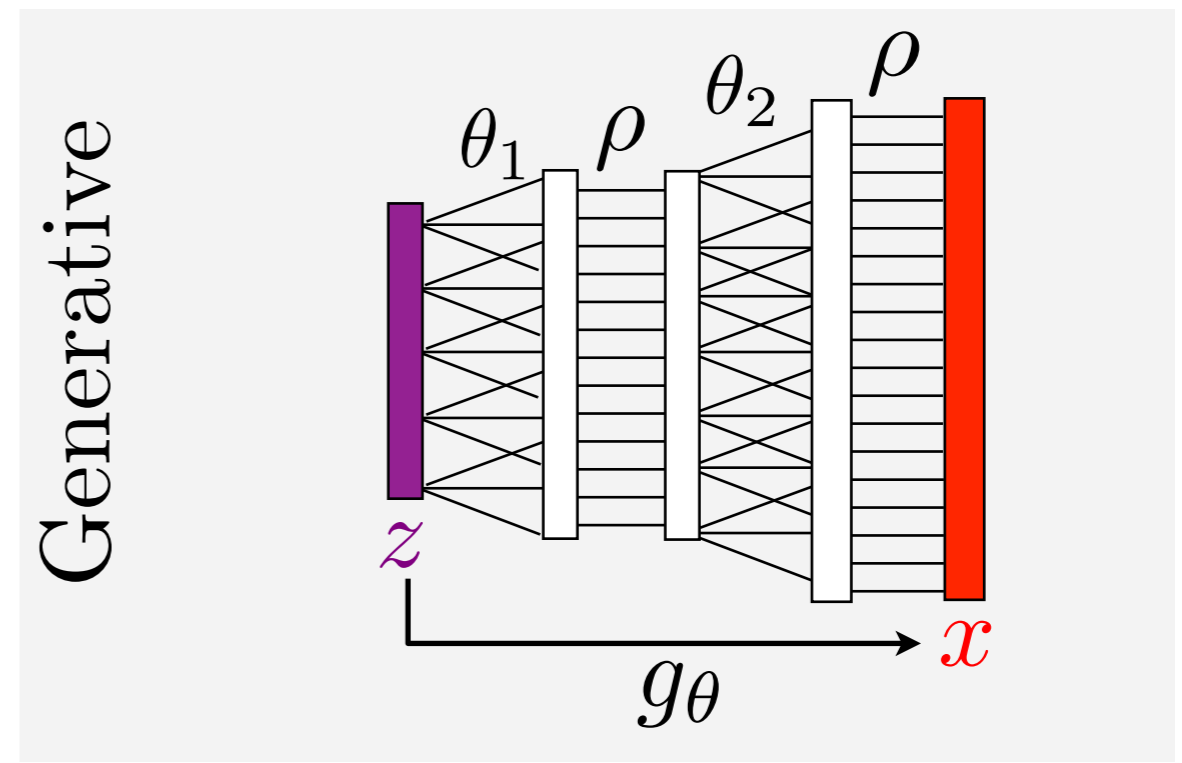
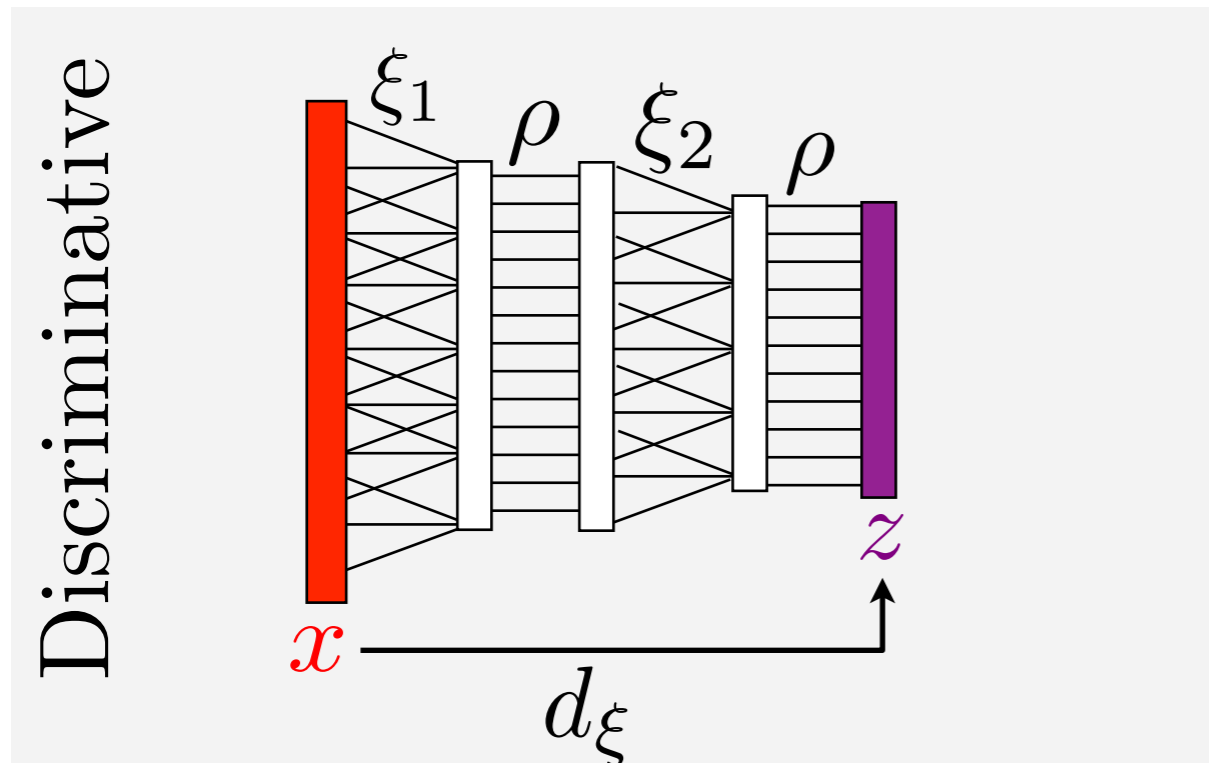
Generative



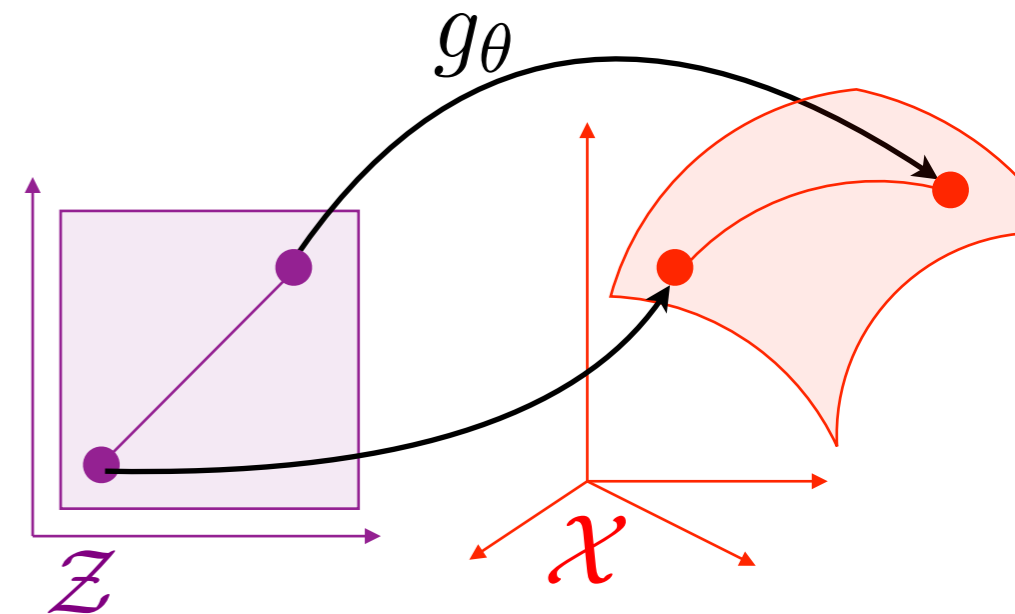


# Deep Discriminative vs Generative Models

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



# Examples of Image Generation



[Credit ArXiv:1511.06434]



# Conclusion: Toward High-dimensional OT

Monge

Kantorovich

Dantzig

Brenier

Otto

McCann

Villani

