

A Review of Regularized Optimal Transport

Marco Cuturi



Joint work with many people, including:

G. Peyré, A. Genevay (*ENS*), A. Doucet (*Oxford*) J. Solomon (*MIT*),
J.D. Benamou, N. Bonneel, F. Bach, L. Nenna (*INRIA*),
G. Carlier (*Dauphine*).

What is Optimal Transport?

A geometric toolbox to
compare probability measures
supported on a metric space.



Monge



Kantorovich



Dantzig



Wasserstein



Brenier



Otto



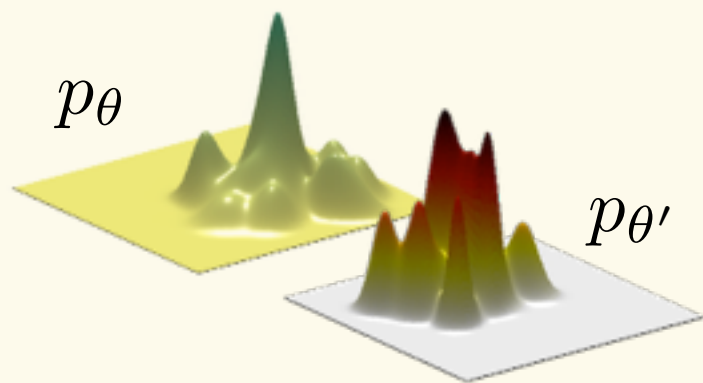
McCann



Villani

What is Optimal Transport?

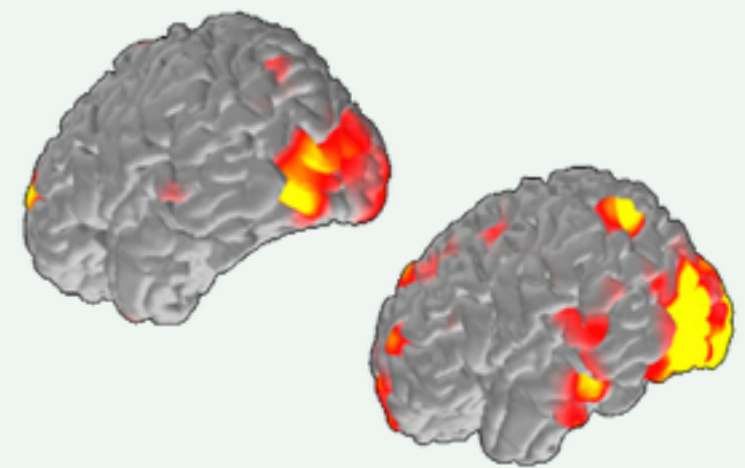
A geometric toolbox to
compare **probability measures**
supported on a metric space.



Statistical Models

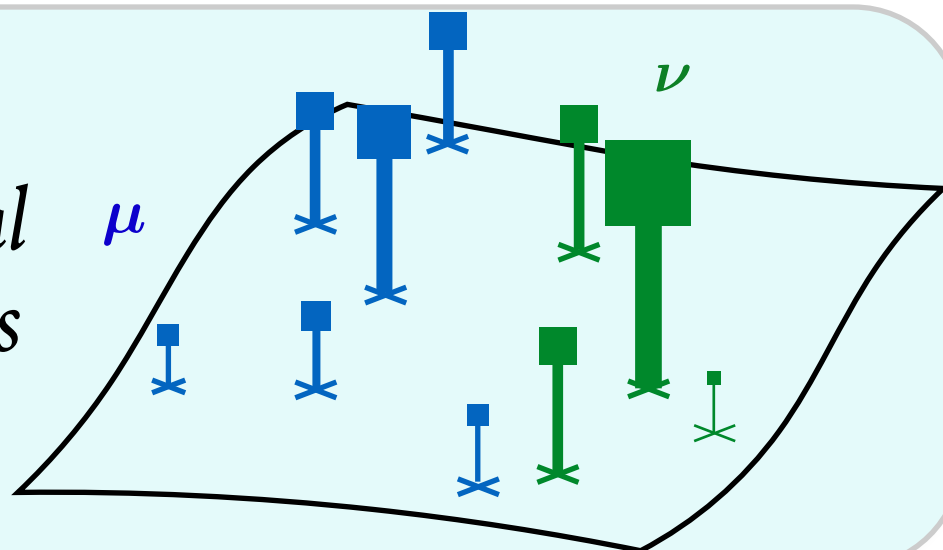


*Bags
of features*



Brain Activation Maps

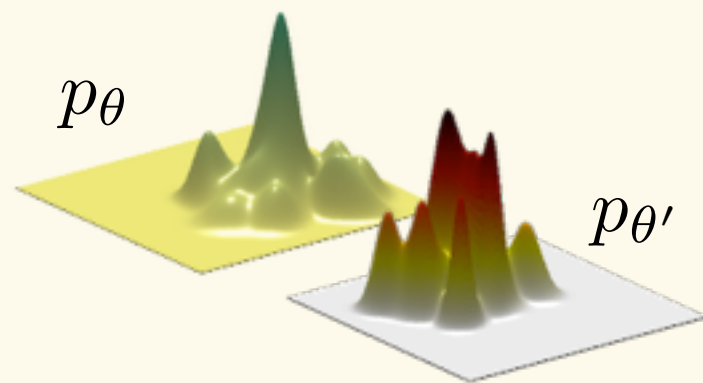
*Empirical
Measures*



Color Histograms

What is Optimal Transport?

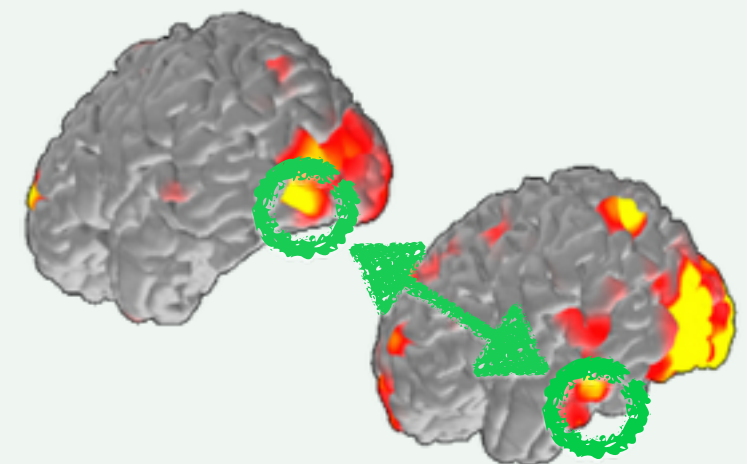
A geometric toolbox to
compare probability measures
supported on a metric space.



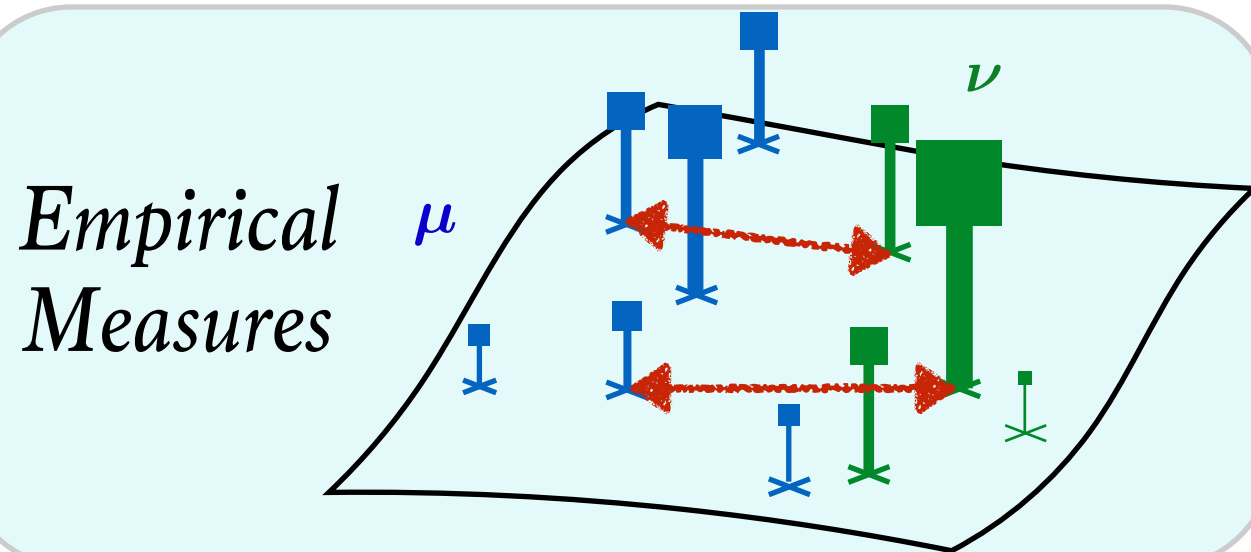
Statistical Models



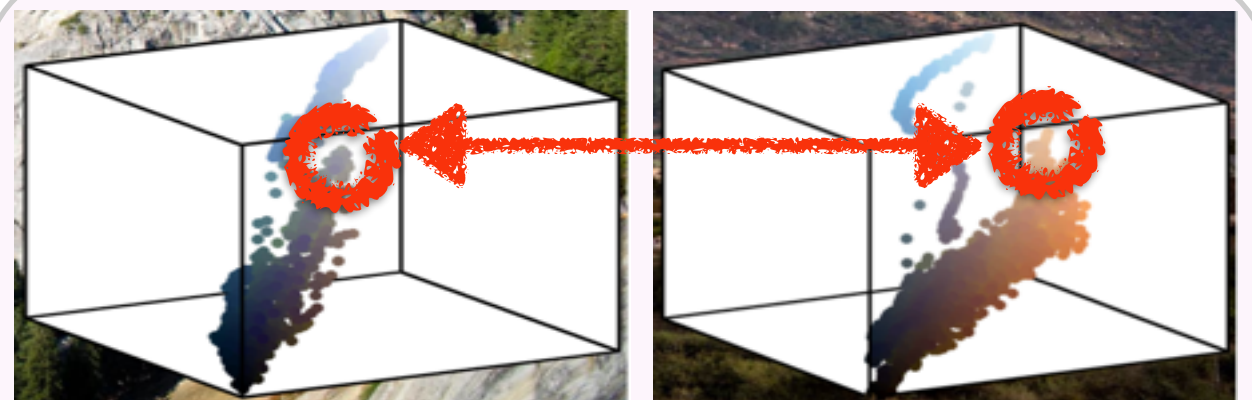
Bags of features



Brain Activation Maps



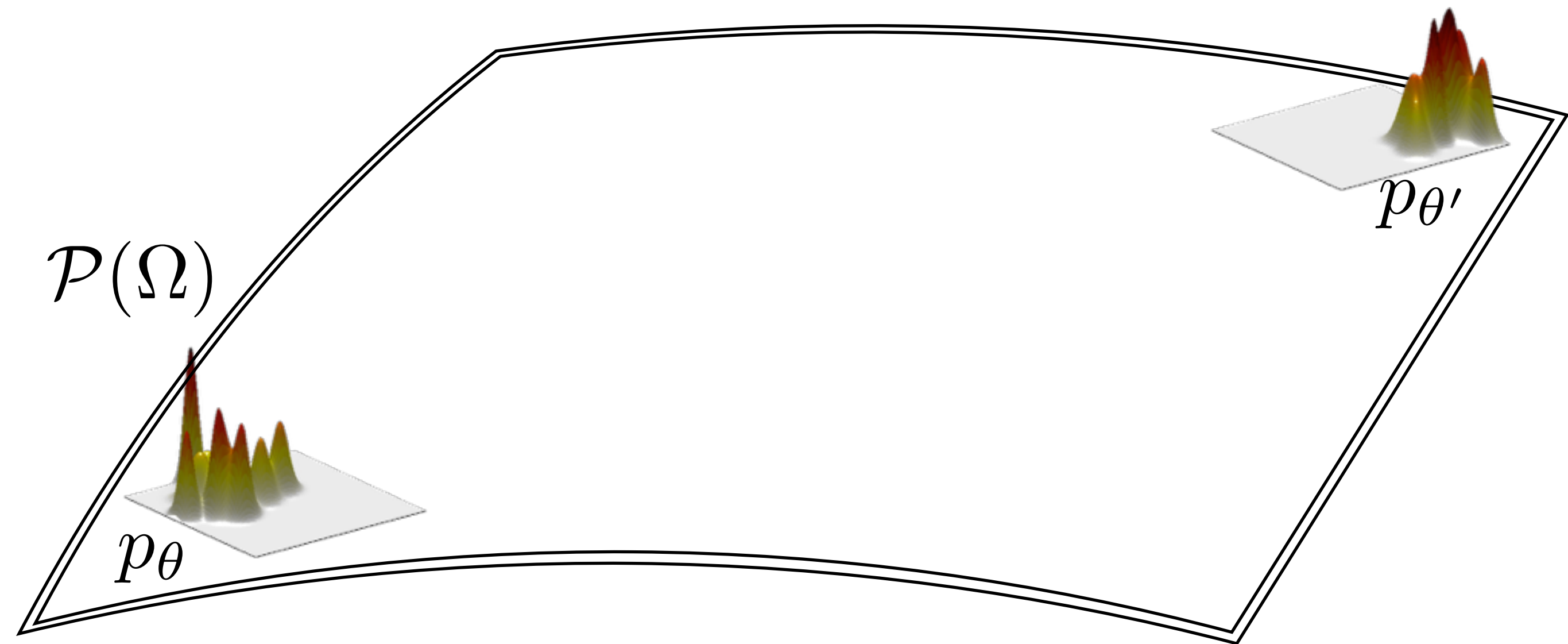
Empirical Measures



Color Histograms

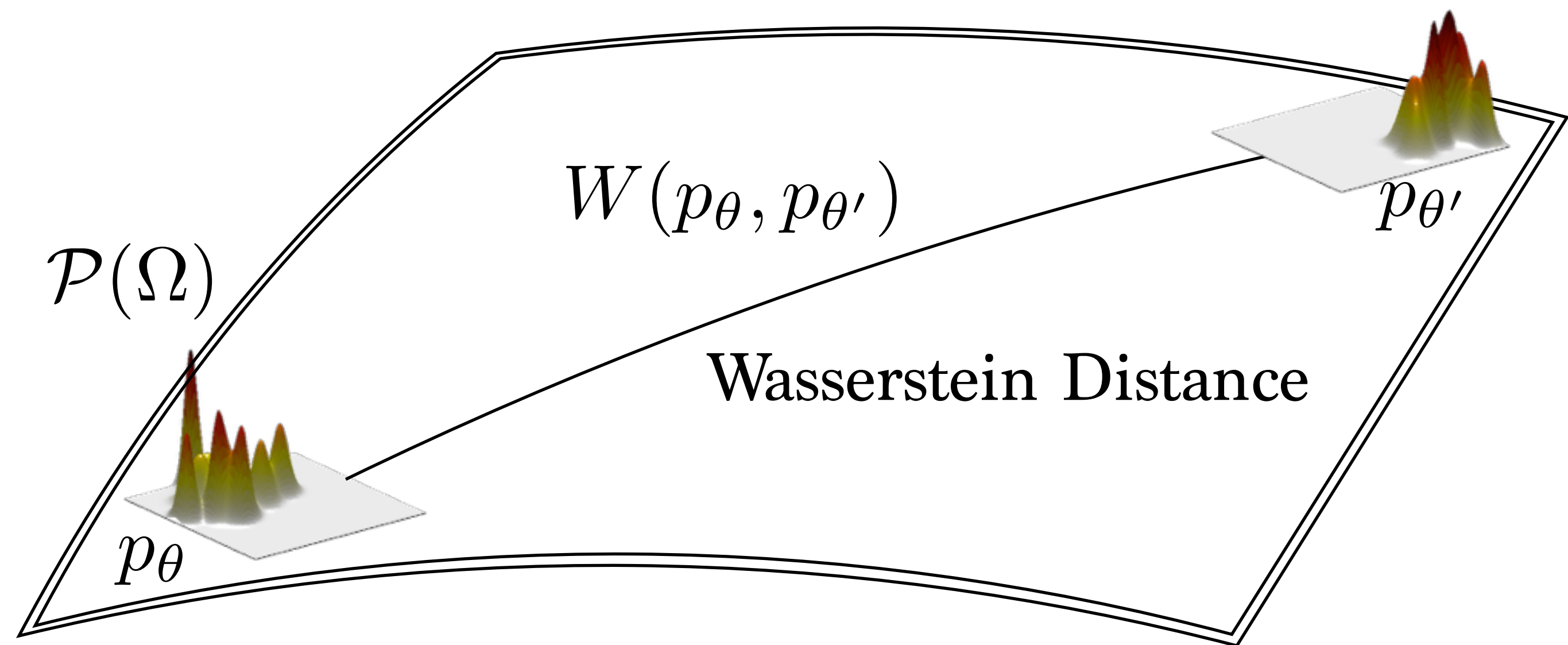
What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



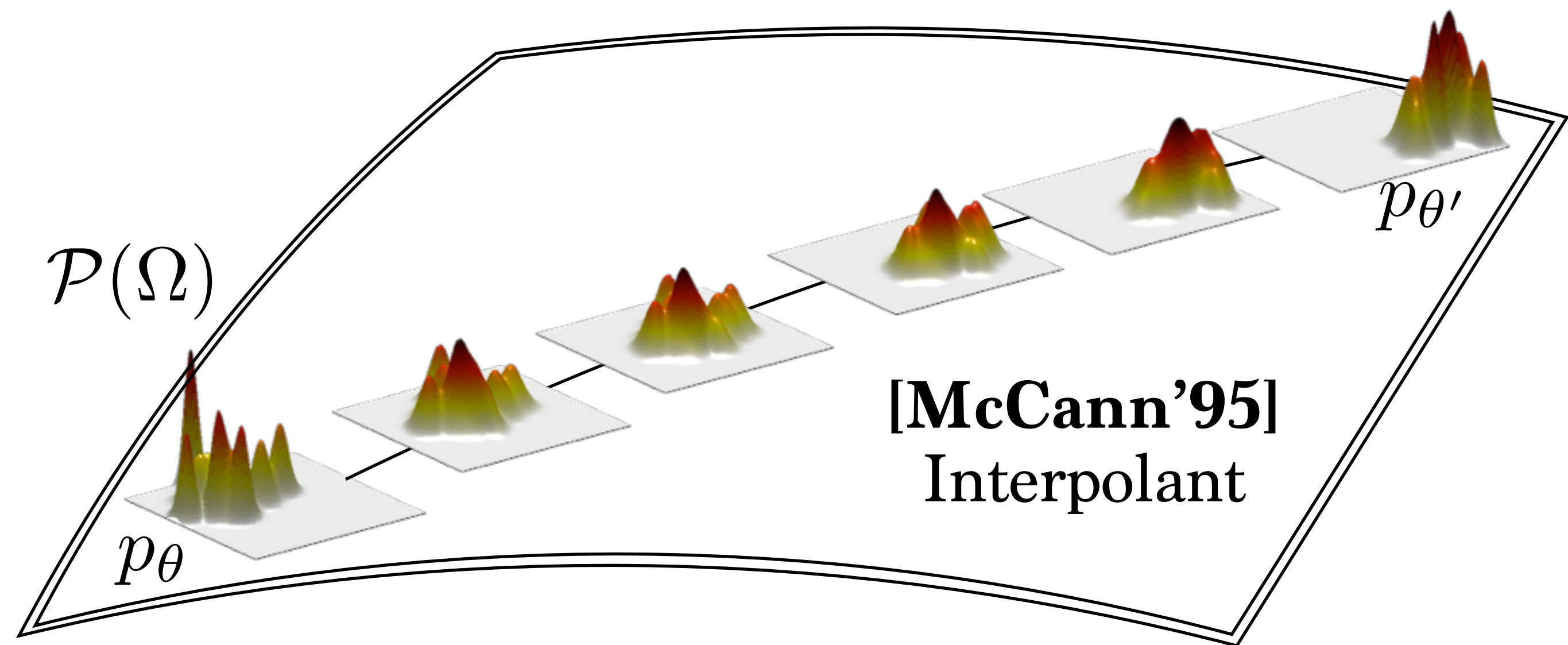
What is Optimal Transport?

A **geometric toolbox** to compare probability measures supported on a metric space.



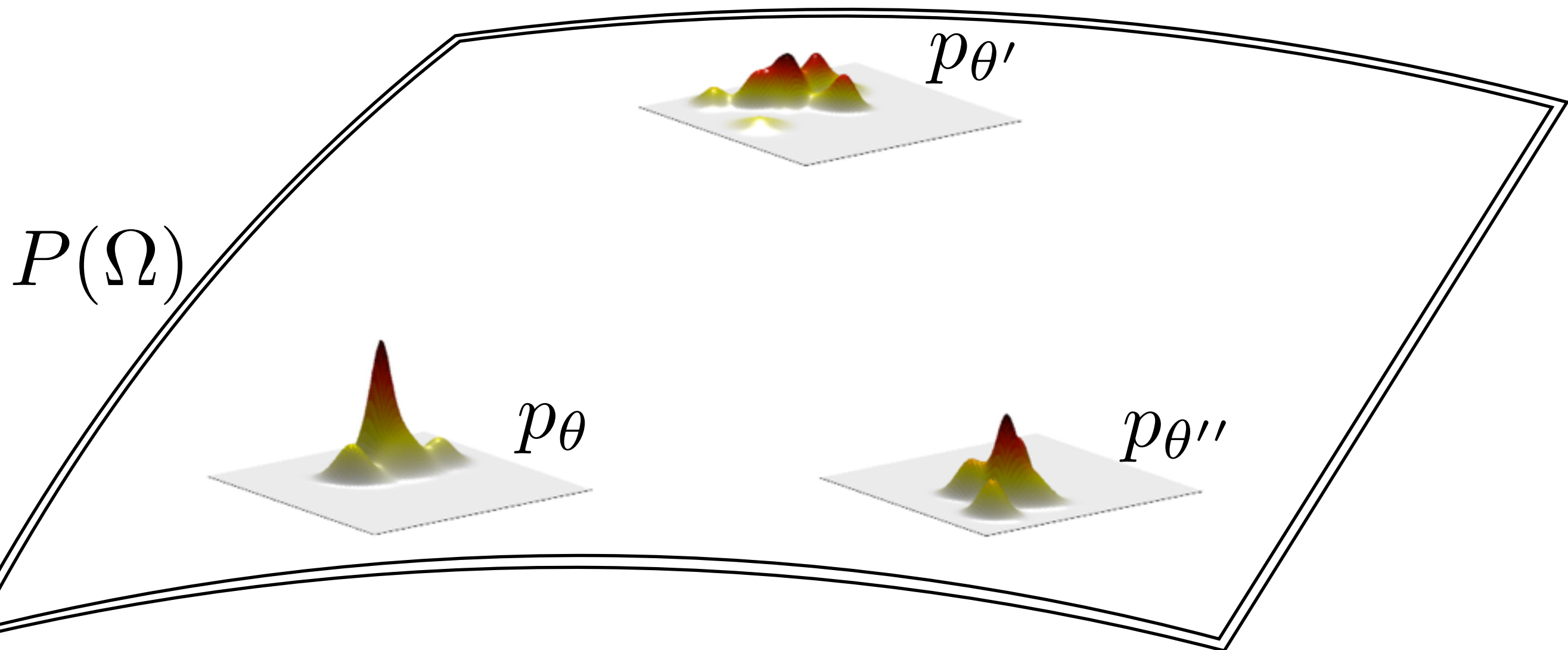
What is Optimal Transport?

A **geometric toolbox** to compare probability measures supported on a metric space.



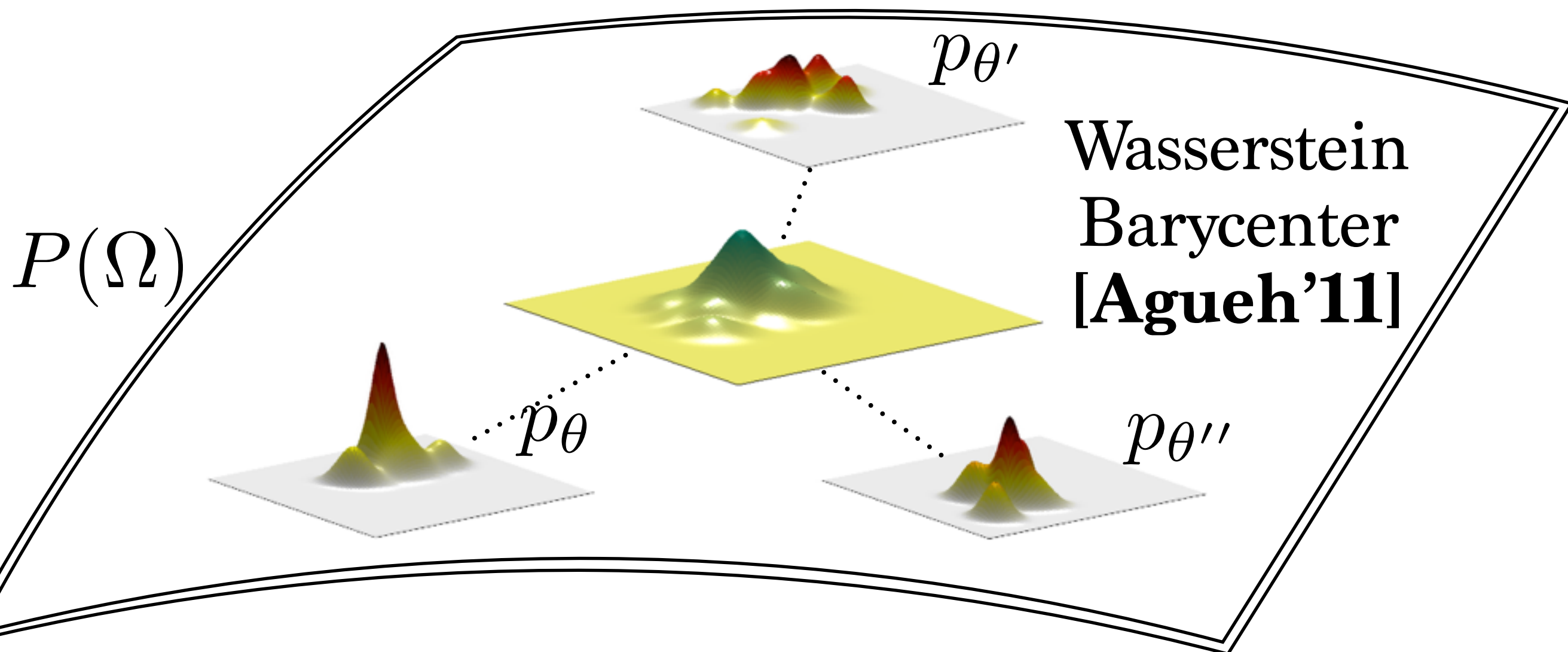
What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



What is Optimal Transport?

A **geometric toolbox** to
compare probability measures
supported on a metric space.



OT and data-analysis

- Key developments in (applied) maths ~'90s
[McCann'95], [JKO'98], [Benamou'98], [Gangbo'98],
[Ambrosio'06], [Villani'03/'09].
 - Key developments in TCS / graphics since '00s
[Rubner'98], [Indyk'03], [Naor'07], [Andoni'15].
- ◎ Small to *no-impact* in large-scale data analysis:
- ♦ computationally heavy;
 - ♦ Wasserstein distance is not differentiable

OT and data-analysis

Today's talk: *Entropy Regularized OT*

- **Very fast** compared to usual approaches, GPGPU parallel.
- **Differentiable**, important if we want to use OT distances as **loss functions**.
- Can be **automatically differentiated**, simple iterative process, *DL*-toolboxes compatible.
- OT can become a building block in ML.

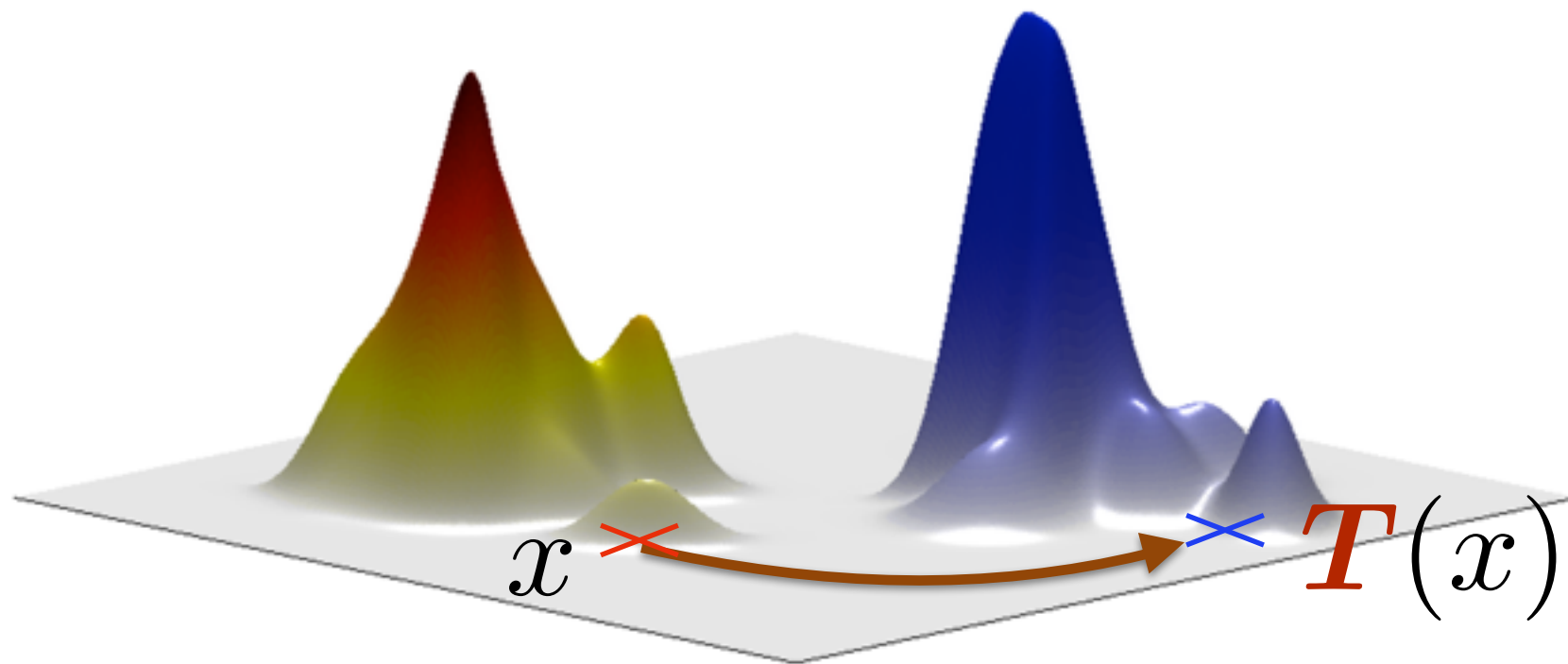
✦ Wasserstein distance is not differentiable

Background: OT Geometry

Consider (Ω, \mathbf{D}) , a metric probability space. Let μ, ν be probability measures in $\mathcal{P}(\Omega)$.

- [Monge'81] problem: find a map $T : \Omega \rightarrow \Omega$

$$\inf_{T \# \mu = \nu} \int_{\Omega} \mathbf{D}(x, T(x)) \mu(dx)$$

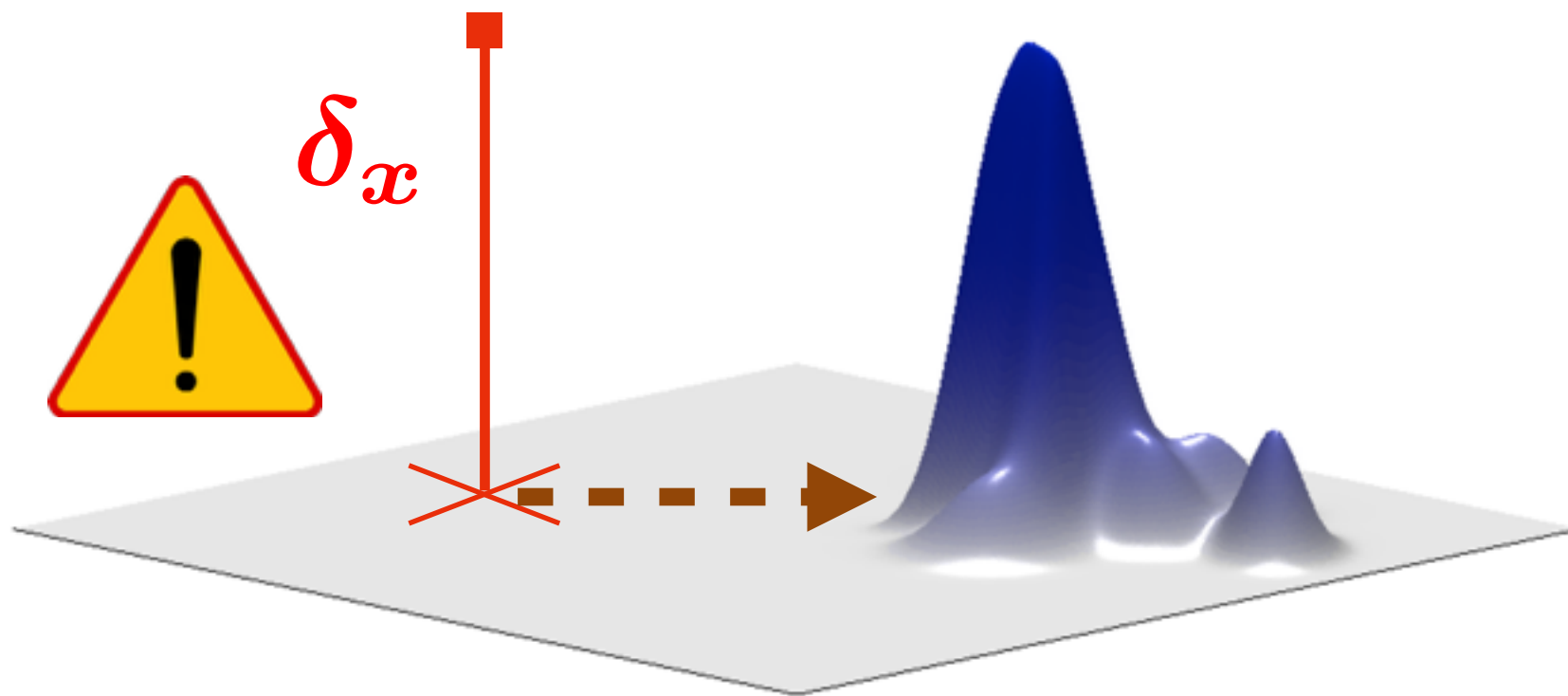


Background: OT Geometry

Consider (Ω, \mathbf{D}) , a metric probability space. Let μ, ν be probability measures in $\mathcal{P}(\Omega)$.

- [Monge'81] problem: find a map $T : \Omega \rightarrow \Omega$

$$\inf_{T \# \mu = \nu} \int_{\Omega} \mathbf{D}(x, T(x)) \mu(dx)$$



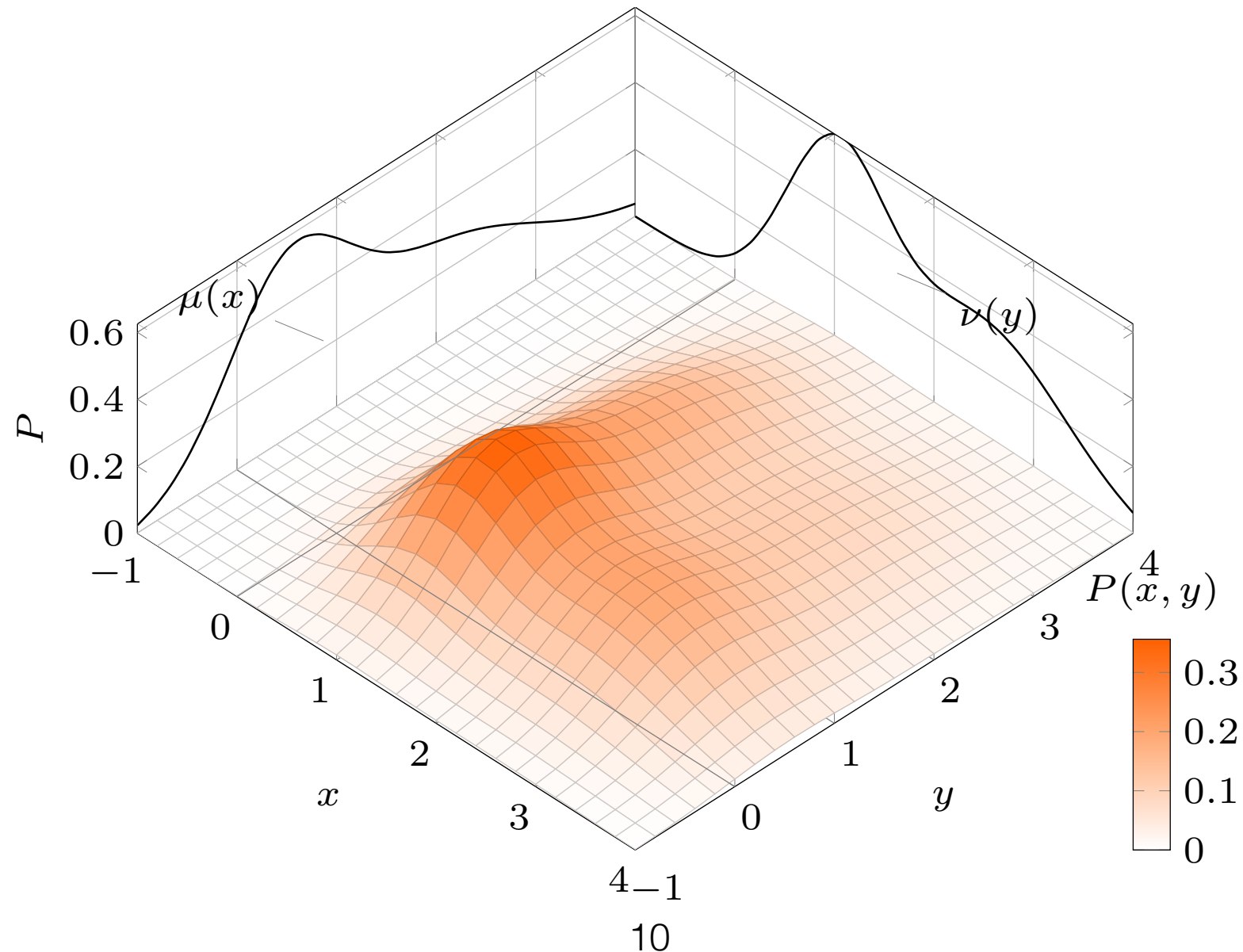
[Kantorovich'42] Relaxation

- Instead of maps $T : \Omega \rightarrow \Omega$, consider probabilistic maps, i.e. **couplings** $P \in \mathcal{P}(\Omega \times \Omega)$:

$$\Pi(\mu, \nu) \stackrel{\text{def}}{=} \{P \in \mathcal{P}(\Omega \times \Omega) \mid \forall A, B \subset \Omega, \\ P(A \times \Omega) = \mu(A), \\ P(\Omega \times B) = \nu(B)\}$$

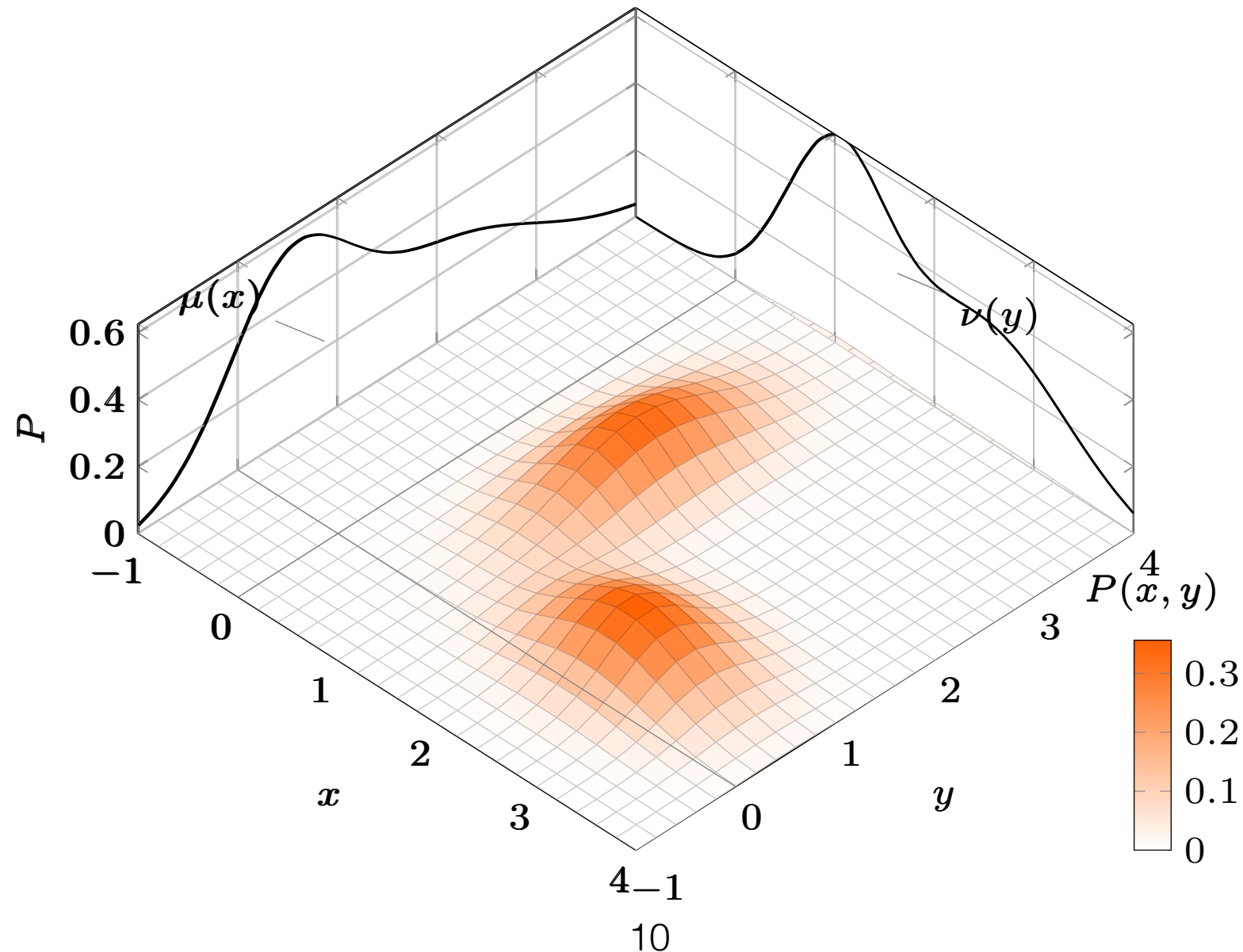
[Kantorovich'42] Relaxation

$$\Pi(\mu, \nu) \stackrel{\text{def}}{=} \{P \in \mathcal{P}(\Omega \times \Omega) \mid \forall A, B \subset \Omega, \\ P(A \times \Omega) = \mu(A), P(\Omega \times B) = \nu(B)\}$$



[Kantorovich'42] Relaxation

$$\Pi(\mu, \nu) \stackrel{\text{def}}{=} \{P \in \mathcal{P}(\Omega \times \Omega) \mid \forall A, B \subset \Omega, \\ P(A \times \Omega) = \mu(A), P(\Omega \times B) = \nu(B)\}$$

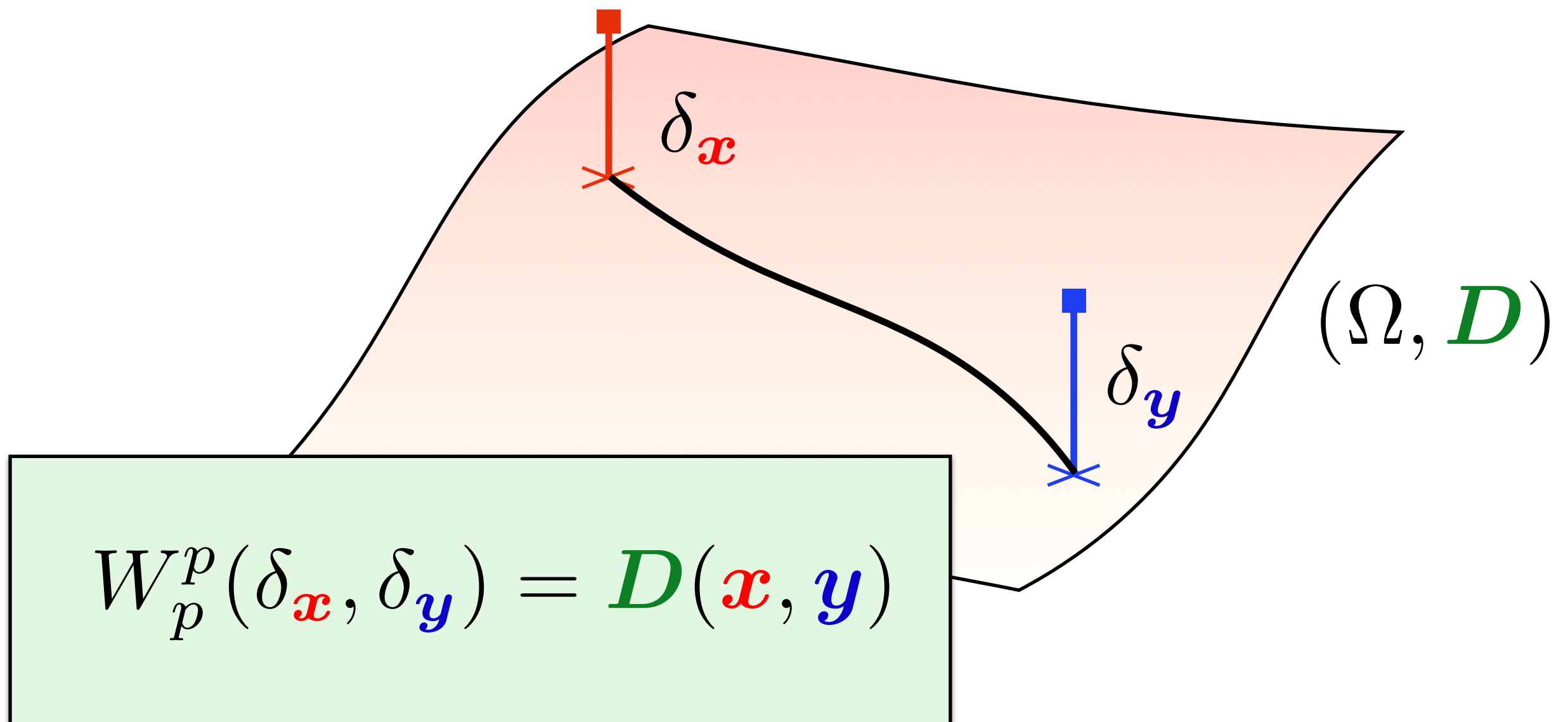


Wasserstein Distance

Def. For $p \geq 1$, the p -Wasserstein distance between μ, ν in $\mathcal{P}(\Omega)$ is

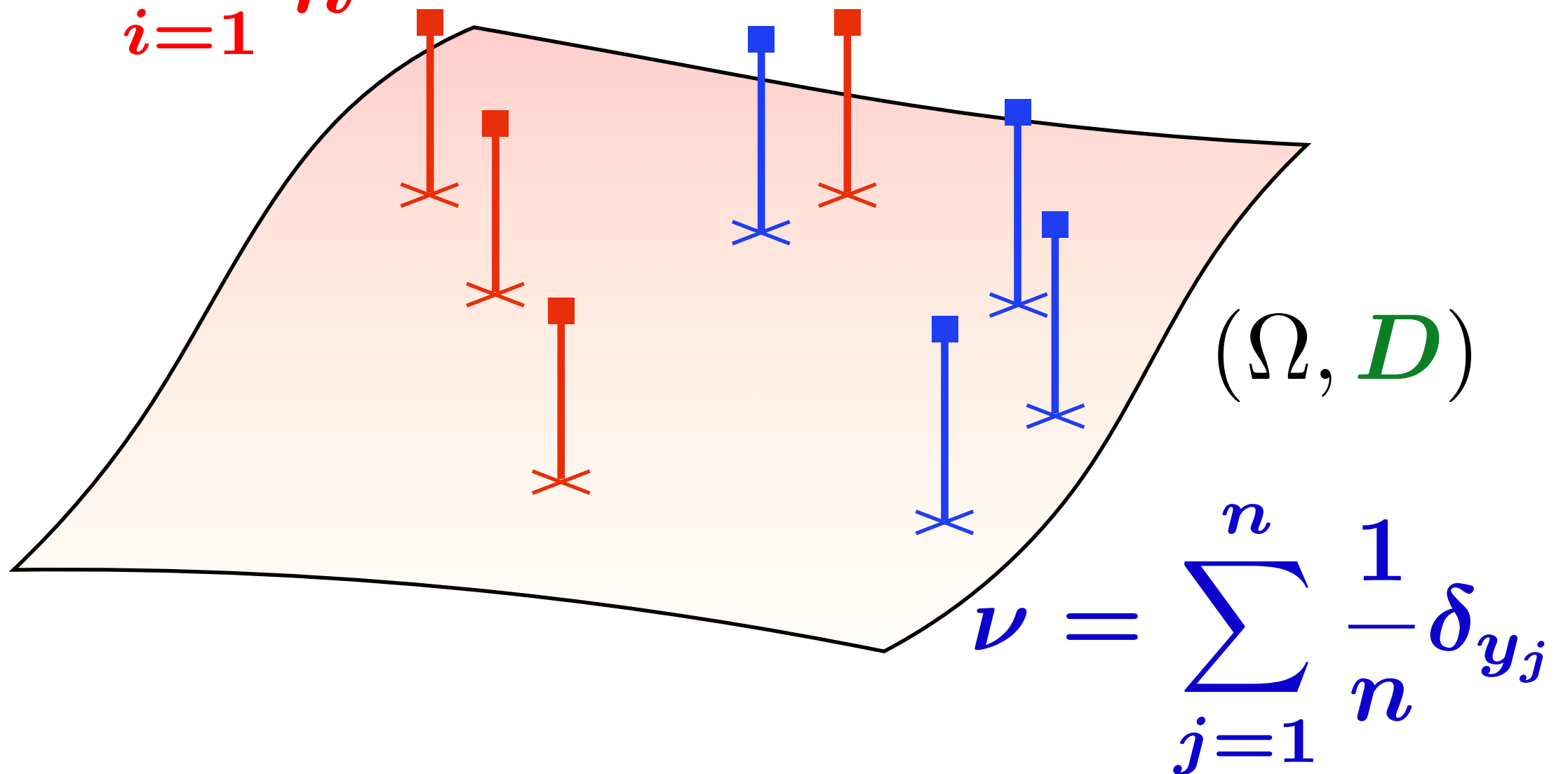
$$W_p(\mu, \nu) \stackrel{\text{def}}{=} \left(\inf_{P \in \Pi(\mu, \nu)} \mathbb{E}_P [D(X, Y)^p] \right)^{1/p}.$$

Wasserstein between 2 Diracs



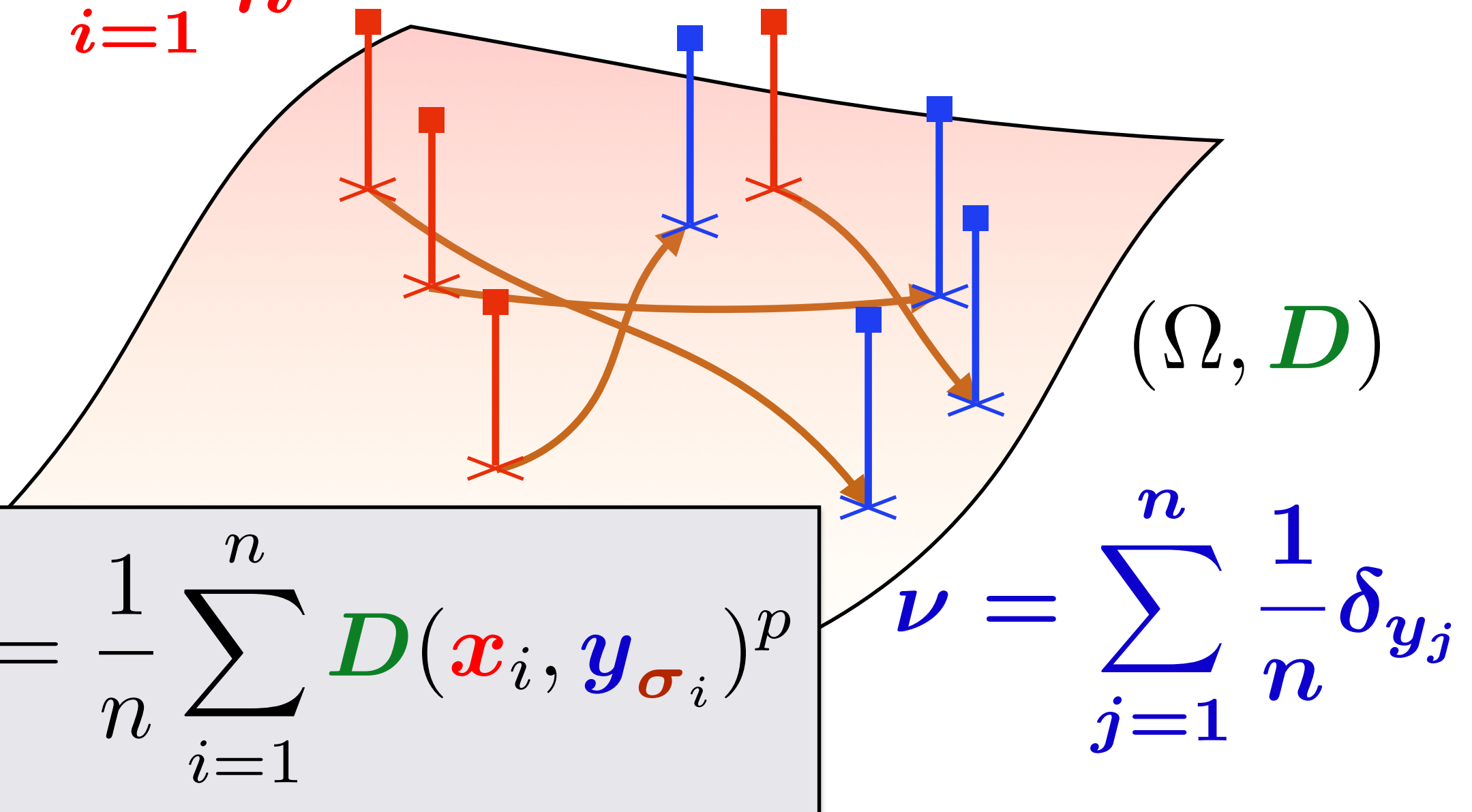
Wasserstein on Uniform Measures

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



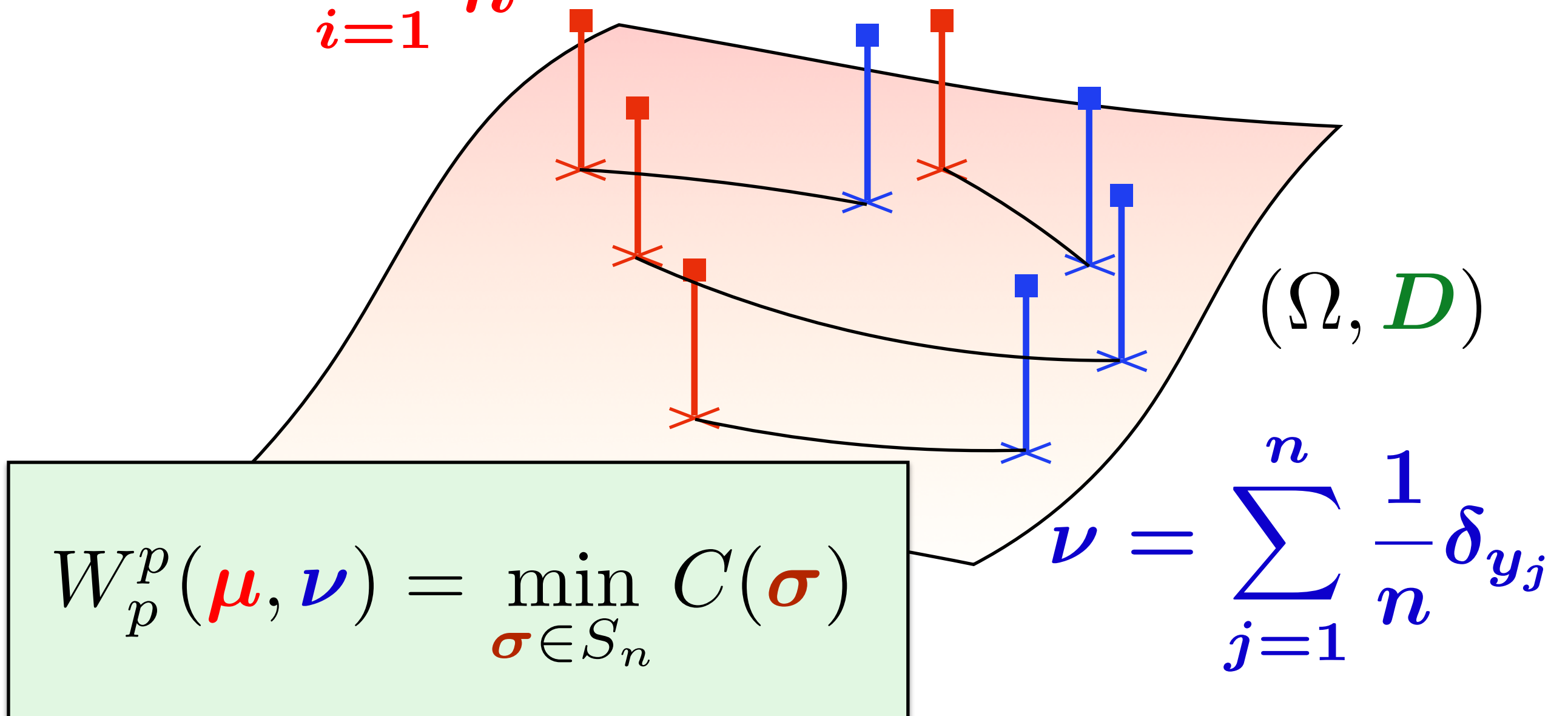
Wasserstein on Uniform Measures

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



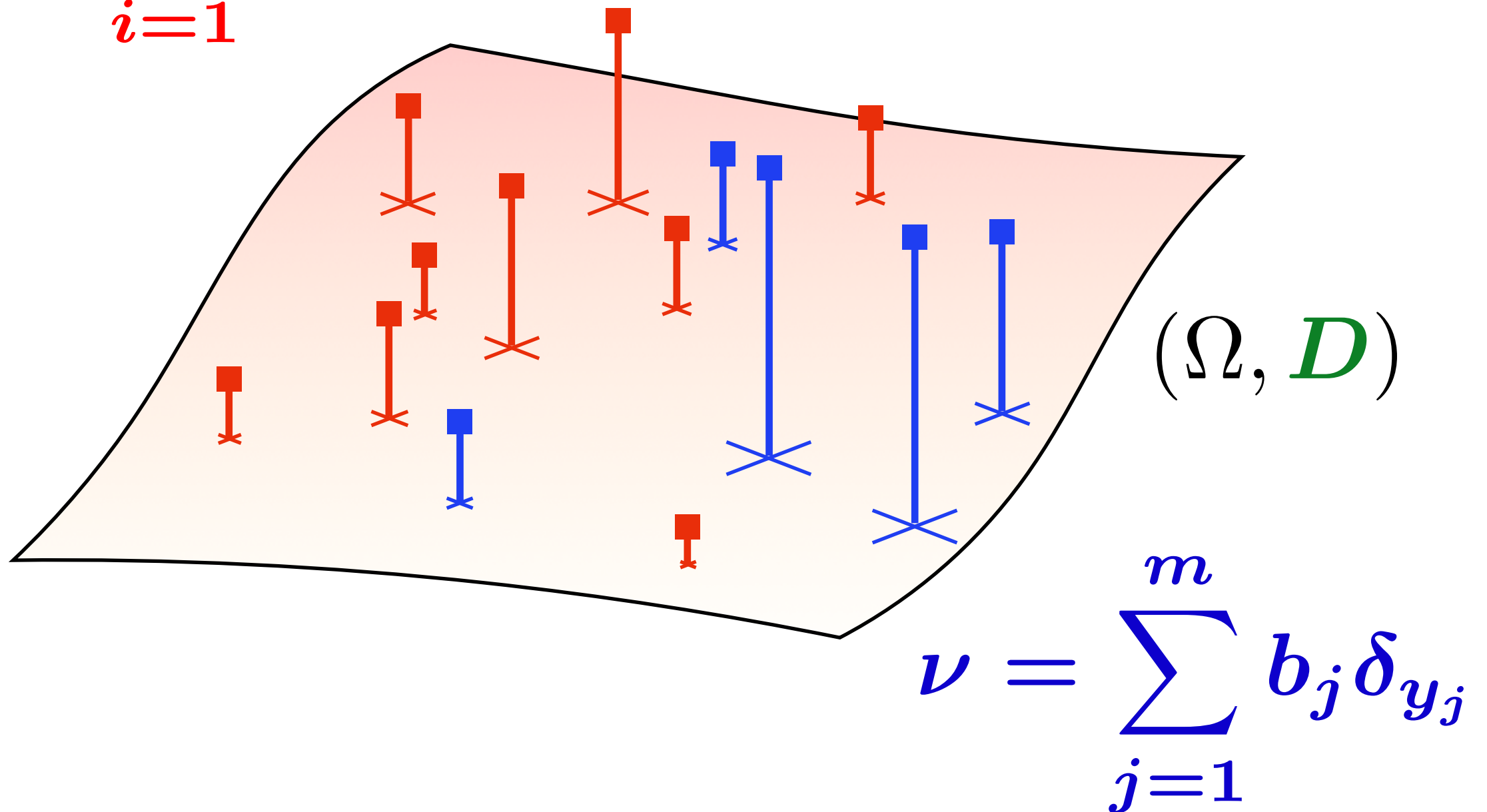
Optimal Assignment \subset Wasserstein

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



Wasserstein on Empirical Measures

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



Wasserstein on Empirical Measures

Consider $\mu = \sum_{i=1}^n a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$.

$$M_{\mathbf{x}\mathbf{y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij}$$

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

$$\begin{array}{c}
 \mathbf{x}_1 \\
 \vdots \\
 \mathbf{x}_n
 \end{array}
 \begin{array}{c}
 \mathbf{y}_1 \quad \dots \quad \mathbf{y}_m \\
 \left[\begin{array}{ccc}
 \cdot & \cdot & \cdot \\
 \cdot & D(\mathbf{x}_i, \mathbf{y}_j)^p & \cdot \\
 \cdot & \cdot & \cdot
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 \mathbf{a}_1 \\
 \vdots \\
 \mathbf{a}_n
 \end{array}
 \begin{array}{c}
 \mathbf{b}_1 \quad \dots \quad \mathbf{b}_m \\
 \left[\begin{array}{ccc}
 \cdot \cdot \cdot & \cdot \cdot \cdot & \cdot \cdot \cdot \\
 \cdot \cdot \cdot & \mathbf{P} \mathbf{1}_m = \mathbf{a} & \cdot \cdot \cdot \\
 \cdot \cdot \cdot & \cdot \cdot \cdot & \cdot \cdot \cdot
 \end{array} \right]
 \end{array}$$

Wasserstein on Empirical Measures

Consider $\mu = \sum_{i=1}^n a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$.

$$M_{\mathbf{x}\mathbf{y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij}$$

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

$$\begin{array}{c}
 \mathbf{x}_1 \\
 \vdots \\
 \mathbf{x}_n
 \end{array}
 \begin{array}{c}
 \mathbf{y}_1 \quad \dots \quad \mathbf{y}_m \\
 \left[\begin{array}{ccc}
 \cdot & & \cdot \\
 \cdot & D(\mathbf{x}_i, \mathbf{y}_j)^p & \cdot \\
 \cdot & & \cdot
 \end{array} \right]
 \end{array}
 \begin{array}{c}
 \mathbf{a}_1 \\
 \vdots \\
 \mathbf{a}_n
 \end{array}
 \begin{array}{c}
 b_1 \quad \dots \quad b_m \\
 \left[\begin{array}{ccc}
 \vdots & & \vdots \\
 \vdots & \mathbf{P}^T \mathbf{1}_n = \mathbf{b} & \vdots \\
 \vdots & & \vdots
 \end{array} \right]
 \end{array}$$

Wasserstein on Empirical Measures

Consider $\mu = \sum_{i=1}^n a_i \delta_{x_i}$ and $\nu = \sum_{j=1}^m b_j \delta_{y_j}$.

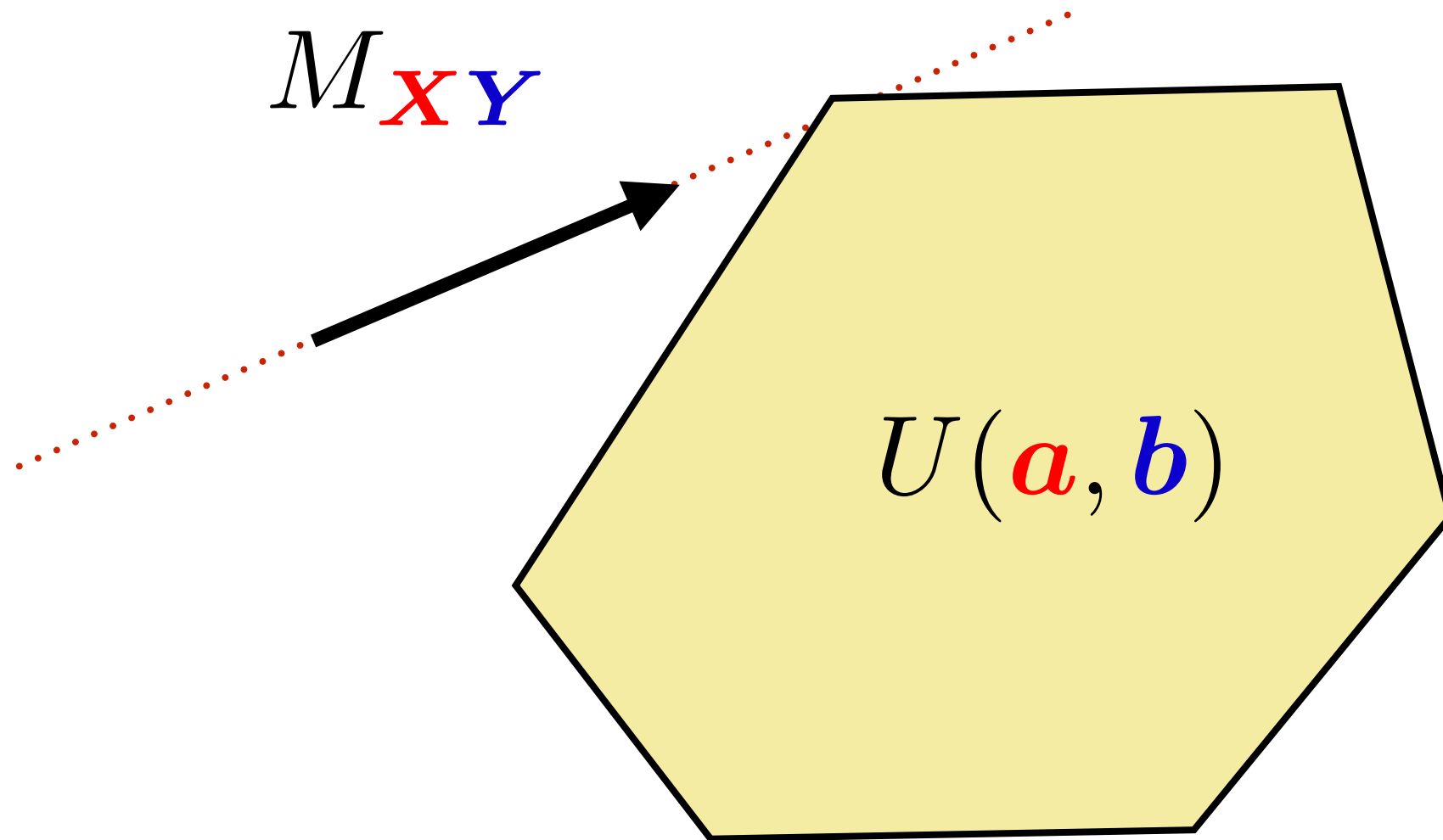
$$M_{\mathbf{x}\mathbf{y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij}$$

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

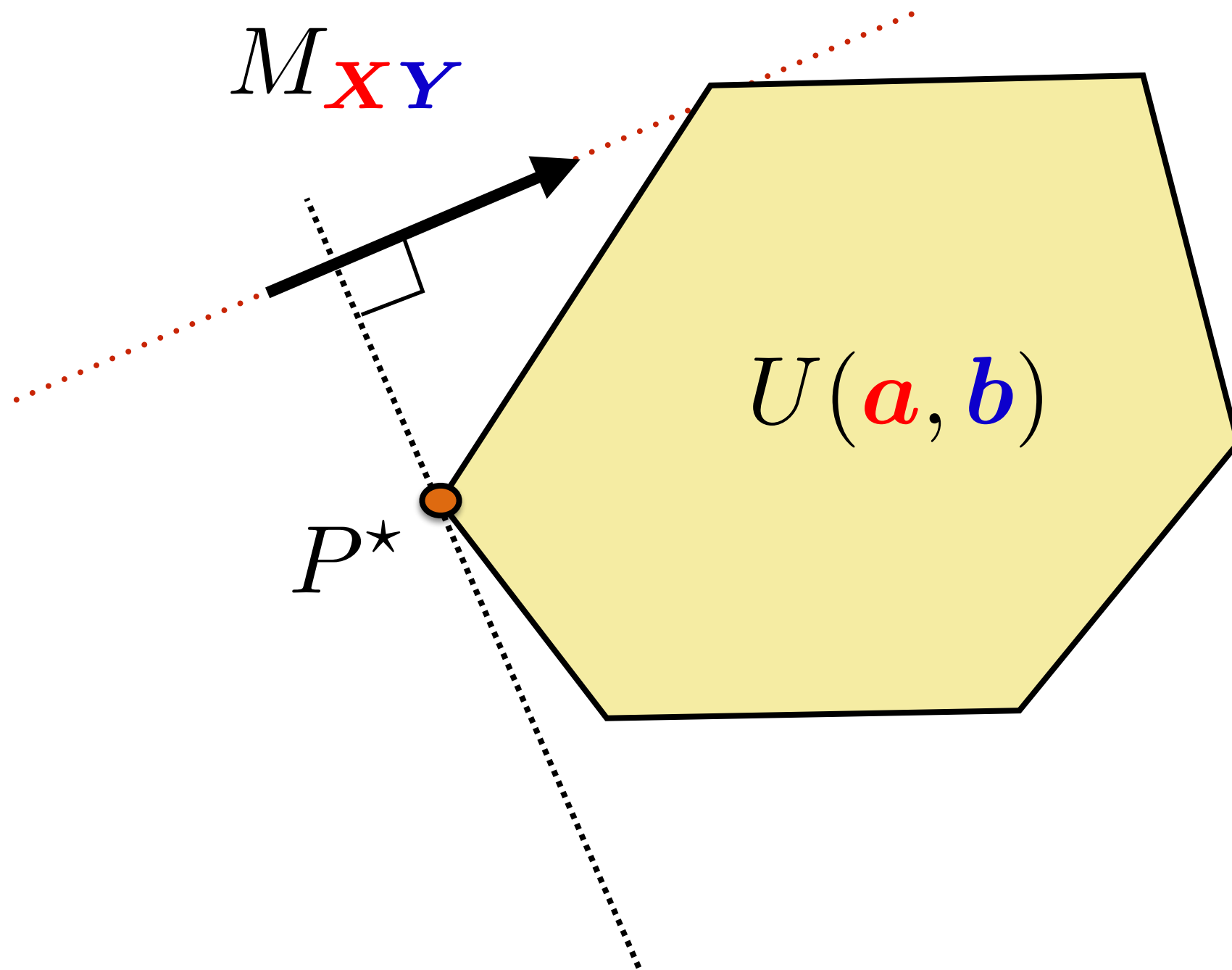
Def. Optimal Transport Problem

$$W_p^p(\mu, \nu) = \min_{\mathbf{P} \in U(\mathbf{a}, \mathbf{b})} \langle \mathbf{P}, M_{\mathbf{x}\mathbf{y}} \rangle$$

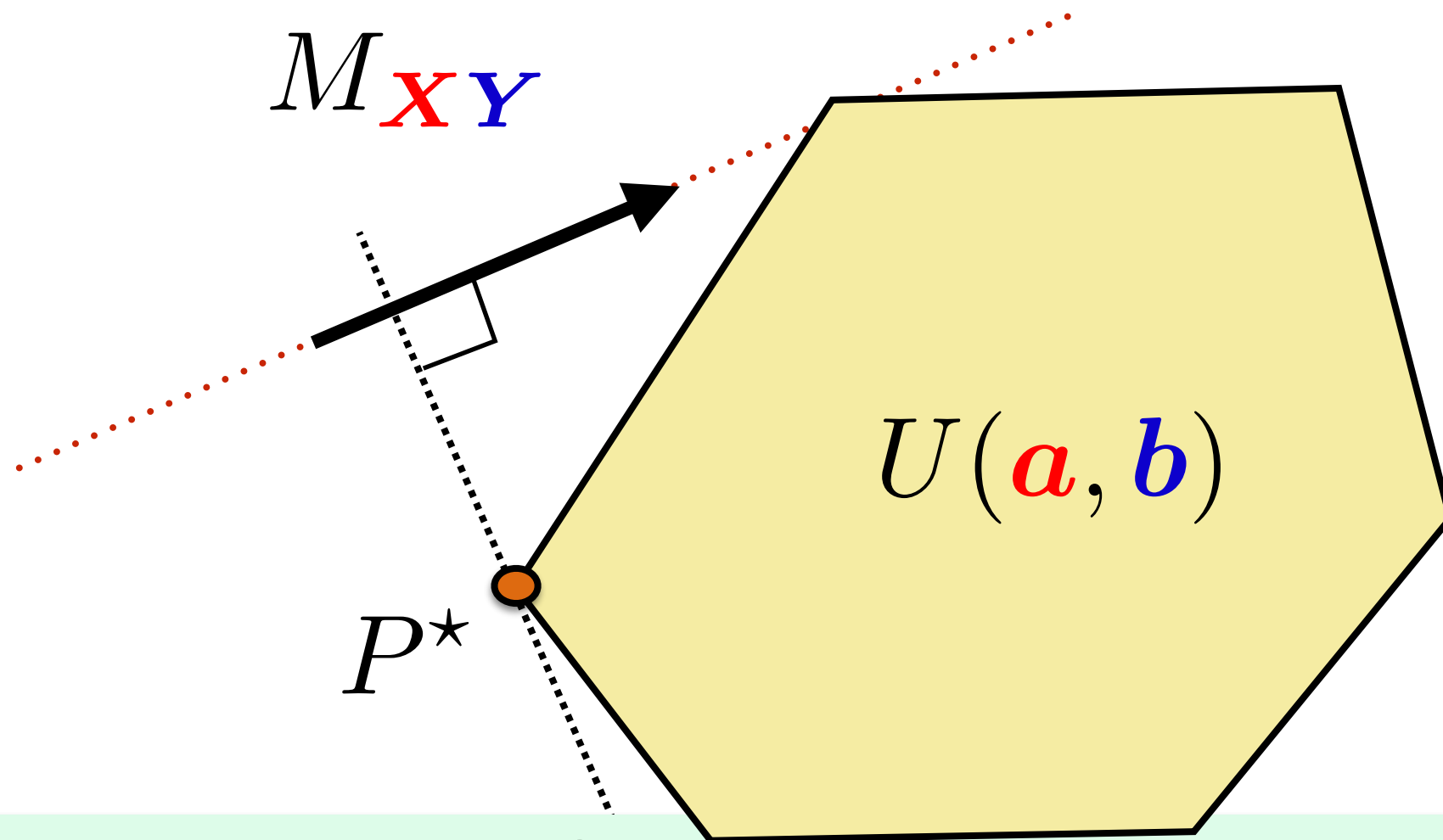
Discrete OT Problem



Discrete OT Problem



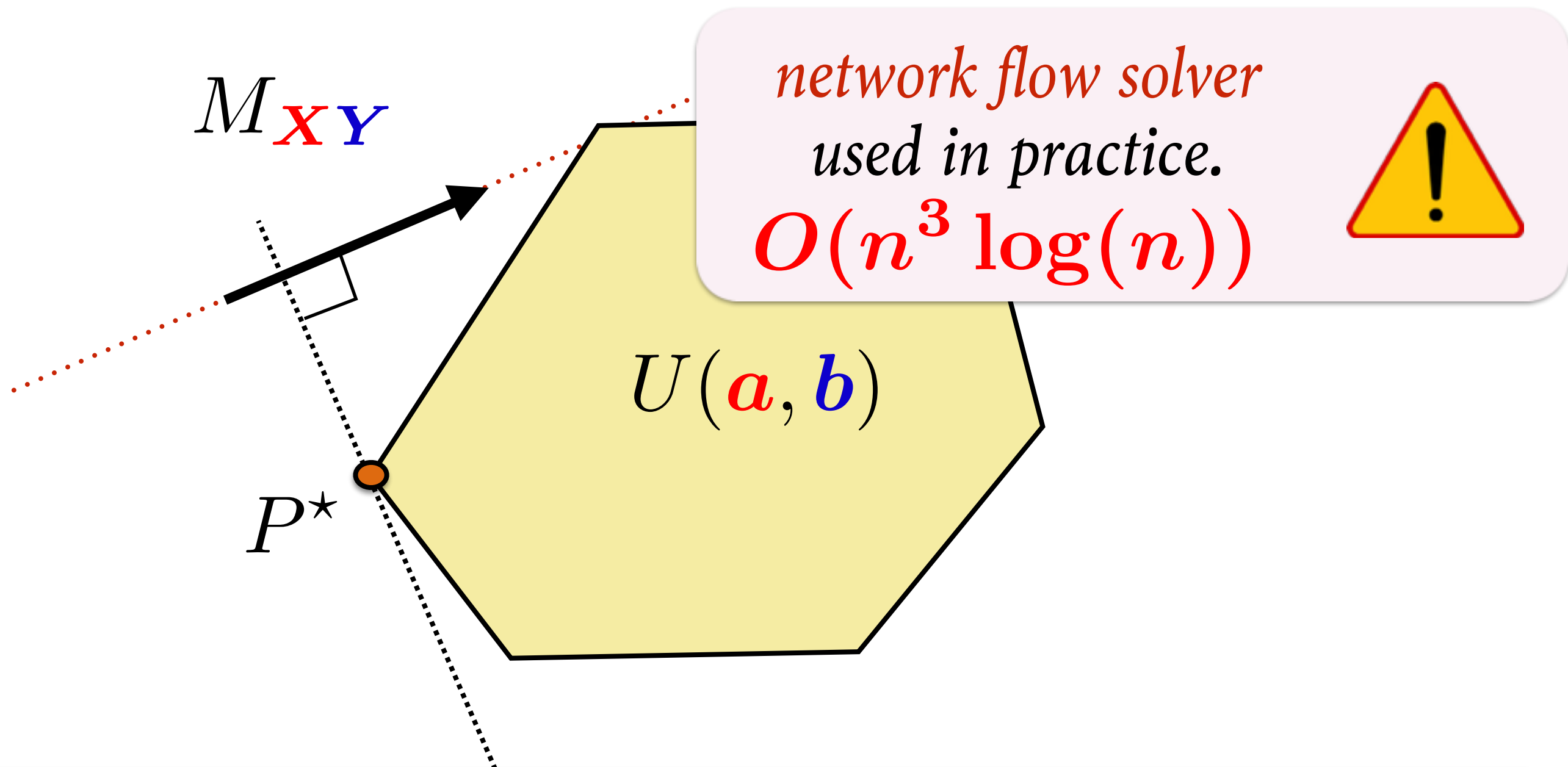
Discrete OT Problem



Def. Dual OT problem

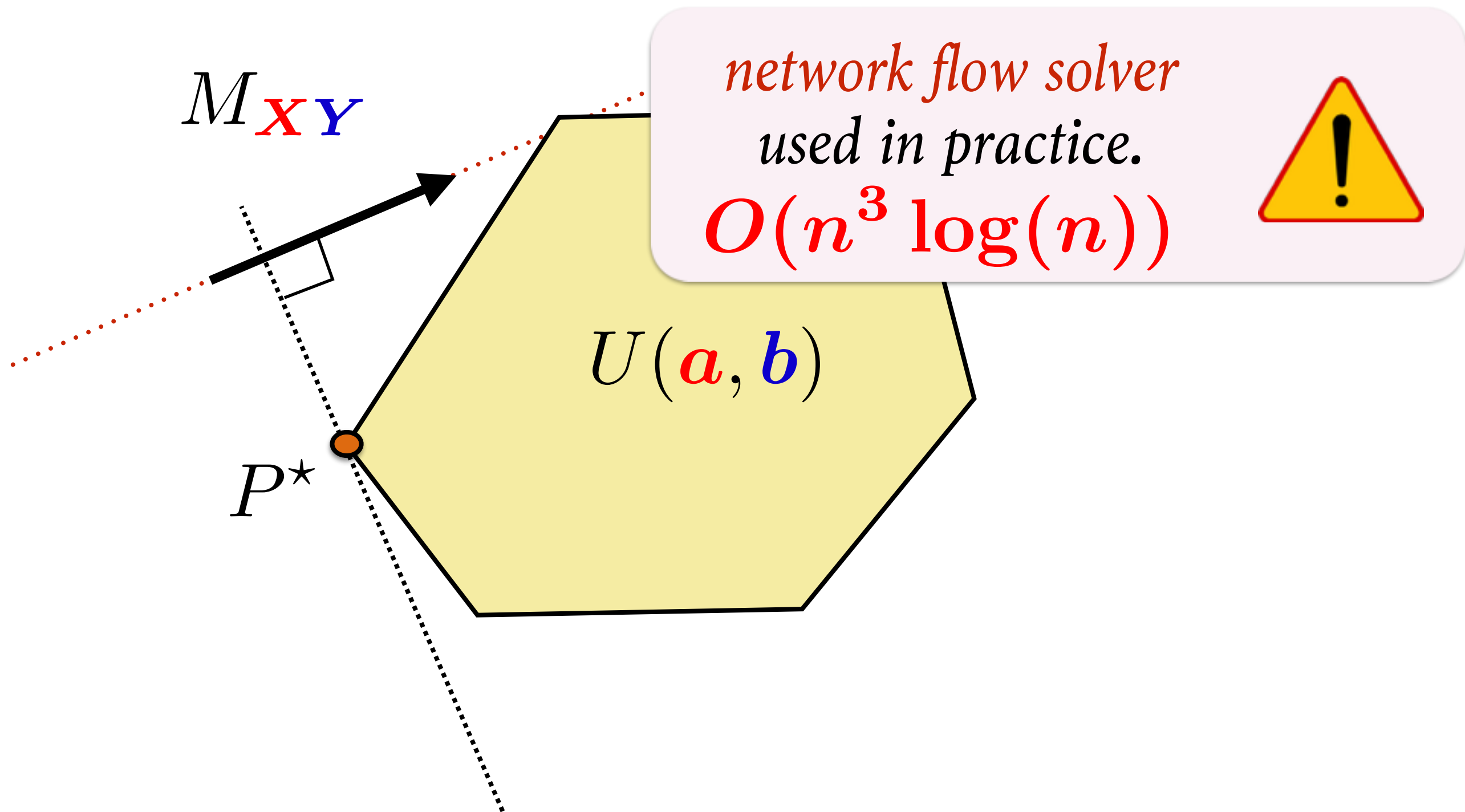
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \alpha_i + \beta_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \alpha^T \mathbf{a} + \beta^T \mathbf{b}$$

Discrete OT Problem

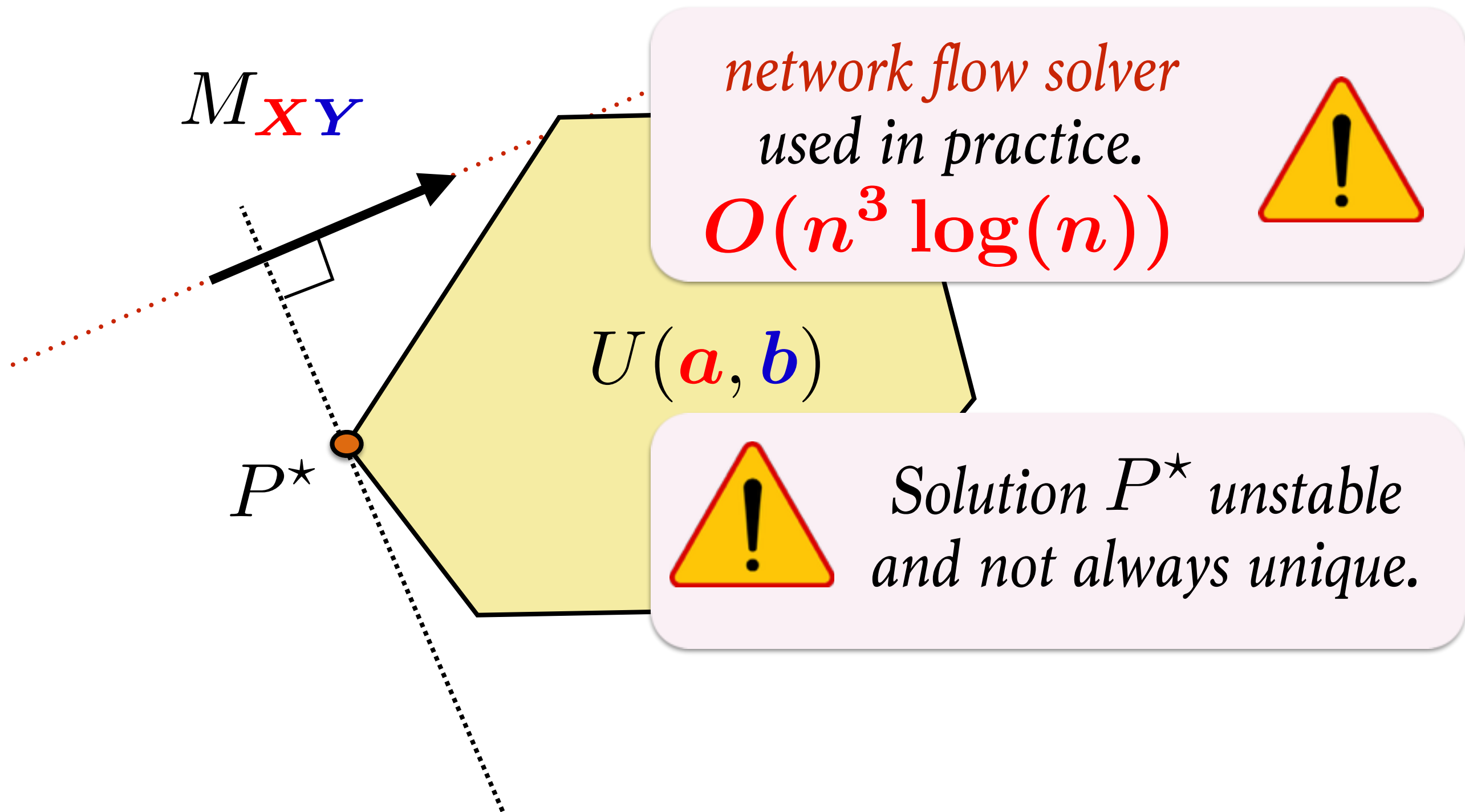


Note: flow/PDE formulations [Beckman'61]/[Benamou'98] can be used for $p=1/p=2$ for a sparse-graph metric/Euclidean metric.

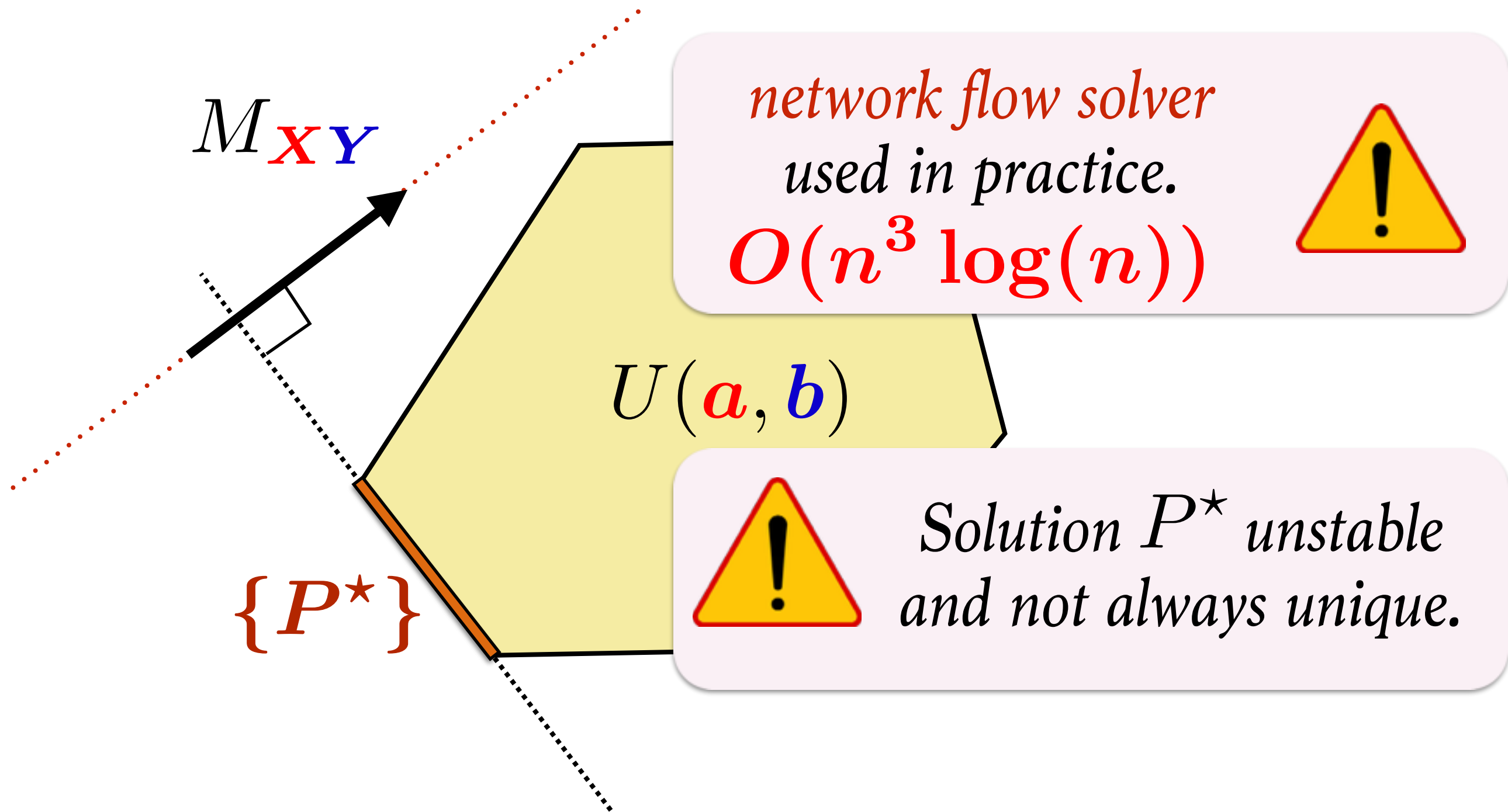
Discrete OT Problem



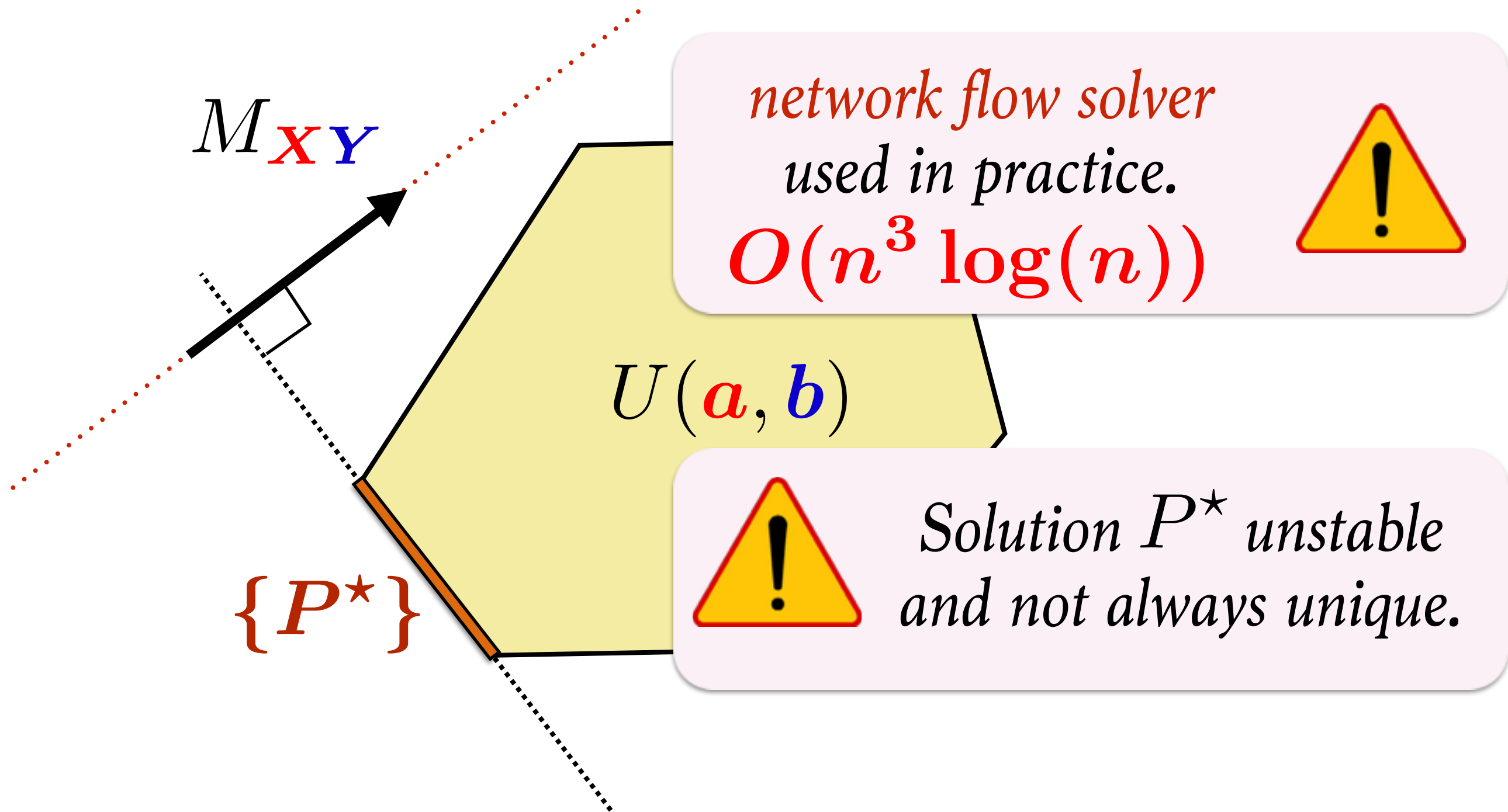
Discrete OT Problem



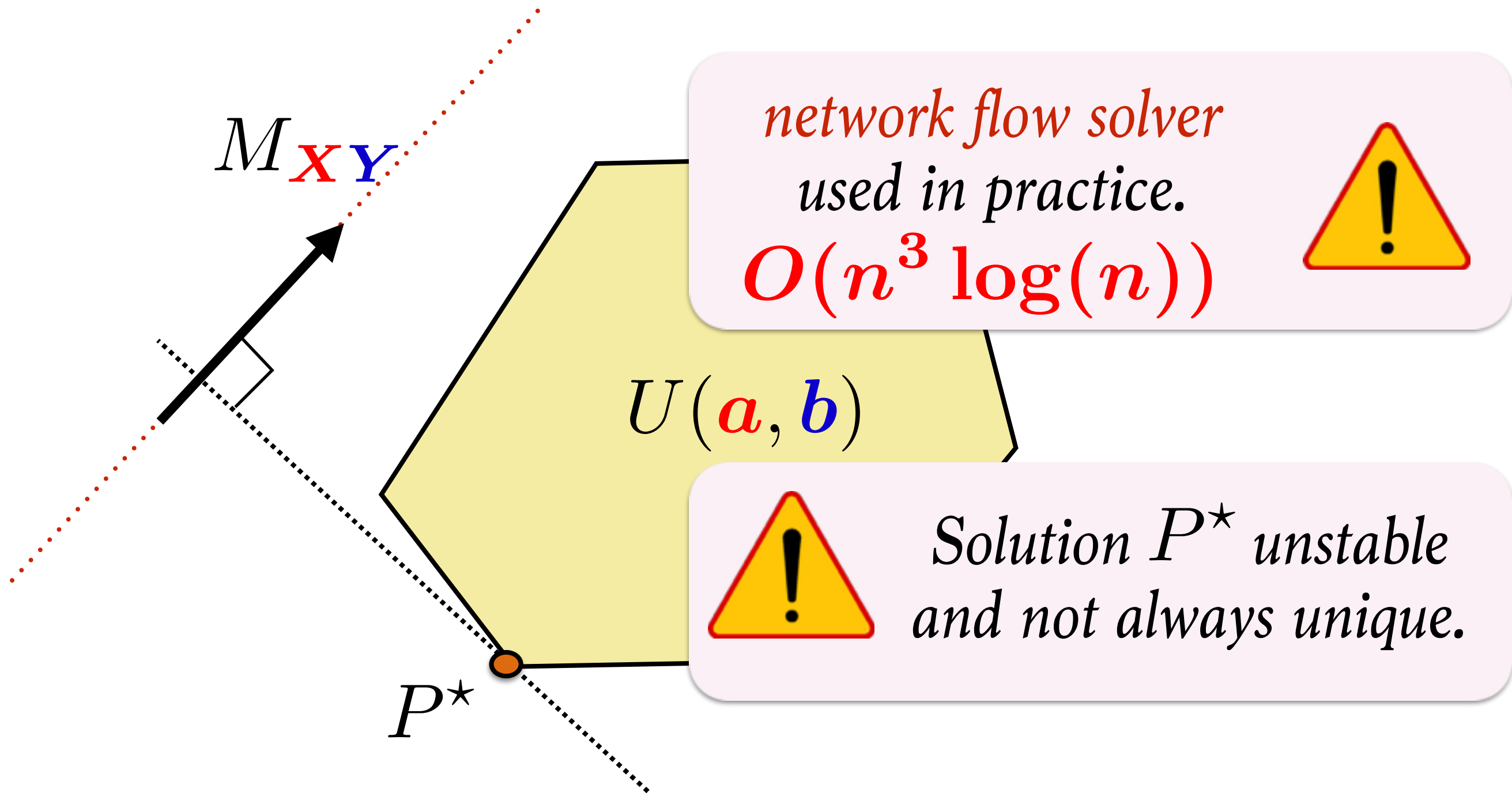
Discrete OT Problem



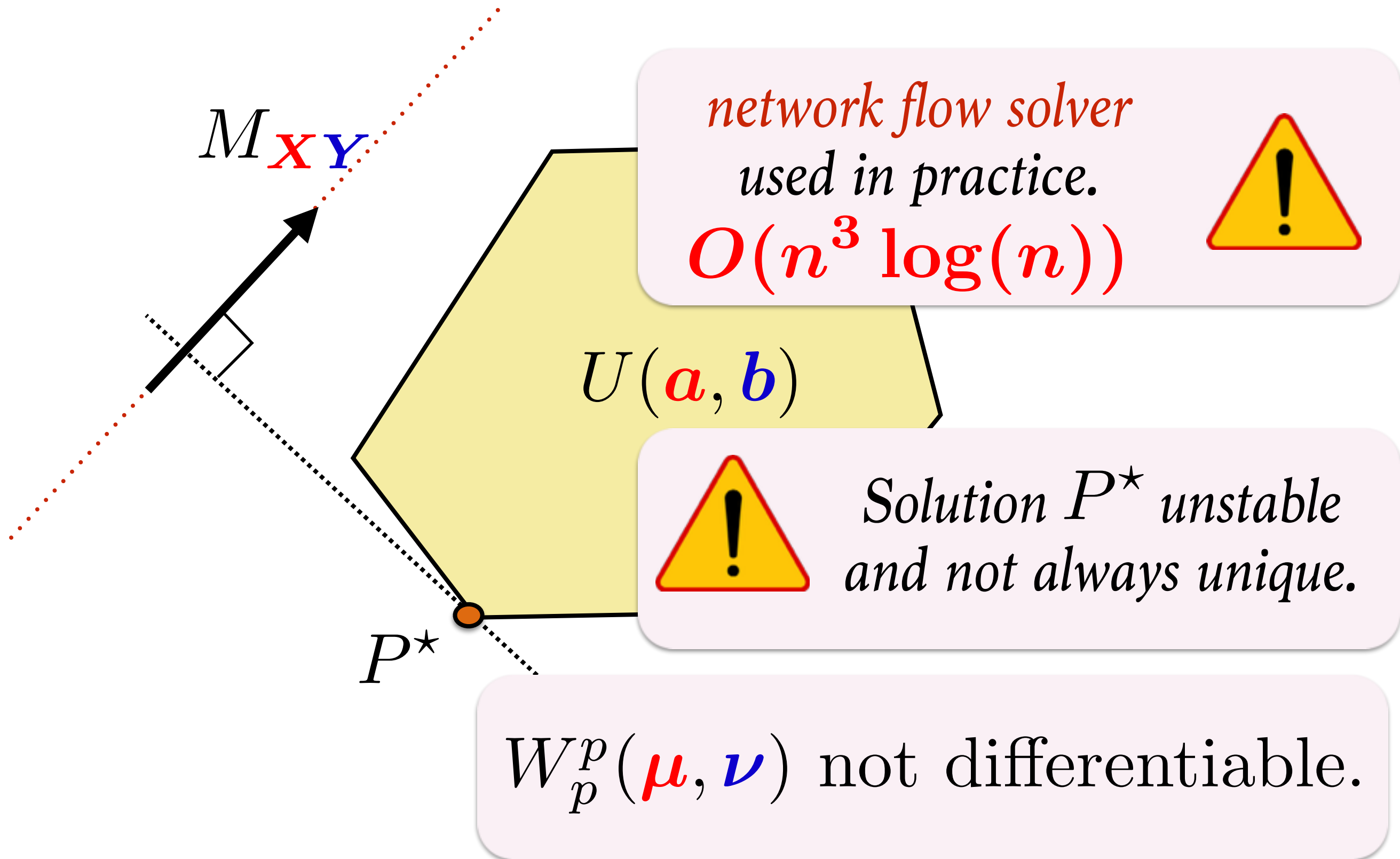
Discrete OT Problem



Discrete OT Problem



Discrete OT Problem



Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein, $\gamma \geq 0$

$$W_\gamma(\mu, \nu) \stackrel{\text{def}}{=} \min_{P \in U(a, b)} \langle P, M_{XY} \rangle - \gamma E(P)$$

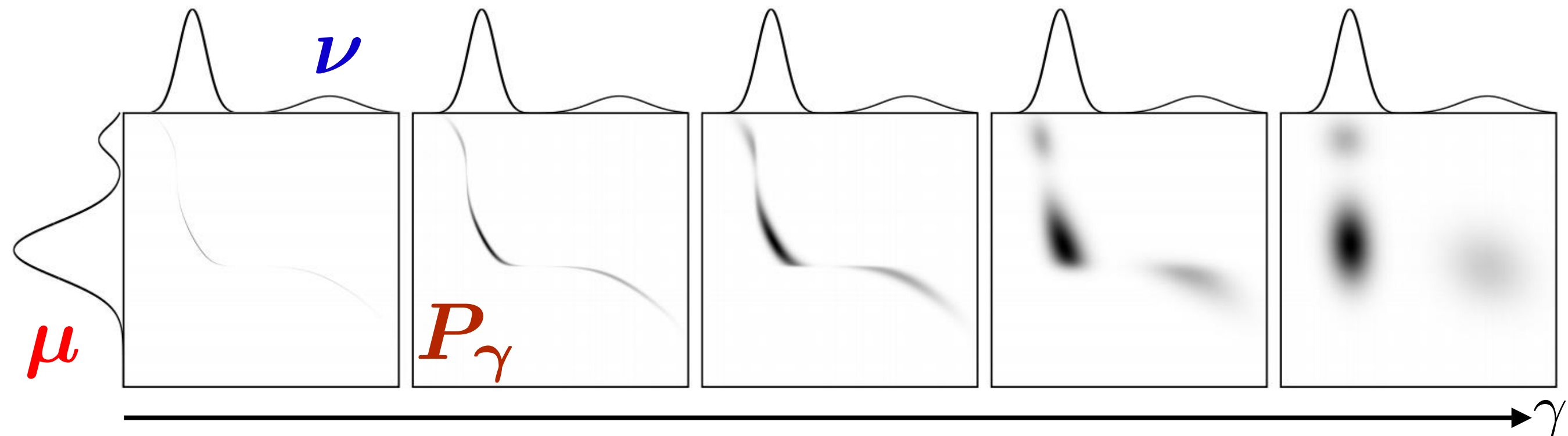
$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij} (\log P_{ij})$$

Note: Unique optimal solution because of strong concavity of Entropy

Entropic Regularization [Wilson'62]

Def. Regularized Wasserstein, $\gamma \geq 0$

$$W_\gamma(\mu, \nu) \stackrel{\text{def}}{=} \min_{P \in U(a, b)} \langle P, M_{XY} \rangle - \gamma E(P)$$



Note: Unique optimal solution because of strong concavity of Entropy

Fast & Scalable Algorithm

Prop. If $P_\gamma \stackrel{\text{def}}{=} \underset{P \in U(\mathbf{a}, \mathbf{b})}{\operatorname{argmin}} \langle \mathbf{P}, M_{\mathbf{x}\mathbf{y}} \rangle - \gamma E(\mathbf{P})$

then $\exists! \mathbf{u} \in \mathbb{R}_+^n, \mathbf{v} \in \mathbb{R}_+^m$, such that

$$P_\gamma = \operatorname{diag}(\mathbf{u}) \mathbf{K} \operatorname{diag}(\mathbf{v}), \quad \mathbf{K} \stackrel{\text{def}}{=} e^{-M_{\mathbf{x}\mathbf{y}} / \gamma}$$

Fast & Scalable Algorithm

Prop. If $P_\gamma \stackrel{\text{def}}{=} \underset{P \in U(\mathbf{a}, \mathbf{b})}{\operatorname{argmin}} \langle \mathbf{P}, M_{\mathbf{x}\mathbf{y}} \rangle - \gamma E(\mathbf{P})$

then $\exists! \mathbf{u} \in \mathbb{R}_+^n, \mathbf{v} \in \mathbb{R}_+^m$, such that

$$P_\gamma = \operatorname{diag}(\mathbf{u}) \mathbf{K} \operatorname{diag}(\mathbf{v}), \quad \mathbf{K} \stackrel{\text{def}}{=} e^{-M_{\mathbf{x}\mathbf{y}} / \gamma}$$

$$L(P, \alpha, \beta) = \sum_{ij} P_{ij} M_{ij} + \gamma P_{ij} \log P_{ij} + \alpha^T (P \mathbf{1} - \mathbf{a}) + \beta^T (P^T \mathbf{1} - \mathbf{b})$$

$$\partial L / \partial P_{ij} = M_{ij} + \gamma (\log P_{ij} + 1) + \alpha_i + \beta_j$$

$$(\partial L / \partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma} + \frac{1}{2}} e^{-\frac{M_{ij}}{\gamma}} e^{\frac{\beta_j}{\gamma} + \frac{1}{2}} = \mathbf{u}_i \mathbf{K}_{ij} \mathbf{v}_j$$

Fast & Scalable Algorithm

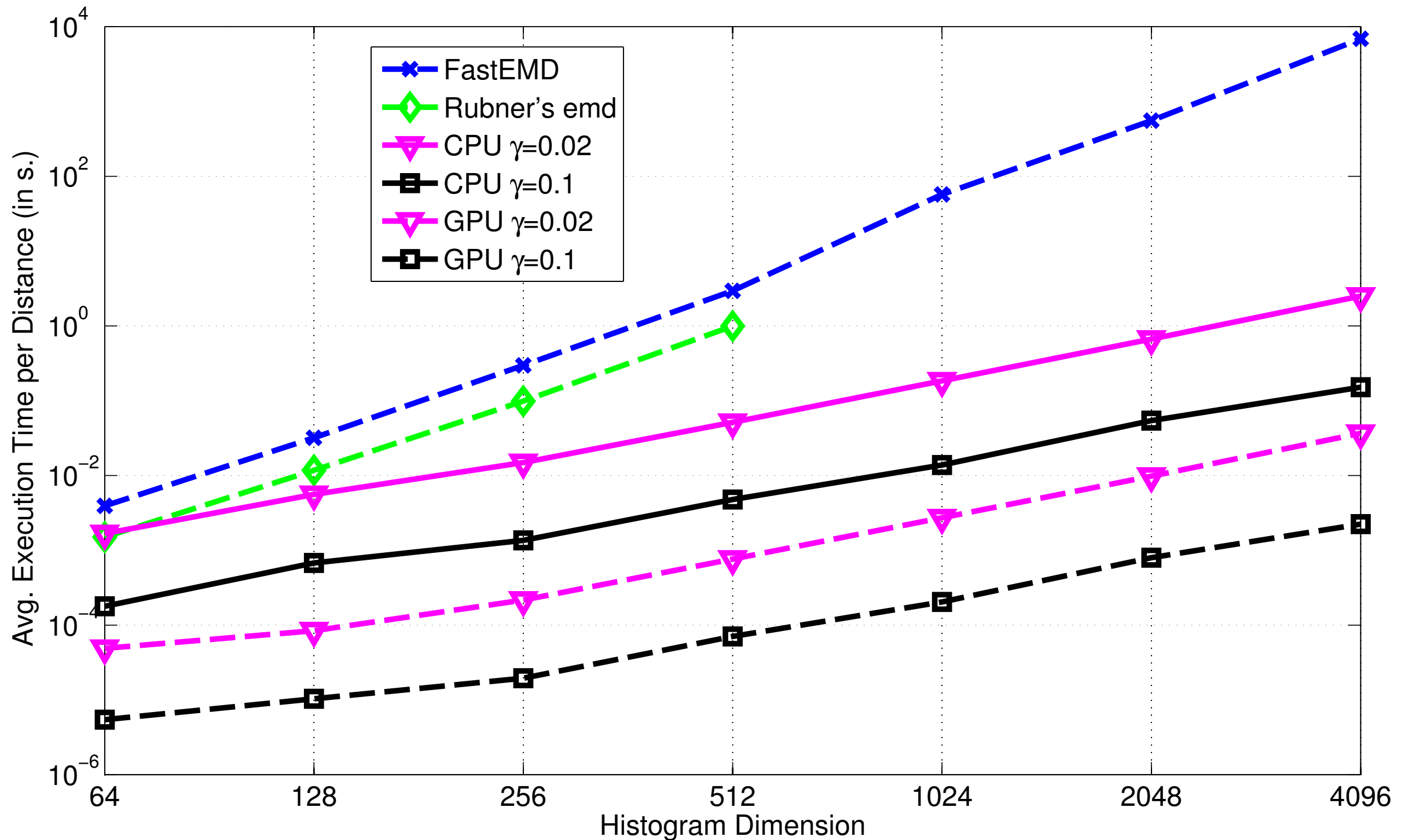
Prop. If $P_\gamma \stackrel{\text{def}}{=} \underset{P \in U(\mathbf{a}, \mathbf{b})}{\operatorname{argmin}} \langle \mathbf{P}, M_{\mathbf{x}\mathbf{y}} \rangle - \gamma E(\mathbf{P})$

then $\exists! \mathbf{u} \in \mathbb{R}_+^n, \mathbf{v} \in \mathbb{R}_+^m$, such that

$$P_\gamma = \operatorname{diag}(\mathbf{u}) \mathbf{K} \operatorname{diag}(\mathbf{v}), \quad \mathbf{K} \stackrel{\text{def}}{=} e^{-M_{\mathbf{x}\mathbf{y}} / \gamma}$$

- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})
$$\mathbf{u} \leftarrow \mathbf{a} / \mathbf{K} \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / \mathbf{K}^T \mathbf{u}$$
- $O(nm)$ complexity, GPGPU parallel [C'13].
- $O(n^{d+1})$ if $\Omega = \{1, \dots, n\}^d$ and \mathbf{D}^p separable.
[S..C..'15]

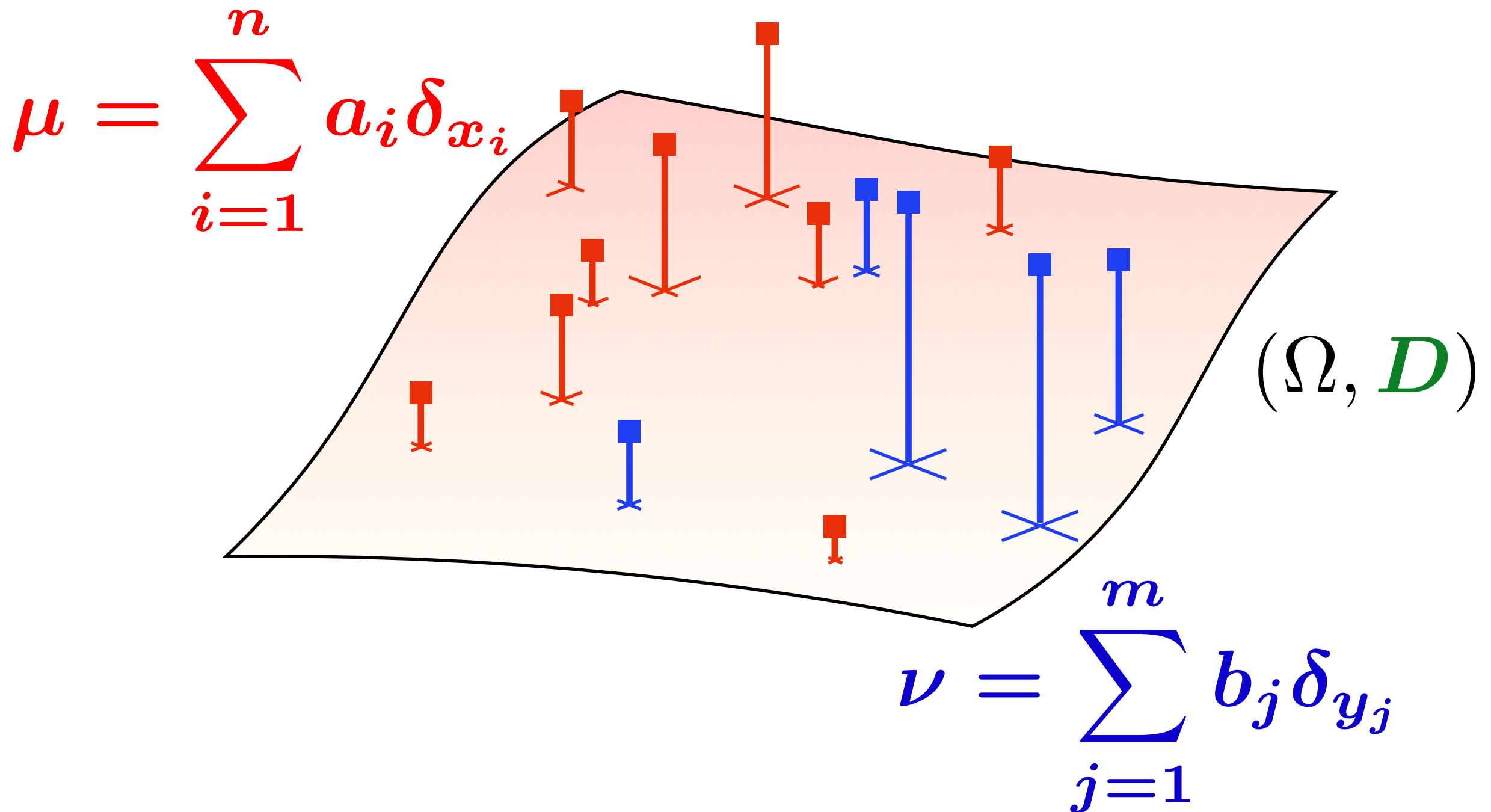
Very Fast EMD Approx. Solver



Note. (Ω, \mathbf{D}) is a random graph with shortest path metric, histograms sampled uniformly on simplex, Sinkhorn tolerance 10^{-2} .

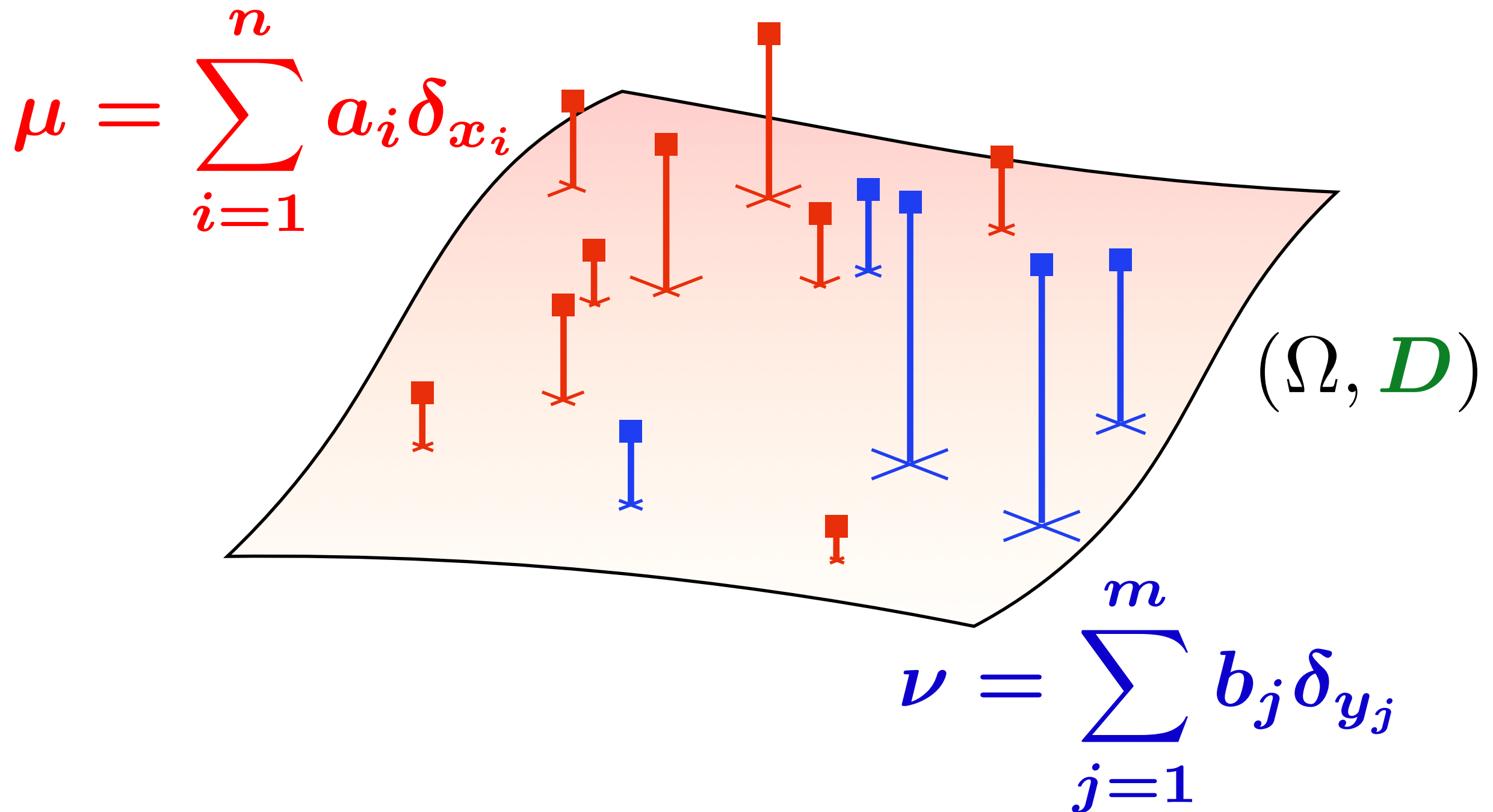
Regularization \rightsquigarrow *Differentiability*

$$W_\gamma((a, X), (b, Y)) = \min_{P \in U(a, b)} \langle P, M_{XY} \rangle - \gamma E(P)$$



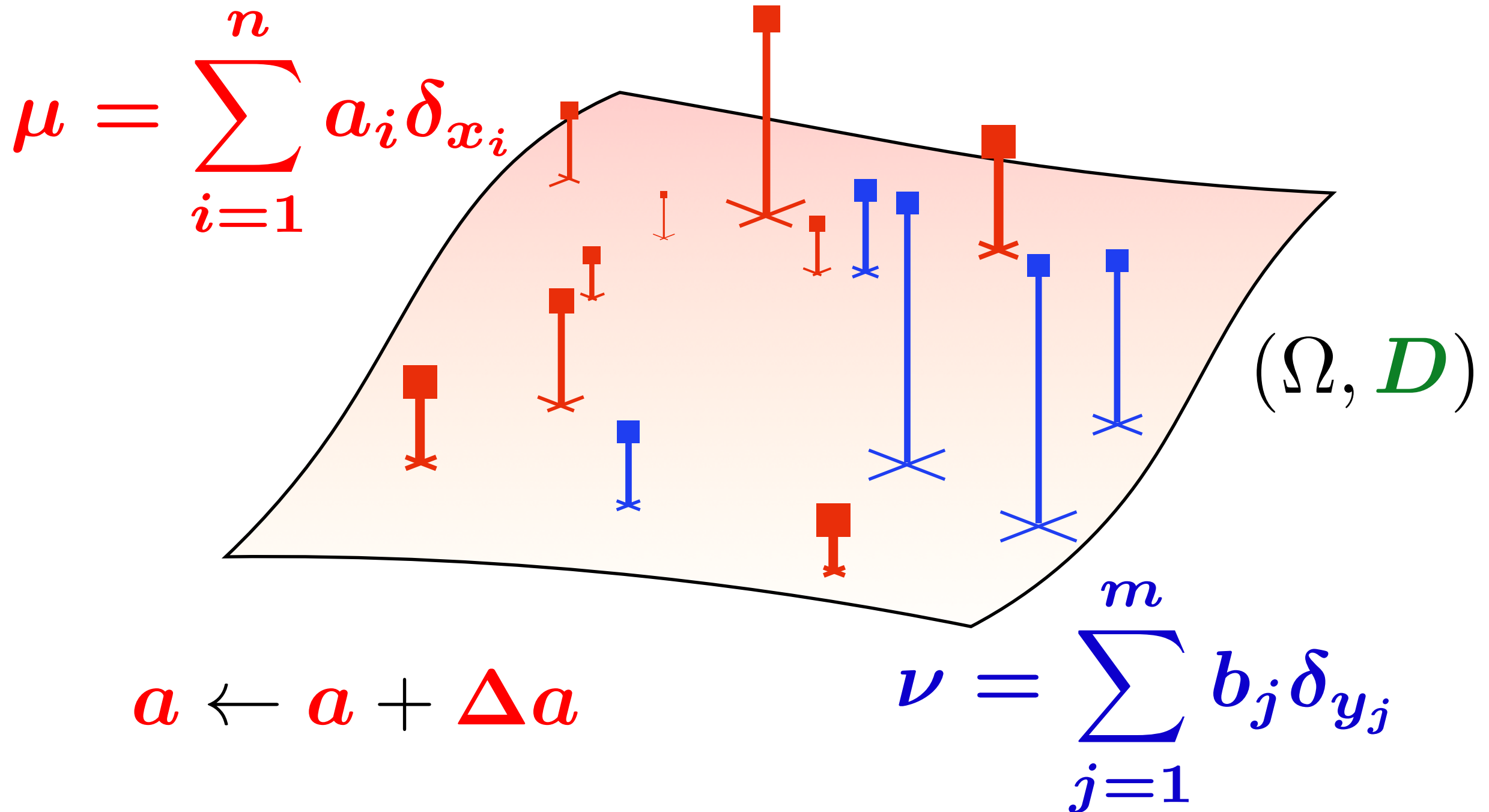
Regularization \rightsquigarrow *Differentiability*

$$W_\gamma((a + \Delta a, X), (b, Y)) = W_\gamma((a, X), (b, Y)) + ??$$



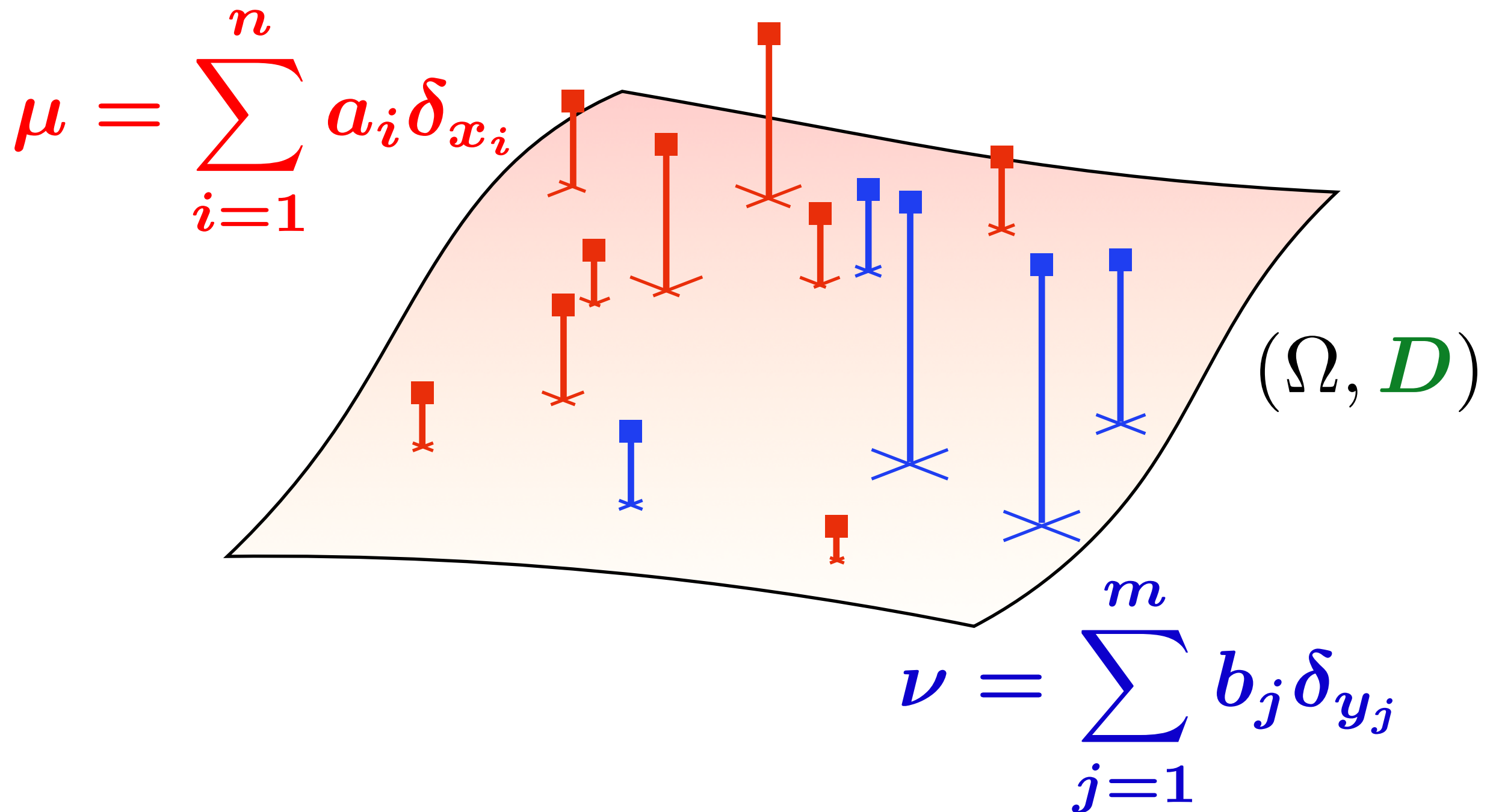
Regularization \rightsquigarrow *Differentiability*

$$W_\gamma((a + \Delta a, X), (b, Y)) = W_\gamma((a, X), (b, Y)) + ??$$



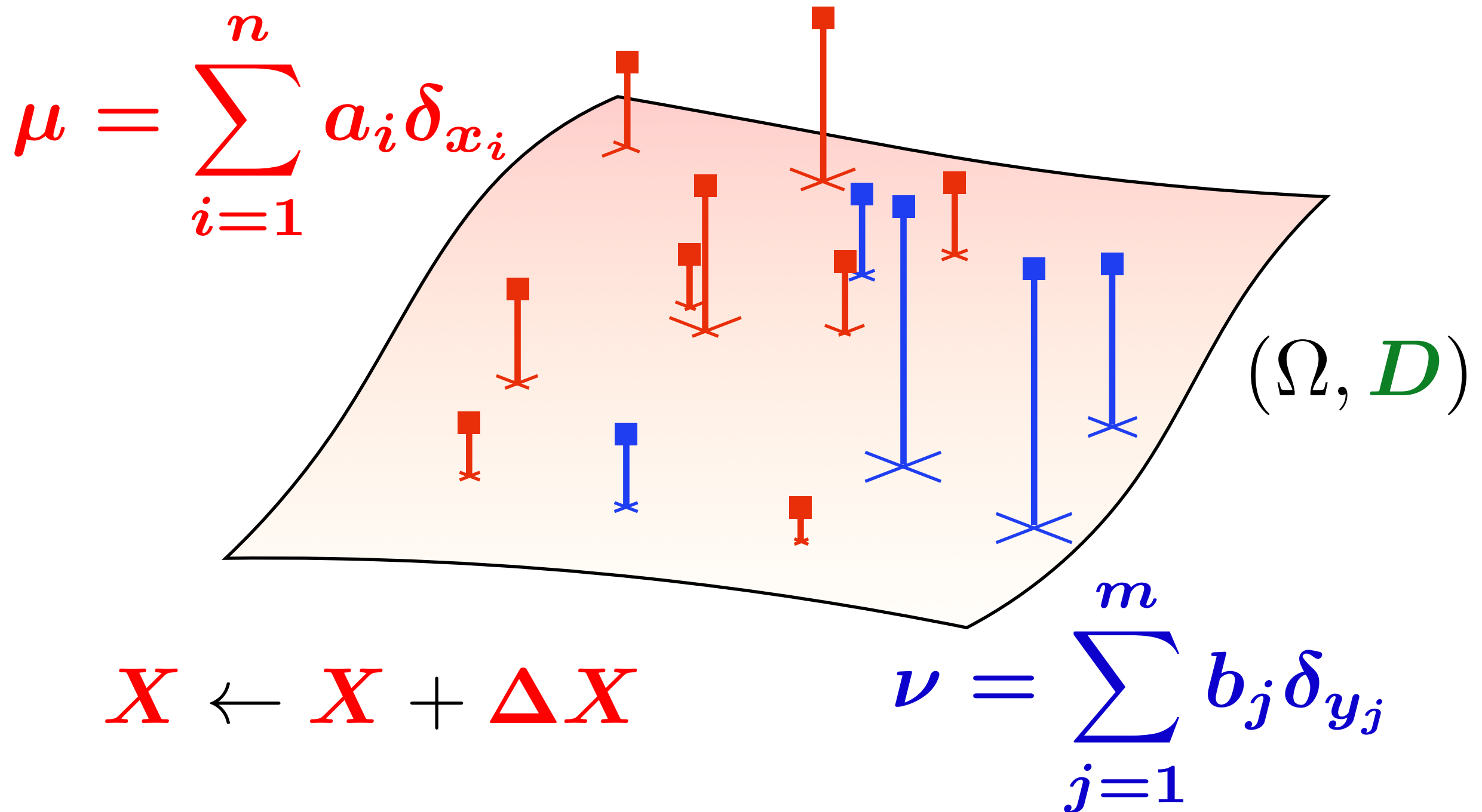
Regularization \rightsquigarrow *Differentiability*

$$W_\gamma((a, X + \Delta X), (b, Y)) = W_\gamma((a, X), (b, Y)) + ??$$



Regularization \rightsquigarrow *Differentiability*

$$W_\gamma((a, X + \Delta X), (b, Y)) = W_\gamma((a, X), (b, Y)) + ??$$



Crucial for “min *data* + *W*” problems

- Quantization, *k*-means problem [Lloyd'82]

$$\min_{\substack{\mu \in \mathcal{P}(\mathbb{R}^d) \\ |\text{supp } \mu| = k}} W_2^2(\mu, \nu_{\text{data}})$$

- [McCann'95] Interpolant

$$\min_{\mu \in \mathcal{P}(\Omega)} (1 - t)W_2^2(\mu, \nu_1) + tW_2^2(\mu, \nu_2)$$

- [JKO'98] PDE's as gradient flows in $(\mathcal{P}(\Omega), W)$.

$$\mu_{t+1} = \operatorname{argmin}_{\mu \in \mathcal{P}(\Omega)} J(\mu) + \lambda_t W_p^p(\mu, \mu_t)$$

Crucial for “min *data* + *W*” problems

- Quantization,

$$\min_{\mu \in \mathcal{P}(\mathbb{R}^d)} W_2^2(\mu, \nu_{\text{data}})$$

Any (ML) problem involving a **KL** or **L2** loss between (parameterized) histograms or probability measures can be easily *Wasserstein-ized* if we can differentiate *W* efficiently.

1. Differentiability of Regularized OT

Def. Dual regularized OT Problem

$$W_\gamma(\mu, \nu) = \max_{\alpha, \beta} \alpha^T \mathbf{a} + \beta^T \mathbf{b} - \frac{1}{\gamma} (e^{\alpha/\gamma})^T \mathbf{K} e^{\beta/\gamma}$$

Prop. $W_\gamma(\mu, \nu)$ is

[CP'16]

1. convex w.r.t. \mathbf{a} (Danskin),

$$\nabla_{\mathbf{a}} W_\gamma = \alpha^* = \gamma \log(\mathbf{u}).$$

2. decreased, when $p = 2, \Omega = \mathbb{R}^d$, using

$$\mathbf{X} \leftarrow \mathbf{Y} P_\gamma^T \mathbf{D}(\mathbf{a}^{-1}).$$

2. Duality for Regularized OT's

Prop. Writing $H_{\nu} : \mathbf{a} \mapsto W_{\gamma}(\mu, \nu)$, [CP'16]

1. H_{ν} has simple Legendre transform:

$$H_{\nu}^* : \mathbf{g} \in \mathbb{R}^n \mapsto \gamma \left(E(\mathbf{b}) + \mathbf{b}^T \log(\mathbf{K} e^{\mathbf{g}/\gamma}) \right)$$

2. If $A \in \mathbb{R}^{n \times d}$, f convex on \mathbb{R}^d ,

$$\min_{\mathbf{a} \in \Sigma_n} H_{\nu}(\mathbf{a}) + f(A\mathbf{a}) = \max_{\mathbf{g} \in \mathbb{R}^d} -H_{\nu}^*(A^T \mathbf{g}) - f^*(-\mathbf{g})$$

3. Stochastic Formulation

$$\begin{aligned} W_\gamma(\mu, \nu) &= \max_{\alpha, \beta} \alpha^T \mathbf{a} + \beta^T \mathbf{b} - \frac{1}{\gamma} (e^{\alpha/\gamma})^T \mathbf{K} e^{\beta/\gamma} \\ &= \max_{\alpha} \alpha^T \mathbf{a} - \gamma (\log \mathbf{K} e^{\alpha/\gamma})^T \mathbf{b} \\ &= \max_{\alpha} \sum_{j=1}^m \mathbf{b}_j \left(\alpha^T \mathbf{a} - \gamma \log \mathbf{K}_{\cdot j}^T e^{\alpha/\gamma} \right) \\ &= \max_{\alpha} \sum_{j=1}^m f_j(\alpha) \end{aligned}$$

- **[GCPB'16]** shows how incremental gradient methods can be used to scale this further.

4. Algorithmic Formulation

Def. For $L \geq 1$, define

$$W_L(\mu, \nu) \stackrel{\text{def}}{=} \langle P_L, M_{\mathbf{x}\mathbf{y}} \rangle,$$

where $P_L \stackrel{\text{def}}{=} \text{diag}(\mathbf{u}_L) K \text{diag}(\mathbf{v}_L)$,

$$\mathbf{v}_0 = \mathbf{1}_m; l \geq 0, \mathbf{u}_l \stackrel{\text{def}}{=} \mathbf{a} / K \mathbf{v}_l, \mathbf{v}_{l+1} \stackrel{\text{def}}{=} \mathbf{b} / K^T \mathbf{u}_l.$$

Prop. $\frac{\partial W_L}{\partial \mathbf{x}}, \frac{\partial W_L}{\partial \mathbf{a}}$ can be computed recursively, in $O(L)$ kernel $K \times$ vector products.

Algorithmic Formulation of Reg. OT

Example: Differentiability w.r.t. a

$$\left(\frac{\partial \mathbf{v}_0}{\partial a} \right)^T = \mathbf{0}_{m \times n},$$

$$\left(\frac{\partial \mathbf{u}_l}{\partial a} \right)^T \mathbf{x} = \frac{\mathbf{x}}{\mathbf{K} \mathbf{v}_l} - \left(\frac{\partial \mathbf{v}_l}{\partial a} \right)^T \mathbf{K}^T \frac{\mathbf{x} \circ a}{(\mathbf{K} \mathbf{v}_l)^2},$$

$$\left(\frac{\partial \mathbf{v}_{l+1}}{\partial a} \right)^T \mathbf{y} = - \left(\frac{\partial \mathbf{u}_l}{\partial a} \right)^T \mathbf{K} \frac{\mathbf{y} \circ b}{(\mathbf{K}^T \mathbf{u}_l)^2}.$$

Algorithmic Formulation of Reg. OT

Example: Differentiability w.r.t. a

$$\textcolor{blue}{N} = \textcolor{blue}{K} \circ M_{\textcolor{red}{X}\textcolor{blue}{Y}}$$

$$\nabla_{\textcolor{red}{a}} W_L(\textcolor{red}{\mu}, \textcolor{blue}{\nu}) = \left(\frac{\partial \textcolor{brown}{u}_L}{\partial a} \right)^T \textcolor{blue}{N} \textcolor{brown}{v}_L + \left(\frac{\partial \textcolor{brown}{v}_L}{\partial a} \right)^T \textcolor{blue}{N}^T \textcolor{brown}{u}_L$$


```

function [d,grad_a,grad_b,hess_a,hess_b] = sinkhornObjGradHess(a,b,K,M,niter)

u_update = @(v,a) a./(K*v);
v_update = @(u,b) b./(K'*u);

% DuDa = @(eps,dvda,a,v) (eps./(K*v)) - (a./((K*v).^2)).*(K*dvda(eps));
%
% DvDa = @(eps,duda,b,u) -(b./((K'*u).^2)).*(K'*duda(eps));
%
% DuDb = @(eps,dvdb,a,v) -(a./((K*v).^2)).*(K*dvdb(eps));
%
% DvDb = @(eps,dudb,b,u) (eps./(K'*u)) - (b./((K'*u).^2)).*(K'*dudb(eps));

DuDat = @(x,dvdat,a,v) bsxfun(@rdivide,x,K*v)... (x./(K*v))
        -dvdat(K'*( bsxfun(@times,x,(a./((K*v).^2)))));...-dvdat(K'*( (a./((K*v).^2)).*x));

DvDat = @(x,dudat,b,u) -dudat(K*(bsxfun(@times,x,(b./((K'*u).^2))))); ... (b./((K'*u).^2)).*x))

JDuDat= @(x,Jdvdat,dvdat,a,v) -diag((x'*dvdat(K'))'./((K*v).^2)) ... (K*dvda(x))
        - Jdvdat(x)*K'*diag(a./((K*v).^2))...
        - dvdat(K'* ...
        ( diag(a.*( (-2*(x'*dvdat(K'))'./((K*v).^3)))+...
        diag(x./((K*v).^2)) )); %1

JDvDat = @(x,Jdudat,dudat,b,u) ...
        -Jdudat(x)*K*diag(b./((K'*u).^2))...
        - dudat(K)* ( ...
        diag(b.*( (-2*(x'*dudat(K))'./((K'*u).^3)))) ;...

```

```
DuDbt = @(x,dvdbt,a,v) -dvdbt(K'*(bsxfun(@times,x,(a./((K*v).^2))))); ... (a./((K*v).^2)).*x);
```

```
DvDbt = @(x,dudbt,b,u) bsxfun(@rdivide,x,K'*u) ... (x./((K'*u)))...  
-dudbt(K*( bsxfun(@times,x,(b./((K'*u).^2))))); ... ( b./((K'*u).^2)) .*x);
```

```
JDvDbt= @(x,Jdudbt,dudbt,b,u) -diag((x'*dudbt(K))'./((K'*u).^2)) ... (K'*dudb(x))  
- Jdudbt(x)*K*diag(b./((K'*u).^2))...  
- dudbt(K)* ( ...  
diag(b.*( (-2*(x'*dudbt(K))' )./((K'*u).^3)))+...  
diag(x./((K'*u).^2)) ) ;
```

```
JDvDbt = @(x,Jdvdbt,dvdbt,a,v) ...  
-Jdvdbt(x)*K'*diag(a./((K*v).^2))...  
- dvdbt(K')* ( ...  
diag(a.*( (-2* (x'*dvdbt(K'))' )./((K*v).^3)))) ;
```

```

n=size(a,1);
m=size(b,1);

DV DAT = @(eps) zeros(n,size(eps,2));
DV DBT = @(eps) zeros(m,size(eps,2));

JDV DAT = @(eps) zeros(n,m);
JDV DBT = @(eps) zeros(m,m);

v=ones(m,size(b,2));

for j=1:niter,
    u=u_update(v,a);
    DU DAT = @(x) DuDat(x,DV DAT,a,v);
    DU DBT = @(x) DuDbt(x,DV DBT,a,v);

    if nargout>3
        JDU DAT = @(x) JDuDat(x,JDV DAT,DV DAT,a,v);
        JDU DBT = @(x) JDuDbt(x,JDV DBT,DV DBT,a,v);
    end

    v=v_update(u,b);
    DV DAT = @(x) DvDat(x,DU DAT,b,u);
    DV DBT = @(x) DvDbt(x,DU DBT,b,u);

    if nargout>3
        JDV DAT = @(x) JvDat(x,JDU DAT,DU DAT,b,u);
        JDV DBT = @(x) JvDbt(x,JDU DBT,DU DBT,b,u);
    end
end
end

```

```

U=K.*M;
d=diag(u'*U*v);

grad_a=(DUDAT(U*v)+DV DAT(U'*u));
grad_b=(DU DBT(U*v)+DV DBT(U'*u));

if nargout>3
    hess_a= @(eps) JDUDAT(eps)*(U*v)+DUDAT((eps'*DV DAT(U'))')+...
            JDV DAT(eps)*(U'*u)+DV DAT((eps'*DUDAT(U))')+...
end

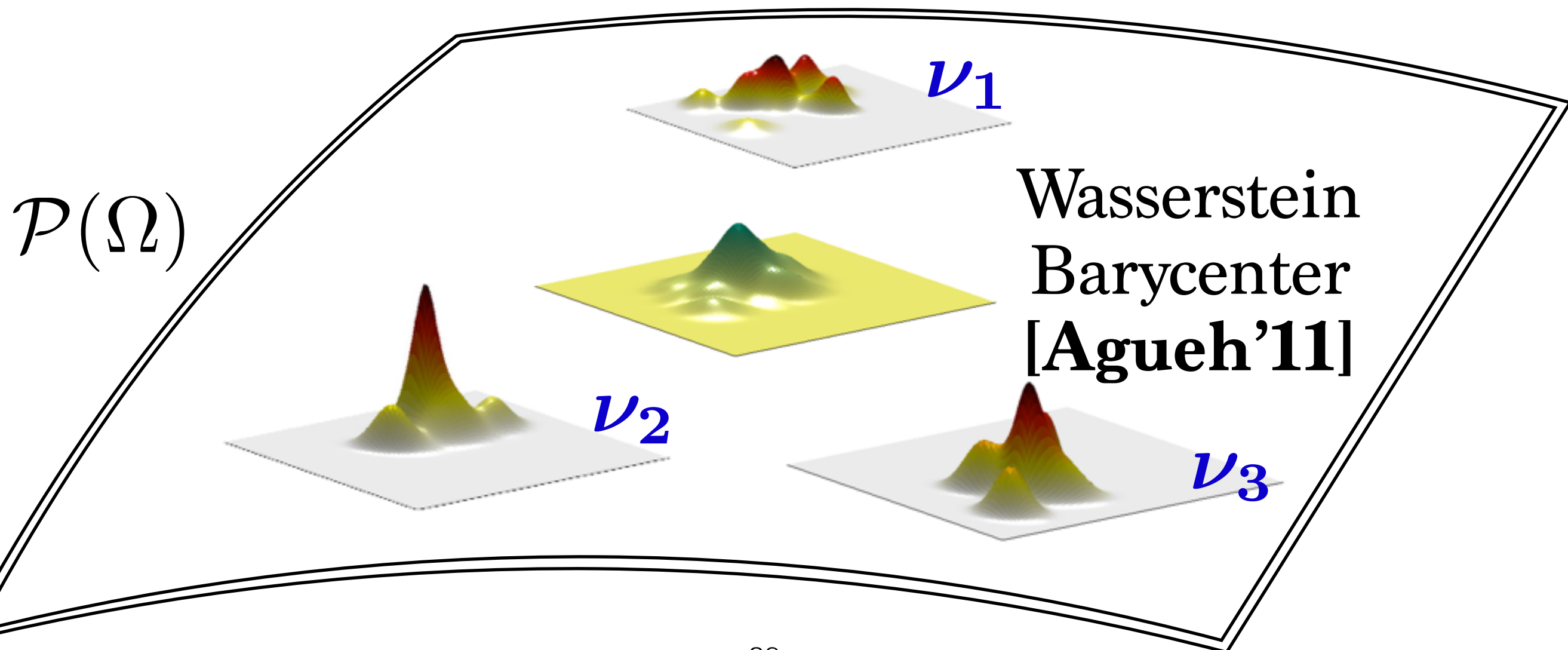
```

Thanks to these tricks...

- **[Agueh'11]** Barycenters **[CD'14][BCCNP'15]**
[GCP'15][S..C..'15]
- **[Burger'12]** TV gradient flow using duality **[CP'16]**
- Dictionary Learning / Latent Factors **[RCP'16]**
- **[Bigot'15]** W-PCA **[SC'15]**
- Density fitting / parameter estimation **[MMC'16]**
- **Inverse problems / Wasserstein regression** **[BPC'16]**

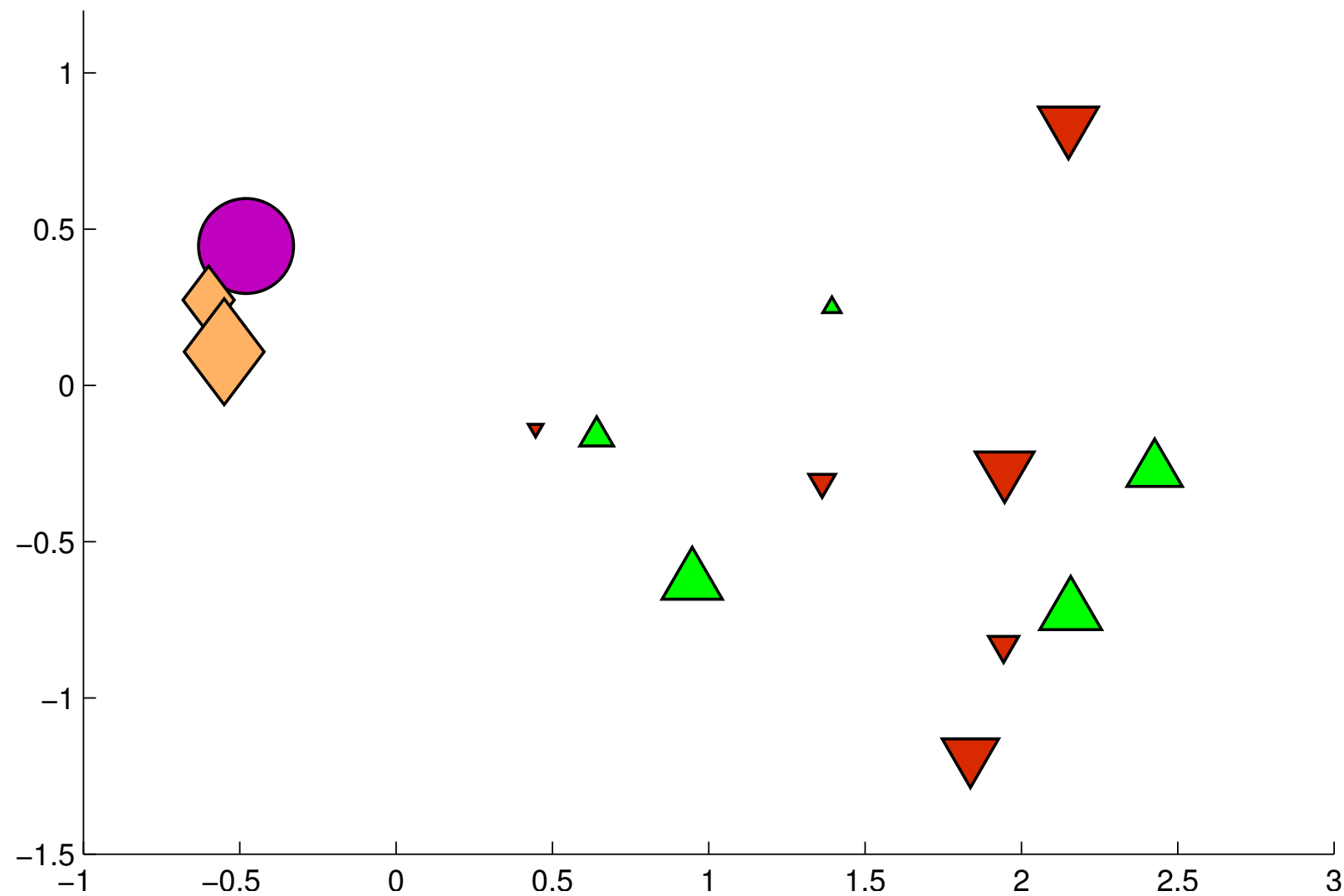
Wasserstein Barycenters

$$\min_{\boldsymbol{\mu} \in \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}_i)$$



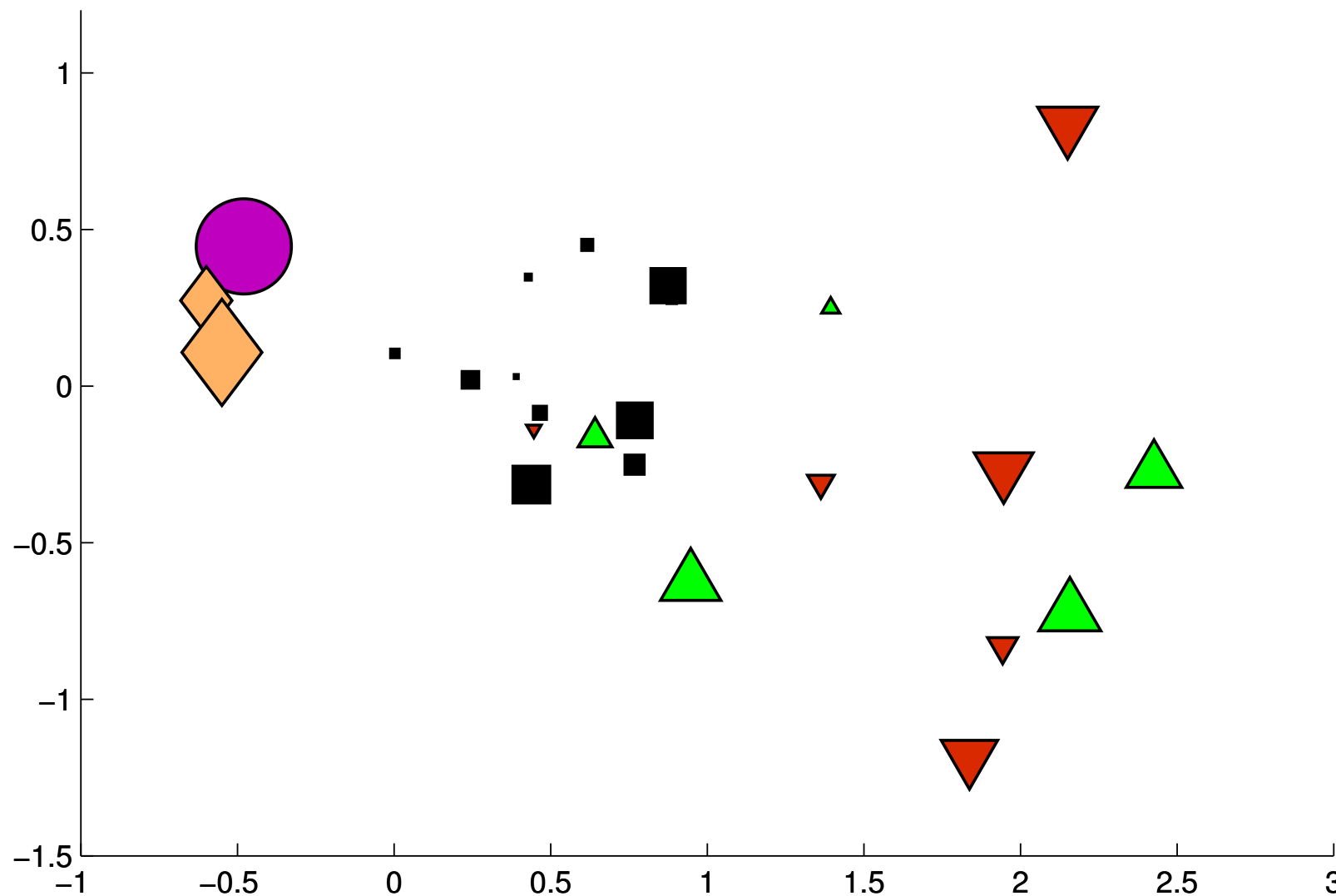
Multimarginal Formulation

- Exact solution (W_2) using MM-OT. [Agueh'11]



Multimarginal Formulation

- Exact solution (W_2) using MM-OT. [Agueh'11]



If $|\text{supp } \nu_i| = n_i$, LP of size $(\prod_i n_i, \sum_i n_i)$

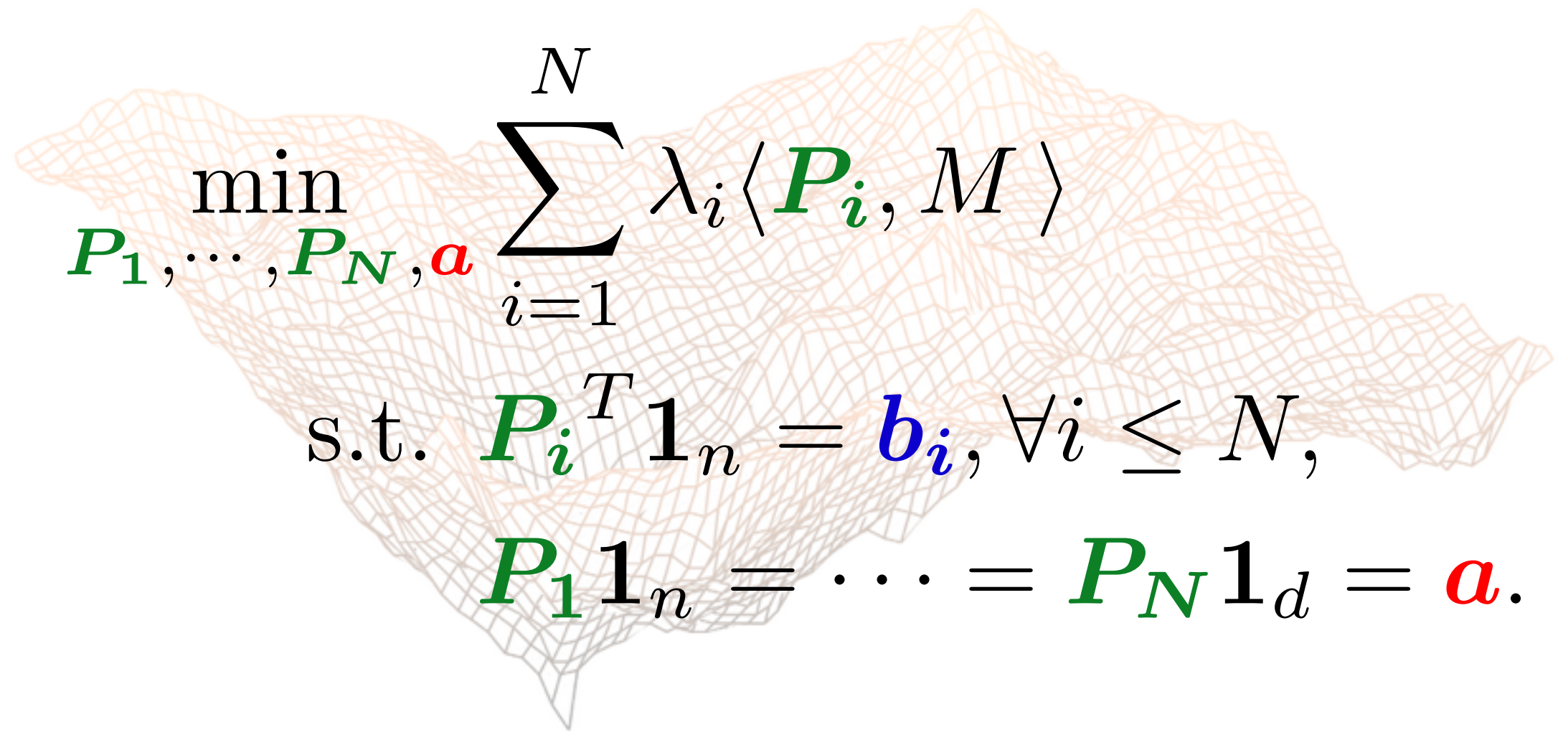
Finite Case, LP Formulation

- When Ω is a **finite set**, metric M , another LP.


$$\min_{\mu} \sum_i \lambda_i W_p^p(\mu, \nu_i)$$

Finite Case, LP Formulation

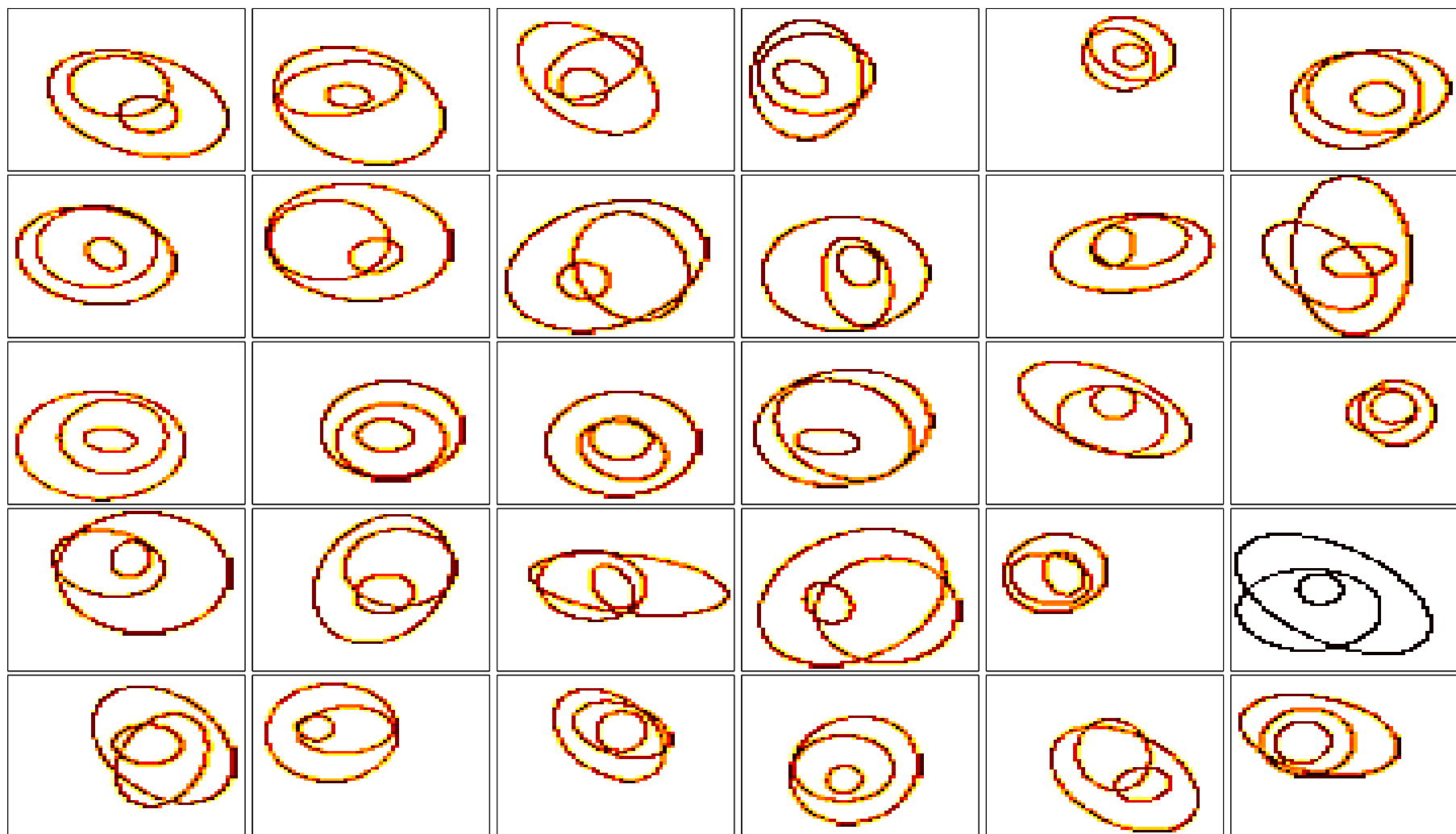
- When Ω is a **finite set**, metric M , another LP.


$$\begin{aligned} \min_{\mathbf{P}_1, \dots, \mathbf{P}_N, \mathbf{a}} \quad & \sum_{i=1}^N \lambda_i \langle \mathbf{P}_i, M \rangle \\ \text{s.t.} \quad & \mathbf{P}_i^T \mathbf{1}_n = \mathbf{b}_i, \forall i \leq N, \\ & \mathbf{P}_1 \mathbf{1}_n = \dots = \mathbf{P}_N \mathbf{1}_d = \mathbf{a}. \end{aligned}$$

If $|\Omega| = n$, LP of size $(Nn^2, (2N - 1)n)$; unstable

Primal Descent on Regularized W

$$\min_{\mu \in Q \subset \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_{\gamma}(\mu, \nu_i)$$

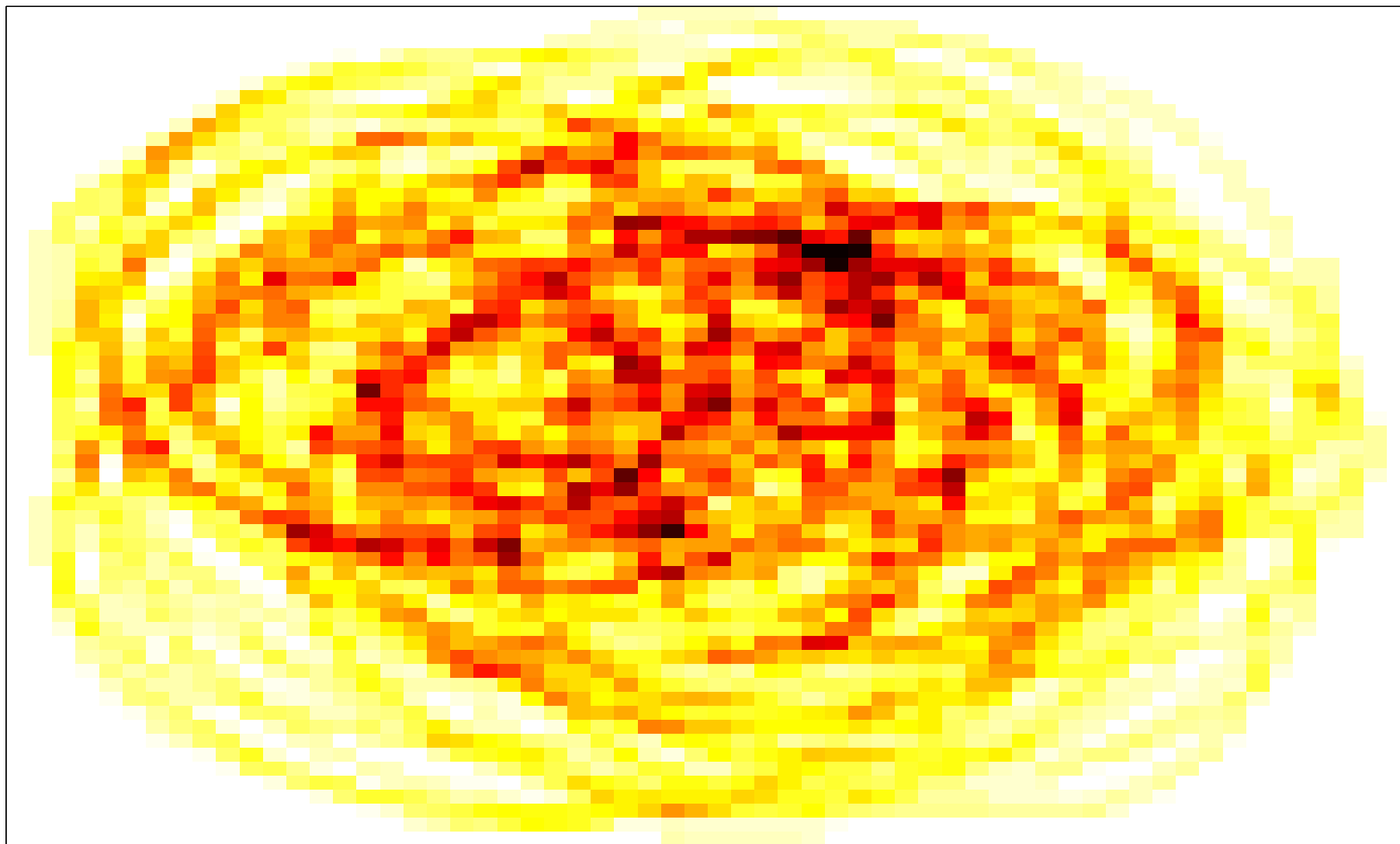


Fast Computation of Wasserstein Barycenters
International Conference on Machine Learning 2014

[CD'14]

Primal Descent on Regularized W

$$\min_{\mu \in Q \subset \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_{\gamma}(\mu, \nu_i)$$

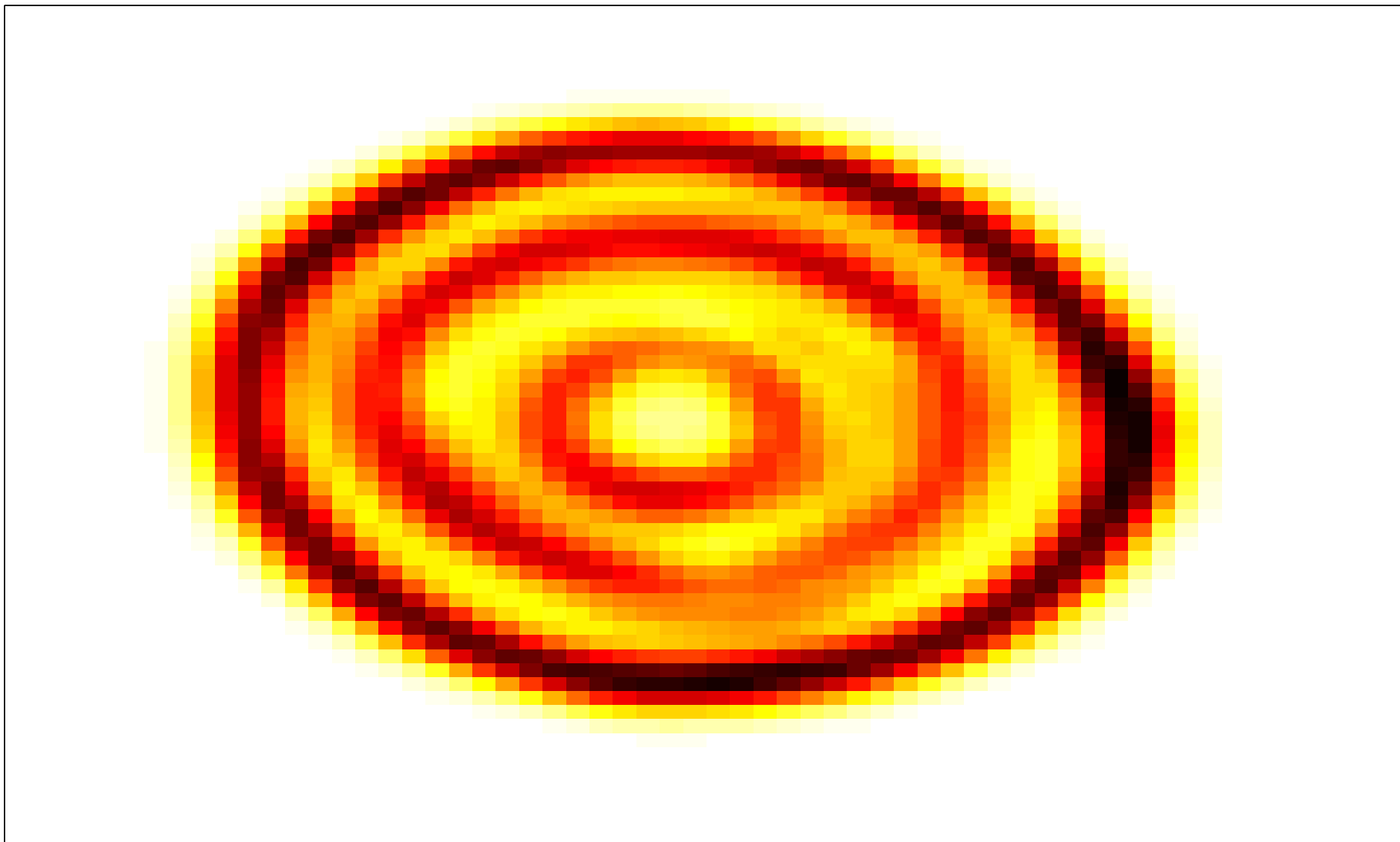


Fast Computation of Wasserstein Barycenters
International Conference on Machine Learning 2014

[CD'14]

Primal Descent on Regularized W

$$\min_{\mu \in Q \subset \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_{\gamma}(\mu, \nu_i)$$



Fast Computation of Wasserstein Barycenters
International Conference on Machine Learning 2014

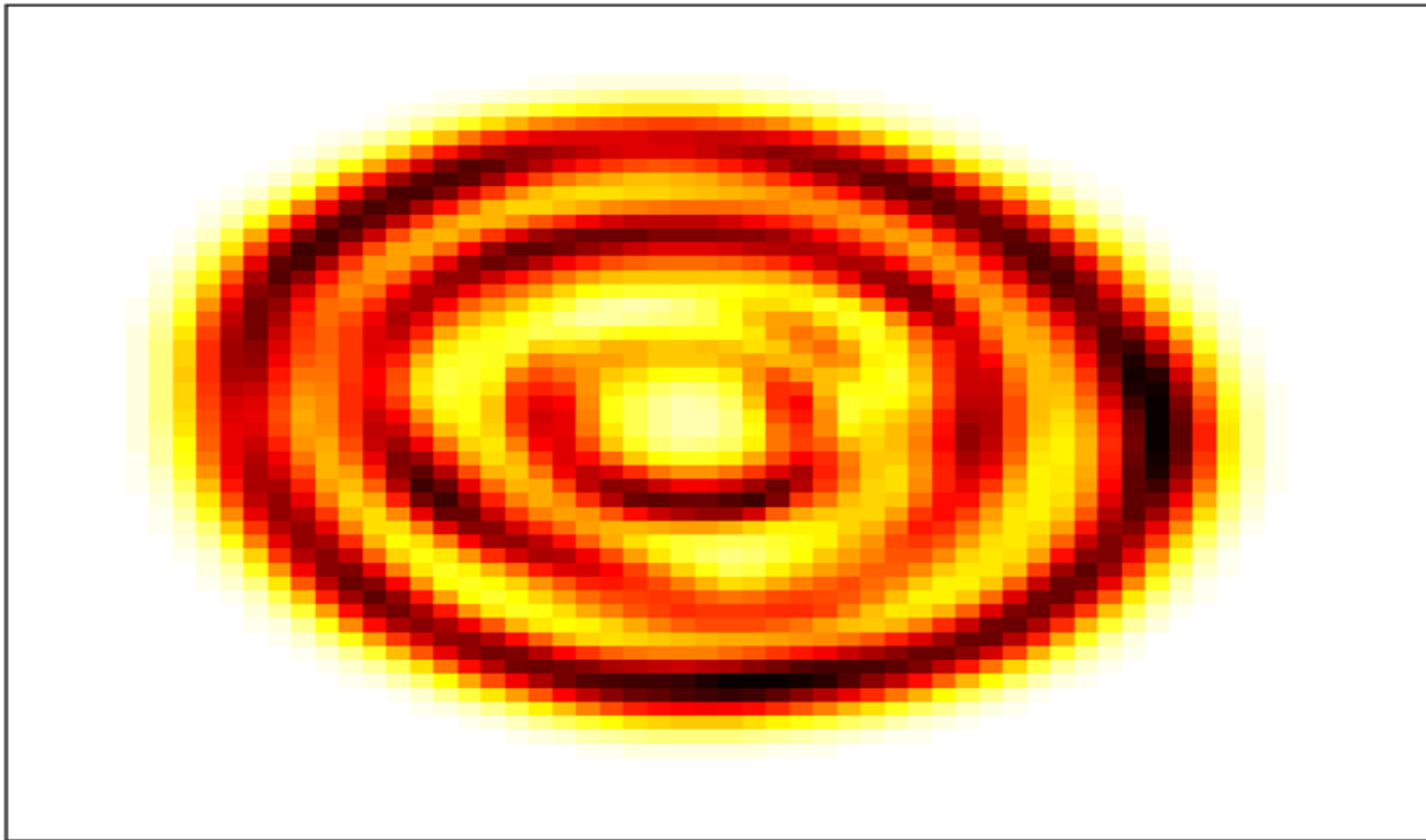
[CD'14]

Primal Descent on Algorithmic W

$$\min_{\mu \in Q \subset \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_{\mathbf{L}}(\mu, \nu_i)$$

Primal Descent on Algorithmic W

$$\min_{\mu \in Q \subset \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_{\mathbf{L}}(\mu, \nu_i)$$



Wasserstein Barycenter = KL Projections

$$\langle P, M_{\mathbf{x}\mathbf{y}} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P | \mathbf{K})$$

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} = [\mathbf{P}_1, \dots, \mathbf{P}_N] \\ \mathbf{P} \in \mathbf{C}_1 \cap \mathbf{C}_2}} \sum_{i=1}^N \lambda_i \mathbf{KL}(\mathbf{P}_i | \mathbf{K})$$

$$\mathbf{C}_1 = \{ \mathbf{P} | \exists \mathbf{a}, \forall i, \mathbf{P}_i \mathbf{1}_m = \mathbf{a} \}$$

$$\mathbf{C}_2 = \{ \mathbf{P} | \forall i, \mathbf{P}_i^T \mathbf{1}_n = \mathbf{b}_i \}$$

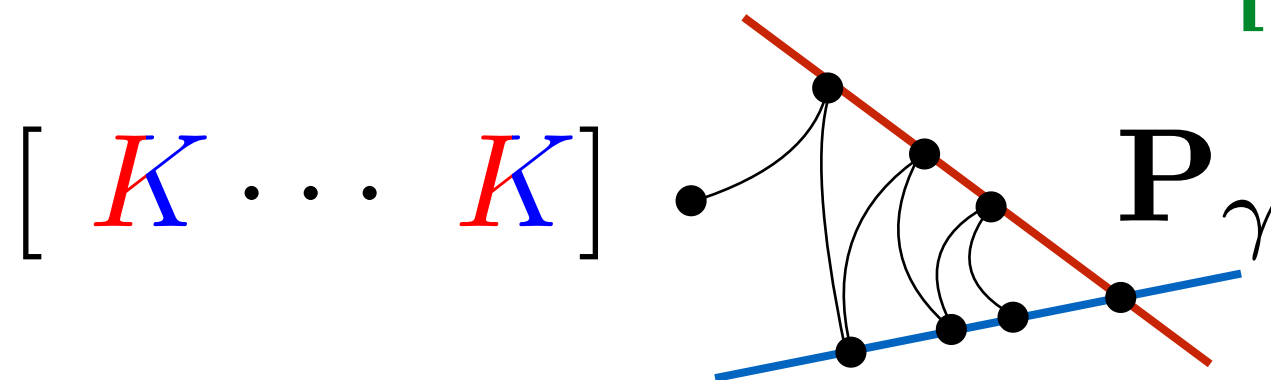
Wasserstein Barycenter = KL Projections

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} = [\mathbf{P}_1, \dots, \mathbf{P}_N] \\ \mathbf{P} \in \mathbf{C}_1 \cap \mathbf{C}_2}} \sum_{i=1}^N \lambda_i \text{KL}(\mathbf{P}_i | \mathbf{K})$$

$$\mathbf{C}_1 = \{ \mathbf{P} | \exists \mathbf{a}, \forall i, \mathbf{P}_i \mathbf{1}_m = \mathbf{a} \}$$

$$\mathbf{C}_2 = \{ \mathbf{P} | \forall i, \mathbf{P}_i^T \mathbf{1}_n = \mathbf{b}_i \}$$

[BCCNP'15]



Wasserstein Barycenter = KL Projections

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} = [\mathbf{P}_1, \dots, \mathbf{P}_N] \\ \mathbf{P} \in \mathcal{C}_1 \cap \mathcal{C}_2}} \sum_{i=1}^N \lambda_i \text{KL}(\mathbf{P}_i | \mathbf{K})$$

$$\mathcal{C}_1 = \{ \mathbf{P} | \exists \mathbf{a}, \forall i, \mathbf{P}_i \mathbf{1}_m = \mathbf{a} \}$$

$$\mathcal{C}_2 = \{ \mathbf{P} | \forall i, \mathbf{P}_i^T \mathbf{1}_n = \mathbf{b}_i \}$$

```
u=ones(size(B)); % d x N matrix
```

```
while not converged
```

```
    v=u.*(K'*(B./(K*u))); % 2(Nd^2) cost
```

```
    u=bsxfun(@times,u,exp(log(v)*weights))./v;
```

```
end
```

```
a=mean(v,2);
```

[BCCNP'15]

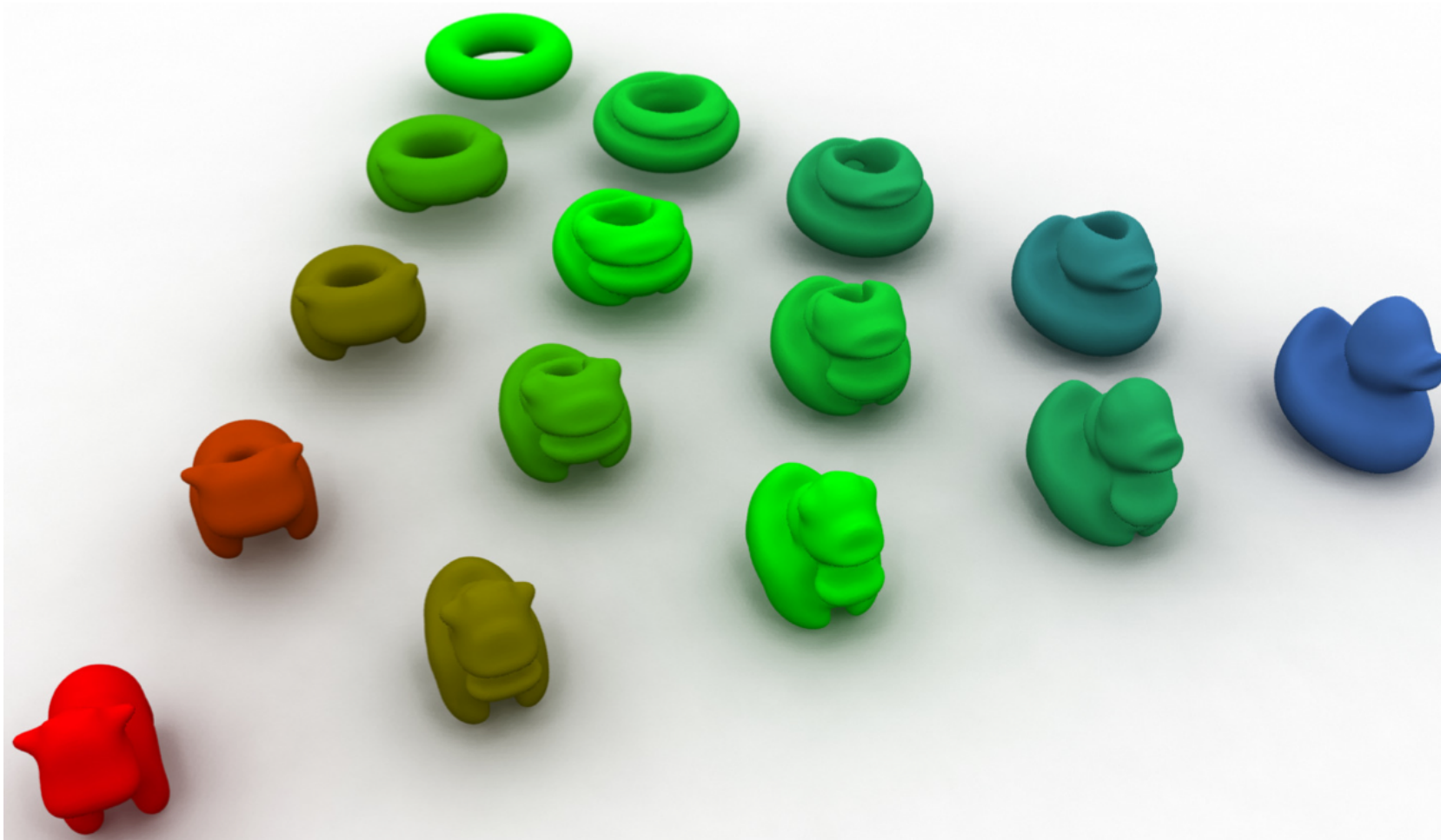
*Iterative Bregman Projections for
Regularized Transportation Problems*
SIAM J. on Sci. Comp. 2015

Application: Graphics



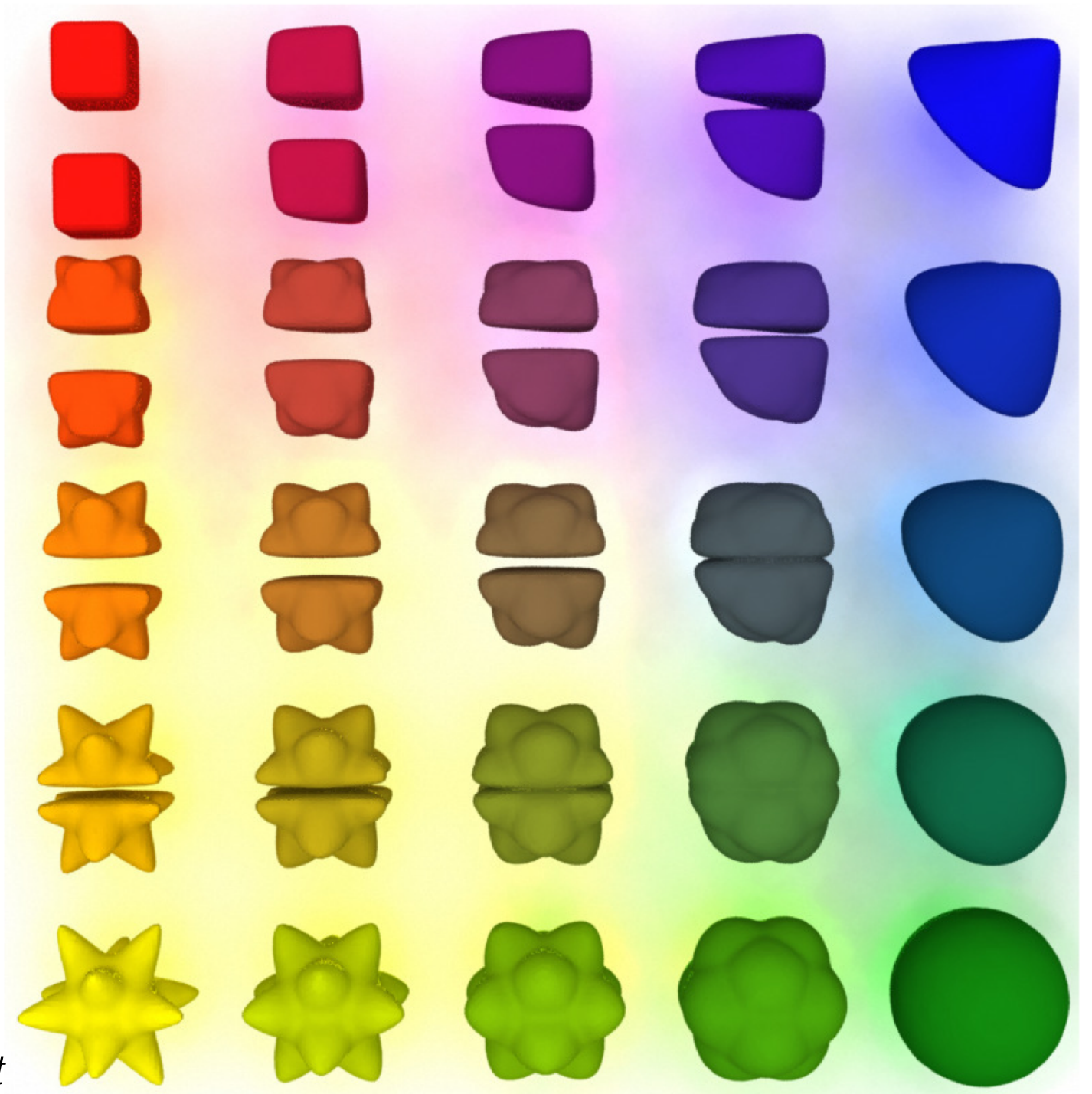
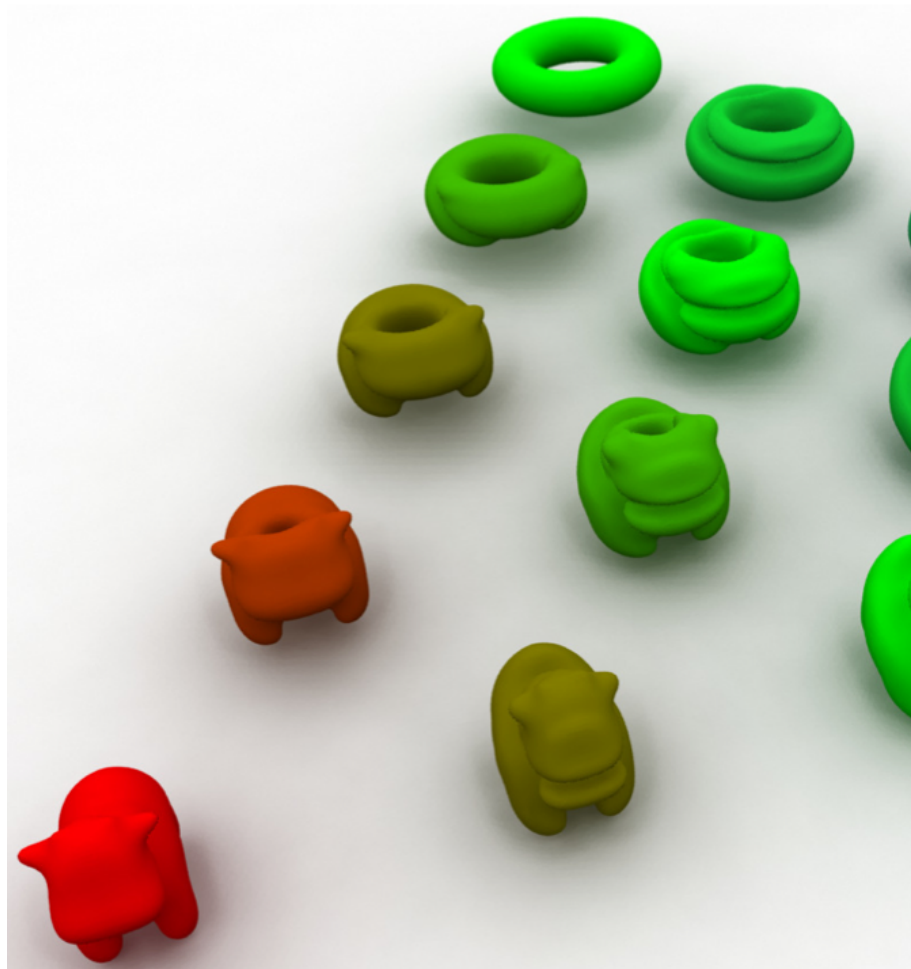
*Convolutional Wasserstein Distances: Efficient
Optimal Transportation on Geometric Domains,*
SIGGRAPH'15 [S..C..'15]

Application: Graphics



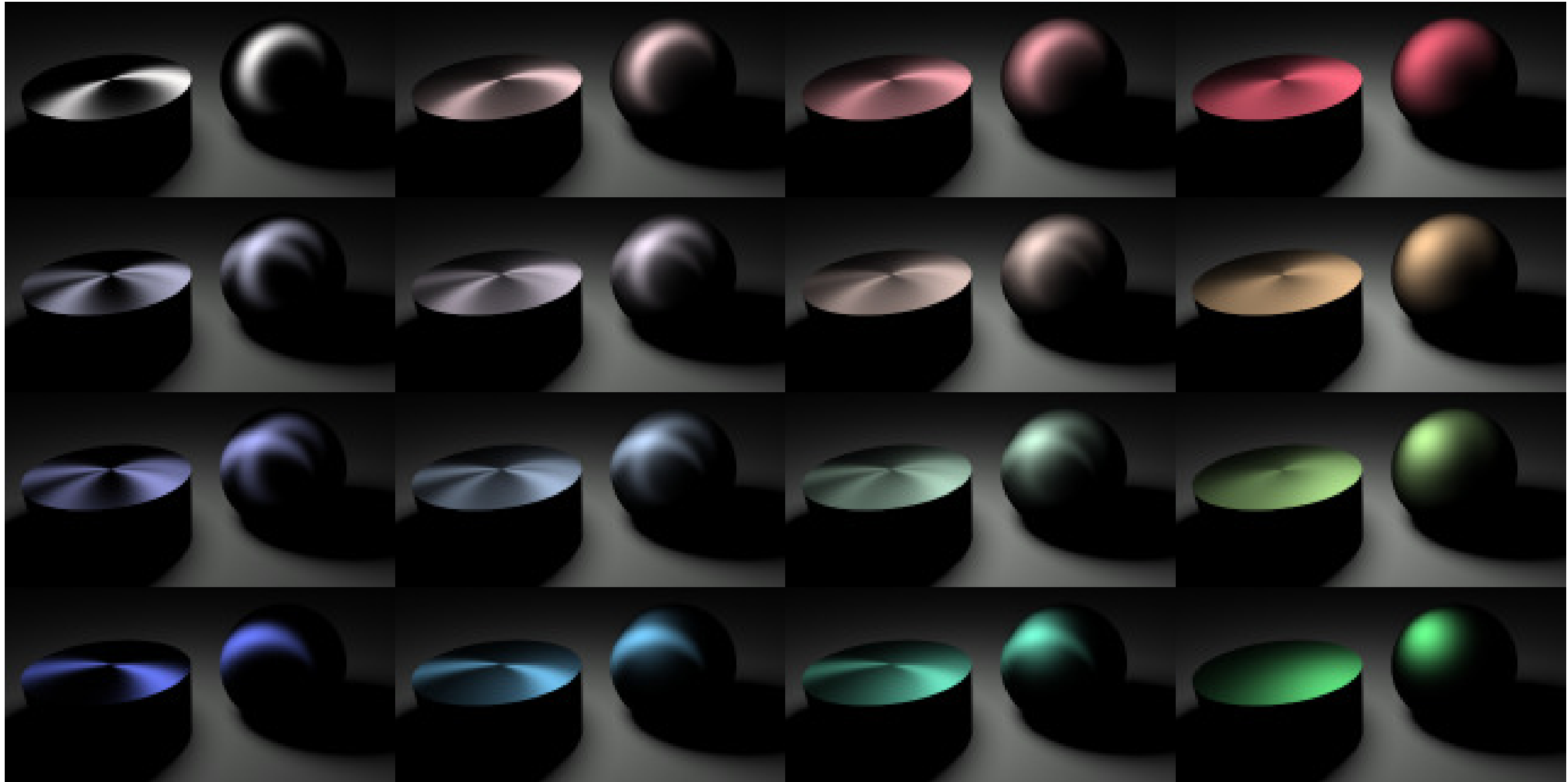
*Convolutional Wasserstein Distances: Efficient
Optimal Transportation on Geometric Domains,*
SIGGRAPH'15 [S..C..'15]

Application: Graphics



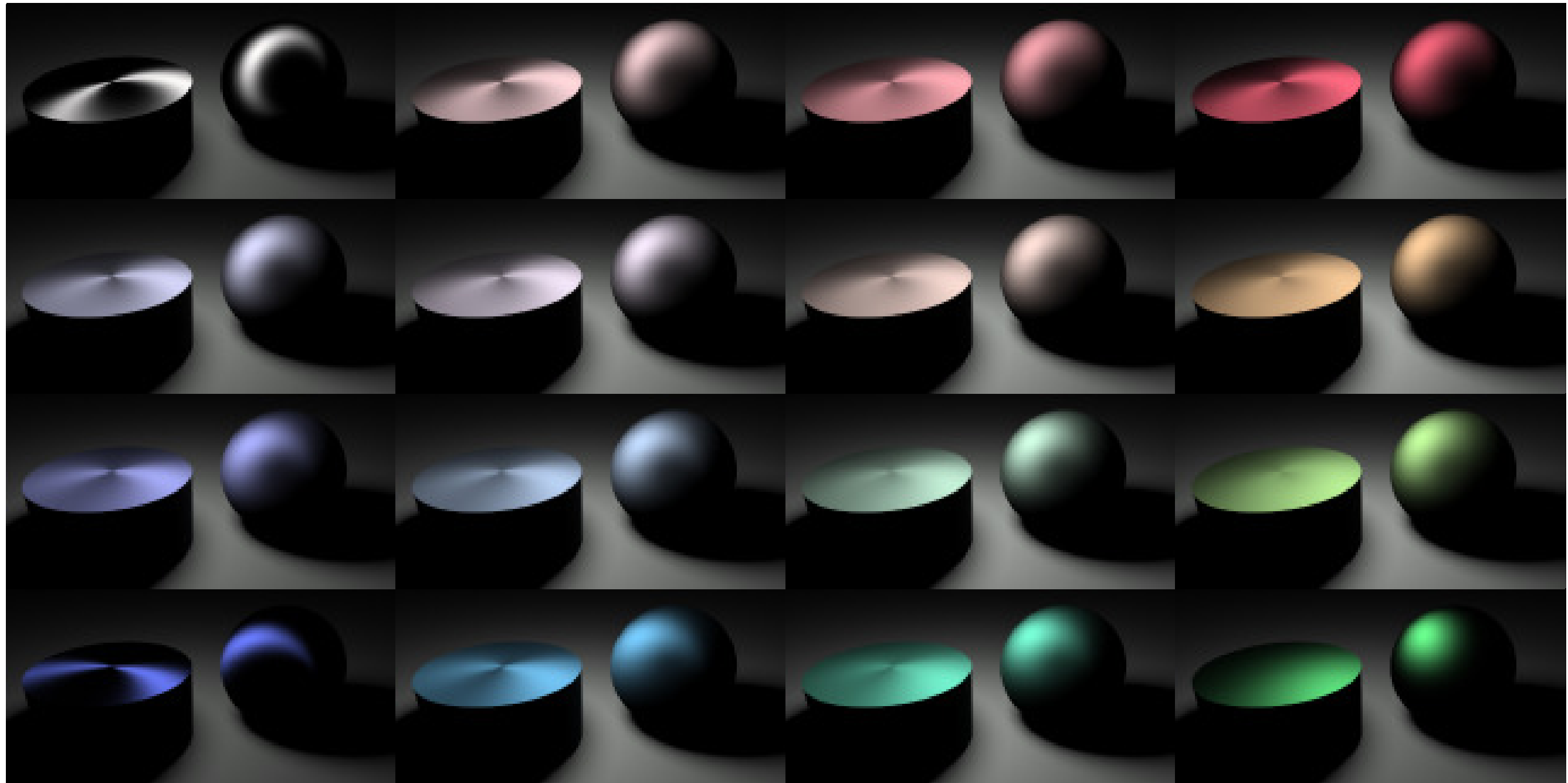
*Convolutional Wasserstein Distances: Efficient
Optimal Transportation on Geometric Domains,*
SIGGRAPH'15 [S..C..'15]

Application: Graphics



*Convolutional Wasserstein Distances: Efficient
Optimal Transportation on Geometric Domains,*
SIGGRAPH'15 [S..C..'15]

Application: Graphics



*Convolutional Wasserstein Distances: Efficient
Optimal Transportation on Geometric Domains,*
SIGGRAPH'15 [S..C..'15]

Inverse Wasserstein Problems

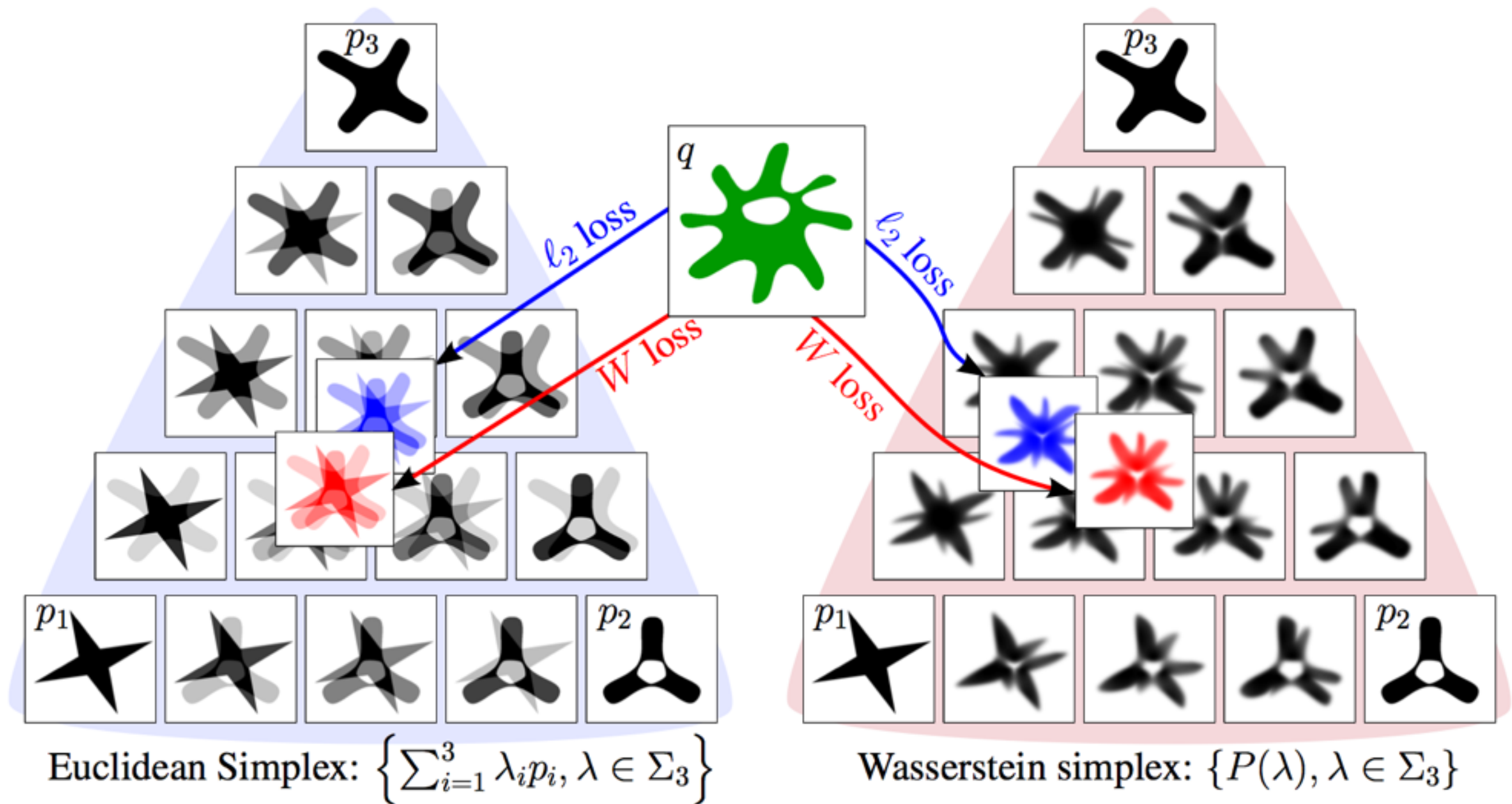
- consider Barycenter operator:

$$\mathbf{b}(\lambda) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$

- address now **Wasserstein inverse problems**:

Given \mathbf{a} , find $\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \text{Loss}(\mathbf{a}, \mathbf{b}(\lambda))$

The Wasserstein Simplex



Barycenters = Fixed Points

Prop. [BCCNP'15] Consider $\mathbf{B} \in \Sigma_d^N$ and let $\mathbf{U}_0 = \mathbf{1}_{d \times N}$, and then for $l \geq 0$:

$$\mathbf{b}^l \stackrel{\text{def}}{=} \exp \left(\log \left(K^T \mathbf{U}_l \right) \lambda \right) ; \begin{cases} \mathbf{V}_{l+1} \stackrel{\text{def}}{=} \frac{\mathbf{b}^l \mathbf{1}_N^T}{K^T \mathbf{U}_l}, \\ \mathbf{U}_{l+1} \stackrel{\text{def}}{=} \frac{\mathbf{B}}{K \mathbf{V}_{l+1}}. \end{cases}$$

Using Truncated Barycenters

- instead of using the exact barycenter

$$\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}(\lambda) \stackrel{\text{def}}{=} \text{Loss}(\textcolor{green}{a}, \textcolor{blue}{b}(\lambda))$$

- use instead the L-iterate barycenter

$$\operatorname{argmin}_{\lambda \in \Sigma_N} \mathcal{E}^{(L)}(\lambda) \stackrel{\text{def}}{=} \text{Loss}(\textcolor{green}{a}, \textcolor{blue}{b}^{(L)}(\lambda))$$

- Differentiate using **the chain rule**.

$$\nabla \mathcal{E}^{(L)}(\lambda) = [\partial \textcolor{blue}{b}^{(L)}]^T(\textcolor{brown}{g}), \quad \textcolor{brown}{g} \stackrel{\text{def}}{=} \nabla \text{Loss}(\textcolor{green}{a}, \cdot) \big|_{\textcolor{blue}{b}^{(L)}(\lambda)}.$$

Gradient / Barycenter Computation

```
function SINKHORN-DIFFERENTIATE( $(p_s)_{s=1}^S, q, \lambda$ )  
   $\forall s, b_s^{(0)} \leftarrow \mathbb{1}$   
   $(w, r) \leftarrow (0^S, 0^{S \times N})$   
  for  $\ell = 1, 2, \dots, L$  // Sinkhorn loop  
     $\forall s, \varphi_s^{(\ell)} \leftarrow K^\top \frac{p_s}{K b_s^{(\ell-1)}}$   
     $p \leftarrow \prod_s \left( \varphi_s^{(\ell)} \right)^{\lambda_s}$   
     $\forall s, b_s^{(\ell)} \leftarrow \frac{p}{\varphi_s^{(\ell)}}$   
     $g \leftarrow \nabla \mathcal{L}(p, q) \odot p$   
    for  $\ell = L, L-1, \dots, 1$  // Reverse loop  
       $\forall s, w_s \leftarrow w_s + \langle \log \varphi_s^{(\ell)}, g \rangle$   
       $\forall s, r_s \leftarrow -K^\top \left( K \left( \frac{\lambda_s g - r_s}{\varphi_s^{(\ell)}} \right) \odot \frac{p_s}{(K b_s^{(\ell-1)})^2} \right) \odot b_s^{(\ell-1)}$   
       $g \leftarrow \sum_s r_s$   
  return  $P^{(L)}(\lambda) \leftarrow p, \nabla \mathcal{E}_L(\lambda) \leftarrow w$ 
```

Application: Volume Reconstruction



Shape database
 (p_1, \dots, p_5)



Input shape q



Projection
 $P(\lambda)$



Iso-surface

*Wasserstein Barycentric Coordinates: Histogram
Regression using Optimal Transport, SIGGRAPH'16*

[BPC'16]

Application: Color Grading



Application: Color Grading



$$\lambda_0 = 0.03$$



$$\lambda_1 = 0.12$$



$$\lambda_2 = 0.40$$

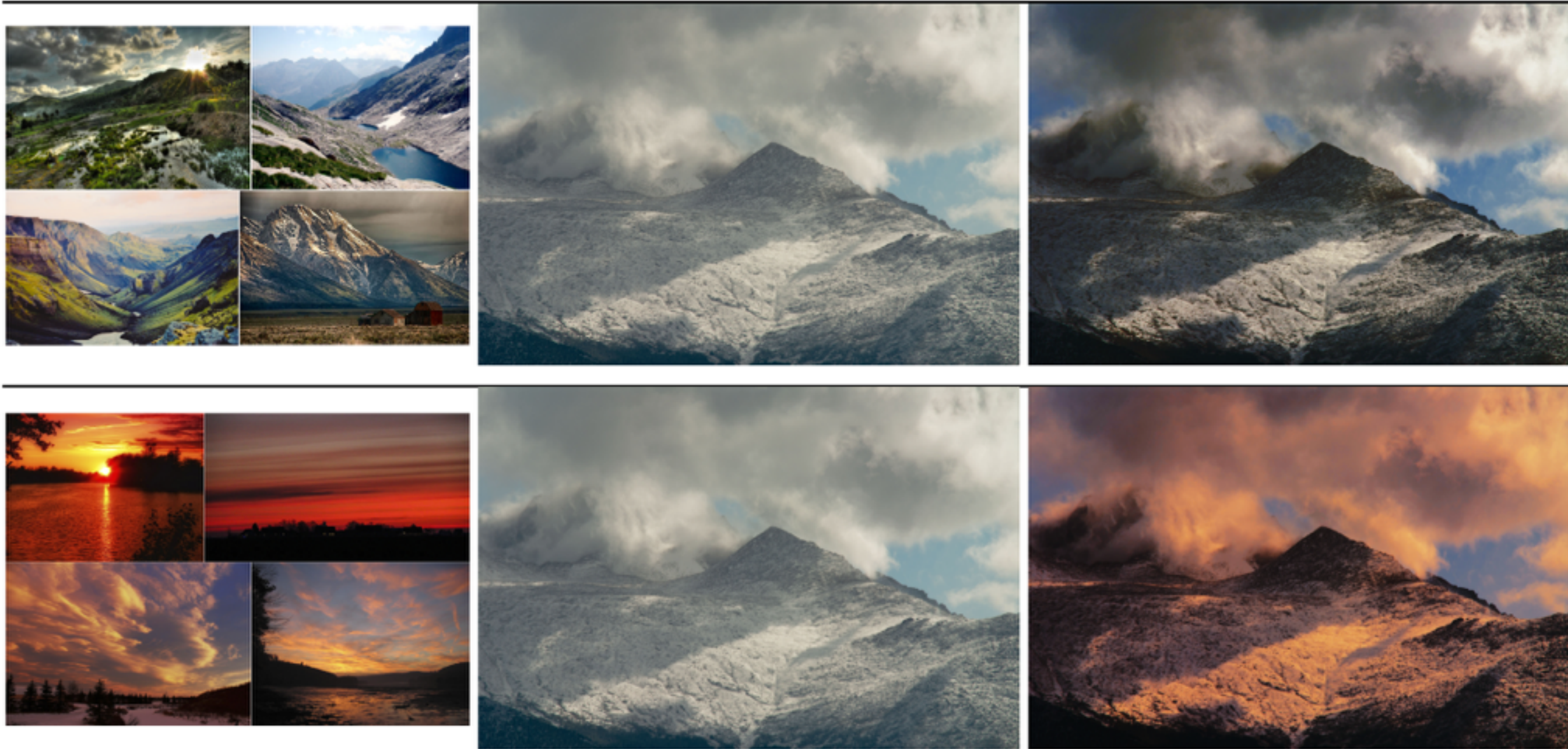


$$\lambda_3 = 0.43$$

Application: Color Grading



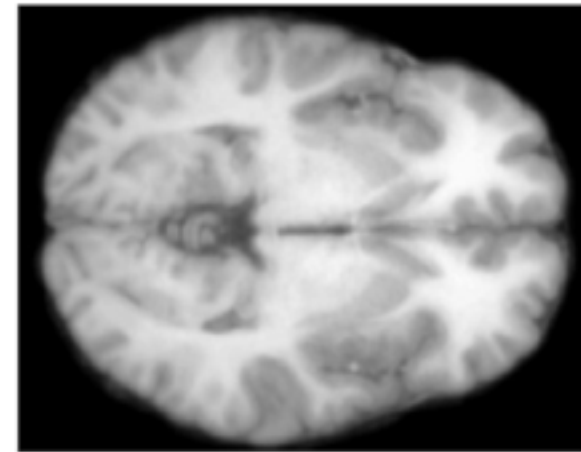
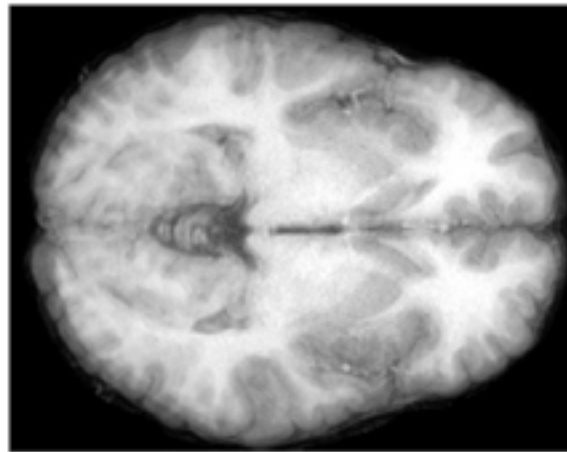
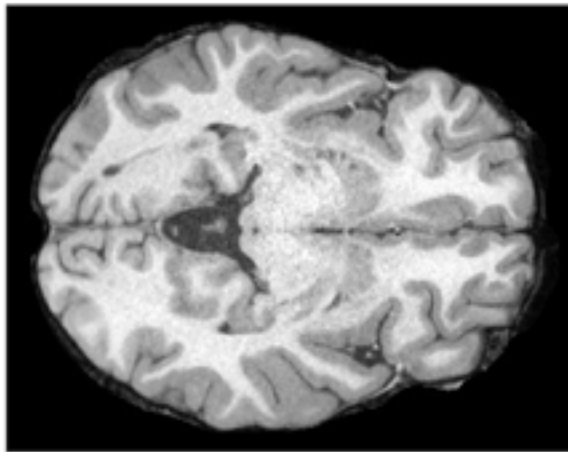
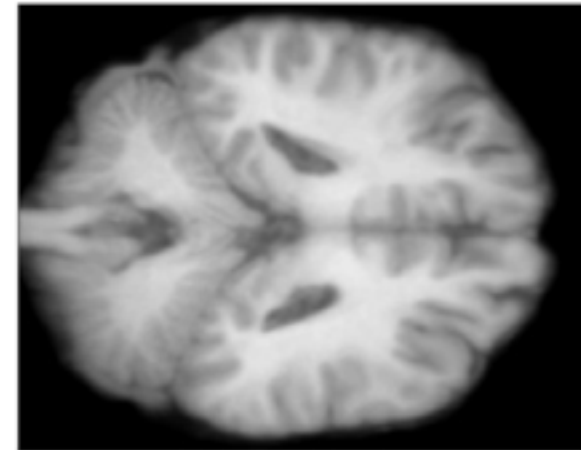
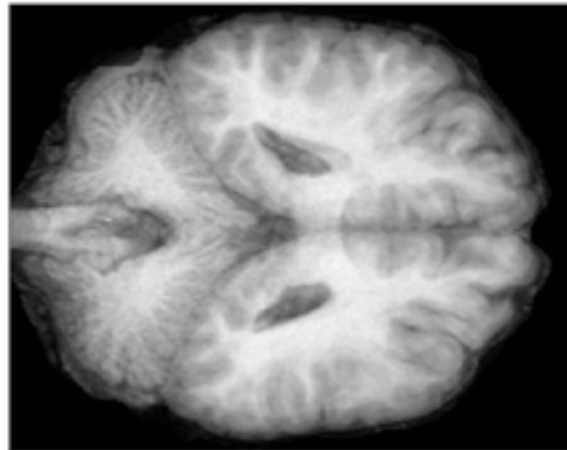
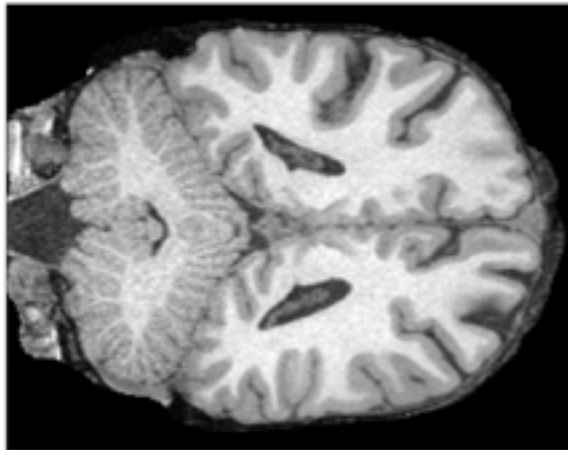
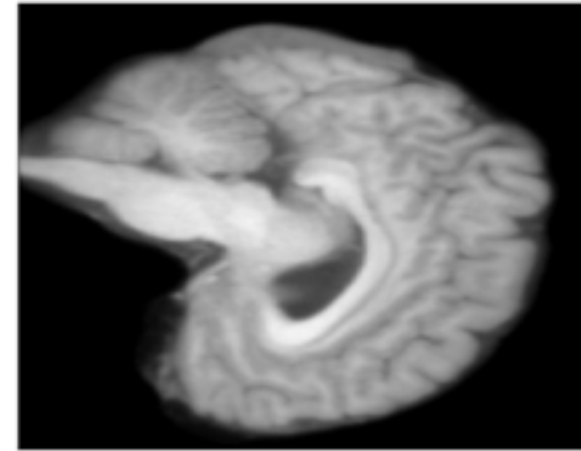
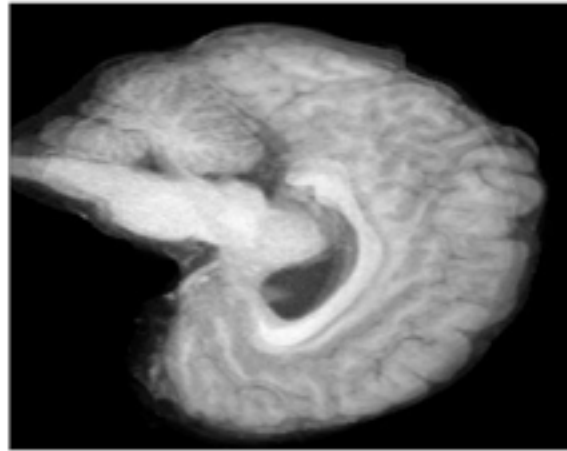
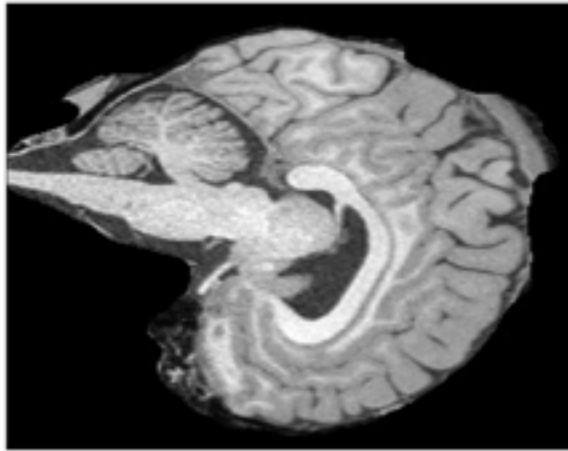
Application: Color Grading



*Wasserstein Barycentric Coordinates: Histogram
Regression using Optimal Transport, **SIGGRAPH'16***

[BPC'16]

Application: Brain Mapping



Original

Euclidean
projection

Wasserstein
projection

To conclude

- *Entropy* regularization is a very effective way to get OT to work as a generic loss.
- Many recent extensions:
 - [Schmitzer'16]: fast multiscale approaches
 - [ZFMAP'15] [CSPV'16]: Unbalanced transport
 - [SPKS'16] [PCS'16] extensions to *Gromov-W.*
 - [FCTR'15] Domain adaptation in ML