

# Generative Models and Optimal Transport

Marco Cuturi

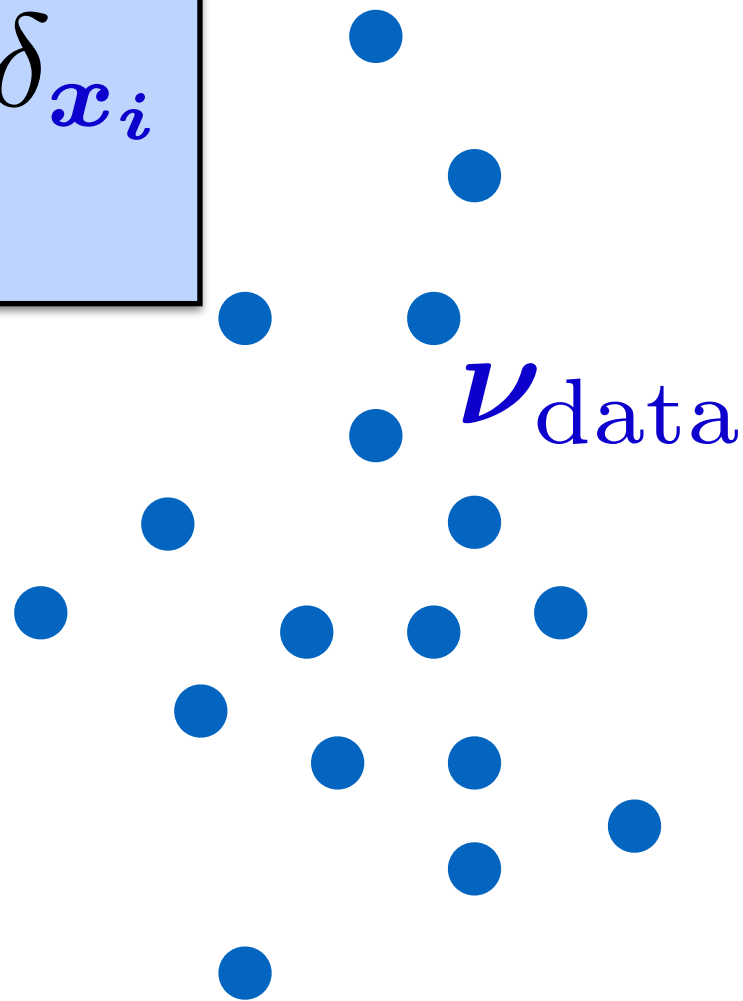


*Joint work / work in progress with*  
G. Peyré, A. Genevay (*ENS*), F. Bach (*INRIA*),  
G. Montavon, K-R Müller (*TU Berlin*)

# Statistics 0.1 : Density Fitting

We collect data

$$\nu_{\text{data}} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_i}$$





# Statistics 0.1 : Density Fitting

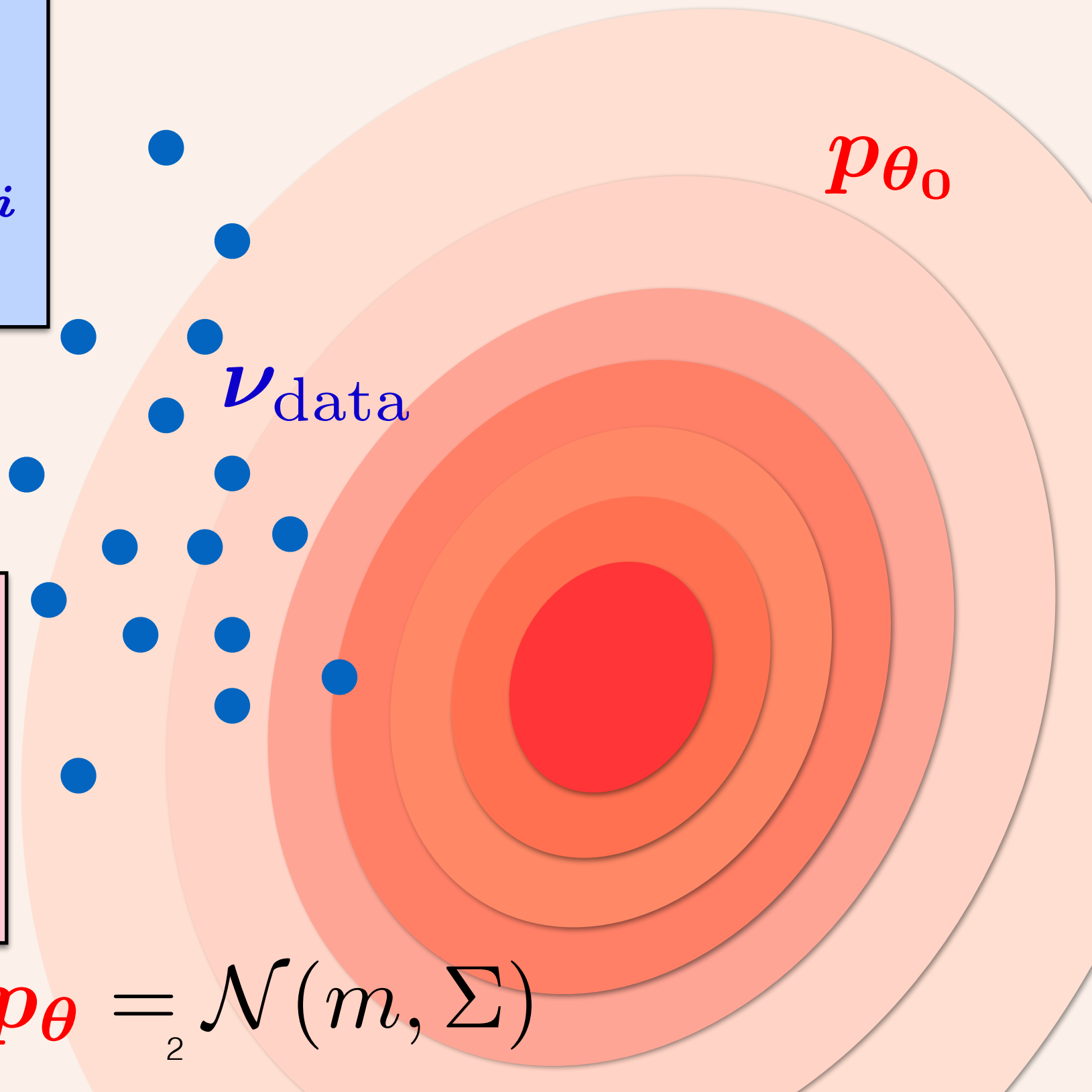
We collect data

$$\nu_{\text{data}} = \frac{1}{N} \sum_{i=1}^N \delta_{\mathbf{x}_i}$$

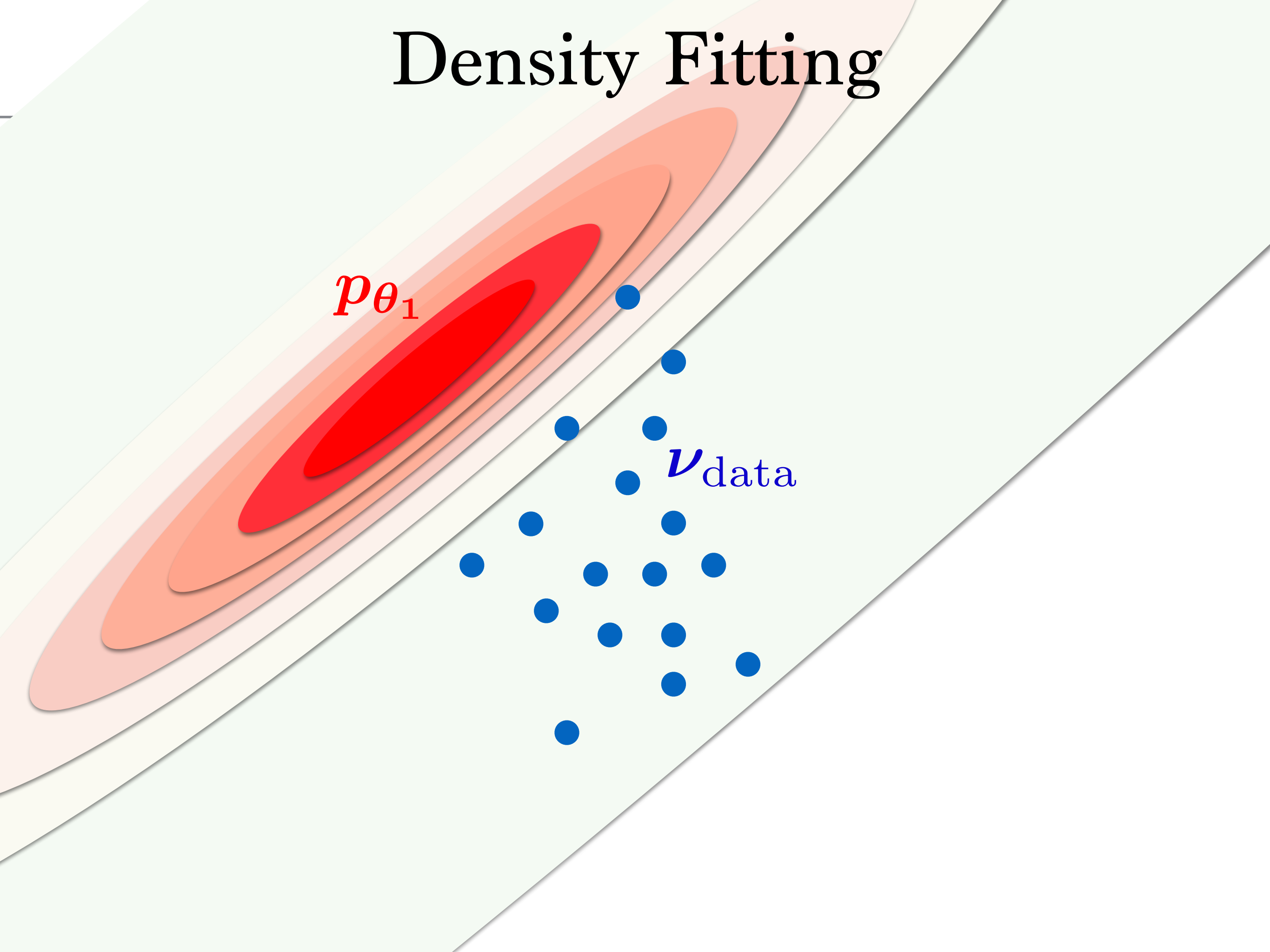
We fit a parametric family of densities

$$\{p_{\theta}, \theta \in \Theta\}$$

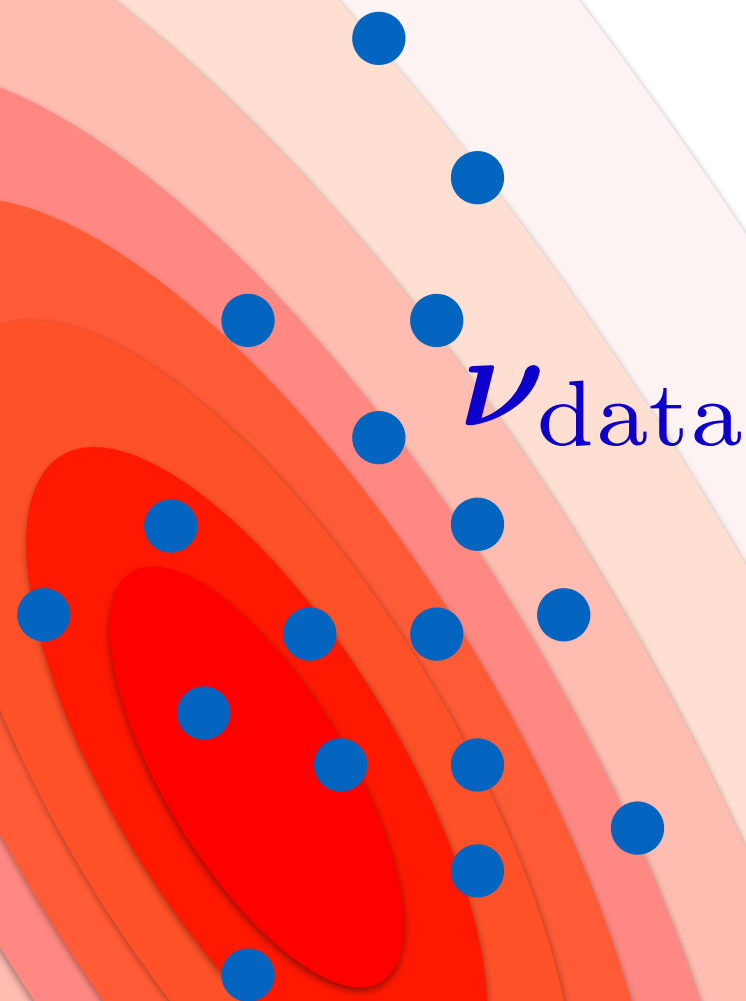
e.g.  $\theta = (m, \Sigma); p_{\theta} = \mathcal{N}_2(m, \Sigma)$



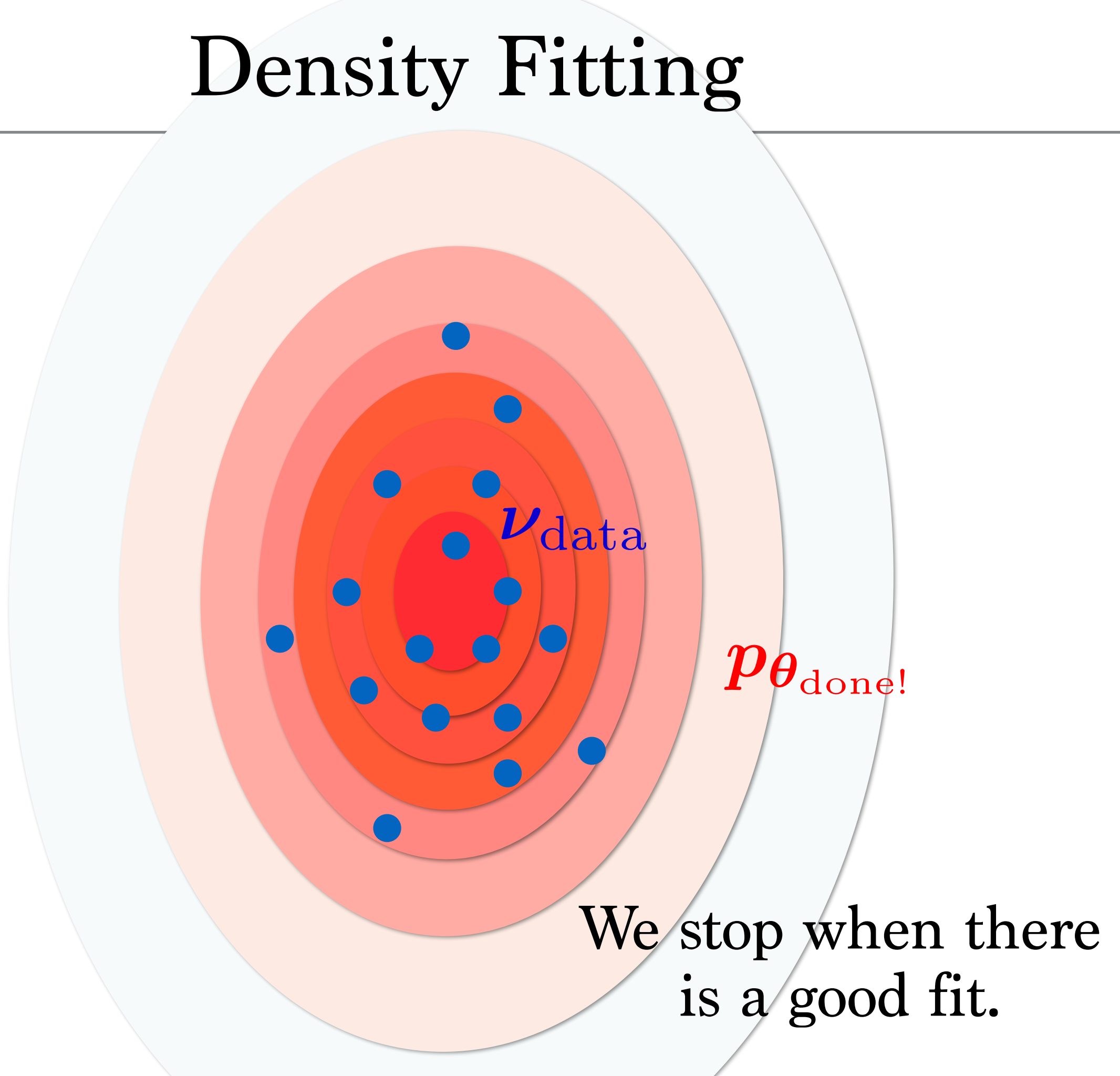
# Density Fitting



# Density Fitting



# Density Fitting



# Maximum Likelihood Estimation

## ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By *R. A. Fisher*, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma



$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i)$$

$p_{\theta}$  done!

$\nu_{\text{data}}$

# Maximum Likelihood Estimation

## ON AN ABSOLUTE CRITERION FOR FITTING FREQUENCY CURVES.

By *R. A. Fisher*, Gonville and Caius College, Cambridge.

1. IF we set ourselves the problem, in its frequent occurrence, of finding the arbitrary function of known form, which best suit a observations, we are met at the outset by an which appears to invalidate any results we ma

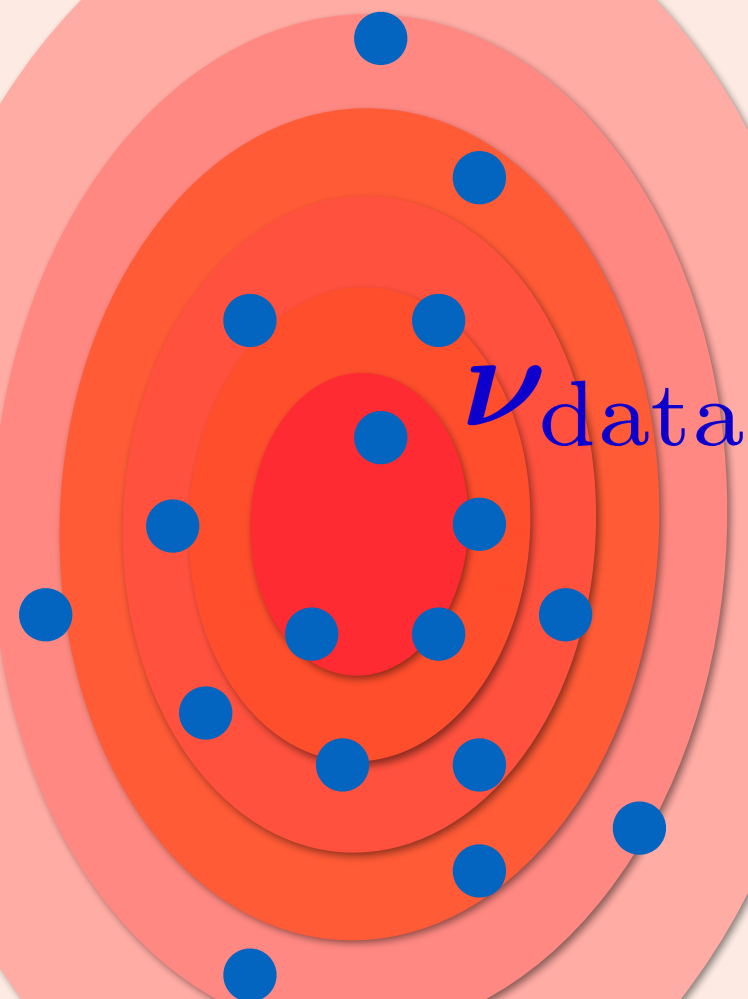


$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i)$$



$$\log 0 = -\infty$$

$p_{\theta}(x_i)$  must be  $> 0$

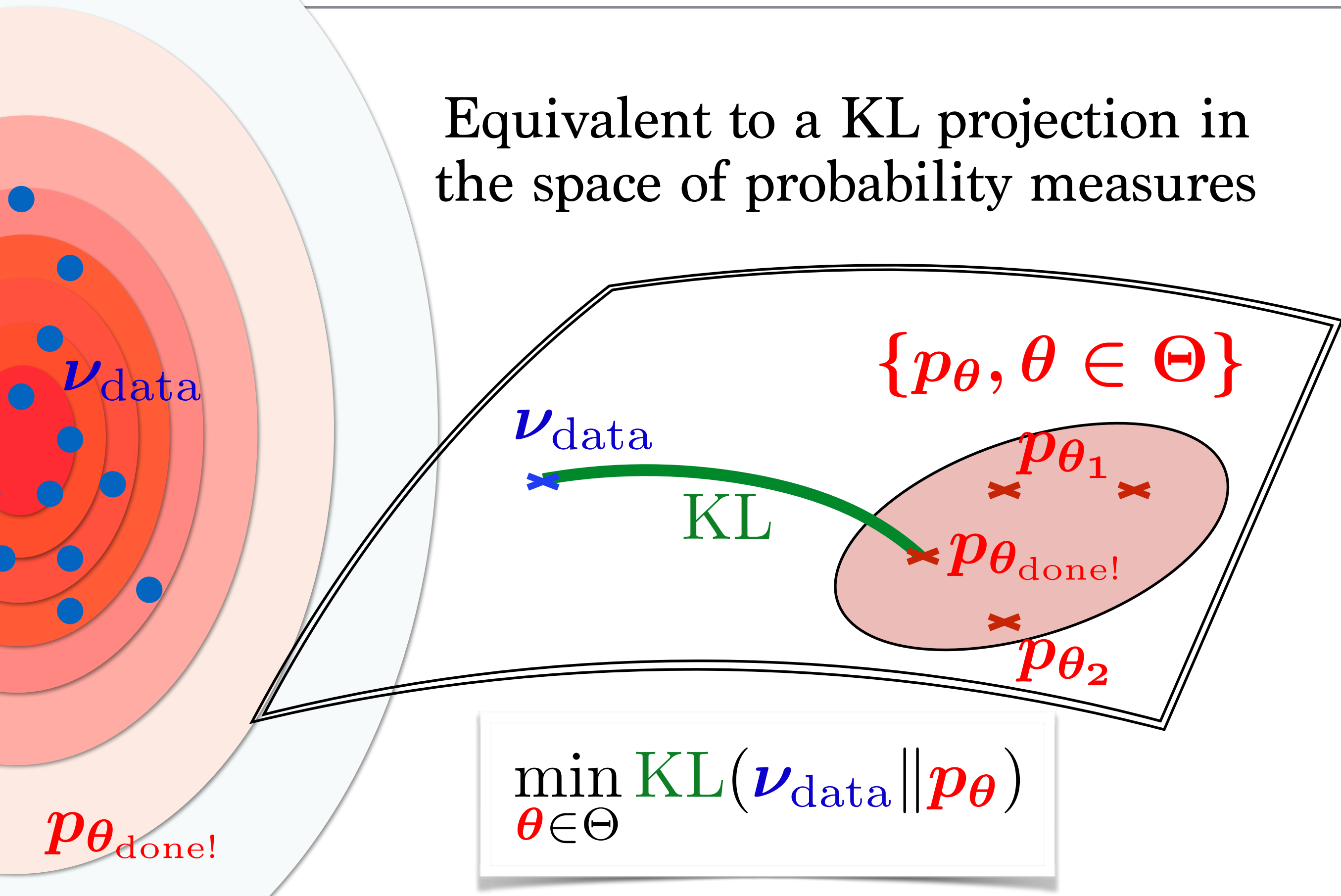


$p_{\theta}$  done!



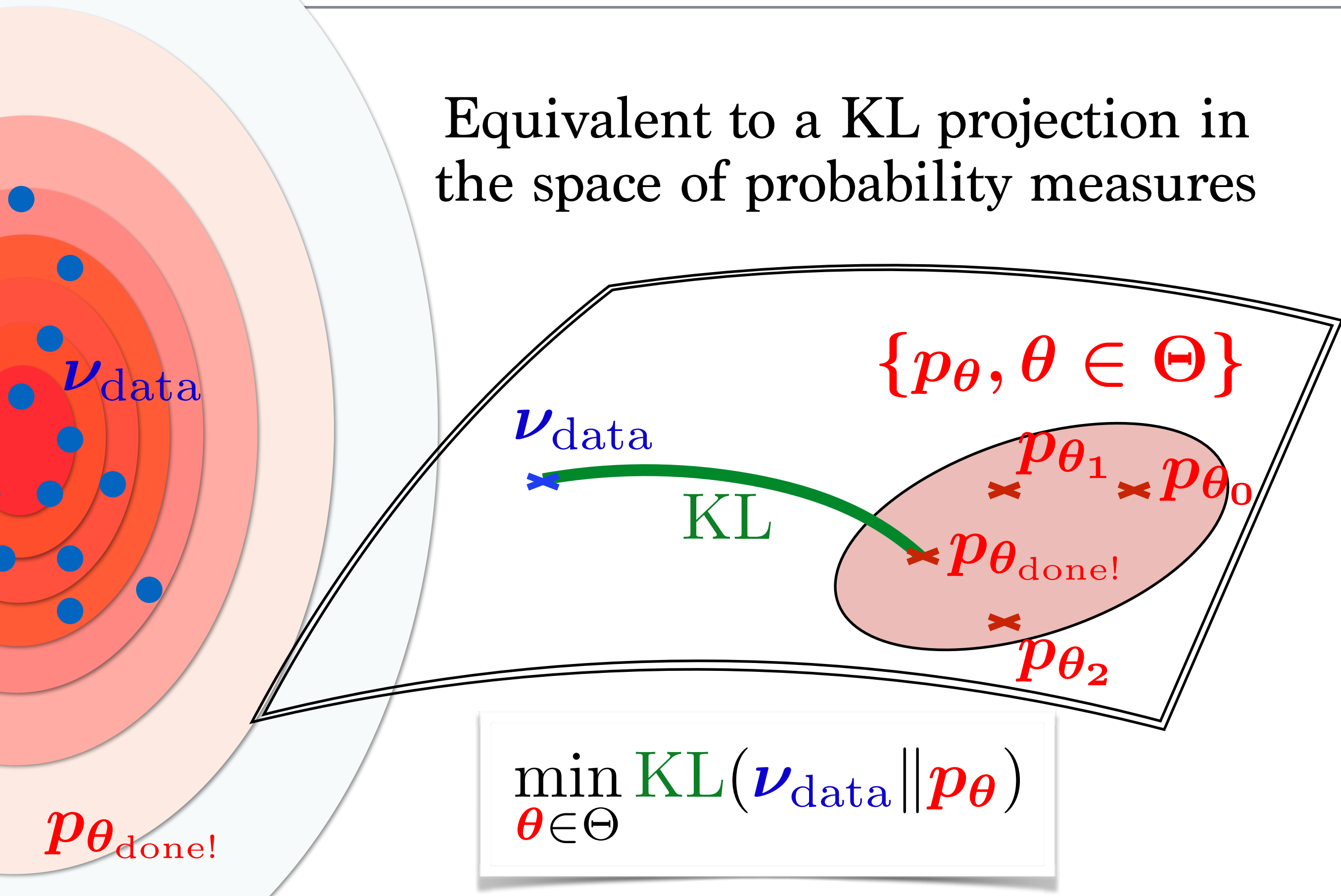
# Maximum Likelihood Estimation

Equivalent to a KL projection in the space of probability measures



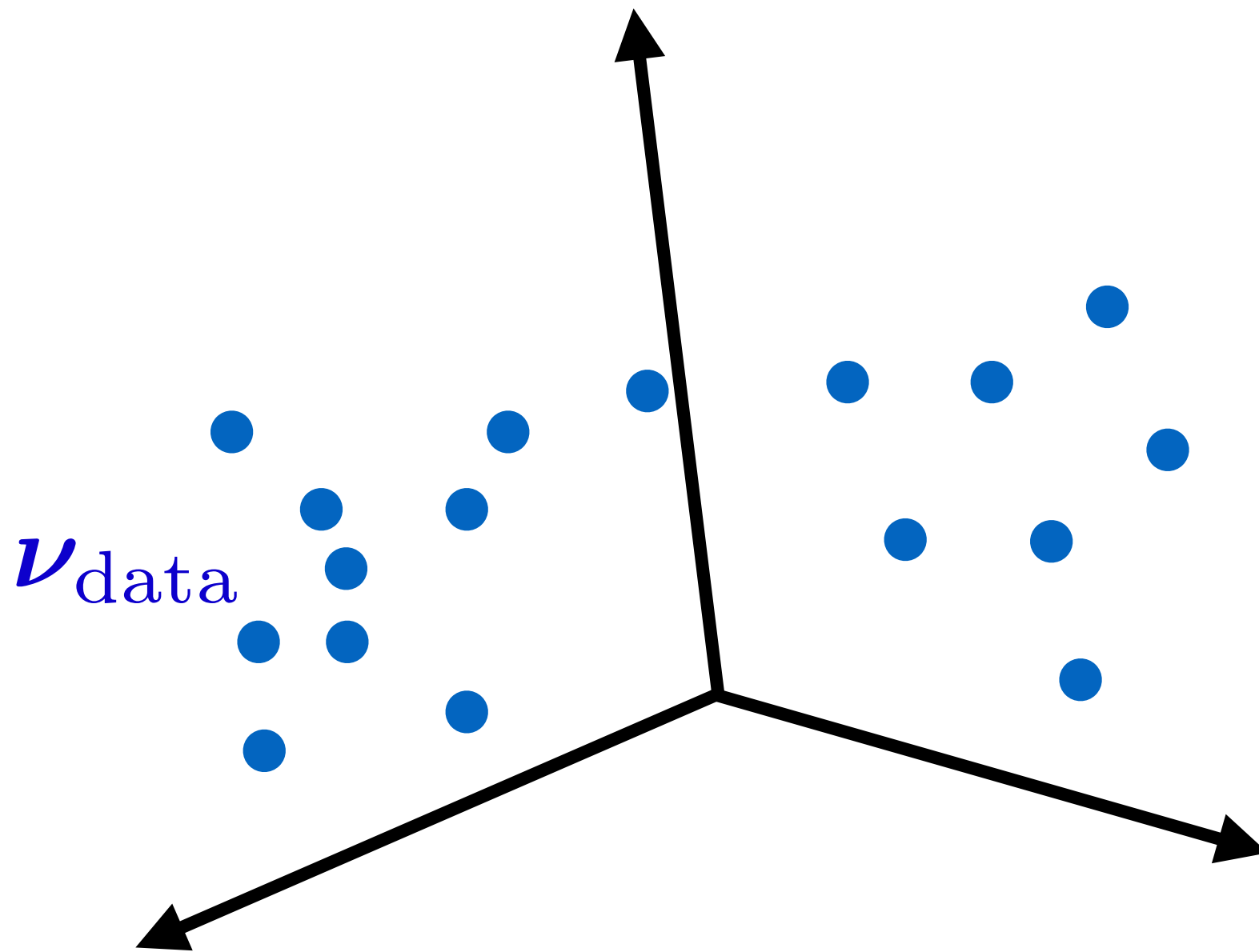
# Maximum Likelihood Estimation

Equivalent to a KL projection in the space of probability measures



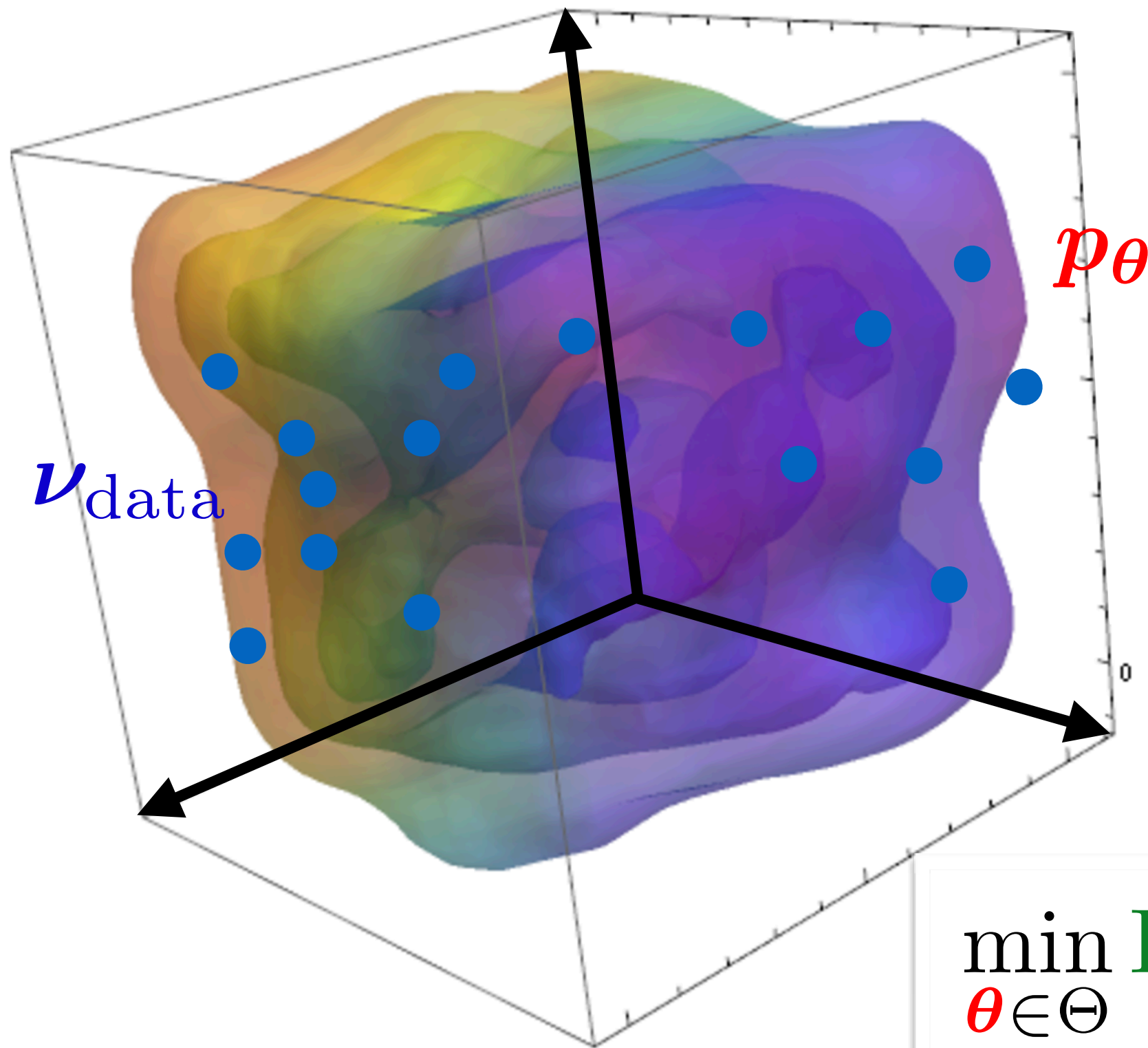


# In higher dimensional spaces...



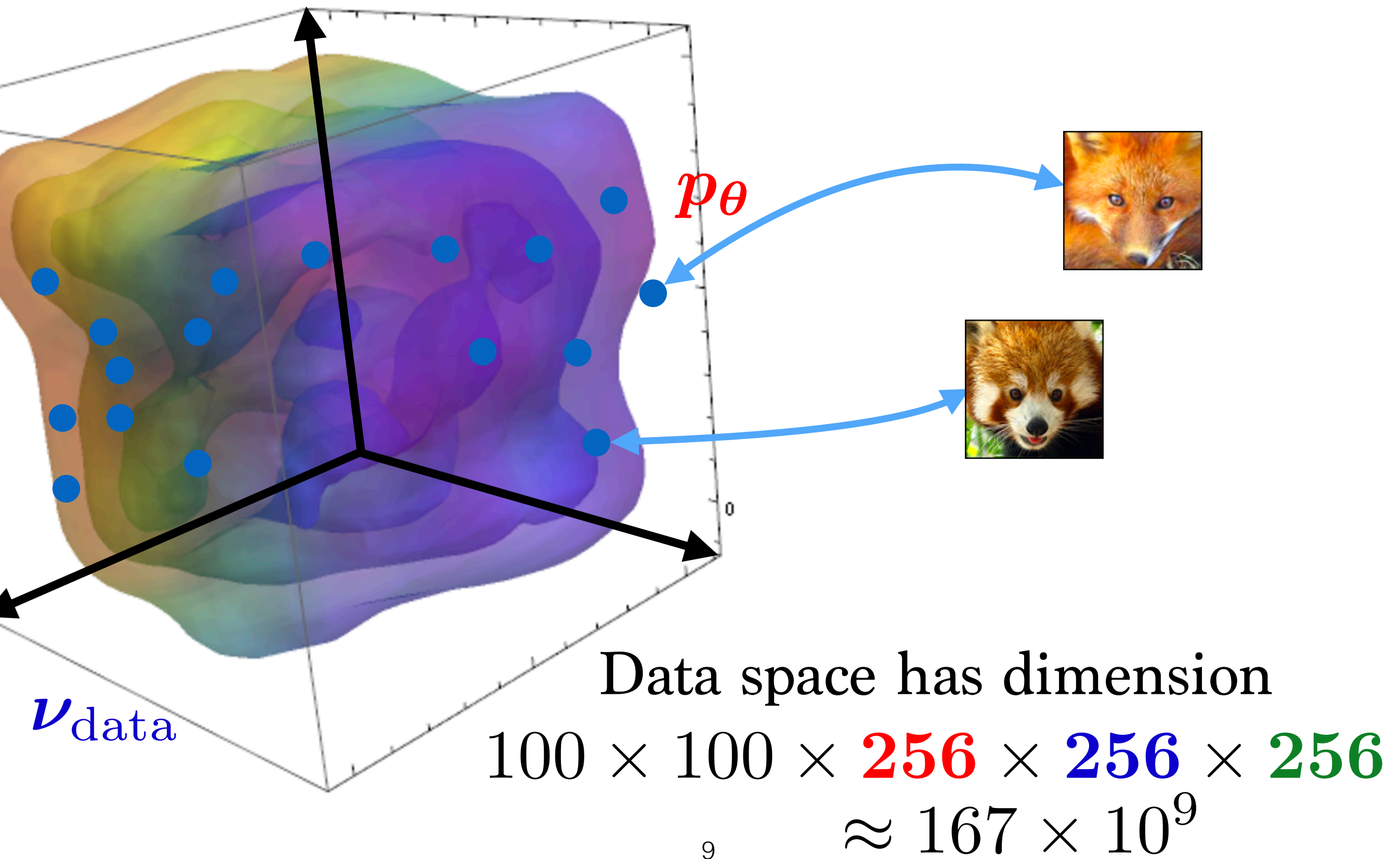
$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} || p_{\theta})$$

# In higher dimensional spaces...



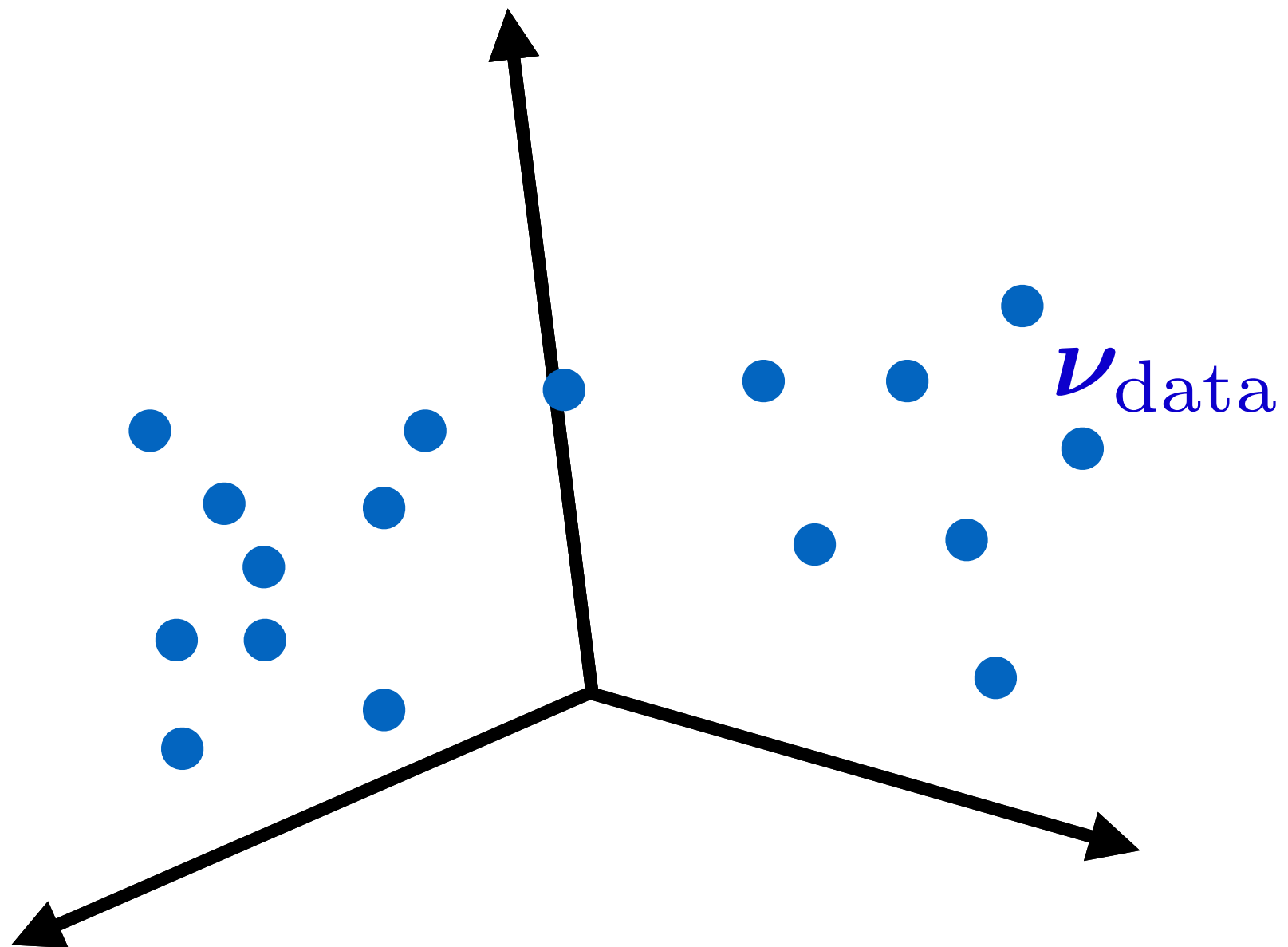
$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} || p_\theta)$$

# In higher dimensional spaces...

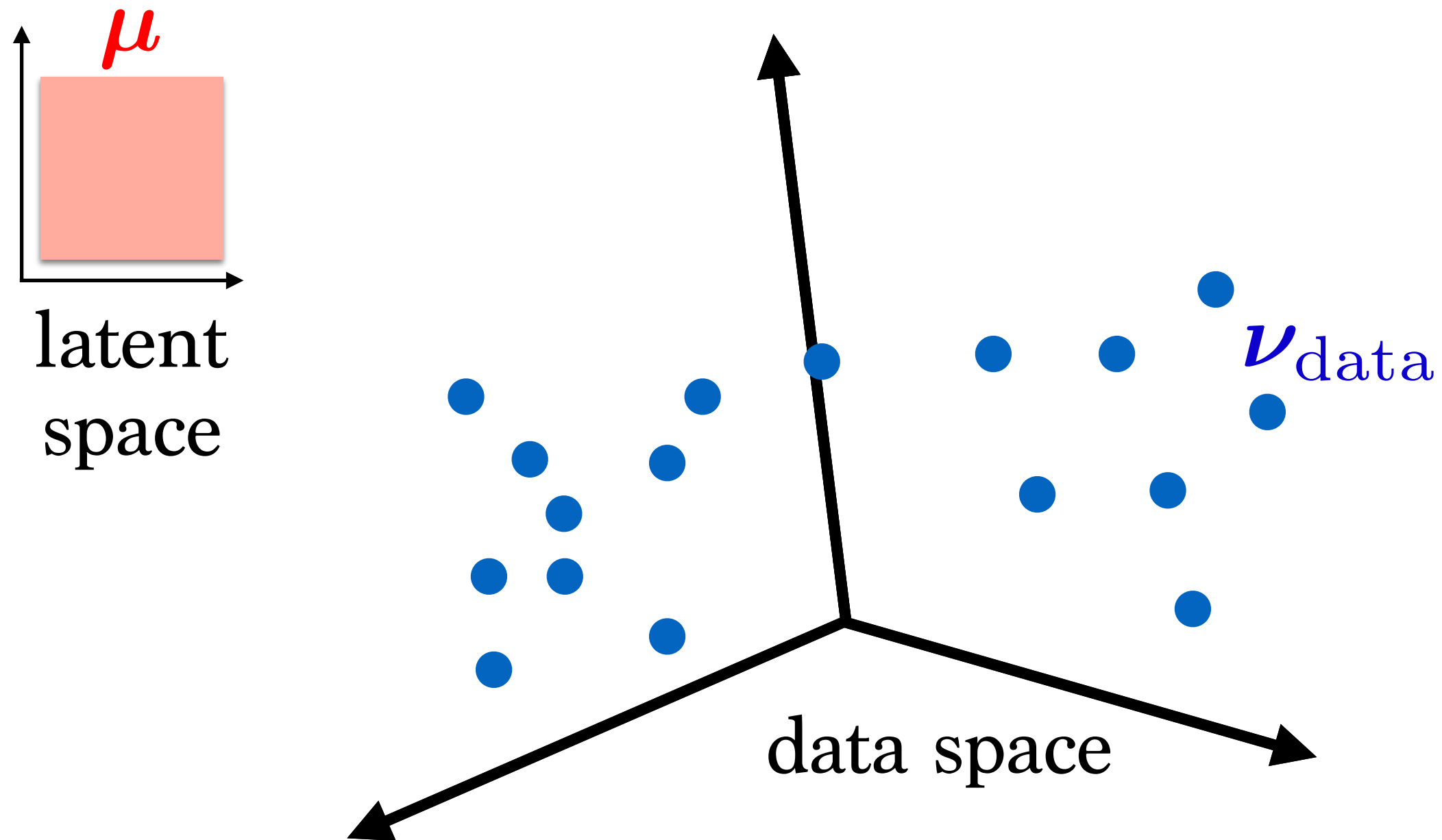


# Generative Models

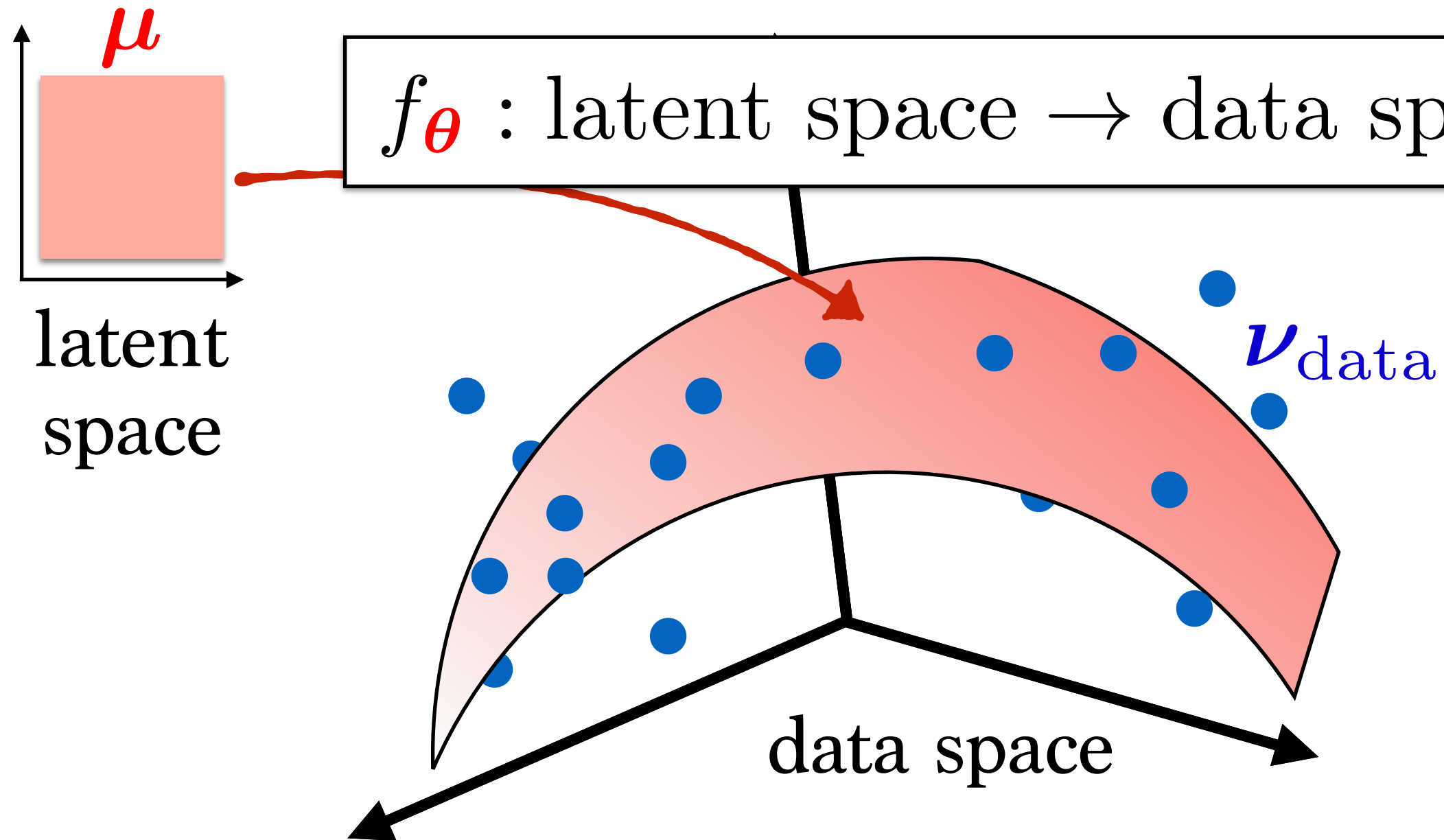
---



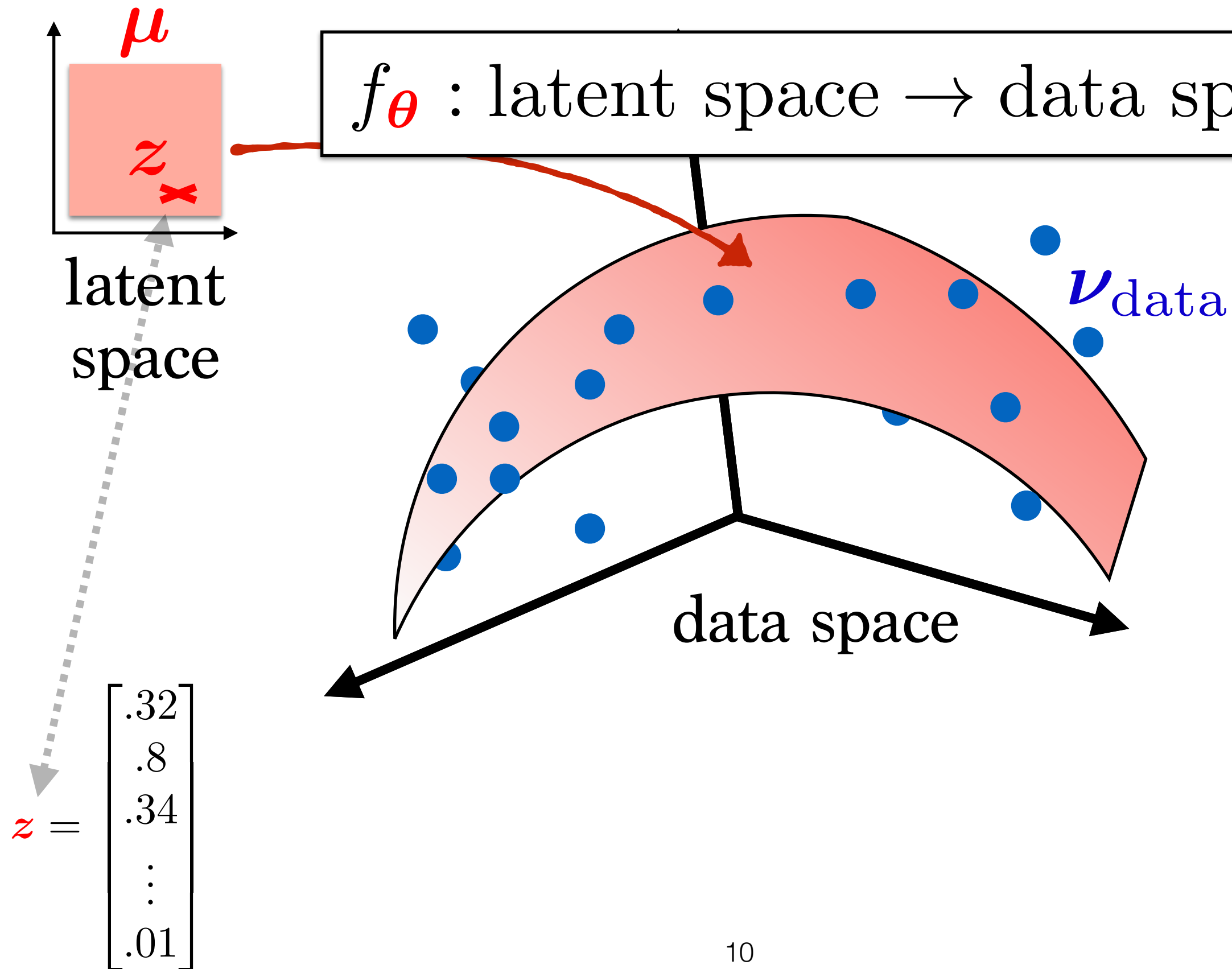
# Generative Models



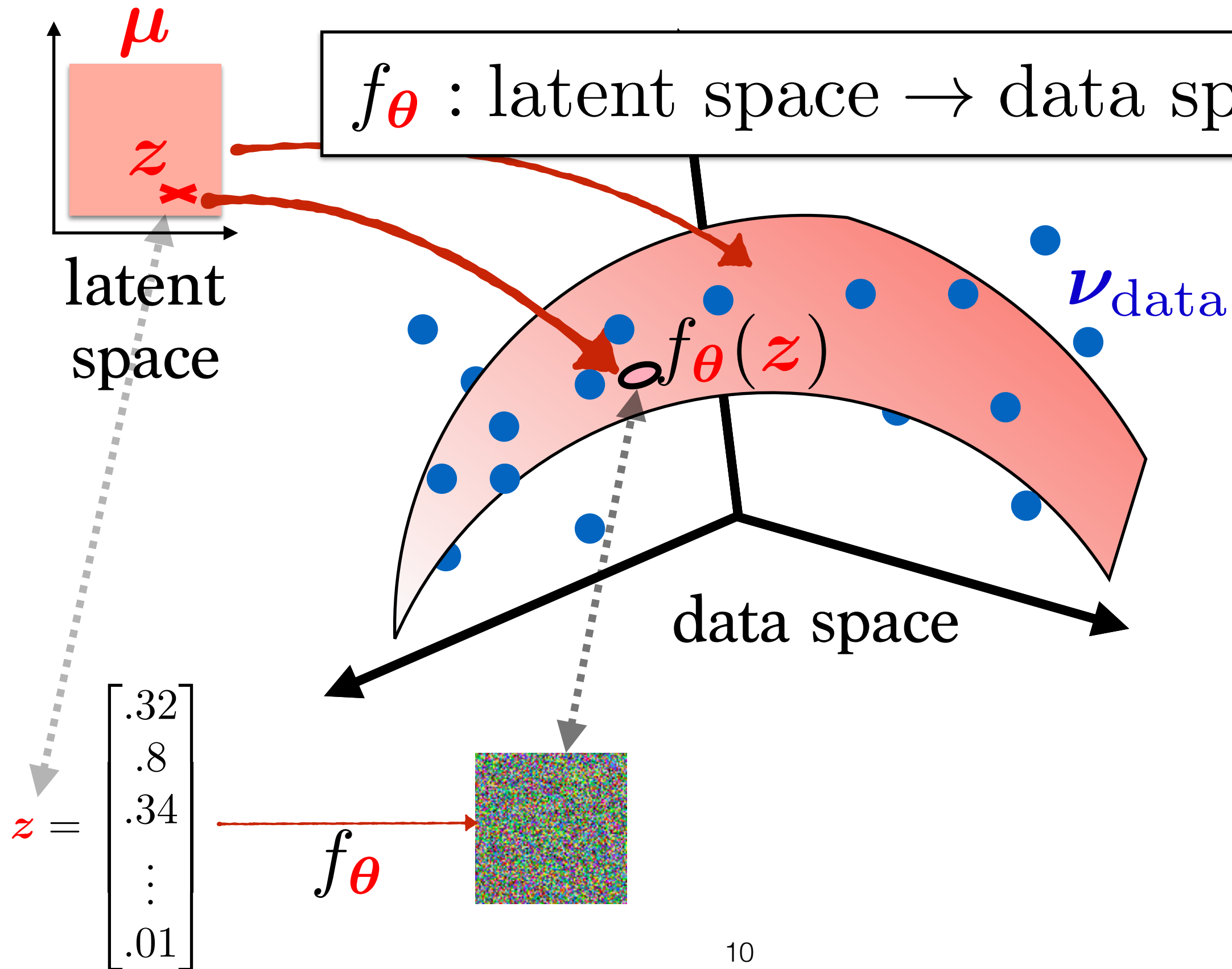
# Generative Models



# Generative Models

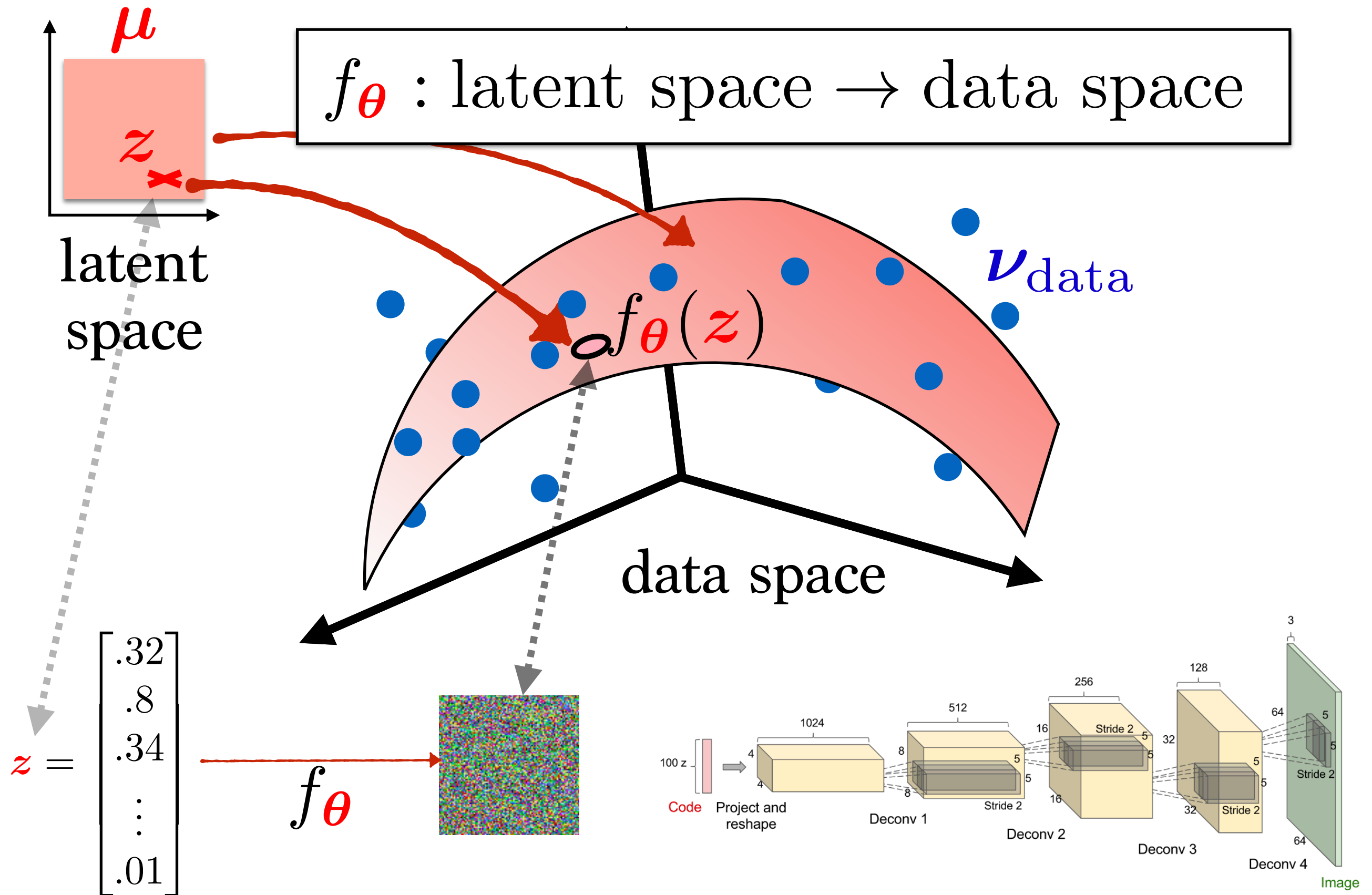


# Generative Models

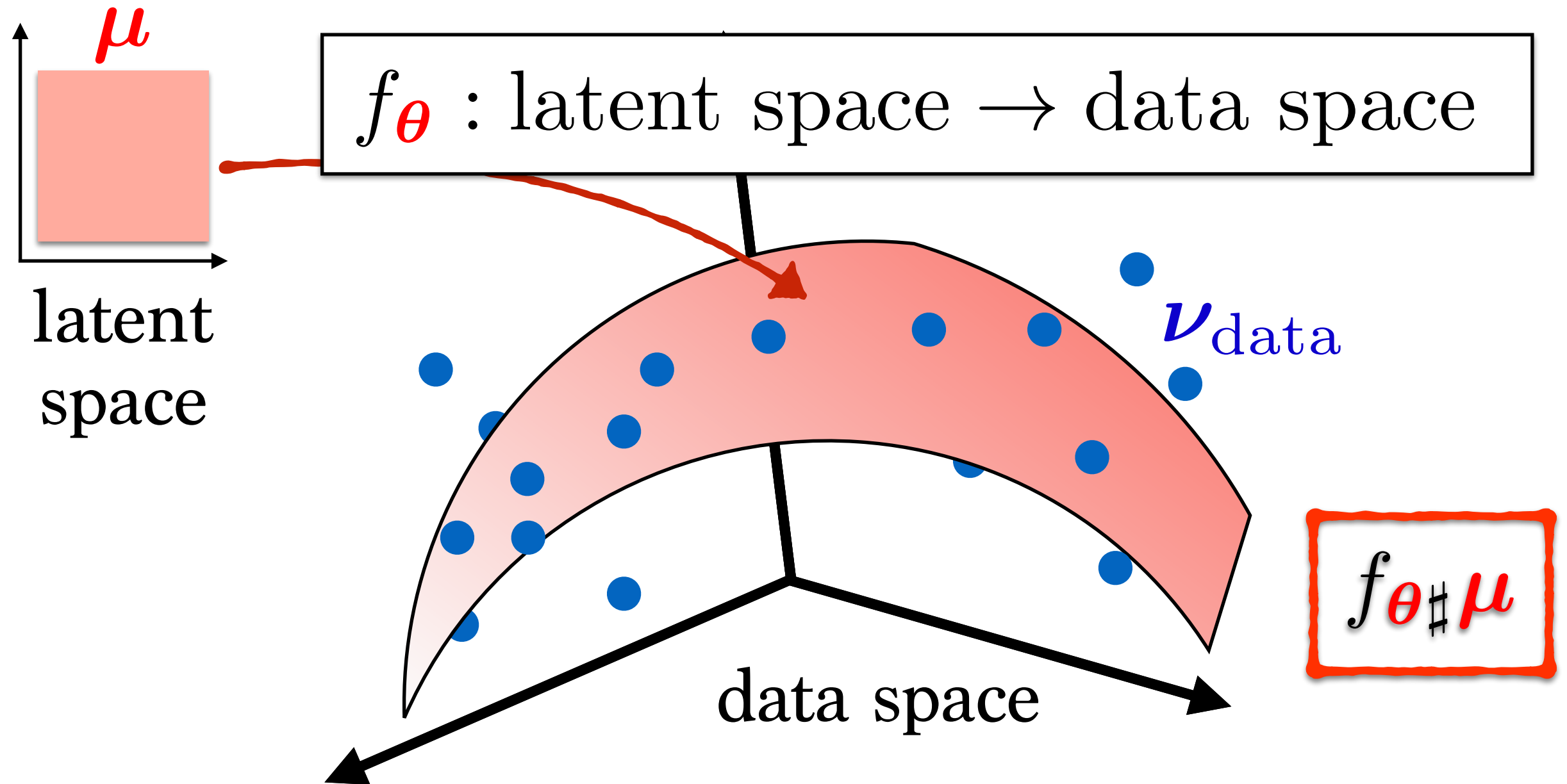




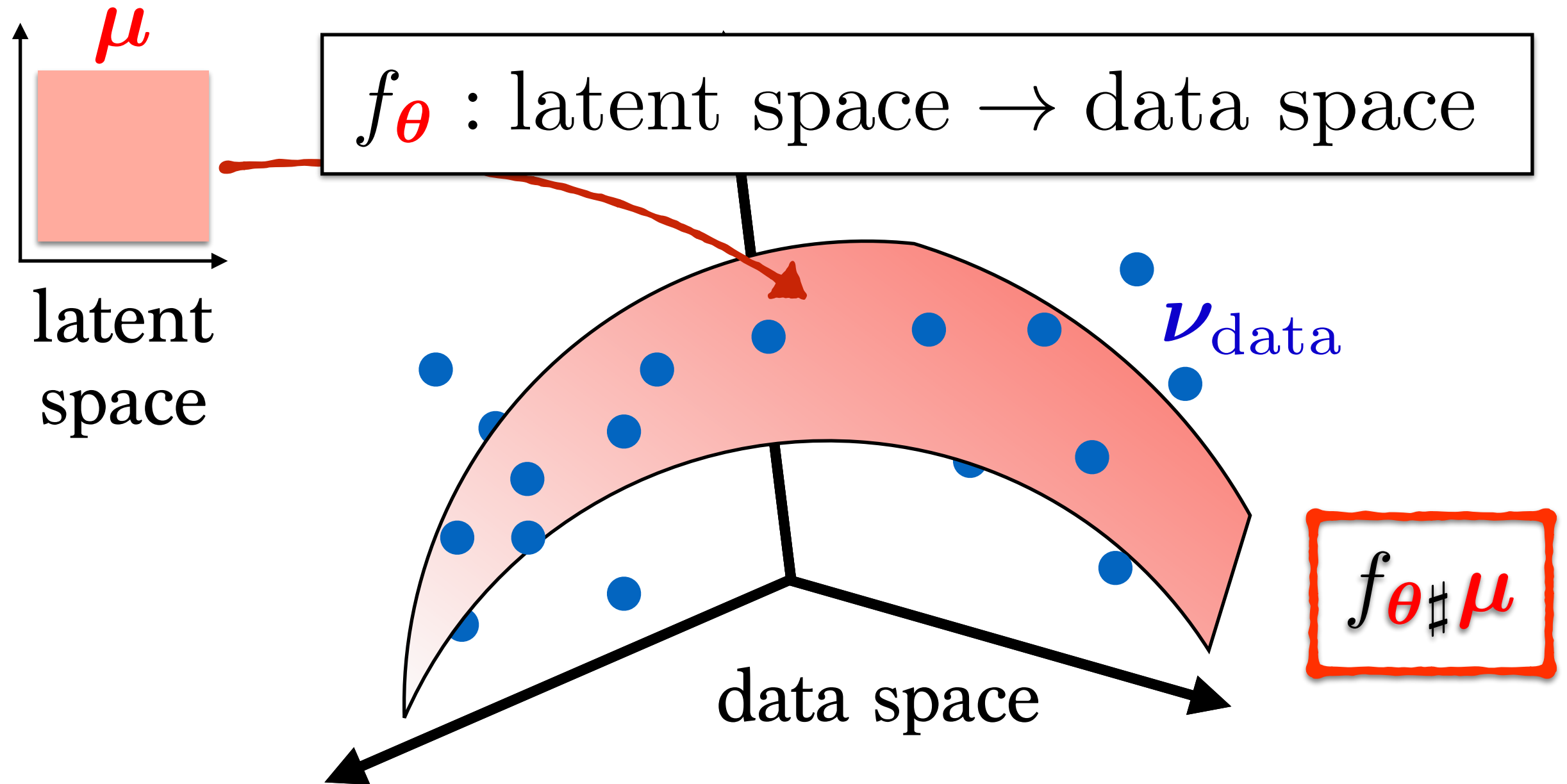
# Generative Models



# Generative Models

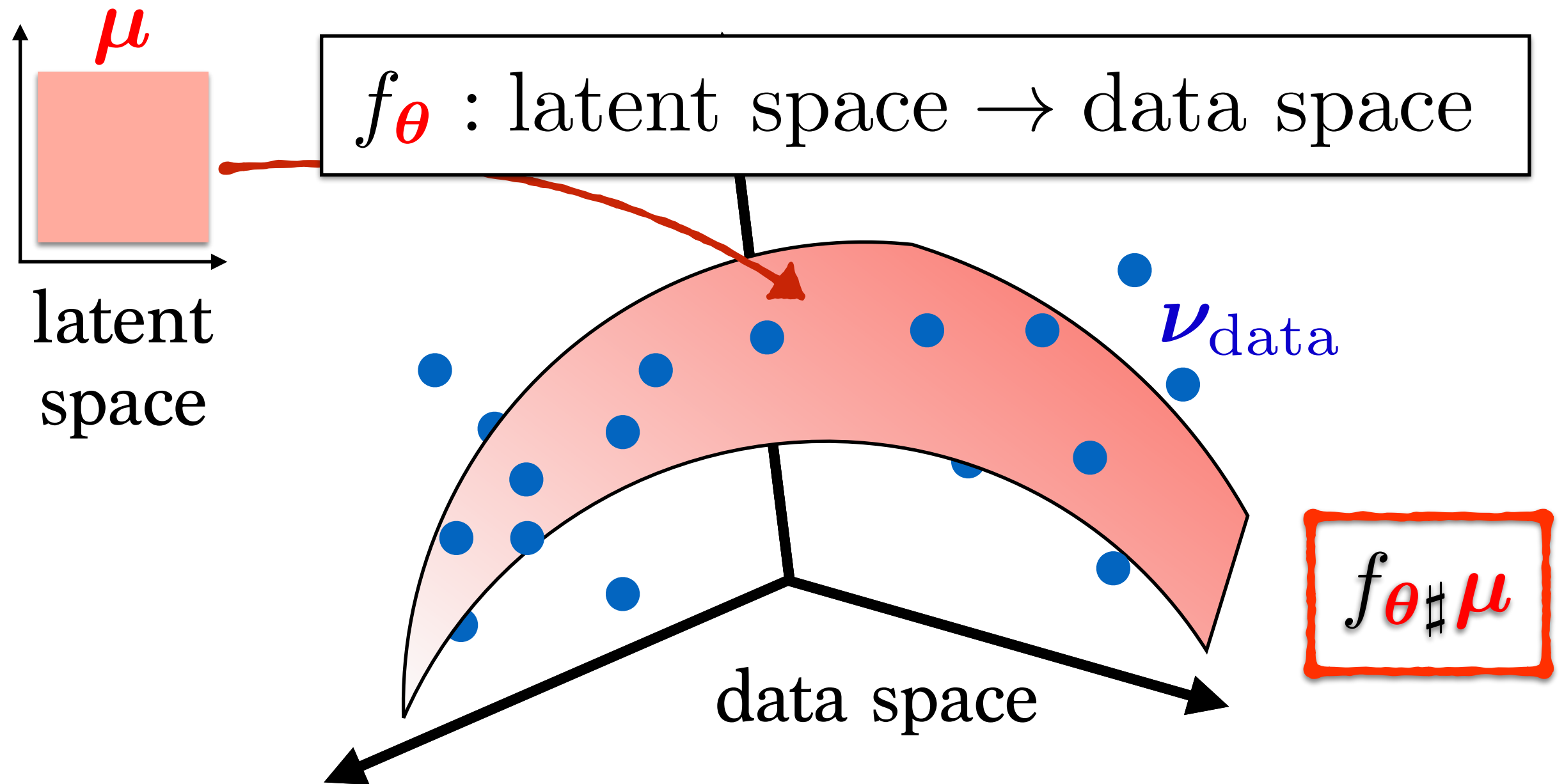


# Generative Models



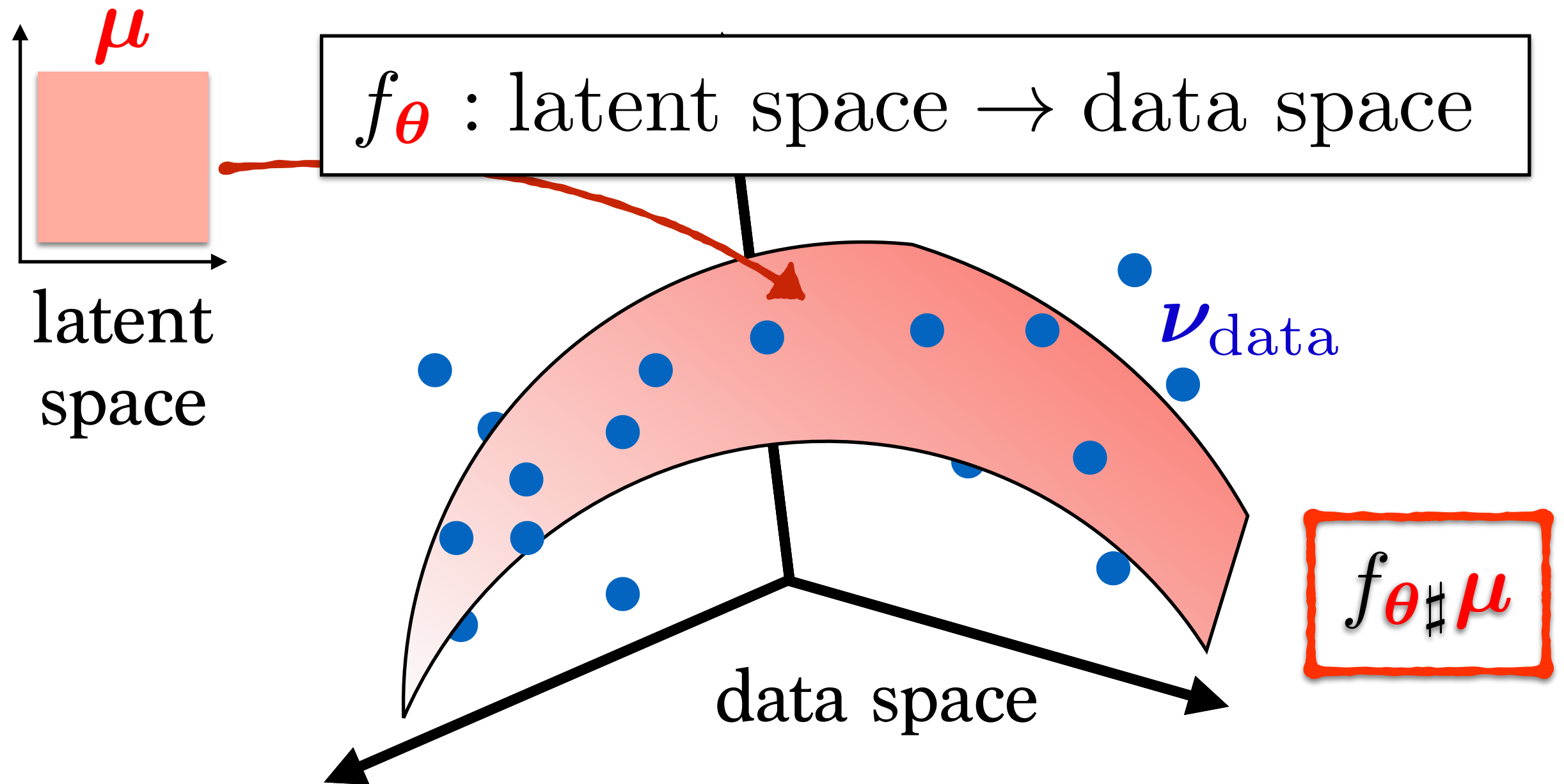
Push-forward:  $\forall B \subset \Omega, f_{\#} \mu(B) := \mu(f^{-1}(B))$

# Generative Models



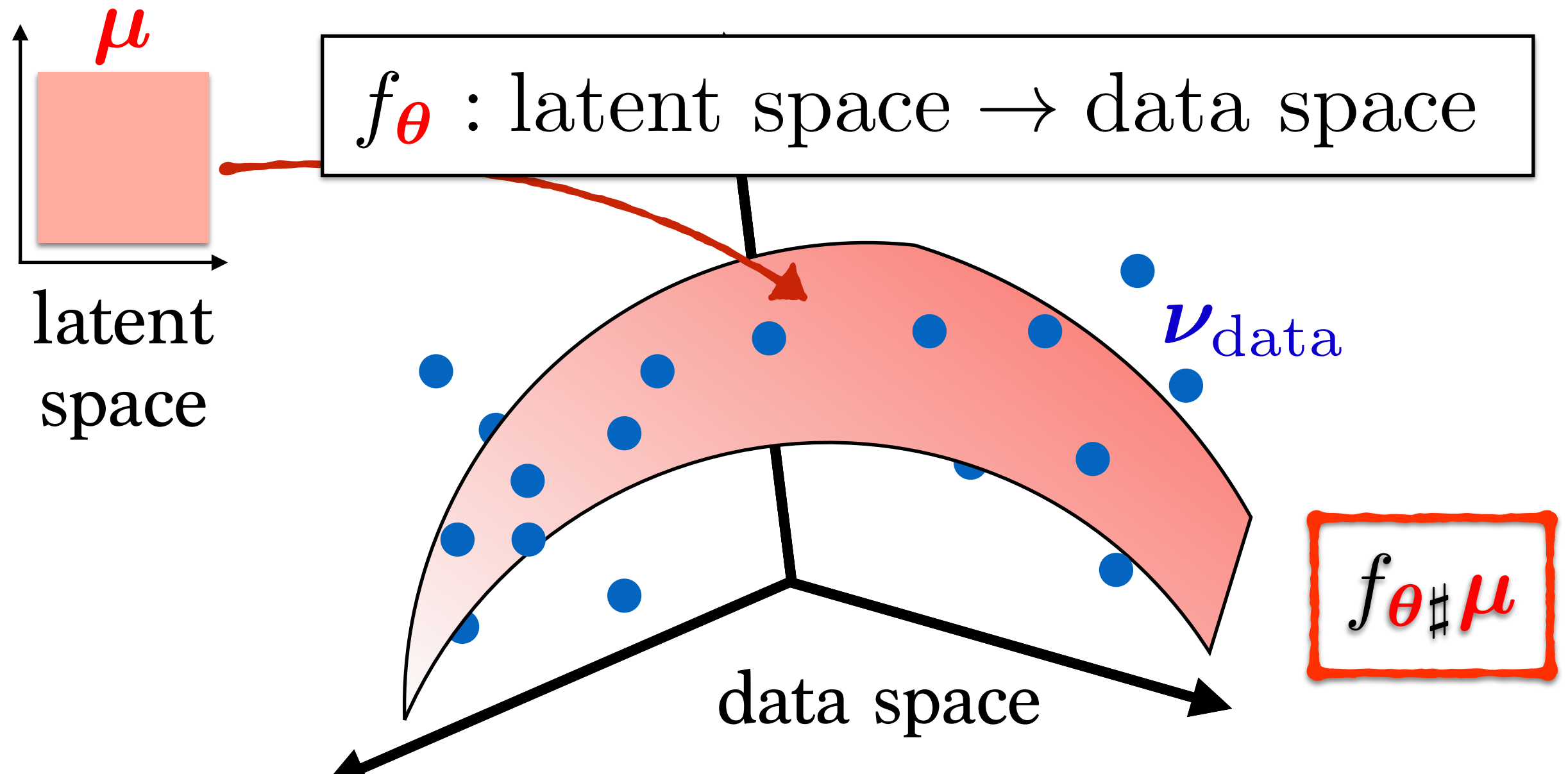
Goal: find  $\theta$  such that  $f_{\theta \# \mu}$  fits  $\nu_{\text{data}}$

# Generative Models



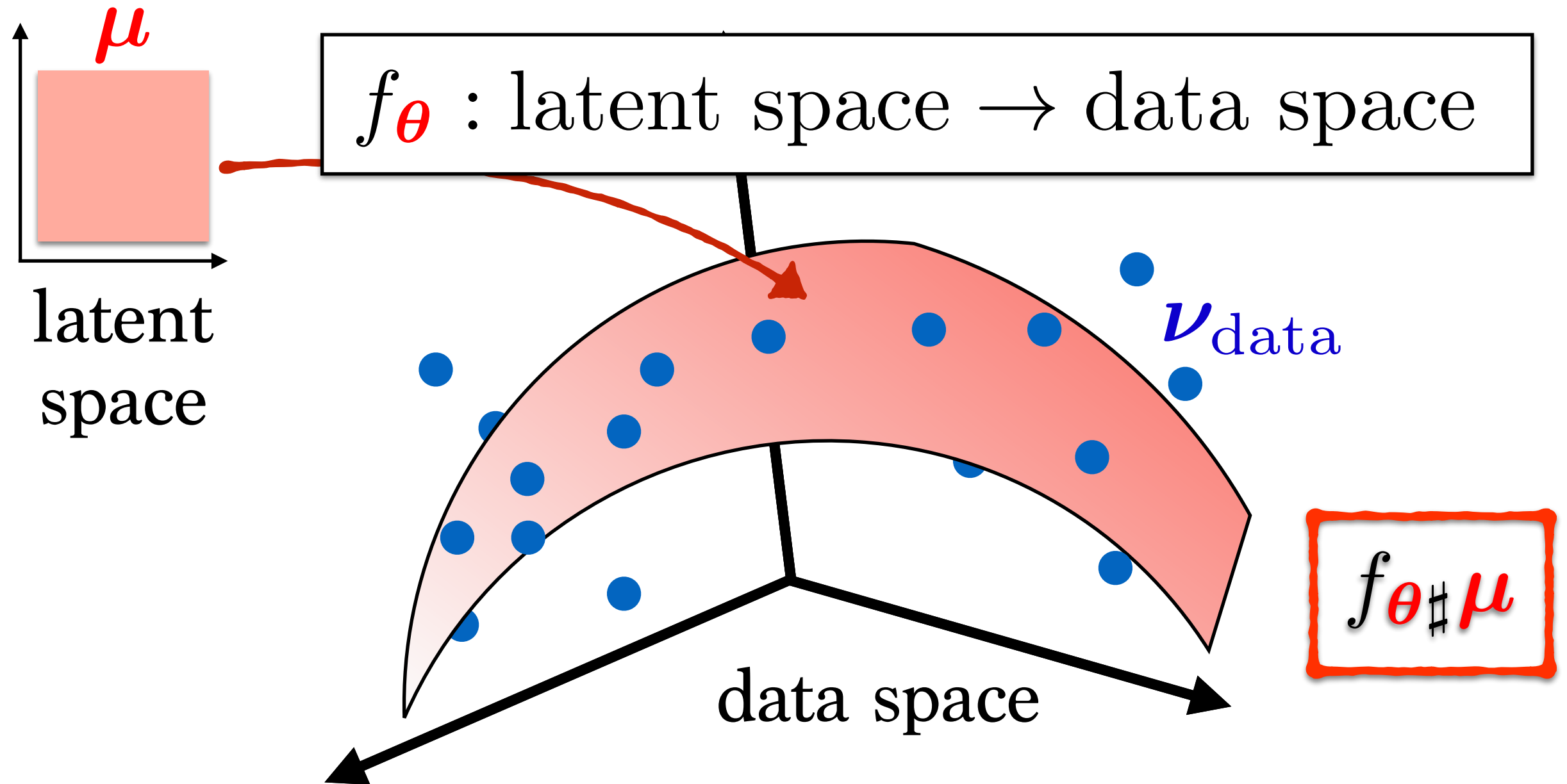
Goal: find  $\theta$  such that  $f_{\theta \# \mu}$  fits  $\nu_{\text{data}}$

# Generative Models



Difference between fitting a push forward measure  $f_{\theta\#}\mu$  vs. a density  $p_{\theta}$ ?

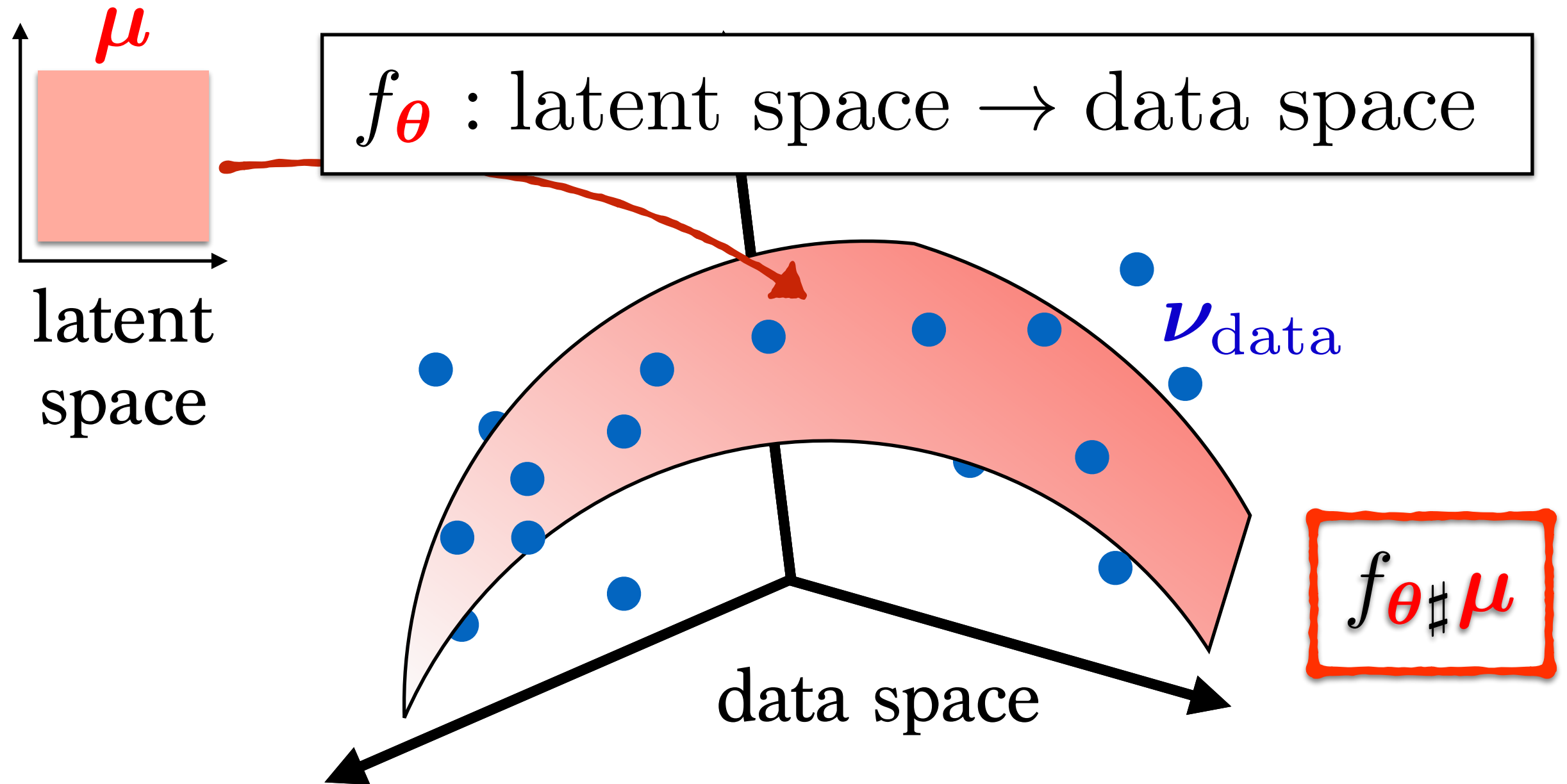
# Generative Models



MLE

$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i) = \min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \| p_{\theta})$$

# Generative Models



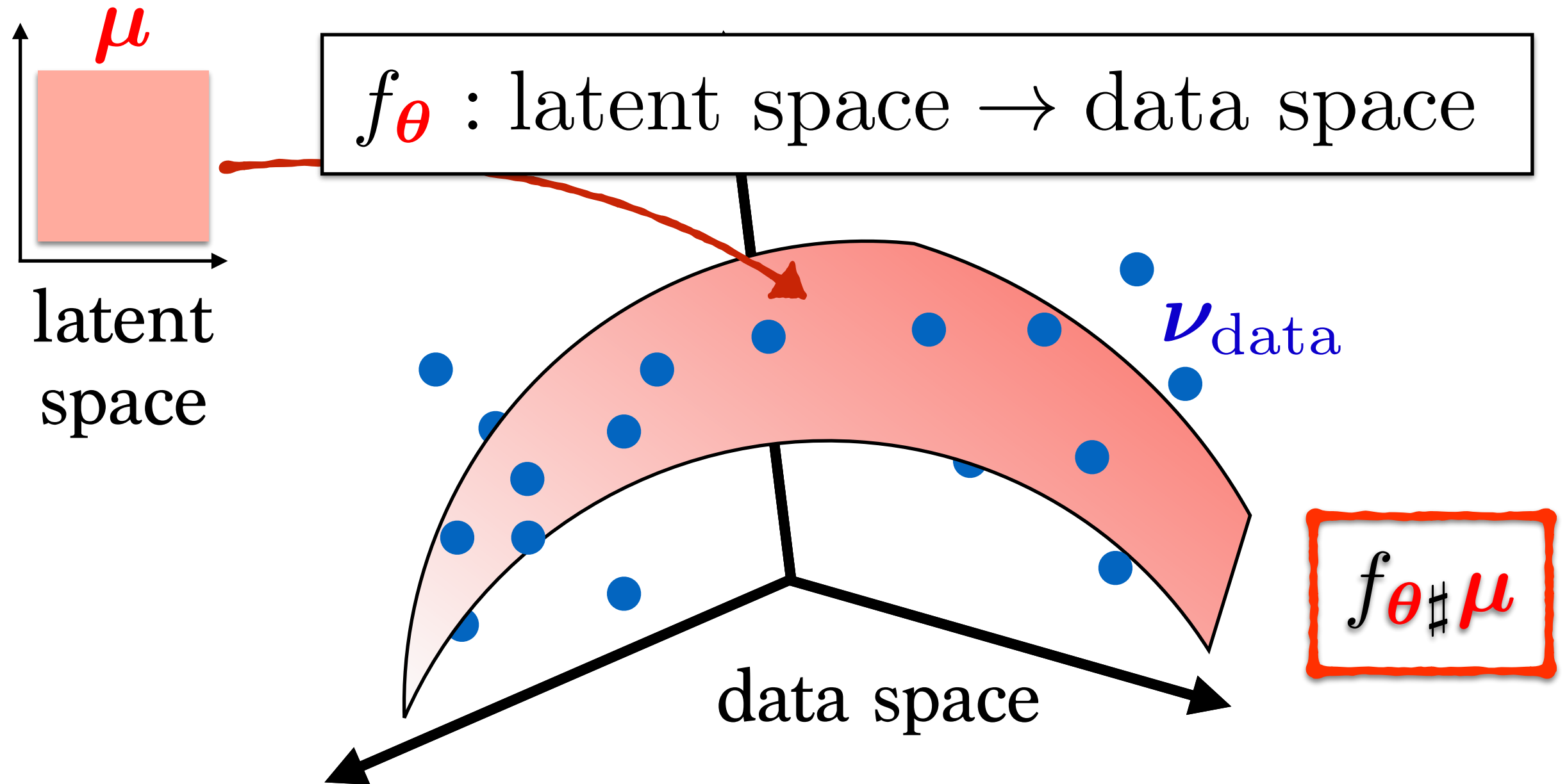
MLE

$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log f_{\theta} \# \mu(x_i)$$

$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \| f_{\theta} \# \mu)$$



# Generative Models

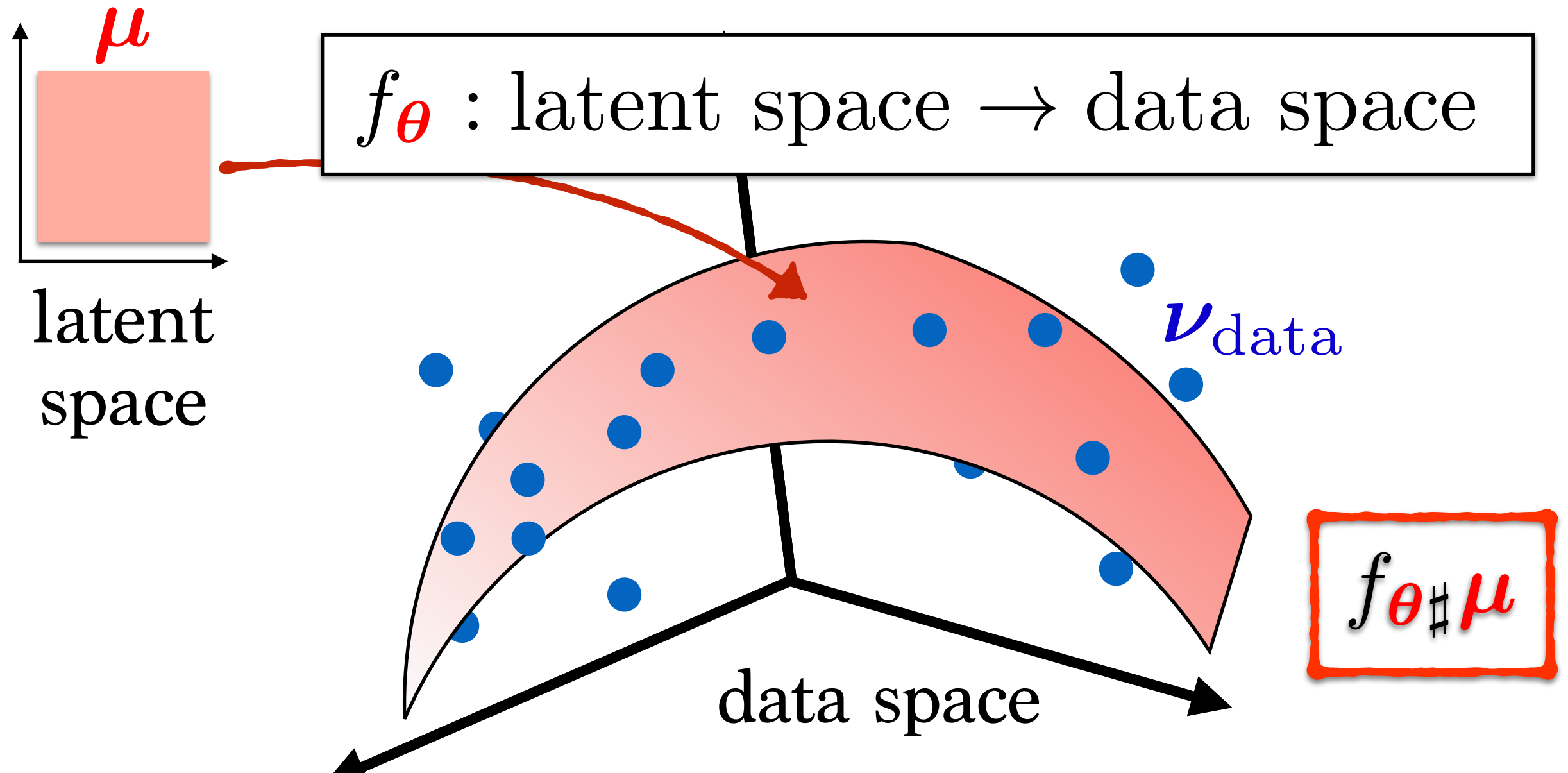


~~MLE~~

$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log f_{\theta \# \mu}(x_i) \quad \min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \| f_{\theta \# \mu})$$

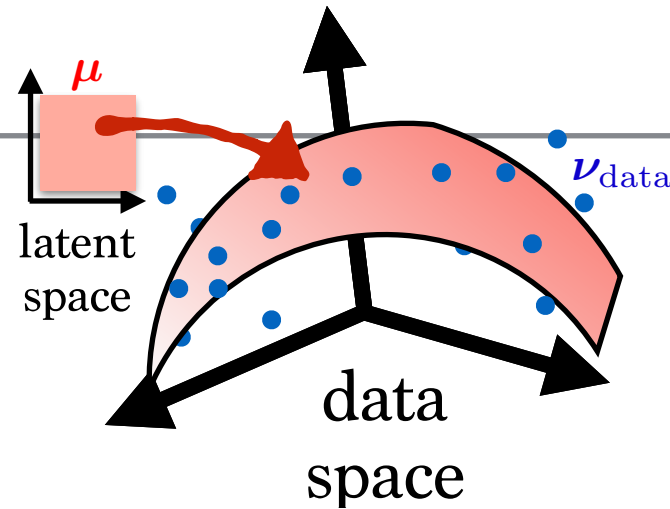


# Generative Models



Need a more flexible **discrepancy function** to compare  $\nu_{\text{data}}$  and  $f_{\theta \# \mu}$

# Workarounds?



- Formulation as adversarial problem [GPM...'14]

$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1))$$

- Use a **richer metric**  $\Delta$  for probability measures, able to handle measures with non-overlapping supports:

$$\min_{\theta \in \Theta} \Delta(\nu_{\text{data}}, p_{\theta}), \quad \text{not } \min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} || p_{\theta})$$

# Minimum $\Delta$ Estimation

*The Annals of Statistics*  
1980, Vol. 8, No. 3, 457–487

## MINIMUM $\chi^2$ CHI-SQUARE, NOT MAXIMUM LIKELIHOOD!

BY JOSEPH BERKSON

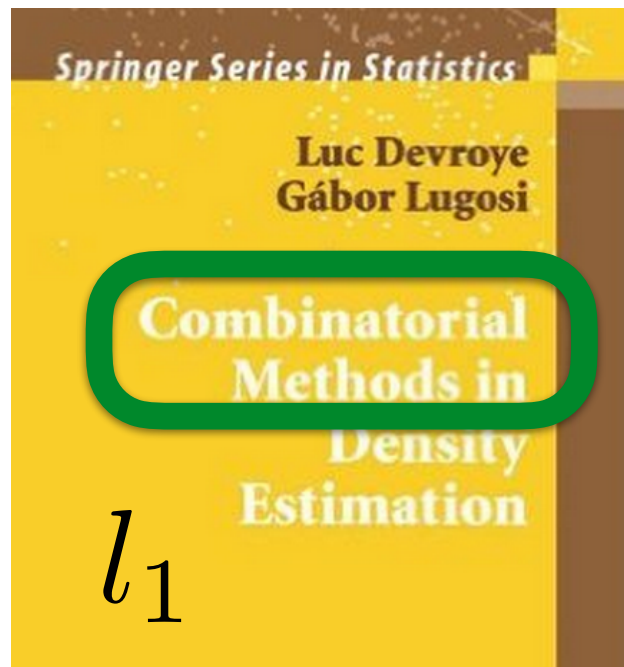
*Mayo Clinic, Rochester, Minnesota*



ELSEVIER

Computational Statistics & Data Analysis 29 (1998) 81–103

COMPUTATIONAL  
STATISTICS  
& DATA ANALYSIS



## Minimum $H$ Hellinger distance estimation for Poisson mixtures

Dimitris Karlis, Evdokia Xekalaki\*

*Department of Statistics, Athens University of Economics and Business, 76 Patission Str., 104 34 Athens, Greece*

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Statistics & Probability Letters 76 (2006) 1298–1302

STATISTICS &  
PROBABILITY  
LETTERS

[www.elsevier.com/locate/stapro](http://www.elsevier.com/locate/stapro)

## On minimum $K$ Kantorovich distance estimators

Federico Bassetti<sup>a</sup>, Antonella Bodini<sup>b</sup>, Eugenio Regazzini<sup>a,\*</sup>

# Minimum Kantorovich Estimation

- Use optimal transport theory, namely *Wasserstein distances* to define discrepancy  $\Delta$ .

$$\min_{\theta \in \Theta} W(\nu_{\text{data}}, f_{\theta\#} \mu)$$

- Optimal transport? fertile field in mathematics.



Monge



Kantorovich



Koopmans



Dantzig



Brenier



Otto



McCann



Villani

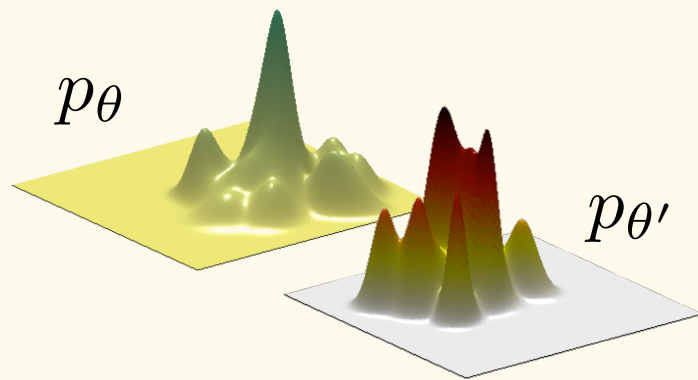
Nobel '75

Fields '10



# What is Optimal Transport?

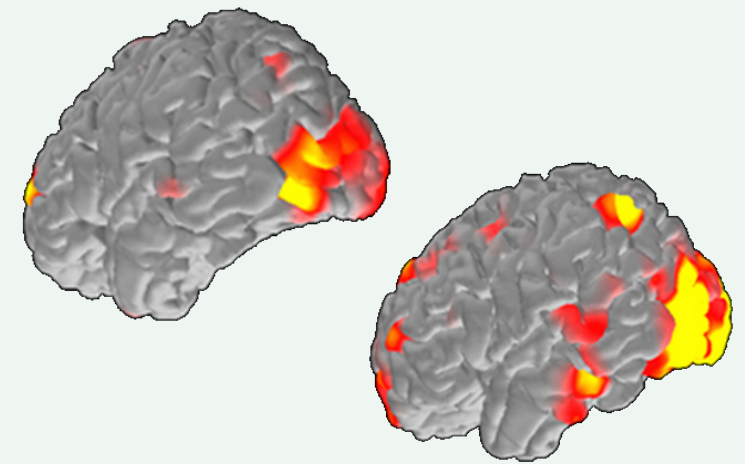
A geometric toolbox to  
compare **probability measures**  
supported on a metric space.



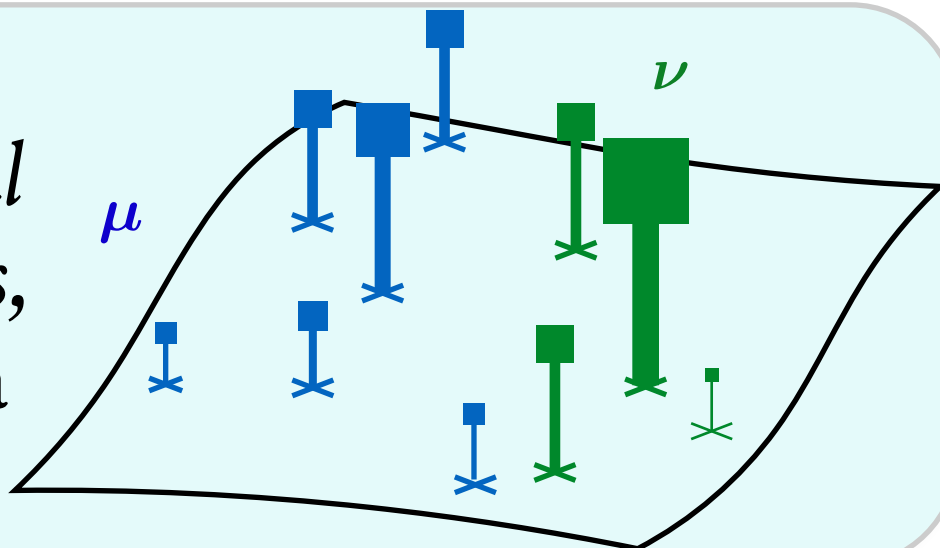
## Statistical Models



# Bags of features



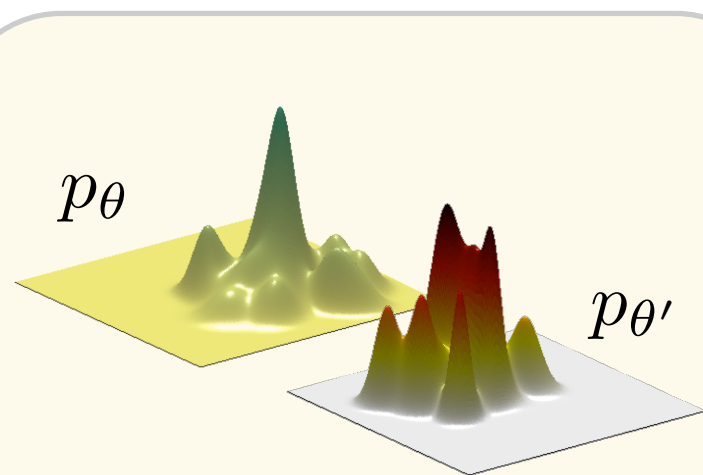
## Brain Activation Maps



## Color Histograms

# What is Optimal Transport?

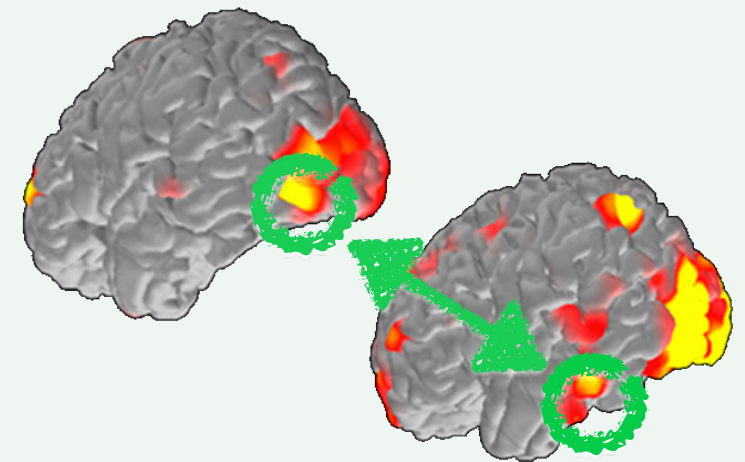
A geometric toolbox to  
compare probability measures  
supported on a metric space.



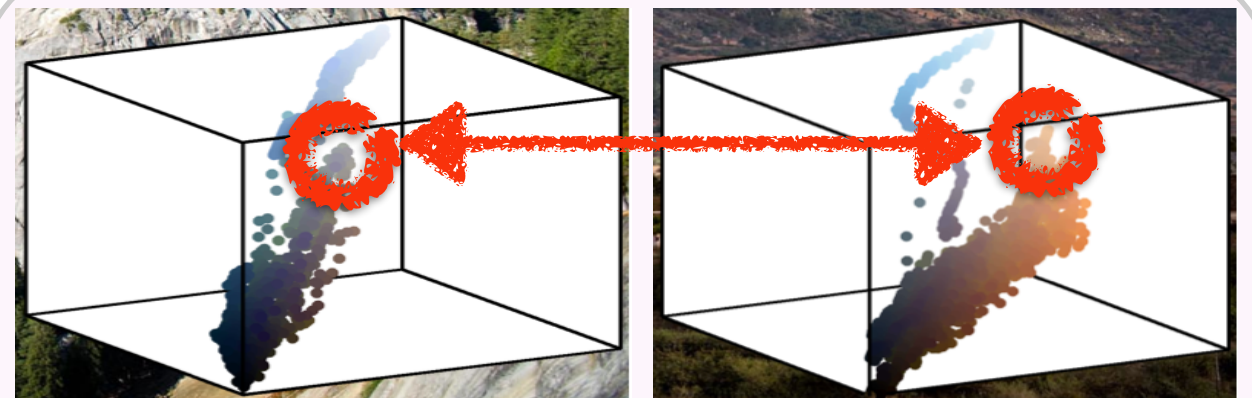
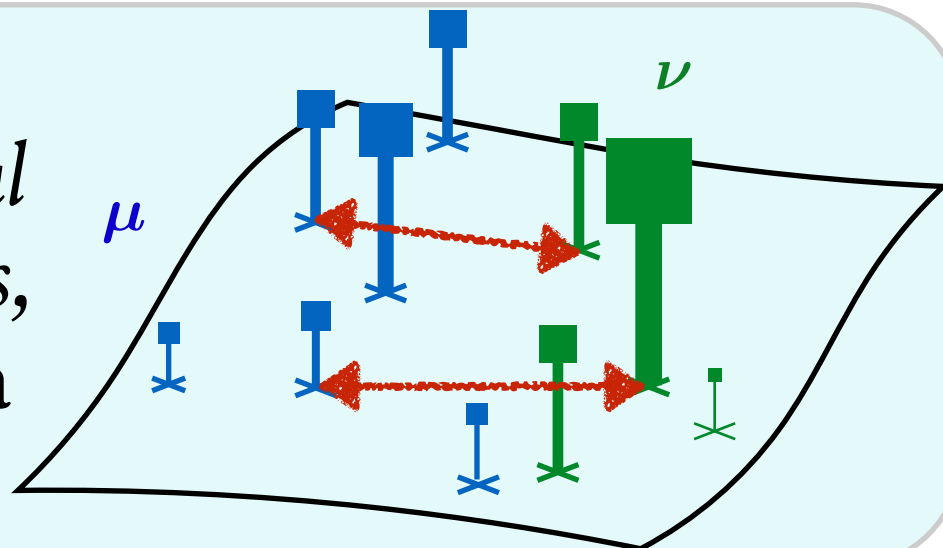
## Statistical Models



# Bags of features



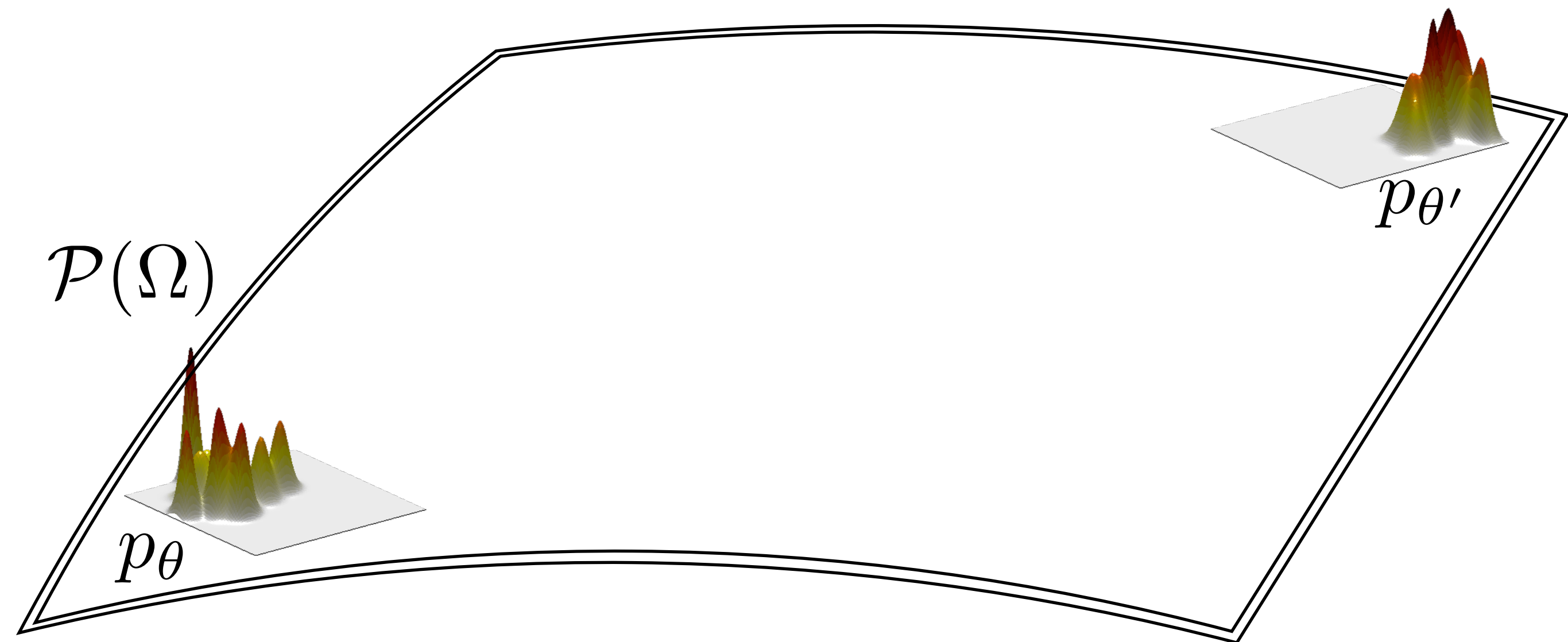
## Brain Activation Maps



## Color Histograms

# Optimal Transport Geometry

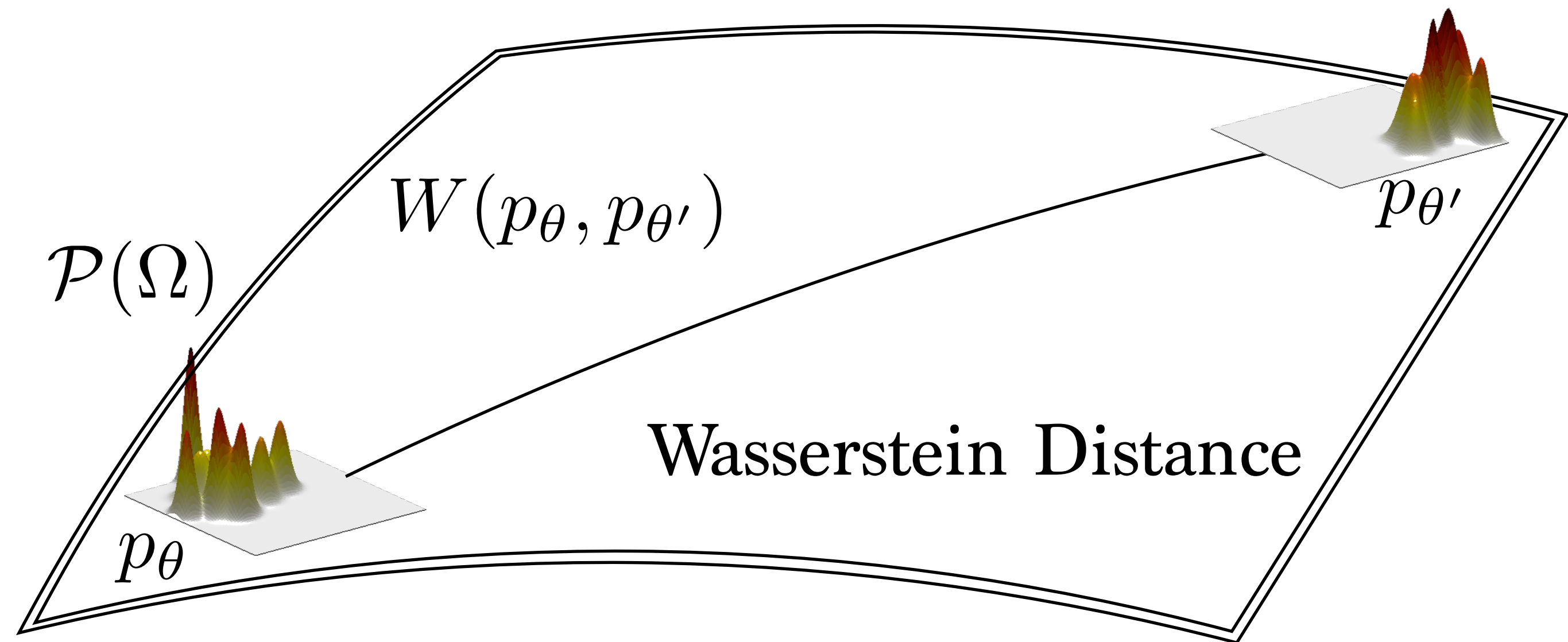
A **geometric toolbox** to  
compare probability measures  
supported on a metric space.





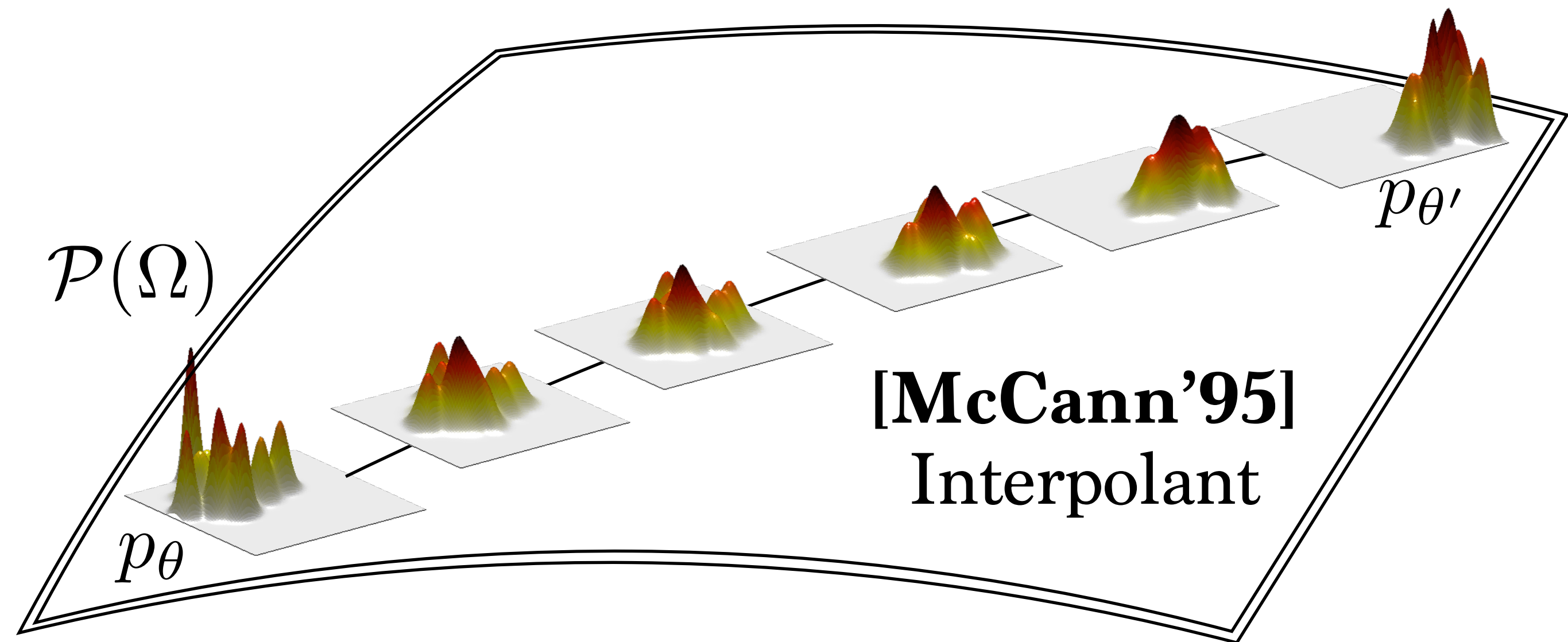
# Optimal Transport Geometry

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



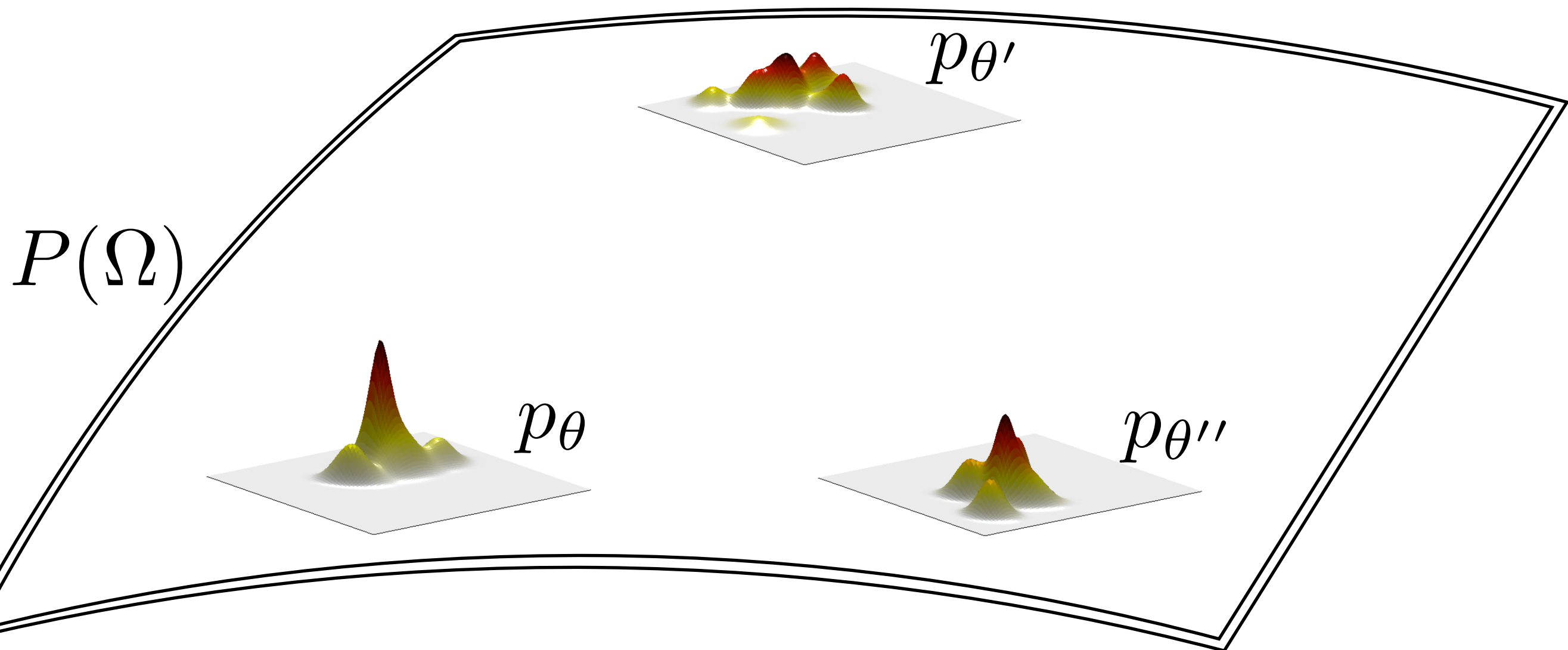
# Optimal Transport Geometry

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



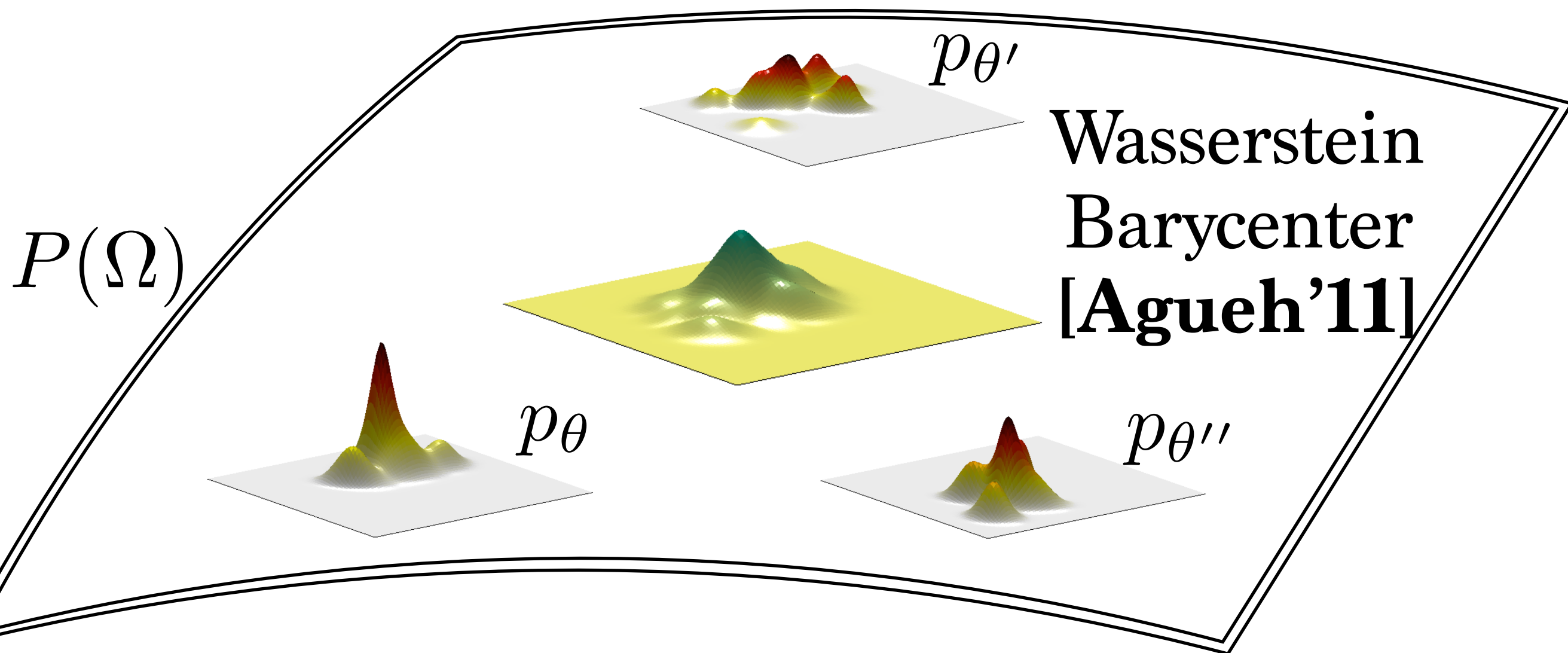
# Optimal Transport Geometry

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



# Optimal Transport Geometry

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



# Optimal Transport Geometry

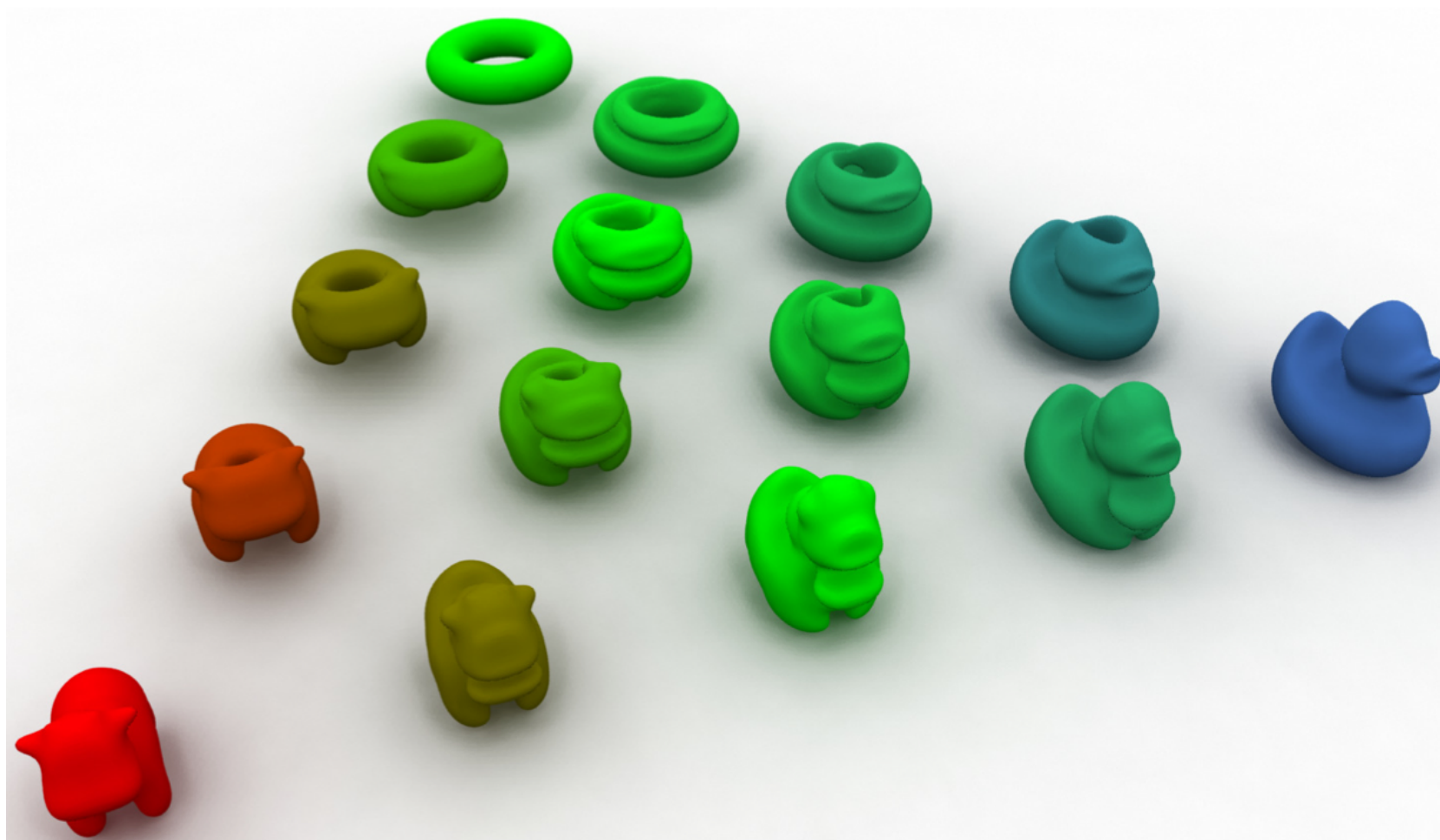
A **geometric toolbox** to compare probability measures supported on a metric space.



[SDPC..'15]

# Optimal Transport Geometry

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

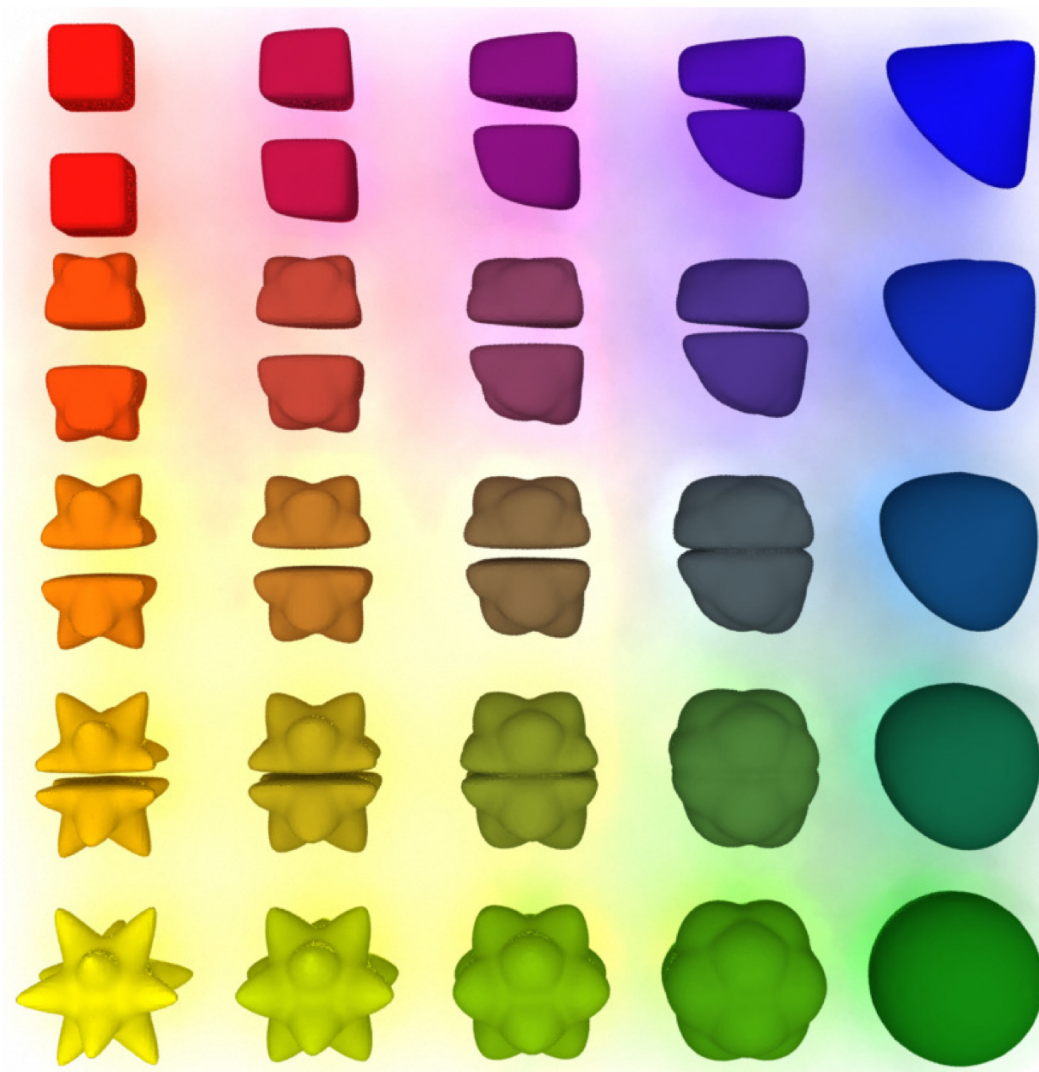


[SDPC..'15]



# Optimal Transport Geometry

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



[SDPC..'15]

# Origins: Monge's Problem

666. MÉMOIRES DE L'ACADÉMIE ROYALE

---

*M É M O I R E*

*S U R L A*

*T H É O R I E D E S D É B L A I S*  
*E T D E S R E M B L A I S.*

Par M. M O N G E.

**L**ORSQU'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.



# Origins: Monge's Problem

---



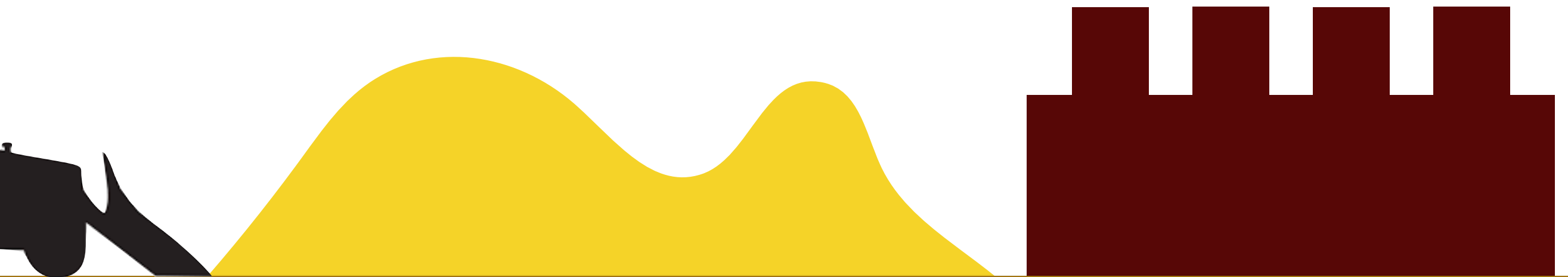
# Origins: Monge's Problem

---



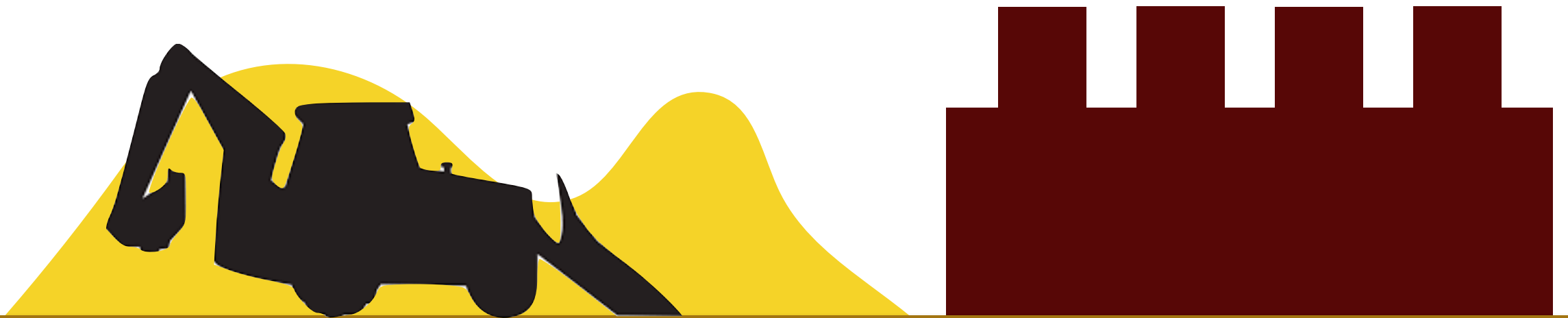
# Origins: Monge's Problem

---



# Origins: Monge's Problem

---



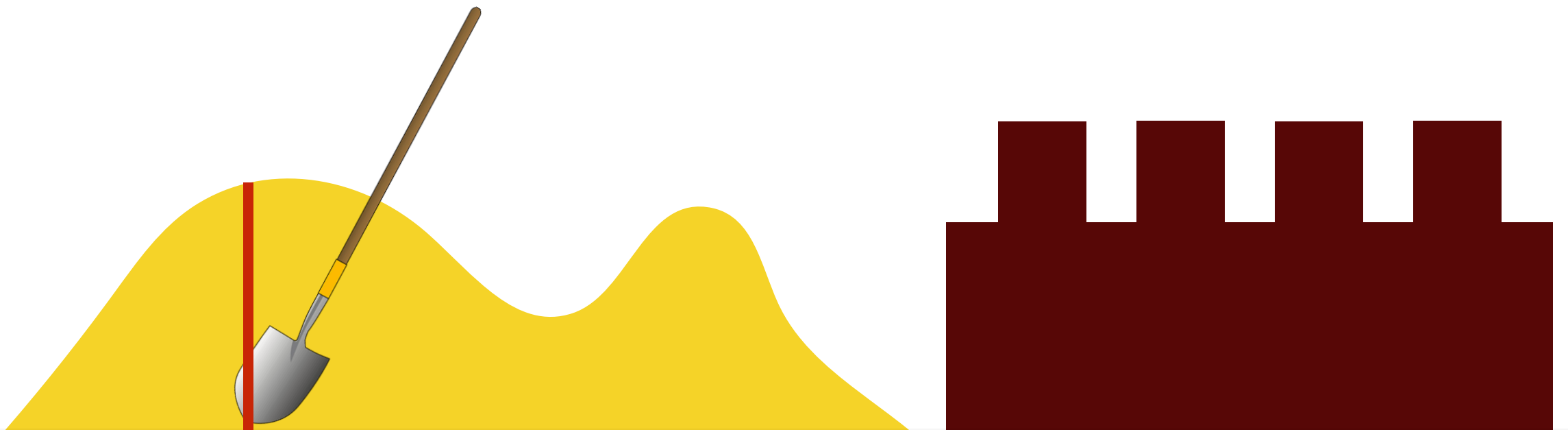
# Origins: Monge's Problem

---



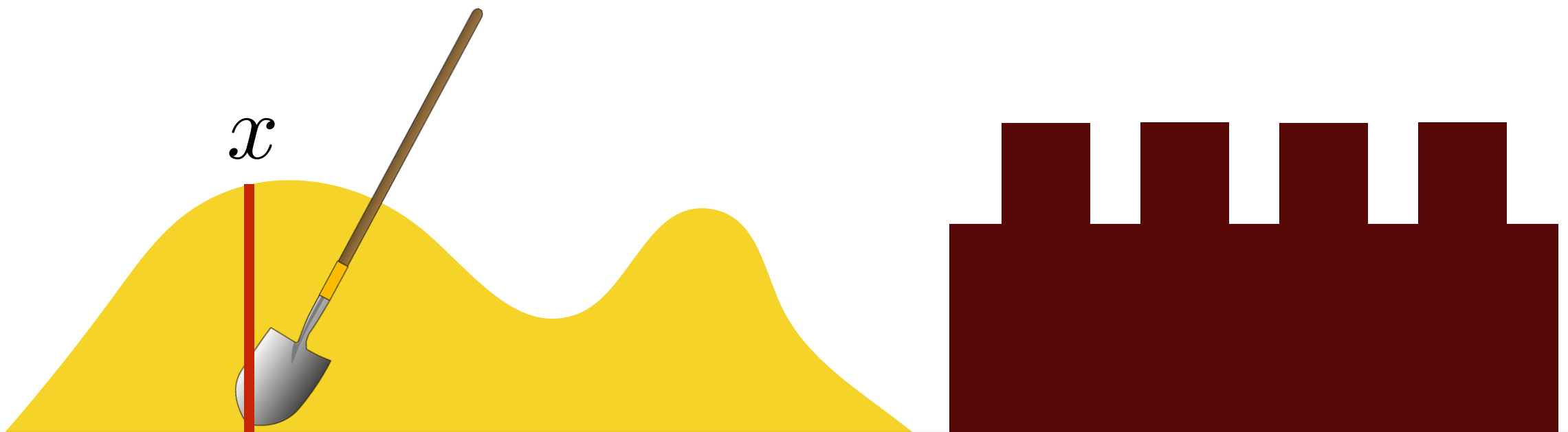
# Origins: Monge's Problem

---



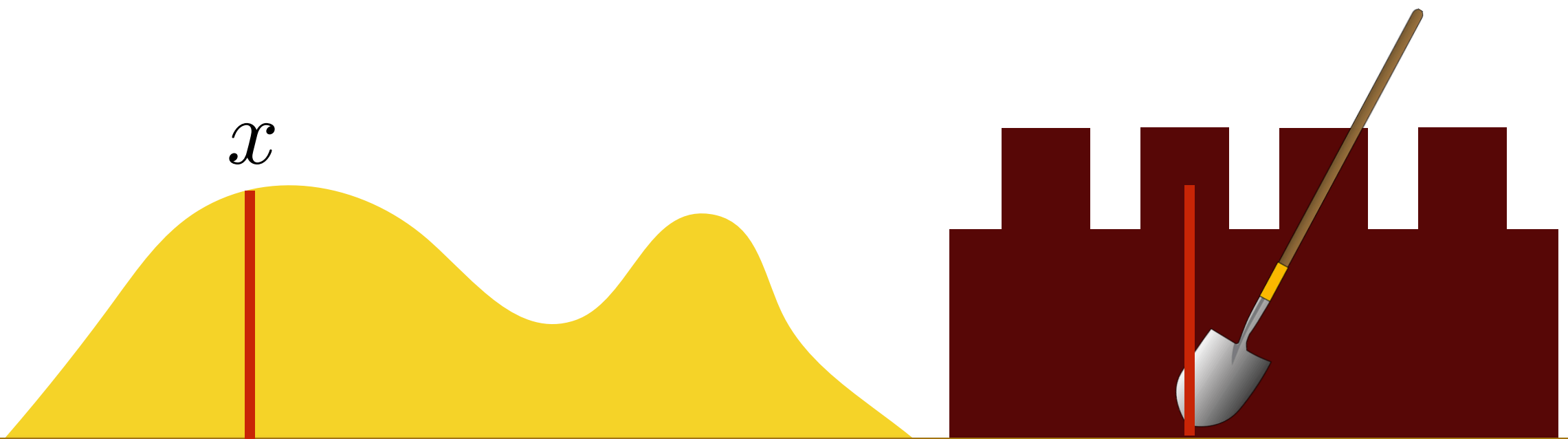
# Origins: Monge's Problem

---



# Origins: Monge's Problem

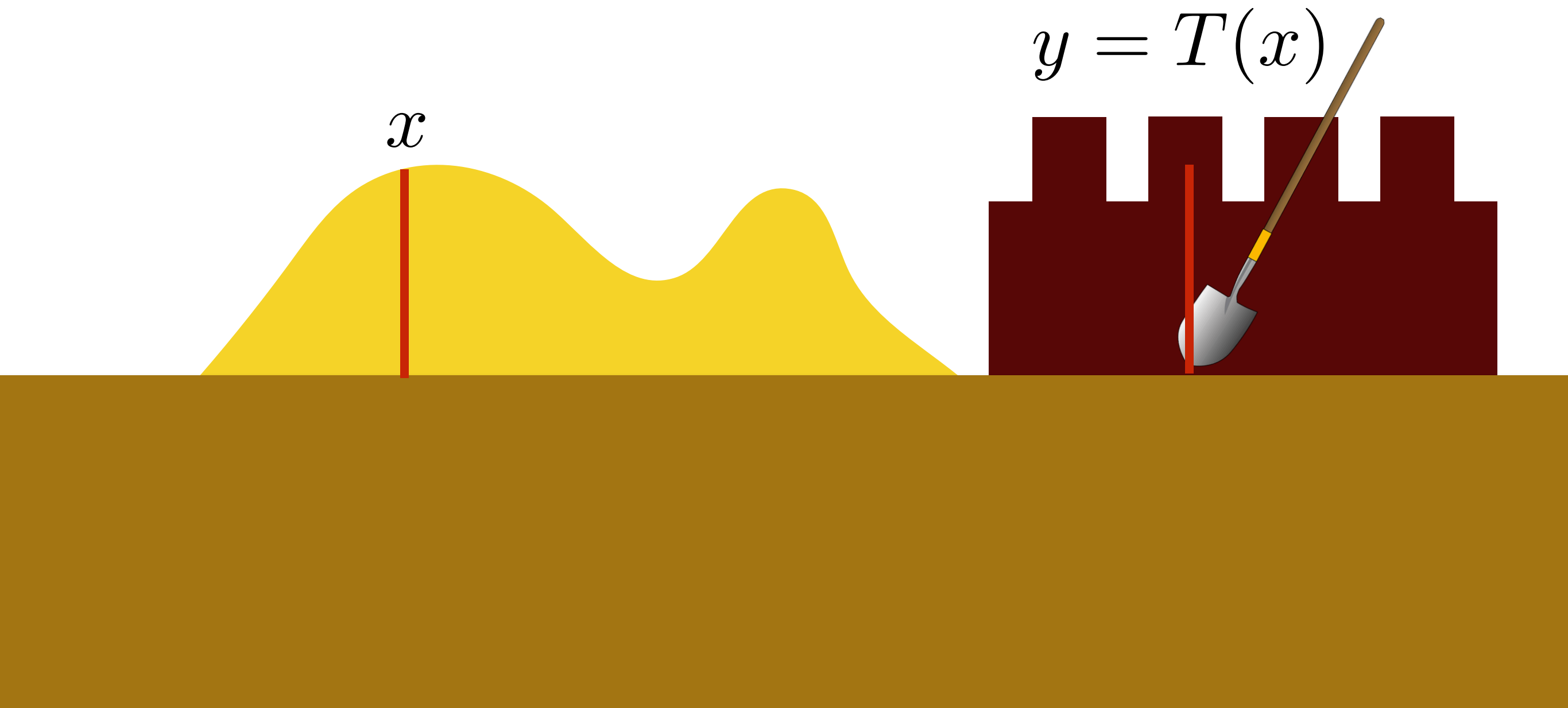
---



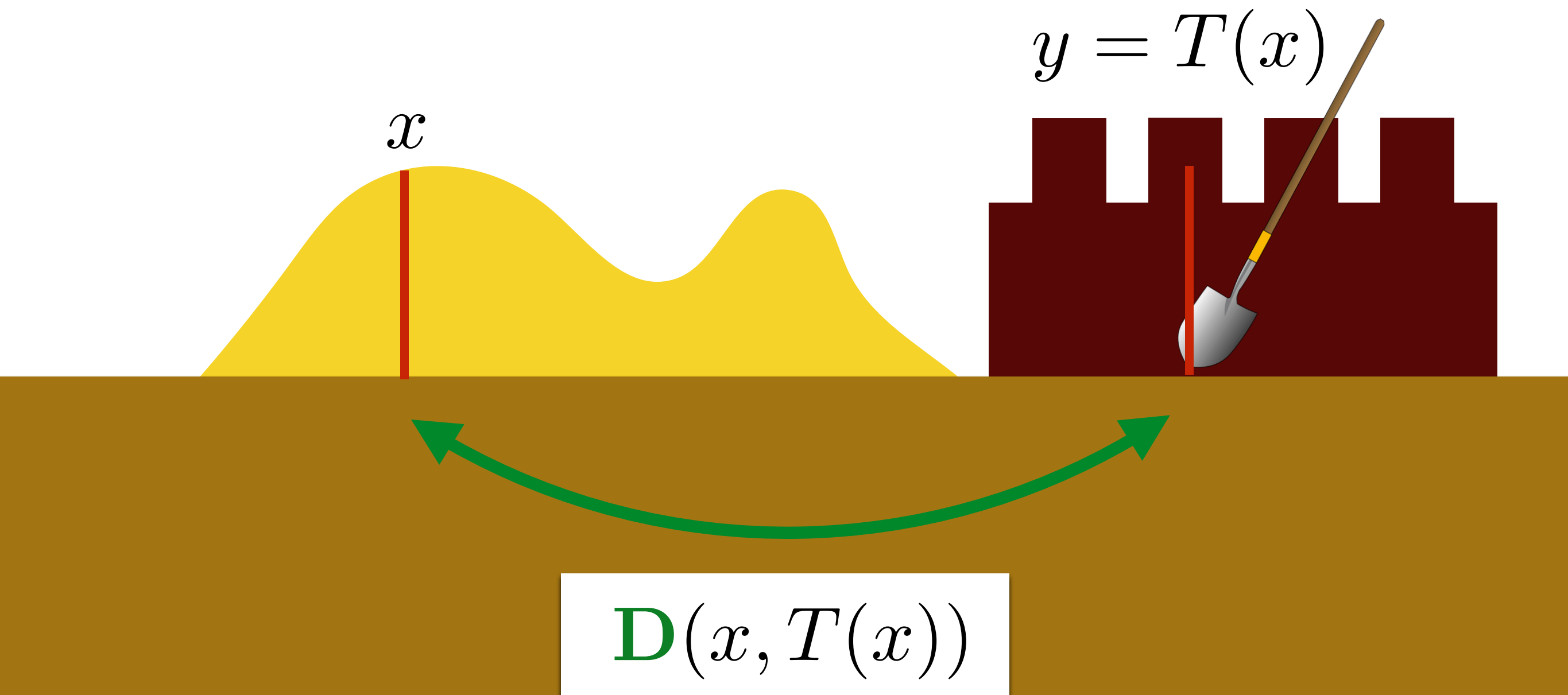


# Origins: Monge's Problem

---



# Origins: Monge's Problem

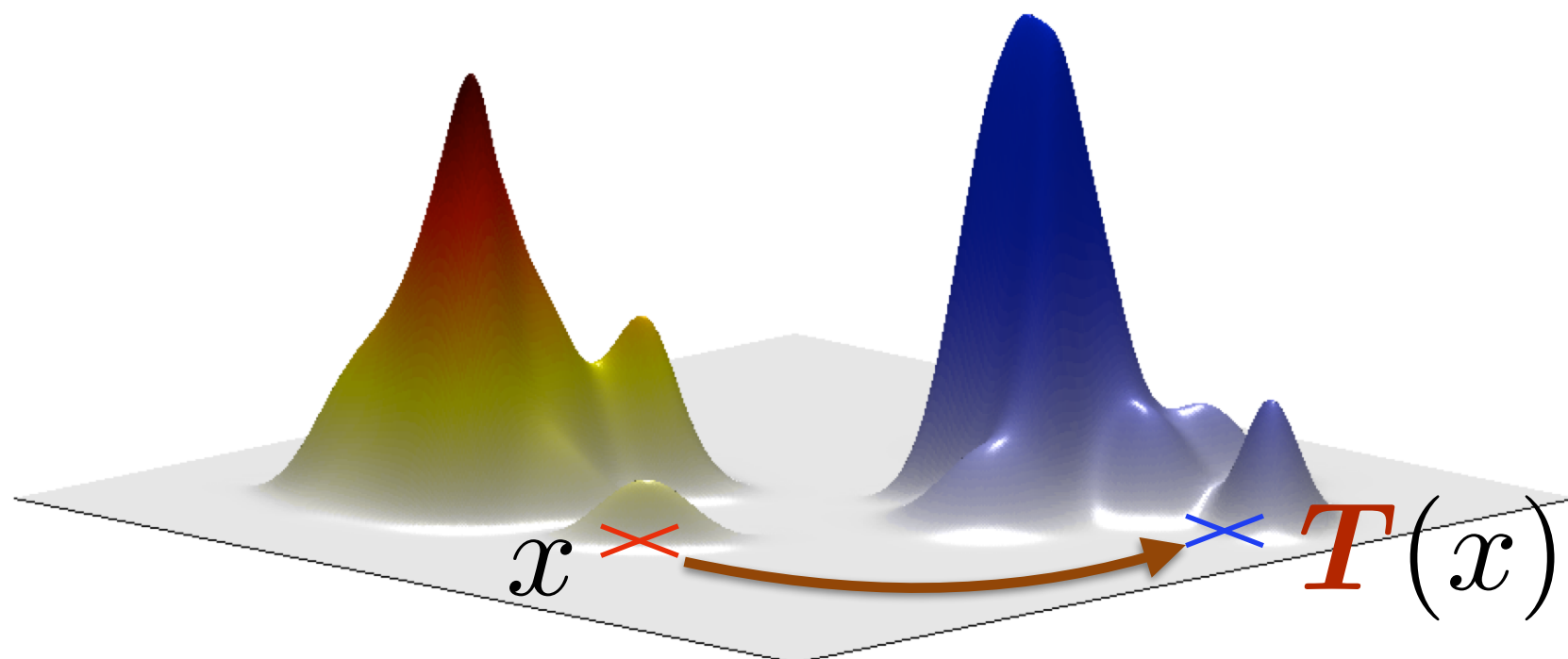


# Origins: Monge's Problem

$\Omega$  a probability space,  $\mathbf{c} : \Omega \times \Omega \rightarrow \mathbb{R}$ .  
 $\mu, \nu$  two probability measures in  $\mathcal{P}(\Omega)$ .

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

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

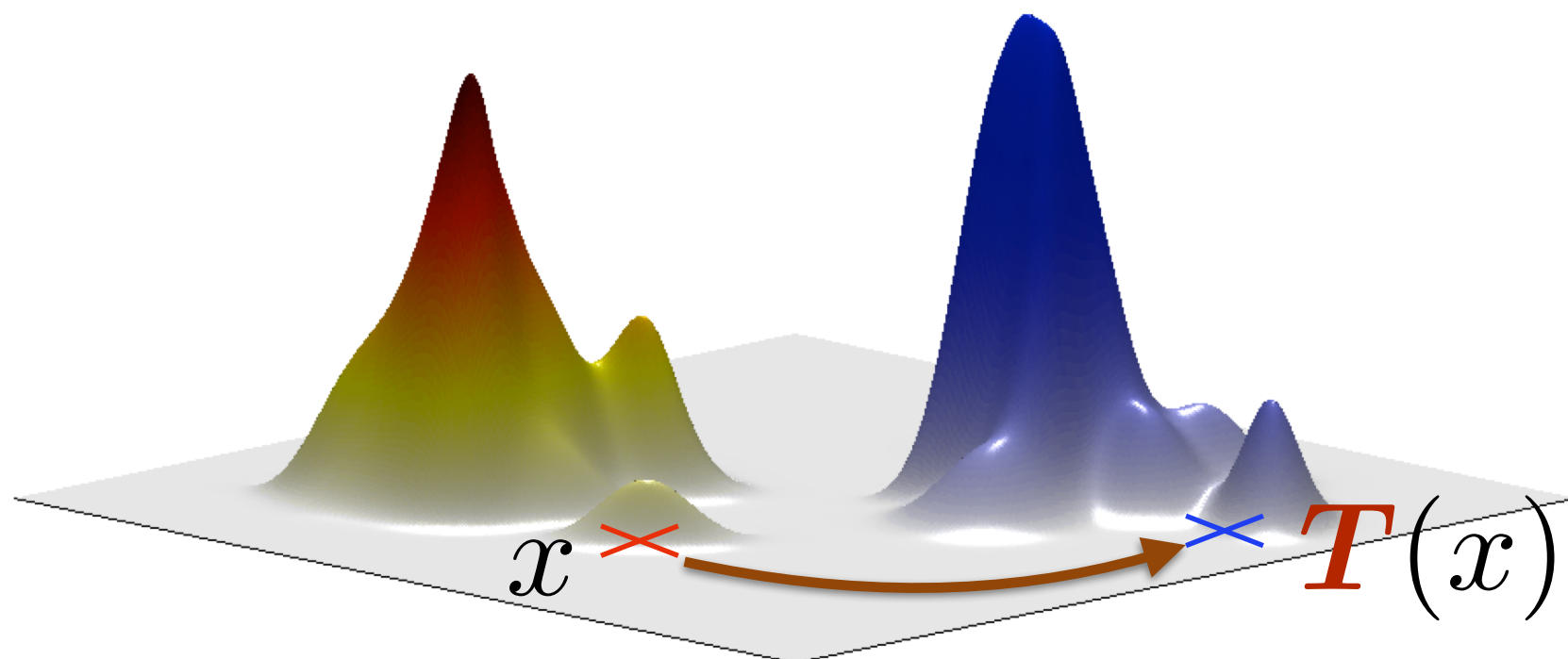


# Origins: Monge's Problem

$\Omega$  a probability space,  $\mathbf{c} : \Omega \times \Omega \rightarrow \mathbb{R}$ .  
 $\mu, \nu$  two probability measures in  $\mathcal{P}(\Omega)$ .

[Monge'81] problem: find a map  $\mathbf{T} : \Omega \rightarrow \Omega$

[Brenier'87] If  $\Omega = \mathbb{R}^d$ ,  $\mathbf{c} = \|\cdot - \cdot\|^2$ ,  
 $\mu, \nu$  a.c., then  $\mathbf{T} = \nabla u$ ,  $u$  convex.

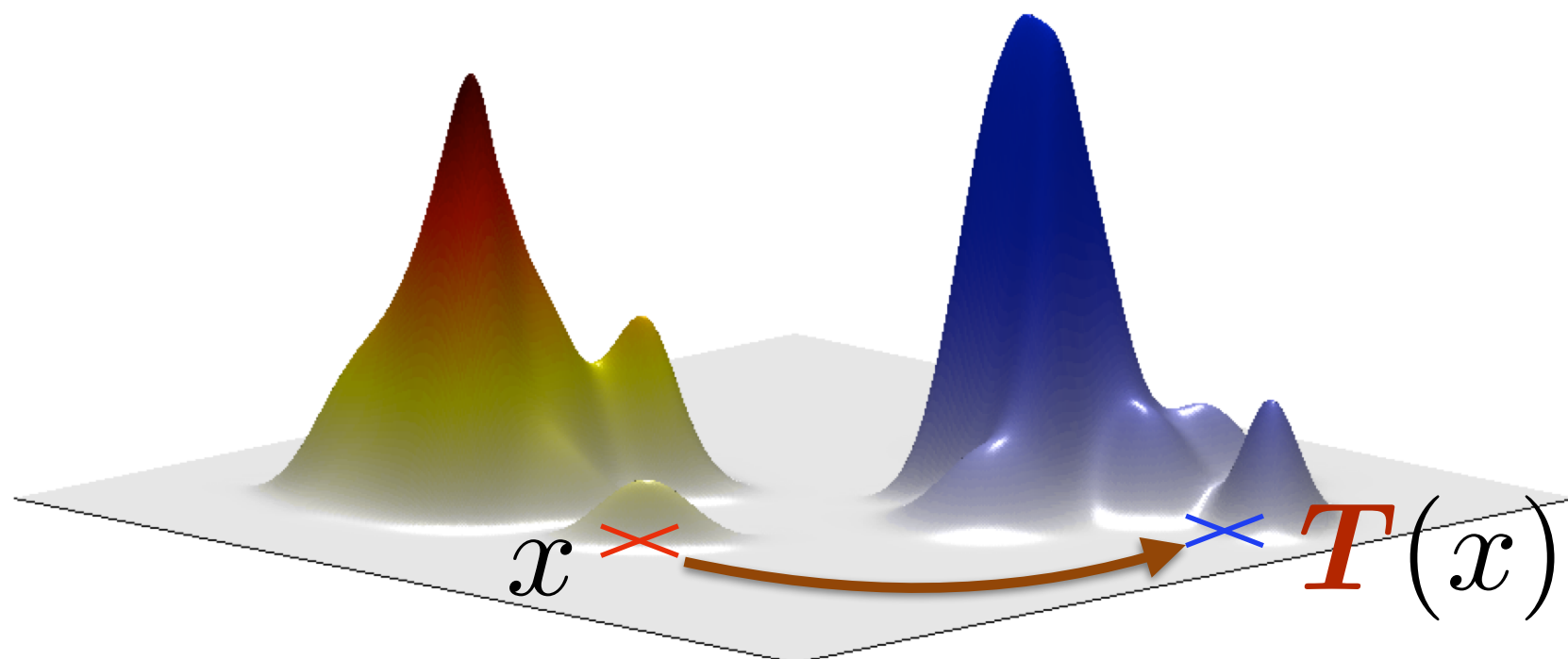


# Monge's Problem

$\Omega$  a probability space,  $\mathbf{c} : \Omega \times \Omega \rightarrow \mathbb{R}$ .  
 $\mu, \nu$  two probability measures in  $\mathcal{P}(\Omega)$ .

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

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

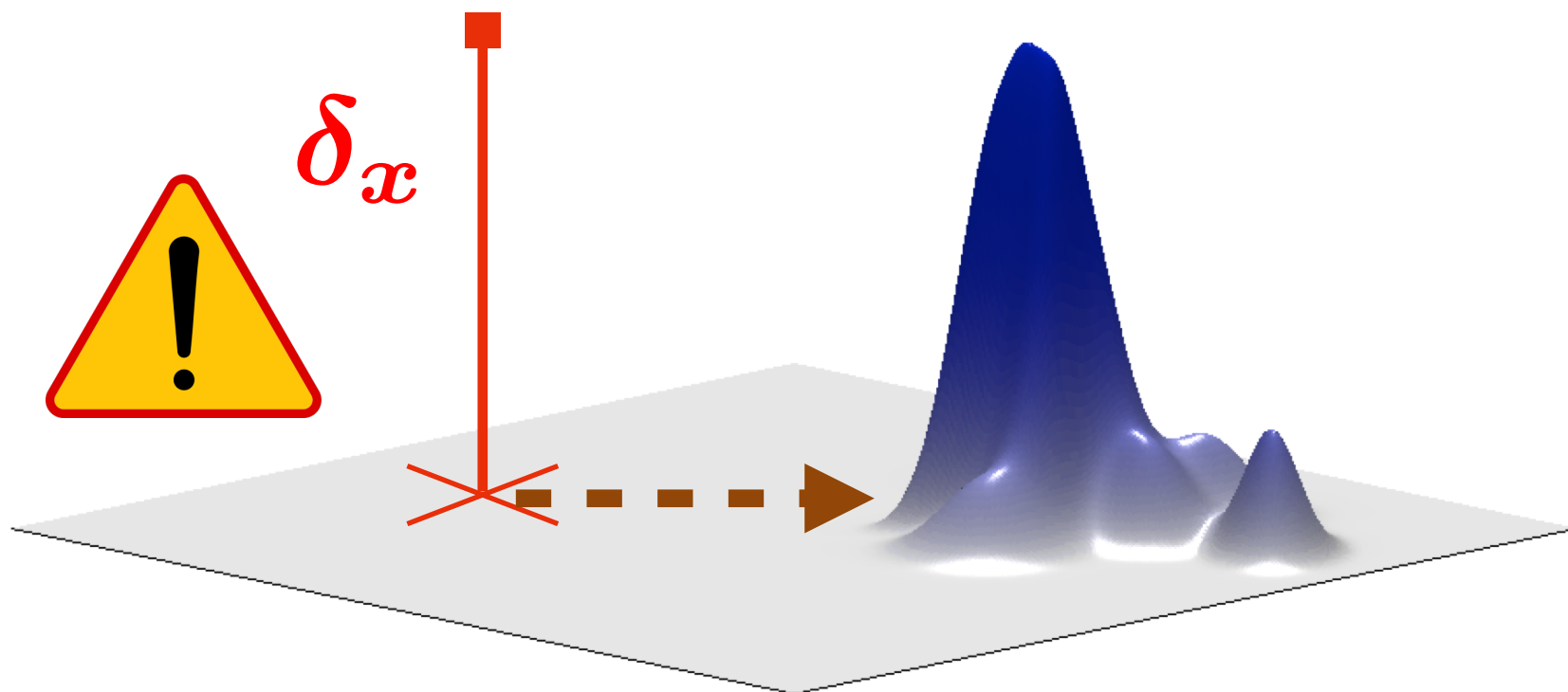


# Monge's Problem

$\Omega$  a probability space,  $\mathbf{c} : \Omega \times \Omega \rightarrow \mathbb{R}$ .  
 $\mu, \nu$  two probability measures in  $\mathcal{P}(\Omega)$ .

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

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



# [Kantorovich'42] Relaxation

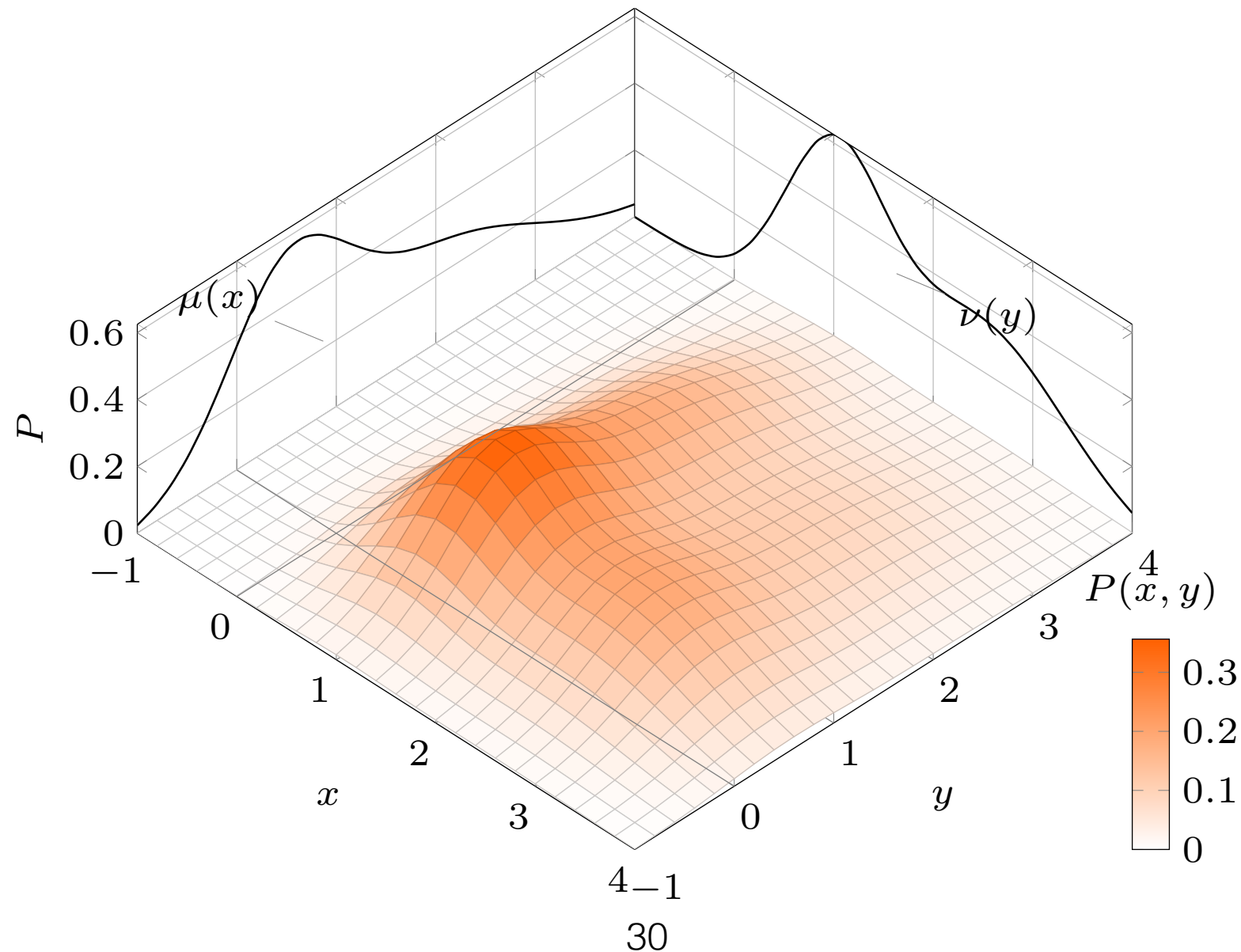
- Instead of maps  $T : \Omega \rightarrow \Omega$ , consider probabilistic maps, i.e. **couplings**  $P \in \mathcal{P}(\Omega \times \Omega)$ :

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



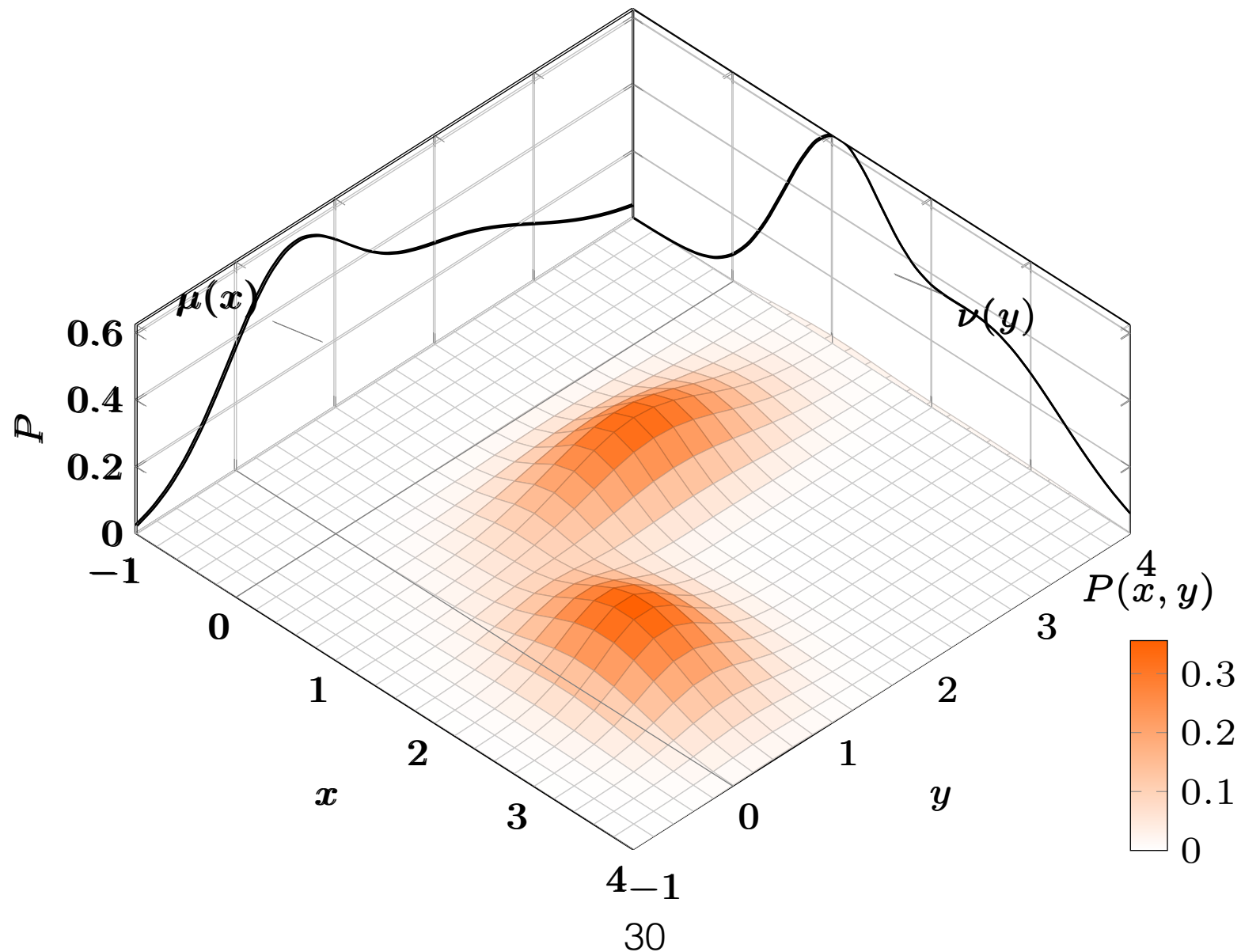
# [Kantorovich'42] Relaxation

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



# [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 Distances

**Def.** For  $p \geq 1$ , the  $p$ -Wasserstein distance between  $\mu, \nu$  in  $\mathcal{P}(\Omega)$ , defined by a metric  $D$  on  $\Omega$ ,

$$W_p^p(\mu, \nu) \stackrel{\text{def}}{=} \inf_{P \in \Pi(\mu, \nu)} \iint D(x, y)^p P(dx, dy).$$

PRIMAL

# Wasserstein Distances

**Def.** For  $p \geq 1$ , the  $p$ -Wasserstein distance between  $\mu, \nu$  in  $\mathcal{P}(\Omega)$ , defined by a metric  $D$  on  $\Omega$ ,

$$W_p^p(\mu, \nu) \stackrel{\text{def}}{=} \inf_{P \in \Pi(\mu, \nu)} \iint D(x, y)^p P(dx, dy).$$

PRIMAL

## THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES

BY FRANK L. HITCHCOCK

**1. Statement of the problem.** When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary from city to city.

## МАТЕМАТИЧЕСКИЕ МЕТОДЫ

ОРГАНИЗАЦИИ  
И ПЛАНИРОВАНИЯ  
ПРОИЗВОДСТВА

# Wasserstein Distances

**Def.** For  $p \geq 1$ , the  $p$ -Wasserstein distance between  $\mu, \nu$  in  $\mathcal{P}(\Omega)$ , defined by a metric  $D$  on  $\Omega$ ,

$$W_p^p(\mu, \nu) \stackrel{\text{def}}{=} \inf_{P \in \Pi(\mu, \nu)} \iint D(x, y)^p P(dx, dy).$$

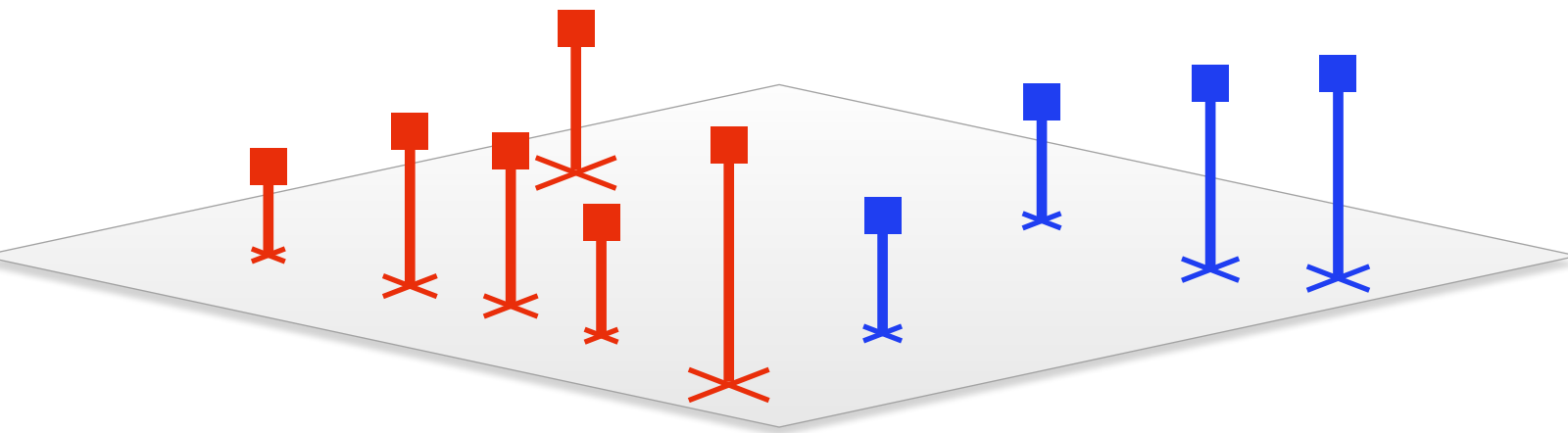
PRIMAL

$$W_p^p(\mu, \nu) = \sup_{\substack{\varphi \in L_1(\mu), \psi \in L_1(\nu) \\ \varphi(x) + \psi(y) \leq D^p(x, y)}} \int \varphi d\mu + \int \psi d\nu.$$

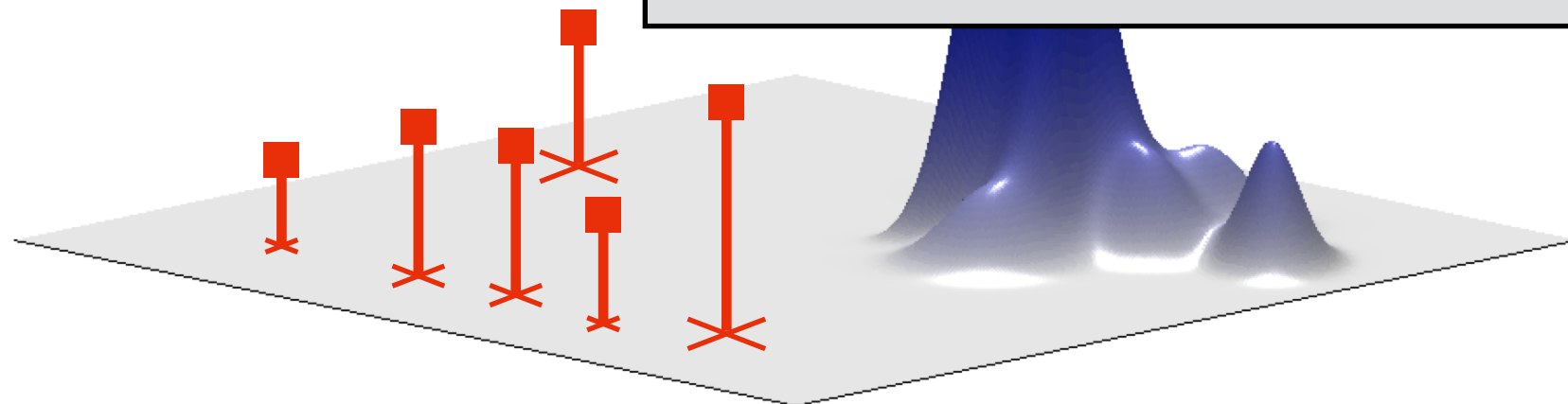
DUAL

# W is versatile

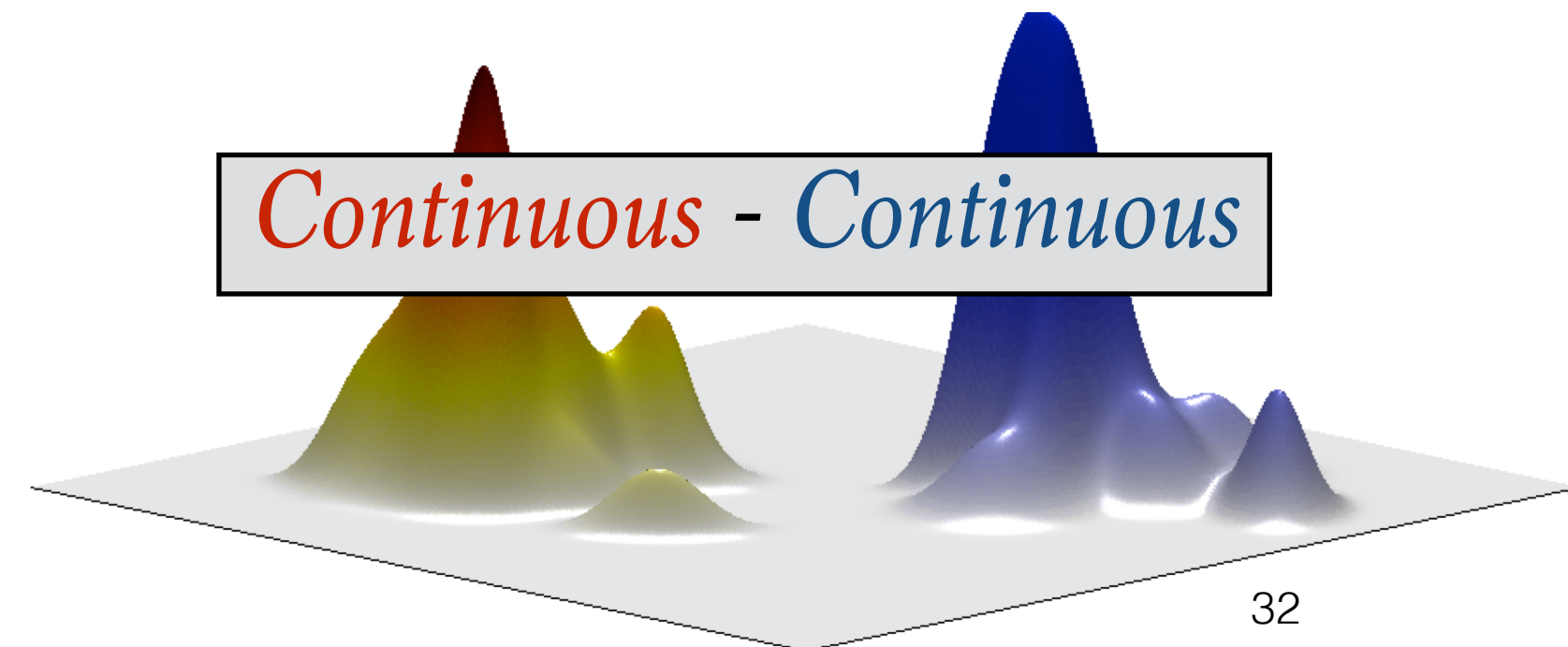
*Discrete* - *Discrete*



*Discrete* - *Continuous*



*Continuous* - *Continuous*

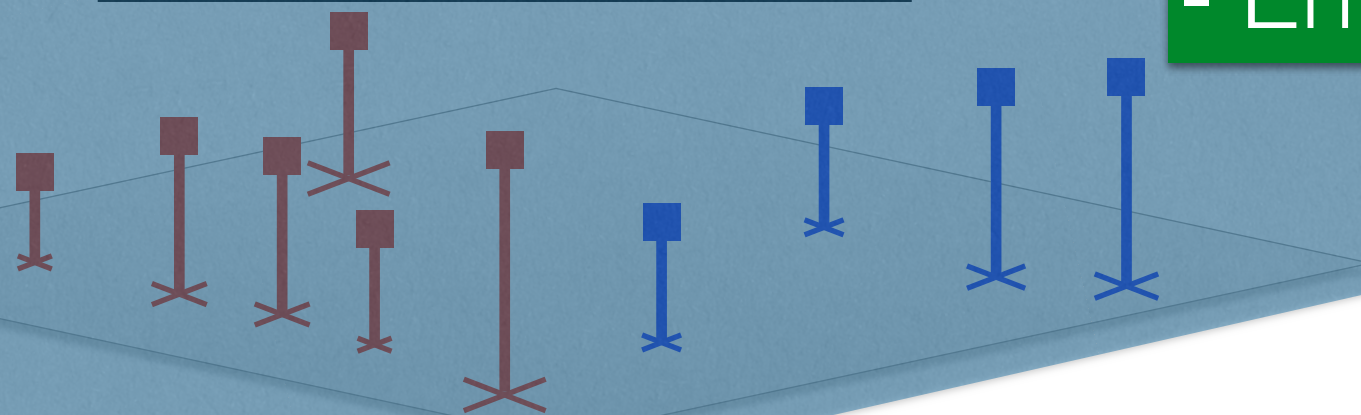




# W is versatile

*Discrete - Discrete*

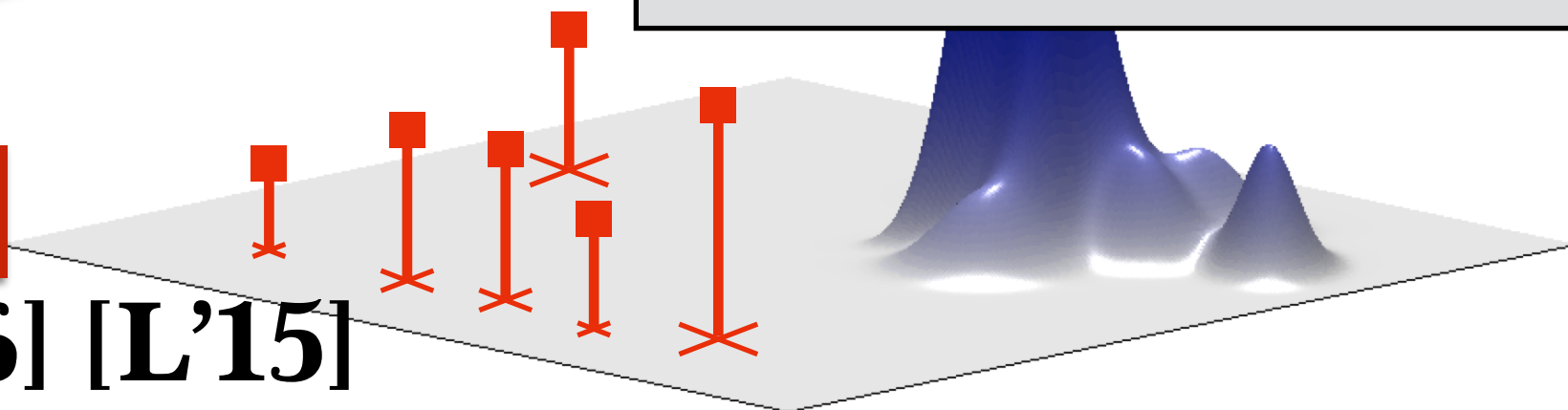
- Network flow solvers
- Entropic regularization



*Discrete - Continuous*

low dim.

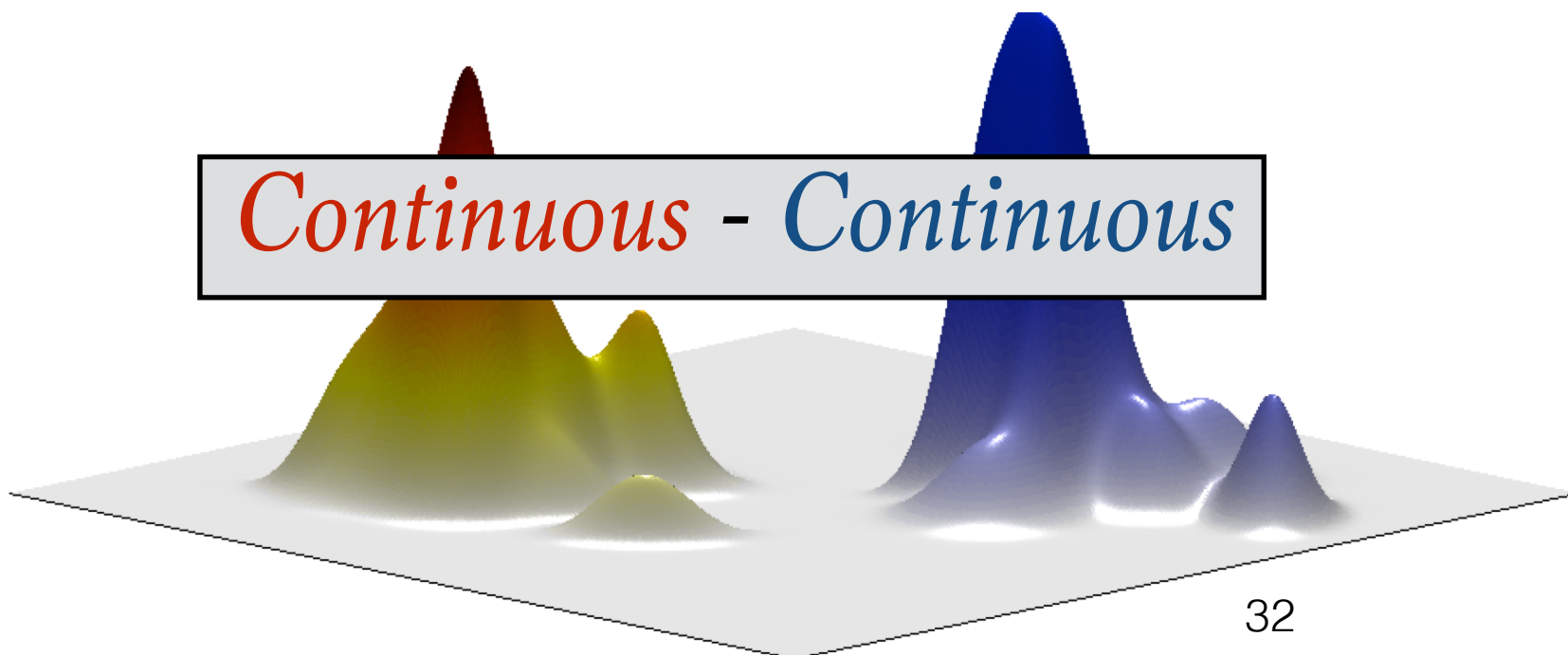
[M'11][KMB'16] [L'15]



*Continuous - Continuous*

Stochastic  
Optimization

[GCPB'16]



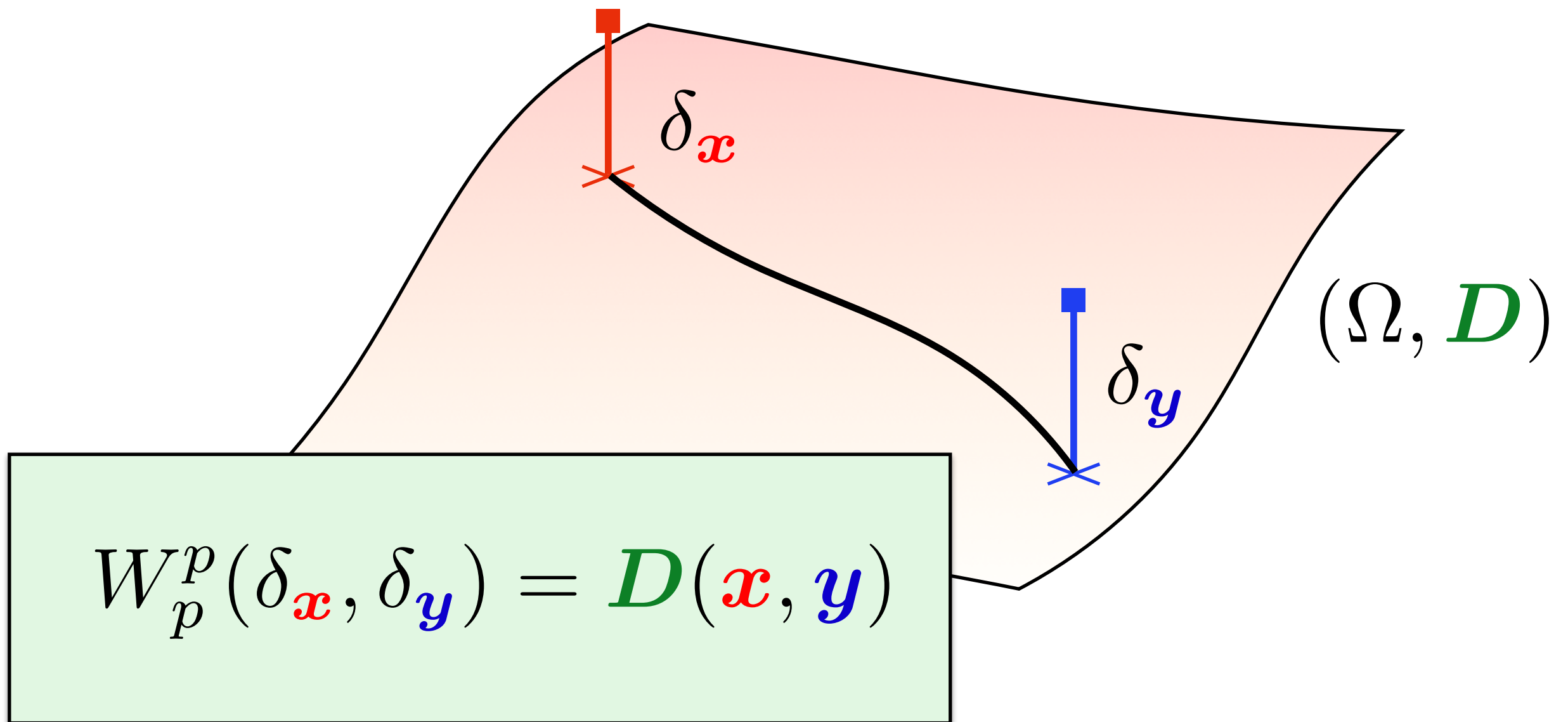
# Minimum Kantorovich Estimators

$$\min_{\theta \in \Theta} W(\nu_{\text{data}}, f_{\theta\#} \mu)$$

- **[Bassetti'06]** 1st reference discussing this approach.
- **[MMC'16]** use regularization in a finite setting.
- **[ACB'17]** (WGAN) **[BJGR'17]** (Wasserstein ABC).
- **Hot topics:** *approximate & differentiate  $W$  efficiently.*
- Today: ideas from our recent preprint **[GPC'17]**

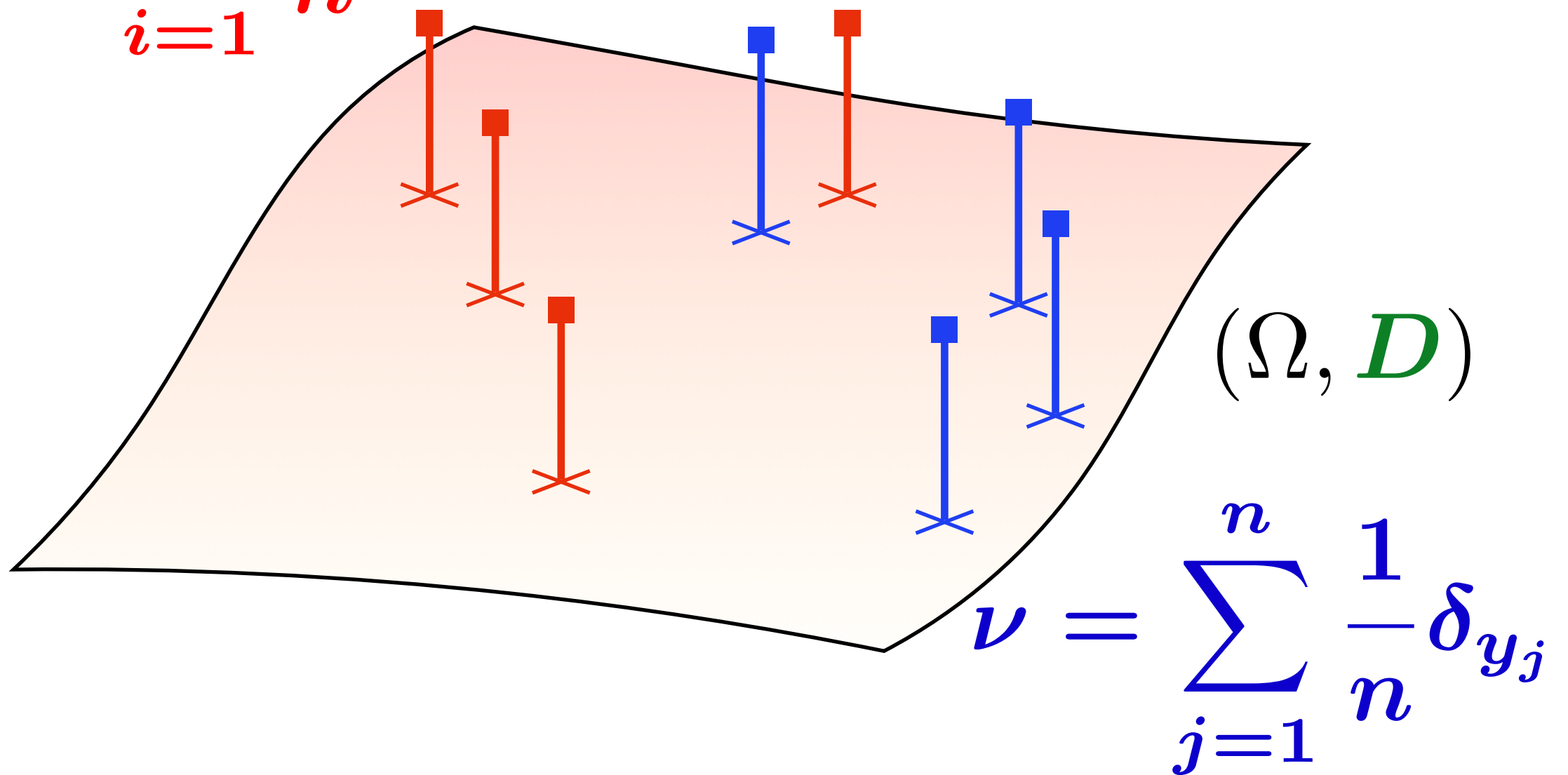


# Wasserstein between 2 Diracs



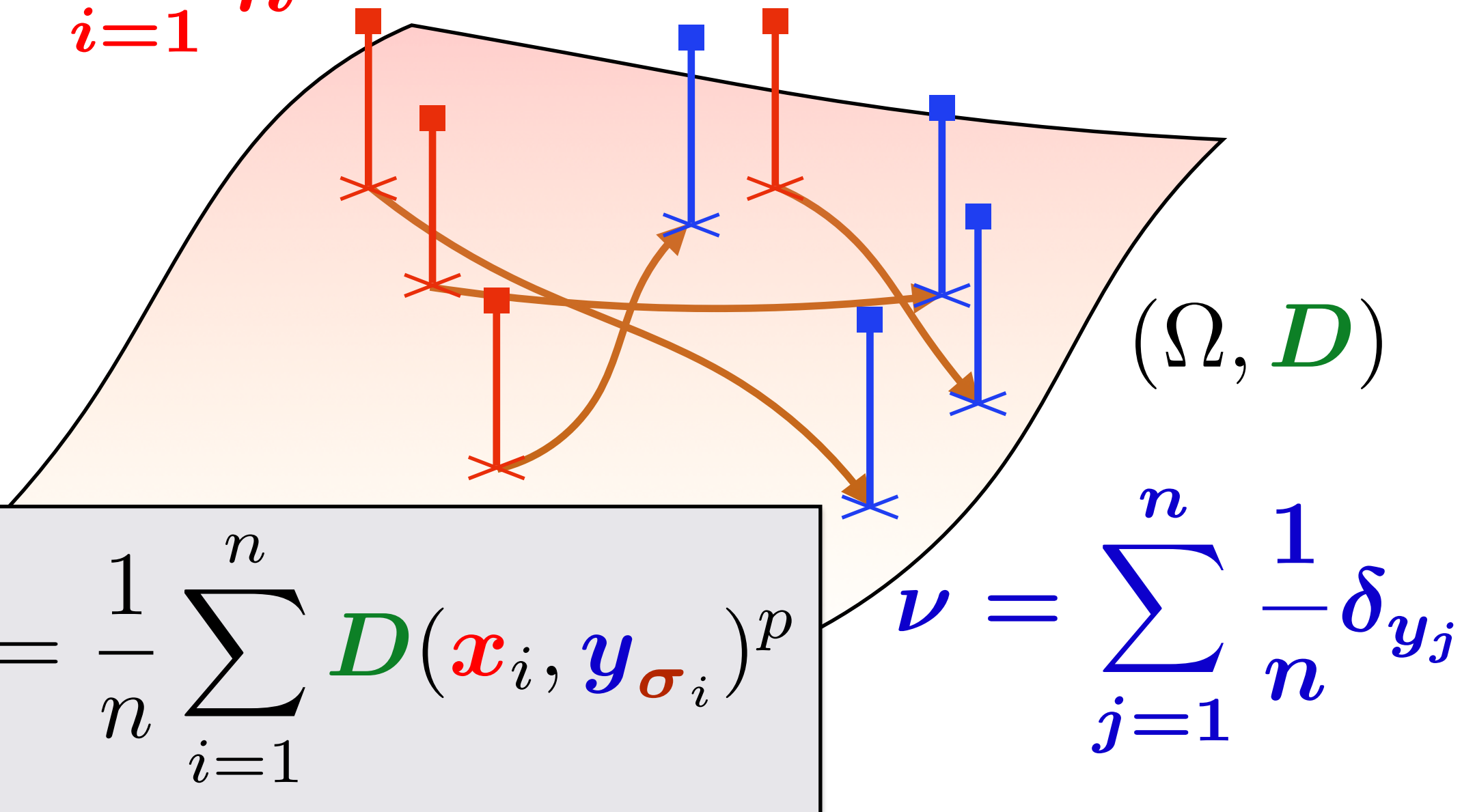
# Wasserstein on Uniform Measures

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



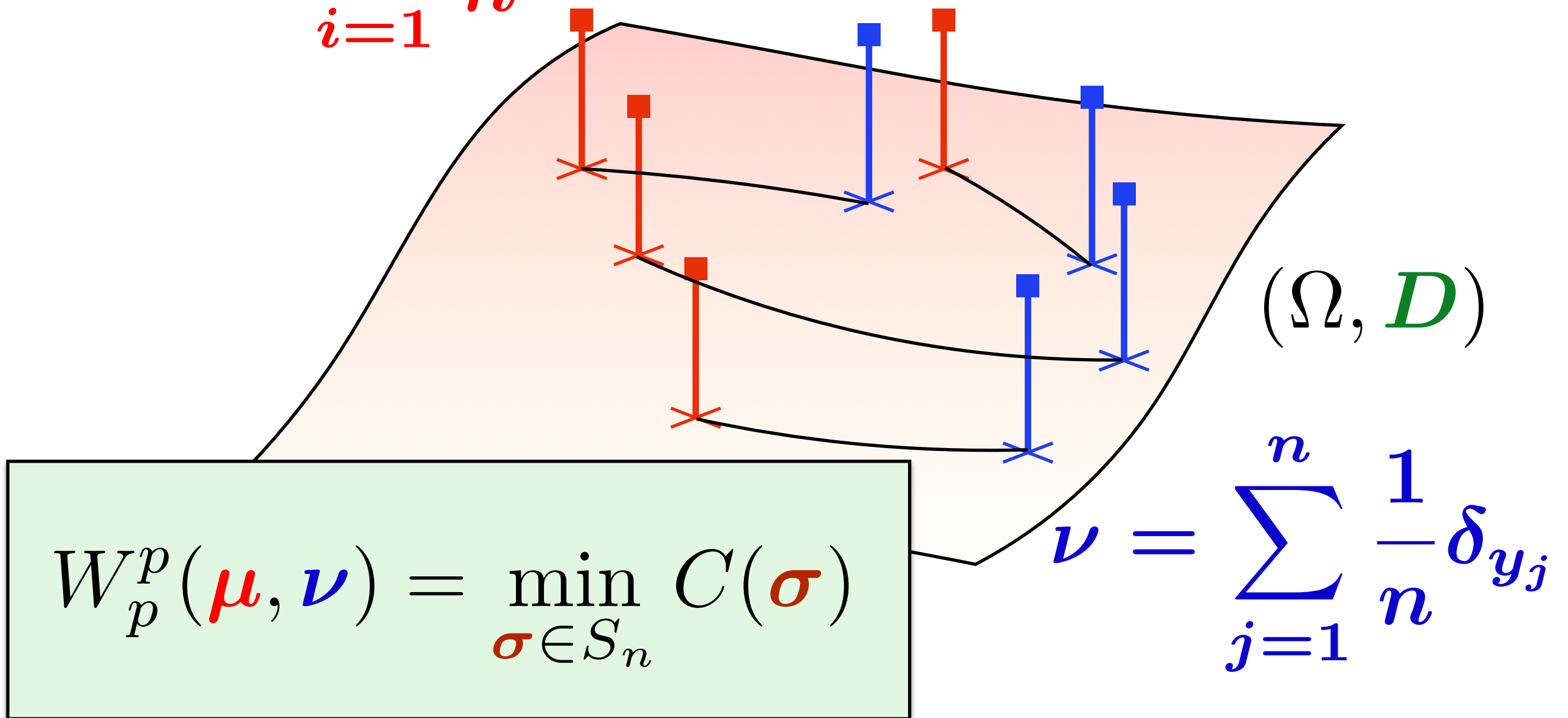
# Wasserstein on Uniform Measures

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



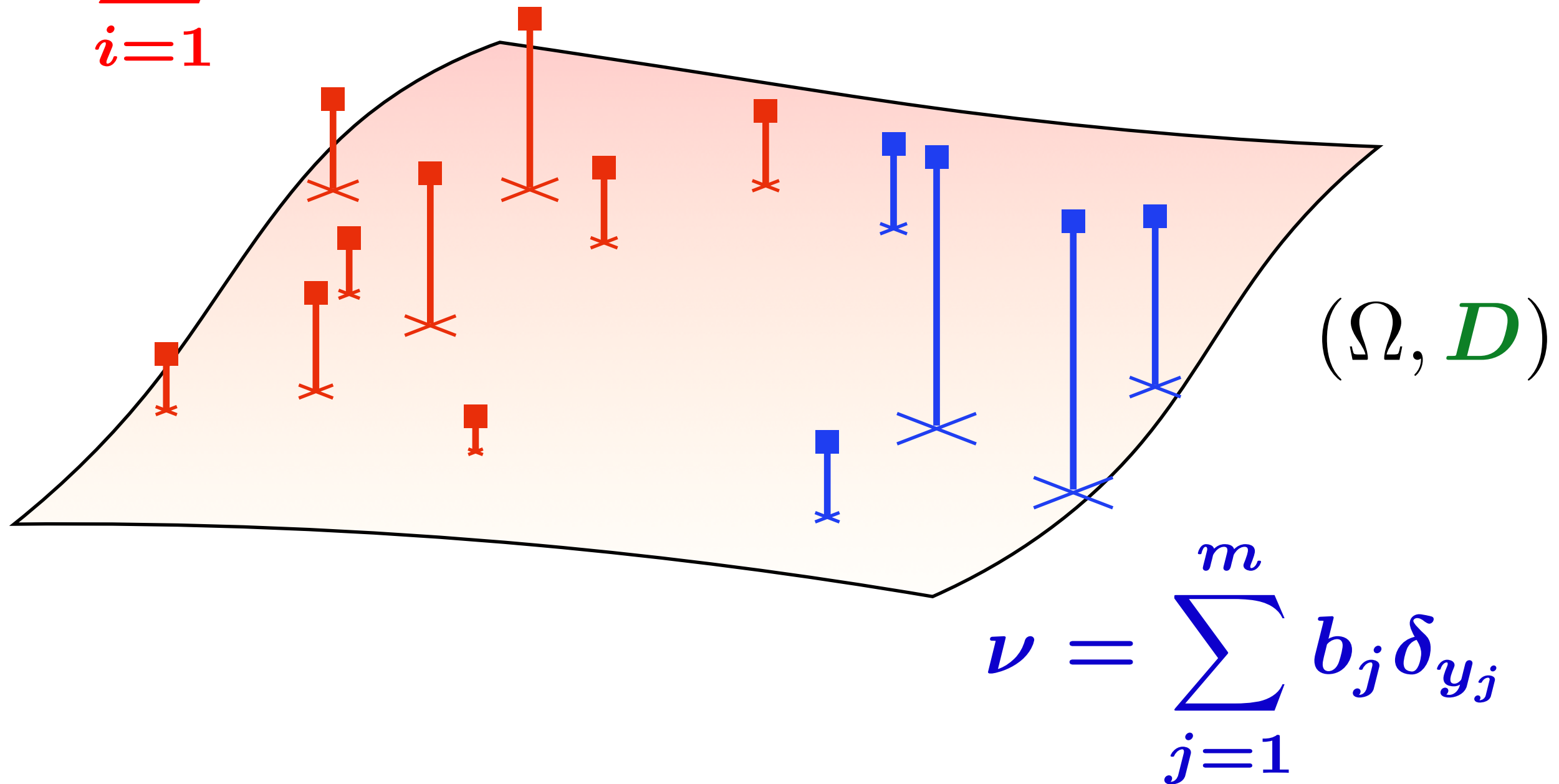
# Optimal Assignment $\subset$ Wasserstein

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



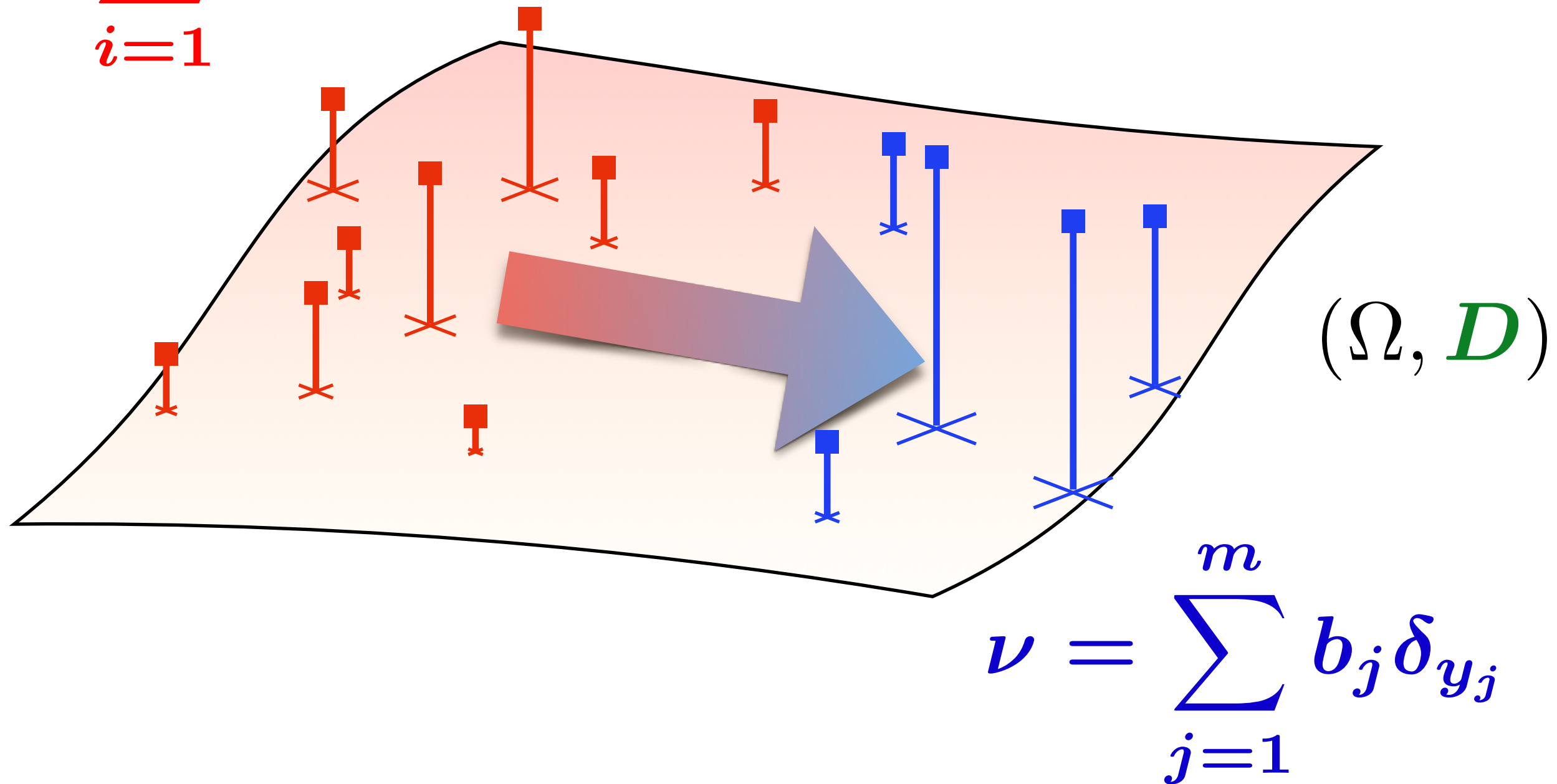
# OT on Two Empirical Measures

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



# OT on Two 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  $\nu = \sum_{j=1}^m b_j \delta_{y_j}$ .

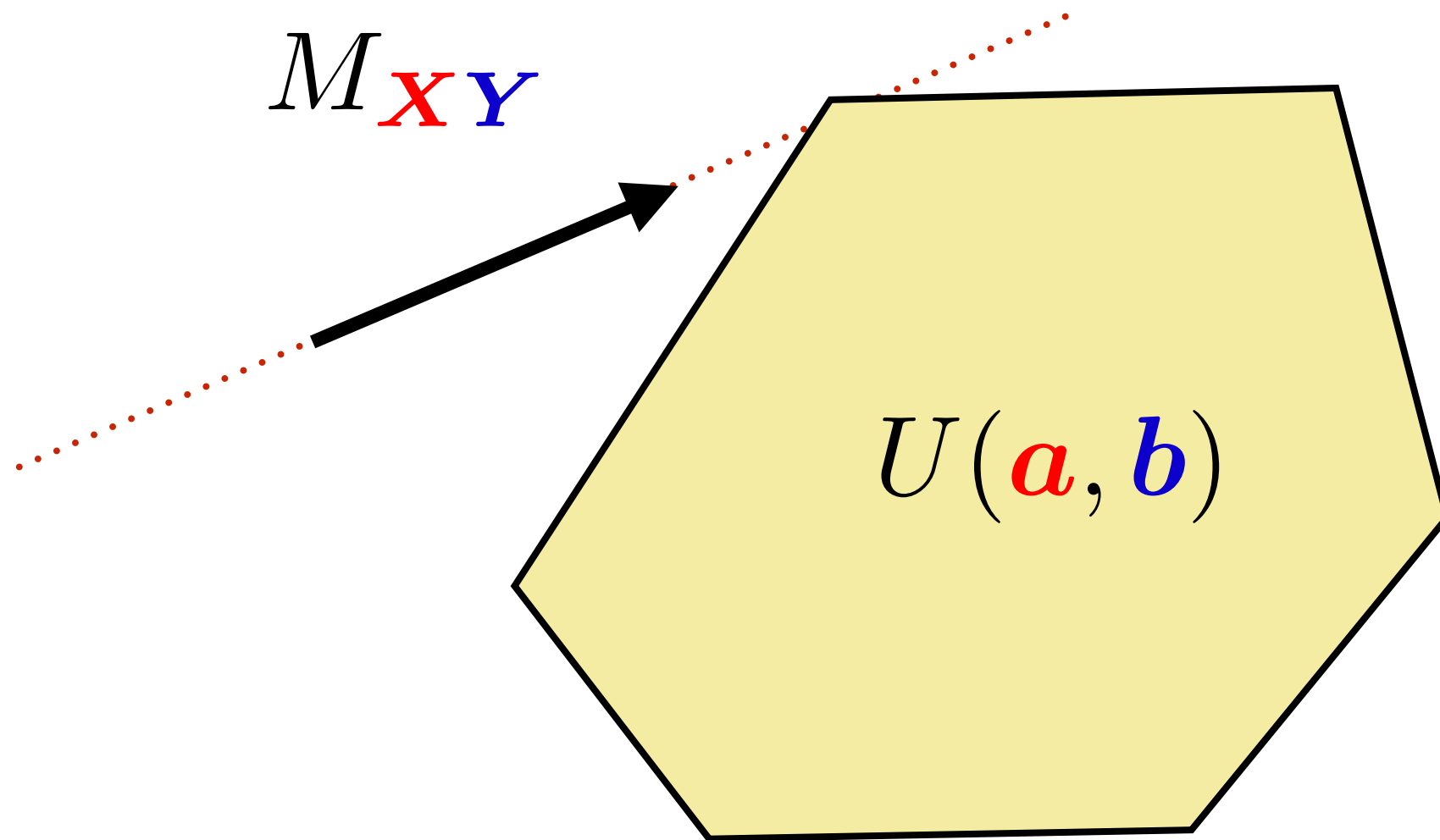
$$M_{\mathbf{x}\mathbf{y}} \stackrel{\text{def}}{=} [D(\mathbf{x}_i, \mathbf{y}_j)^p]_{ij}$$

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

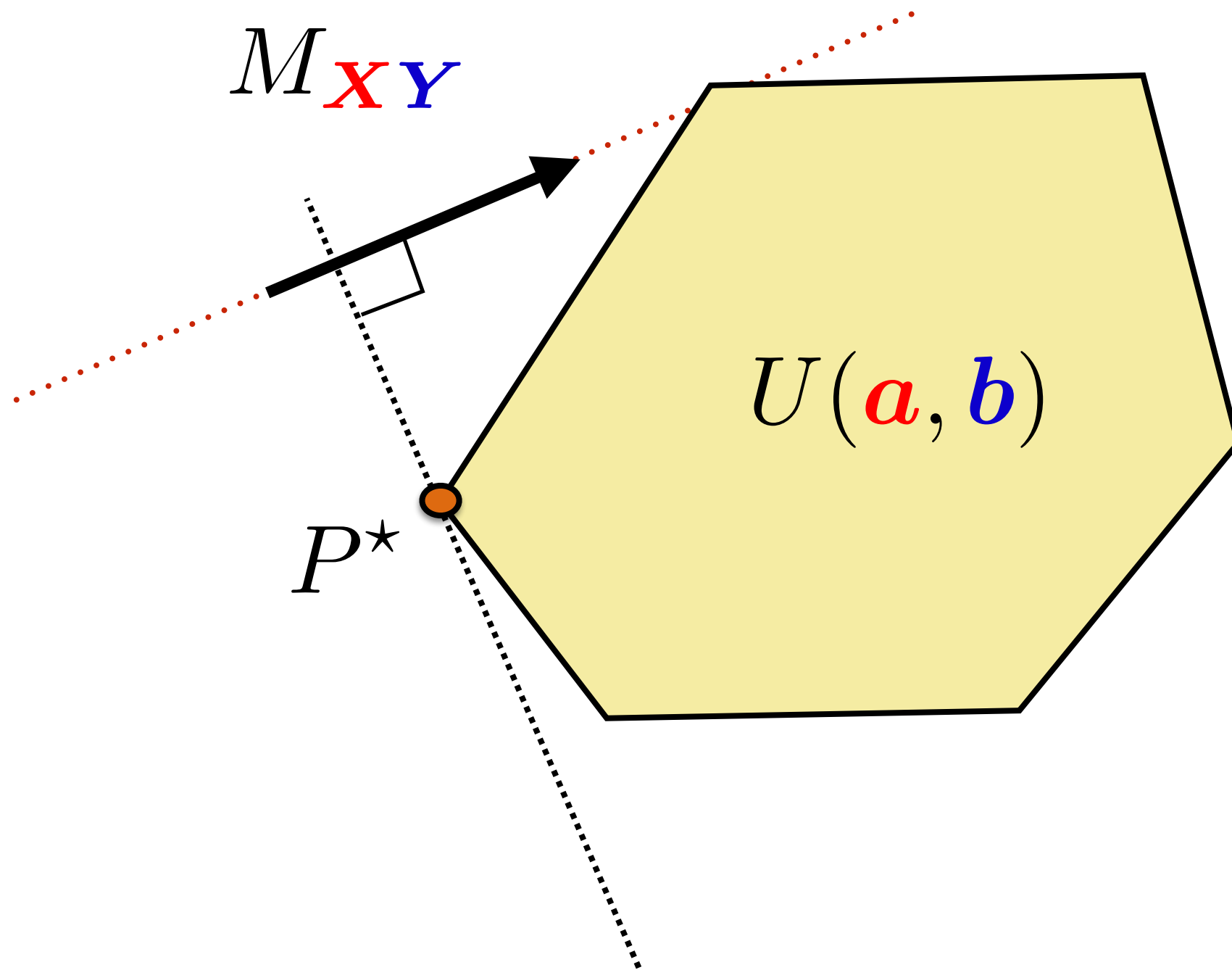
**Def.** Optimal Transport Problem

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

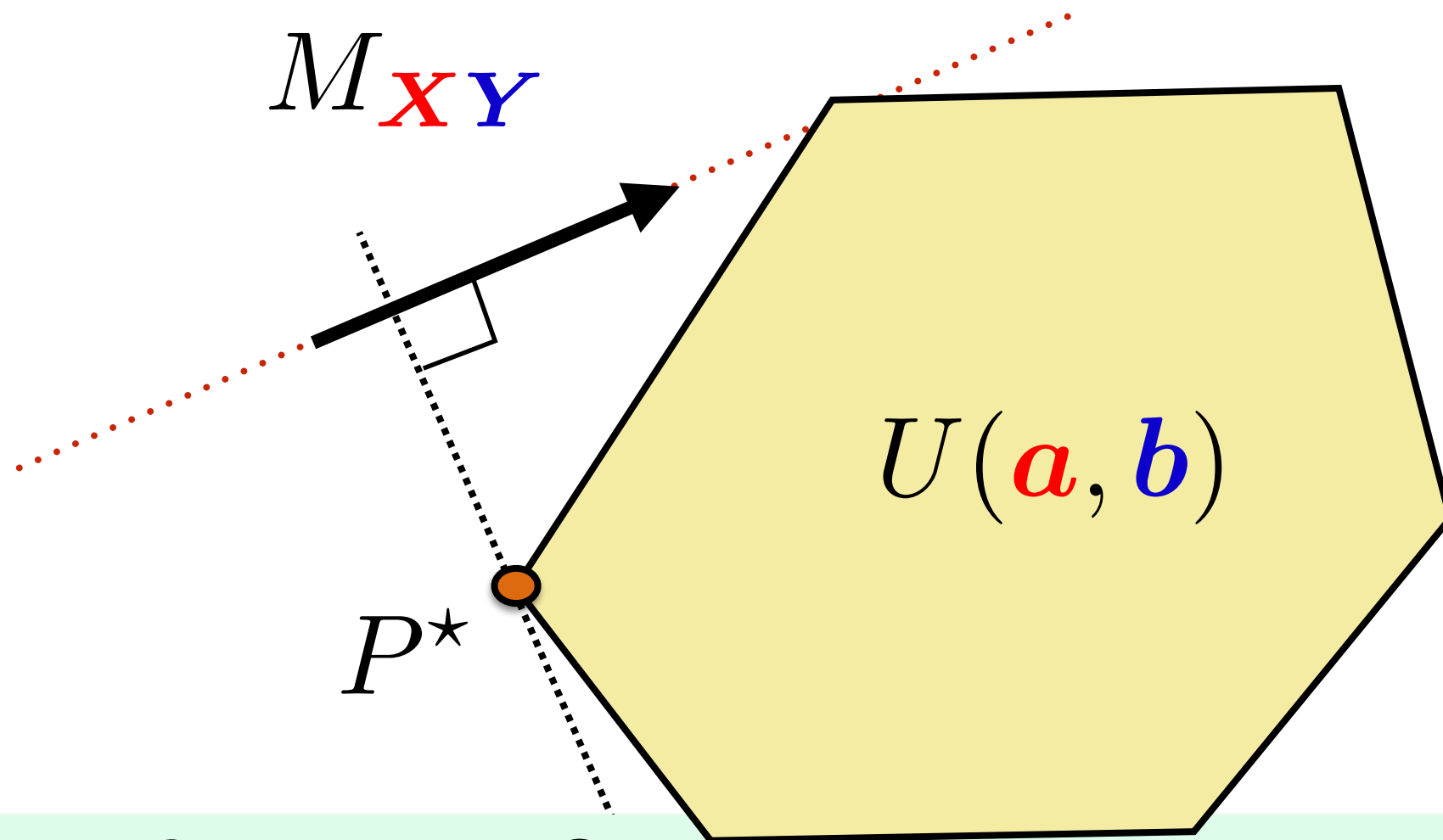
# Discrete OT Problem



# Discrete OT Problem



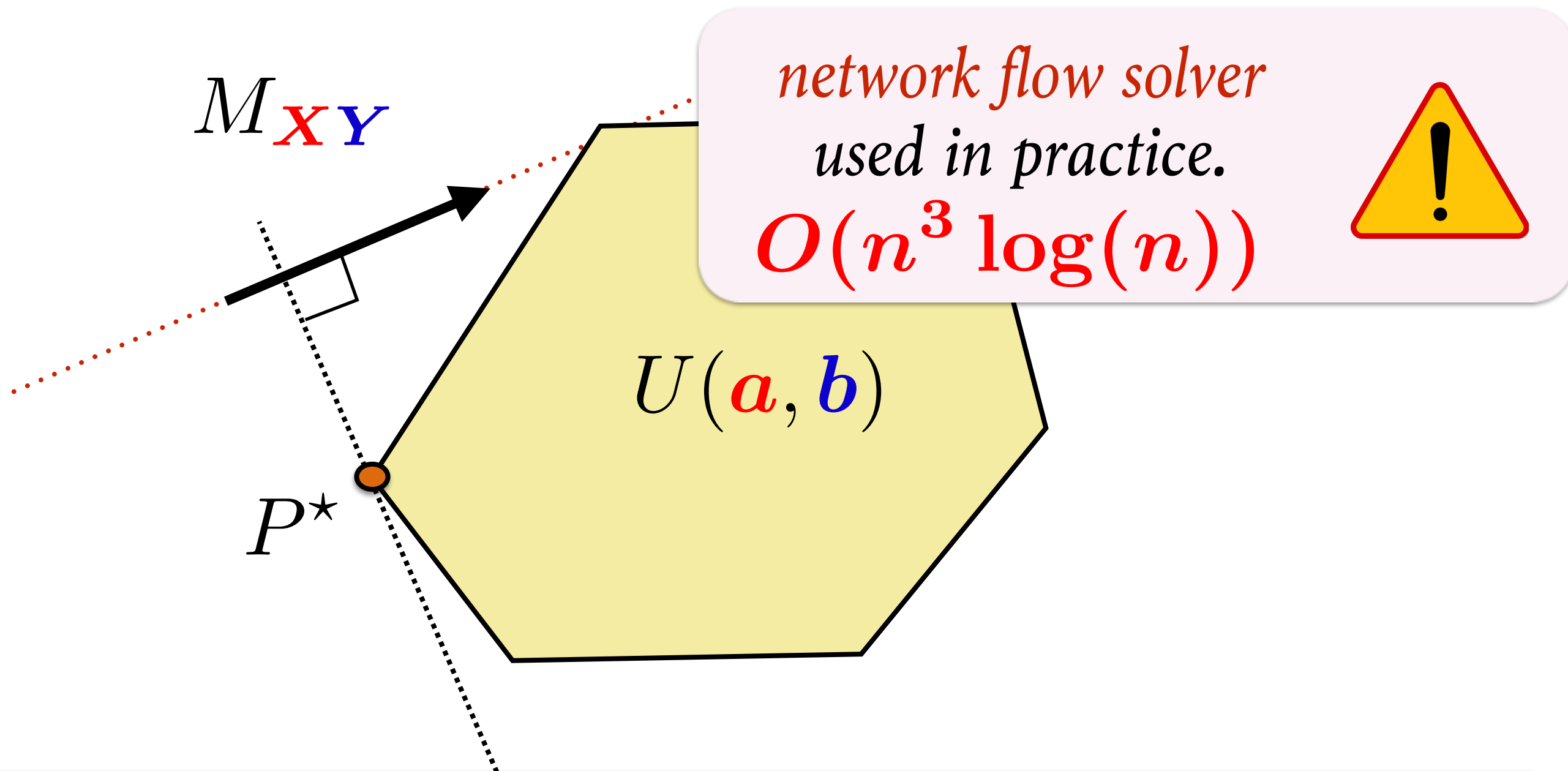
# Discrete OT Problem



**Def.** Dual OT problem

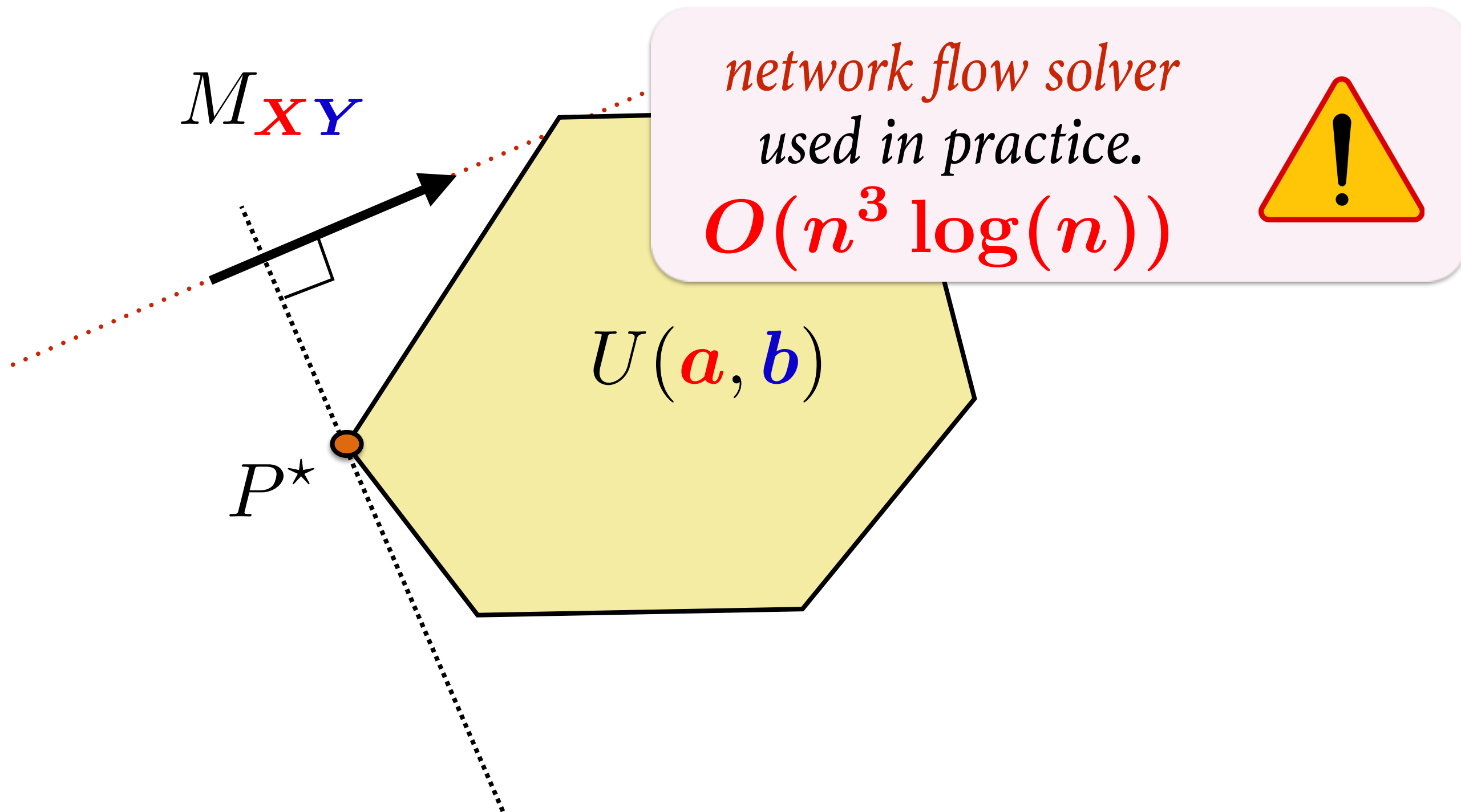
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \boldsymbol{\alpha}_i + \boldsymbol{\beta}_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$

# Discrete OT Problem

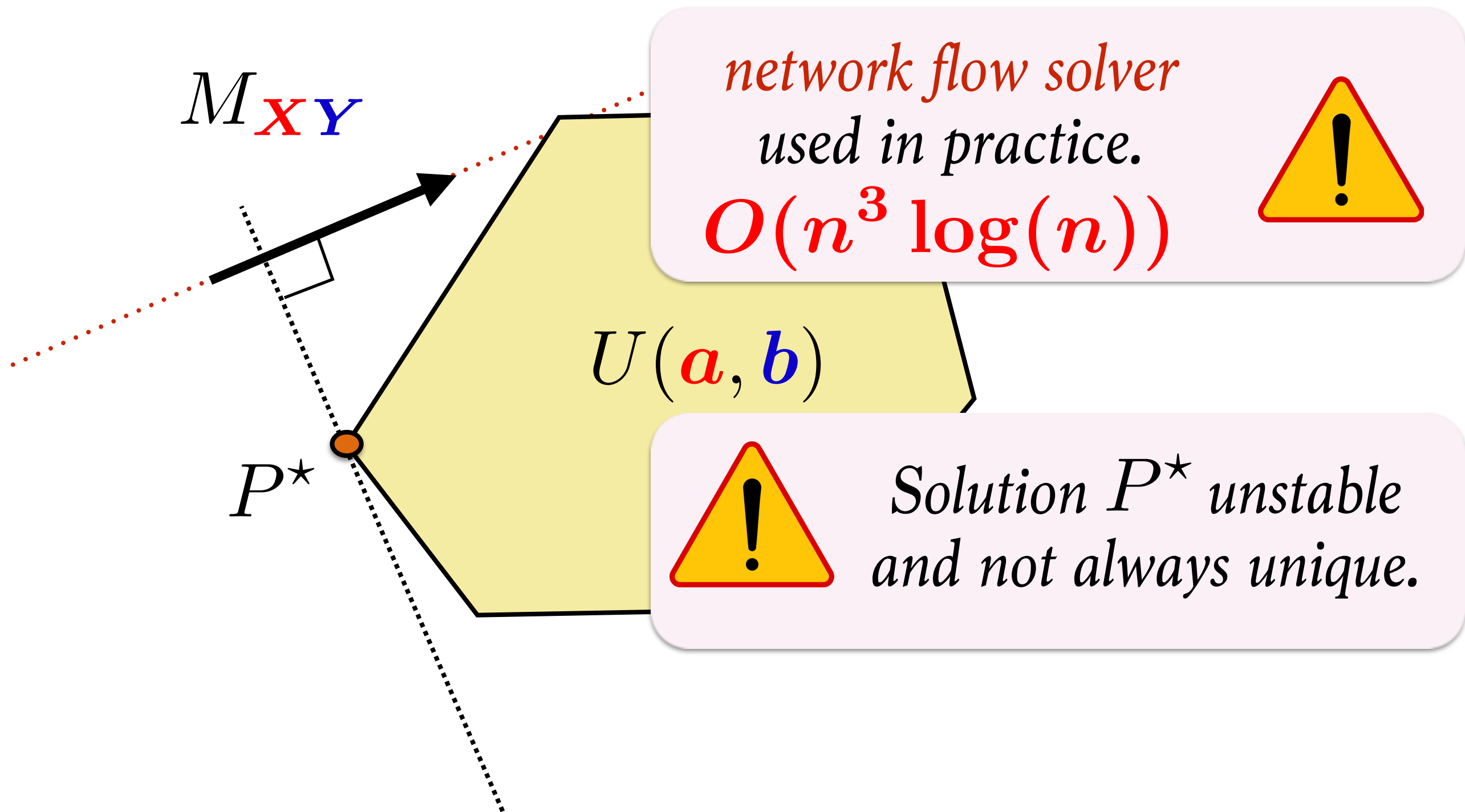


**Note:** flow/PDE formulations [Beckman'61]/[Benamou'98] can be used for  $p=1/p=2$  for a sparse-graph metric/Euclidean metric.

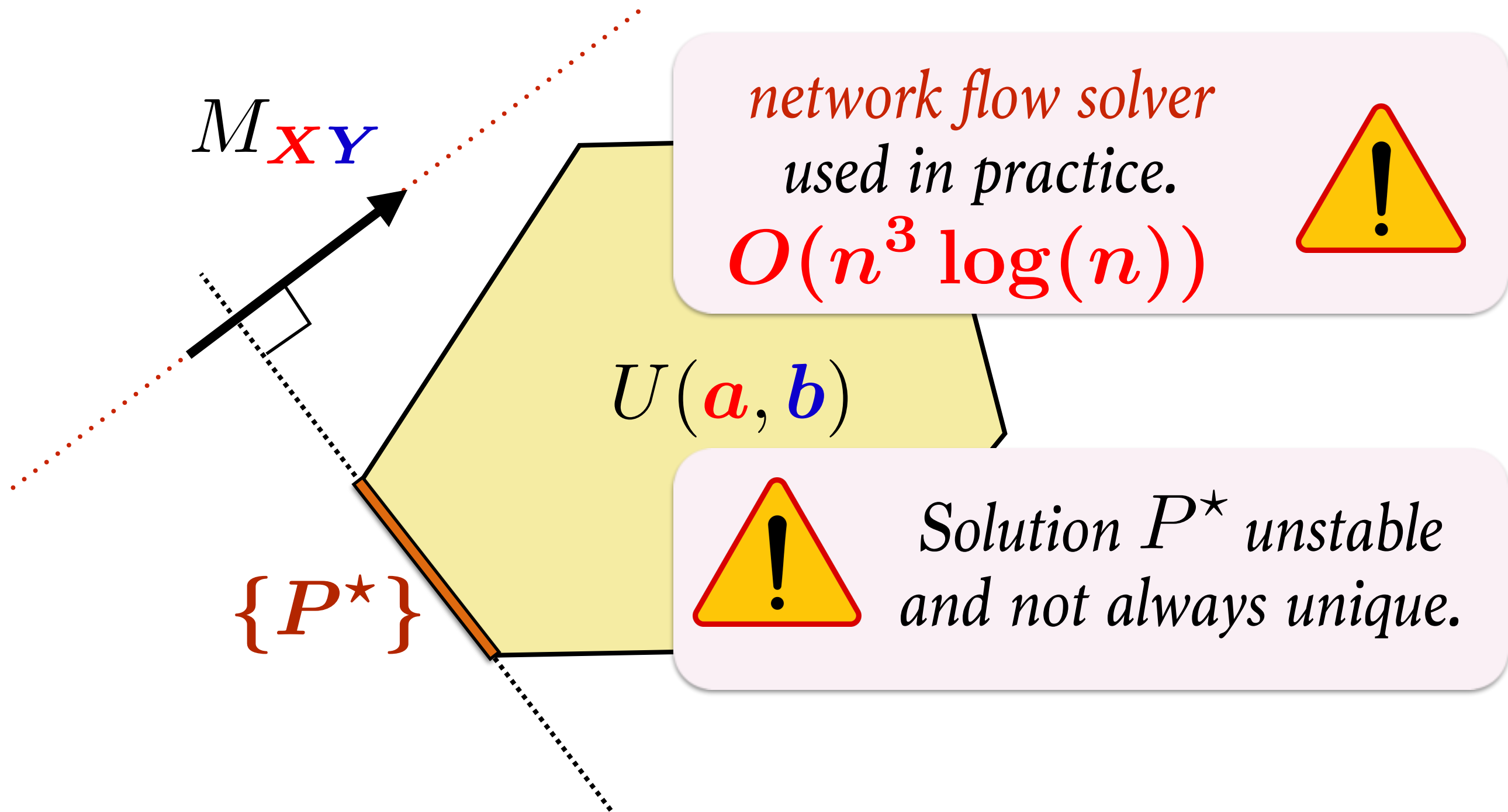
# Discrete OT Problem



# Discrete OT Problem

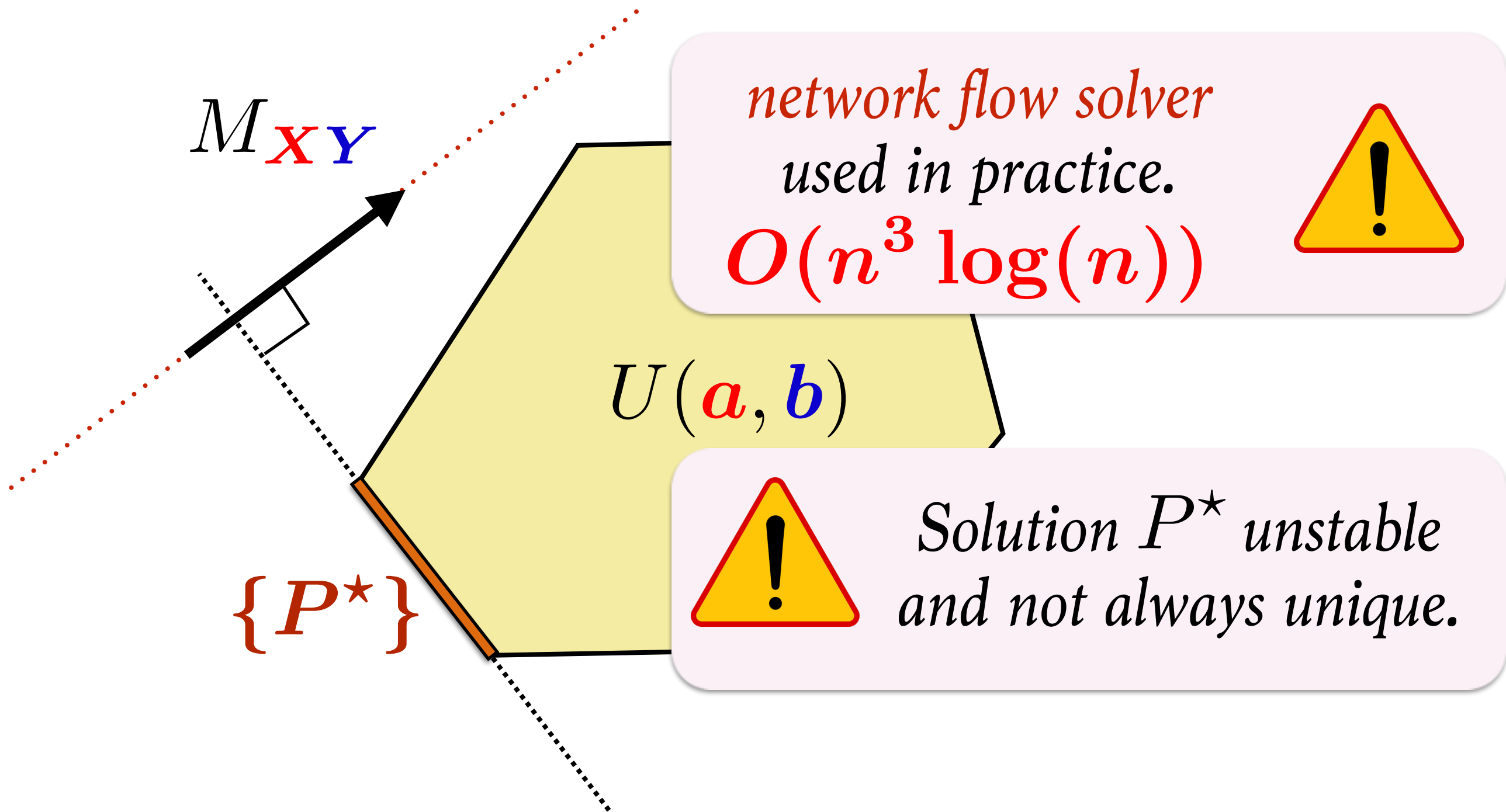


# Discrete OT Problem

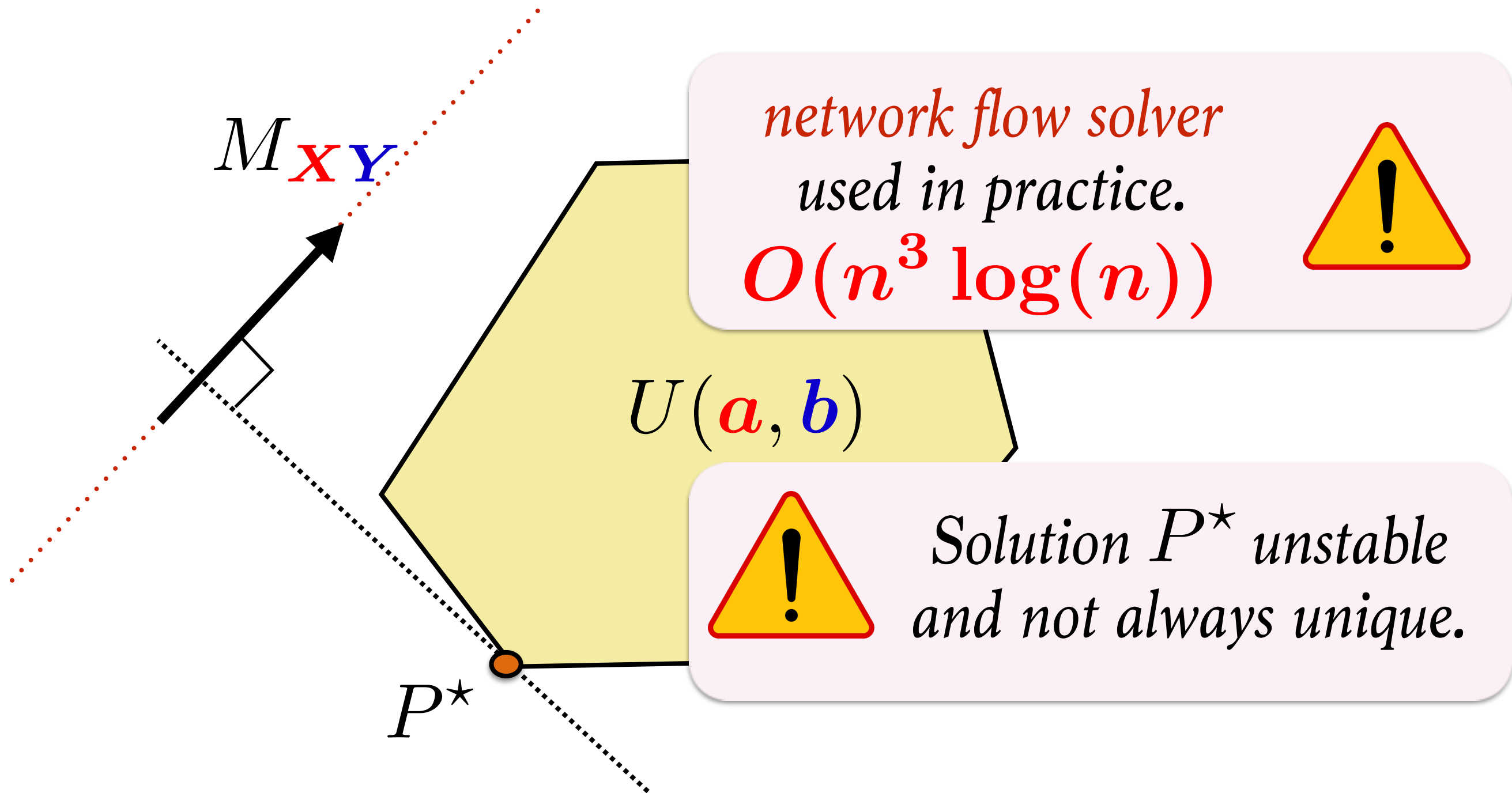




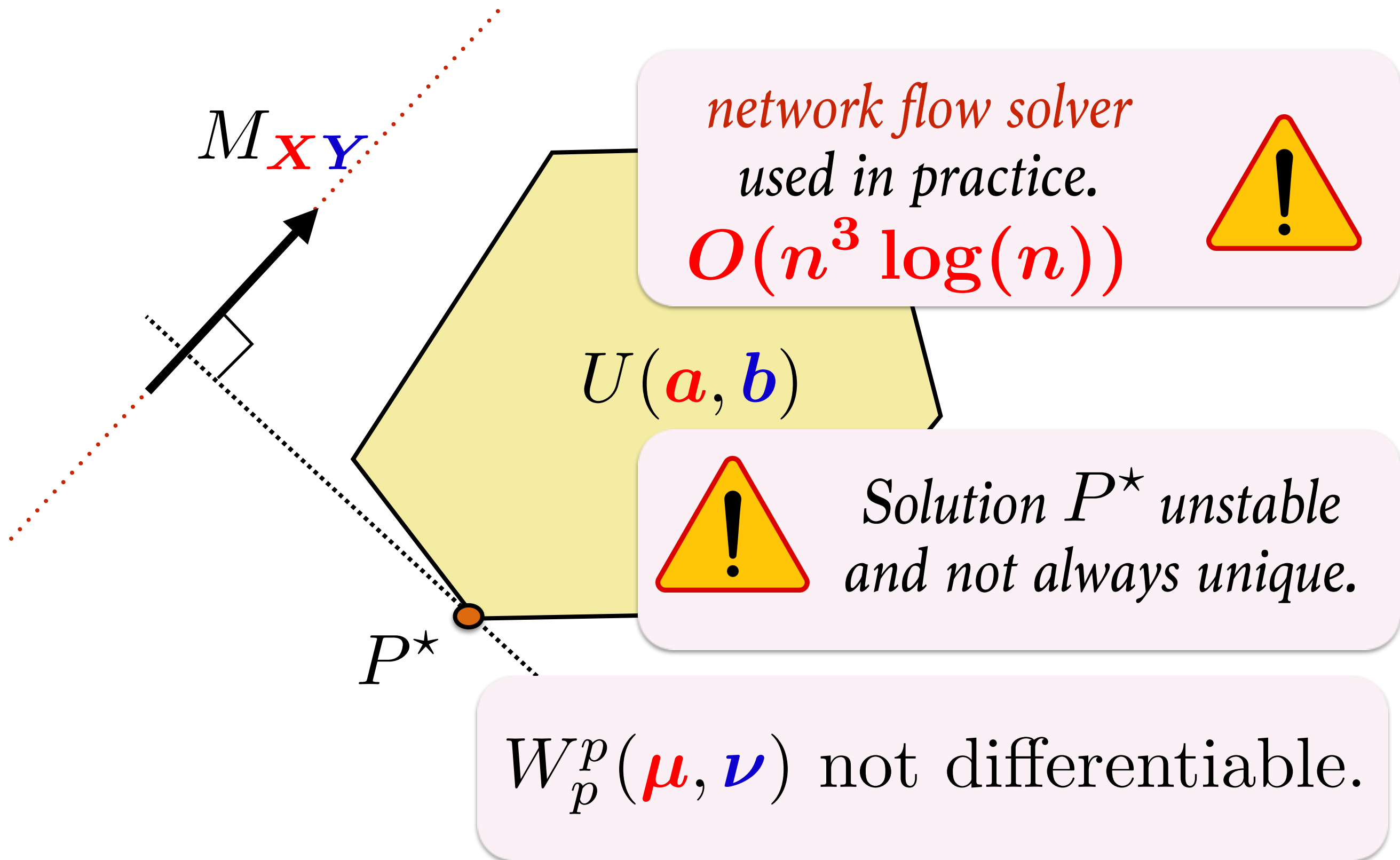
# Discrete OT Problem



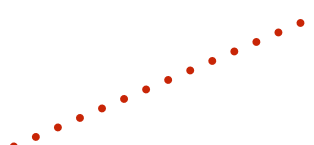
# Discrete OT Problem



# Discrete OT Problem



# Discrete OT Problem



```
emd.c
Last update: 3/14/98

An implementation of the Earth Movers Distance.
Based of the solution for the Transportation problem as described in
"Introduction to Mathematical Programming" by F. S. Hillier and
G. J. Lieberman, McGraw-Hill, 1990.

Copyright (C) 1998 Yossi Rubner
Computer Science Department, Stanford University
E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner

/*
#include <stdio.h>
#include <stdlib.h>
#include <math.h>

#include "emd.h"

#define DEBUG_LEVEL 0
/*
DEBUG_LEVEL:
0 = NO MESSAGES
1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
2 = PRINT THE RESULT AFTER EVERY ITERATION
3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
*/

#define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSSIBLE DUMMY FEATURE */

/* NEW TYPES DEFINITION */

/* node1_t IS USED FOR SINGLE-LINKED LISTS */
typedef struct node1_t {
    int i;
    double val;
    struct node1_t *Next;
} node1_t;

/* node2_t IS USED FOR DOUBLE-LINKED LISTS */
typedef struct node2_t {
    int i, j;
    double val;
    struct node2_t *NextC; /* NEXT COLUMN */
    struct node2_t *NextR; /* NEXT ROW */
} node2_t;

/* GLOBAL VARIABLE DECLARATION */
static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1]; /* THE COST MATRIX */
static node2_t _X[MAX_SIG_SIZE1*2]; /* THE BASIC VARIABLES VECTOR */
```

# Discrete OT Problem

```
emd.c
Last update: 3/14/98

An implementation of the Earth Movers Distance.
Based of the solution for the Transportation problem as described in
"Introduction to Mathematical Programming" by F. S. Hillier and
G. J. Lieberman, McGraw-Hill, 1990.

Copyright (C) 1998 Yossi Rubner
Computer Science Department, Stanford University
E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner

/*
#include <stdio.h>
#include <stdlib.h>
#include <math.h>

#include "emd.h"

#define DEBUG_LEVEL 0
/*
DEBUG_LEVEL:
0 = NO MESSAGES
1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
2 = PRINT THE RESULT AFTER EVERY ITERATION
3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
*/

#define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSSIBLE DUMMY FEATURE */

/* NEW TYPES DEFINITION */

/* node1_t IS USED FOR SINGLE-LINKED LISTS */
typedef struct node1_t {
    int i;
    double val;
    struct node1_t *Next;
} node1_t;

/* node2_t IS USED FOR DOUBLE-LINKED LISTS */
typedef struct node2_t {
    int i, j;
    double val;
    struct node2_t *NextC; /* NEXT COLUMN */
    struct node2_t *NextR; /* NEXT ROW */
} node2_t;

/* GLOBAL VARIABLE DECLARATION */
static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1]; /* THE COST MATRIX */
static node2_t _X[MAX_SIG_SIZE1*2]; /* THE BASIC VARIABLES VECTOR */
```





# Discrete OT Problem

```
emd.c
Last update: 3/14/98

An implementation of the Earth Movers Distance.
Based of the solution for the Transportation problem as described in
"Introduction to Mathematical Programming" by F. S. Hillier and
G. J. Lieberman, McGraw-Hill, 1990.

Copyright (C) 1998 Yossi Rubner
Computer Science Department, Stanford University
E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner

/*
#include <stdio.h>
#include <stdlib.h>
#include <math.h>

#include "emd.h"

#define DEBUG_LEVEL 0
/*
DEBUG_LEVEL:
0 = NO MESSAGES
1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
2 = PRINT THE RESULT AFTER EVERY ITERATION
3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
*/

#define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSSIBLE DUMMY FEATURE */

/* NEW TYPES DEFINITION */

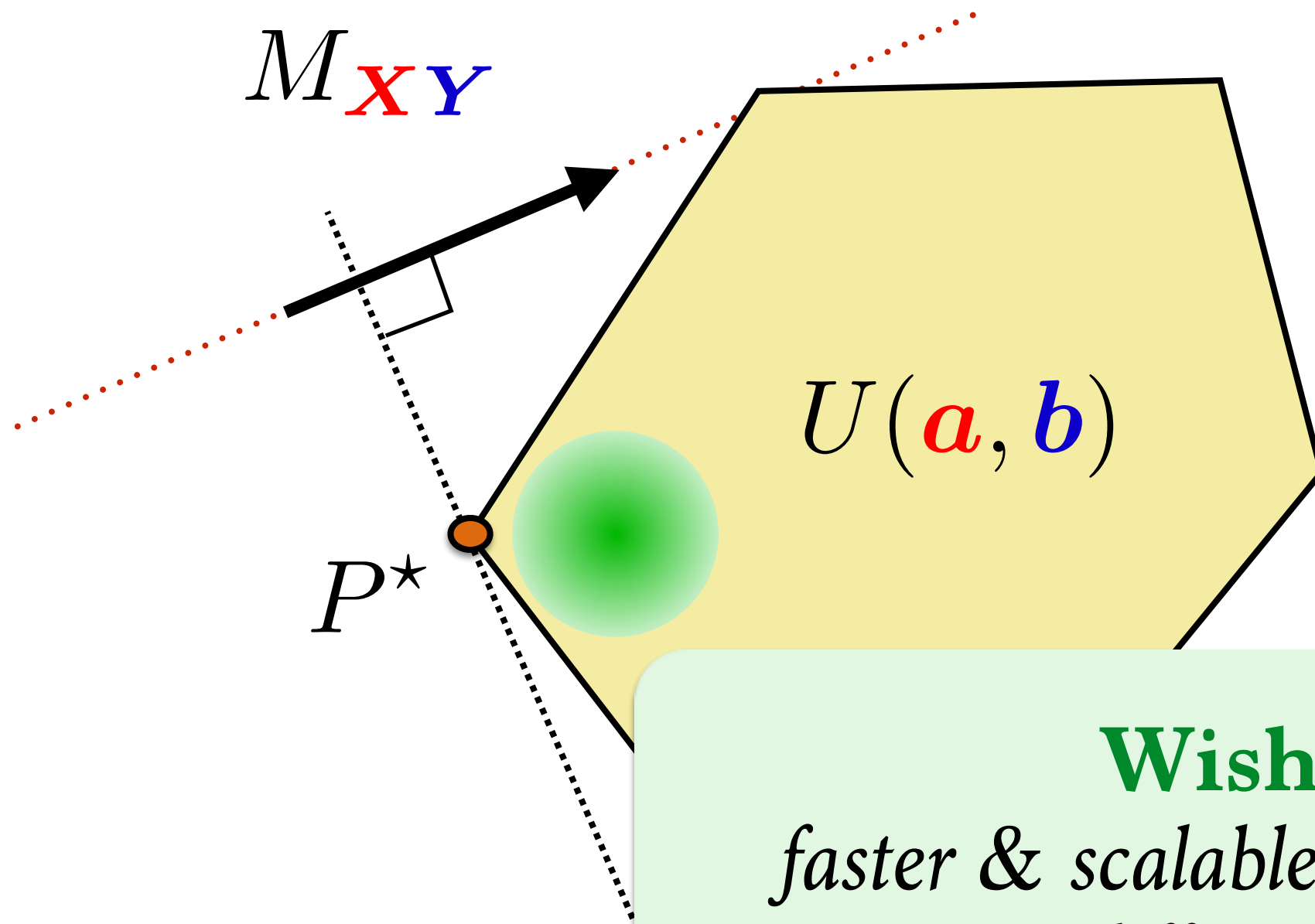
/* node1_t IS USED FOR SINGLE-LINKED LISTS */
typedef struct node1_t {
    int i;
    double val;
    struct node1_t *Next;
} node1_t;

/* node2_t IS USED FOR DOUBLE-LINKED LISTS */
typedef struct node2_t {
    int i, j;
    double val;
    struct node2_t *NextC; /* NEXT COLUMN */
    struct node2_t *NextR; /* NEXT ROW */
} node2_t;

/* GLOBAL VARIABLE DECLARATION */
static int _n1, _n2; /* SIGNATURES SIZES */
static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1]; /* THE COST MATRIX */
static node2_t _X[MAX_SIG_SIZE1*2]; /* THE BASIC VARIABLES VECTOR */
```



# Solution: Modify OT Problem



**Wishlist:**  
*faster & scalable, more stable,  
differentiable*

# Entropic Regularization [Wilson'62]

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

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

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

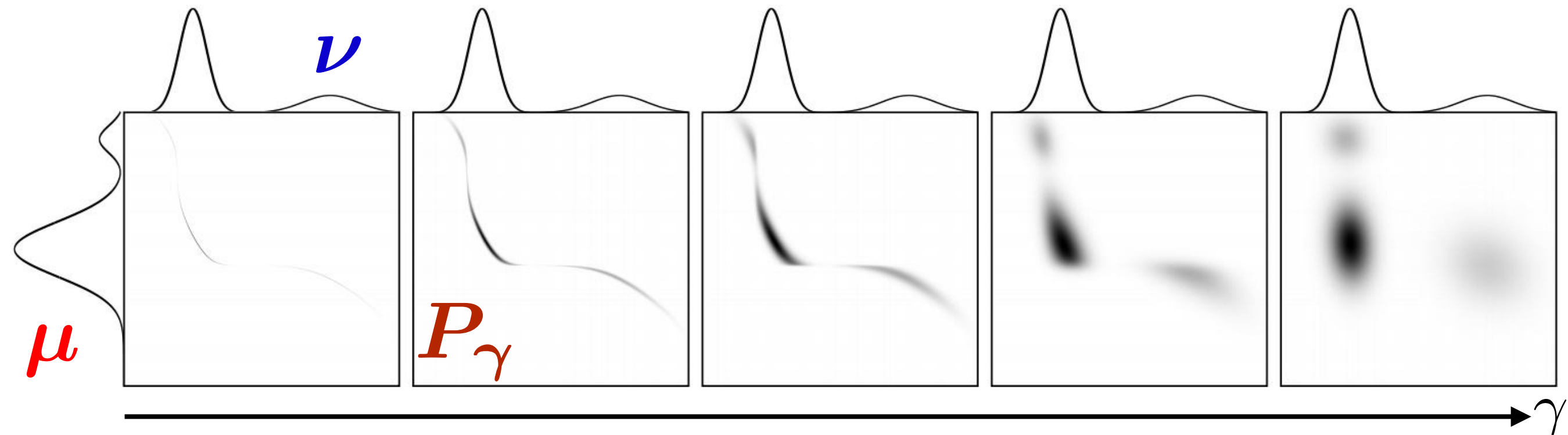
**Note:** Unique optimal solution because of strong concavity of Entropy



# Entropic Regularization [Wilson'62]

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

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



**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

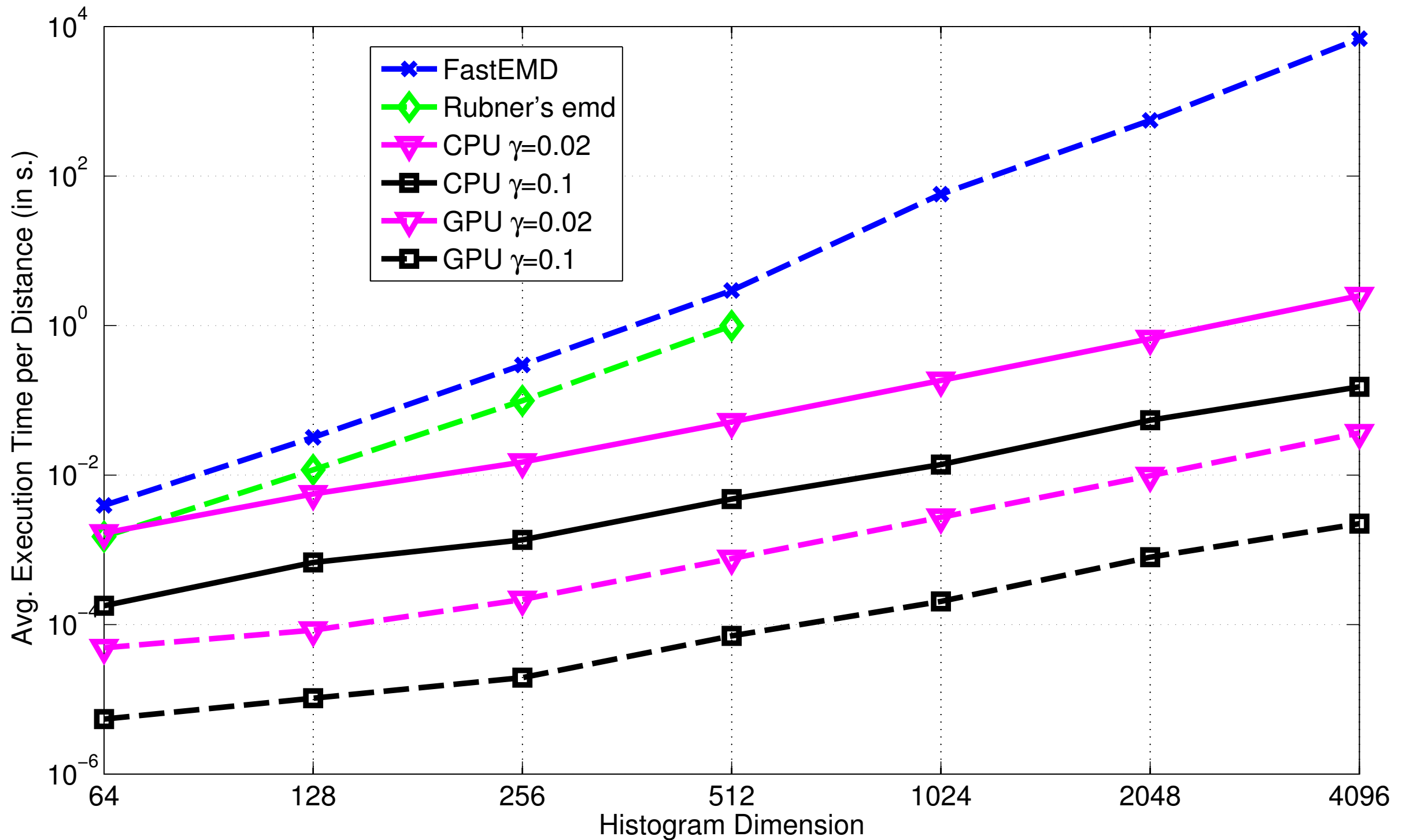
**Prop.** If  $P_\gamma \stackrel{\text{def}}{=} \underset{P \in U(\mathbf{a}, \mathbf{b})}{\operatorname{argmin}} \langle \mathbf{P}, M_{\mathbf{x}\mathbf{y}} \rangle - \gamma E(\mathbf{P})$

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  $D^p$  separable.  
[S..C..'15]

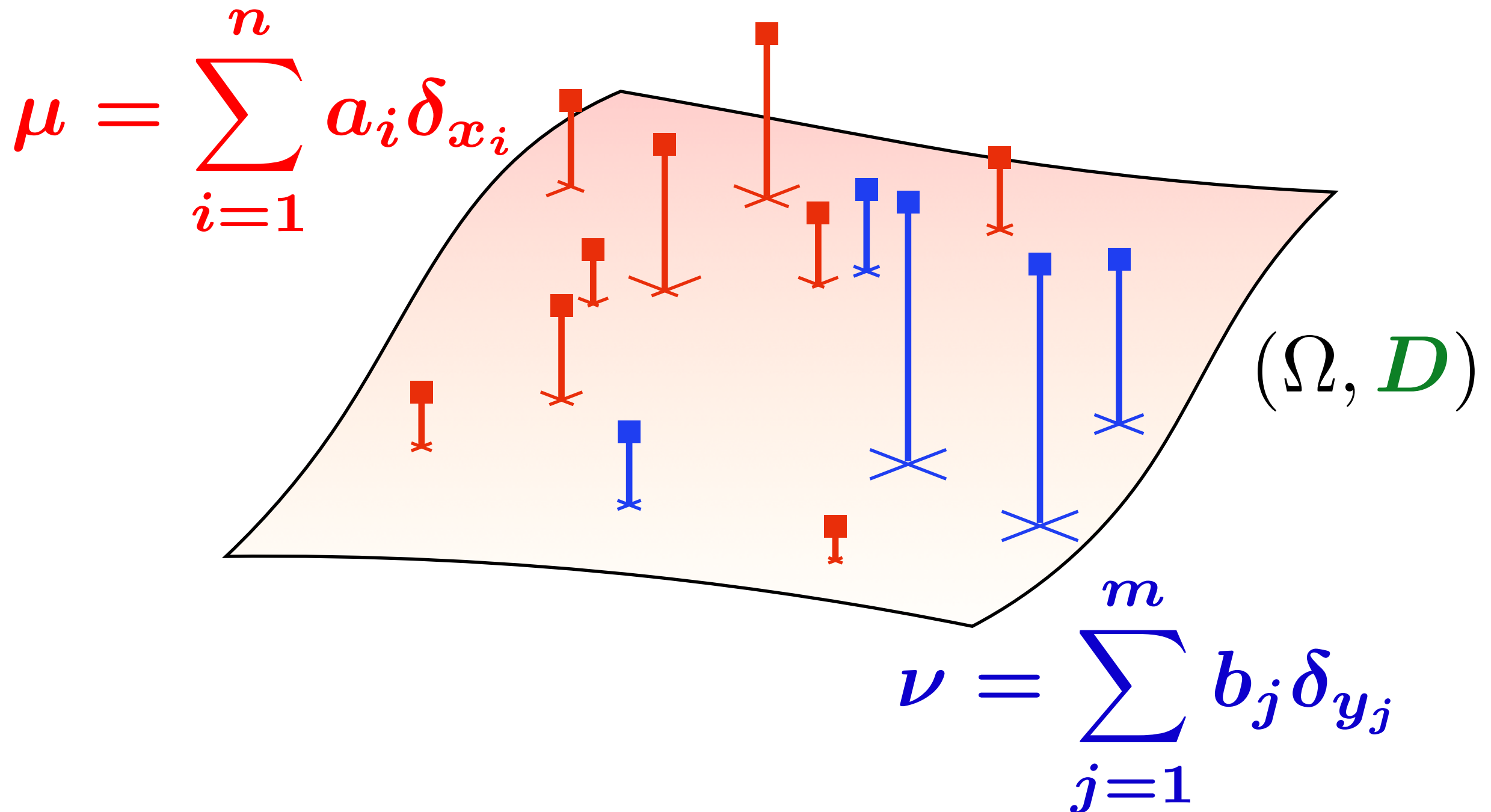
# Very Fast EMD Approx. Solver



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

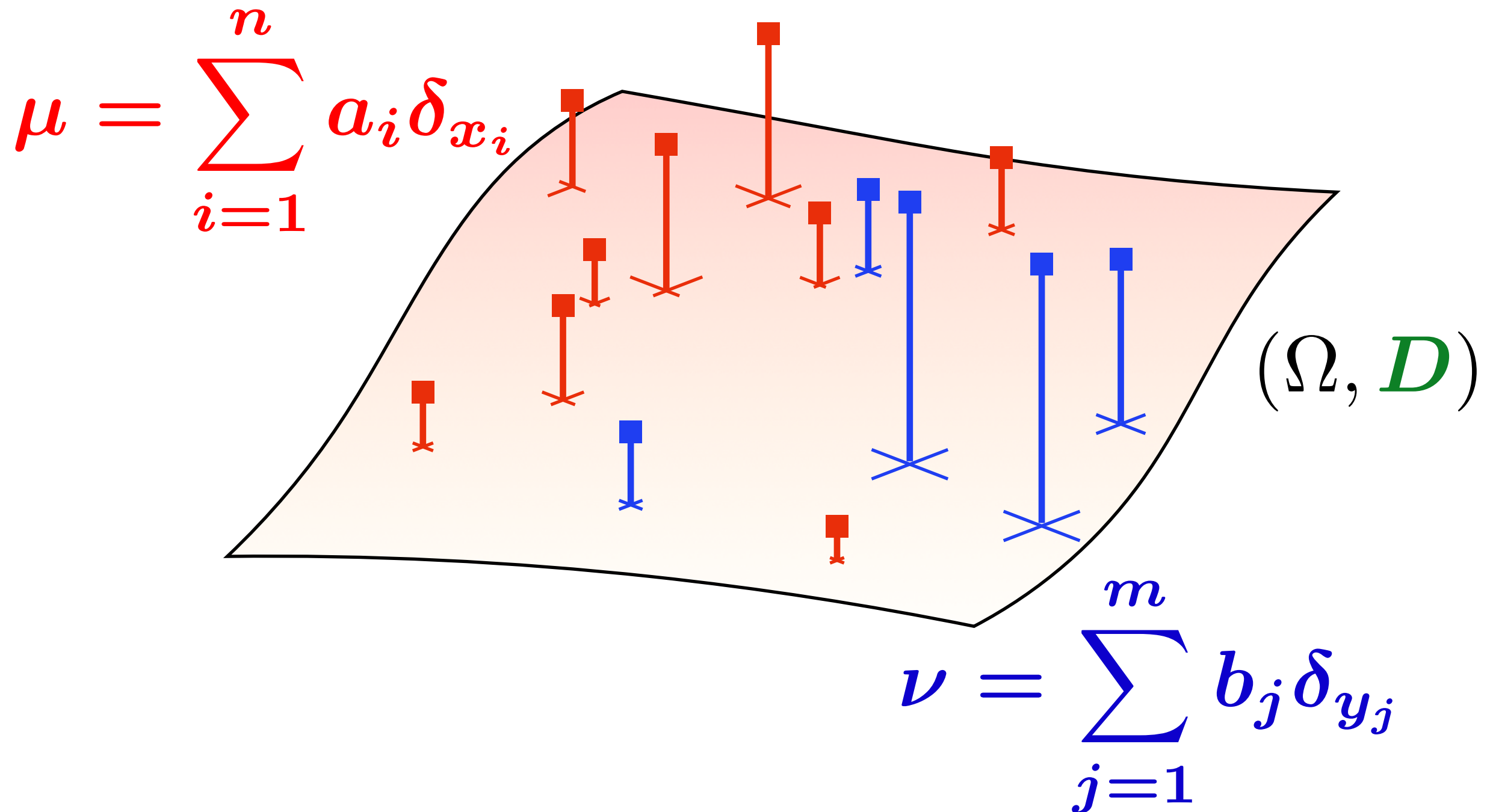
# Regularization $\rightsquigarrow$ *Differentiability*

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



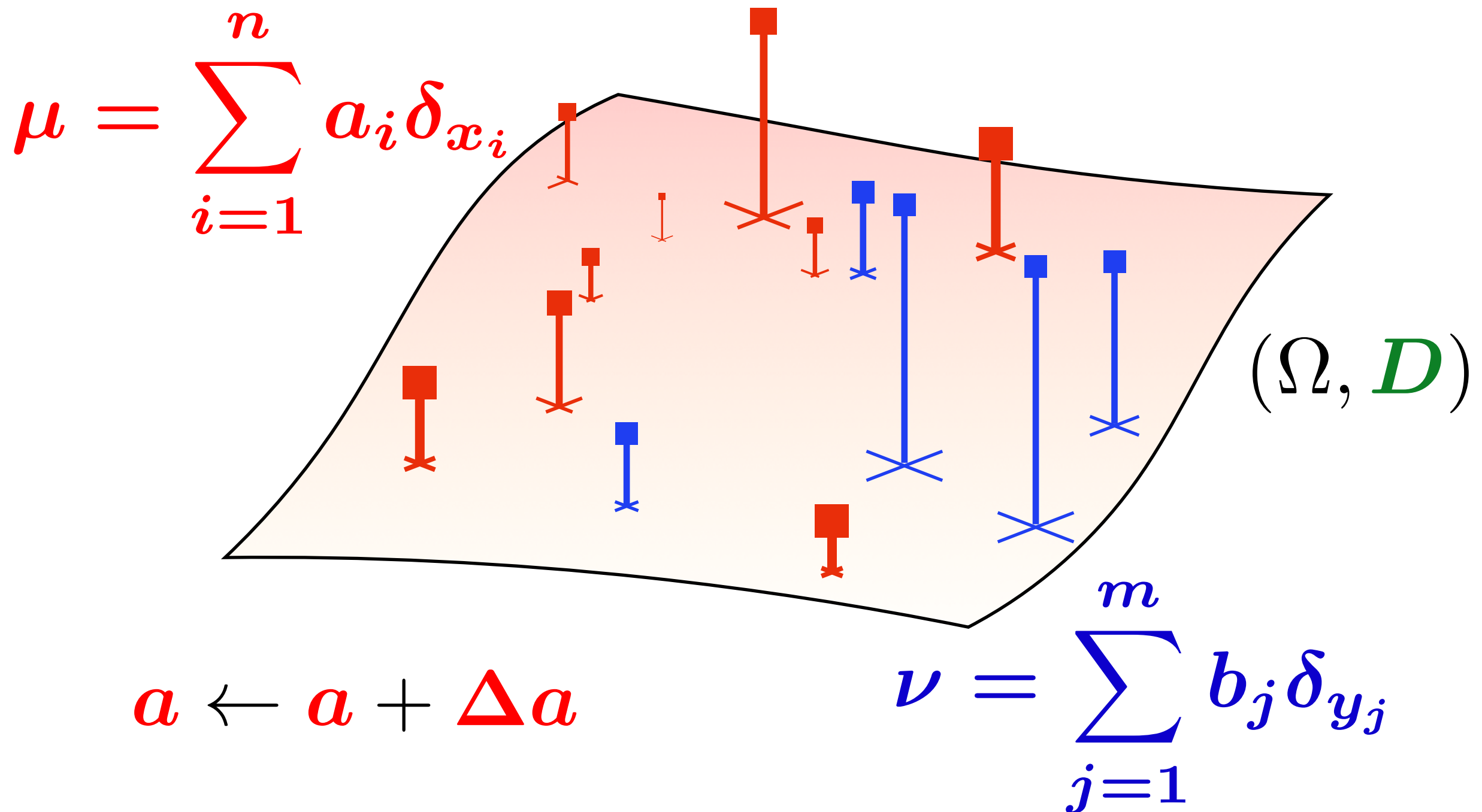
# Regularization $\rightsquigarrow$ *Differentiability*

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



# Regularization $\rightsquigarrow$ *Differentiability*

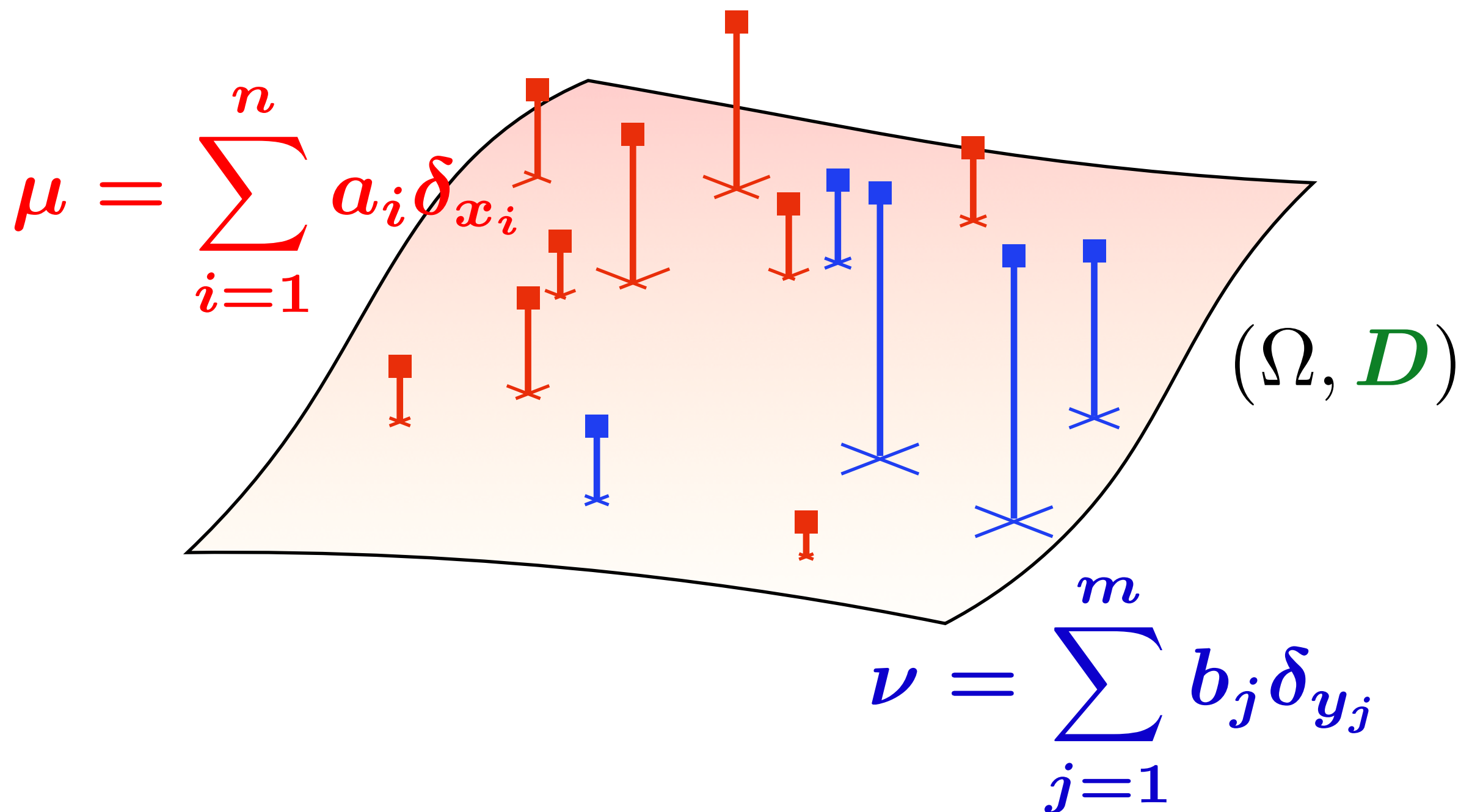
$$W_\gamma((a + \Delta a, X), (b, Y)) = W_\gamma((a, X), (b, Y)) + ??$$





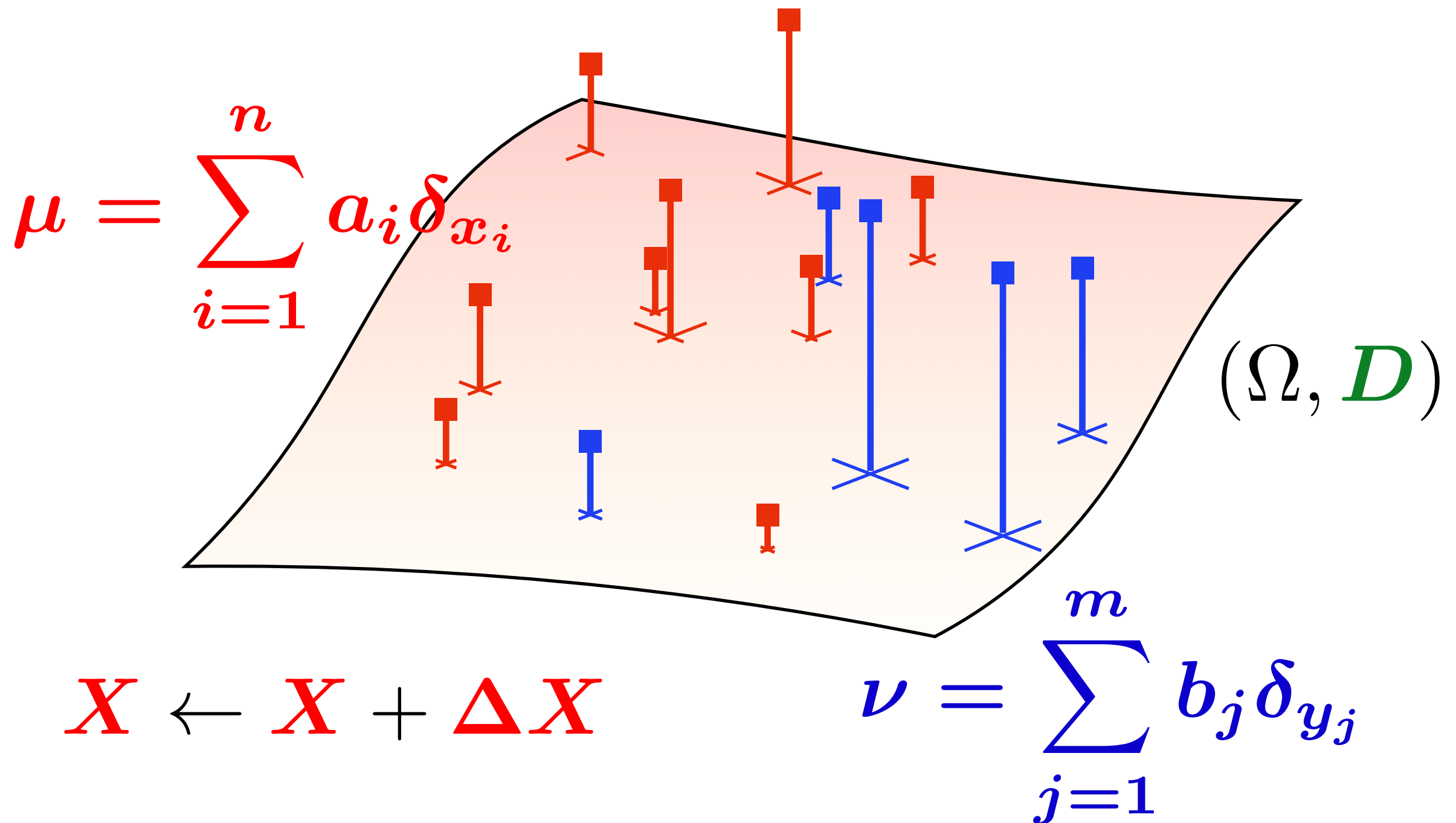
# Regularization $\rightsquigarrow$ *Differentiability*

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



# Regularization $\rightsquigarrow$ *Differentiability*

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



# 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

[CD'14]

1. convex w.r.t.  $\mathbf{a}$ ,

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

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 Discrete Reg. OT's

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

1.  $H_{\nu}$  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

$$W_p^p(\mu, \nu) = \sup_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu - \iota_C(\varphi, \psi) \quad [\text{GCPB'16}]$$

$$C = \{(\varphi, \psi) \mid \varphi \oplus \psi \leq D^p\}$$

DUAL

*regularizing dual*  *constraints*  $\gamma > 0$

$$W_\gamma(\mu, \nu) = \sup_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu - \iota_C^\gamma(\varphi, \psi)$$

$$\iota_C^\gamma(\varphi, \psi) = \gamma \iint e^{(\varphi \oplus \psi - D^p)/\gamma} d\mu d\nu$$

REGULARIZED DUAL

### 3. Stochastic Formulation

$$W_p^p(\mu, \nu) = \sup_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu - \iota_C(\varphi, \psi) \quad [\text{GCPB'16}]$$

$$C = \{(\varphi, \psi) \mid \varphi \oplus \psi \leq D^p\}$$

DUAL

regularizing dual  constraints  $\gamma > 0$

$$W_\gamma(\mu, \nu) = \sup_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu - \iota_C^\gamma(\varphi, \psi)$$

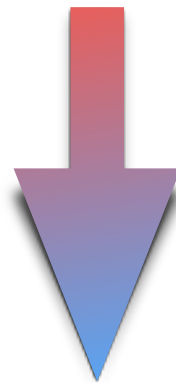
$$\iota_C^\gamma(\varphi, \psi) = \gamma \iint e^{(\varphi \oplus \psi - D^p)/\gamma} d\mu d\nu$$

REGULARIZED DUAL

# Smoothed $D$ transforms

$$W_p^p(\mu, \nu) = \sup_{\varphi} \int \varphi d\mu + \int \varphi^D d\nu.$$

SEMI-DUAL



$$\gamma > 0$$

$$W_\gamma(\mu, \nu) = \sup_{\varphi} \int \varphi d\mu + \int \varphi^{D,\gamma} d\nu.$$
$$\varphi^{D,\gamma} = -\gamma \log \int e^{\frac{\varphi(x) - D(x, \cdot)^p}{\gamma}} d\mu(x)$$

REGULARIZED SEMI-DUAL

# Regularized Semidual Wasserstein

$$W_\gamma(\mu, \nu) = \sup_{\varphi} \int \varphi d\mu + \int \varphi^{D, \gamma} d\nu.$$

$$\varphi^{D, \gamma} = -\gamma \log \int e^{\frac{\varphi(x) - D(x, \cdot)^p}{\gamma}} d\mu(x)$$

REGULARIZED SEMI-DUAL

substituting

$$\sup_{\varphi} \int_y \left[ \int_x \varphi(x) d\mu(x) - \gamma \log \int_x e^{\frac{\varphi(x) - D(x, y)^p}{\gamma}} d\mu(x) \right] d\nu(y).$$

REGULARIZED SEMI-DUAL



# Stochastic Regularized Semidual

$$\sup_{\varphi} \int_y \left[ \int_x \varphi(x) d\mu(x) - \gamma \log \int_x e^{\frac{\varphi(x) - D(x,y)^p}{\gamma}} d\mu(x) \right] d\nu(y).$$

REGULARIZED SEMI-DUAL

# Stochastic Regularized Semidual

$$\sup_{\varphi} \int_y \left[ \int_x \varphi(x) d\mu(x) - \gamma \log \int_x e^{\frac{\varphi(x) - D(x,y)^p}{\gamma}} d\mu(x) \right] d\nu(y).$$

REGULARIZED SEMI-DUAL

What if  $\mu$  is a discrete measure?

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

$\varphi \in L_1(\mu)$  is now just a vector  $\alpha \in \mathbb{R}^n$ !

# Stochastic Regularized Semidual

$$\sup_{\varphi} \int_y \left[ \int_x \varphi(x) d\mu(x) - \gamma \log \int_x e^{\frac{\varphi(x) - D(x,y)^p}{\gamma}} d\mu(x) \right] d\nu(y).$$

REGULARIZED SEMI-DUAL

What if  $\mu$  is a discrete measure?

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

$\varphi \in L_1(\mu)$  is now just a vector  $\alpha \in \mathbb{R}^n$ !

$$\sup_{\alpha \in \mathbb{R}^n} \int_y \left[ \sum_{i=1}^n \alpha_i a_i - \gamma \log \sum_{i=1}^n e^{\frac{\alpha_i - D(x_i, y)^p}{\gamma}} a_i \right] d\nu(y)$$

$$= \sup_{\alpha \in \mathbb{R}^n} \mathbb{E}_{\nu} [f(\alpha, y)]$$

STOCHASTIC REGULARIZED SEMI-DUAL

## 4. Sinkhorn Divergence

**Def.** For  $\gamma > 0$ , let  $W_\gamma(\mu, \nu) \stackrel{\text{def}}{=} \langle P_\gamma, M_{\mathbf{x}\mathbf{y}} \rangle$

**Prop.**  $W_\gamma(\mu, \mu) > 0$

**Def.** Normalized Sinkhorn Divergence

$$\bar{W}_\gamma(\mu, \nu) \stackrel{\text{def}}{=} W_\gamma(\mu, \nu) - \frac{1}{2} (W_\gamma(\mu, \mu) + W_\gamma(\nu, \nu))$$

**Prop.** If  $p = 1$ ,  $\bar{W}_\gamma(\mu, \nu) \xrightarrow{\gamma \rightarrow \infty} \text{ED}(\mu, \nu)$

# 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

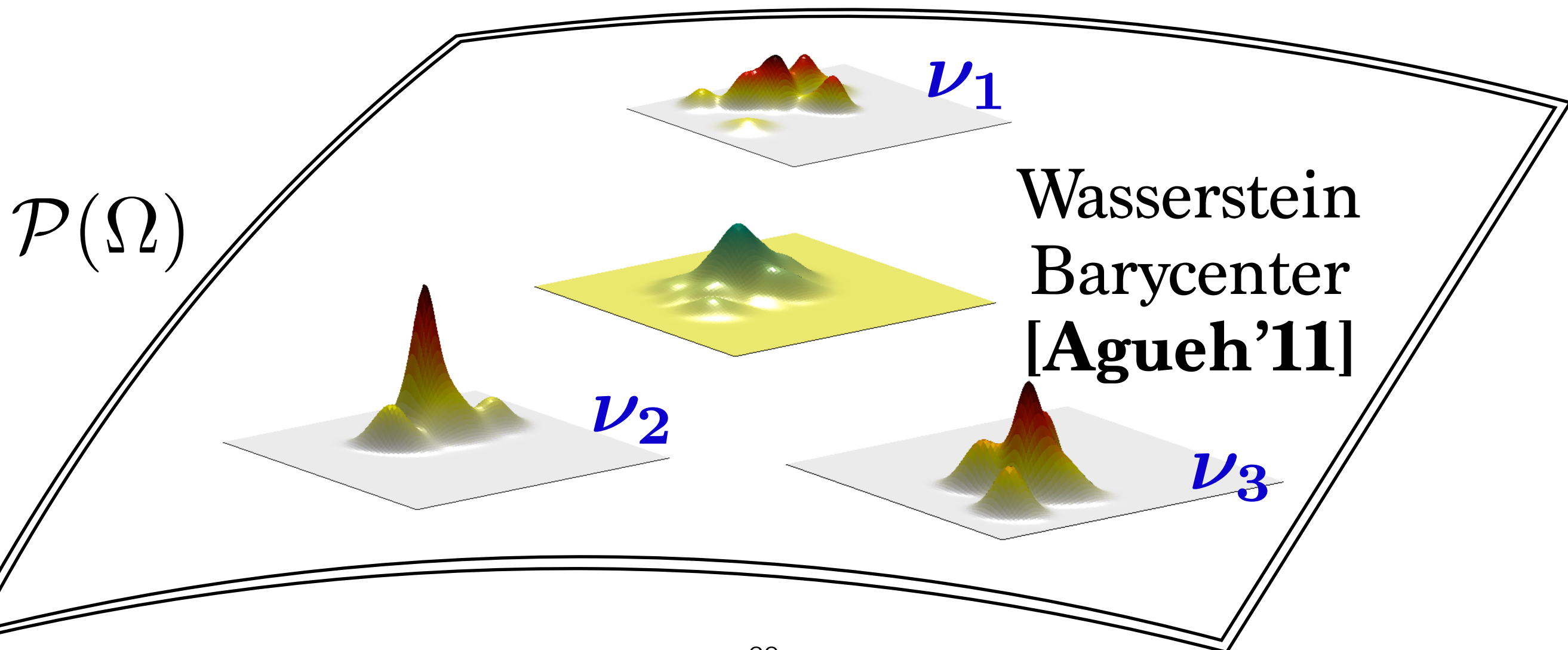
**Example:** Differentiability w.r.t.  $a$

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

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

# Wasserstein Barycenters

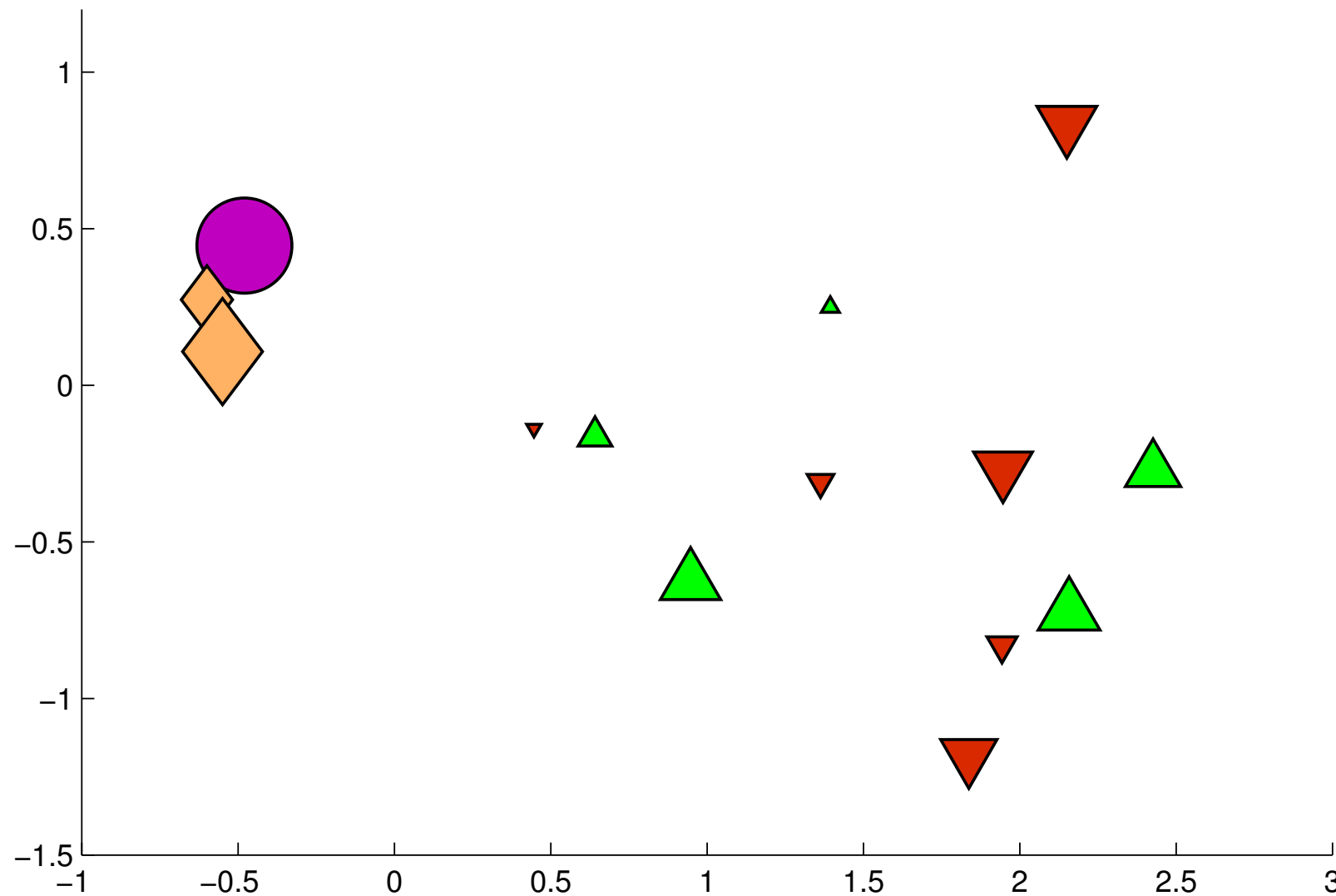
$$\min_{\mu \in \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_p^p(\mu, \nu_i)$$





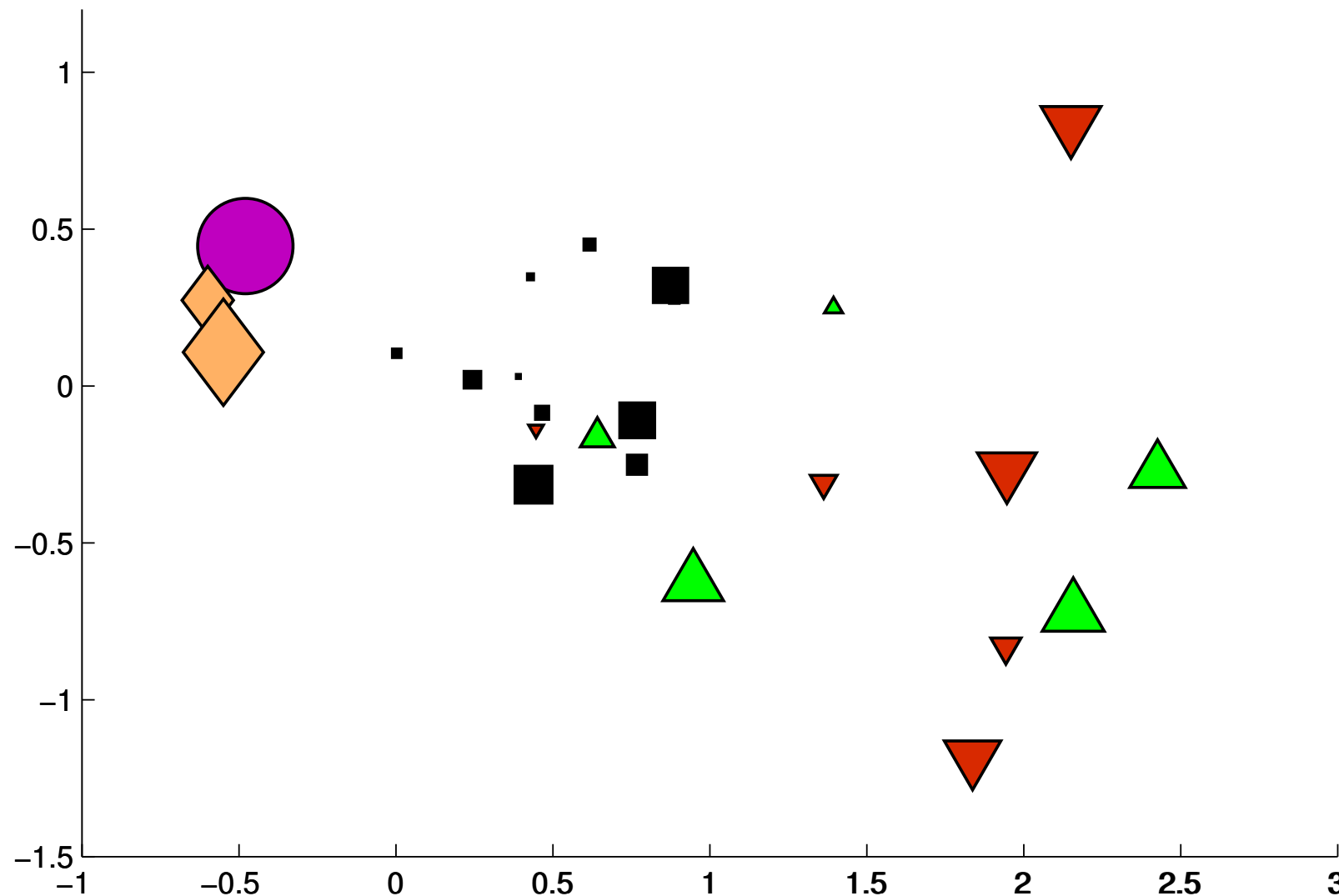
# Multimarginal Formulation

- 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  $(\prod_i n_i, \sum_i n_i)$

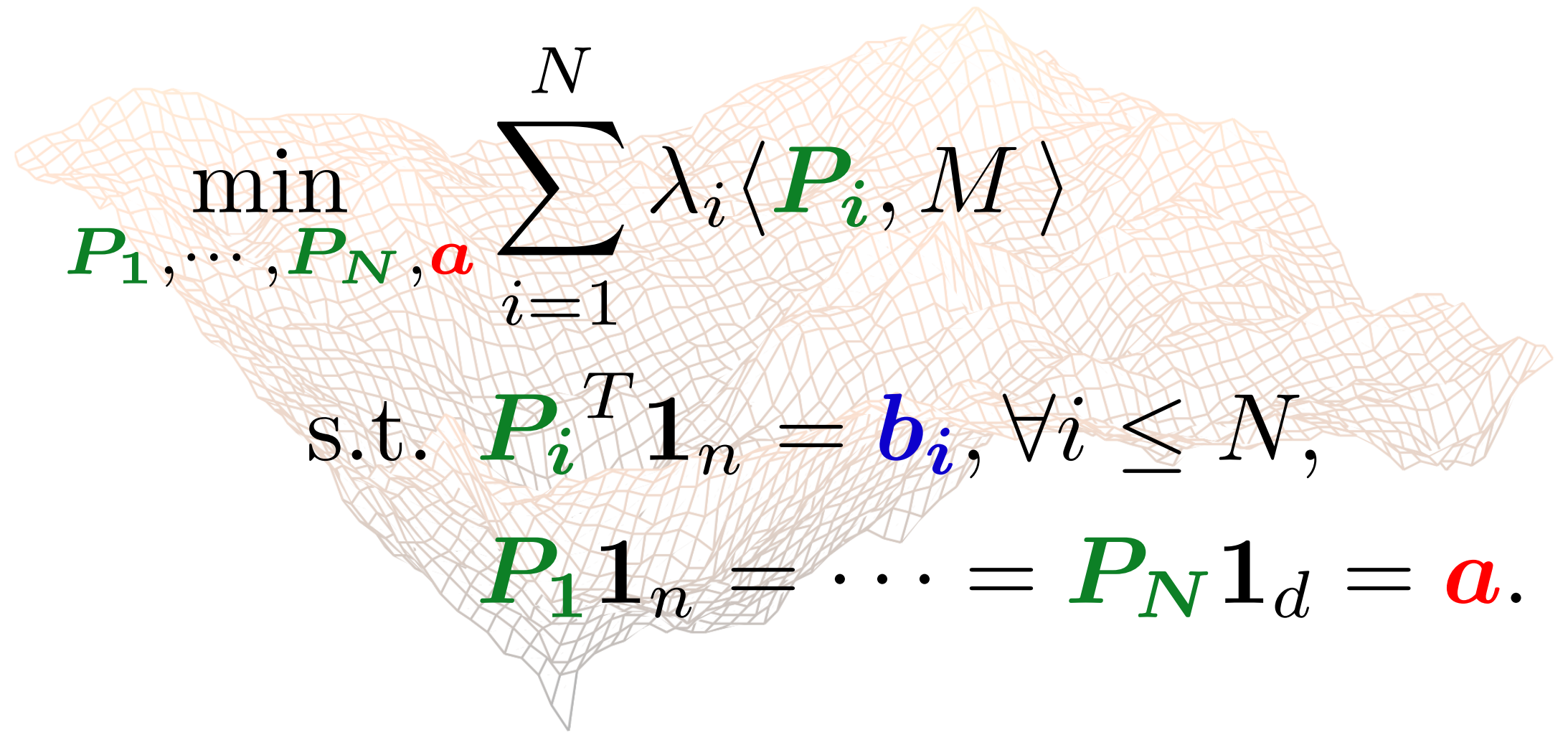
# Finite Case, LP Formulation

- When  $\Omega$  is a **finite set**, metric  $M$ , another LP.


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

# Finite Case, LP Formulation

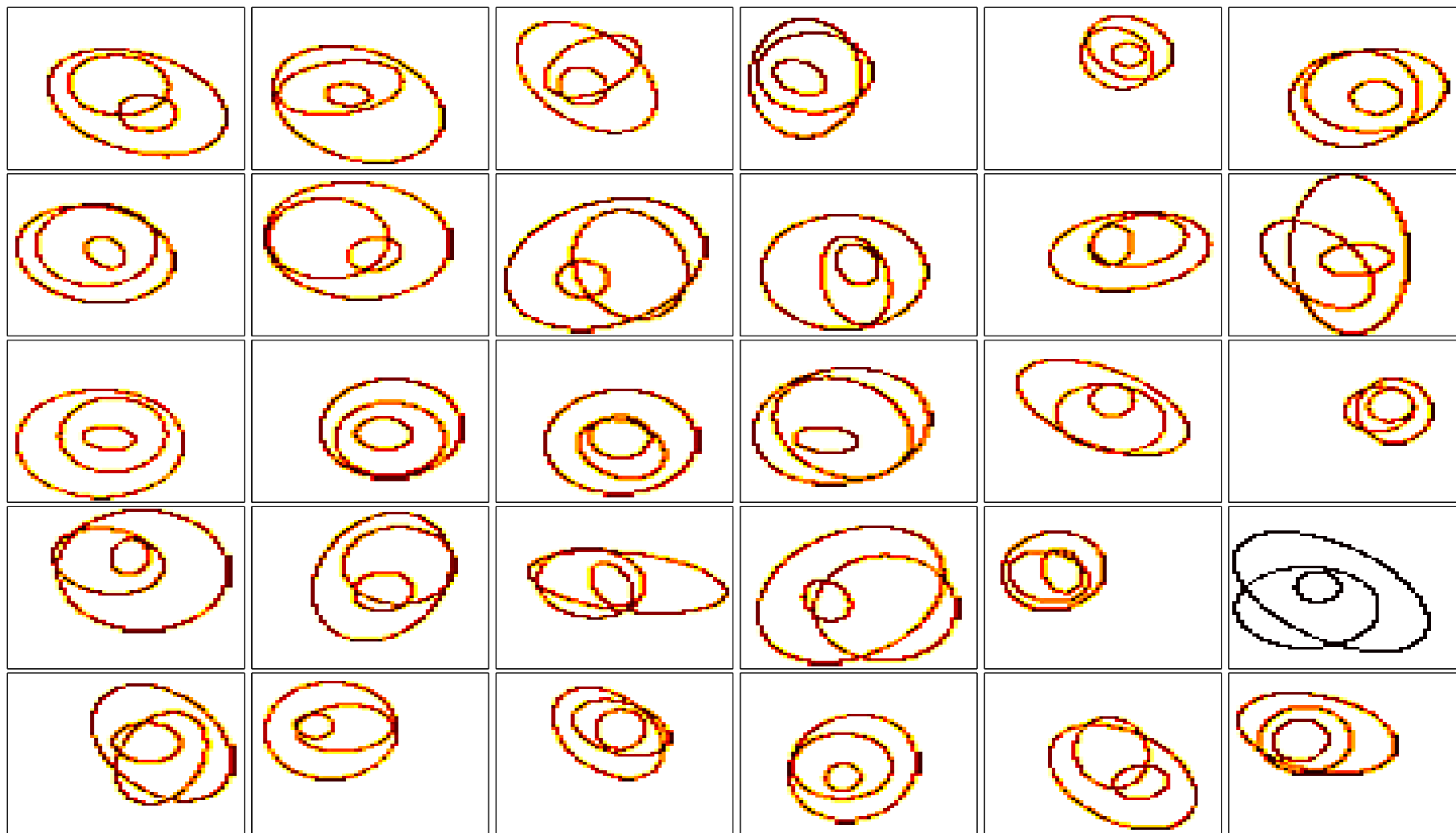
- When  $\Omega$  is a **finite set**, metric  $M$ , another LP.


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

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

# Primal Descent on Regularized $W$

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

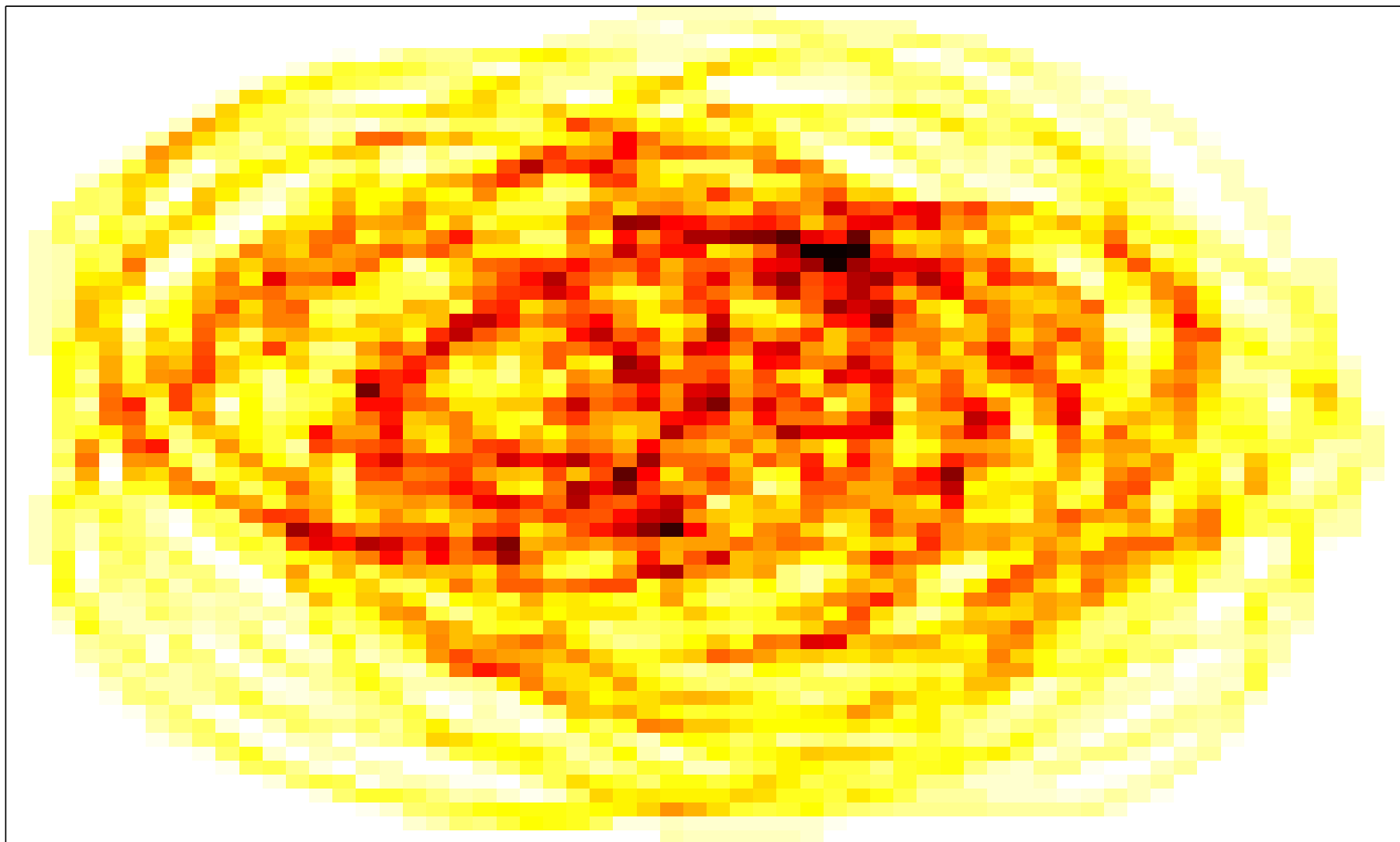


*Fast Computation of Wasserstein Barycenters*  
**International Conference on Machine Learning 2014**

**[CD'14]**

# Primal Descent on Regularized $W$

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

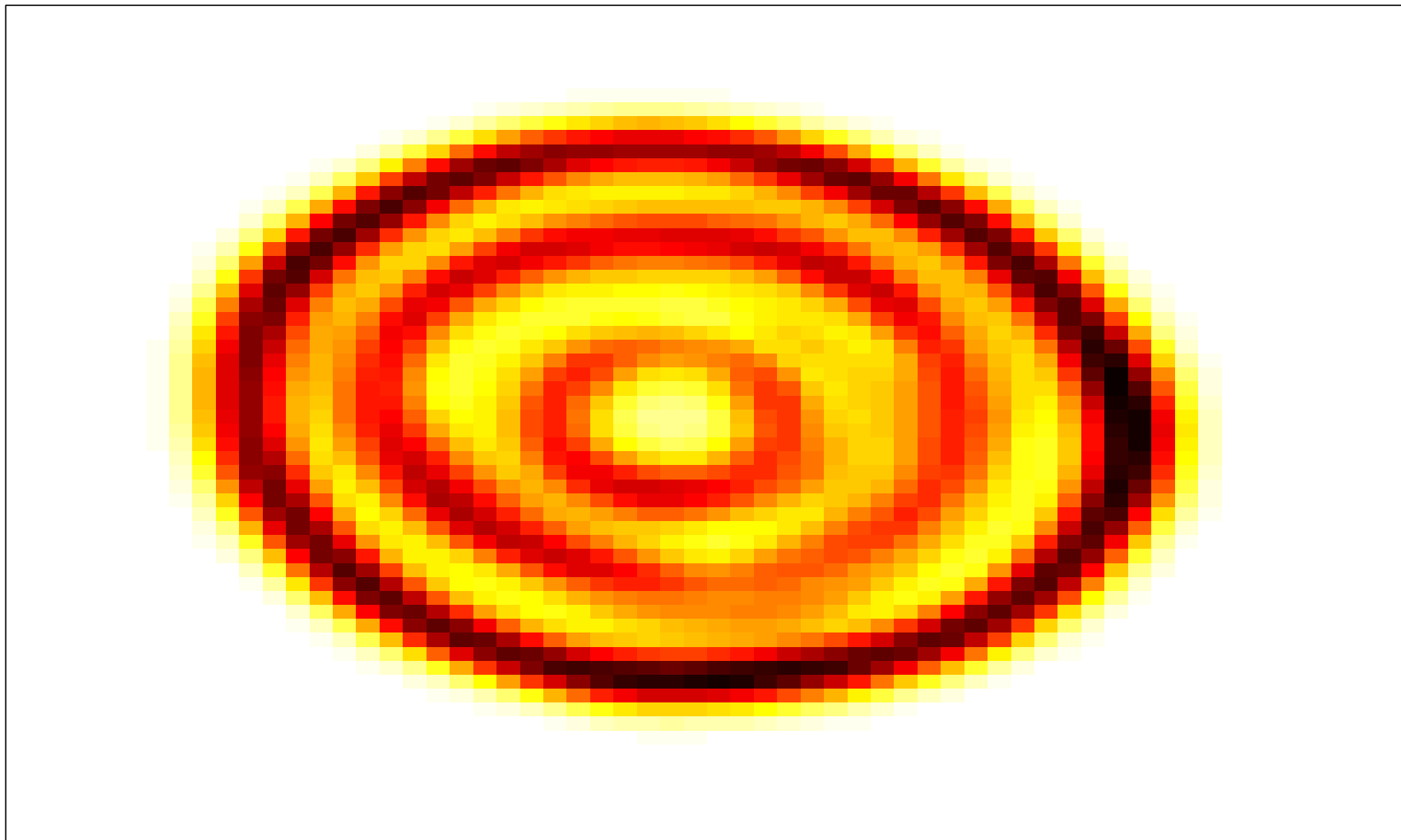


*Fast Computation of Wasserstein Barycenters*  
**International Conference on Machine Learning 2014**

**[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 $W$

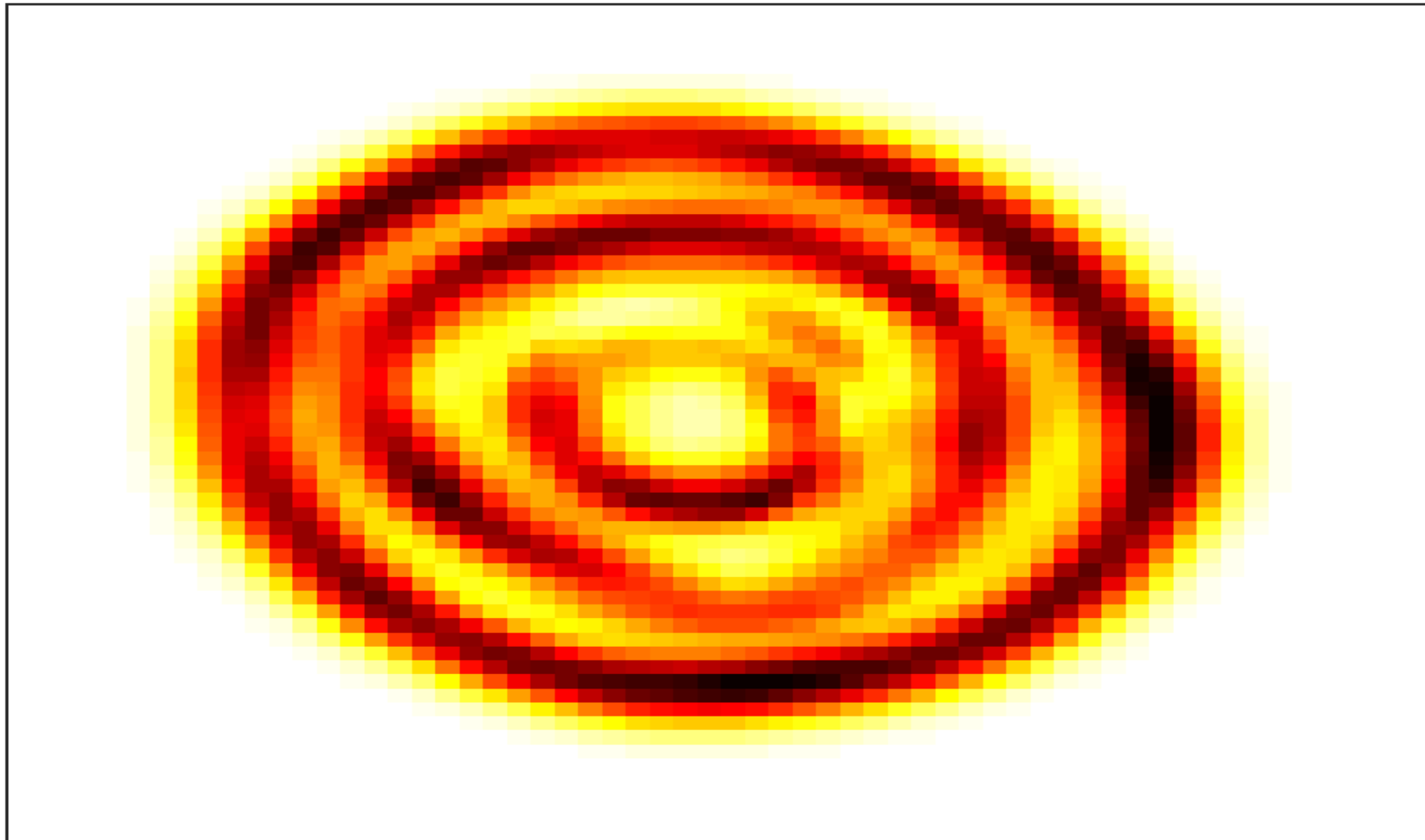
---

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



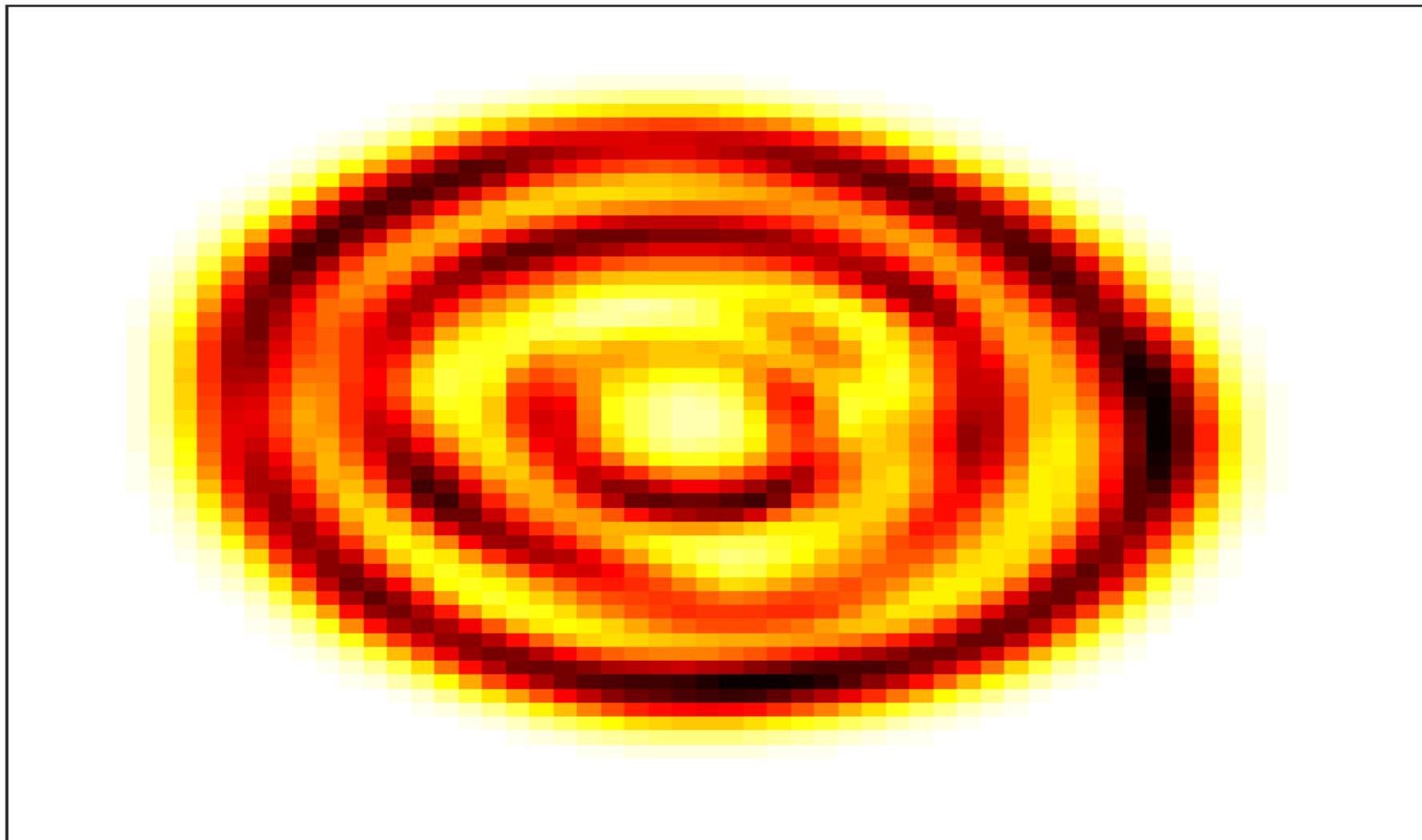
# Primal Descent on Algorithmic $W$

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



# Primal Descent on Algorithmic $W$

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



**not a convex problem**

# Inverse Wasserstein Problems

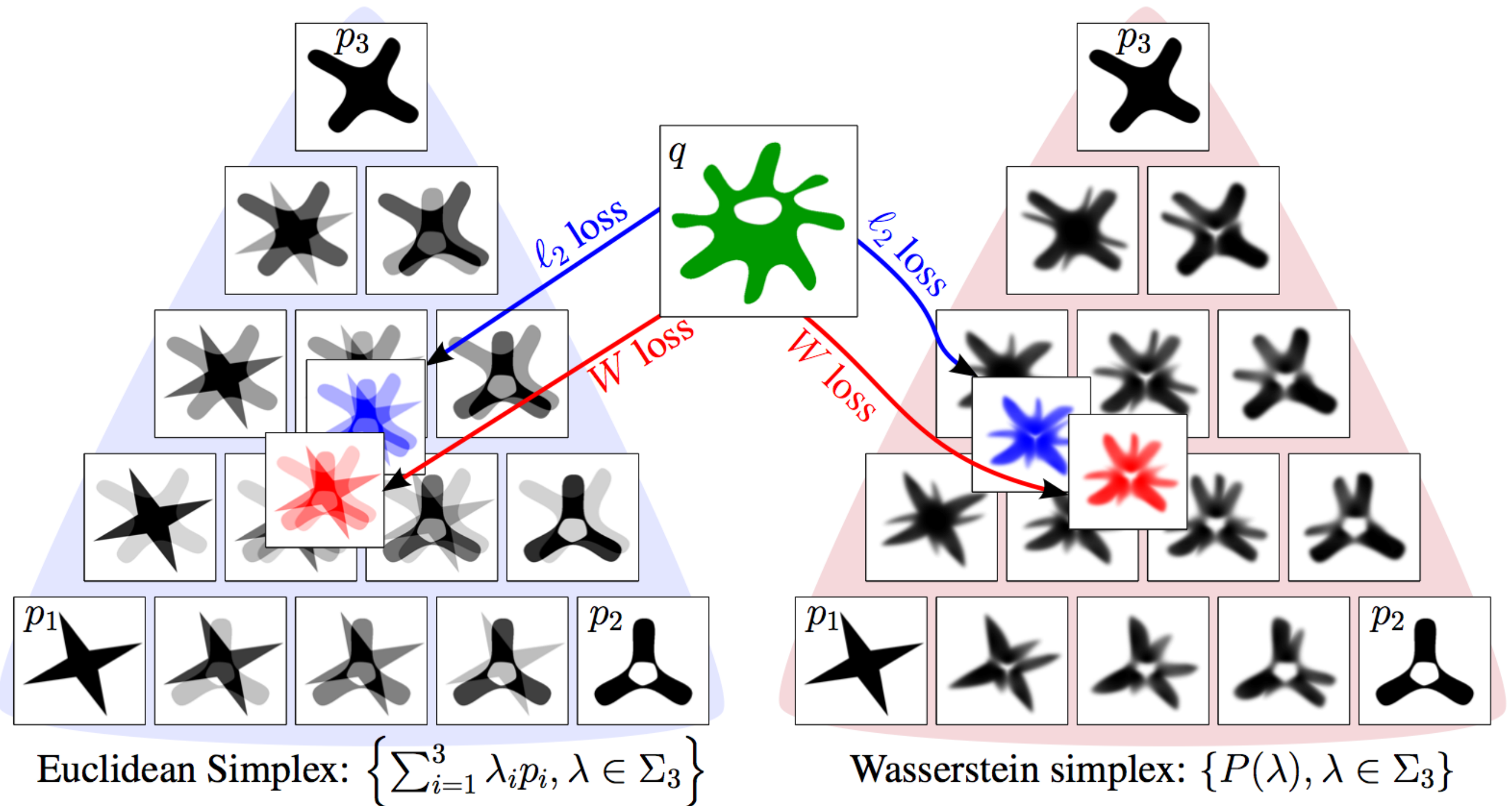
- consider Barycenter operator:

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

- 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**.

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

# 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



$$\lambda_0 = 0.03$$



$$\lambda_1 = 0.12$$



$$\lambda_2 = 0.40$$



$$\lambda_3 = 0.43$$

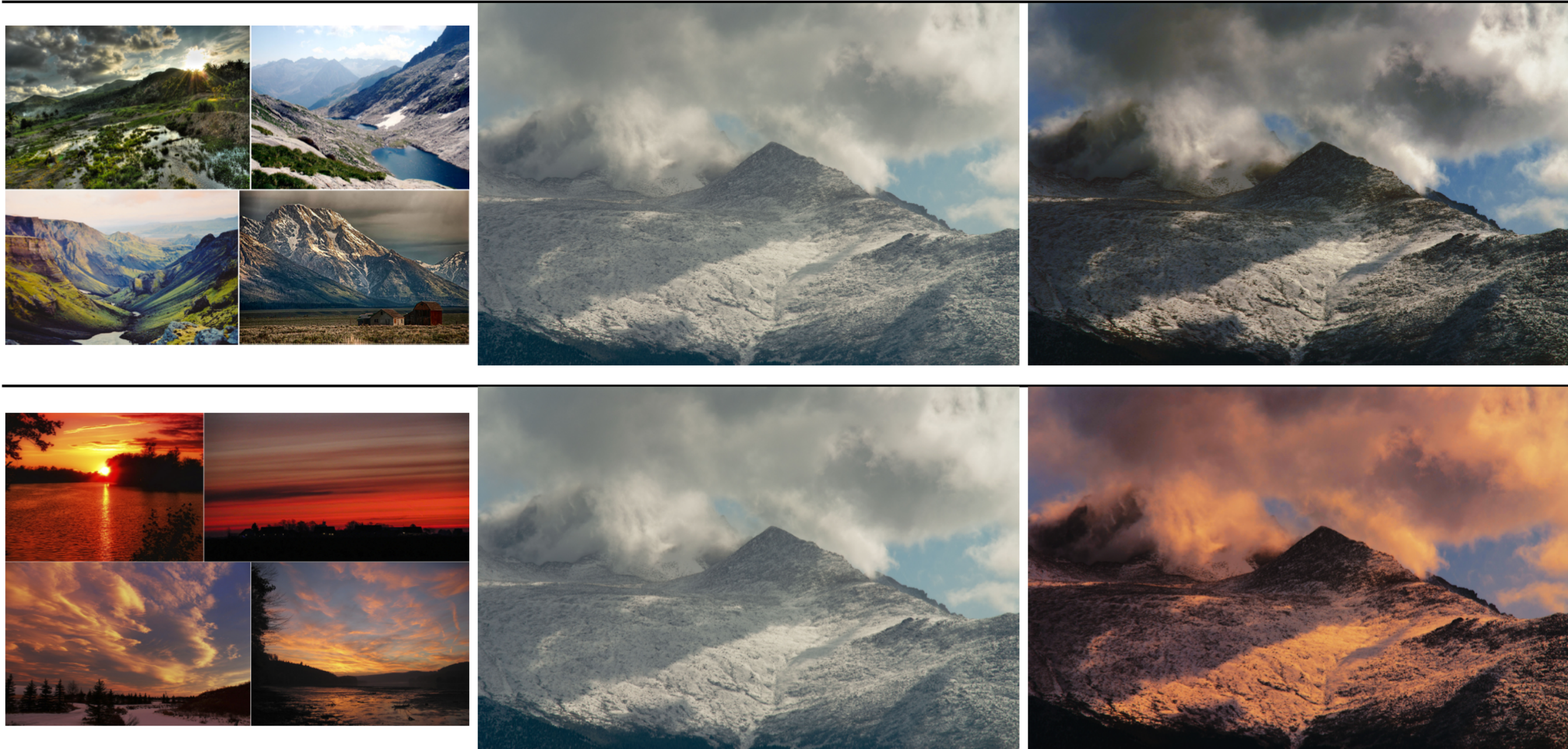


# Application: Color Grading





# Application: Color Grading

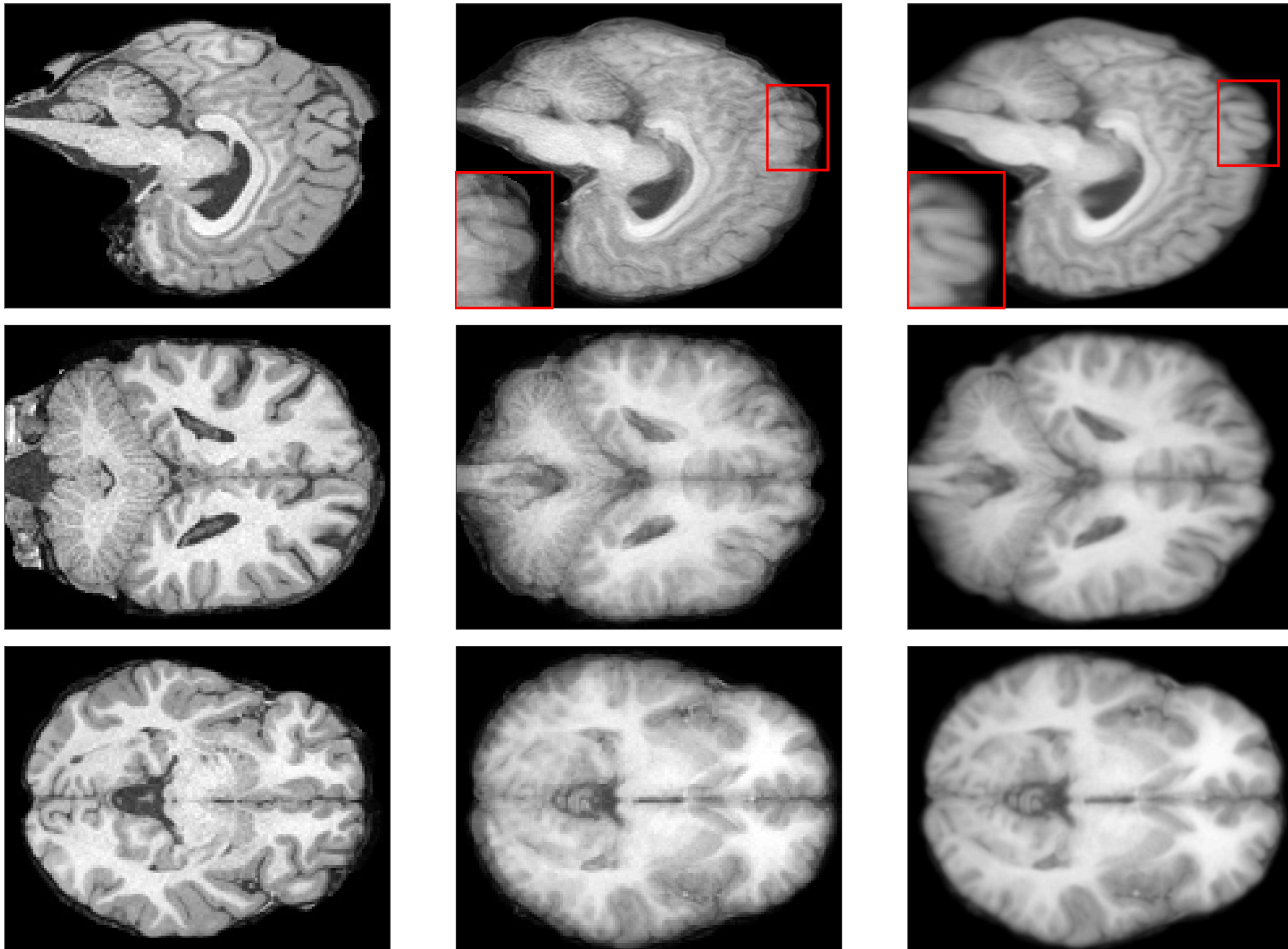


*Wasserstein Barycentric Coordinates: Histogram  
Regression using Optimal Transport, SIGGRAPH'16*

**[BPC'16]**



# Application: Brain Mapping

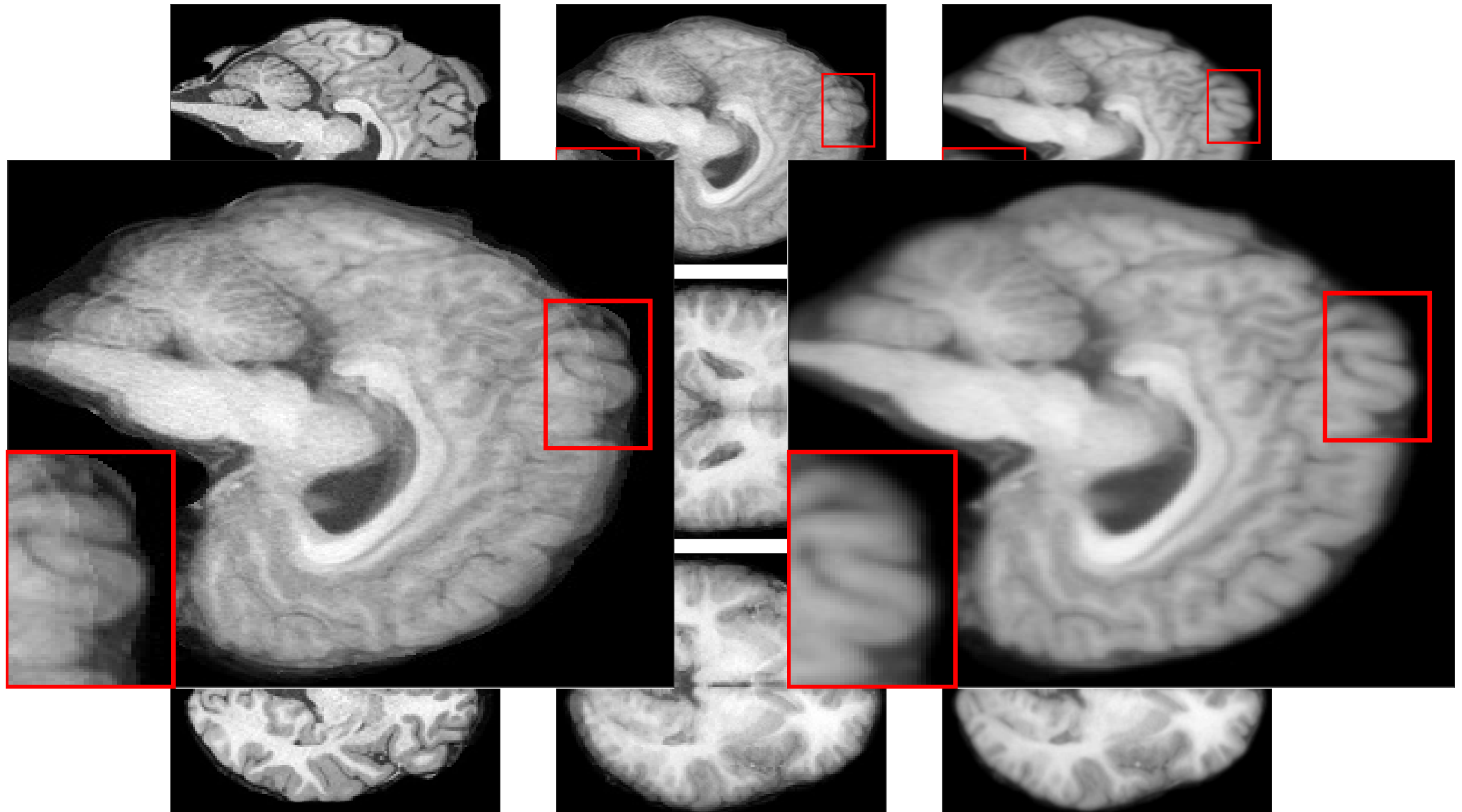


Original

Euclidean  
projection

Wasserstein  
projection

# Application: Brain Mapping



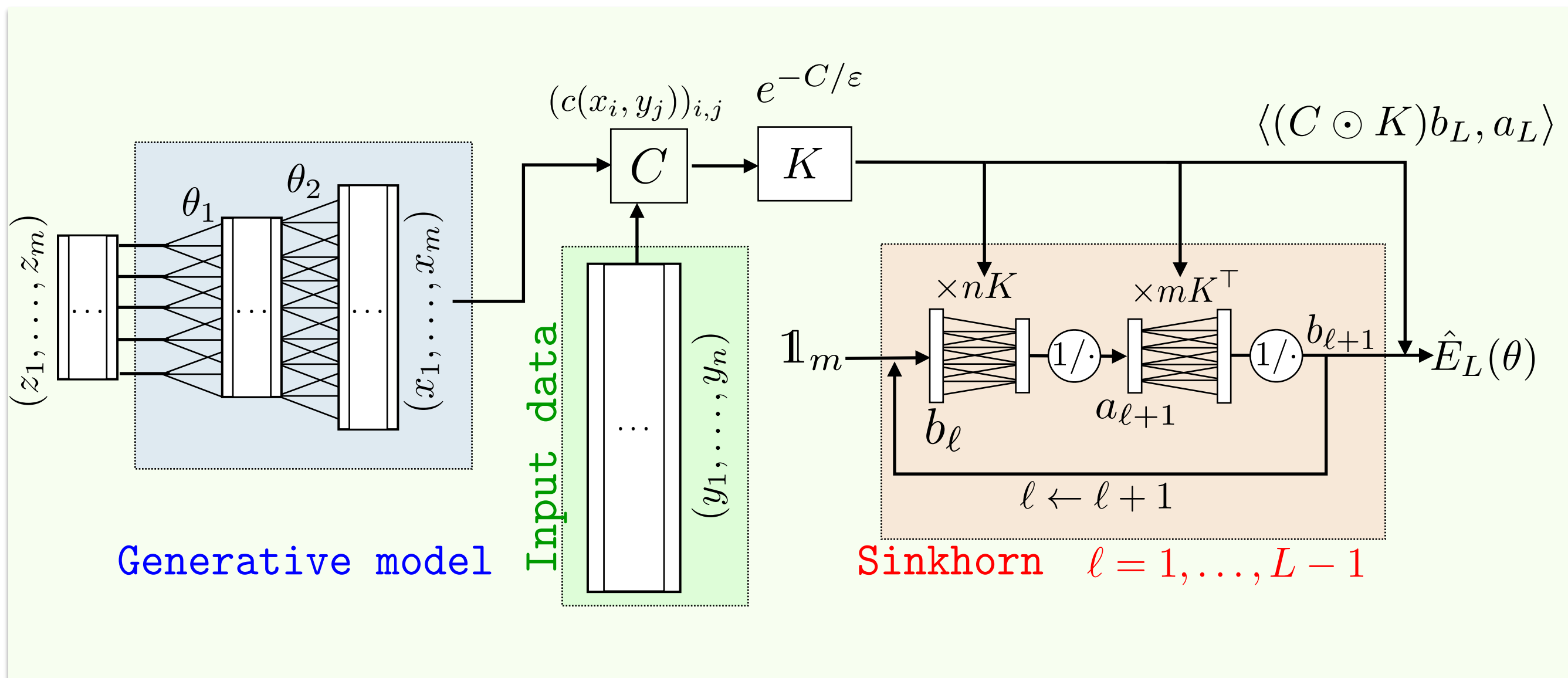
Original

Euclidean  
projection

Wasserstein  
projection

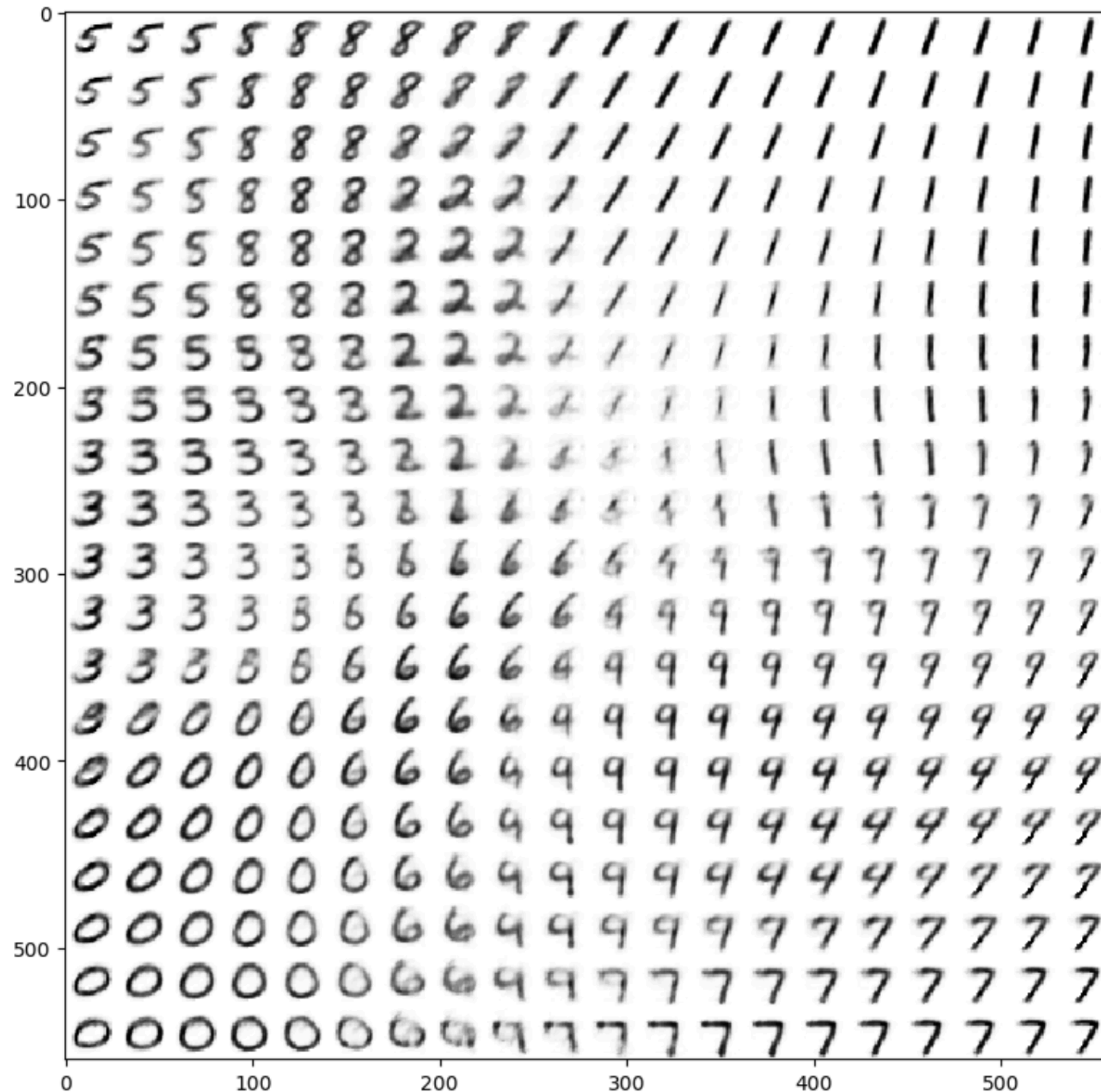
# At Last: Application to Generative Models

Approximate  $W$  loss by the transport cost  $\bar{W}_L$  after  $L$  Sinkhorn iterations.



[GPC'17]

# Example: MNIST, Learning $f_{\theta}$

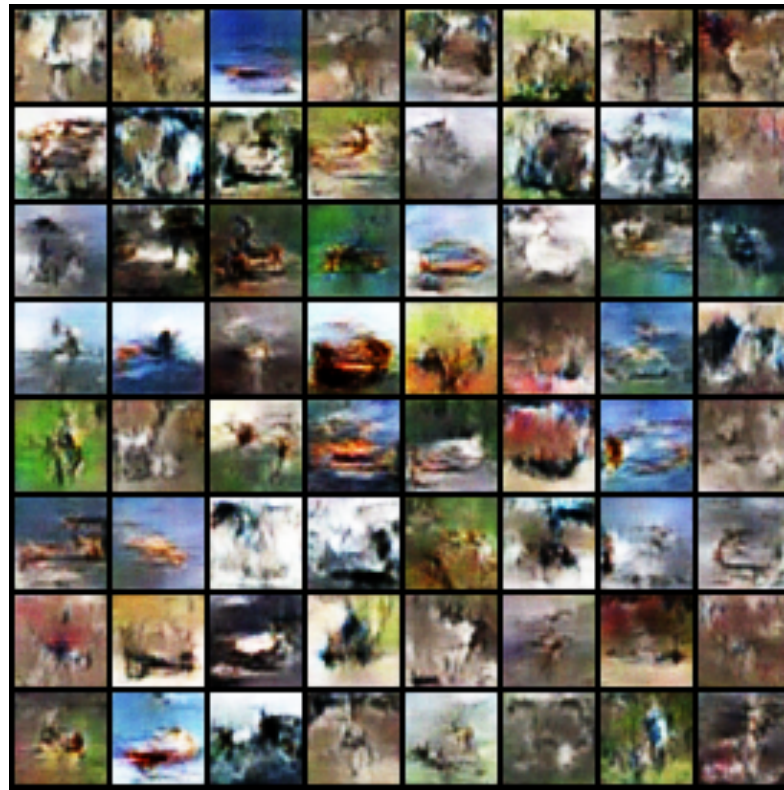




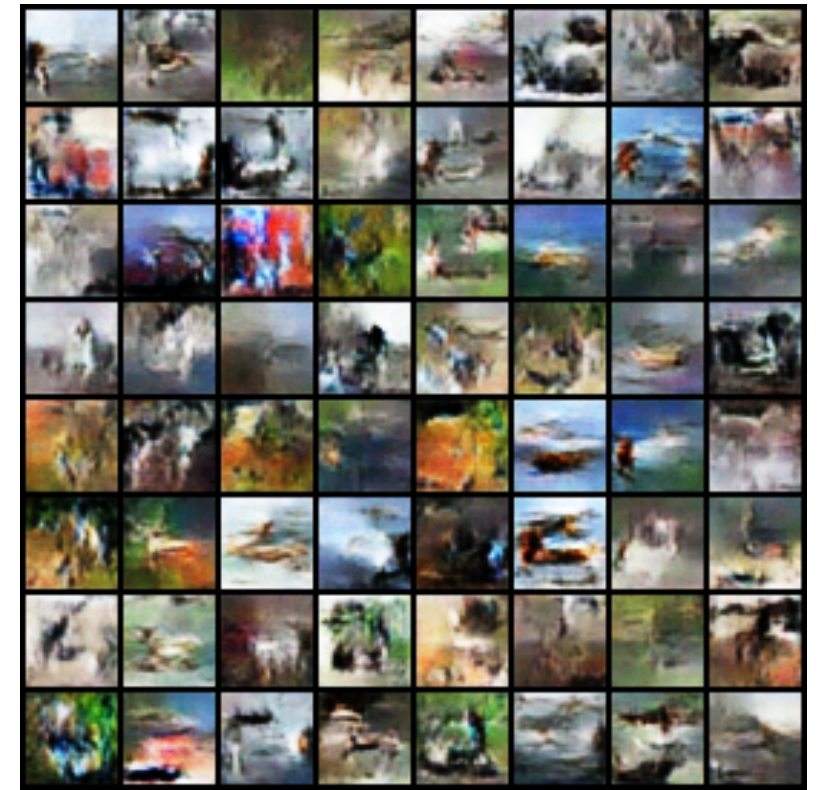
# Example: Generation of Images



MMD-GAN



gamma = 1000



gamma=10

- CIFAR-10 images
- In these examples the cost function is also learned adversarially, as a NN mapping onto feature vectors.

# Concluding Remarks

- *Regularized OT* is much faster than OT.
- *Regularized OT* can interpolate between  $W$  and the *MMD / Energy distance* metrics.
- The solution of *regularized OT* is “*auto-differentiable*”.
- **Many open problems remain!**

Sat Dec 9th 08:00 AM -- 06:30 PM @ None

***Optimal Transport and Machine Learning***

Olivier Bousquet · Marco Cuturi · Gabriel Peyré · Fei Sha · Justin Solomon

NIPS'17 WORKSHOP

Workshop

Dates n/a. @ TBA

***A Primer on Optimal Transport***

Marco Cuturi · Justin M Solomon

NIPS'17 TUTORIAL

Tutorial