

A Primer on
Optimal Transport

Marco Cuturi



École nationale
de la statistique
et de l'administration
économique

université
PARIS-SACLAY

Justin Solomon



What is Optimal Transport?

The natural geometry for **probability measures**



Monge



Kantorovich



Koopmans



Dantzig



Brenier



Otto



McCann



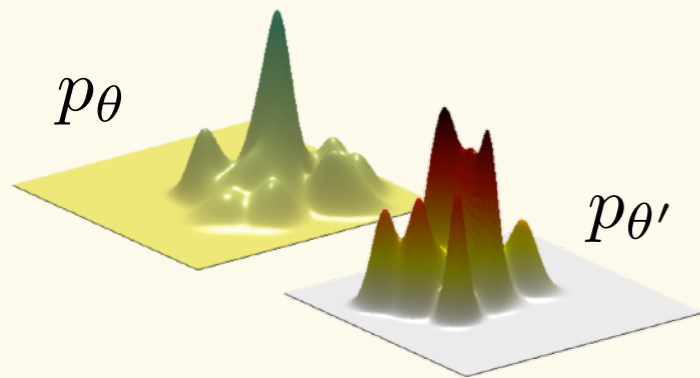
Villani

Nobel '75

Fields '10

What is Optimal Transport?

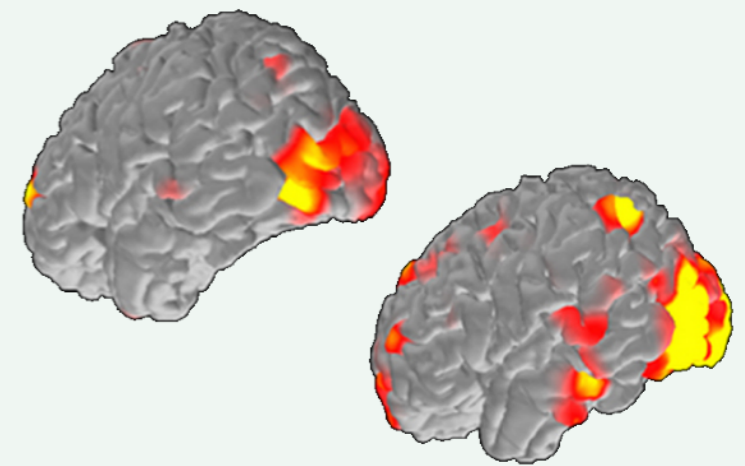
The natural geometry for **probability measures**



Statistical Models

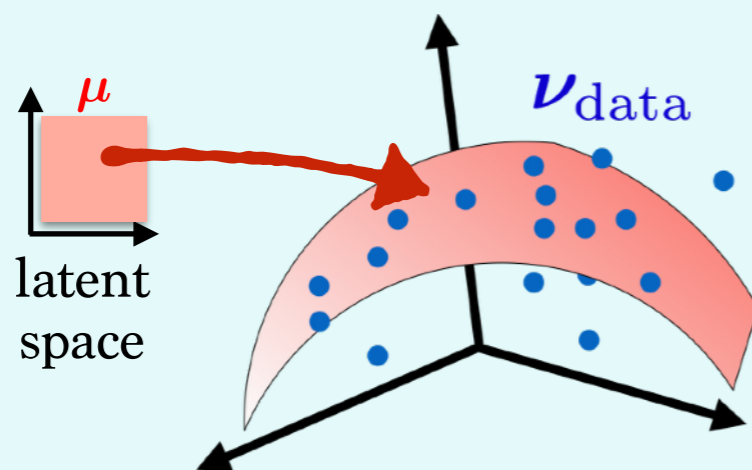


Bags of features



Brain Activation Maps

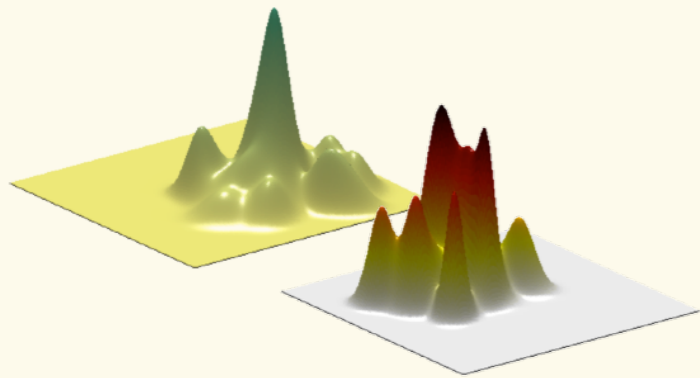
Generative Models vs. data



Color Histograms

What is Optimal Transport?

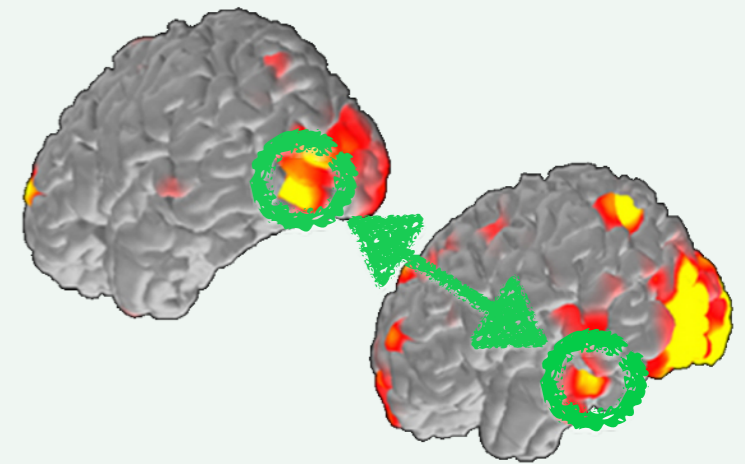
The natural geometry for **probability measures** supported on a geometric space.



Statistical Models

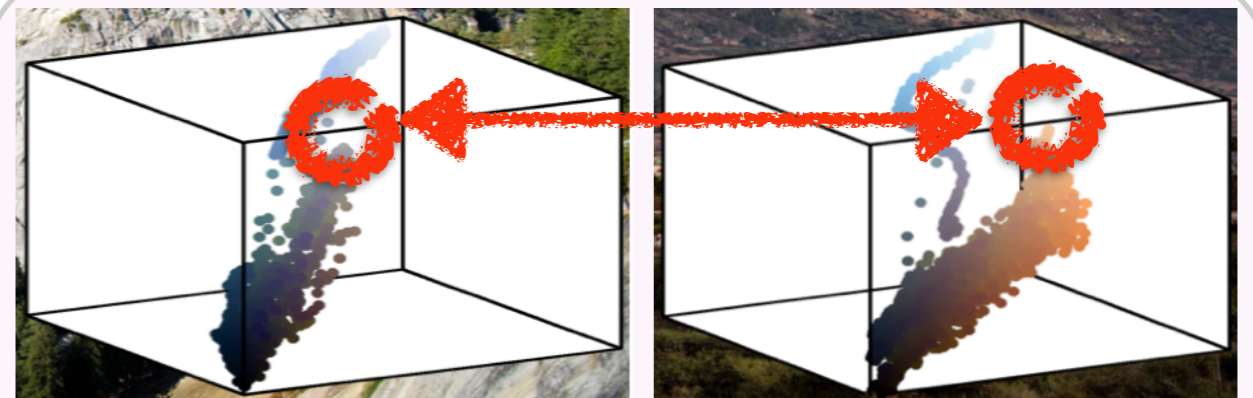
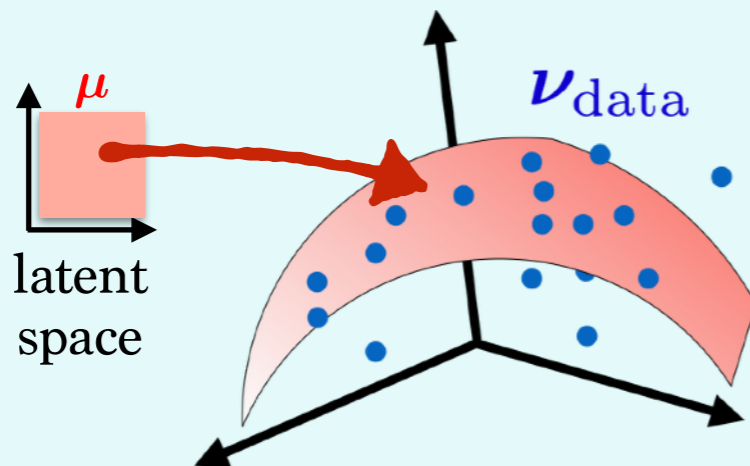


Bags of Features



Brain Activation Maps

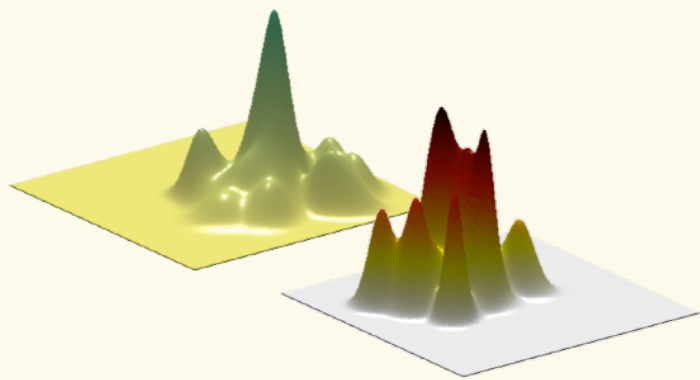
Generative Models vs. Data



Color Histograms

What is Optimal Transport?

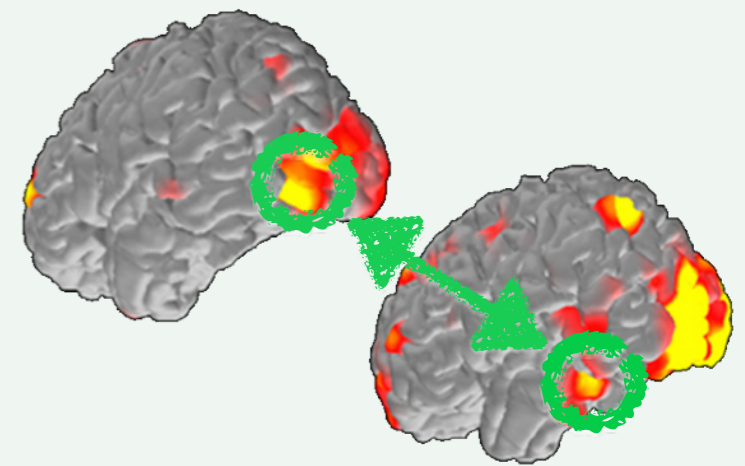
The natural geometry for **probability measures** supported on a geometric space.



Statistical Models

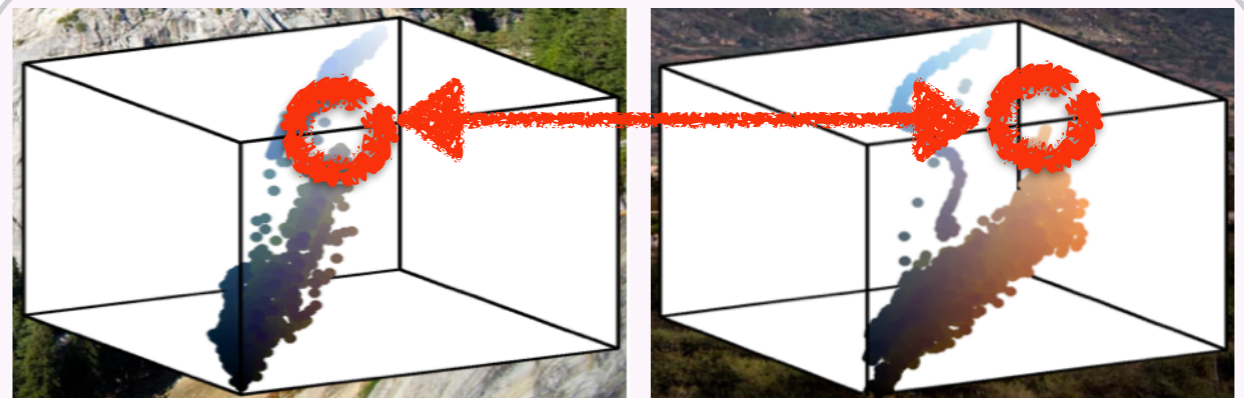
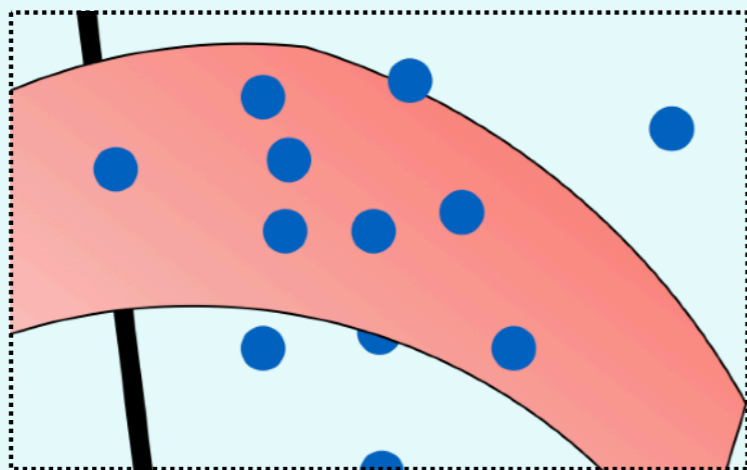


Bags of Features



Brain Activation Maps

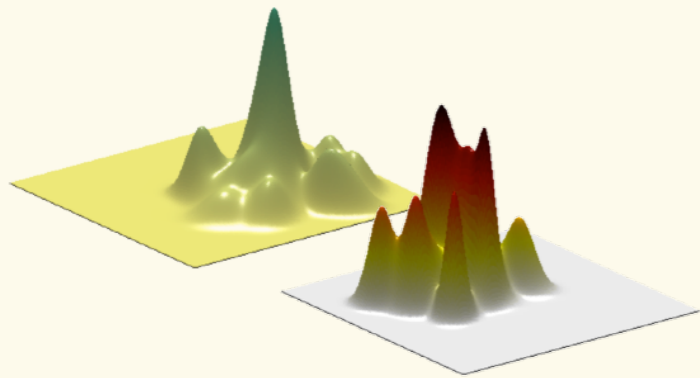
Generative Models vs. Data



Color Histograms

What is Optimal Transport?

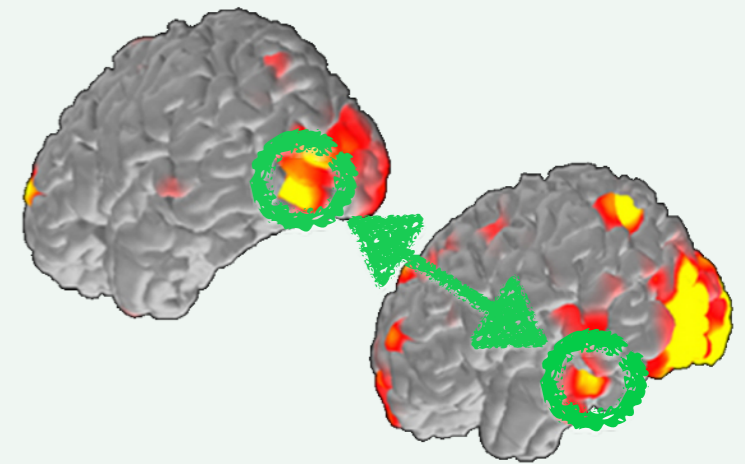
The natural geometry for **probability measures** supported on a geometric space.



Statistical Models

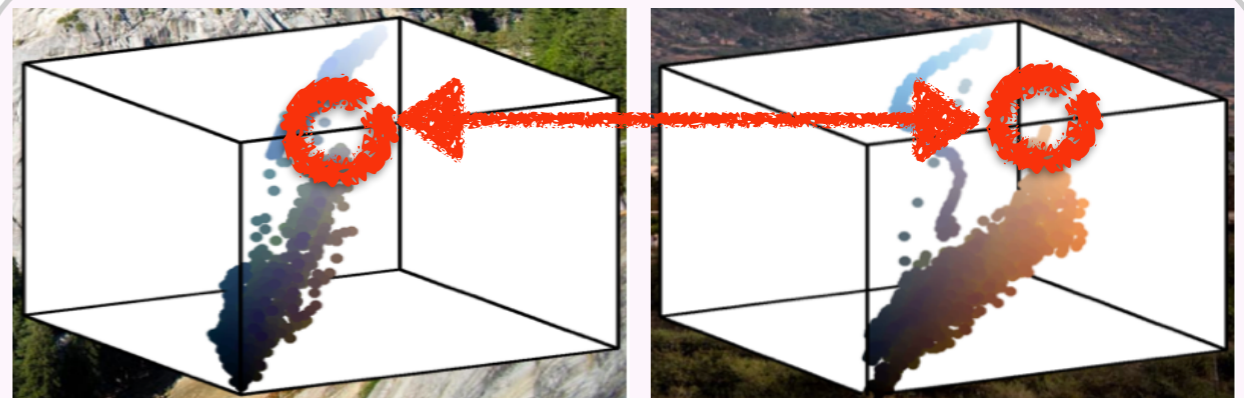
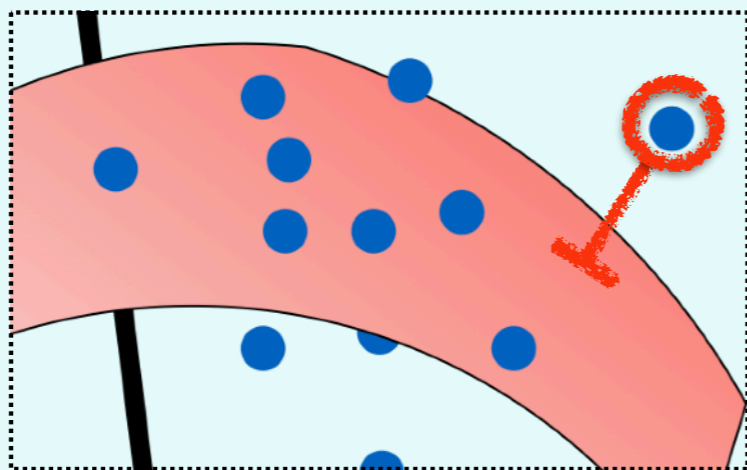


Bags of Features



Brain Activation Maps

Generative Models vs. Data



Color Histograms

Today's Outline

1. Introduction to optimal transport
2. Optimal transport algorithms
3. Applications (W as a loss)
4. Applications (W for estimation)

More...

This Saturday: OT & ML Workshop

7 Talks by: Jacob, Kraig, Andoni, Gangbo, Bottou, Flamary, Bach

17 posters and spotlight presentations.

Organizers: Bousquet, Cuturi, Peyré, Sha, Solomon

Survey and slides: <https://optimaltransport.github.io/>

Introduction to OT

- Two examples: moving earth & soldiers
- Monge problem, Kantorovich problem
- OT as geometry, OT as a loss function

Origins: Monge Problem (1781)



60 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.

LORSQU'ON doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

Origins: Monge Problem (1781)



60 MÉMOIRES DE L'ACADÉMIE ROYALE

M É M O I R E

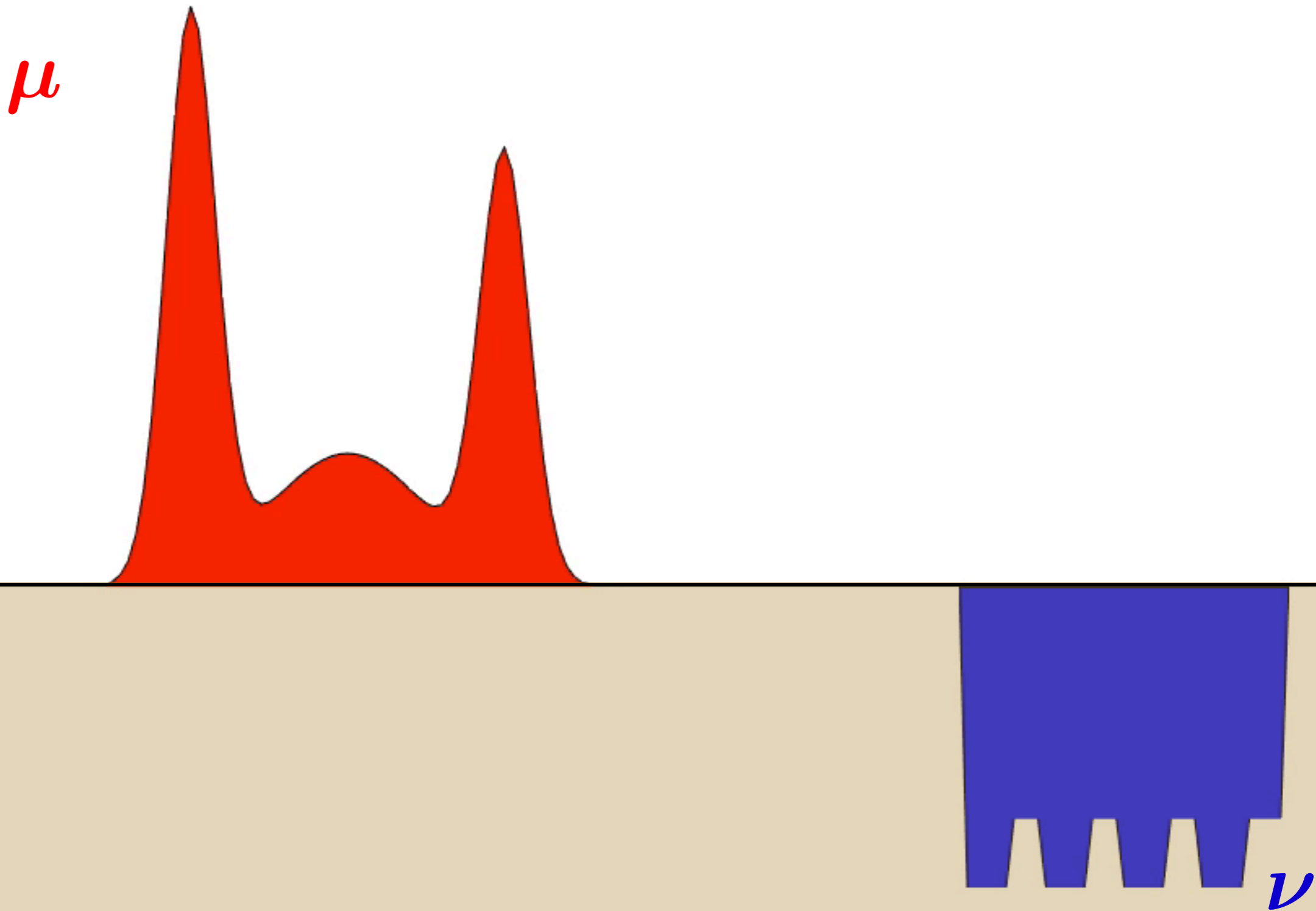
S U R L A

T H É O R I E D E S D É B L A I S

*When one has to bring earth
from one place to another...*

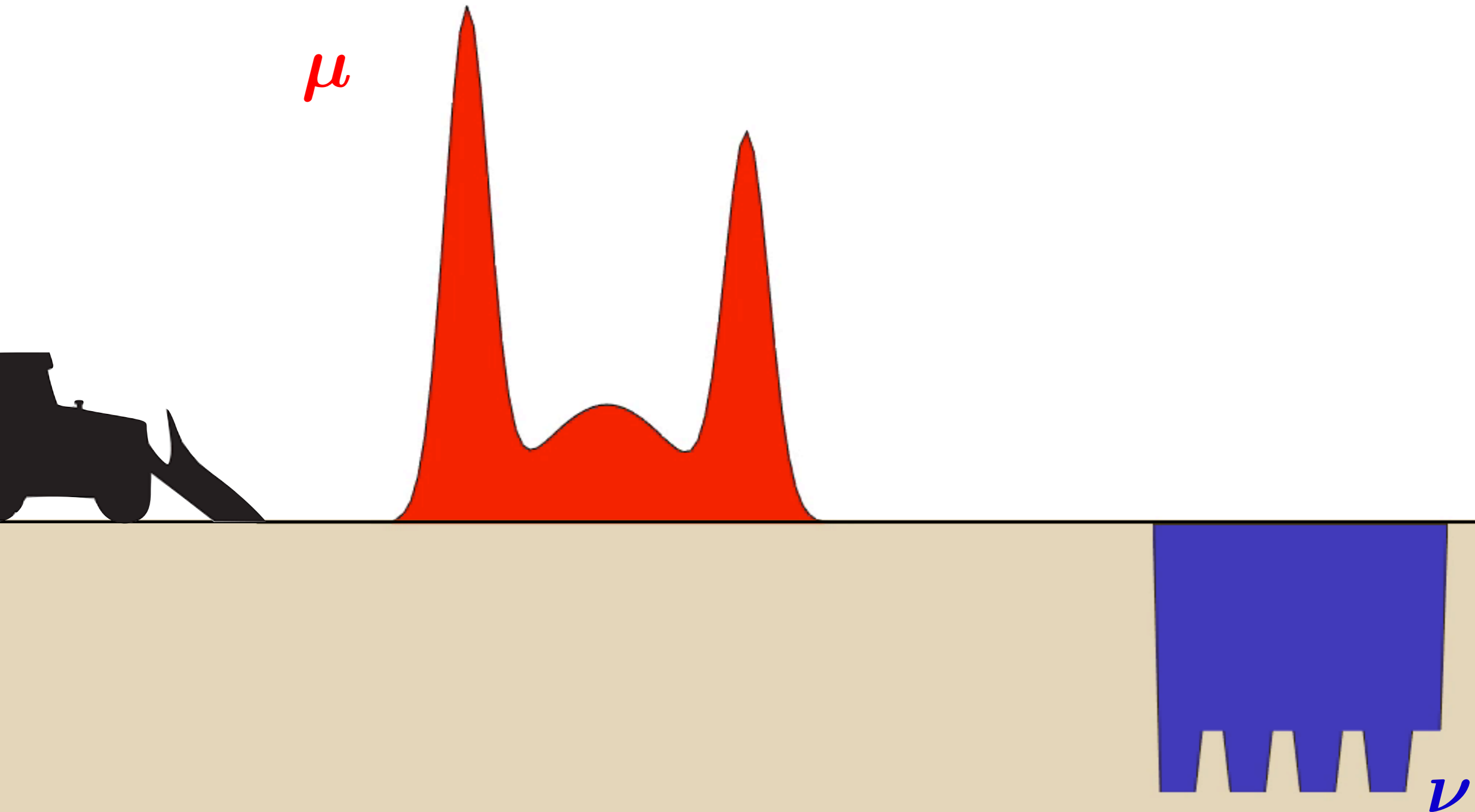
LORSQU'ON doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

Origins: Monge Problem



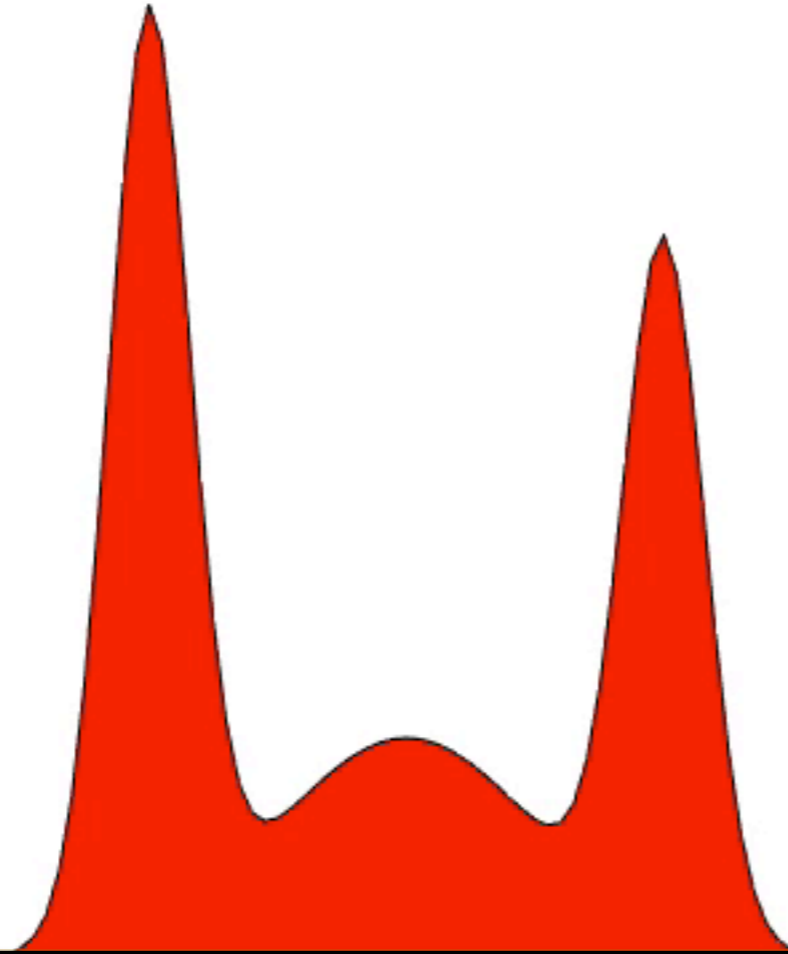
Origins: Monge Problem

In the 21st Century...



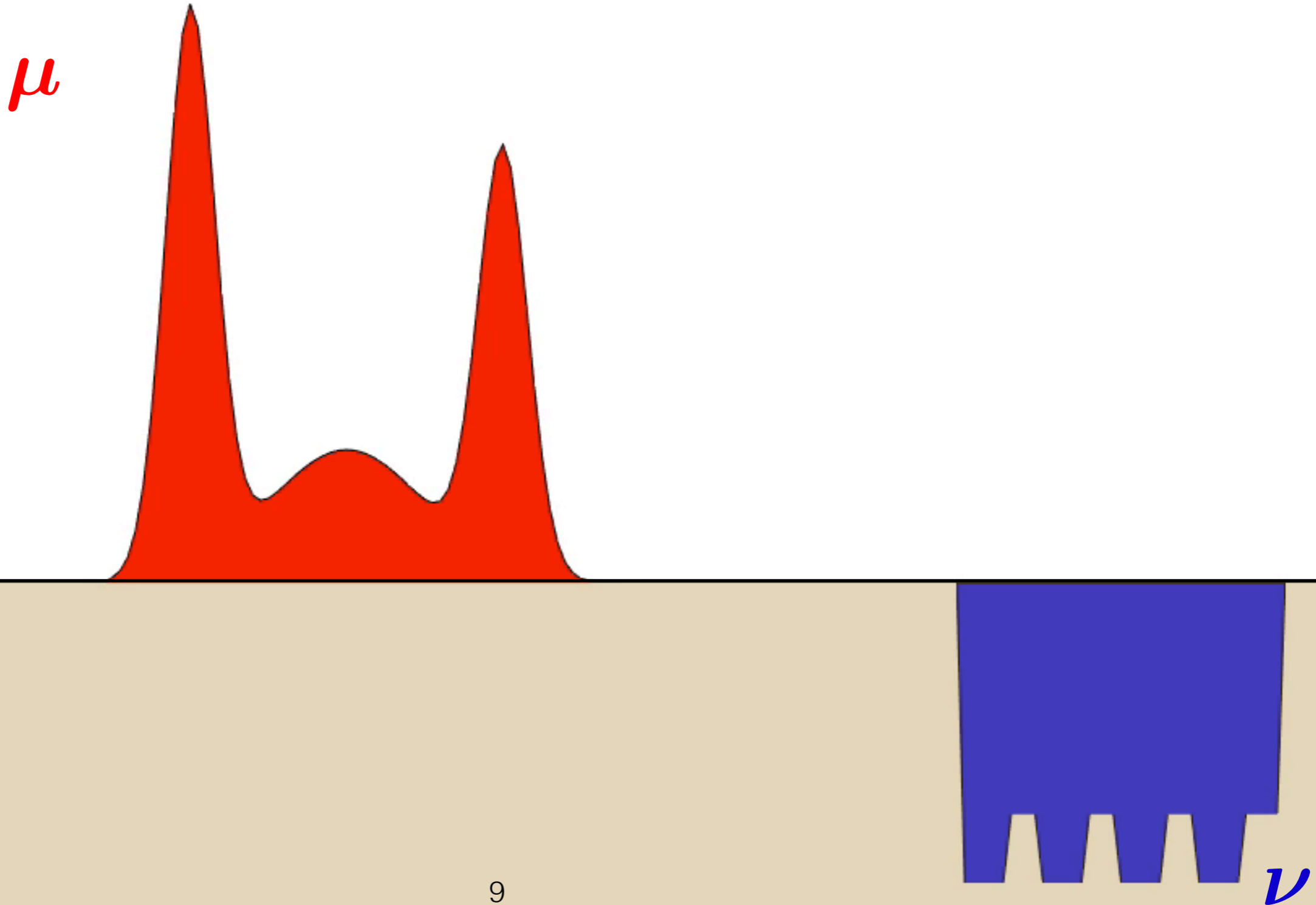
Origins: Monge Problem

In the 21st Century...



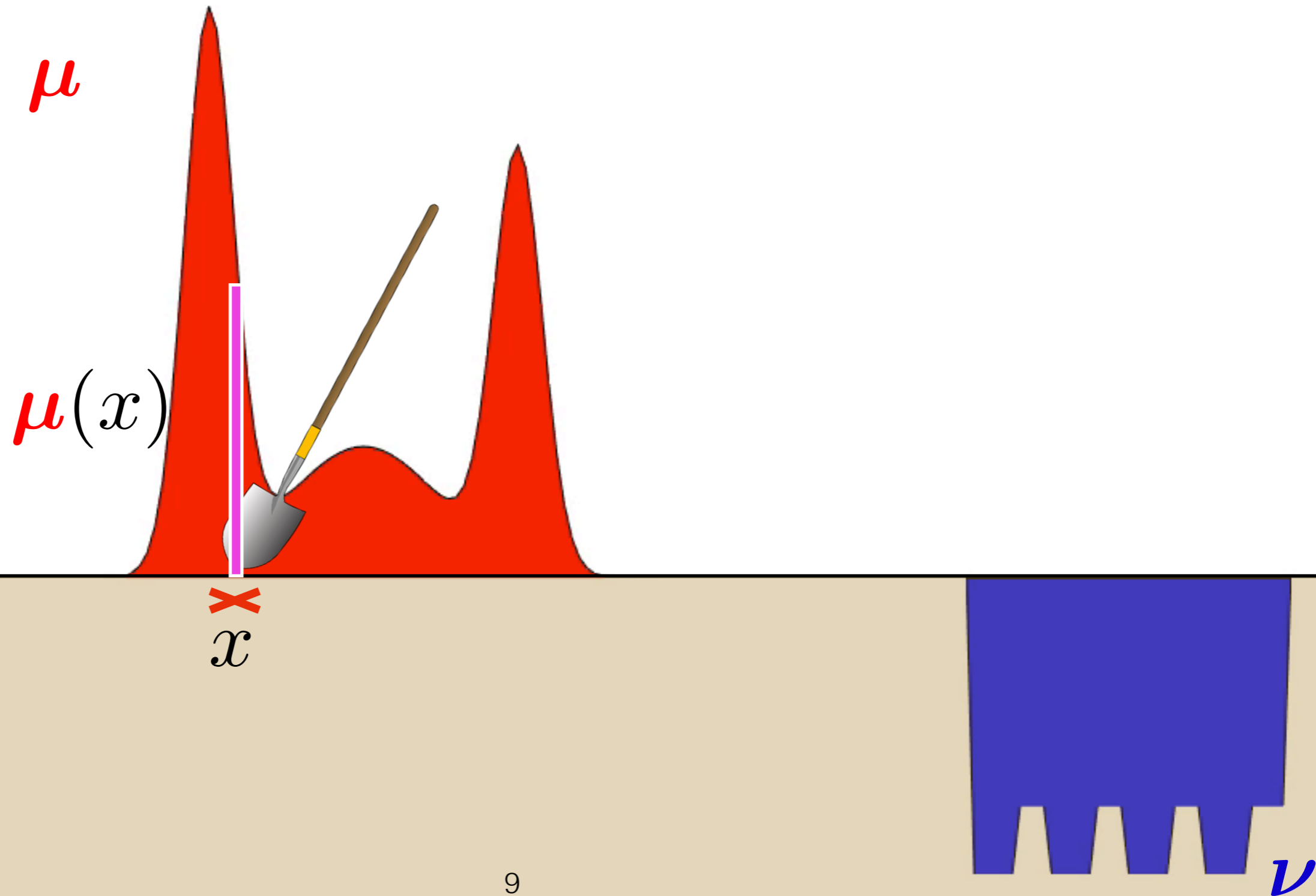
Origins: Monge's Problem

In 1781 however...



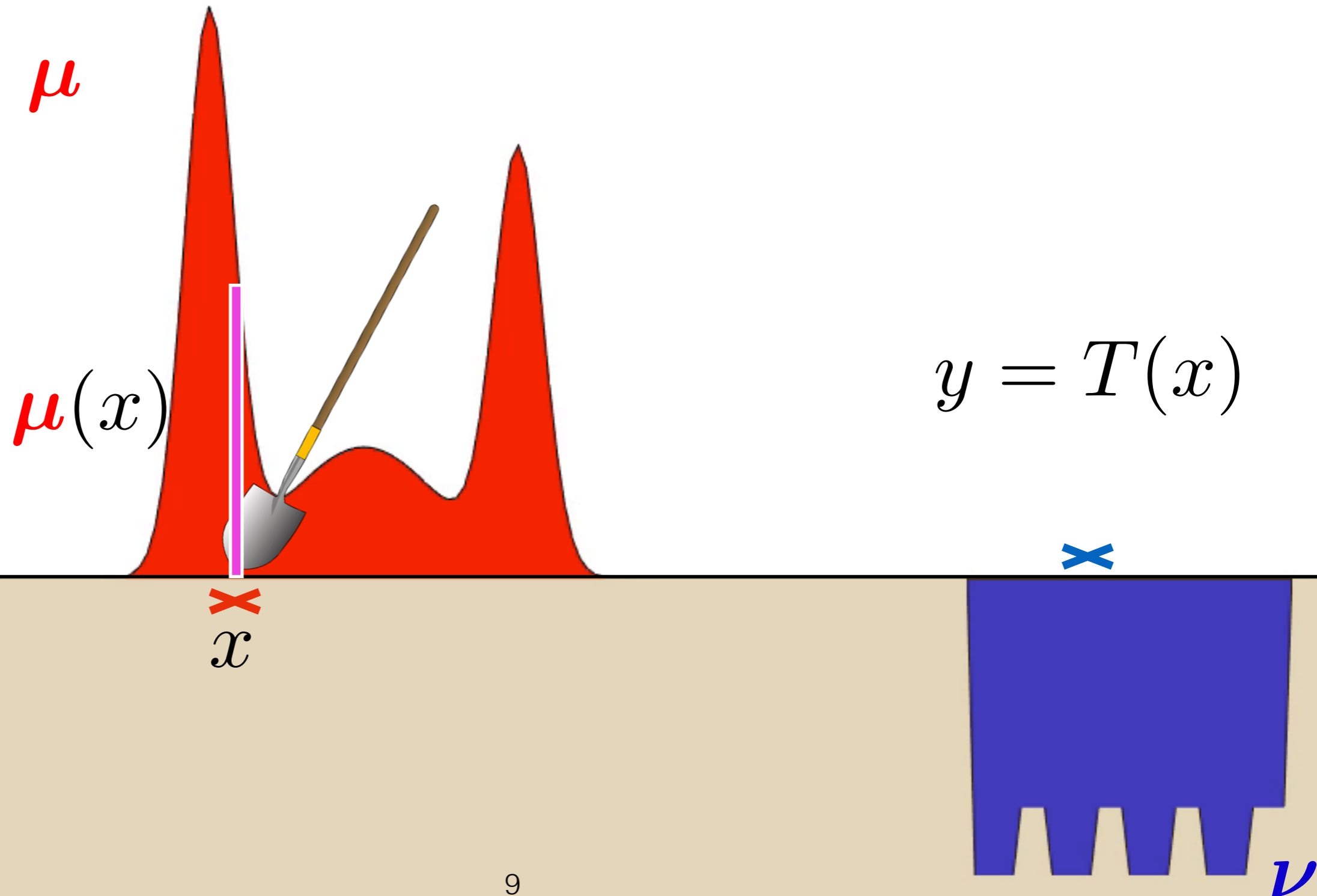
Origins: Monge's Problem

In 1781 however...



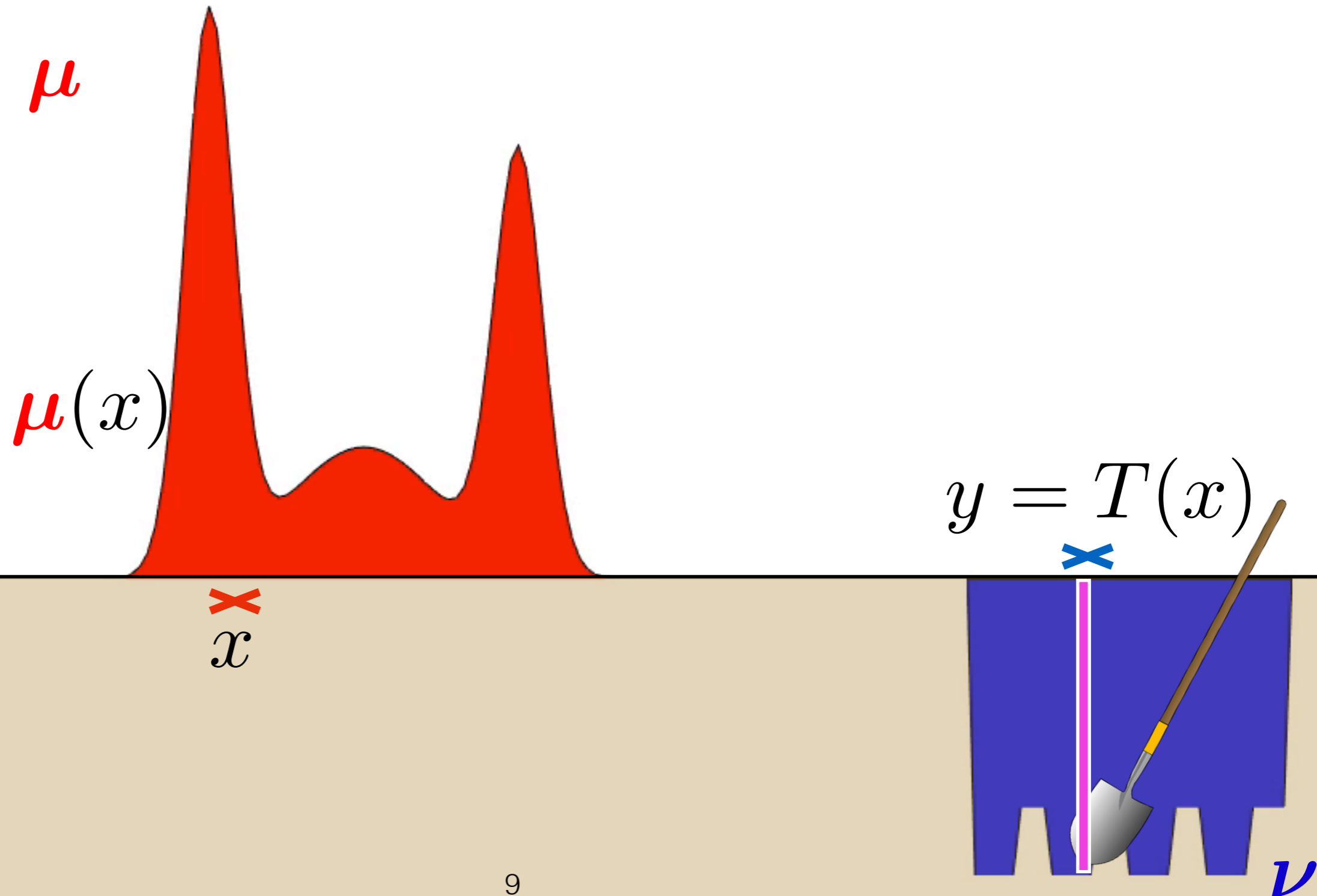
Origins: Monge's Problem

In 1781 however...



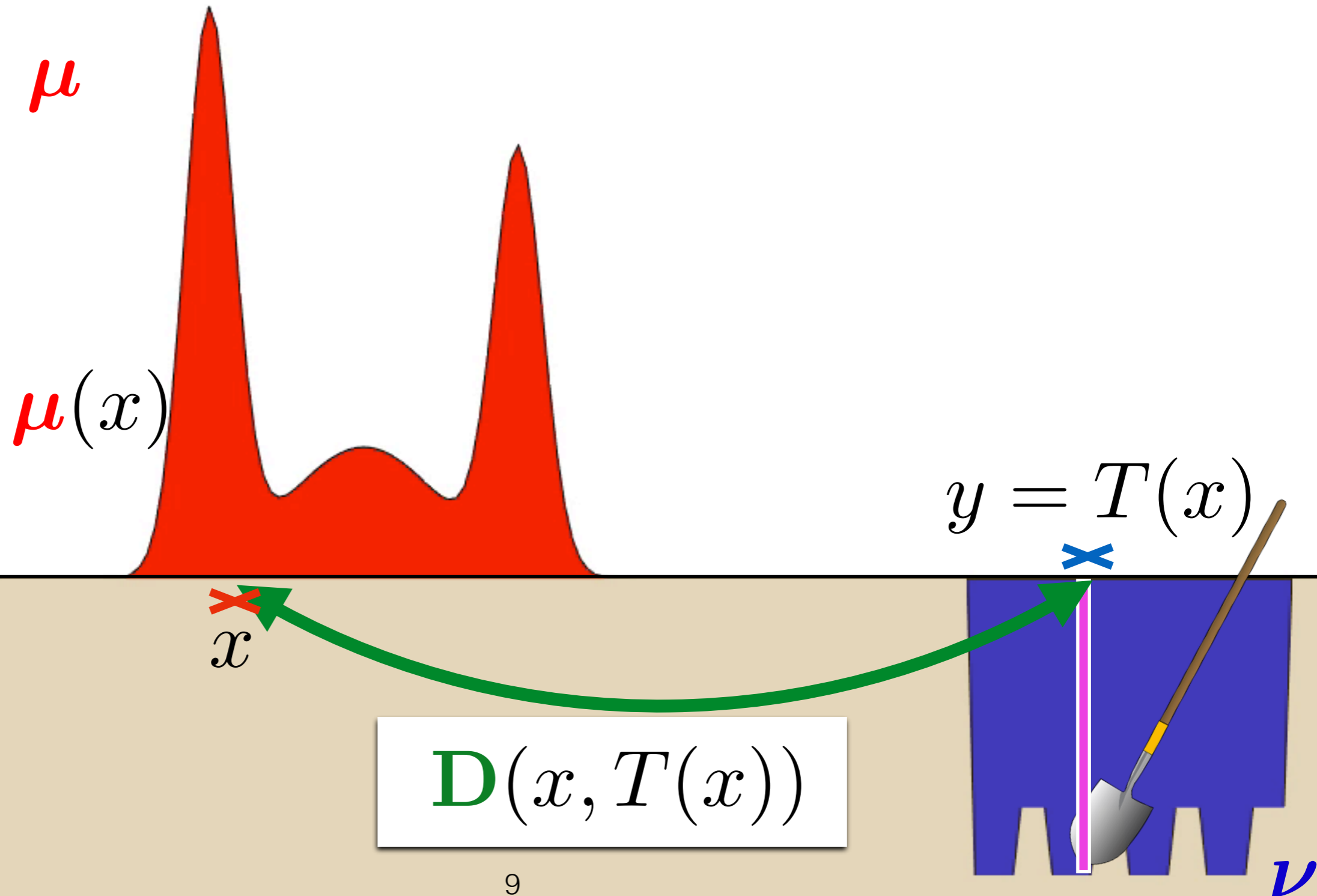
Origins: Monge's Problem

In 1781 however...



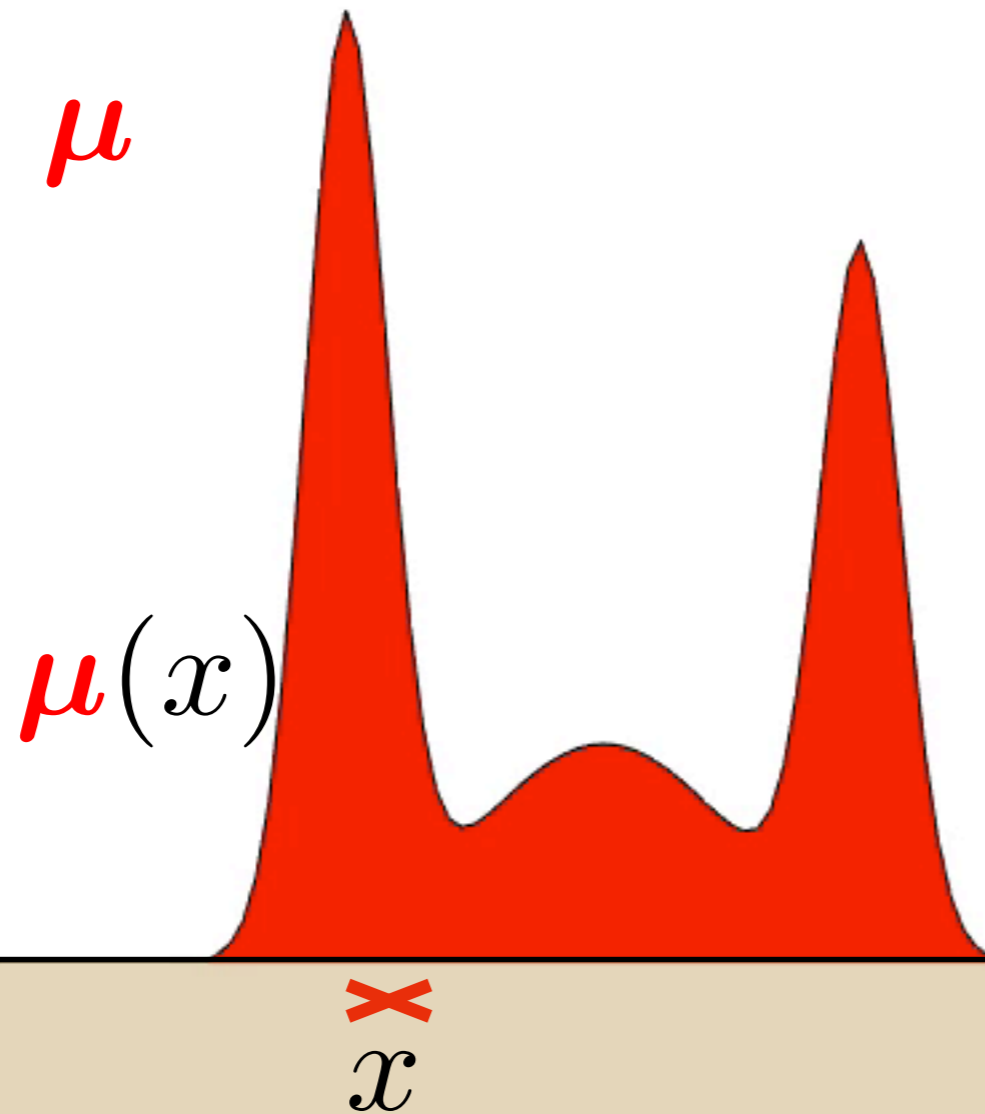
Origins: Monge's Problem

In 1781 however...

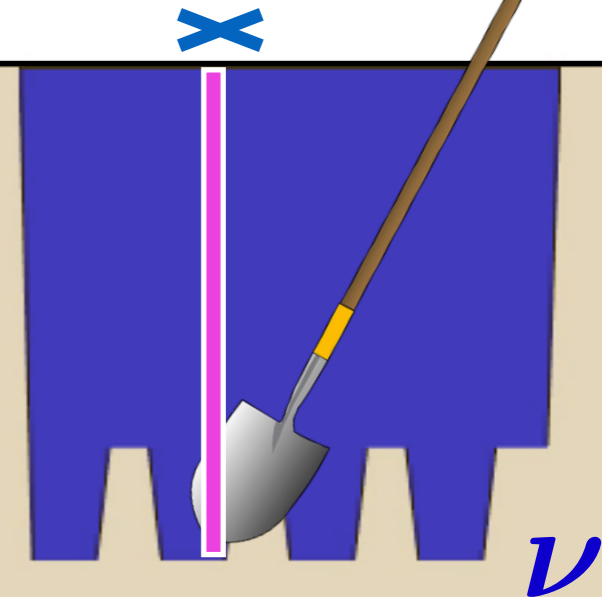


Origins: Monge's Problem

In 1781 however...



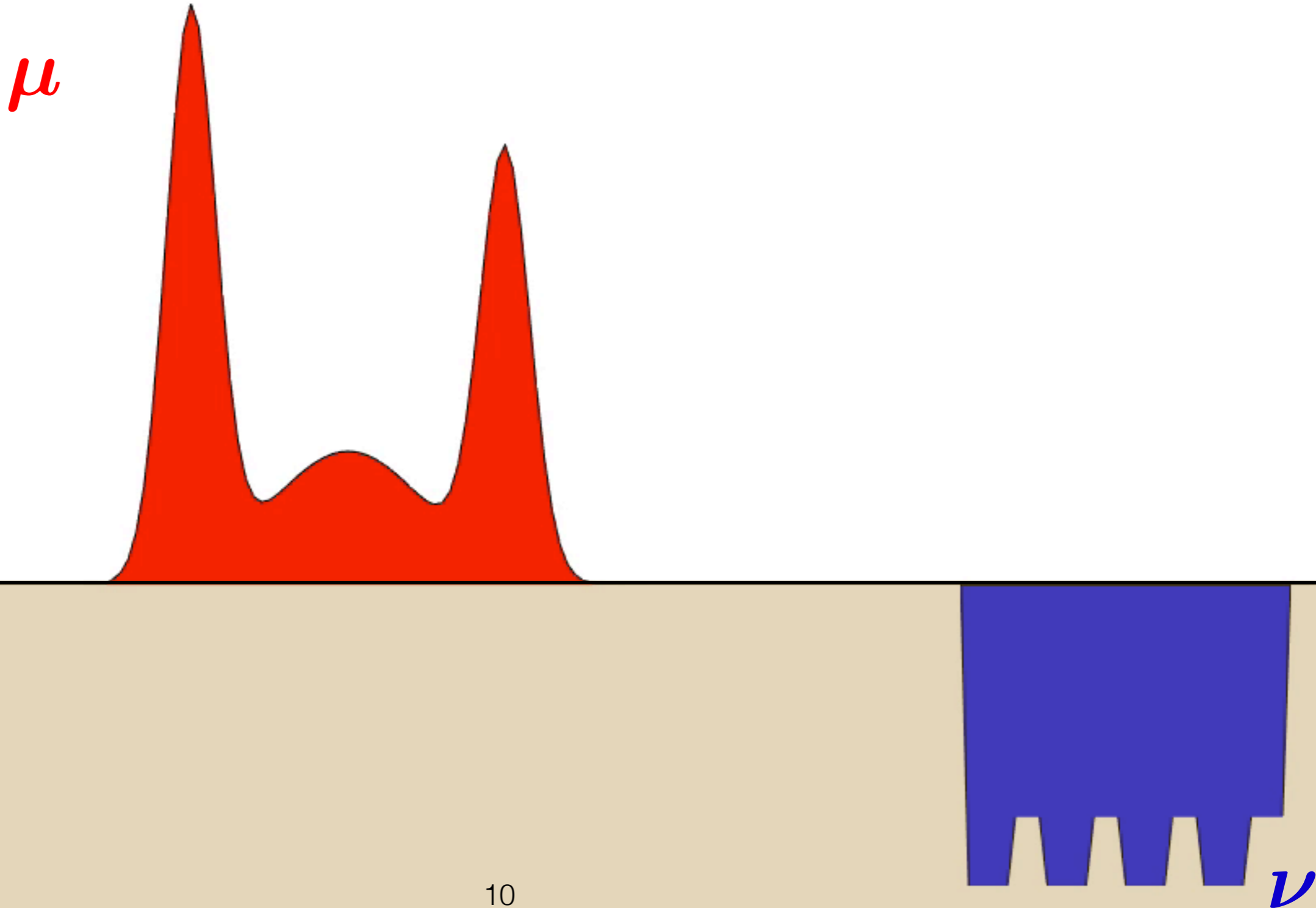
$$y = T(x)$$



work: $\mu(x) D(x, T(x))$

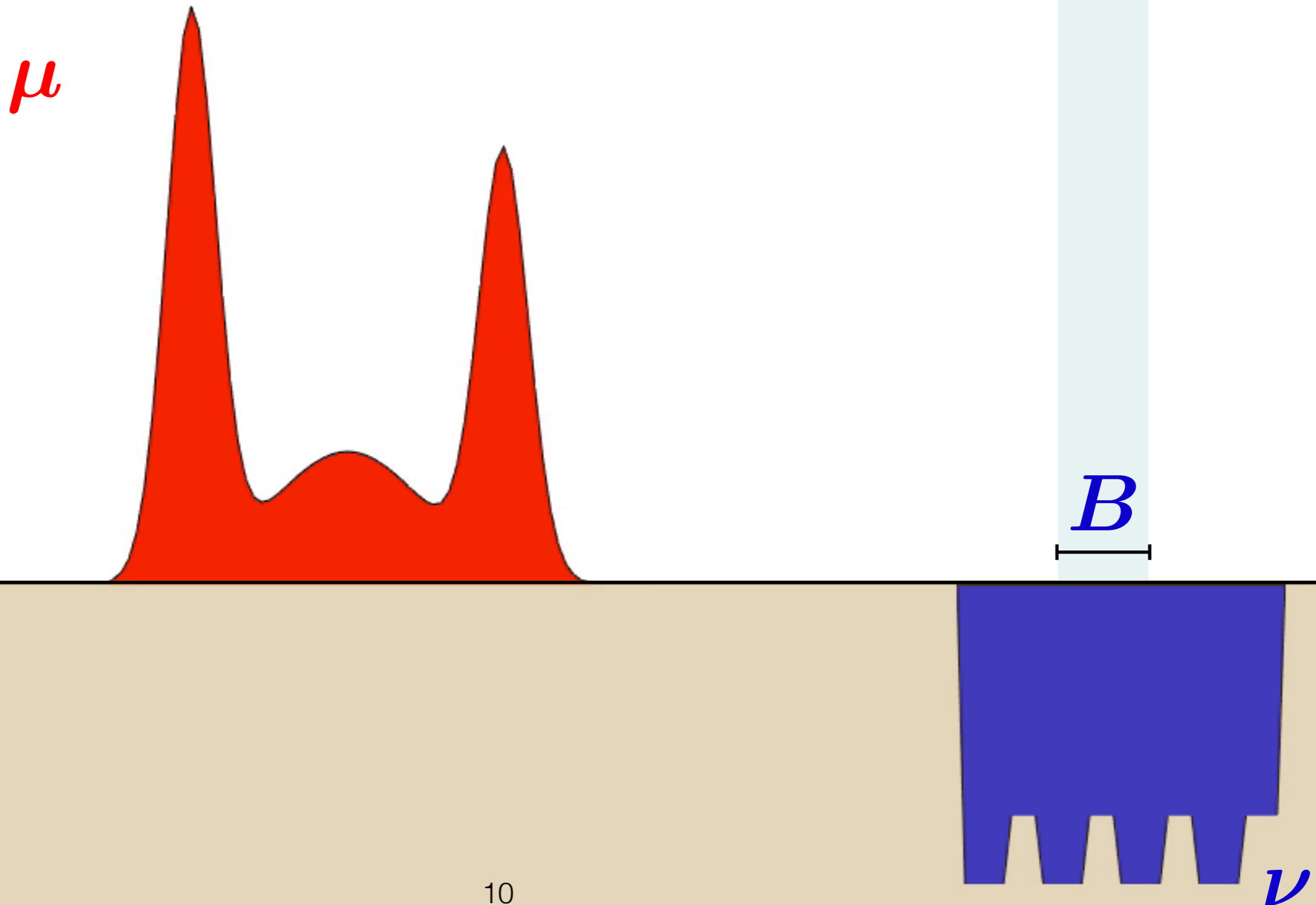
Origins: Monge's Problem

T must map red to blue.



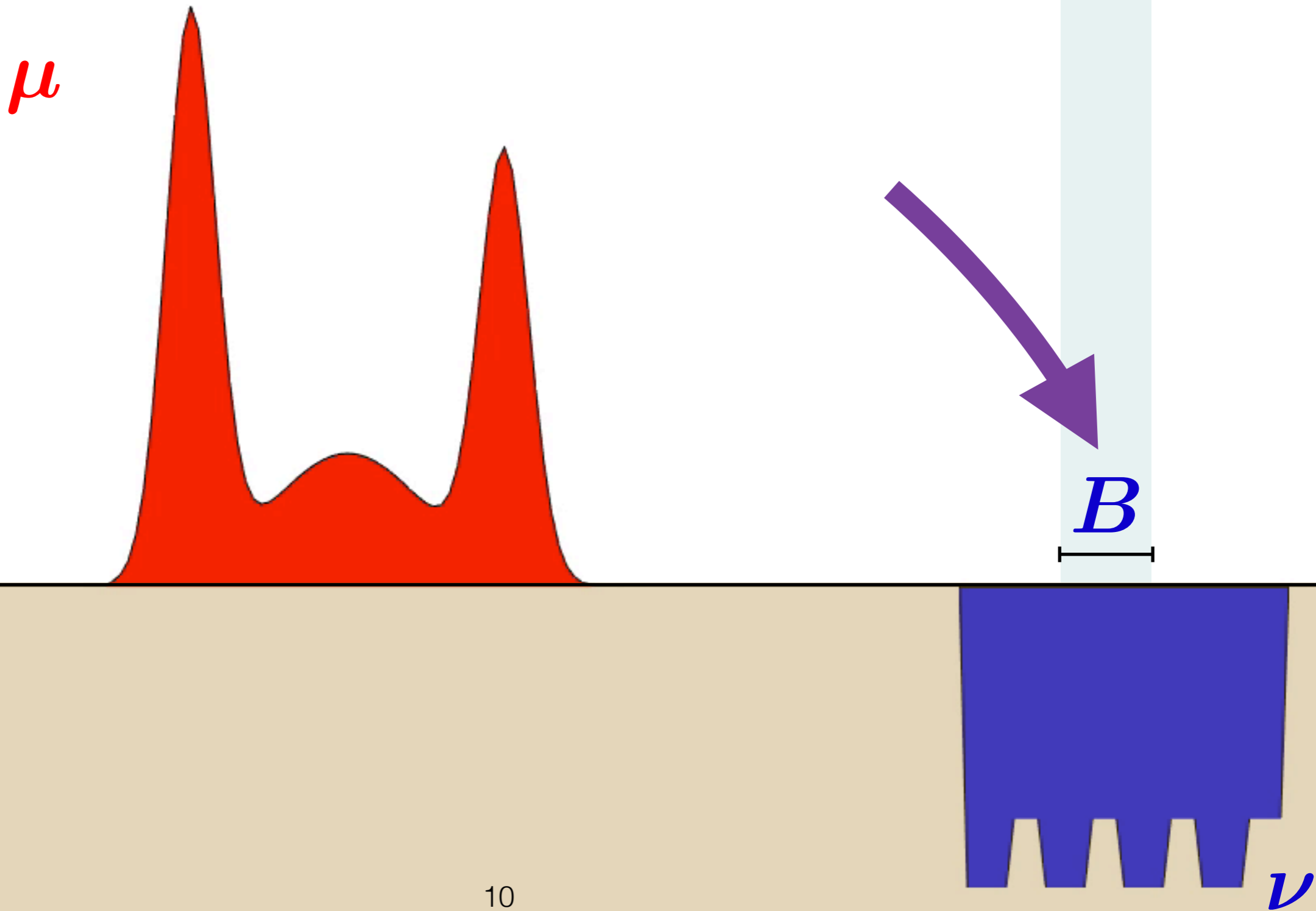
Origins: Monge's Problem

T must map red to blue.



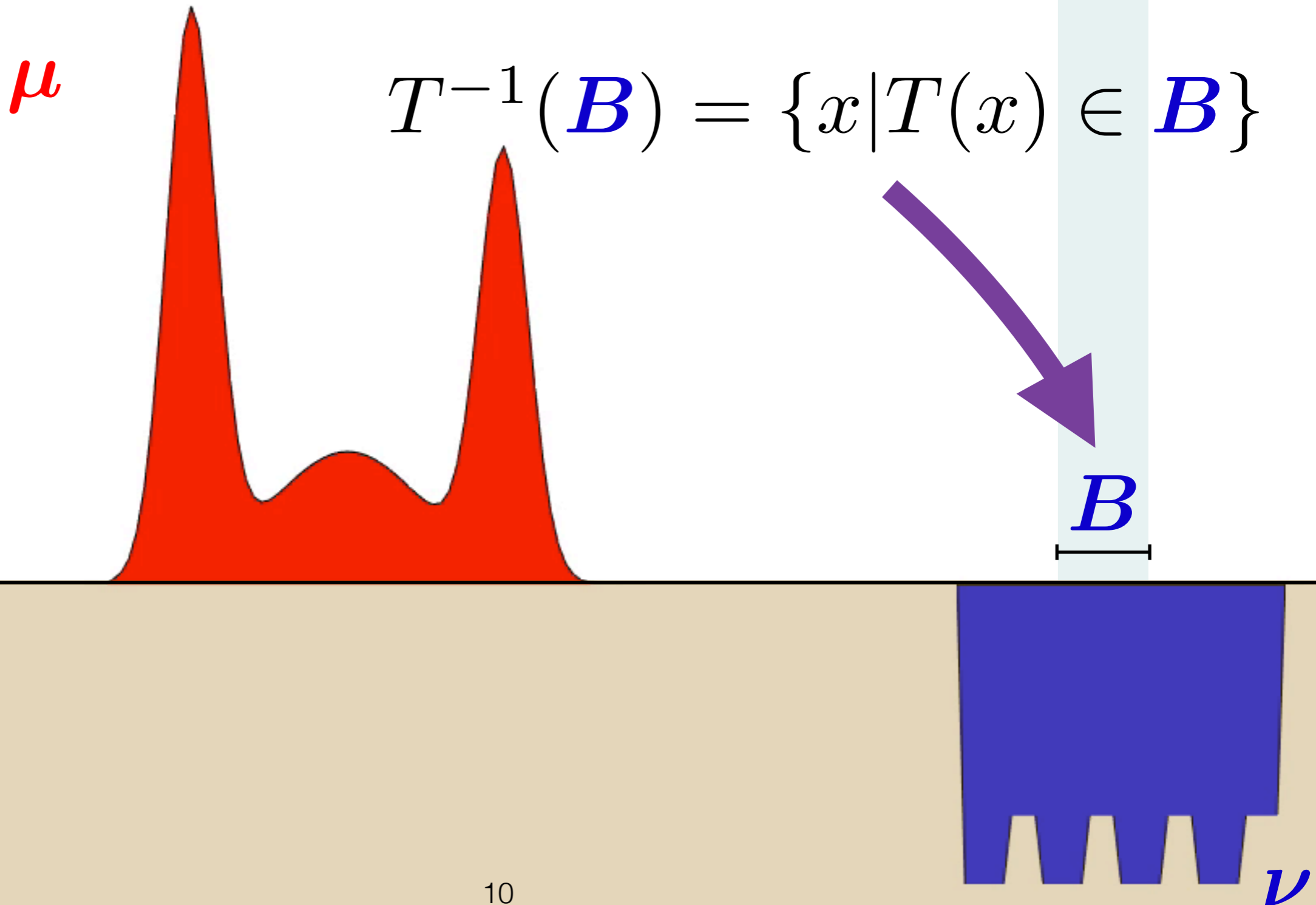
Origins: Monge's Problem

T must map red to blue.



Origins: Monge's Problem

T must map red to blue.

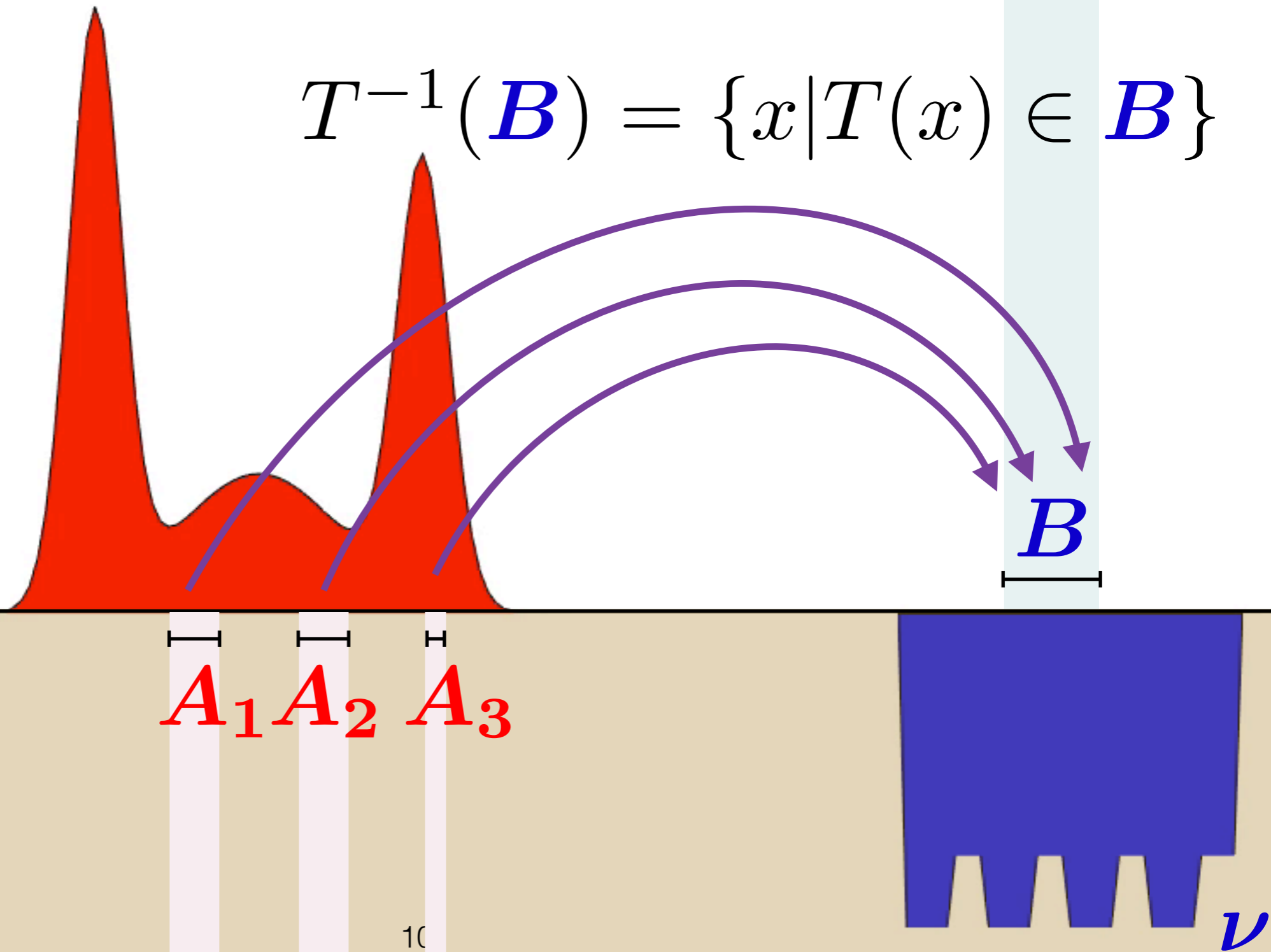


Origins: Monge's Problem

T must map red to blue.

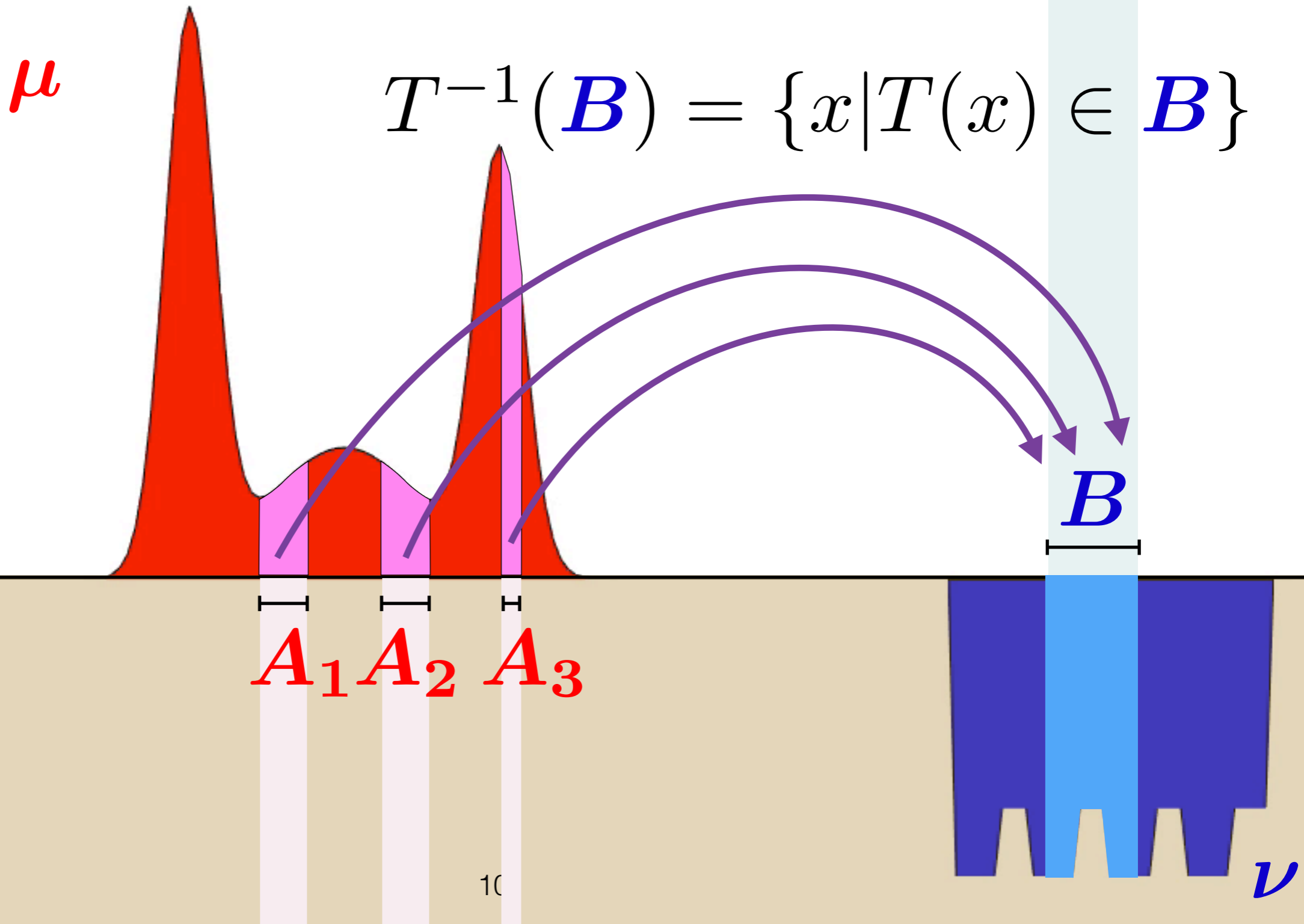
μ

$$T^{-1}(B) = \{x \mid T(x) \in B\}$$



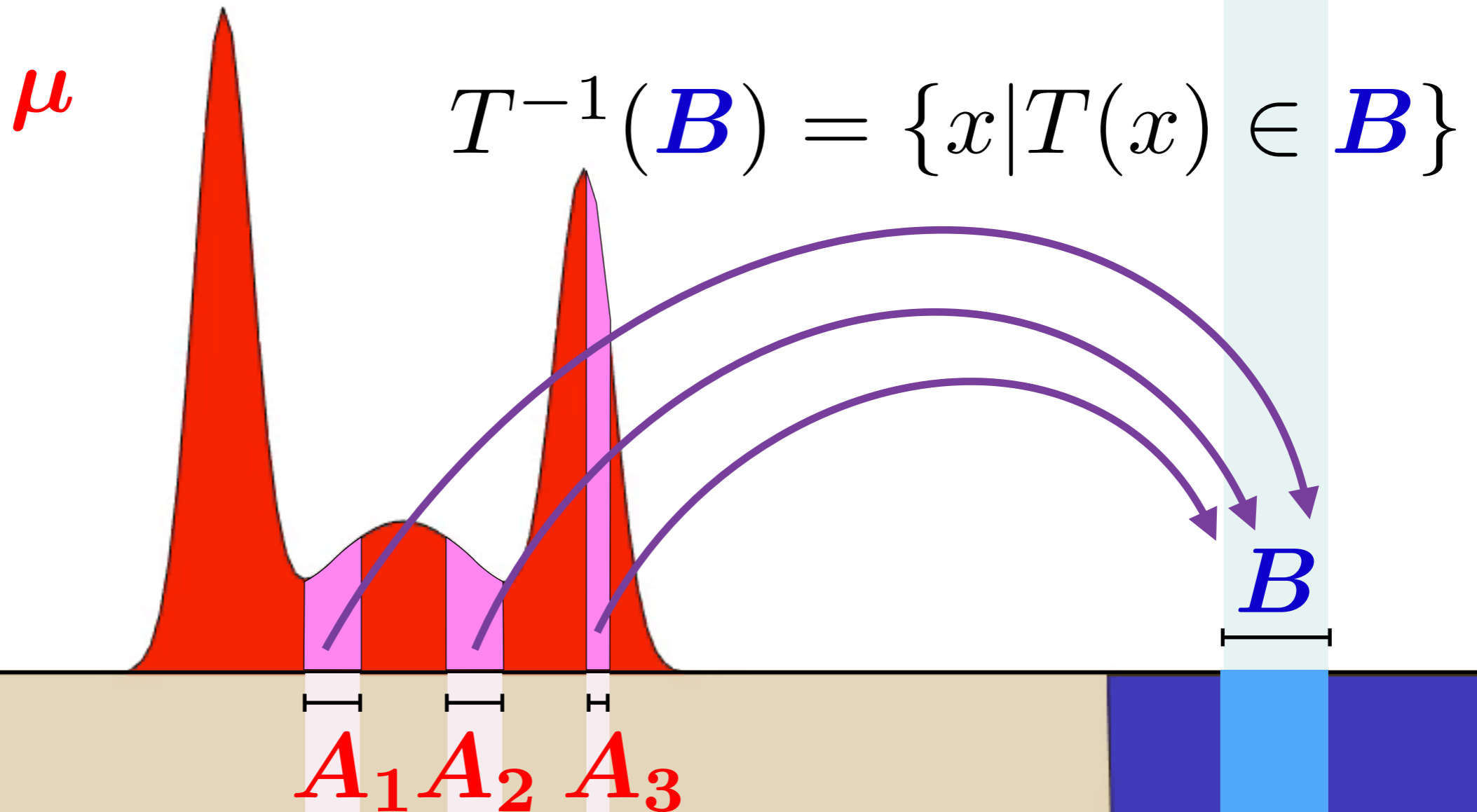
Origins: Monge's Problem

T must map red to blue.



Origins: Monge's Problem

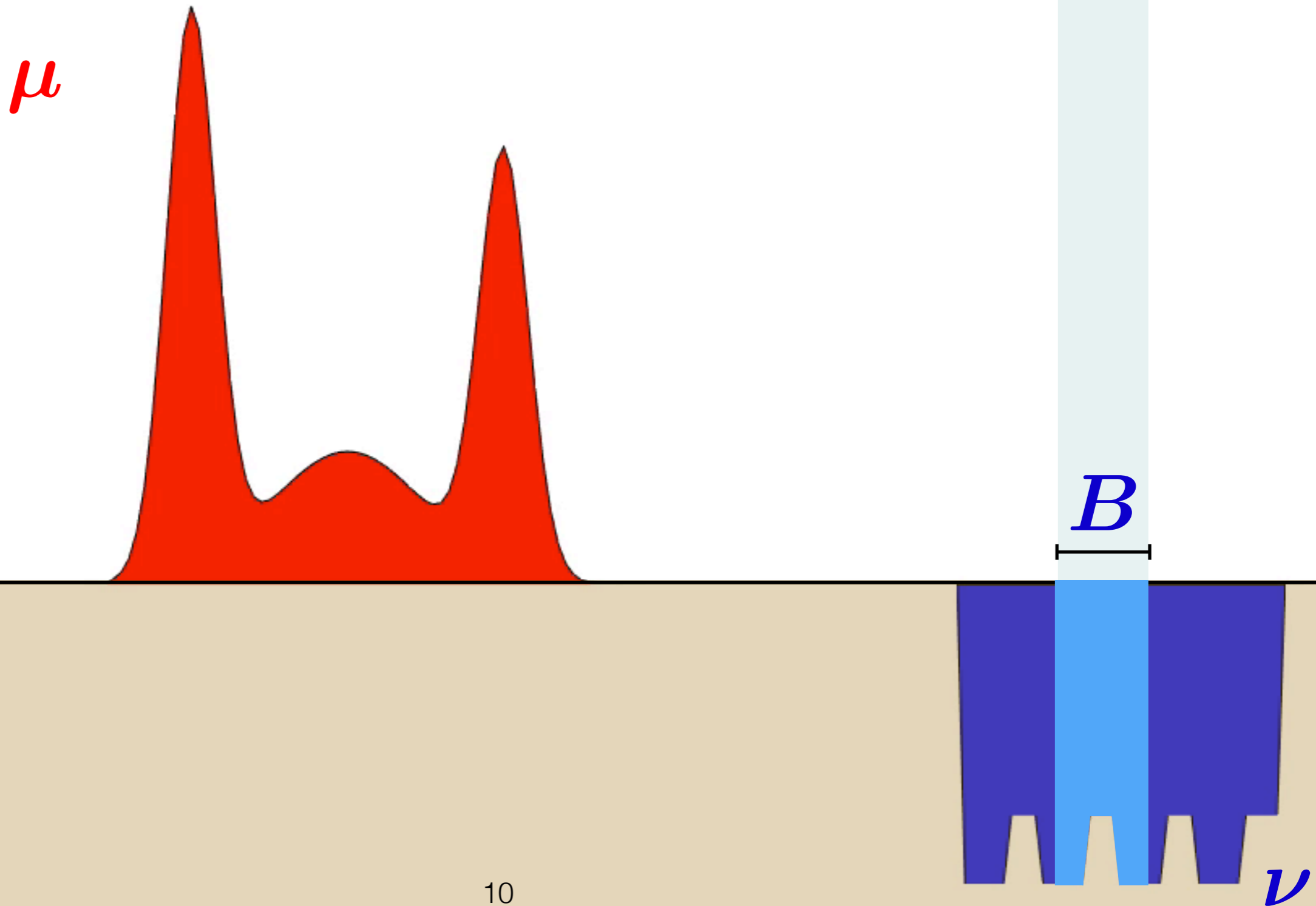
T must map red to blue.



$$\mu(A_1) + \mu(A_2) + \mu(A_3) = \nu(B)$$

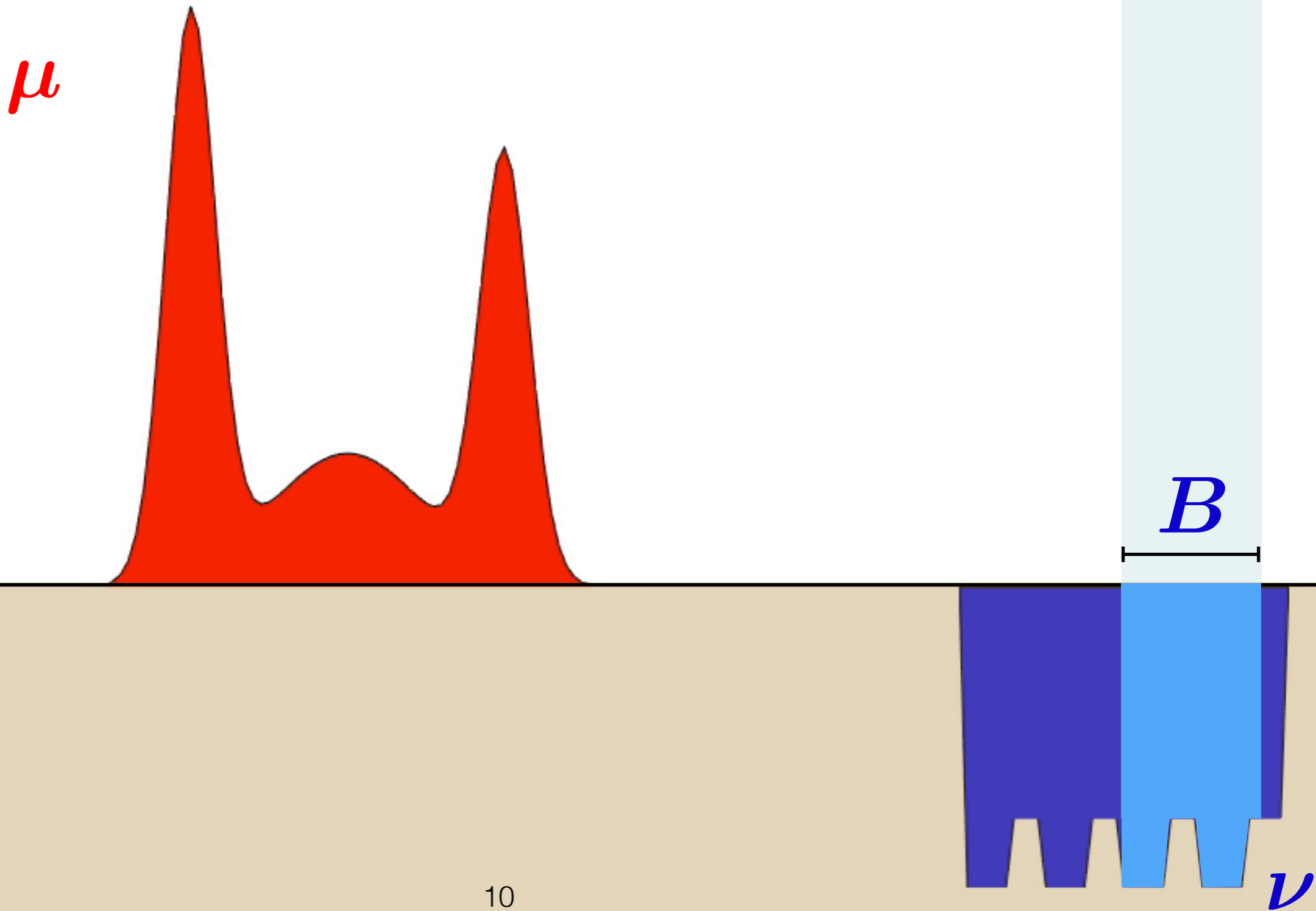
Origins: Monge's Problem

T must map red to blue.



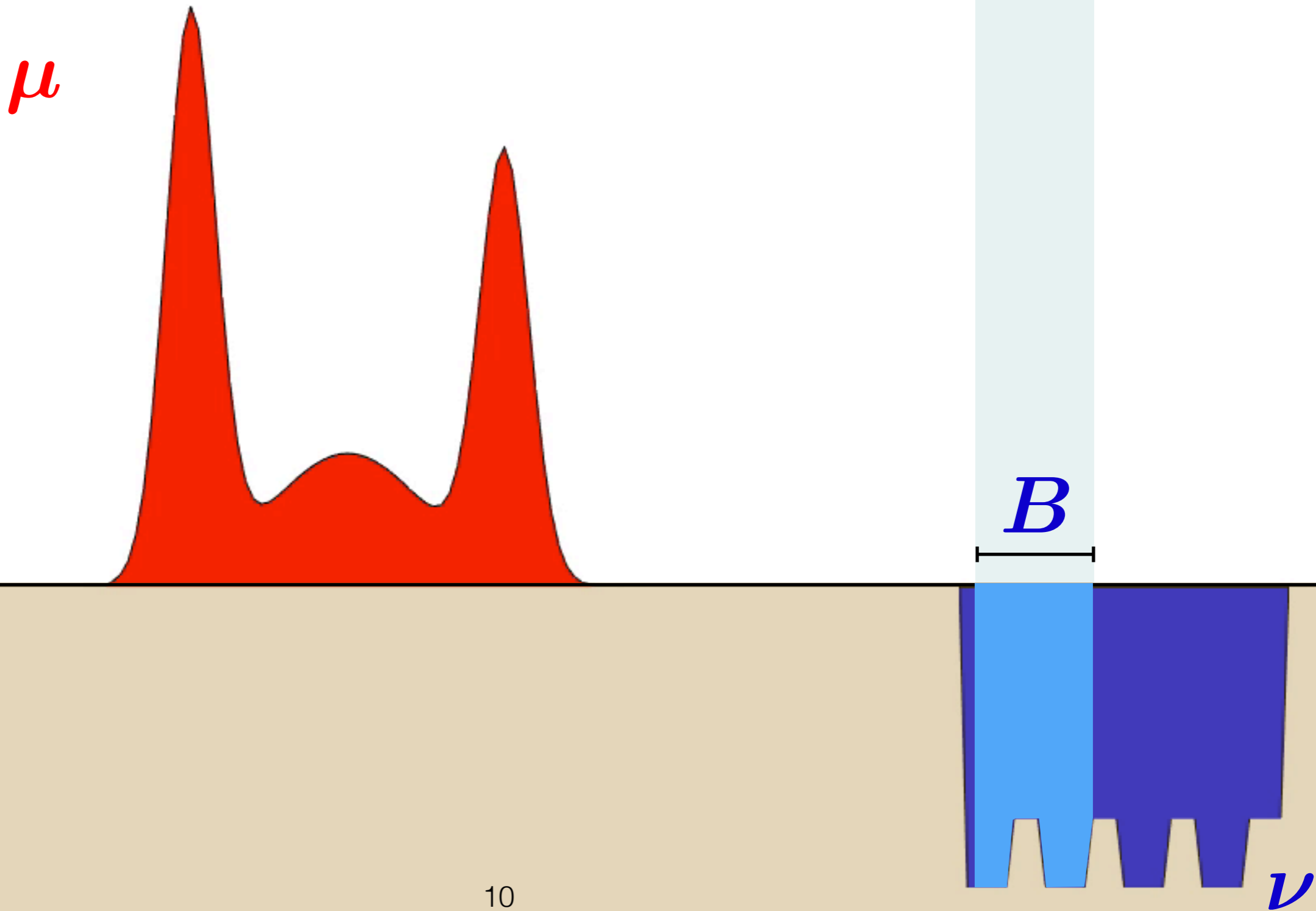
Origins: Monge's Problem

T must map red to blue.



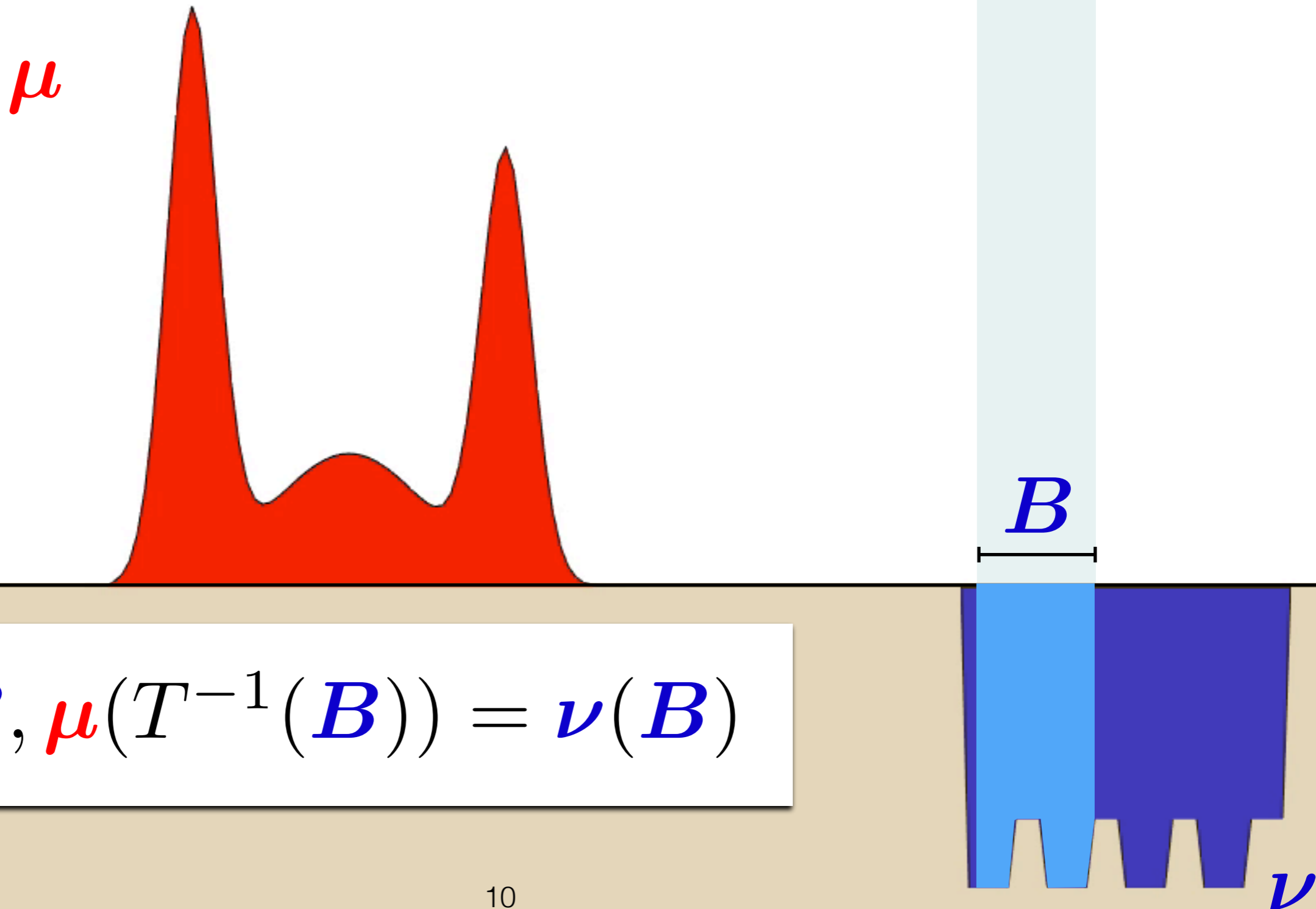
Origins: Monge's Problem

T must map red to blue.



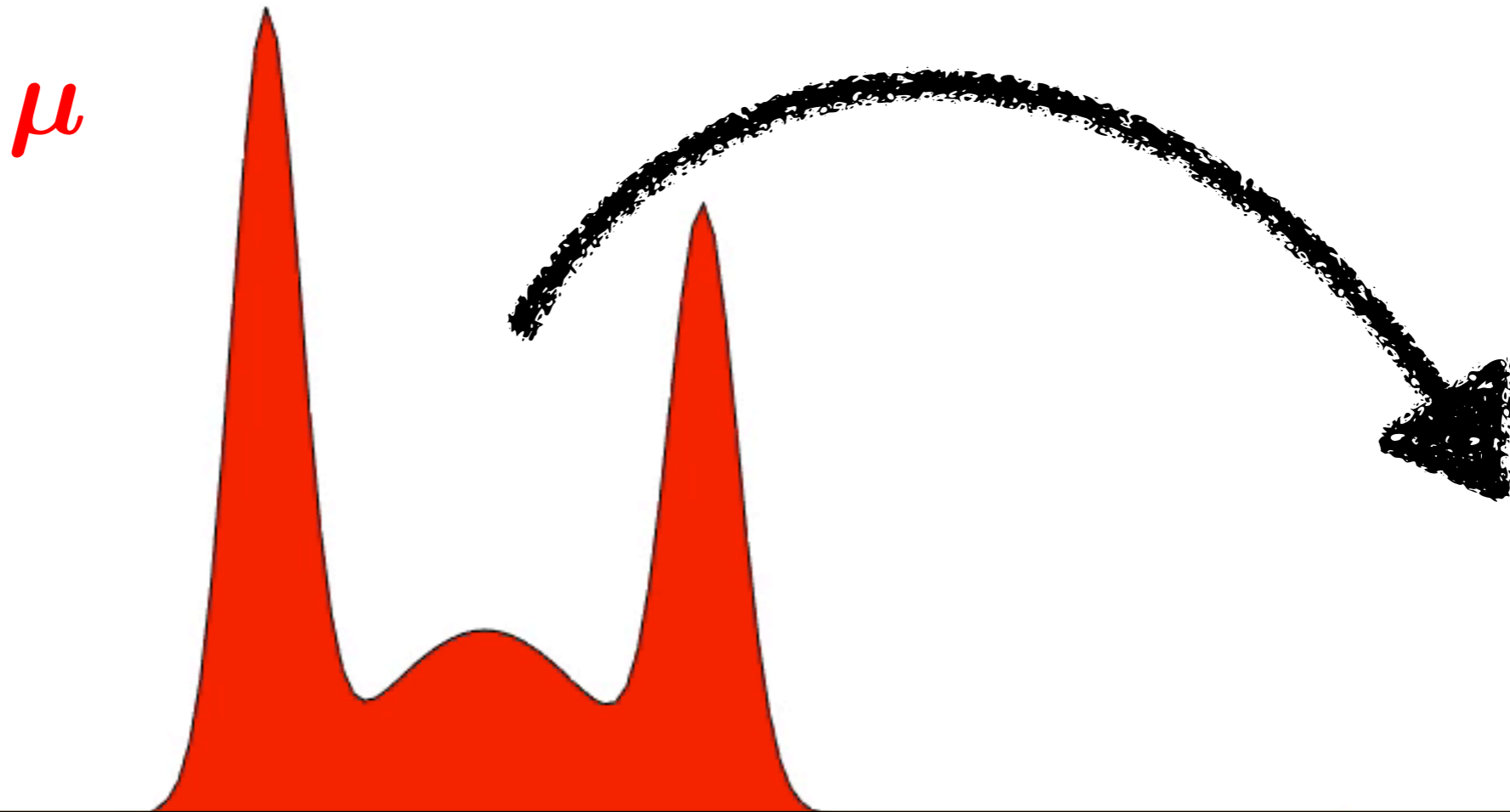
Origins: Monge's Problem

T must map red to blue.



Origins: Monge's Problem

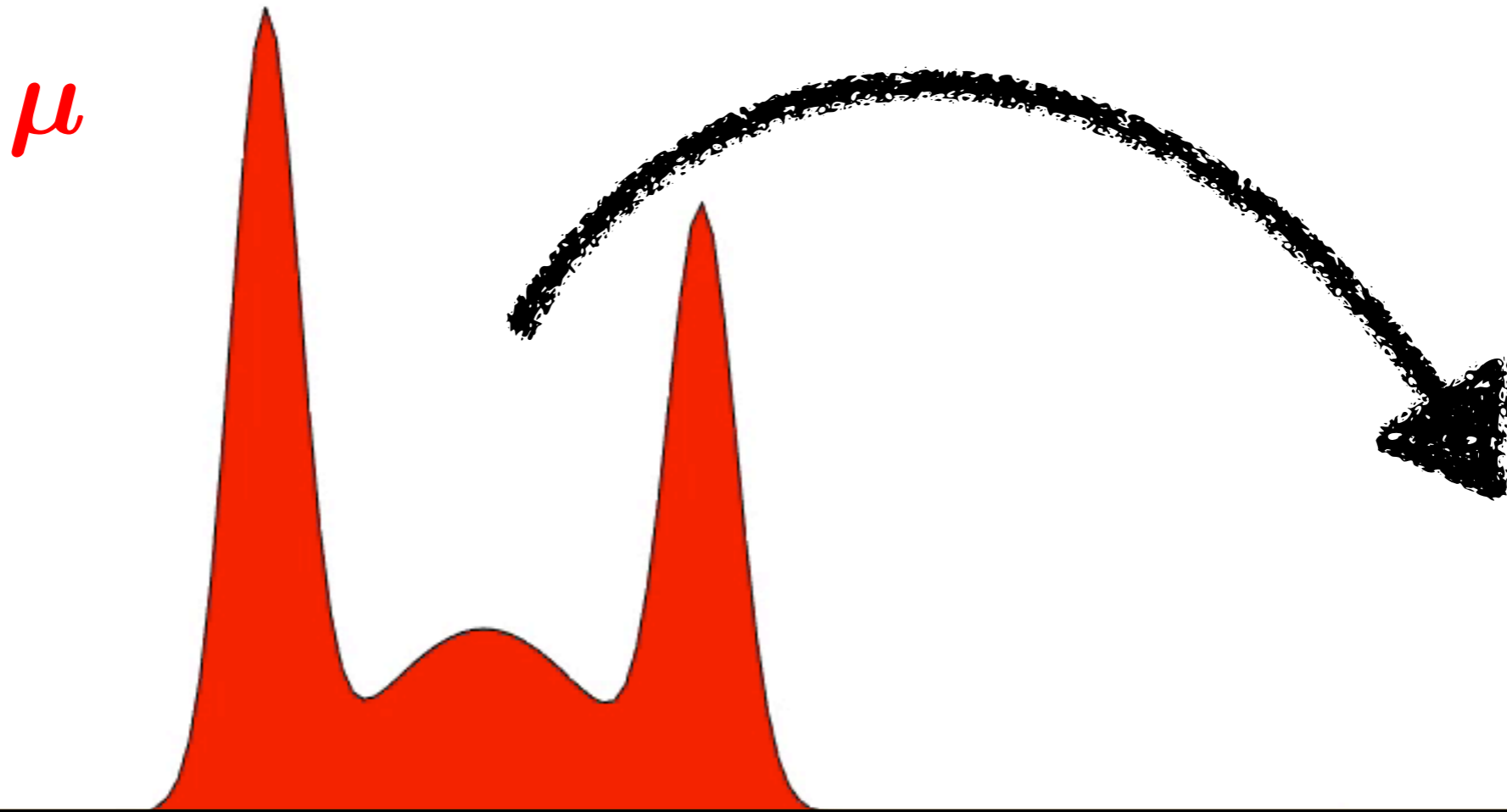
T must **push-forward** the red measure towards the blue



$$T_{\#} \mu = \nu$$

Origins: Monge's Problem

T must **push-forward** the red measure towards the blue



What T s.t. $T_{\#}\mu = \nu$
minimizes $\int D(x, T(x)) \mu(dx)$?

Kantorovich Problem



Kantorovich



Tolstoi
1930



Hitchcock



1939

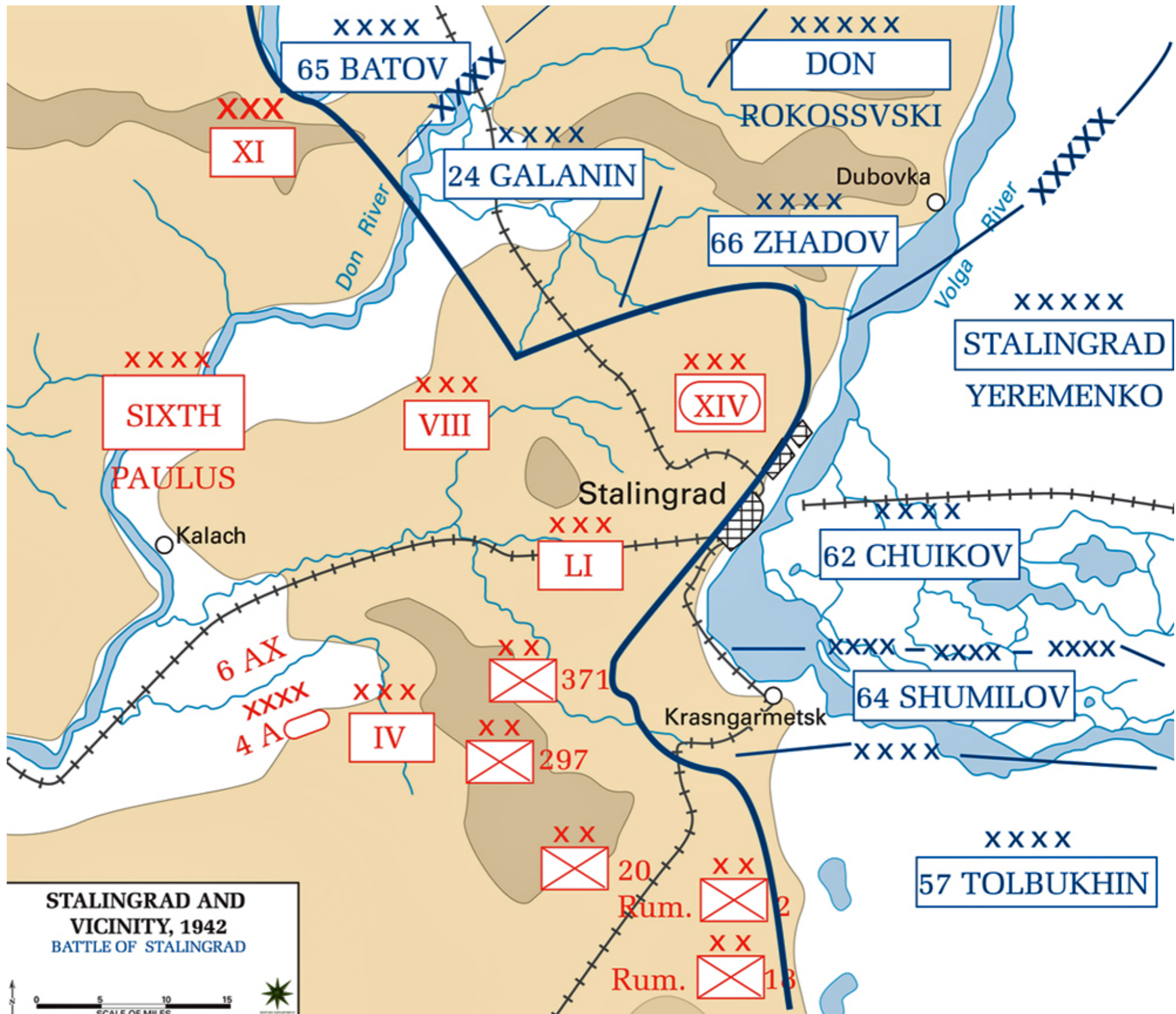
THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES

BY FRANK L. HITCHCOCK

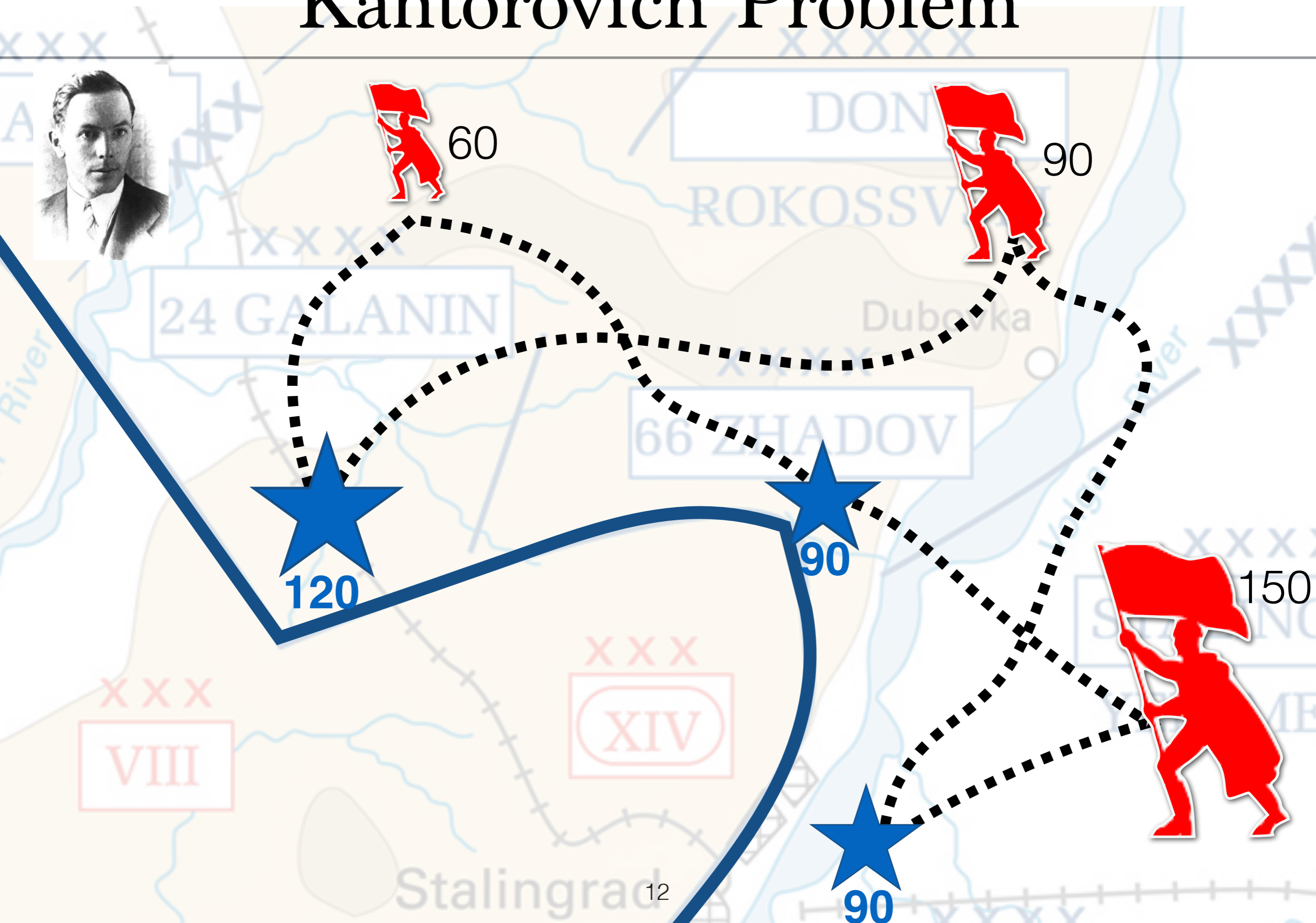
1. Statement of the problem. When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary from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

1941

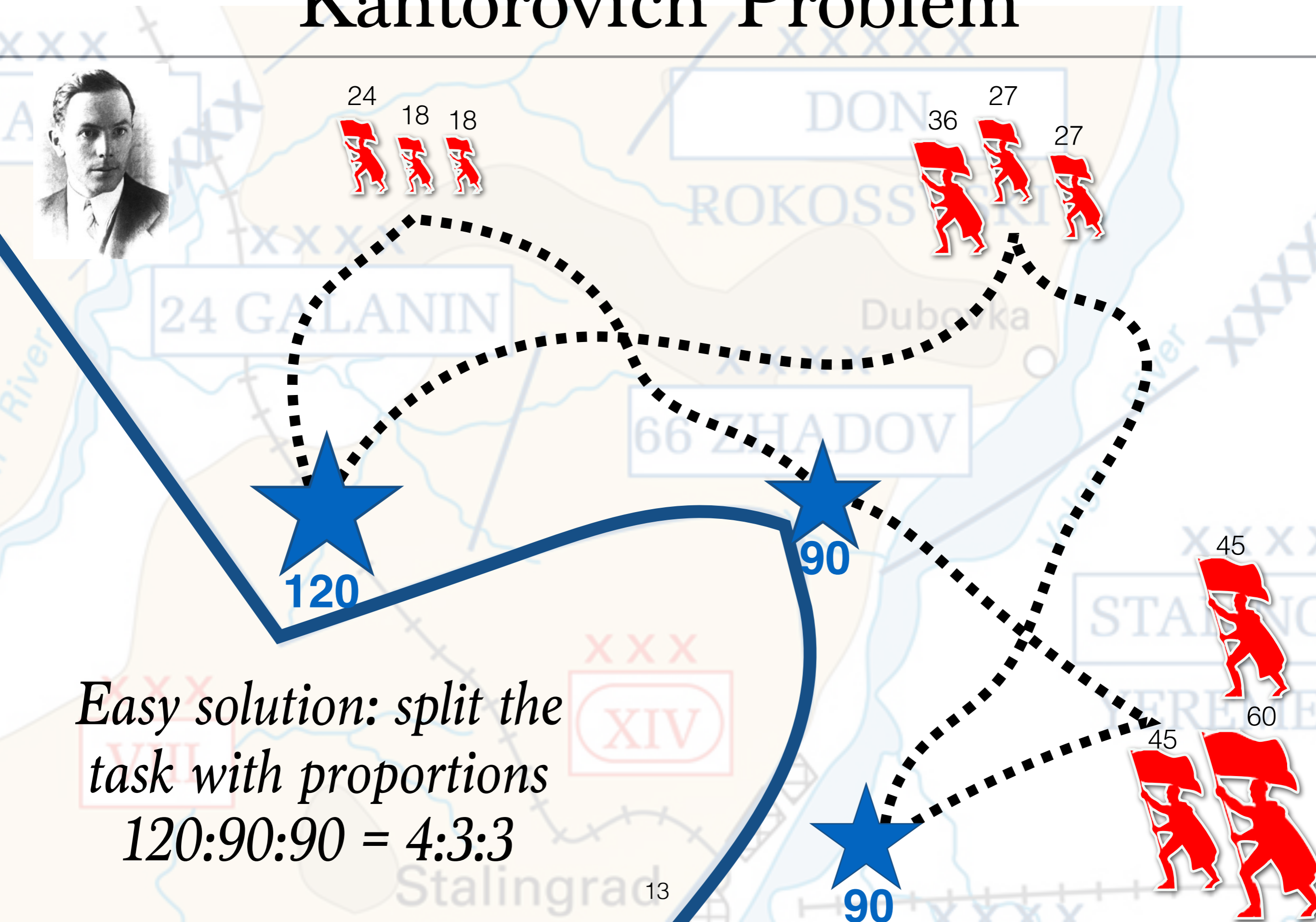
Kantorovich Problem



Kantorovich Problem

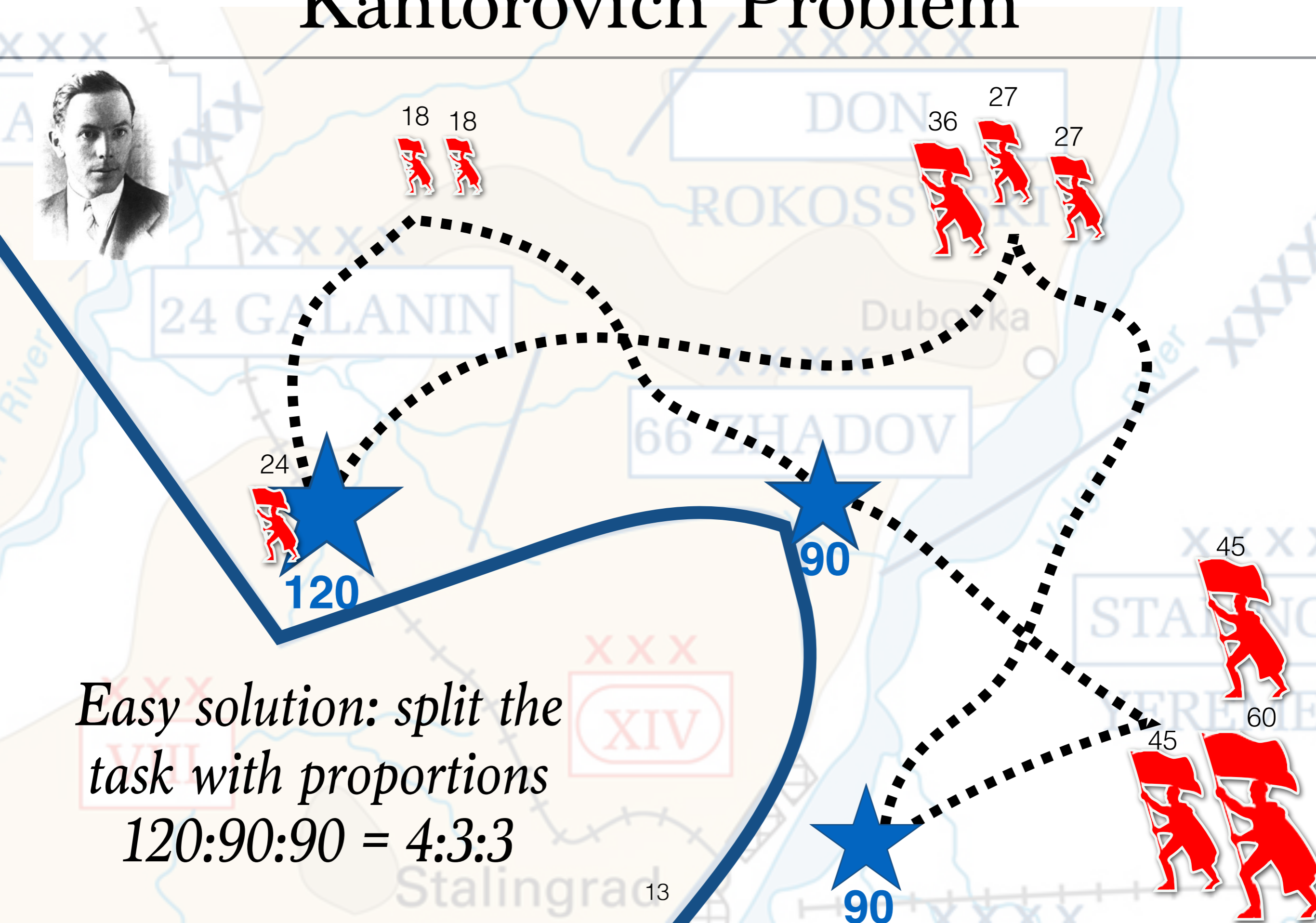


Kantorovich Problem



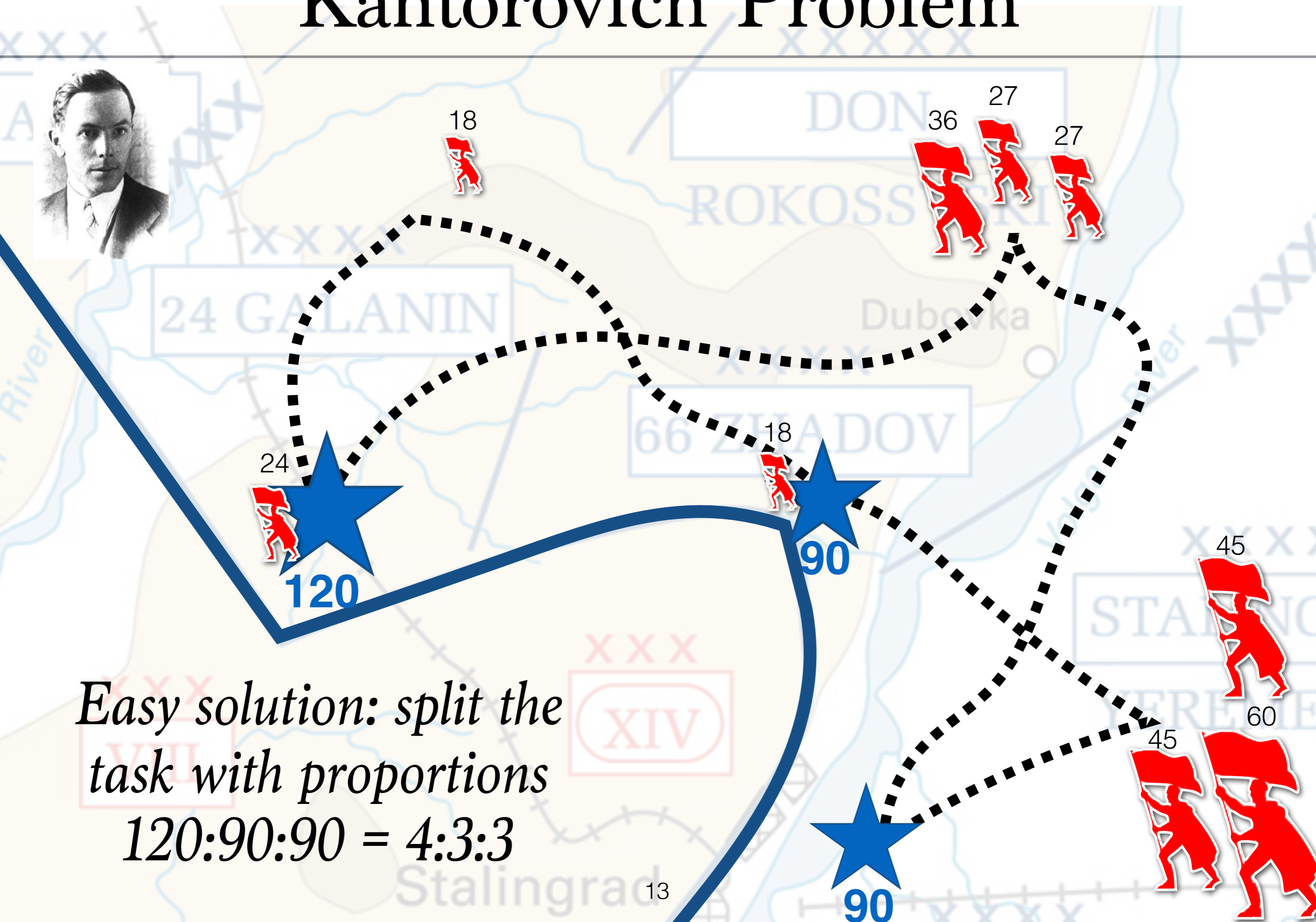
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



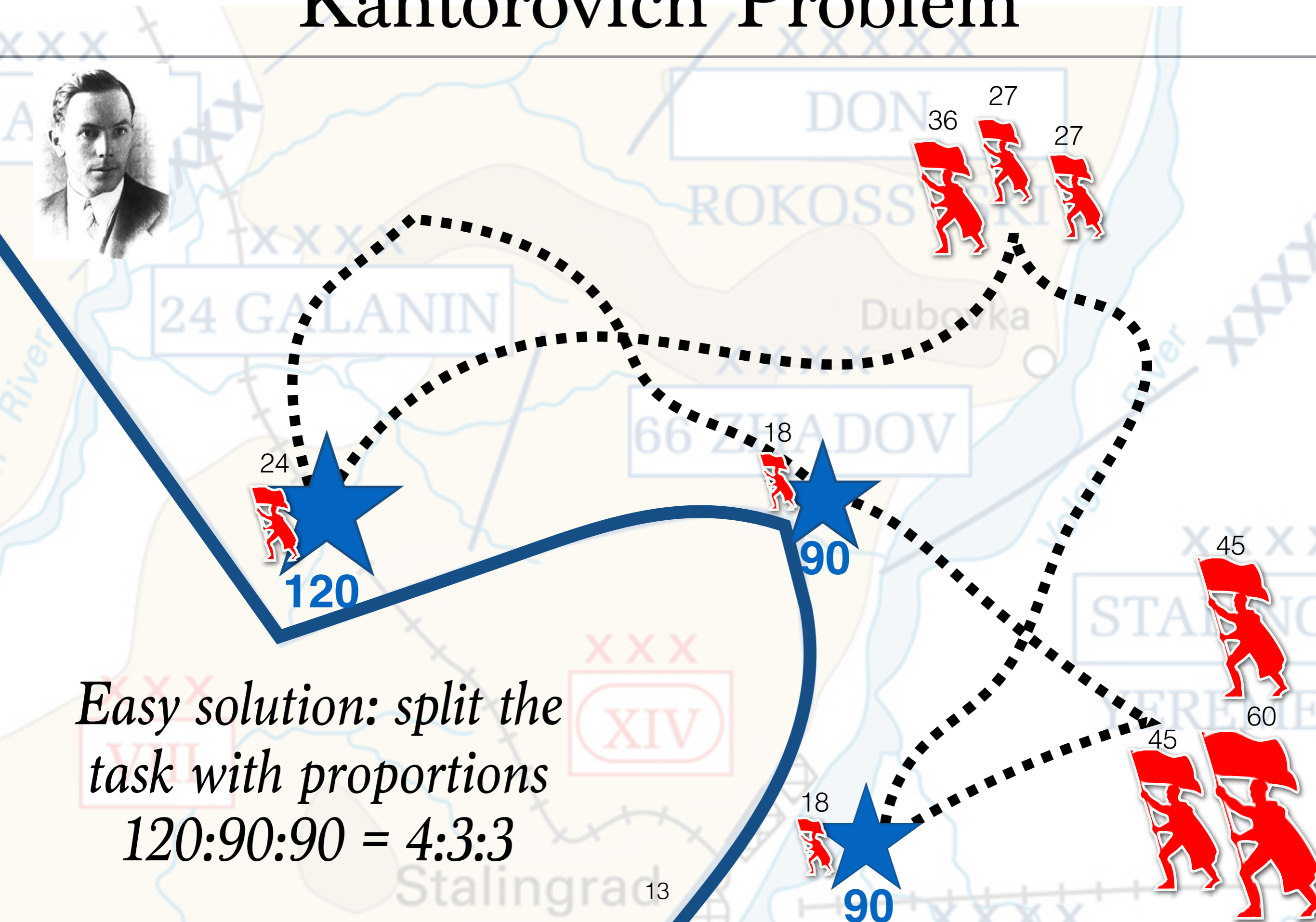
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



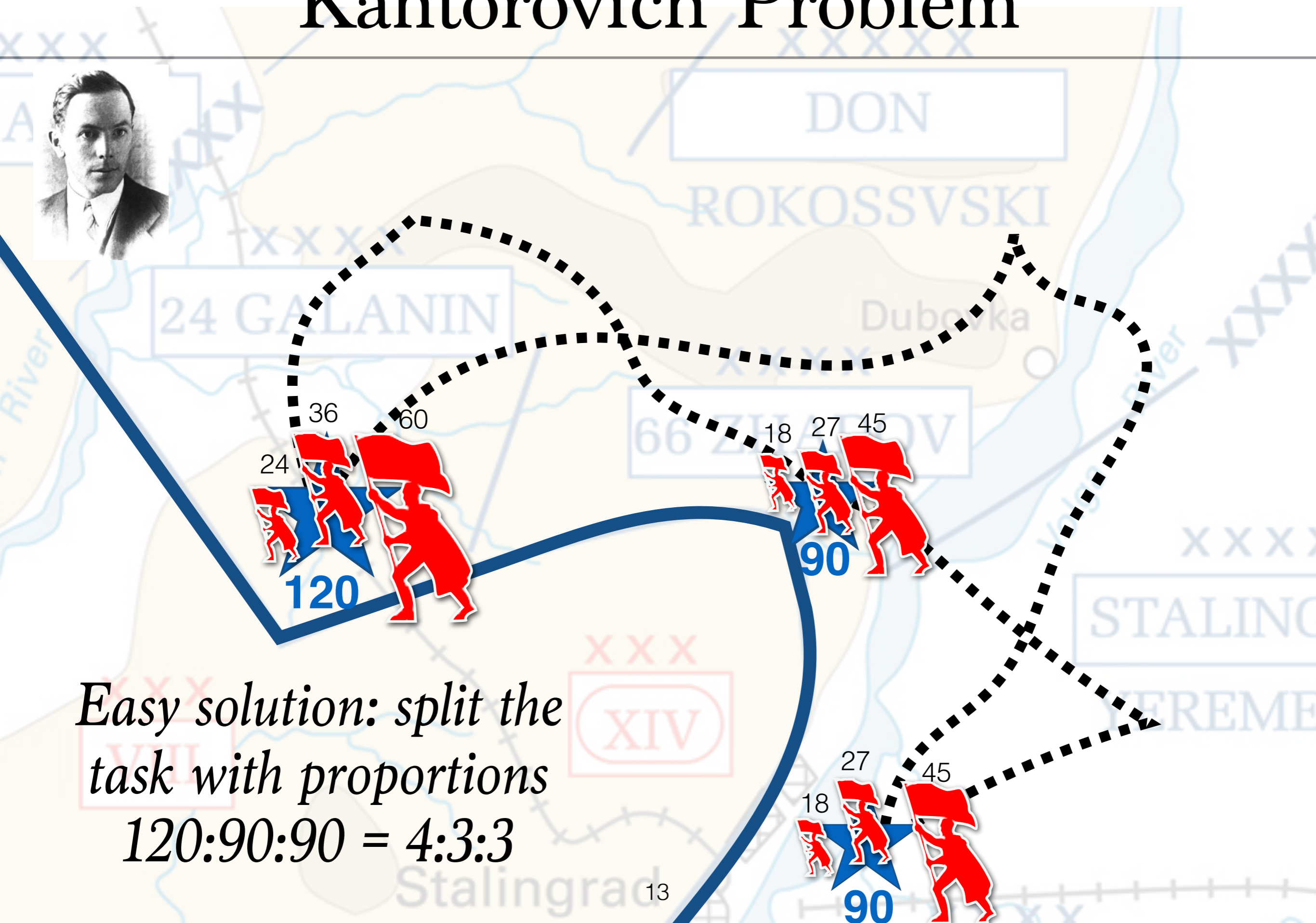
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem

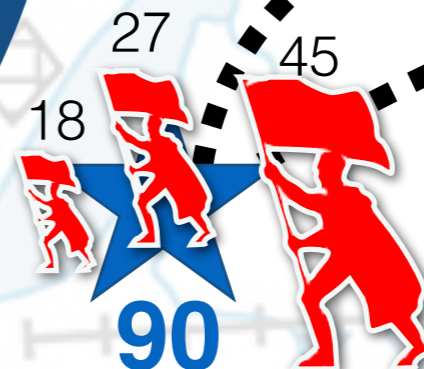


Naive approach results in too many displacements.

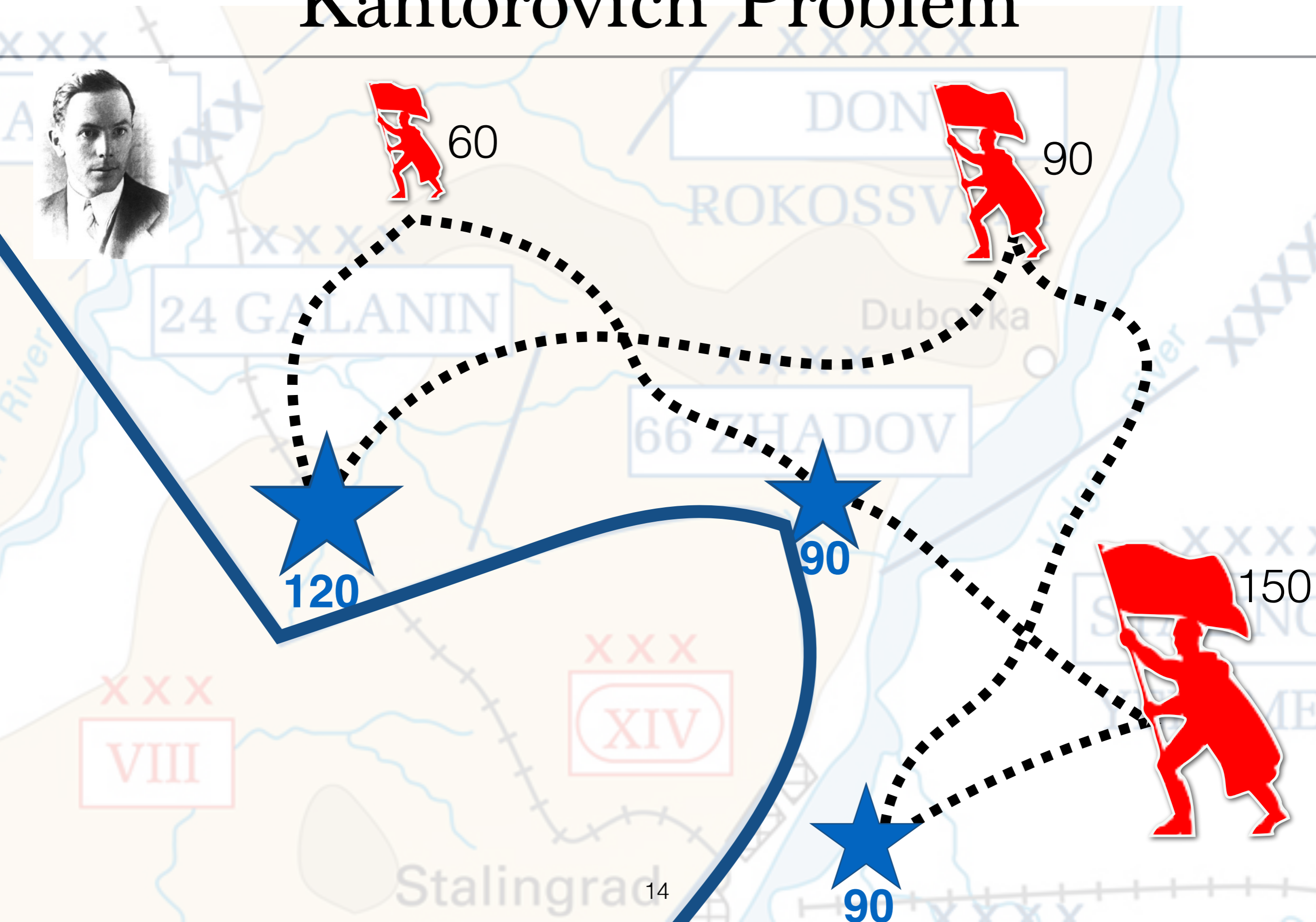
Goal: find a cheaper alternative

Easy solution: split the task with proportions

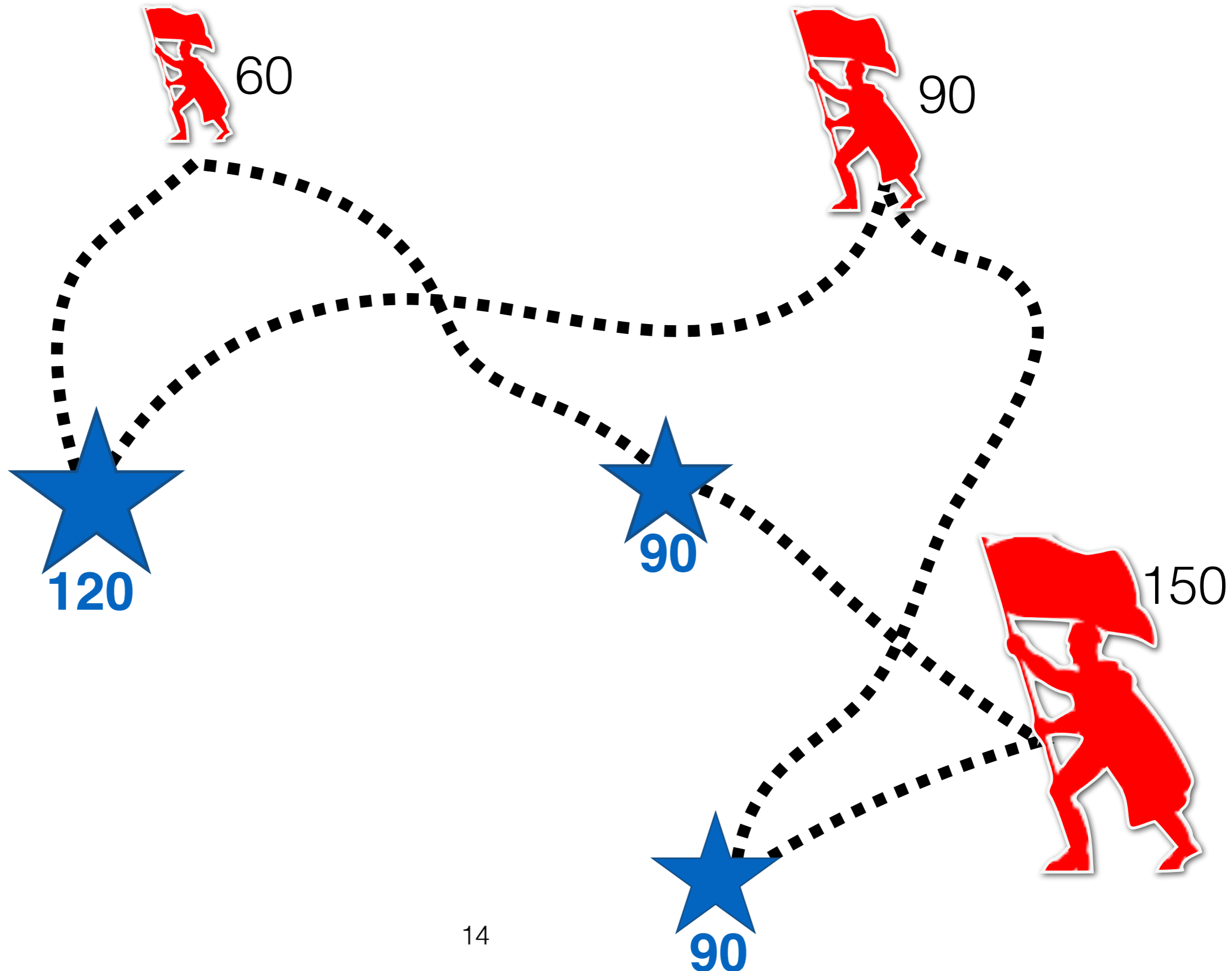
$$120:90:90 = 4:3:3$$



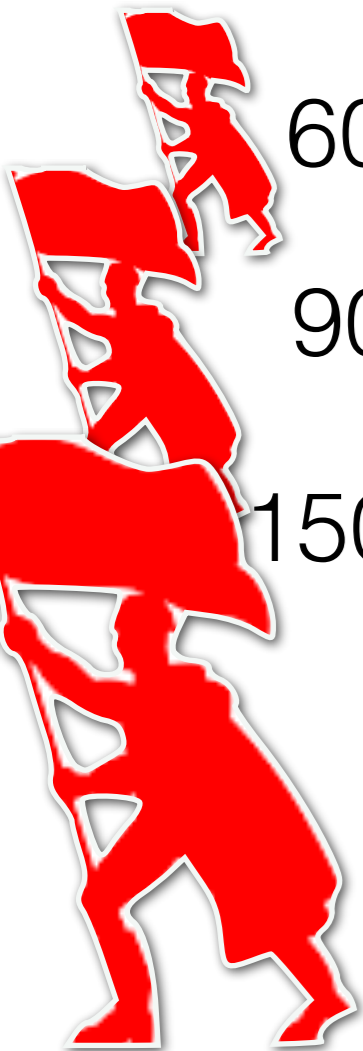
Kantorovich Problem



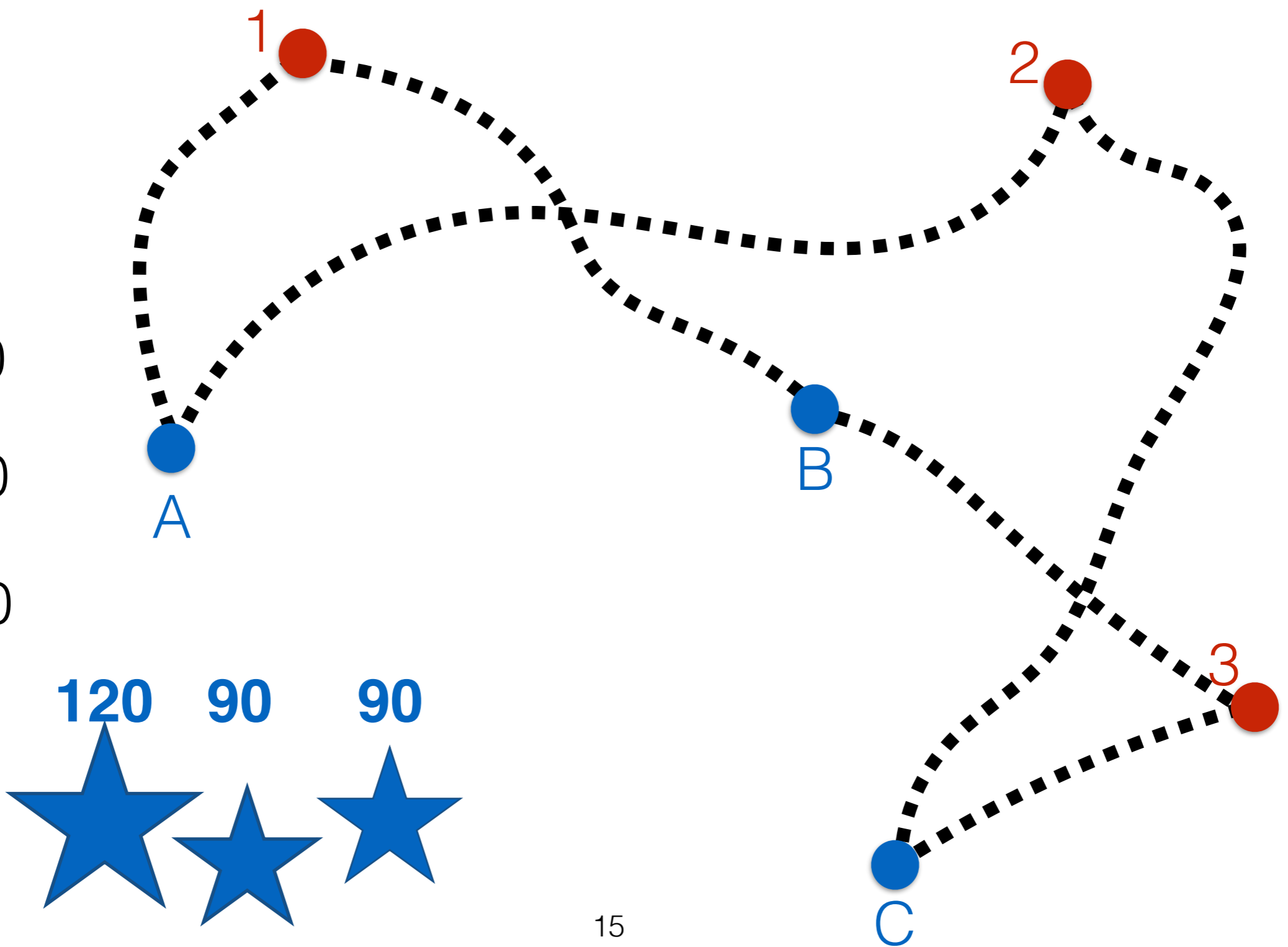
Kantorovich Problem



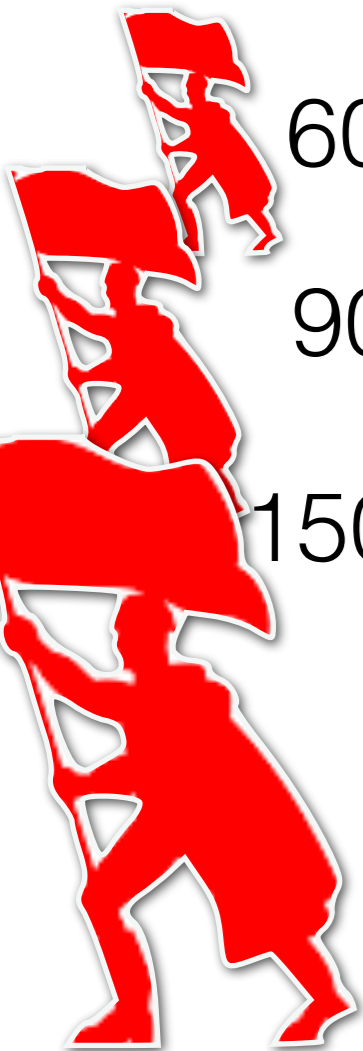
Kantorovich Problem



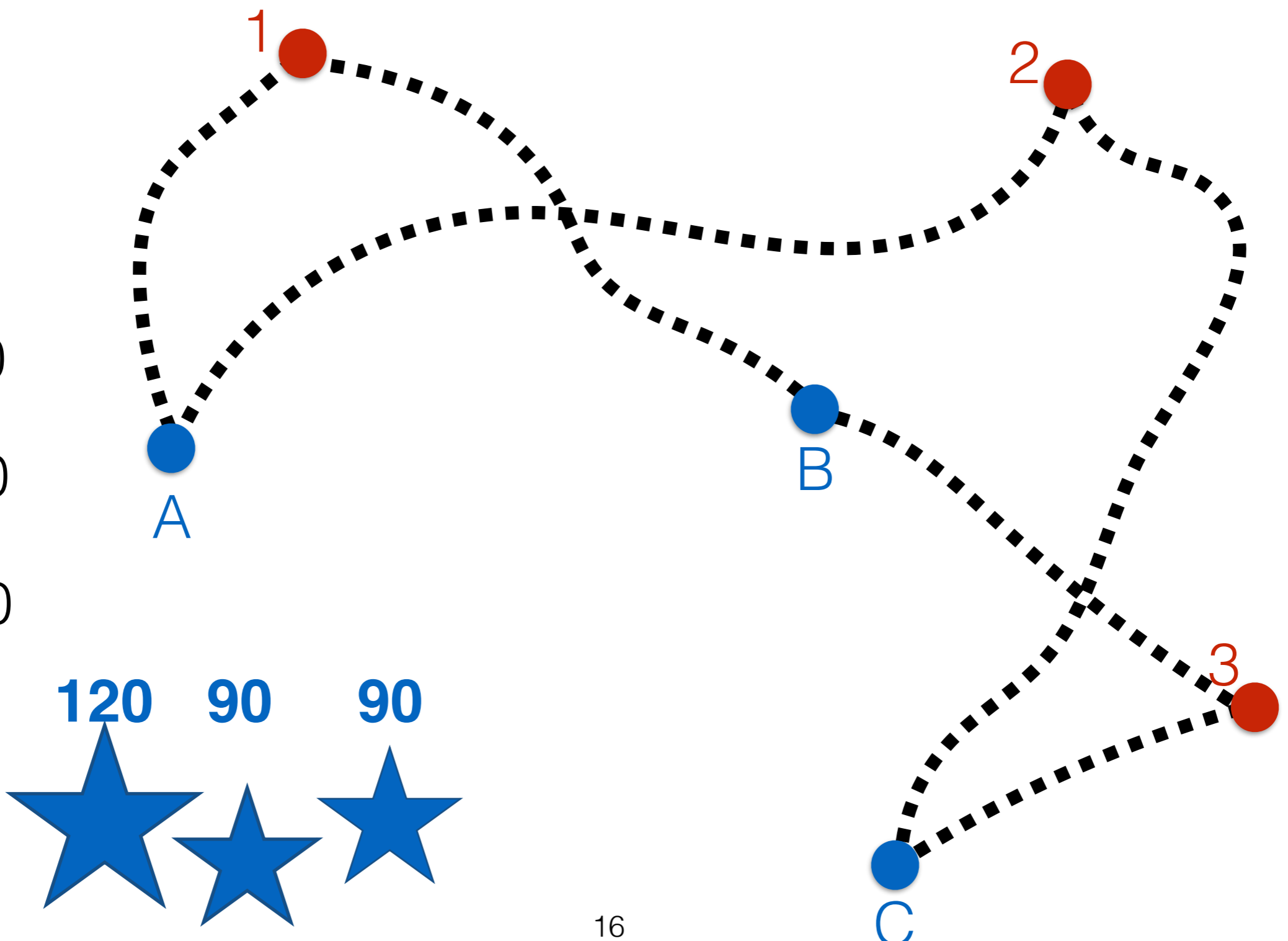
60
90
150



Kantorovich Problem



60
90
150



Kantorovich Problem



60	?	?	?
90	?	?	?
150	?	?	?

120 90 90



1 ●	d_{1A}	d_{1B}	d_{1C}
2 ●	d_{2A}	d_{2B}	d_{2C}
3 ●	d_{3A}	d_{3B}	d_{3C}
	● A	● B	● C

Kantorovich Problem



Transportation matrix

60	?	?	?
90	?	?	?
150	?	?	?

120 90 90



Distance matrix


1 ●	d_{1A}	d_{1B}	d_{1C}
2 ●	d_{2A}	d_{2B}	d_{2C}
3 ●	d_{3A}	d_{3B}	d_{3C}
	● A	● B	● C

Kantorovich Problem



The problem is entirely described by
counts and a cost/distance matrix

Transportation matrix



60	?	?	?
90	?	?	?
150	?	?	?
	120	90	90



Distance matrix

1 ●	d_{1A}	d_{1B}	d_{1C}
2 ●	d_{2A}	d_{2B}	d_{2C}
3 ●	d_{3A}	d_{3B}	d_{3C}
	● A	● B	● C

Kantorovich Problem

Transportation matrix

60	?	?	?
90	?	?	?
150	?	?	?
	120	90	90

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Kantorovich Problem

Transportation matrix

60	p_{1A}	p_{1B}	p_{1C}
90	p_{2A}	p_{2B}	p_{2C}
150	p_{3A}	p_{3B}	p_{3C}
	120	90	90

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Kantorovich Problem

Transportation matrix

a_1	p_{1A}	p_{1B}	p_{1C}
a_2	p_{2A}	p_{2B}	p_{2C}
a_3	p_{3A}	p_{3B}	p_{3C}
	b_A	b_B	b_C

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Kantorovich Problem

Transportation matrix

a_1	p_{1A}	p_{1B}	p_{1C}
a_2	p_{2A}	p_{2B}	p_{2C}
a_3	p_{3A}	p_{3B}	p_{3C}
	b_A	b_B	b_C

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Constraints

$$\forall i \in \{1, 2, 3\}, \quad \sum_{j \in \{A, B, C\}} p_{ij} = a_i$$

$$\forall j \in \{A, B, C\}, \quad \sum_{i \in \{1, 2, 3\}} p_{ij} = b_j$$

$$p_{ij} \geq 0$$

Kantorovich Problem

Transportation matrix

a_1	p_{1A}	p_{1B}	p_{1C}
a_2	p_{2A}	p_{2B}	p_{2C}
a_3	p_{3A}	p_{3B}	p_{3C}
	b_A	b_B	b_C

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Constraints

$$\forall i \in \{1, 2, 3\}, \quad \sum_{j \in \{A, B, C\}} p_{ij} = a_i$$

$$\forall j \in \{A, B, C\}, \quad \sum_{i \in \{1, 2, 3\}} p_{ij} = b_j$$

$$p_{ij} \geq 0$$

Cost function

$$C(P) = \sum_{j \in \{A, B, C\}} \sum_{i \in \{1, 2, 3\}} p_{ij} d_{ij}$$

Kantorovich Problem

Transportation matrix

a_1	p_{1A}	p_{1B}	p_{1C}
a_2	p_{2A}	p_{2B}	p_{2C}
a_3	p_{3A}	p_{3B}	p_{3C}
	b_A	b_B	b_C

Distance matrix

1	d_{1A}	d_{1B}	d_{1C}
2	d_{2A}	d_{2B}	d_{2C}
3	d_{3A}	d_{3B}	d_{3C}
	A	B	C

Constraints

$$\forall i \in \{1, 2, 3\}, \quad \sum_{j \in \{A, B, C\}} p_{ij} = a_i$$

$$\forall j \in \{A, B, C\}, \quad \sum_{i \in \{1, 2, 3\}} p_{ij} = b_j$$

$$p_{ij} \geq 0$$

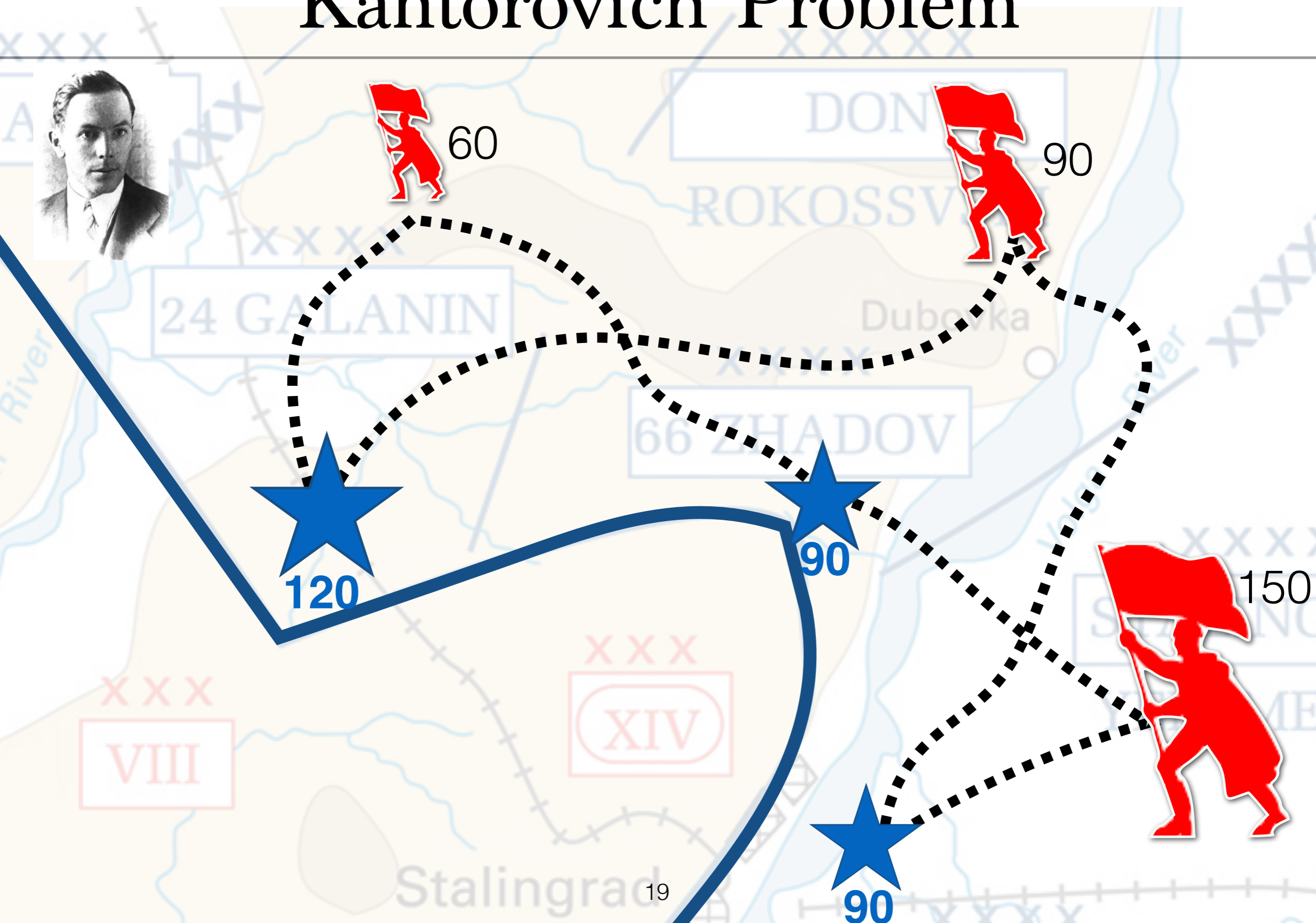
Cost function

$$C(P) = \sum_{j \in \{A, B, C\}} \sum_{i \in \{1, 2, 3\}} p_{ij} d_{ij}$$

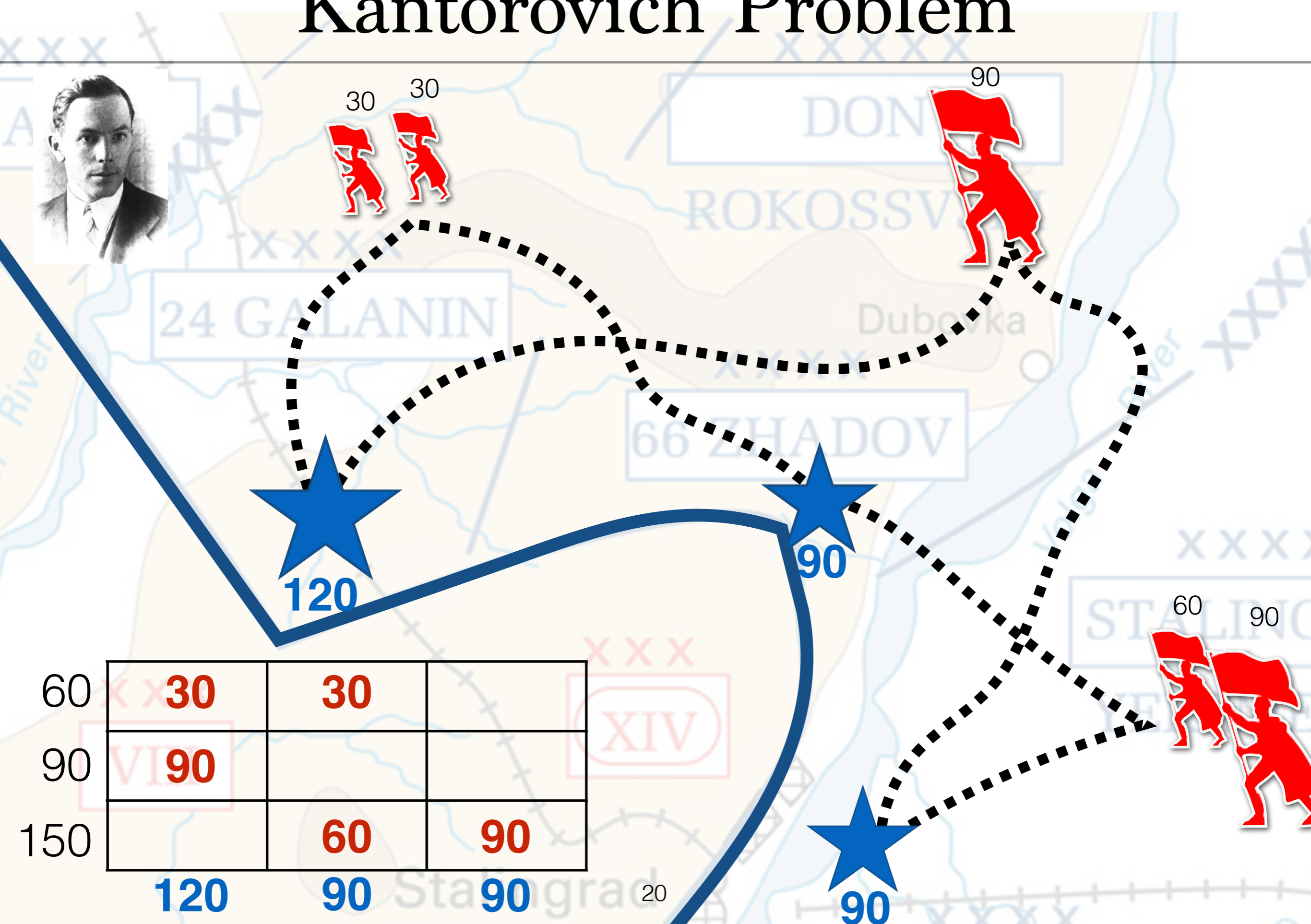
Problem

$$\min_{\text{all valid } P} C(P)$$

Kantorovich Problem

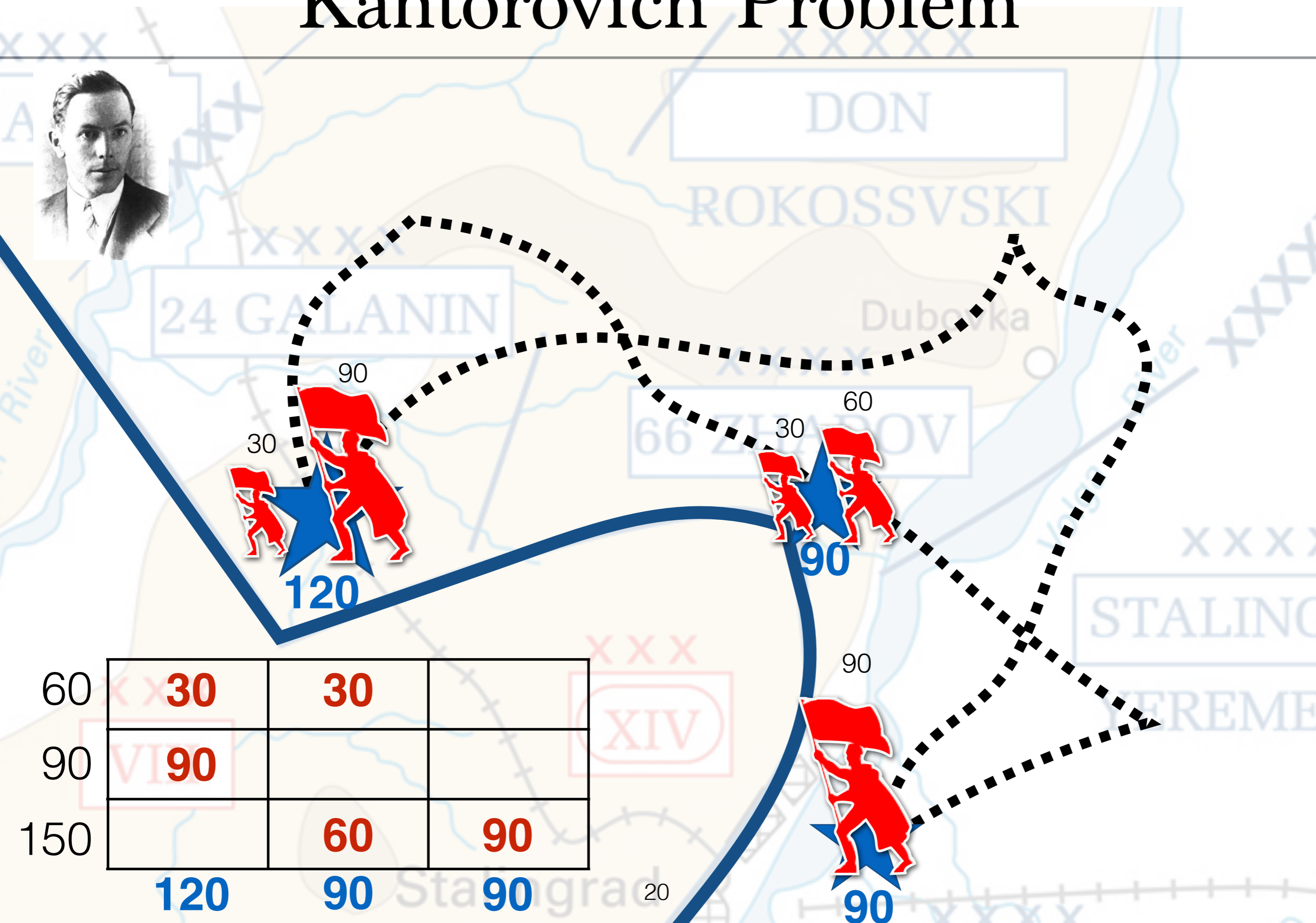


Kantorovich Problem



60	30	30	
90	90		
150	60	90	
	120	90	90

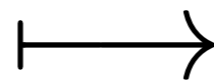
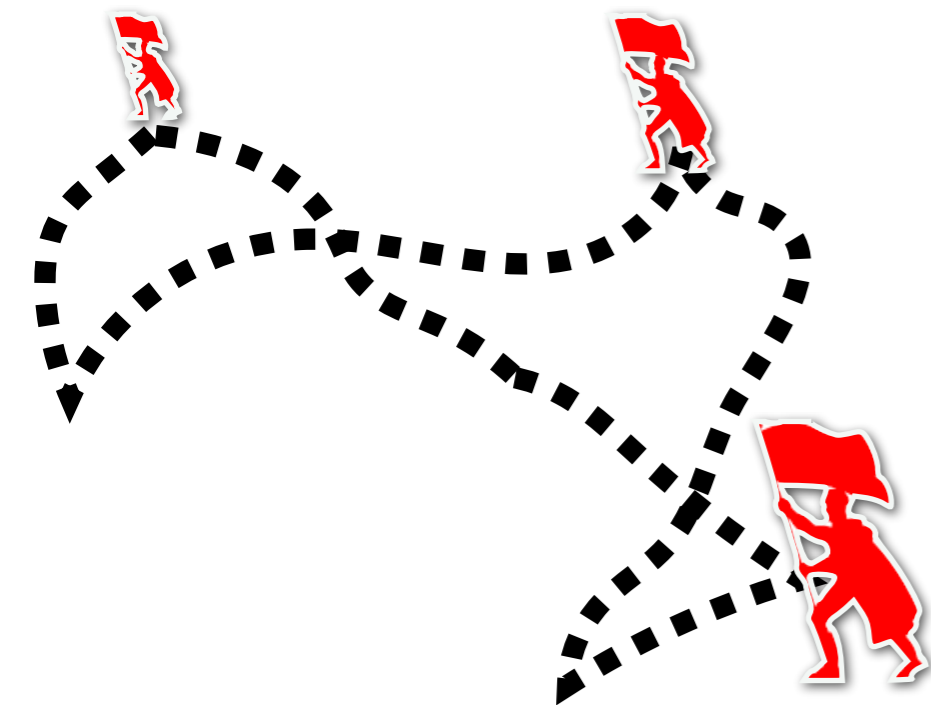
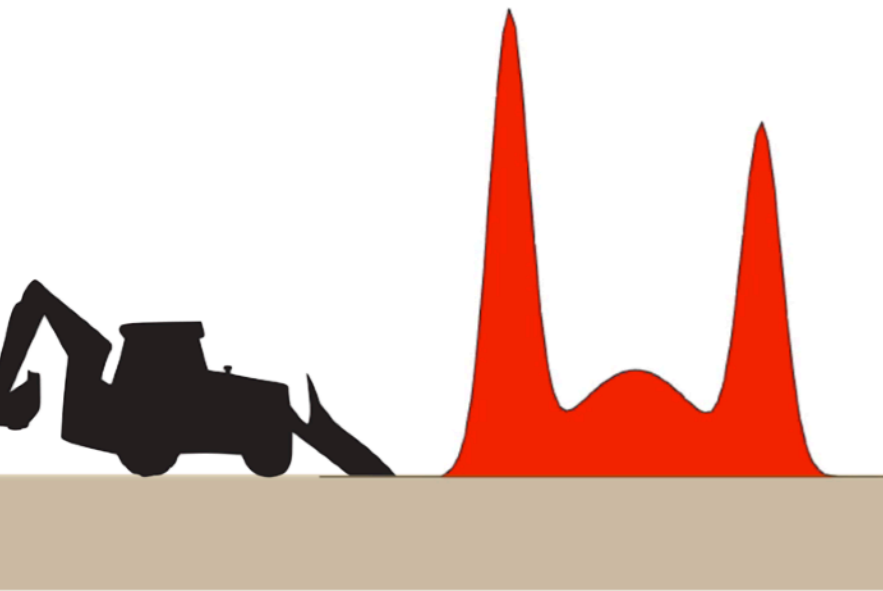
Kantorovich Problem



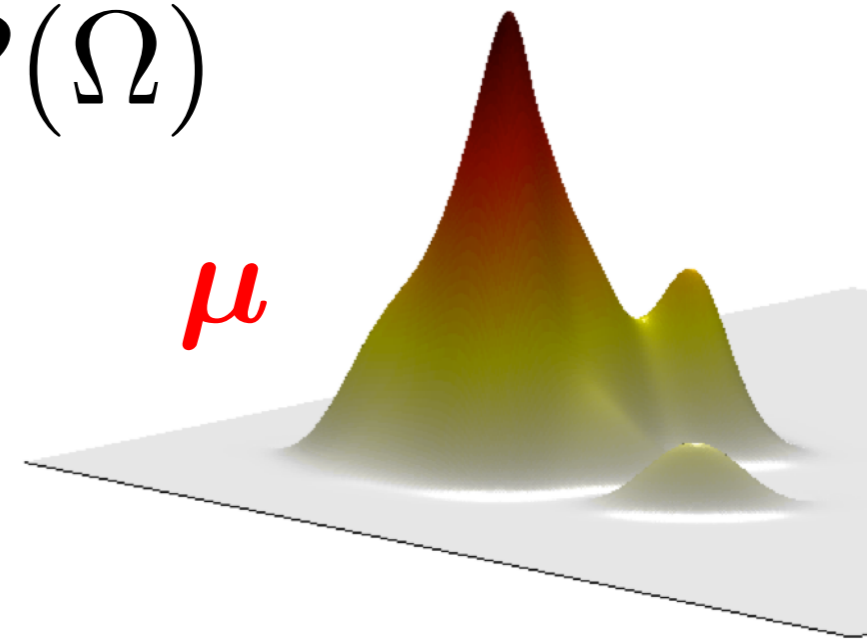
60	30	30	
90	90		
150		60	90
	120	90	90

Mathematical Formalism

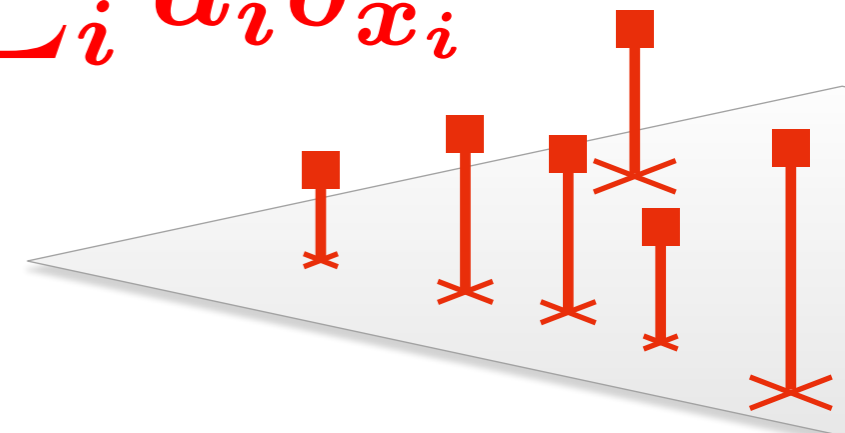
These problems involve discrete and continuous **probability measures** on a geometric space



$$\mathcal{P}(\Omega)$$



$$\sum_i a_i \delta_{x_i}$$

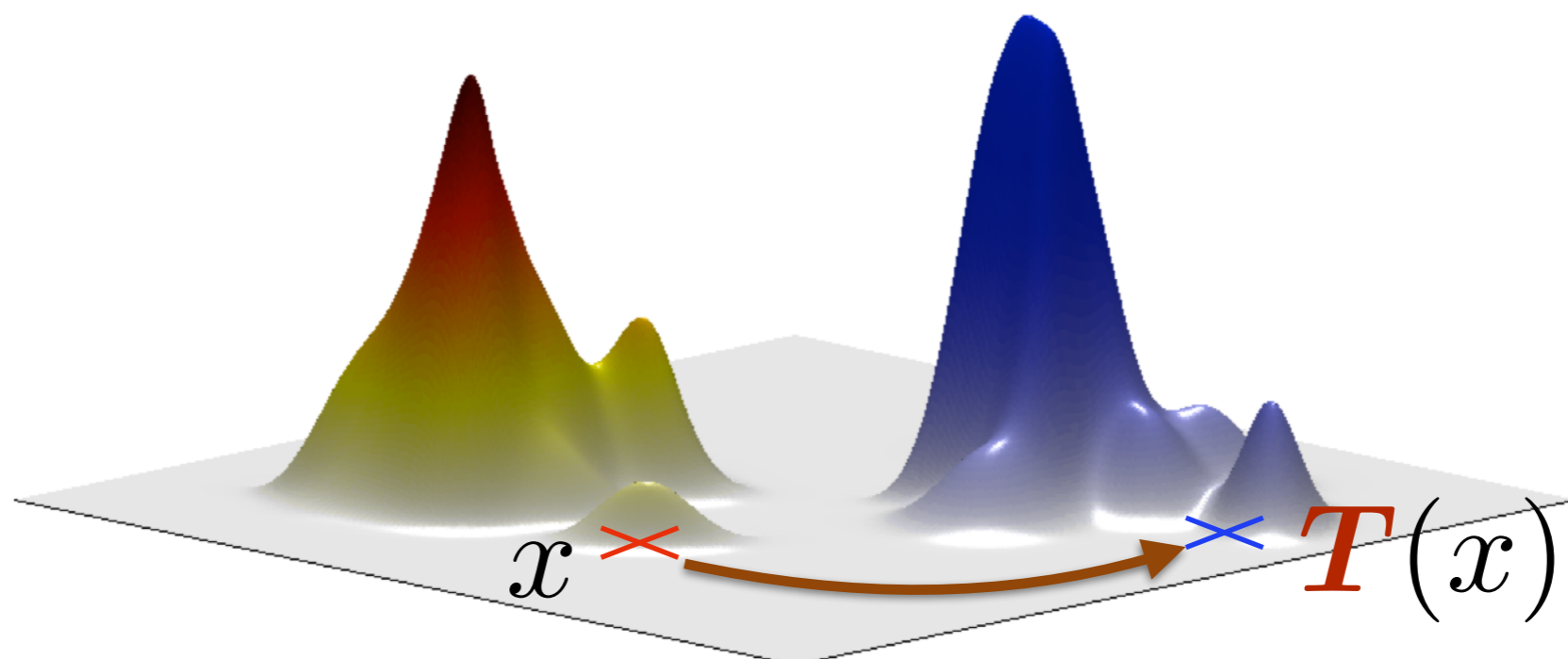


Monge Problem

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

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

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

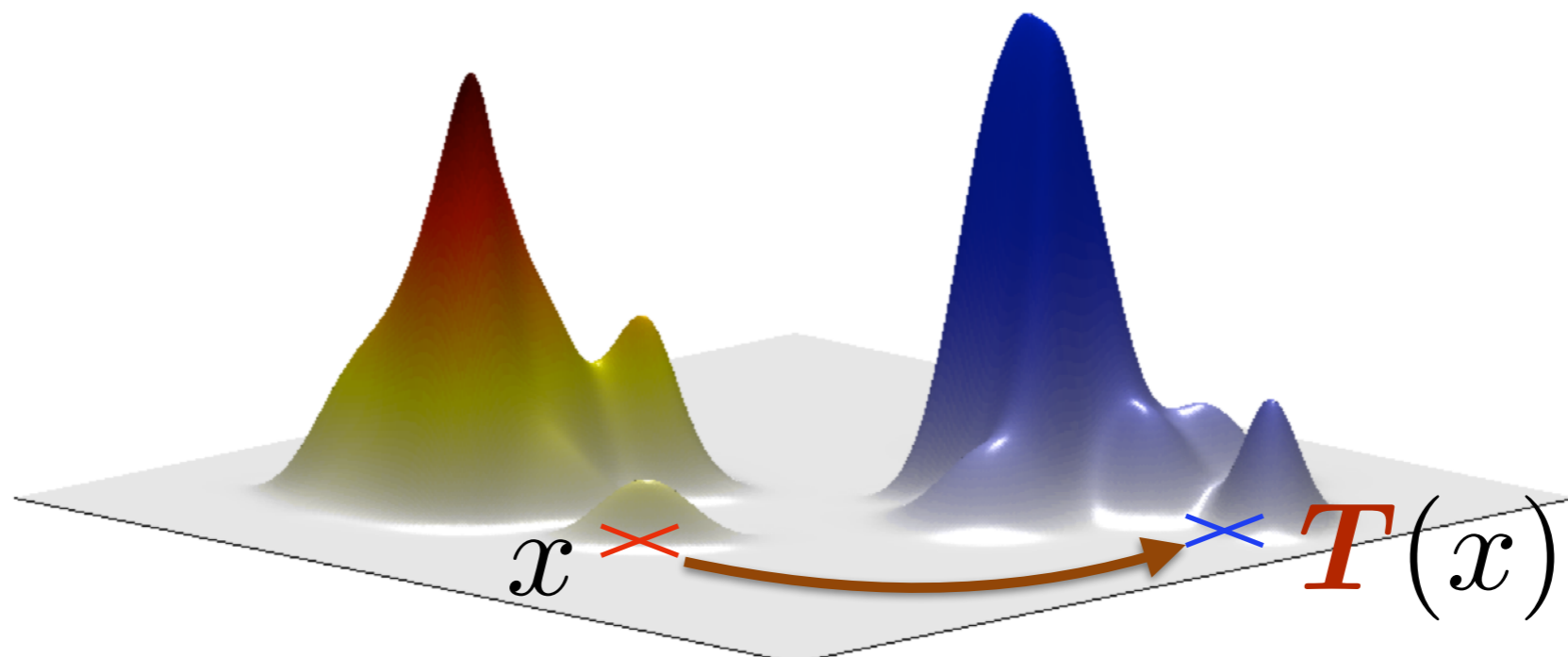


Monge Problem

Ω a probability space, $c : \Omega \times \Omega \rightarrow \mathbb{R}$.
 μ, ν two probability measures in $\mathcal{P}(\Omega)$.

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

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

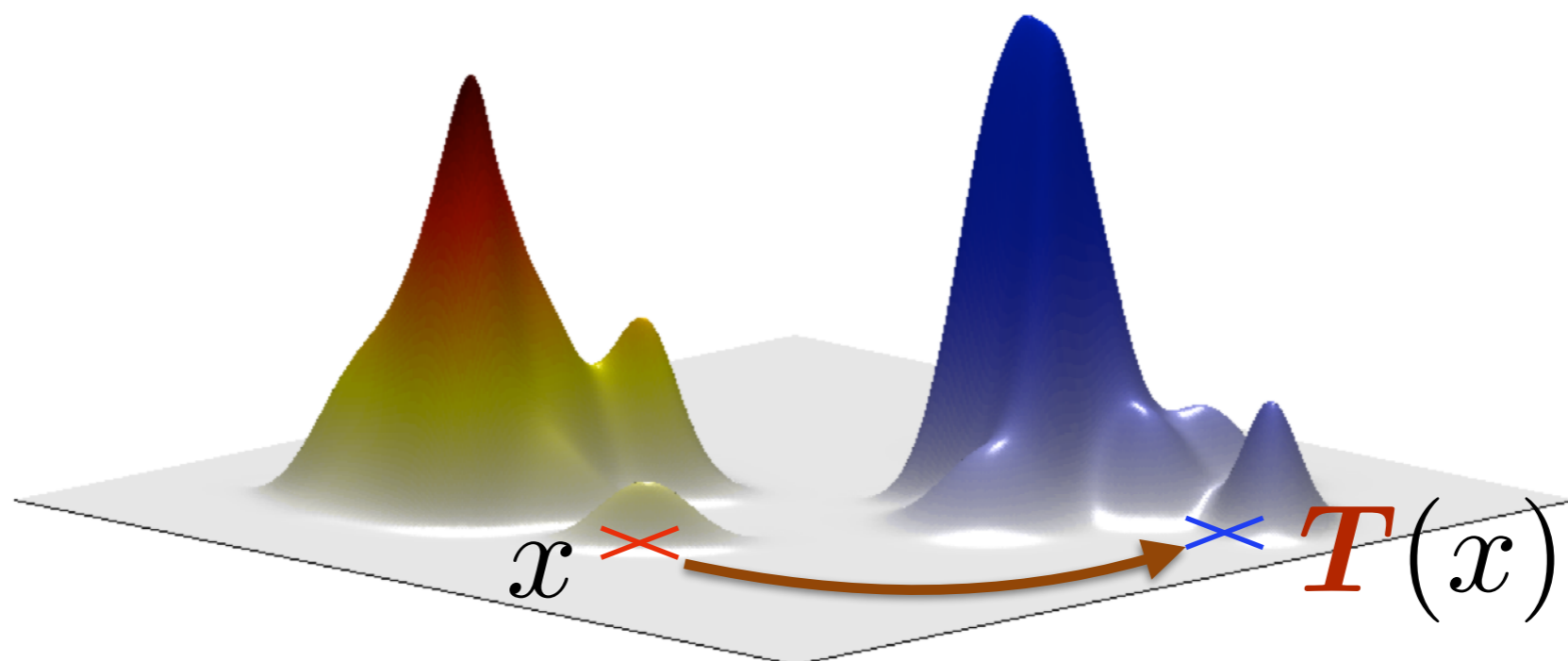


Monge Problem

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

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

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

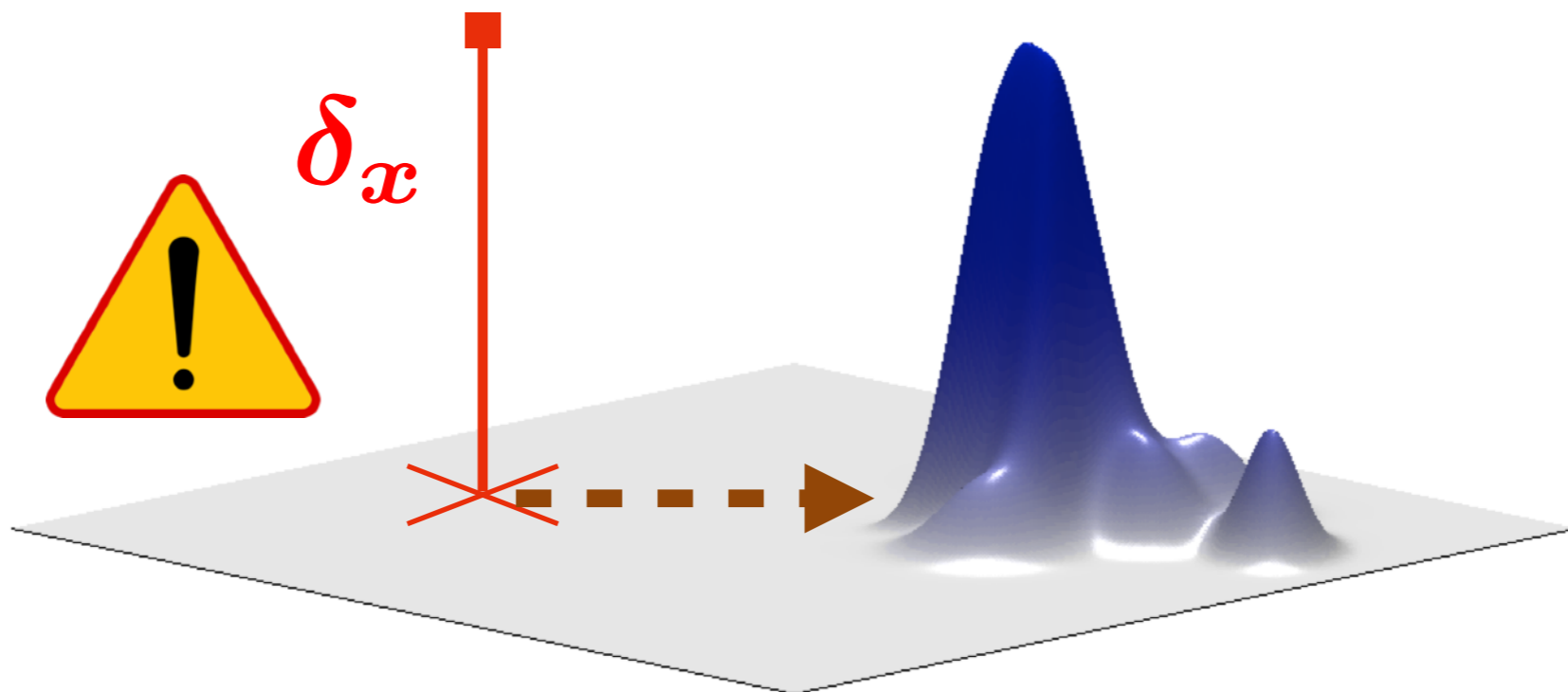


Monge Problem

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

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

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



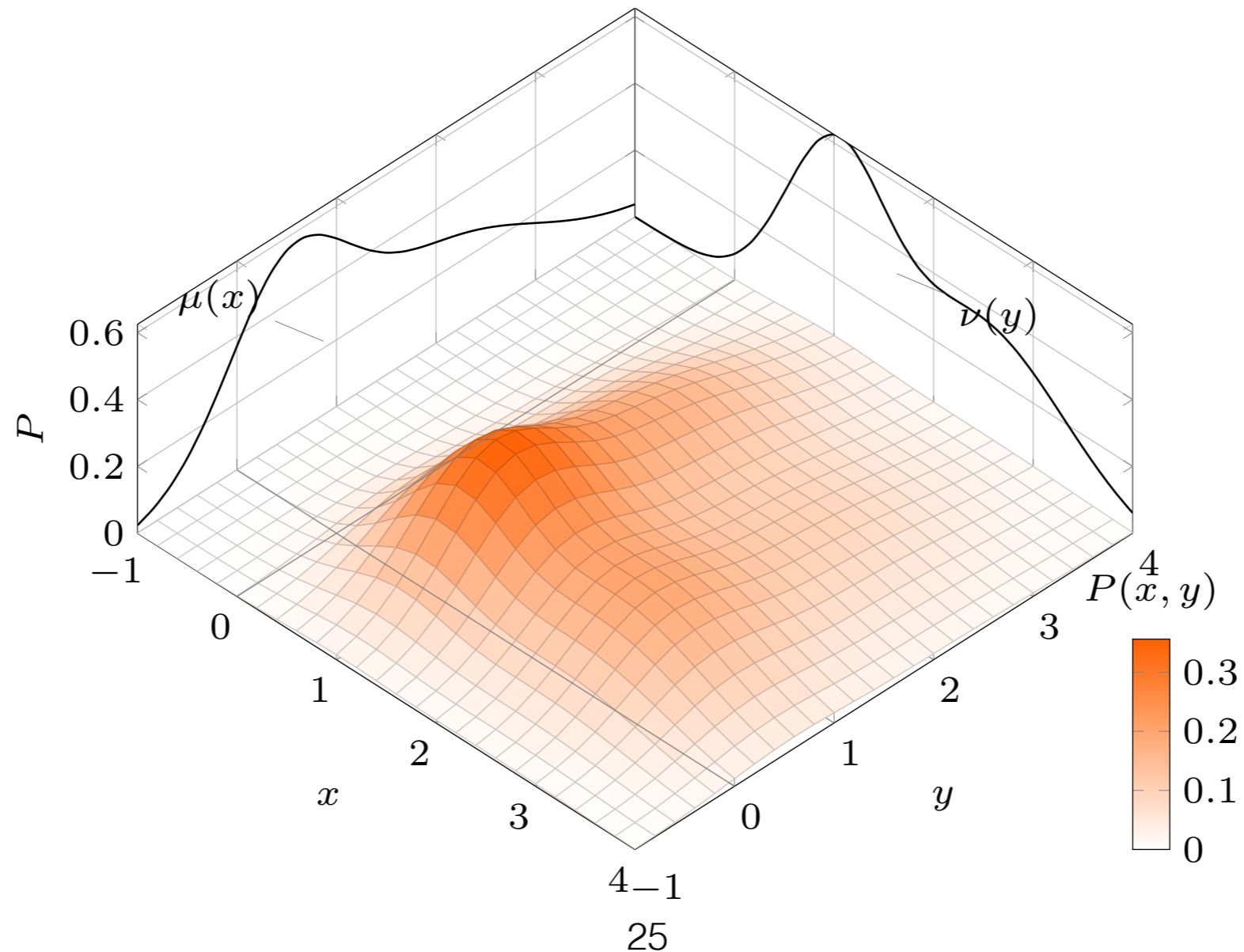
Kantorovich Relaxation

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

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

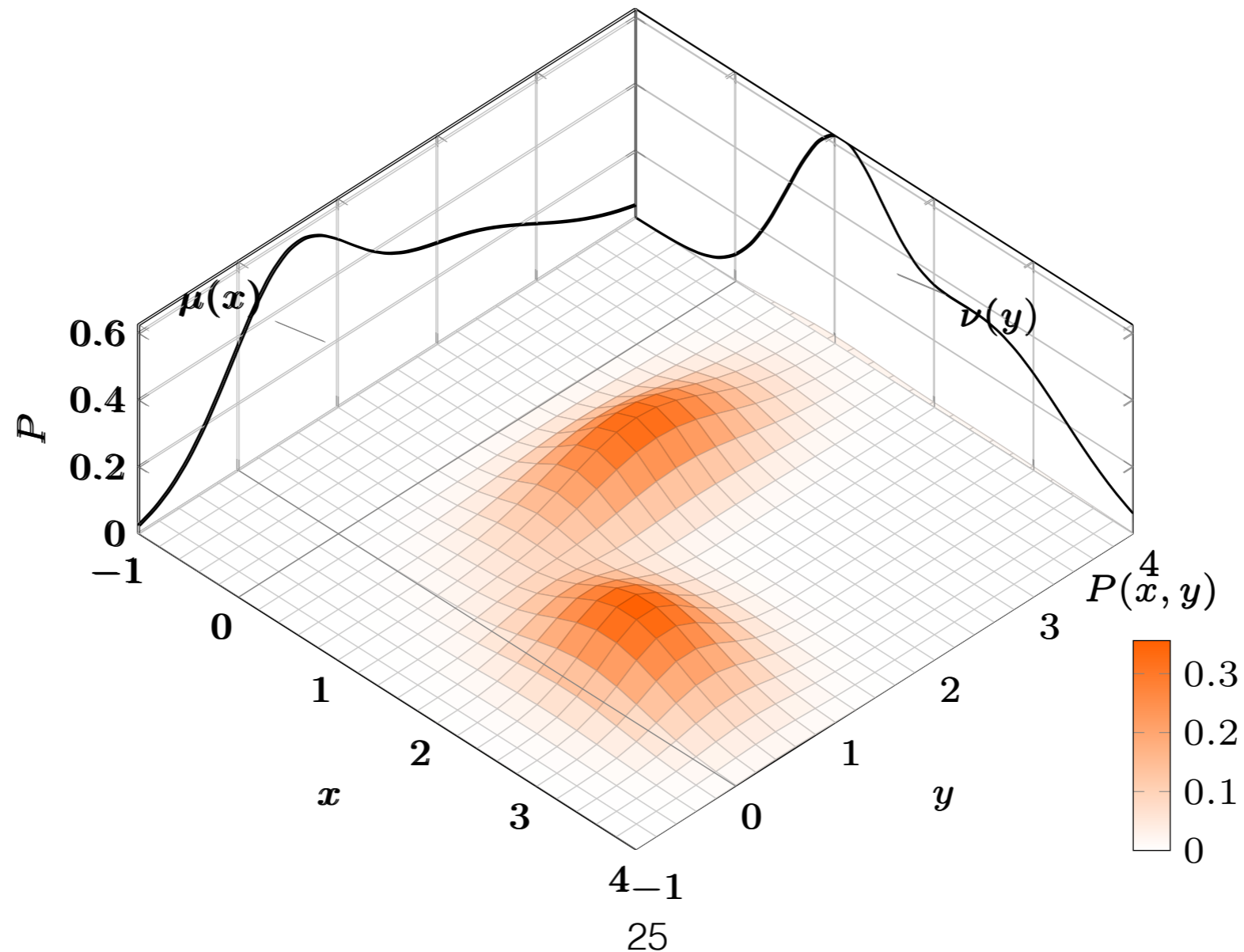
Kantorovich Relaxation

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



Kantorovich 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 Problem

Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function c on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{P \in \Pi(\mu, \nu)} \iint c(x, y) P(dx, dy).$$

PRIMAL

Kantorovich Problem

Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function c on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{P \in \Pi(\mu, \nu)} \iint c(x, y) P(dx, dy).$$

PRIMAL

$$\inf_{P \in \Pi(\mu, \nu)} \mathbb{E}_P [c(X, Y)]$$

Kantorovich Problem

Def. Given μ, ν in $\mathcal{P}(\Omega)$; a cost function c on $\Omega \times \Omega$, the Kantorovich problem is

$$\inf_{P \in \Pi(\mu, \nu)} \iint c(x, y) P(dx, dy).$$

PRIMAL

$$\sup_{\substack{\varphi \in L_1(\mu), \psi \in L_1(\nu) \\ \varphi(x) + \psi(y) \leq c(x, y)}} \int \varphi d\mu + \int \psi d\nu.$$

DUAL

Links between Monge & Kantorovich

Prop. For “well behaved” costs c , if μ has a density then an *optimal* Monge map T^* between μ and ν must exist.

Prop. In that case

$$P^* := (\text{Id}, T^*)\# \mu \in \Pi(\mu, \nu)$$

is also *optimal* for the Kantorovich problem.

[Brenier'91] [Smith&Knott'87] [McCann'01]

(Kantorovich) Wasserstein Distances

Let $p \geq 1$.

Let $c := D$, a metric.

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

$$W_p(\mu, \nu) \stackrel{\text{def}}{=} \left(\inf_{P \in \Pi(\mu, \nu)} \iint D(x, y)^p P(dx, dy) \right)^{1/p}.$$

(Kantorovich) Wasserstein Distances

Let $p \geq 1$.

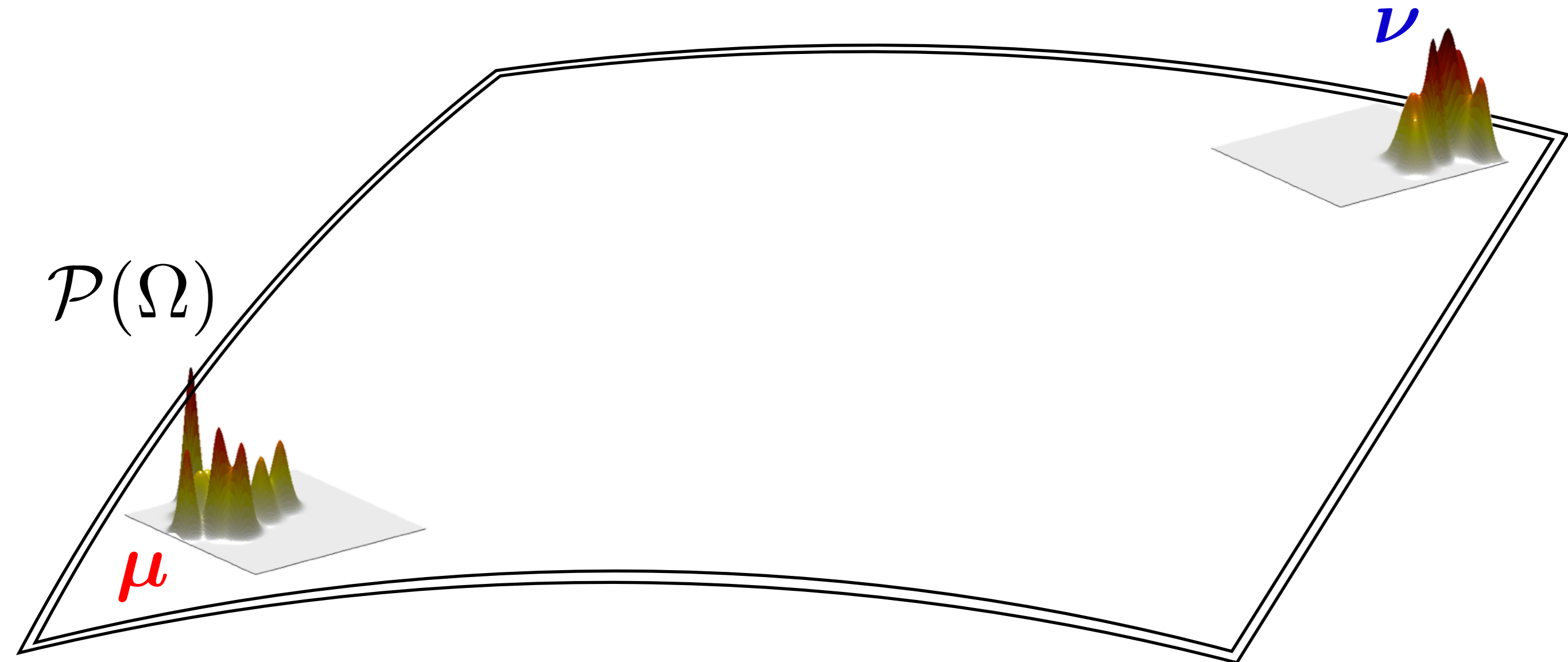
Let $c := D$, a metric.

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

$$W_p(\mu, \nu) \stackrel{\text{def}}{=} \left(\inf_{P \in \Pi(\mu, \nu)} \iint D(x, y)^p P(dx, dy) \right)^{1/p}.$$

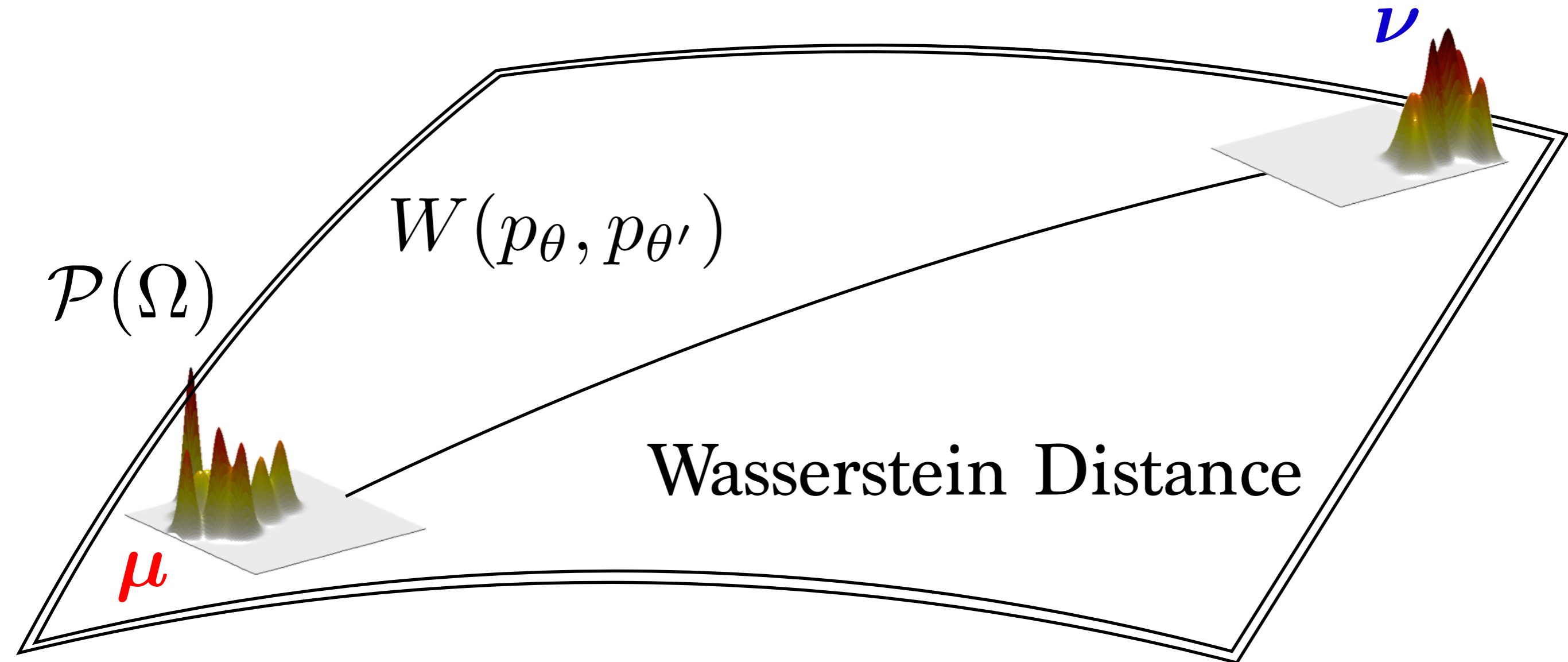
Optimal Transport Geometry

Very different geometry than standard information divergences (*KL*, Euclidean)



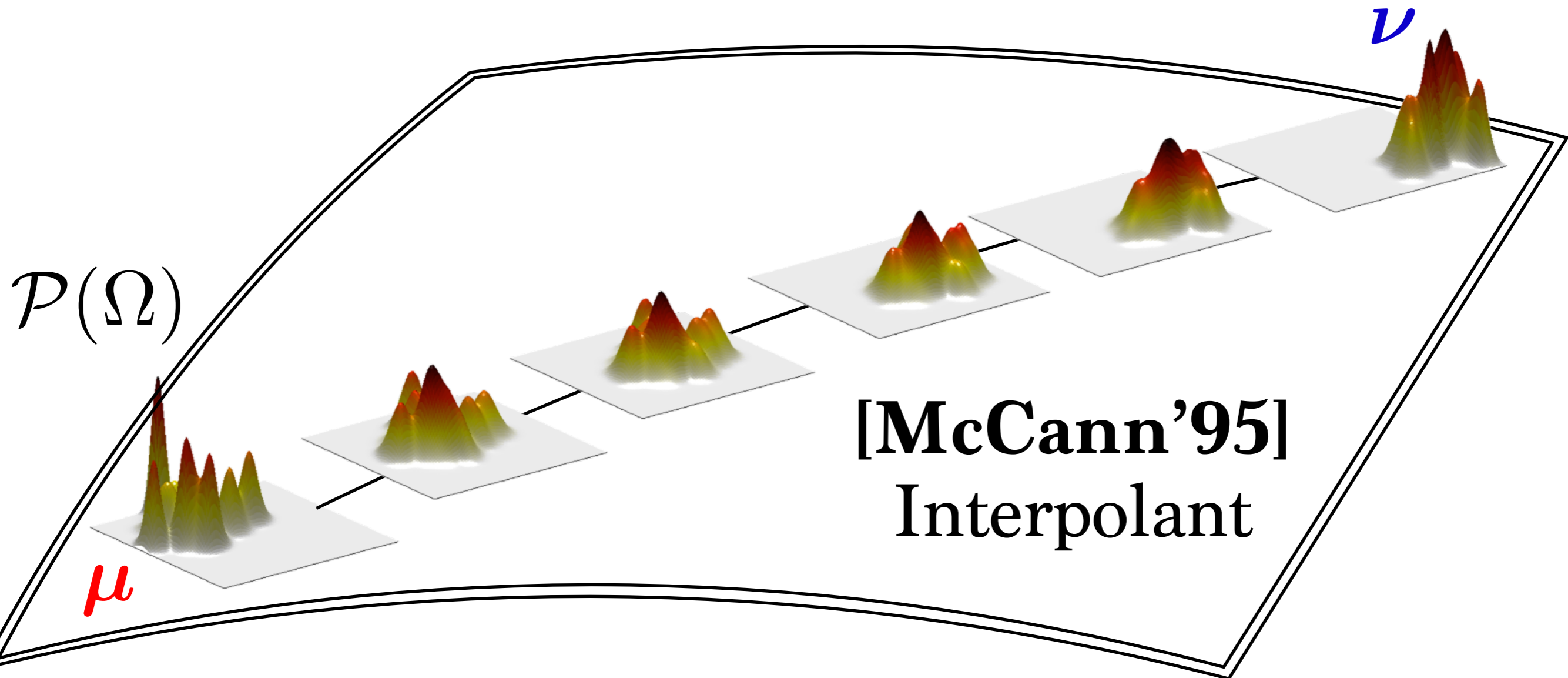
Optimal Transport Geometry

Very different geometry than standard information divergences (*KL*, Euclidean)



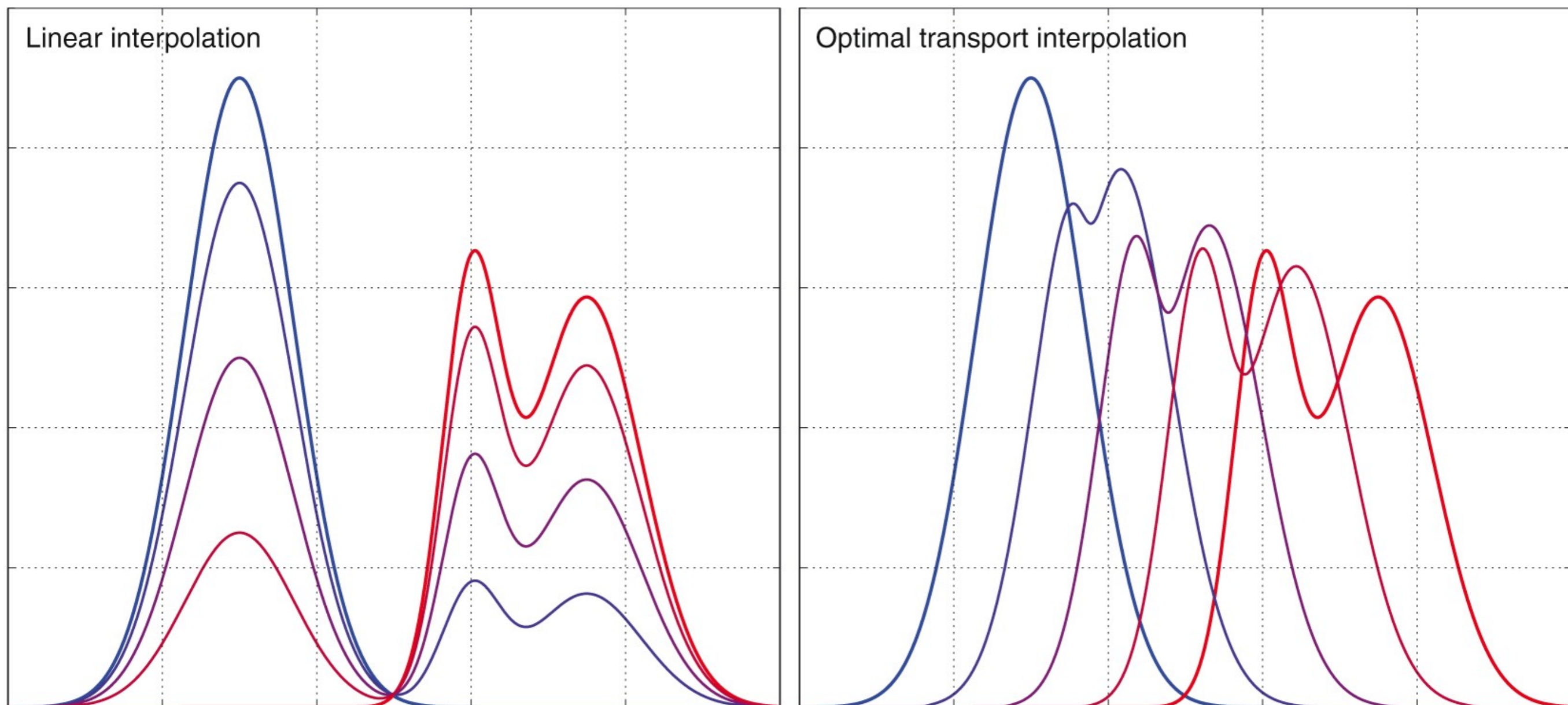
Optimal Transport Geometry

Very different geometry than standard information divergences (*KL*, Euclidean)



Optimal Transport Geometry

Very different geometry than standard information divergences (KL , Euclidean)



Optimal Transport Geometry

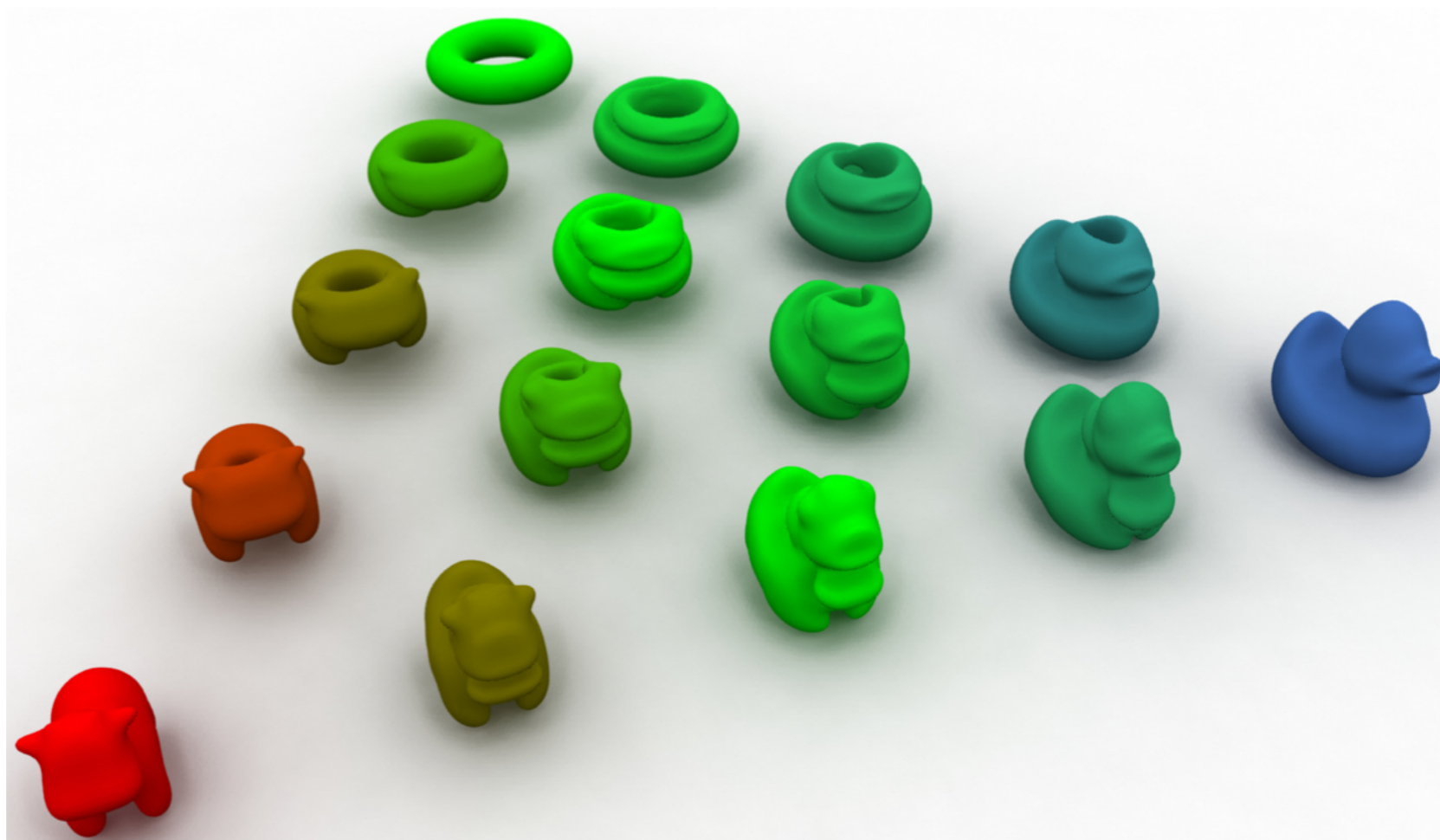
Very different geometry than standard information divergences (KL , Euclidean)



[Solomon'15]

Optimal Transport Geometry

Very different geometry than standard information divergences (KL , Euclidean)



[Solomon'15]

Computational OT

Up to 2010: OT solvers $W_p(\mu, \nu) = ?$

Goal now: use OT as a **loss or fidelity** term

$\operatorname{argmin}_{\mu \in \mathcal{P}(\Omega)} F(W_p(\mu, \nu_1), W_p(\mu, \nu_2), \dots, \mu) = ?$

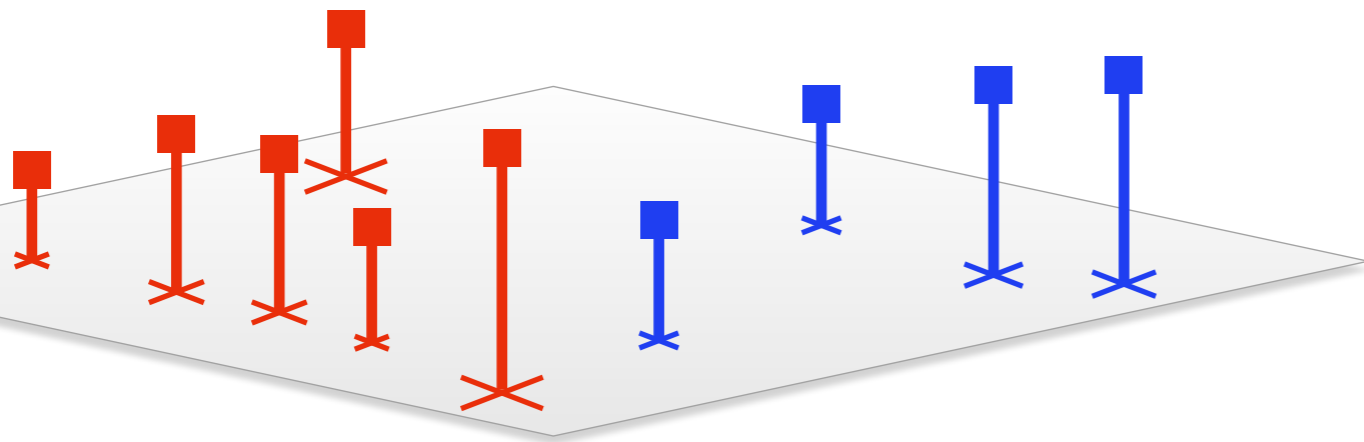
$\nabla_{\mu} W_p(\mu, \nu_1) = ?$

2. How to compute OT

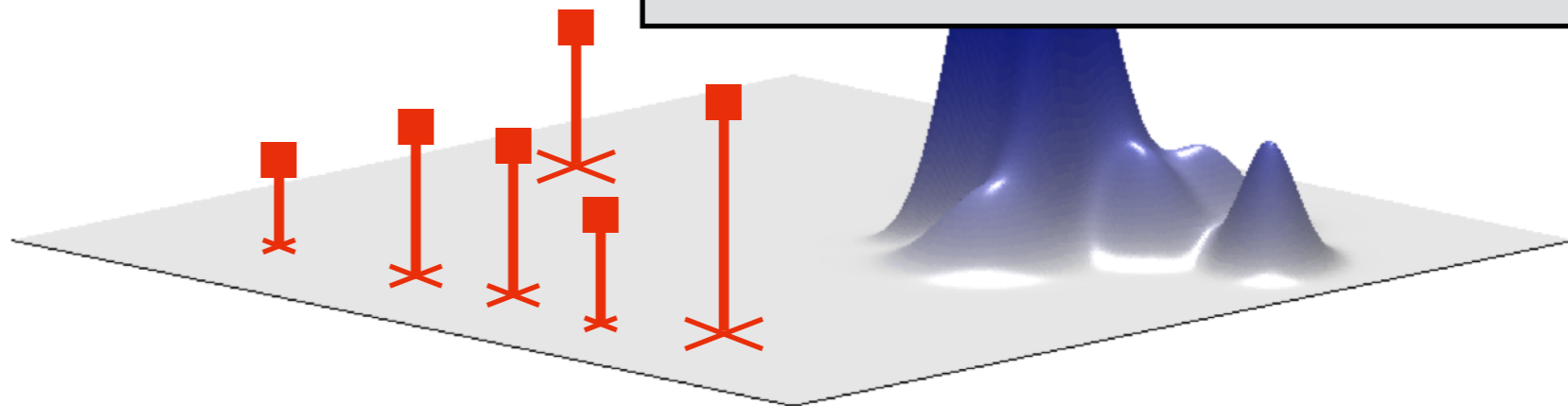
- Typology: discrete/continuous problems
- Easy cases, zoo of solvers
- Entropic regularization
- Differentiability of the W distance

How can we compute OT?

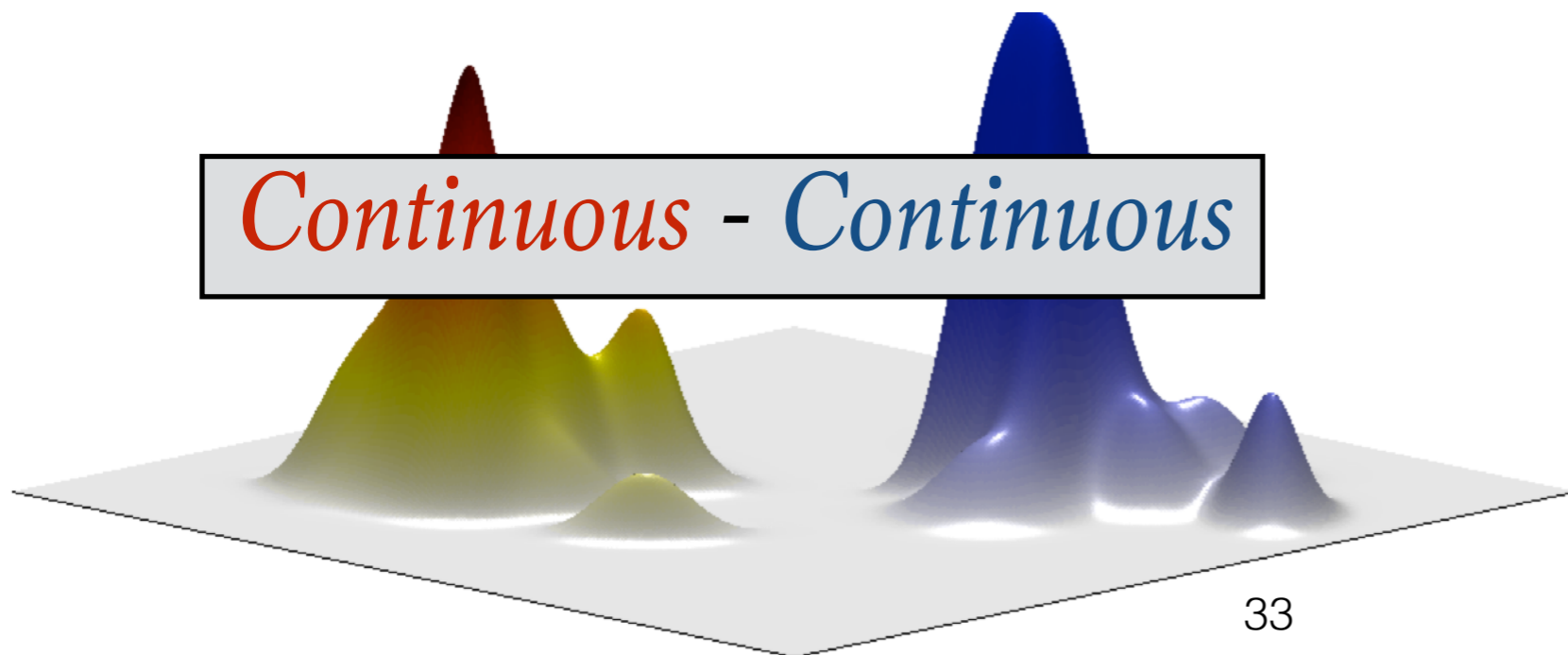
Discrete - Discrete



Discrete - Continuous



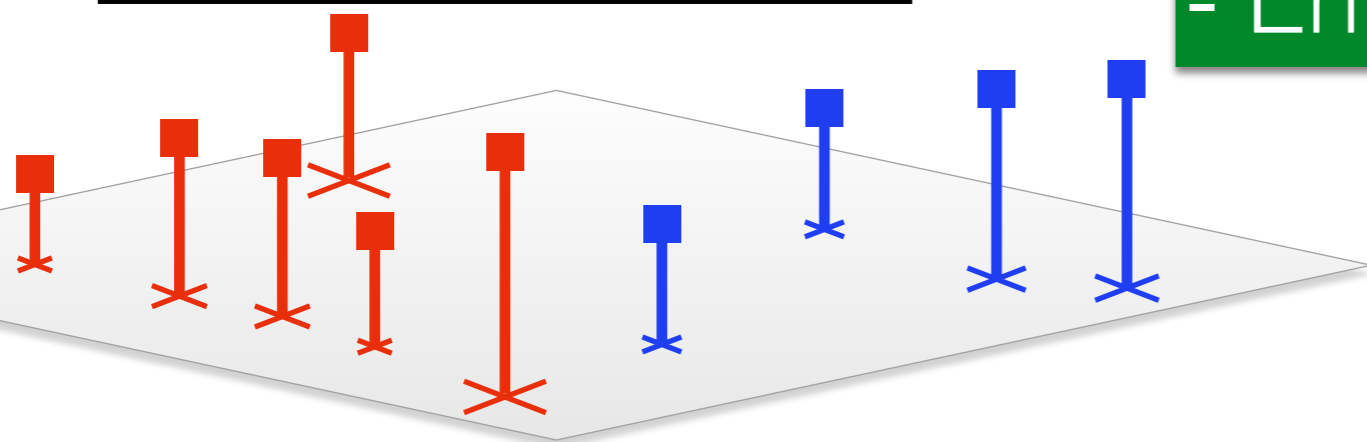
Continuous - Continuous



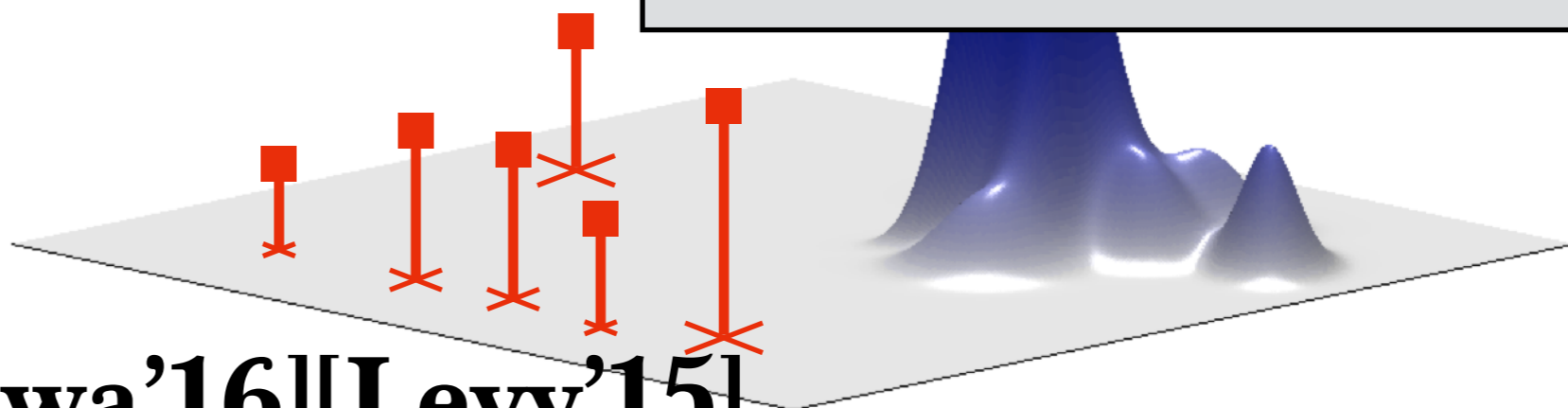
How can we compute OT?

Discrete - Discrete

- Network flow solvers
- Entropic regularization



Discrete - Continuous



low dim.

[Mérigot'11][Kitagawa'16][Levy'15]

Continuous - Continuous

Stochastic
Optimization

[Genevay'16]

PDE's

[Benamou'98]

Easy (1): Univariate Measures

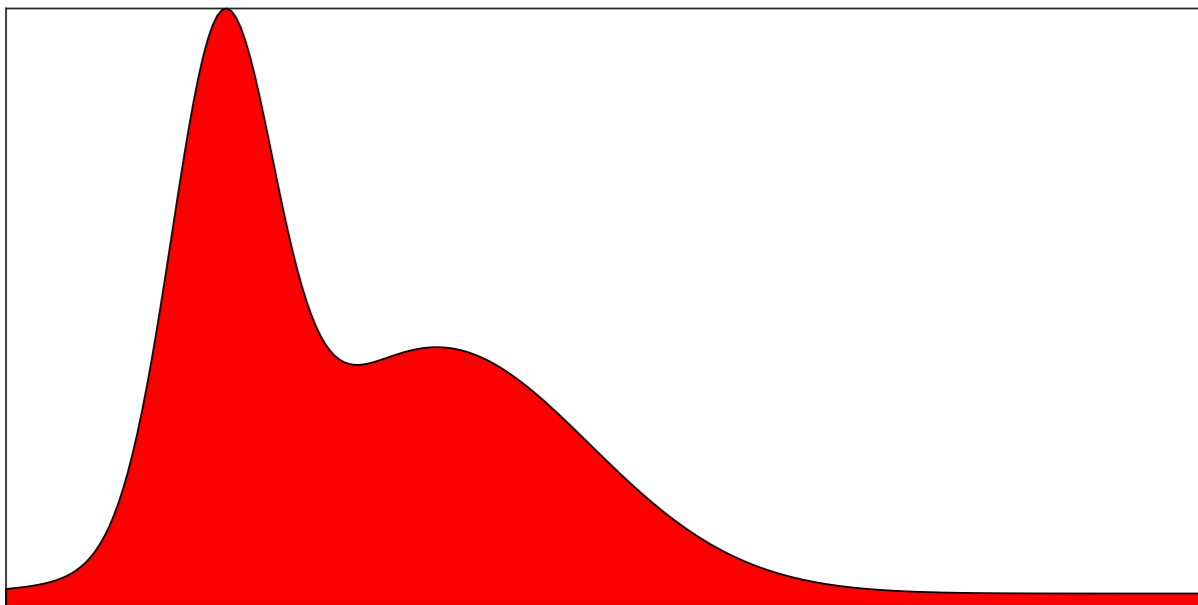
Remark. If $\Omega = \mathbb{R}$, $\mathbf{c}(x, y) = \mathbf{c}(|x - y|)$, \mathbf{c} convex, $F_{\mu}^{-1}, F_{\nu}^{-1}$ quantile functions,

$$W(\mu, \nu) = \int_0^1 \mathbf{c}(|F_{\mu}^{-1}(x) - F_{\nu}^{-1}(x)|) dx$$

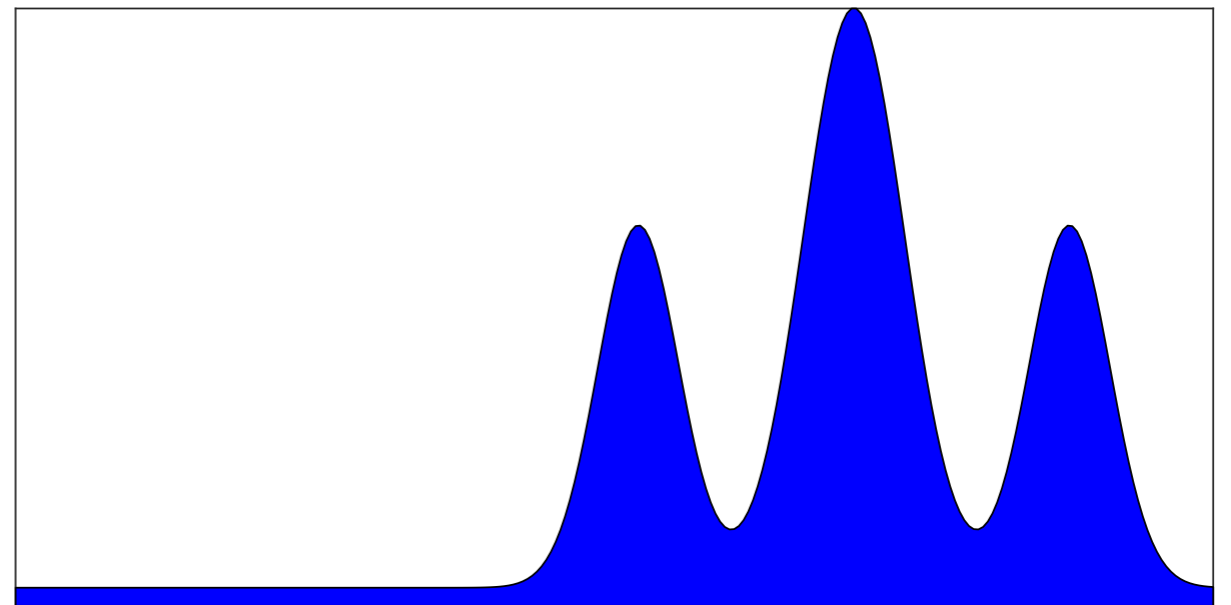
Easy (1): Univariate Measures

Remark. If $\Omega = \mathbb{R}$, $\mathbf{c}(x, y) = \mathbf{c}(|x - y|)$, \mathbf{c} convex, $F_{\mu}^{-1}, F_{\nu}^{-1}$ quantile functions,

$$W(\mu, \nu) = \int_0^1 \mathbf{c}(|F_{\mu}^{-1}(x) - F_{\nu}^{-1}(x)|) dx$$



μ

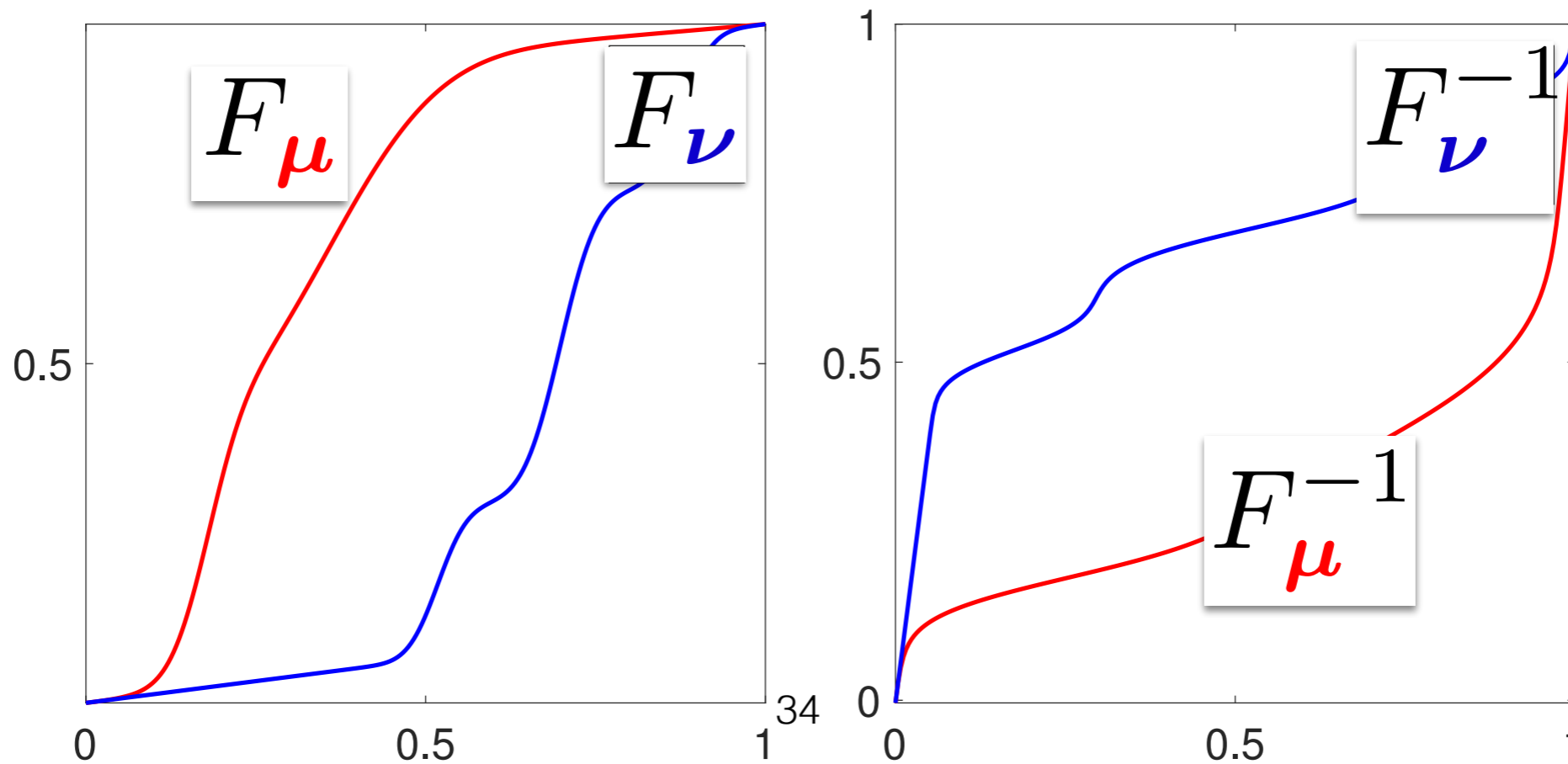


ν

Easy (1): Univariate Measures

Remark. If $\Omega = \mathbb{R}$, $\mathbf{c}(x, y) = \mathbf{c}(|x - y|)$, \mathbf{c} convex, $F_{\mu}^{-1}, F_{\nu}^{-1}$ quantile functions,

$$W(\mu, \nu) = \int_0^1 \mathbf{c}(|F_{\mu}^{-1}(x) - F_{\nu}^{-1}(x)|) dx$$



Easy (2): Gaussian Measures

Remark. If $\Omega = \mathbb{R}^d$, $\mathbf{c}(x, y) = \|x - y\|^2$, and $\mu = \mathcal{N}(\mathbf{m}_\mu, \Sigma_\mu)$, $\nu = \mathcal{N}(\mathbf{m}_\nu, \Sigma_\nu)$ then

$$W_2^2(\mu, \nu) = \|\mathbf{m}_\mu - \mathbf{m}_\nu\|^2 + B(\Sigma_\mu, \Sigma_\nu)^2$$

where B is the Bures metric

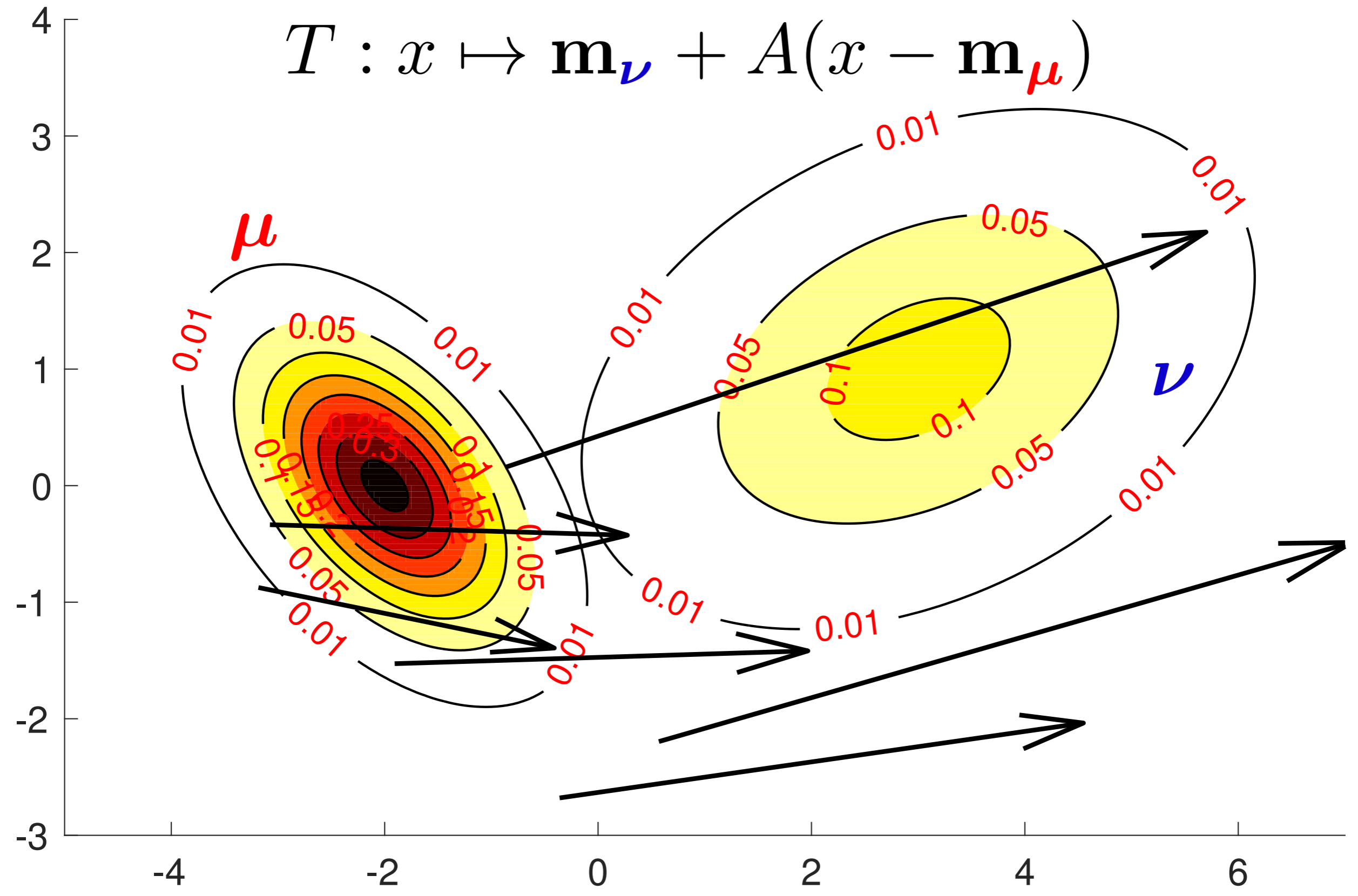
$$B(\Sigma_\mu, \Sigma_\nu)^2 = \text{trace}(\Sigma_\mu + \Sigma_\nu - 2(\Sigma_\mu^{1/2} \Sigma_\nu \Sigma_\mu^{1/2})^{1/2}).$$

The map $T : x \mapsto \mathbf{m}_\nu + A(x - \mathbf{m}_\mu)$ is **optimal**,

where $A = \Sigma_\mu^{-\frac{1}{2}} \left(\Sigma_\mu^{\frac{1}{2}} \Sigma_\nu \Sigma_\mu^{\frac{1}{2}} \right)^{\frac{1}{2}} \Sigma_\mu^{-\frac{1}{2}}$.

Easy (2): Gaussian Measures

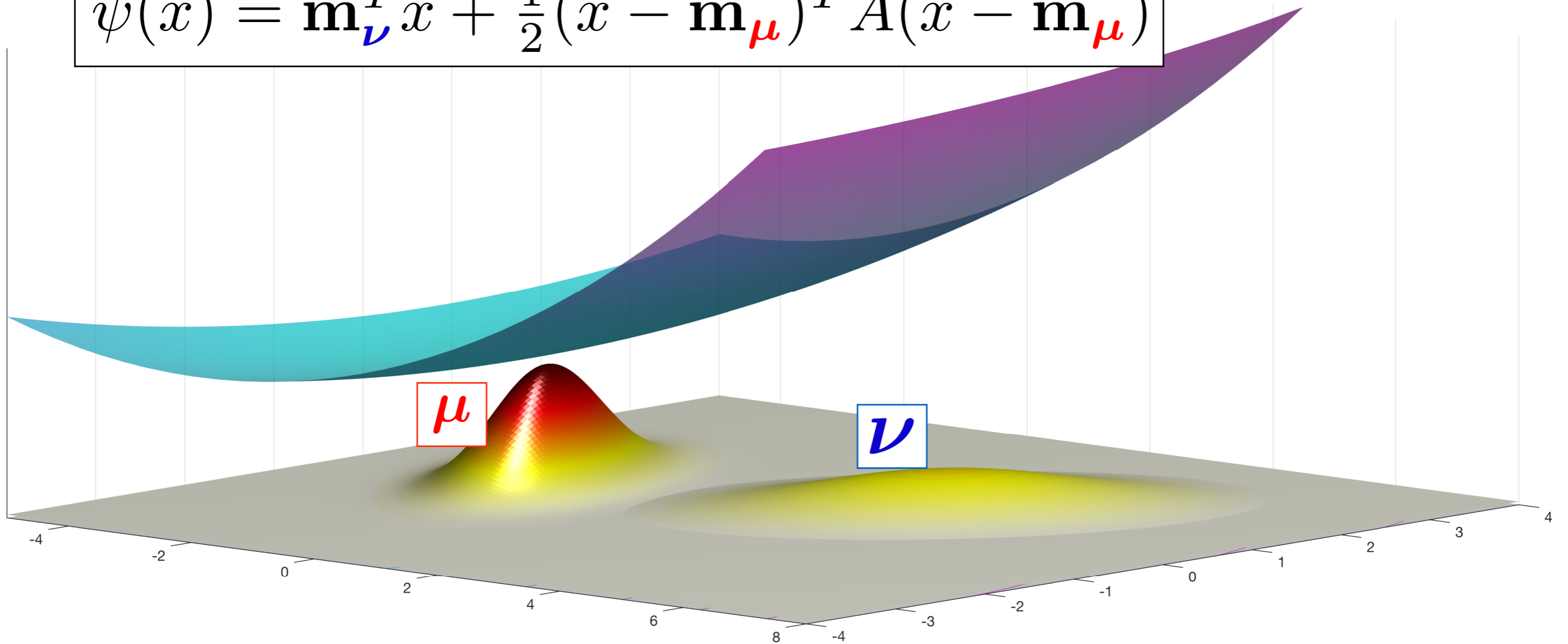
$$T : x \mapsto \mathbf{m}_\nu + A(x - \mathbf{m}_\mu)$$



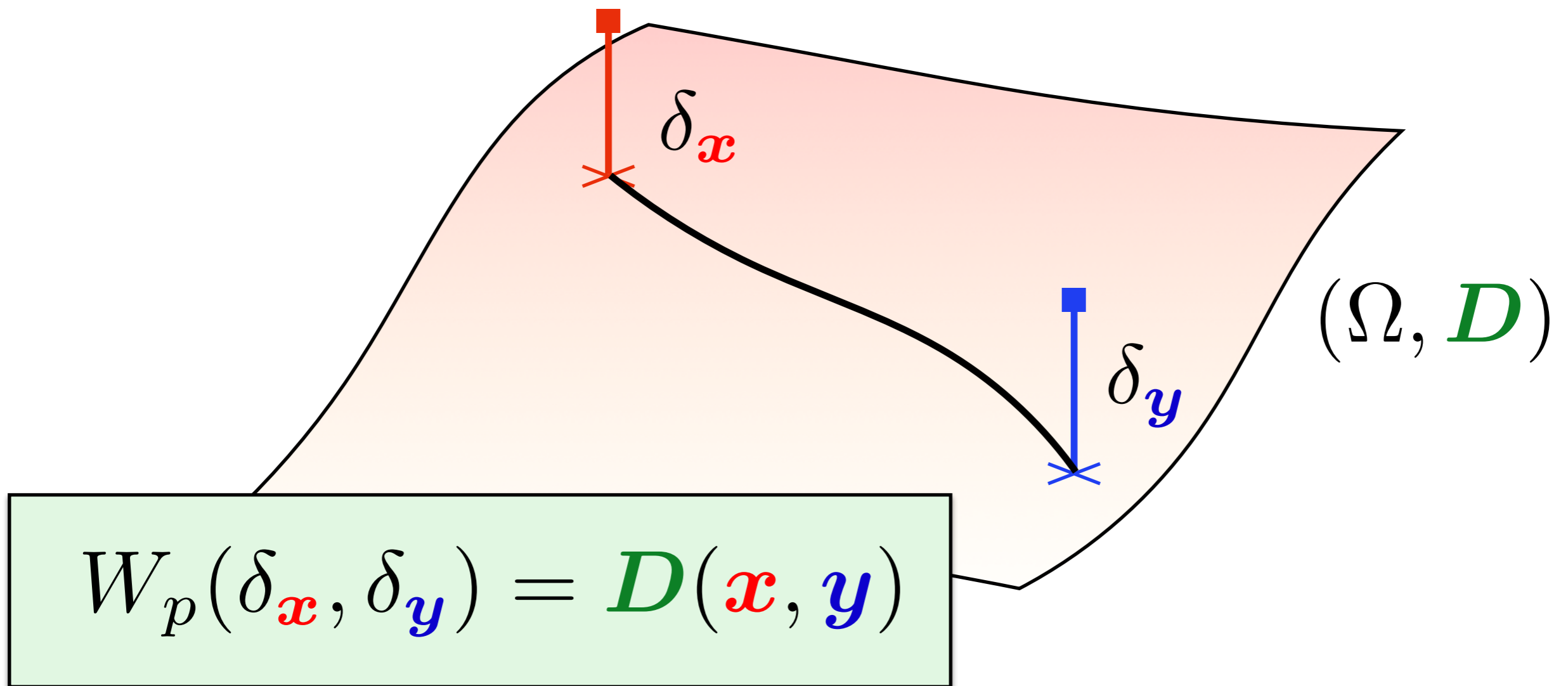
Easy (2): Gaussian Measures

$$T = \nabla \psi : x \mapsto \mathbf{m}_\nu + A(x - \mathbf{m}_\mu)$$

$$\psi(x) = \mathbf{m}_\nu^T x + \frac{1}{2} (x - \mathbf{m}_\mu)^T A (x - \mathbf{m}_\mu)$$



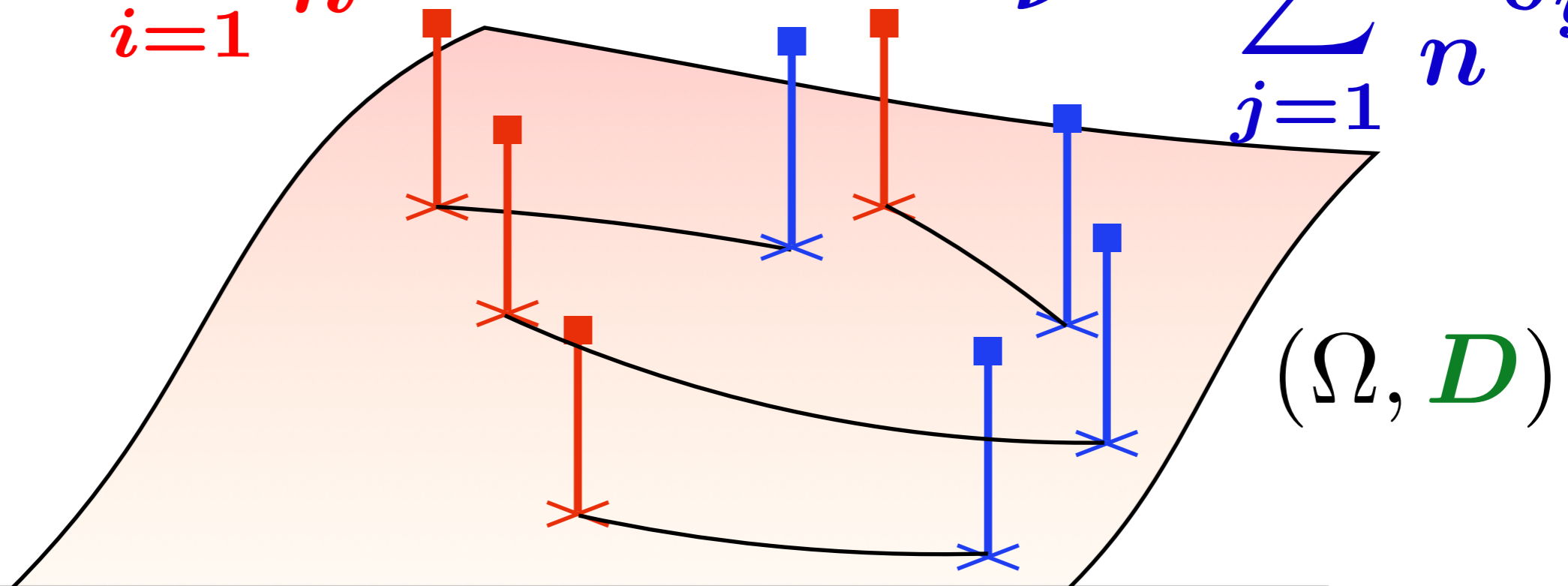
Wasserstein Between Two Diracs



Linear Assignment \subset Wasserstein

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

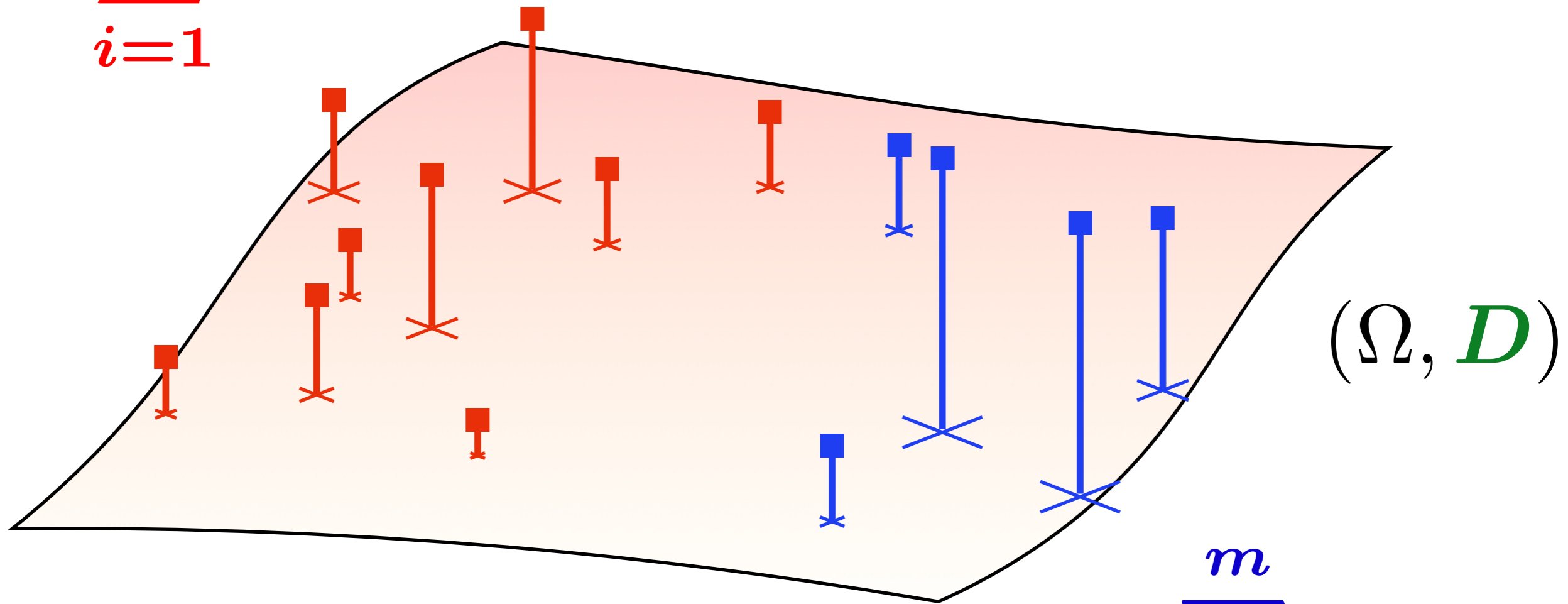
$$\nu = \sum_{j=1}^n \frac{1}{n} \delta_{y_j}$$



$$W_p^p(\mu, \nu) = \min_{\sigma \in S_n} \frac{1}{n} \sum_{i=1}^n D(x_i, y_{\sigma_i})^p$$

OT on Two Empirical Measures

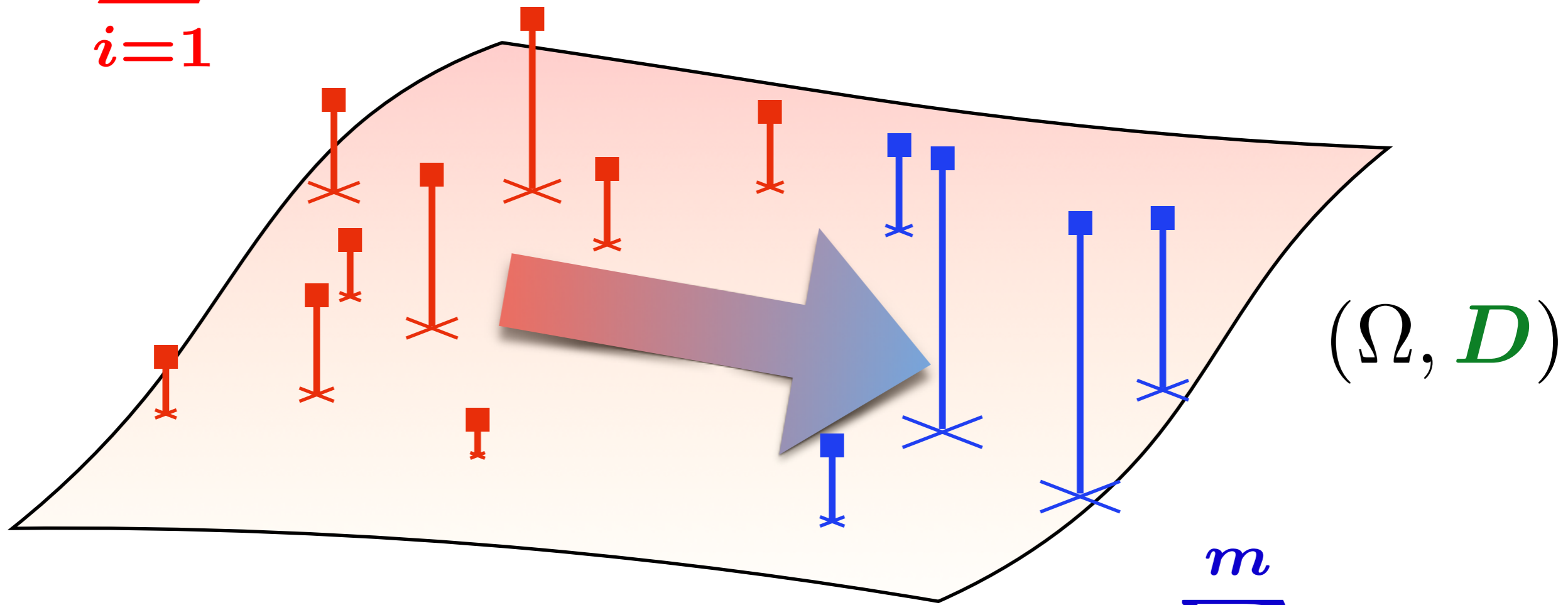
$$\mu = \sum_{i=1}^n a_i \delta_{x_i}$$



$$\nu = \sum_{j=1}^m b_j \delta_{y_j}$$

OT on Two Empirical Measures

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



$$\nu = \sum_{j=1}^m b_j \delta_{y_j}$$

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

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

Dual Kantorovich Problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{X}\mathbf{Y}} \rangle$$

Dual Kantorovich Problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{x}\mathbf{y}} \rangle$$

Def. Dual OT problem

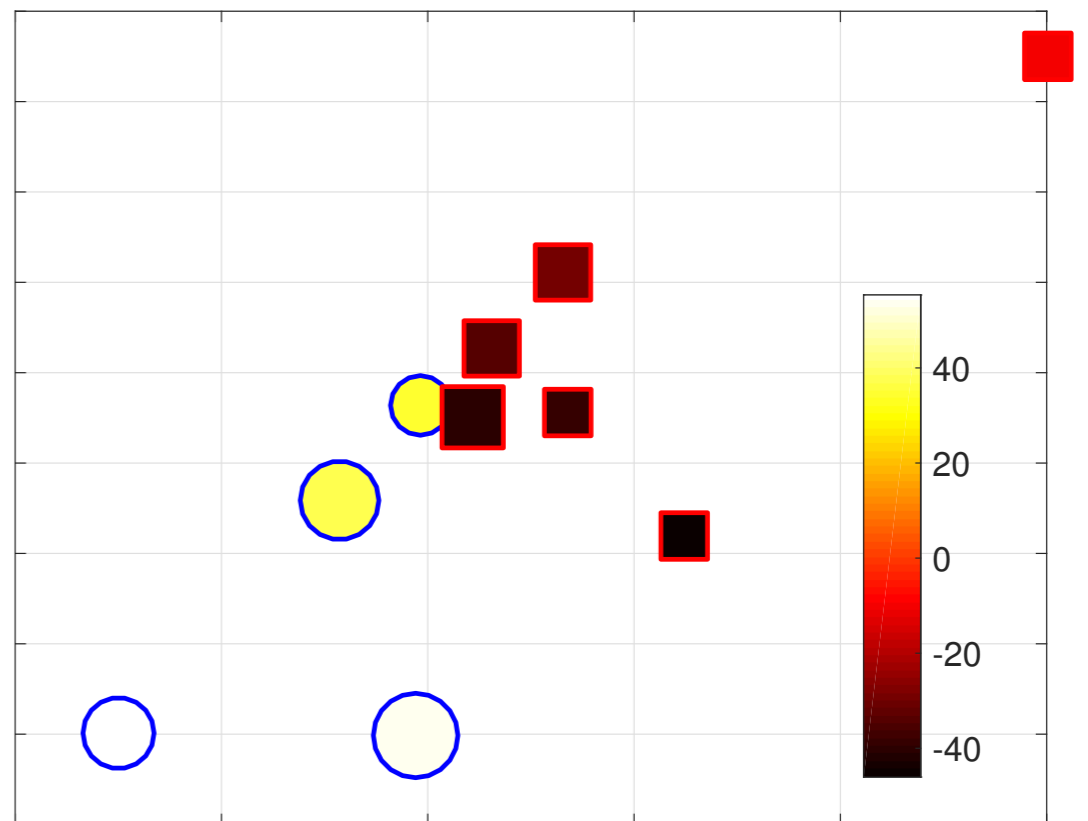
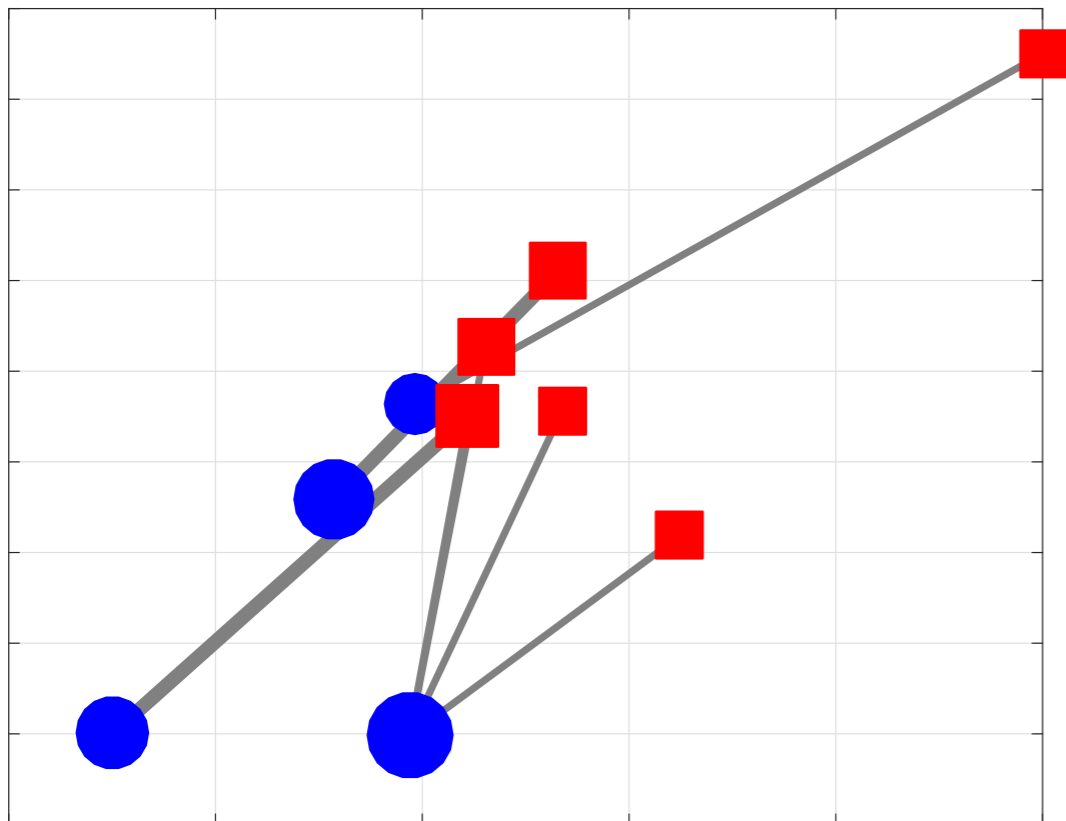
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \alpha_i + \beta_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$

Dual Kantorovich Problem

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{x}\mathbf{y}} \rangle$$

Def. Dual OT problem

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

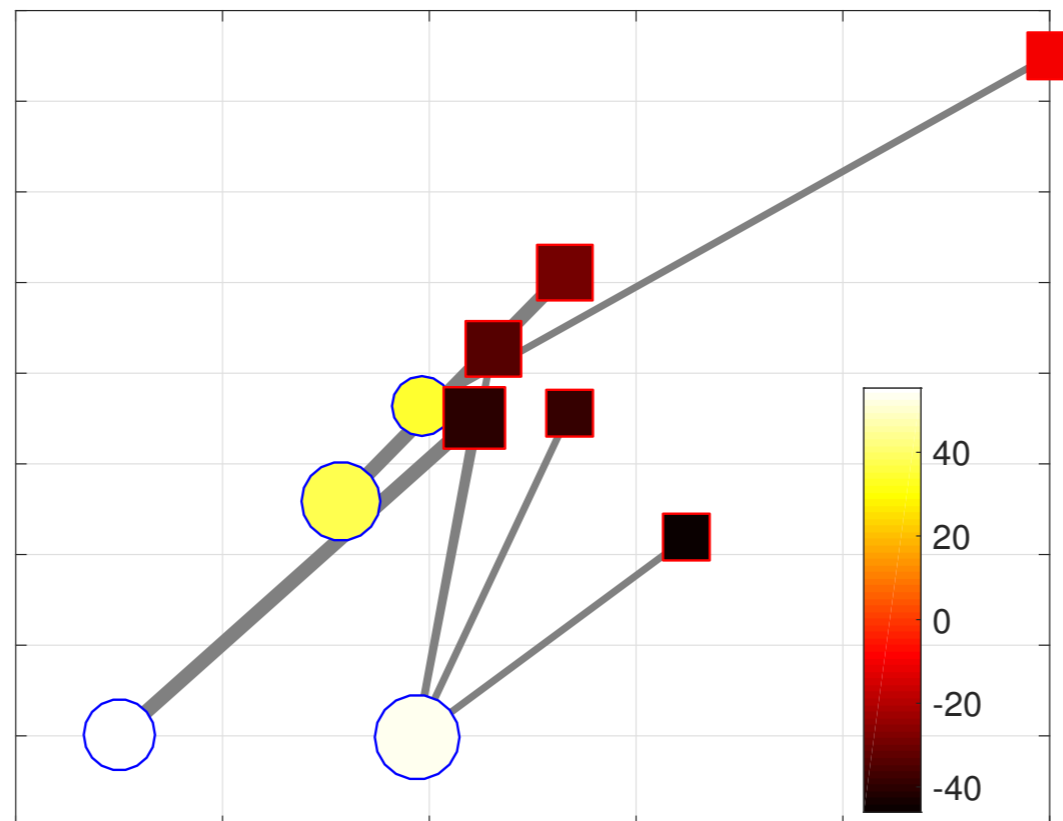


Dual Kantorovich Problem

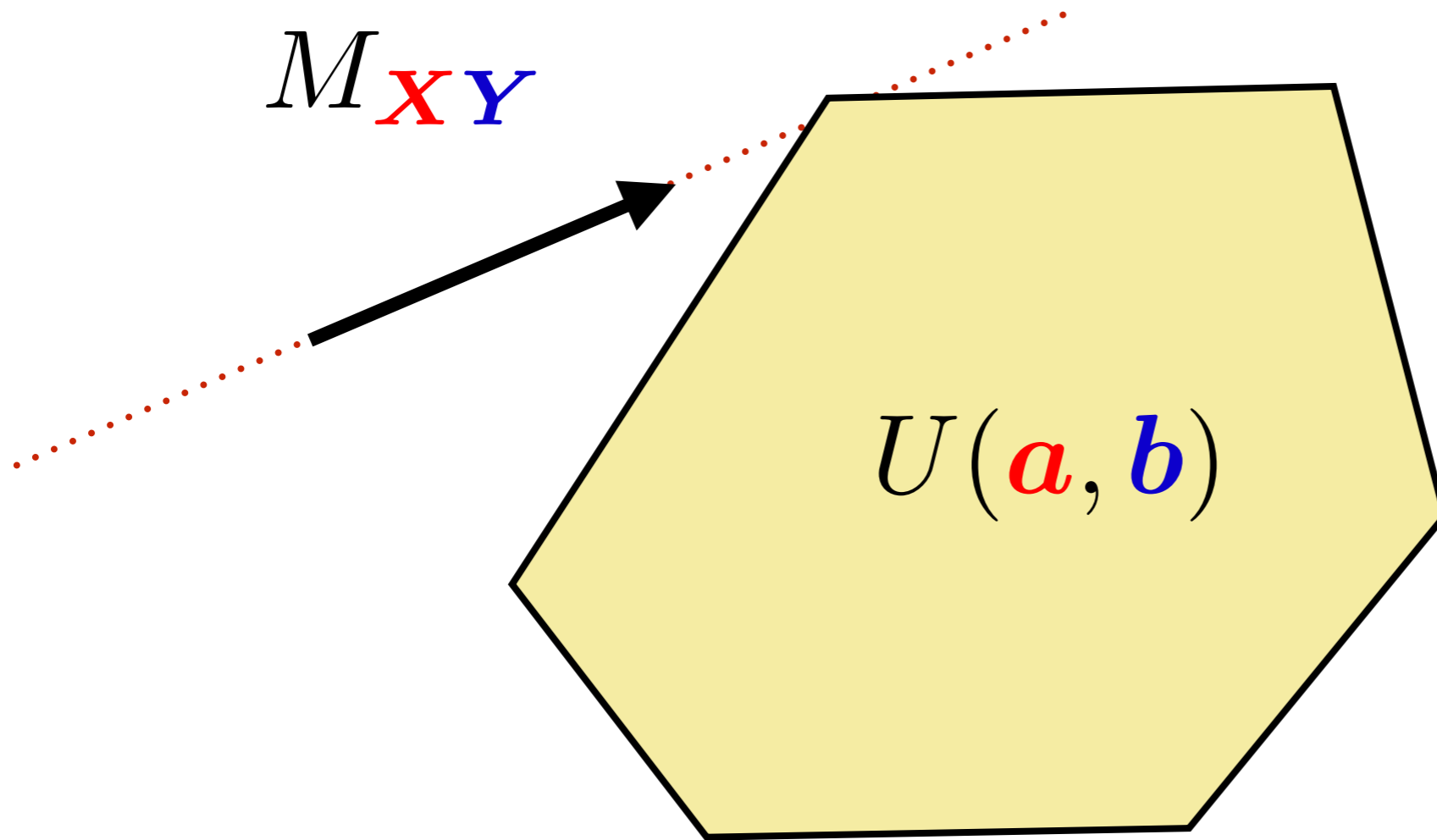
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{x}\mathbf{y}} \rangle$$

Def. Dual OT problem

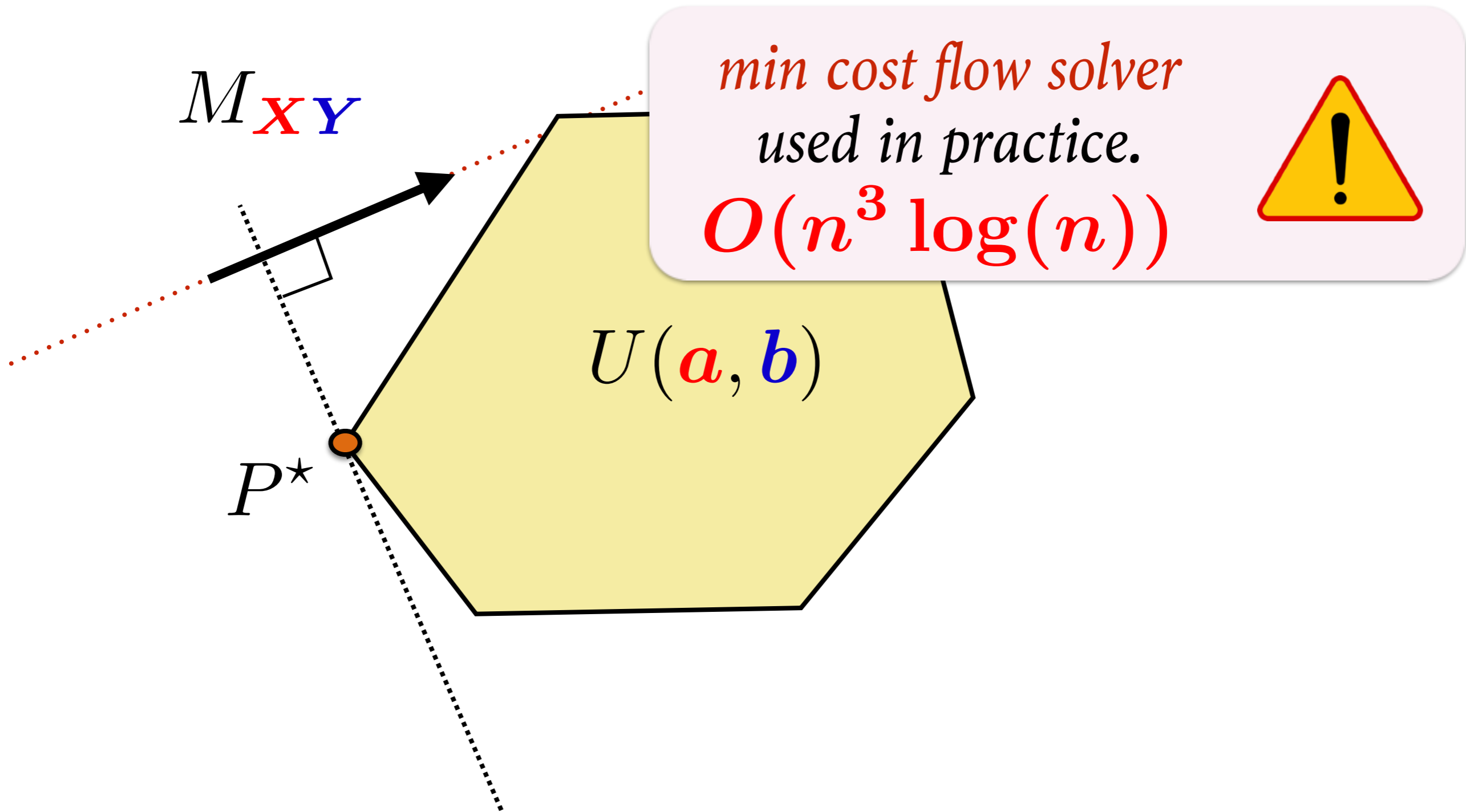
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \alpha_i + \beta_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$



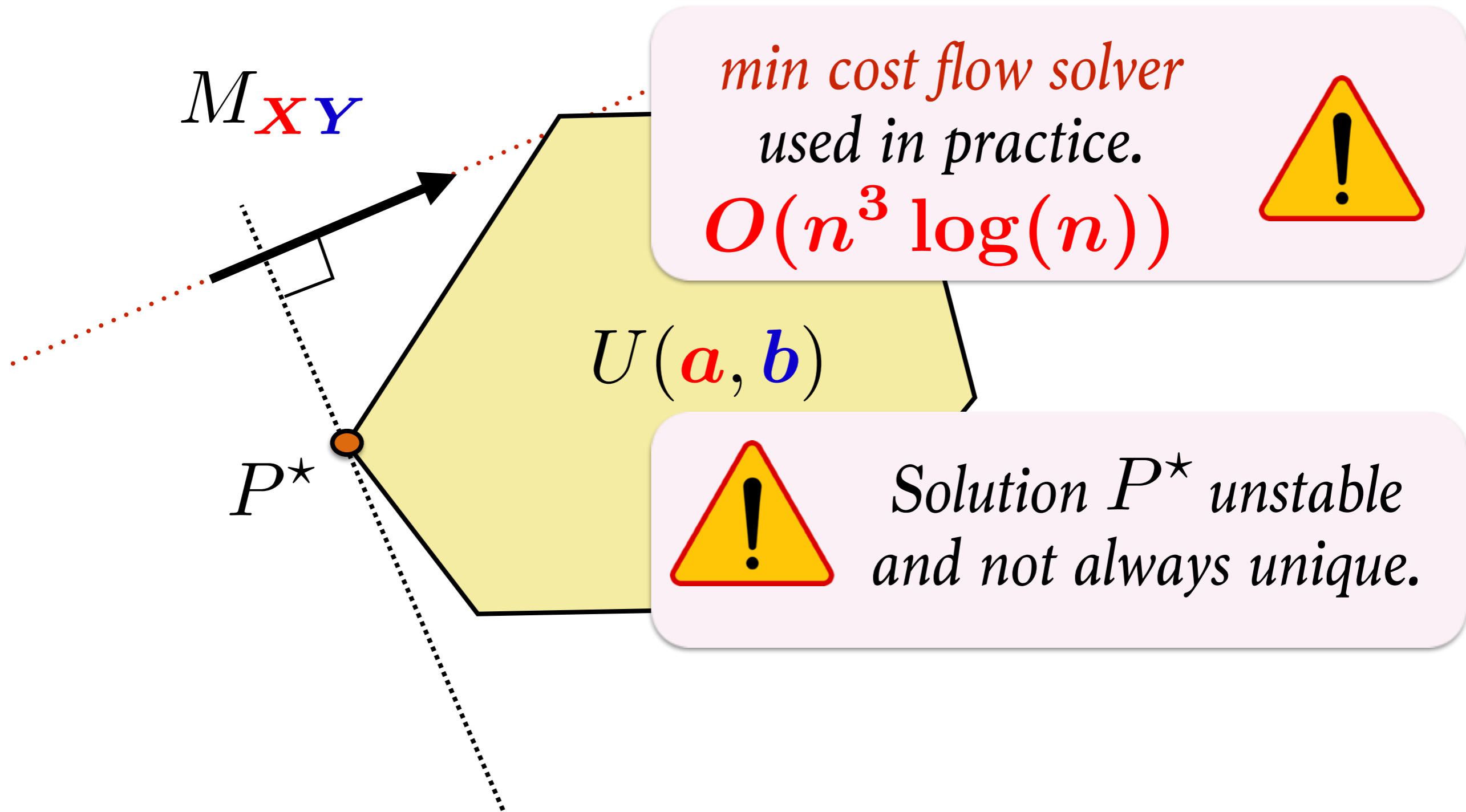
Solving the OT Problem



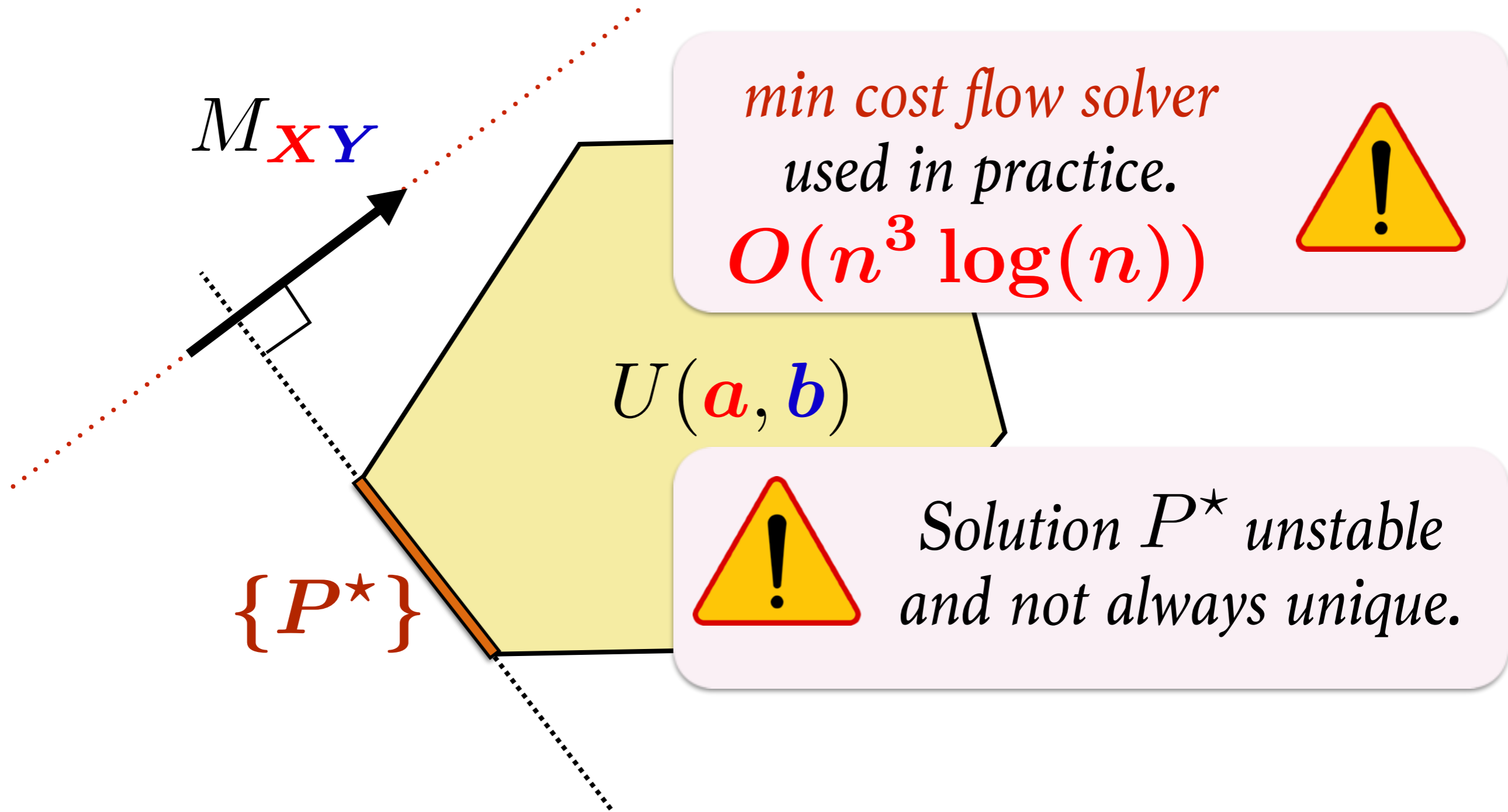
Solving the OT Problem



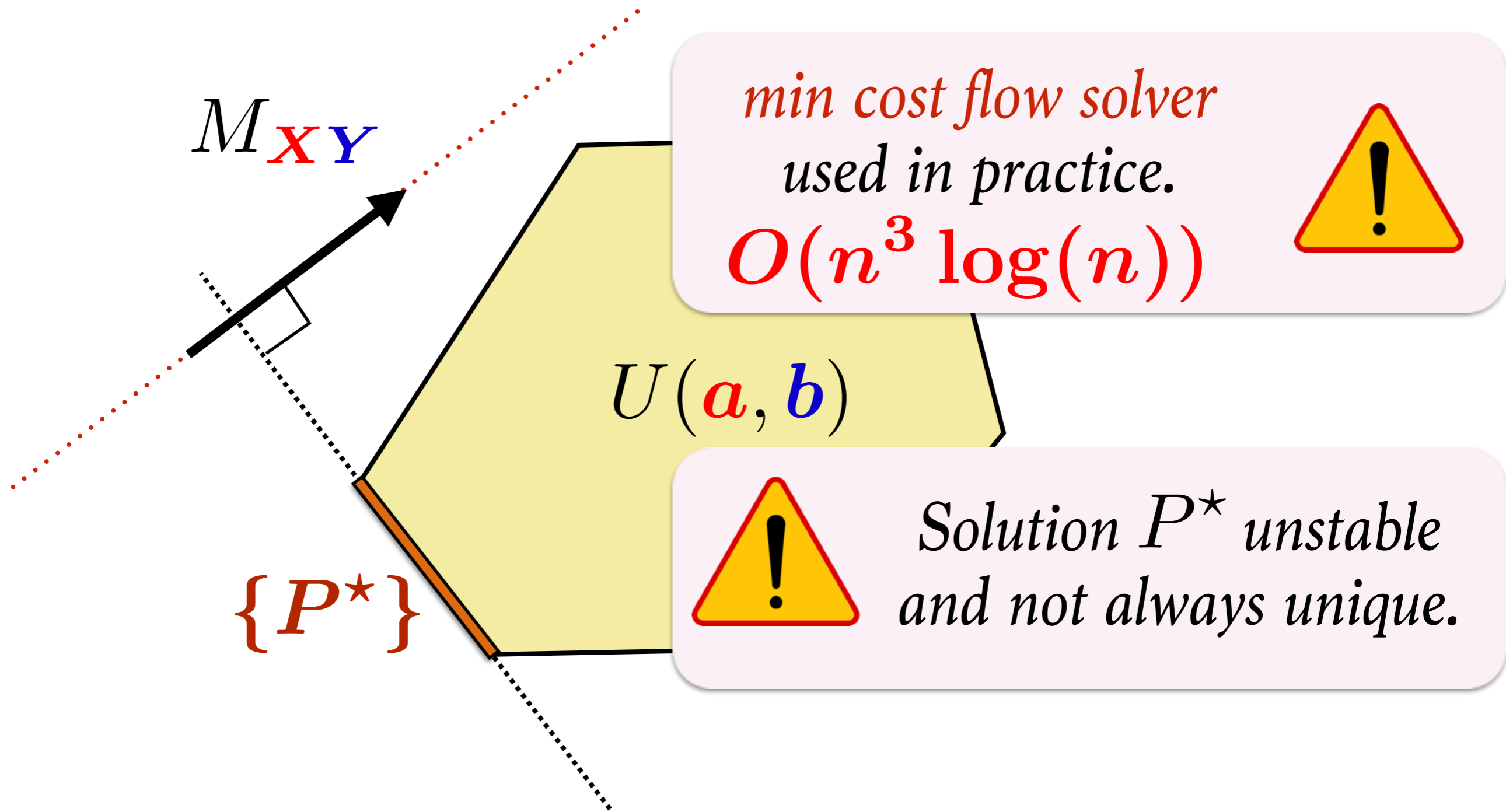
Solving the OT Problem



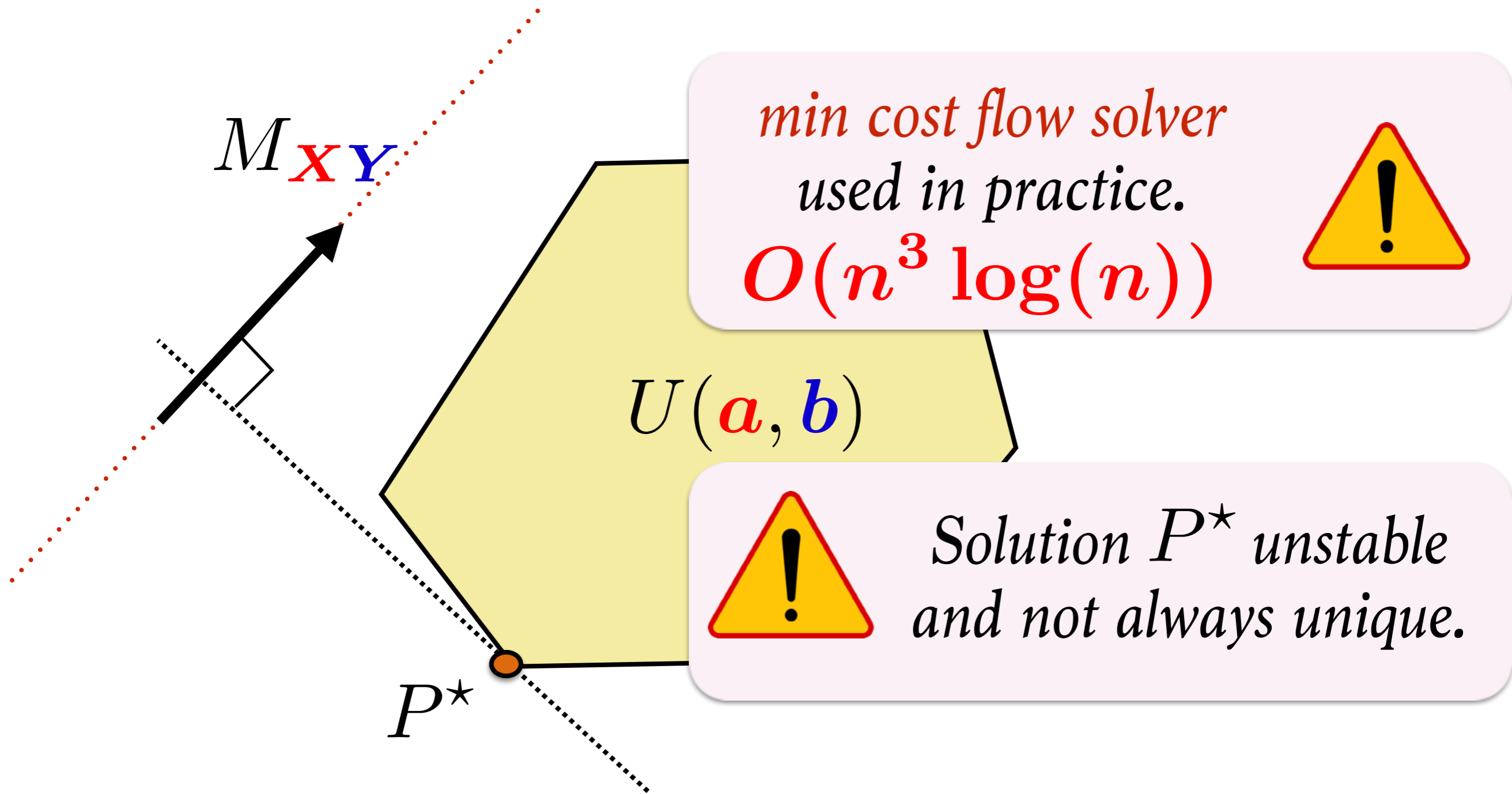
Solving the OT Problem



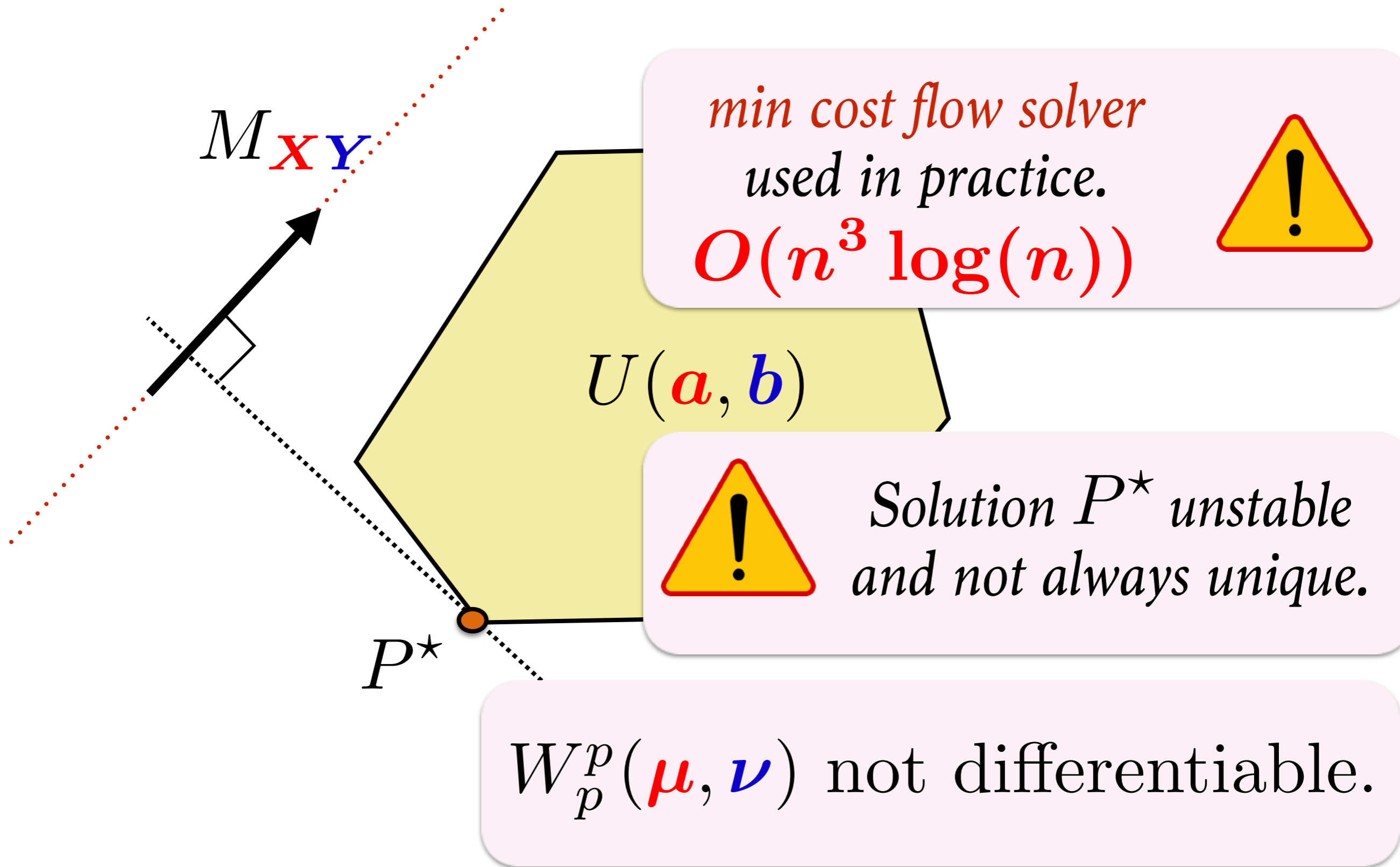
Solving the OT Problem



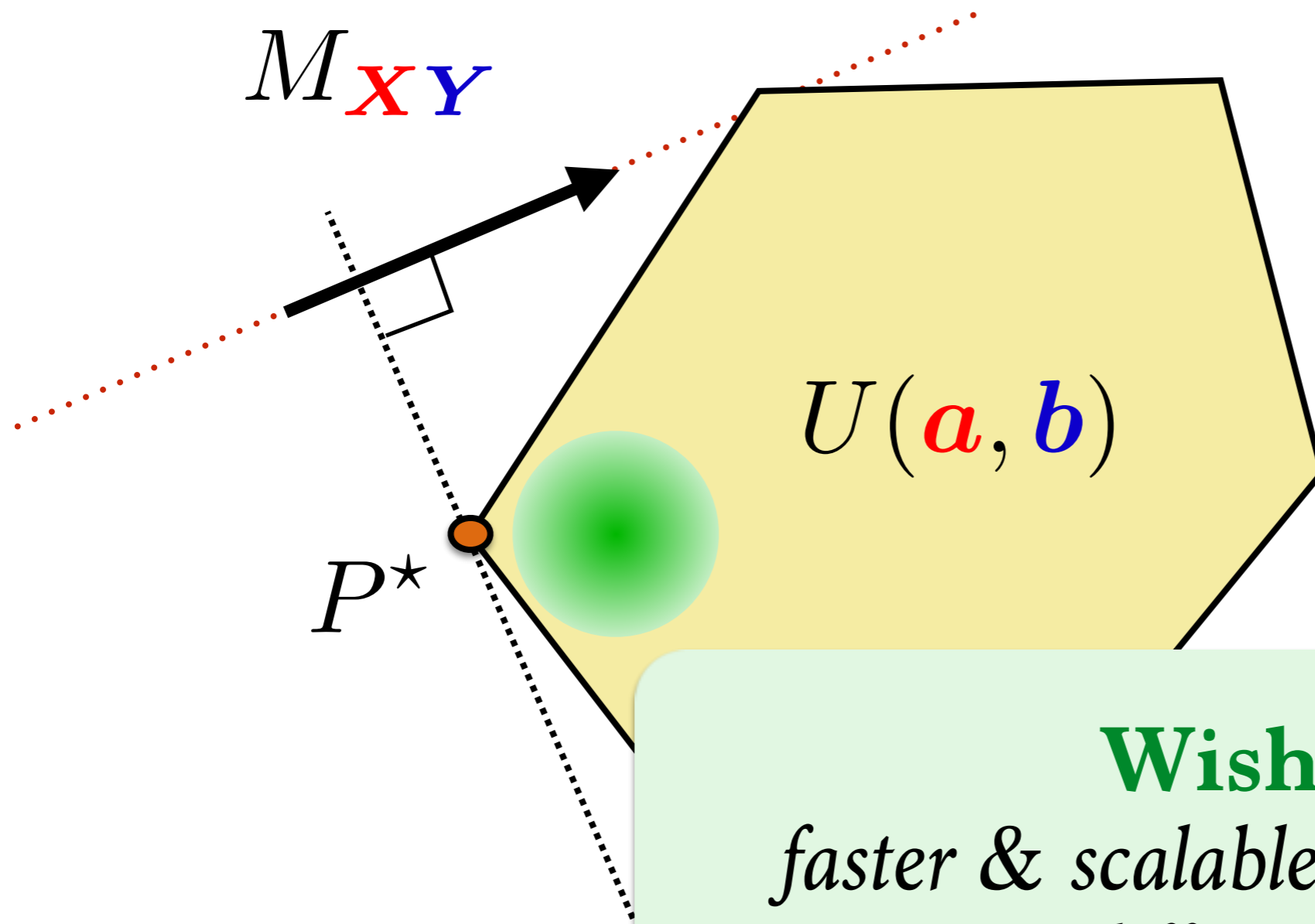
Solving the OT Problem



Solving the OT Problem



Solution: Regularization



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

Entropic Regularization [Wilson'62]

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

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

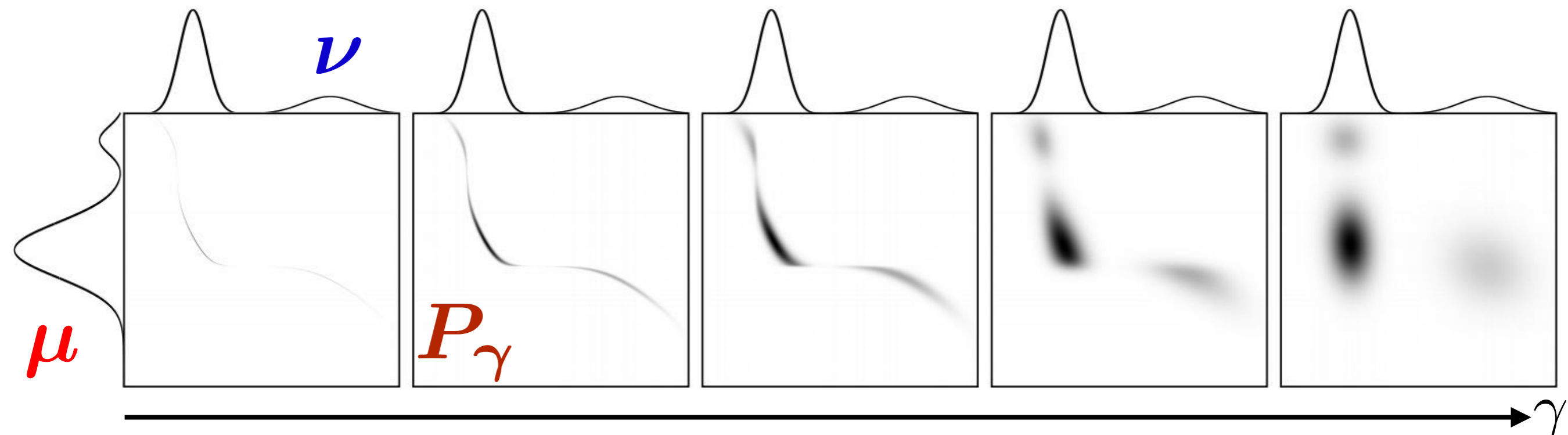
$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij} (\log P_{ij} - 1)$$

Note: Unique optimal solution because of strong concavity of entropy

Entropic Regularization [Wilson'62]

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

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



Note: Unique optimal solution because of strong concavity of entropy

Fast & Scalable Algorithm

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

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

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

Fast & Scalable Algorithm

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

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

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

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

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

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

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \text{diag}(\mathbf{u}) K \text{diag}(\mathbf{v}) \mathbf{1}_m & = \mathbf{a} \\ \text{diag}(\mathbf{v}) K^T \text{diag}(\mathbf{u}) \mathbf{1}_n & = \mathbf{b} \end{cases}$$

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \operatorname{diag}(\mathbf{u}) K \operatorname{diag}(\mathbf{v}) \mathbf{1}_m & = \mathbf{a} \\ \operatorname{diag}(\mathbf{v}) K^T \underbrace{\operatorname{diag}(\mathbf{u}) \mathbf{1}_n}_{\mathbf{u}} & = \mathbf{b} \end{cases}$$

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \text{diag}(\mathbf{u}) K \text{diag}(\mathbf{v}) \mathbf{1}_m & = \mathbf{a} \\ \text{diag}(\mathbf{v}) K^T \text{diag}(\mathbf{u}) \mathbf{1}_n & = \mathbf{b} \end{cases}$$

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \text{diag}(\mathbf{u}) \mathbf{K} \mathbf{v} & = \mathbf{a} \\ \text{diag}(\mathbf{v}) \mathbf{K}^T \mathbf{u} & = \mathbf{b} \end{cases}$$

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \mathbf{u} \odot \mathbf{K} \mathbf{v} & = \mathbf{a} \\ \mathbf{v} \odot \mathbf{K}^T \mathbf{u} & = \mathbf{b} \end{cases}$$

Fast & Scalable Algorithm

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

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

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

$$P_\gamma \in U(\mathbf{a}, \mathbf{b}) \Leftrightarrow \begin{cases} \mathbf{u} = \mathbf{a} / \mathbf{K} \mathbf{v} \\ \mathbf{v} = \mathbf{b} / \mathbf{K}^T \mathbf{u} \end{cases}$$

Fast & Scalable Algorithm

Sinkhorn's Algorithm : Repeat

1. $\mathbf{u} = \mathbf{a} / K \mathbf{v}$

2. $\mathbf{v} = \mathbf{b} / K^T \mathbf{u}$

Fast & Scalable Algorithm

Sinkhorn's Algorithm : Repeat

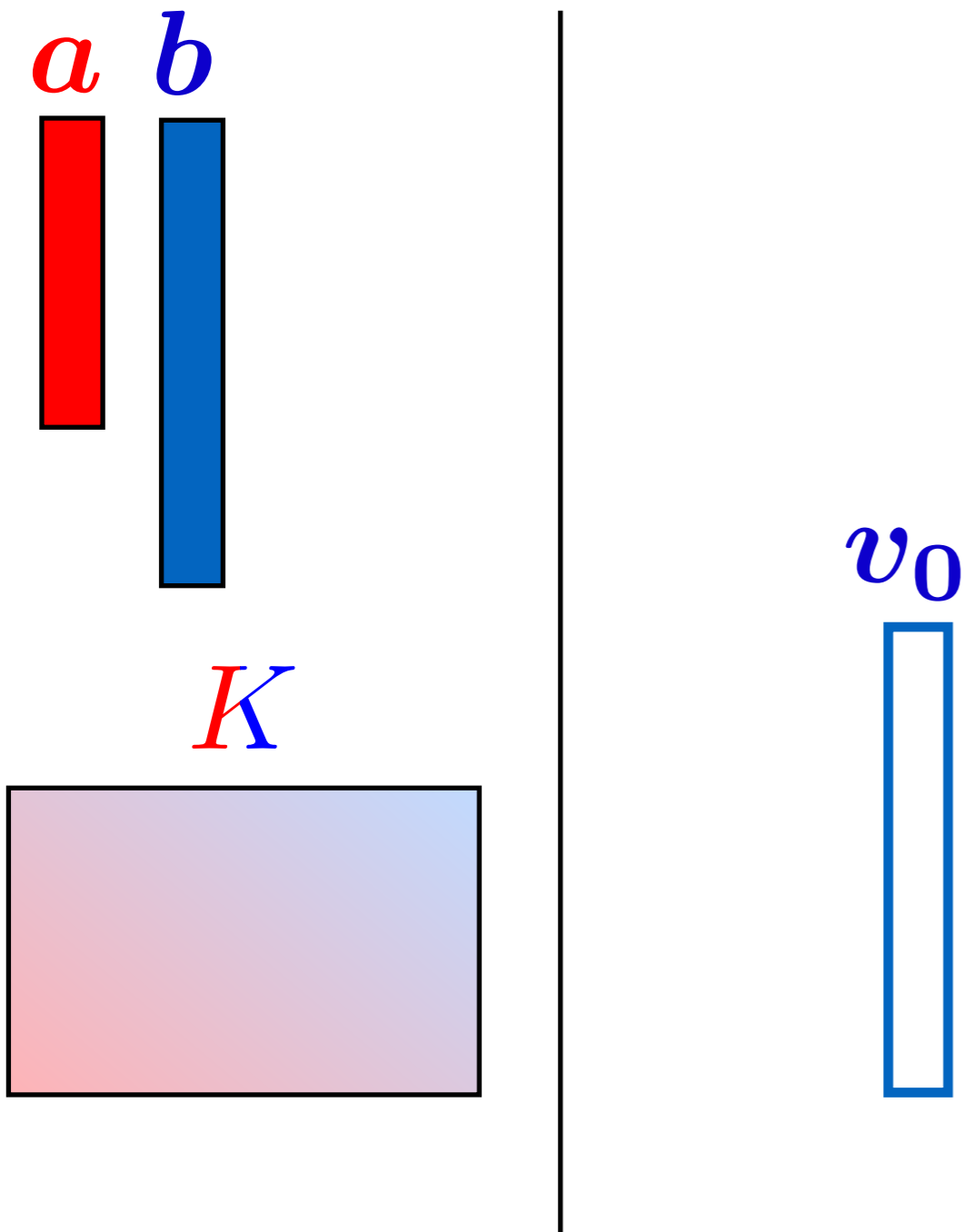
$$1. \quad \mathbf{u} = \mathbf{a} / \mathbf{K} \mathbf{v}$$

$$2. \quad \mathbf{v} = \mathbf{b} / \mathbf{K}^T \mathbf{u}$$

- [Sinkhorn'64] proved convergence for the first time.
- [Lorenz'89] linear convergence, see [Altschuler'17]
- $O(nm)$ complexity, GPGPU parallel [Cuturi'13].
- $O(n \log n)$ on gridded spaces using convolutions.
[Solomon'15]

Fast & Scalable Algorithm

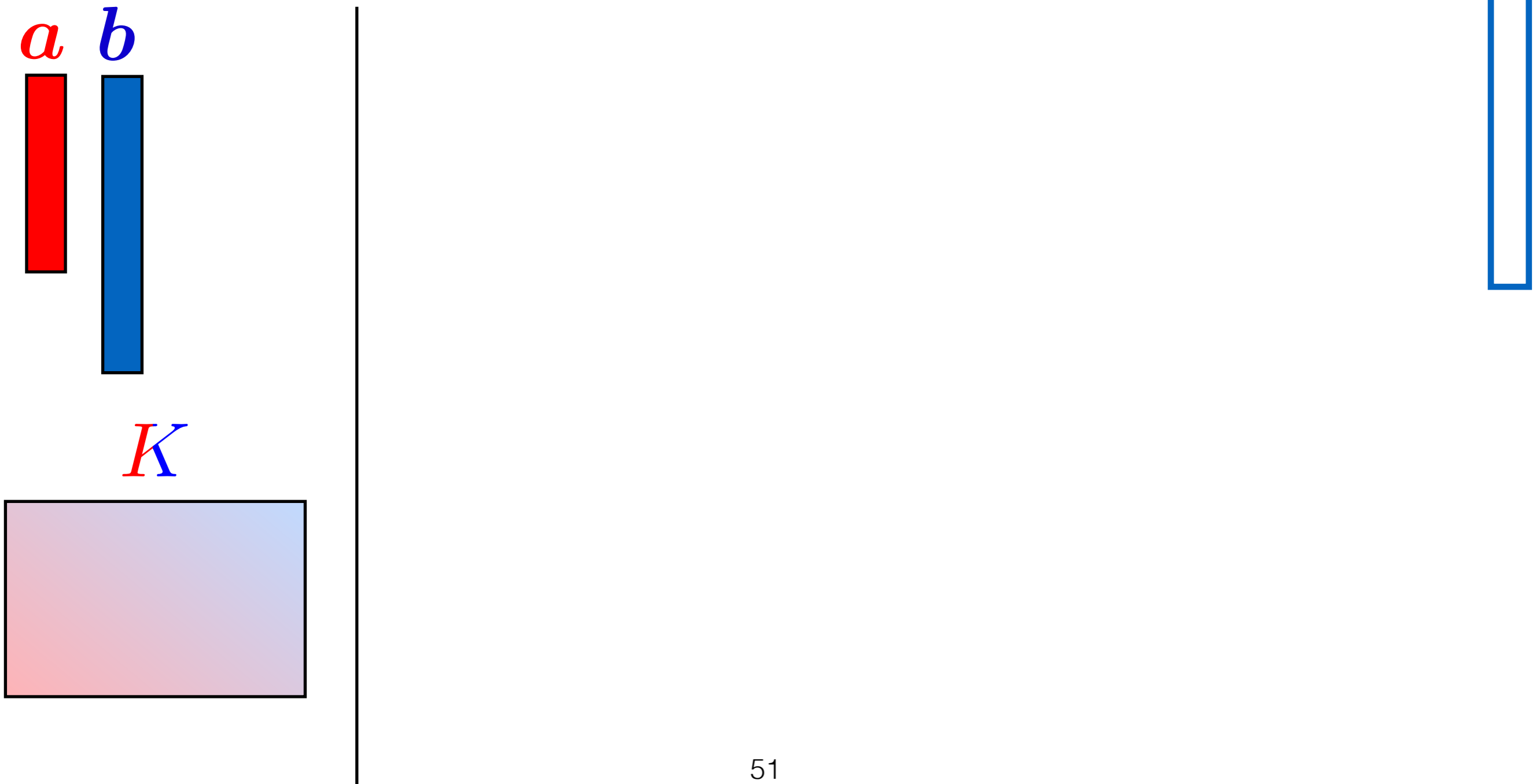
- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})

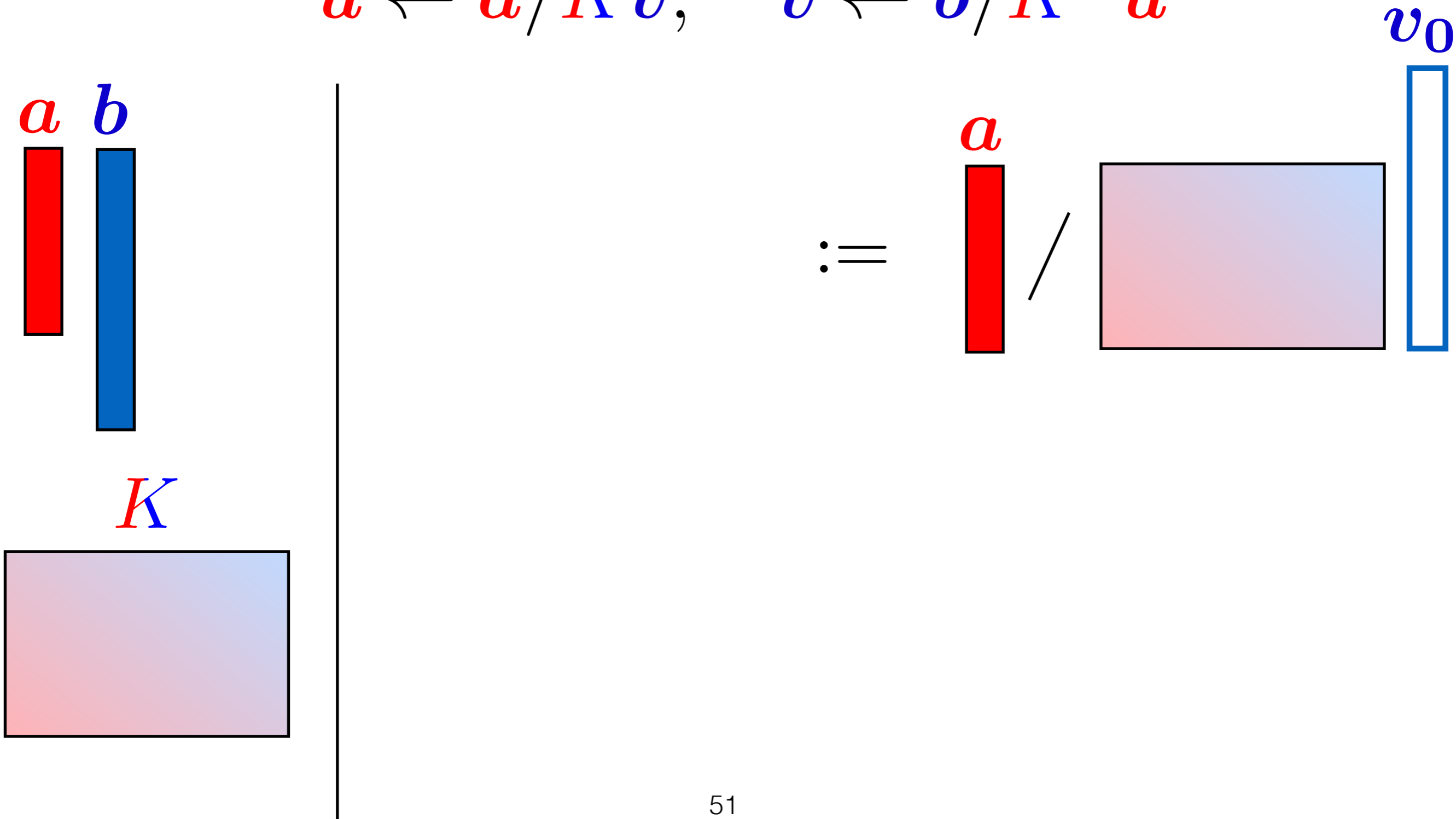
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})

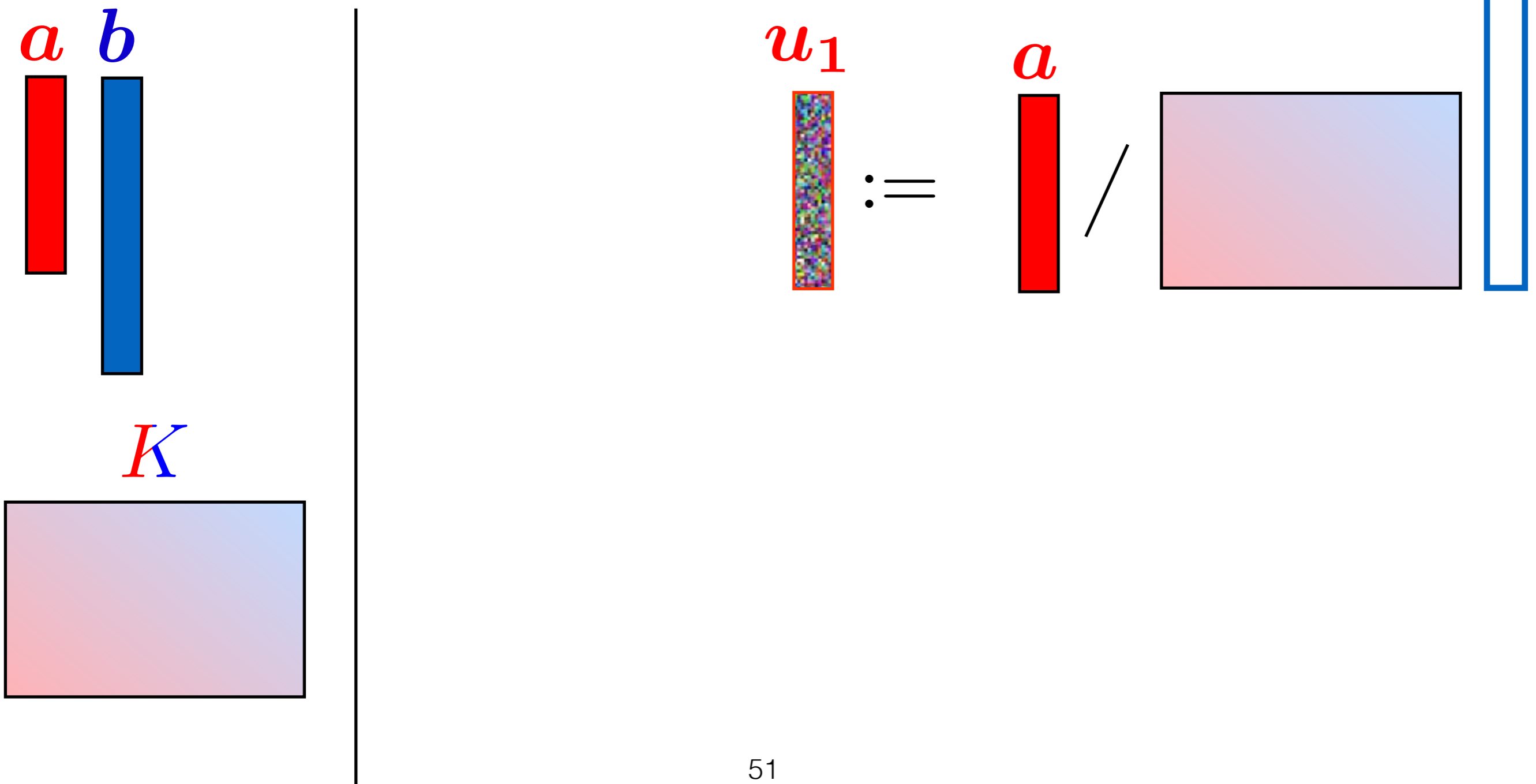
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

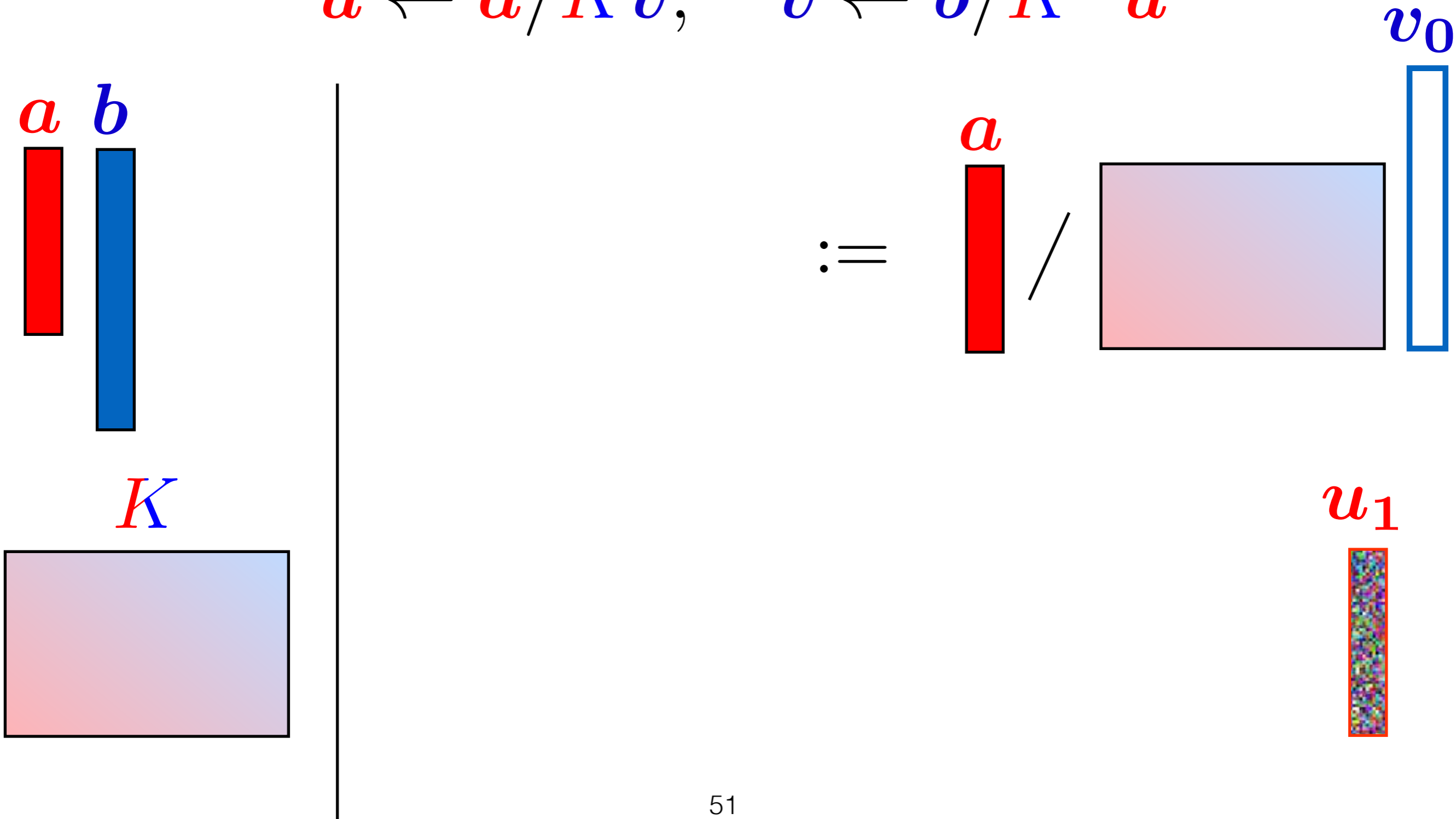
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})

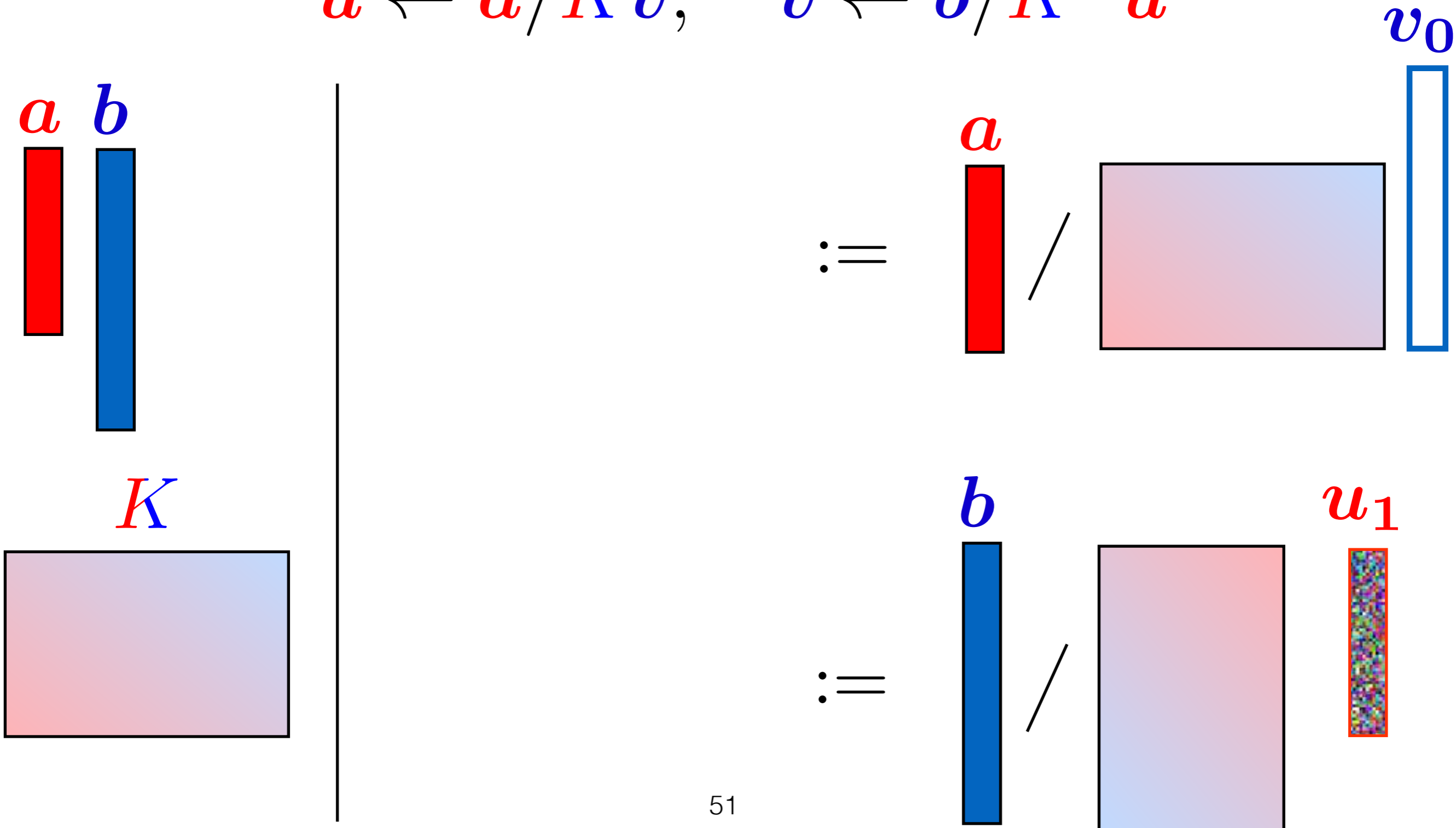
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

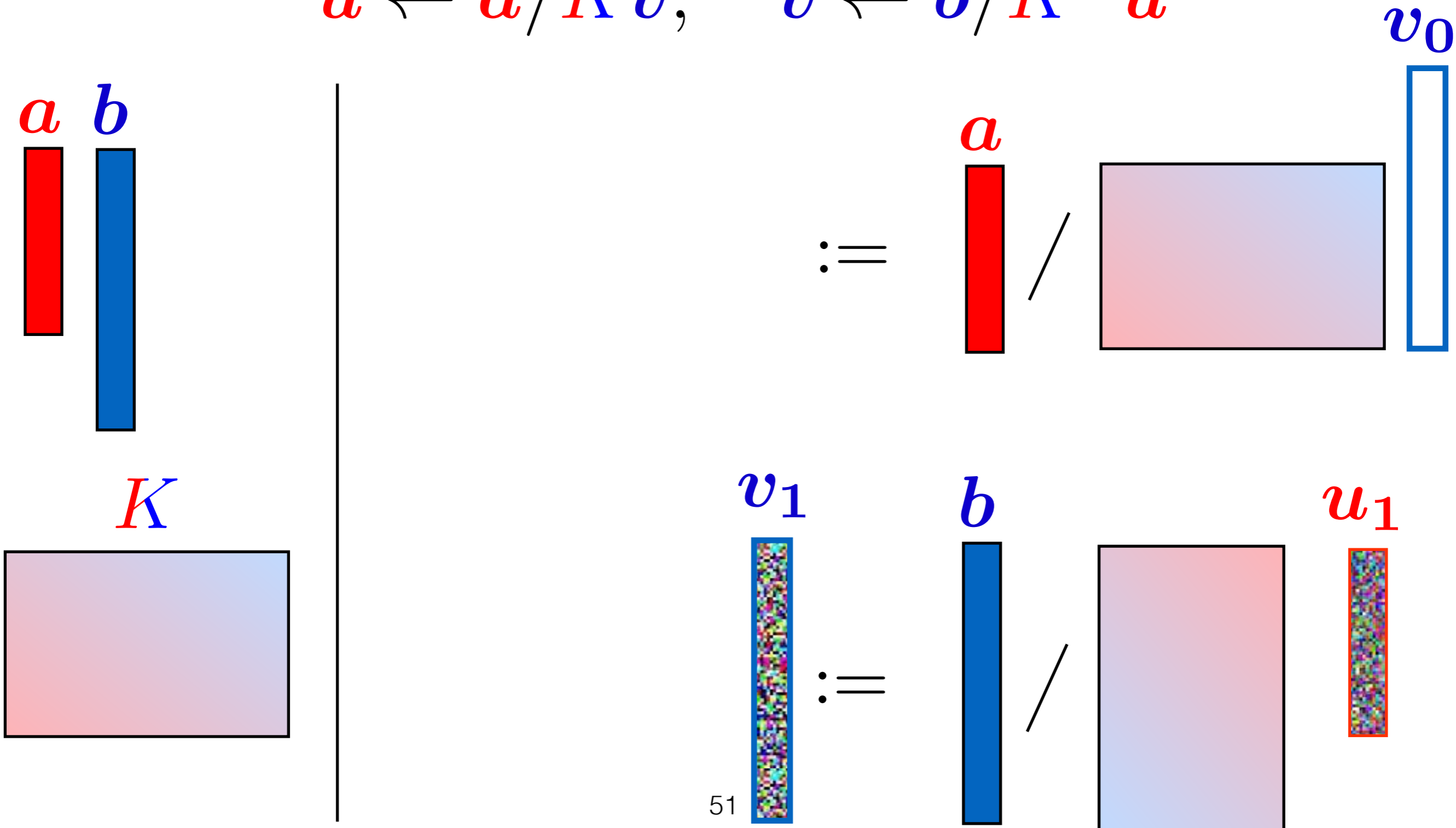
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

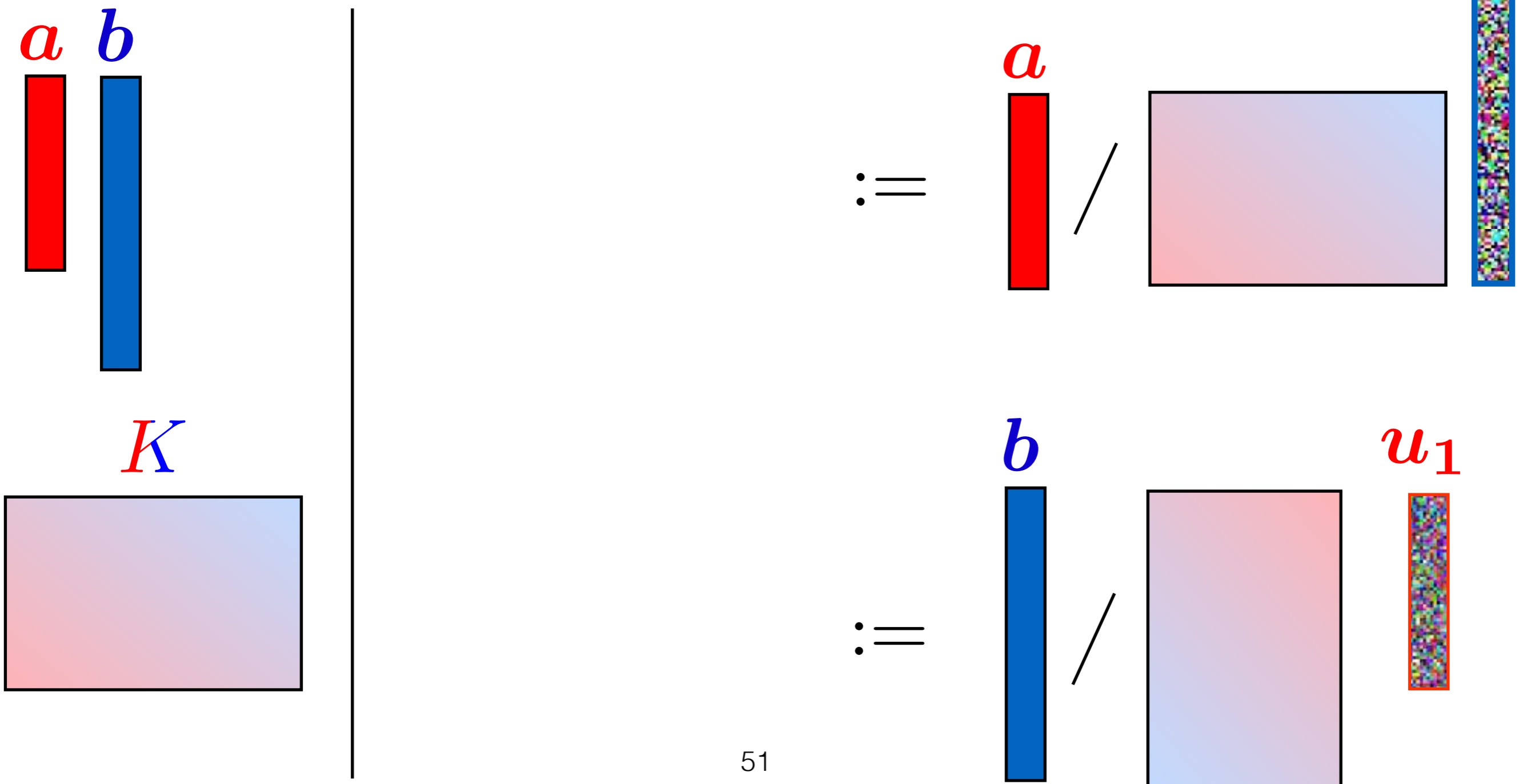
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

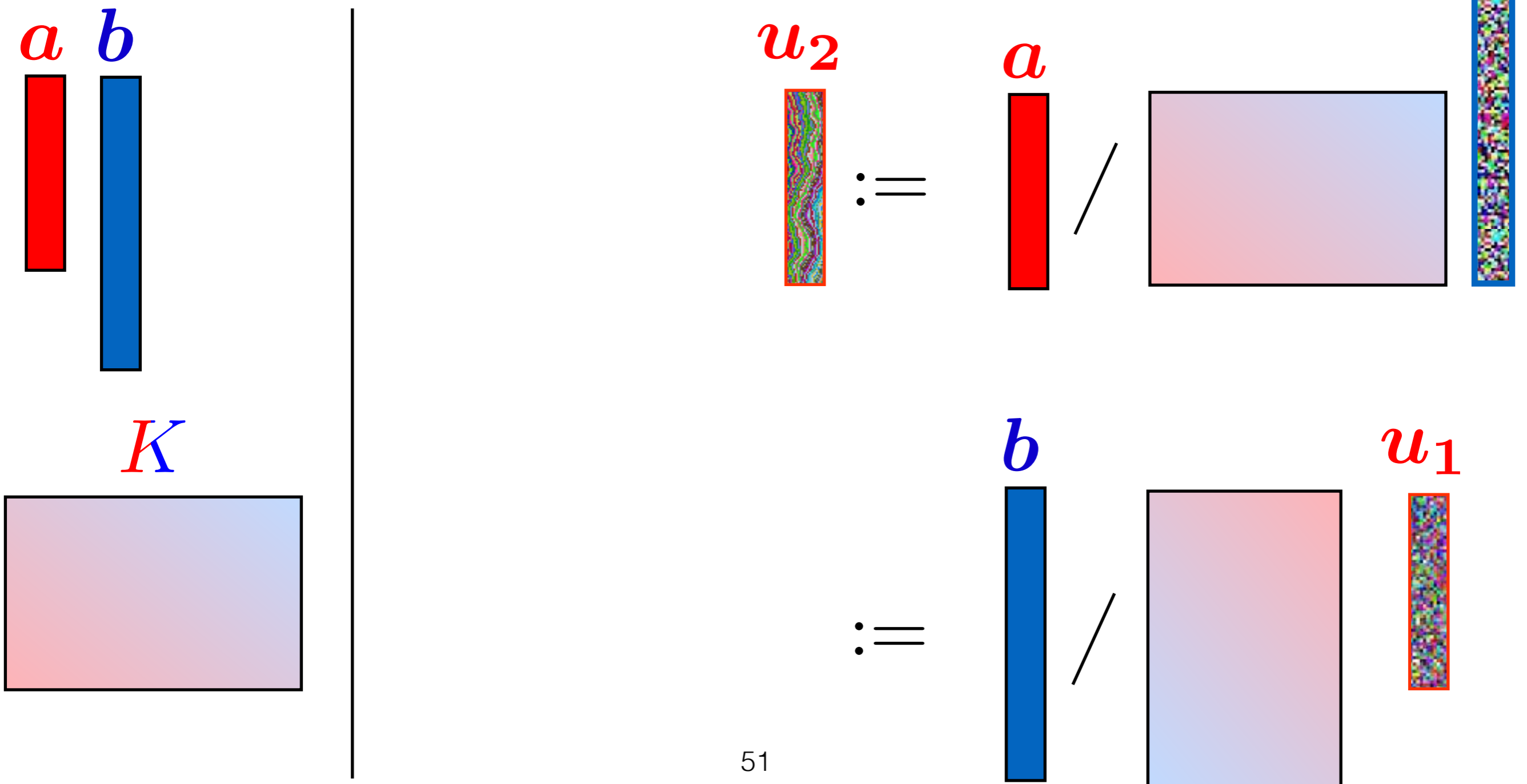
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

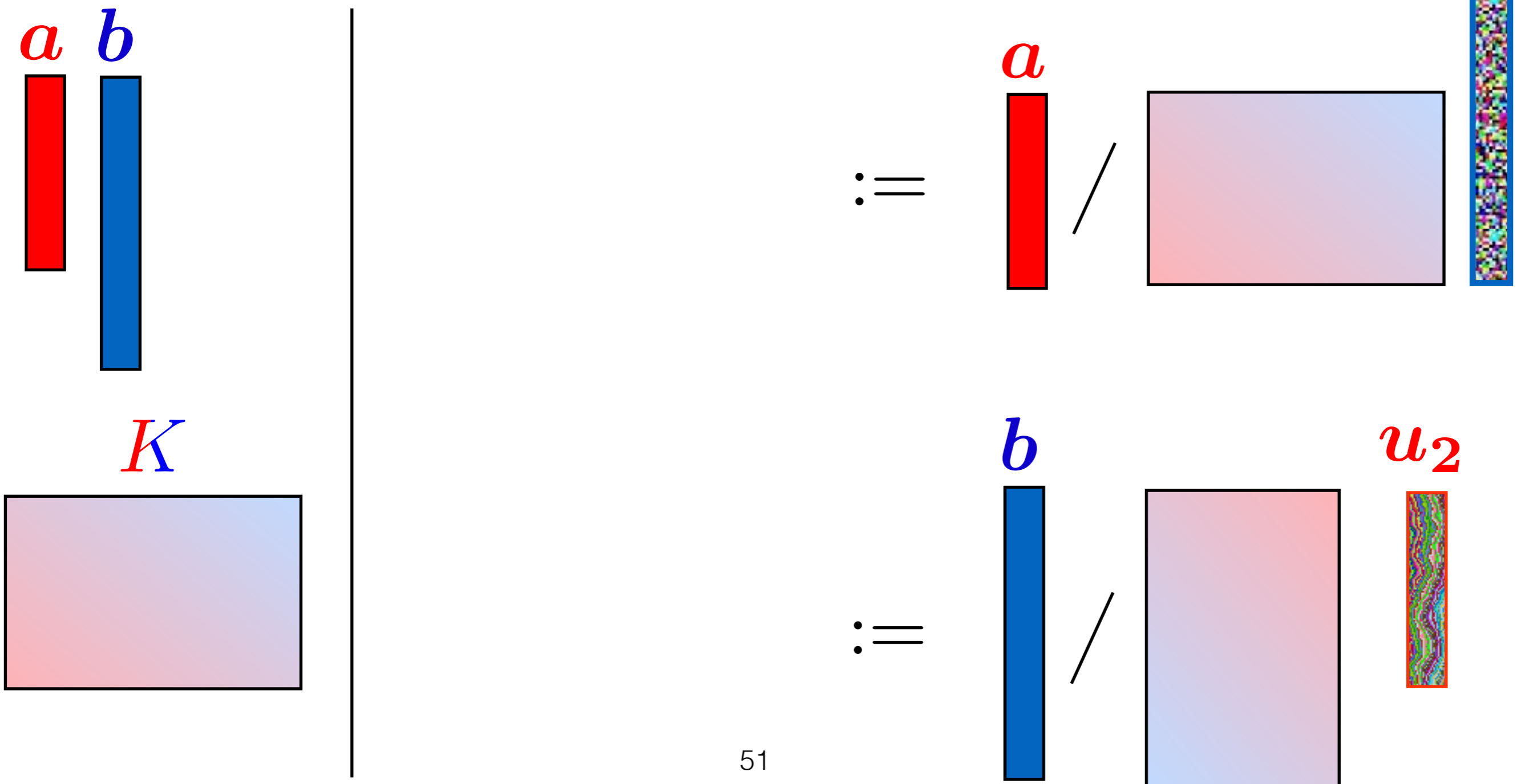
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

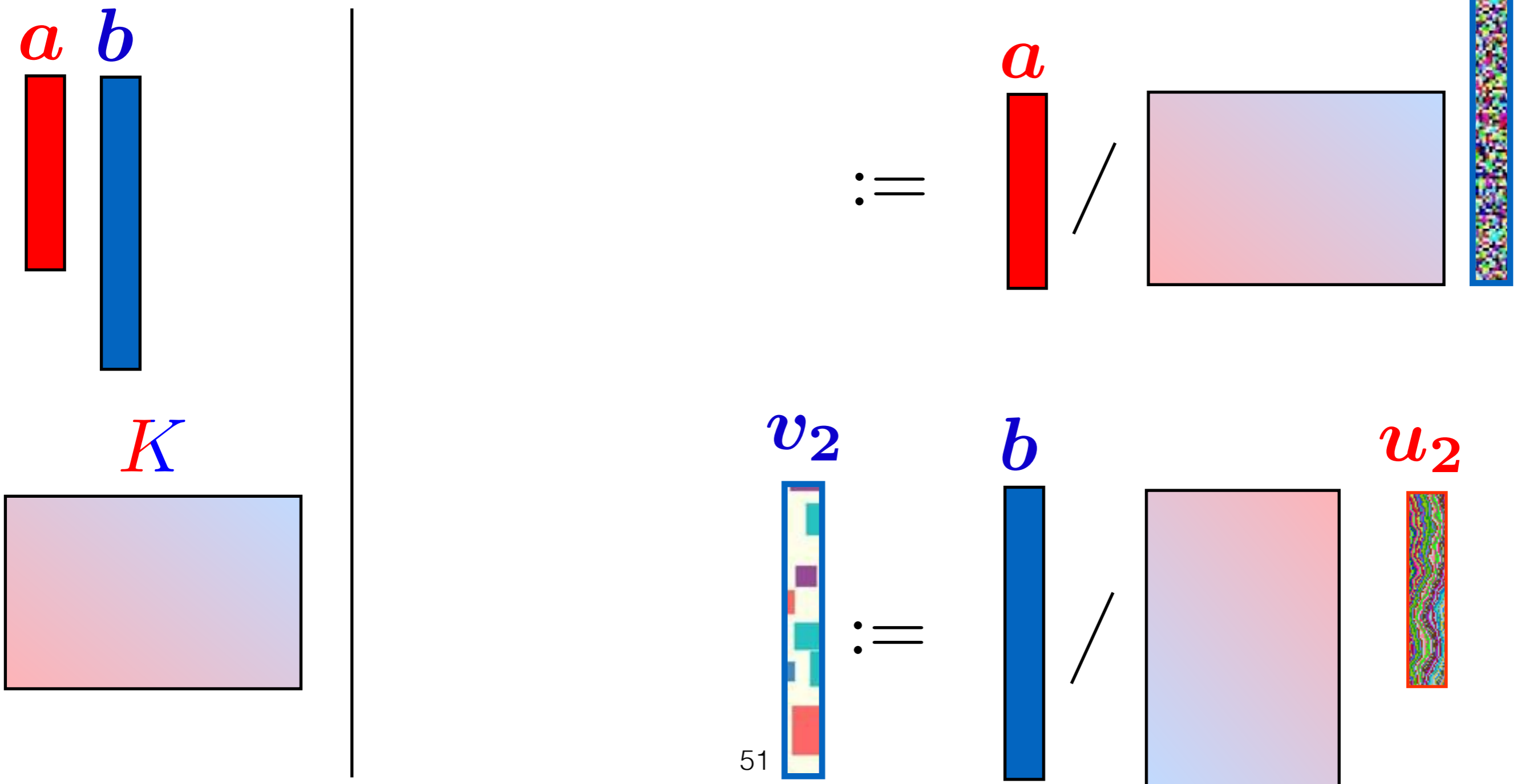
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

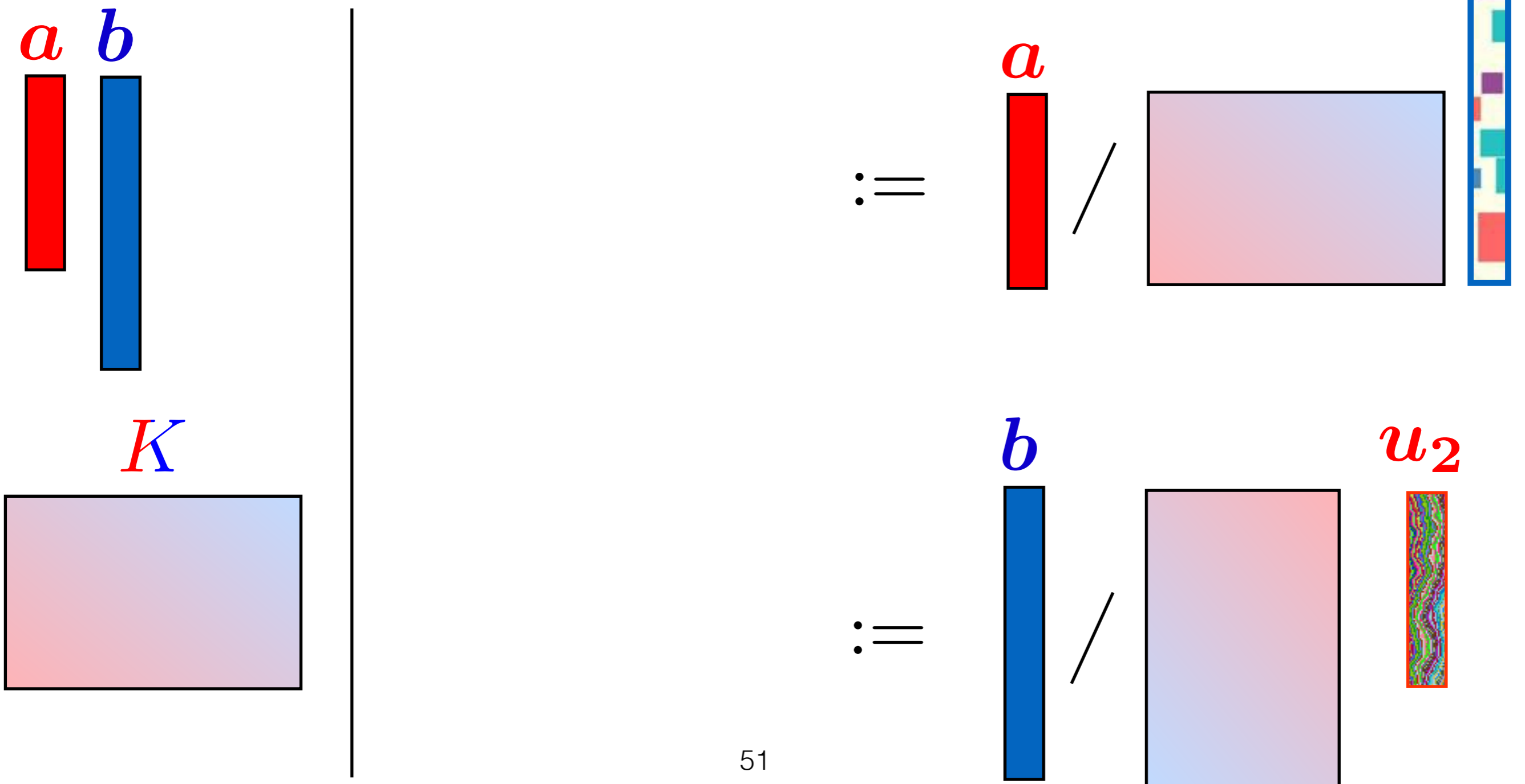
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

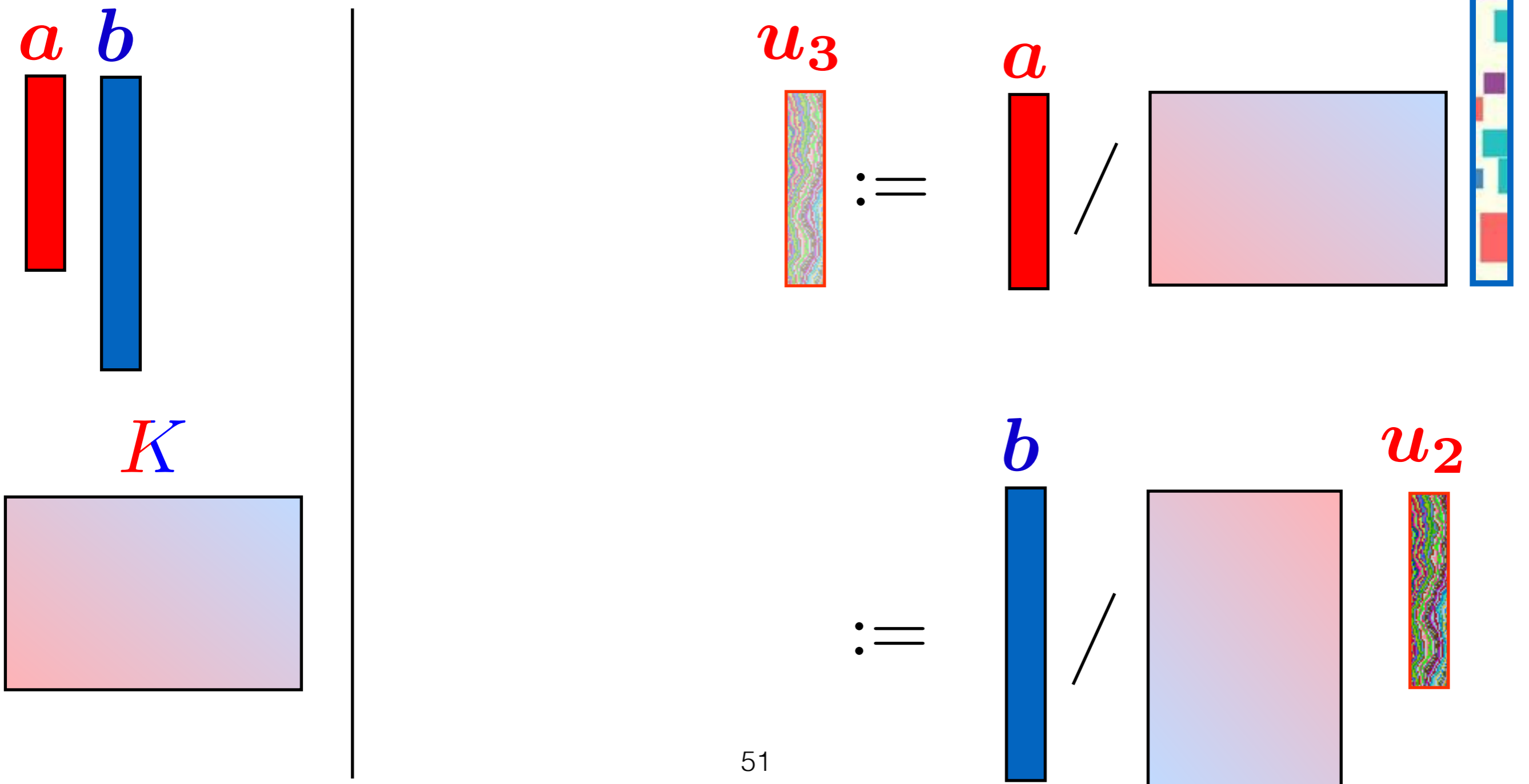
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

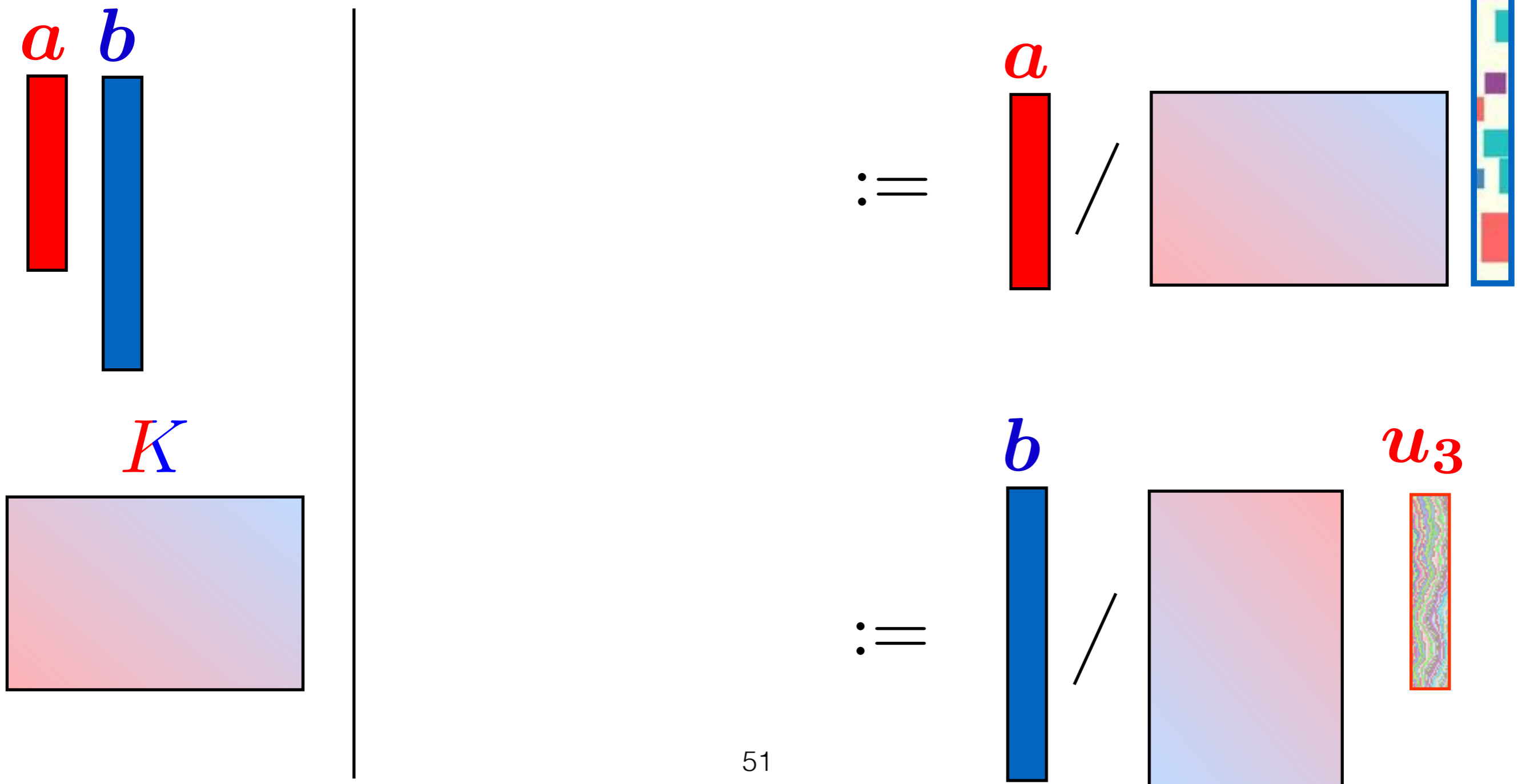
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations for (u, v)

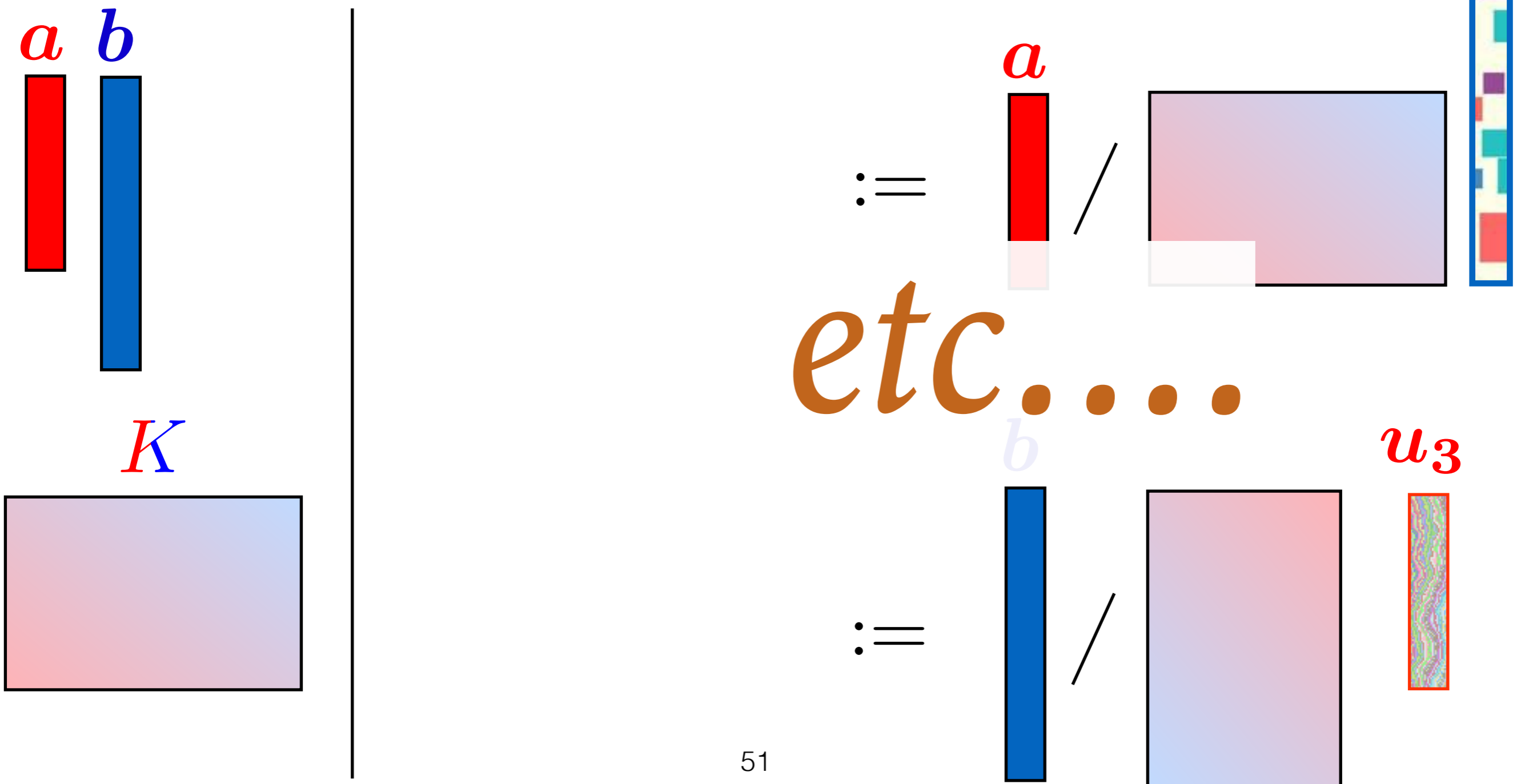
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

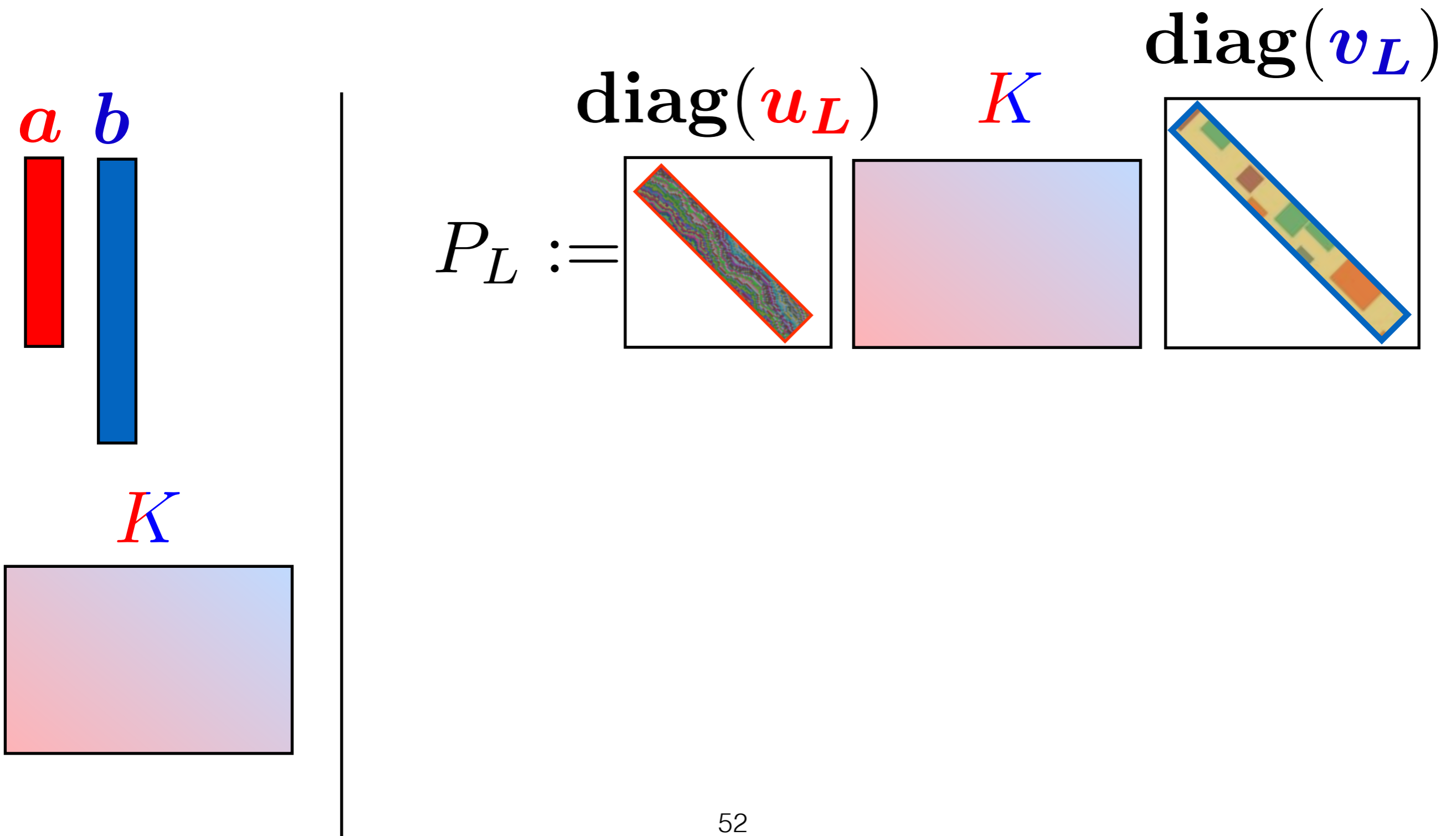
- [Sinkhorn'64] fixed-point iterations for (u, v)

$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



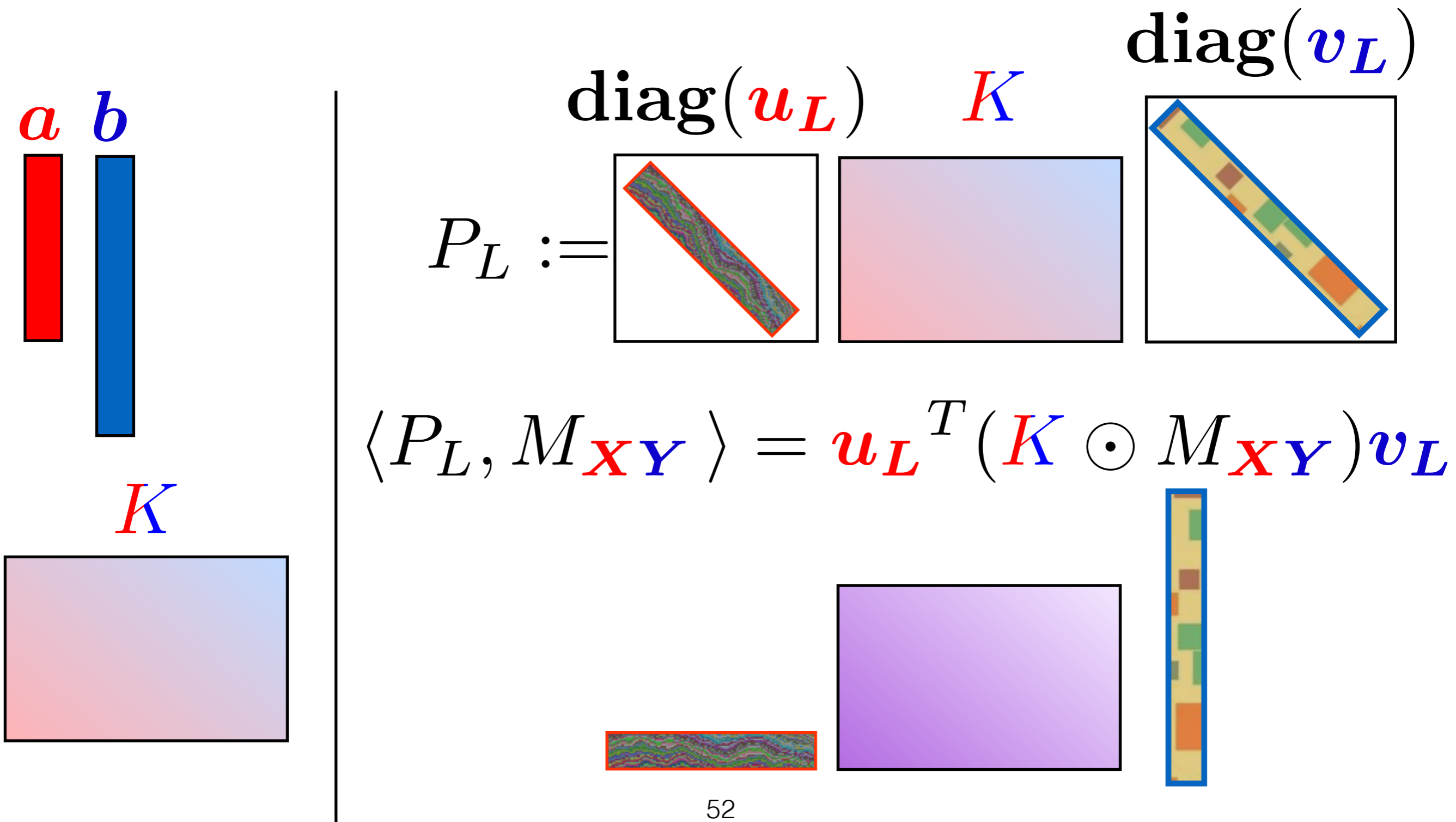
Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations.



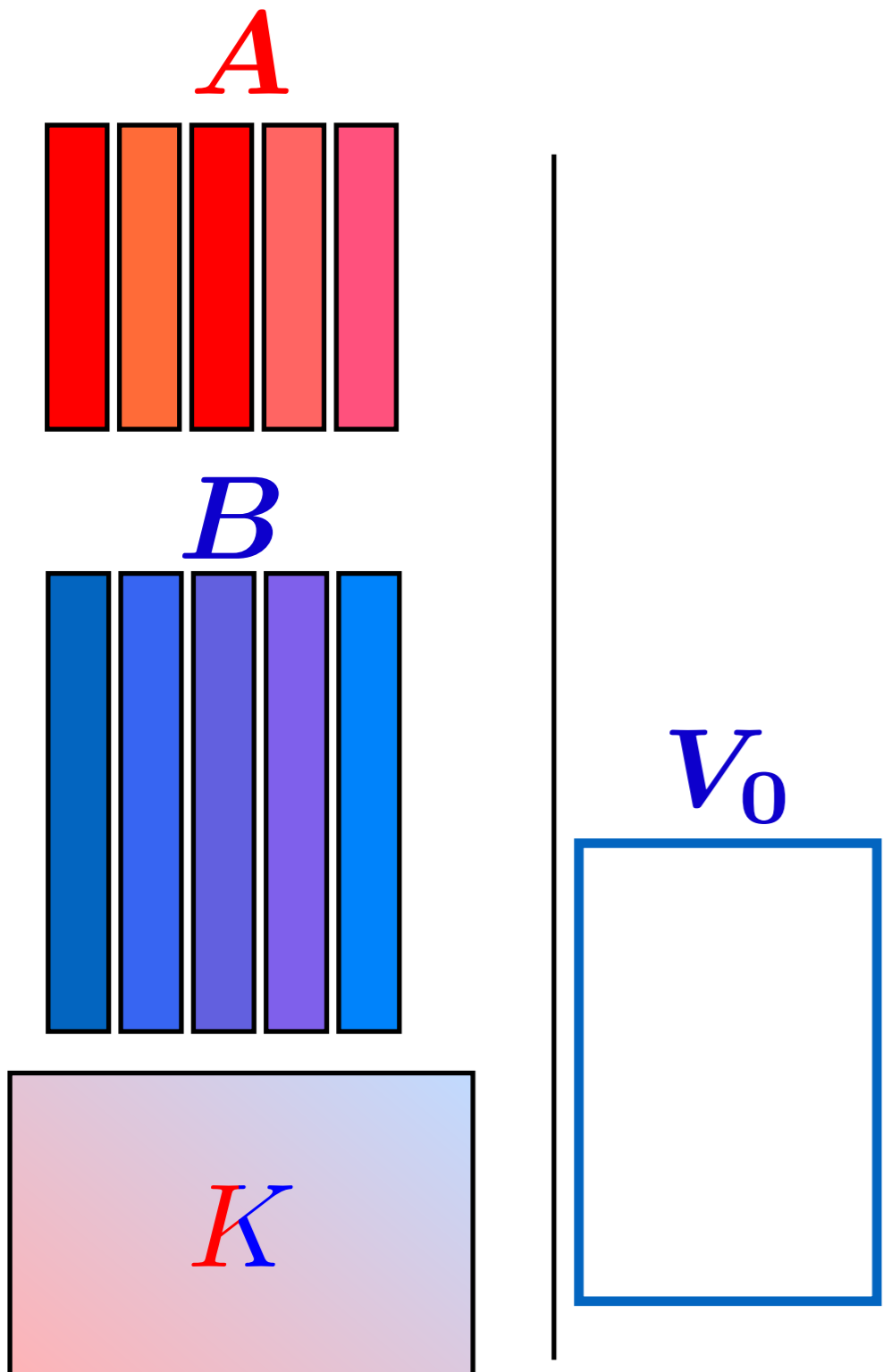
Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations.



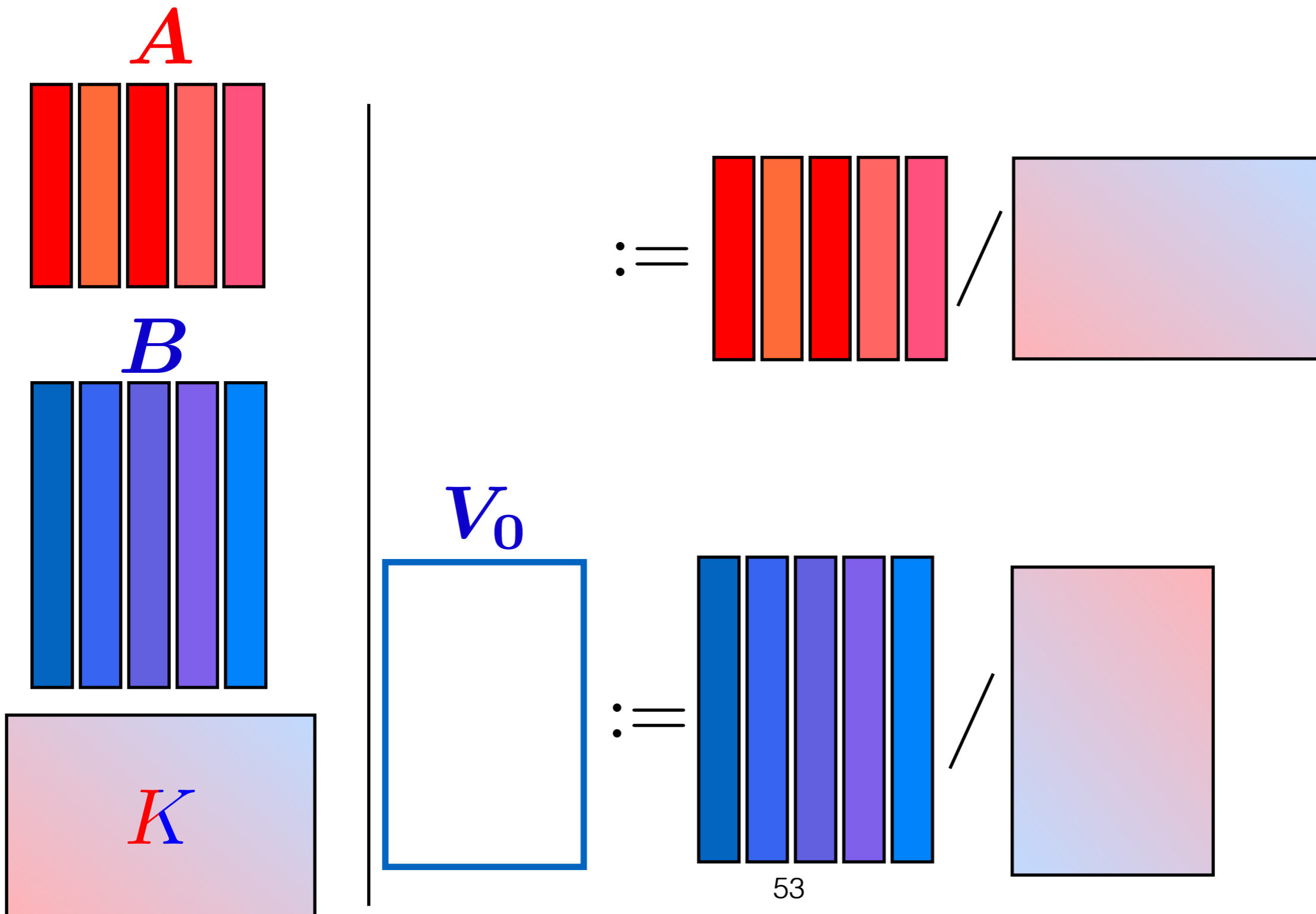
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



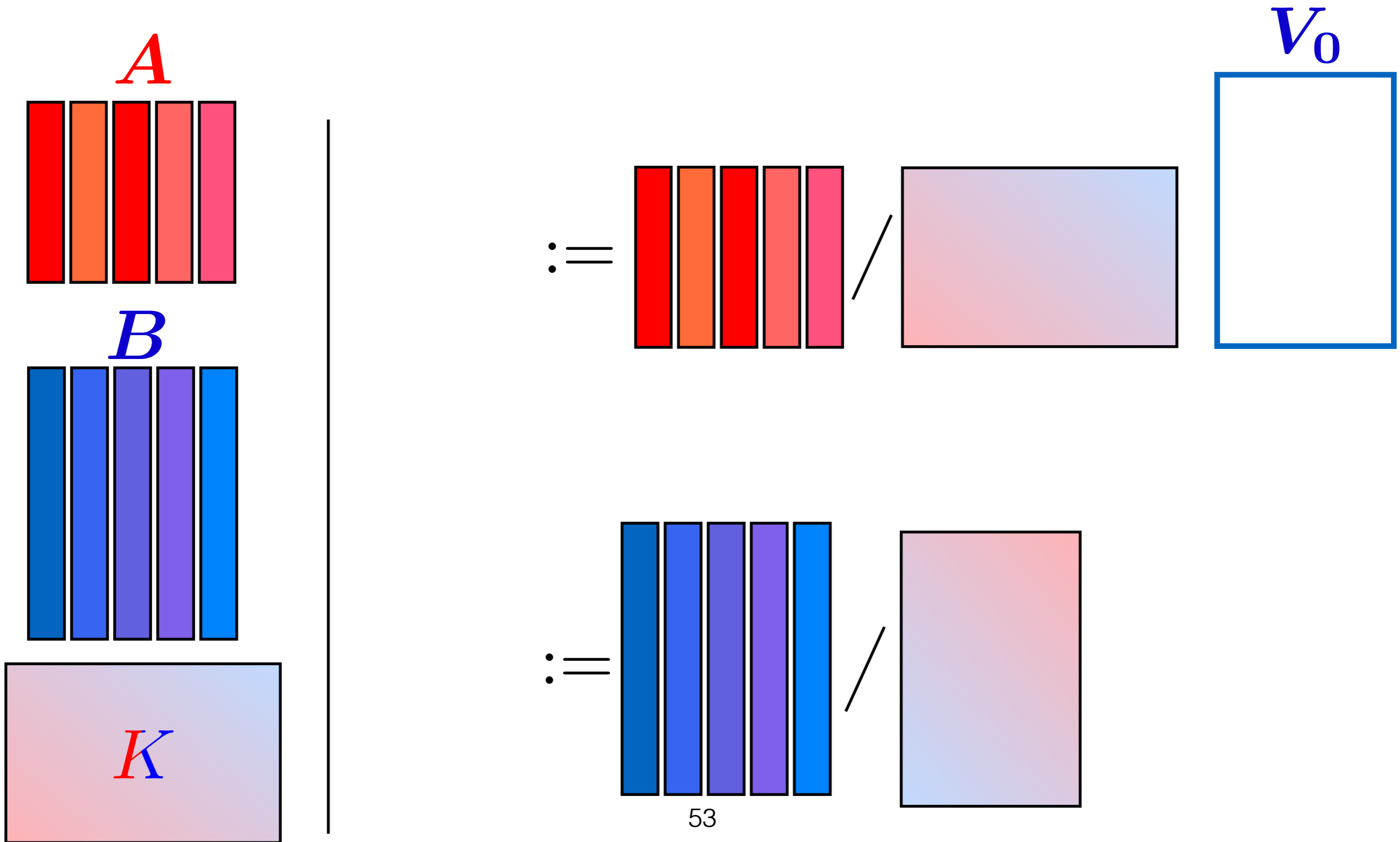
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



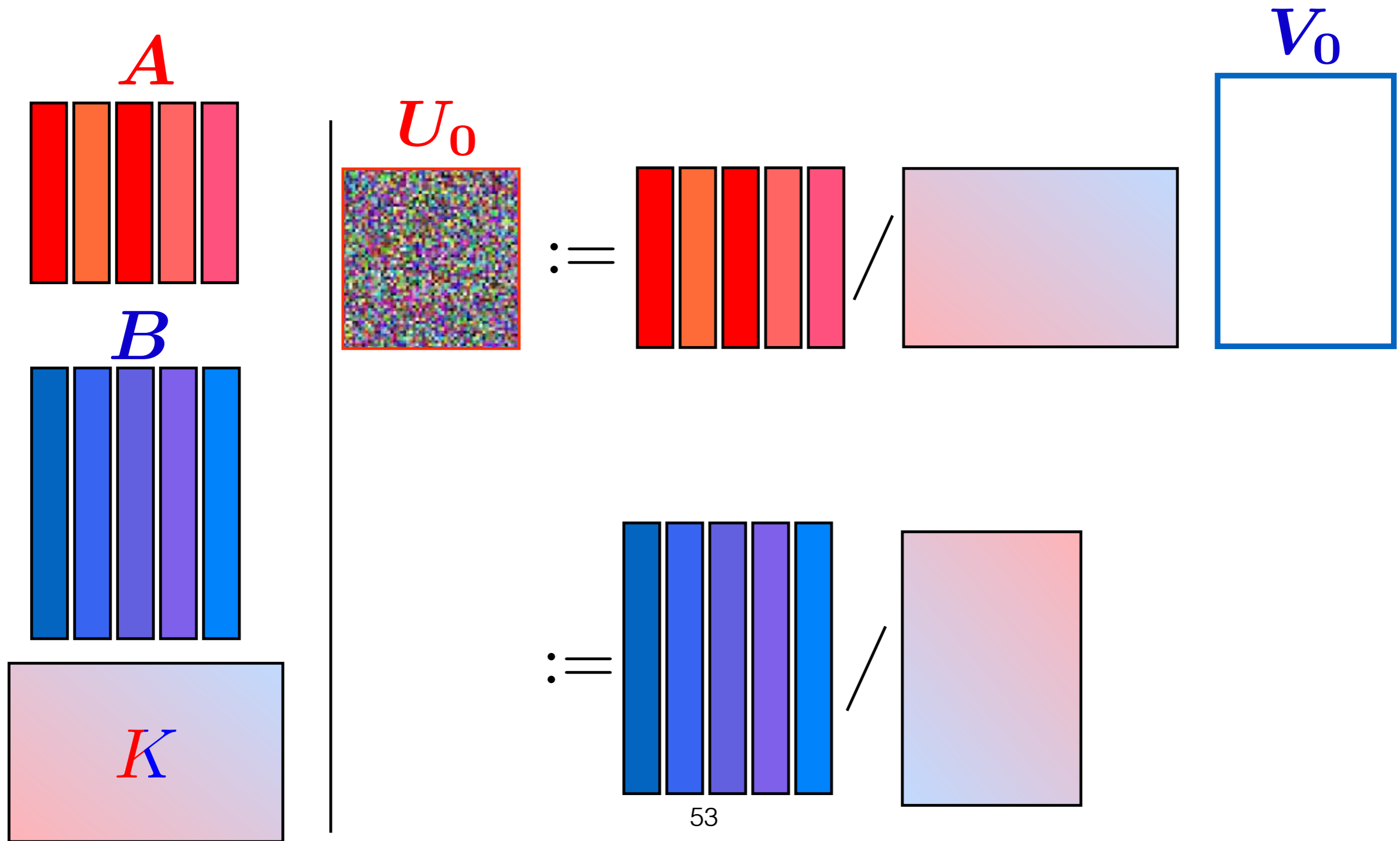
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



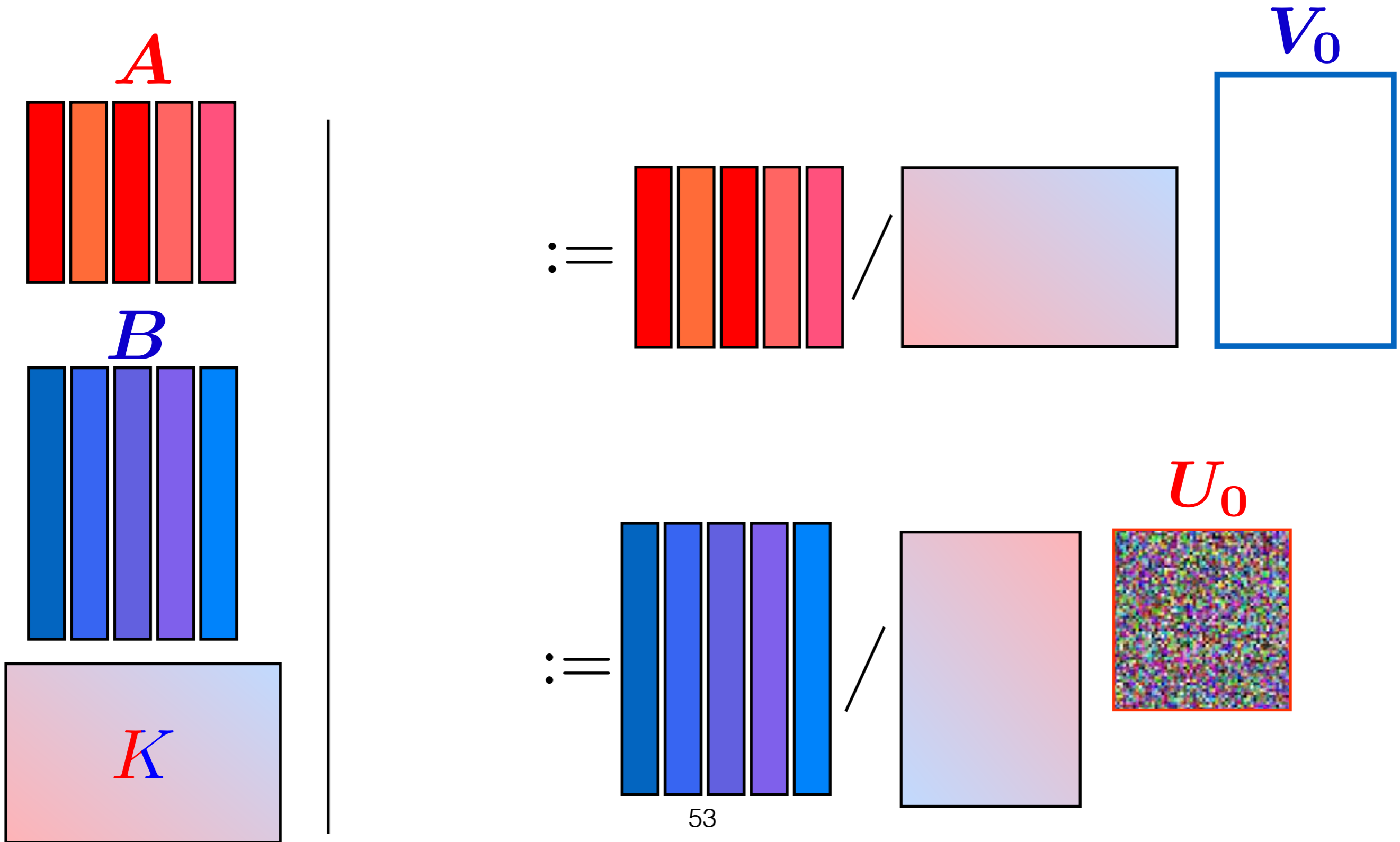
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



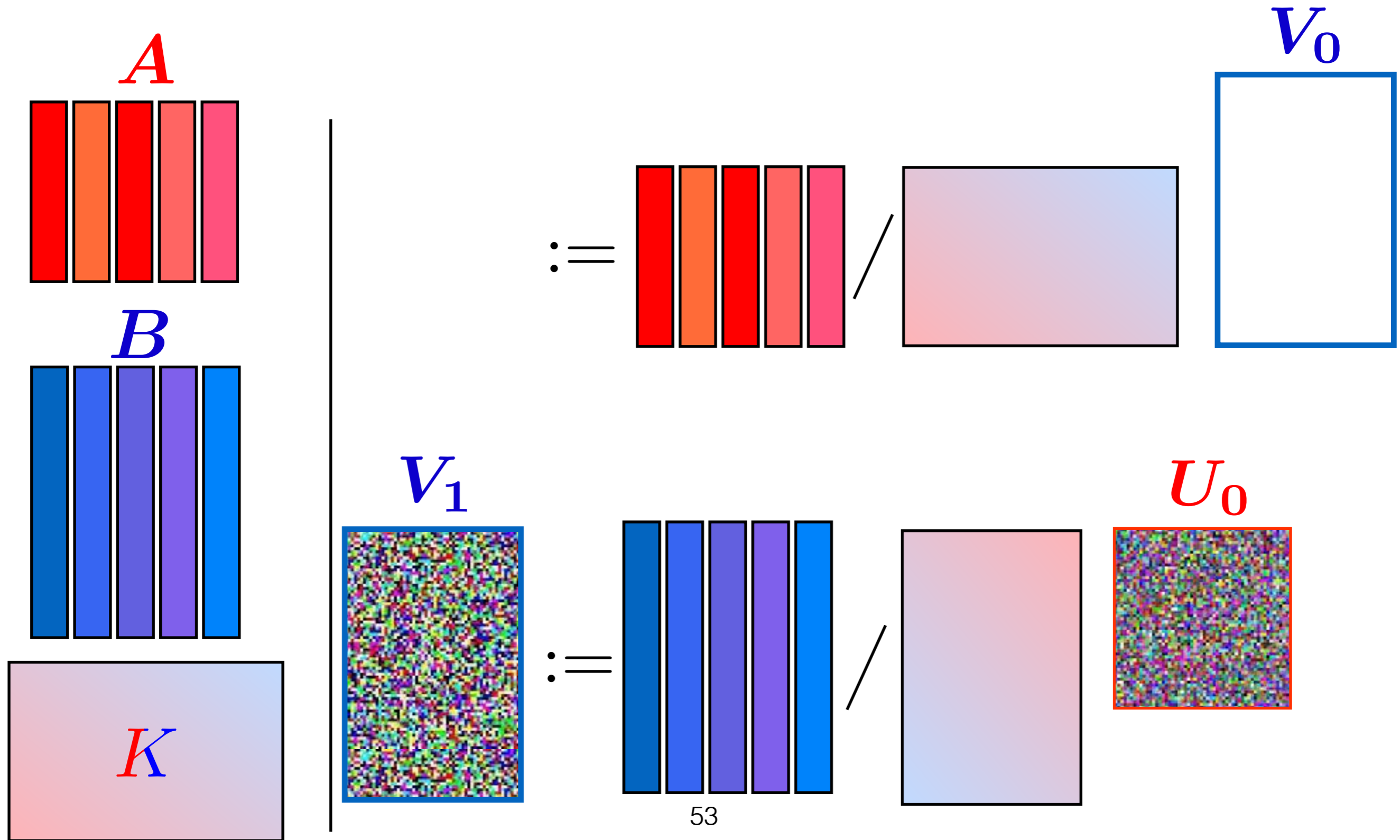
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



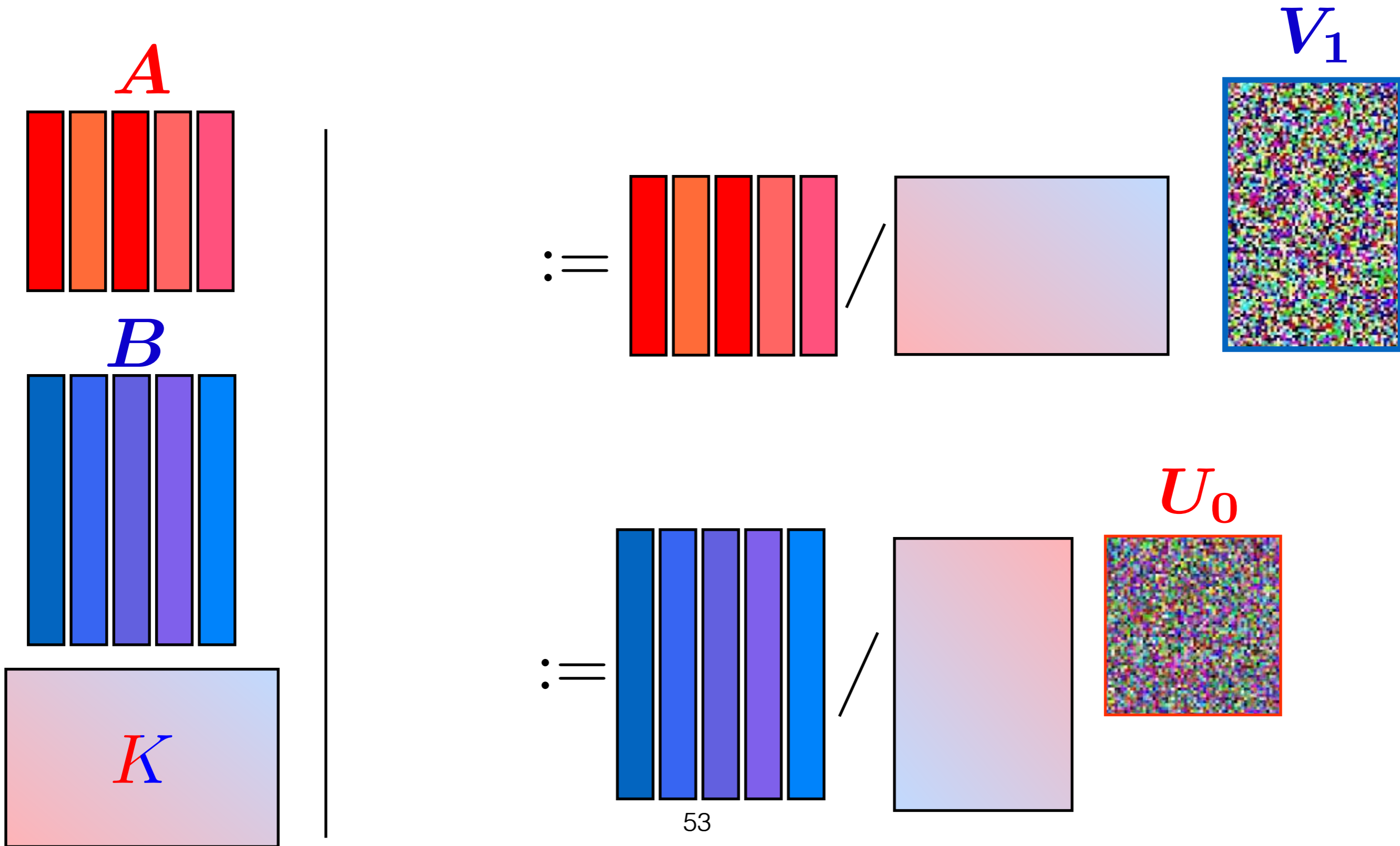
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



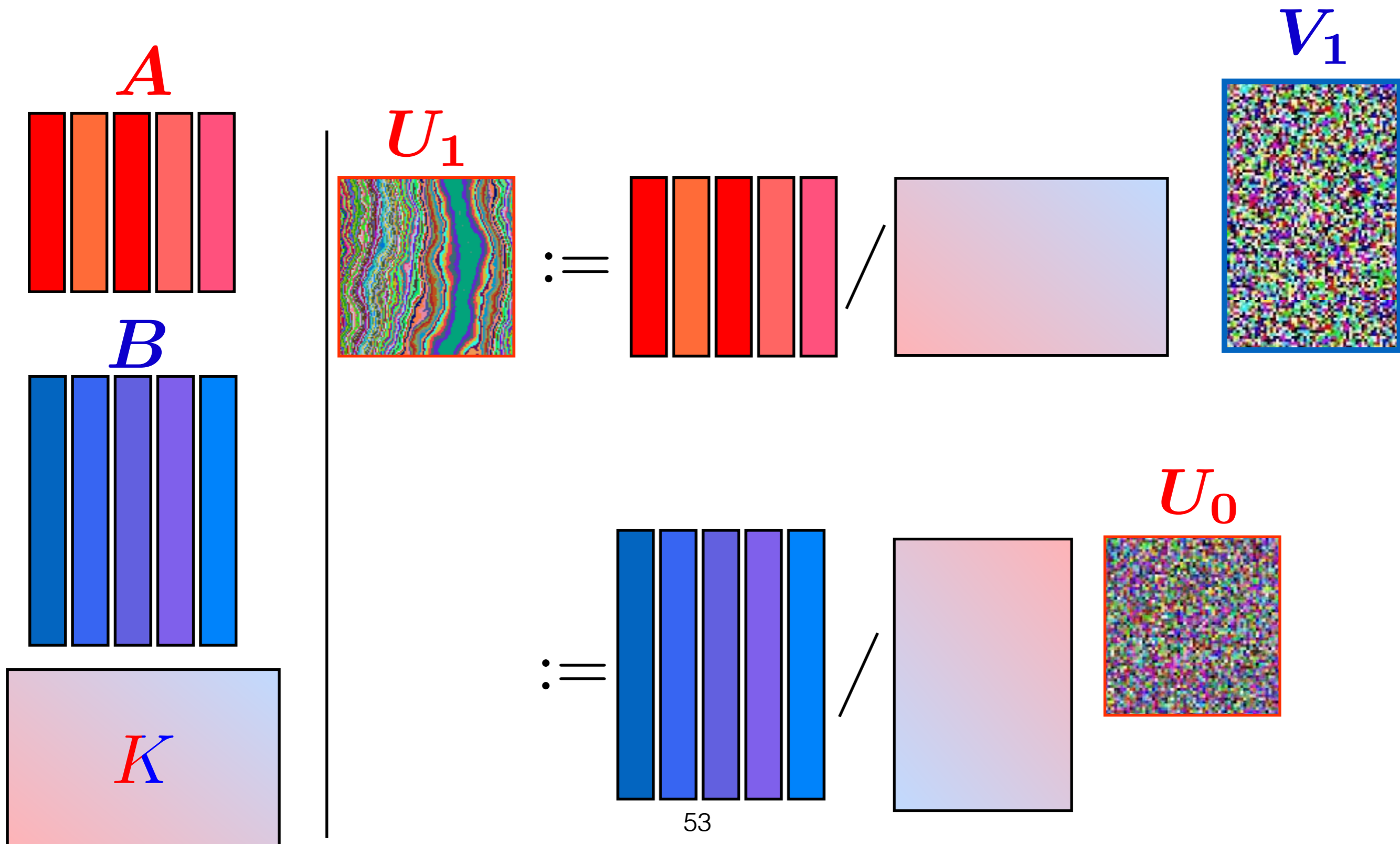
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



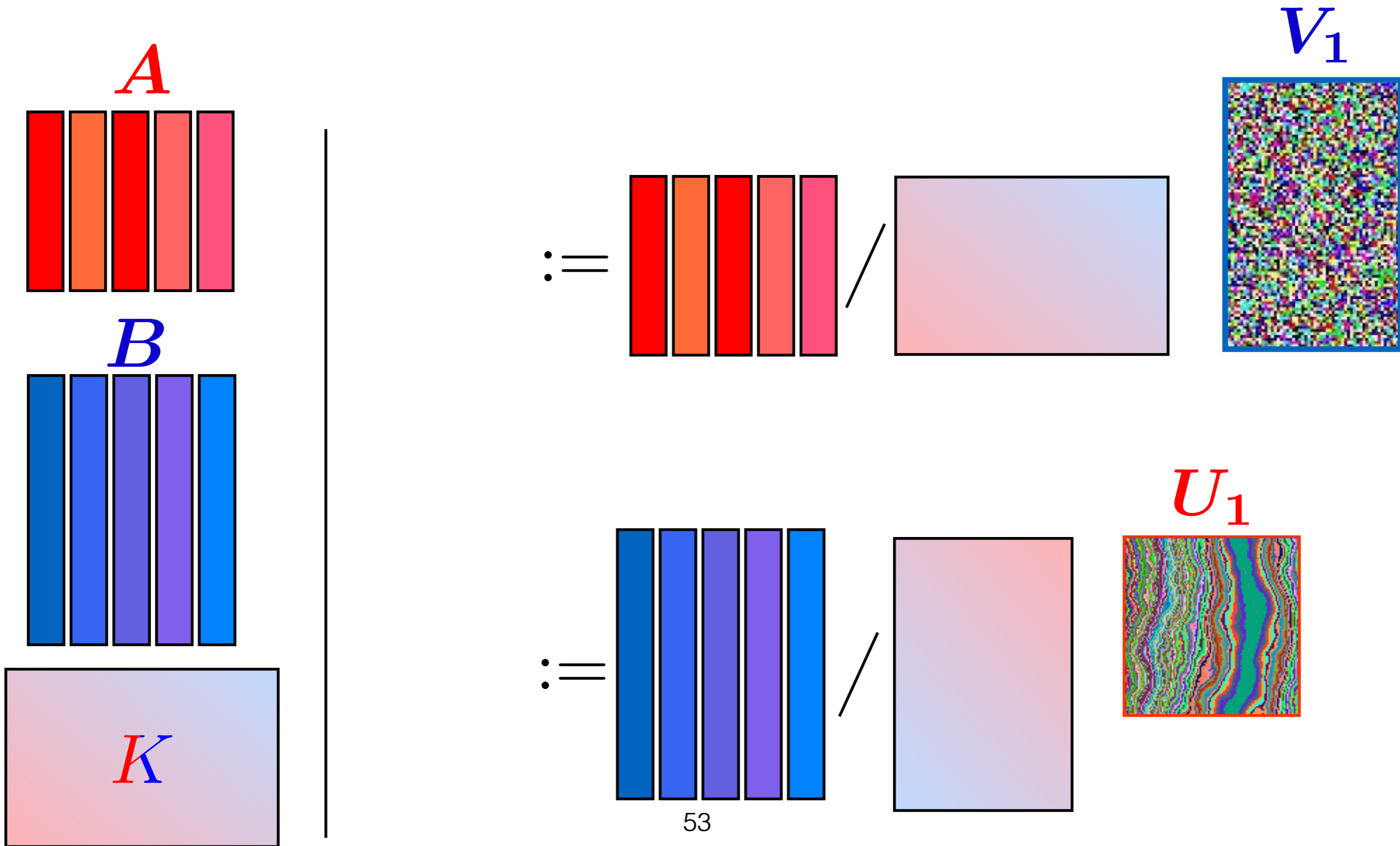
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



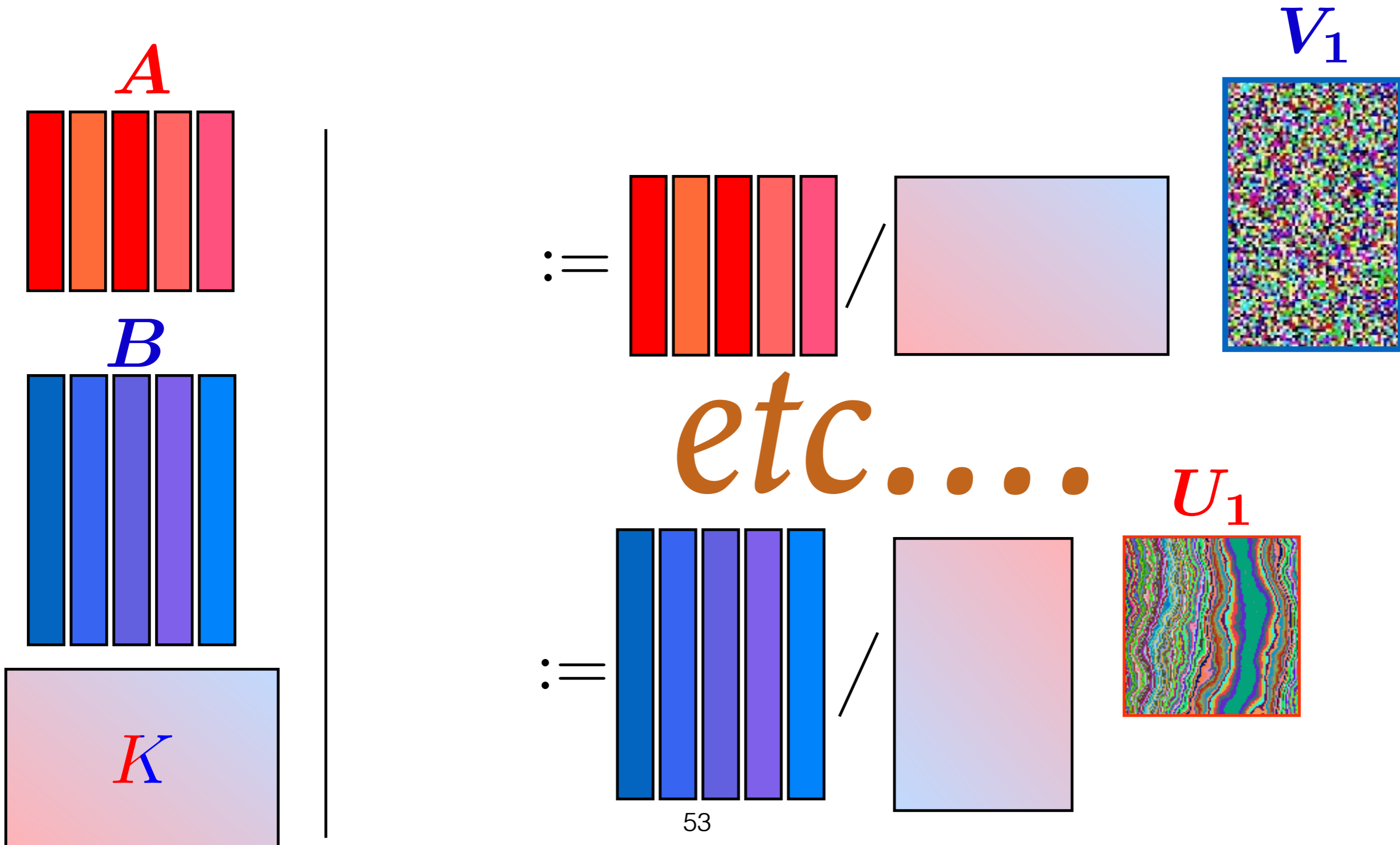
Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations

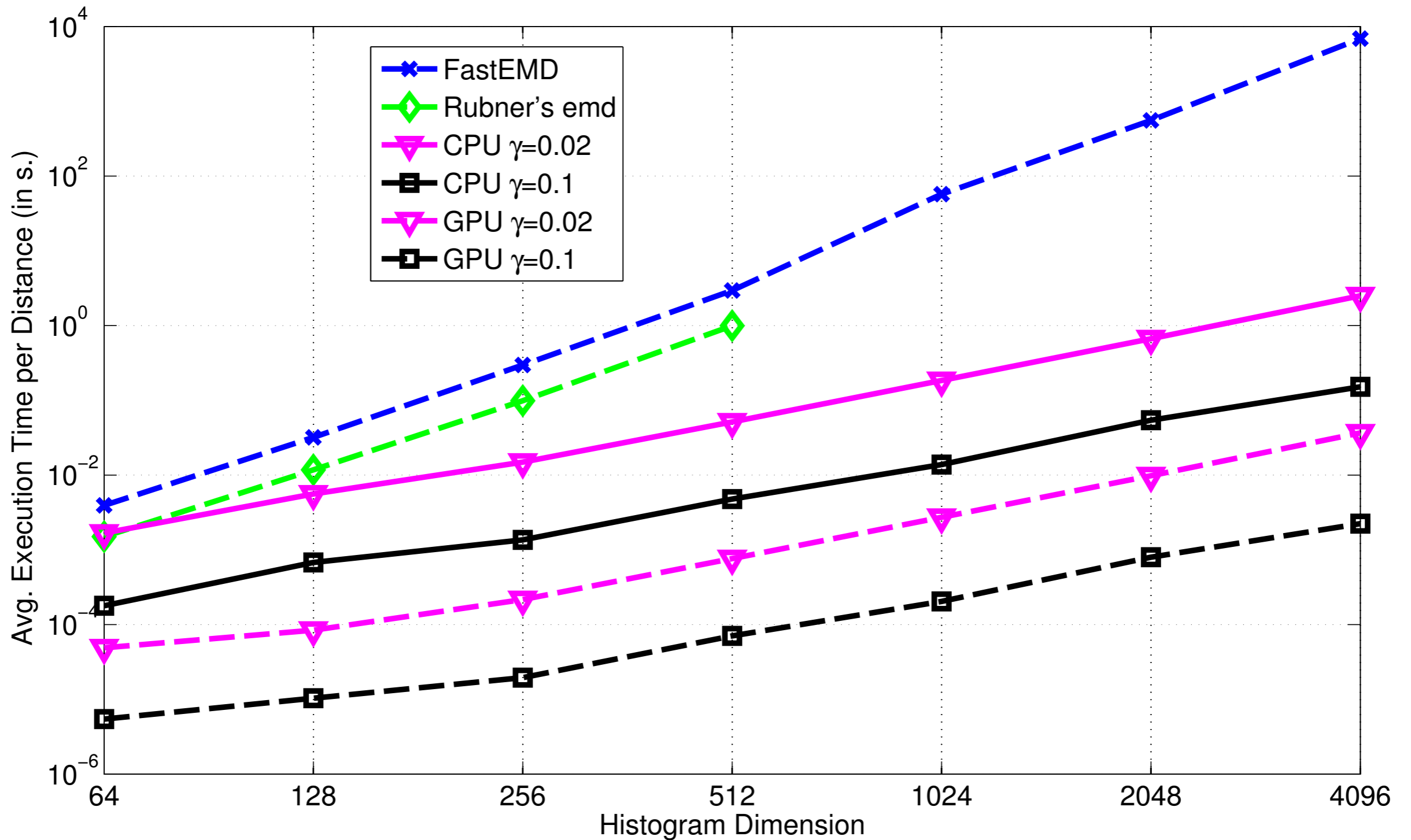


Also embarrassingly parallel

- [Sinkhorn'64] with *matrix* fixed-point iterations



Very Fast EMD Approx. Solver



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

Sinkhorn as a Dual Algorithm

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

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

REGULARIZED DISCRETE PRIMAL

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

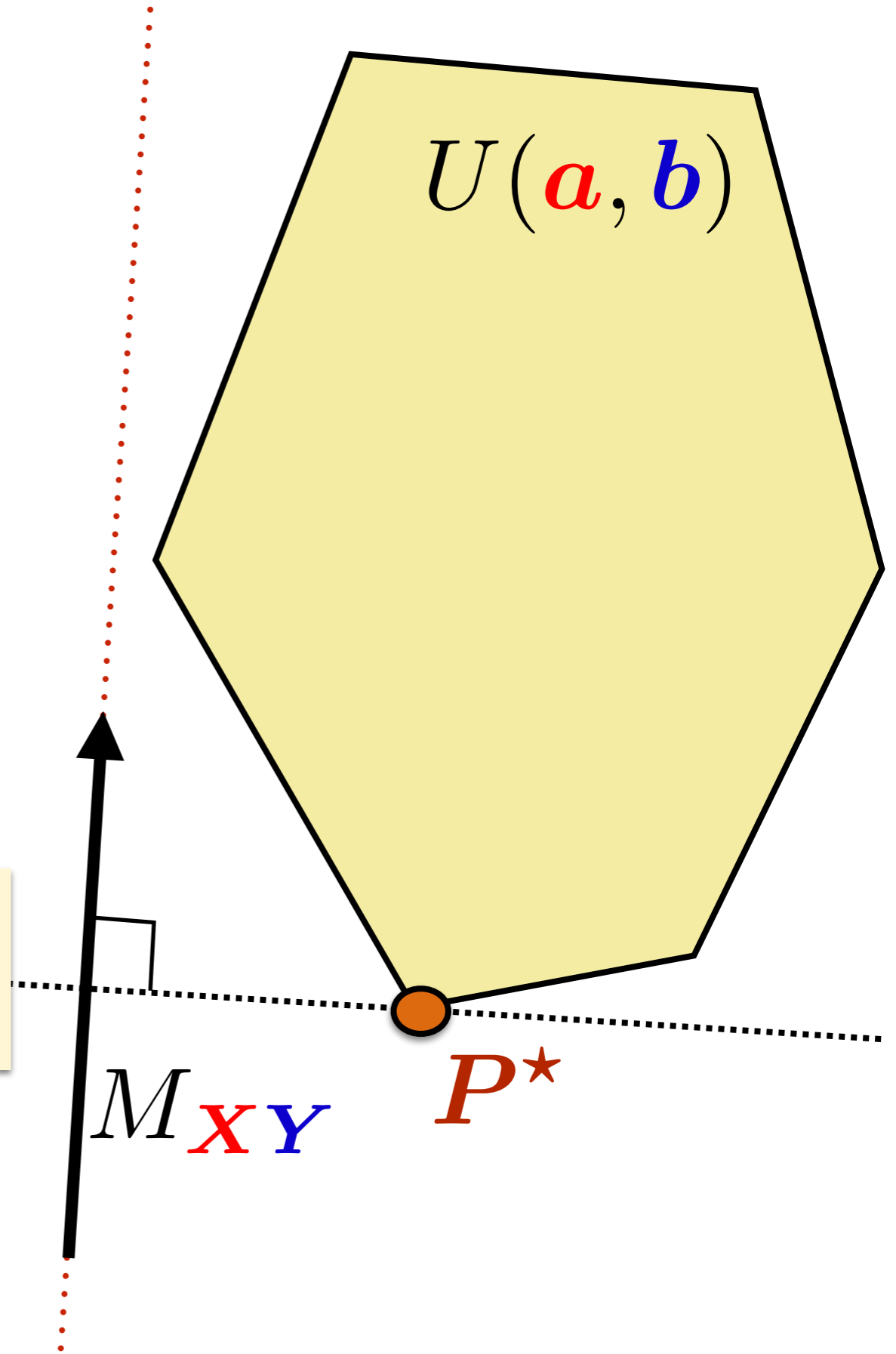
REGULARIZED DISCRETE DUAL

Sinkhorn = *Block Coordinate Ascent* on Dual

Sinkhorn in between W and MMD

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

$$W^p(\mu, \nu) = \langle P^*, M_{XY} \rangle$$

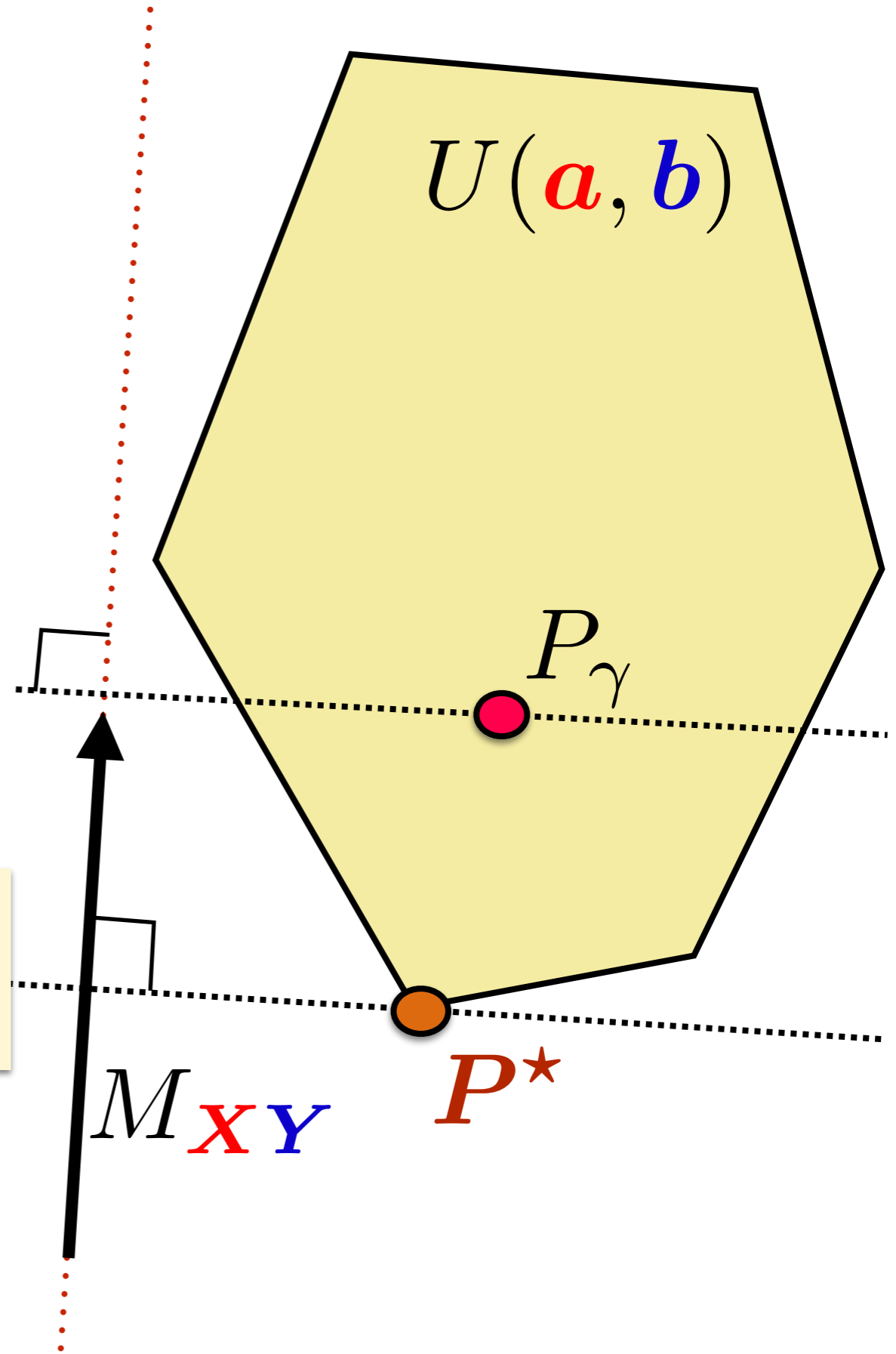


Sinkhorn in between W and MMD

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

$$W_\gamma(\mu, \nu) = \langle P_\gamma, M_{XY} \rangle$$

$$W^p(\mu, \nu) = \langle P^*, M_{XY} \rangle$$



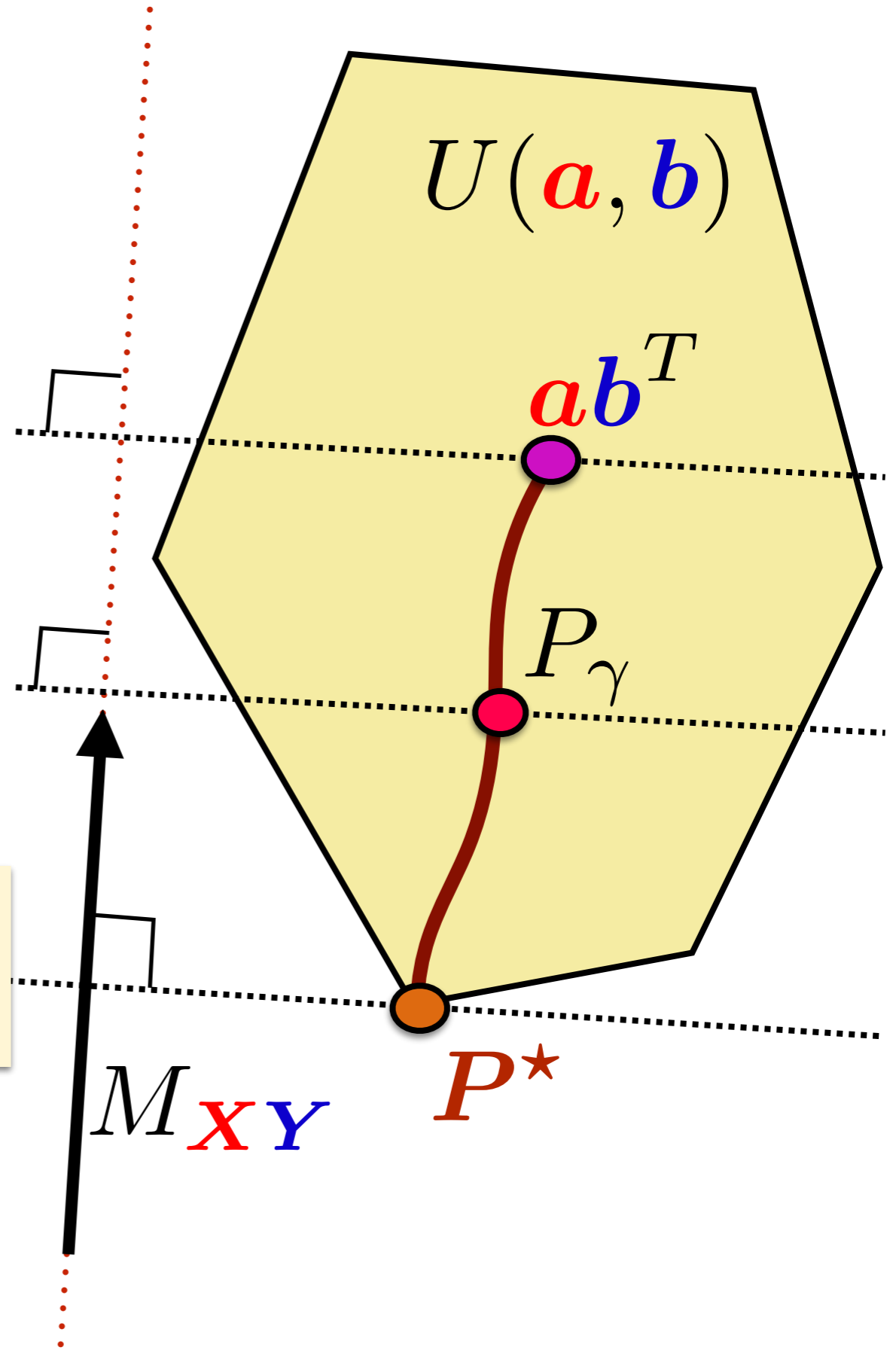
Sinkhorn in between W and MMD

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

$$\mathcal{E}(\mu, \nu) = \langle ab^T, M_{XY} \rangle$$

$$W_\gamma(\mu, \nu) = \langle P_\gamma, M_{XY} \rangle$$

$$W^p(\mu, \nu) = \langle P^*, M_{XY} \rangle$$



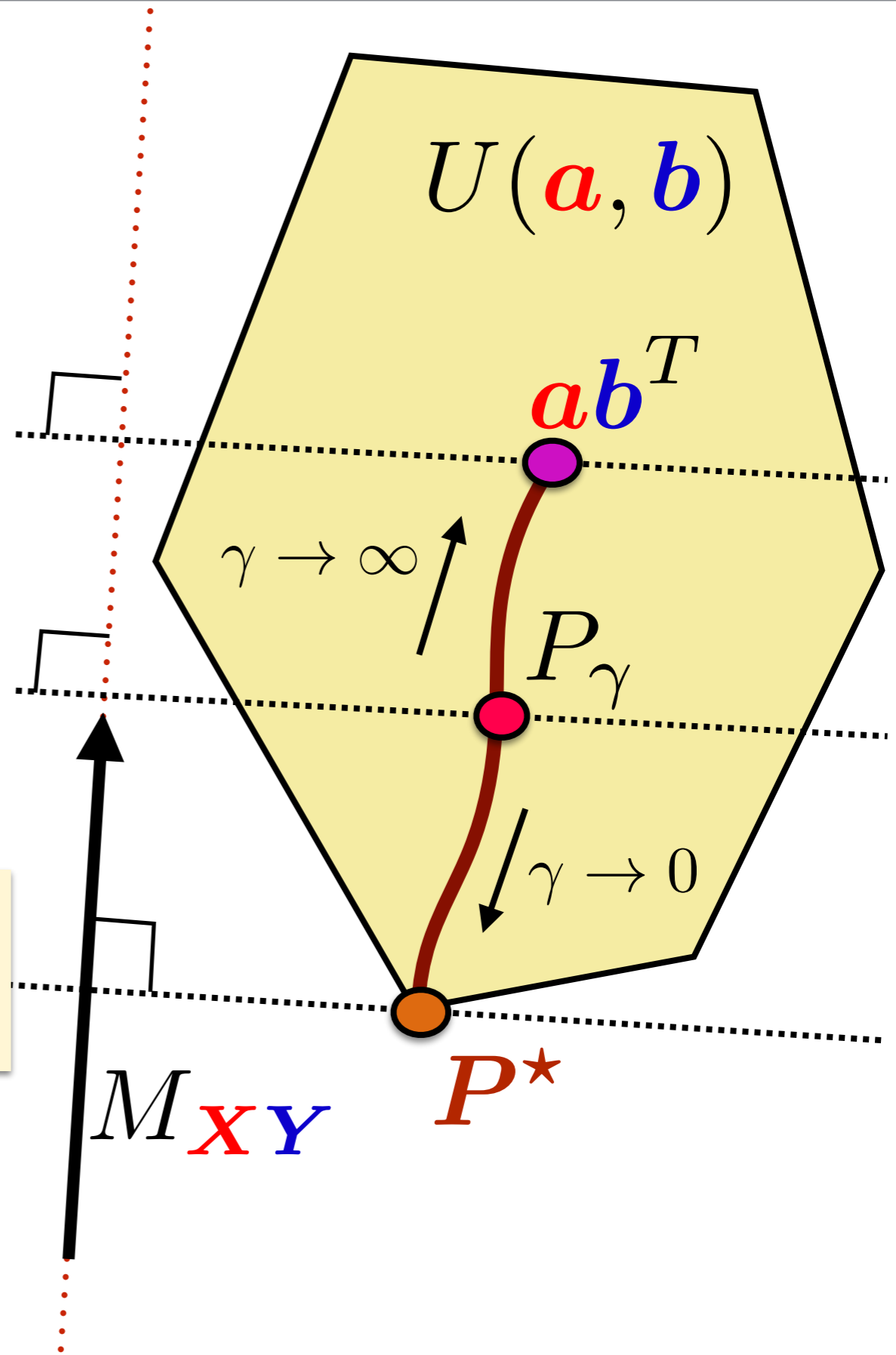
Sinkhorn in between W and MMD

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

$$\mathcal{E}(\mu, \nu) = \langle ab^T, M_{XY} \rangle$$

$$W_\gamma(\mu, \nu) = \langle P_\gamma, M_{XY} \rangle$$

$$W^p(\mu, \nu) = \langle P^*, M_{XY} \rangle$$



Sinkhorn in between W and MMD

$$\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \mathbf{ab}^T, M_{\mathbf{XY}} \rangle$$

$$MMD(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2} (\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

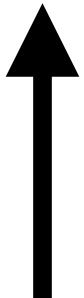
$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle P_\gamma, M_{\mathbf{XY}} \rangle$$

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

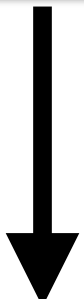
$$W^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \mathbf{P}^*, M_{\mathbf{XY}} \rangle$$

Sinkhorn in between W and MMD

$$MMD(\boldsymbol{\mu}, \boldsymbol{\nu}) = \mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(\mathcal{E}(\boldsymbol{\mu}, \boldsymbol{\mu}) + \mathcal{E}(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$\gamma \rightarrow \infty$ 

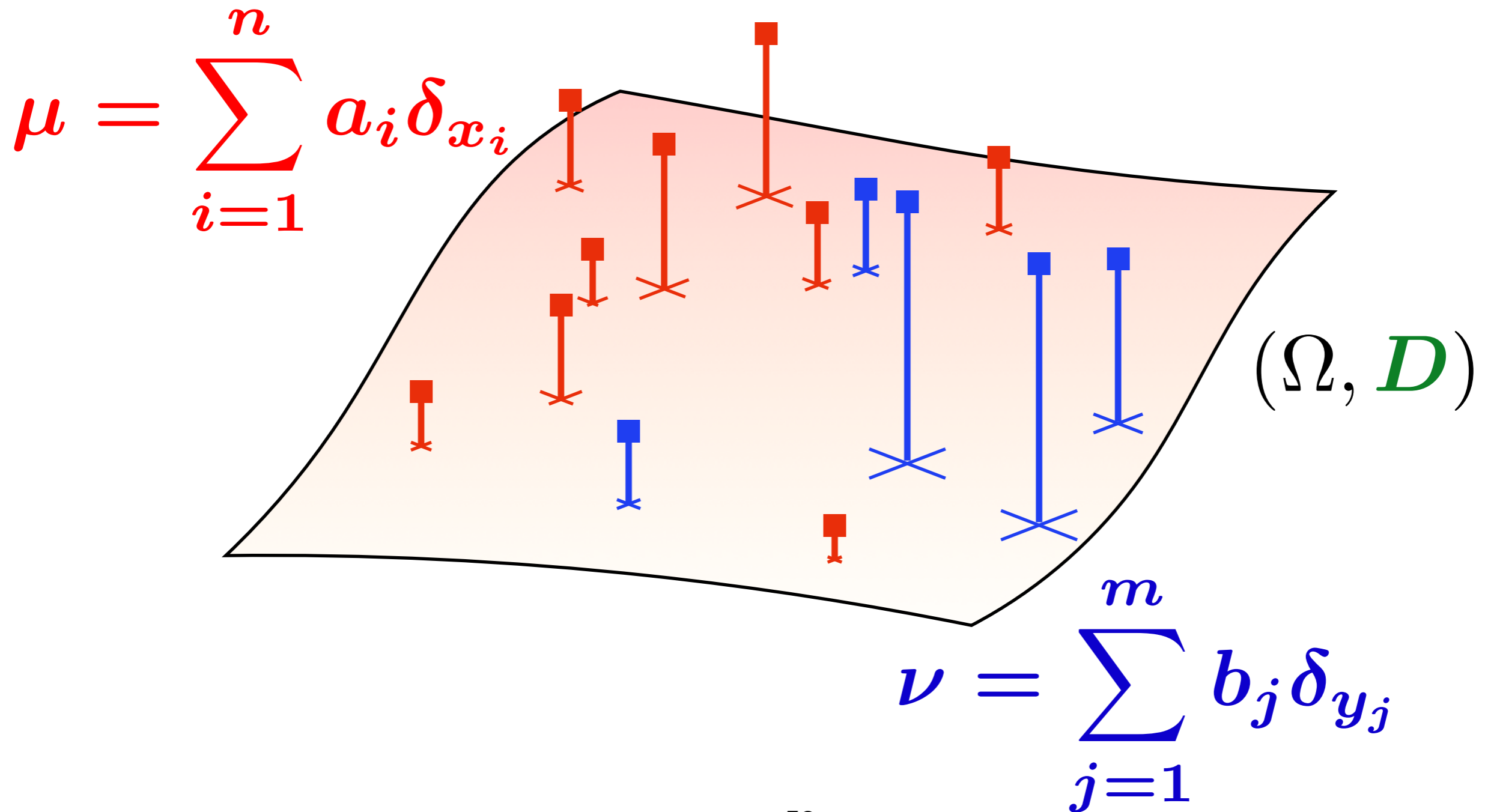
$$\bar{W}_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) - \frac{1}{2}(W_\gamma(\boldsymbol{\mu}, \boldsymbol{\mu}) + W_\gamma(\boldsymbol{\nu}, \boldsymbol{\nu}))$$

$\gamma \rightarrow 0$ 

$$W^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \boldsymbol{P}^*, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle$$

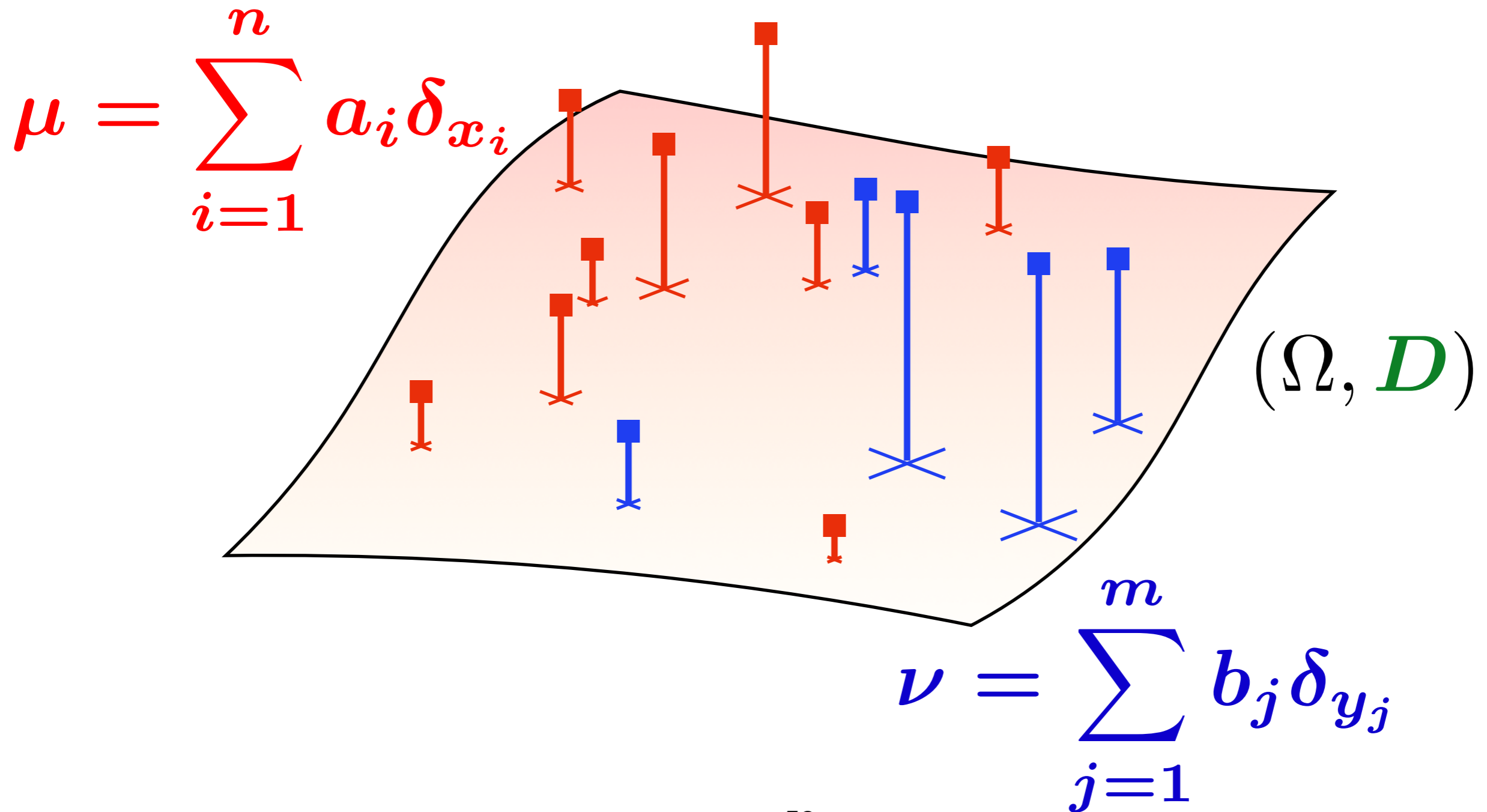
Differentiability of W

$$W((a, X), (b, Y))$$



Differentiability of W

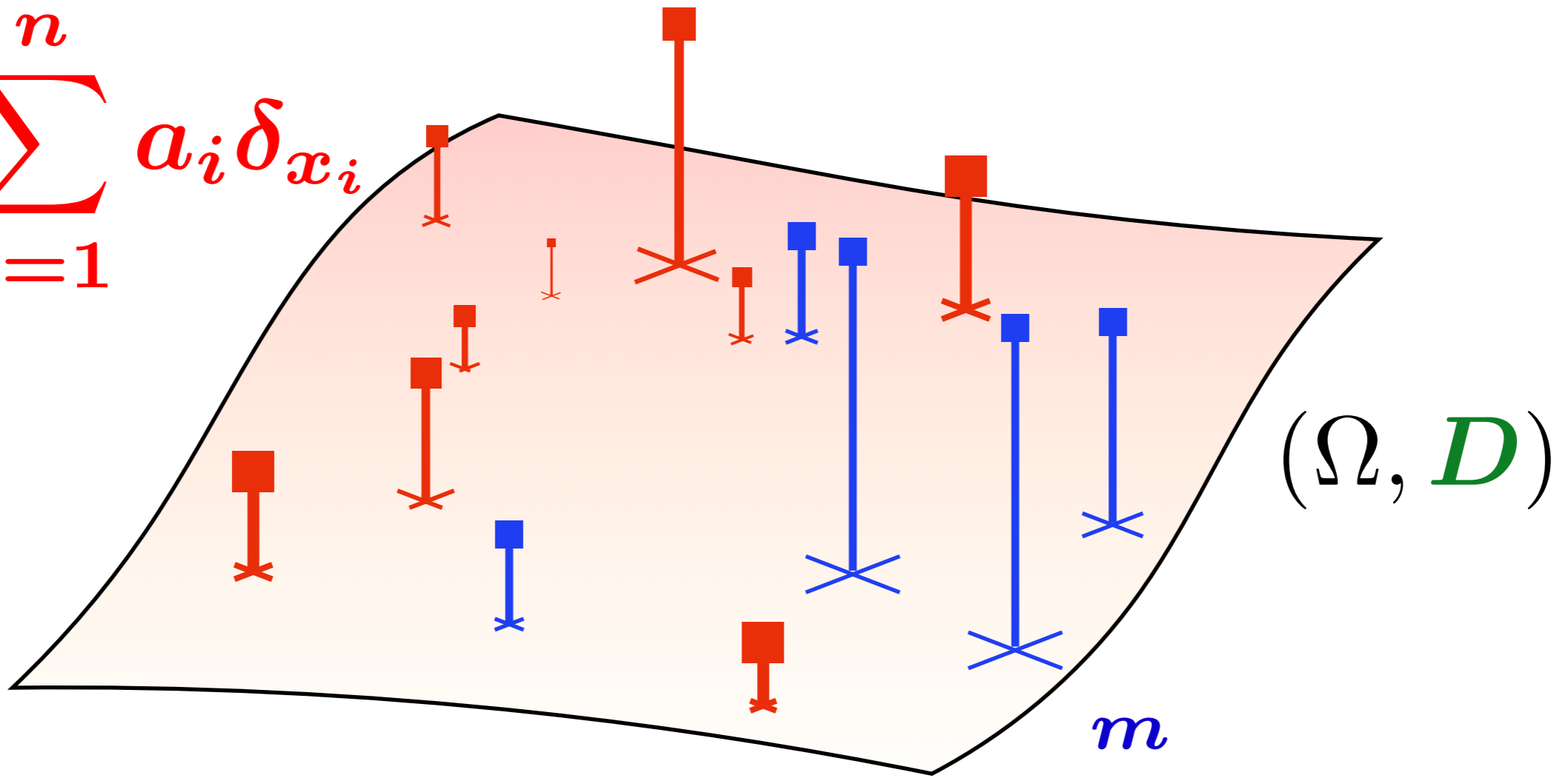
$$W((a + \Delta a, X), (b, Y)) = W((a, X), (b, Y)) + ??$$



Differentiability of W

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

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

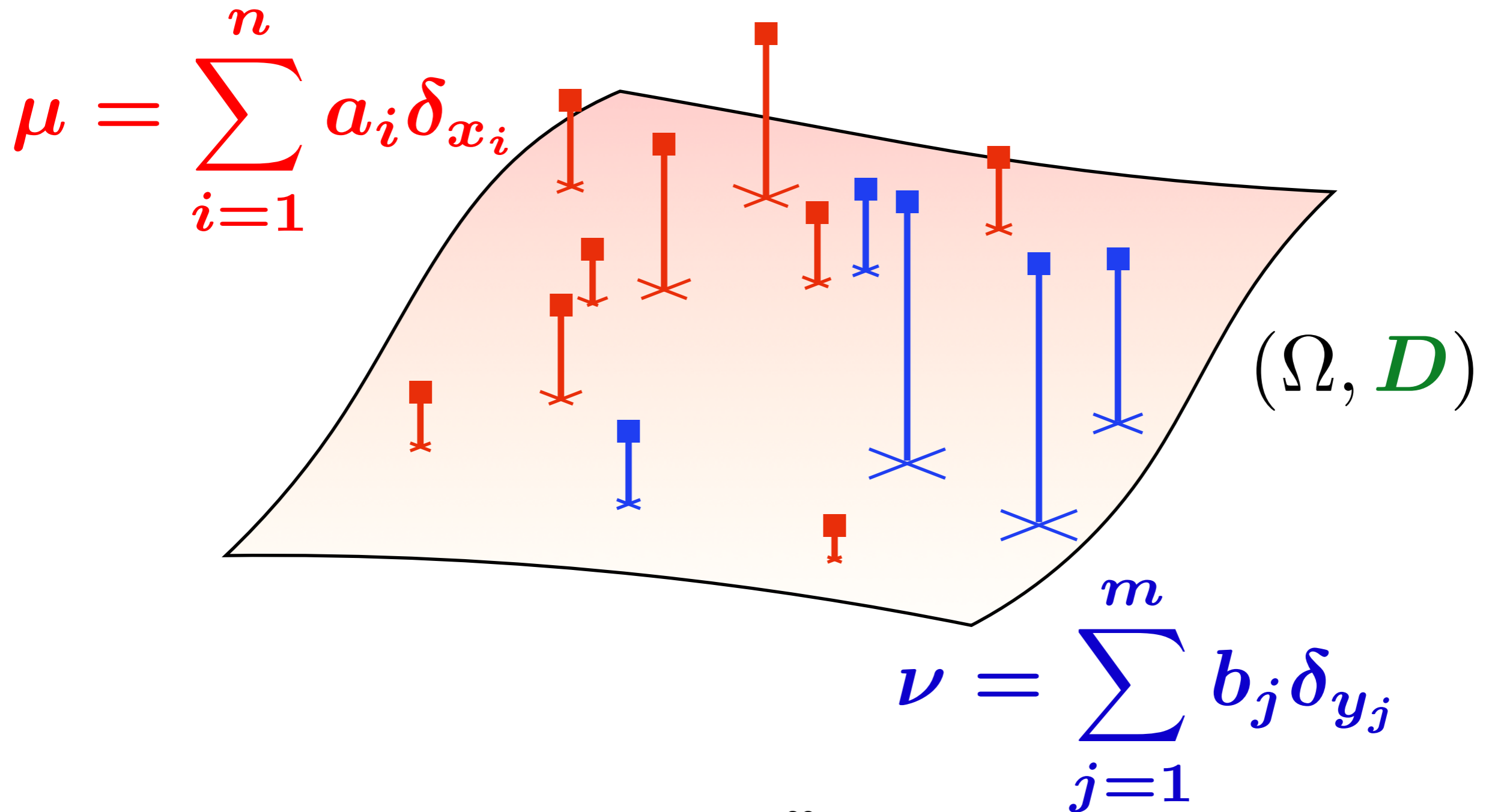


$$a \leftarrow a + \Delta a$$

$$\nu = \sum_{j=1}^m b_j \delta y_j$$

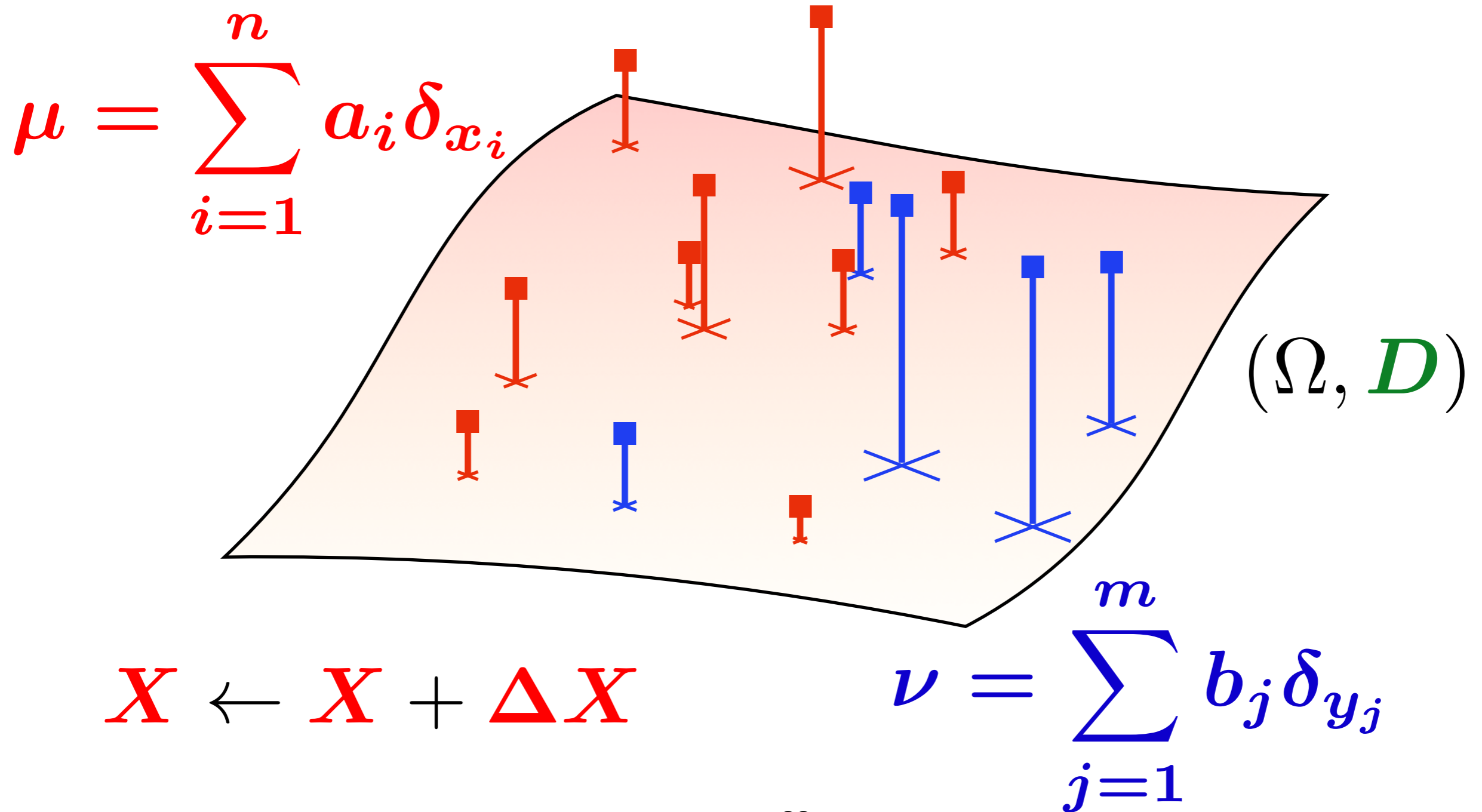
Sinkhorn \rightsquigarrow *Differentiability*

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



Sinkhorn \rightsquigarrow *Differentiability*

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



How to decrease W ? change weights

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

DUAL

Prop. $W(\boldsymbol{\mu}, \boldsymbol{\nu})$ is convex w.r.t. \mathbf{a} , $\partial_{\mathbf{a}} W = \boldsymbol{\alpha}^*$

Prop. $W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu})$ is convex and differentiable w.r.t. \mathbf{a} , $\nabla_{\mathbf{a}} W_\gamma = \gamma \log \mathbf{u}$

How to decrease W ? change locations

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P\mathbf{1}_m = \mathbf{a}, P^T\mathbf{1}_n = \mathbf{b}}} \langle P, \mathbf{1}_n \mathbf{1}_d^T X^2 + Y^{2T} \mathbf{1}_d \mathbf{1}_m - 2X^T Y \rangle$$

PRIMAL

Prop. $p = 2, \Omega = \mathbb{R}^d$. $W(\boldsymbol{\mu}, \boldsymbol{\nu})$ decreases if
 $X \leftarrow Y P^{*T} \mathbf{D}(\mathbf{a}^{-1})$.

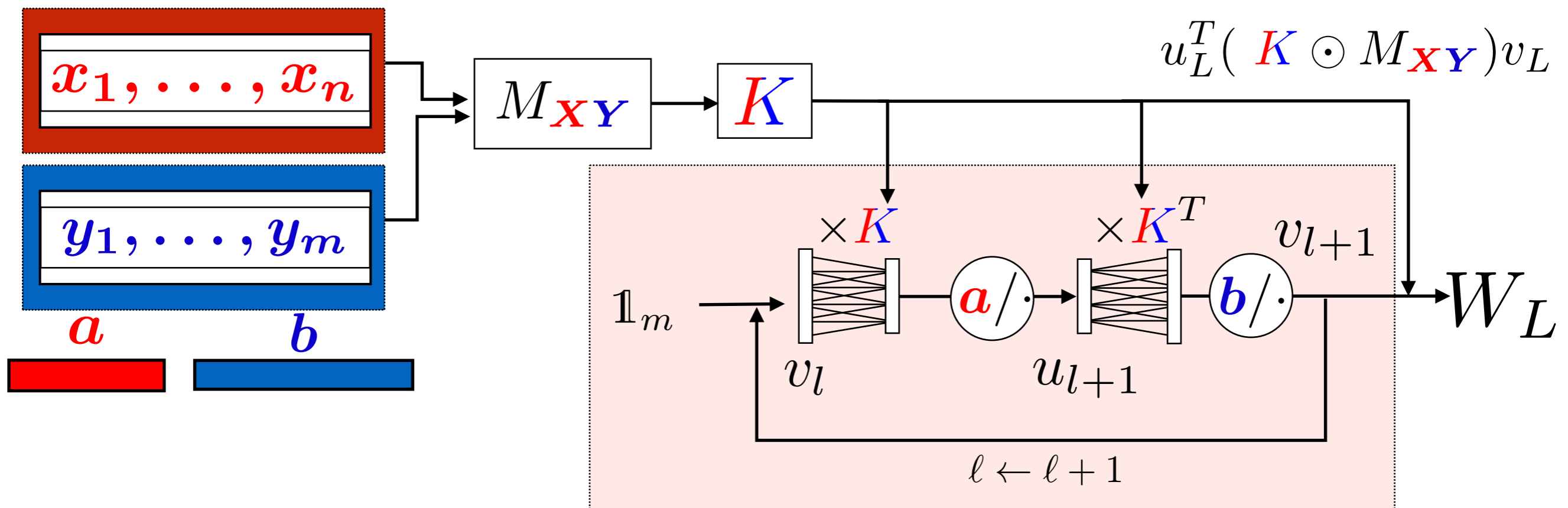
Prop. $p = 2, \Omega = \mathbb{R}^d$. $W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu})$ is differentiable w.r.t. X , with

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

Sinkhorn: A Programmer View

Def. For $L \geq 1$, define

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



Sinkhorn $l = 1, \dots, L - 1$

Sinkhorn: A Programmer View

Def. For $L \geq 1$, define

$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P}_L, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle,$$

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

The Programmer's Way

Def. For $L \geq 1$, define

$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P}_L, M_{\boldsymbol{X}\boldsymbol{Y}} \rangle,$$

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

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

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

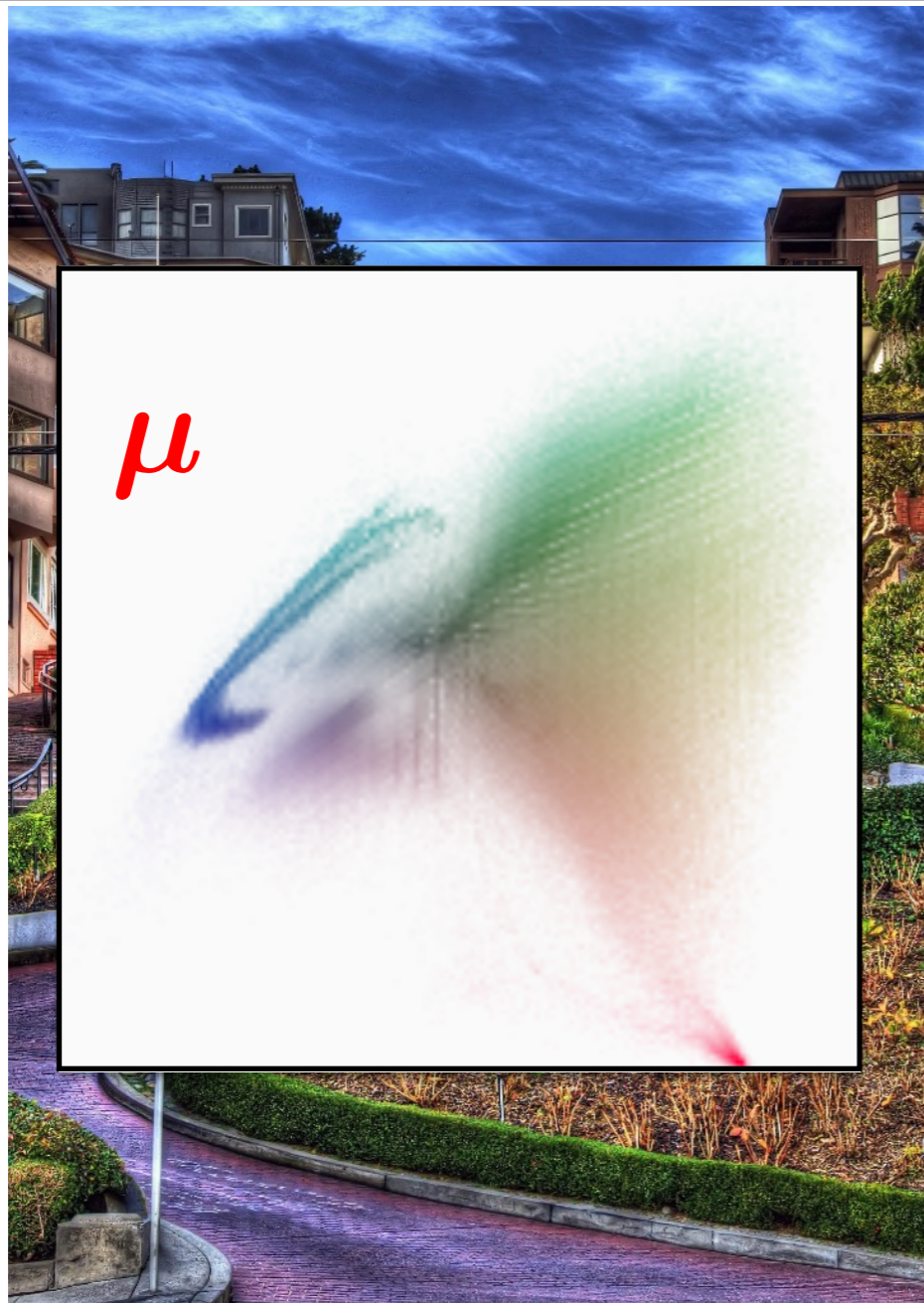
3. Applications

- Wasserstein distances for retrieval
- Wasserstein barycenters
- W for unsupervised learning
- W inverse problems
- W to learn parameters and generative models

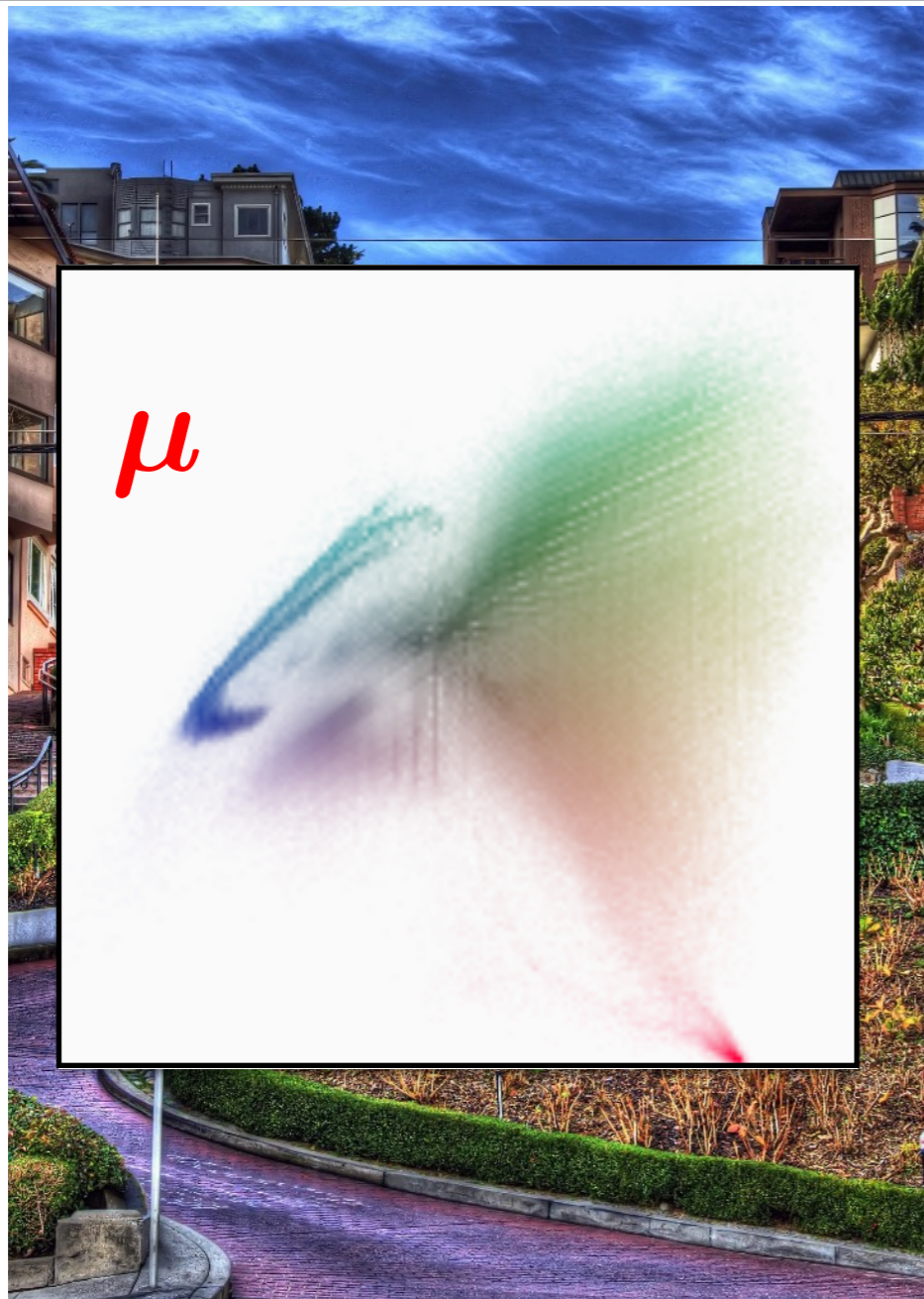
The Earth Mover's Distance



The Earth Mover's Distance



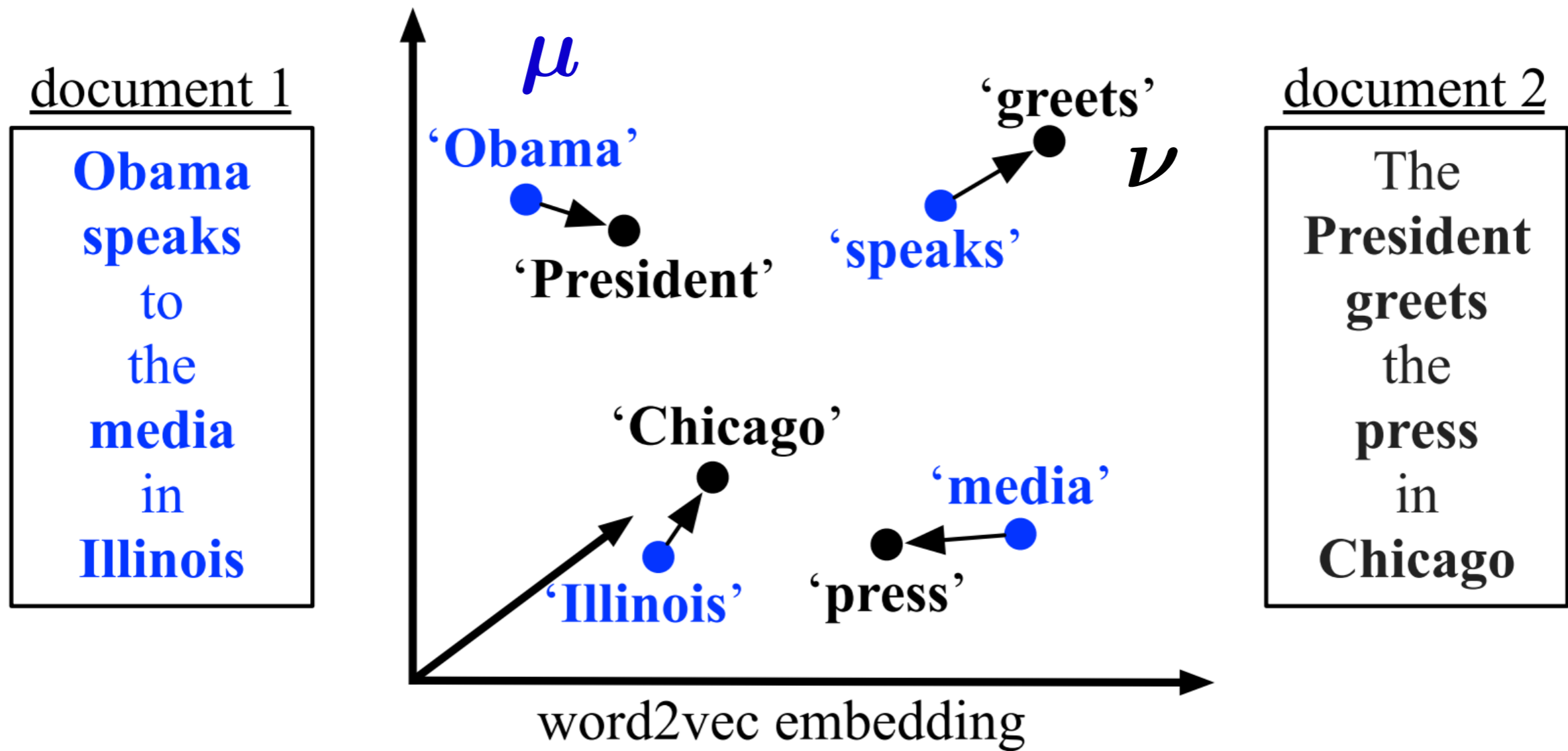
The Earth Mover's Distance



[Rubner'98]

$$\text{dist}(I_1, I_2) = W_1(\mu, \nu)$$

The Word Mover's Distance



[Kusner'15] $\text{dist}(D_1, D_2) = W_2(\mu, \nu)$

Recall

Up to 2010: OT solvers $W_p(\mu, \nu) = ?$

Goal now: use OT as a **loss or fidelity** term

$\operatorname{argmin}_{\mu \in \mathcal{P}(\Omega)} F(W_p(\mu, \nu_1), W_p(\mu, \nu_2), \dots, \mu) = ?$

$\nabla_{\mu} W_p(\mu, \nu_1) = ?$

Wassersteinization

[wos-ur-stahyn-ahy-sey-shuh-n]

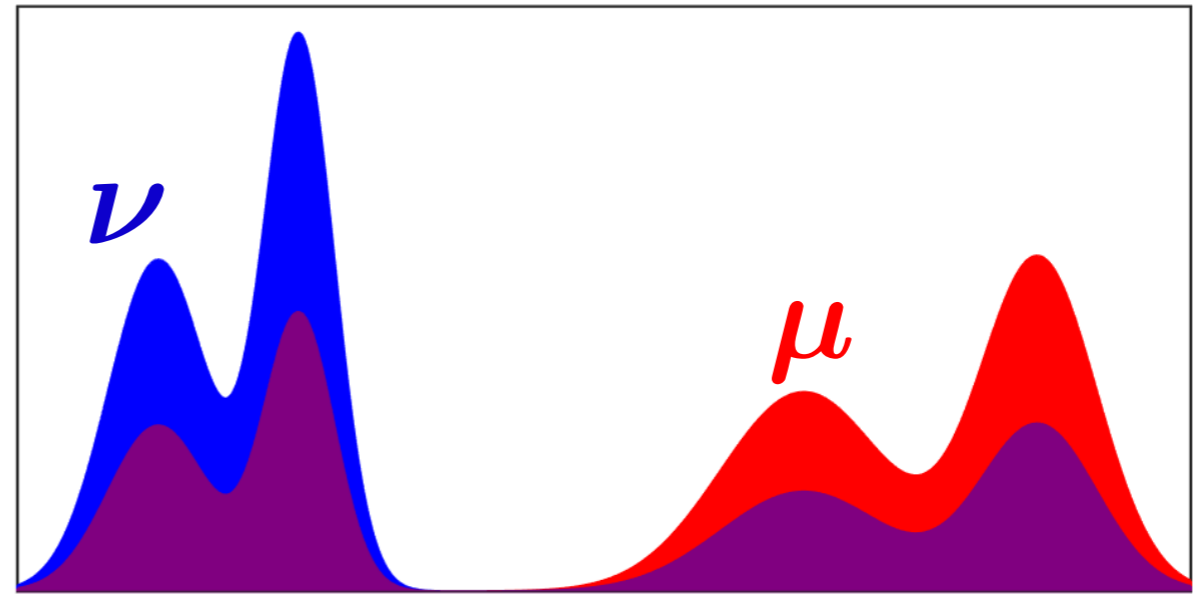
noun.

Introduction of optimal transport into an optimization or learning problem.

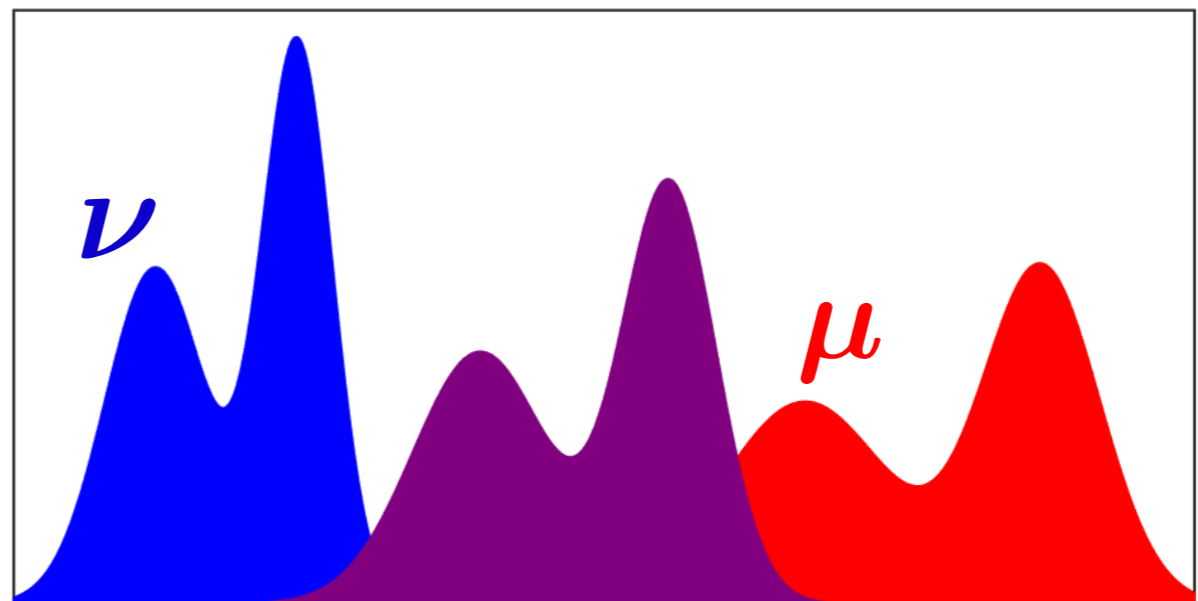
cf. least-squarification, L_1 ification, deep-netification, kernelization

Averaging Measures

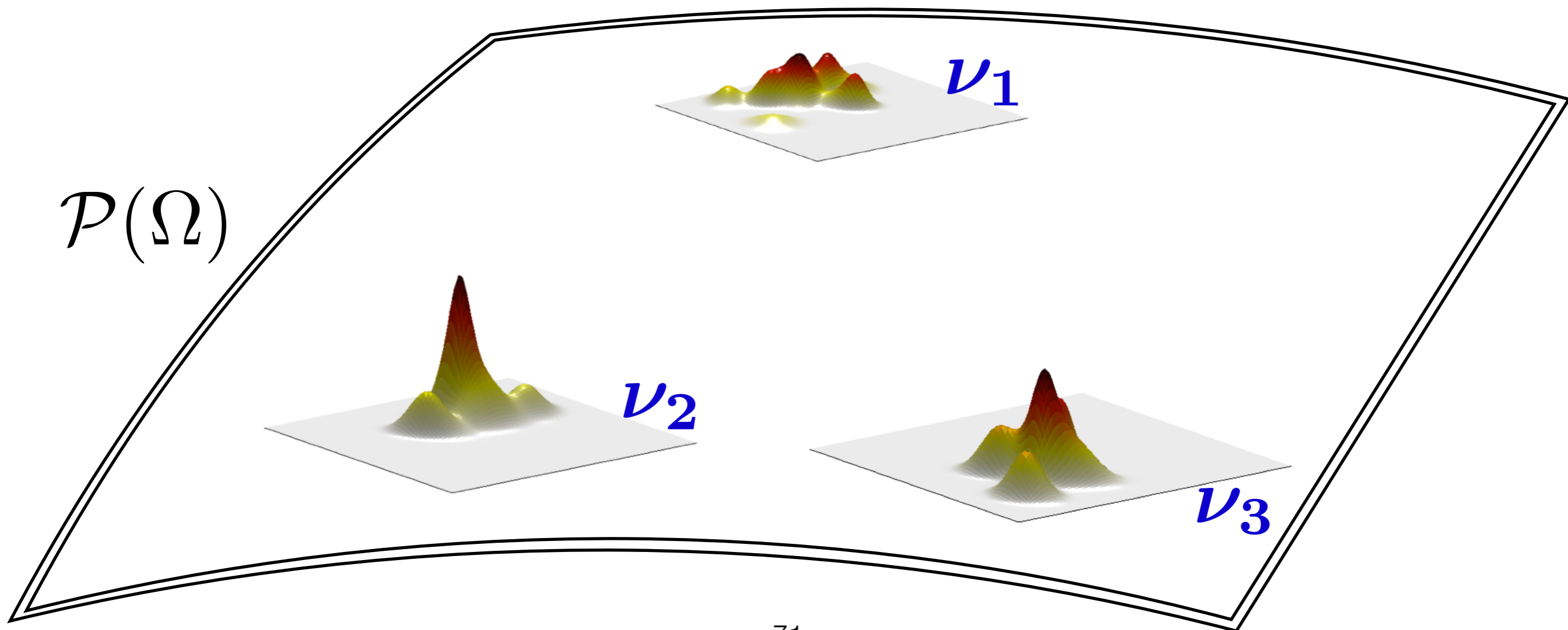
L_2 average



W average

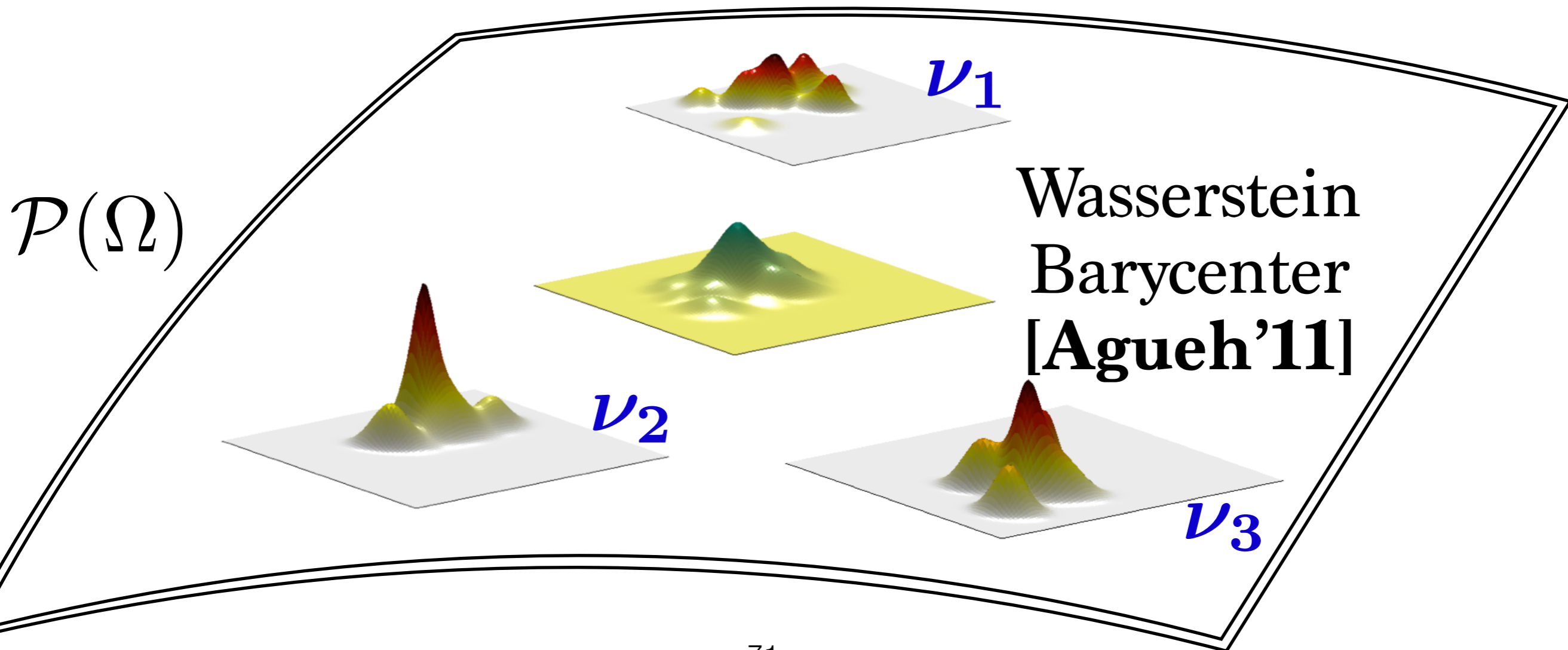


Barycenter for Measures?

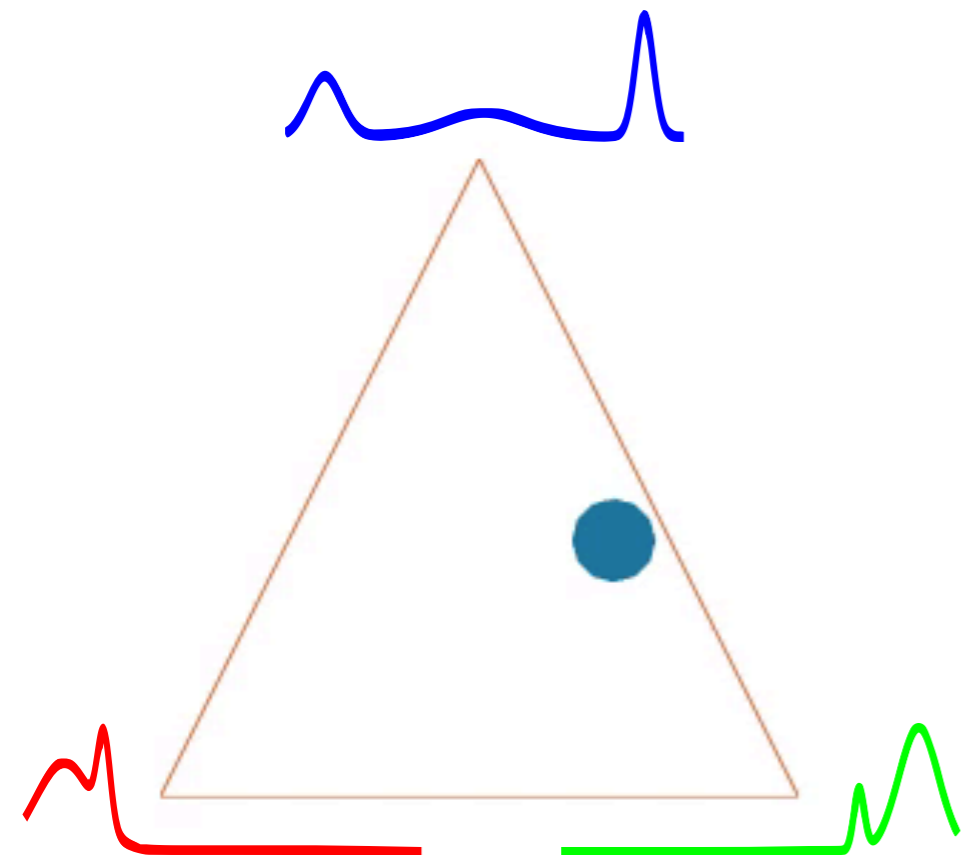


Barycenter for Measures?

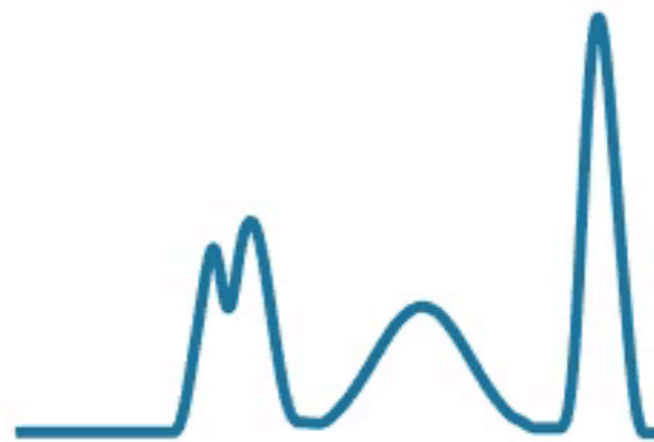
$$\min_{\mu \in \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_p^p(\mu, \nu_i)$$



Barycenter for Measures?



$\lambda \in \Sigma_3$

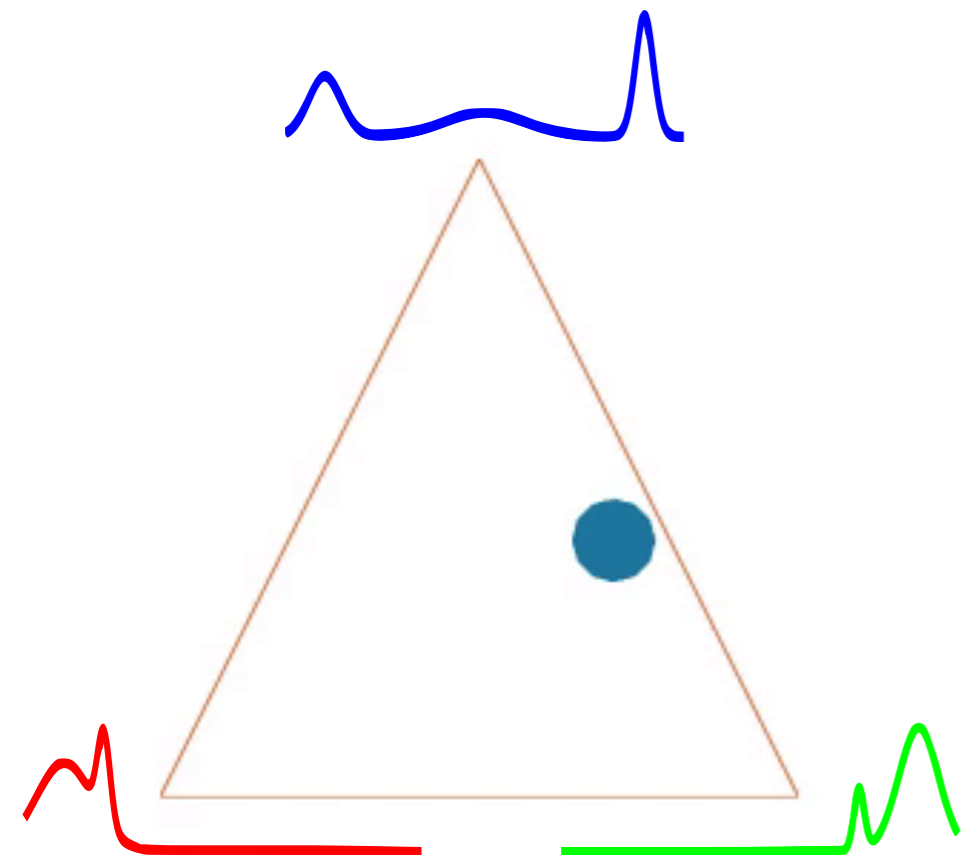


Wasserstein mean



L_2 mean

Barycenter for Measures?



$$\lambda \in \Sigma_3$$

Wasserstein mean

L_2 mean

Averaging Histograms is a LP

When Ω is a **finite metric space** defined by M .

$$\min_{\mathbf{a} \in \Sigma_n} \sum_i \lambda_i W_M(\mathbf{a}, \mathbf{b}_i)$$

Averaging Histograms is a LP

When Ω is a **finite metric space** defined by M .

$$\min_{P_1, \dots, P_N, \mathbf{a}} \sum_{i=1}^N \lambda_i \langle P_i, M \rangle$$

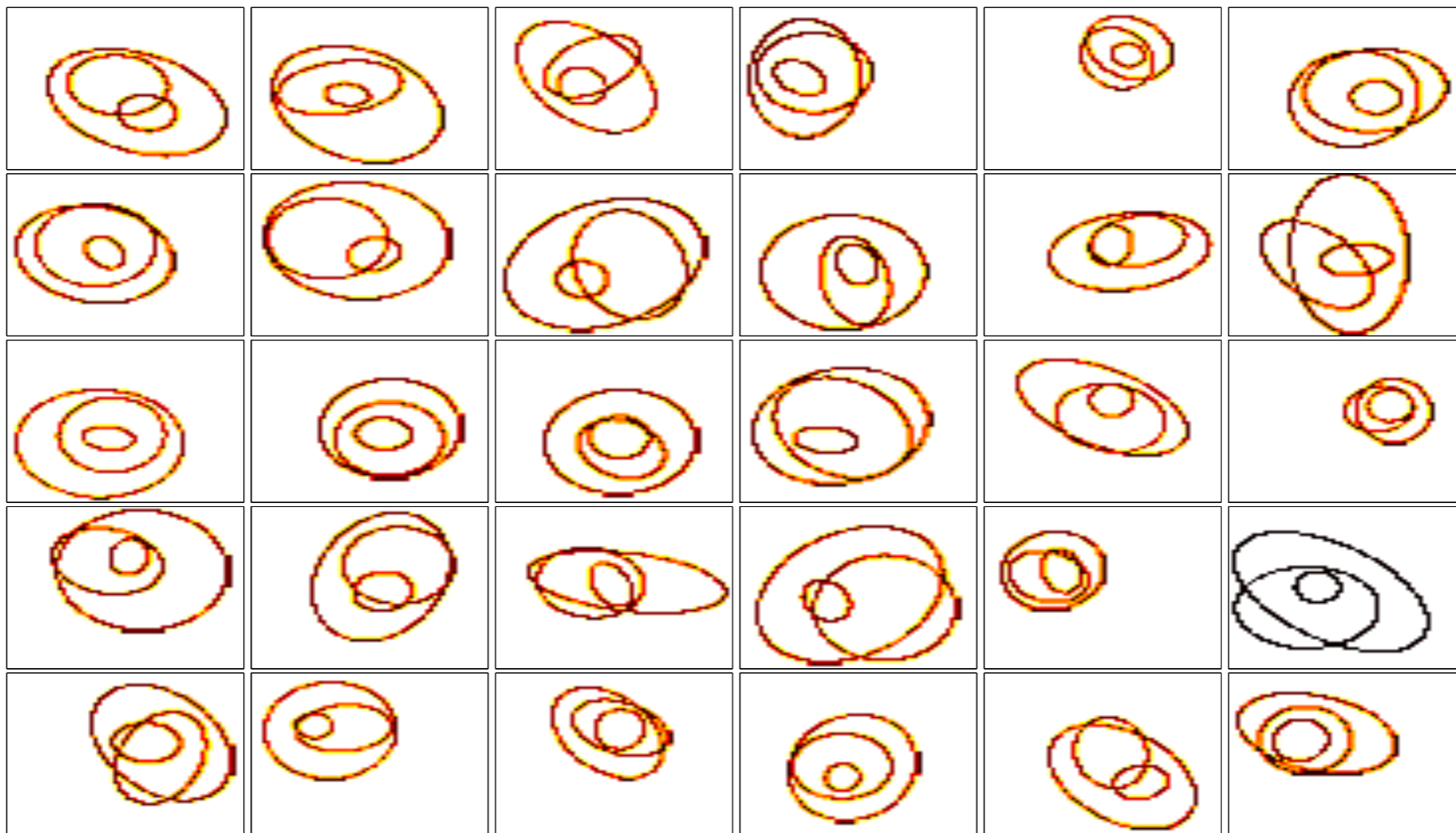
$$\text{s.t. } P_i^T \mathbf{1}_n = \mathbf{b}_i, \forall i \leq N,$$

$$P_1 \mathbf{1}_n = \dots = P_N \mathbf{1}_d = \mathbf{a}.$$

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

Primal Descent on Regularized W

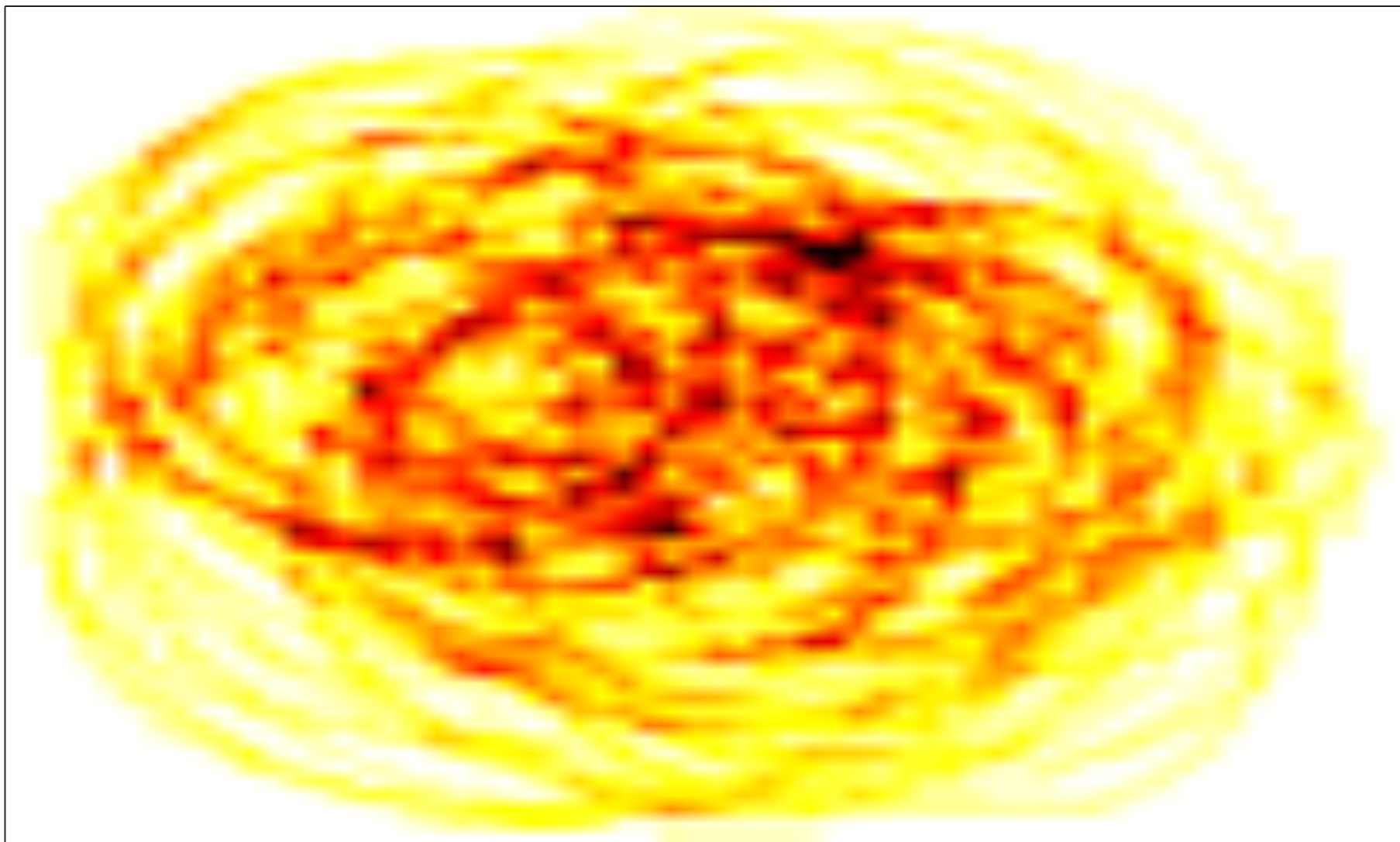
$$\min_{\mathbf{a} \in \Sigma_{h \times h}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$



[Cuturi'14]

Primal Descent on Regularized W

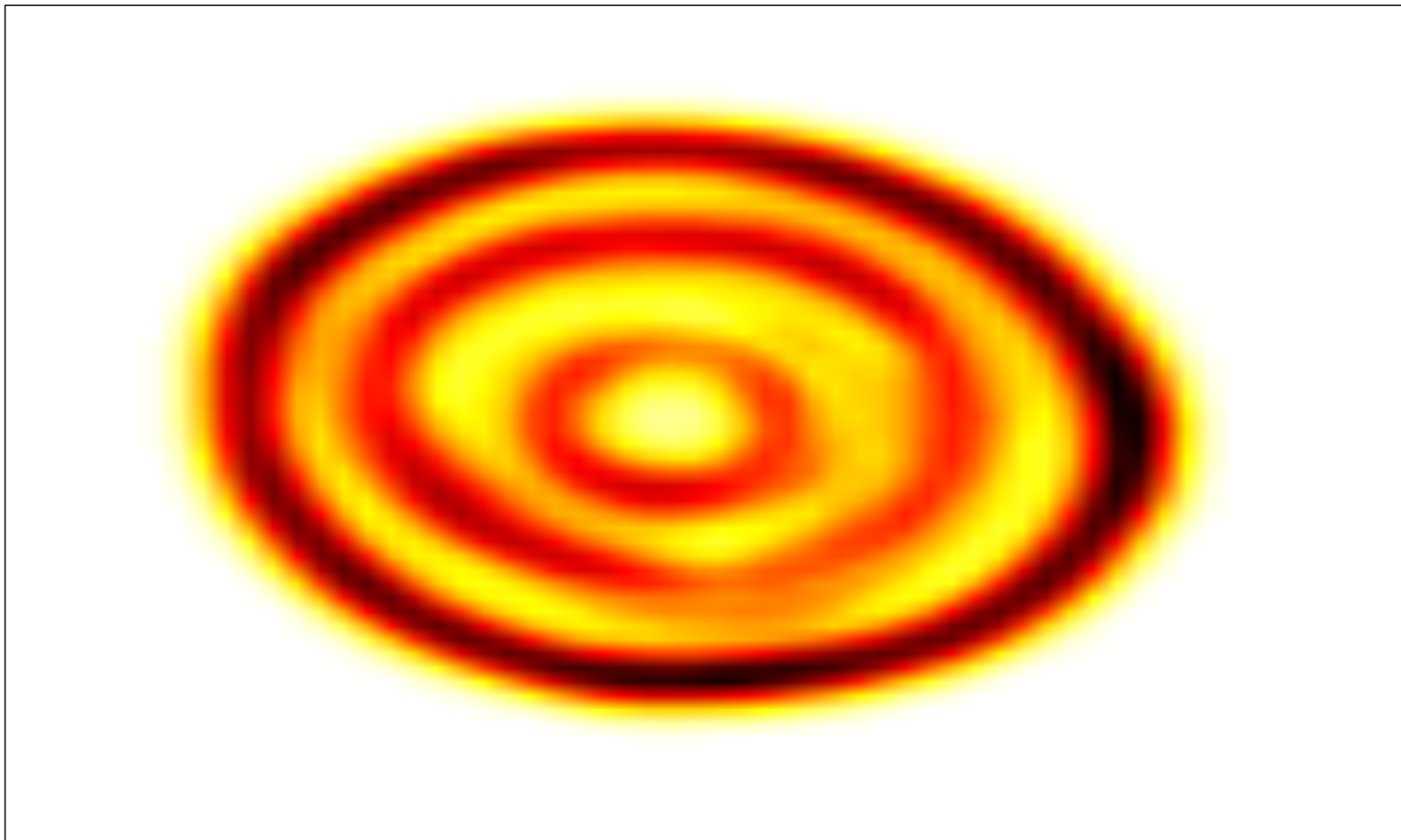
$$\min_{\mathbf{a} \in \Sigma_{h \times h}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$



[Cuturi'14]

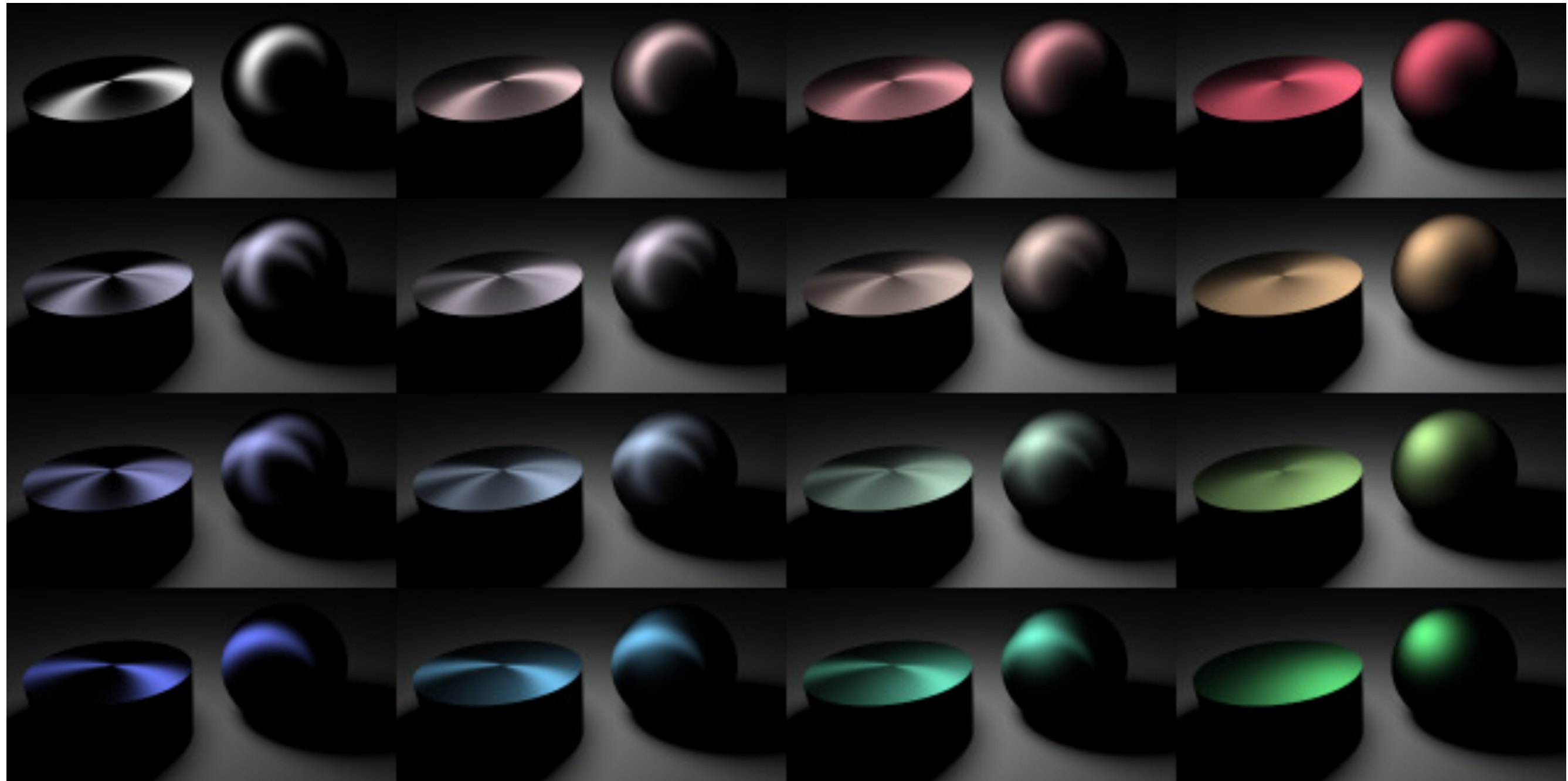
Primal Descent on Regularized W

$$\min_{\mathbf{a} \in \Sigma_{h \times h}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$



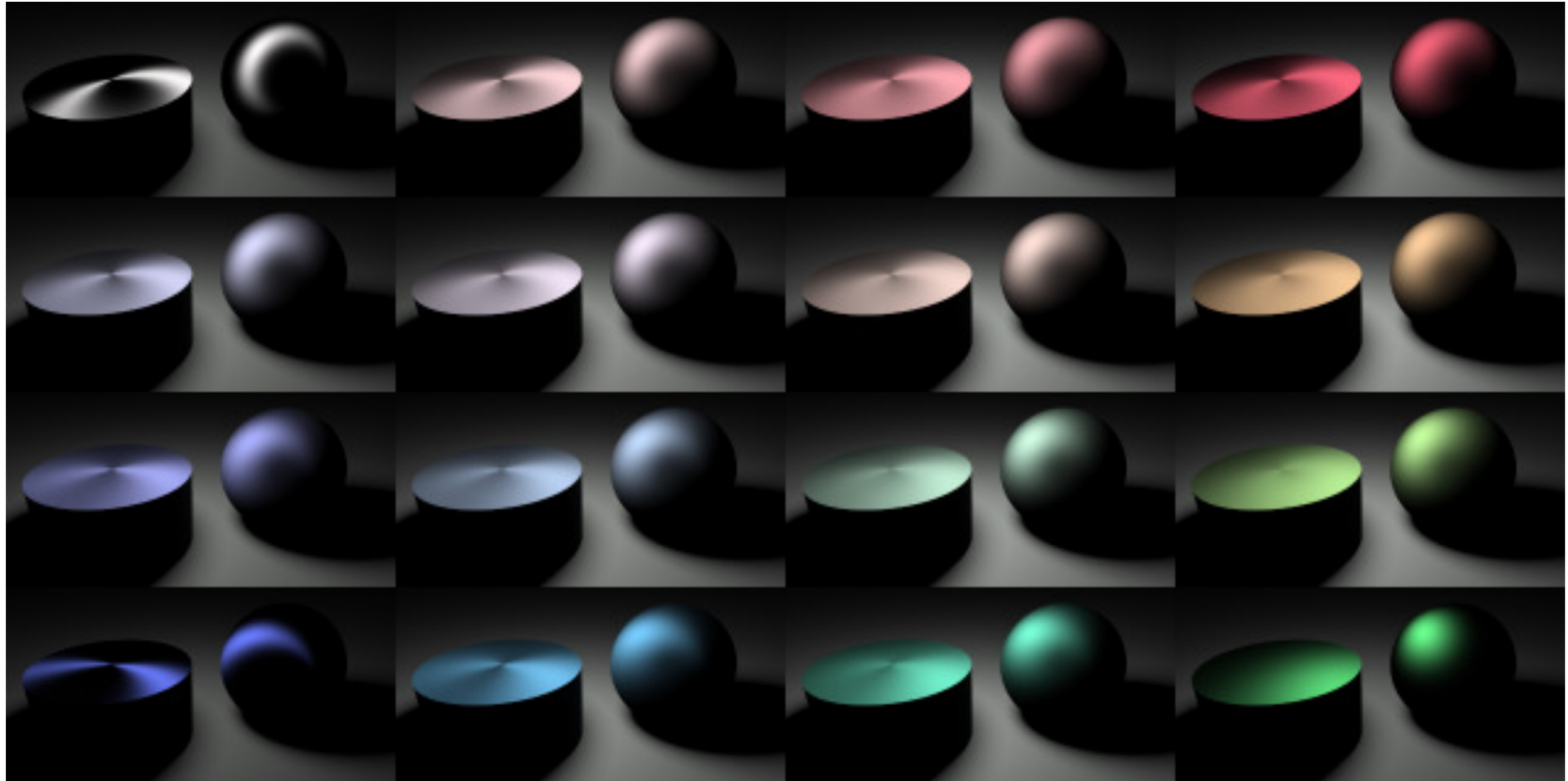
[Cuturi'14]

Applications in Imaging



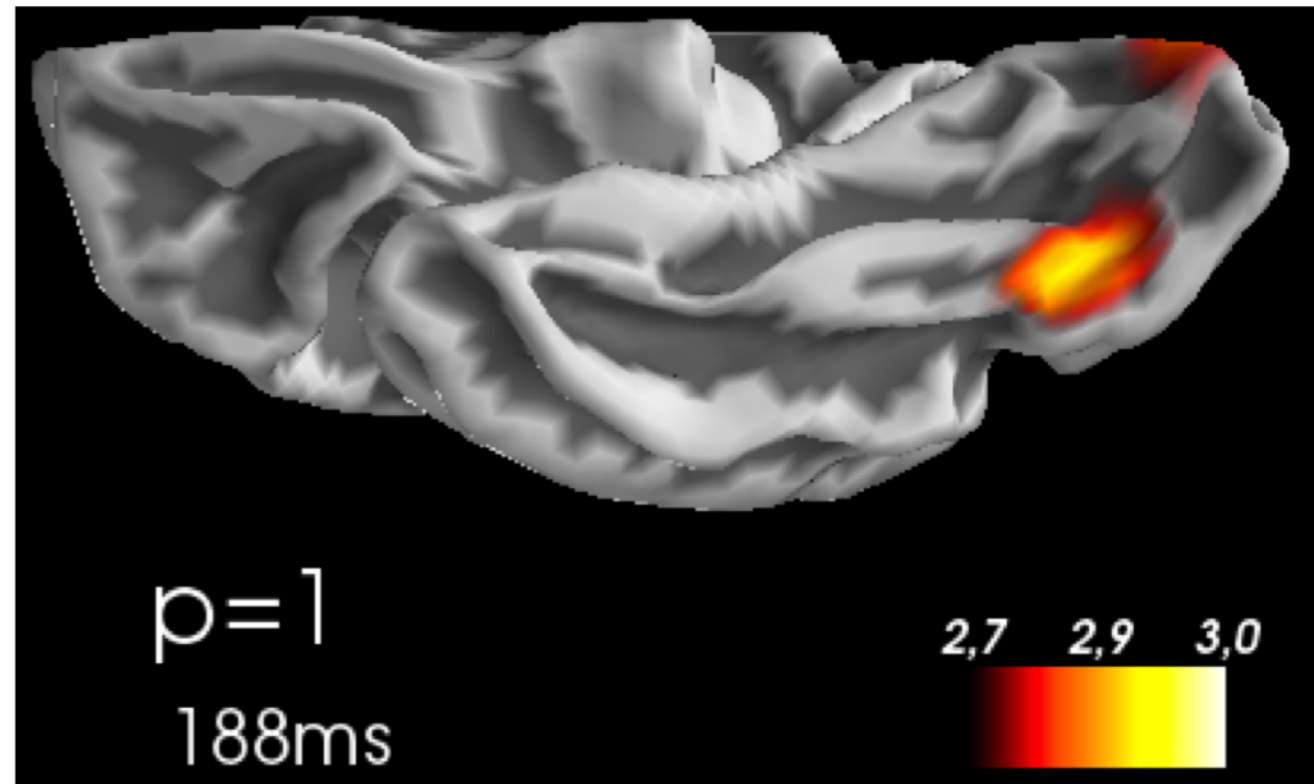
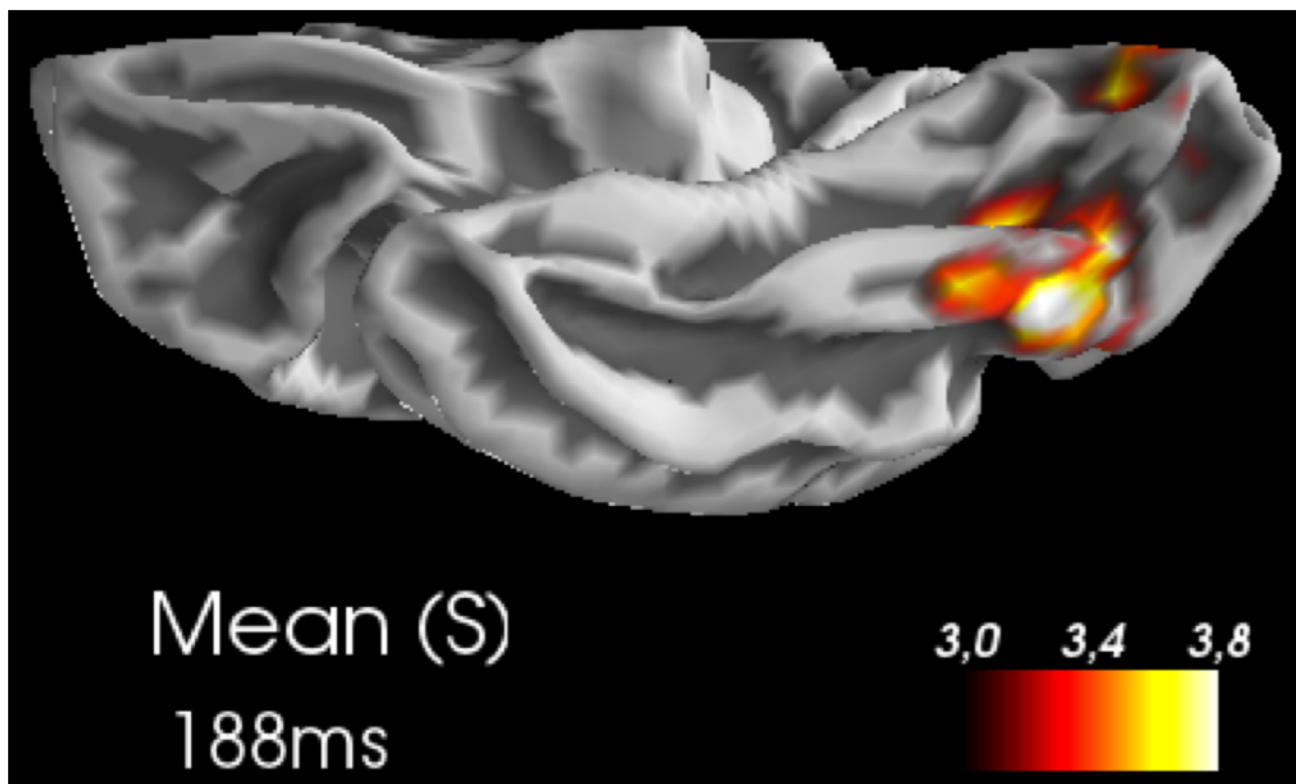
[Solomon'15]

Applications in Imaging



[Solomon'15]

Applications: Brain Imaging



Extension to **non-normalized** data!
Applied to MEG and fMRI.

[Gramfort'16]

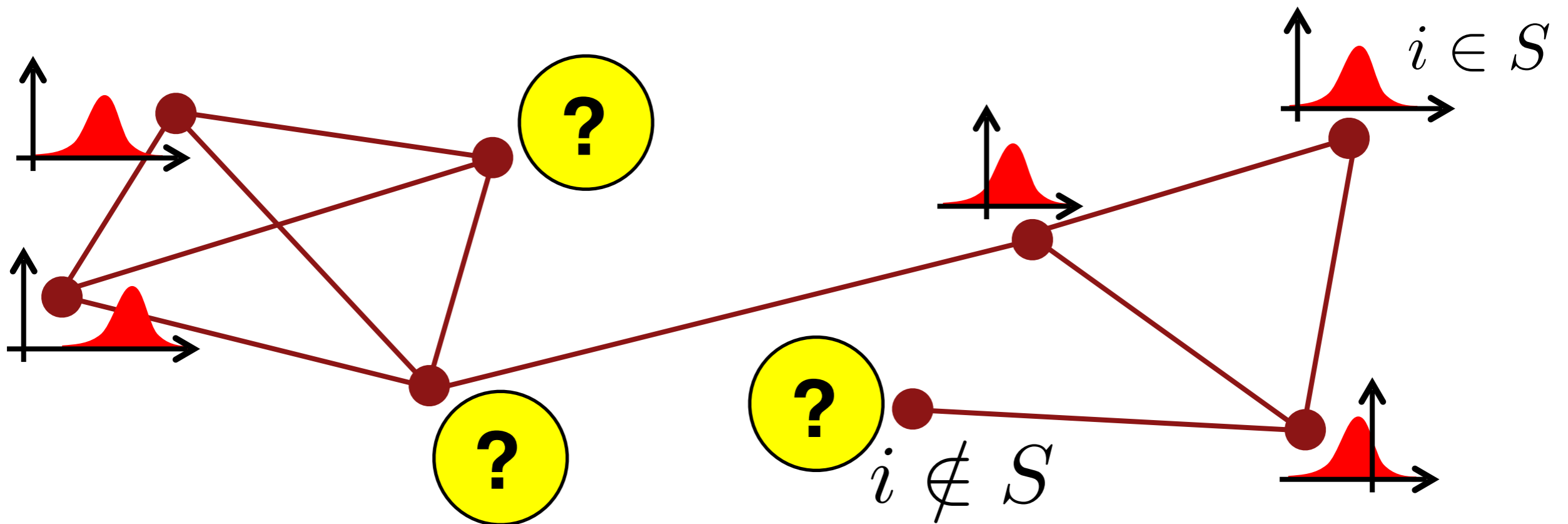
Wasserstein Posterior (WASP)

Merge Bayesian subset posteriors.

1. **Split** data into J subsets S_1, \dots, S_J
2. **Distribute** to J machines.
3. In parallel, **sample** from $\{p(\theta) | S_i\}_{i=1}^J$ using MCMC.
4. **Aggregate** using Wasserstein barycenters

[Srivastava'15]

Wasserstein Propagation



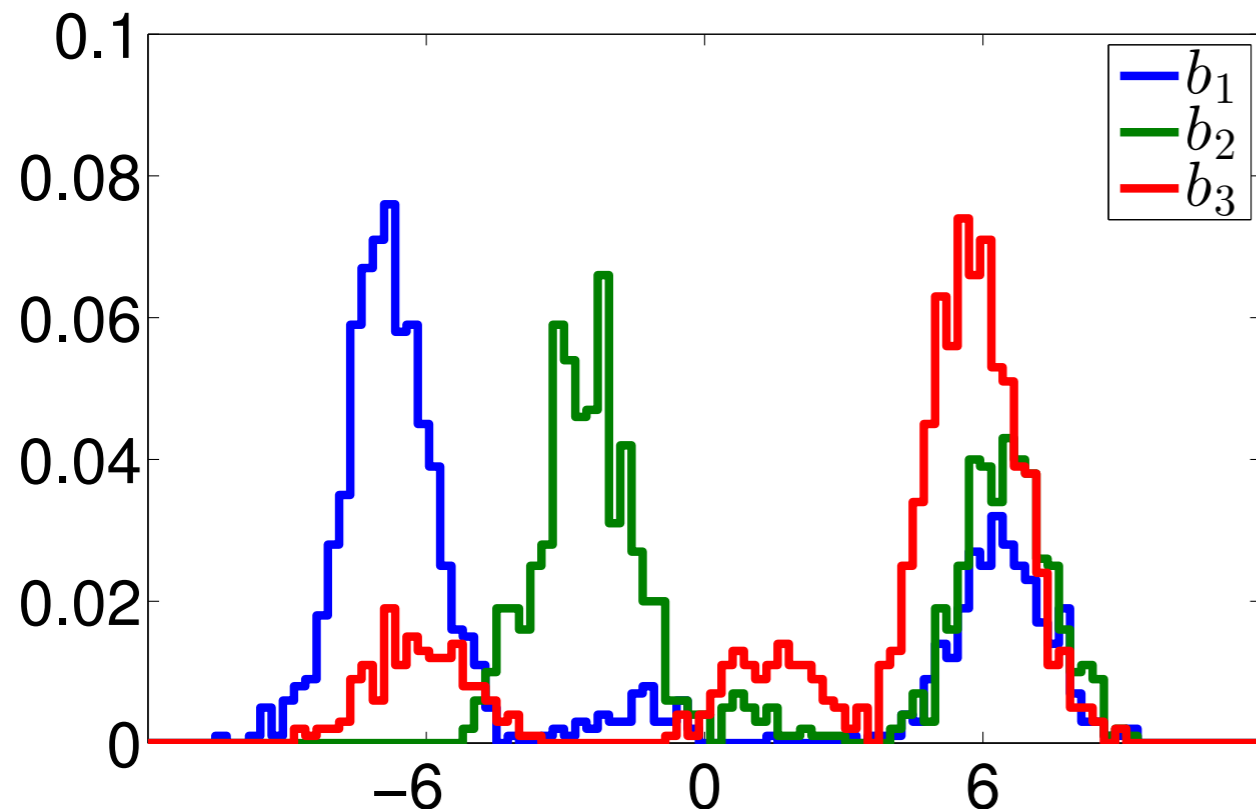
$$\min_{\substack{\mu_i \in \mathcal{P}(\Omega) \\ \mu_i \text{ fixed for } i \in S}} \sum_{(e_1, e_2) \in E} W_2^2(\mu_{e_1}, \mu_{e_2})$$

[Solomon'14]

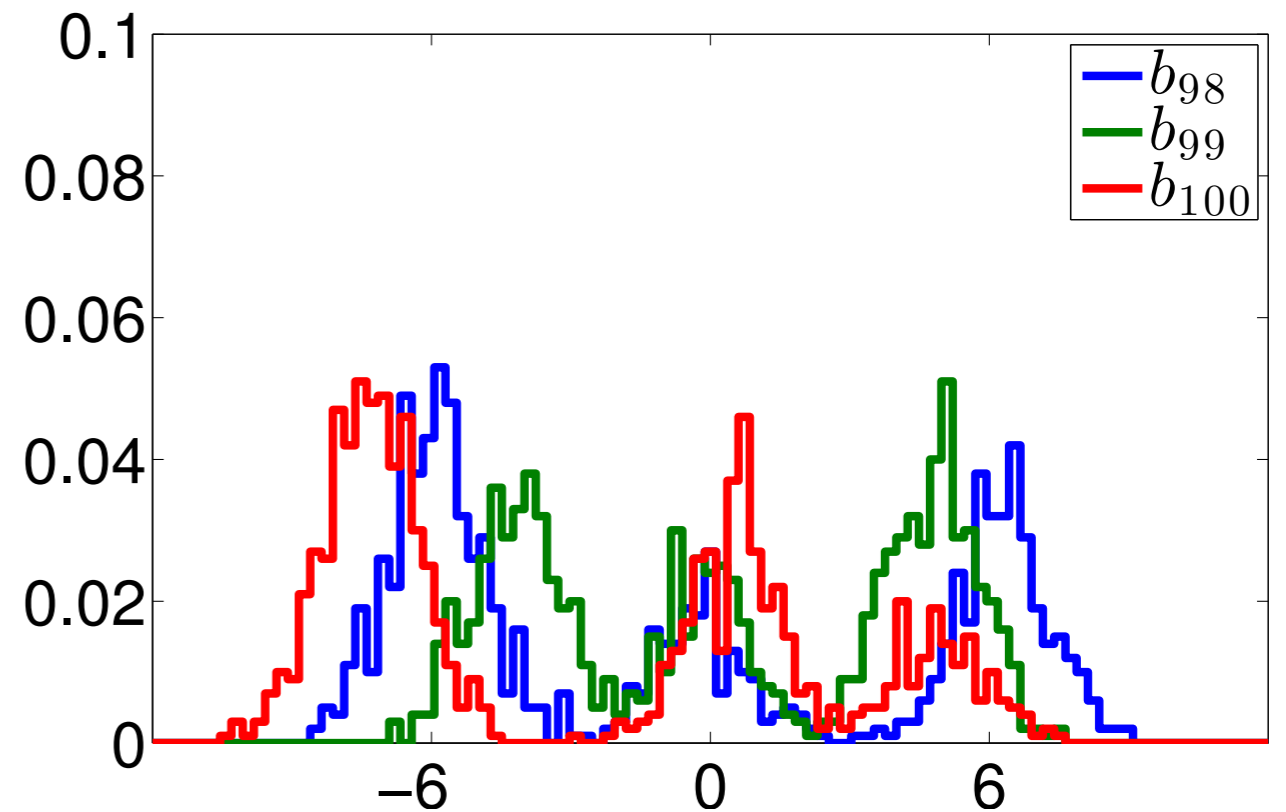
Dictionary Learning

$$\mathbf{A} \in (\Sigma_n)^K, \mathbf{\Lambda} \in (\Sigma_K)^N \min \sum_{i=1}^N W \left(b_i, \sum_{k=1}^K \Lambda_k^i a_k \right)$$

Data samples



Data samples

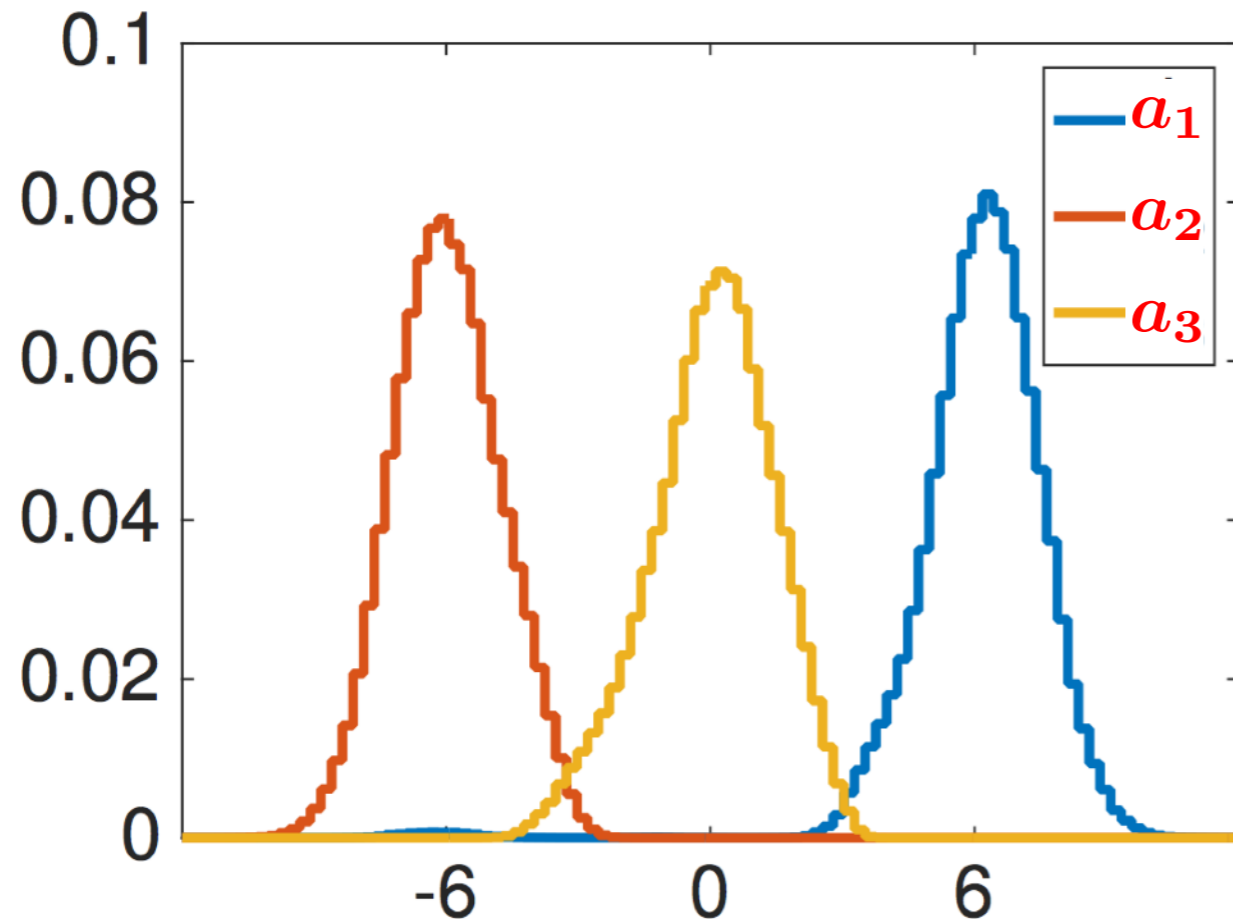


[Sandler'11] [Zen'14] [Rolet'16]

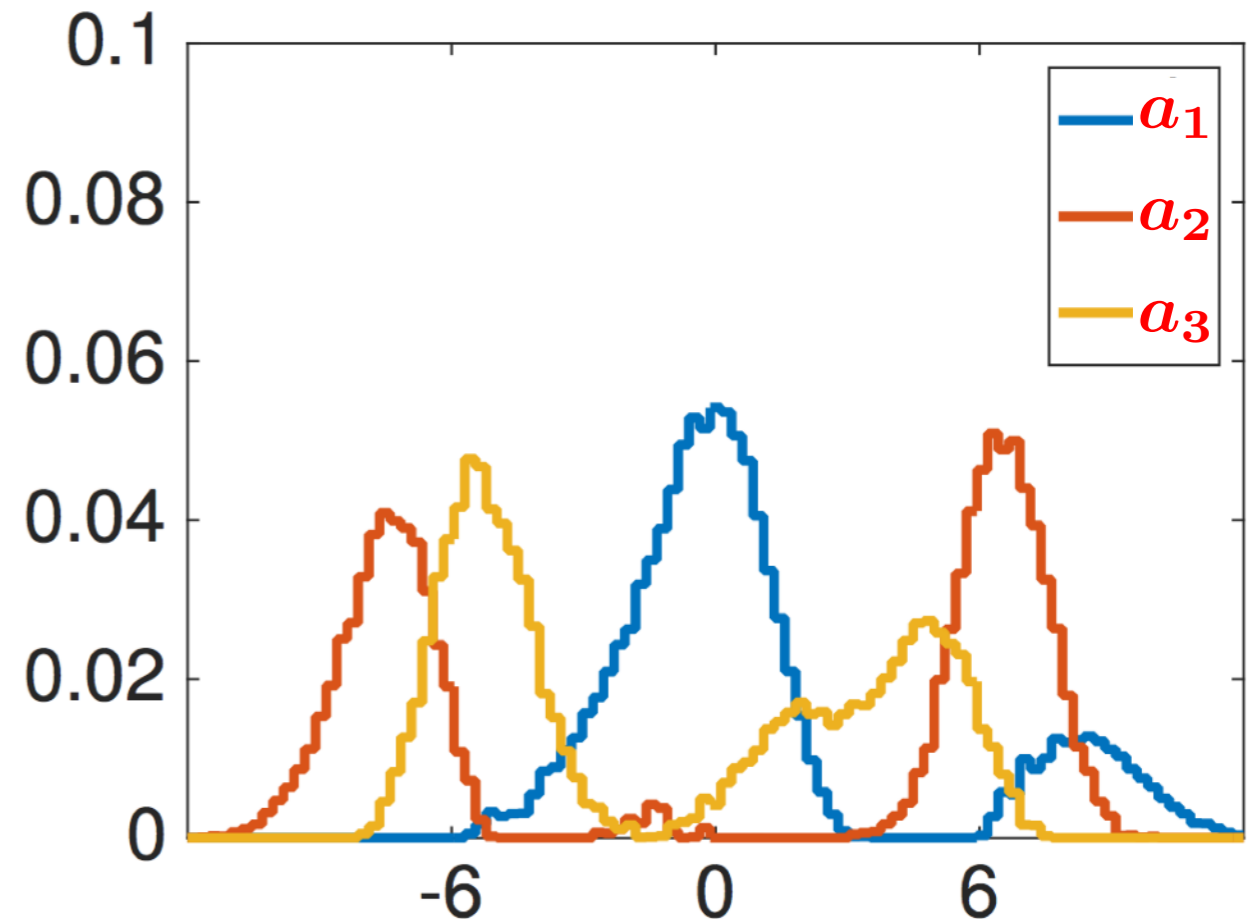
Dictionary Learning

$$\mathbf{A} \in (\Sigma_n)^K, \mathbf{\Lambda} \in (\Sigma_K)^N \min \sum_{i=1}^N W \left(\mathbf{b}_i, \sum_{k=1}^K \Lambda_k^i \mathbf{a}_k \right)$$

Wasserstein NMF



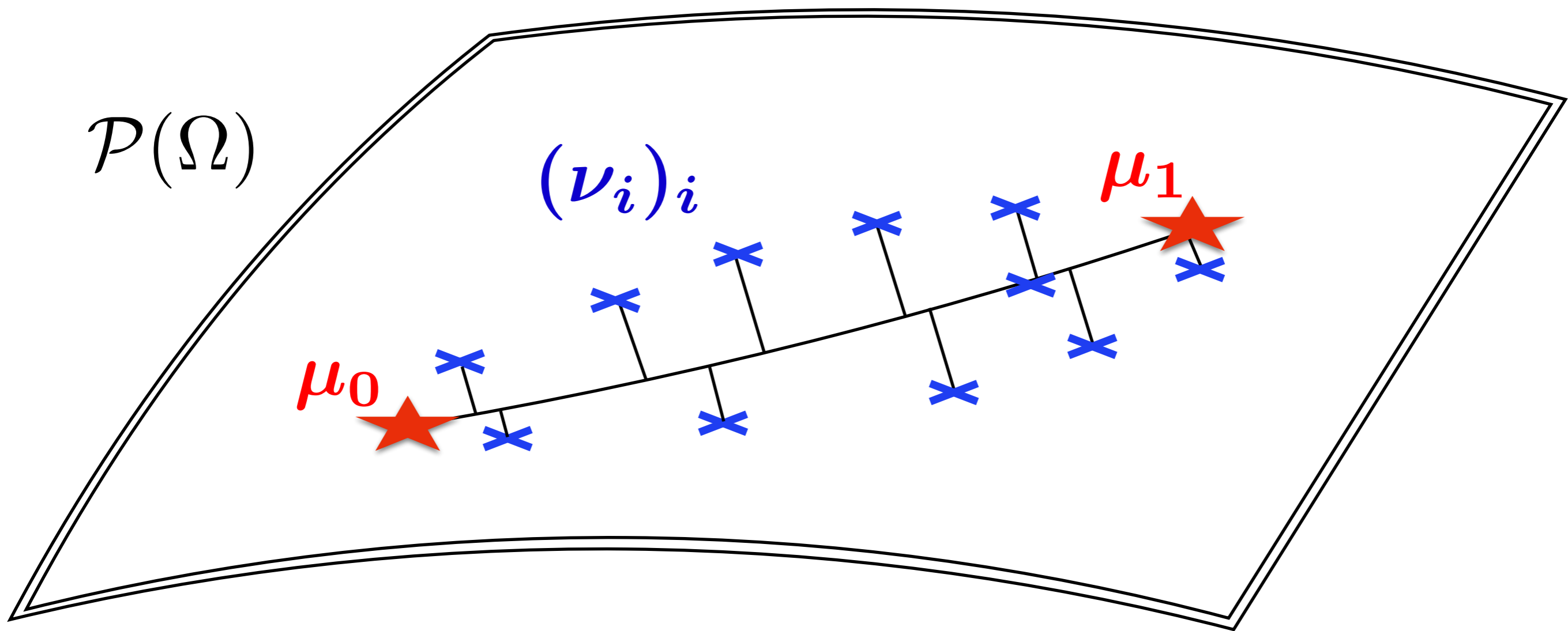
KL NMF



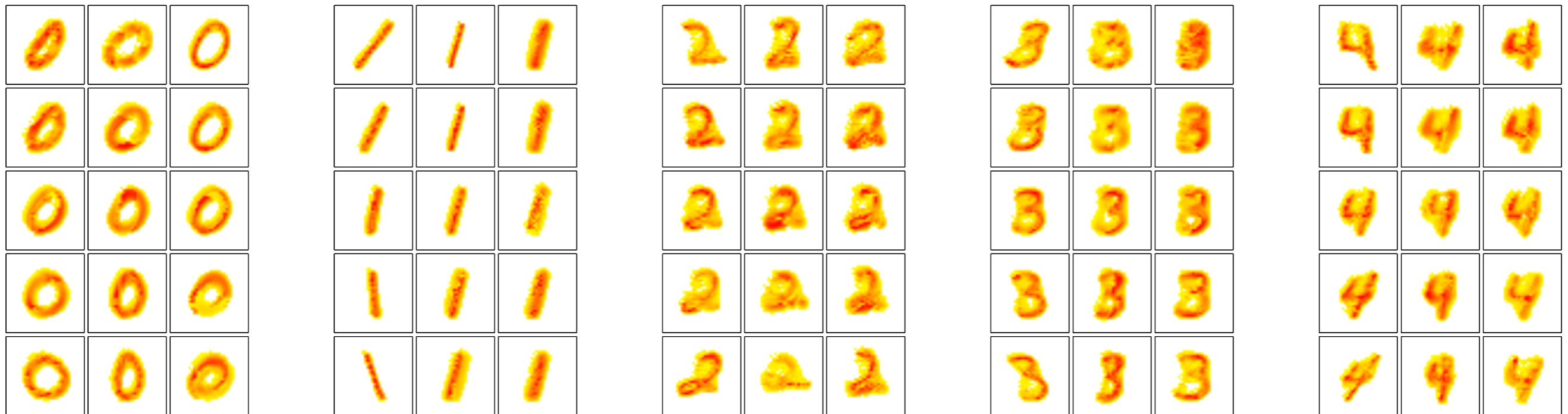
[Sandler'11] [Zen'14] [Rolet'16]

Wasserstein PCA

$$\min_{\mu_0, \mu_1} \sum_{i=1}^N \min_t W_2^2(\rho_{\mu_0 \rightarrow \mu_1}^t, \nu_i)$$



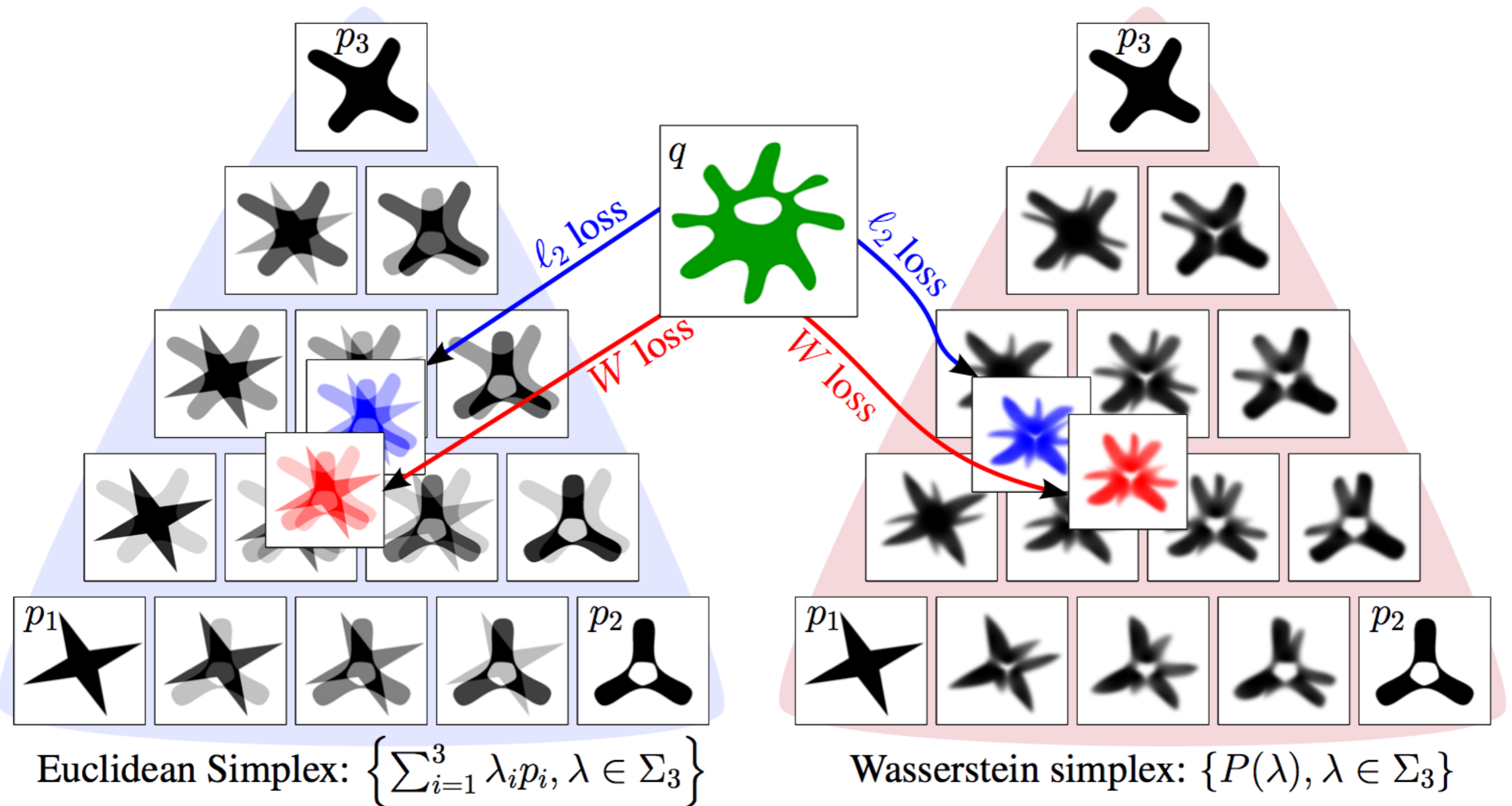
Generalized Principal Geodesics



For each digit, 1,000 MNIST images

[Seguy'15]

Wasserstein Inverse Problems



Application: Volume Reconstruction



Shape database
 (p_1, \dots, p_5)



Input shape q



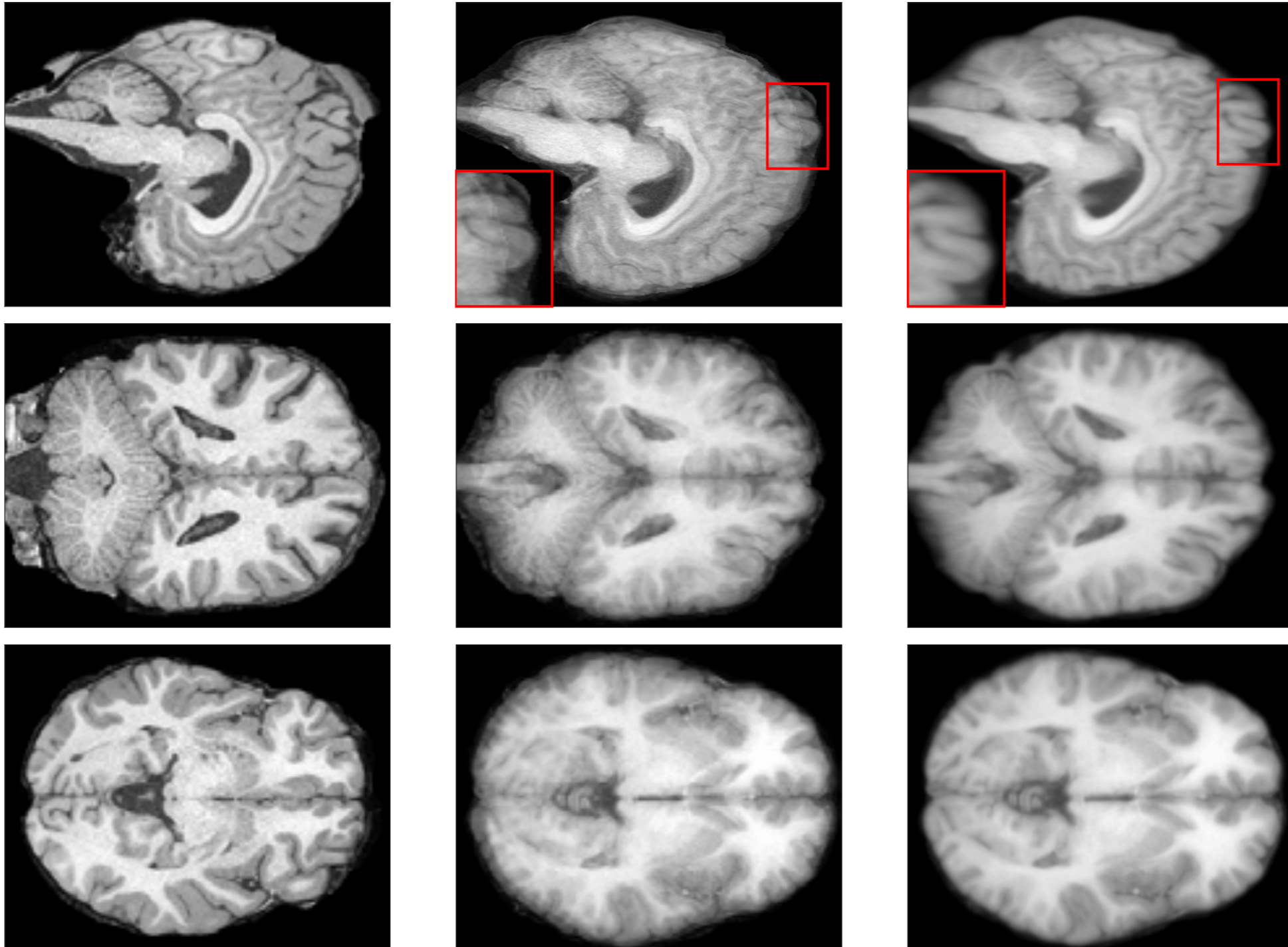
Projection
 $P(\lambda)$



Iso-surface

[Bonneel'16]

Application: Brain Mapping

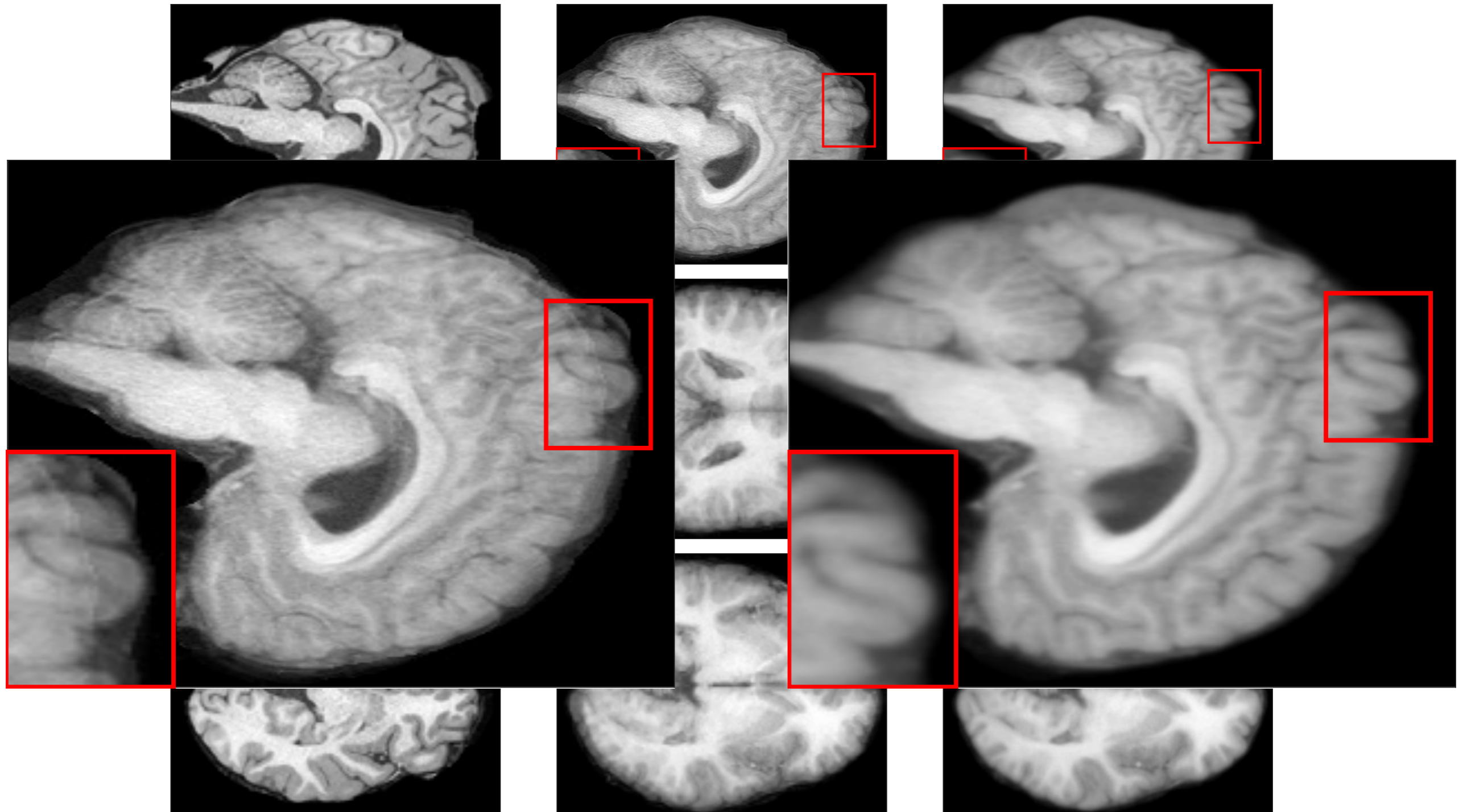


Original

Euclidean
projection

Wasserstein
projection

Application: Brain Mapping



Original

Euclidean
projection

Wasserstein
projection

Distributionally Robust Learning

$$\nu_{\text{data}} = \frac{1}{n} \sum_{i=1}^N \delta_{(x_i, y_i)}$$

Supervised learning

$$\inf_{\theta \in \Theta} \mathbb{E}_{\nu_{\text{data}}} [\mathcal{L}(f_{\theta}(X), Y)]$$

Learning with Wasserstein Ambiguity

$$\inf_{\theta \in \Theta} \sup_{\mu: W_p(\nu_{\text{data}}, \mu) < \varepsilon} \mathbb{E}_{\mu} [\mathcal{L}(f_{\theta}(X), Y)]$$

[Esvahani'17]

Distributionally Robust Learning

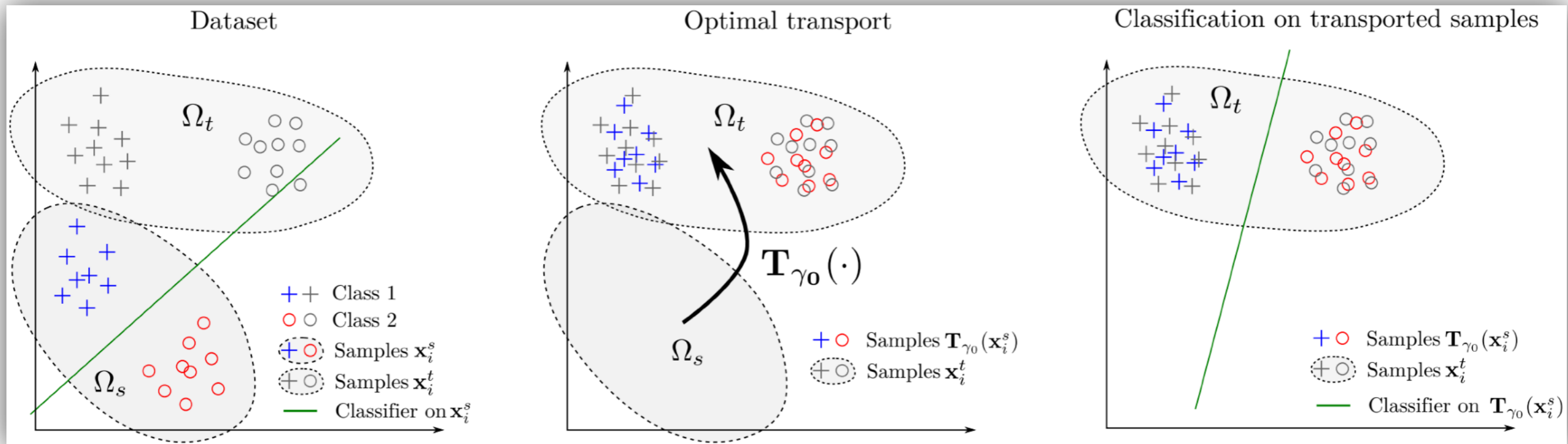
Learning with Wasserstein Ambiguity

$$\inf_{\theta \in \Theta} \sup_{\mu: W_p(\nu_{\text{data}}, \mu) < \varepsilon} \mathbb{E}_{\mu} [\mathcal{L}(f_{\theta}(X), Y)]$$

Advantages:

- Bound on out-of-sample performance
- Converges as size of dataset increases
- Often reduces to a finite convex program (e.g. when f is element-wise max over elementary concave functions)

Domain Adaptation



1. **Estimate** transport map
2. **Transport** labeled samples to new domain
3. **Train** classifier on transported labeled samples

Learning with a Wasserstein Loss

Dataset $\{(x_i, y_i)\}$, $x_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}_+^n$



x_i

husky
snow
sled
slope
men

y_i

Goal is to find $f_{\theta} : \text{Images} \mapsto \text{Labels}$

Learning with a Wasserstein Loss

$$\min_{\theta \in \Theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(x_i), y_i)$$



x_i

husky
snow
sled
slope
men

y_i

Which loss \mathcal{L} could we use?

Learning with a Wasserstein Loss

$$\min_{\theta \in \Theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(x_i), y_i)$$

dog
driver
winter
ice

$f_{\theta}(x_i)$

husky
snow
sled
slope
men

y_i

Which loss \mathcal{L} could we use?

Learning with a Wasserstein Loss

$$\min_{\theta \in \Theta} \sum_{i=1}^N \mathcal{L}(f_{\theta}(x_i), y_i)$$

$$\begin{aligned} \mathcal{L}(\mathbf{a}, \mathbf{b}) = & \min_{P \in \mathbb{R}^{nm}} \langle P, M \rangle + \varepsilon \text{KL}(P \mathbf{1}, \mathbf{a}) \\ & + \varepsilon \text{KL}(P^T \mathbf{1}, \mathbf{b}) - \gamma E(P) \end{aligned}$$

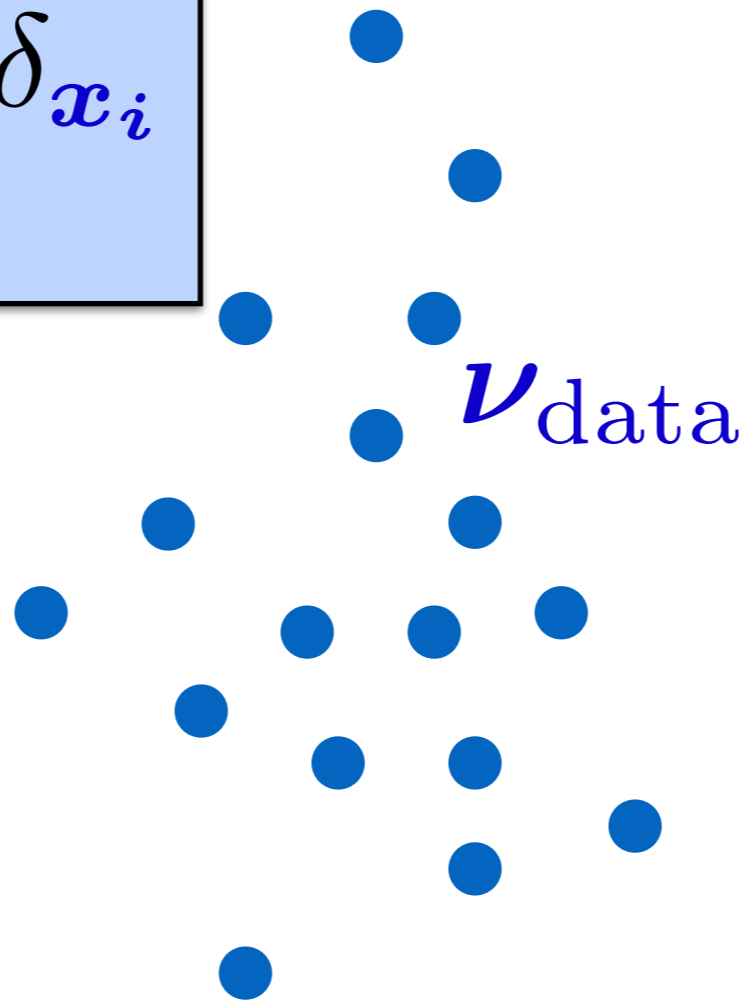
1. Generalizes Word Mover's to label clouds
2. Sinkhorn algorithm can be generalized

[Frogner'15] [Chizat'15][Chizat'16]

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

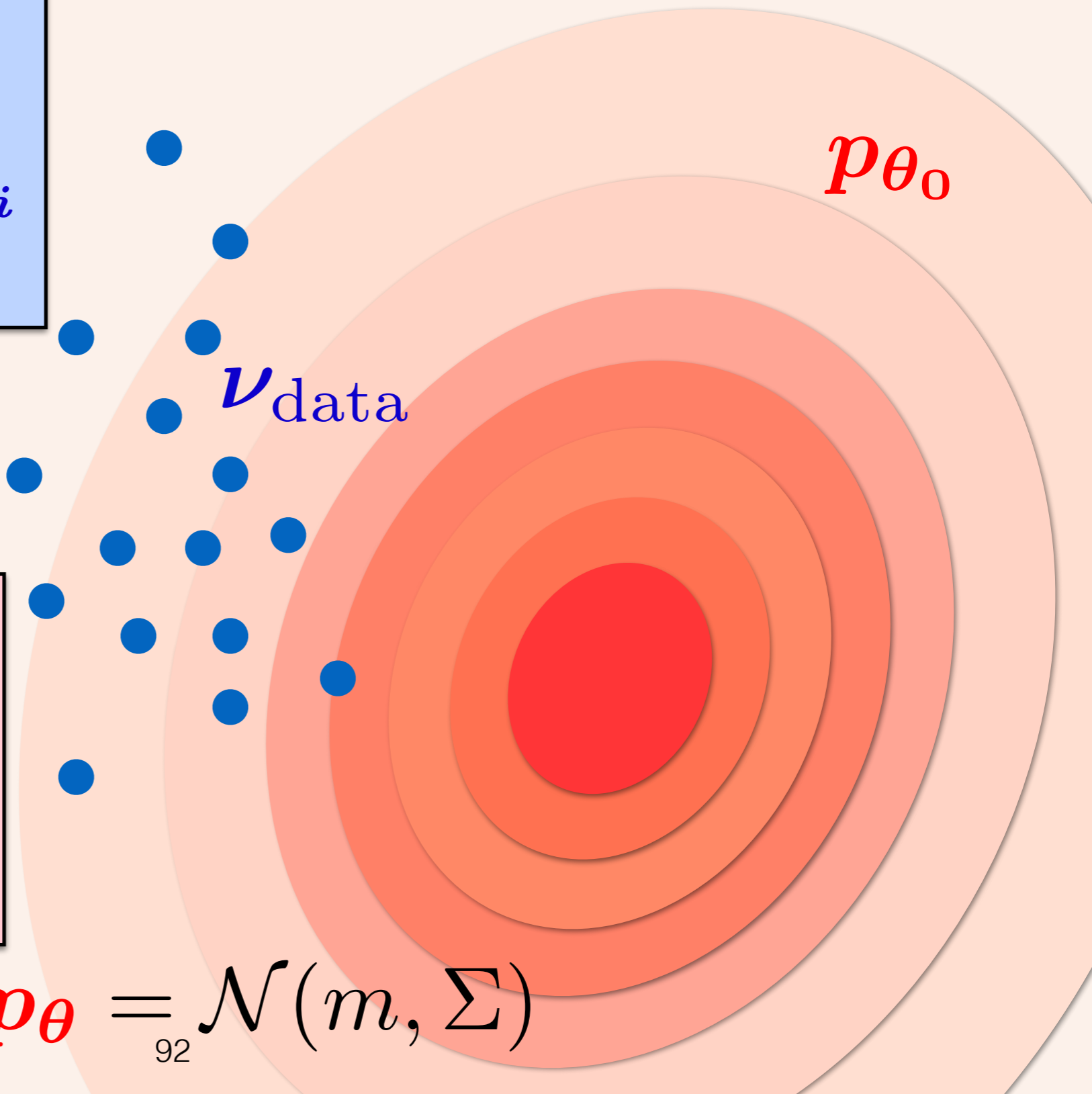
We collect data

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

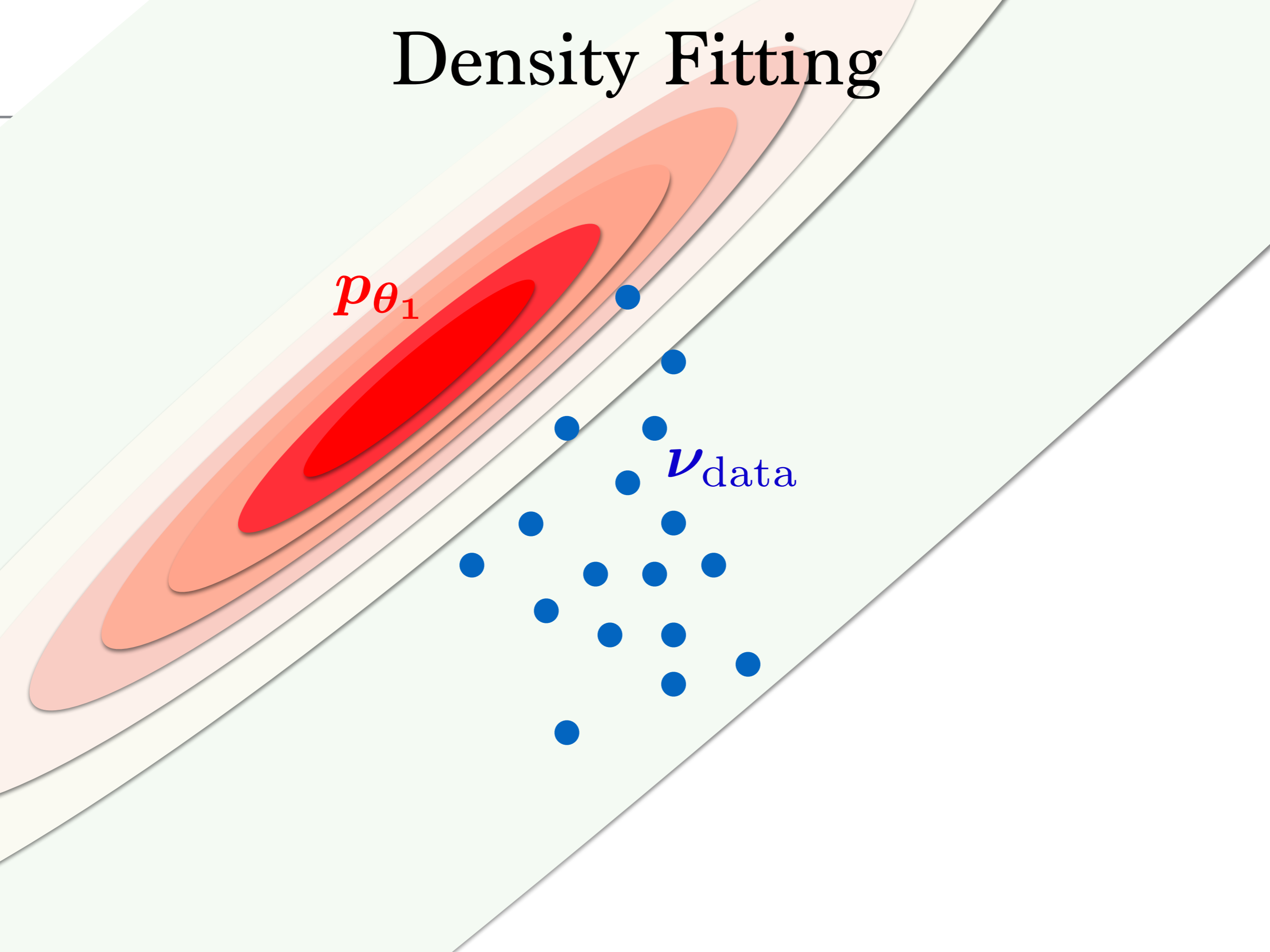
We fit a parametric family of densities

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

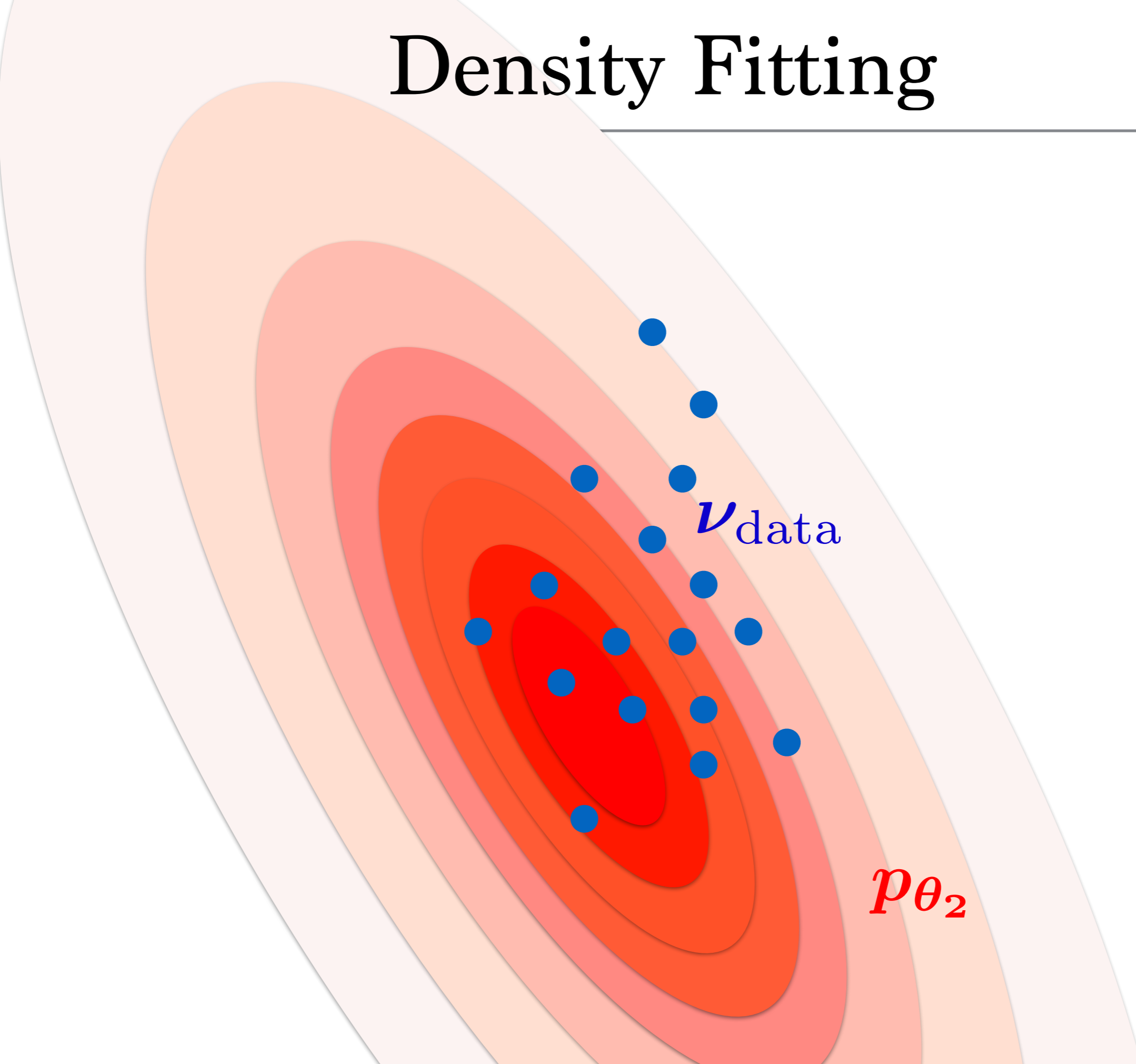
e.g. $\theta = (m, \Sigma); p_{\theta} = \mathcal{N}(m, \Sigma)$



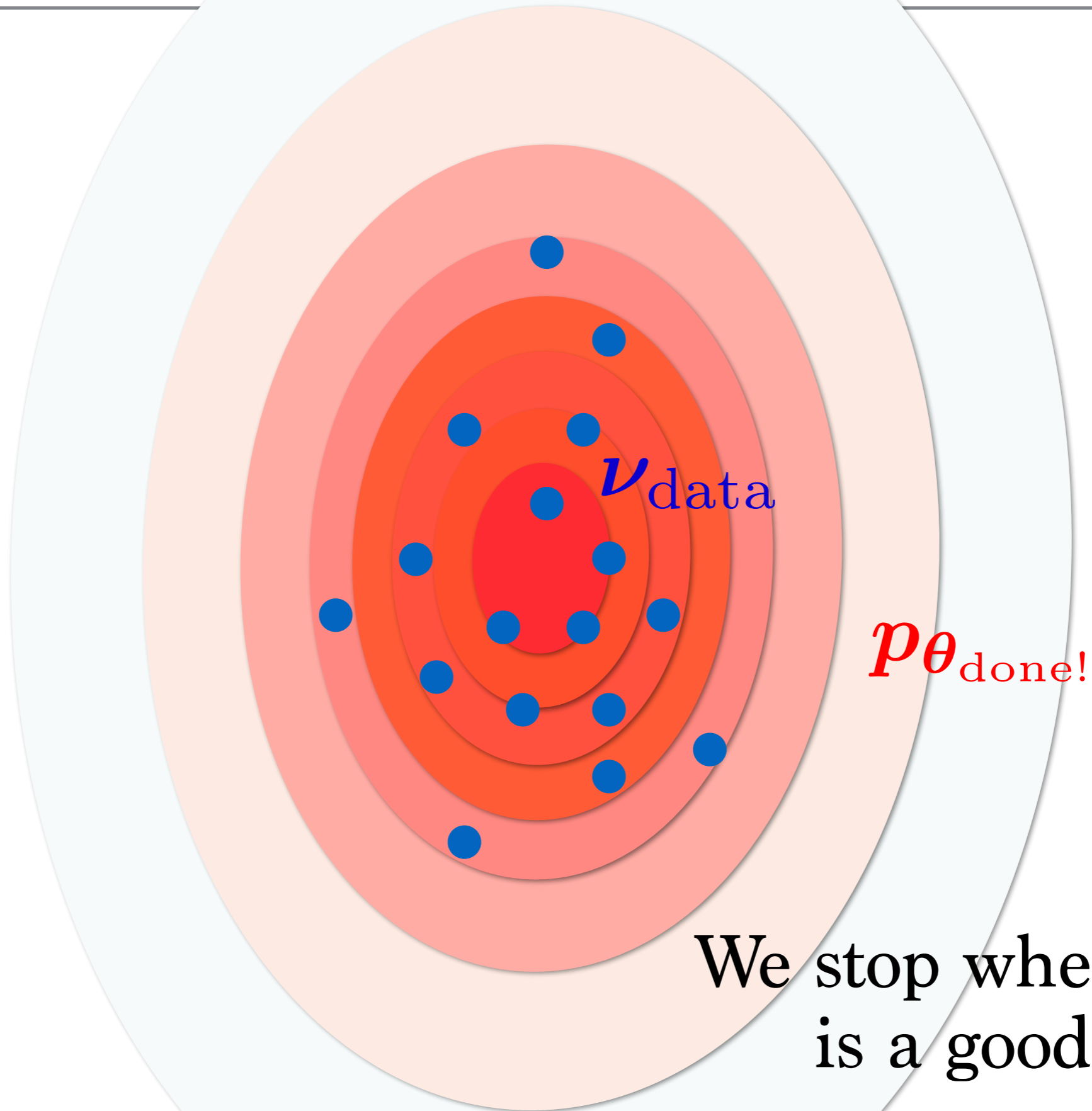
Density Fitting



Density Fitting



Density Fitting

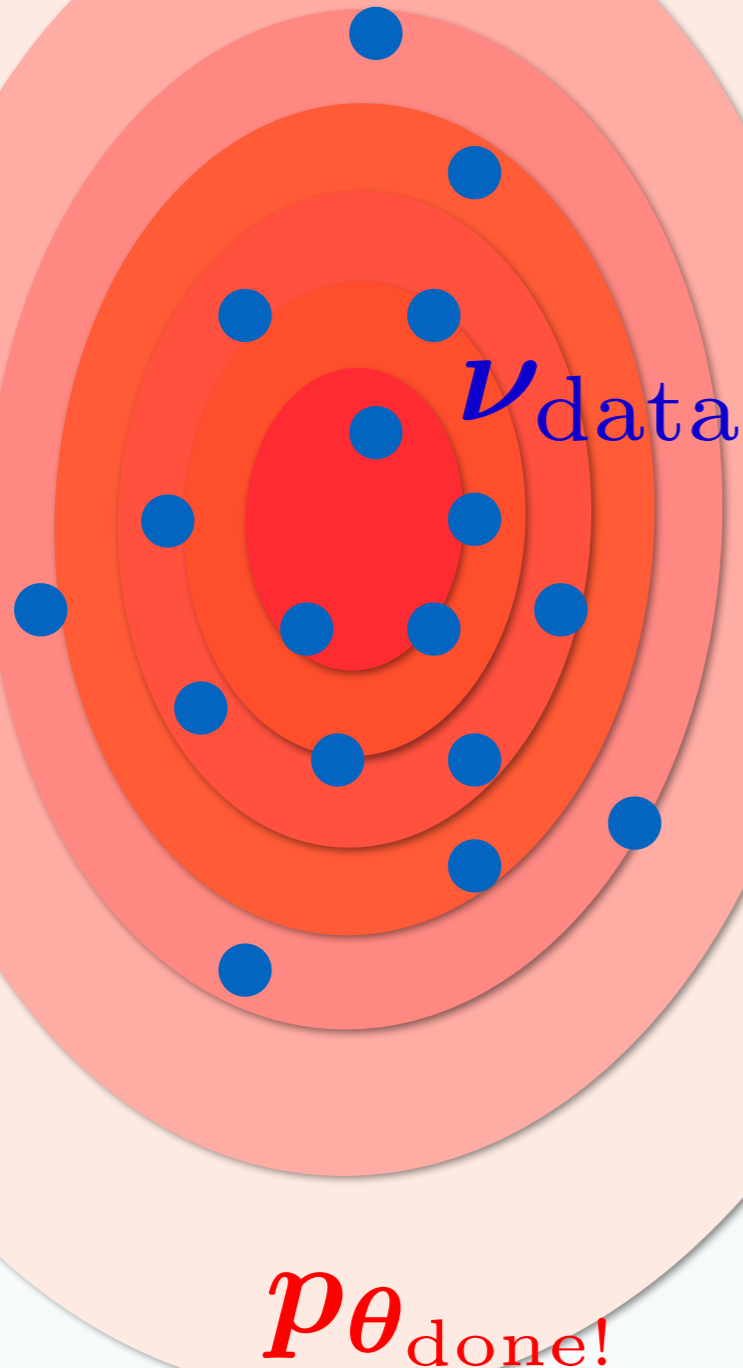


Maximum Likelihood Estimation

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



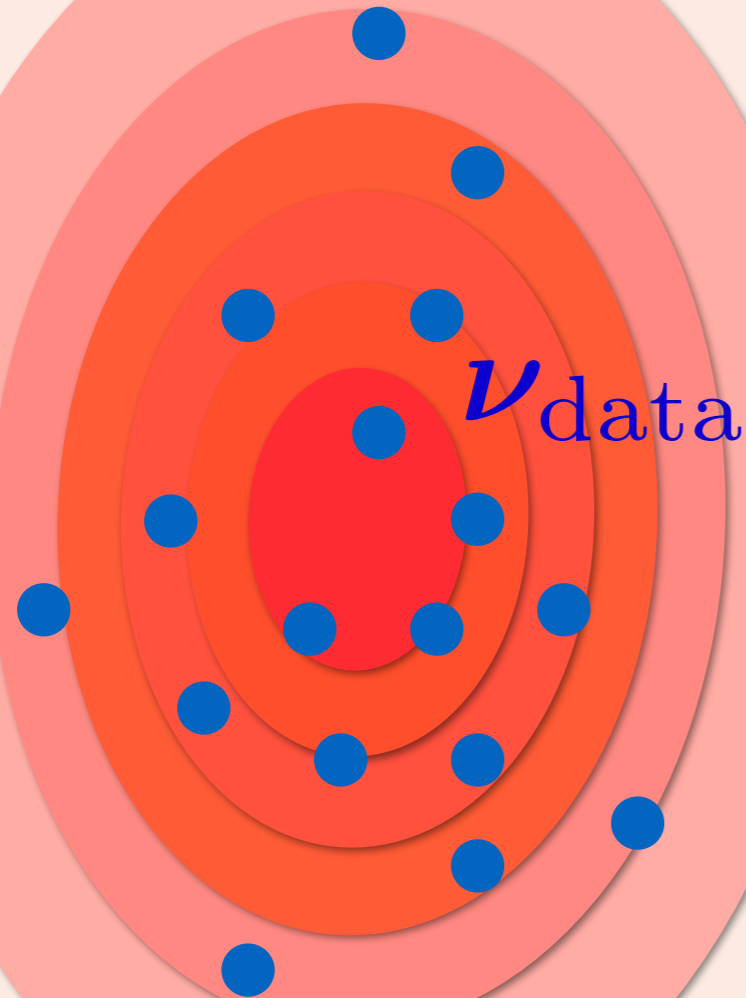
$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i)$$

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



p_{θ} done!

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

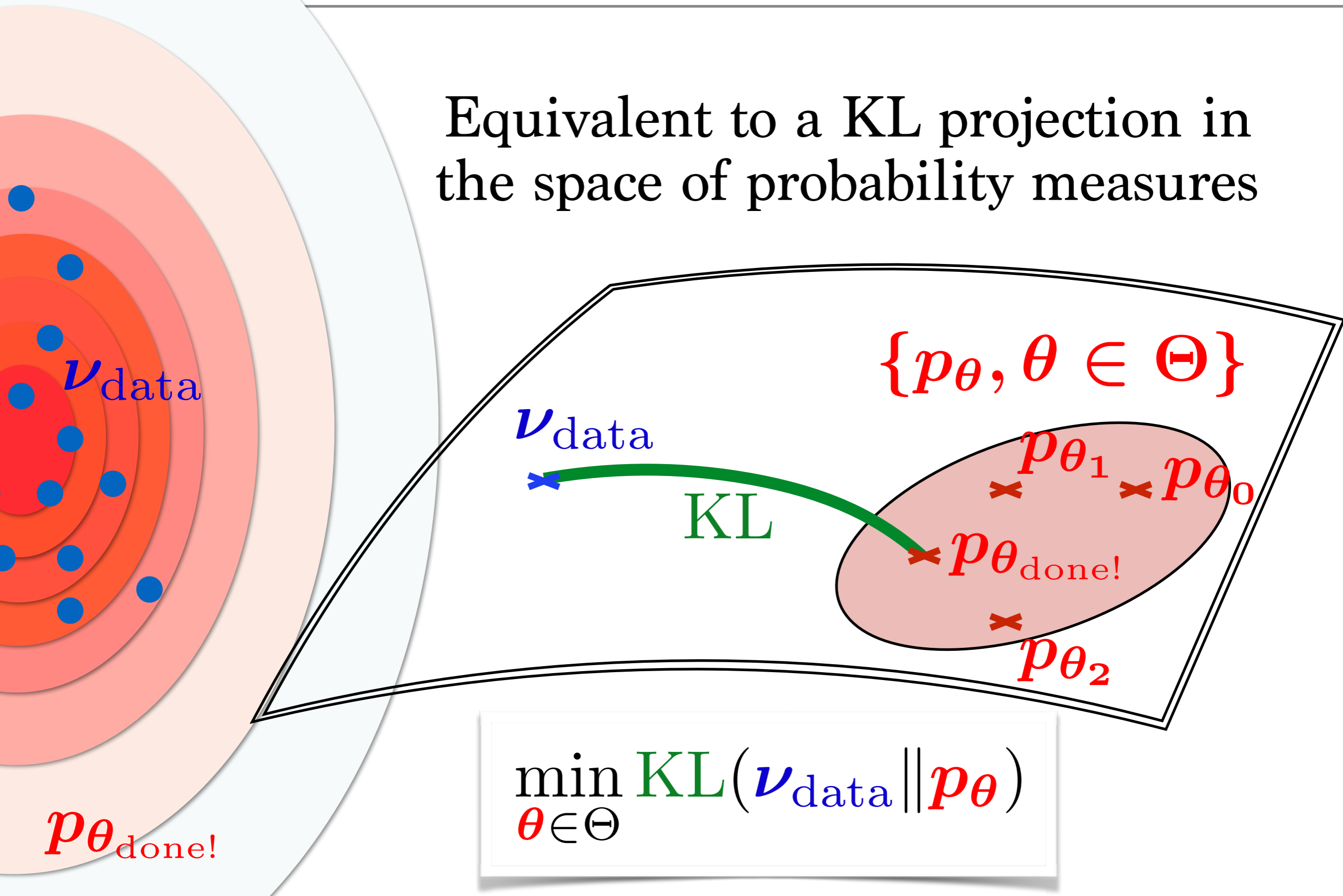


$\log 0 = -\infty$

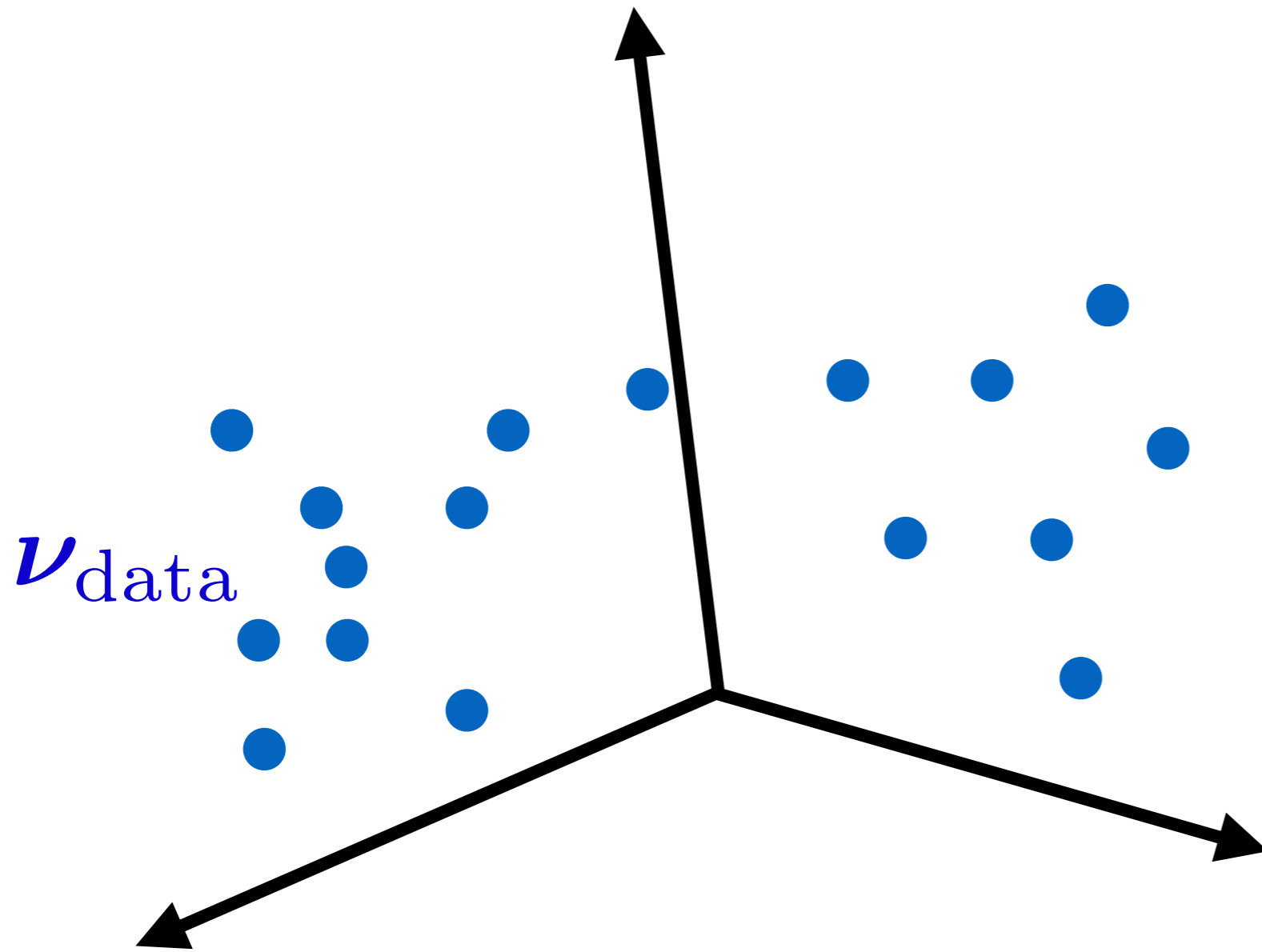
$p_{\theta}(x_i)$ must be > 0

Maximum Likelihood Estimation

Equivalent to a KL projection in the space of probability measures

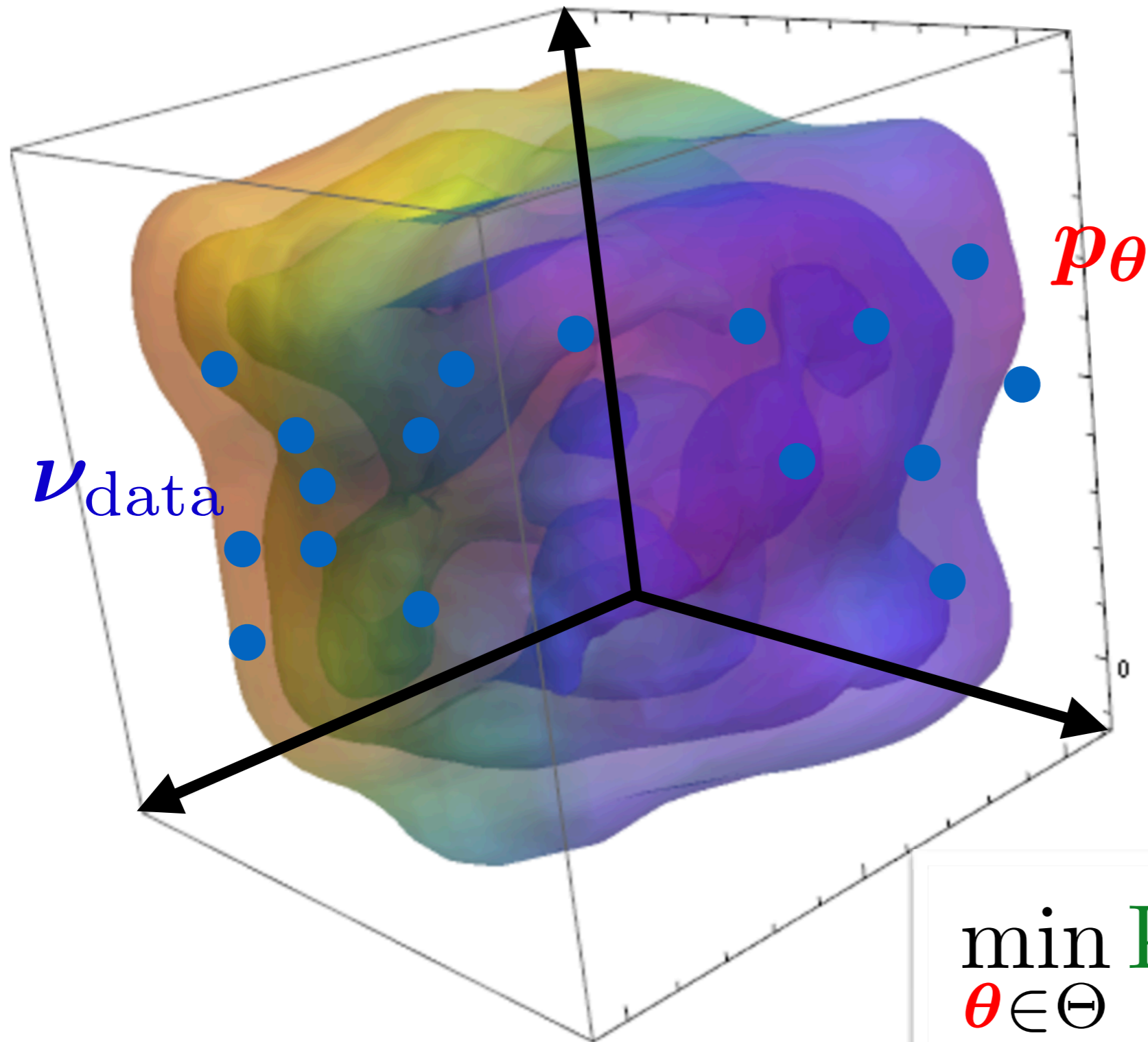


In higher dimensional spaces...



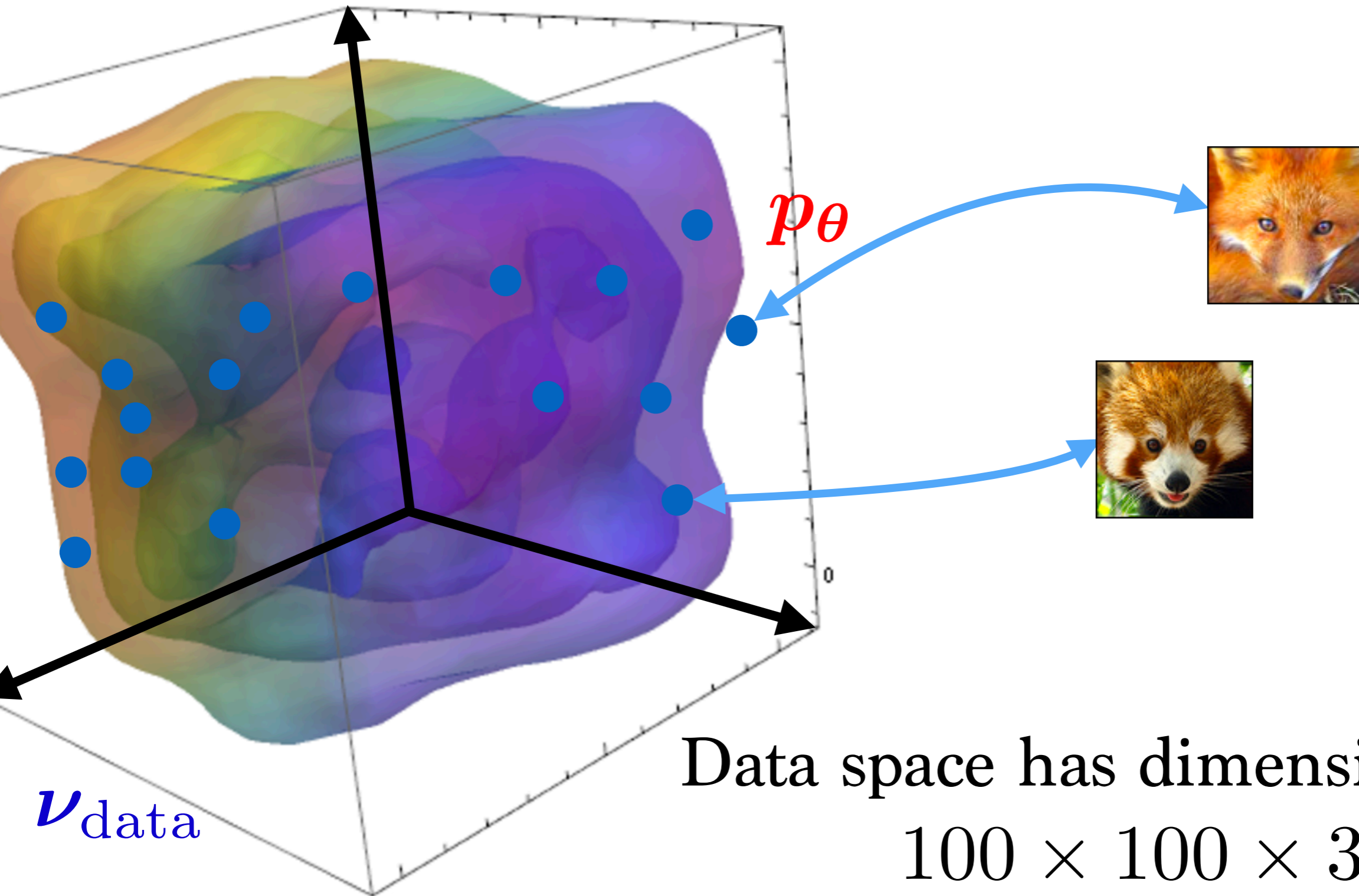
$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \parallel p_{\theta})$$

In higher dimensional spaces...

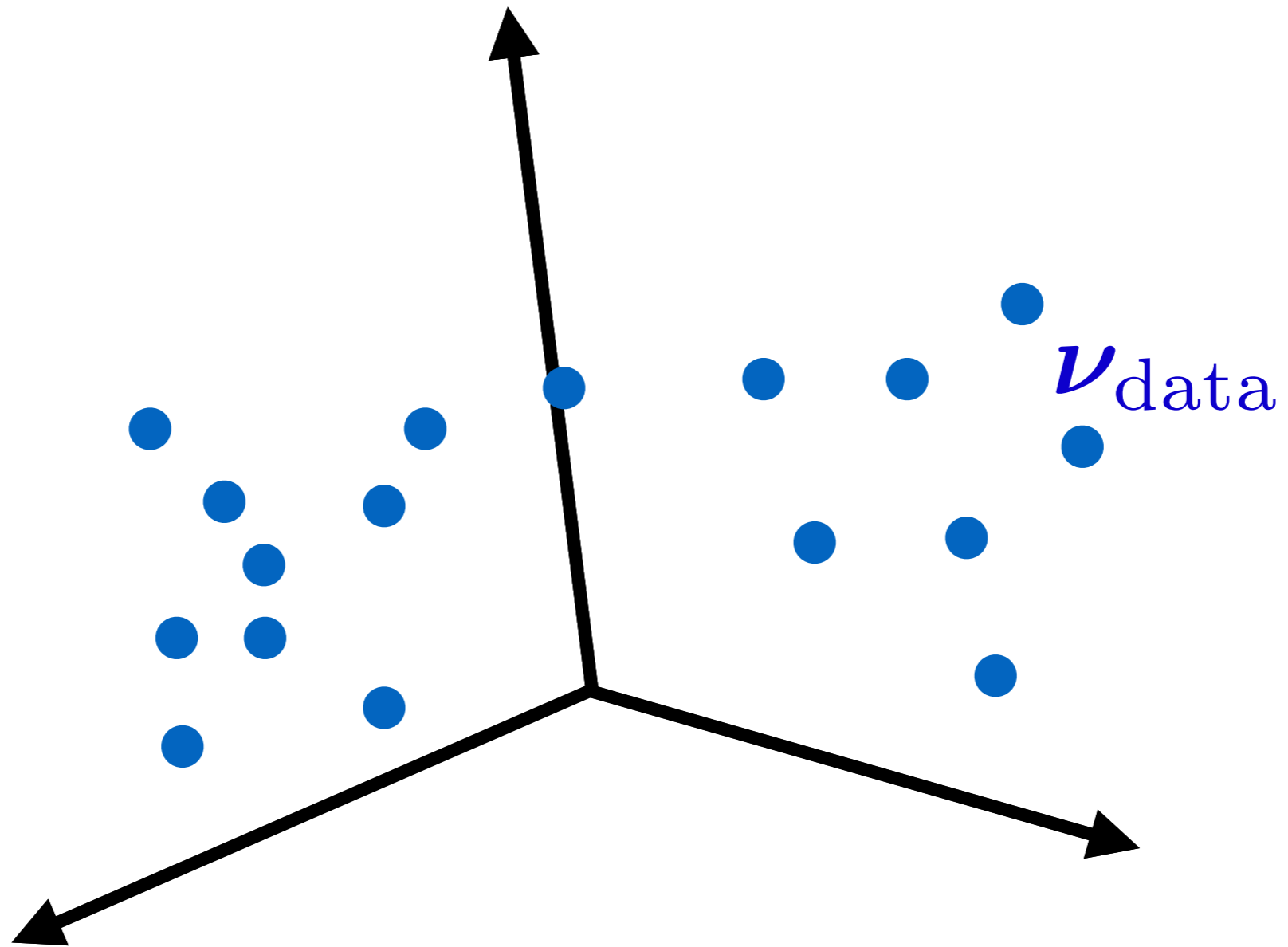


$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \parallel p_{\theta})$$

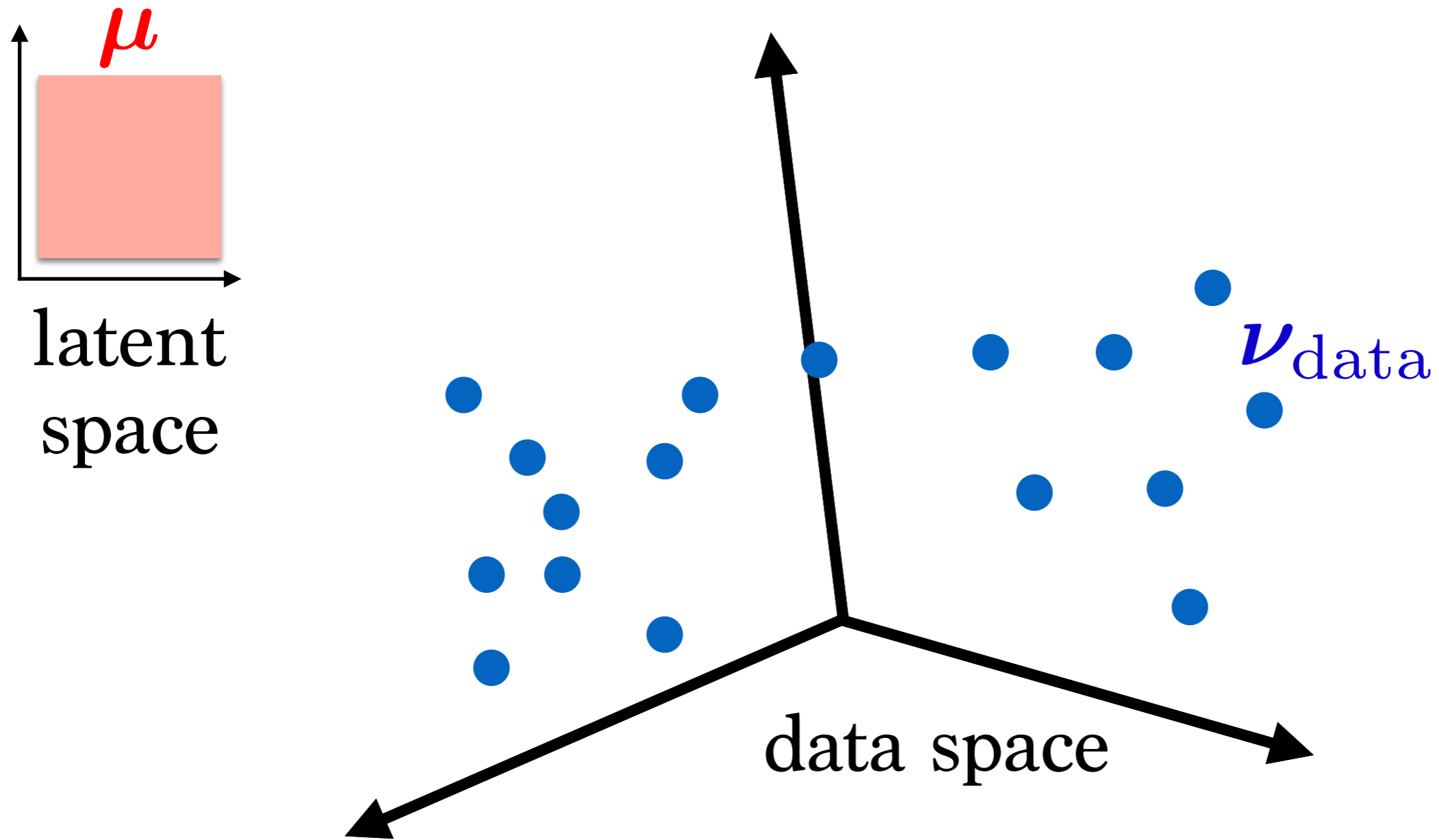
In higher dimensional spaces...



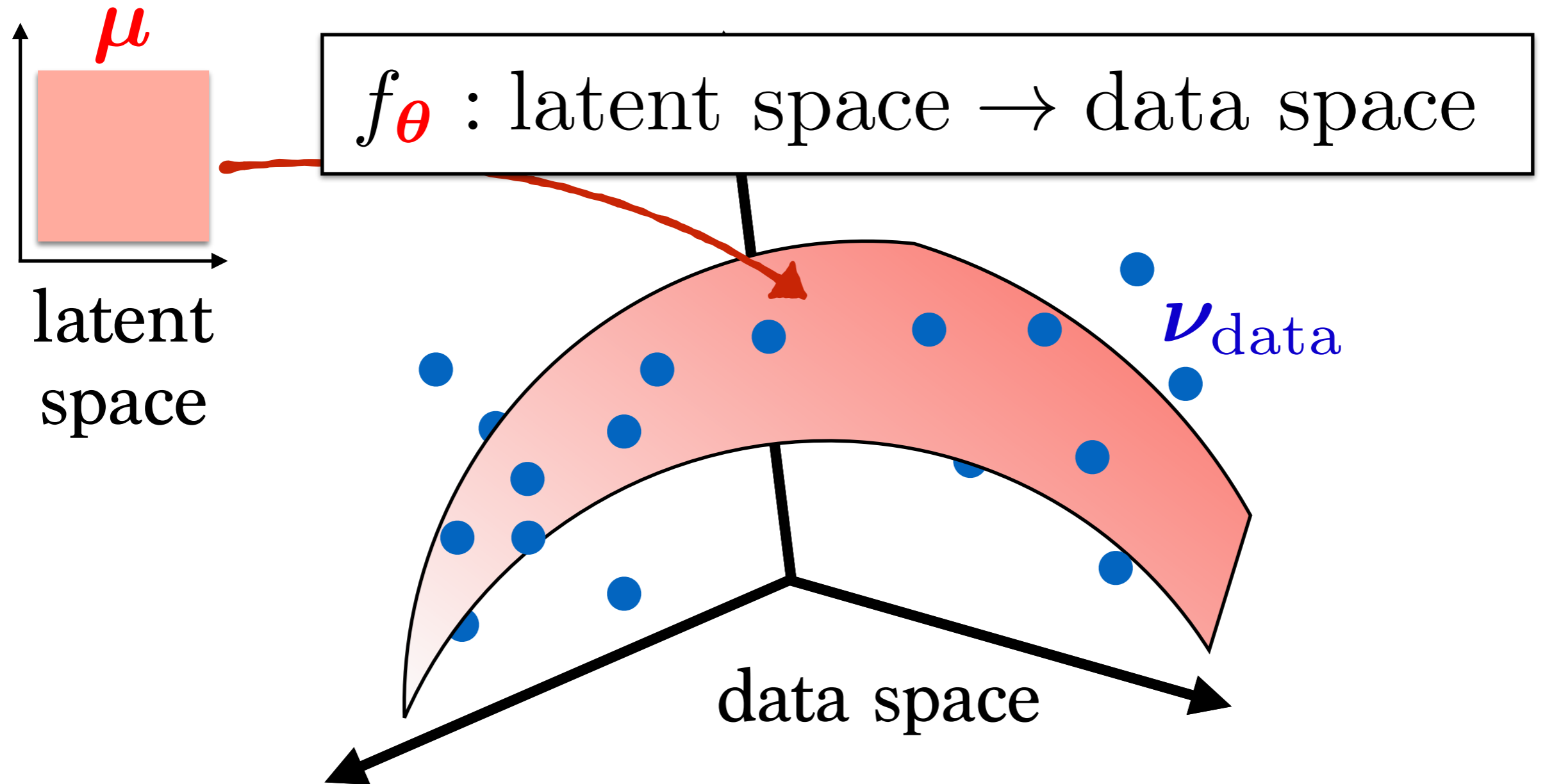
Generative Models



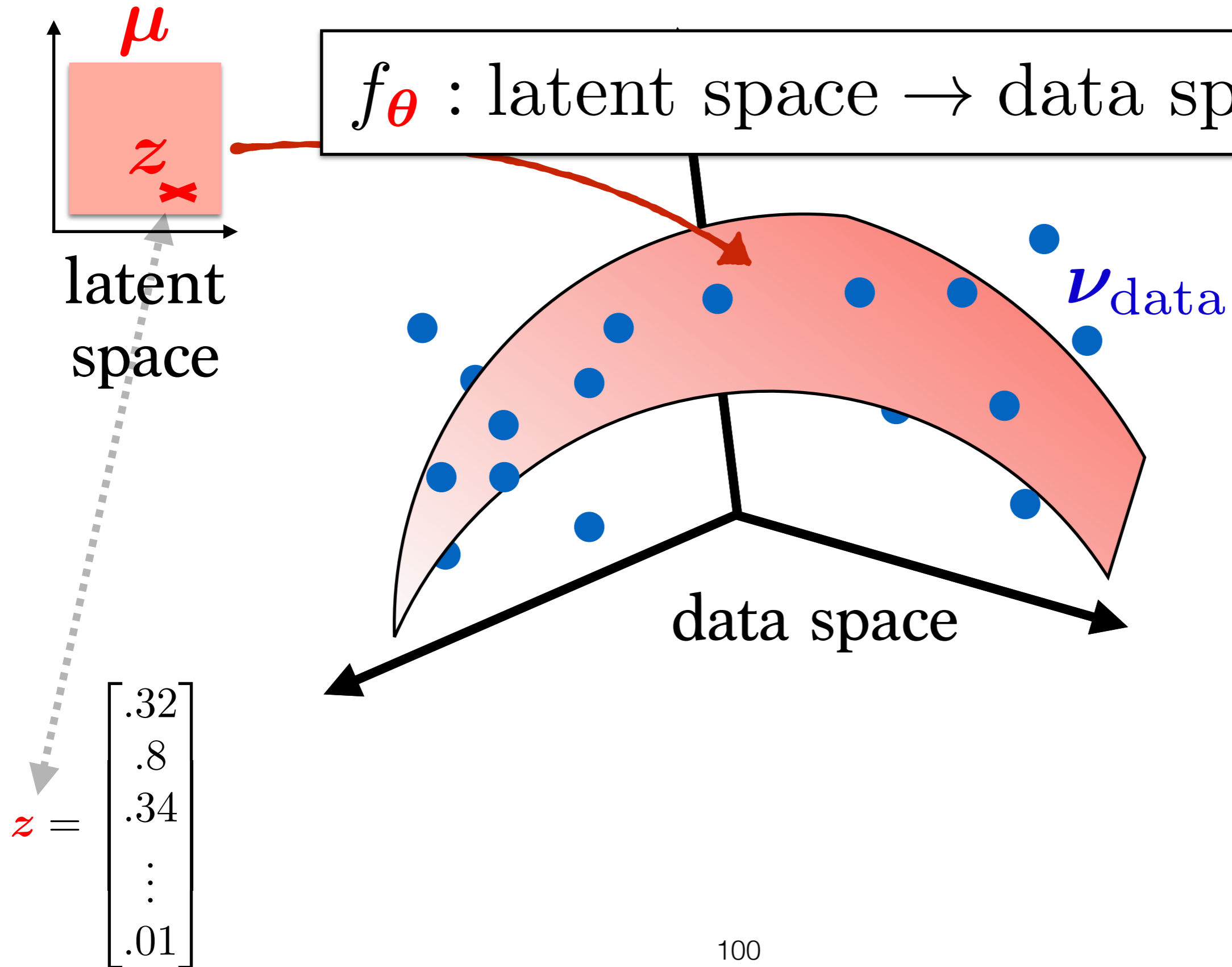
Generative Models



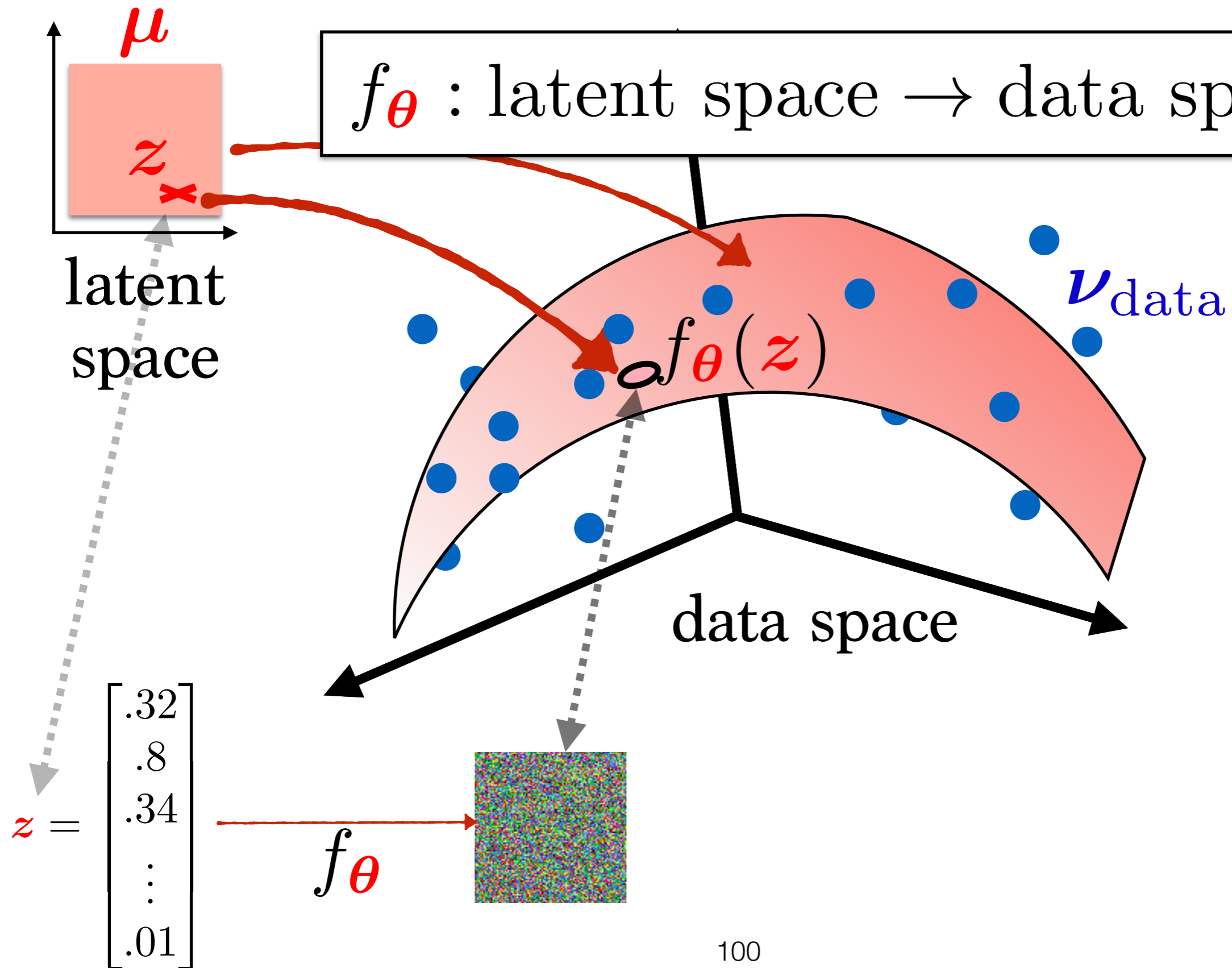
Generative Models



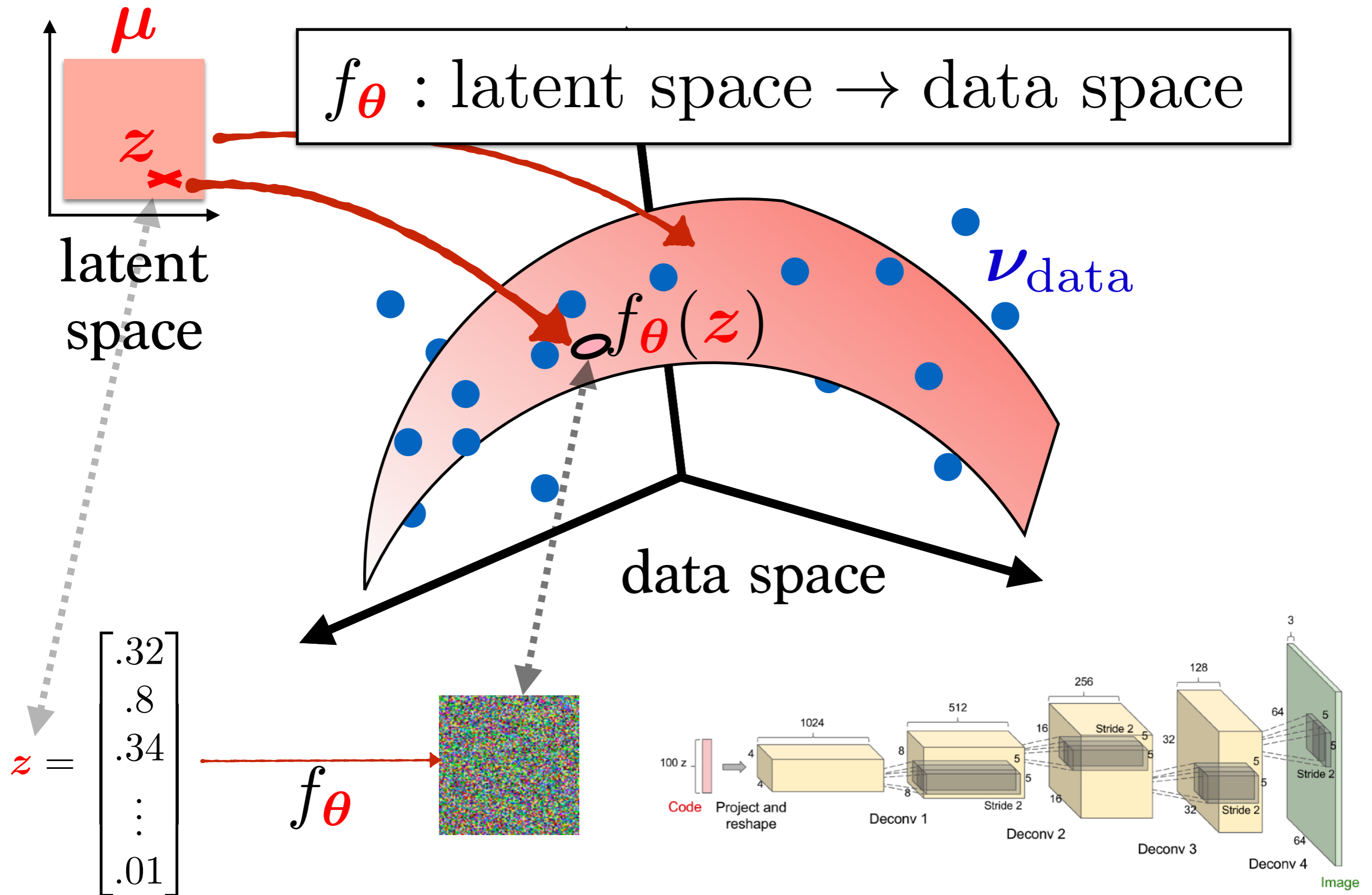
Generative Models



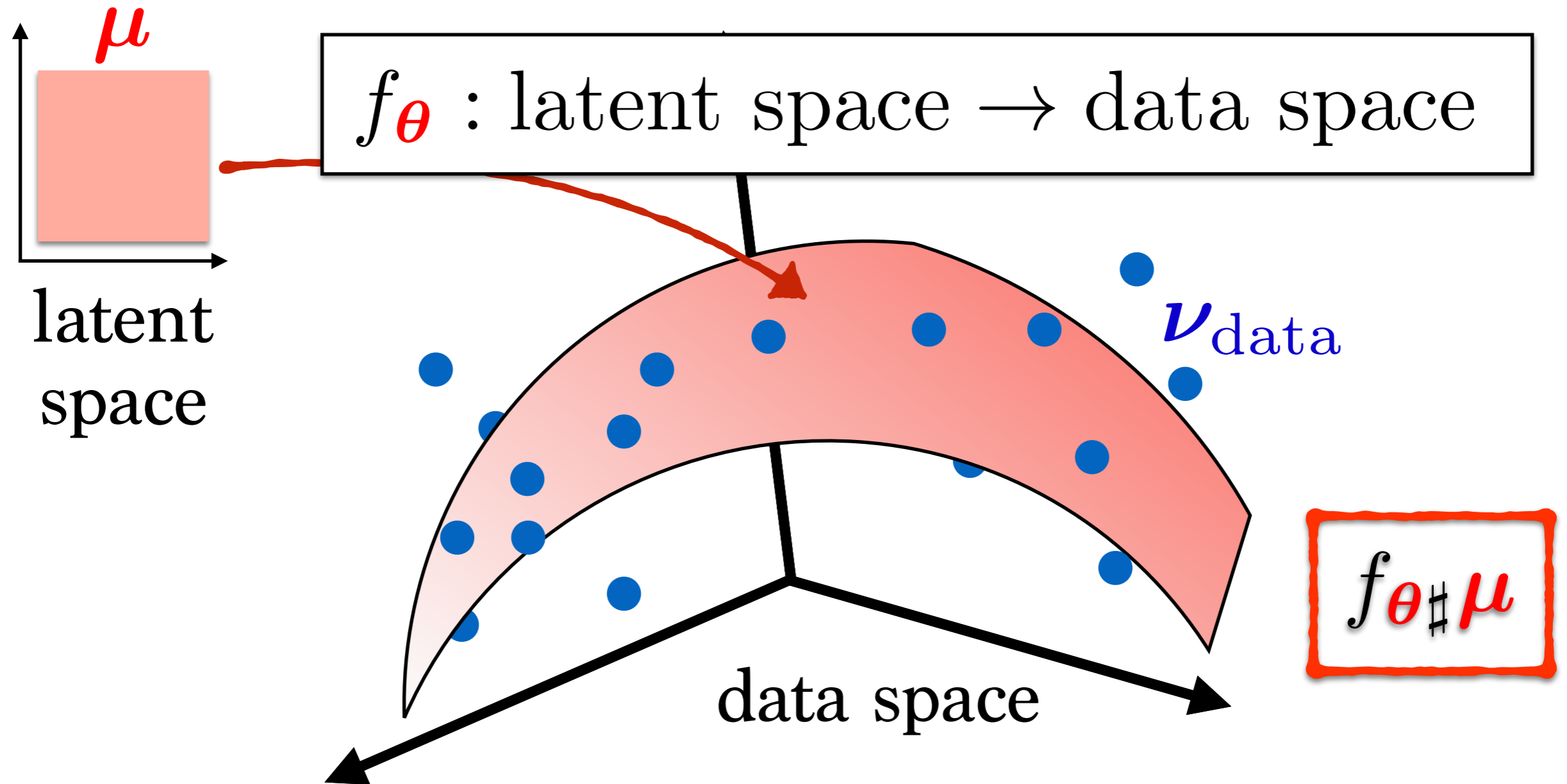
Generative Models



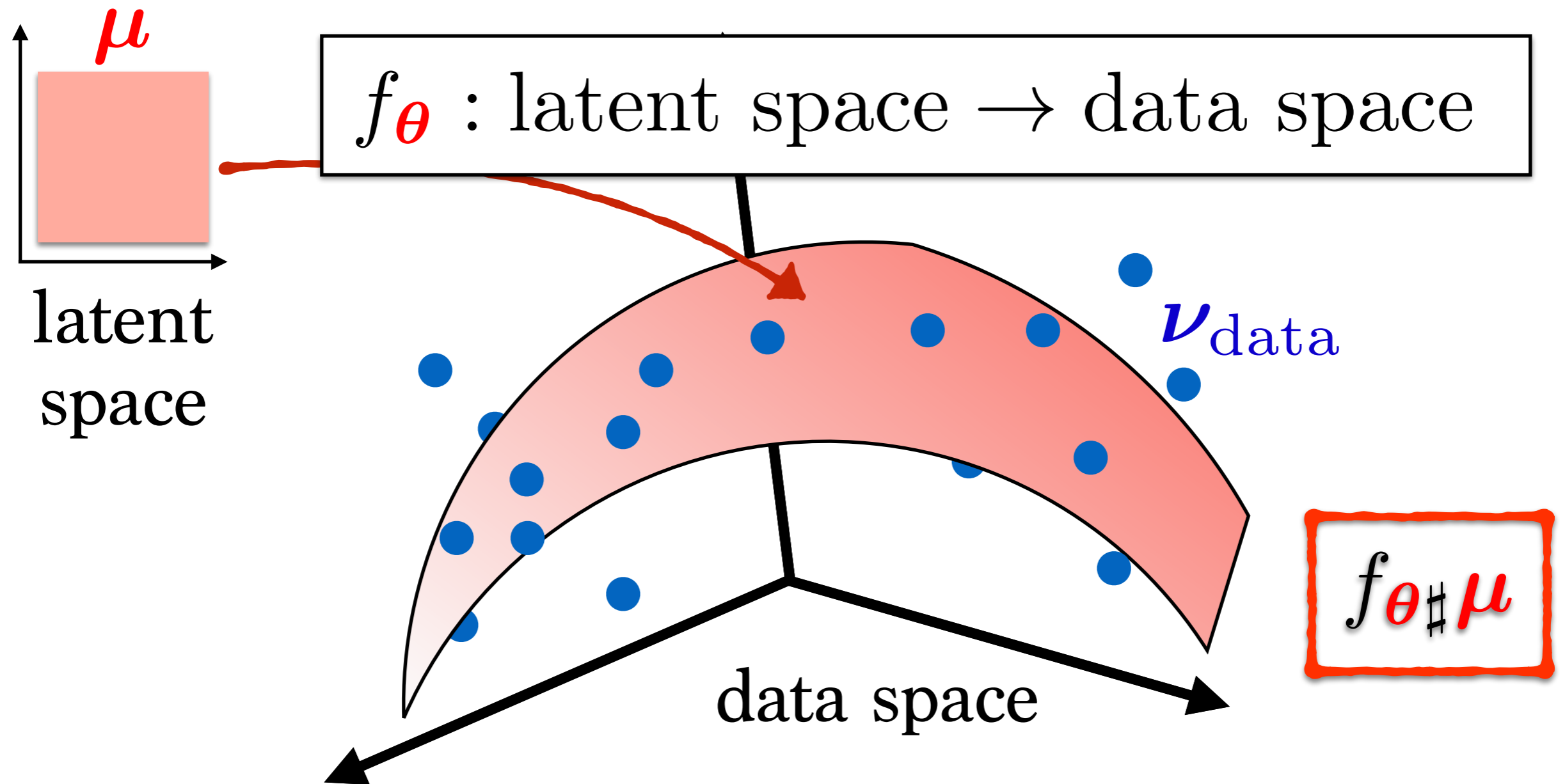
Generative Models



Generative Models

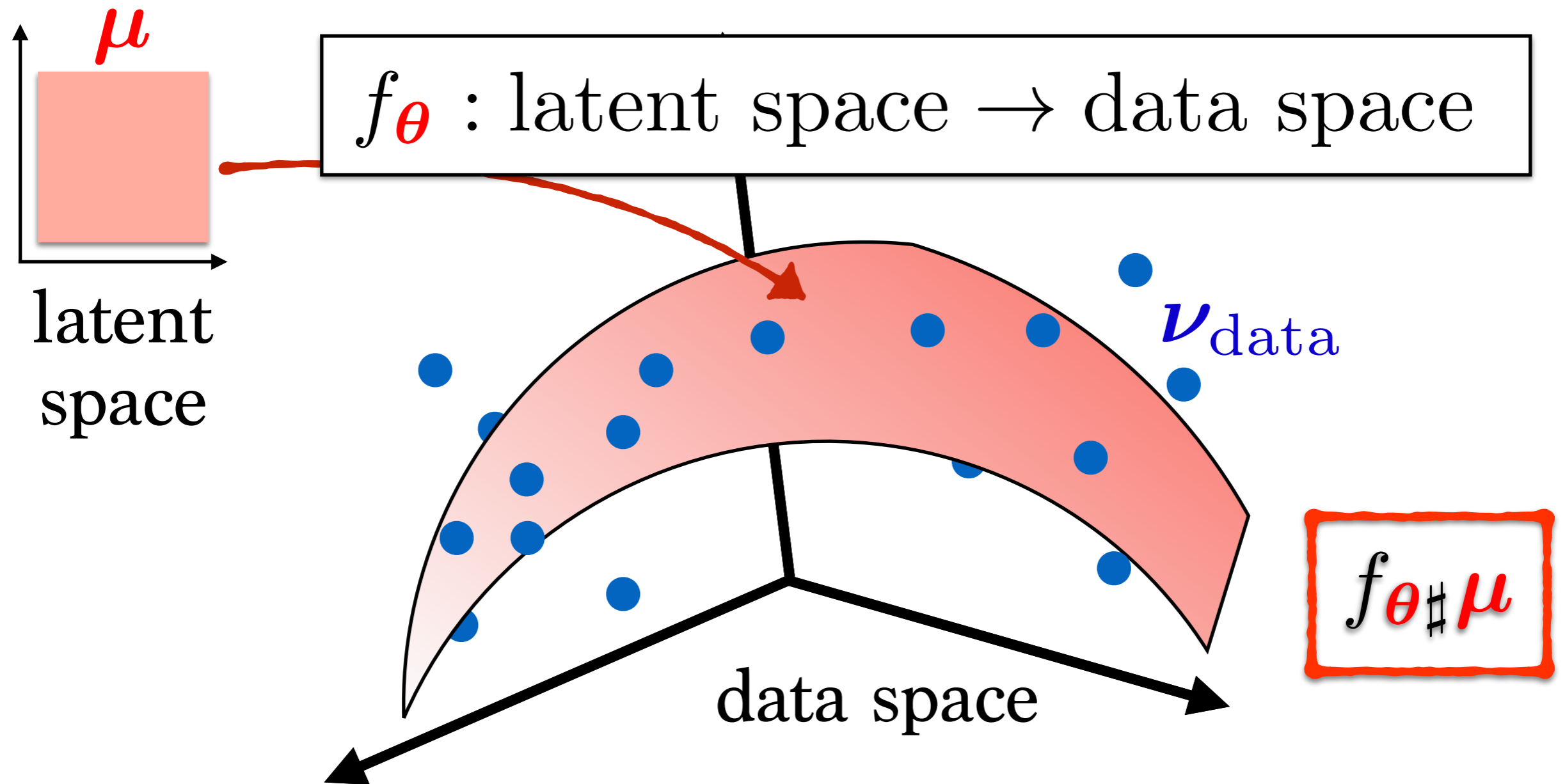


Generative Models



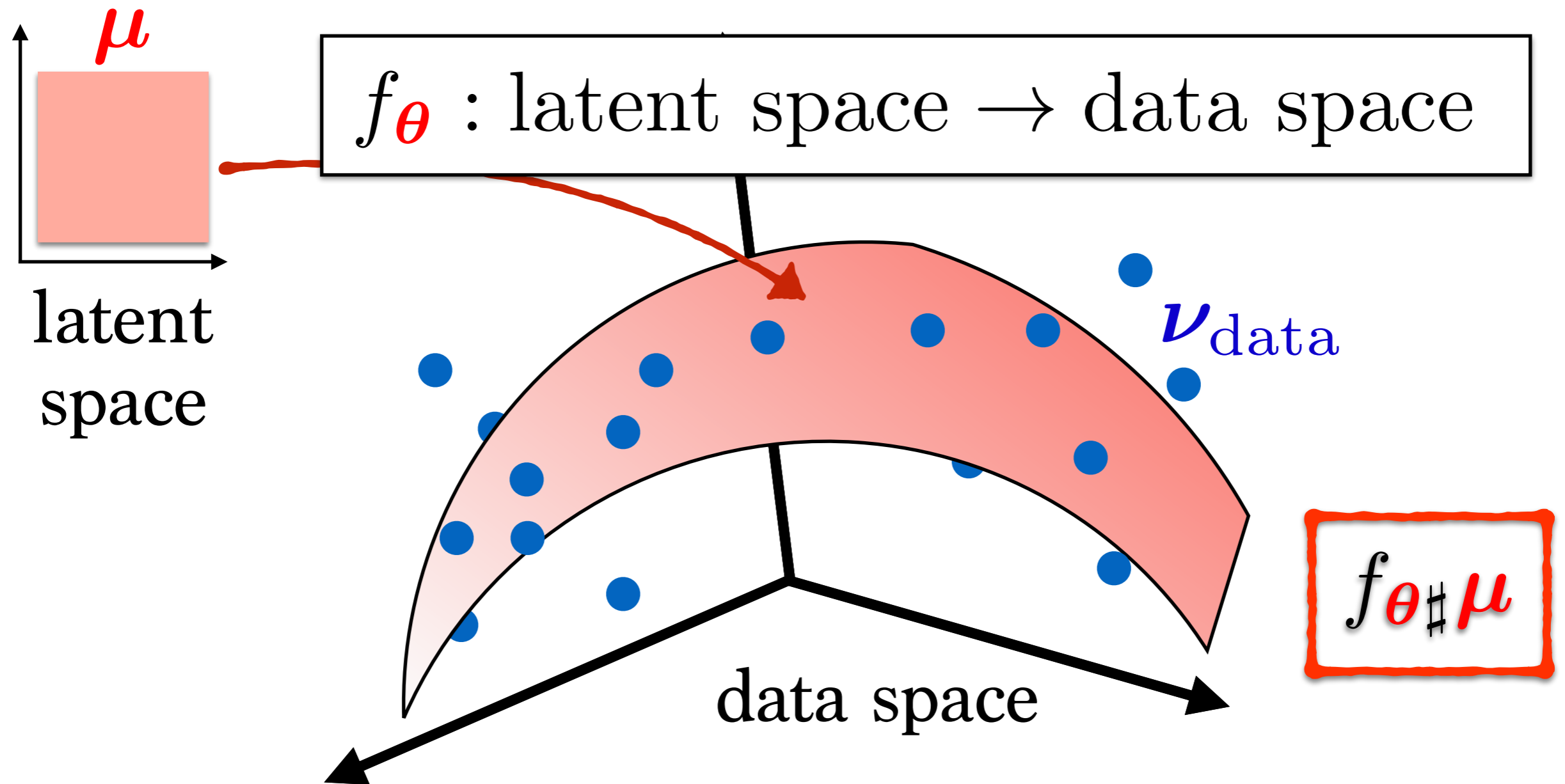
Push-forward: $\forall B \subset \Omega, f_{\#}\mu(B) := \mu(f^{-1}(B))$

Generative Models



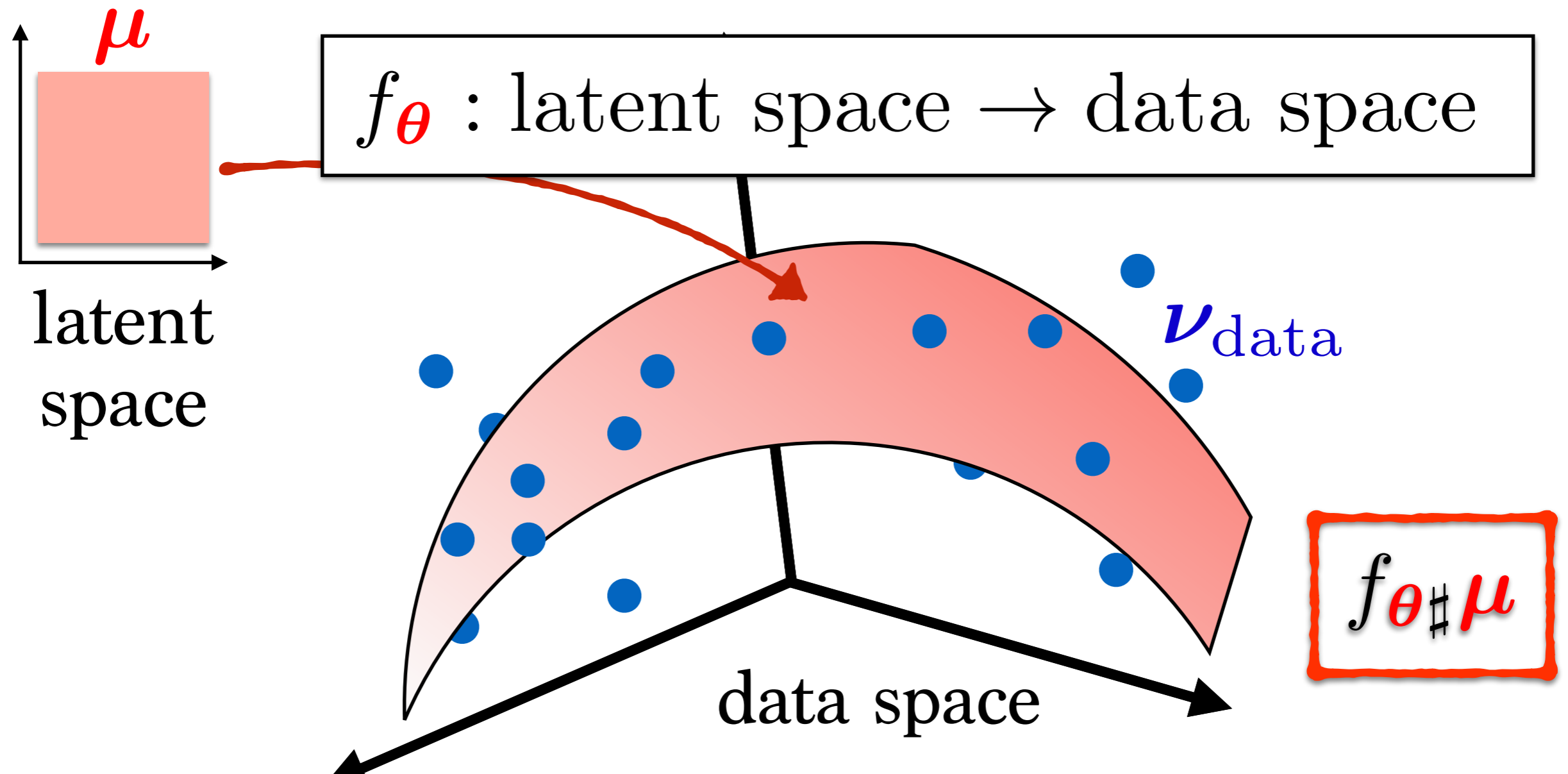
Goal: find θ such that $f_{\theta \# \mu}$ fits $\mathcal{V}_{\text{data}}$

Generative Models



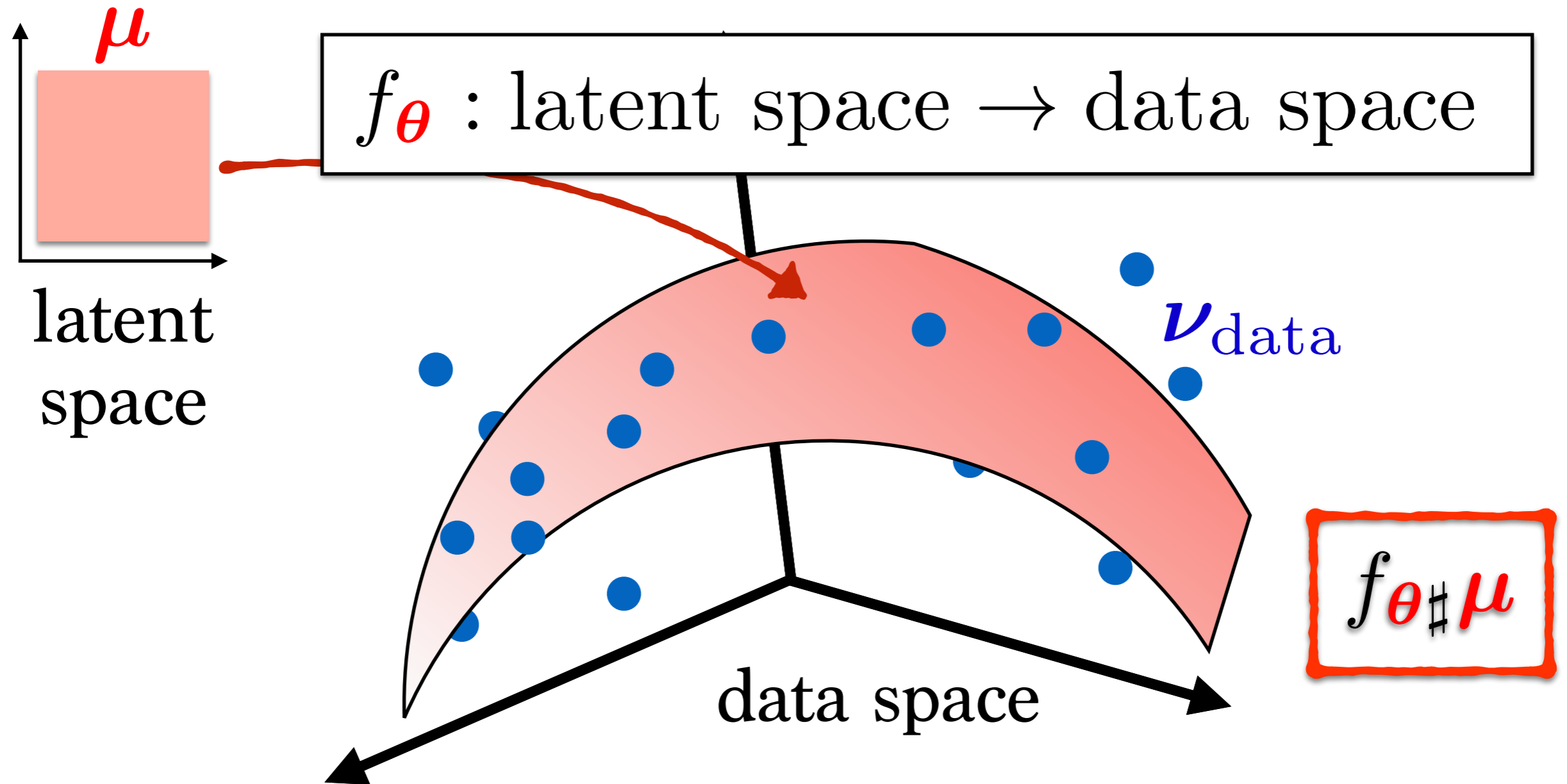
Goal: find θ such that $f_{\theta \# \mu}$ fits $\mathcal{V}_{\text{data}}$

Generative Models



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

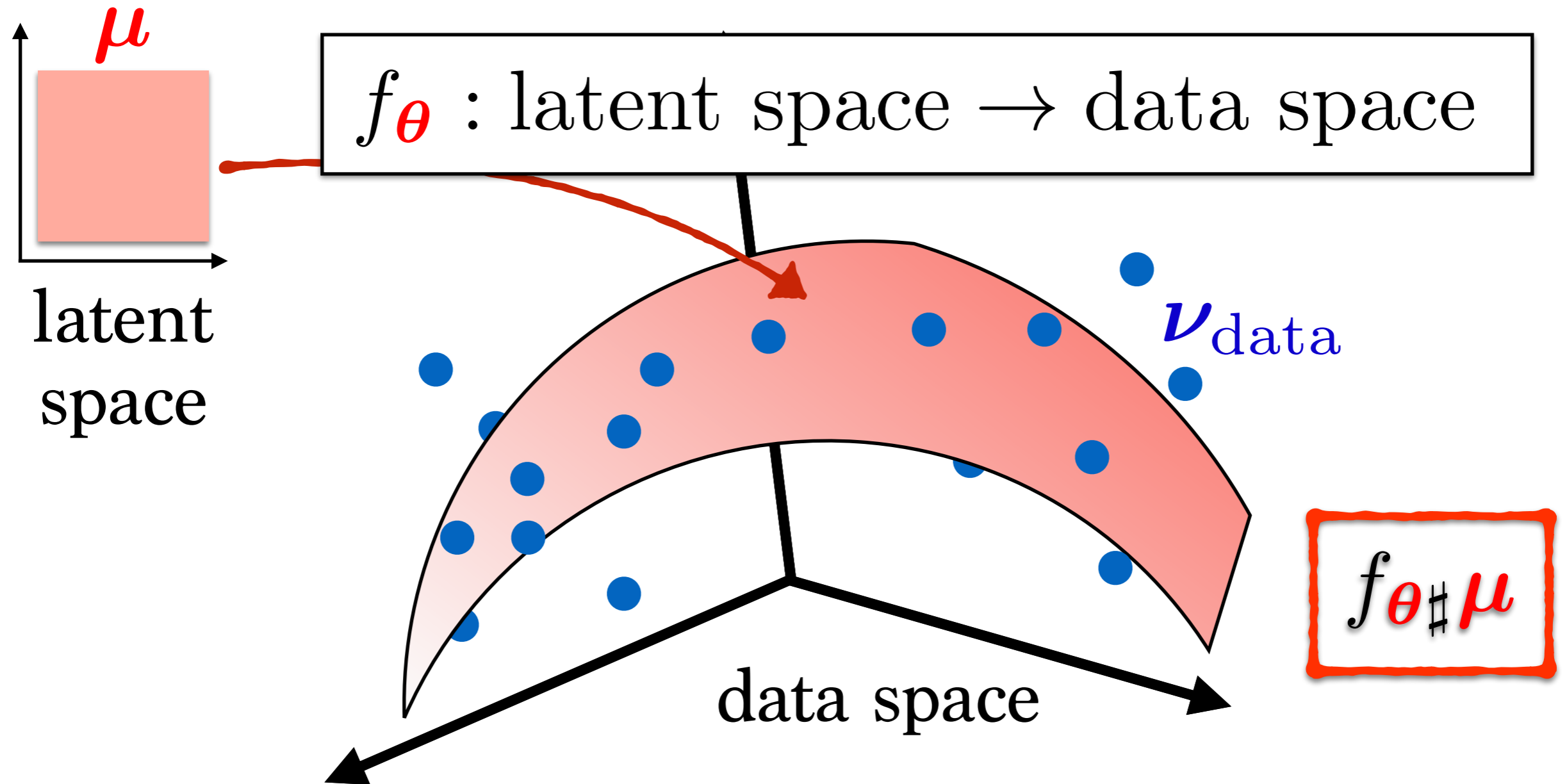
Generative Models



MLE

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

Generative Models

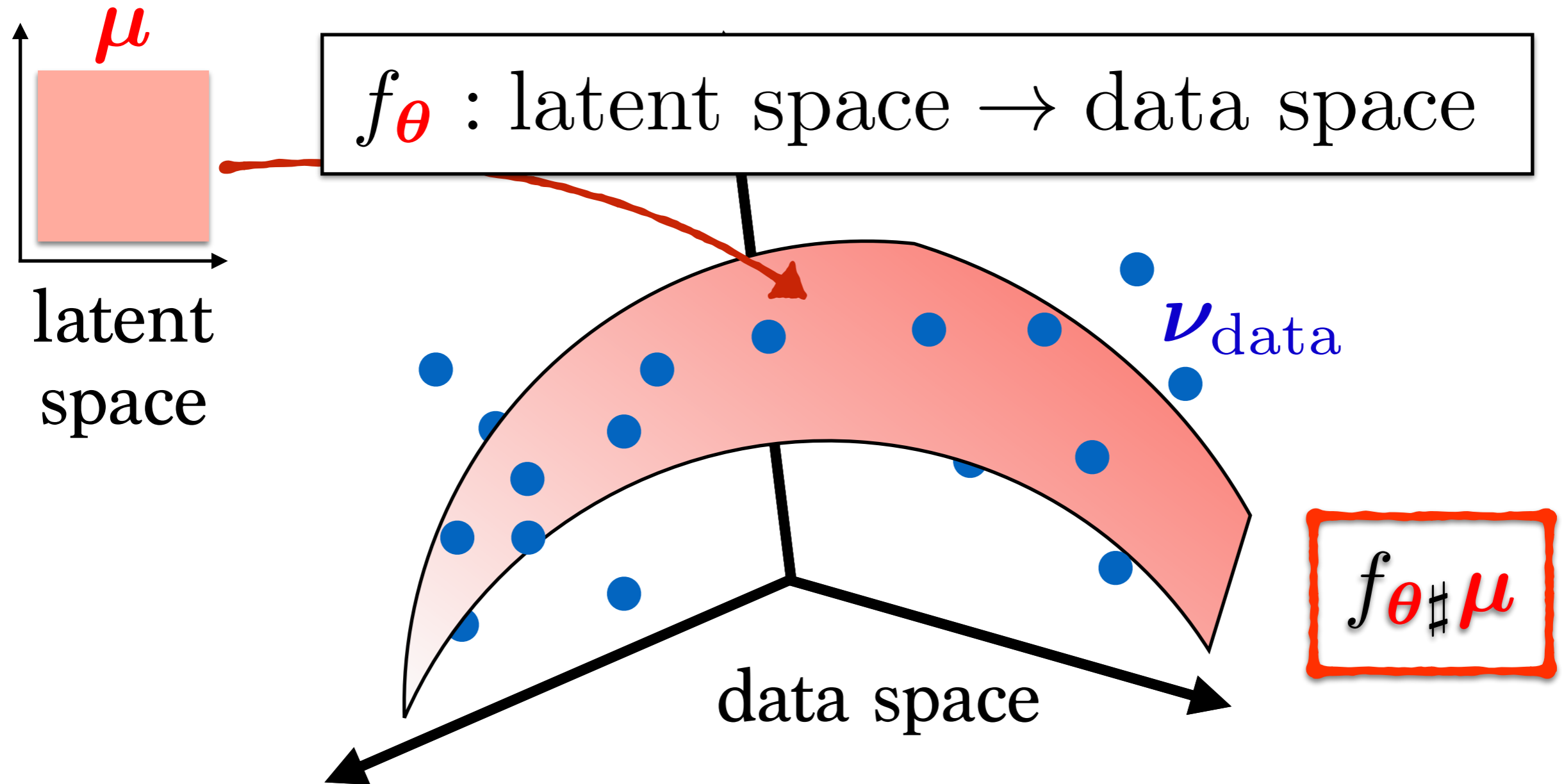


MLE

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

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

Generative Models

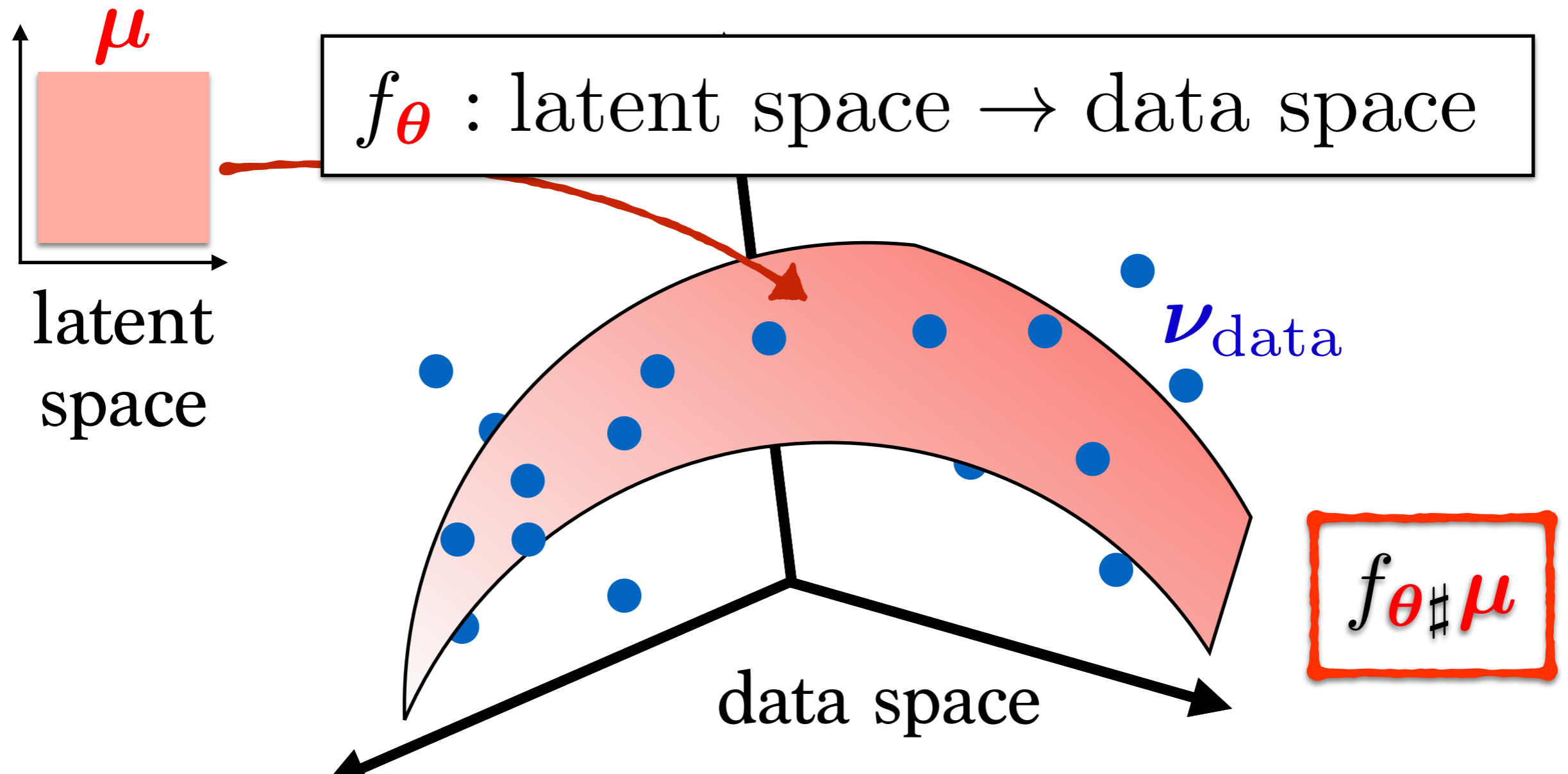


~~MLE~~

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

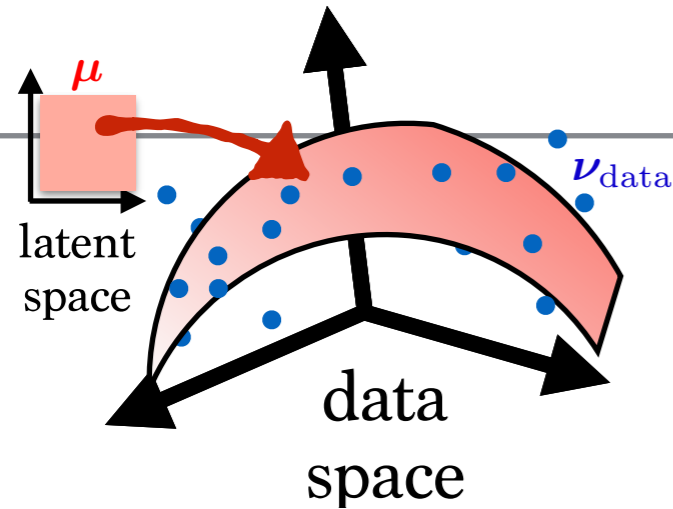


Generative Models



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

Workarounds?



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

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

- Use a richer metric Δ for measures, able to handle measures with non-overlapping supports.

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

Minimum Δ Estimation

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

MINIMUM l_1 CHI-SQUARE, NOT MAXIMUM LIKELIHOOD!

BY JOSEPH BERKSON

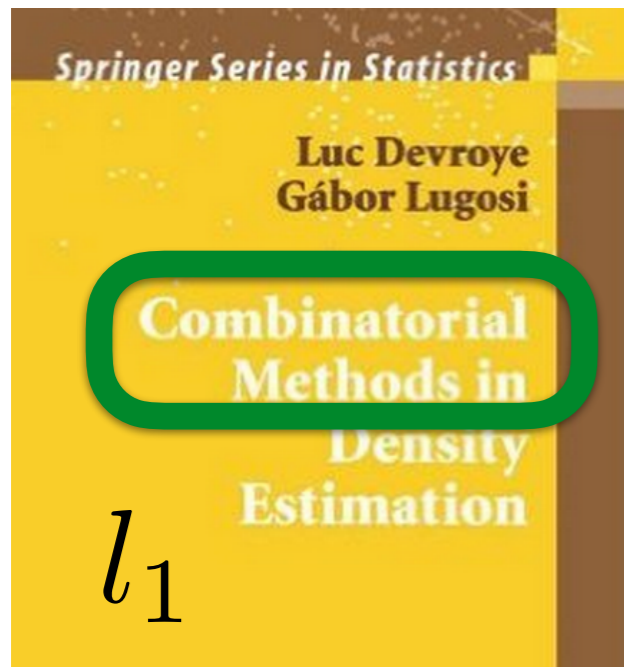
Mayo Clinic, Rochester, Minnesota



ELSEVIER

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

COMPUTATIONAL
STATISTICS
& DATA ANALYSIS



Minimum 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

SCIENCE @ DIRECT®



ELSEVIER

Statistics & Probability Letters 76 (2006) 1298–1302

STATISTICS &
PROBABILITY
LETTERS

www.elsevier.com/locate/stapro

On minimum Kantorovich distance estimators

Federico Bassetti^a, Antonella Bodini^b, Eugenio Regazzini^{a,*}

Minimum Kantorovich Estimation



ELSEVIER

Available online at www.sciencedirect.com

SCIENCE @ DIRECT®

Statistics & Probability Letters 76 (2006) 1298–1302

STATISTICS &
PROBABILITY
LETTERS

www.elsevier.com/locate/stapro

On minimum Kantorovich distance estimators

Federico Bassetti^a, Antonella Bodini^b, Eugenio Regazzini^{a,*}

Use *Wasserstein distances* to define a loss between data and model.

$$\min_{\theta \in \Theta} W(\nu_{\text{data}}, p_{\theta})$$

Minimum Kantorovich Estimators

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

[Bassetti'06] 1st reference discussing this approach.

Challenge: $\nabla_{\theta} W(\nu_{\text{data}}, f_{\theta\#} \mu)$?

[Montavon'16] use regularized OT in a finite setting.

[Arjovsky'17] (*WGAN*) uses a NN to approximate dual solutions and recover gradient w.r.t. parameter

[Bernton'17] (*Wasserstein ABC*)

[Genevay'17, Salimans'17] (*Sinkhorn approach*)

Concluding Remarks

1. Introduction to optimal transport
2. Optimal transport algorithms
3. Applications (W as a loss)
4. Applications (W for estimation)

This Saturday: OT & ML Workshop

7 Talks by: Jacob, Kraig, Andoni, Gangbo, Bottou, Flamary, Bach

17 posters and spotlight presentations.

Organizers: Bousquet, Cuturi, Peyré, Sha, Solomon

What we could not talk about...

- almost infinite supply of **maths**...
- **Statistical** challenges to compute W .
- If **linear assignment** = Wasserstein, then **quadratic assignment** = Gromov-Wasserstein.
- Wasserstein gradient flows (a.k.a. **JKO** flow).
- **Dynamical** aspects of optimal transport
- Transporting vectors and matrices

<https://optimaltransport.github.io/>