

A Primer on Optimal Transport

Marco Cuturi



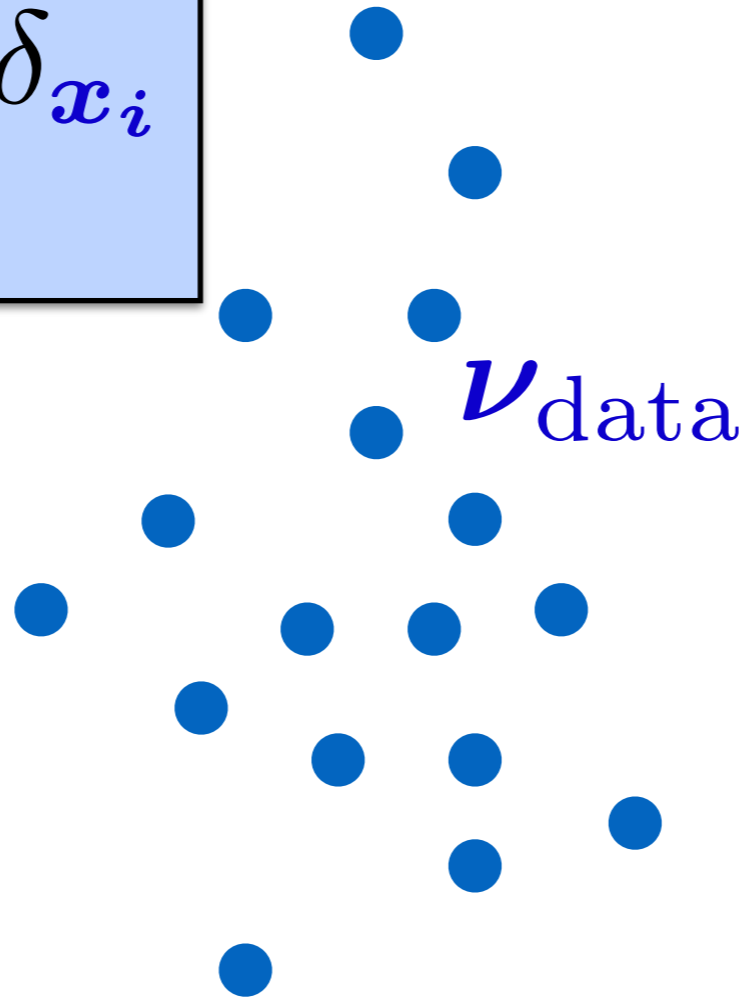
book with Gabriel Peyré

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

A Motivating Example

We collect data

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



A Motivating Example

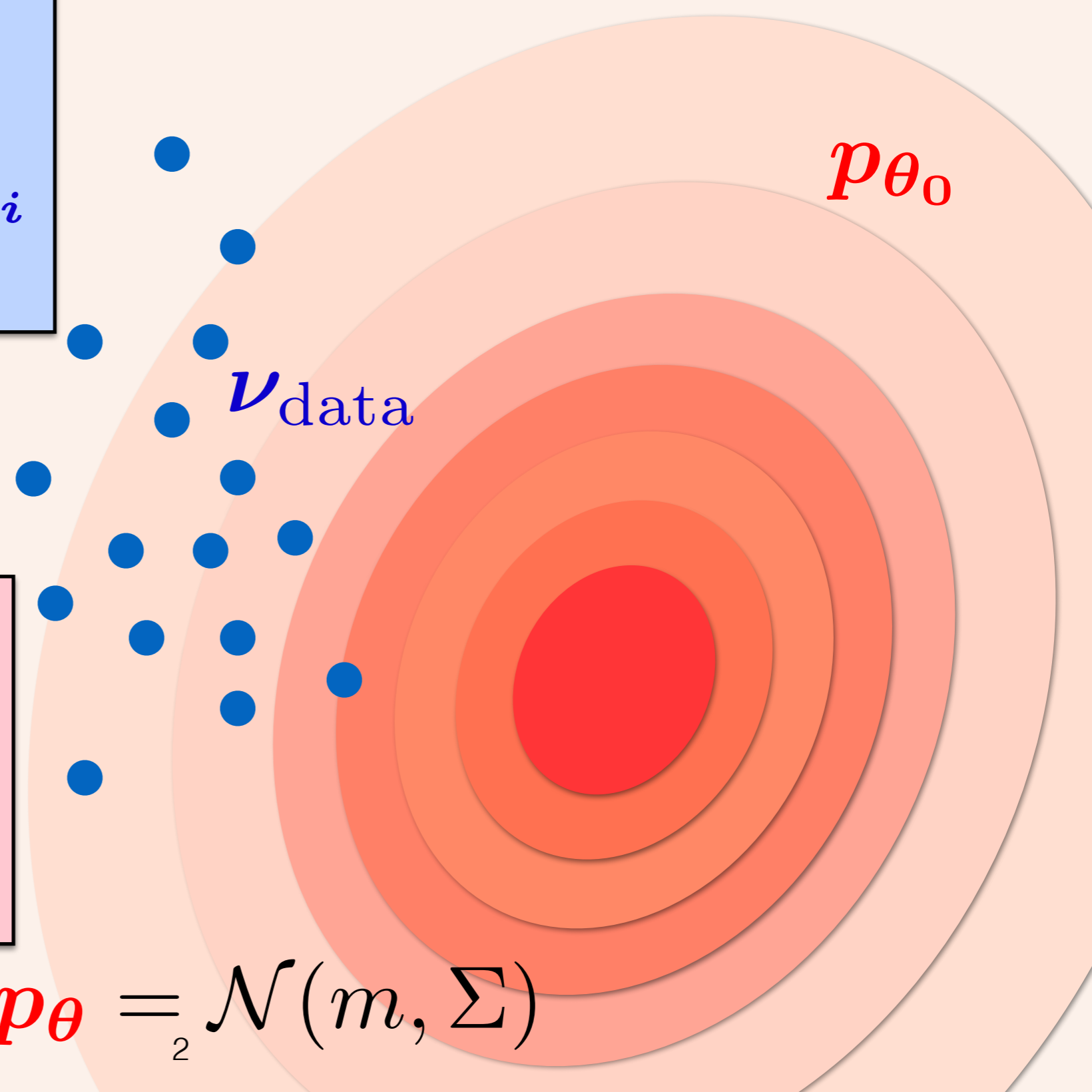
We collect data

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

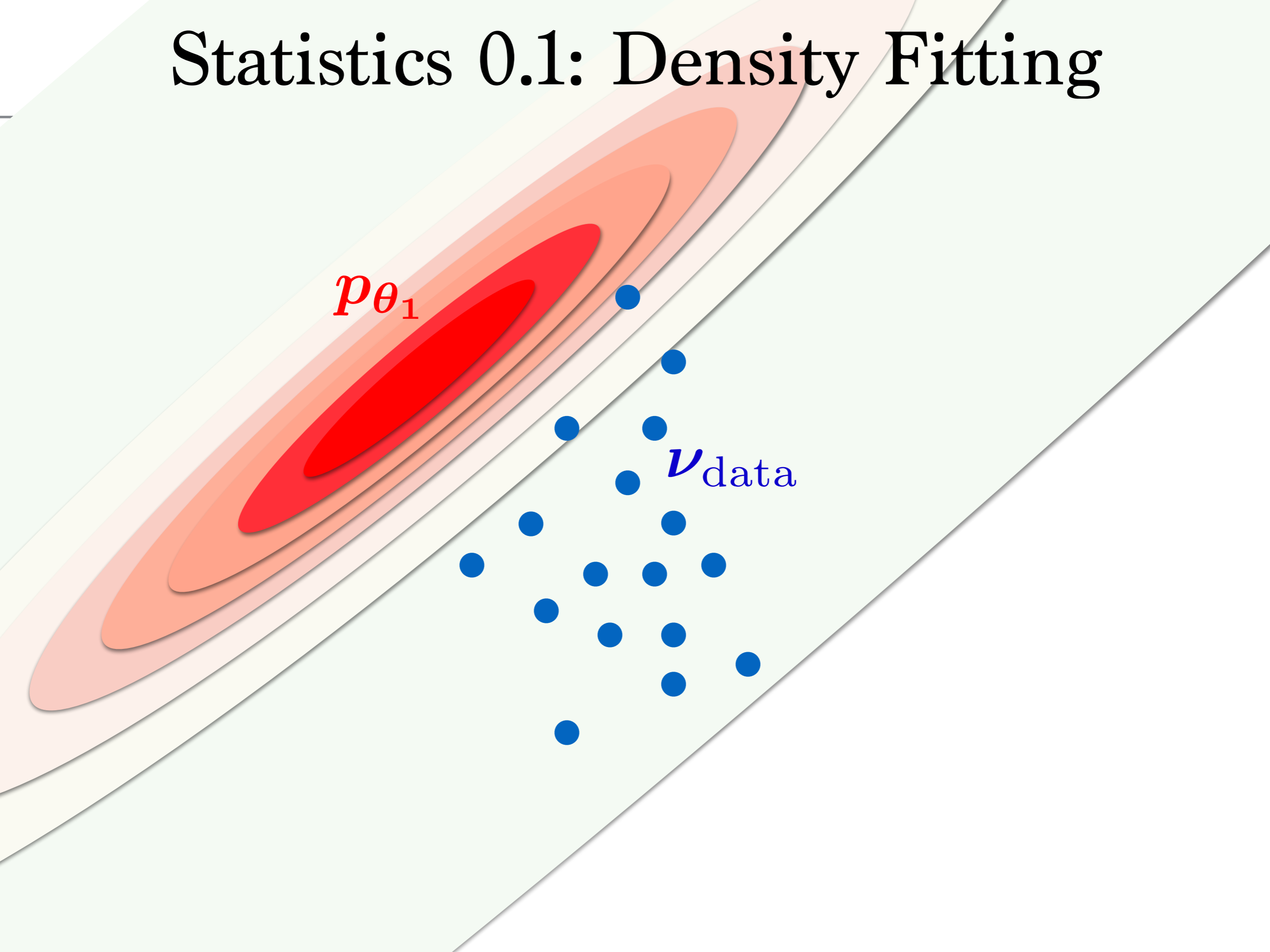
We fit a parametric family of densities

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

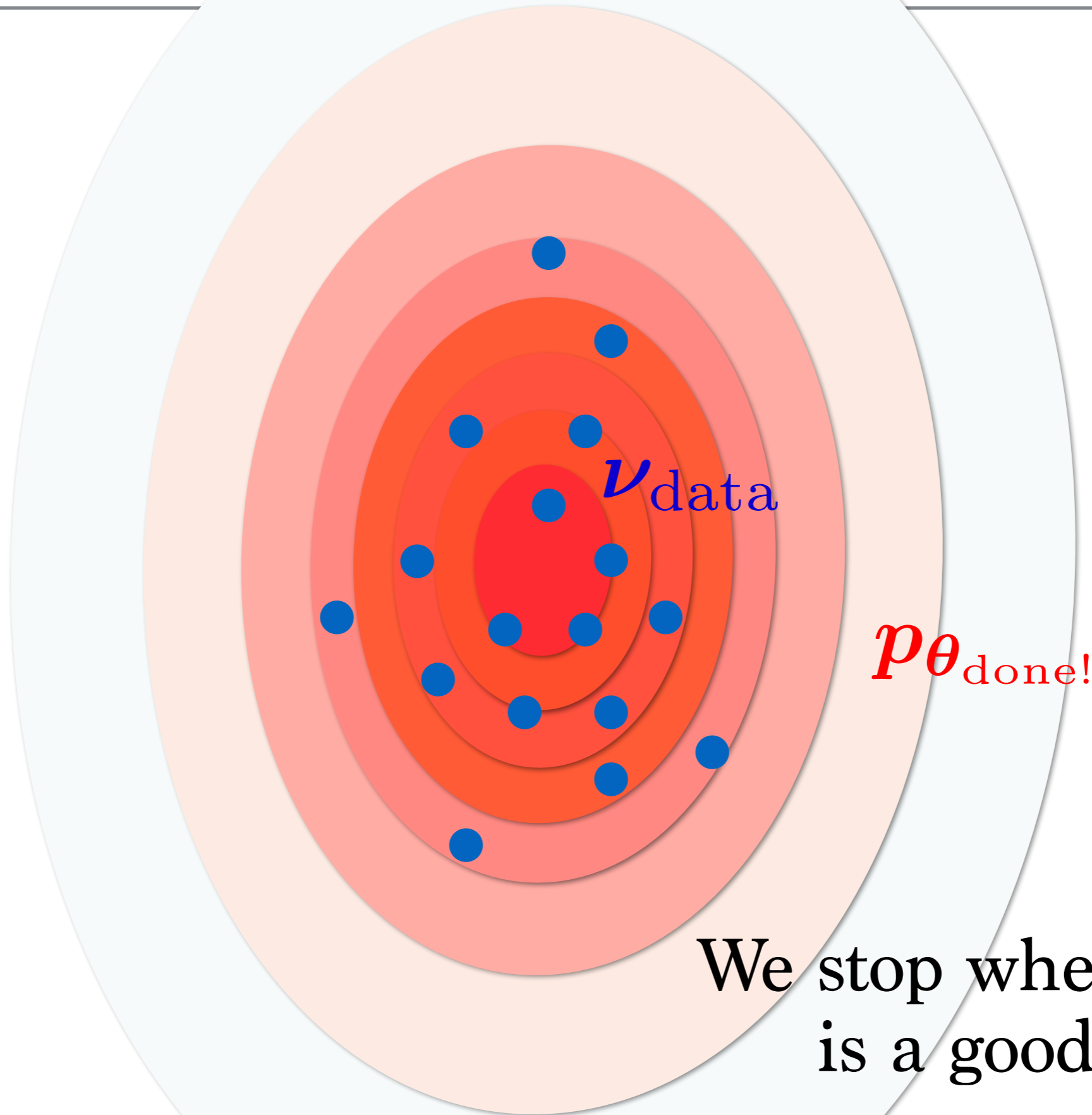
e.g. $\theta = (m, \Sigma); p_{\theta} = \mathcal{N}(m, \Sigma)$



Statistics 0.1: Density Fitting



Statistics 0.1: Density Fitting

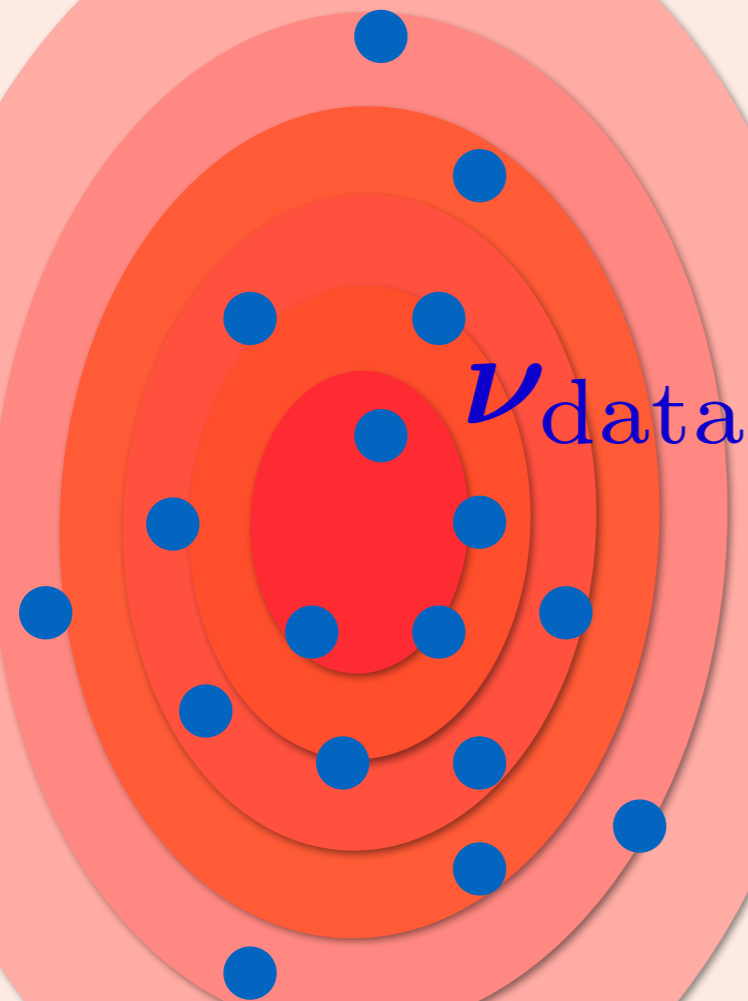


Maximum Likelihood Estimation

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



p_{θ} done!

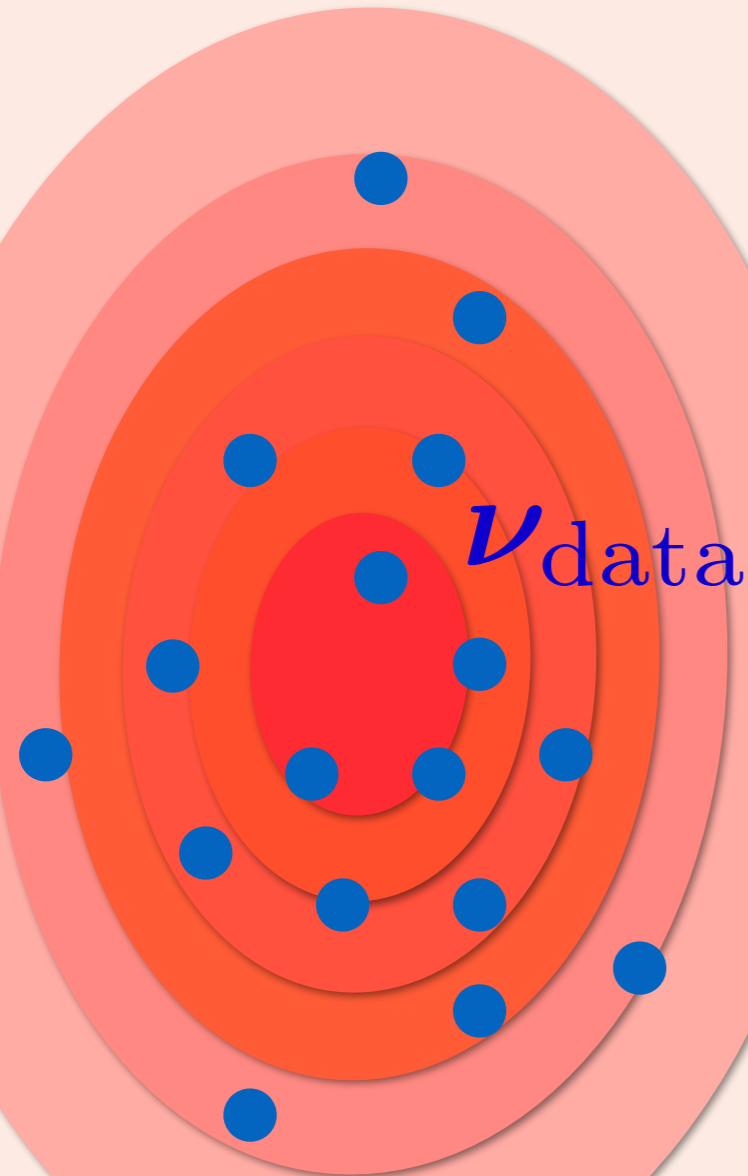
$$\max_{\theta \in \Theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(x_i)$$

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!

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

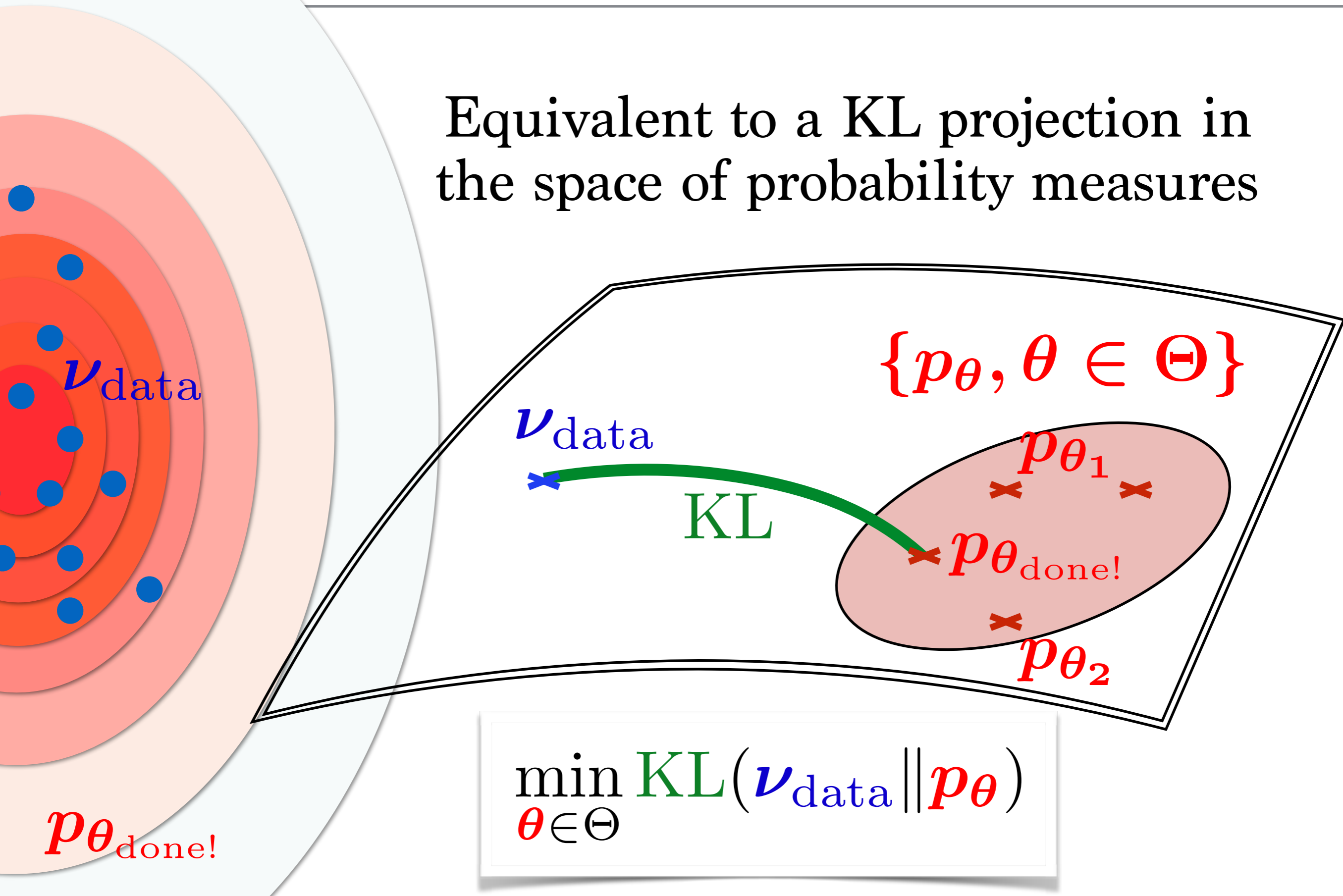


$\log 0 = -\infty$

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

Maximum Likelihood Estimation

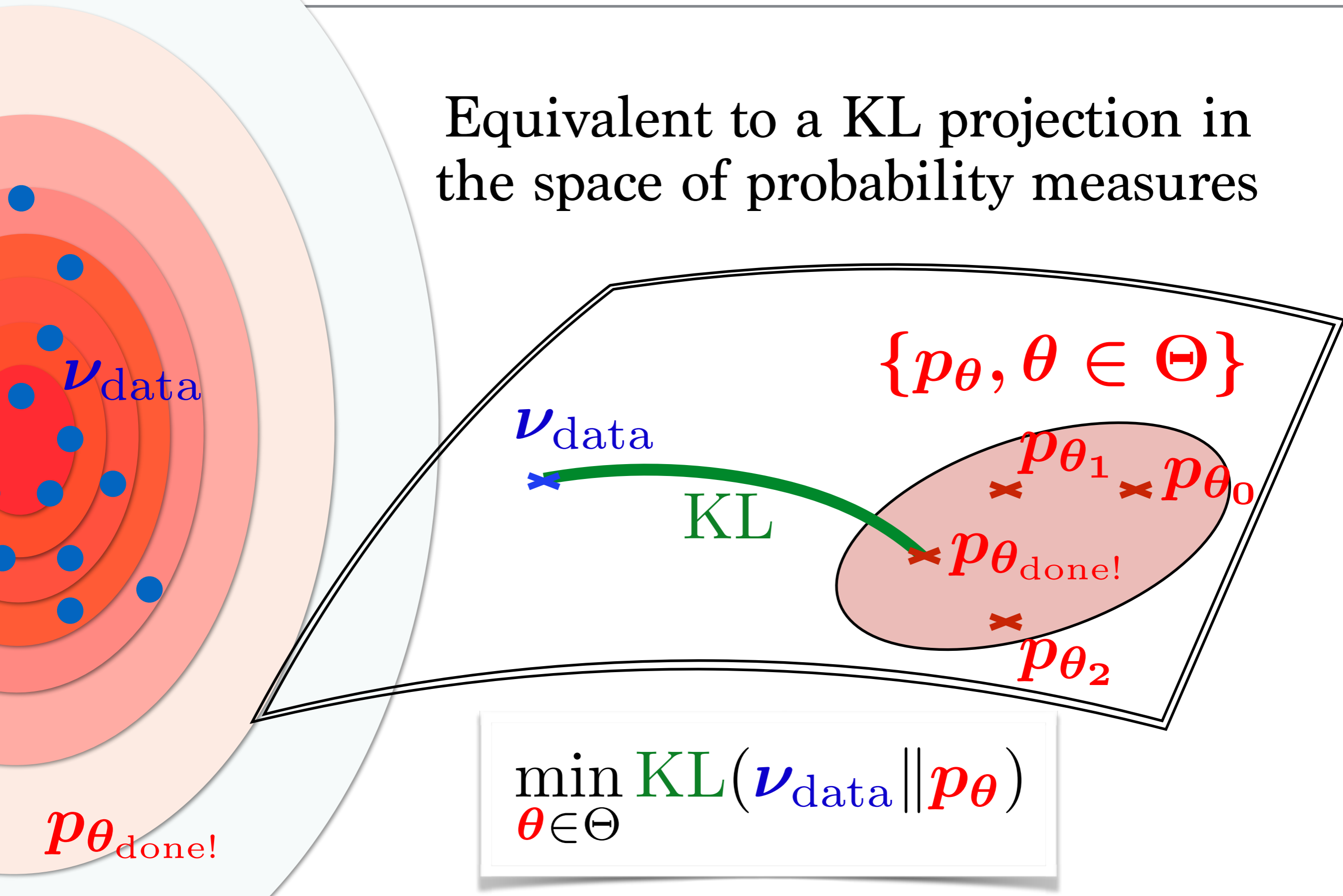
Equivalent to a KL projection in the space of probability measures



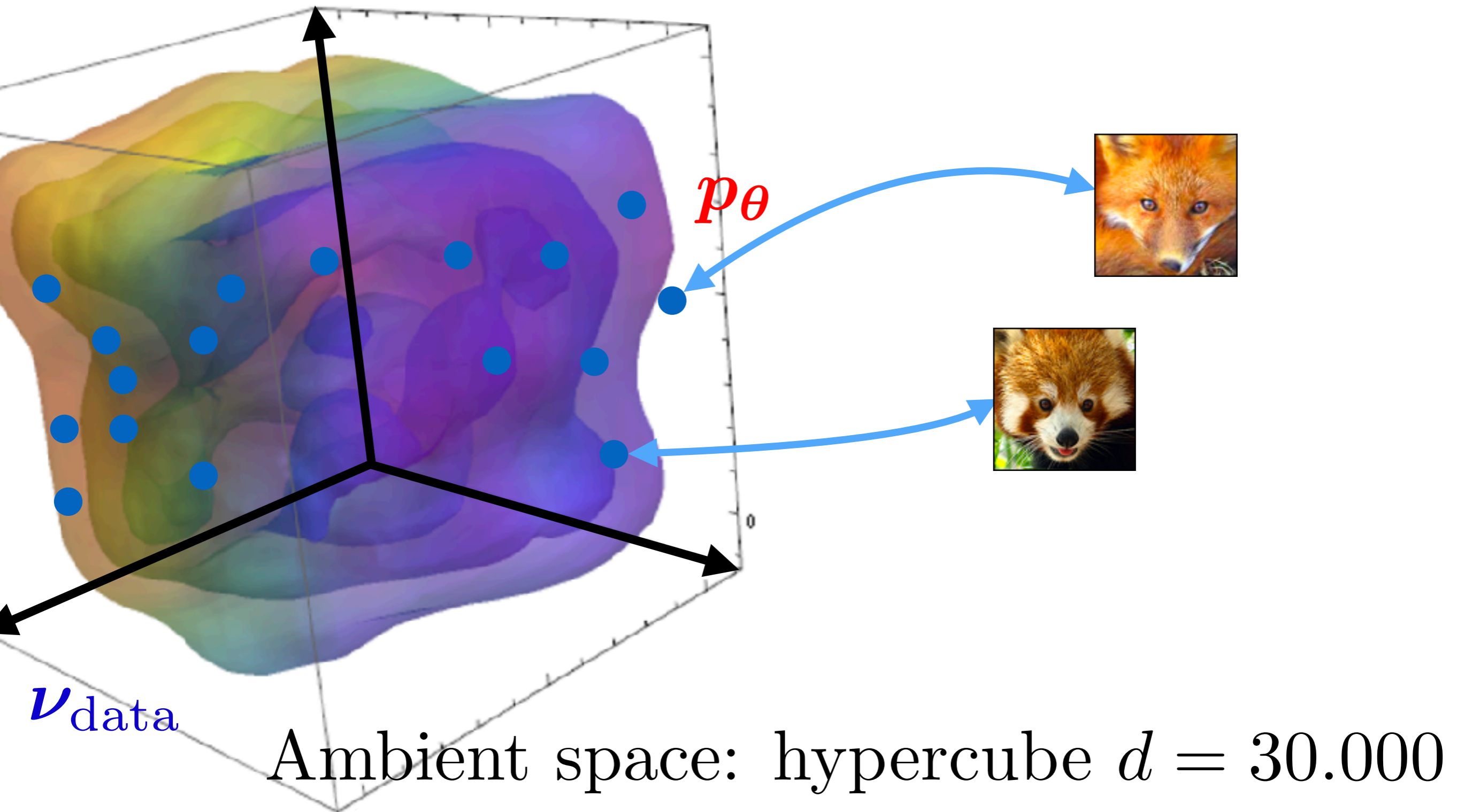
$$\min_{\theta \in \Theta} \text{KL}(\nu_{\text{data}} \parallel p_{\theta})$$

Maximum Likelihood Estimation

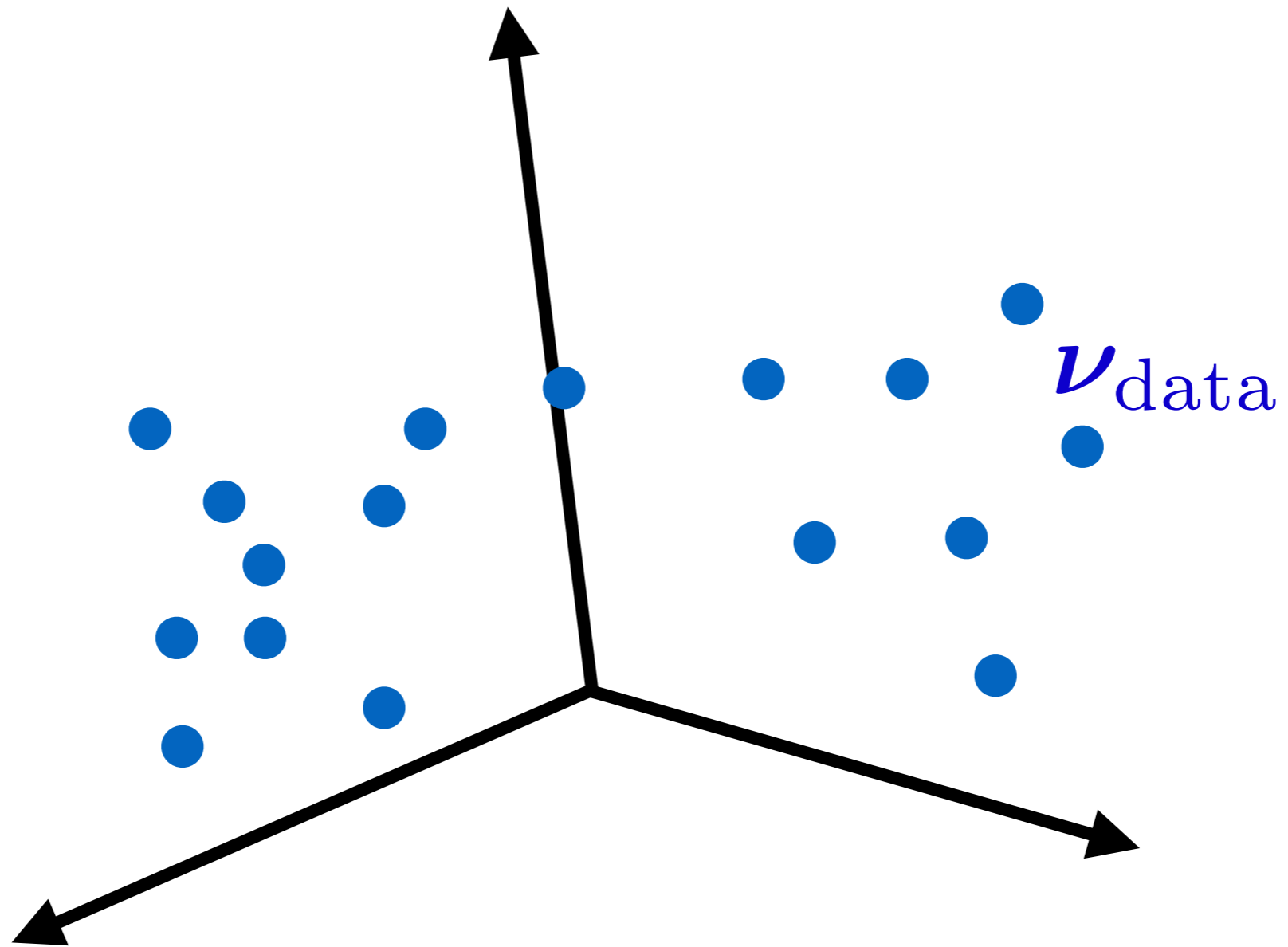
Equivalent to a KL projection in the space of probability measures



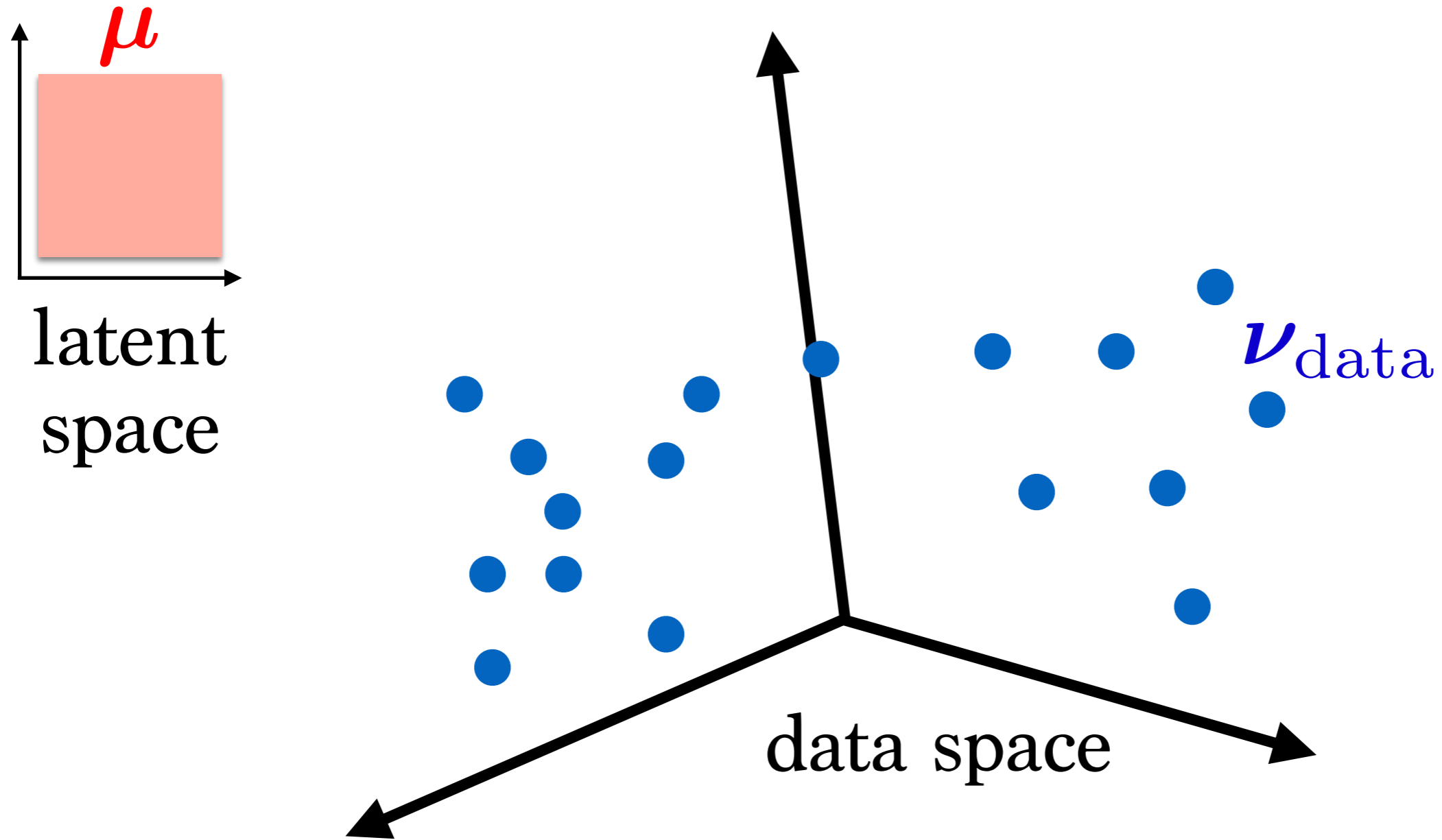
In higher dimensional spaces...



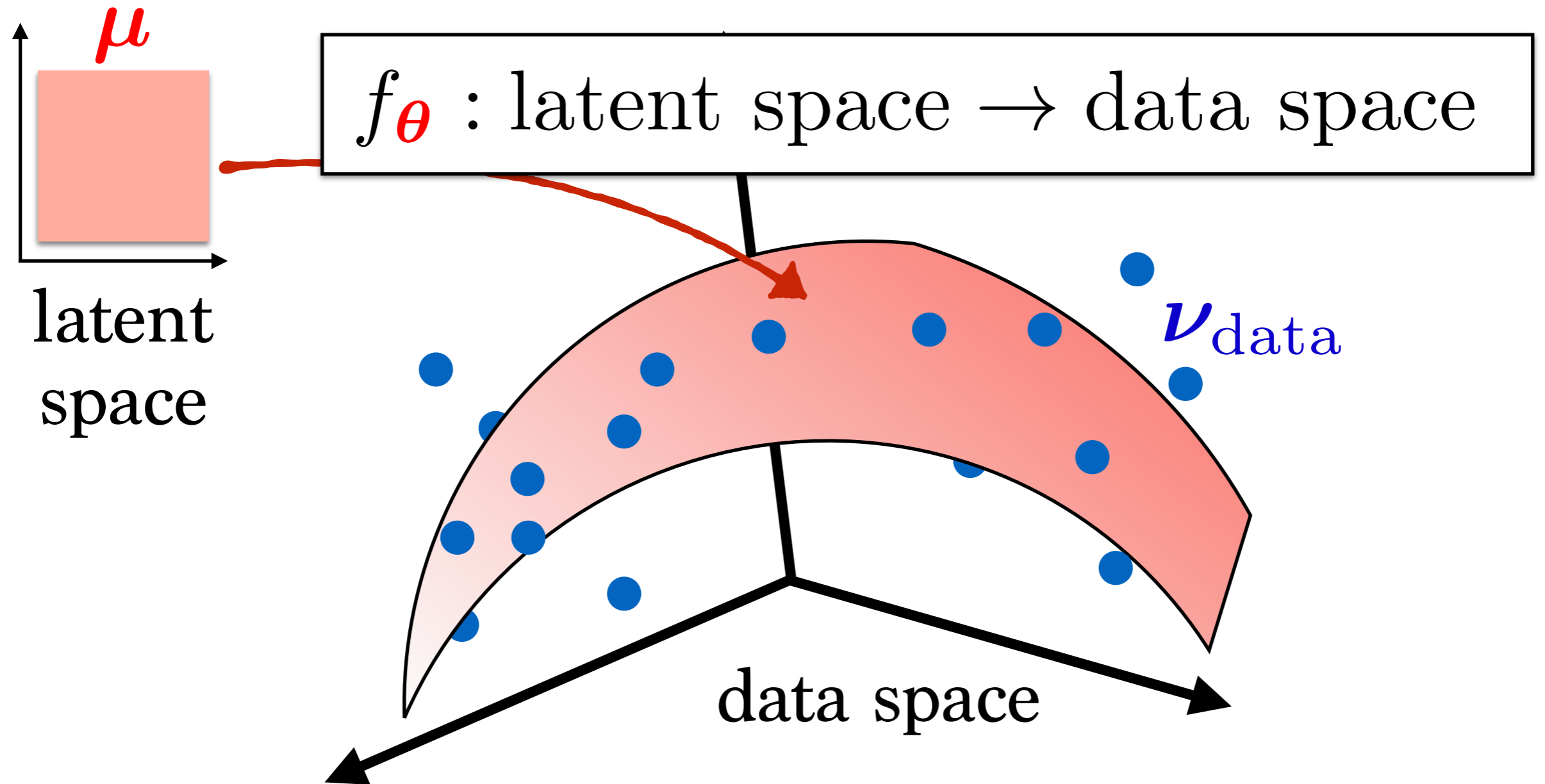
Generative Models



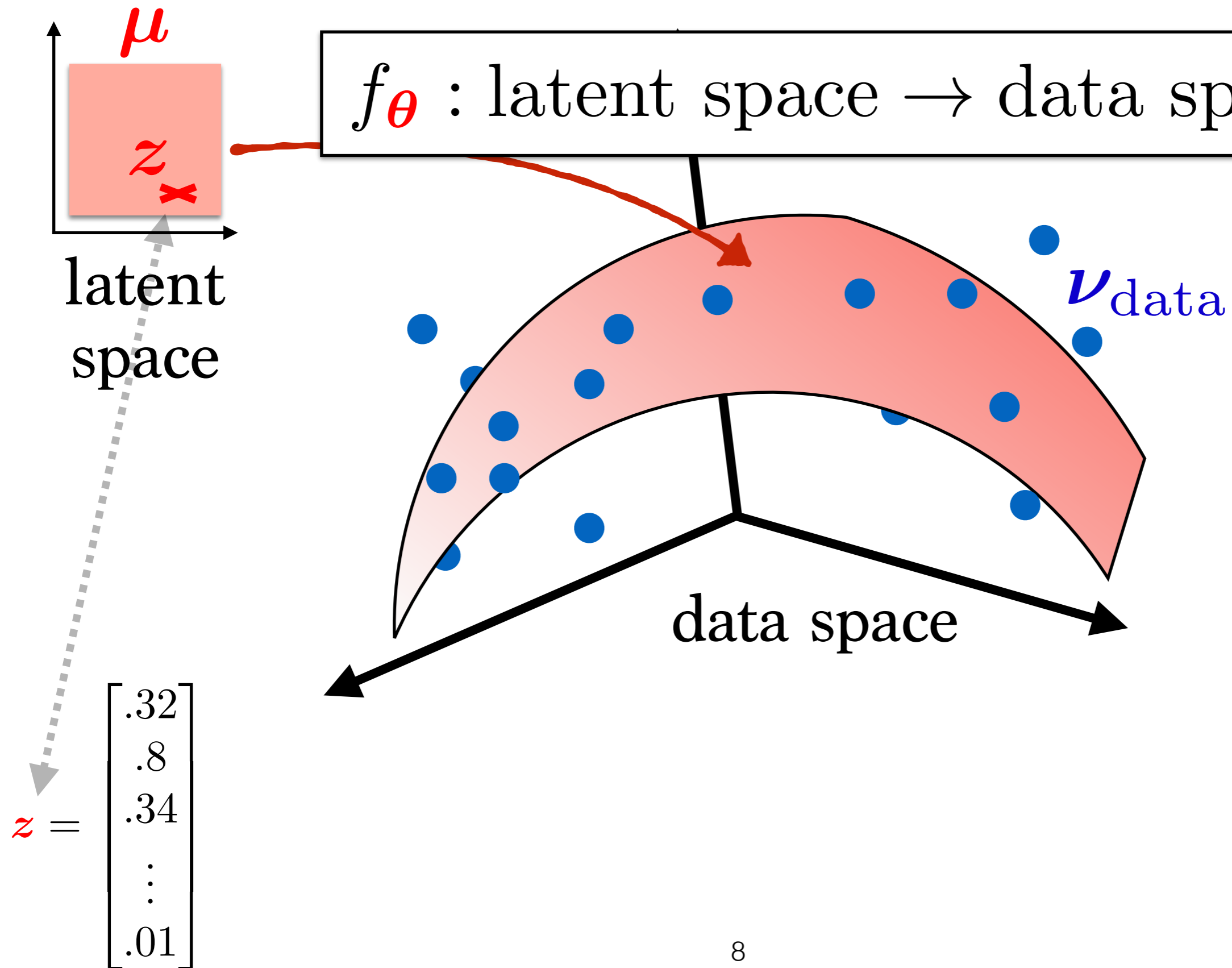
Generative Models



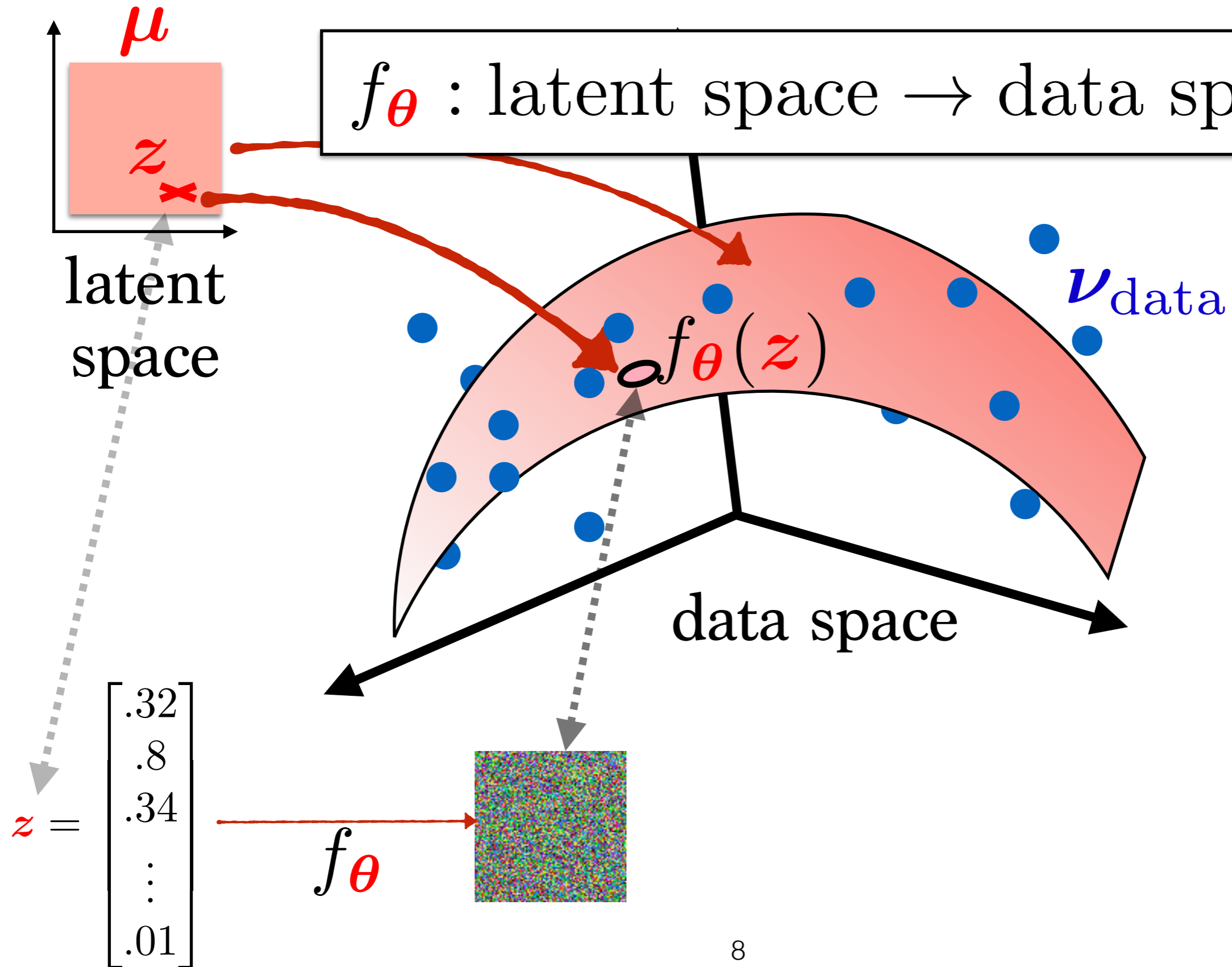
Generative Models



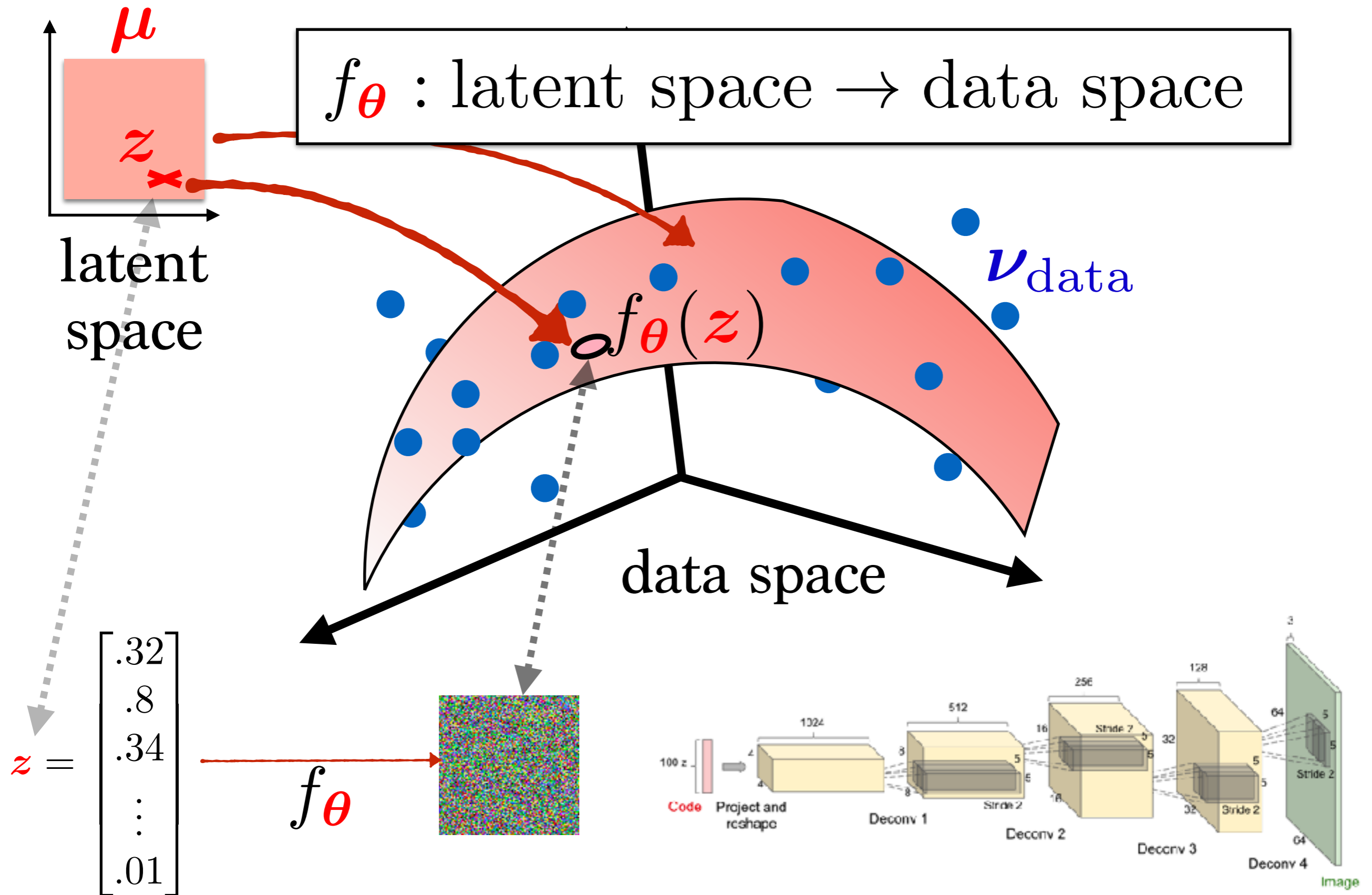
Generative Models



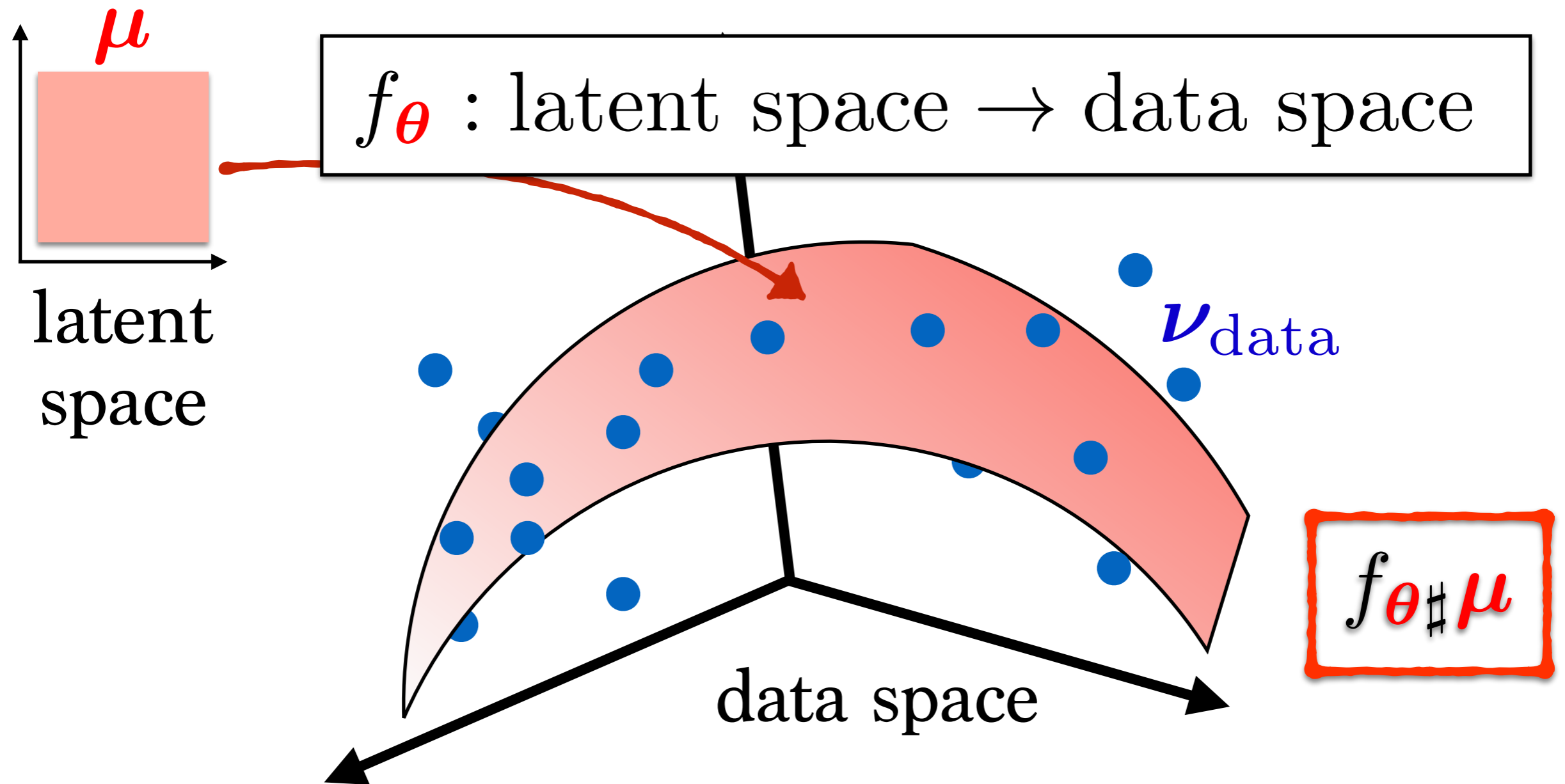
Generative Models



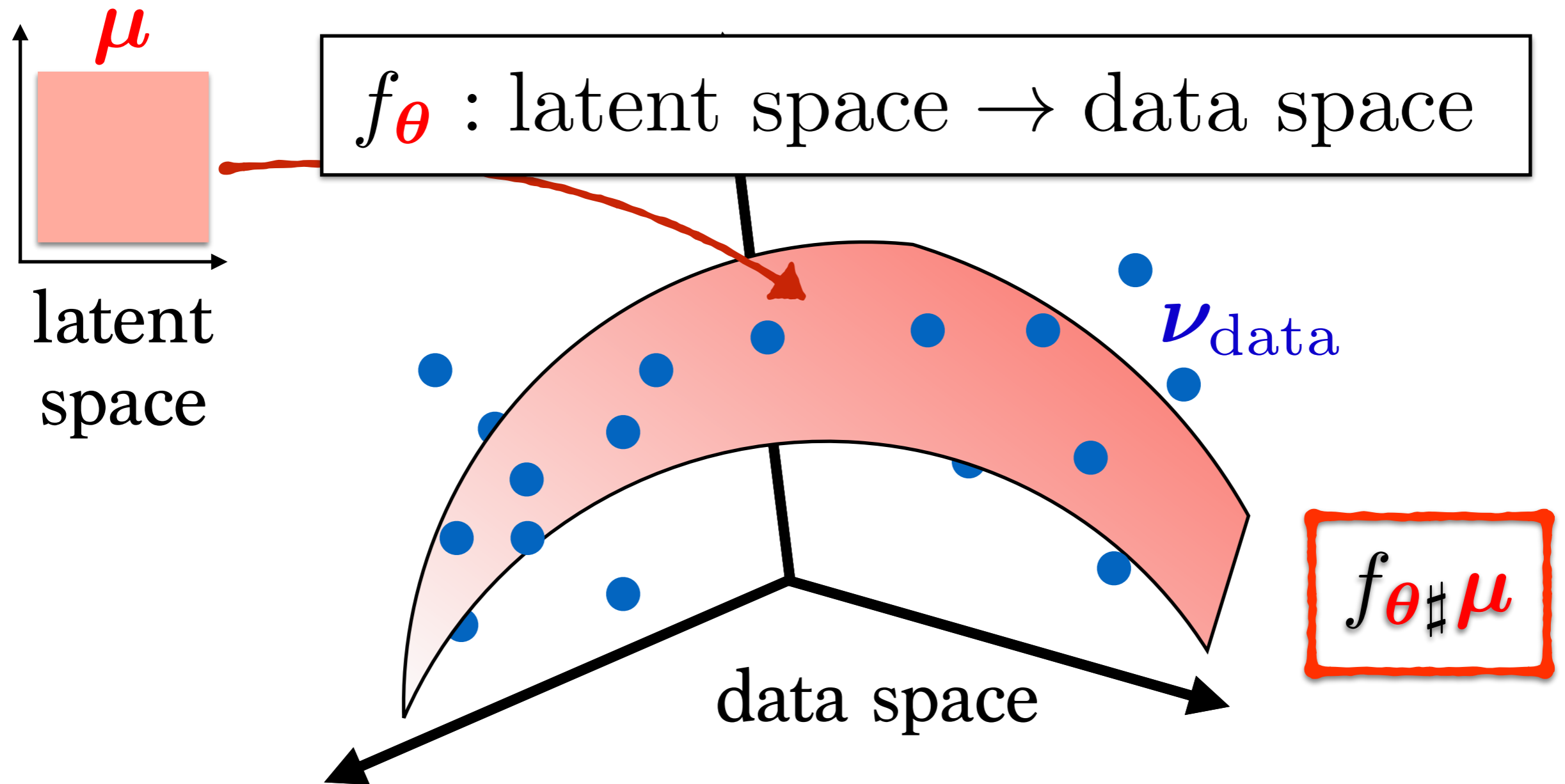
Generative Models



Generative Models

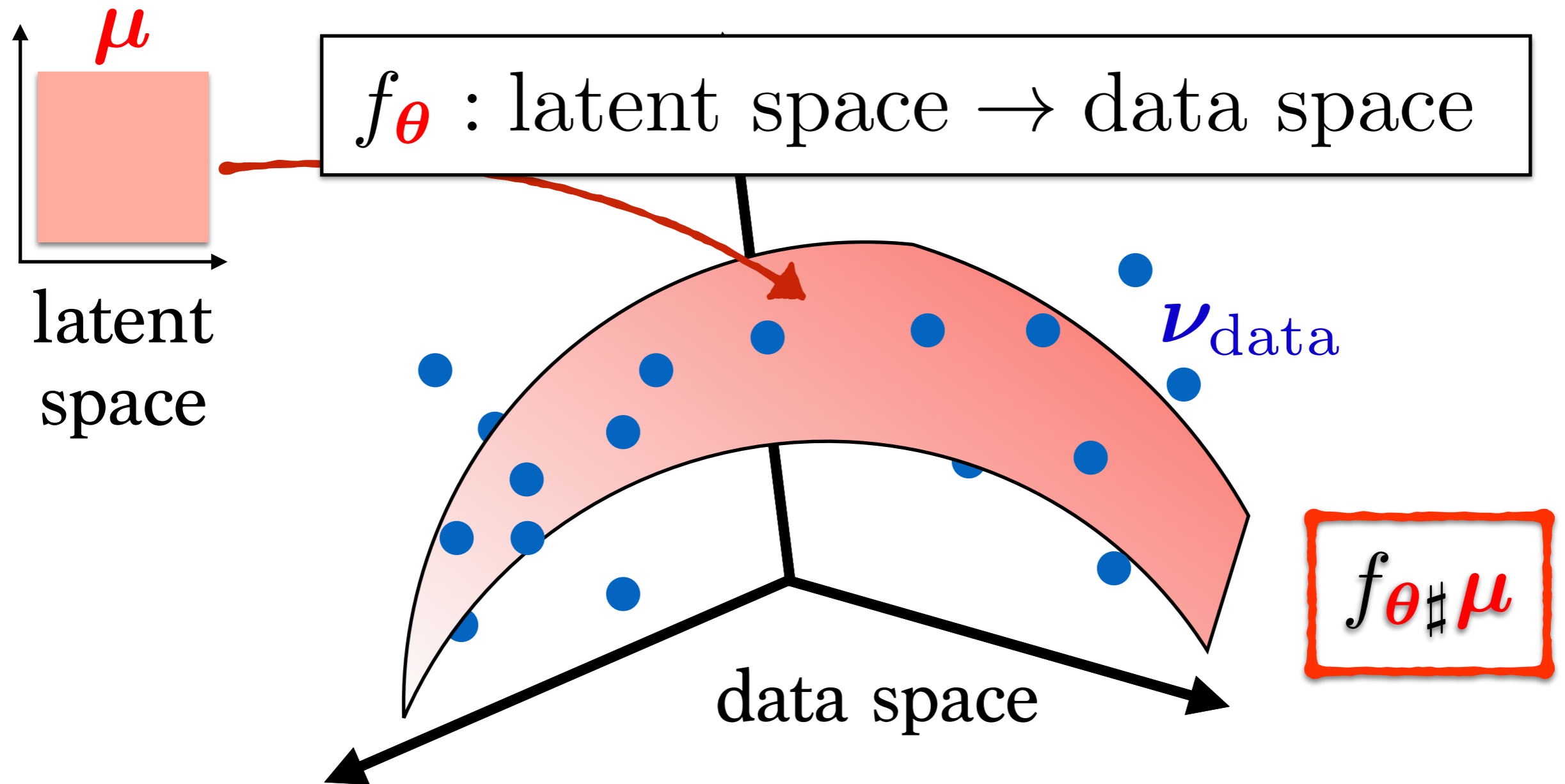


Generative Models



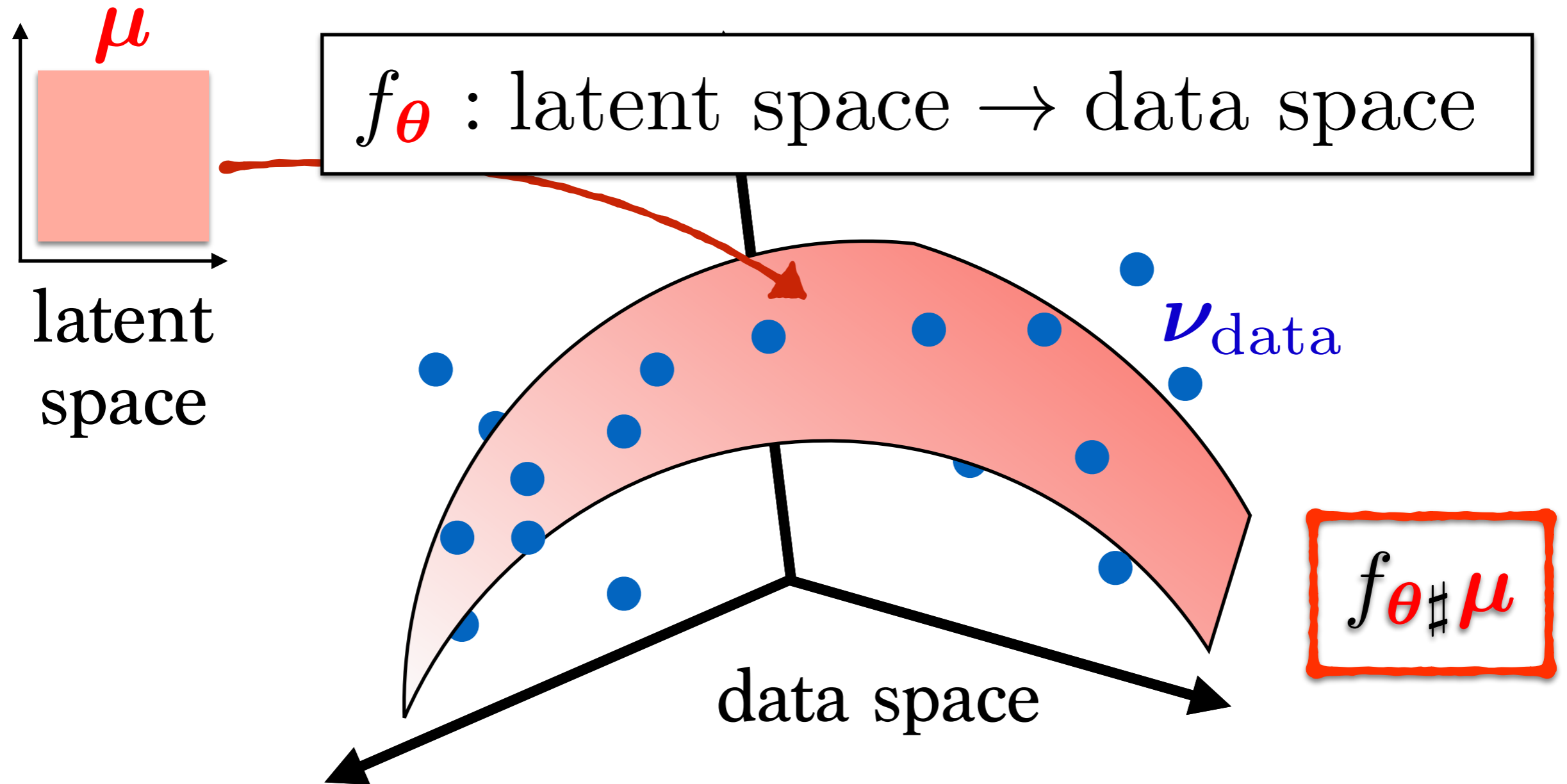
Goal: find θ such that $f_{\theta \# \mu}$ fits $\mathcal{V}_{\text{data}}$

Generative Models



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

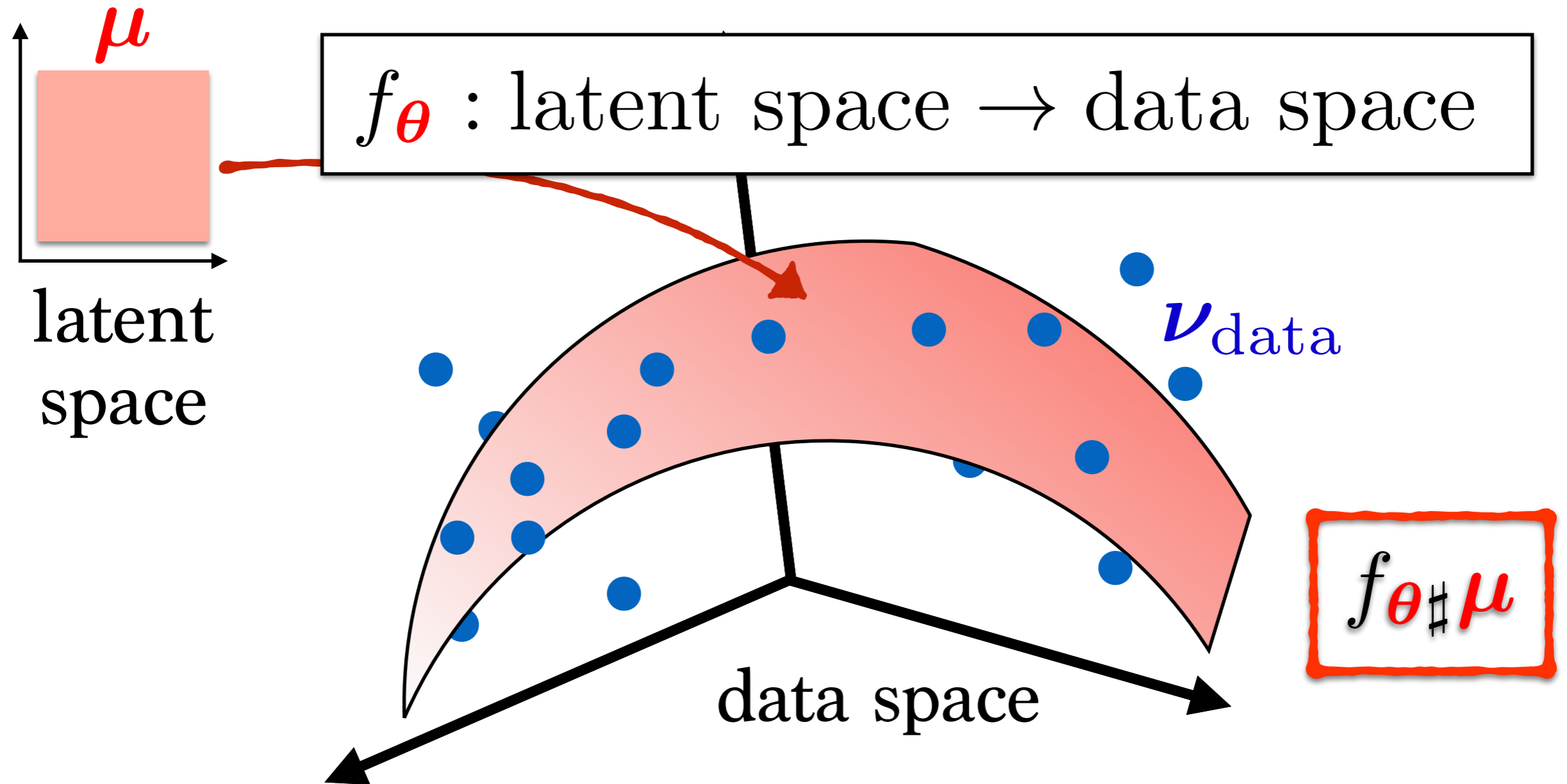
Generative Models



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

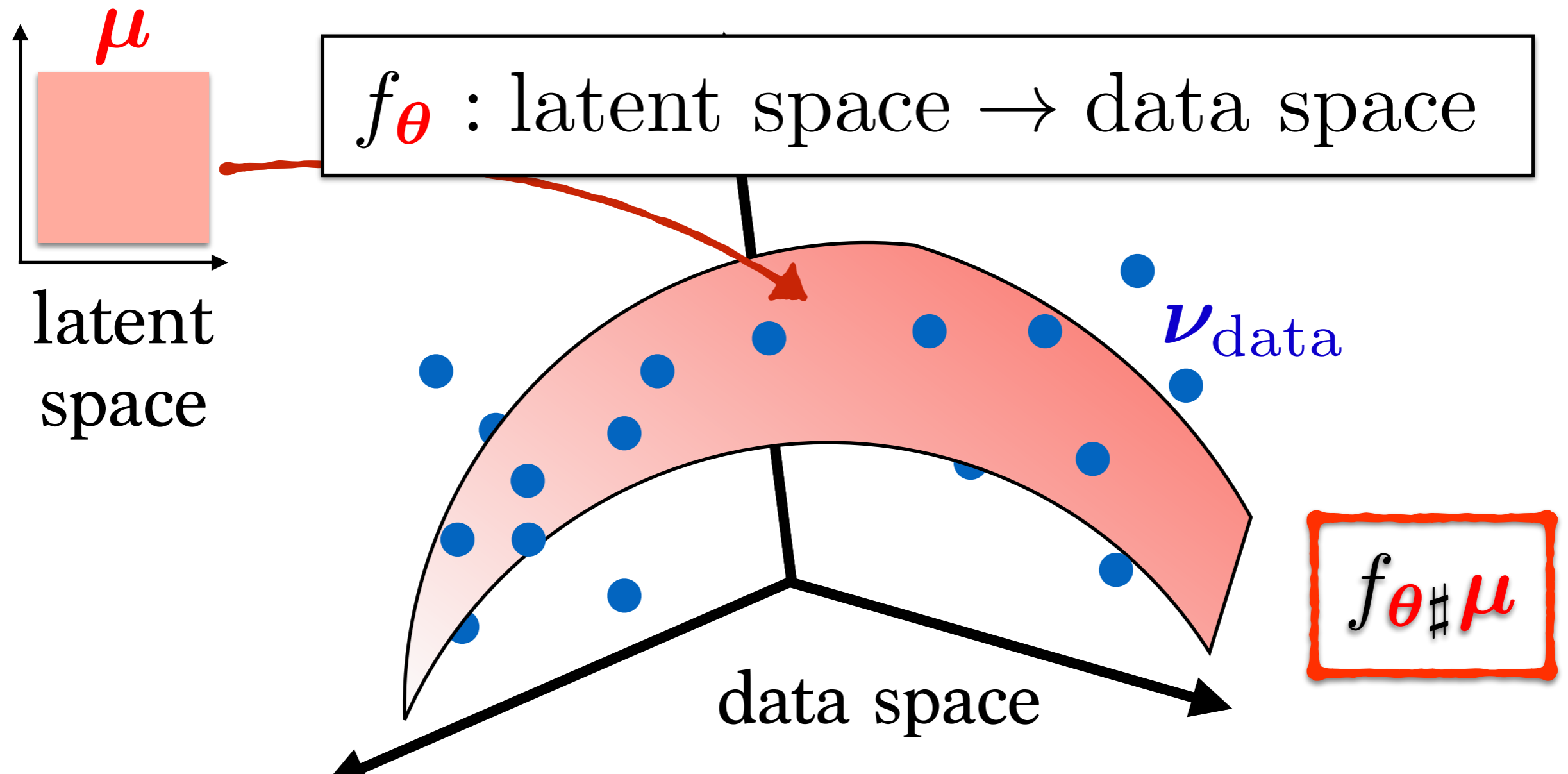


~~MLE~~

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



Generative Models

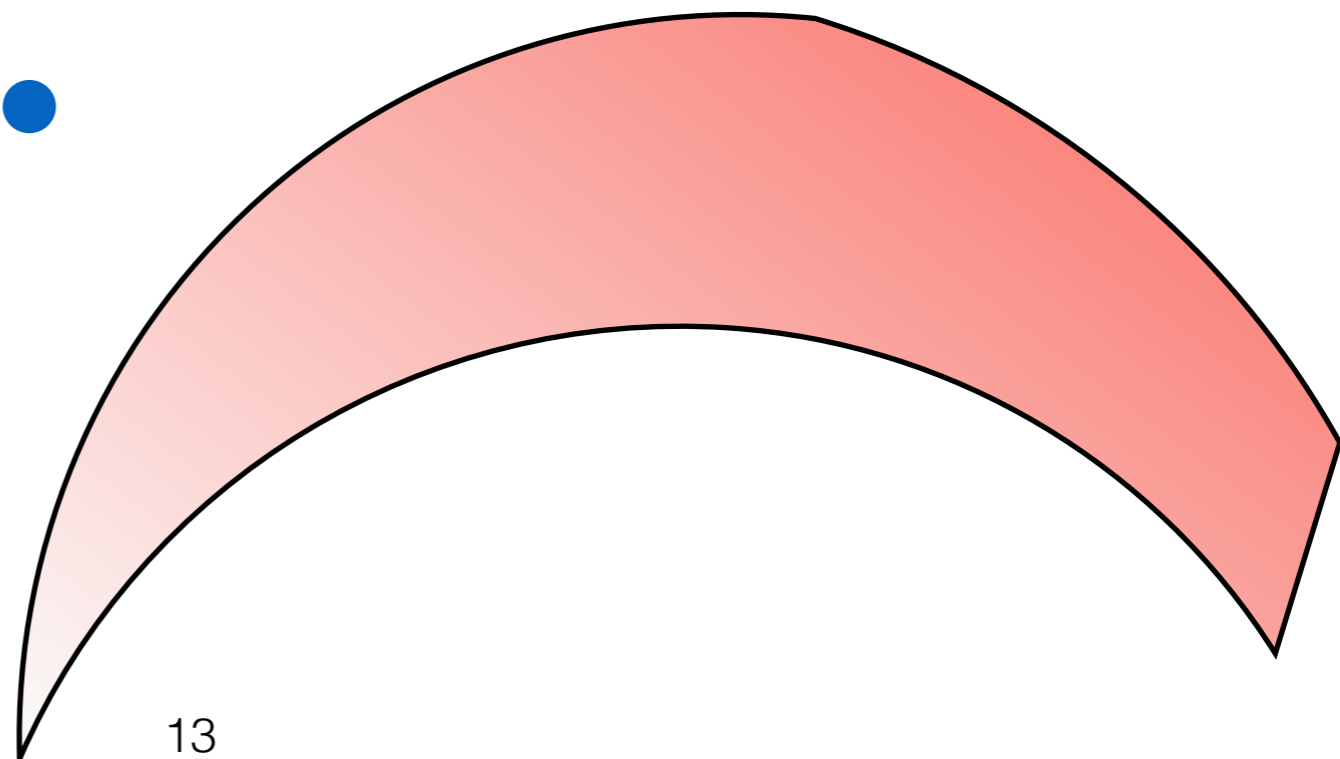
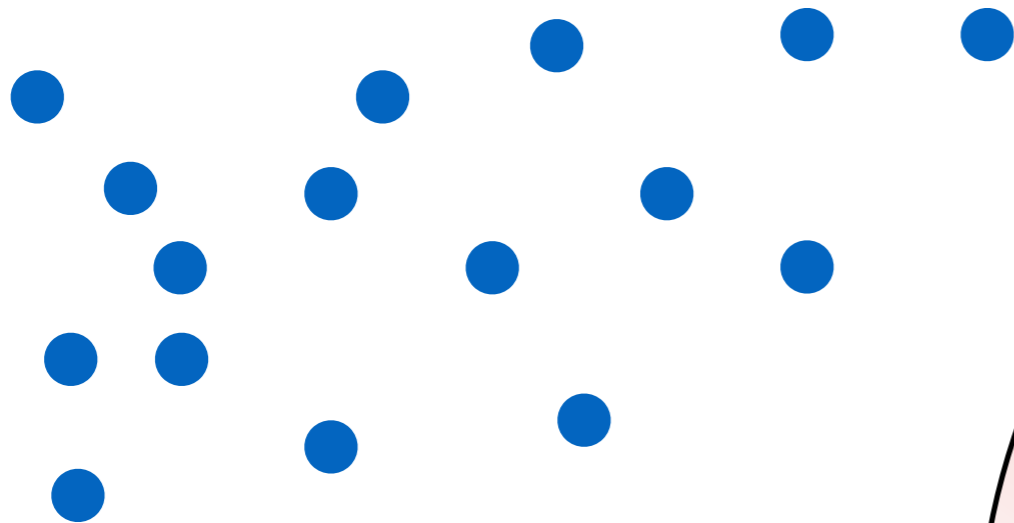


Need a more flexible **discrepancy function** to compare $\mathcal{V}_{\text{data}}$ and $f_{\theta \# \mu}$

Workarounds?

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

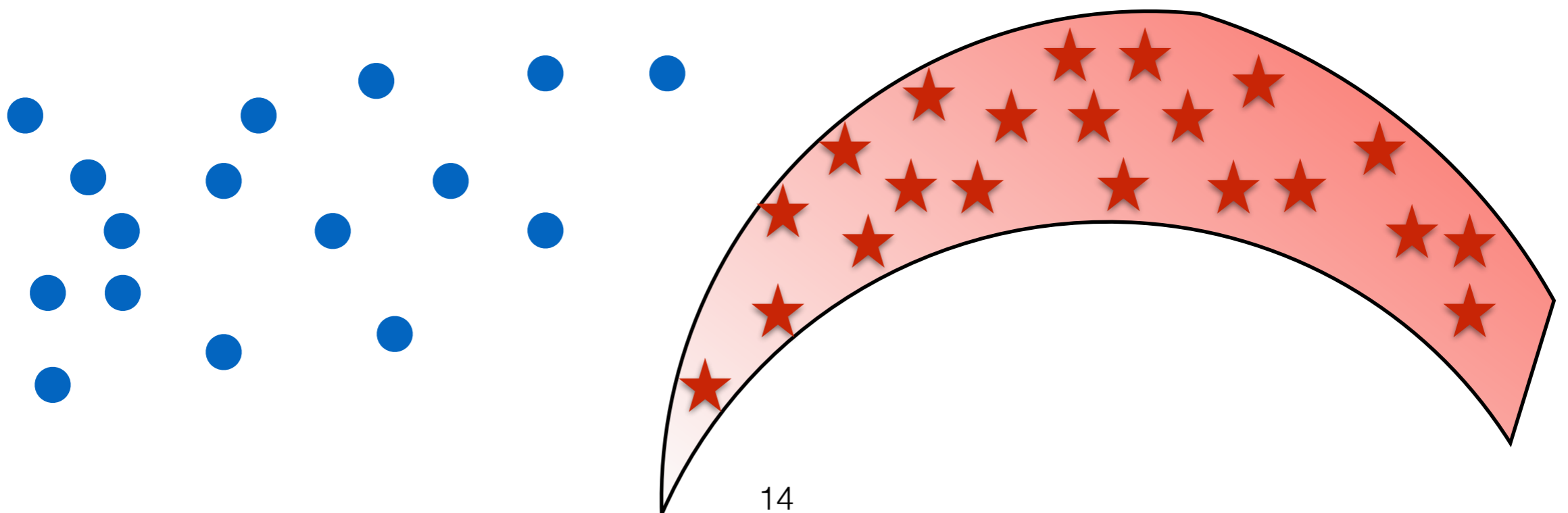
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

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

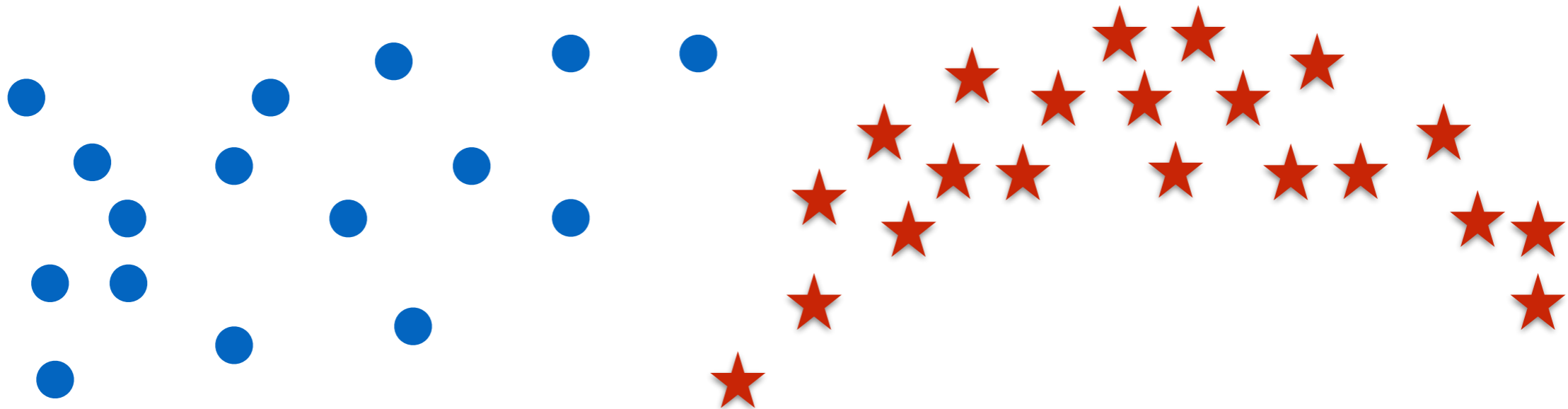
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

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

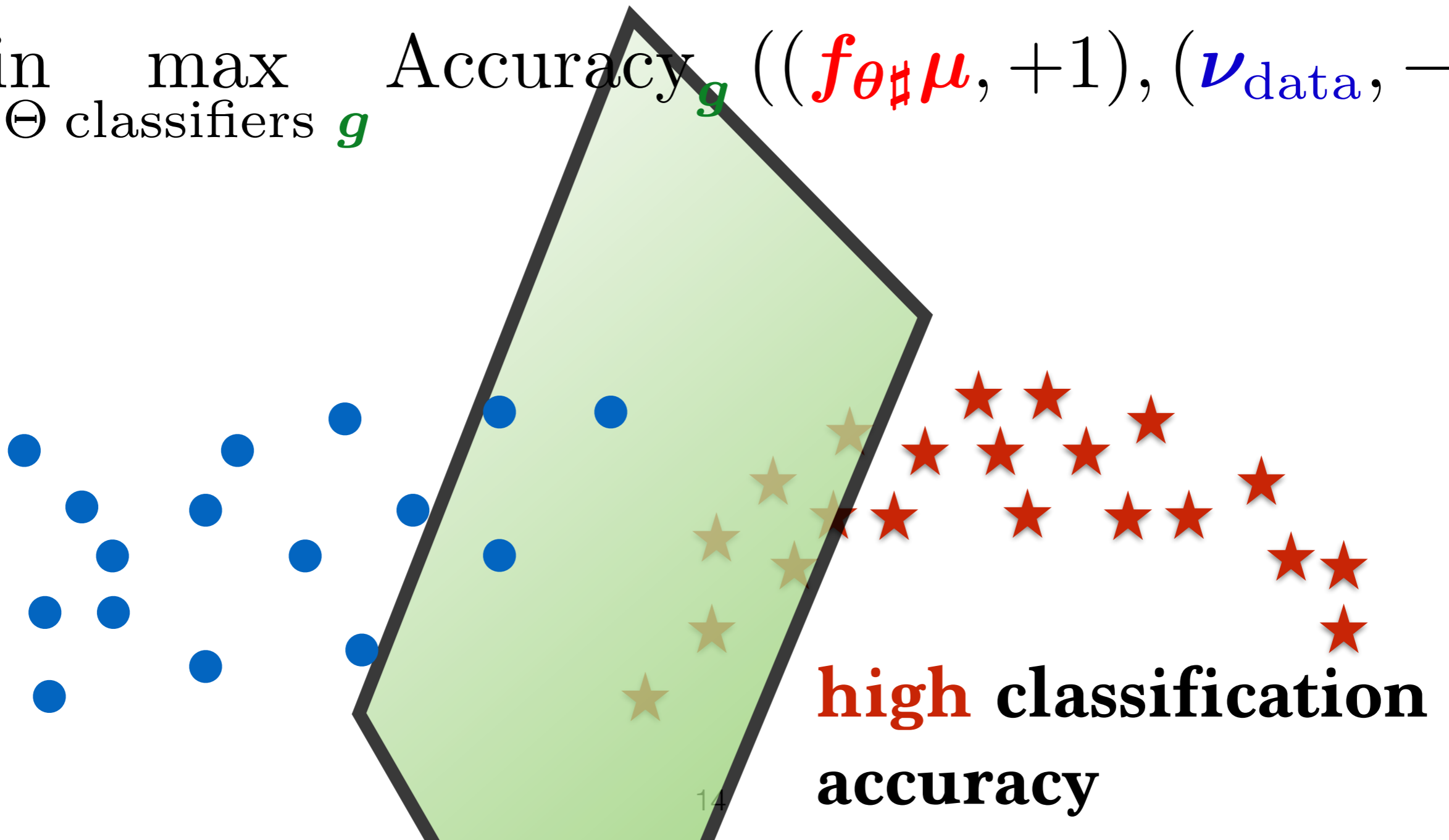
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

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

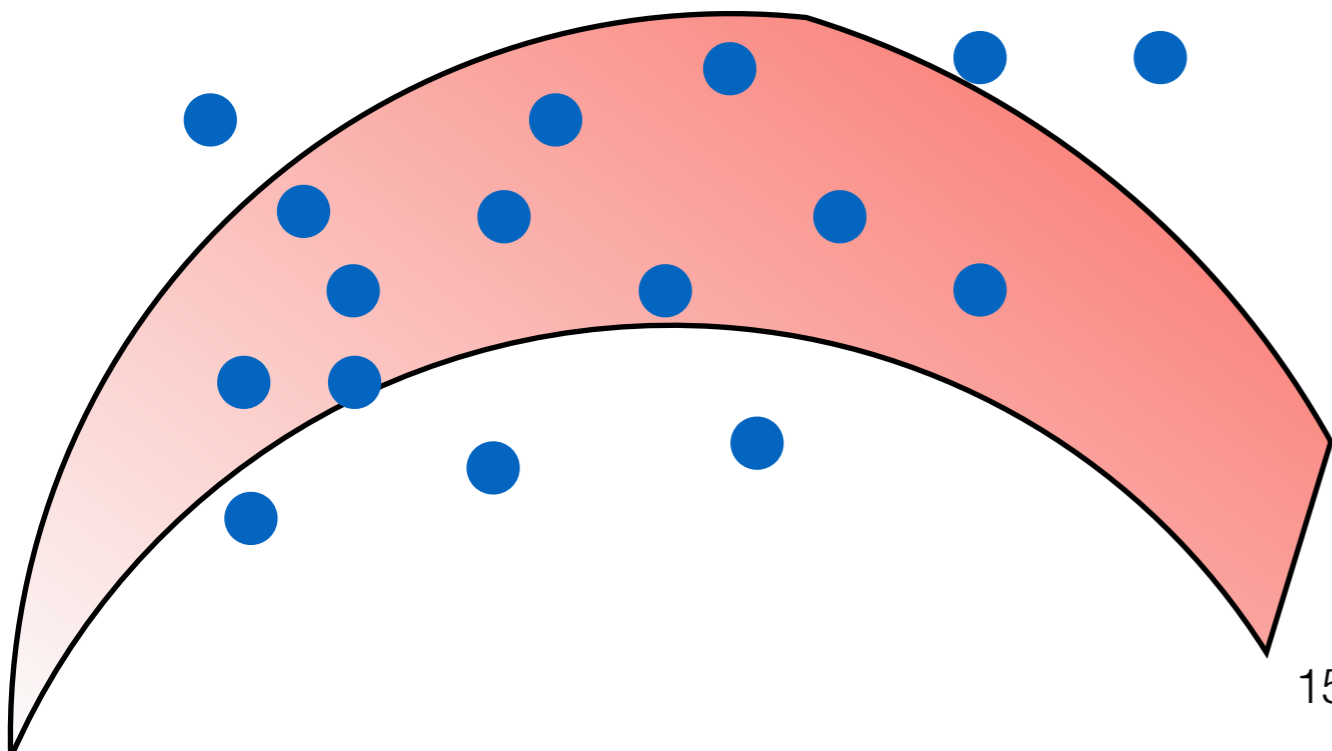
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

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

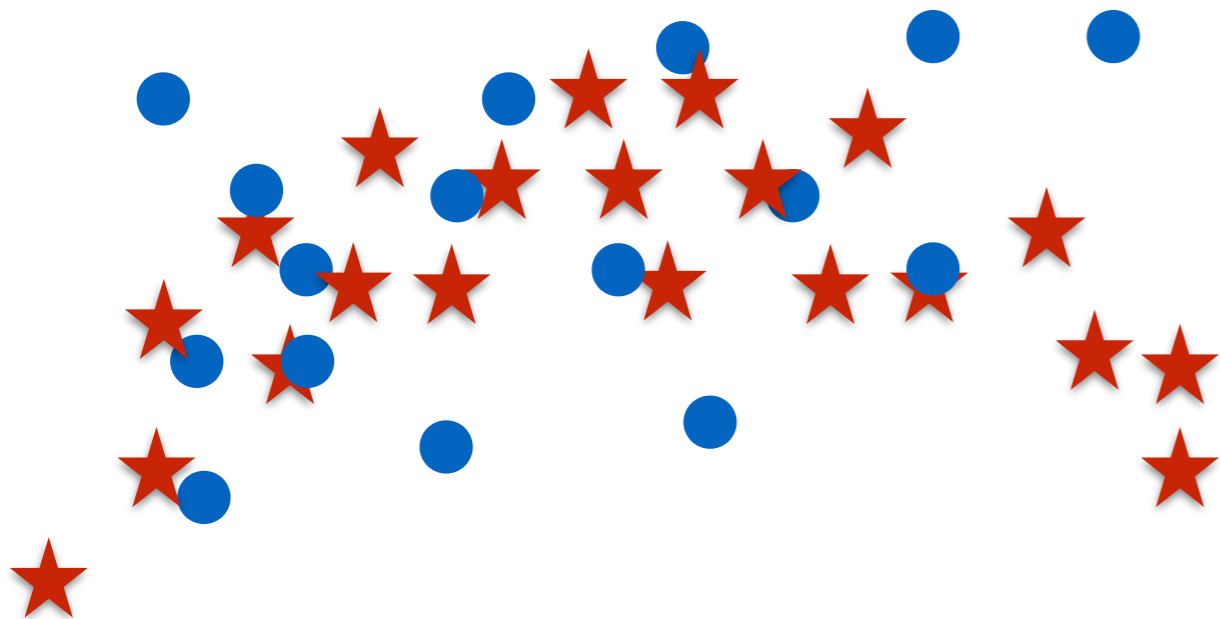
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

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

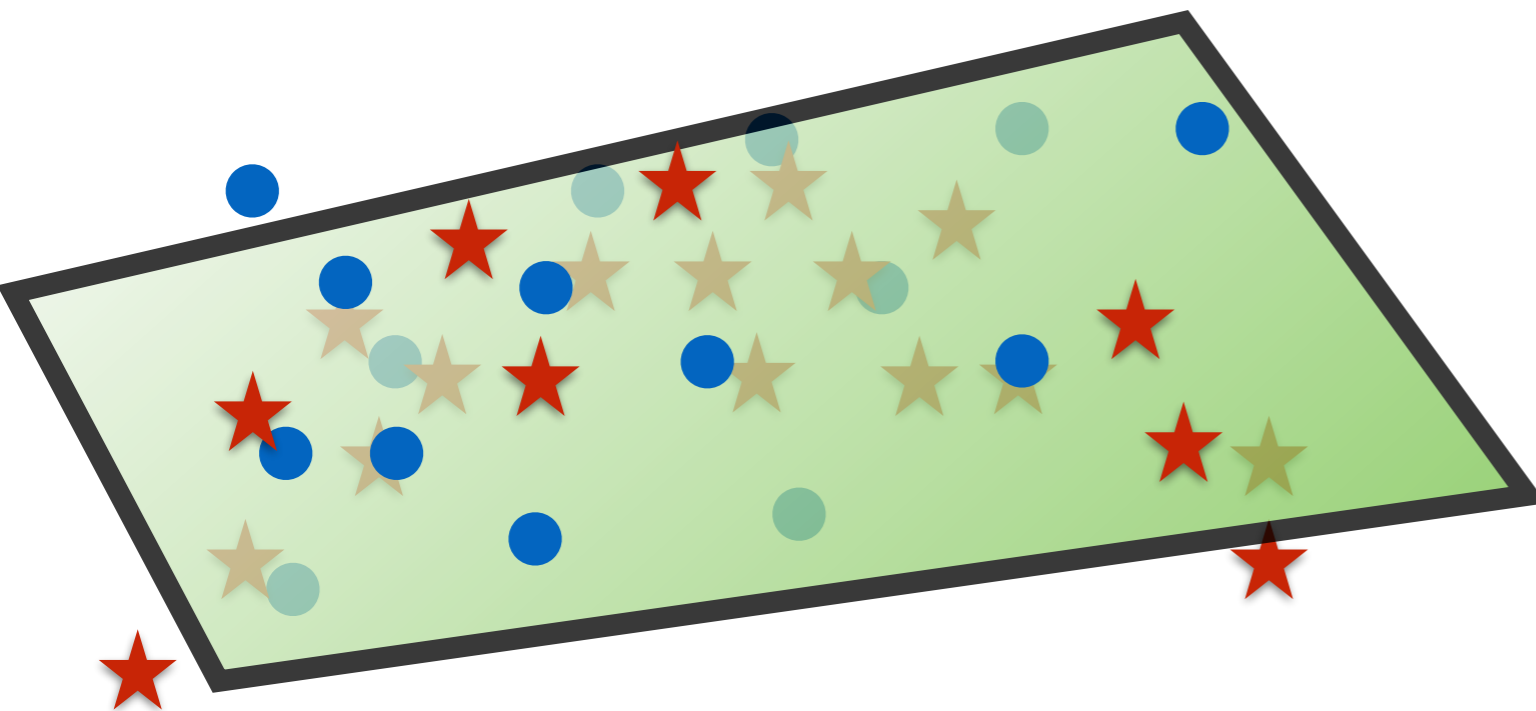
$$\min_{\theta \in \Theta} \max_{\text{classifiers } g} \text{Accuracy}_g \left((f_{\theta} \# \mu, +1), (\nu_{\text{data}}, -1) \right)$$



Workarounds?

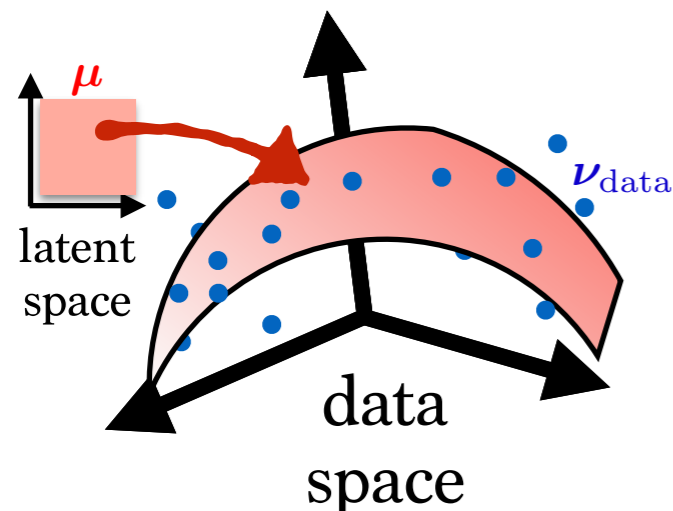
- Formulation as adversarial problem [GPM...'14]

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



**low classification
accuracy...
is the goal.**

Another idea?



- Use a **metric** Δ for probability measures, that can 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

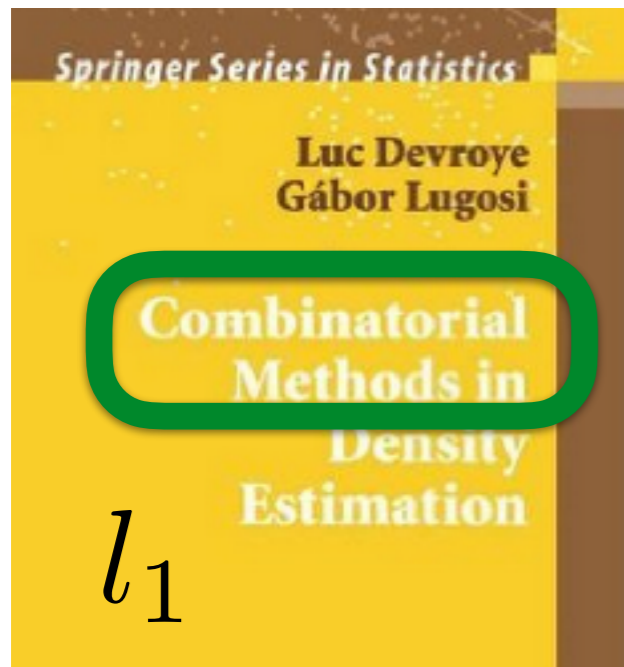
Mayo Clinic, Rochester, Minnesota



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,*}

△ Generative Model Estimation

Generative **Moment Matching** Networks

Yujia Li¹
Kevin Swersky¹
Richard Zemel^{1,2}

YUJIALI@CS.TORONTO.EDU
KSWERSKY@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU

¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA

²Canadian Institute for Advanced Research, Toronto, ON, CANADA

MMD GAN: Towards Deeper Understanding of
Moment Matching Network

Chun-Liang Li^{1,*} Wei-Cheng Chang^{1,*} Yu Cheng² Yiming Yang¹ Barnabás Póczos¹
¹ Carnegie Mellon University, ²IBM Research
{chunlial,wchang2,yiming,bapoczos}@cs.cmu.edu chengyu@us.ibm.com

Training generative neural networks via **Maximum Mean Discrepancy**
optimization

Gintare Karolina Dziugaite
University of Cambridge

Daniel M. Roy
University of Toronto

Zoubin Ghahramani
University of Cambridge

△ Generative Model Estimation

Generative **Moment Matching** Networks

Yujia Li¹
Kevin Swersky¹
Richard Zemel^{1,2}

YUJIALI@CS.TORONTO.EDU
KSWERSKY@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU

¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA

²Canadian Institute for Advanced Research, Toronto, ON, CANADA

MMD GAN: Towards Deeper Understanding of Moment Matching Network

Chun-Liang Li^{1,*} Wei-Cheng Chang^{1,*} Yu Cheng² Yiming Yang¹ Barnabás Póczos¹
¹ Carnegie Mellon University, ²IBM Research
{chunlial,wchang2,yiming,bapoczos}@cs.cmu.edu chengyu@us.ibm.com

Inference in generative models using the **Wasserstein** distance

Espen Bernton, Mathieu Gerber, Pierre E. Jacob, Christian P. Robert

Wasserstein GAN

Martin Arjovsky¹, Soumith Chintala², and Léon Bottou^{1,2}

¹Courant Institute of Mathematical Sciences

²Facebook AI Research

Training generative neural networks via **Maximum Mean Discrepancy** optimization

Gintare Karolina Dziugaite
University of Cambridge

Daniel M. Roy
University of Toronto

Zoubin Ghahramani
University of Cambridge

Wasserstein Training of Restricted Boltzmann Machines

Grégoire Montavon
Technische Universität Berlin
gregoire.montavon@tu-berlin.de

Klaus-Robert Müller*
Technische Universität Berlin
klaus-robert.mueller@tu-berlin.de

Marco Cuturi
CREST, ENSAE, Université Paris-Saclay
marco.cuturi@ensae.fr

△ Generative Model Estimation

Generative **Moment Matching** Networks

Yujia Li¹
Kevin Swersky¹
Richard Zemel^{1,2}

¹Department of Computer Science, University of Toronto, Toronto, ON, CANADA

²Canadian Institute for Advanced Research, Toronto, ON, CANADA

YUJIALI@CS.TORONTO.EDU
KSWERSKY@CS.TORONTO.EDU
ZEMEL@CS.TORONTO.EDU

MMD GAN: Towards Deeper Understanding of Moment Matching Network

Chun-Liang Li^{1,*} Wei-Cheng Chang^{1,*} Yu Cheng² Yiming Yang¹ Barnabás Póczos¹
¹ Carnegie Mellon University, ²IBM Research
{chunli1,wchang2,yiming,bapoczos}@cs.cmu.edu chengyu@us.ibm.com

Inference in generative models using the **Wasserstein** distance

Espen Bernton, Mathieu Gerber, Pierre E. Jacob, Christian P. Robert

Wasserstein GAN

Martin Arjovsky¹, Soumith Chintala², and Léon Bottou^{1,2}

¹Courant Institute of Mathematical Sciences

²Facebook AI Research

Learning Generative Models with **Sinkhorn** Divergences

Aude Genevay
CEREMADE,
Université Paris-Dauphine

Gabriel Peyré
CNRS and DMA,
École Normale Supérieure

Marco Cuturi
ENSAE CREST
Université Paris-Saclay

Tim Salimans*
OpenAI
tim@openai.com

Han Zhang*
Rutgers University
han.zhang@cs.rutgers.edu

Alec Radford
OpenAI
alec@openai.com

Dimitris Metaxas
Rutgers University
dnm@cs.rutgers.edu

Training generative neural networks via **Maximum Mean Discrepancy** optimization

Gintare Karolina Dziugaite
University of Cambridge

Daniel M. Roy
University of Toronto

Zoubin Ghahramani
University of Cambridge

Wasserstein Training of Restricted Boltzmann Machines

Grégoire Montavon
Technische Universität Berlin
gregoire.montavon@tu-berlin.de

Klaus-Robert Müller*
Technische Universität Berlin
klaus-robert.mueller@tu-berlin.de

Marco Cuturi
CREST, ENSAE, Université Paris-Saclay
marco.cuturi@ensae.fr

Improving GANs Using **Optimal Transport**

Minimum Kantorovich Estimation

- Use optimal transport theory, namely *Wasserstein distances* to define discrepancy Δ .

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

- Optimal transport? fertile field in mathematics.



Monge



Kantorovich



Koopmans



Dantzig



Brenier



Gangbo



Otto



McCann



Villani



Figalli

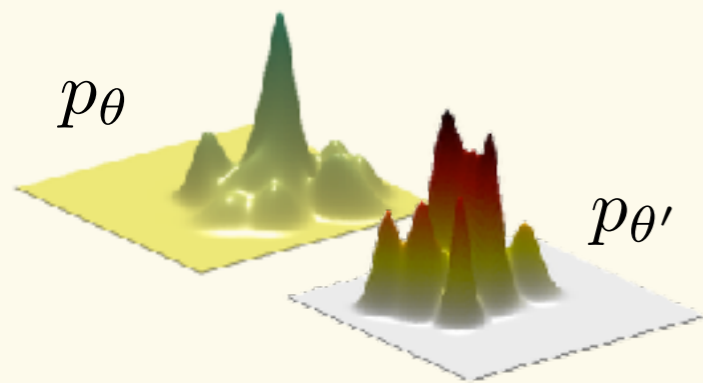
Nobel'75

Fields'10

Fields'18

What is Optimal Transport?

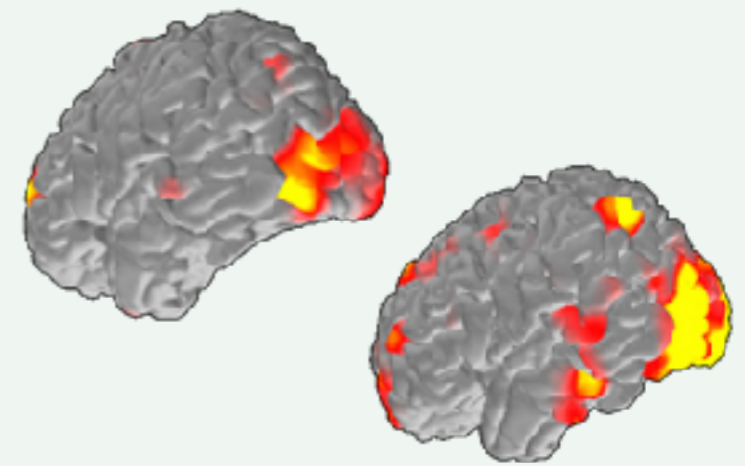
The natural geometry for **probability measures**



Statistical Models

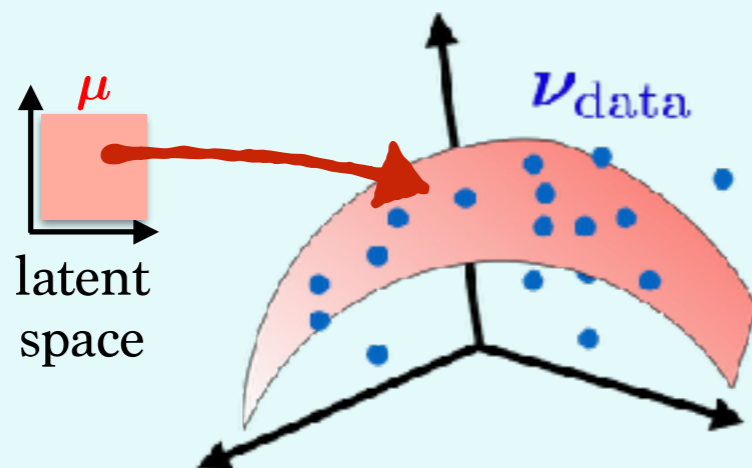


Bags of features



Brain Activation Maps

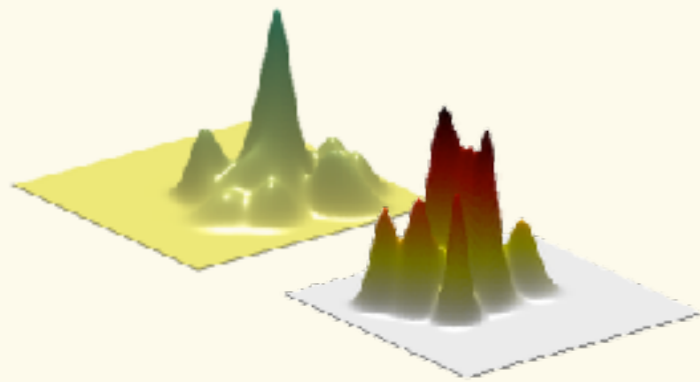
Generative Models vs. data



Color Histograms

What is Optimal Transport?

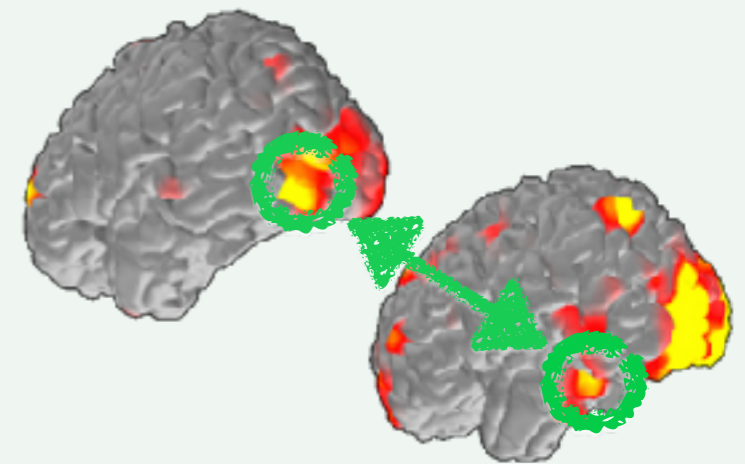
The natural geometry for **probability measures** supported on a metric space.



Statistical Models

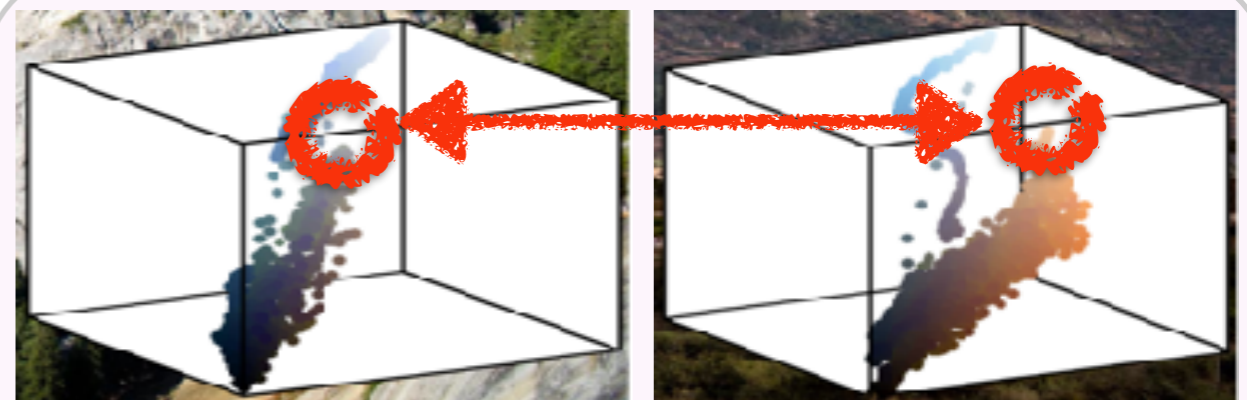
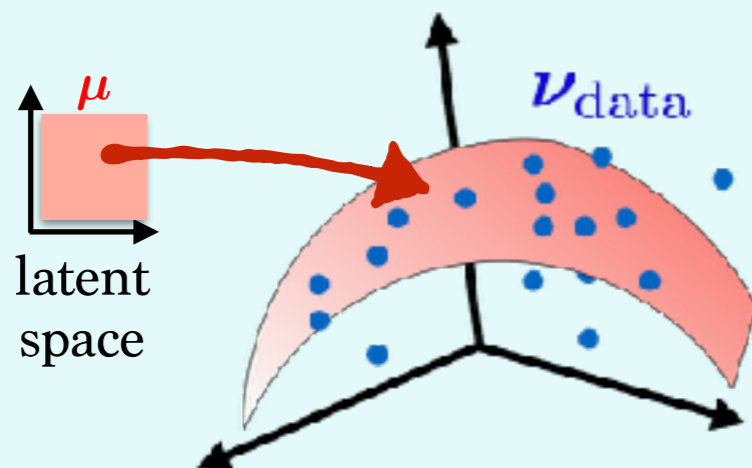


Bags of Features



Brain Activation Maps

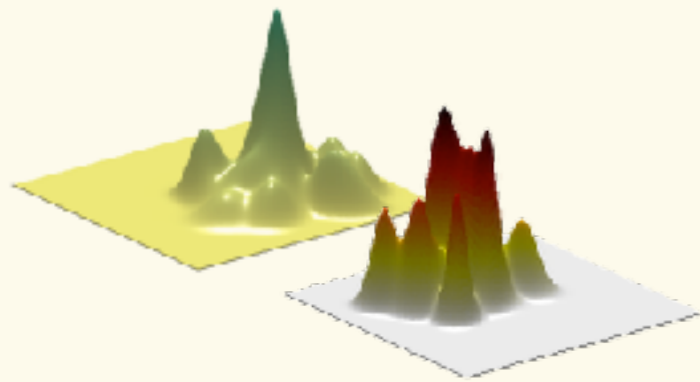
Generative Models vs. Data



Color Histograms

What is Optimal Transport?

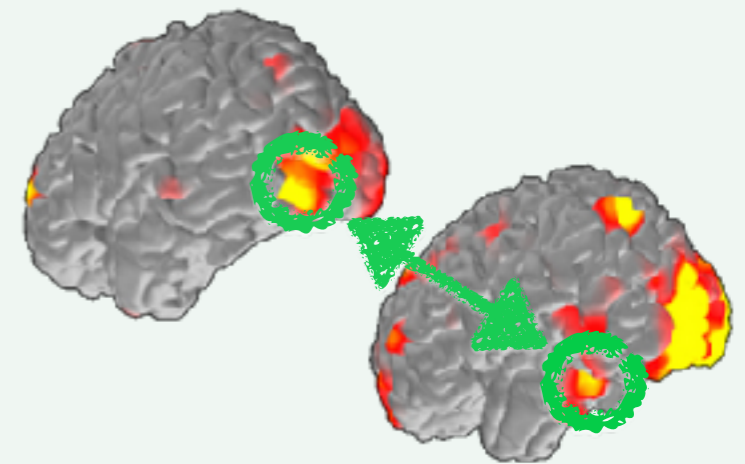
The natural geometry for **probability measures** supported on a metric space.



Statistical Models

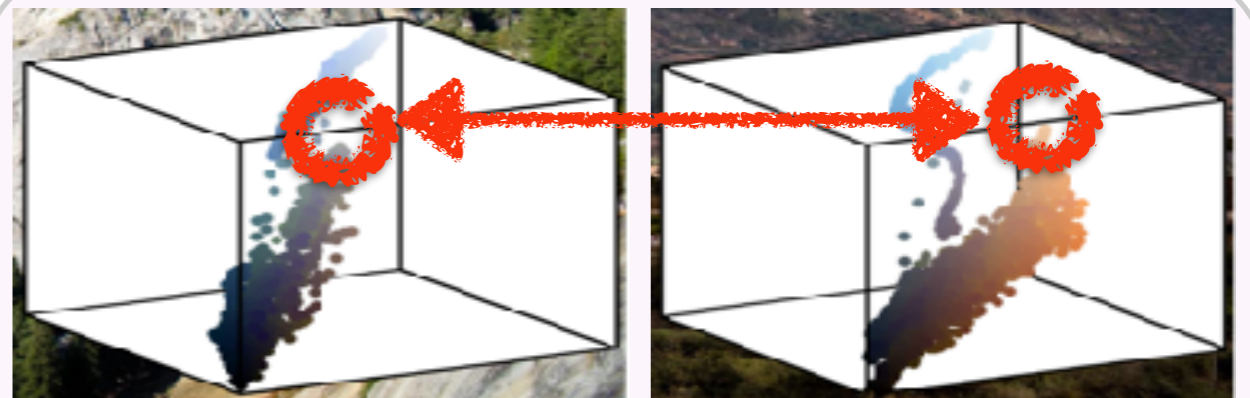
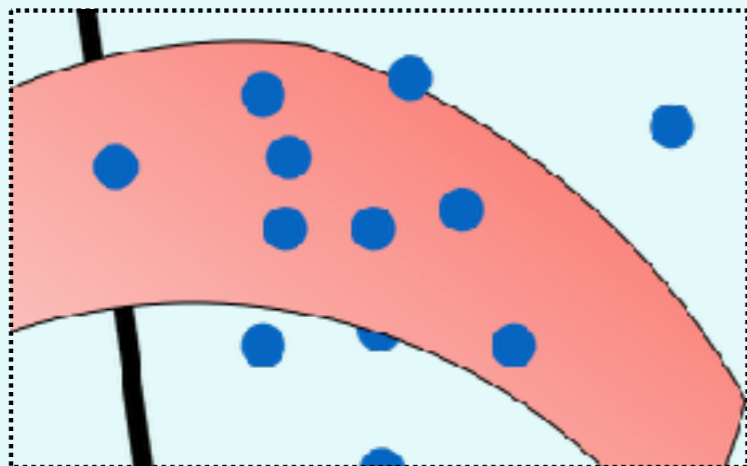


*Bags
of Features*



Brain Activation Maps

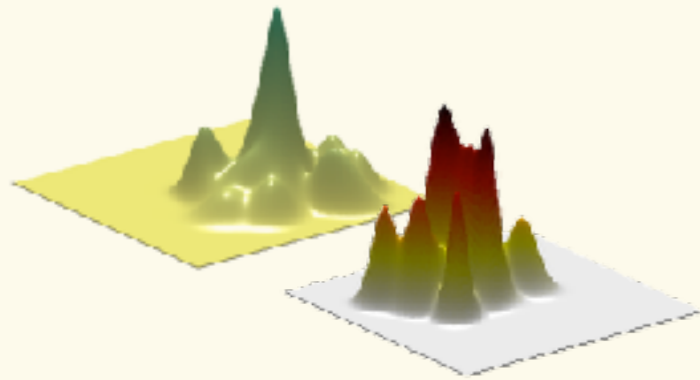
*Generative
Models
vs. Data*



Color Histograms

What is Optimal Transport?

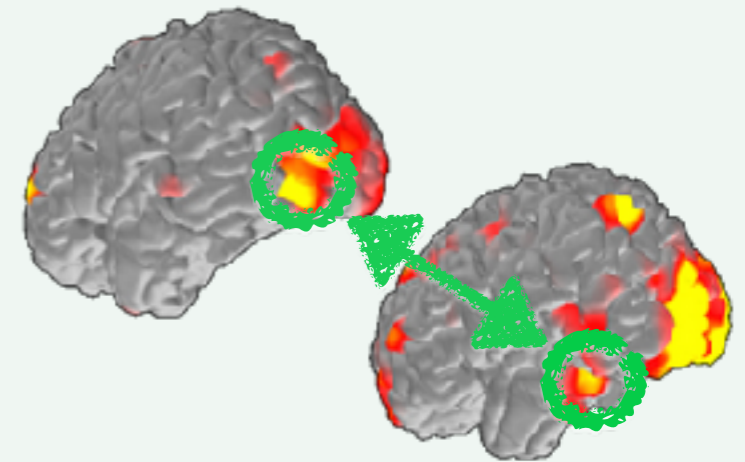
The natural geometry for **probability measures** supported on a metric space.



Statistical Models

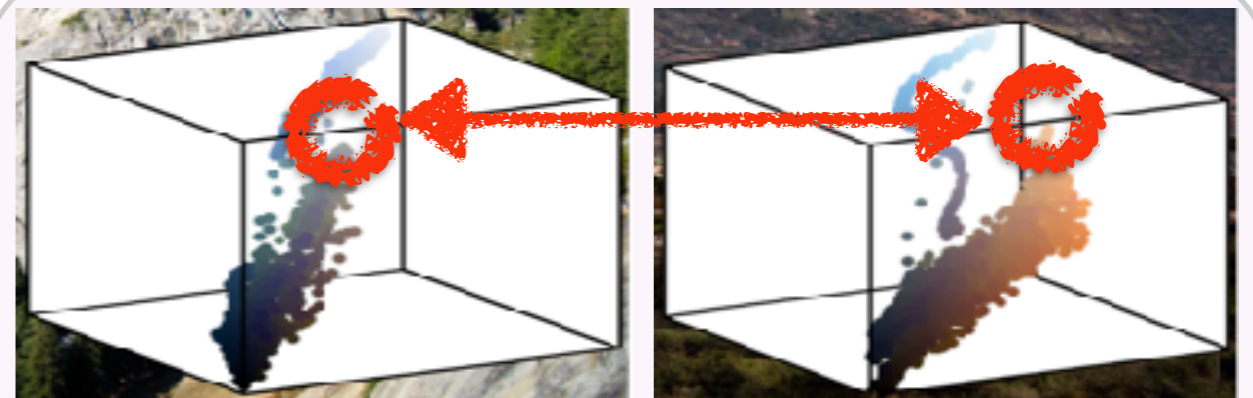
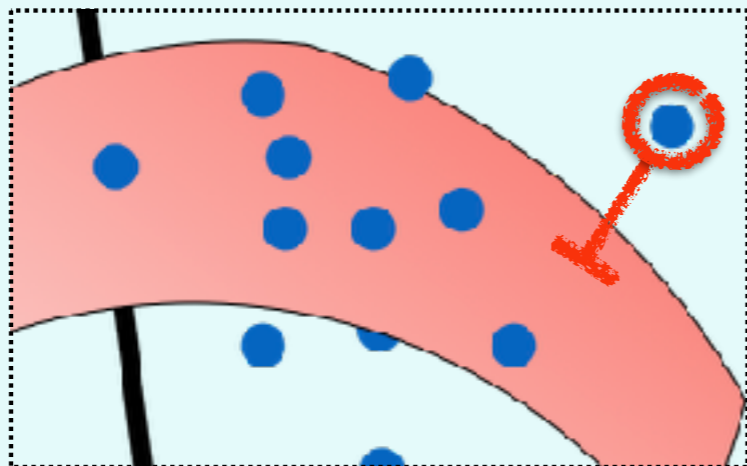


Bags of Features



Brain Activation Maps

Generative Models vs. Data



Color Histograms

Short Course Outline

1. Introduction to optimal transport
2. Optimal transport algorithms
3. Some Applications

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

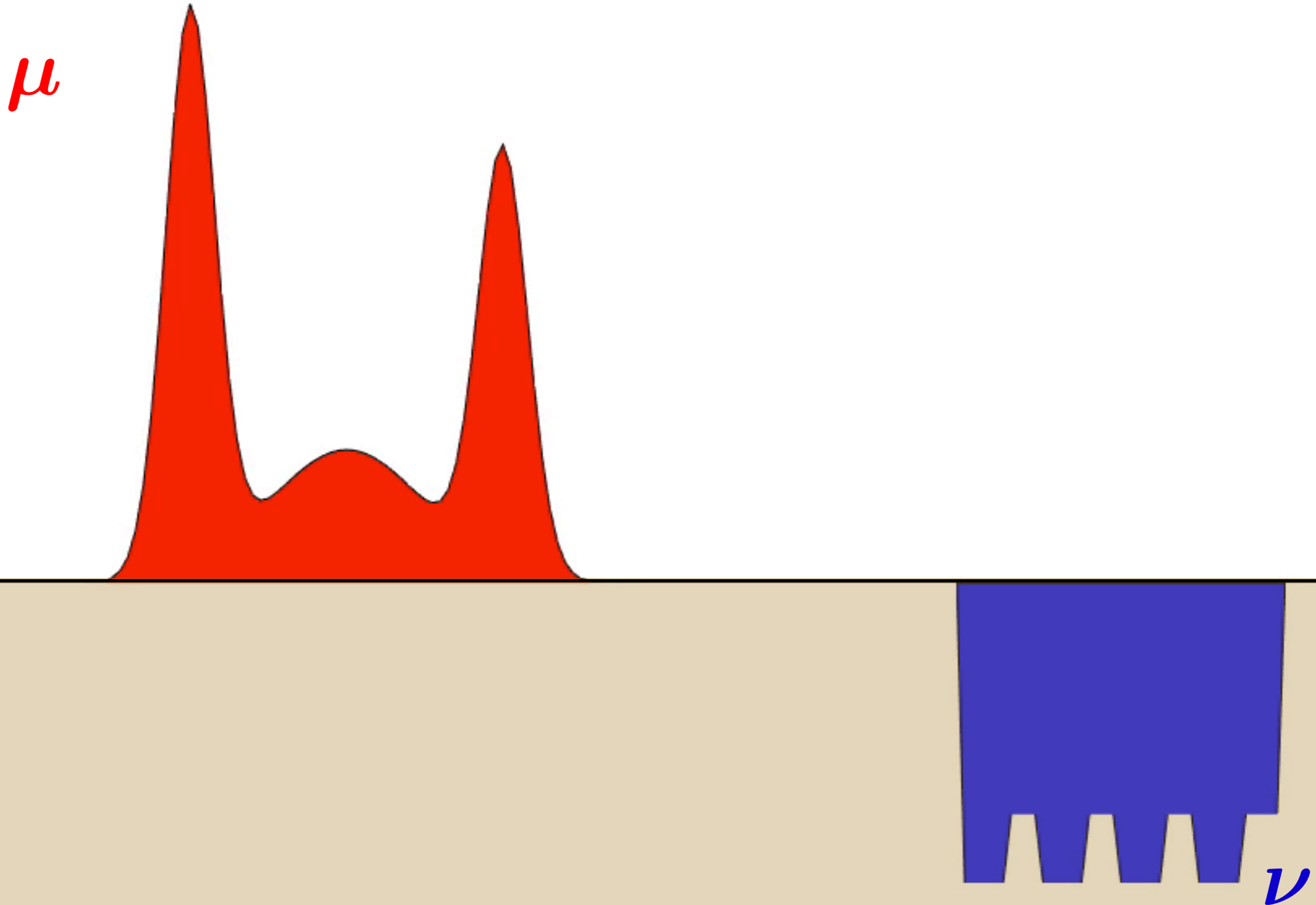
SUR LA

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...*

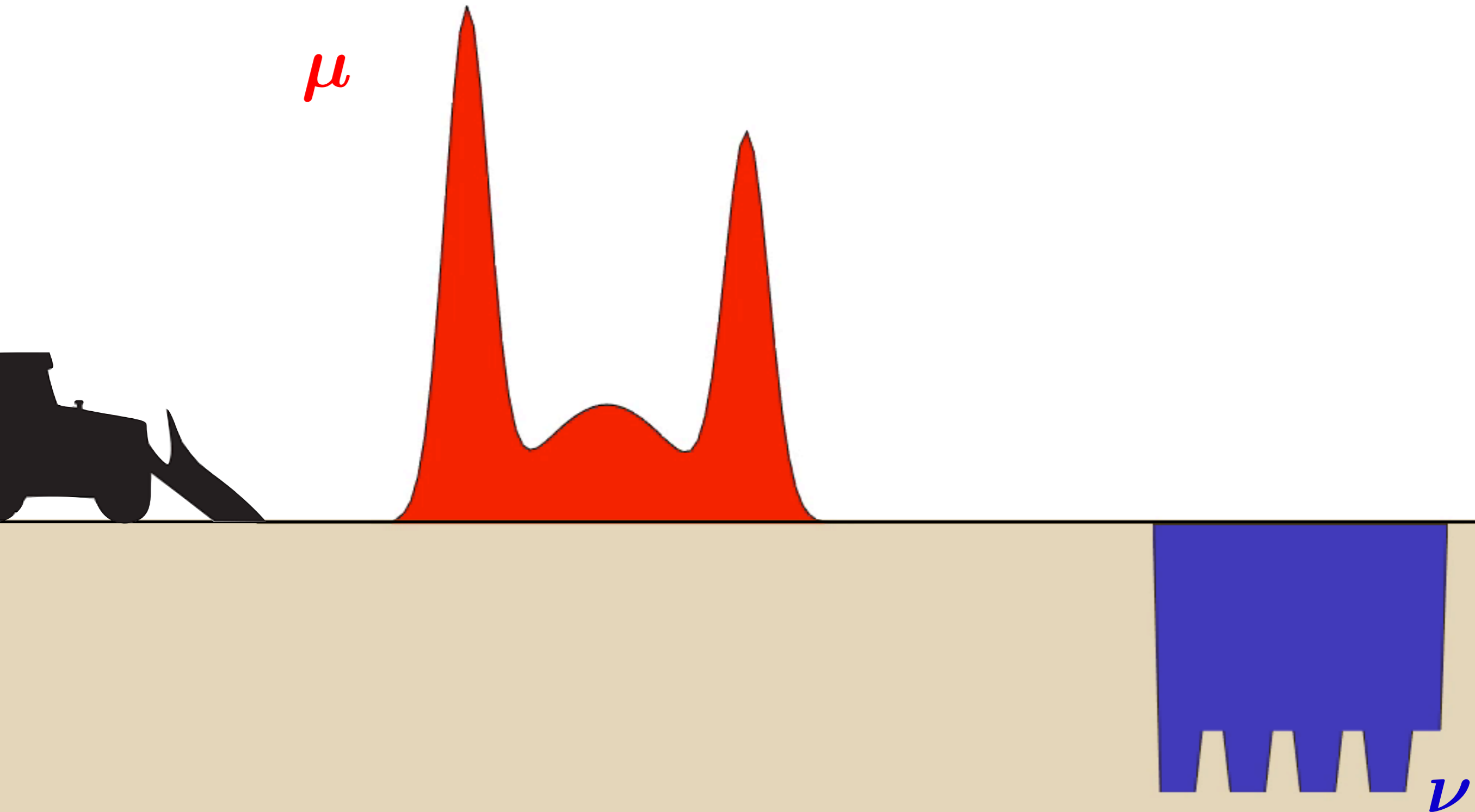
LORSQU'ON doit transporter des terres d'un lieu dans un autre, on a coutume de donner le nom de *Déblai* au volume des terres que l'on doit transporter, & le nom de *Remblai* à l'espace qu'elles doivent occuper après le transport.

Origins: Monge Problem



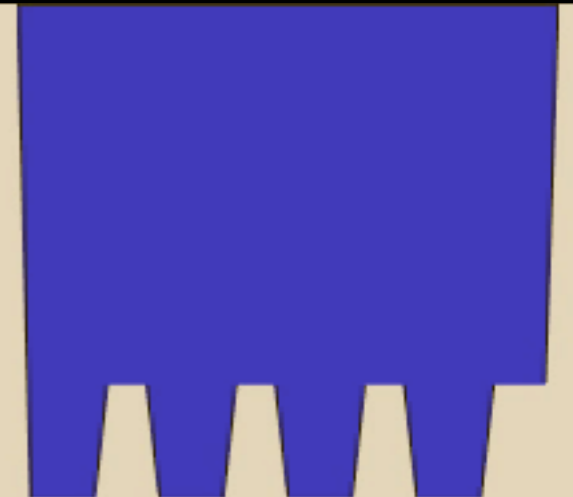
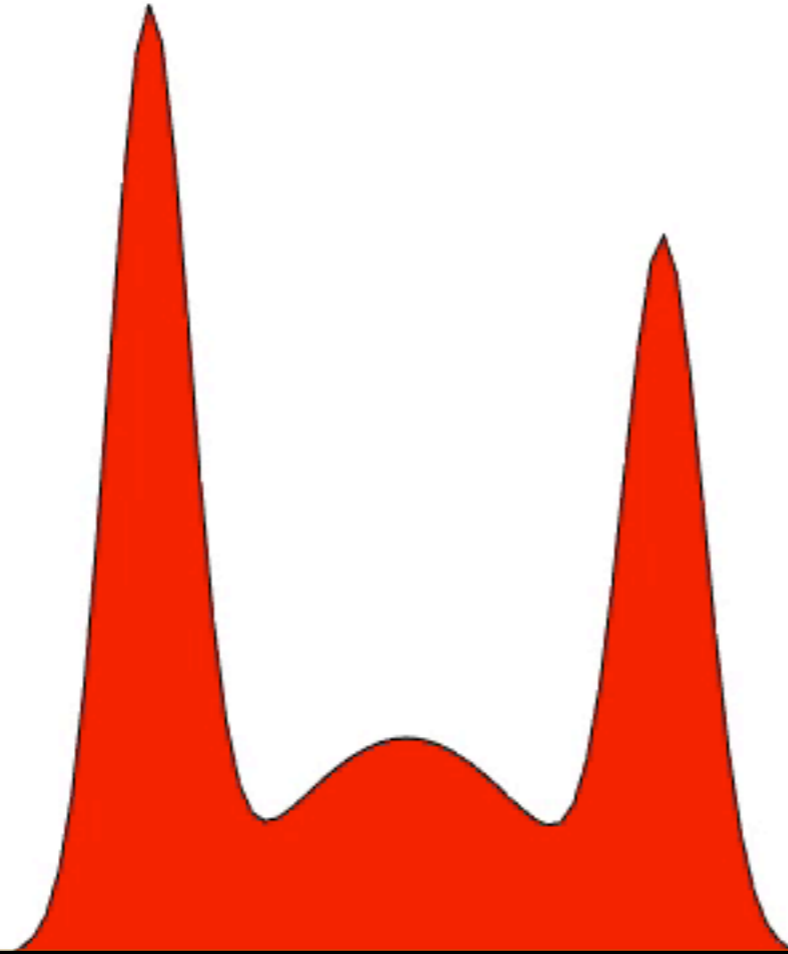
Origins: Monge Problem

In the 21st Century...



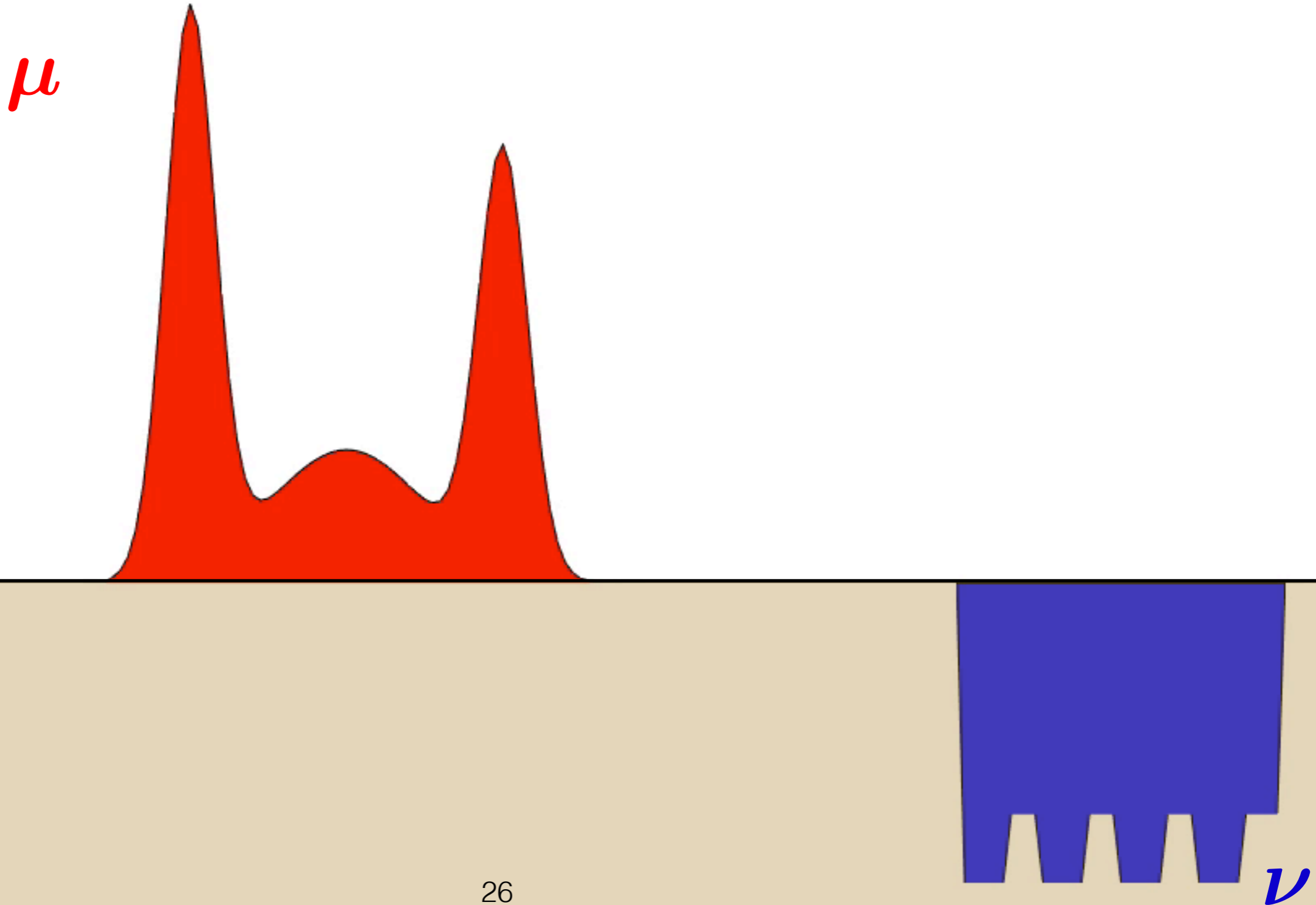
Origins: Monge Problem

In the 21st Century...



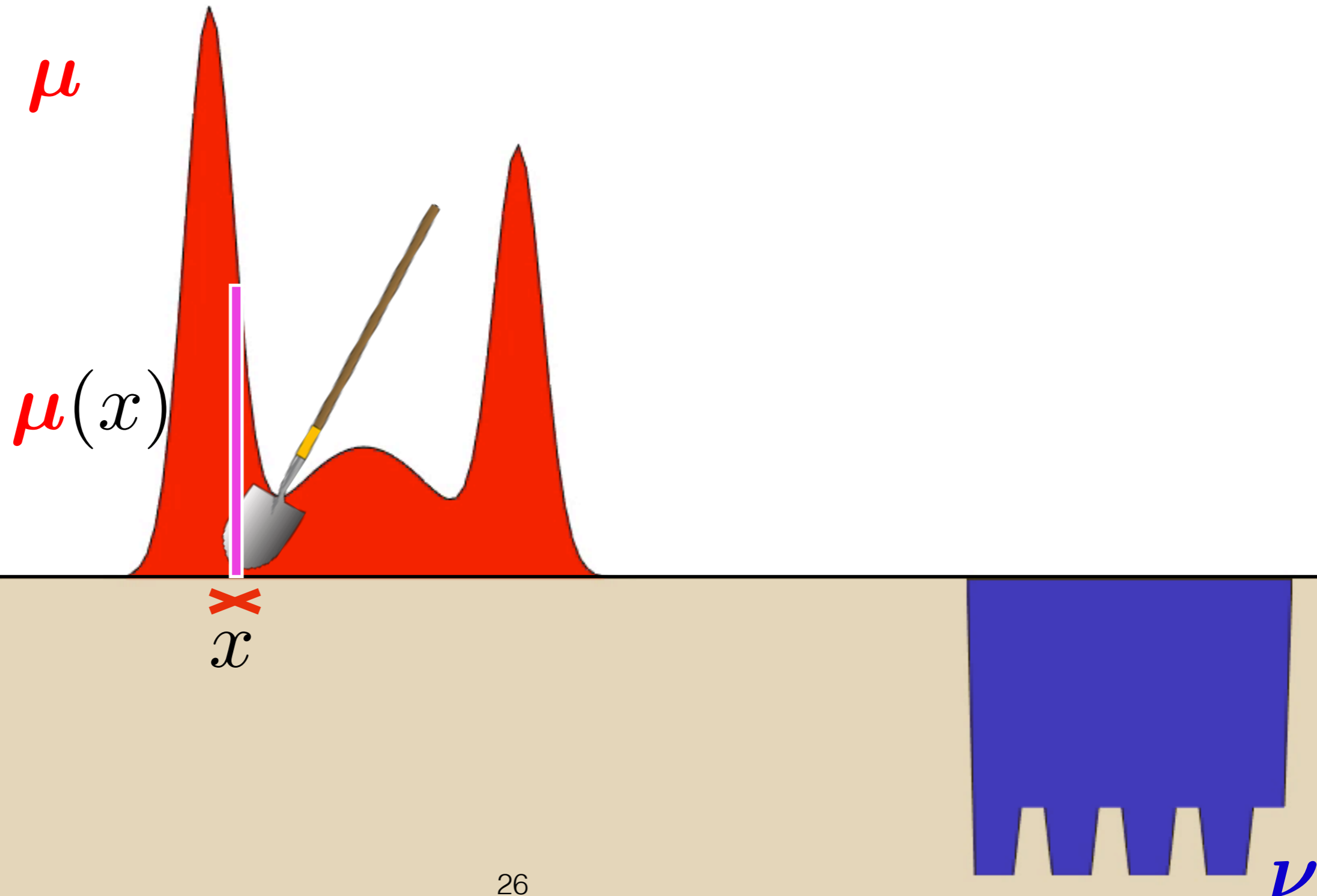
Origins: Monge's Problem

In 1781 however...



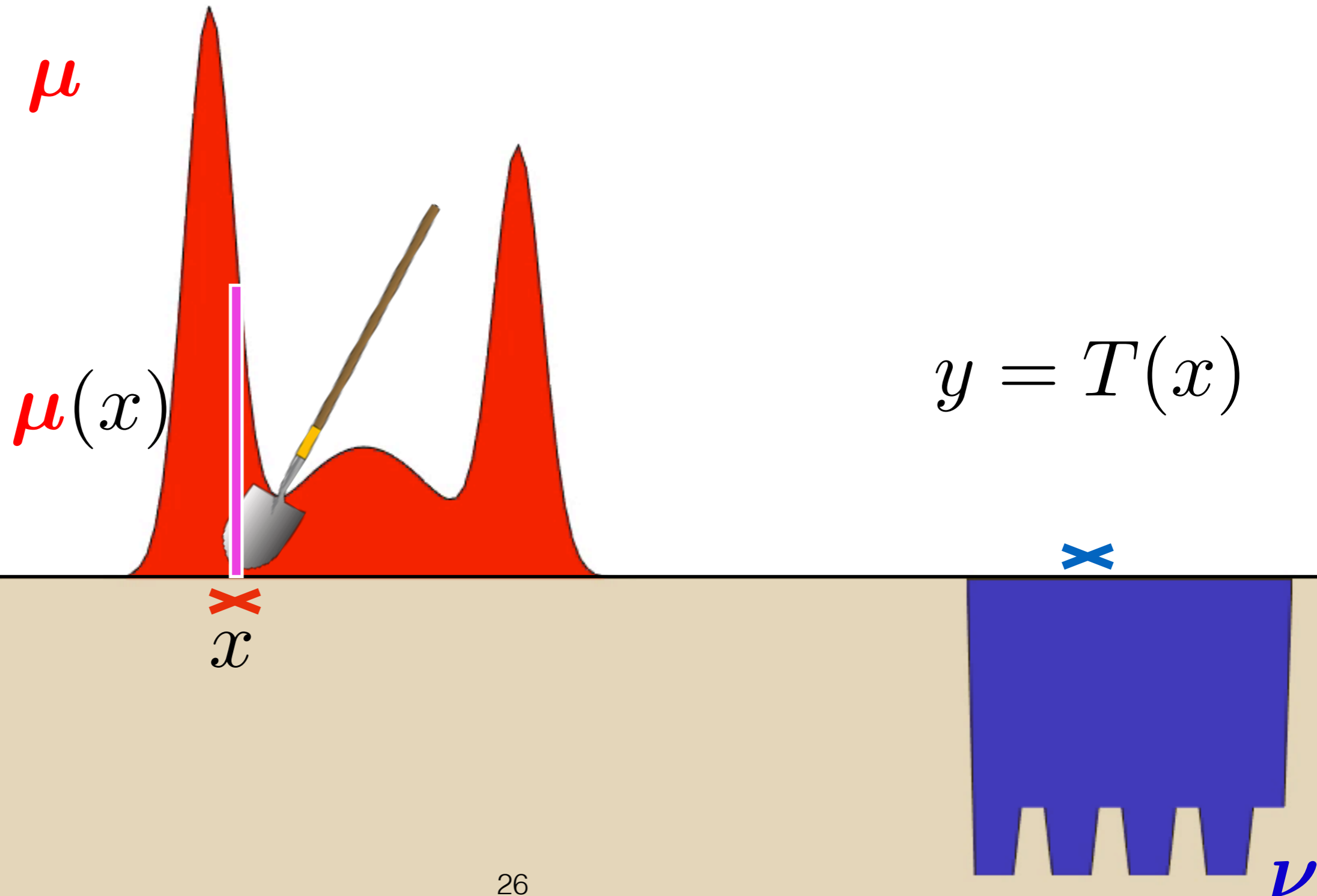
Origins: Monge's Problem

In 1781 however...



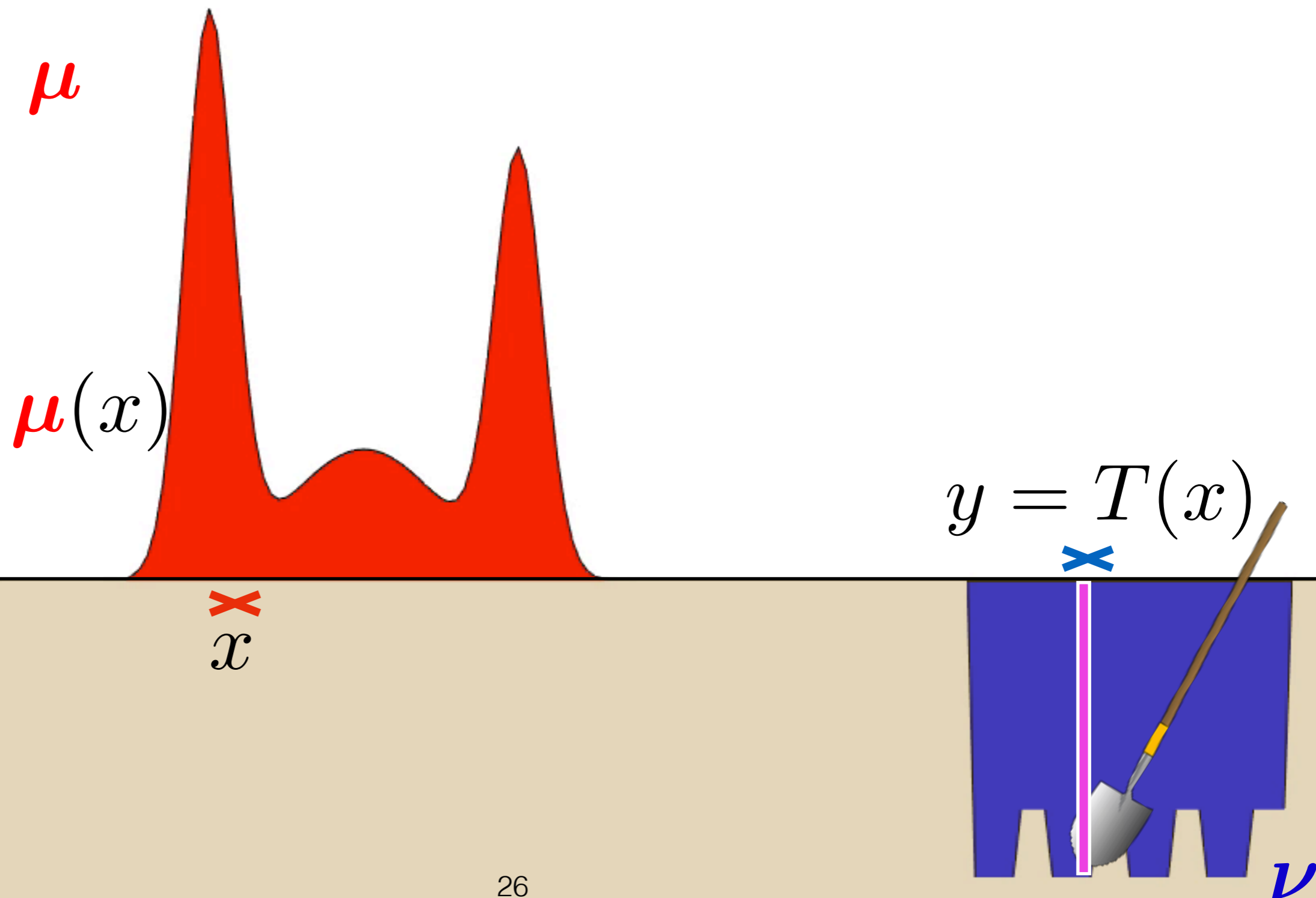
Origins: Monge's Problem

In 1781 however...



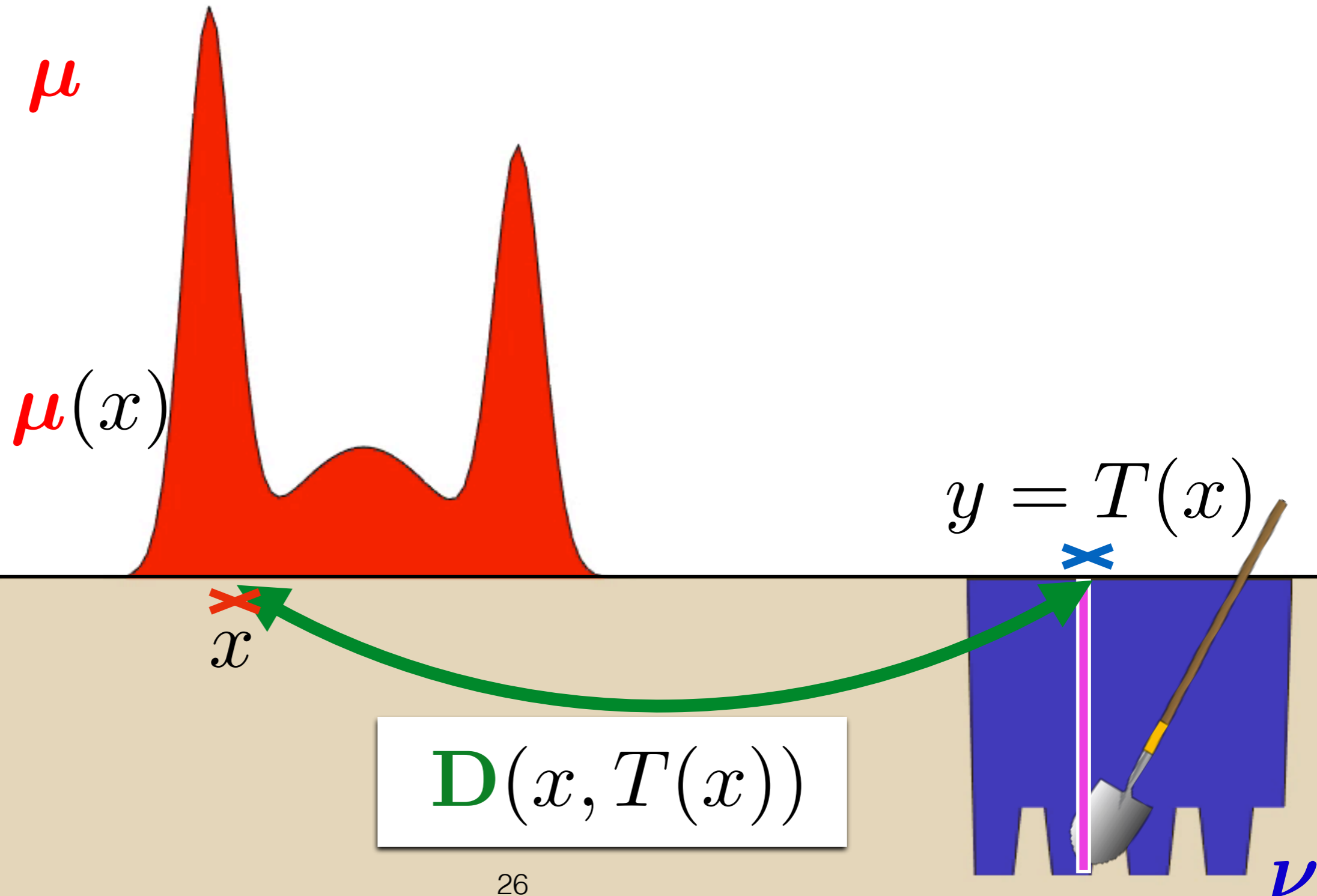
Origins: Monge's Problem

In 1781 however...



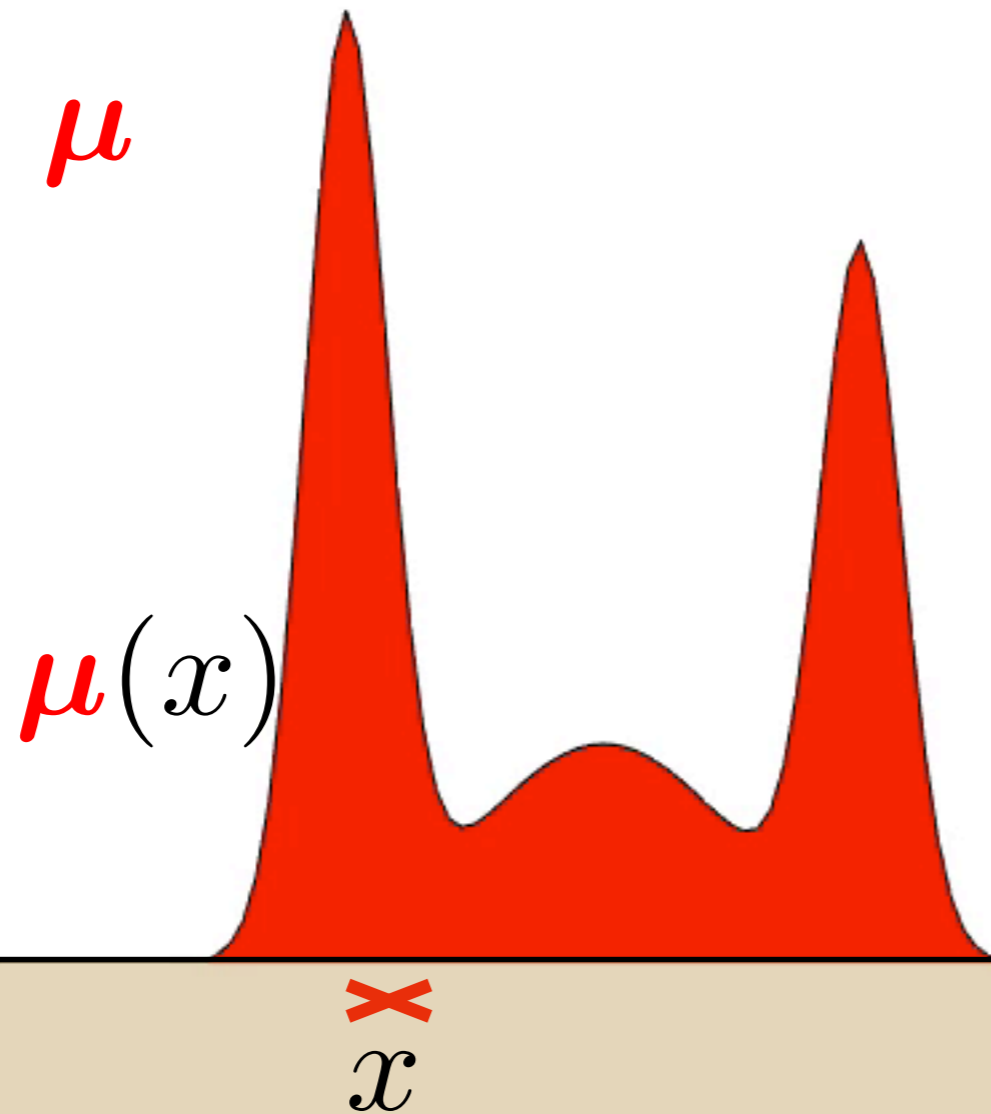
Origins: Monge's Problem

In 1781 however...

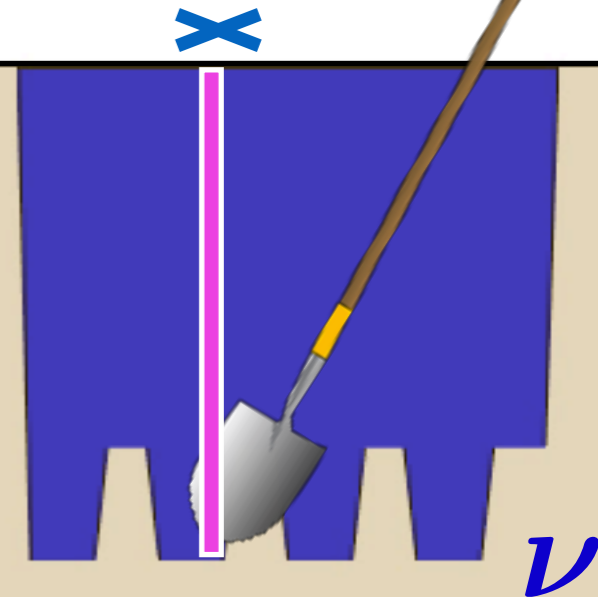


Origins: Monge's Problem

In 1781 however...



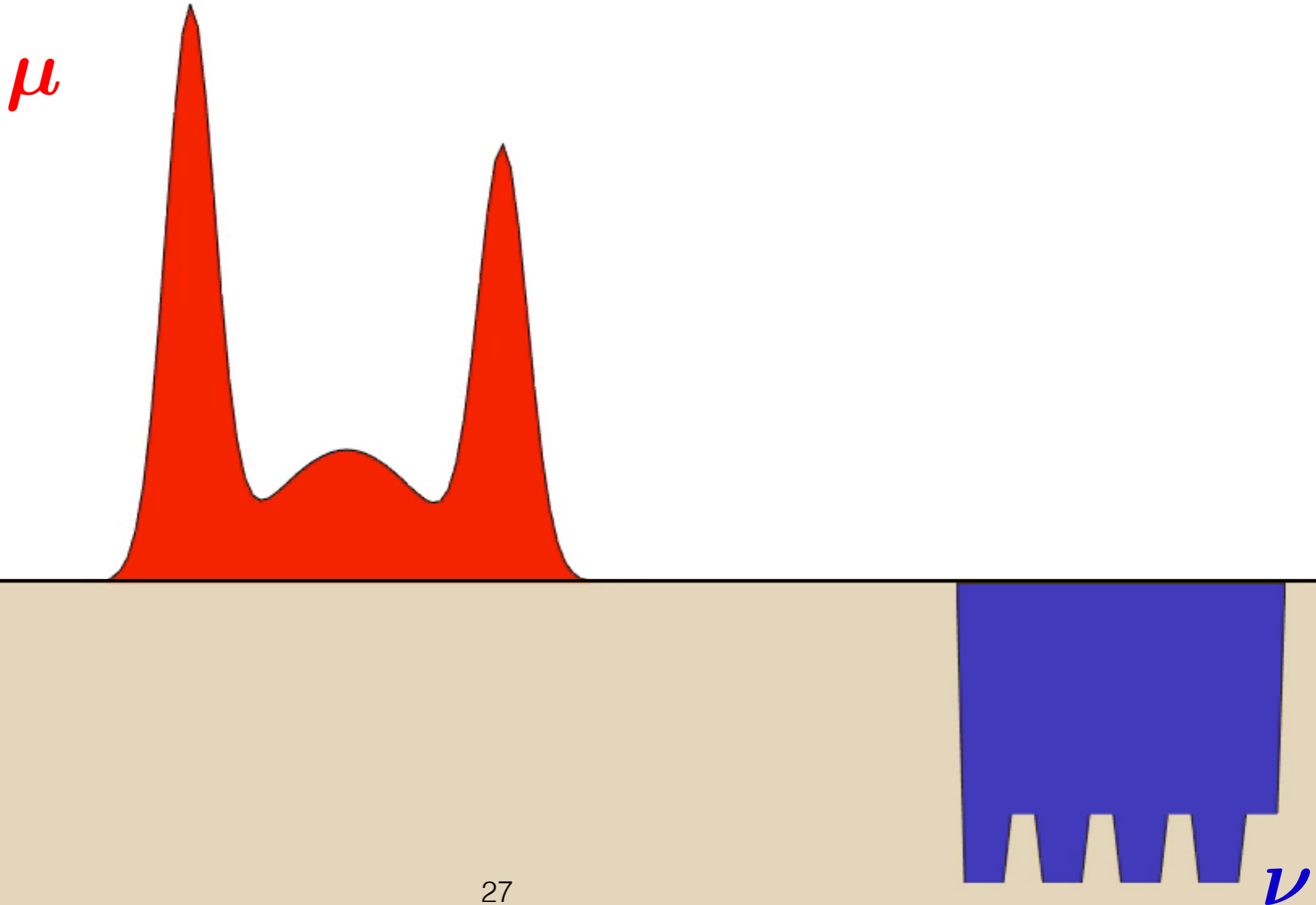
$$y = T(x)$$



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

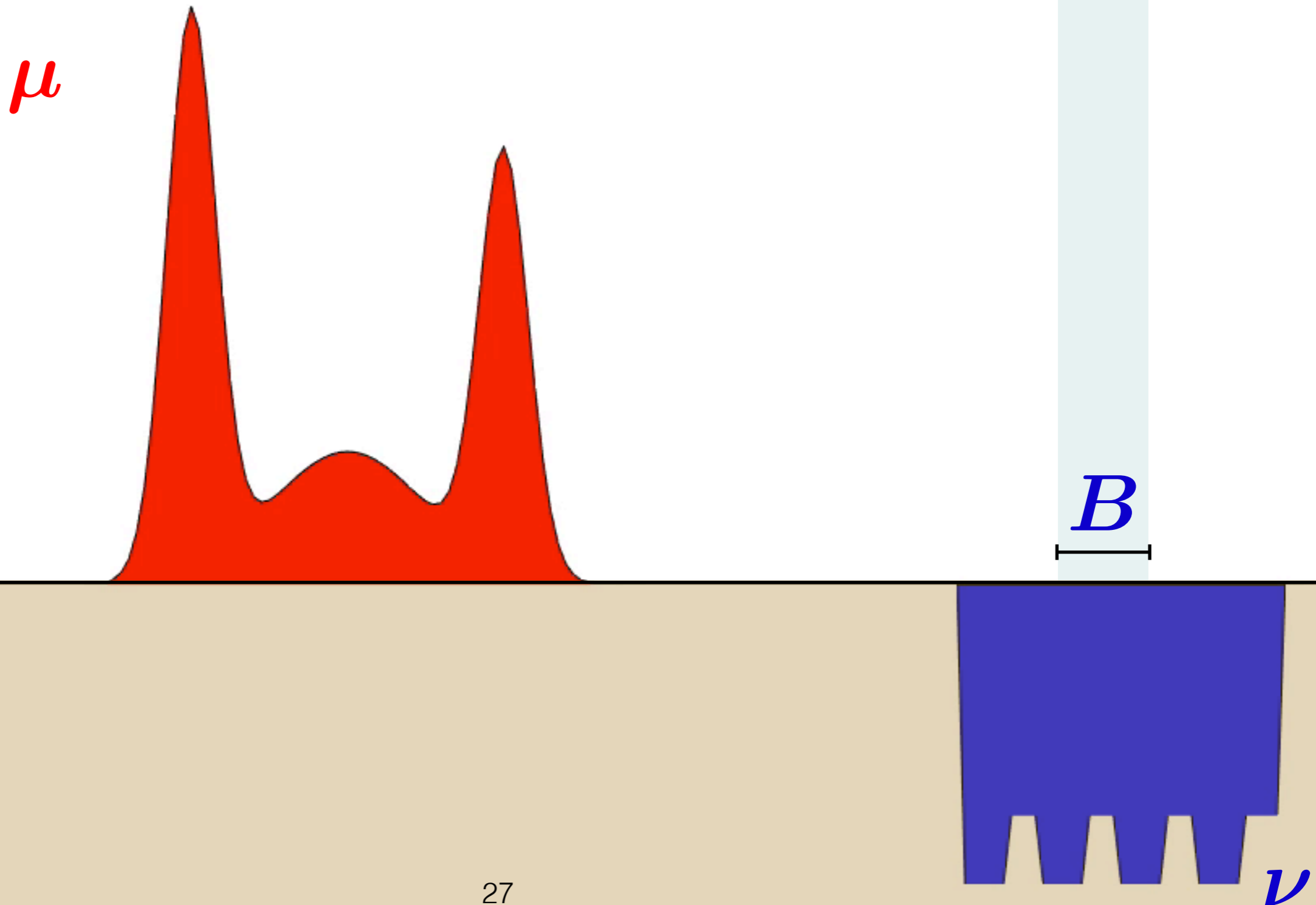
Origins: Monge's Problem

T must map red to blue.



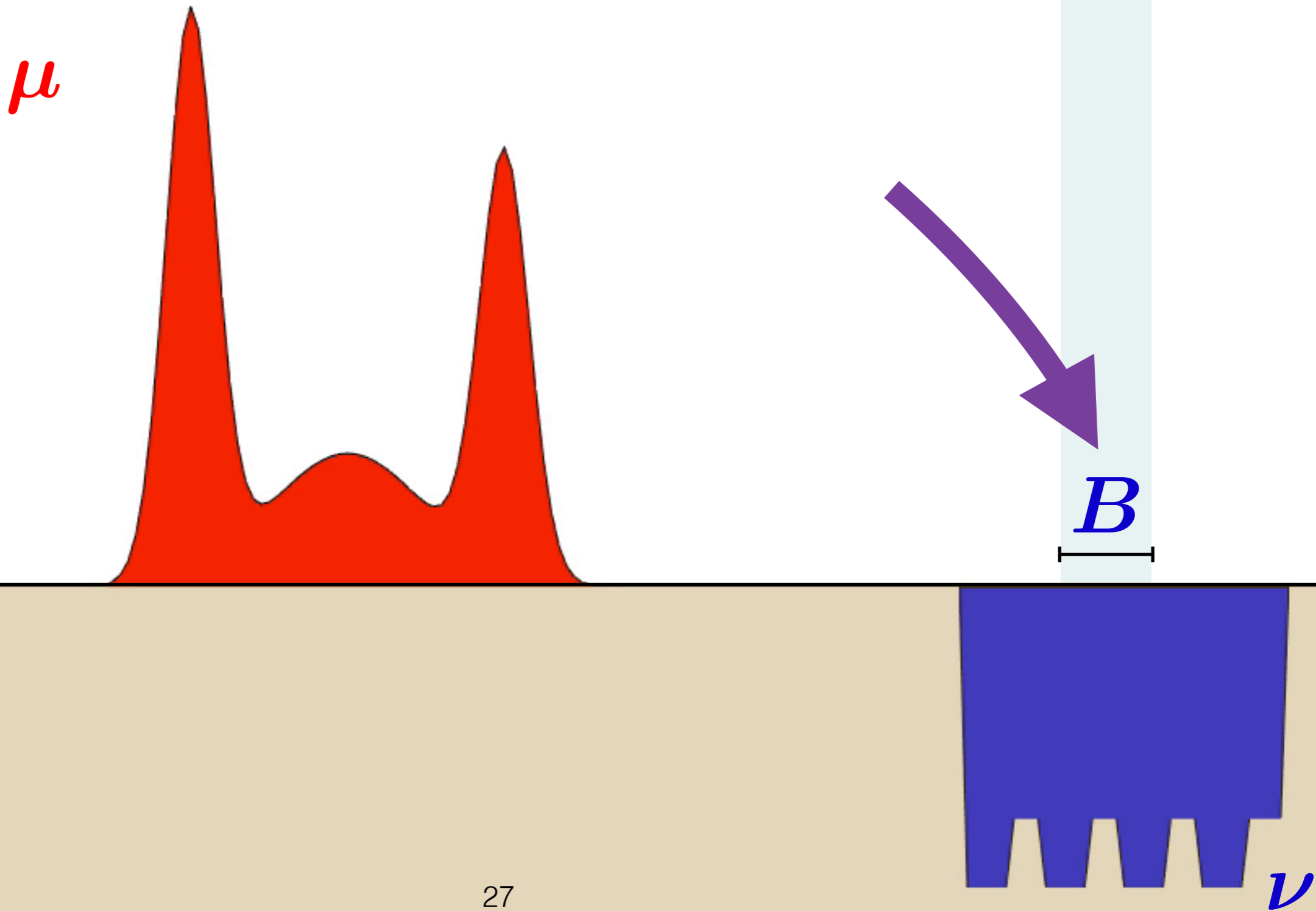
Origins: Monge's Problem

T must map red to blue.



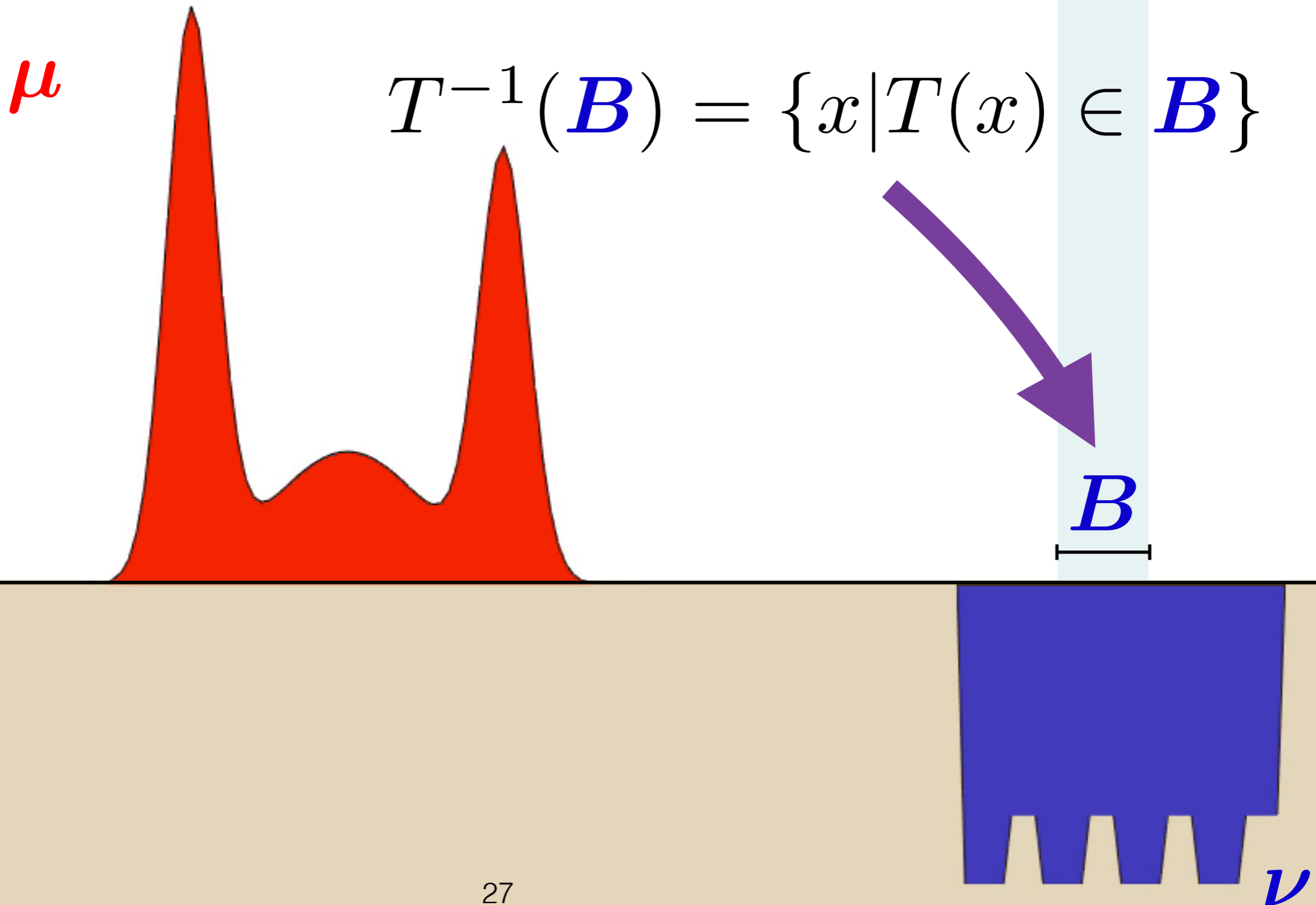
Origins: Monge's Problem

T must map red to blue.



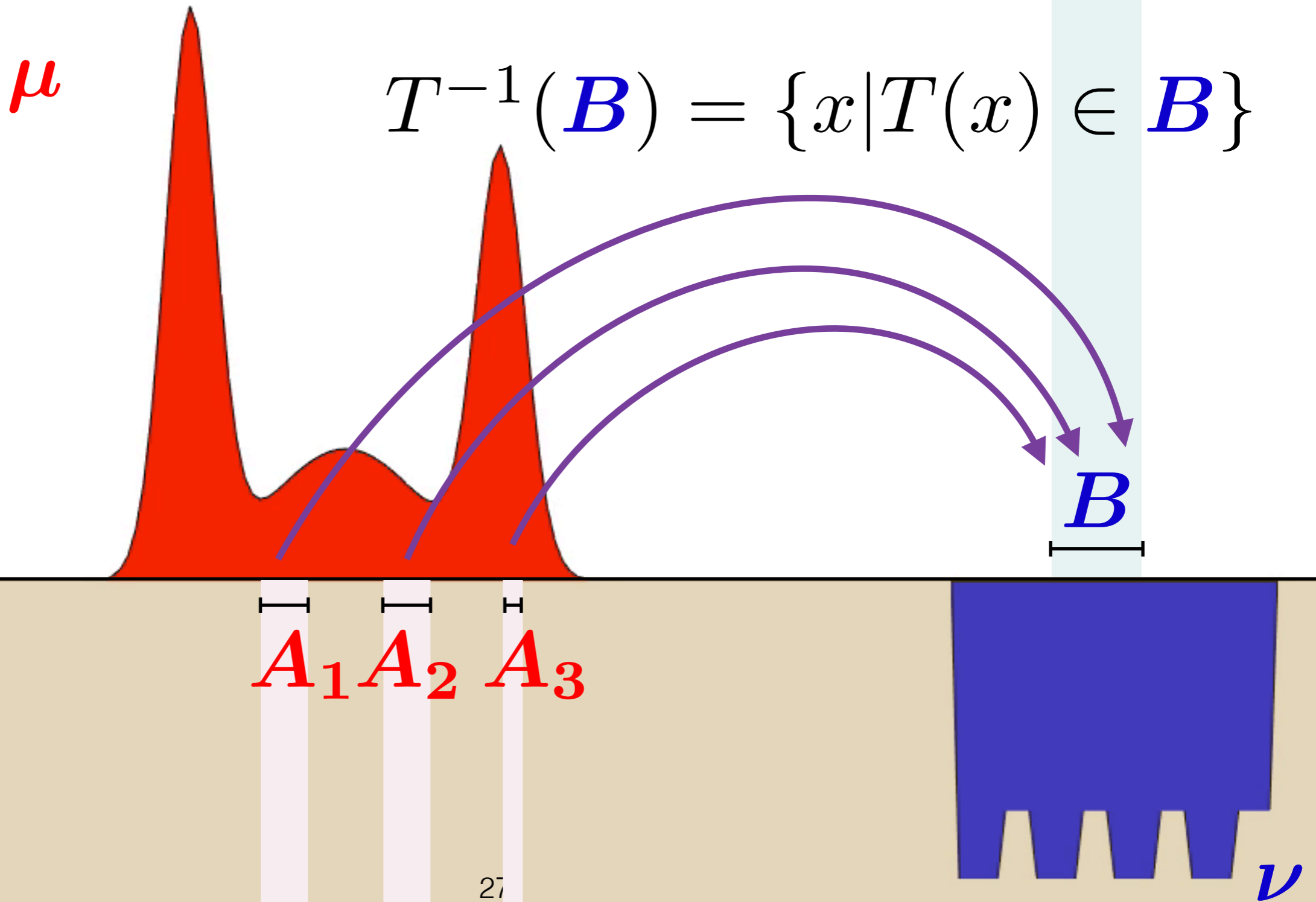
Origins: Monge's Problem

T must map red to blue.



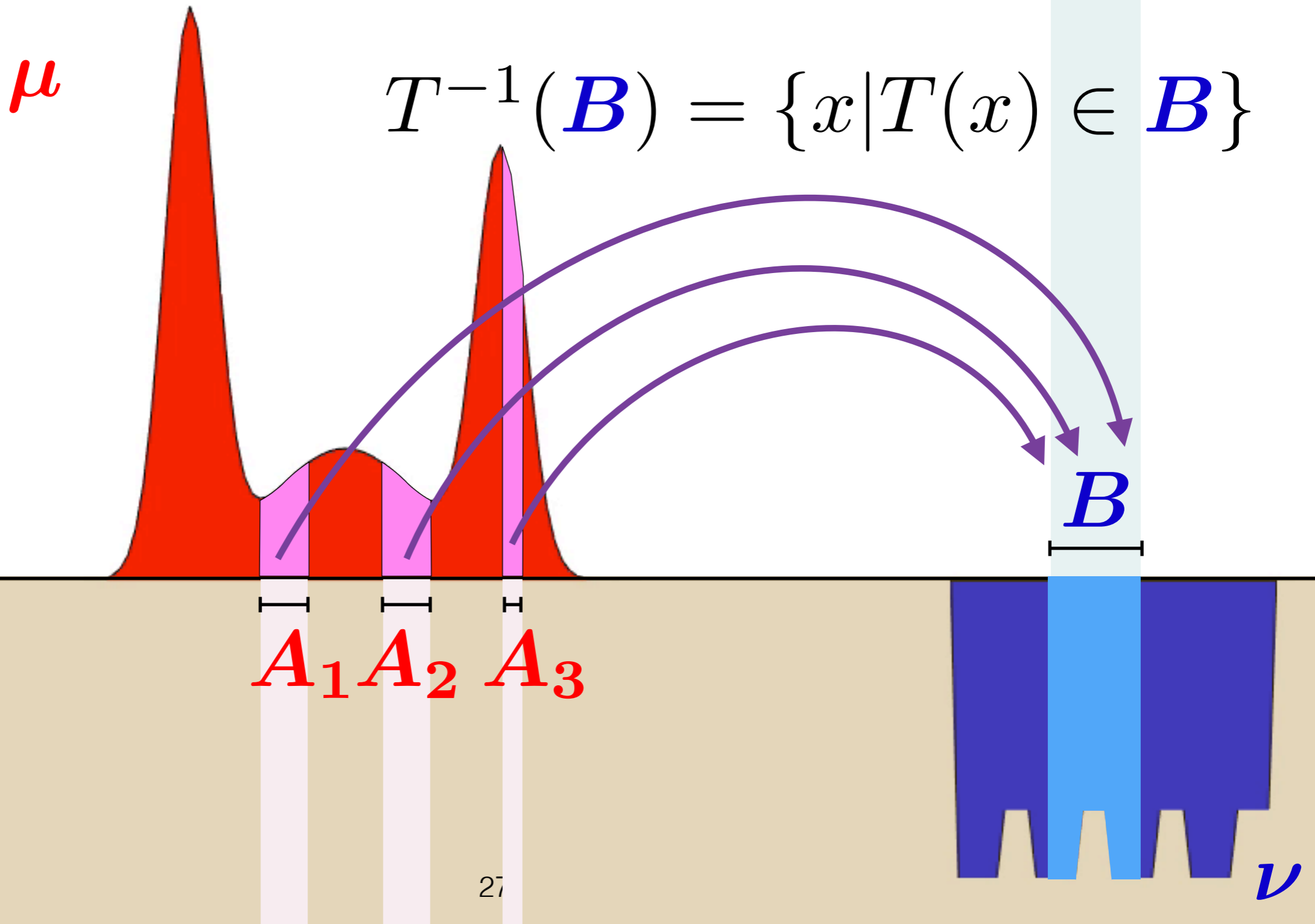
Origins: Monge's Problem

T must map red to blue.



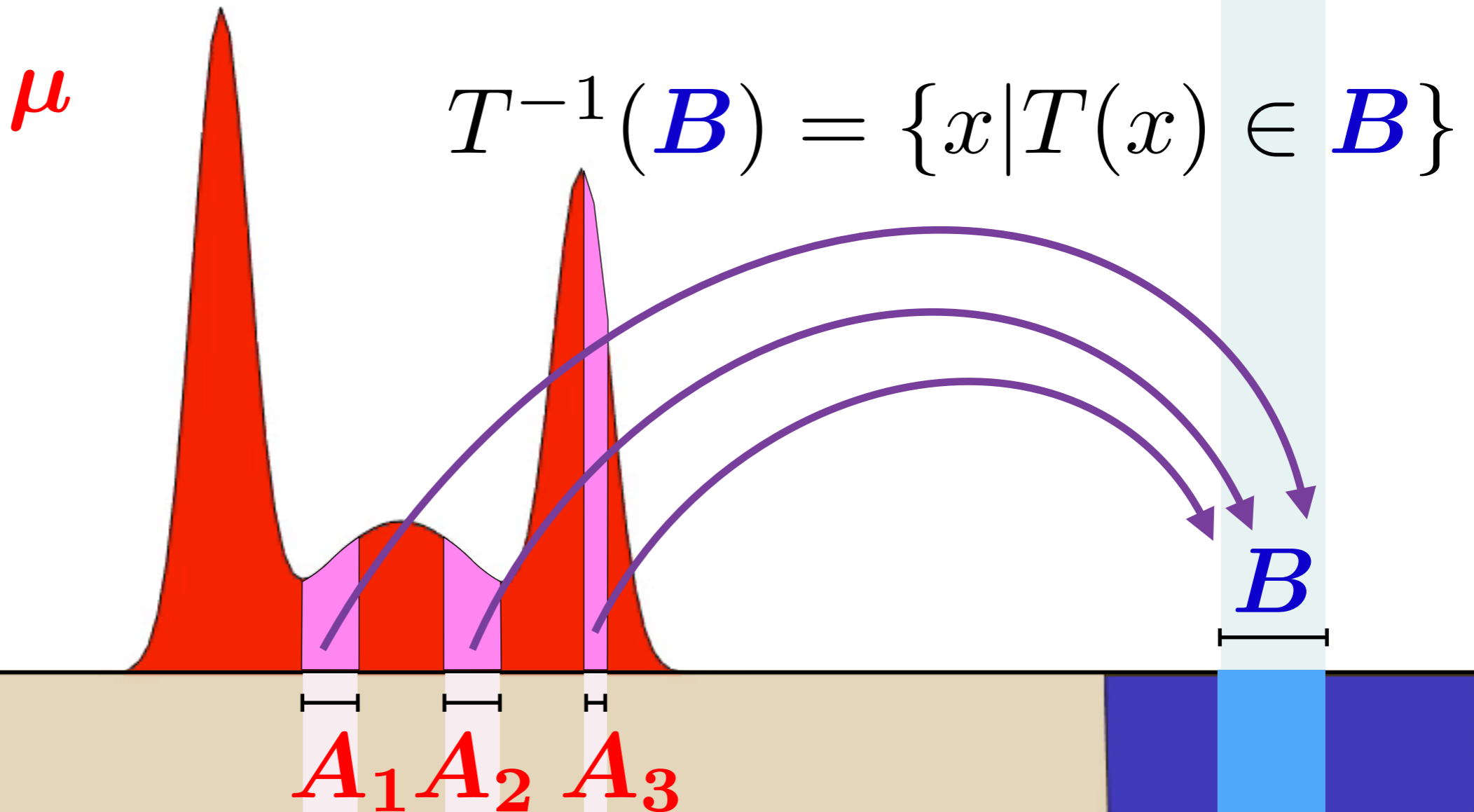
Origins: Monge's Problem

T must map red to blue.



Origins: Monge's Problem

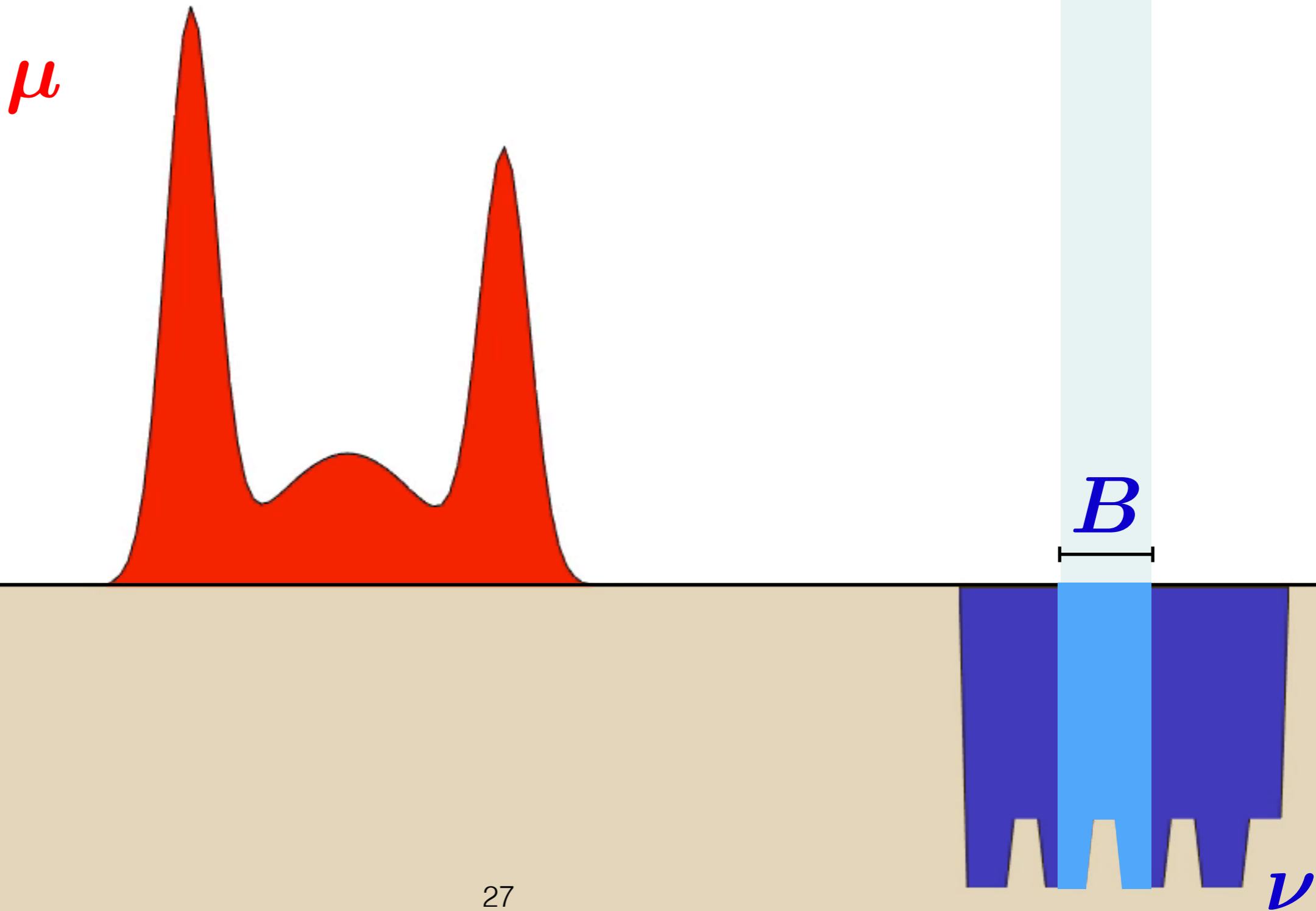
T must map red to blue.



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

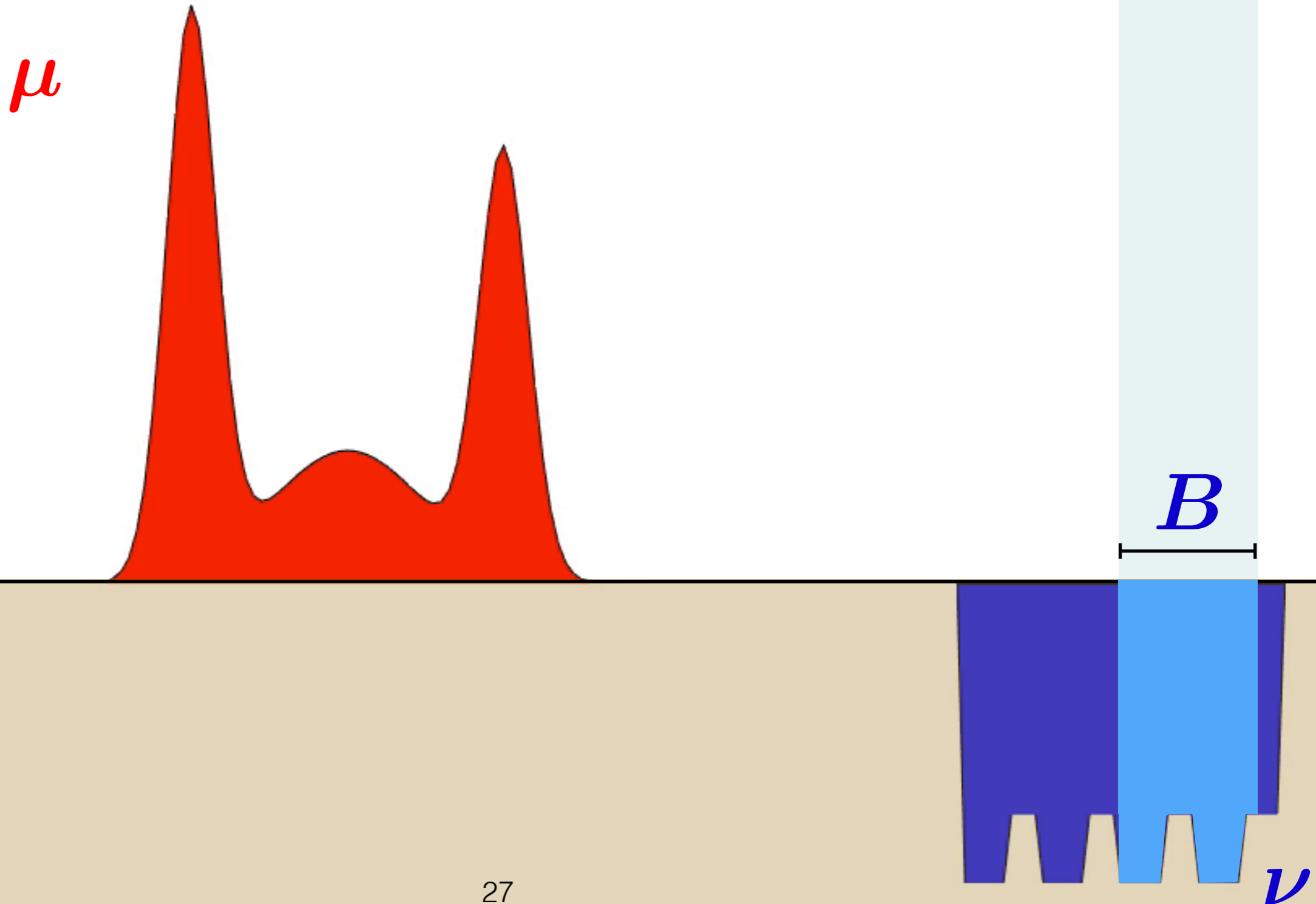
Origins: Monge's Problem

T must map red to blue.



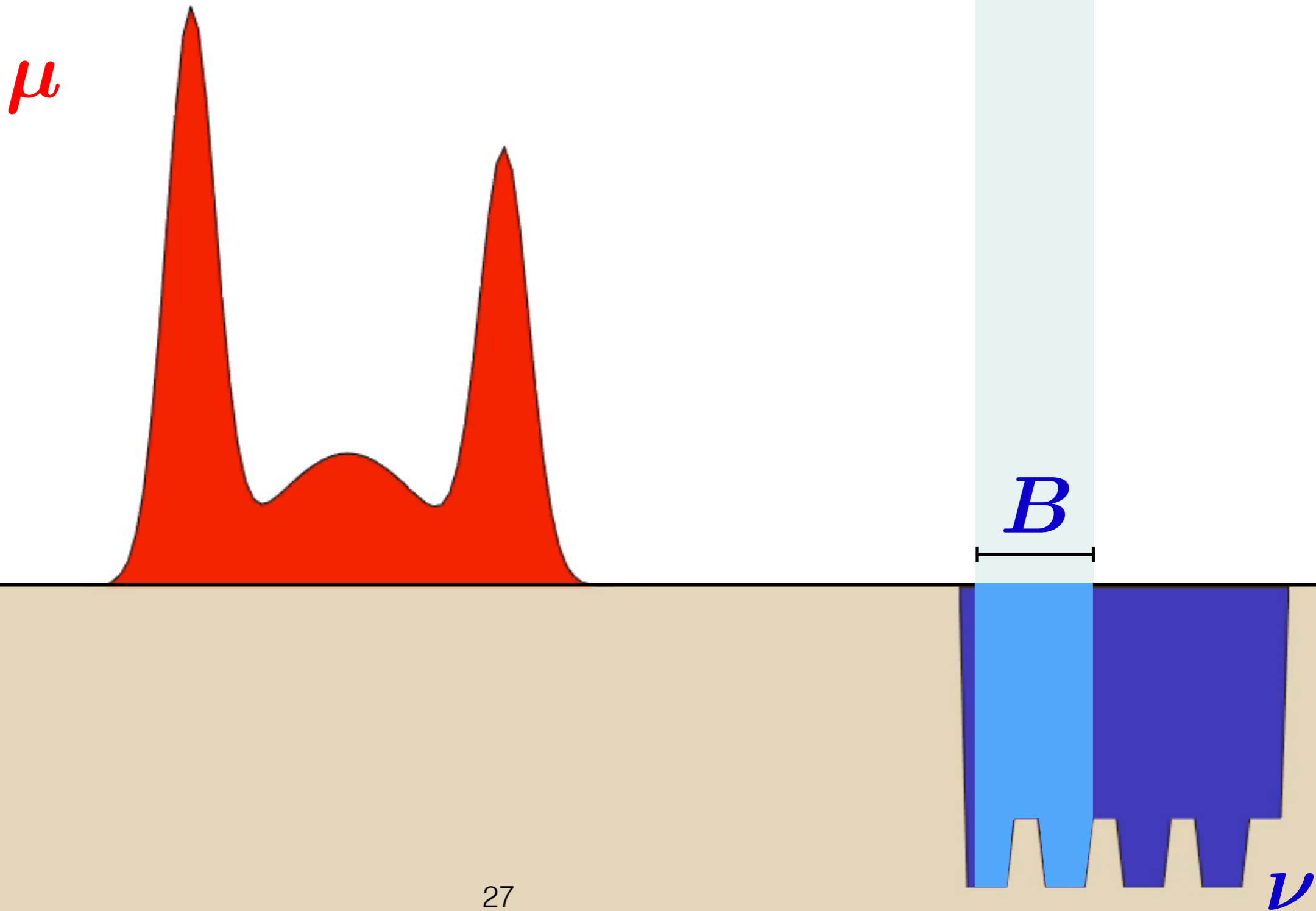
Origins: Monge's Problem

T must map red to blue.



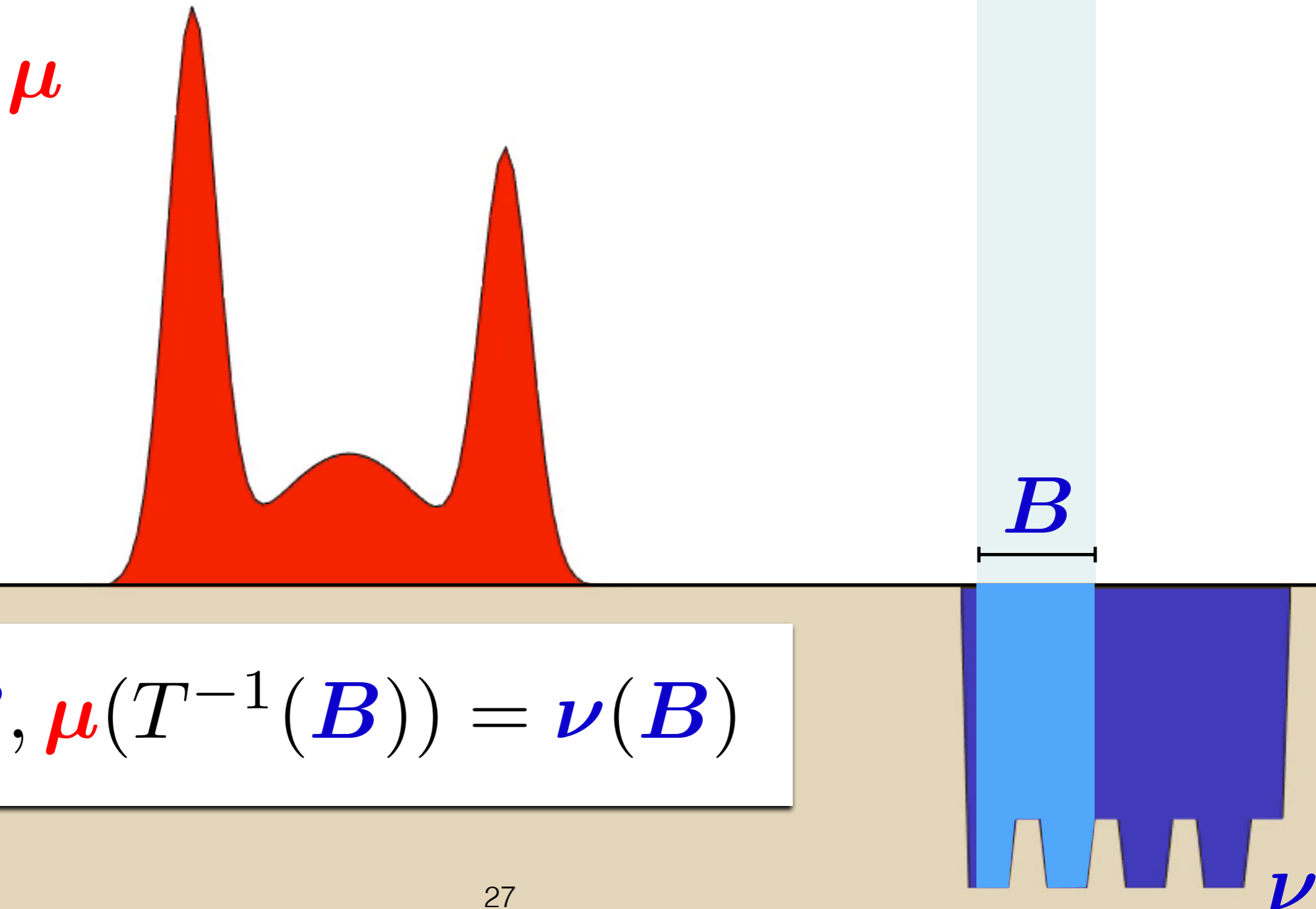
Origins: Monge's Problem

T must map red to blue.



Origins: Monge's Problem

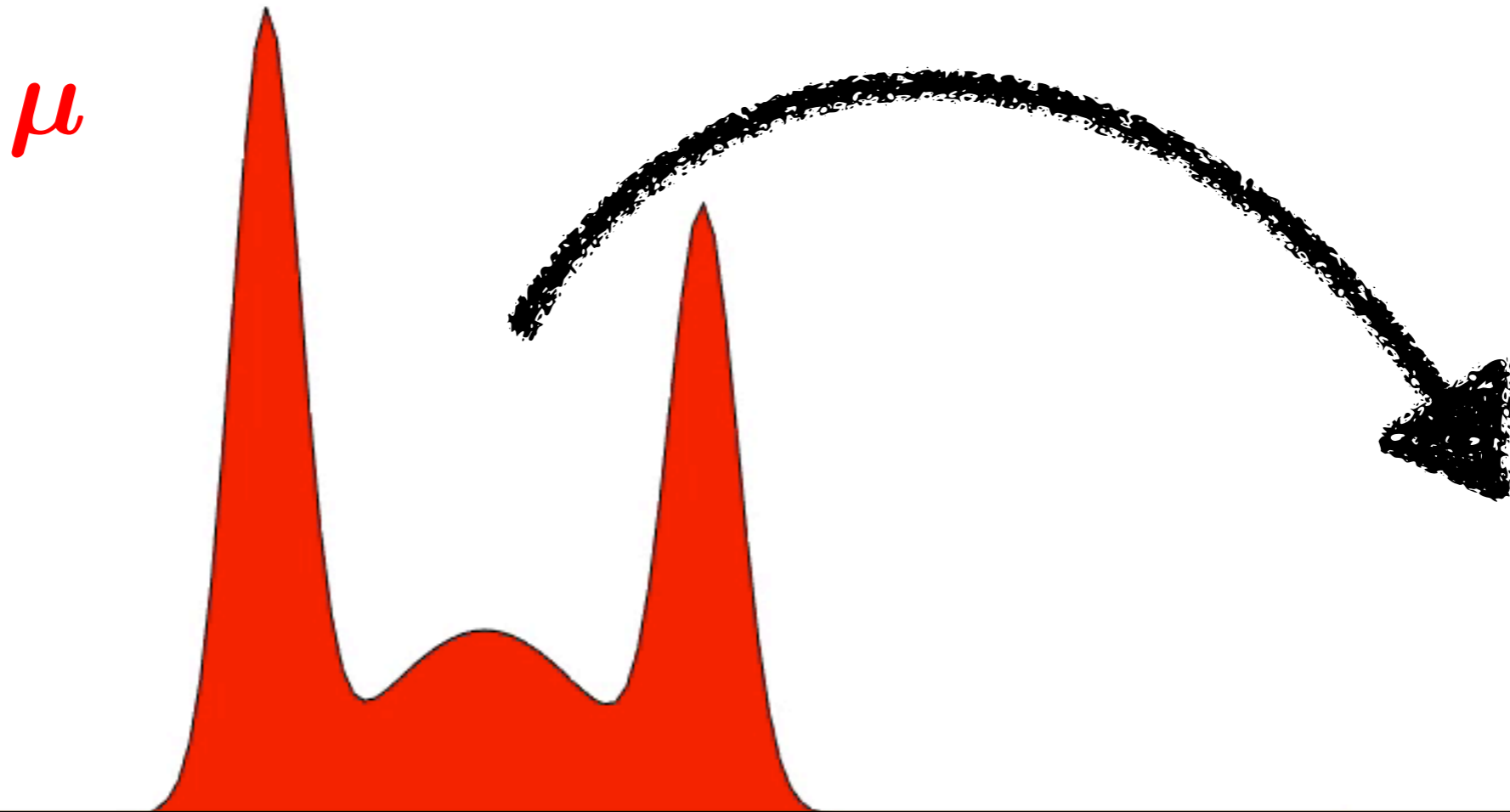
T must map red to blue.



$$\forall B, \mu(T^{-1}(B)) = \nu(B)$$

Origins: Monge's Problem

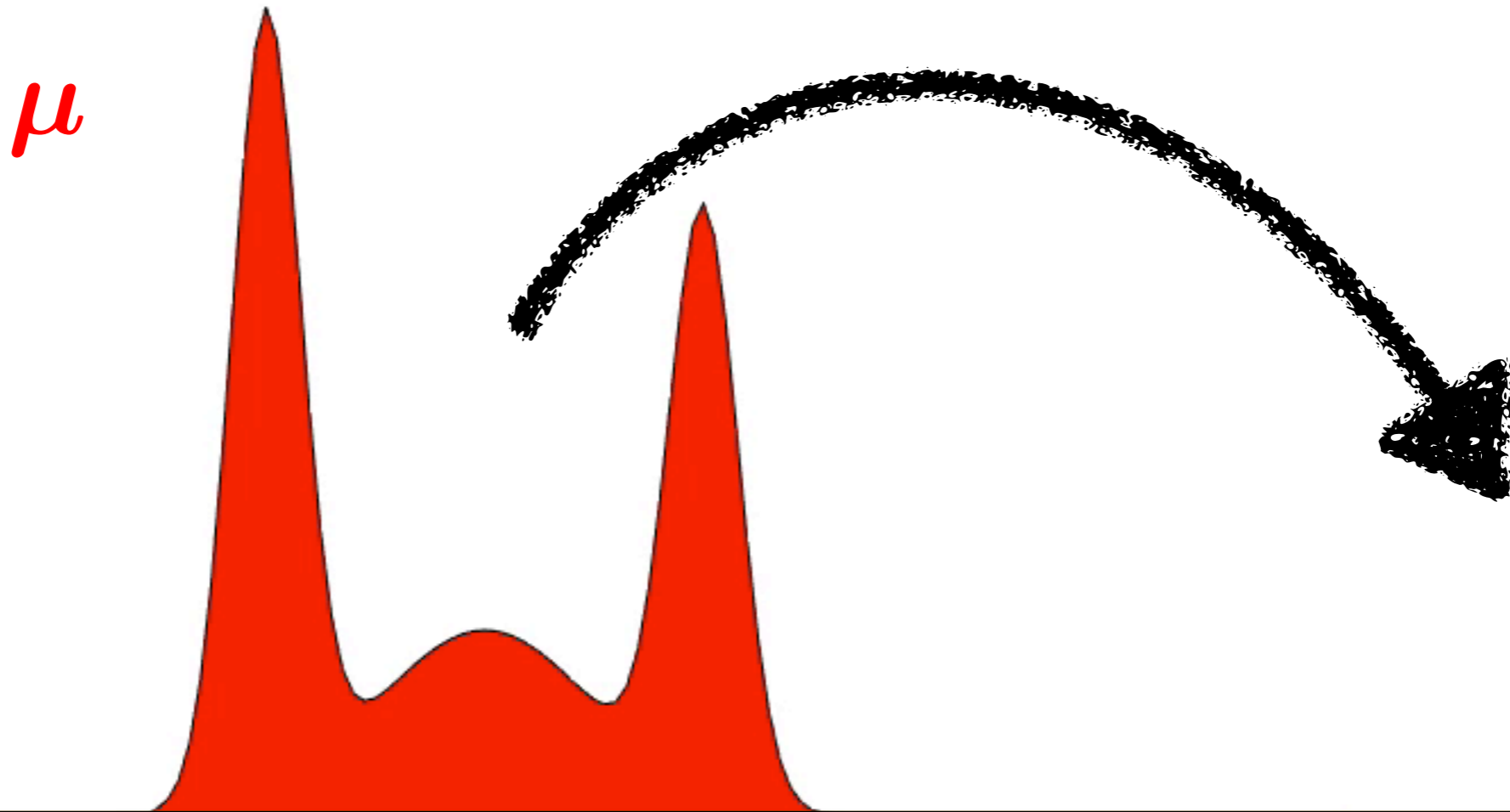
T must **push-forward** the red measure towards the blue



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

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

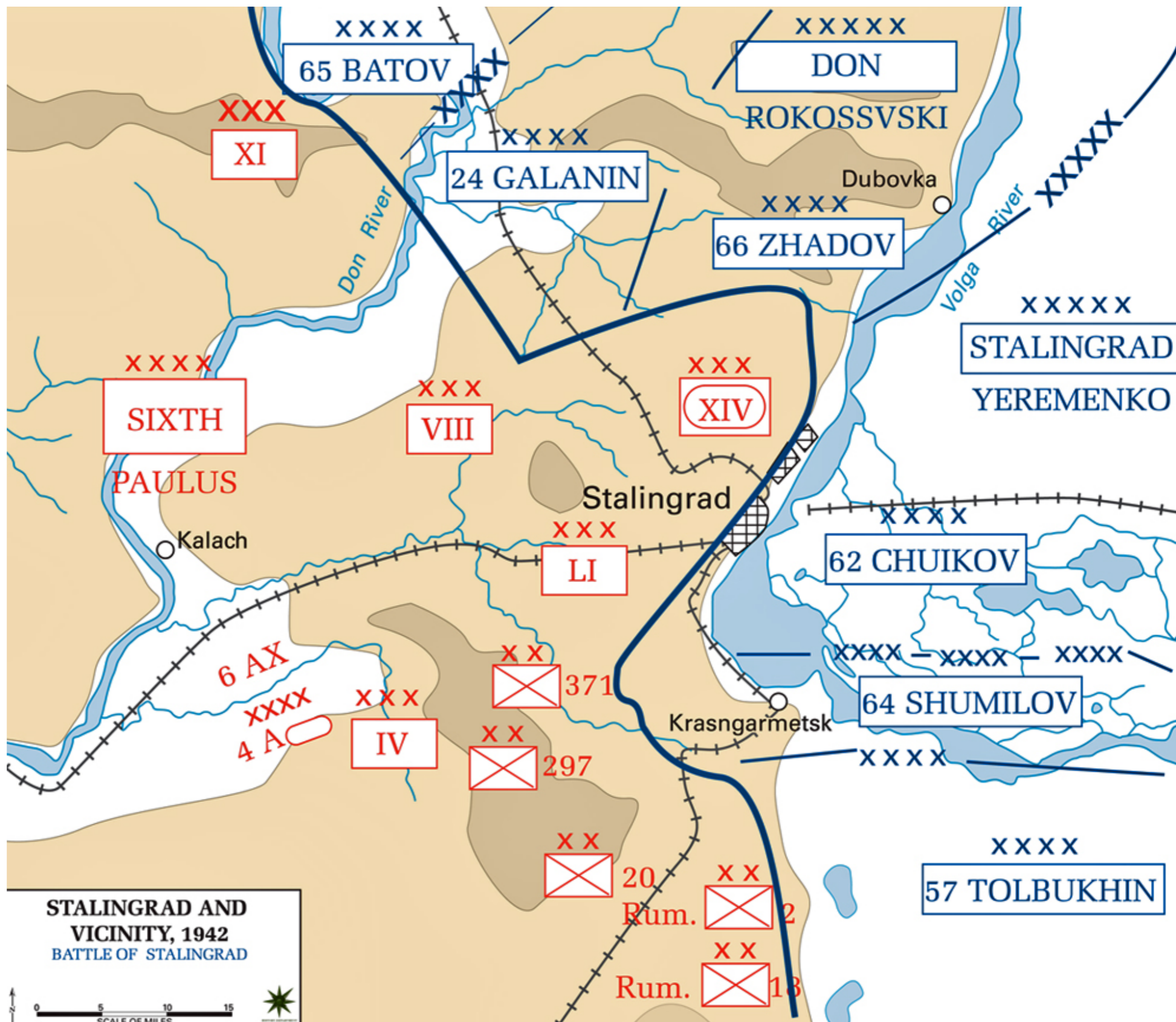
THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES

BY FRANK L. HITCHCOCK

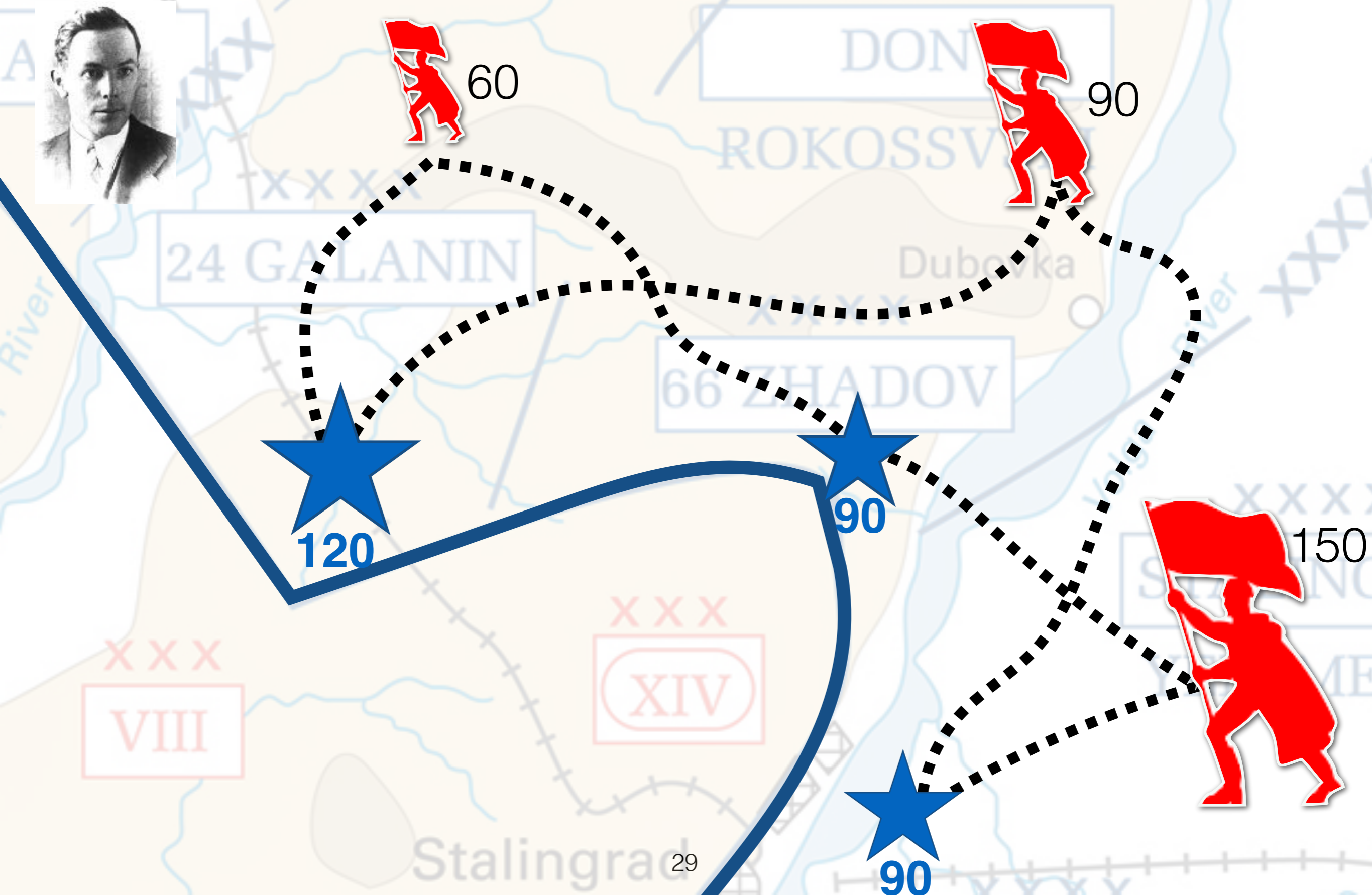
1. Statement of the problem. When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary from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

1941

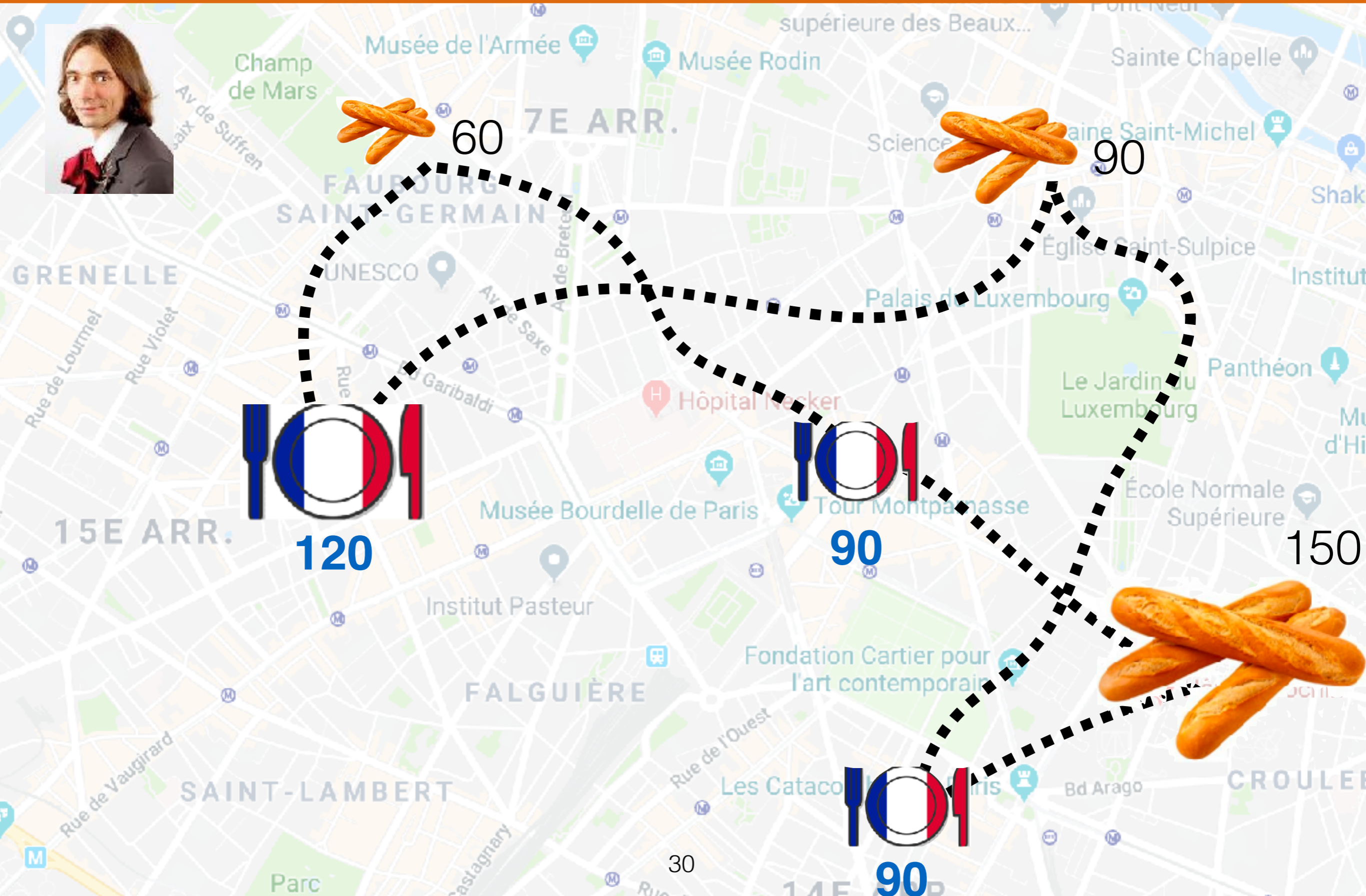
Kantorovich Problem



Kantorovich Problem



Kantorovich Problem *à la française*



120

60

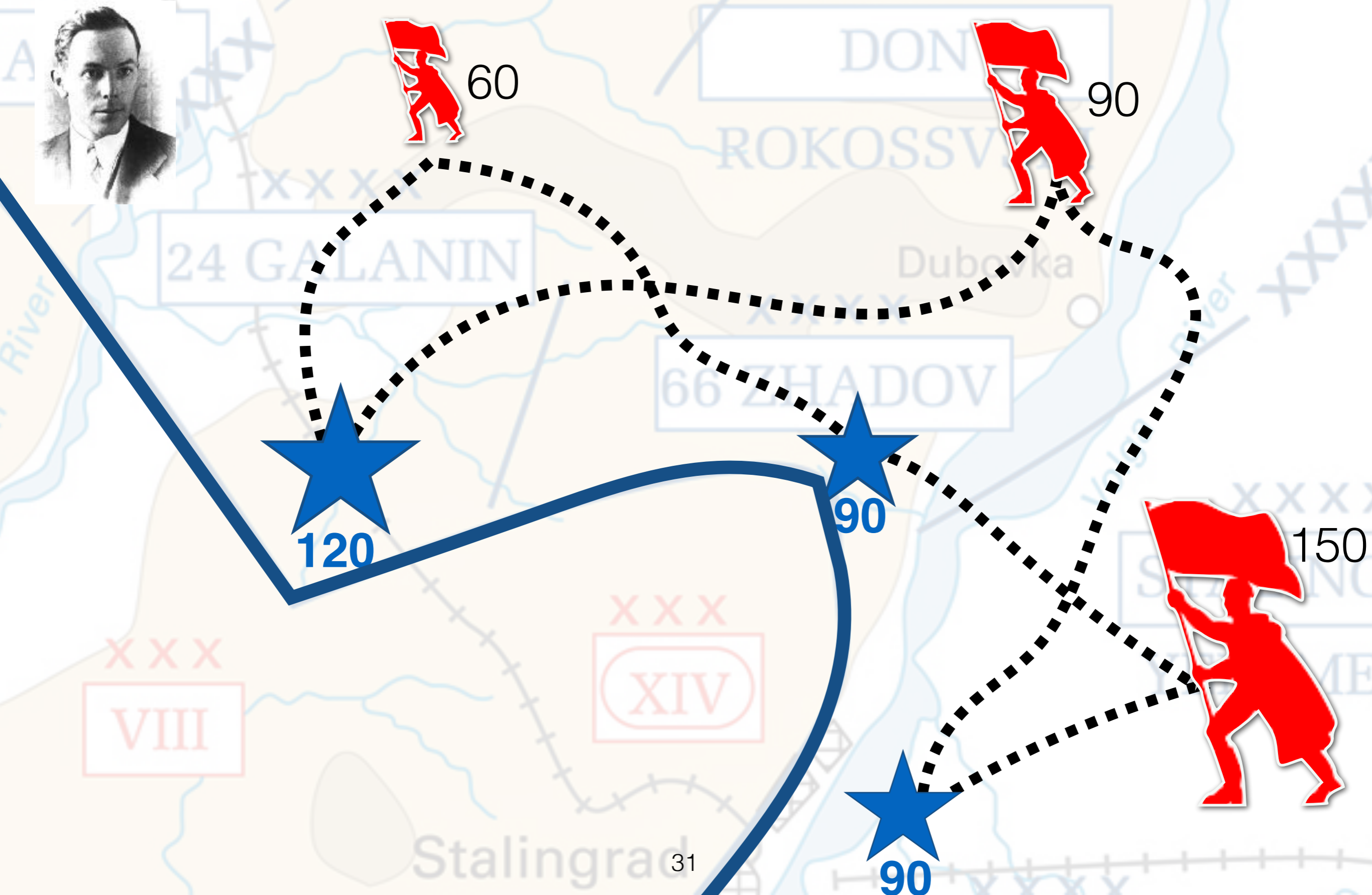
90

90

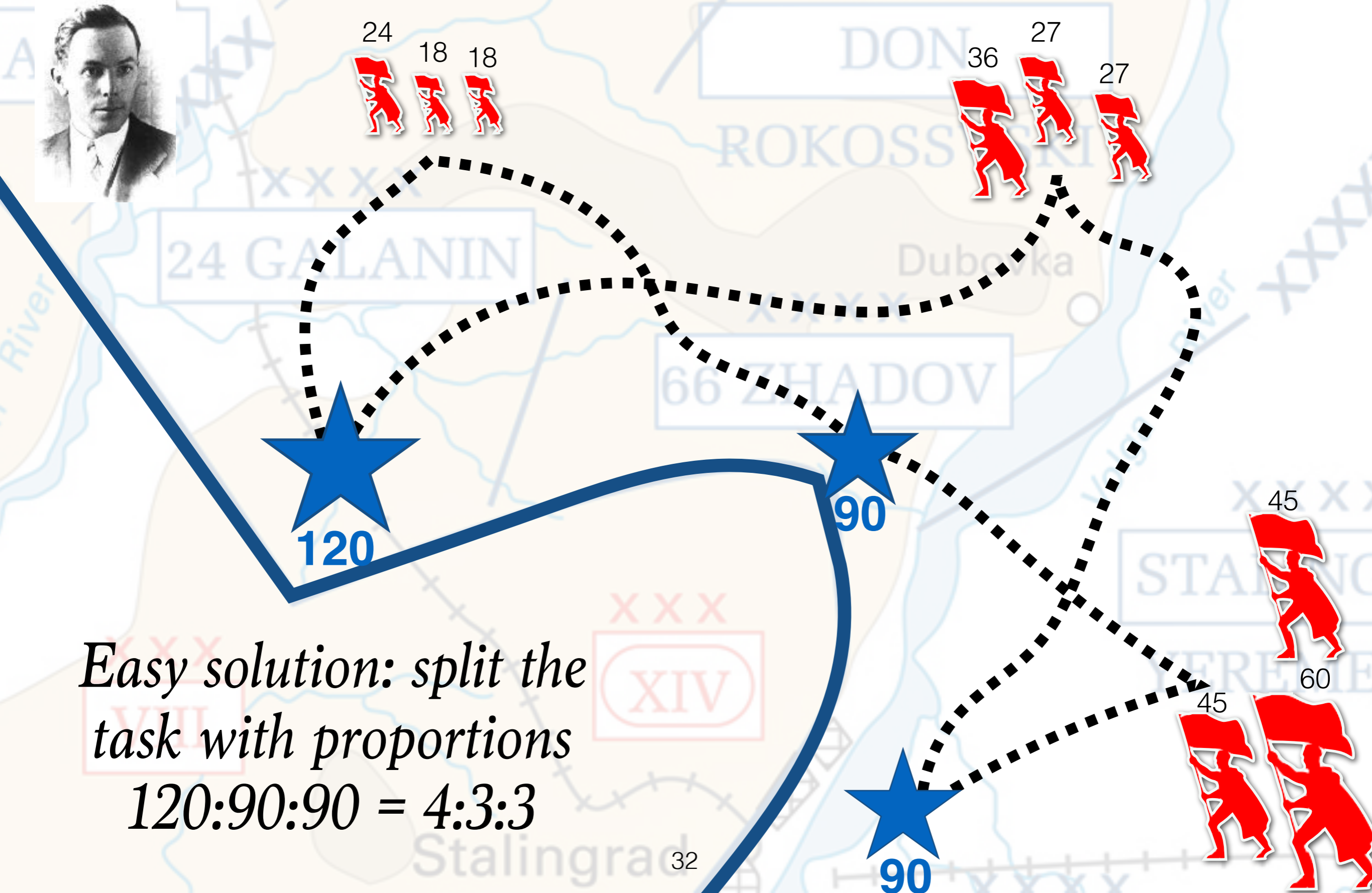
150

90

Kantorovich Problem

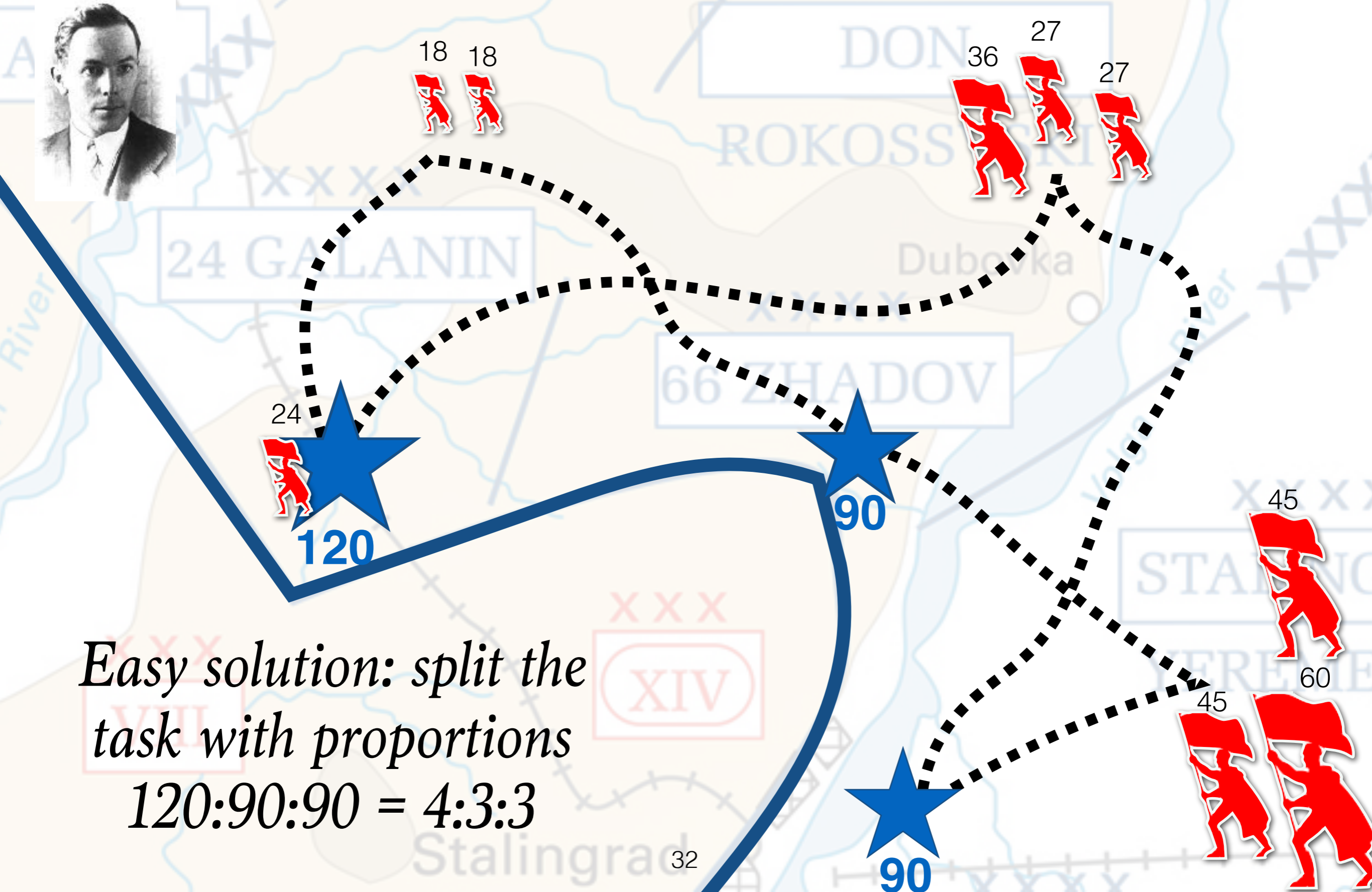


Kantorovich Problem



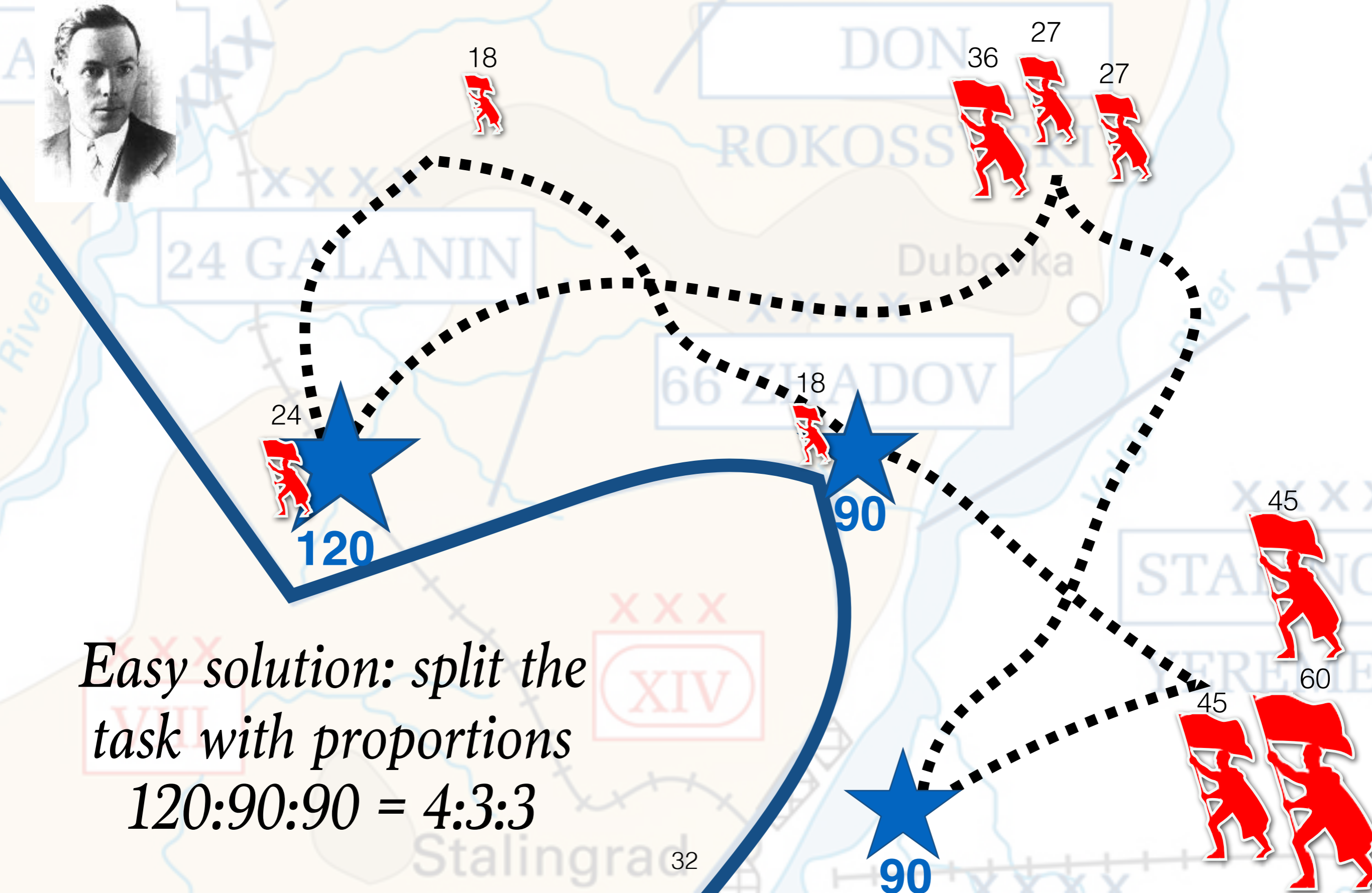
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



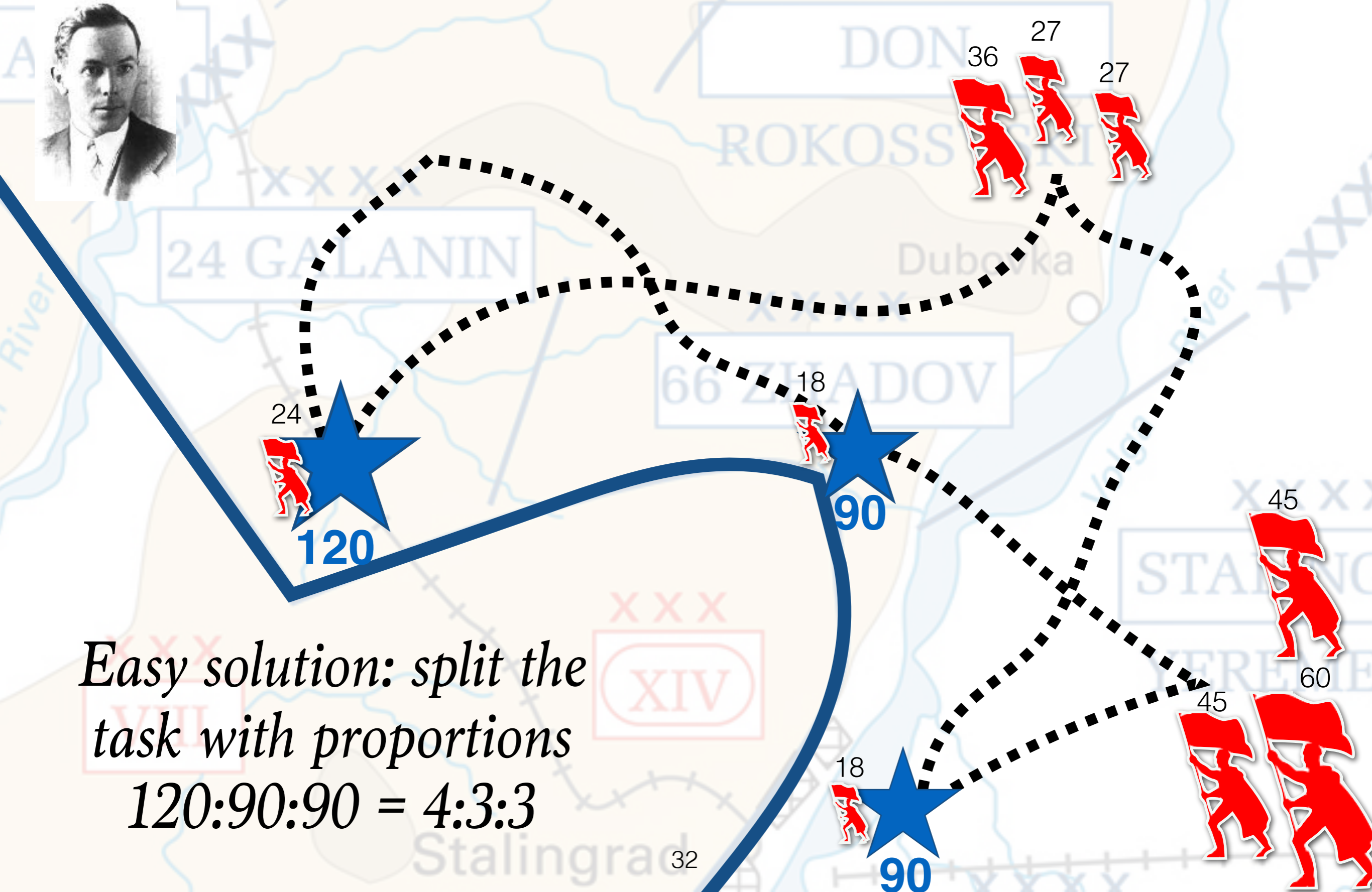
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

Kantorovich Problem



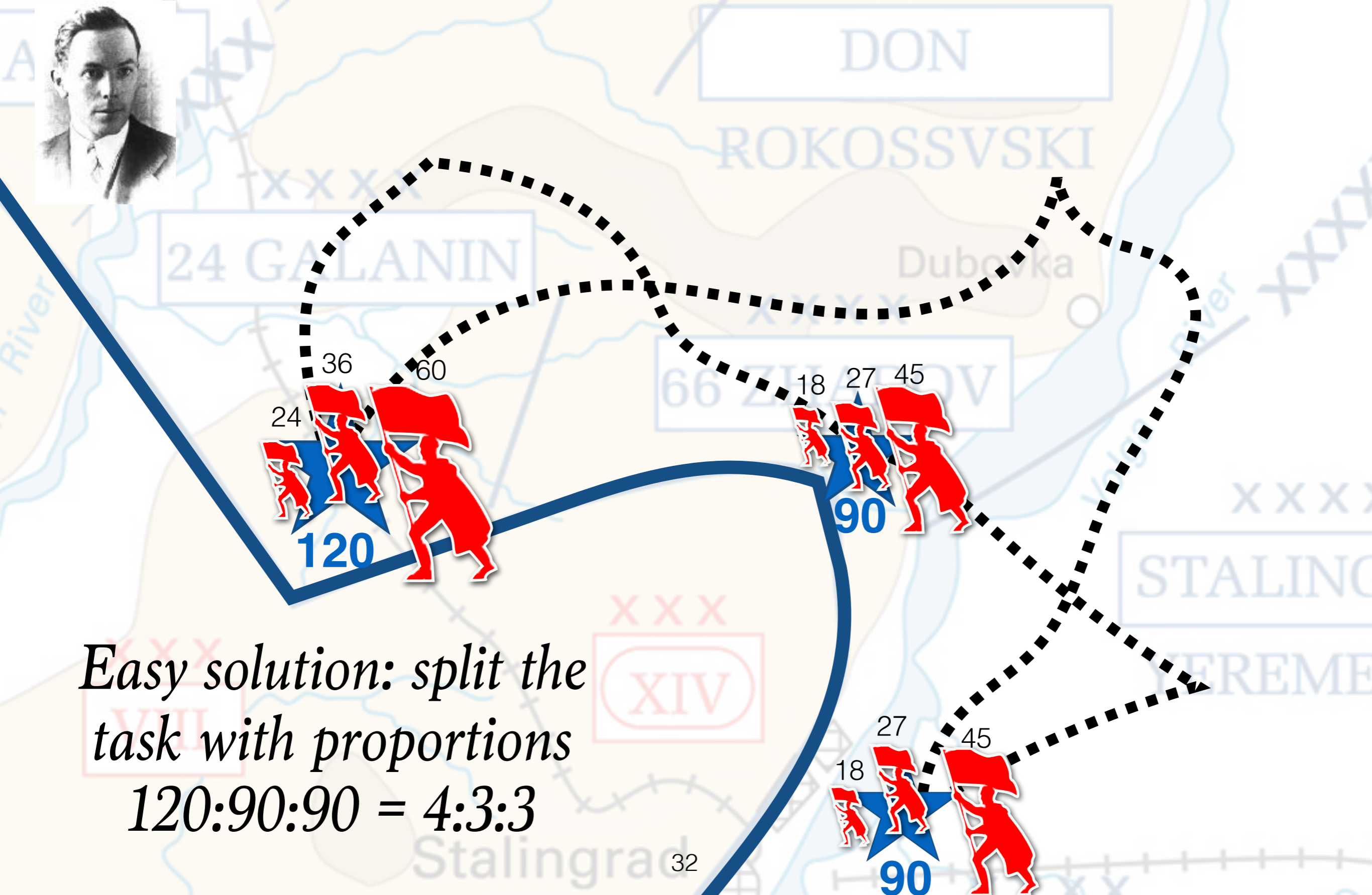
Easy solution: split the task with proportions
 $120:90:90 = 4:3:3$

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 many displacements...

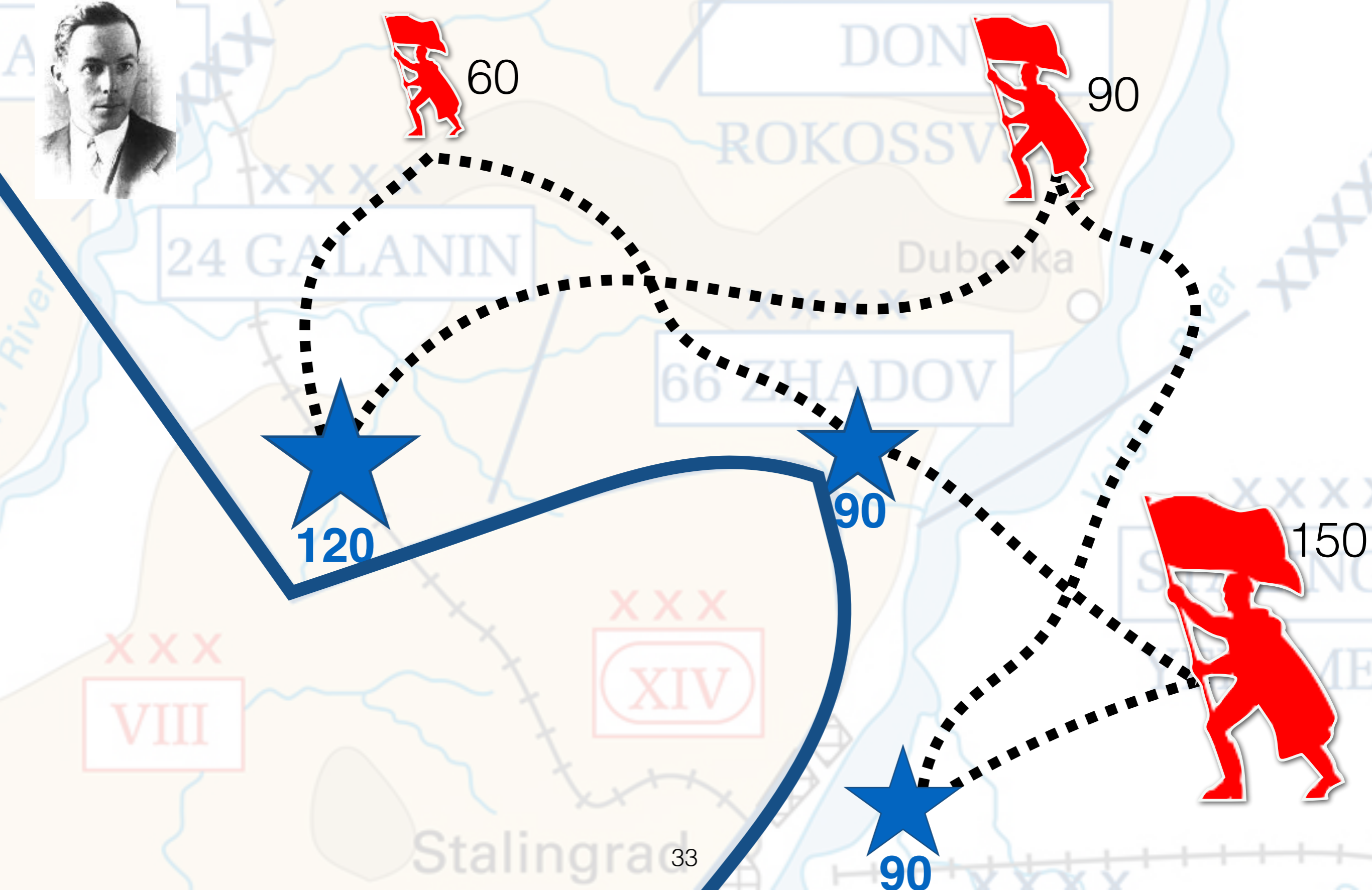
Can we find a cheaper alternative?

Easy solution: split the task with proportions

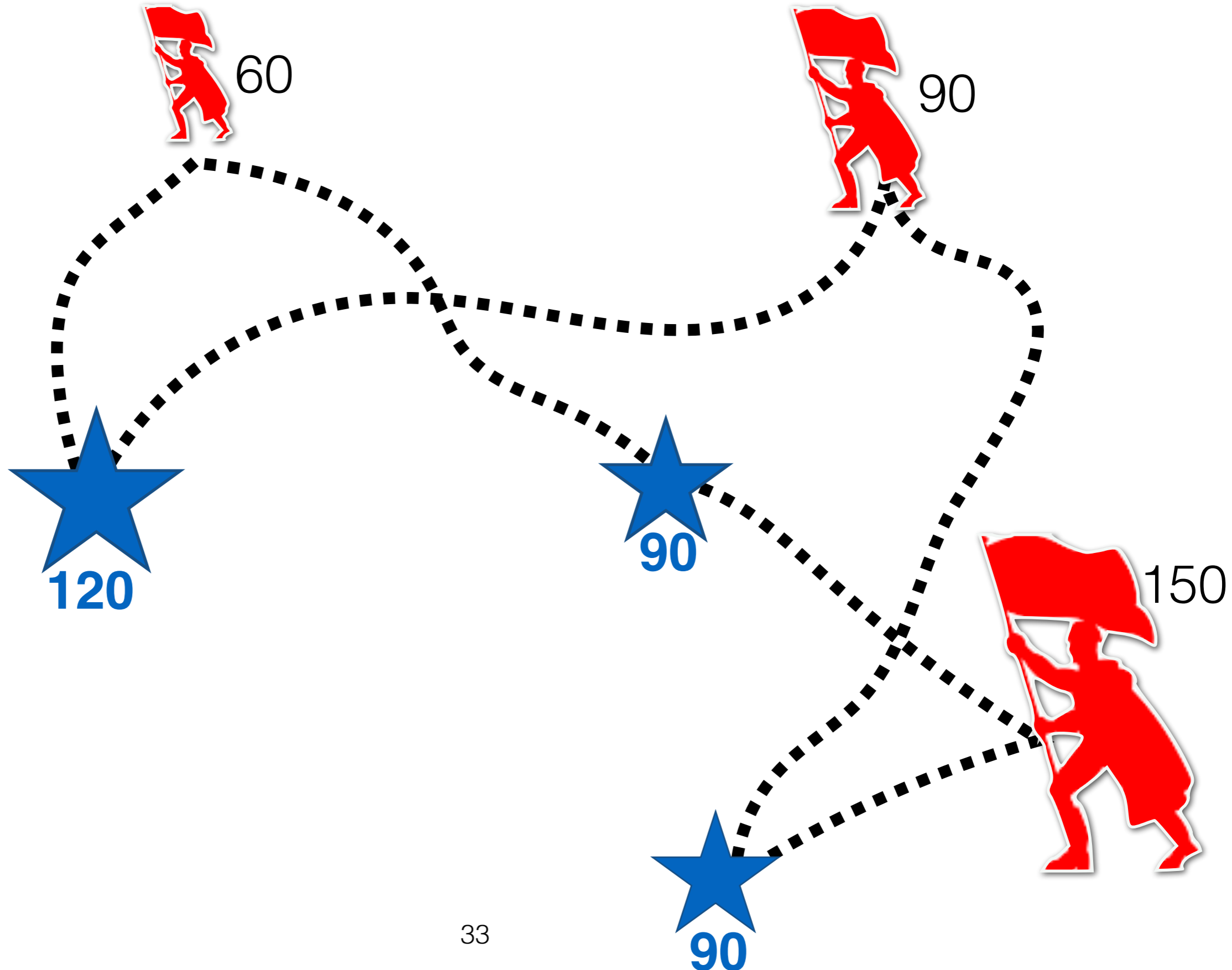
$$120:90:90 = 4:3:3$$



Kantorovich Problem



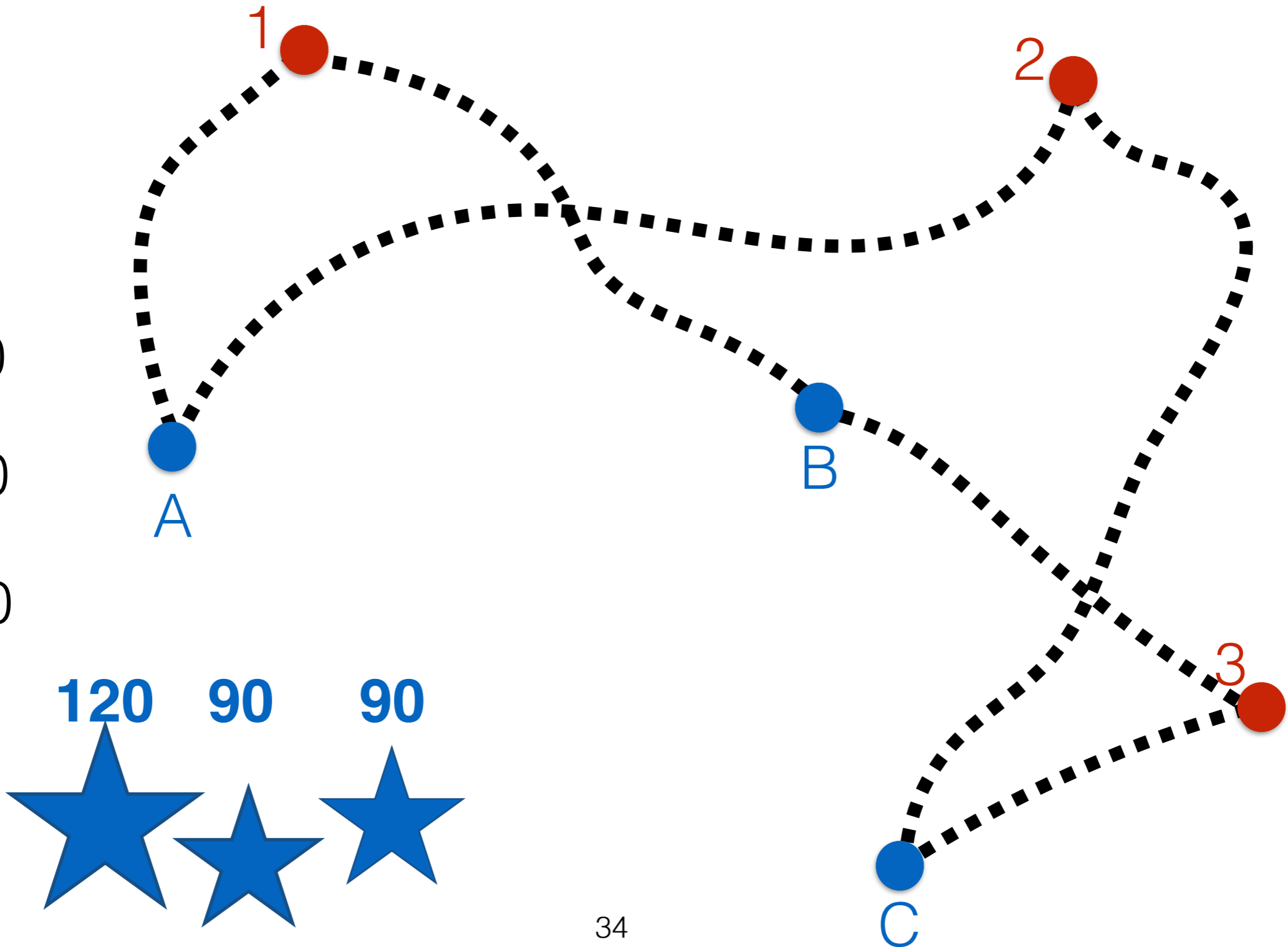
Kantorovich Problem



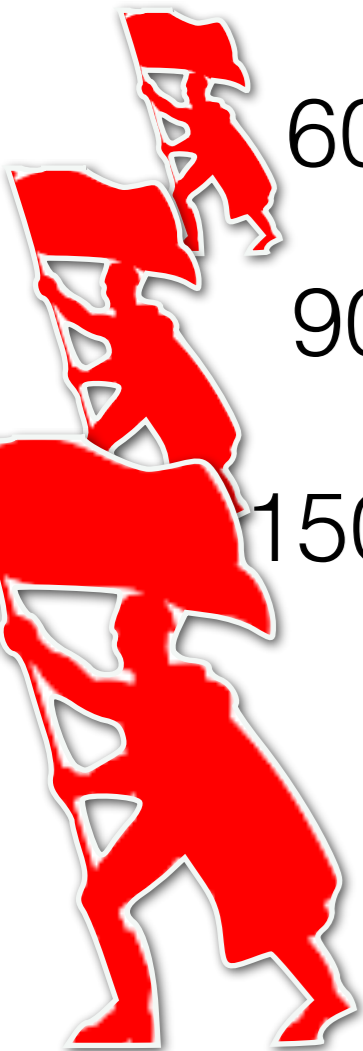
Kantorovich Problem



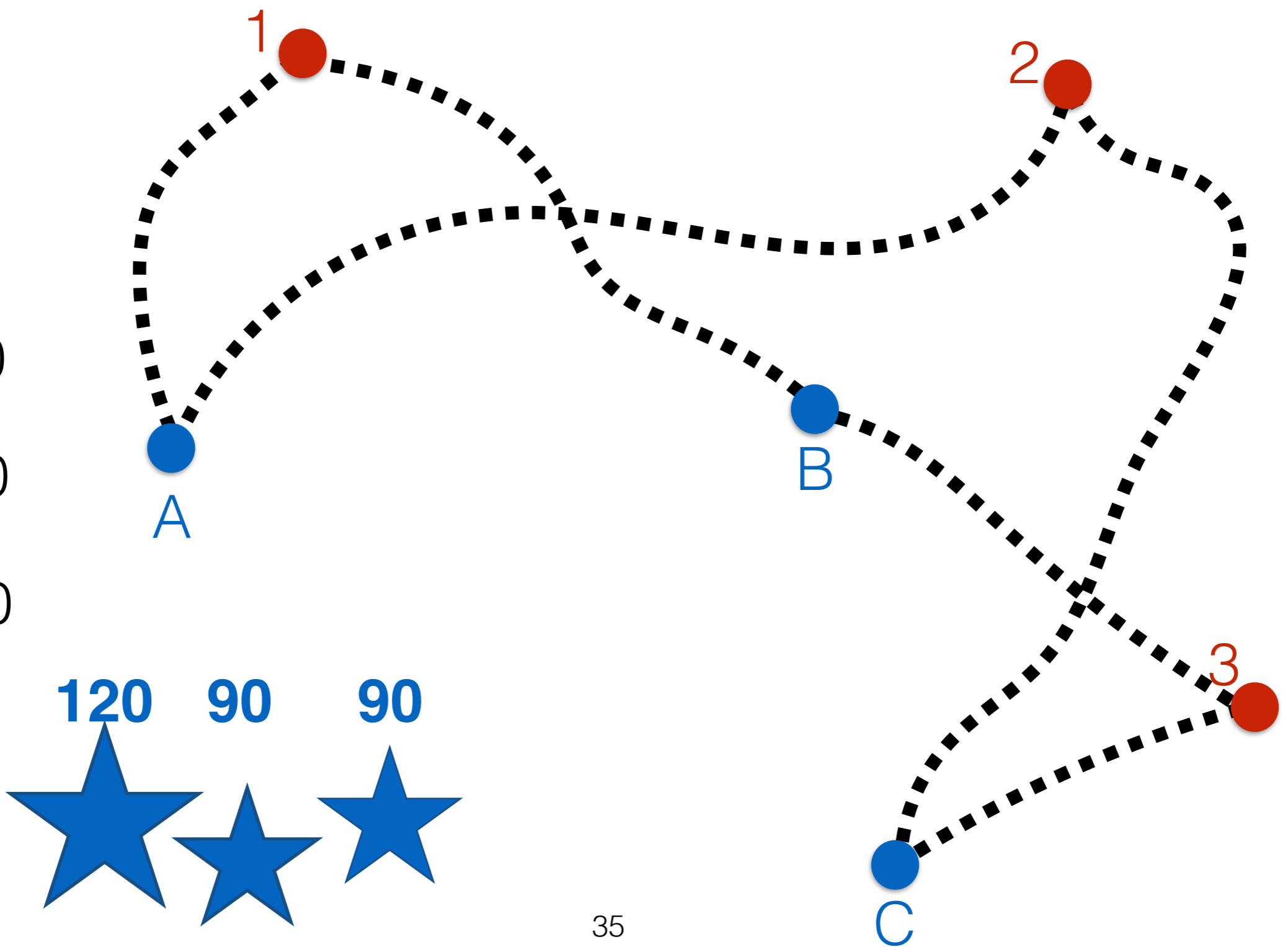
60
90
150



Kantorovich Problem



60
90
150



Kantorovich Problem



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

Below the table are three blue stars of decreasing size, corresponding to the values 120, 90, and 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

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

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

$$p_{ij} \geq 0$$

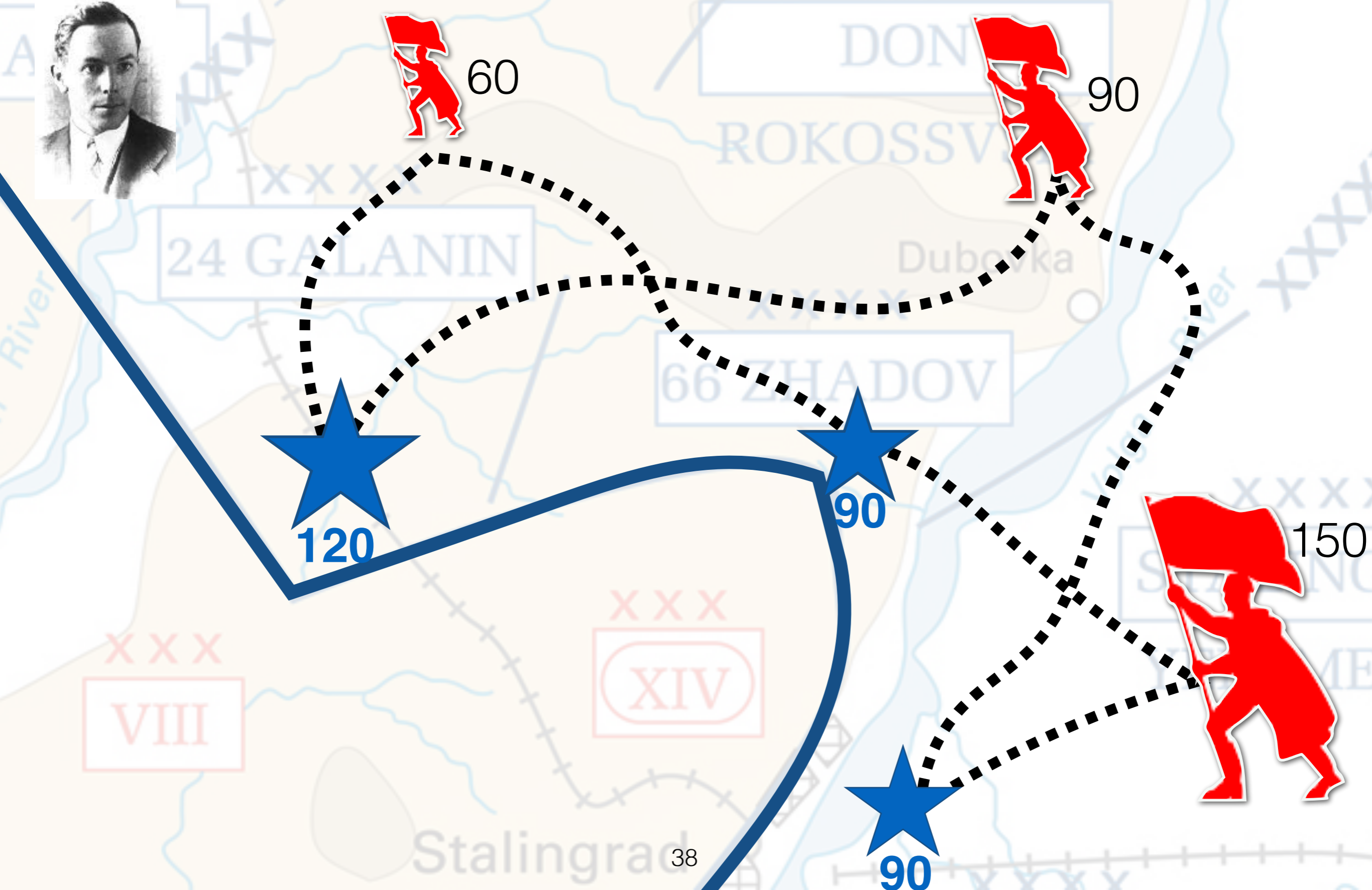
Cost function

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

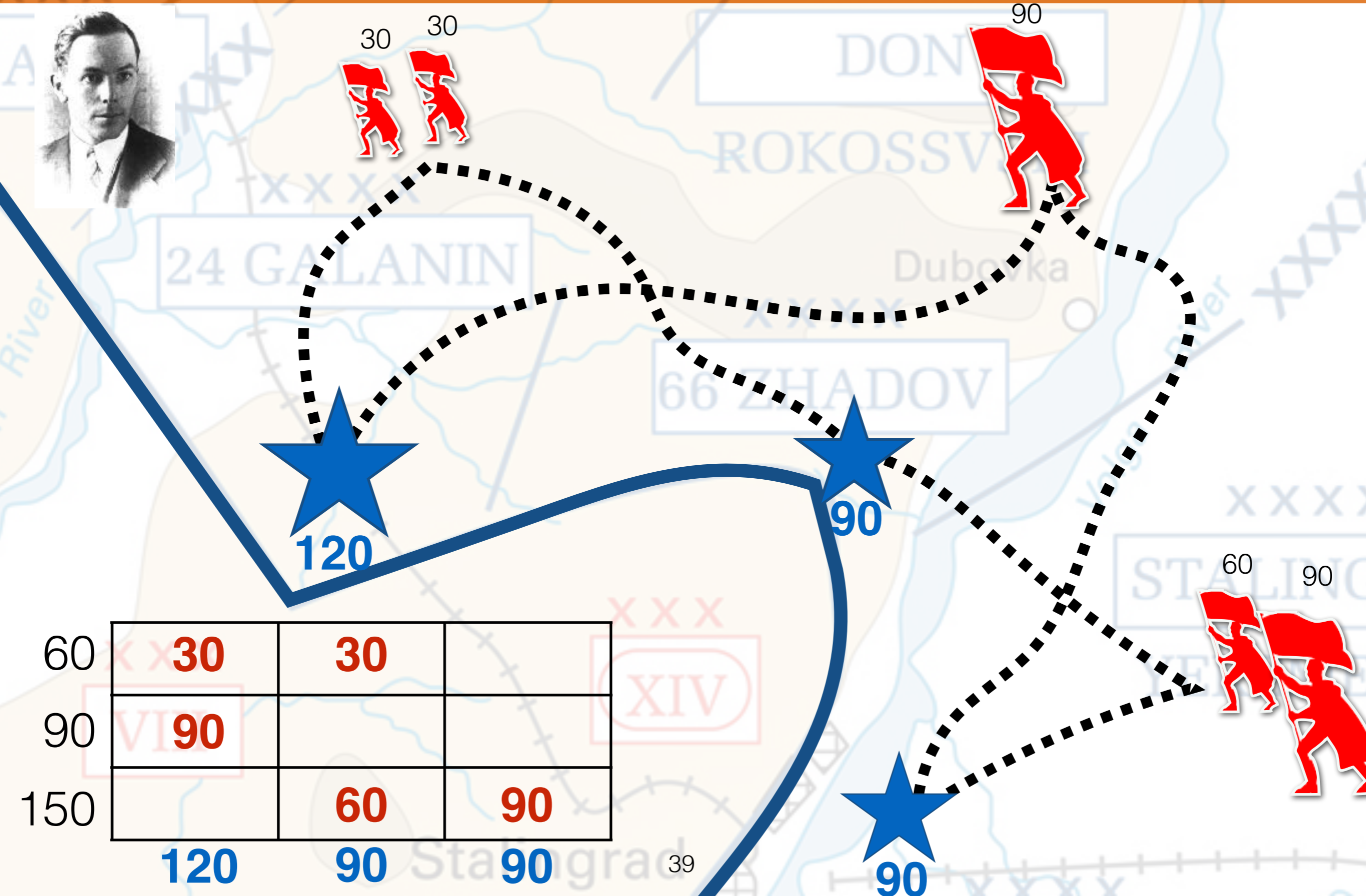
Problem

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

Kantorovich Problem

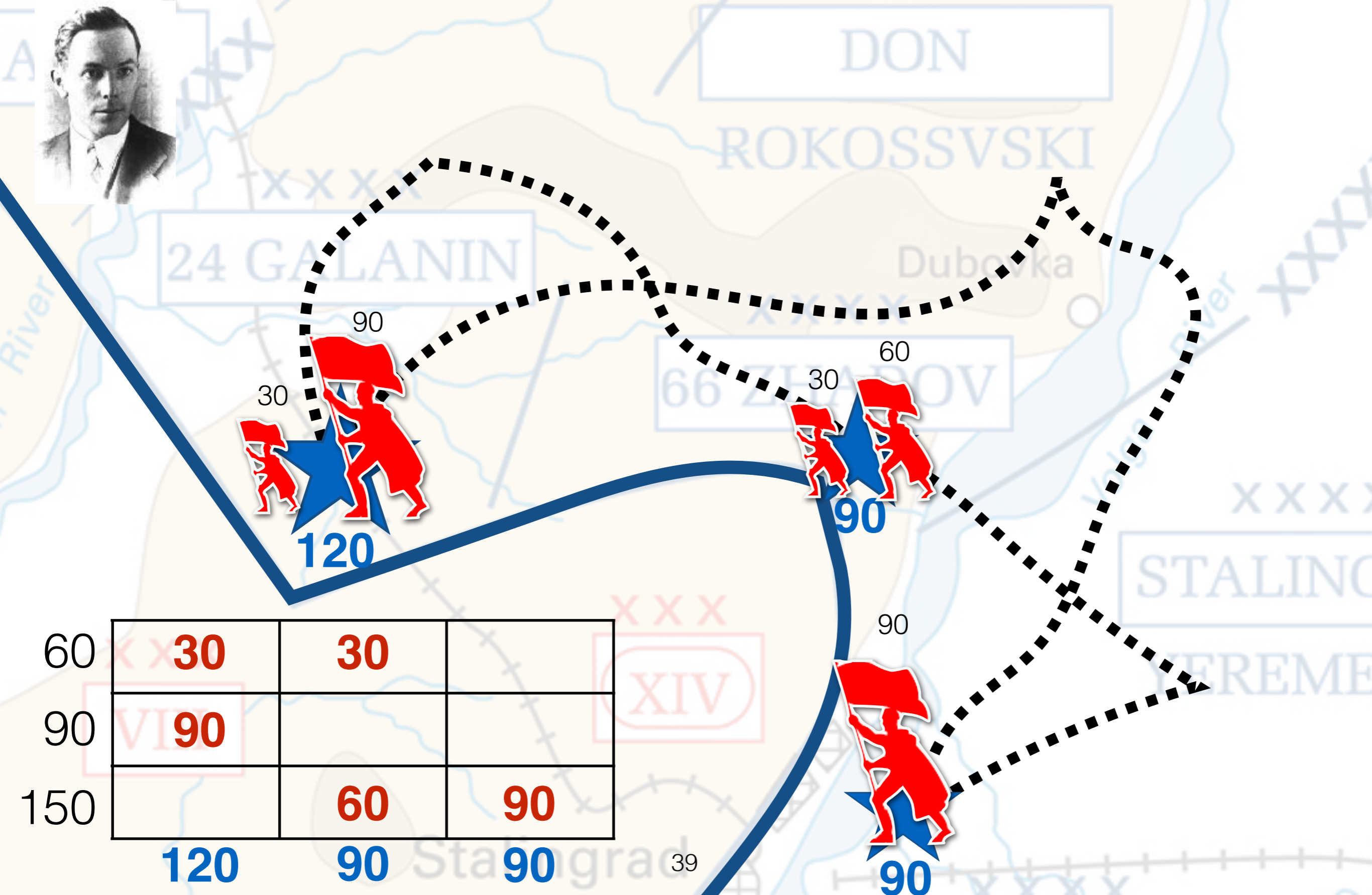


Kantorovich Problem



60	30	30	
90	90		
150		60	90
	120	90	90

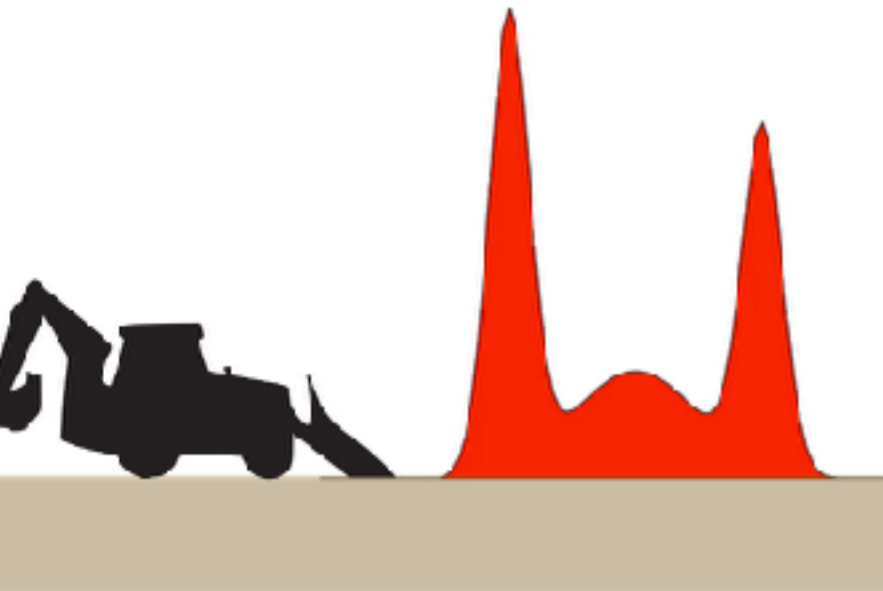
Kantorovich Problem



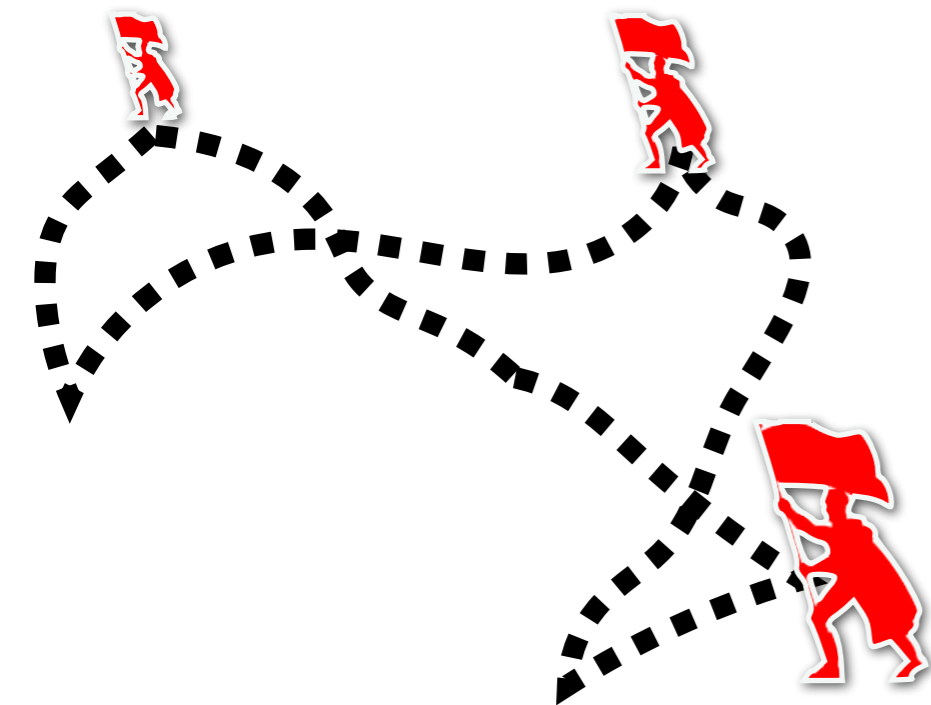
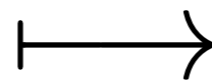
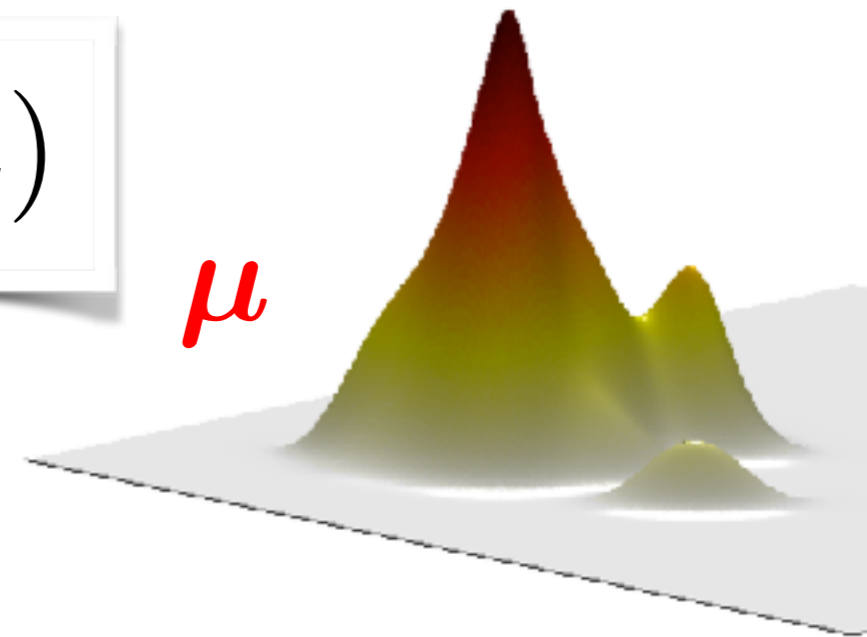
60	30	30	
90	90		
150		60	90
	120	90	90

Mathematical Formalism

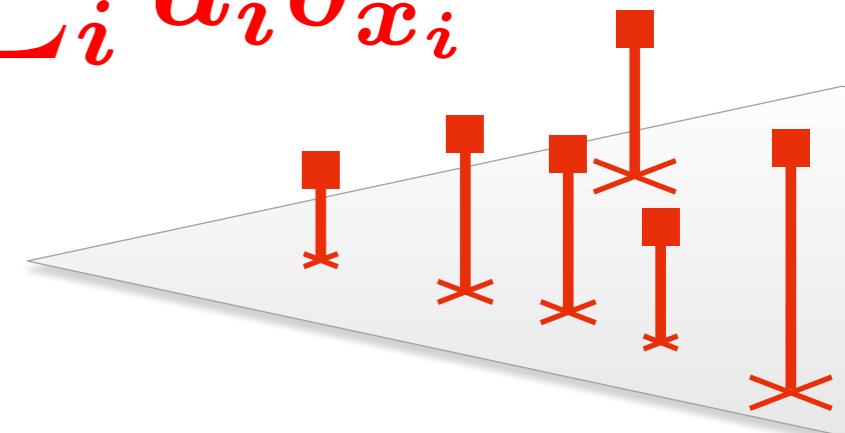
These problems involve discrete and continuous **probability measures** on a geometric space Ω



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



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

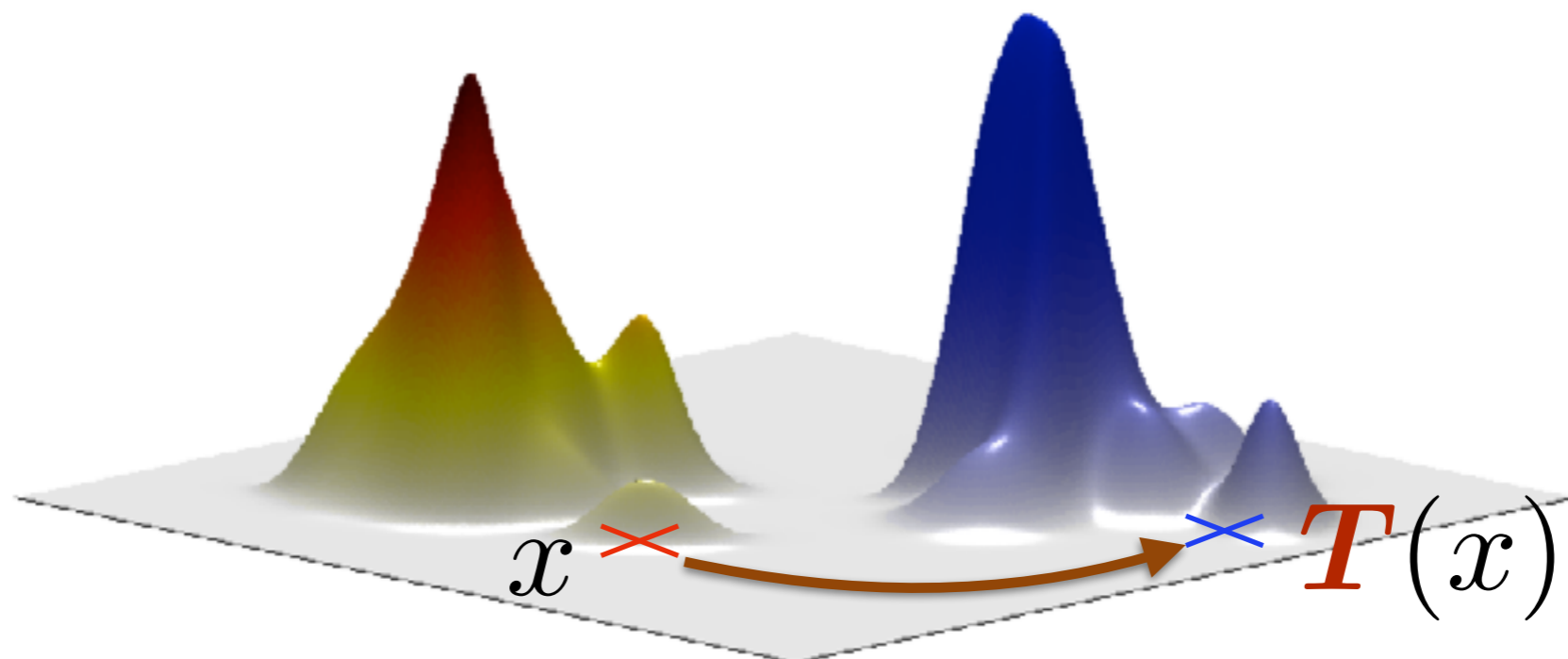


Monge Problem

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

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

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

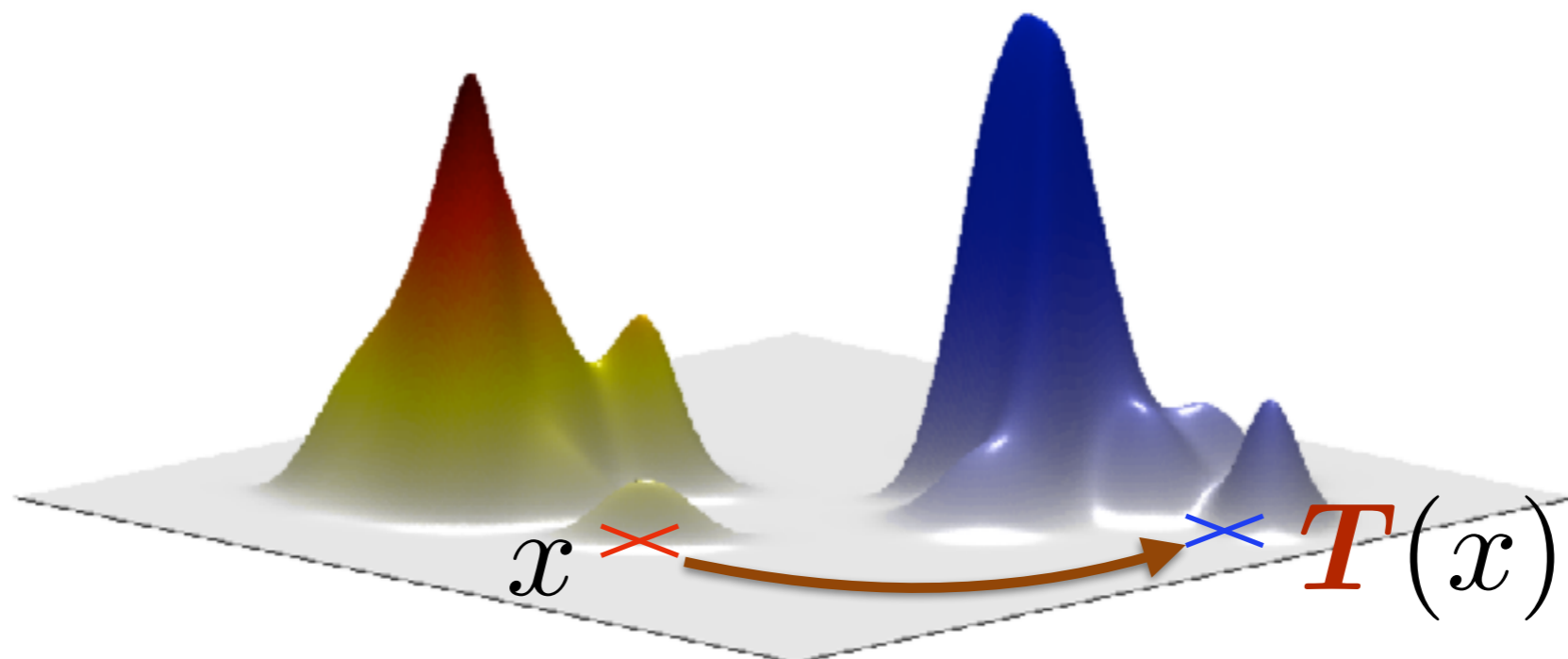


Monge Problem

Ω a measurable 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.

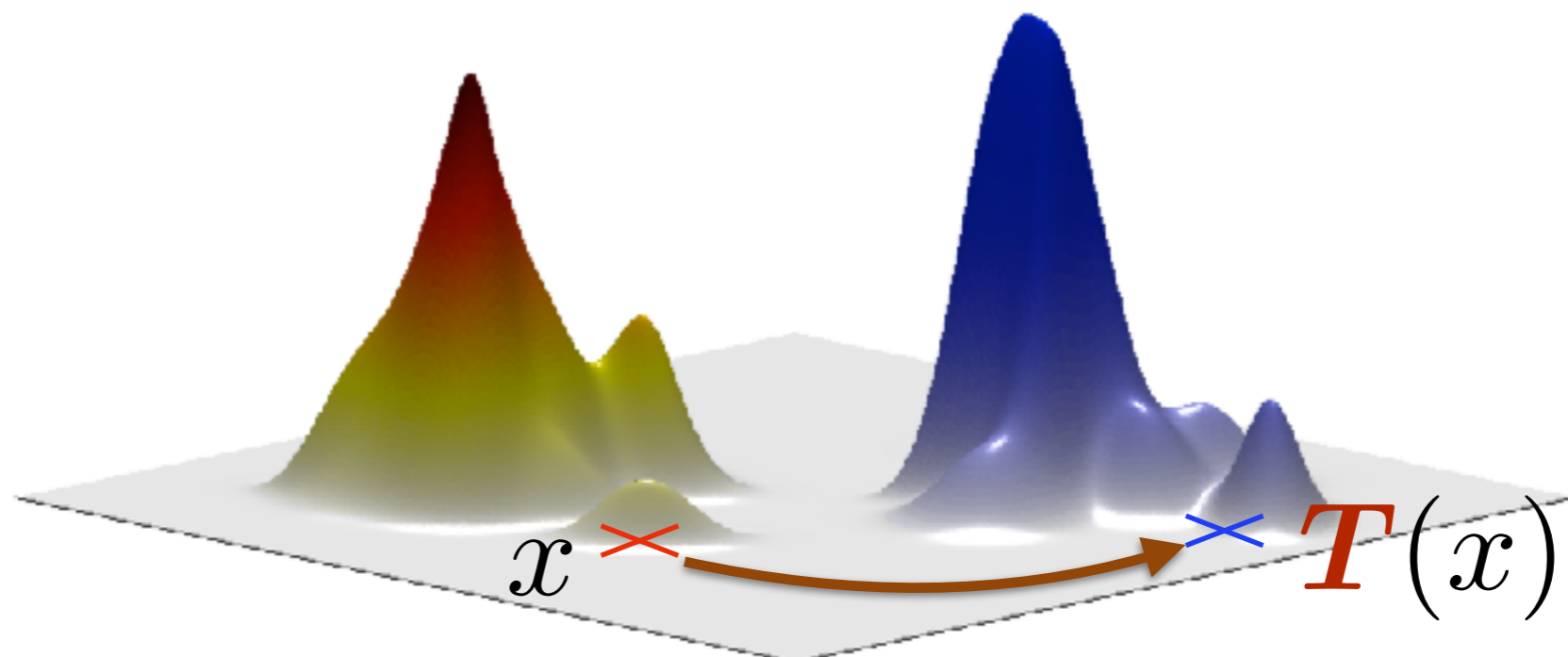


Monge Problem

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

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

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

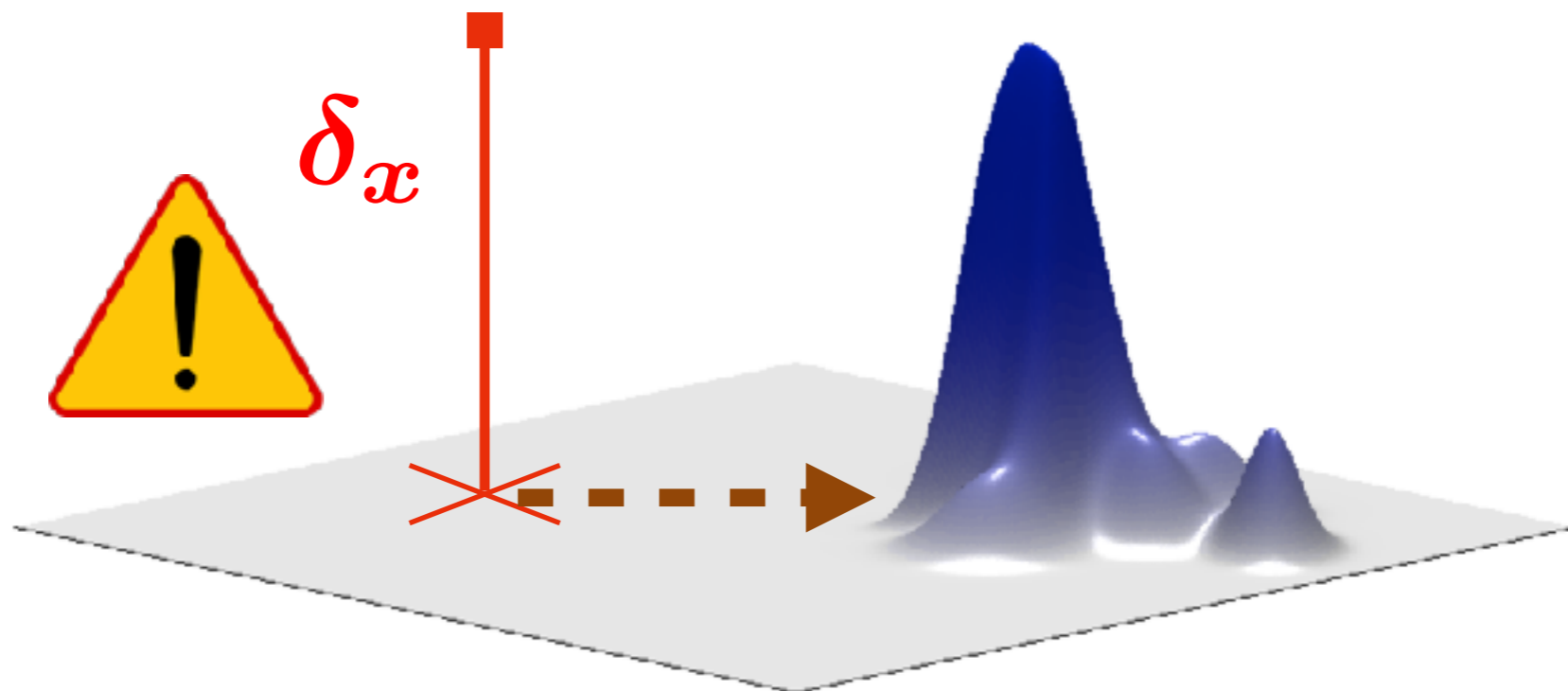


Monge Problem

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

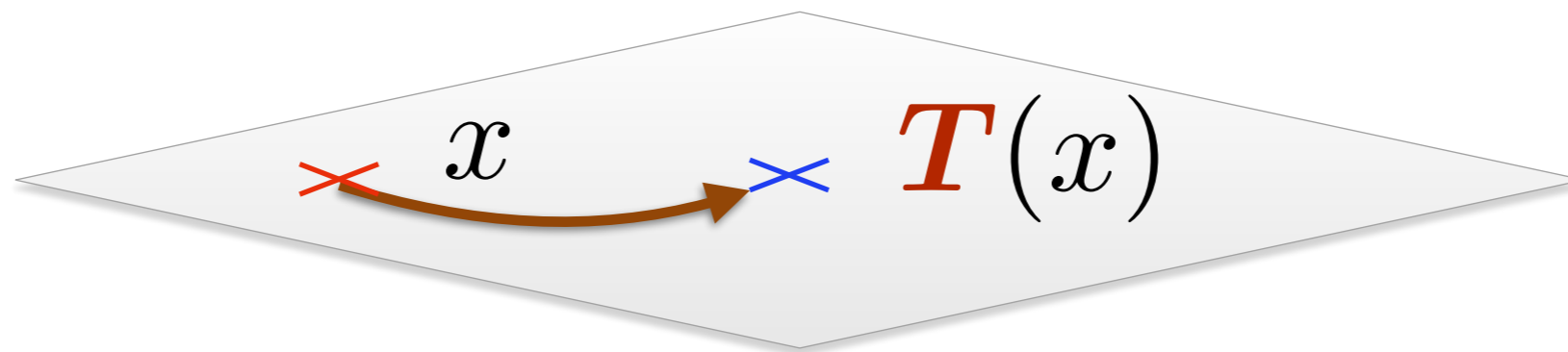
[Monge'81] problem: find a map $T : \Omega \rightarrow \Omega$

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



Kantorovich Relaxation

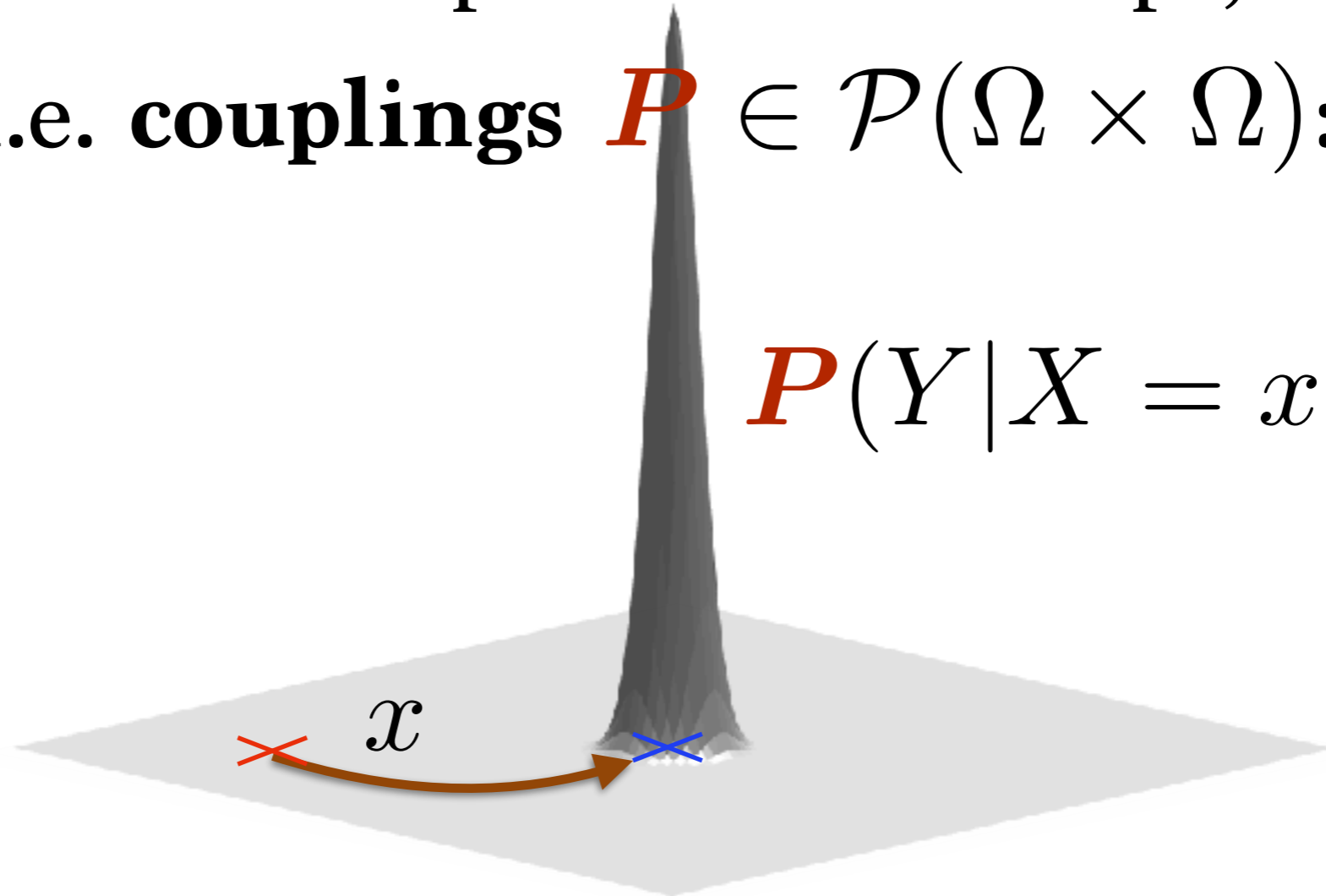
Instead of maps $T : \Omega \rightarrow \Omega$,
consider probabilistic maps,
i.e. couplings $P \in \mathcal{P}(\Omega \times \Omega)$:



Kantorovich Relaxation

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

$$P(Y|X = x)$$



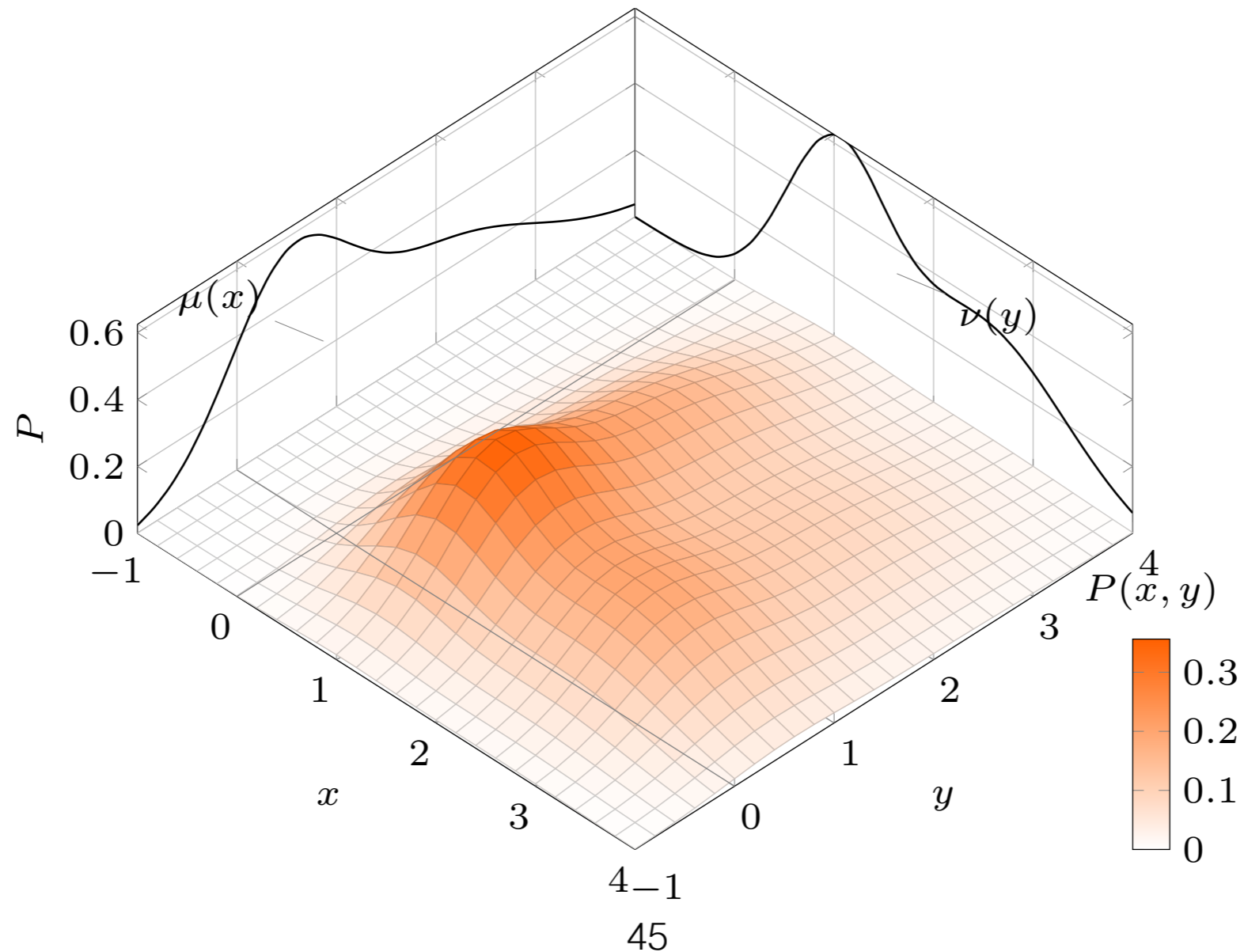
Kantorovich Relaxation

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

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

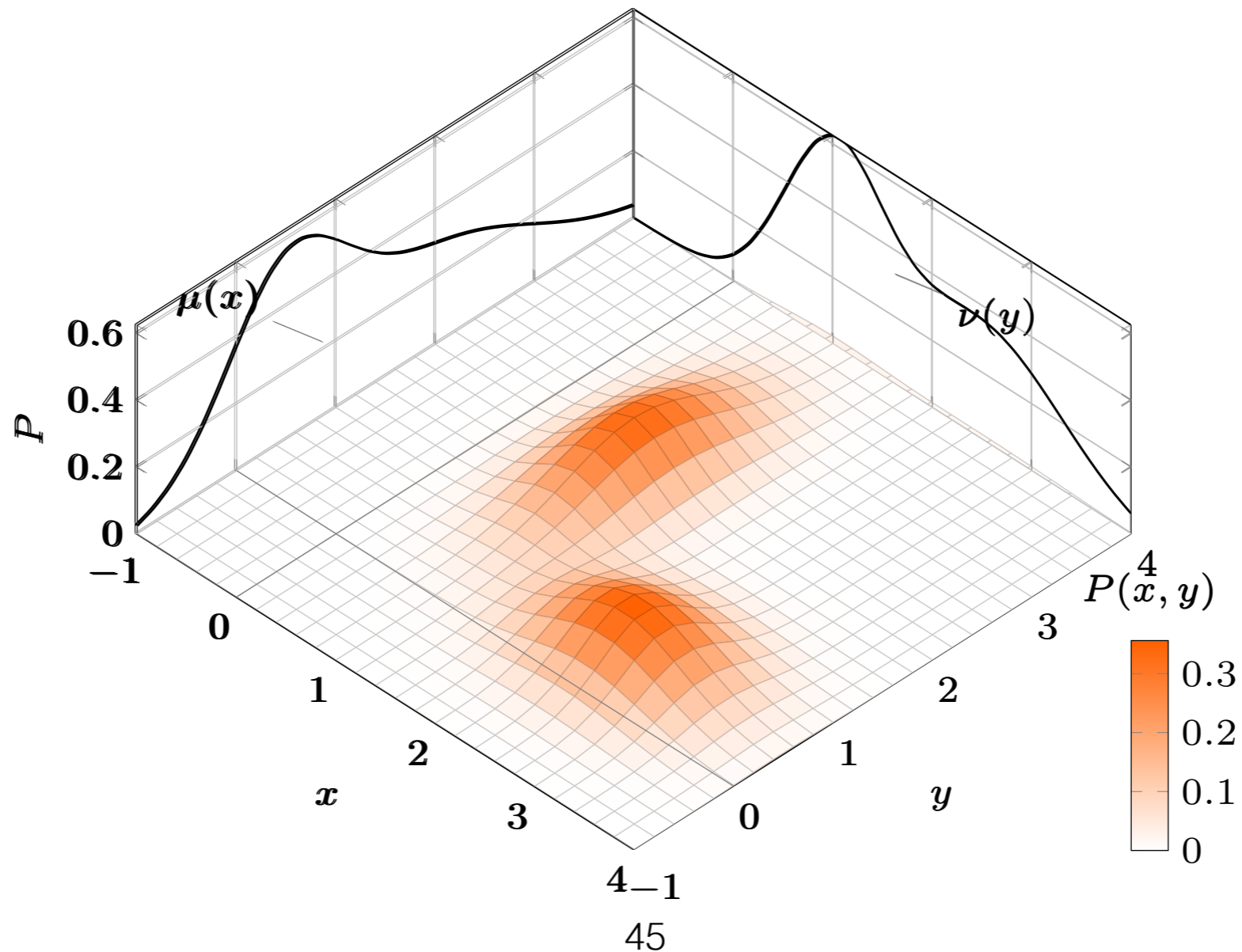
Kantorovich Relaxation

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



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

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

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

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

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

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

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

For two real-valued functions φ, ψ on Ω ,

$$(\varphi \oplus \psi)(x, y) \stackrel{\text{def}}{=} \varphi(x) + \psi(y)$$

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

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

DUAL

Deriving Kantorovich Duality

$$\begin{aligned} \iota_{\Pi}(\mathbf{P}) &= \sup_{\varphi, \psi} \left[\int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iint \varphi \oplus \psi d\mathbf{P} \right] \\ &= \begin{cases} 0 & \text{if } \mathbf{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu}), \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

Deriving Kantorovich Duality

$$\begin{aligned} \iota_{\Pi}(\mathbf{P}) &= \sup_{\varphi, \psi} \left[\int \varphi d\mu + \int \psi d\nu - \iint \varphi \oplus \psi d\mathbf{P} \right] \\ &= \begin{cases} 0 & \text{if } \mathbf{P} \in \Pi(\mu, \nu), \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

$$\inf_{\mathbf{P} \in \Pi(\mu, \nu)} \iint \mathbf{c} d\mathbf{P}$$

Deriving Kantorovich Duality

$$\begin{aligned} \iota_{\Pi}(\mathbf{P}) &= \sup_{\varphi, \psi} \left[\int \varphi d\mu + \int \psi d\nu - \iint \varphi \oplus \psi d\mathbf{P} \right] \\ &= \begin{cases} 0 & \text{if } \mathbf{P} \in \Pi(\mu, \nu), \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

$$\inf_{\mathbf{P} \in \Pi(\mu, \nu)} \iint \mathbf{c} d\mathbf{P}$$

$$\inf_{\mathbf{P} \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} d\mathbf{P} + \boxed{\iota_{\Pi}(\mathbf{P})}$$

Deriving Kantorovich Duality

$$\begin{aligned} \iota_{\Pi}(\mathbf{P}) &= \sup_{\varphi, \psi} \left[\int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iint \varphi \oplus \psi d\mathbf{P} \right] \\ &= \begin{cases} 0 & \text{if } \mathbf{P} \in \Pi(\boldsymbol{\mu}, \boldsymbol{\nu}), \\ +\infty & \text{otherwise.} \end{cases} \end{aligned}$$

$$\inf_{\mathbf{P} \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} d\mathbf{P} + \iota_{\Pi}(\mathbf{P})$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint c dP + \iota_{\Pi}(P)$$

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint c dP + \sup_{\varphi, \psi} \int \varphi d\mu + \int \psi d\nu - \iint \varphi \oplus \psi dP$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint c dP + \iota_{\Pi}(P)$$

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \sup_{\varphi, \psi} \iint c dP + \int \varphi d\mu + \int \psi d\nu - \iint \varphi \oplus \psi dP$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} dP + \iota_{\Pi}(P)$$

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \sup_{\varphi, \psi} \iint \mathbf{c} dP - \iint \varphi \oplus \psi dP + \int \varphi d\mu + \int \psi d\nu$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} dP + \iota_{\Pi}(P)$$

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \sup_{\varphi, \psi} \iint (\mathbf{c} - \varphi \oplus \psi) dP + \int \varphi d\mu + \int \psi d\nu$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} dP + \iota_{\Pi}(P)$$

$$\sup_{\varphi, \psi} \inf_{P \in \mathcal{P}_+(\Omega^2)} \iint (\mathbf{c} - \varphi \oplus \psi) dP + \int \varphi d\mu + \int \psi d\nu$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint c dP + \iota_{\Pi}(P)$$

$$\sup_{\varphi, \psi} \inf_{P \in \mathcal{P}_+(\Omega^2)} \iint (c - \varphi \oplus \psi) dP + \int \varphi d\mu + \int \psi d\nu$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint \mathbf{c} dP + \iota_{\Pi}(P)$$

$$\sup_{\varphi, \psi} \inf_{P \in \mathcal{P}_+(\Omega^2)} \iint (\mathbf{c} - \varphi \oplus \psi) dP + \int \varphi d\mu + \int \psi d\nu$$

$$\inf_{P \in \mathcal{P}_+(\Omega)} \iint (\mathbf{c} - \varphi \oplus \psi) dP = \begin{cases} 0 & \text{if } \mathbf{c} - \varphi \oplus \psi \geq 0. \\ -\infty & \text{otherwise} \end{cases}$$

Deriving Kantorovich Duality

$$\inf_{P \in \mathcal{P}_+(\Omega^2)} \iint c dP + \iota_{\Pi}(P)$$

$$\sup_{\varphi, \psi} \inf_{P \in \mathcal{P}_+(\Omega^2)} \iint (c - \varphi \oplus \psi) dP + \int \varphi d\mu + \int \psi d\nu$$

$$\inf_{P \in \mathcal{P}_+(\Omega)} \iint (c - \varphi \oplus \psi) dP = \begin{cases} 0 & \text{if } c - \varphi \oplus \psi \geq 0. \\ -\infty & \text{otherwise} \end{cases}$$

$$\sup_{\varphi \oplus \psi \leq c} \int \varphi d\mu + \int \psi d\nu.$$

DUAL

Wasserstein Distances

Let $p \geq 1$. Let $\mathbf{c}(x, y) := \mathbf{D}^p(x, y)$, a metric.

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

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

Wasserstein Distances

Let $p \geq 1$. Let $\mathbf{c}(x, y) := \mathbf{D}^p(x, y)$, a metric.

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

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

Kantorovich Duality

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

DUAL

- Kantorovich Duality is **interesting** from a computational perspective: easier to store 2 functions than a whole coupling.
- D transforms: go from **two** to **one** dual potential.

D transforms

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

DUAL

Imagine we choose a φ . Can we find a good ψ ?

D transforms

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

DUAL

Imagine we choose a φ . Can we find a good ψ ?

We need that ψ satisfies for *all* $\boldsymbol{x}, \boldsymbol{y}$

$$\varphi(\boldsymbol{x}) + \psi(\boldsymbol{y}) \leq D^p(\boldsymbol{x}, \boldsymbol{y})$$

D transforms

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

DUAL

Imagine we choose a φ . Can we find a good ψ ?

We need that ψ satisfies for *all* $\boldsymbol{x}, \boldsymbol{y}$

$$\varphi(\boldsymbol{x}) + \psi(\boldsymbol{y}) \leq D^p(\boldsymbol{x}, \boldsymbol{y})$$

$$\psi(\boldsymbol{y}) \leq D^p(\boldsymbol{x}, \boldsymbol{y}) - \varphi(\boldsymbol{x})$$

D transforms

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

DUAL

Imagine we choose a φ . Can we find a good ψ ?

We need that ψ satisfies for *all* $\boldsymbol{x}, \boldsymbol{y}$

$$\varphi(\boldsymbol{x}) + \psi(\boldsymbol{y}) \leq D^p(\boldsymbol{x}, \boldsymbol{y})$$

$$\psi(\boldsymbol{y}) \leq D^p(\boldsymbol{x}, \boldsymbol{y}) - \varphi(\boldsymbol{x})$$

$$\psi(\boldsymbol{y}) \leq \inf_{\boldsymbol{x}} D^p(\boldsymbol{x}, \boldsymbol{y}) - \varphi(\boldsymbol{x})$$

D transforms

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

DUAL

For given φ , cannot get a better ψ than

$$\bar{\varphi}(y) \stackrel{\text{def}}{=} \inf_x D^p(x, y) - \varphi(x).$$

D transforms

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

DUAL

For given φ , cannot get a better ψ than

$$\bar{\varphi}(y) \stackrel{\text{def}}{=} \inf_x D^p(x, y) - \varphi(x).$$

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

SEMI-DUAL

D transforms

$$\overline{\varphi}(\mathbf{y}) \stackrel{\text{def}}{=} \inf_{\mathbf{x}} D^p(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x}).$$

$$\overline{\psi}(\mathbf{x}) = \inf_{\mathbf{y}} D^p(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y}).$$

D transforms

$$\overline{\varphi}(y) \stackrel{\text{def}}{=} \inf_x D^p(x, y) - \varphi(x).$$

$$\overline{\psi}(x) = \inf_y D^p(x, y) - \psi(y).$$

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

D transforms

$$\overline{\varphi}(y) \stackrel{\text{def}}{=} \inf_x D^p(x, y) - \varphi(x).$$

$$\overline{\psi}(x) = \inf_y D^p(x, y) - \psi(y).$$

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



D transforms

$$\overline{\varphi}(\mathbf{y}) \stackrel{\text{def}}{=} \inf_{\mathbf{x}} D^p(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x}).$$

$$\overline{\psi}(\mathbf{x}) = \inf_{\mathbf{y}} D^p(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y}).$$

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

For all φ , we have $\overline{\overline{\varphi}} = \overline{\varphi}$

D transforms

$$\overline{\varphi}(\mathbf{y}) \stackrel{\text{def}}{=} \inf_{\mathbf{x}} D^p(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x}).$$

$$\overline{\psi}(\mathbf{x}) = \inf_{\mathbf{y}} D^p(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y}).$$

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

For all φ , we have $\overline{\overline{\varphi}} = \overline{\varphi}$

φ is D^p -concave if $\exists \phi : \varphi = \overline{\phi}$

φ is D^p -concave $\Rightarrow \overline{\overline{\varphi}} = \varphi$

D transforms

$$\overline{\varphi}(\mathbf{y}) \stackrel{\text{def}}{=} \inf_{\mathbf{x}} D^p(\mathbf{x}, \mathbf{y}) - \varphi(\mathbf{x}).$$

$$\overline{\psi}(\mathbf{x}) = \inf_{\mathbf{y}} D^p(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y}).$$

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

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi \text{ is } D^p\text{-concave}} \int \varphi d\boldsymbol{\mu} + \int \overline{\varphi} d\boldsymbol{\nu}.$$

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

$$\bar{\varphi}_x(y) - \bar{\varphi}_x(y') = D(x, y) - D(x, y') \leq D(y, y')$$

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) = \inf_x \bar{\varphi}_x(y)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) - \bar{\varphi}(x) \leq D(x, y)$

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) = \inf_x \bar{\varphi}_x(y)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) - \bar{\varphi}(x) \leq D(x, y)$

$\Rightarrow -\bar{\varphi}(x) \leq D(x, y) - \bar{\varphi}(y)$

$\Rightarrow -\bar{\varphi}(x) \leq \inf_y D(x, y) - \bar{\varphi}(y)$

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) = \inf_x \bar{\varphi}_x(y)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) - \bar{\varphi}(x) \leq D(x, y)$

$\Rightarrow -\bar{\varphi}(x) \leq D(x, y) - \bar{\varphi}(y)$

$\Rightarrow -\bar{\varphi}(x) \leq \inf_y D(x, y) - \bar{\varphi}(y)$

$\Rightarrow -\bar{\varphi}(x) \leq \inf_y D(x, y) - \bar{\varphi}(y) \leq -\bar{\varphi}(x)$

D transforms, W_1

Prop. If $c = D$, namely $p = 1$, then
 φ is D -concave $\Leftrightarrow \bar{\varphi} = -\varphi$, φ is 1-Lipschitz

For given x , $\bar{\varphi}_x(y) \stackrel{\text{def}}{=} D(x, y) - \varphi(x)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) = \inf_x \bar{\varphi}_x(y)$ is 1-Lipschitz.

$\Rightarrow \bar{\varphi}(y) - \bar{\varphi}(x) \leq D(x, y)$

$\Rightarrow -\bar{\varphi}(x) \leq D(x, y) - \bar{\varphi}(y)$

$\Rightarrow -\bar{\varphi}(x) \leq \inf_y D(x, y) - \bar{\varphi}(y)$

$\Rightarrow -\bar{\varphi}(x) \leq \inf_y D(x, y) - \bar{\varphi}(y) \leq -\bar{\varphi}(x)$

$\Rightarrow -\bar{\varphi}(x) \leq \bar{\varphi}(x) \leq -\bar{\varphi}(x)$ and $\bar{\varphi}(x) = -\varphi(x)$

D transforms, W_1

$$W_1(\mu, \nu) = \sup_{\varphi \text{ is } D\text{-concave}} \int \varphi d\mu + \int \overline{\varphi} d\nu.$$

SEMI-DUAL

Prop. If $c = D$, then

φ is D -concave $\Leftrightarrow \overline{\varphi} = -\varphi$, φ is 1-Lipschitz

$$W_1(\mu, \nu) = \sup_{\varphi \text{ 1-Lipschitz}} \int \varphi (d\mu - d\nu).$$

W1

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

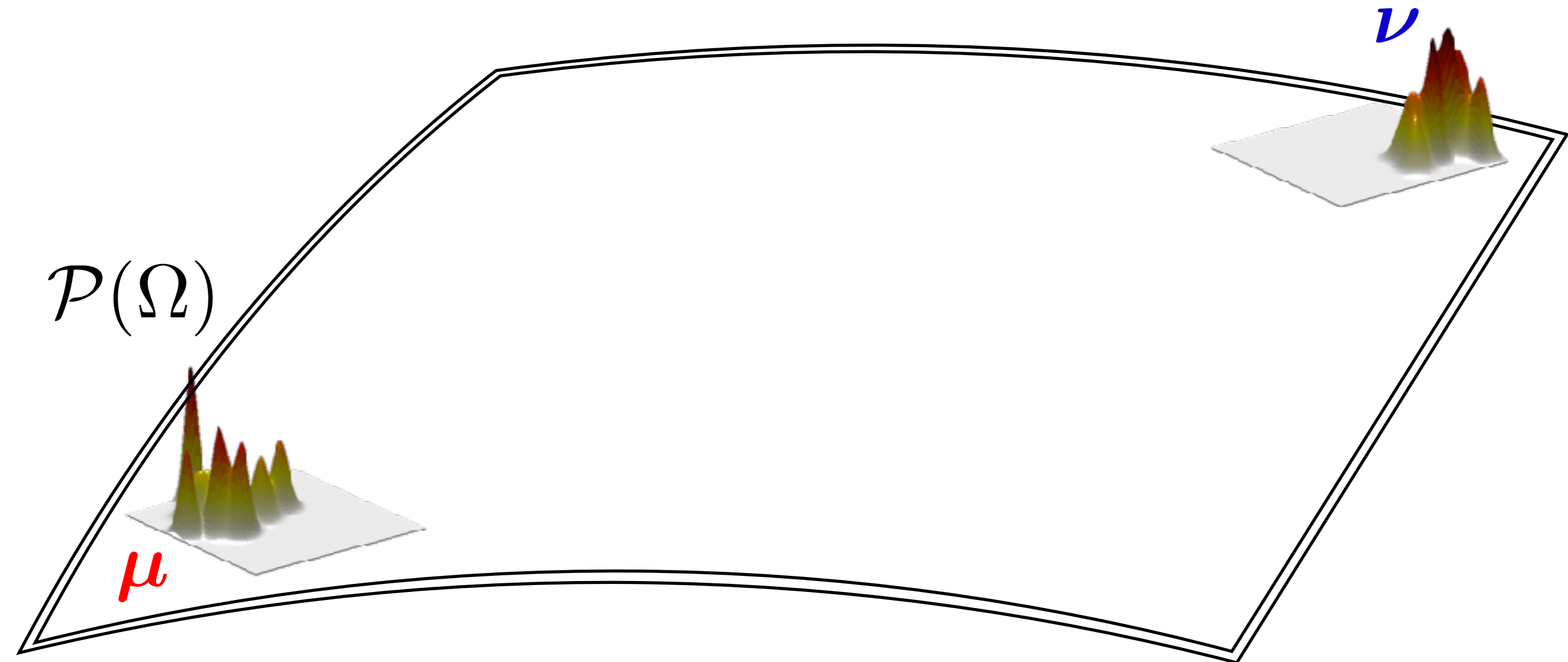
$$P^* := (\text{Id}, T^*)\# \mu \in \Pi(\mu, \nu)$$

is also *optimal* for the Kantorovich problem.

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

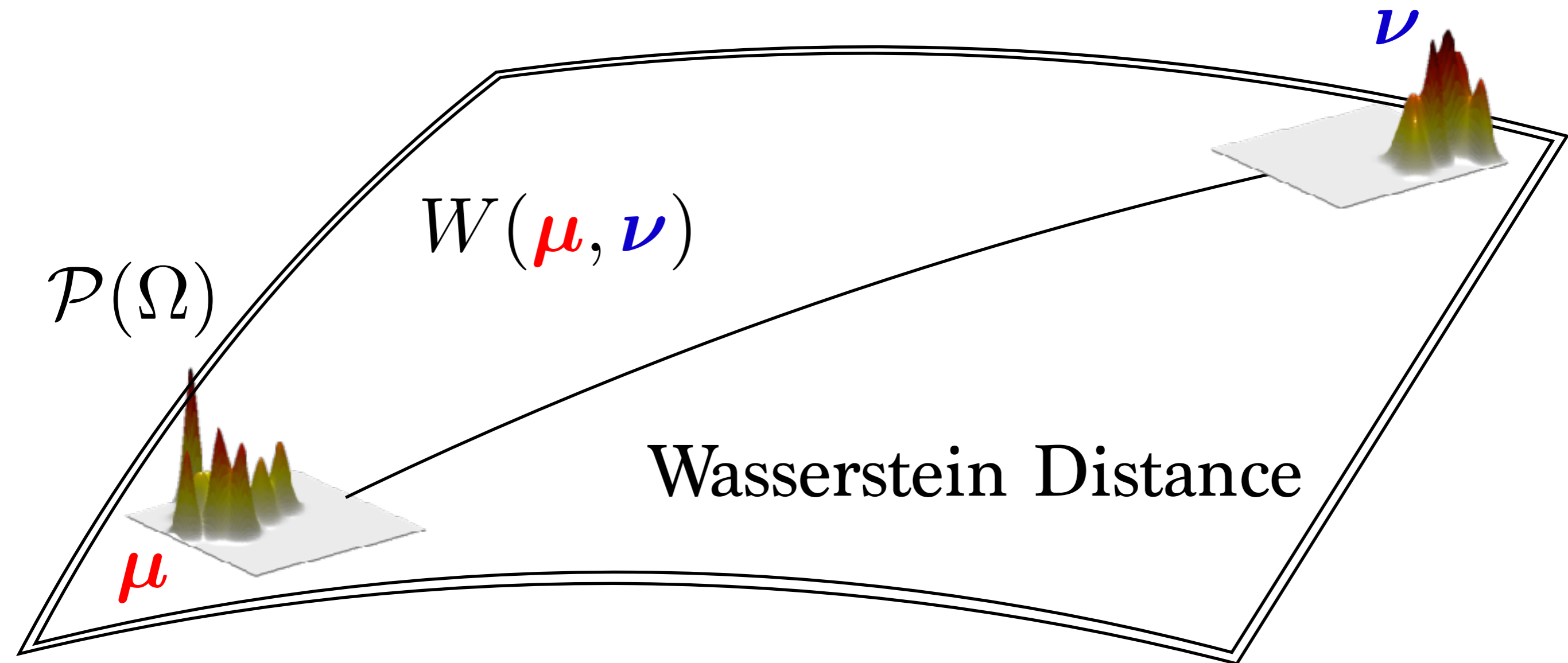
Optimal Transport Geometry

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



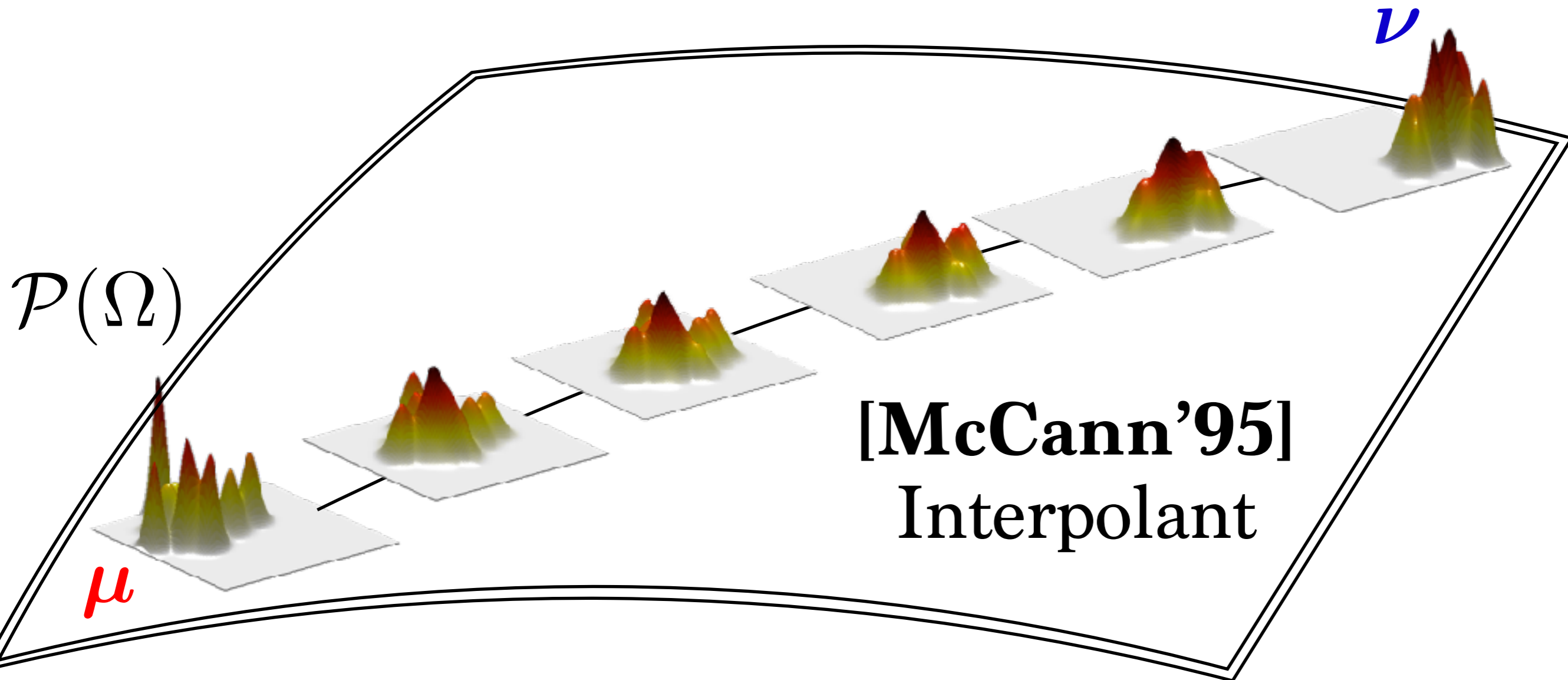
Optimal Transport Geometry

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



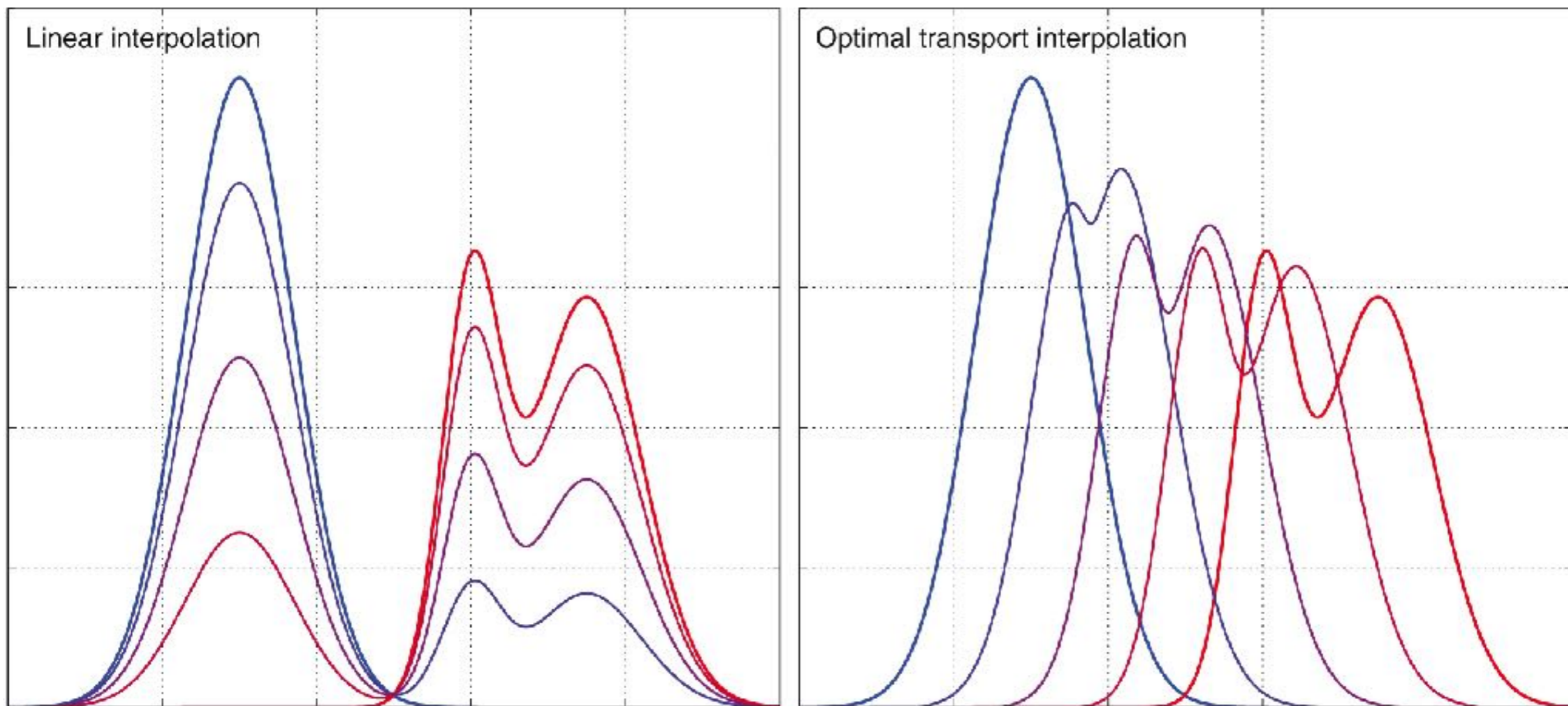
Optimal Transport Geometry

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



Optimal Transport Geometry

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



Optimal Transport Geometry

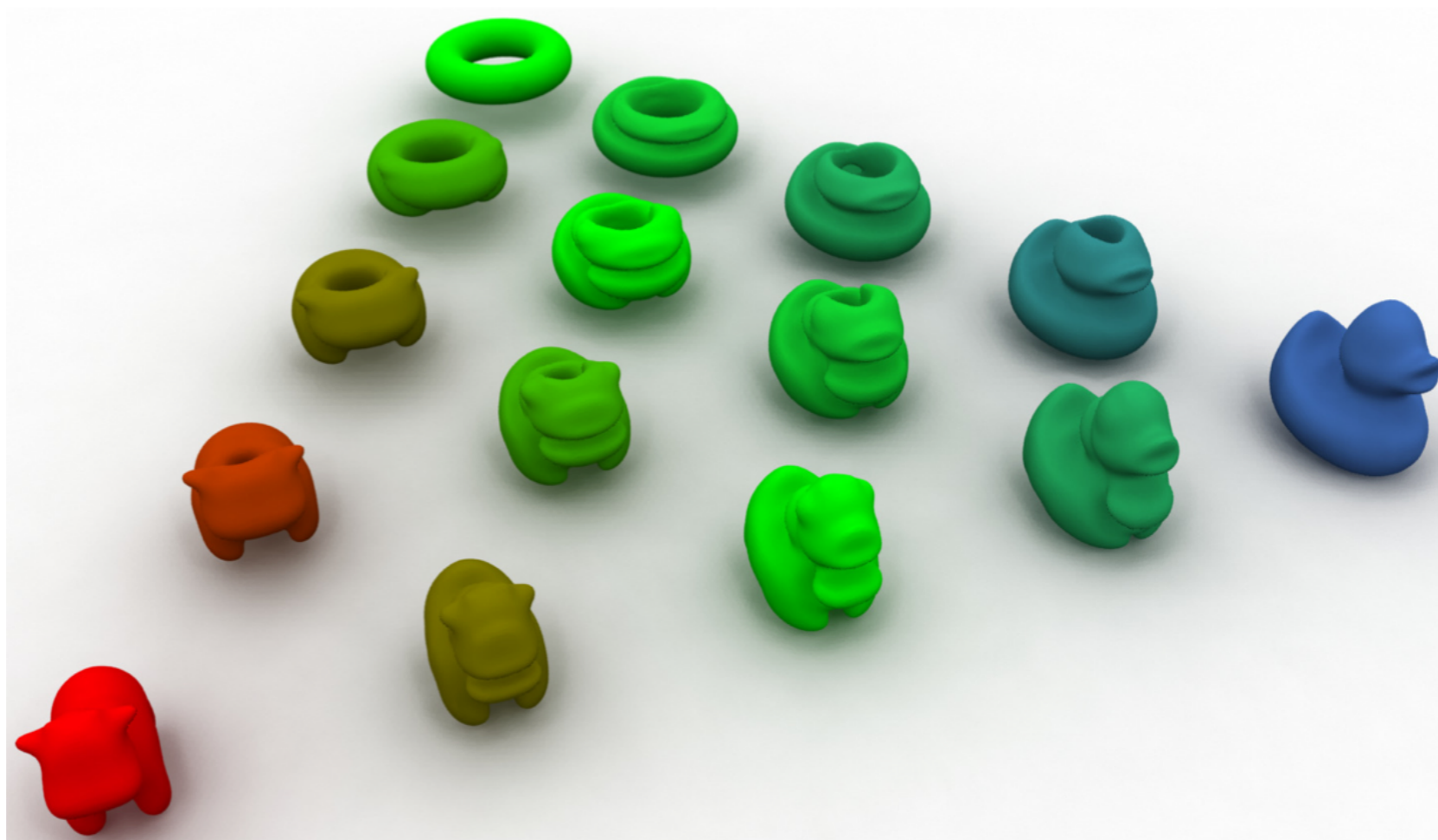
Very different geometry than standard information divergences (KL , Euclidean)



[SDPC.'15]

Optimal Transport Geometry

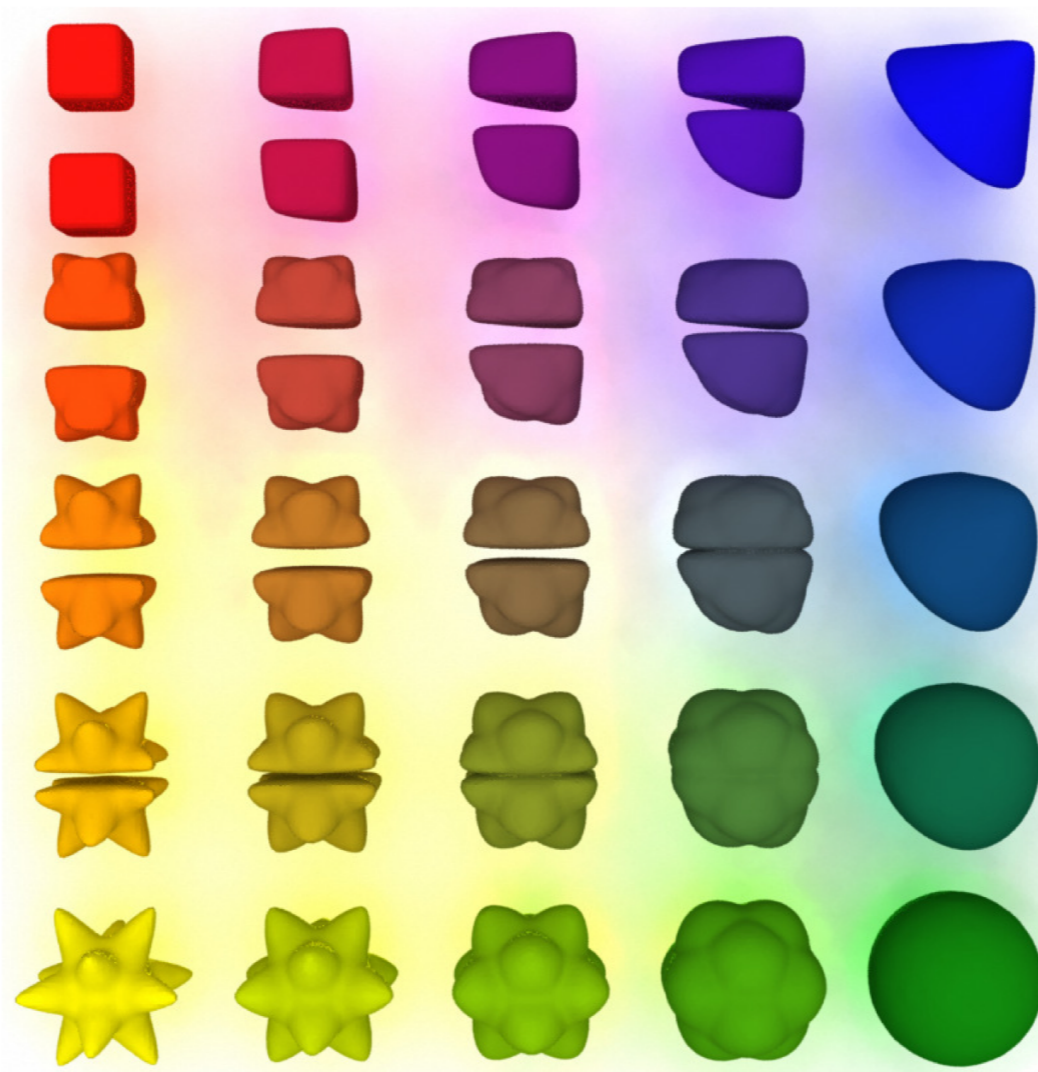
Very different geometry than standard information divergences (KL , Euclidean)



[SDPC.'15]

Optimal Transport Geometry

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



[SDPC.'15]

Computational OT

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

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

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

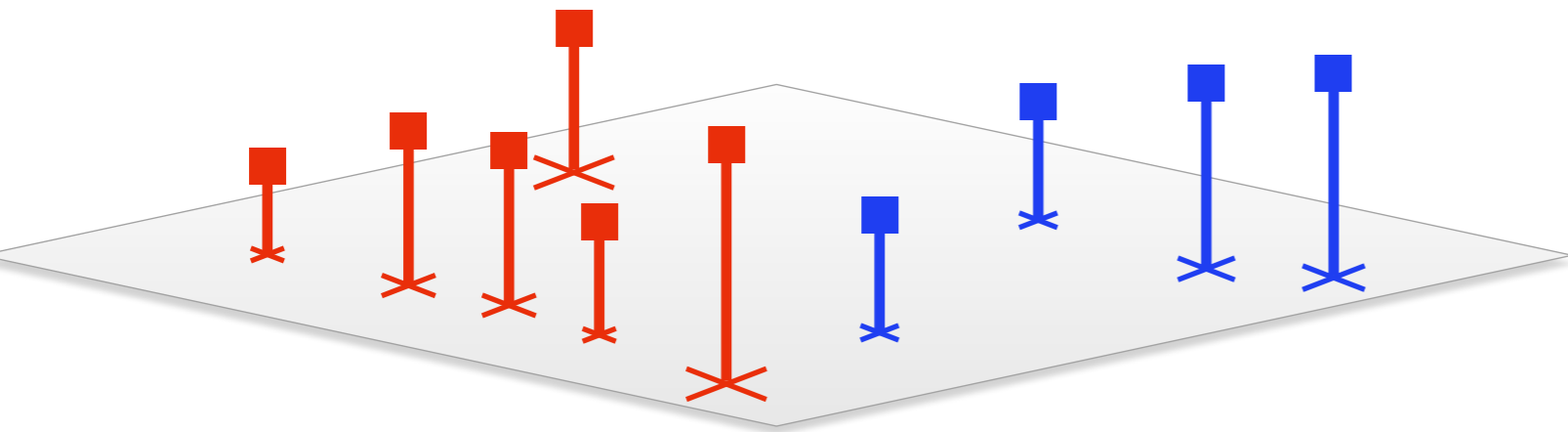
$\nabla_{\mu} W_p(\mu, \nu_1) = ?$

2. How to compute OT

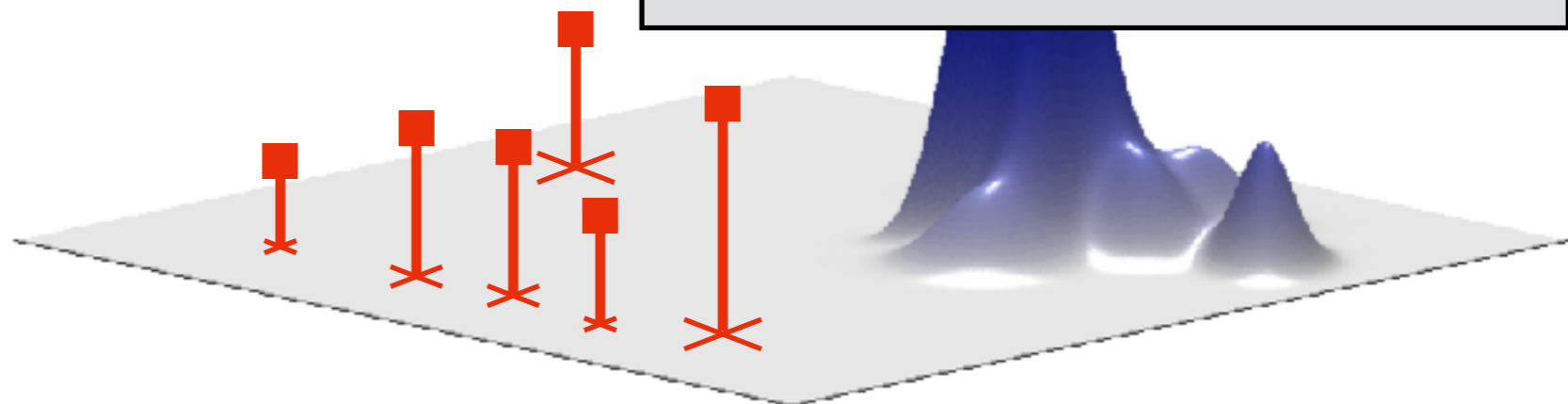
- Typology: discrete/continuous problems
- Easy cases, zoo of solvers
- Entropic regularization
- Differentiability of the W distance

How can we compute OT?

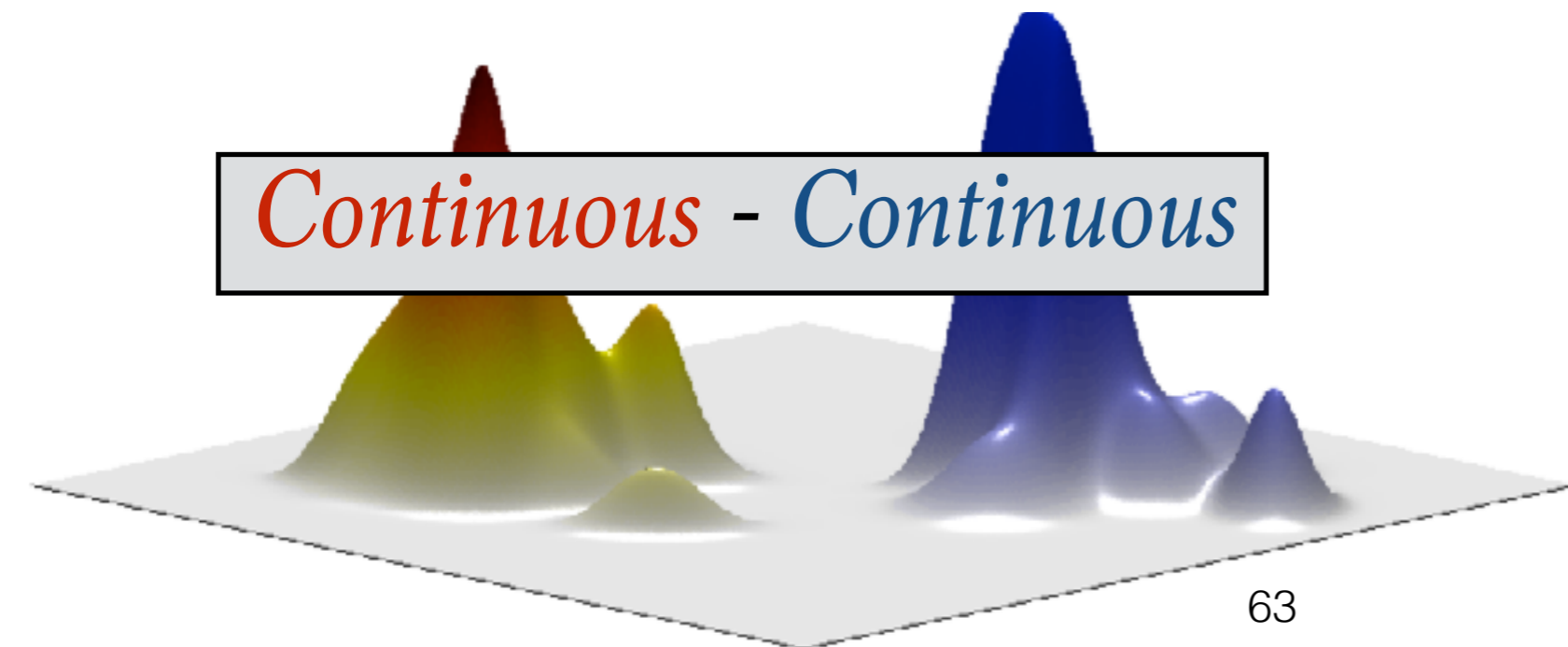
Discrete - Discrete



Discrete - Continuous



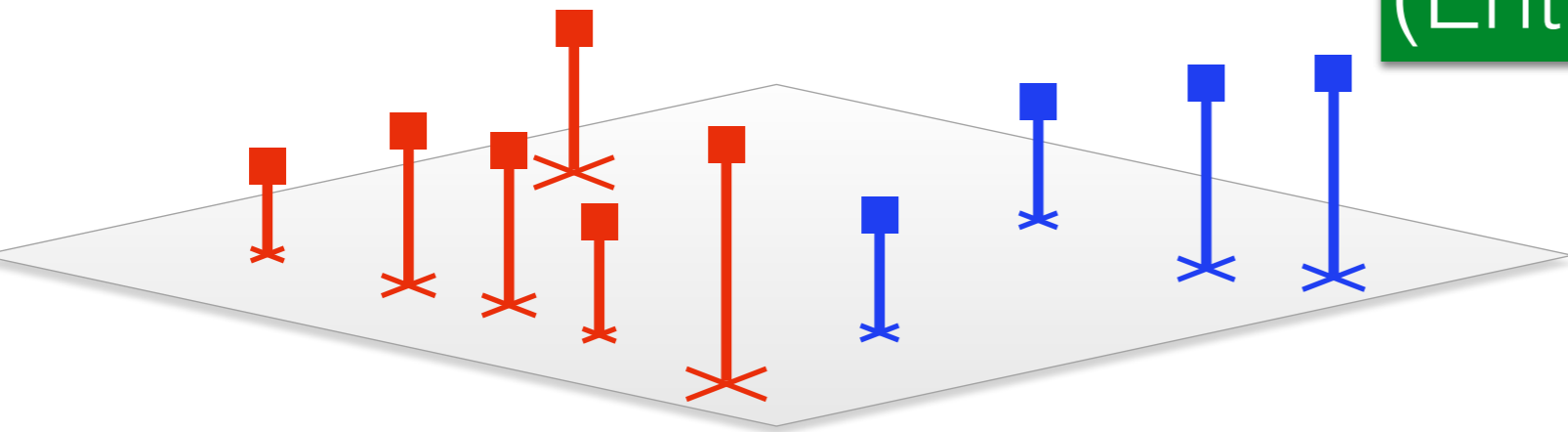
Continuous - Continuous



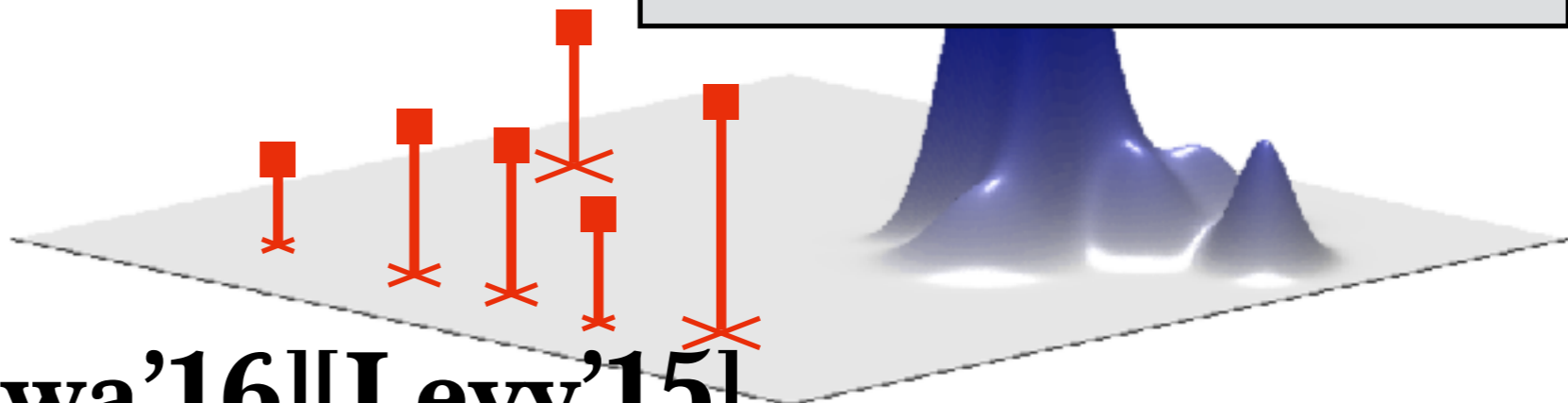
How can we compute OT?

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

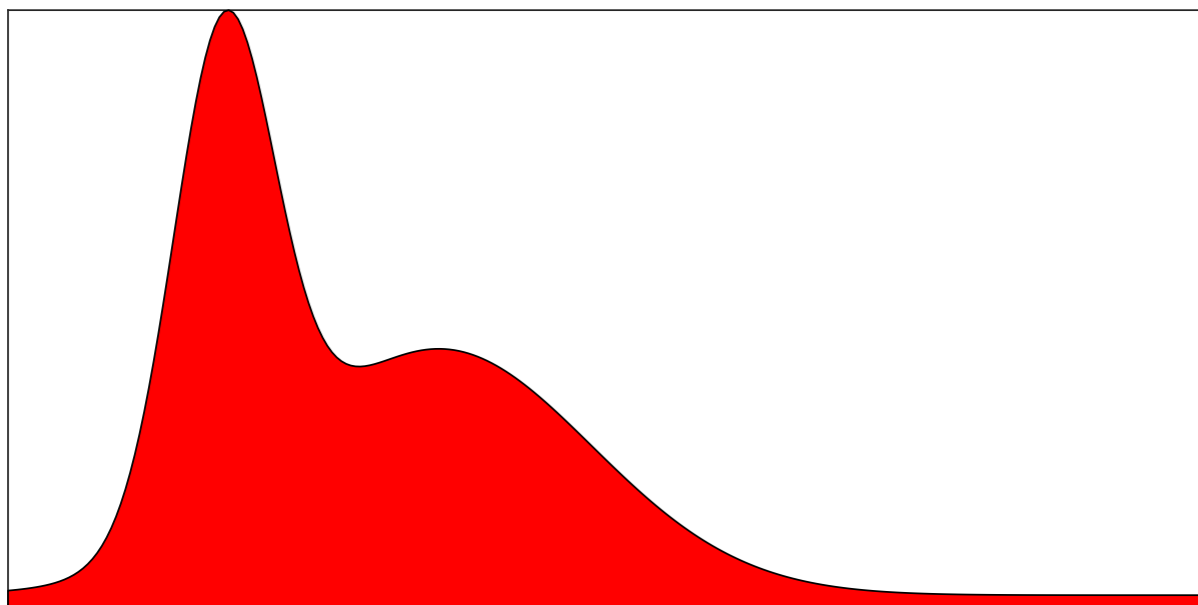
Remark. If $\Omega = \mathbb{R}$, $\mathbf{c}(x, y) = \mathbf{c}(|x - y|)$,
 \mathbf{c} convex, $F_{\mu}^{-1}, F_{\nu}^{-1}$ quantile functions,

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

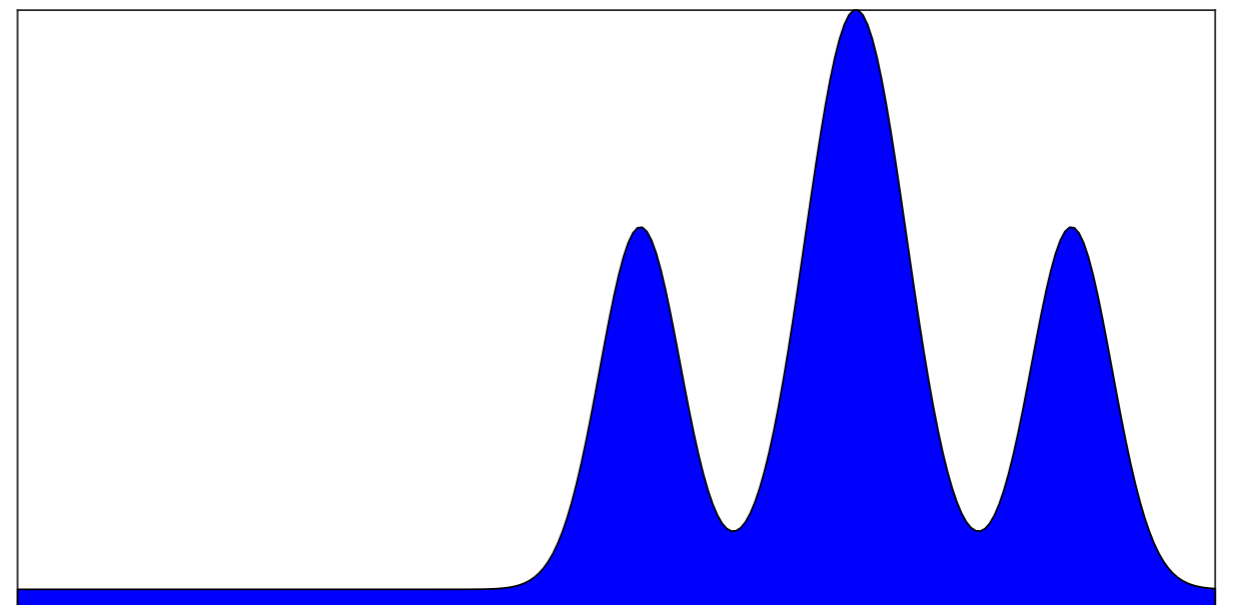
Easy (1): Univariate Measures

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

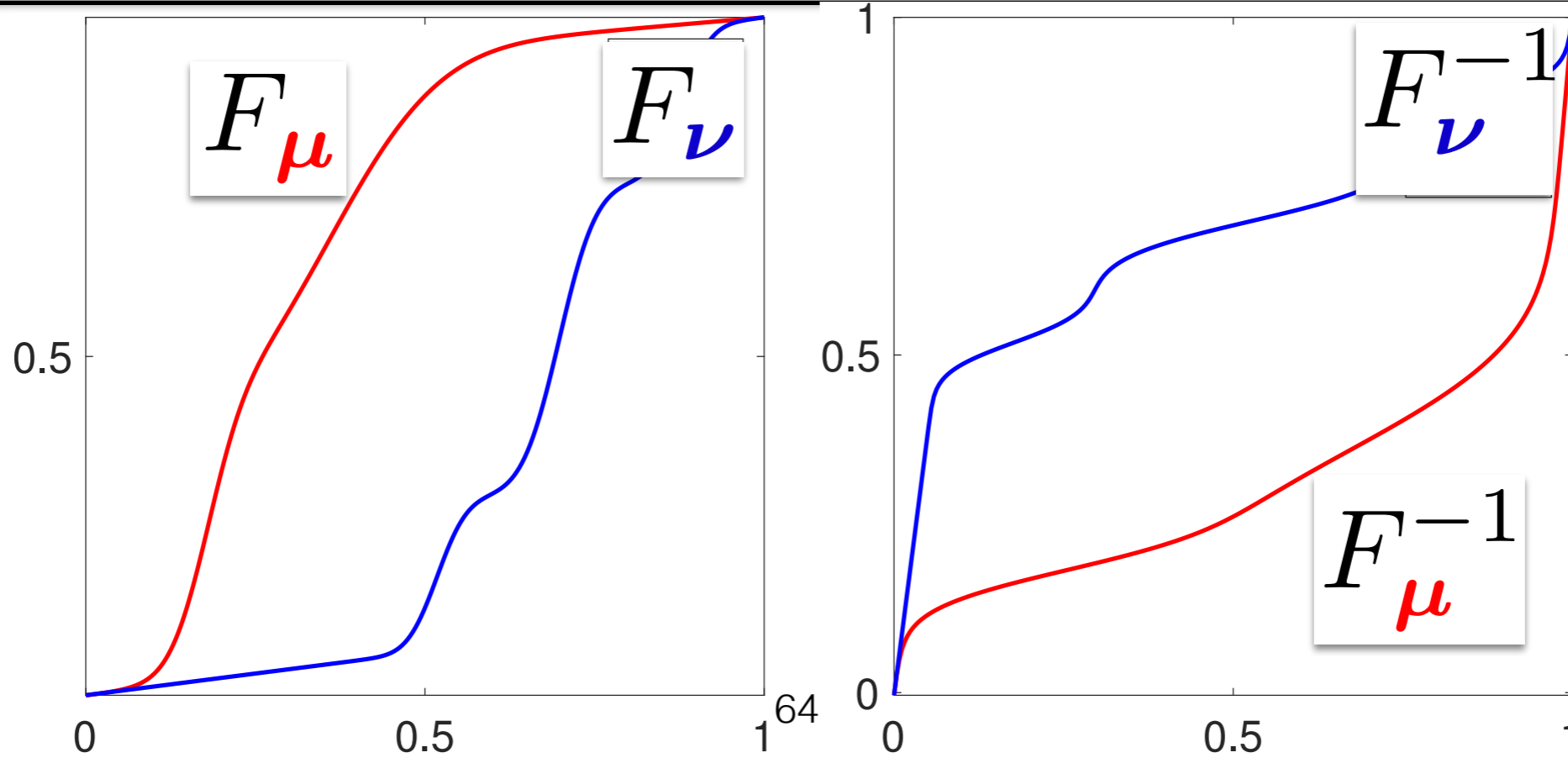


ν

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,

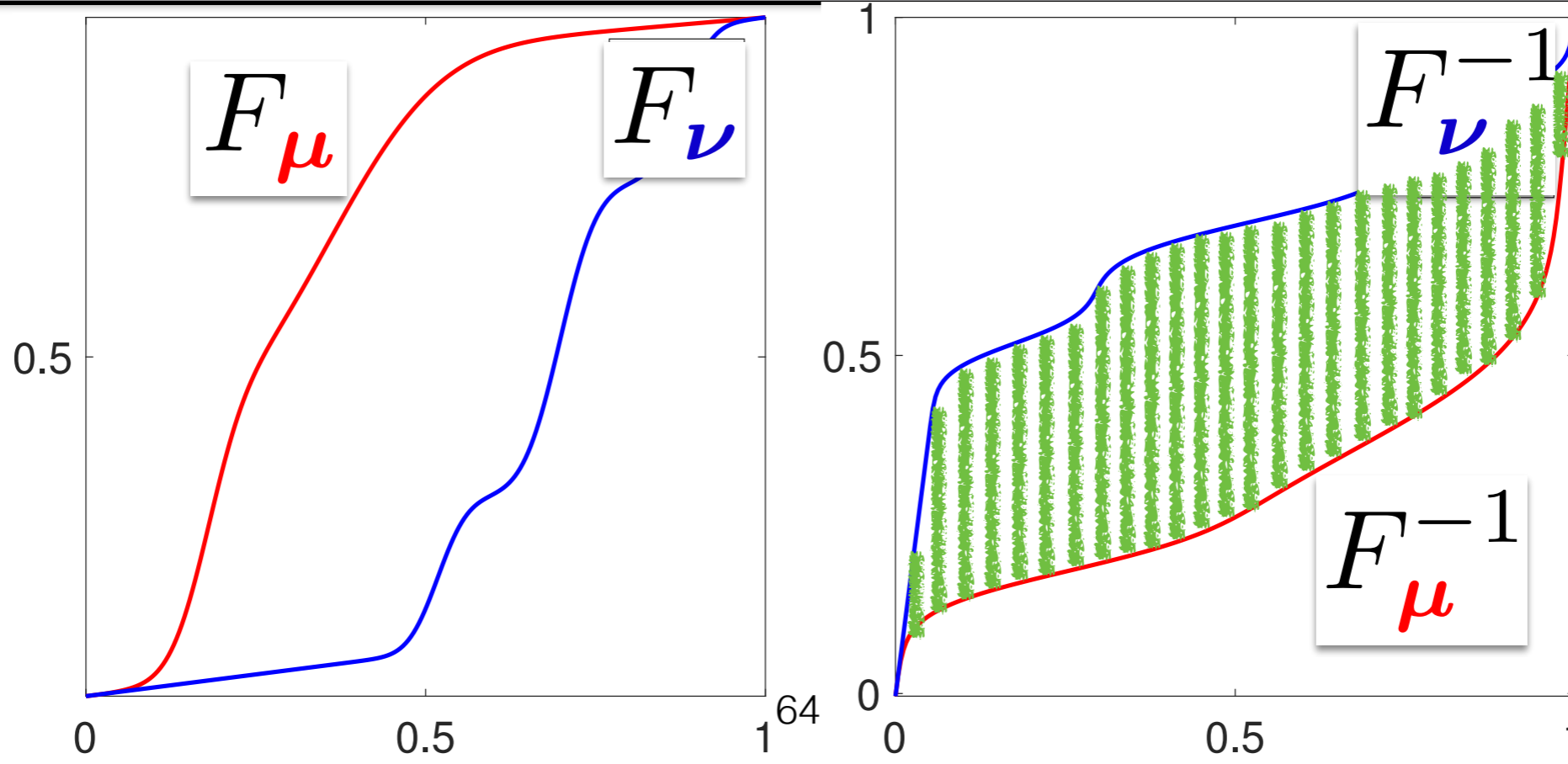
$$W(\mu, \nu) = \int_0^1 \mathbf{c}(|F_{\mu}^{-1}(x) - F_{\nu}^{-1}(x)|) dx$$



Easy (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

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

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

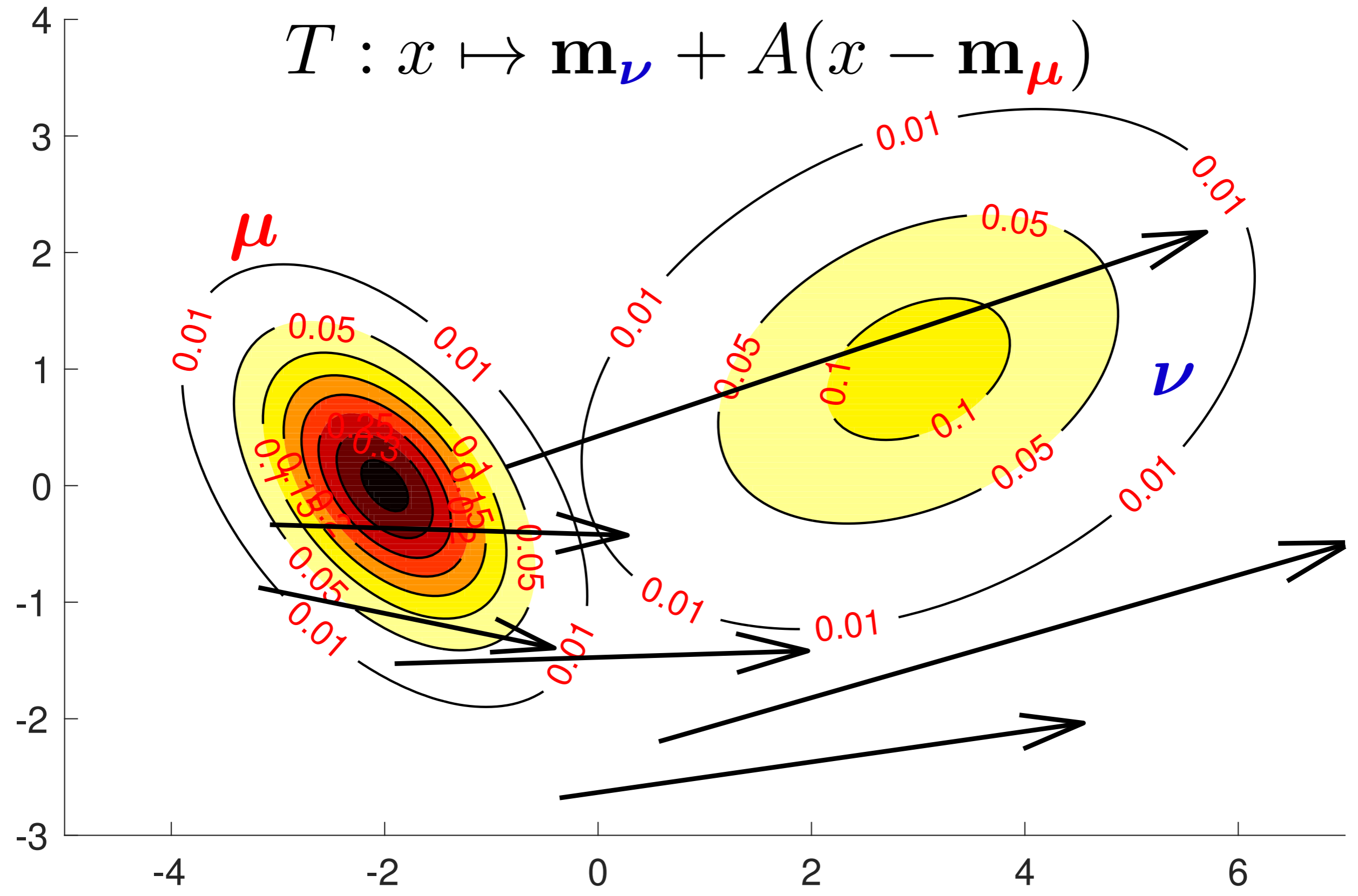
$$B(\Sigma_\mu, \Sigma_\nu)^2 = \text{trace}(\Sigma_\mu + \Sigma_\nu - 2(\Sigma_\mu^{1/2} \Sigma_\nu \Sigma_\mu^{1/2})^{1/2}).$$

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

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

Easy (2): Gaussian Measures

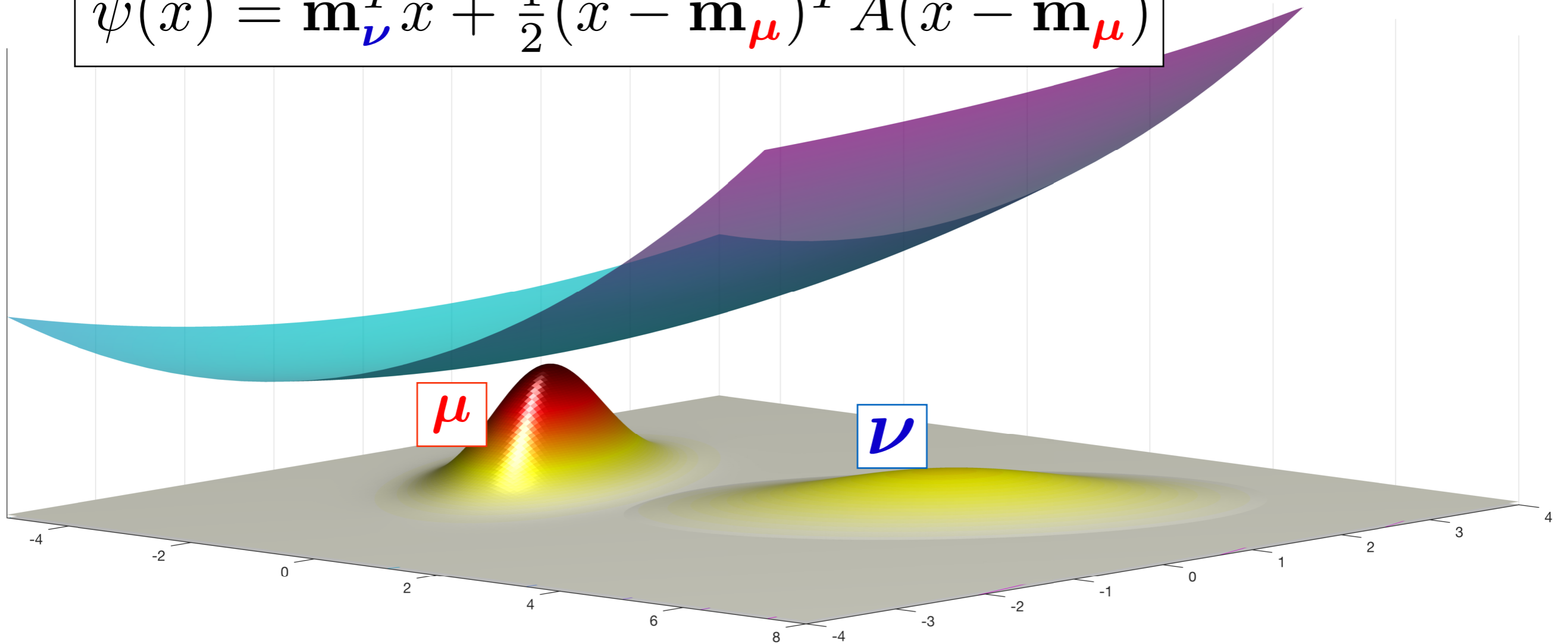
$$T : x \mapsto \mathbf{m}_\nu + A(x - \mathbf{m}_\mu)$$



Easy (2): Gaussian Measures

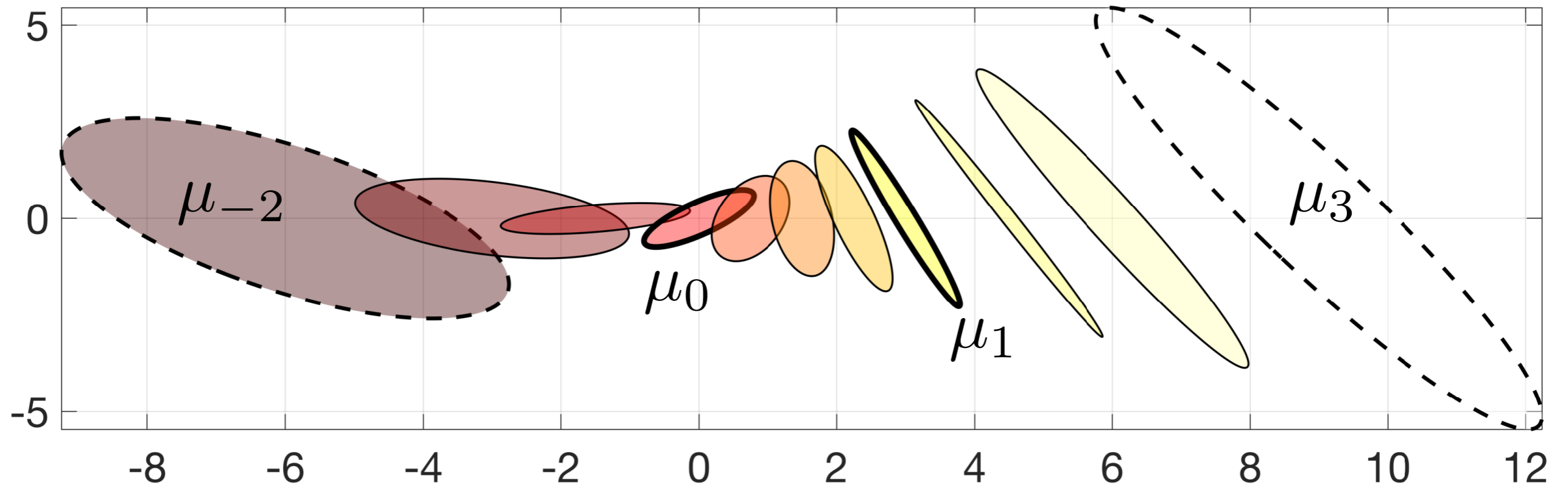
$$T = \nabla \psi : x \mapsto \mathbf{m}_\nu + A(x - \mathbf{m}_\mu)$$

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



Easy (2): Gaussian Measures

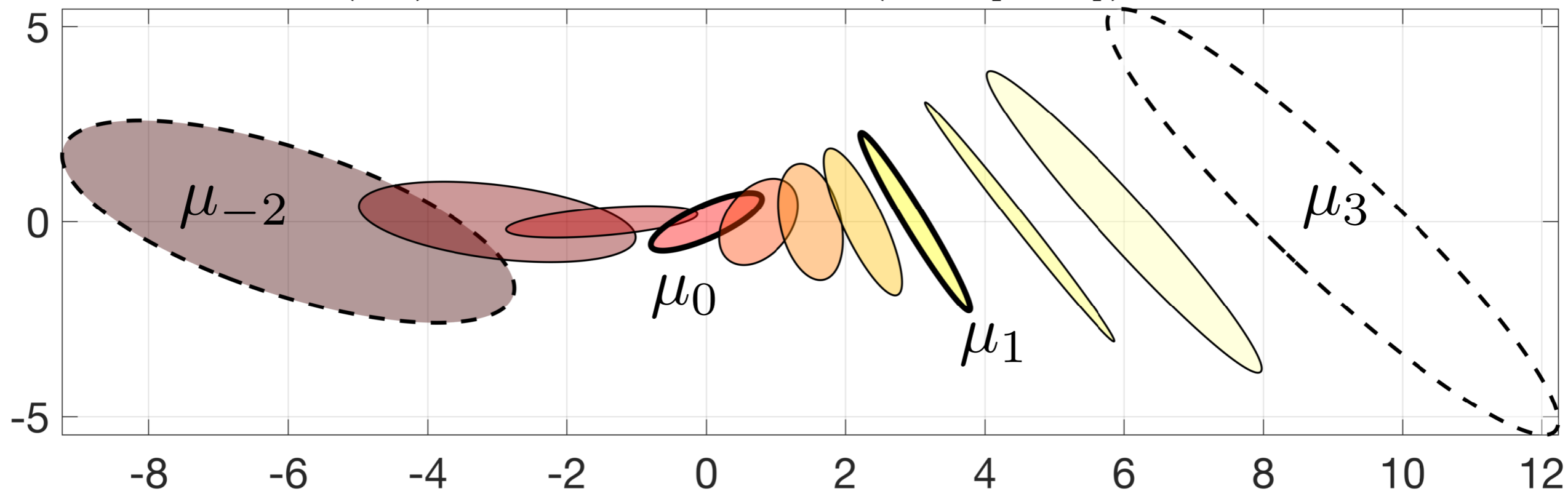
W_2 geodesic $(\mu_t)_t$ from μ_0 to μ_1 ($t \in [0, 1]$) and extrapolation



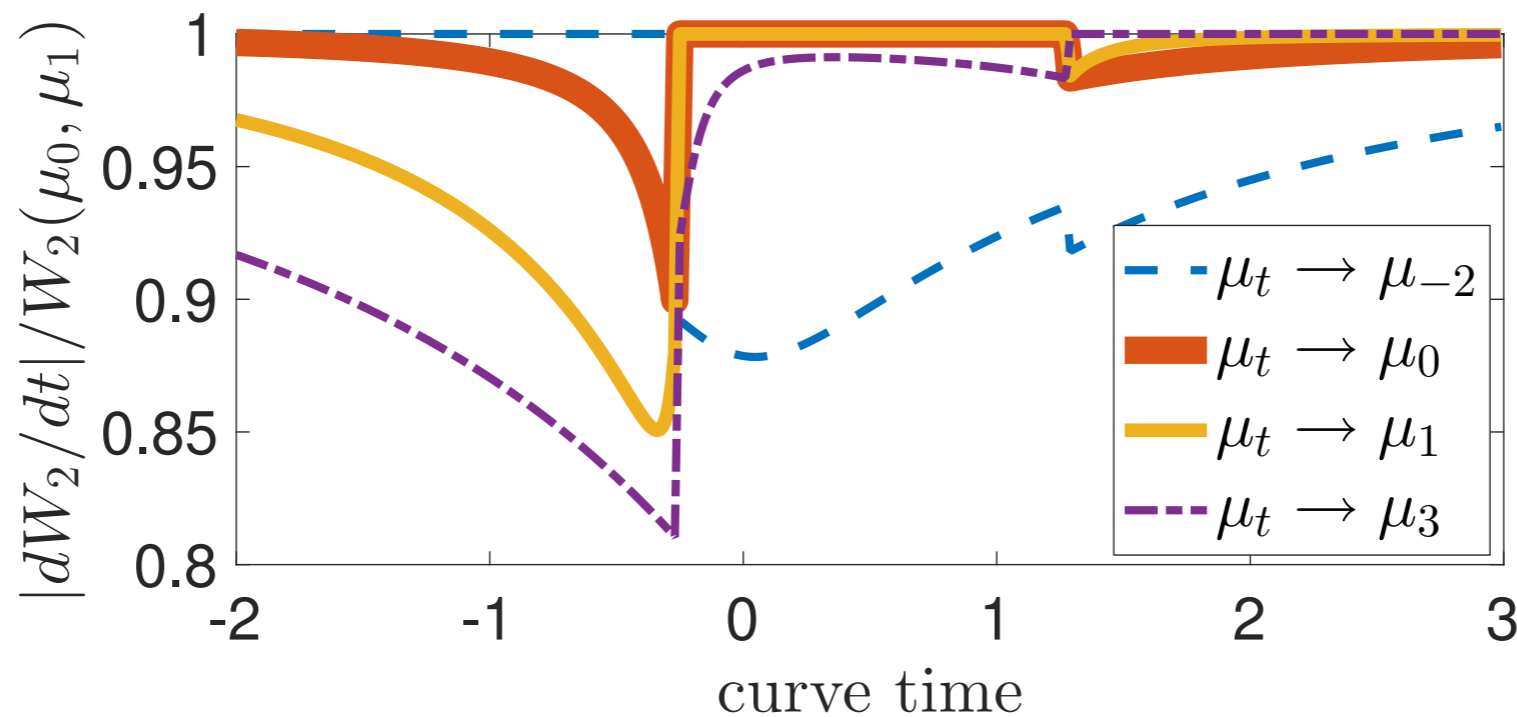
$$\Sigma_t = ((1 - t) I + tA) \Sigma_\mu ((1 - t) I + tA)$$

Easy (2): Gaussian Measures

W_2 geodesic $(\mu_t)_t$ from μ_0 to μ_1 ($t \in [0, 1]$) and extrapolation



Metric derivative on curve



Easy (3): Elliptical Distributions

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

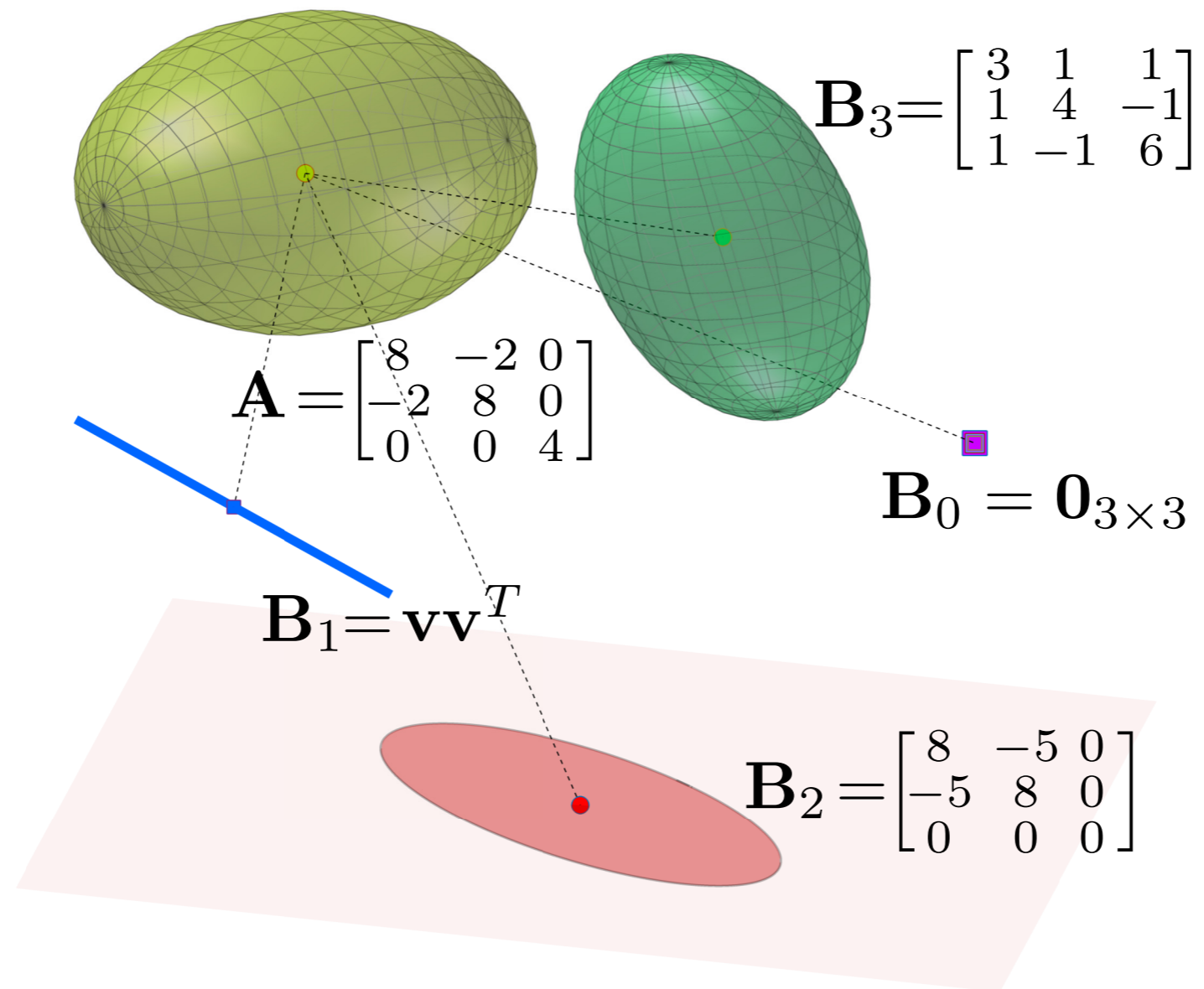
[Gelbrich'92] shows that the linear map T is also **optimal** for elliptically contoured distributions, *i.e.* distributions whose MGF are

$$\phi_X(\mathbf{t}) = \mathbb{E} \left[e^{\sqrt{-1} \mathbf{t}^T X} \right] = e^{\sqrt{-1} \mathbf{t}^T \mathbf{m}} g(\mathbf{t}^T \mathbf{C} \mathbf{t})$$

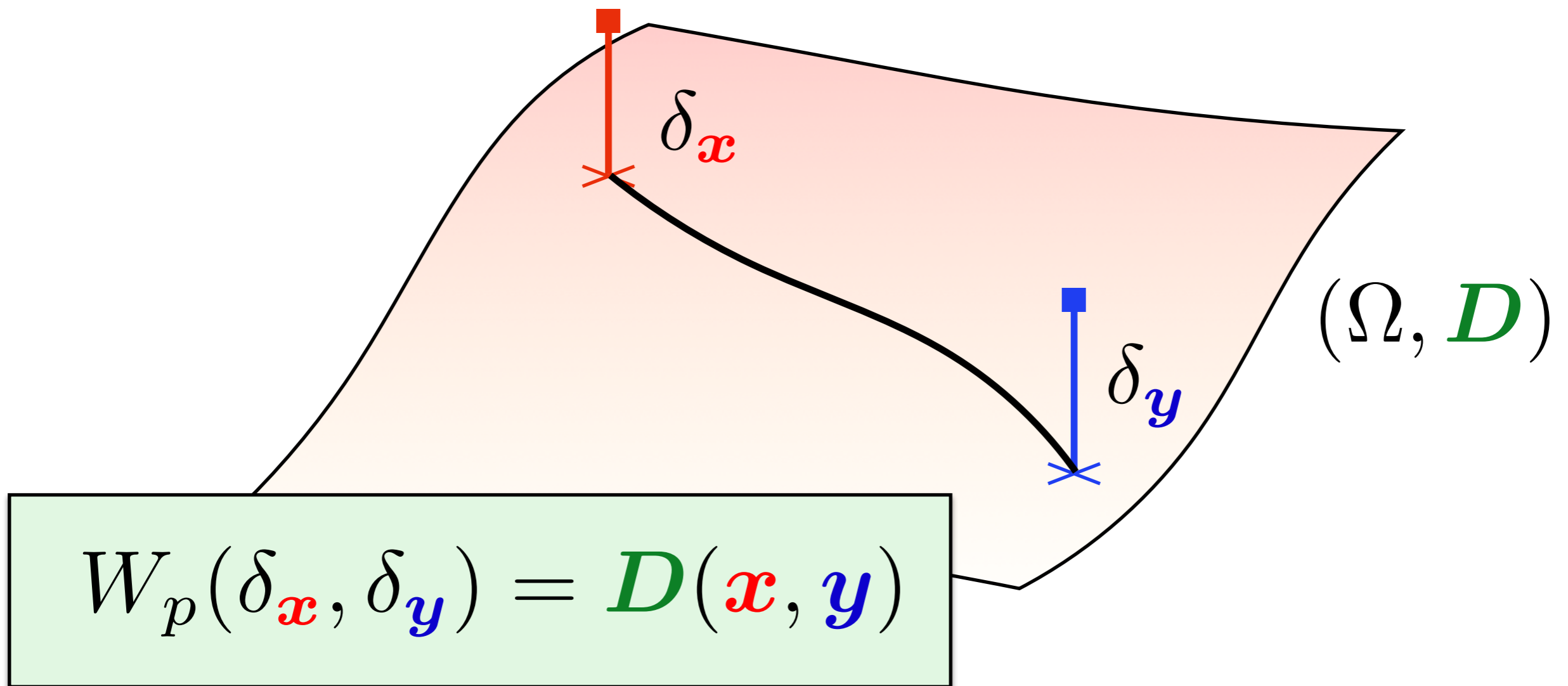
g of positive type.

Same formula applies, but variance is a factor (depends on g) of \mathbf{C} , hence Bures factor is scaled.

Easy (3): Uniform Ellipses



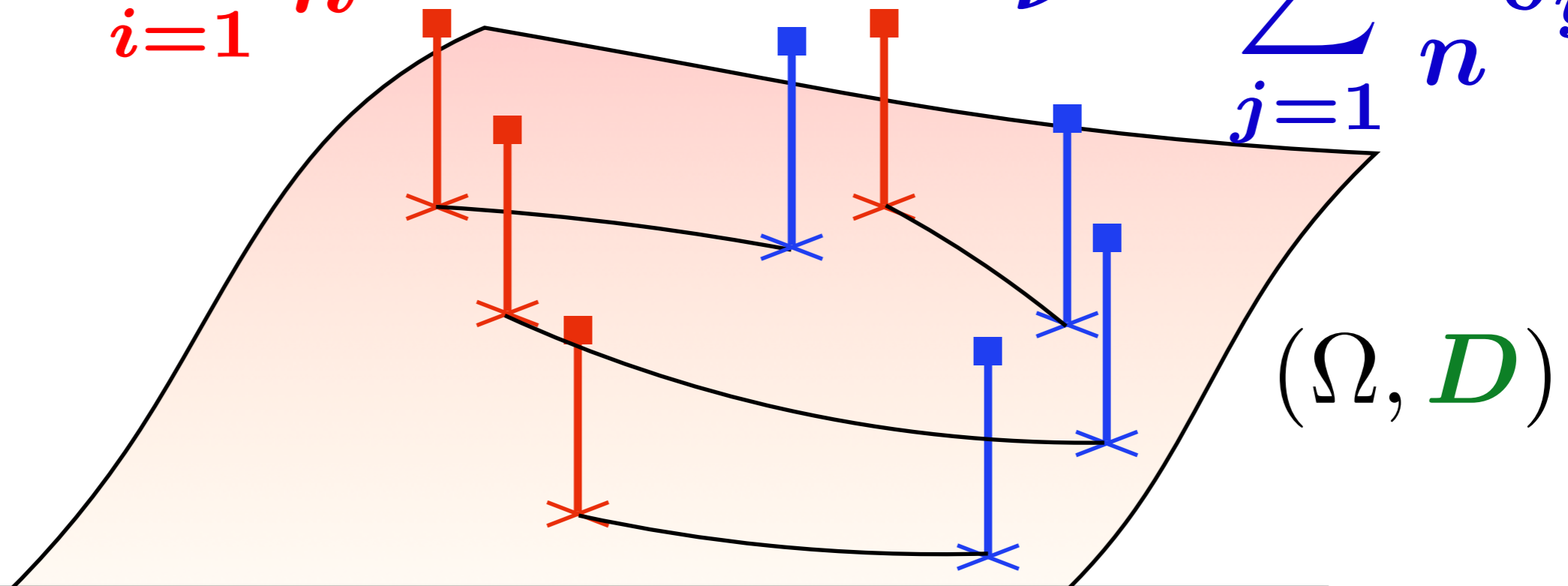
Wasserstein Between Two Diracs



Linear Assignment \subset Wasserstein

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

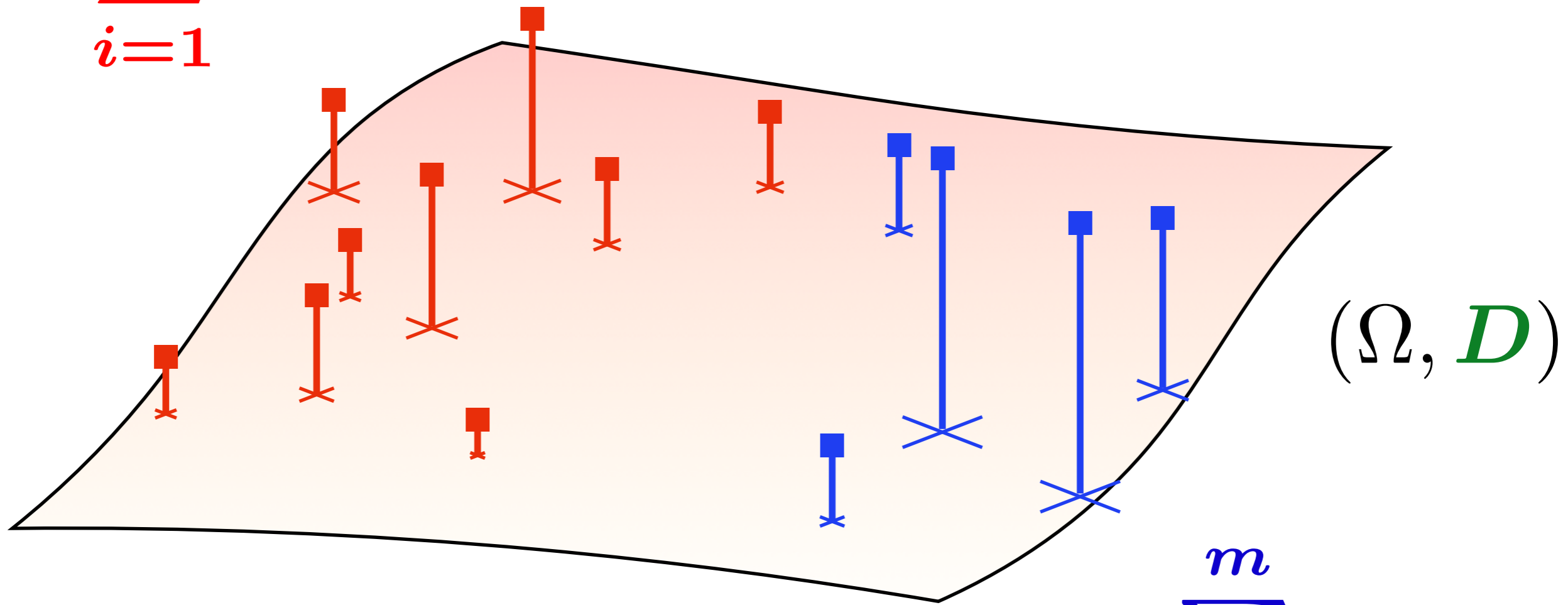
$$\nu = \sum_{j=1}^n \frac{1}{n} \delta_{y_j}$$



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

OT on Two Empirical Measures

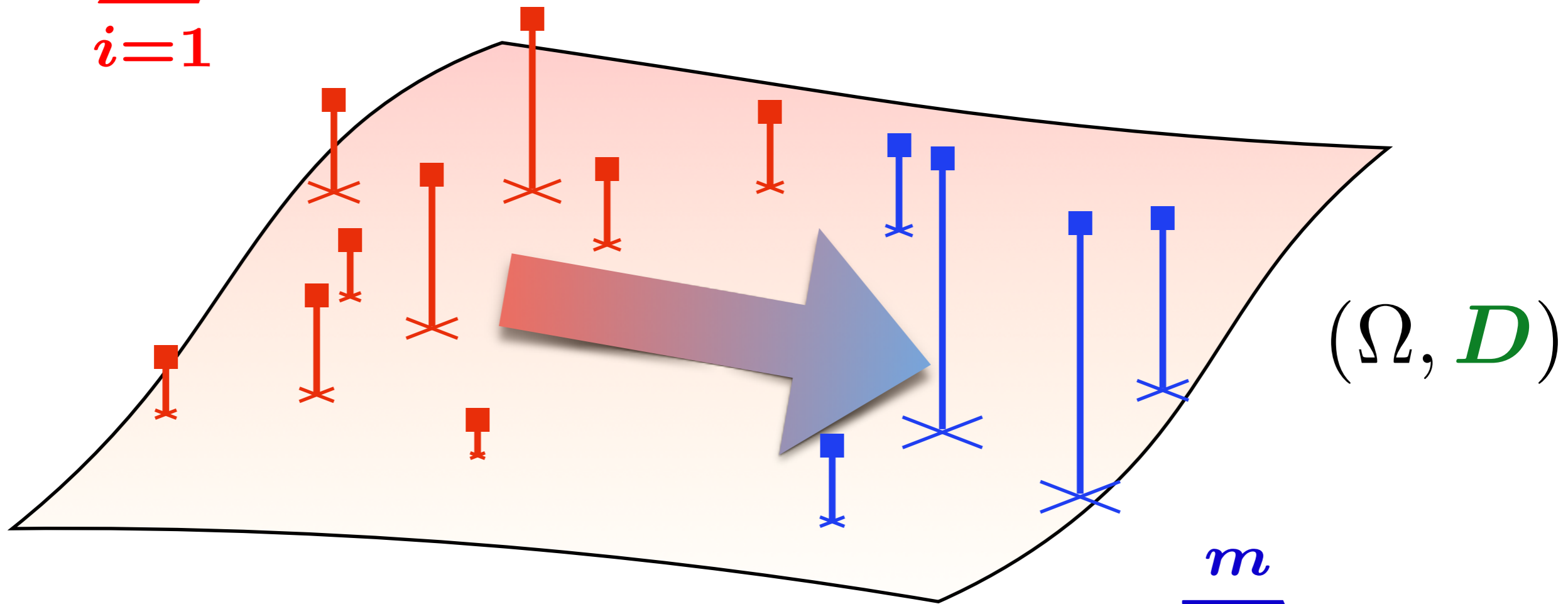
$$\mu = \sum_{i=1}^n a_i \delta_{x_i}$$



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

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

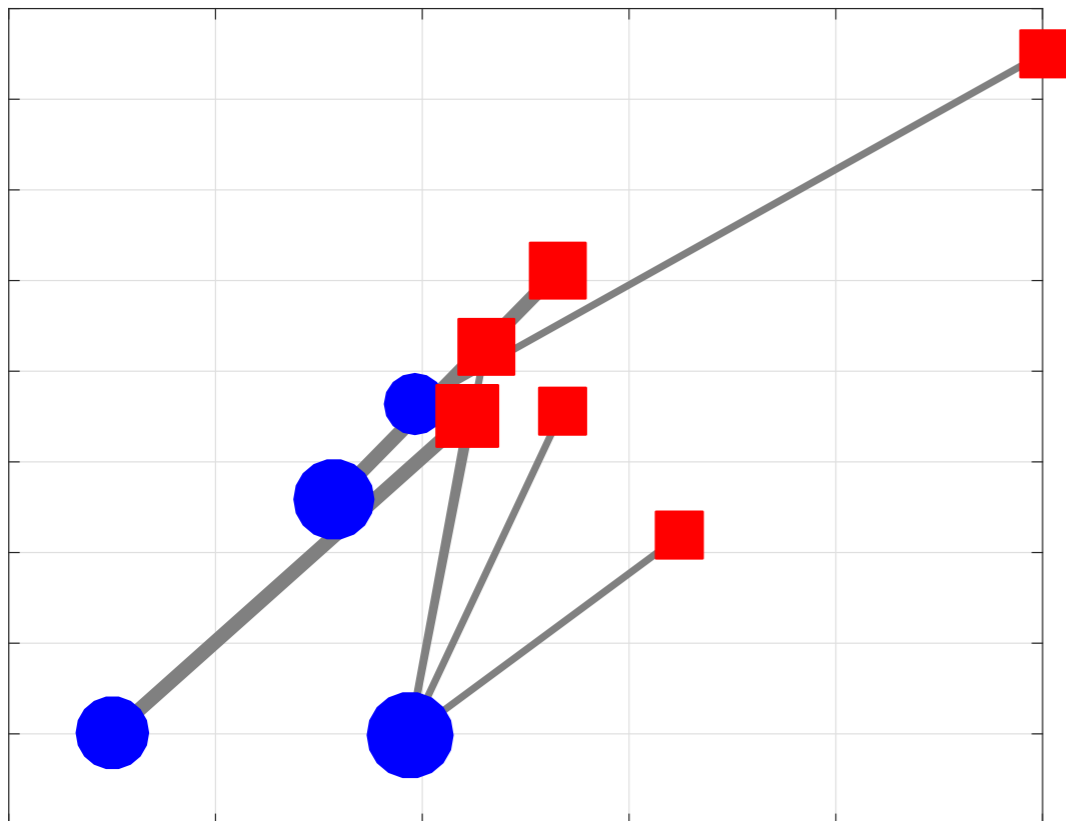
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \alpha_i + \beta_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$

Dual Kantorovich Problem

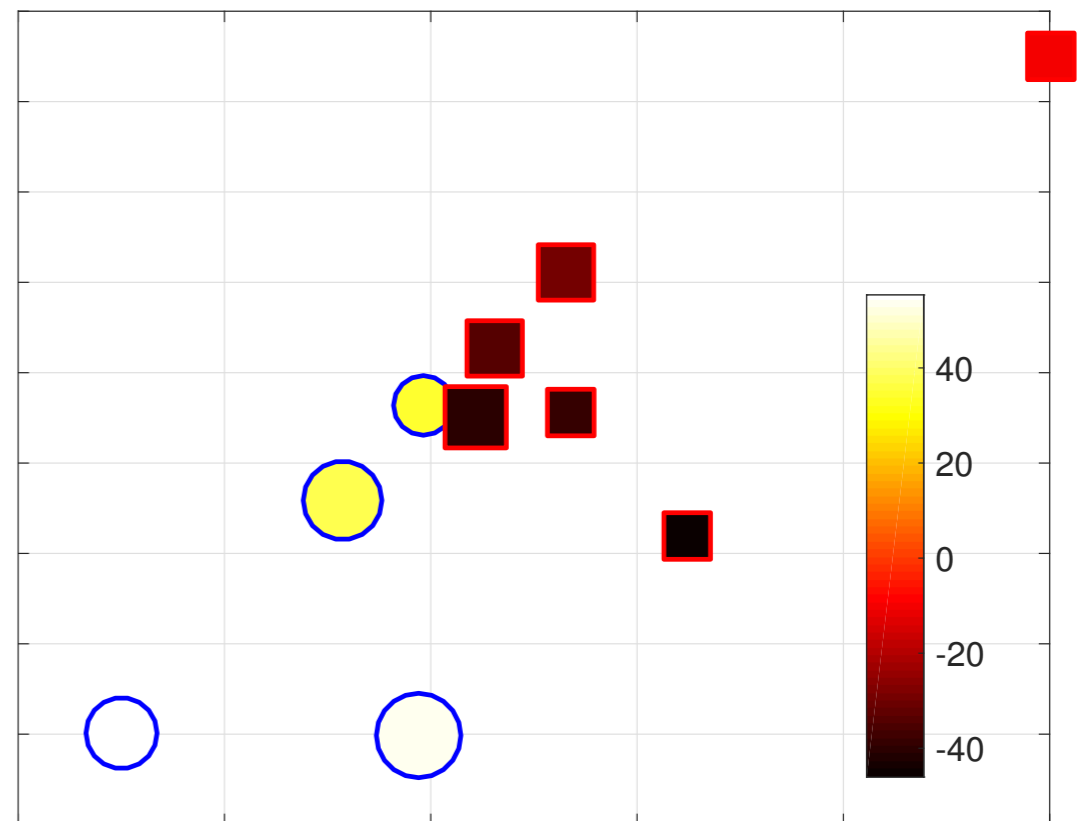
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{x}\mathbf{y}} \rangle$$

Def. Dual OT problem

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



75

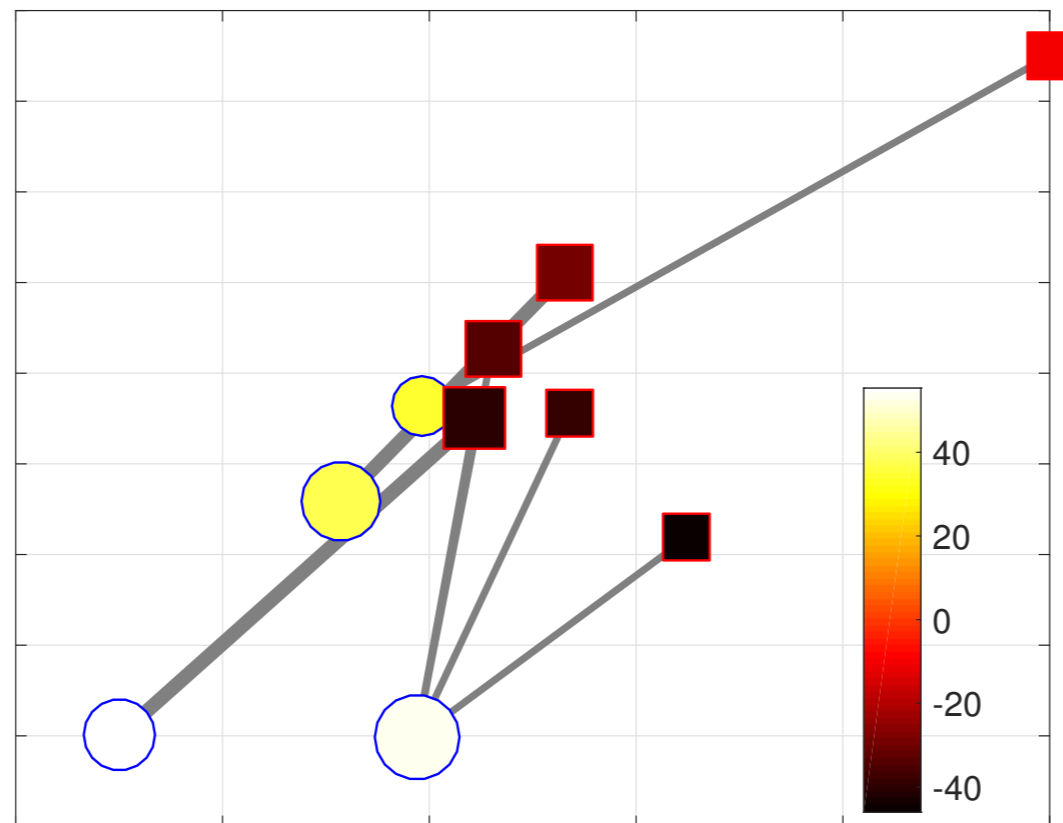


Dual Kantorovich Problem

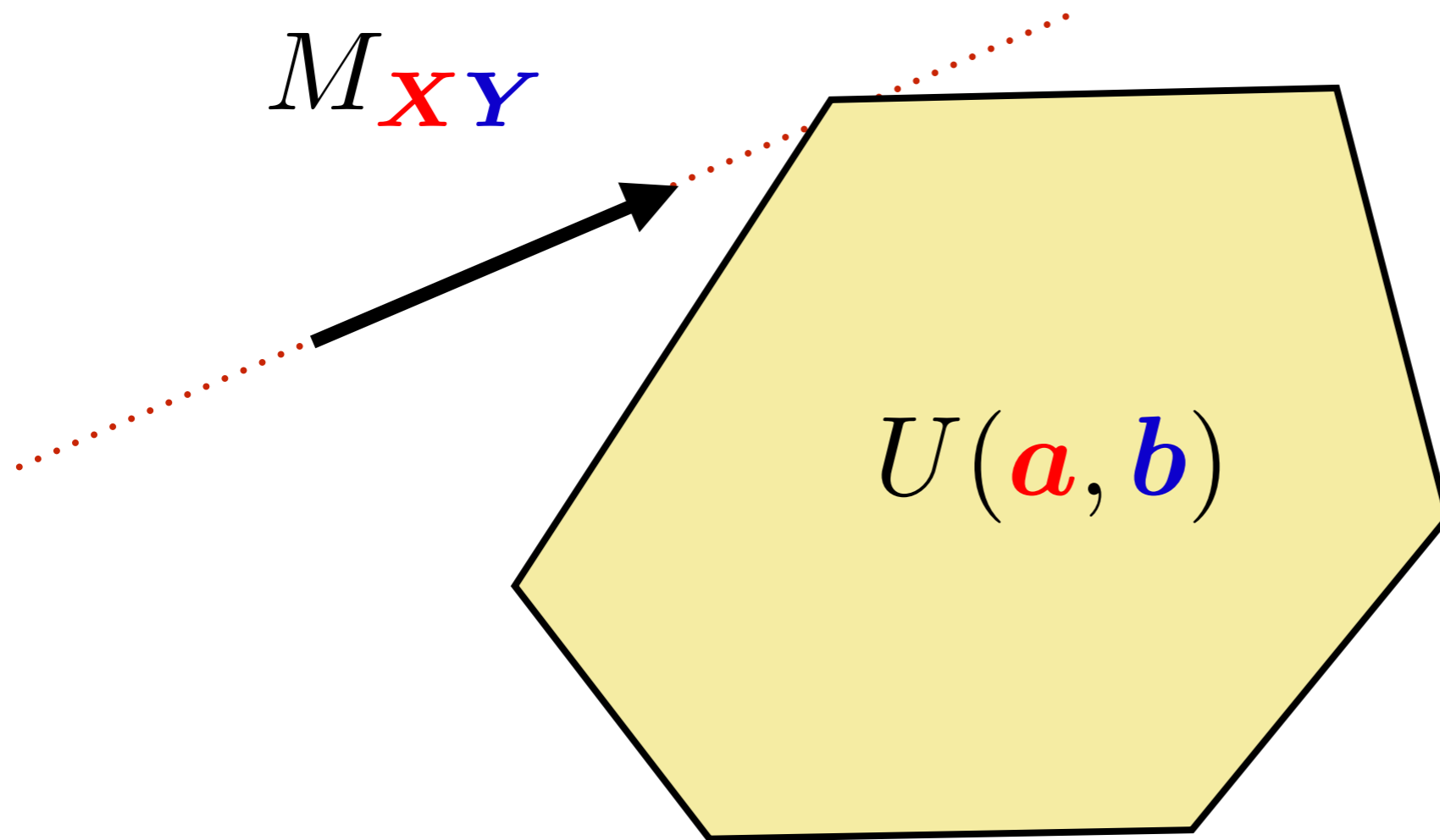
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P \mathbf{1}_m = \mathbf{a}, P^T \mathbf{1}_n = \mathbf{b}}} \langle P, M_{\mathbf{x}\mathbf{y}} \rangle$$

Def. Dual OT problem

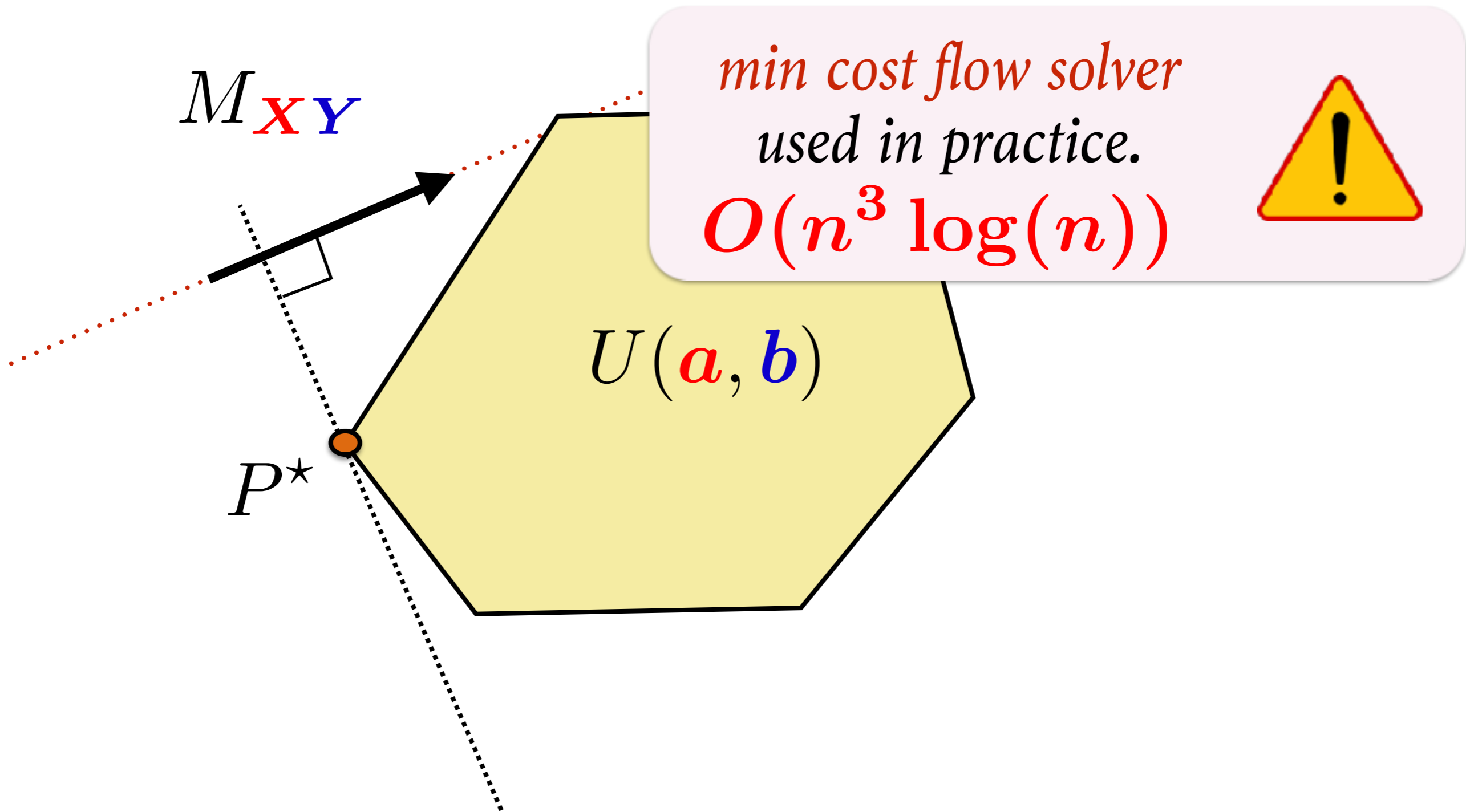
$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \alpha_i + \beta_j \leq D(\mathbf{x}_i, \mathbf{y}_j)^p}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}$$



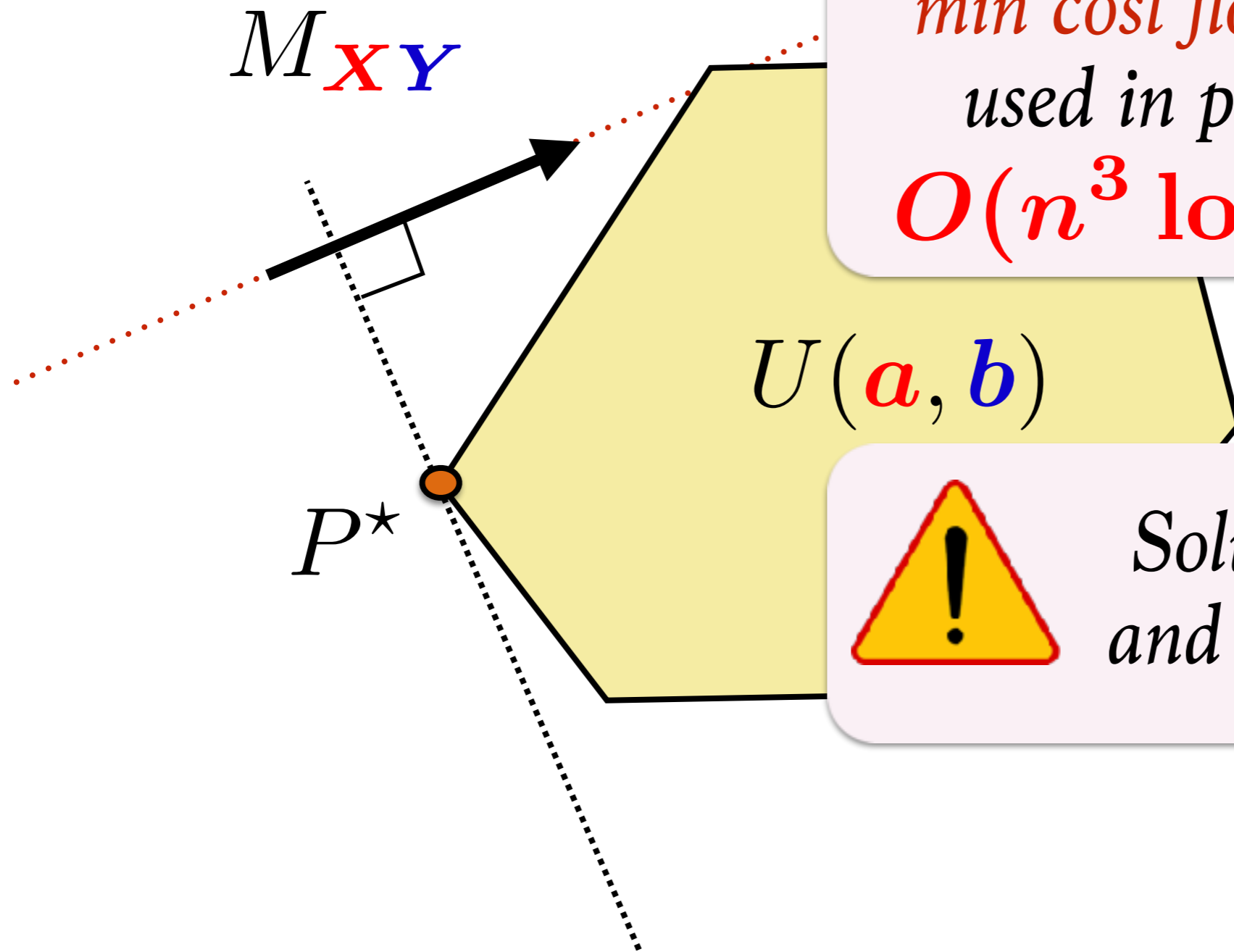
Solving the OT Problem



Solving the OT Problem



Solving the OT Problem

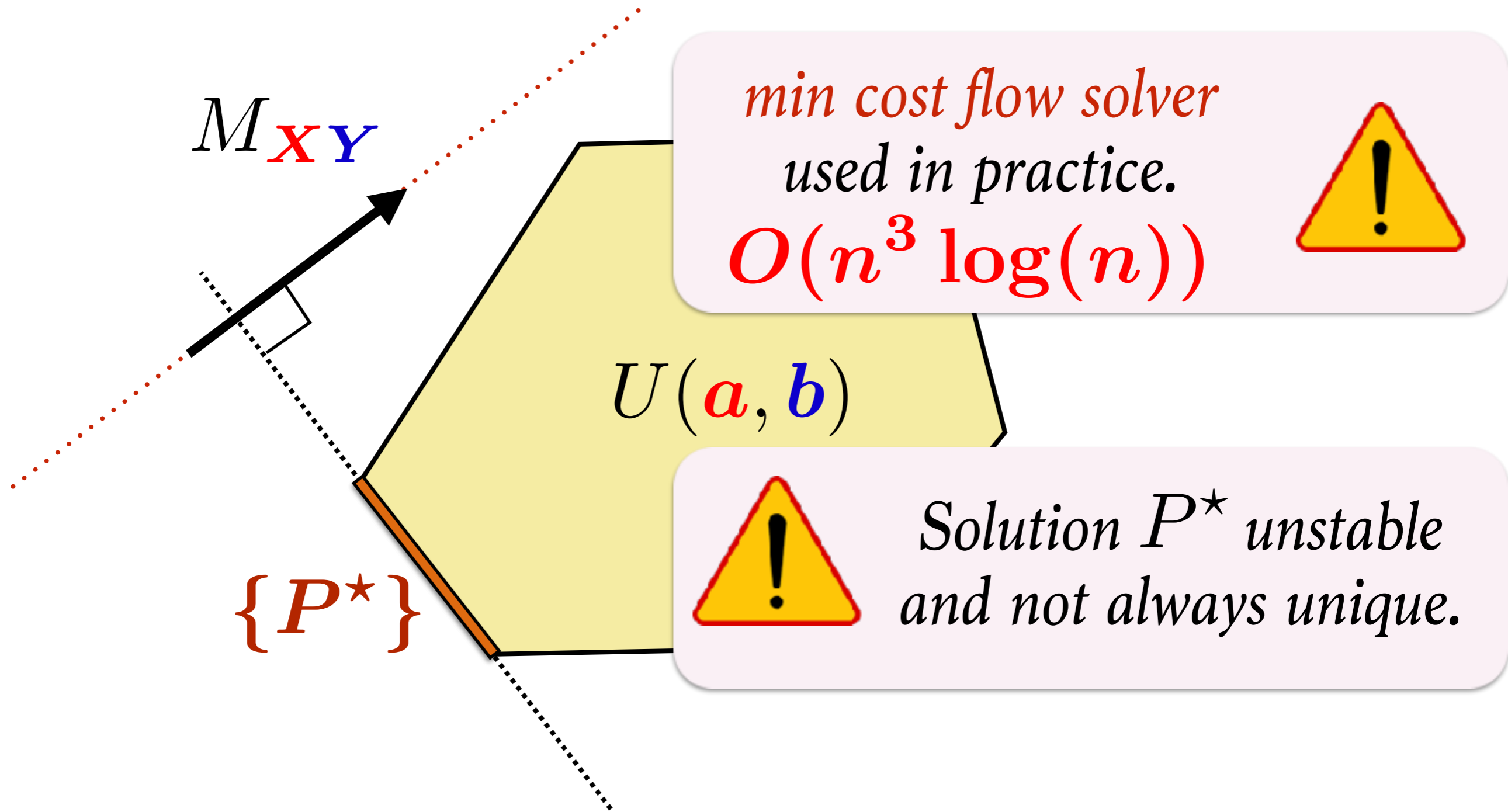


*min cost flow solver
used in practice.
 $O(n^3 \log(n))$*

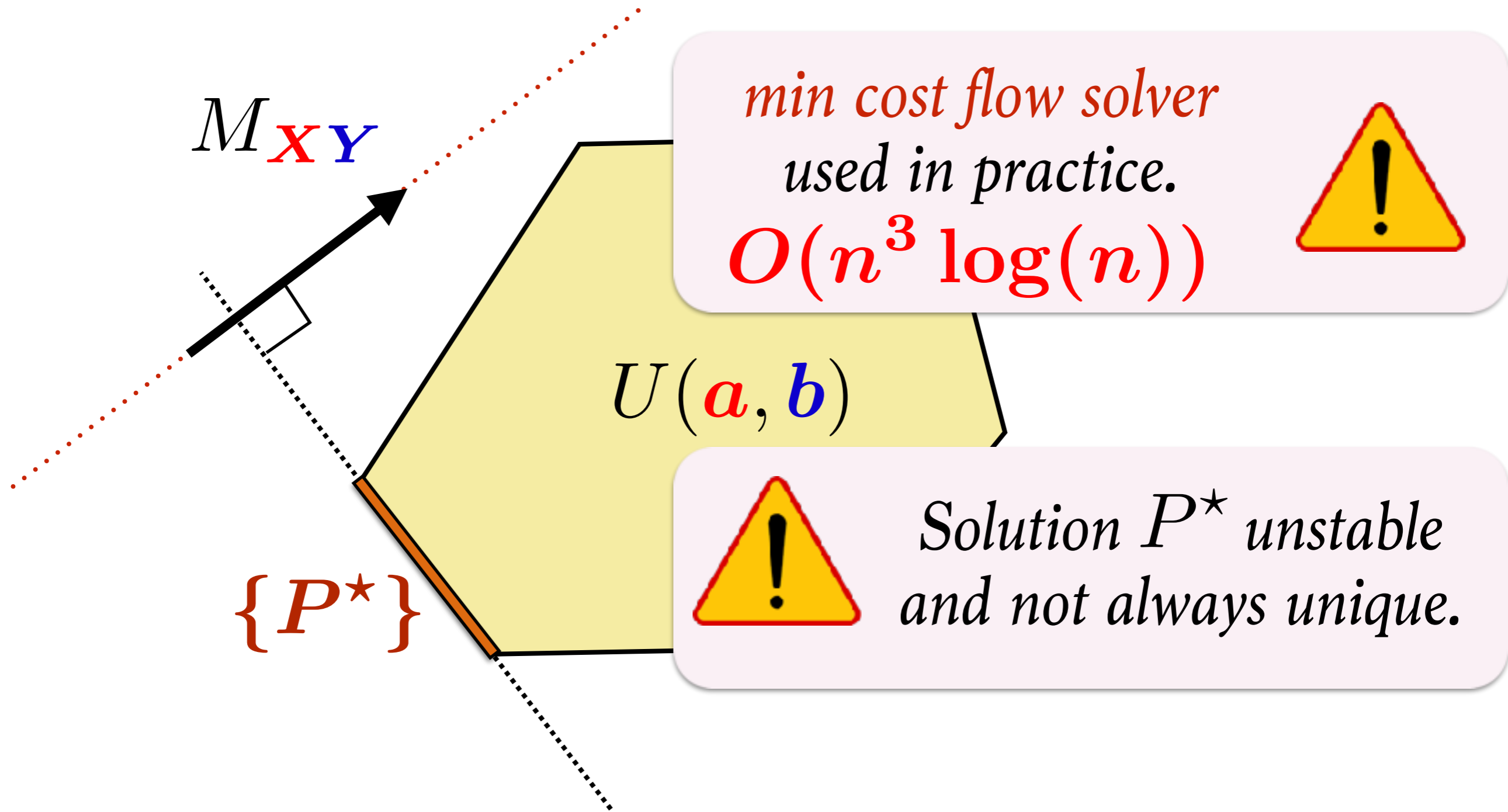


Solution P^ unstable
and not always unique.*

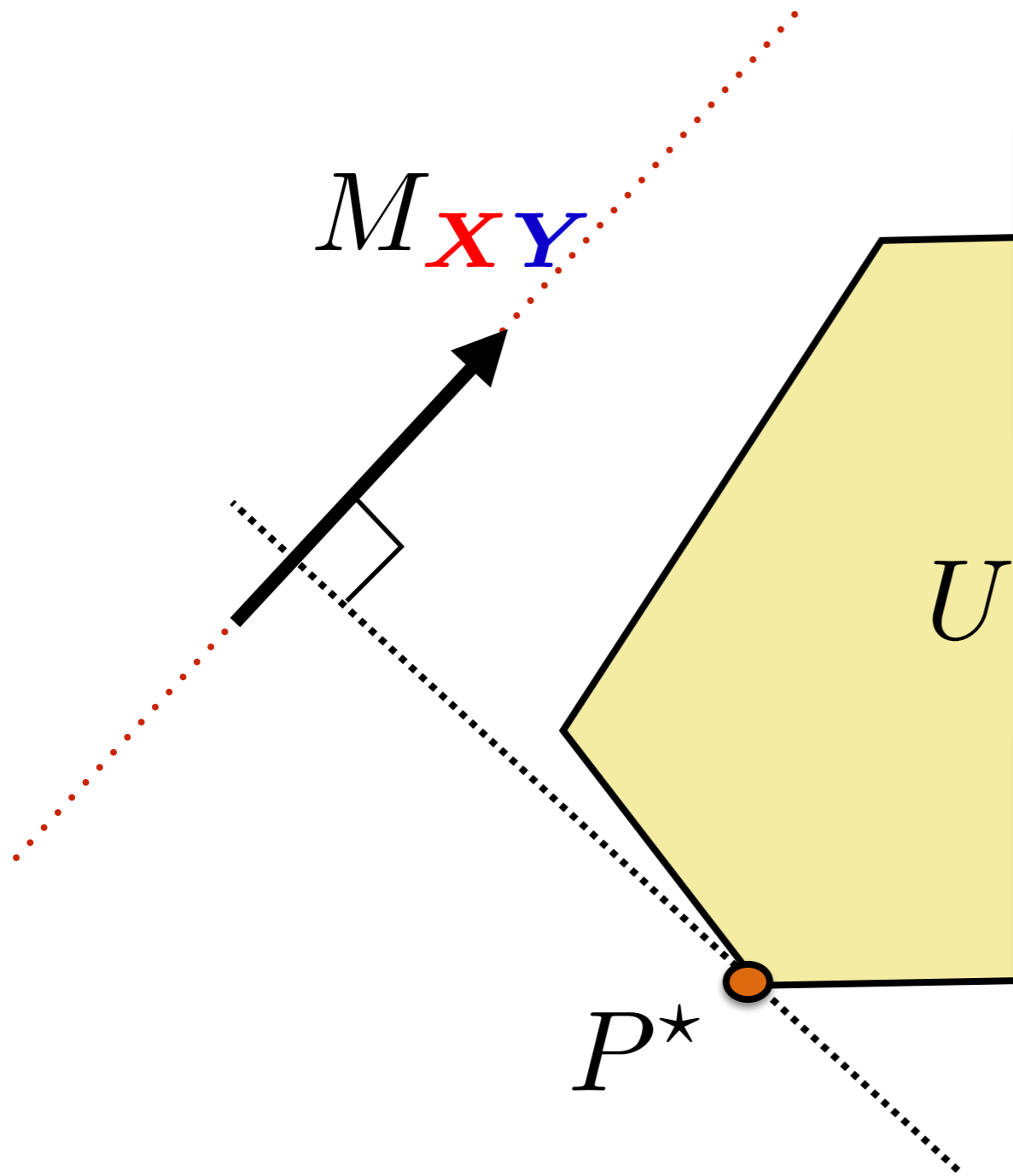
Solving the OT Problem



Solving the OT Problem



Solving the OT Problem

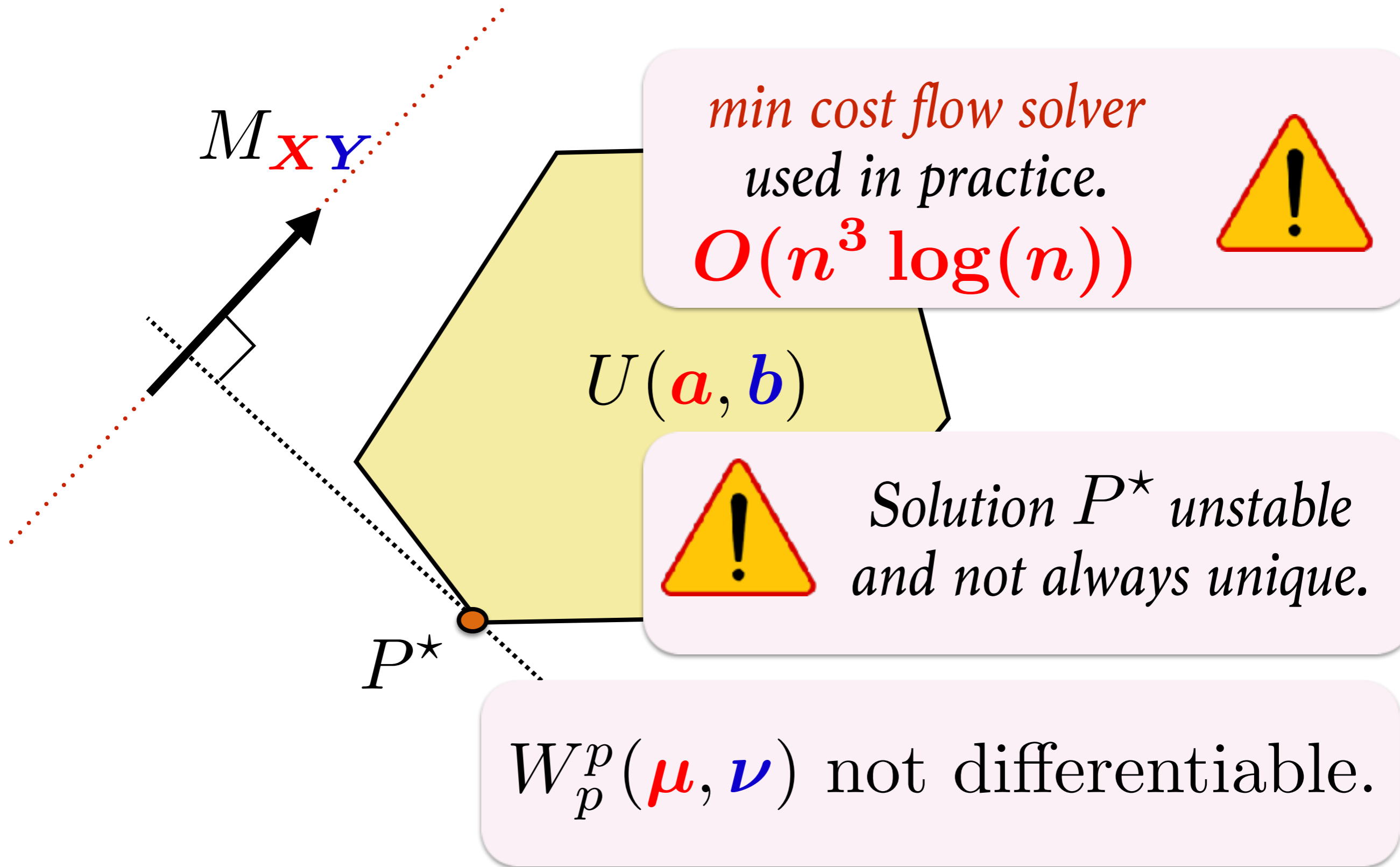


*min cost flow solver
used in practice.
 $O(n^3 \log(n))$*



Solution P^ unstable
and not always unique.*

Solving the 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

```
1  /*
2  end.c
3
4  Last update: 3/14/98
5
6  An implementation of the Earth Movers Distance.
7  Based of the solution for the Transportation problem as described in
8  "Introduction to Mathematical Programming" by F. S. Hillier and
9  G. J. Lieberman, McGraw-Hill, 1990.
10
11  Copyright (C) 1998 Yossi Rubner
12  Computer Science Department, Stanford University
13  E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner
14  */
15
16  /*#include <stdio.h>
17  #include <stdlib.h>*/
18  #include <math.h>
19
20  #include "emd.h"
21
22  #define DEBUG_LEVEL 0
23  /*
24  DEBUG_LEVEL:
25  0 = NO MESSAGES
26  1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
27  2 = PRINT THE RESULT AFTER EVERY ITERATION
28  3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29  4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
30  */
31
32
33  #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSSIBLE DUMMY FEATURE */
34
35  /* NEW TYPES DEFINITION */
36
37  /* node1_t IS USED FOR SINGLE-LINKED LISTS */
38  typedef struct node1_t {
39  int i;
40  double val;
41  struct node1_t *Next;
42  } node1_t;
43
44  /* node2_t IS USED FOR DOUBLE-LINKED LISTS */
45  typedef struct node2_t {
46  int i, j;
47  double val;
48  struct node2_t *NextC; /* NEXT COLUMN */
49  struct node2_t *NextR; /* NEXT ROW */
50  } node2_t;
51
52
53
54  /* GLOBAL VARIABLE DECLARATION */
55  static int _n1, _n2; /* SIGNATURES SIZES */
56  static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1]; /* THE COST MATRIX */
57  static node2_t _X[MAX_SIG_SIZE1+2]; /* THE BASIC VARIABLES VECTOR */
58
```

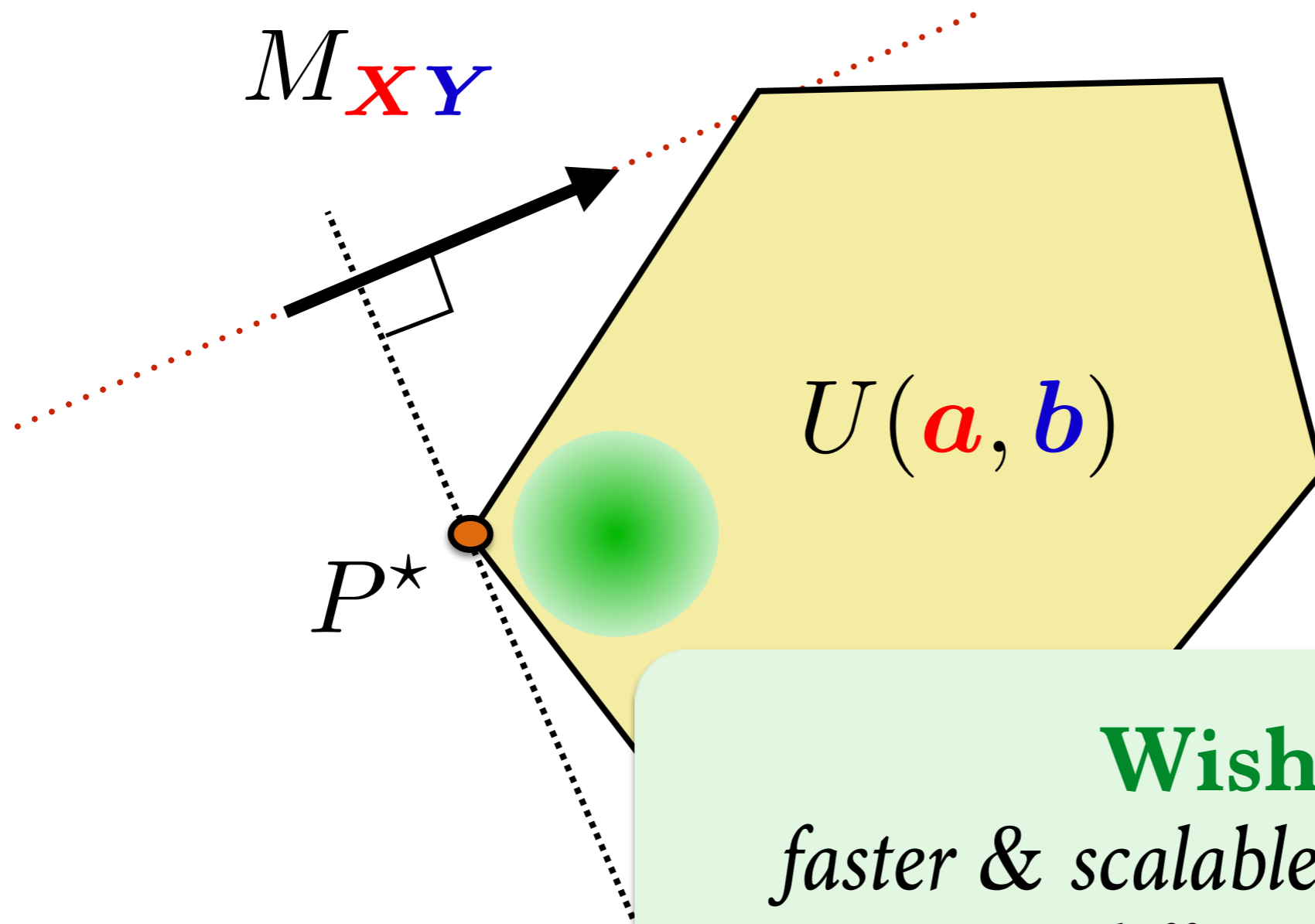


Discrete OT Problem

```
1  /*
2  end.c
3
4  Last update: 3/14/98
5
6  An implementation of the Earth Movers Distance.
7  Based of the solution for the Transportation problem as described in
8  "Introduction to Mathematical Programming" by F. S. Hillier and
9  G. J. Lieberman, McGraw-Hill, 1990.
10
11  Copyright (C) 1998 Yossi Rubner
12  Computer Science Department, Stanford University
13  E-Mail: rubner@cs.stanford.edu URL: http://vision.stanford.edu/~rubner
14  */
15
16  /*#include <stdio.h>
17  #include <stdlib.h>*/
18  #include <math.h>
19
20  #include "emd.h"
21
22  #define DEBUG_LEVEL 0
23  /*
24  DEBUG_LEVEL:
25  0 = NO MESSAGES
26  1 = PRINT THE NUMBER OF ITERATIONS AND THE FINAL RESULT
27  2 = PRINT THE RESULT AFTER EVERY ITERATION
28  3 = PRINT ALSO THE FLOW AFTER EVERY ITERATION
29  4 = PRINT A LOT OF INFORMATION (PROBABLY USEFUL ONLY FOR THE AUTHOR)
30  */
31
32
33  #define MAX_SIG_SIZE1 (MAX_SIG_SIZE+1) /* FOR THE POSSIBLE DUMMY FEATURE */
34
35  /* NEW TYPES DEFINITION */
36
37  /* node1_t IS USED FOR SINGLE-LINKED LISTS */
38  typedef struct node1_t {
39  int i;
40  double val;
41  struct node1_t *Next;
42  } node1_t;
43
44  /* node2_t IS USED FOR DOUBLE-LINKED LISTS */
45  typedef struct node2_t {
46  int i, j;
47  double val;
48  struct node2_t *NextC; /* NEXT COLUMN */
49  struct node2_t *NextR; /* NEXT ROW */
50  } node2_t;
51
52
53
54  /* GLOBAL VARIABLE DECLARATION */
55  static int _n1, _n2; /* SIGNATURES SIZES */
56  static float _C[MAX_SIG_SIZE1][MAX_SIG_SIZE1]; /* THE COST MATRIX */
57  static node2_t _X[MAX_SIG_SIZE1+2]; /* THE BASIC VARIABLES VECTOR */
58
```



Solution: Regularization



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

Entropic Regularization [Wilson'62]

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

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

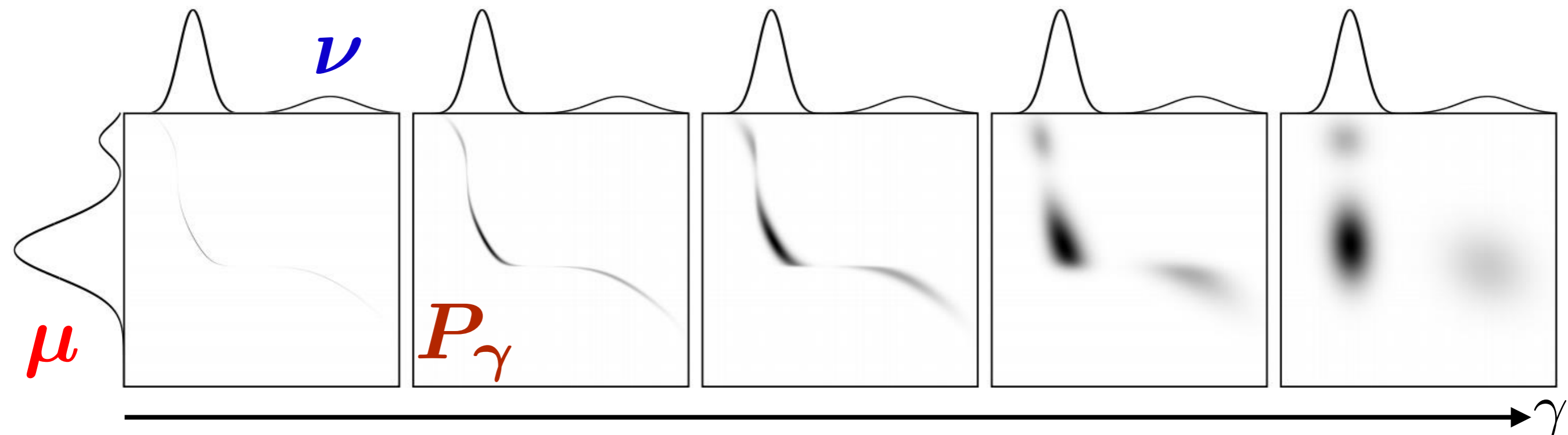
$$E(P) \stackrel{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij} (\log P_{ij} - 1)$$

Note: Unique optimal solution because of strong concavity of entropy

Entropic Regularization [Wilson'62]

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

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

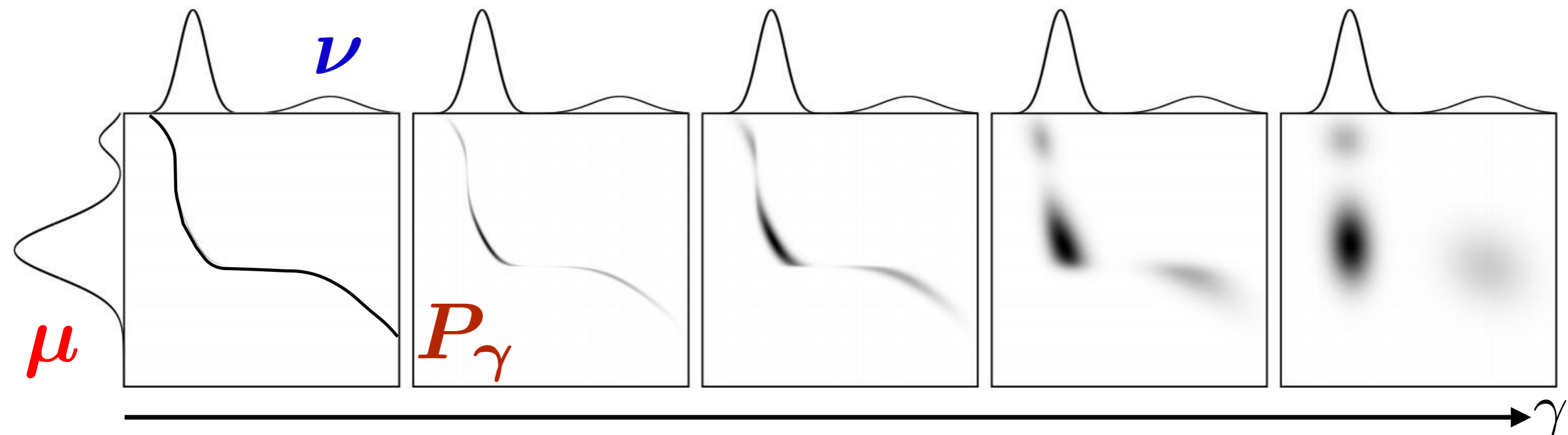


Note: Unique optimal solution because of strong concavity of entropy

Entropic Regularization [Wilson'62]

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

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

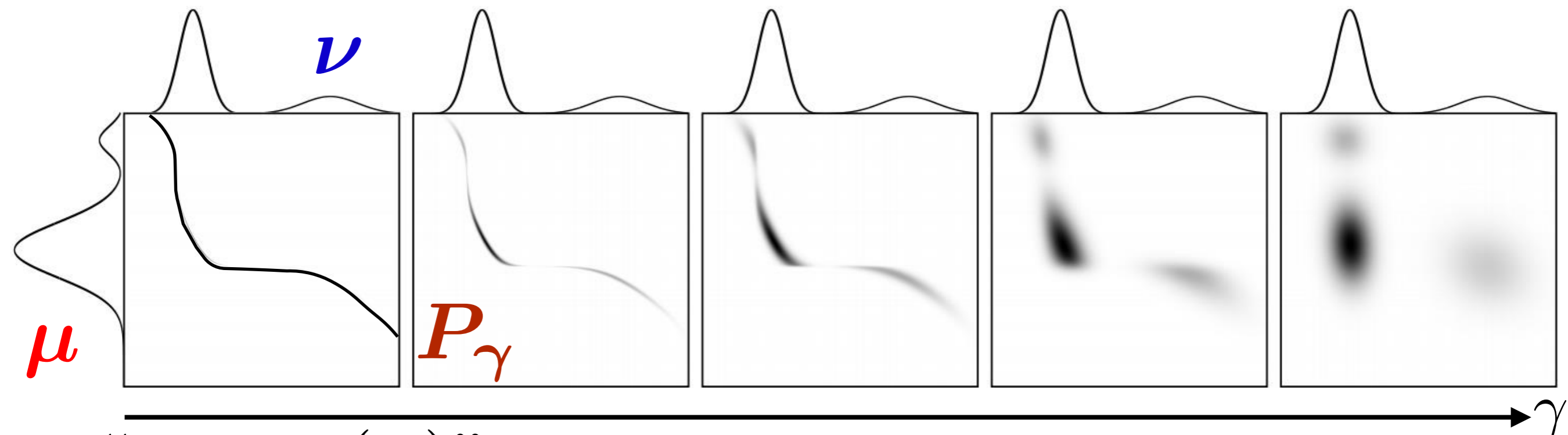


Note: Unique optimal solution because of strong concavity of entropy

Entropic Regularization [Wilson'62]

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

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



\approx “ $\mathbf{y} = T(\mathbf{x})$ ”

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

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

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

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

Fast & Scalable Algorithm

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

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

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

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

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

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

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

Fast & Scalable Algorithm

Sinkhorn's Algorithm : Repeat

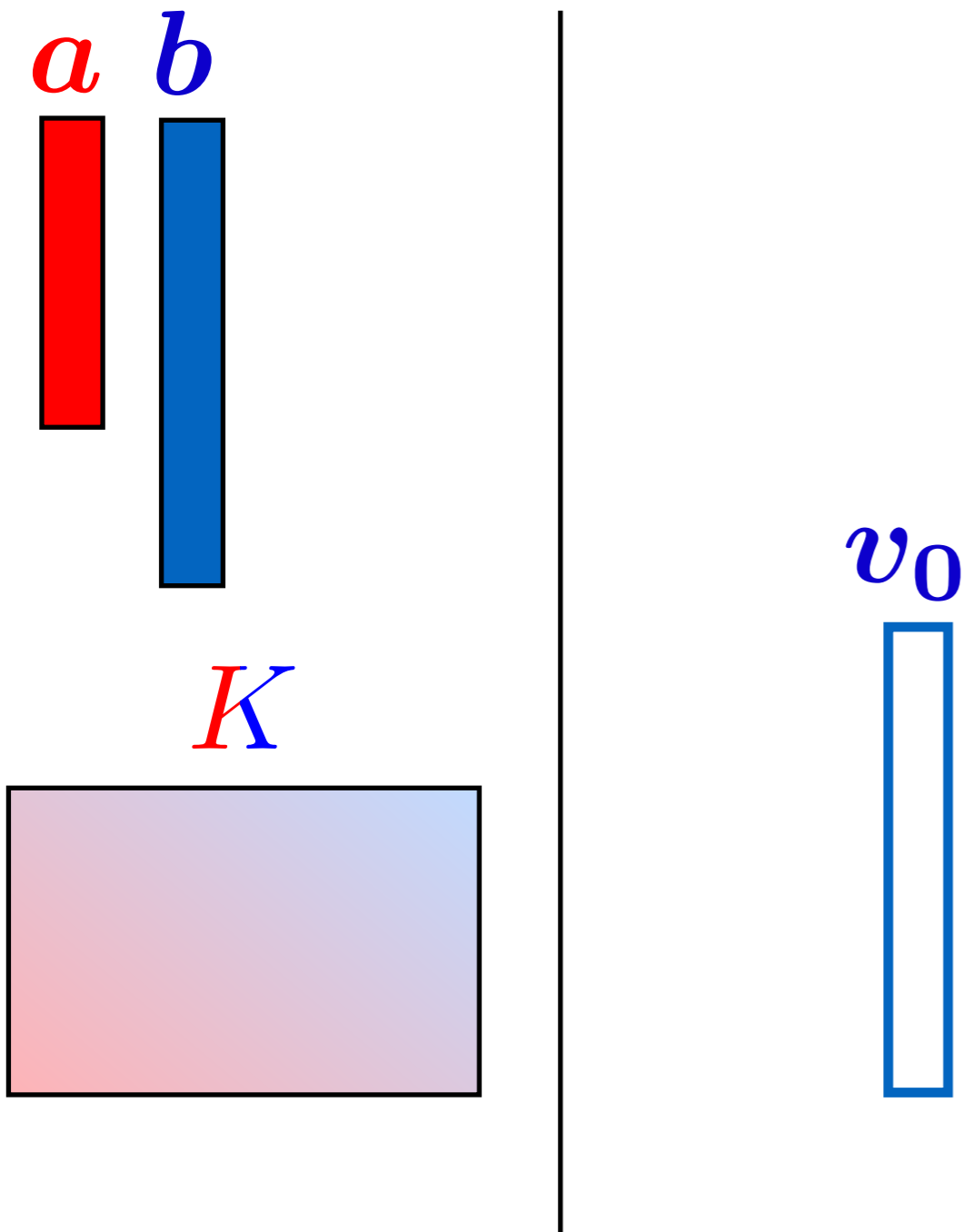
$$1. \quad \mathbf{u} = \mathbf{a} / \mathbf{K} \mathbf{v}$$

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

- [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

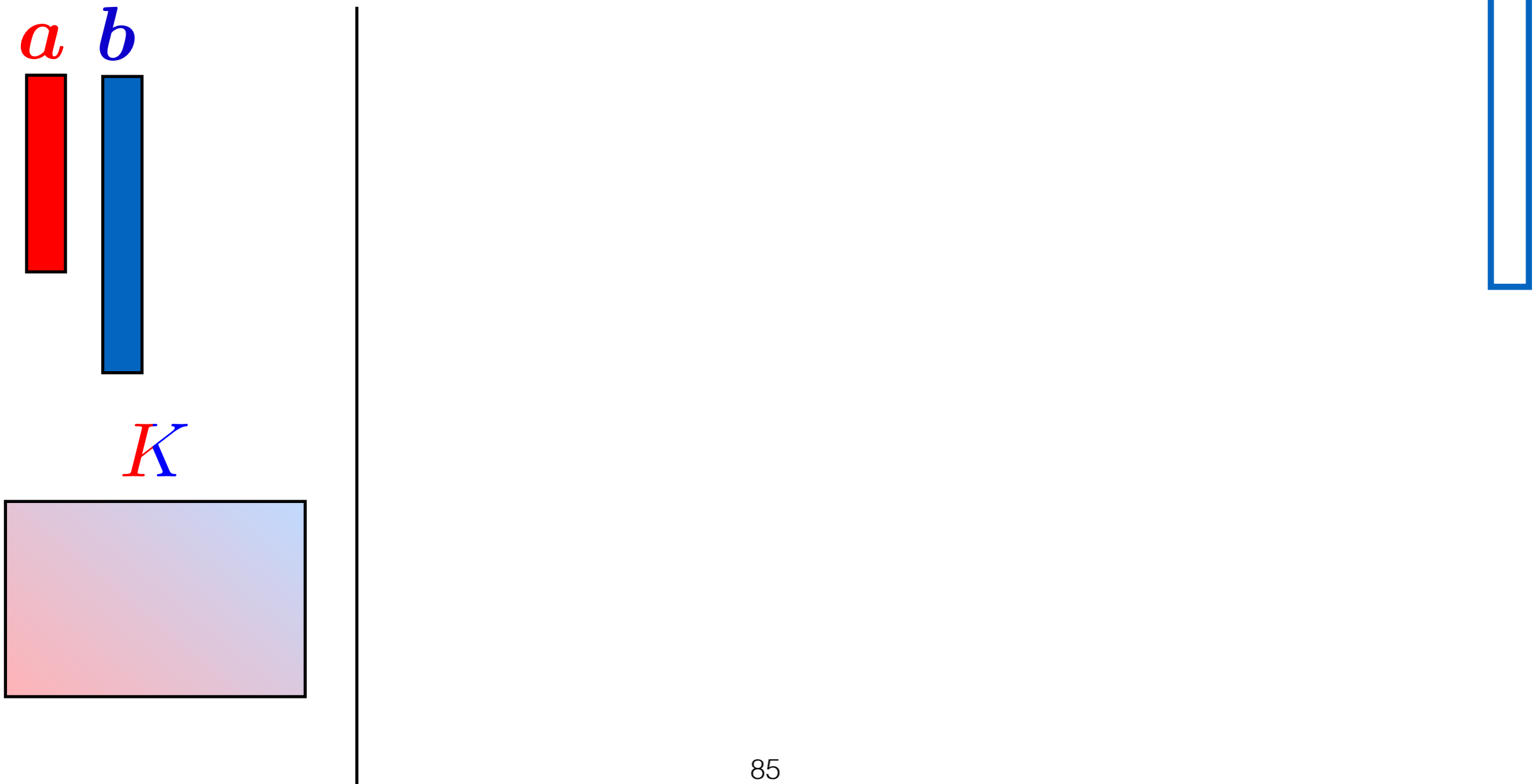
- [Sinkhorn'64] fixed-point iterations for (\mathbf{u}, \mathbf{v})
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

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

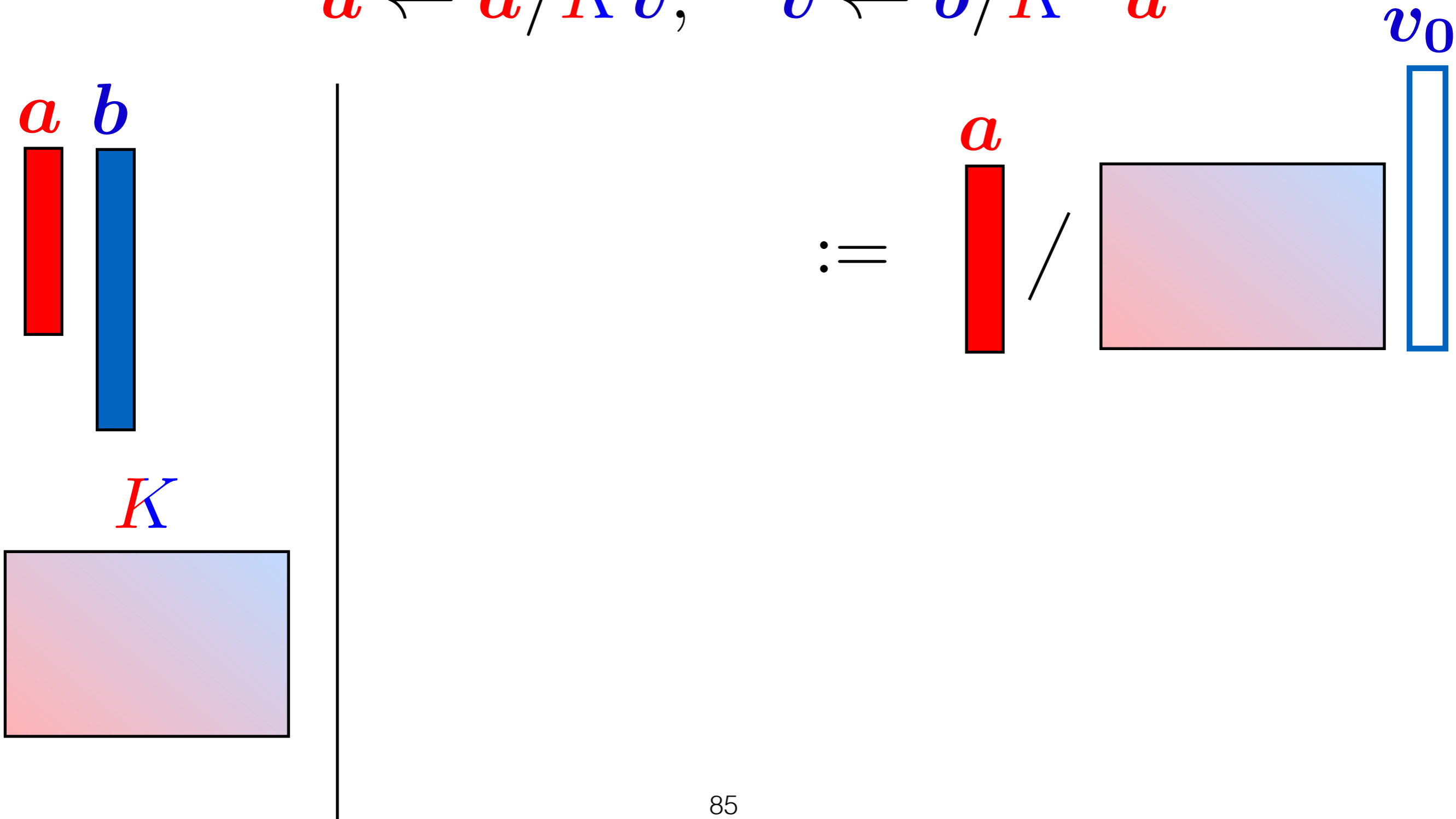
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

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

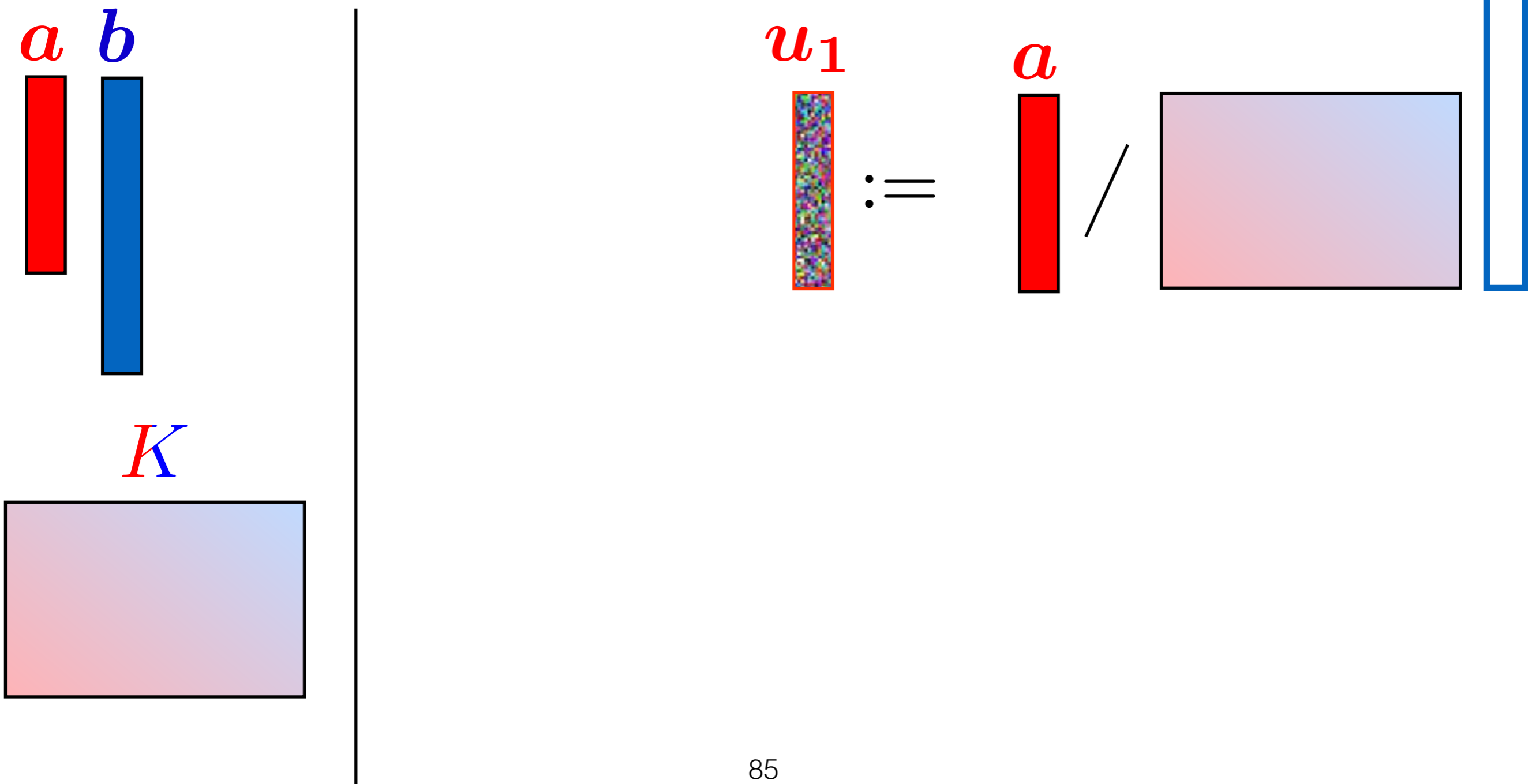
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

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

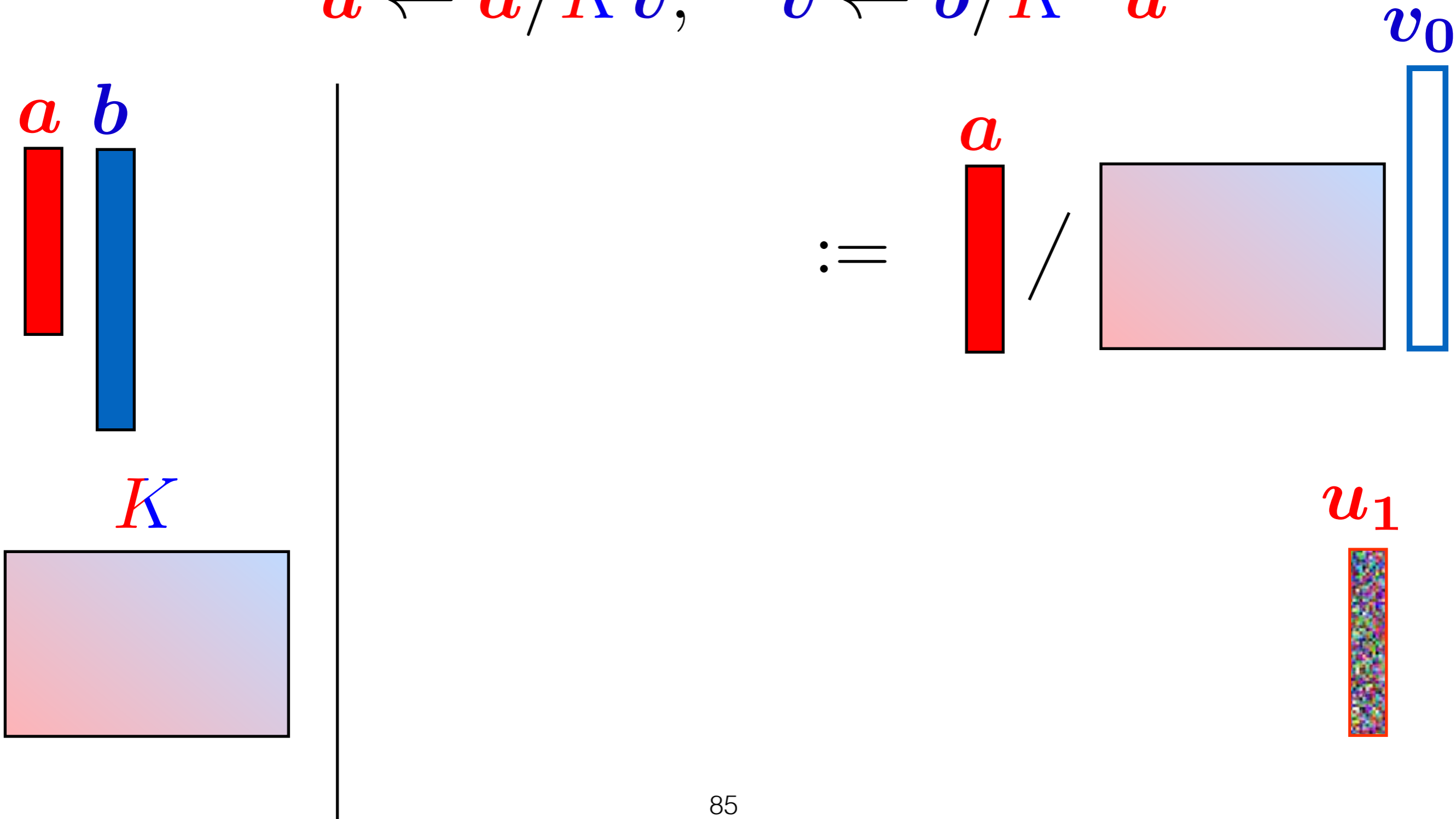
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

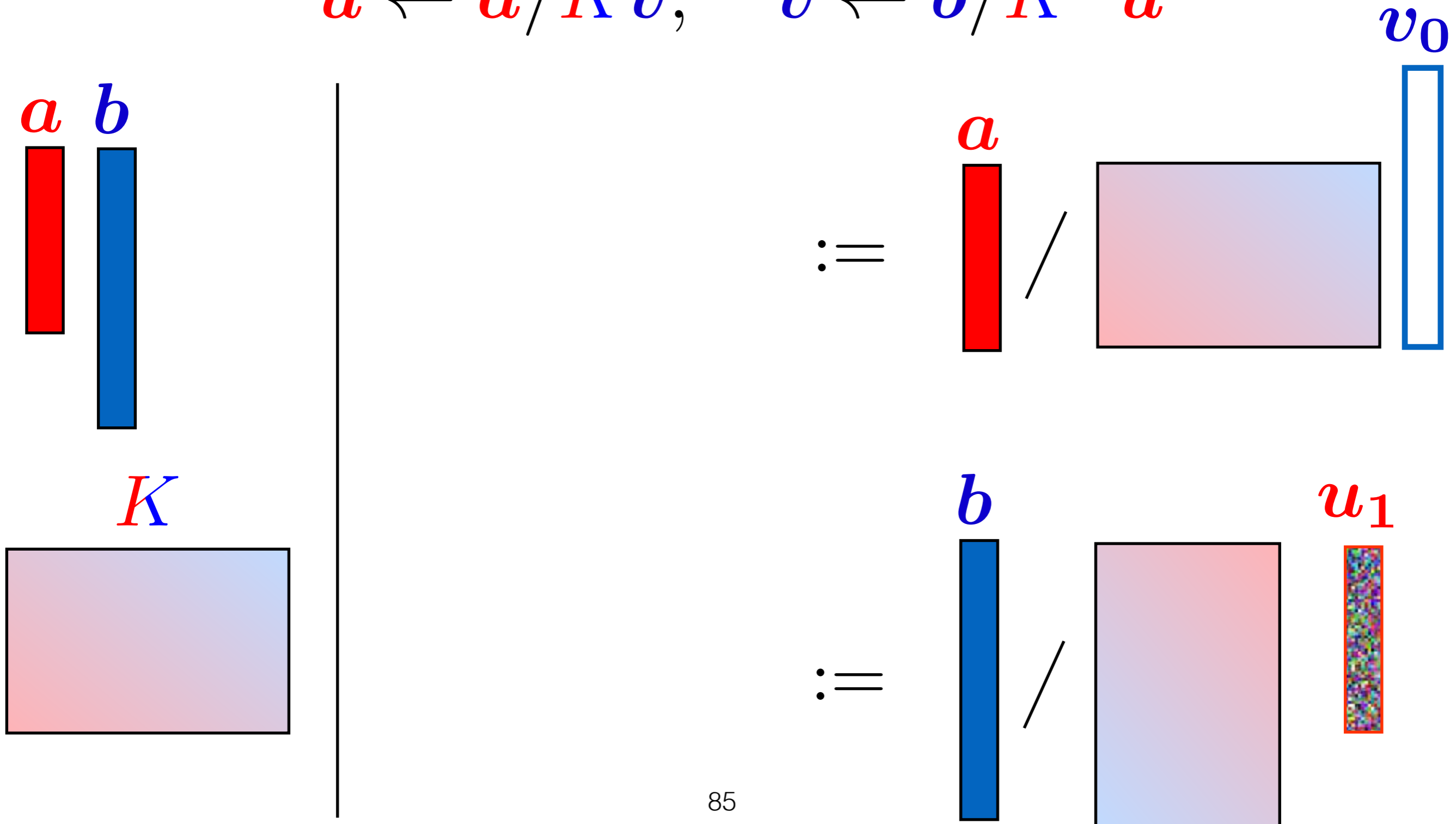
$$\mathbf{u} \leftarrow \mathbf{a} / K \mathbf{v}, \quad \mathbf{v} \leftarrow \mathbf{b} / K^T \mathbf{u}$$



Fast & Scalable Algorithm

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

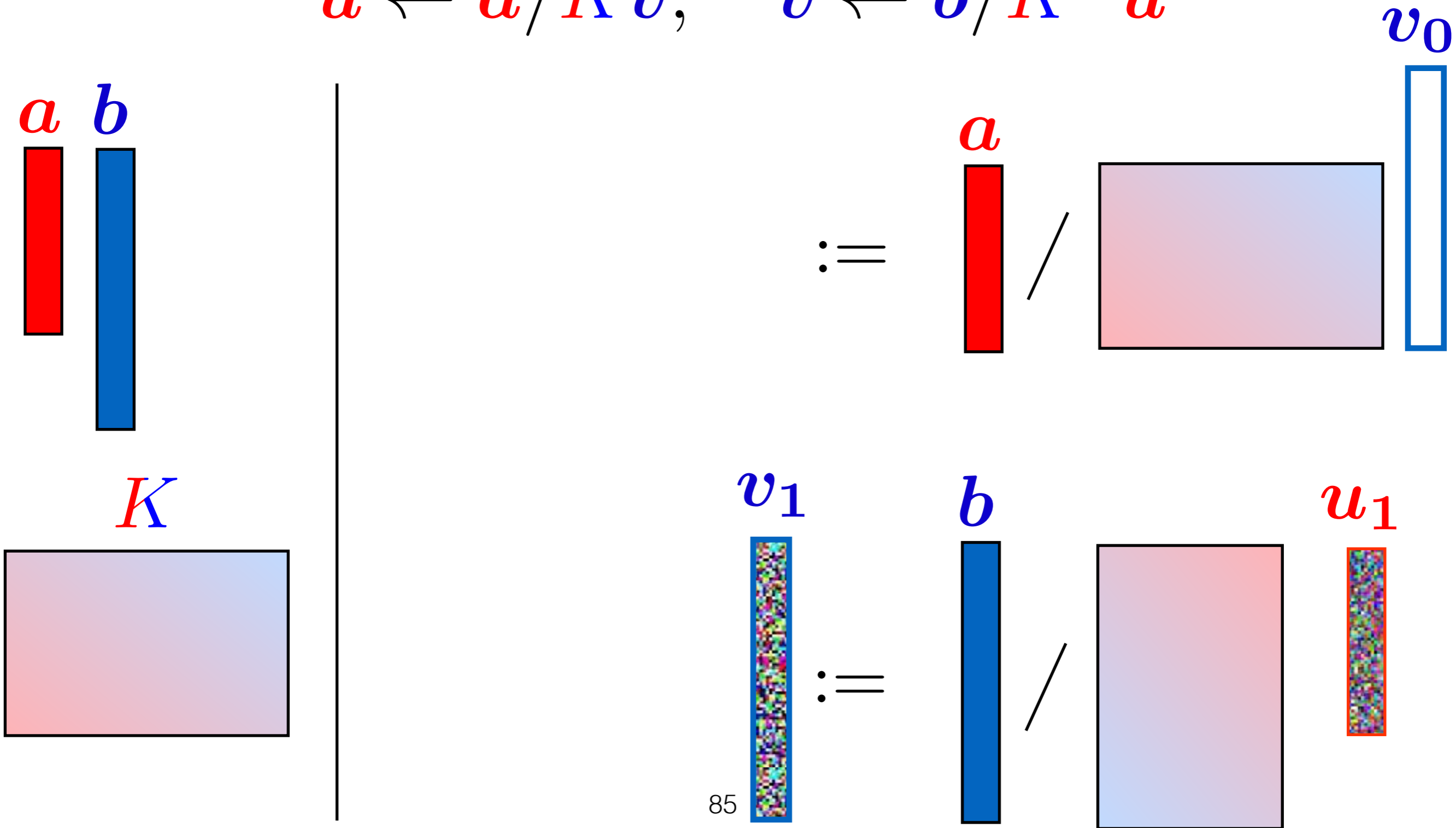
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

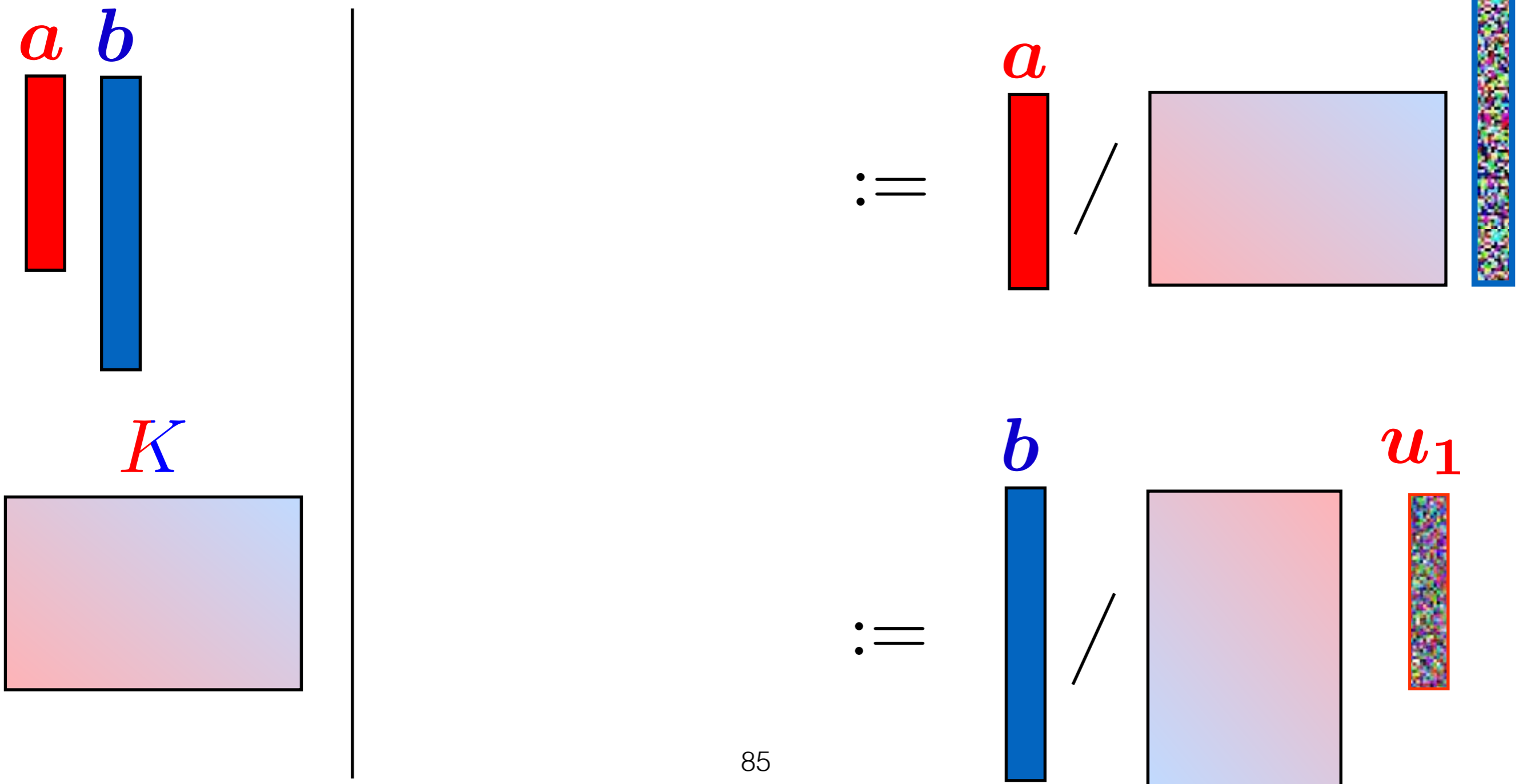
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

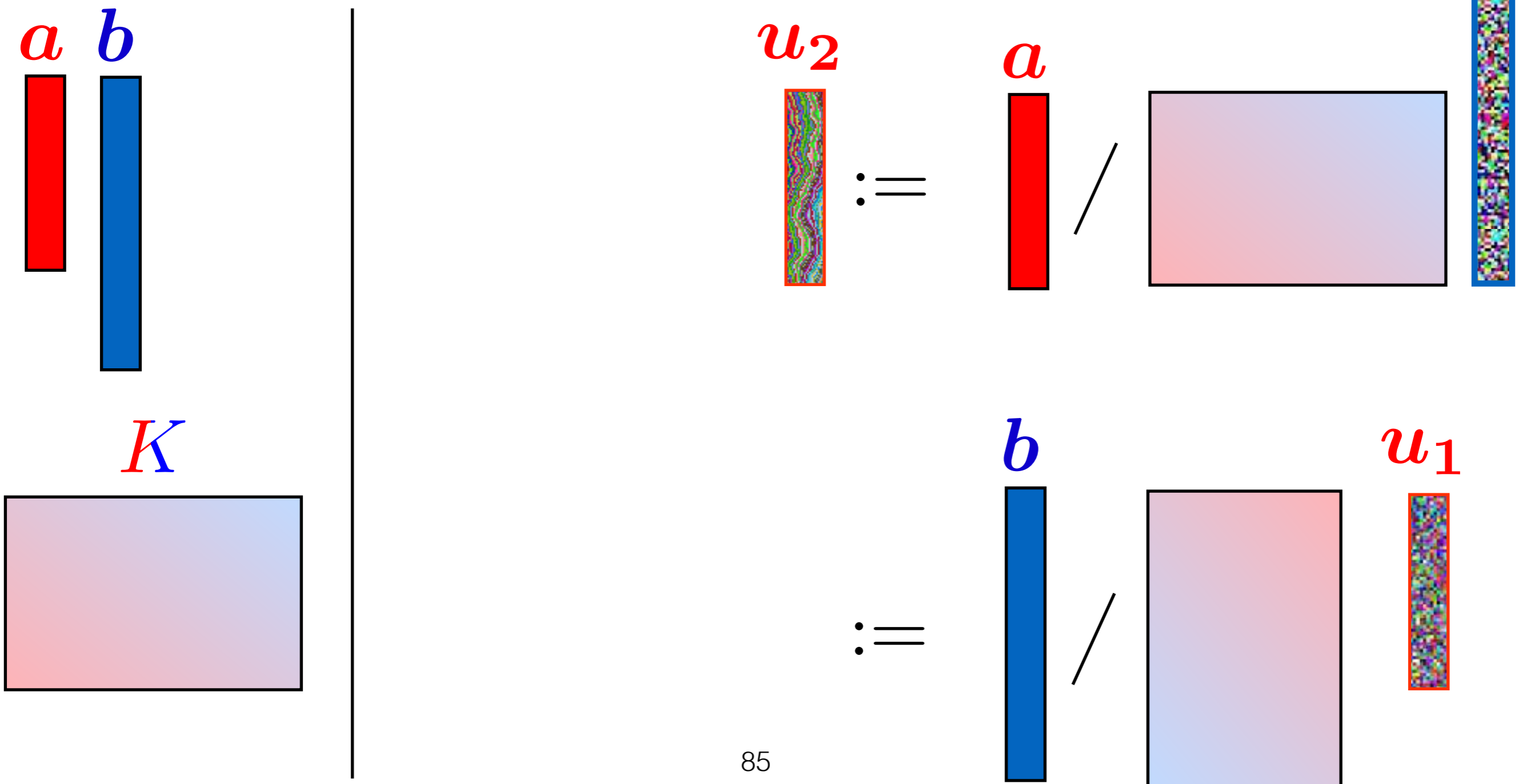
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

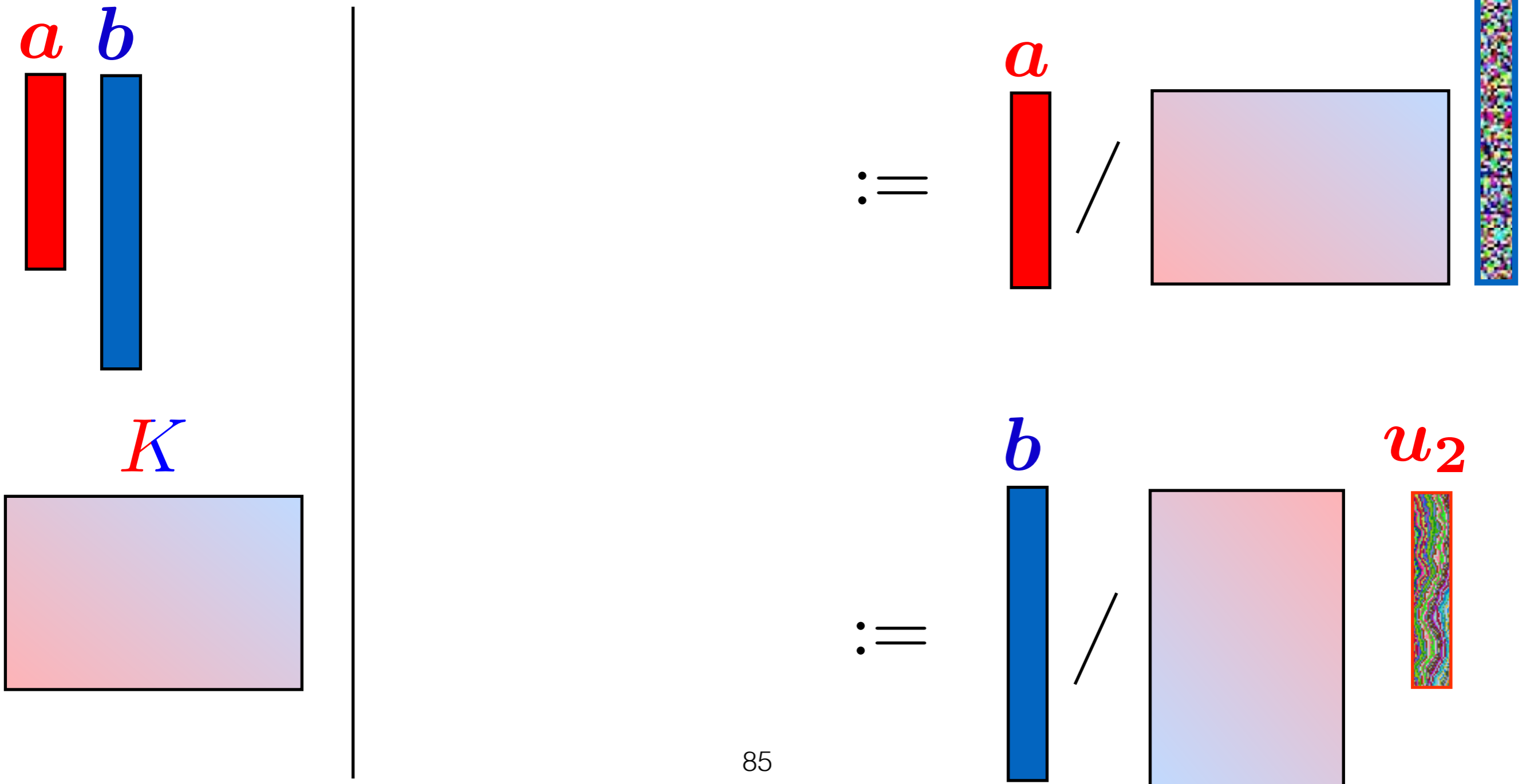
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

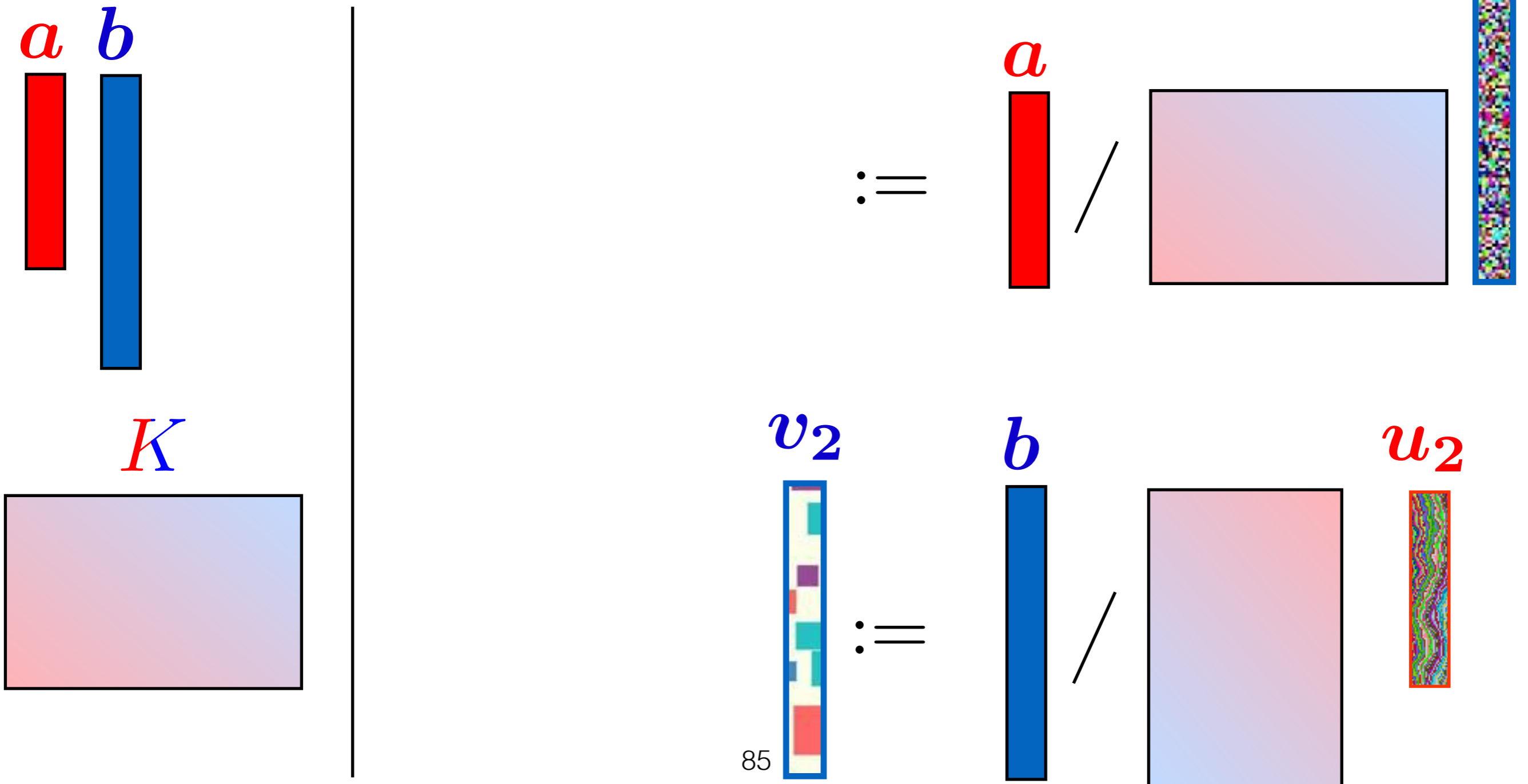
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

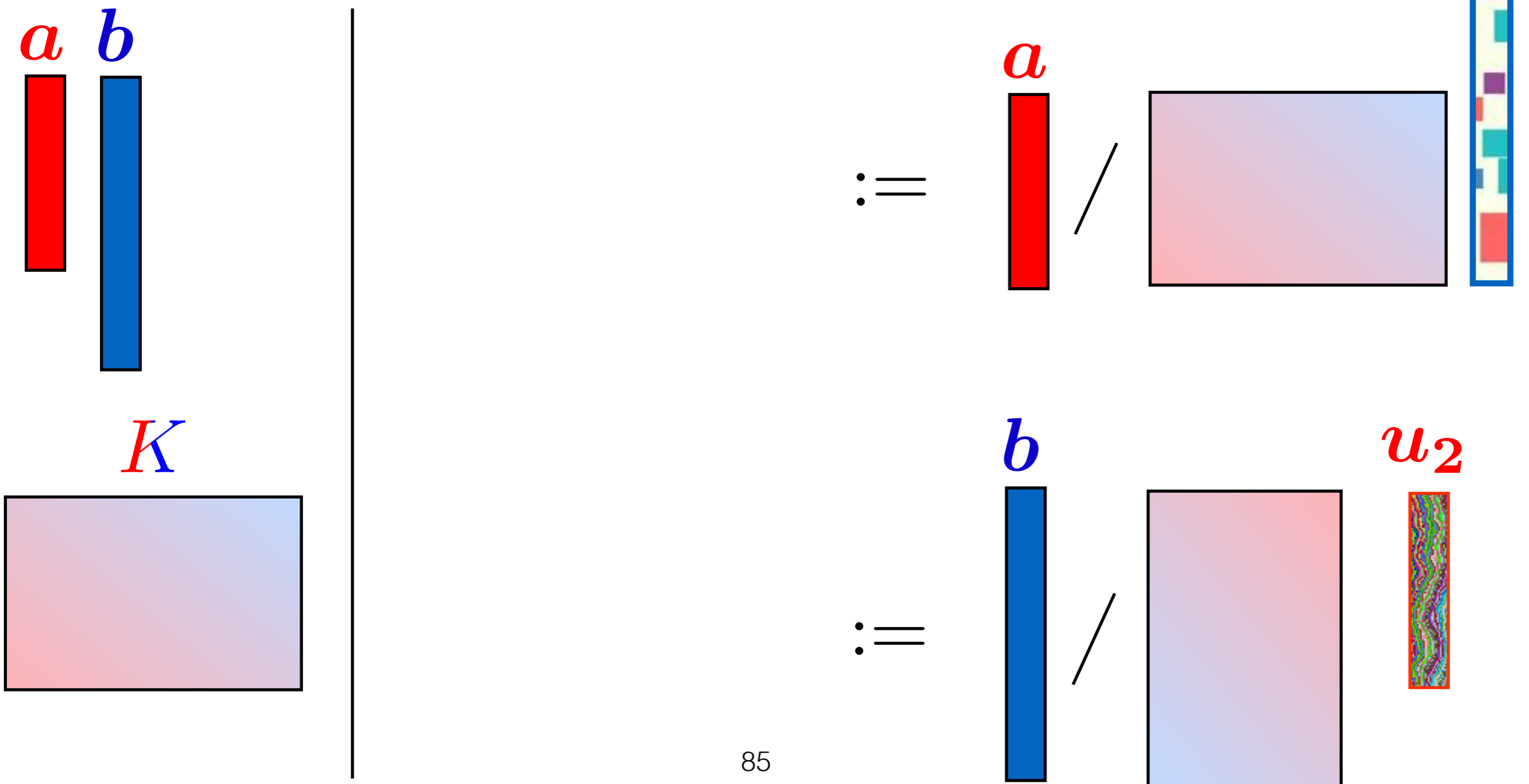
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

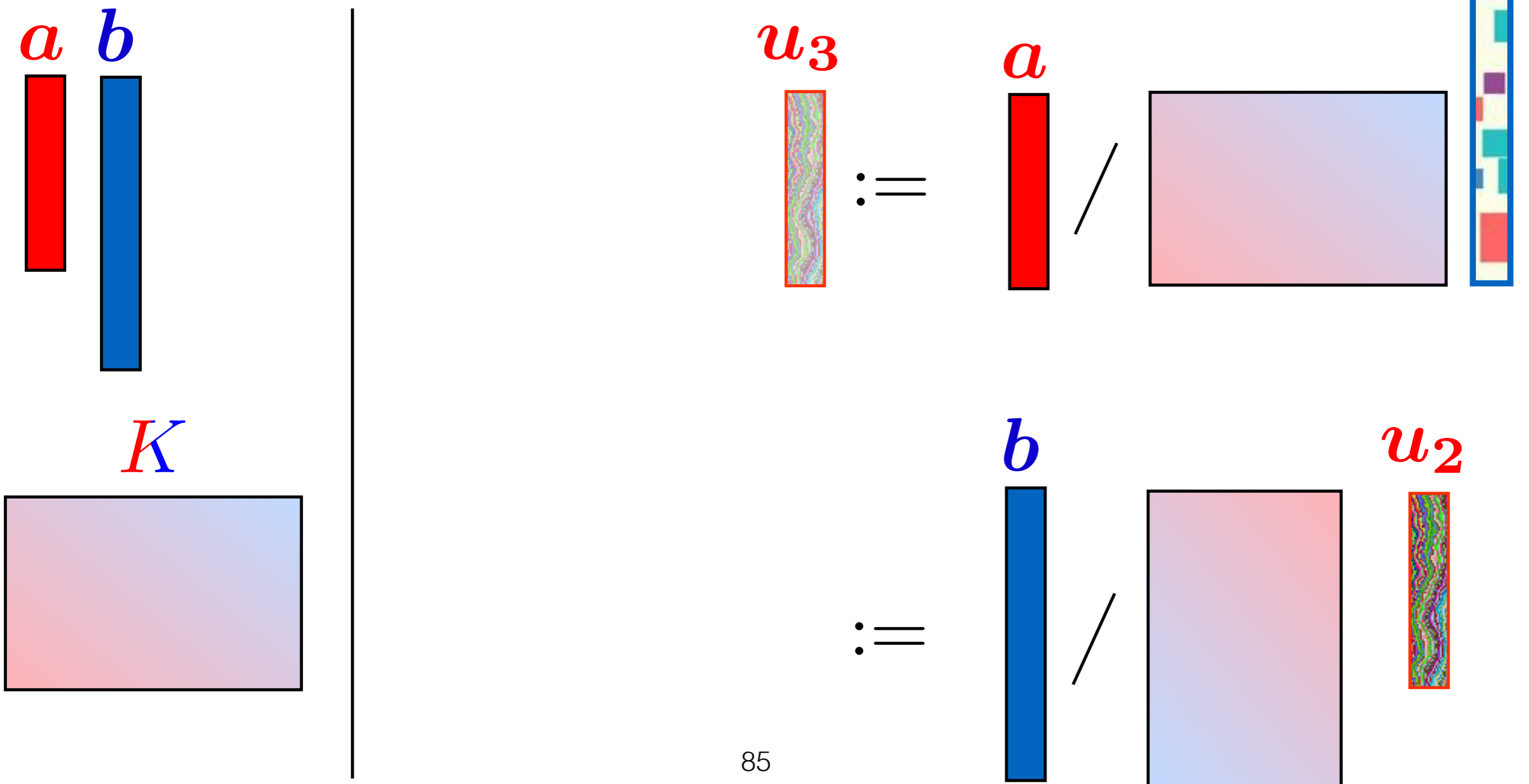
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

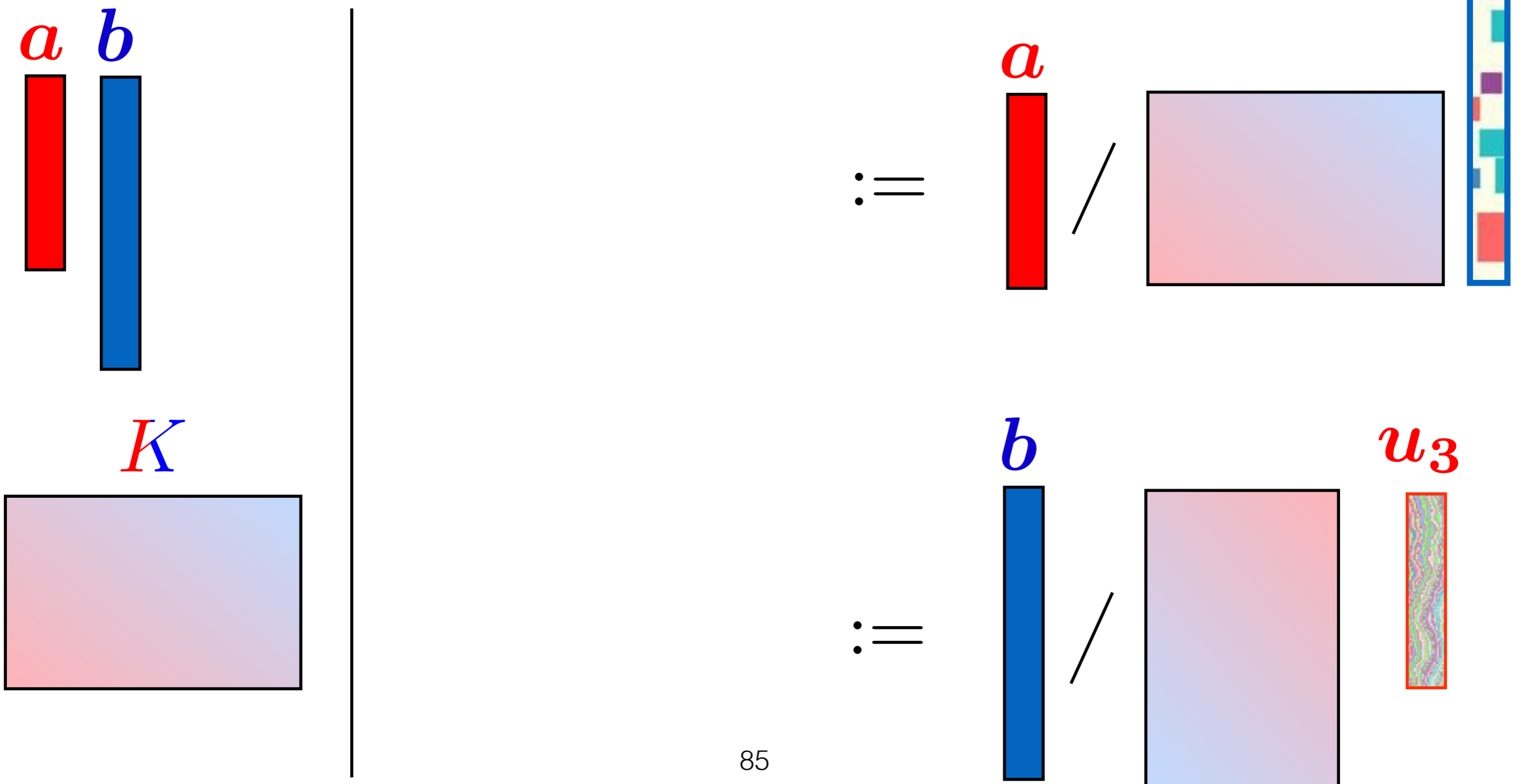
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

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

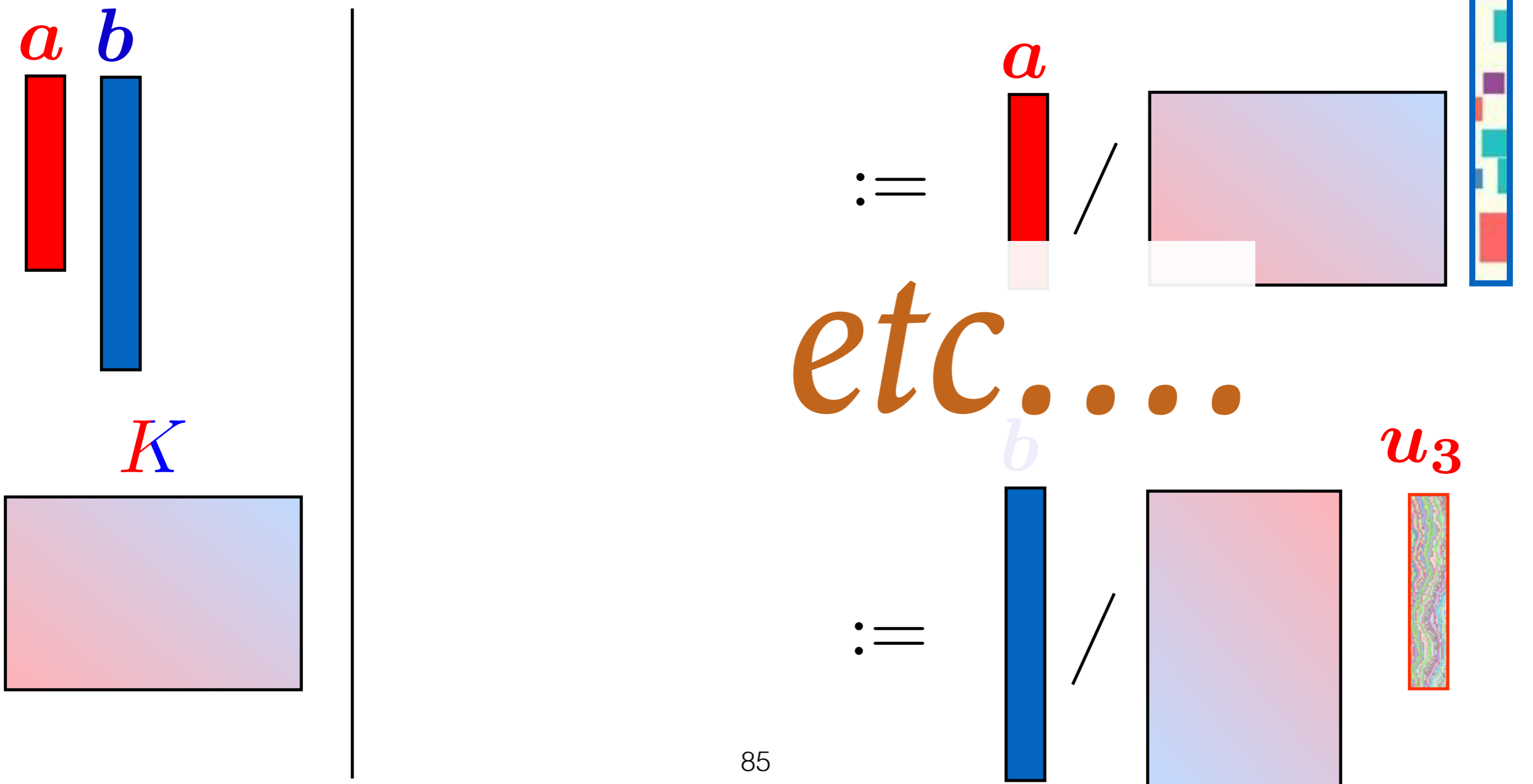
$$u \leftarrow a / K v, \quad v \leftarrow b / K^T u$$



Fast & Scalable Algorithm

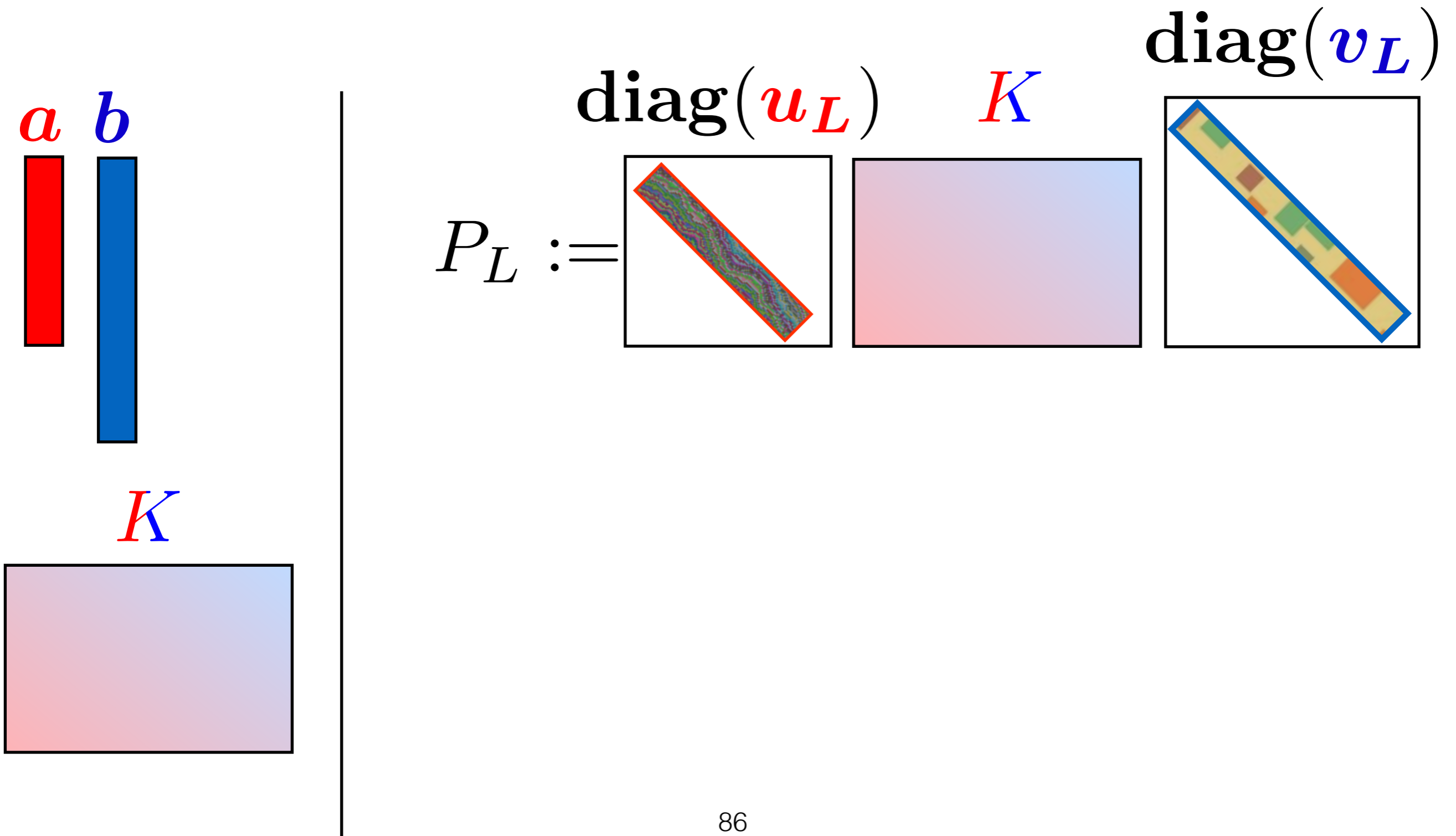
- [Sinkhorn'64] fixed-point iterations for (u, v)

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



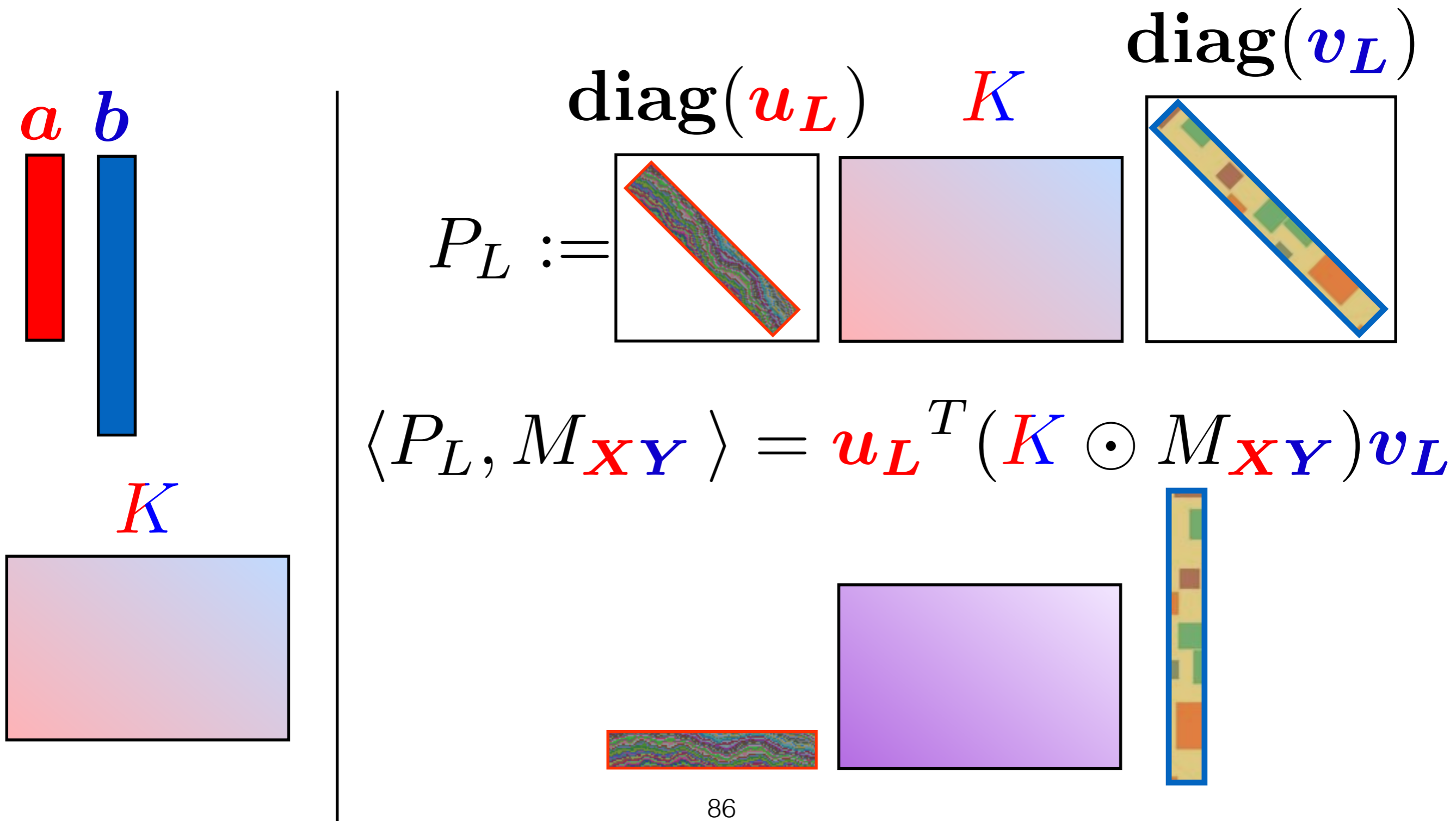
Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations.



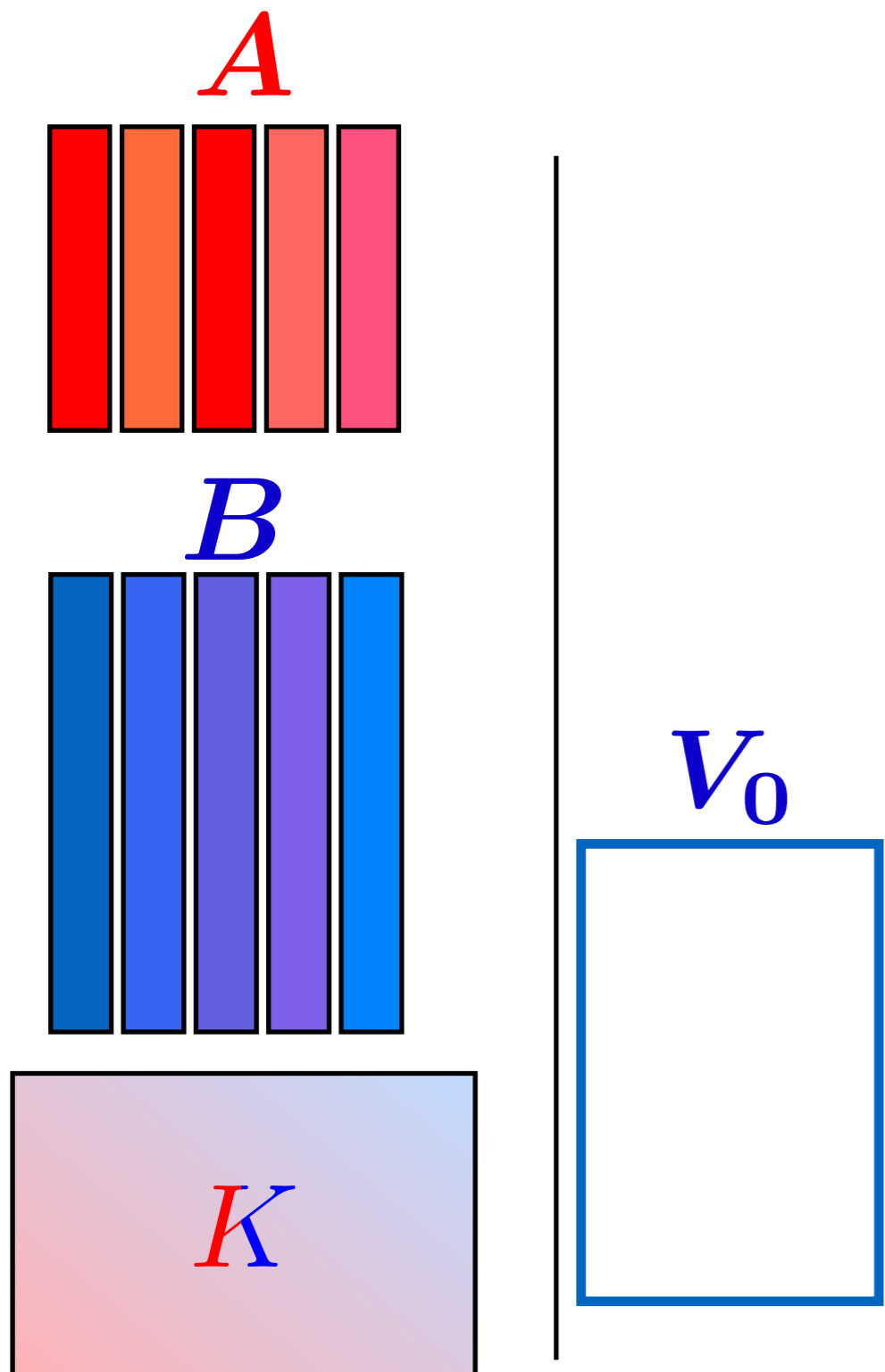
Fast & Scalable Algorithm

- [Sinkhorn'64] fixed-point iterations.



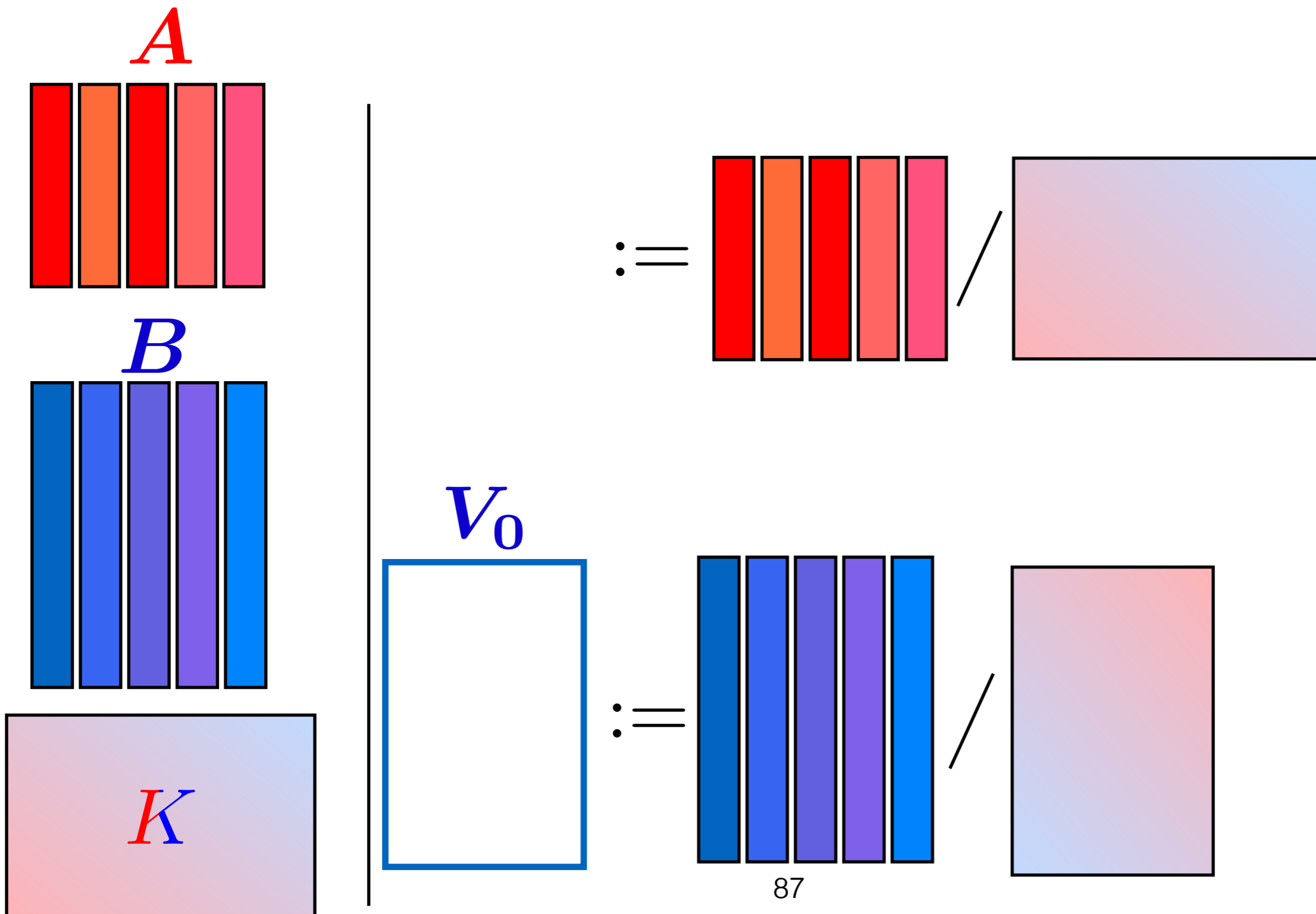
Also embarrassingly parallel

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



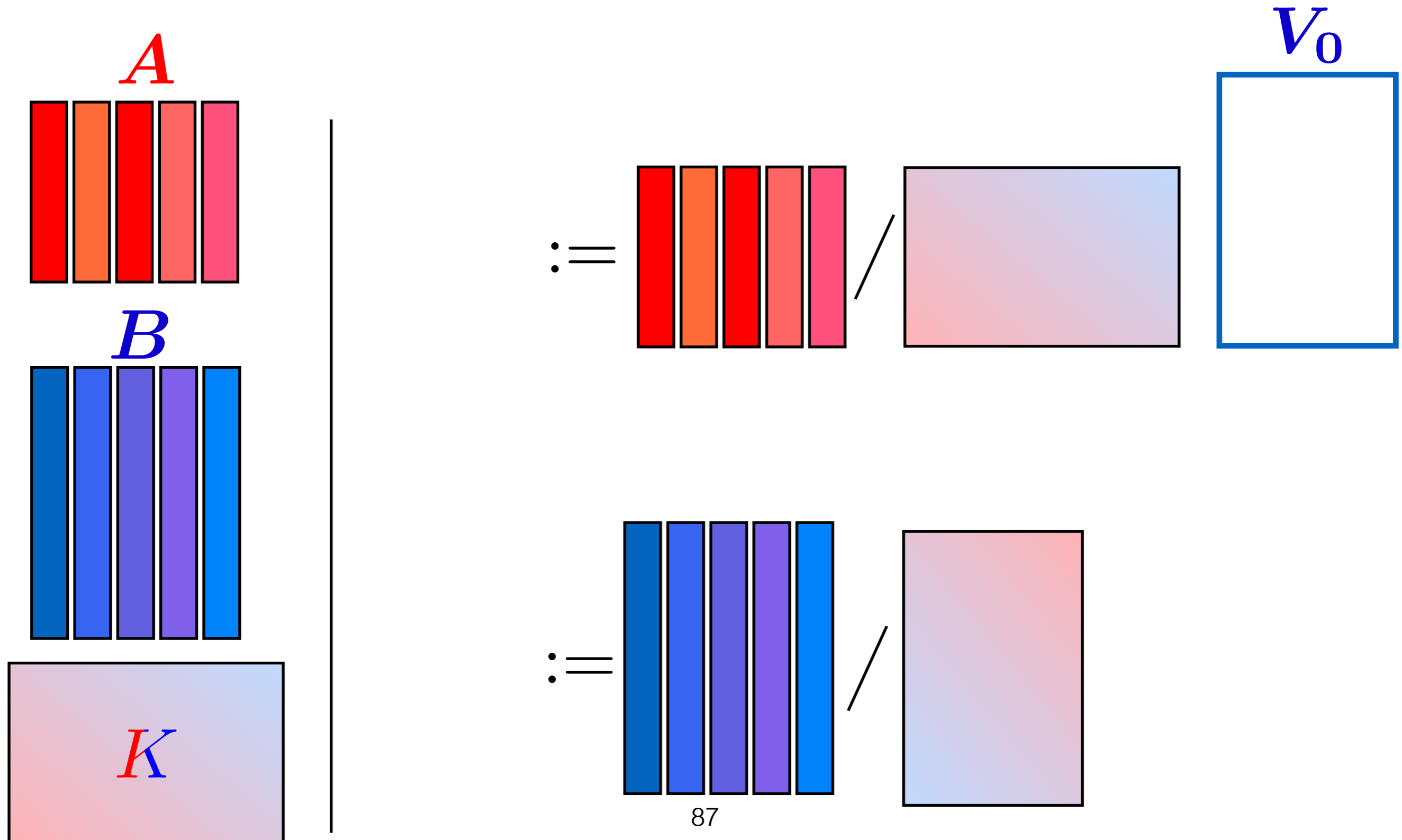
Also embarrassingly parallel

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



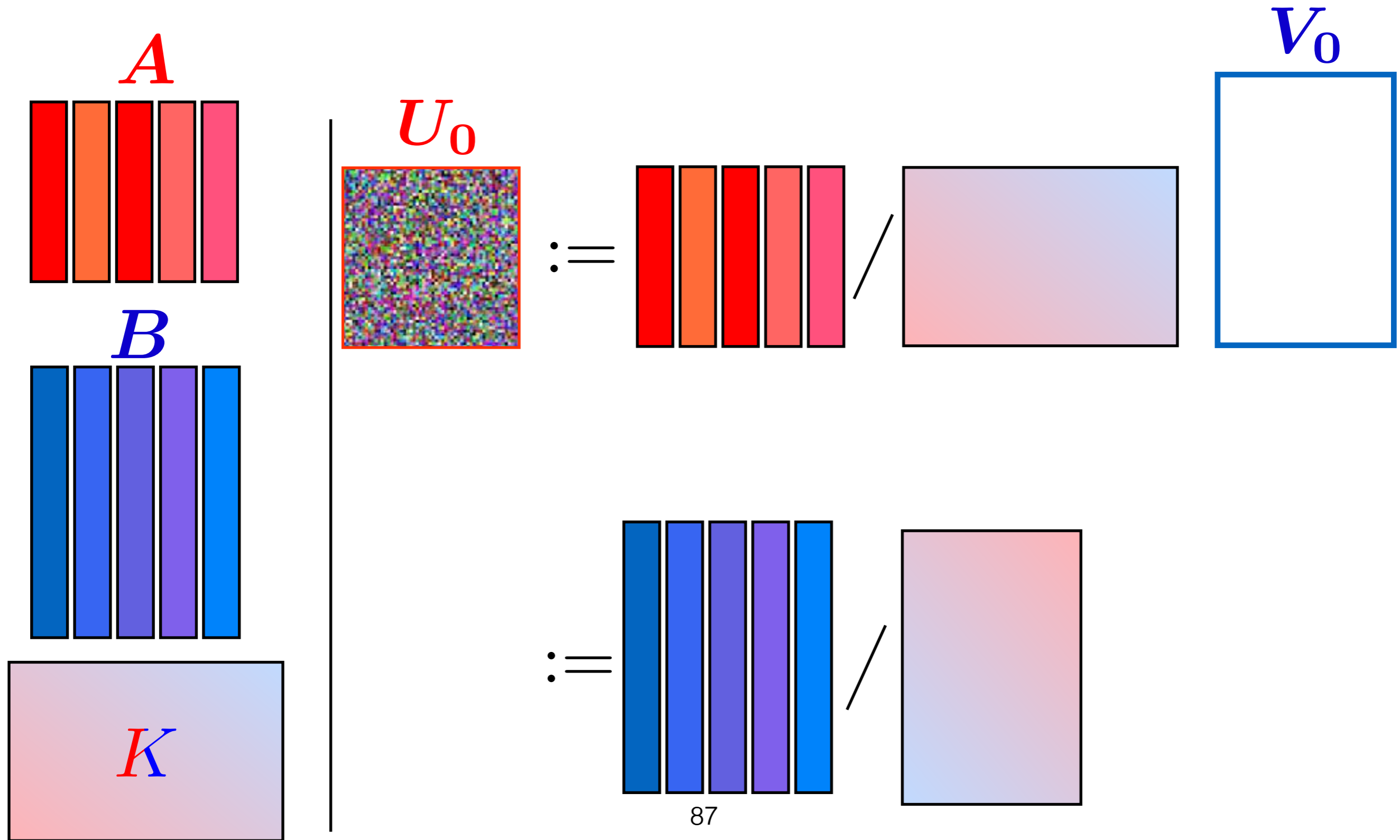
Also embarrassingly parallel

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



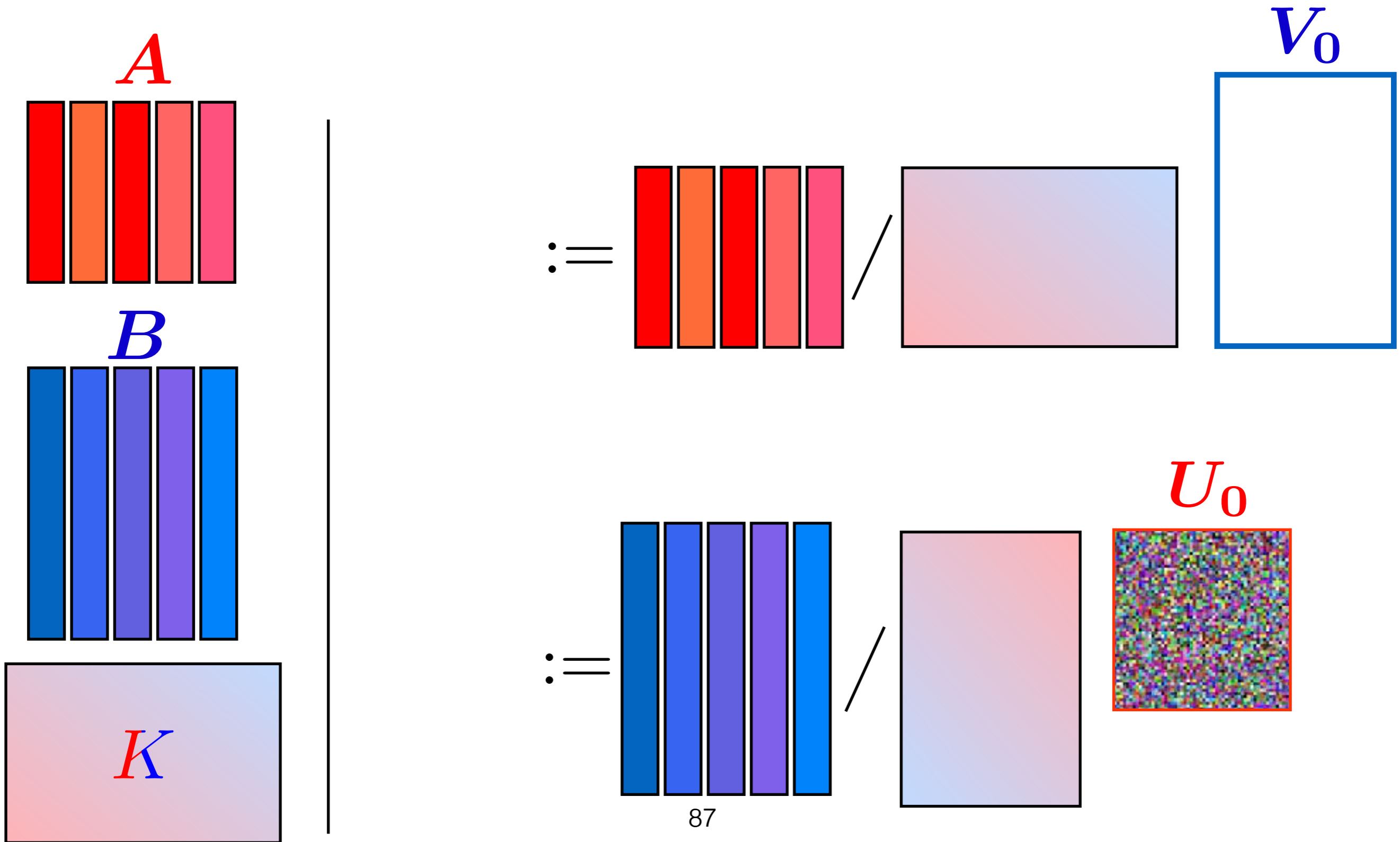
Also embarrassingly parallel

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



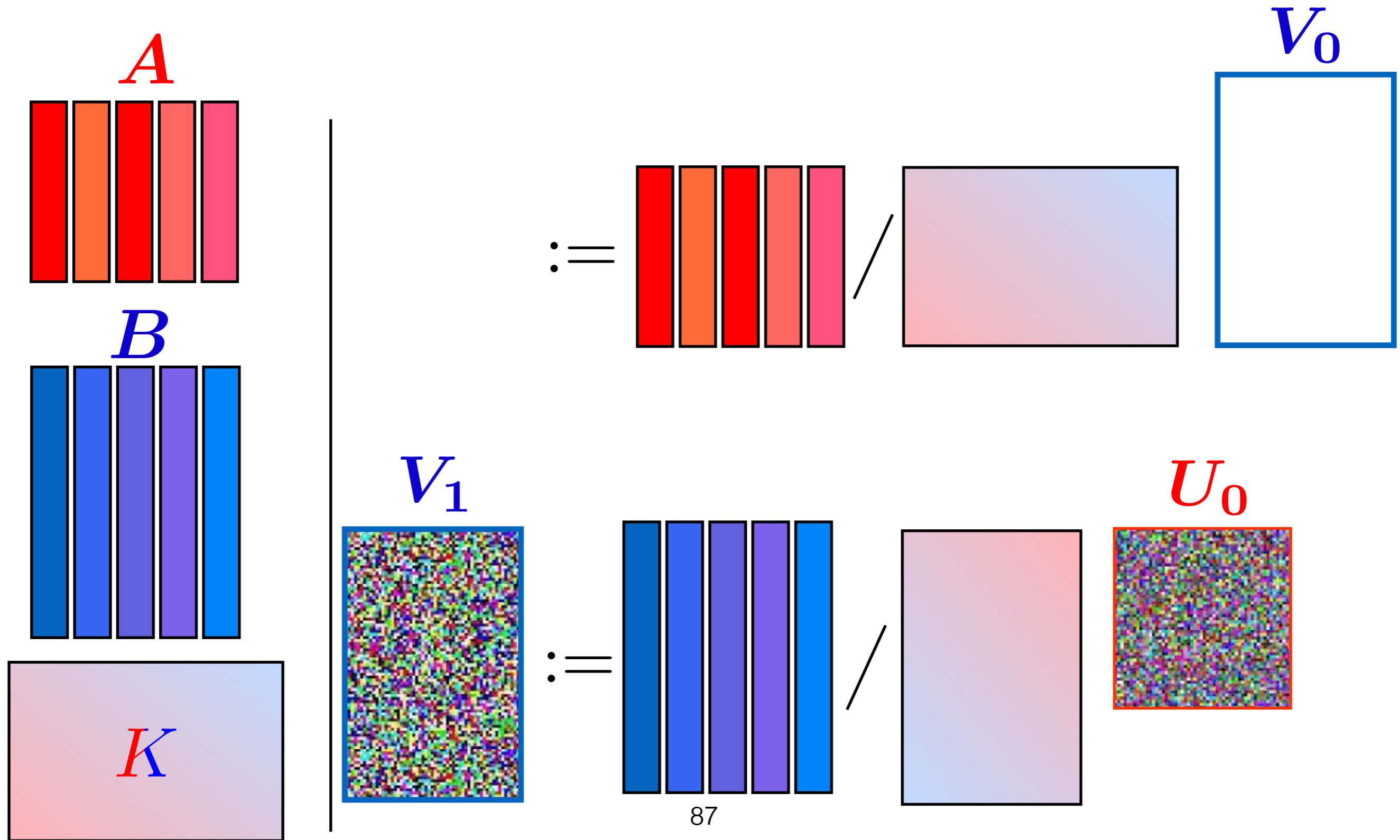
Also embarrassingly parallel

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



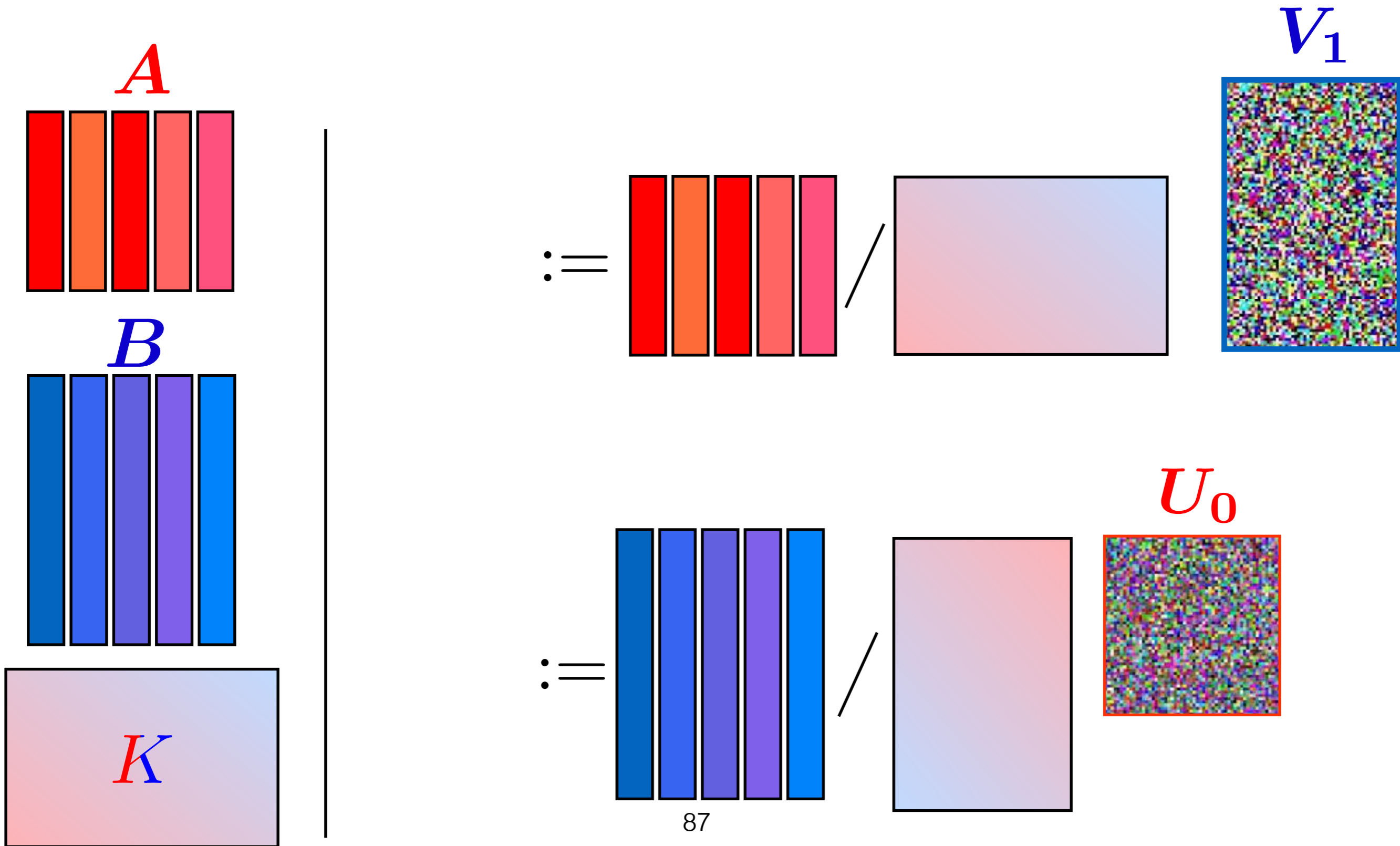
Also embarrassingly parallel

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



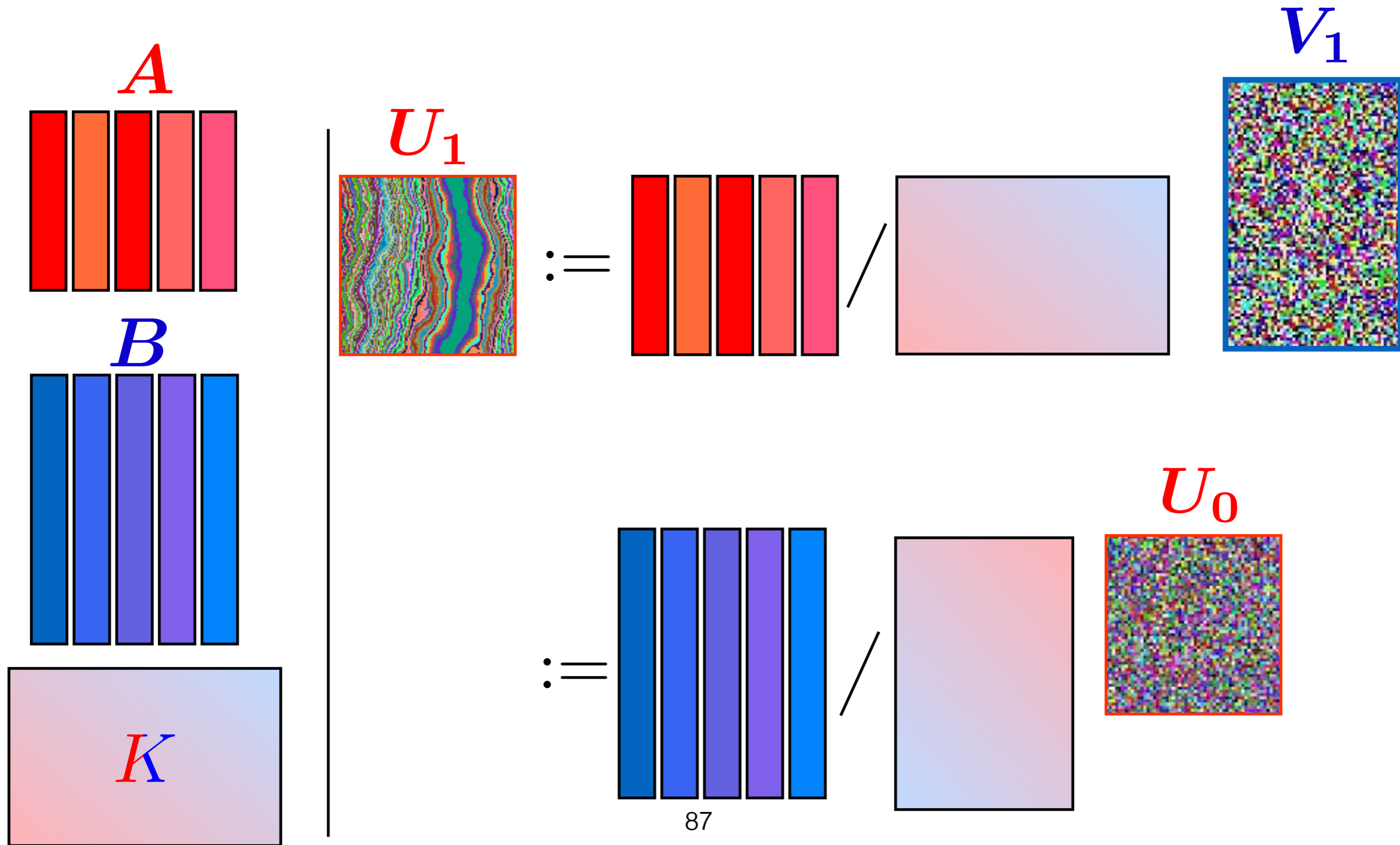
Also embarrassingly parallel

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



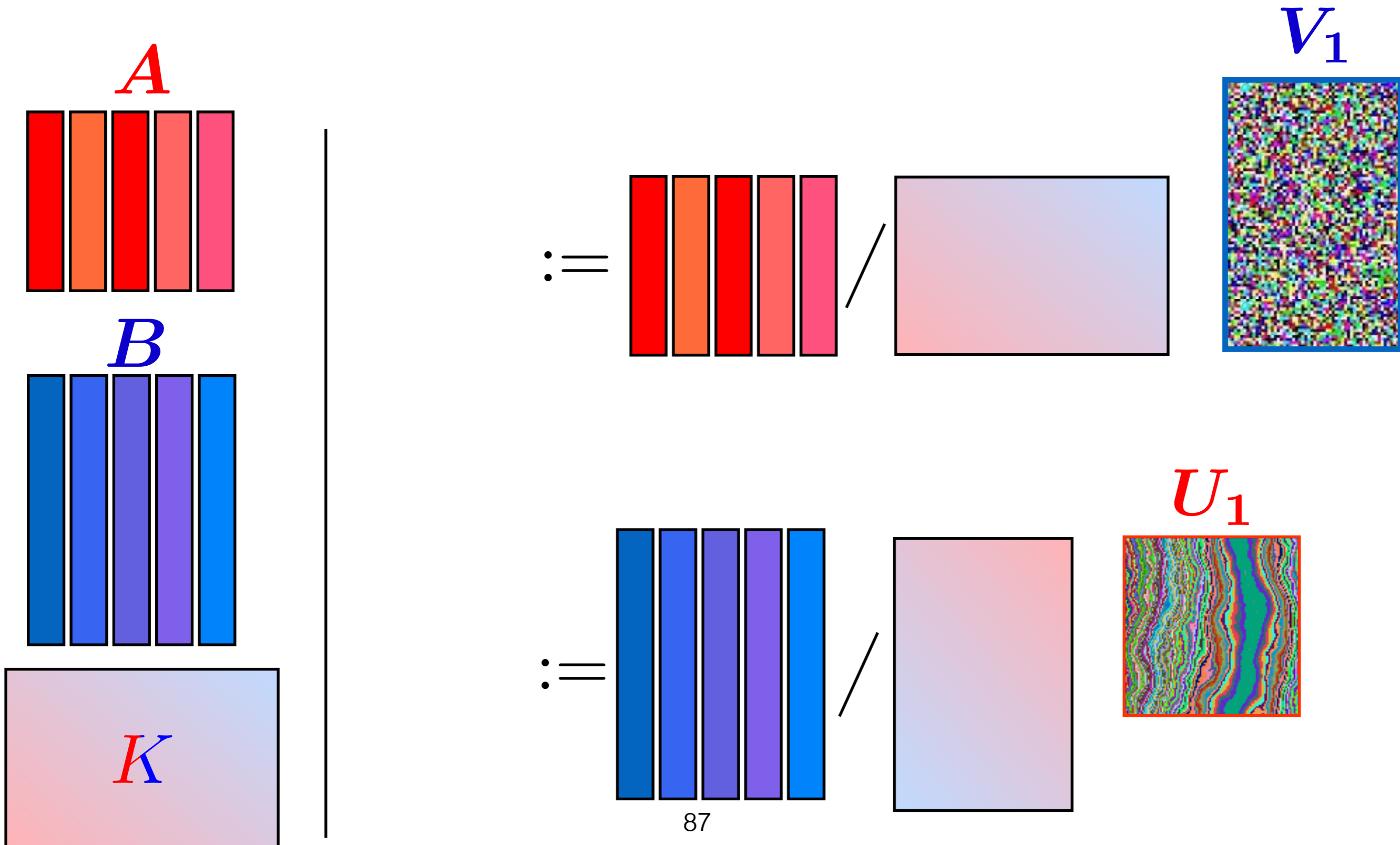
Also embarrassingly parallel

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



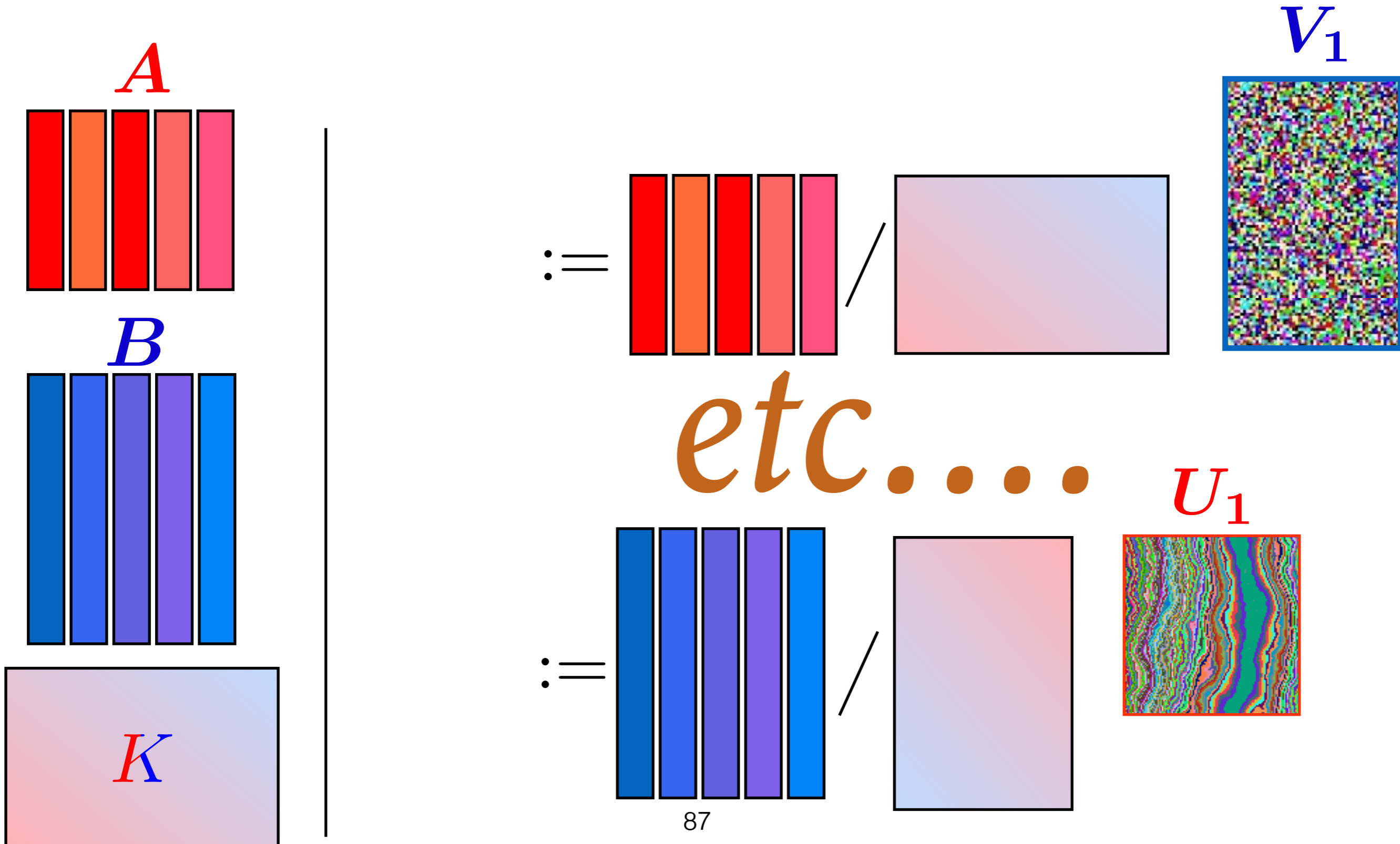
Also embarrassingly parallel

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

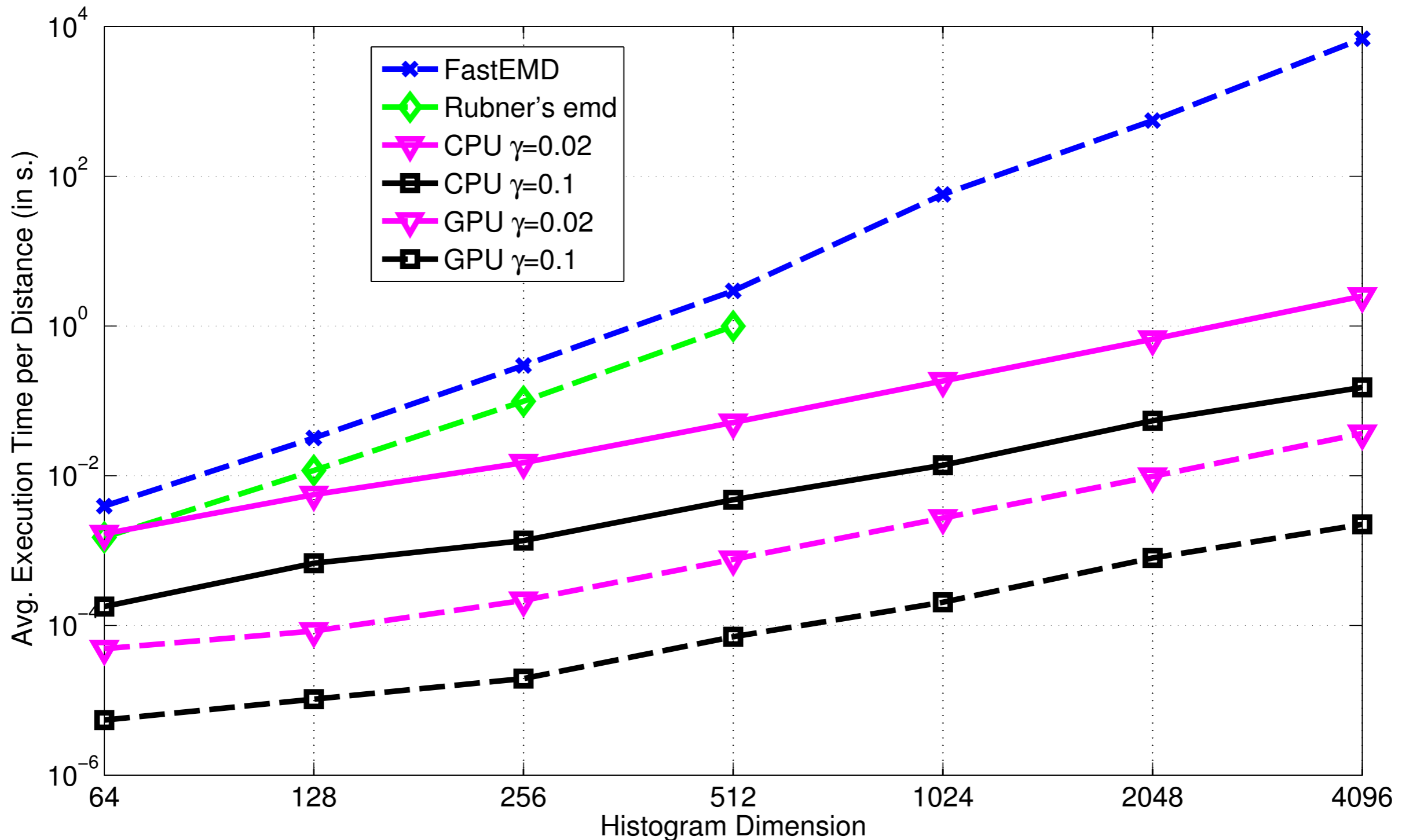


Also embarrassingly parallel

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

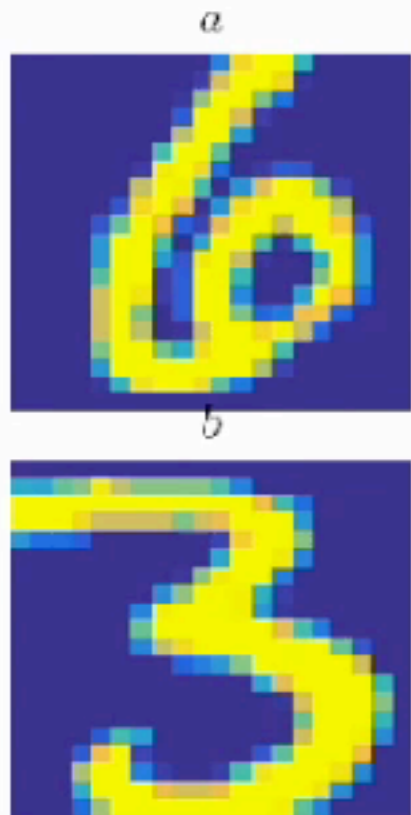


Very Fast EMD Approx. Solver

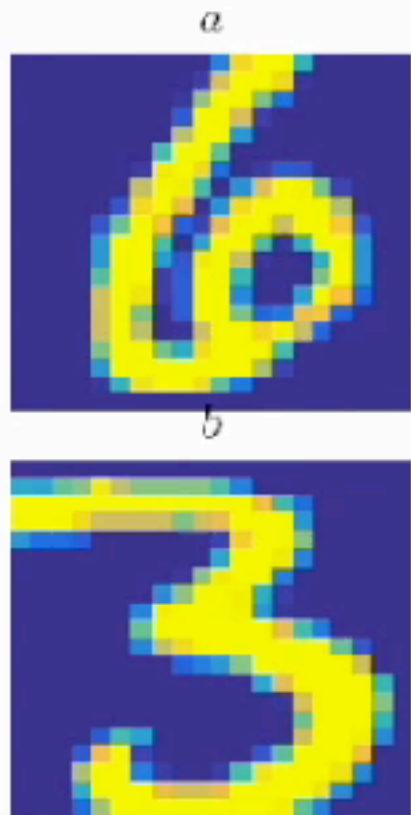


Note. (Ω, D) is a random graph with shortest path metric, histograms sampled uniformly on simplex, Sinkhorn tolerance 10^{-2} .

Very Fast EMD Approx. Solver



Very Fast EMD Approx. Solver



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

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

$$\text{where } K = \left[e^{-\frac{D^p(\mathbf{x}_i, \mathbf{y}_j)}{\gamma}} \right]_{ij}$$

REGULARIZED DISCRETE DUAL

Sinkhorn = *Block Coordinate Ascent* on Dual

Block Coordinate Ascent, *a.k.a* Sinkhorn

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

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^T \mathbf{K} e^{\boldsymbol{\beta}/\gamma}$$

Block Coordinate Ascent, *a.k.a* Sinkhorn

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

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^T \mathbf{K} e^{\boldsymbol{\beta}/\gamma}$$

$$\nabla_{\boldsymbol{\alpha}} \mathcal{E} = \mathbf{a} - e^{\boldsymbol{\alpha}/\gamma} \odot \mathbf{K} e^{\boldsymbol{\beta}/\gamma}$$

$$\nabla_{\boldsymbol{\beta}} \mathcal{E} = \mathbf{b} - e^{\boldsymbol{\beta}/\gamma} \odot \mathbf{K}^T e^{\boldsymbol{\alpha}/\gamma}$$

Block Coordinate Ascent, a.k.a Sinkhorn

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

REGULARIZED DISCRETE DUAL

$$\mathcal{E}(\boldsymbol{\alpha}, \boldsymbol{\beta}) = \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b} - \gamma (e^{\boldsymbol{\alpha}/\gamma})^T \mathbf{K} e^{\boldsymbol{\beta}/\gamma}$$

$$\nabla_{\boldsymbol{\alpha}} \mathcal{E} = \mathbf{a} - e^{\boldsymbol{\alpha}/\gamma} \odot \mathbf{K} e^{\boldsymbol{\beta}/\gamma}$$

$$\boldsymbol{\alpha} \leftarrow \gamma \left(\log \mathbf{a} - \log \mathbf{K} (e^{\boldsymbol{\beta}/\gamma}) \right)$$

$$\nabla_{\boldsymbol{\beta}} \mathcal{E} = \mathbf{b} - e^{\boldsymbol{\beta}/\gamma} \odot \mathbf{K}^T e^{\boldsymbol{\alpha}/\gamma}$$

$$\boldsymbol{\beta} \leftarrow \gamma \left(\log \mathbf{b} - \log \mathbf{K}^T (e^{\boldsymbol{\alpha}/\gamma}) \right)$$

Block Coordinate Ascent, *a.k.a* Sinkhorn

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

REGULARIZED DISCRETE DUAL

Block Coordinate Ascent, *a.k.a* Sinkhorn

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

REGULARIZED DISCRETE DUAL

$$(\boldsymbol{u}, \boldsymbol{v}) \stackrel{\text{def}}{=} (e^{\boldsymbol{\alpha}/\gamma}, e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{u} \leftarrow \frac{\boldsymbol{a}}{K \boldsymbol{v}}$$

$$\boldsymbol{v} \leftarrow \frac{\boldsymbol{b}}{K^T \boldsymbol{u}}$$

Block Coordinate Ascent, *a.k.a* Sinkhorn

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

REGULARIZED DISCRETE DUAL

$$(\boldsymbol{u}, \boldsymbol{v}) \stackrel{\text{def}}{=} (e^{\boldsymbol{\alpha}/\gamma}, e^{\boldsymbol{\beta}/\gamma})$$

$$\boldsymbol{\alpha} \leftarrow \gamma \left(\log \boldsymbol{a} - \log \boldsymbol{K} (e^{\boldsymbol{\beta}/\gamma}) \right)$$

$$\boldsymbol{u} \leftarrow \frac{\boldsymbol{a}}{\boldsymbol{K} \boldsymbol{v}}$$

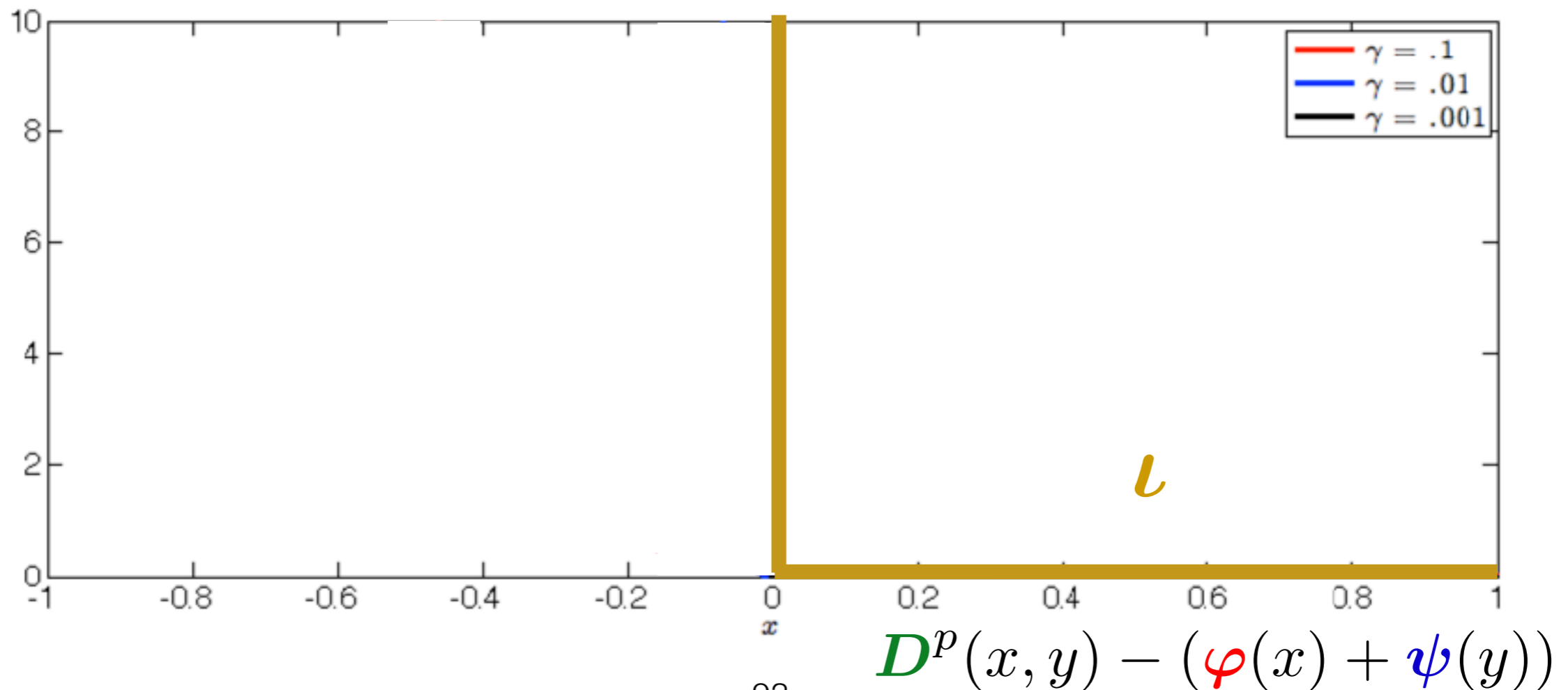
$$\boldsymbol{\beta} \leftarrow \gamma \left(\log \boldsymbol{b} - \log \boldsymbol{K}^T (e^{\boldsymbol{\alpha}/\gamma}) \right)$$

$$\boldsymbol{v} \leftarrow \frac{\boldsymbol{b}}{\boldsymbol{K}^T \boldsymbol{u}}$$

Stochastic Formulation

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi, \psi} \int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iota_C(\varphi, \psi)$$
$$C = \{(\varphi, \psi) \mid \forall x, y, \varphi(x) + \psi(y) \leq D(x, y)^p\}$$

DUAL

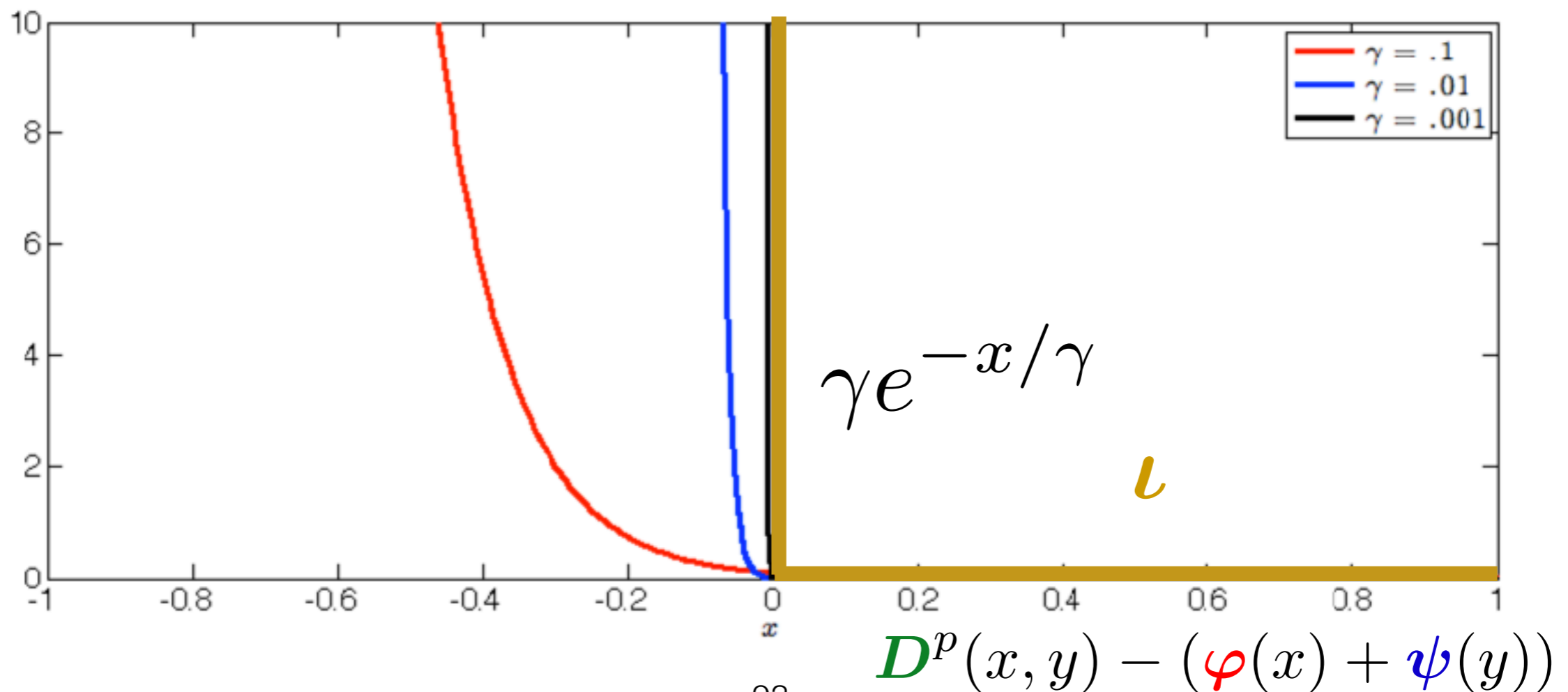


Stochastic Formulation

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi, \psi} \int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iota_C(\varphi, \psi)$$

$$C = \{(\varphi, \psi) \mid \forall x, y, \varphi(x) + \psi(y) \leq D(x, y)^p\}$$

DUAL



Stochastic Formulation

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi, \psi} \int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iota_C(\varphi, \psi)$$
$$C = \{(\varphi, \psi) \mid \forall x, y, \varphi(x) + \psi(y) \leq D(x, y)^p\}$$

DUAL

regularizing dual  *constraints* $\gamma > 0$

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

REGULARIZED DUAL

Stochastic Formulation

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \sup_{\varphi, \psi} \int \varphi d\boldsymbol{\mu} + \int \psi d\boldsymbol{\nu} - \iota_C(\varphi, \psi)$$

$$C = \{(\varphi, \psi) \mid \forall x, y, \varphi(x) + \psi(y) \leq D(x, y)^p\}$$

DUAL

regularizing dual  constraints $\gamma > 0$

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

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

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 μ 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 μ is a discrete measure?

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

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

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

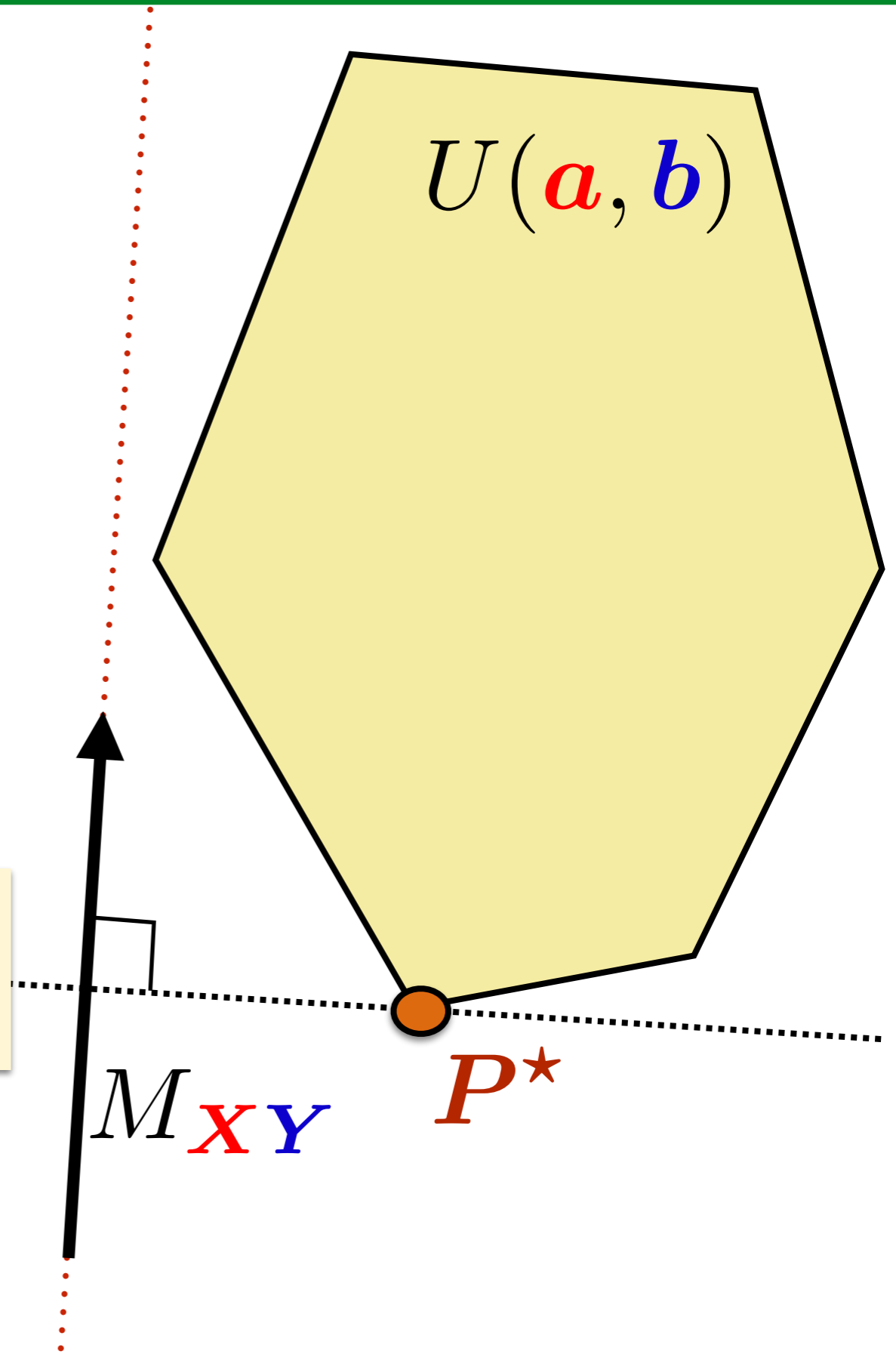
$$= \sup_{\alpha \in \mathbb{R}^n} \mathbb{E}_{\nu} [f(\alpha, y)]$$

STOCHASTIC REGULARIZED SEMI-DUAL

Sinkhorn in between W and MMD

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

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

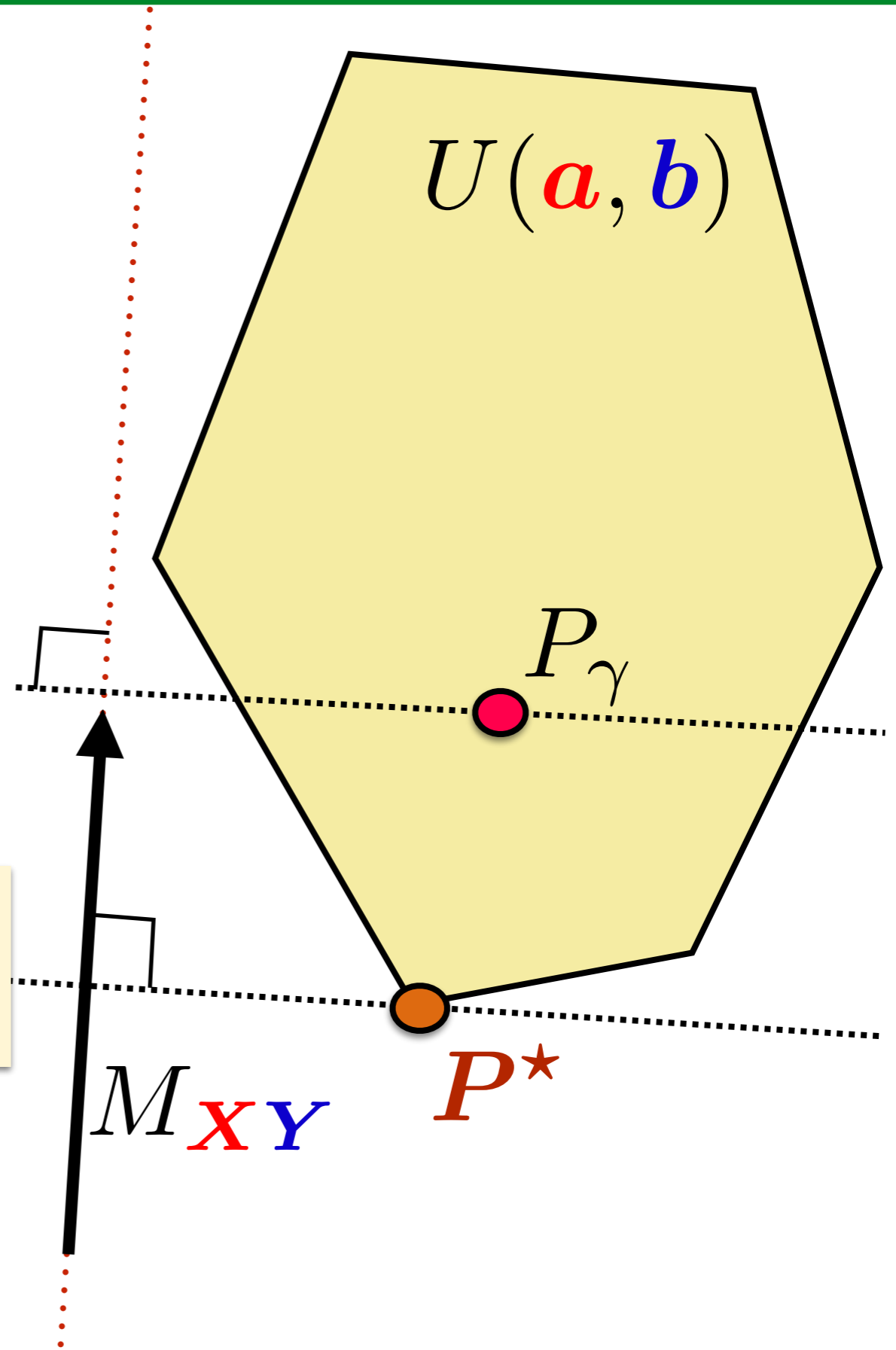


Sinkhorn in between W and MMD

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

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

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



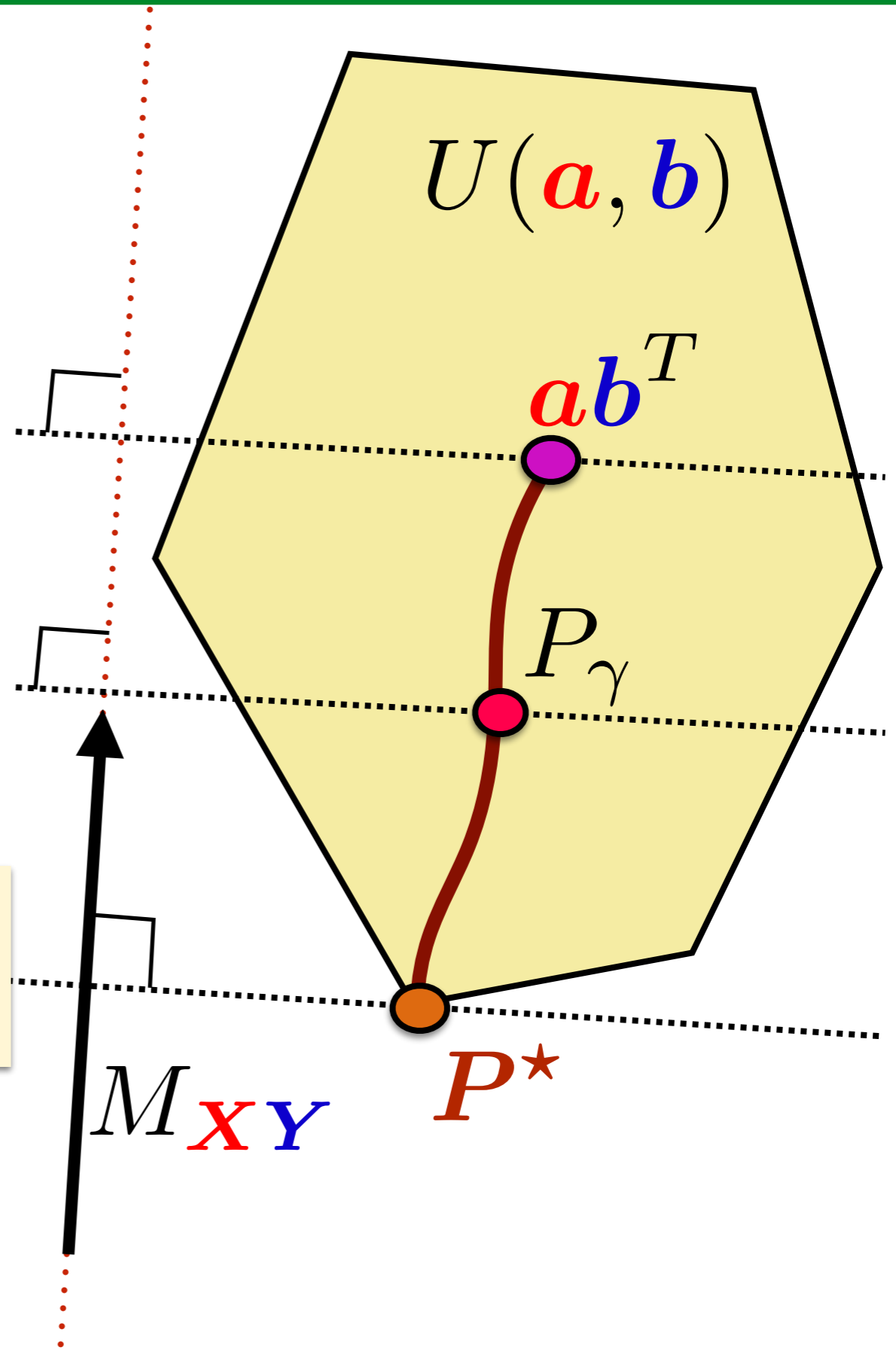
Sinkhorn in between W and MMD

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

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

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

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



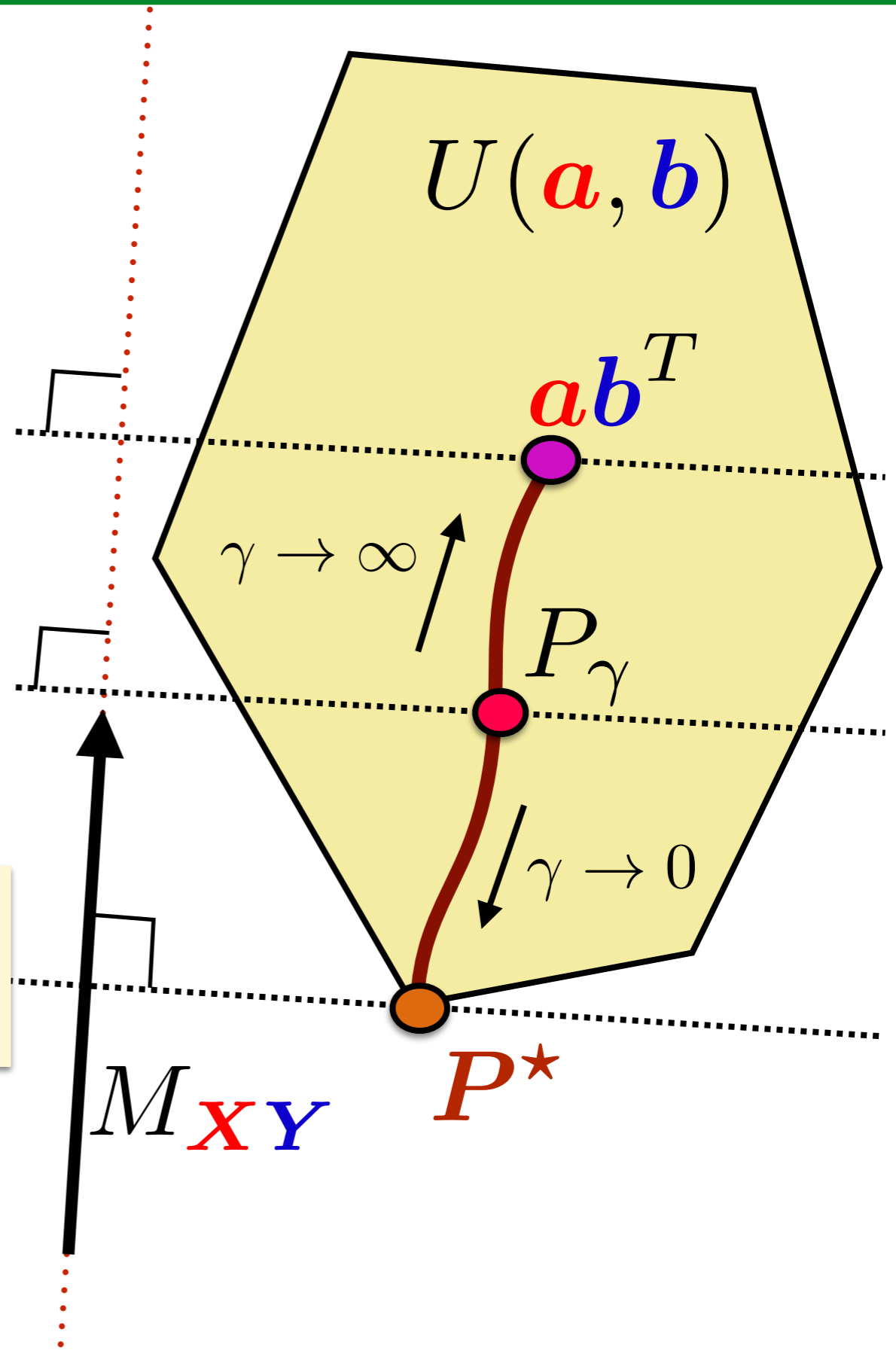
Sinkhorn in between W and MMD

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

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

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

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



Sinkhorn in between W and MMD

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

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

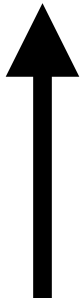
$$W_\gamma(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle P_\gamma, M_{\mathbf{XY}} \rangle$$

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

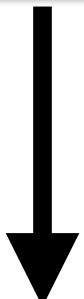
$$W^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \langle \mathbf{P}^*, M_{\mathbf{XY}} \rangle$$

Sinkhorn in between W and MMD

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

$\gamma \rightarrow \infty$ 

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

$\gamma \rightarrow 0$ 

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

How to compare them?

i.i.d samples $\mathbf{x}_1, \dots, \mathbf{x}_n \sim \mu$, $\mathbf{y}_1, \dots, \mathbf{y}_m \sim \nu$,

$$\hat{\mu}_n \stackrel{\text{def}}{=} \frac{1}{n} \sum_i \delta_{\mathbf{x}_i}, \hat{\nu}_m \stackrel{\text{def}}{=} \frac{1}{m} \sum_j \delta_{\mathbf{y}_j}$$

Computational properties

Effort to compute/approximate $\Delta(\hat{\mu}_n, \hat{\nu}_m)$?

Statistical properties

$$|\Delta(\mu, \nu) - \Delta(\hat{\mu}_n, \hat{\nu}_n)| \leq f(n)?$$

Sinkhorn in between W and MMD

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

$$(n + m)^2$$

$$O(1/\sqrt{n})$$

[see Arthur]

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

$$O((n + m)nm \log(n + m))$$

$$O(1/n^{1/d})$$

Sinkhorn in between W and MMD

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

$$(n + m)^2$$

$$O(1/\sqrt{n})$$

[see Arthur]

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

$$O((n + m)^2)$$

$$O\left(\frac{1}{\gamma^{d/2} \sqrt{n}}\right)$$

[GCBCP'18]

[FSVATP'18]

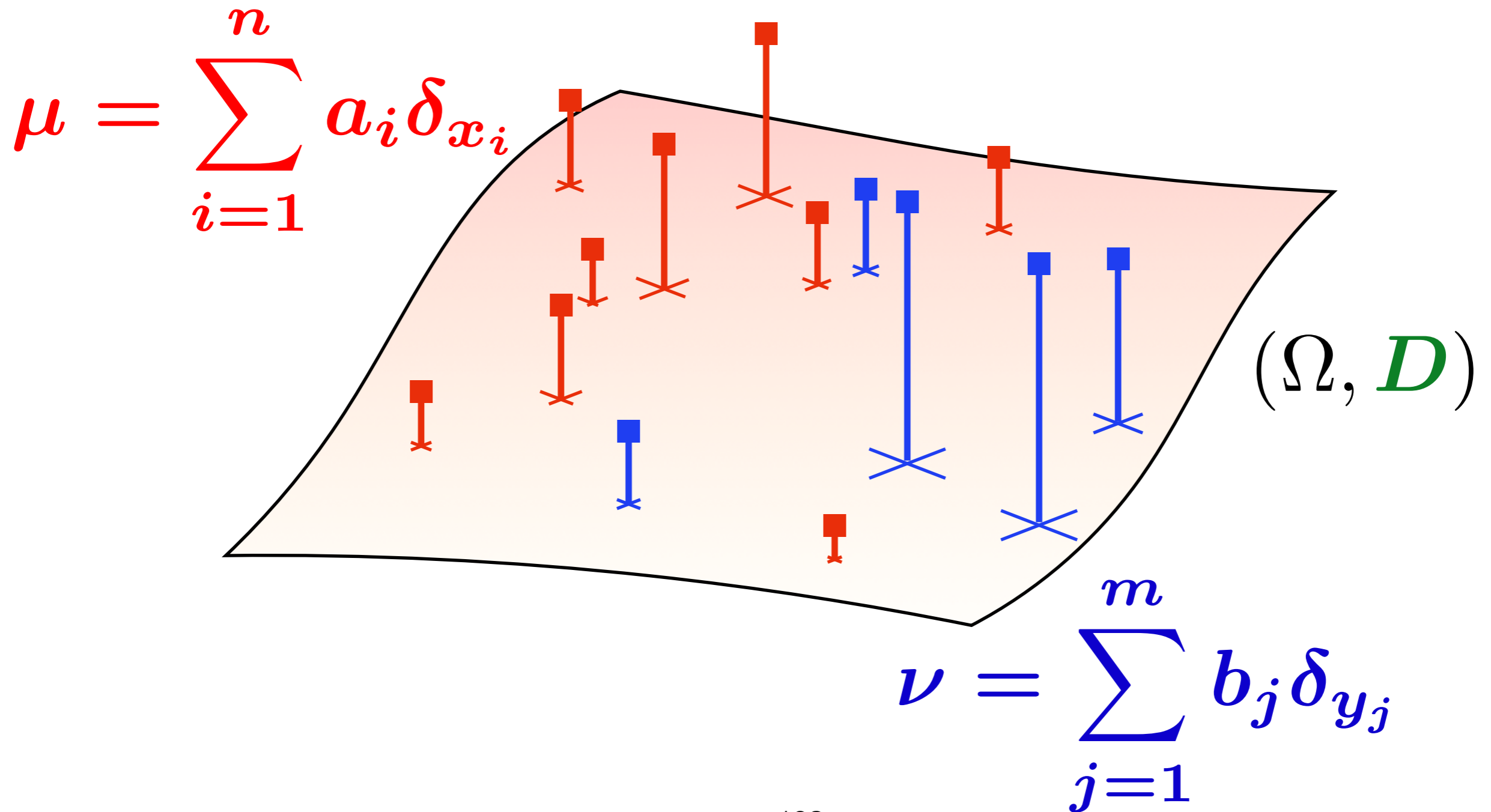
$$W^p(\mu, \nu) = \langle P^*, M_{XY} \rangle$$

$$O((n + m)nm \log(n + m))$$

$$O(1/n^{1/d})$$

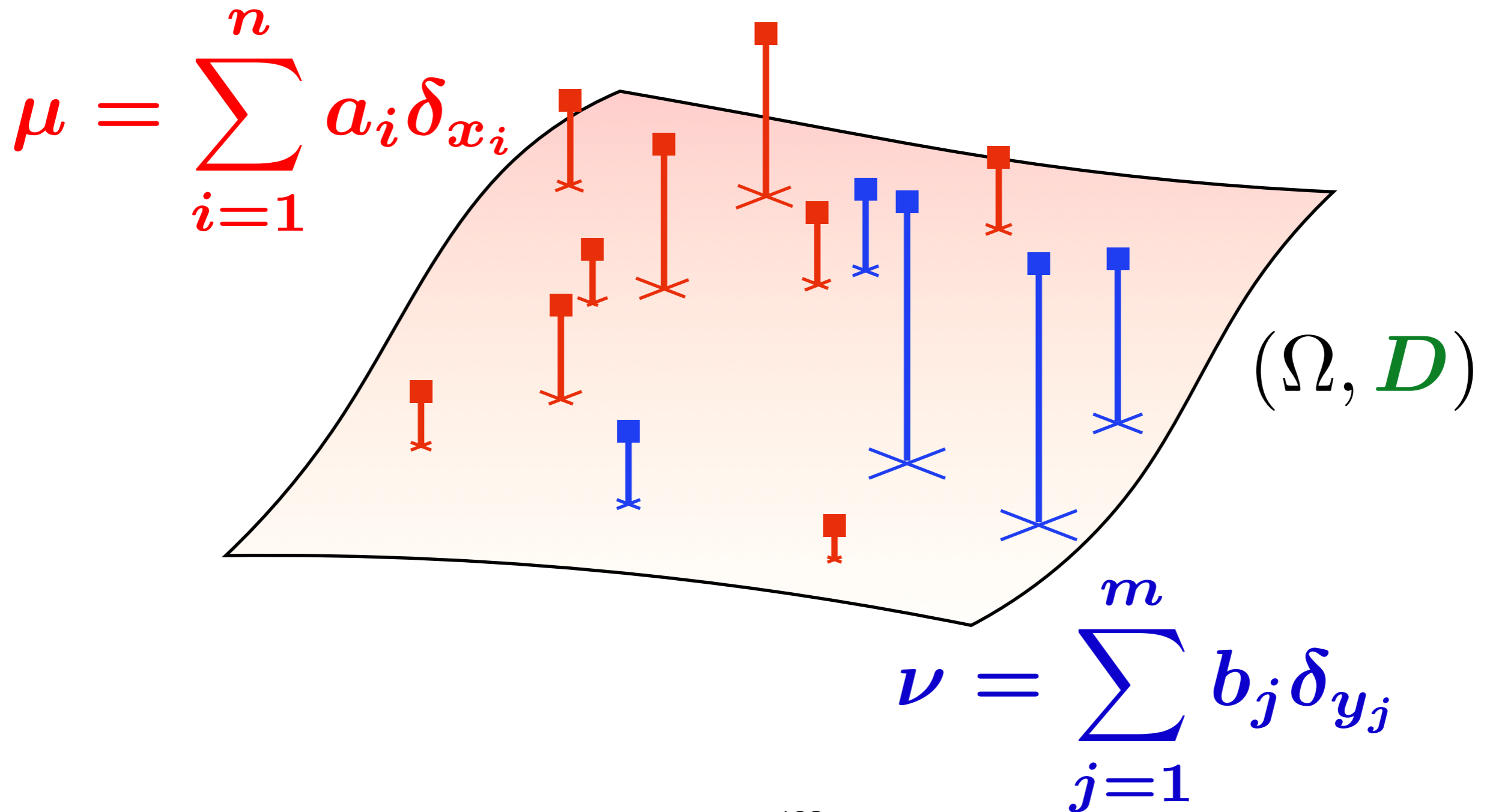
Differentiability of W

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



Differentiability of W

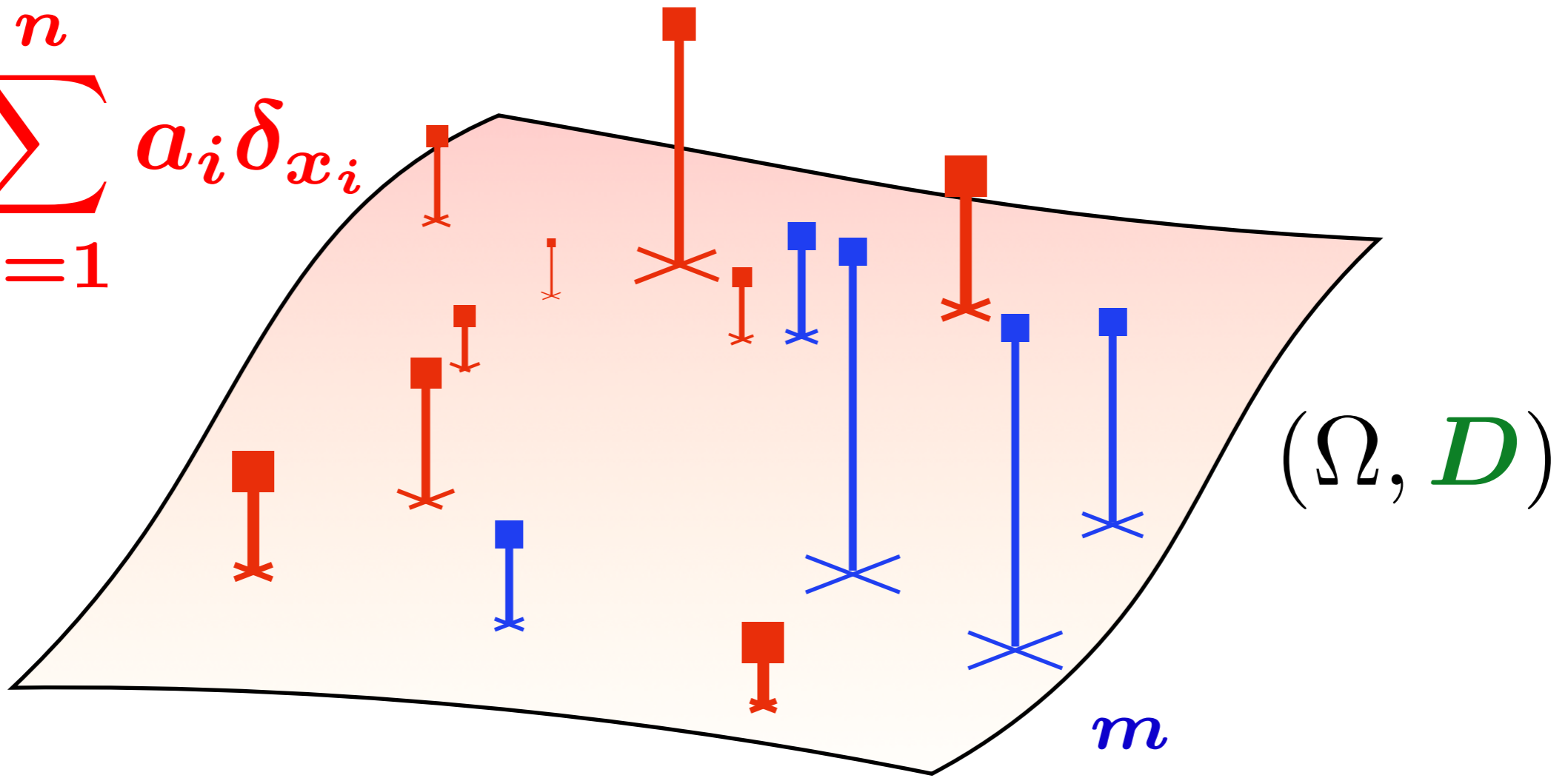
$$W((a + \Delta a, X), (b, Y)) = W((a, X), (b, Y)) + ??$$



Differentiability of W

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

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

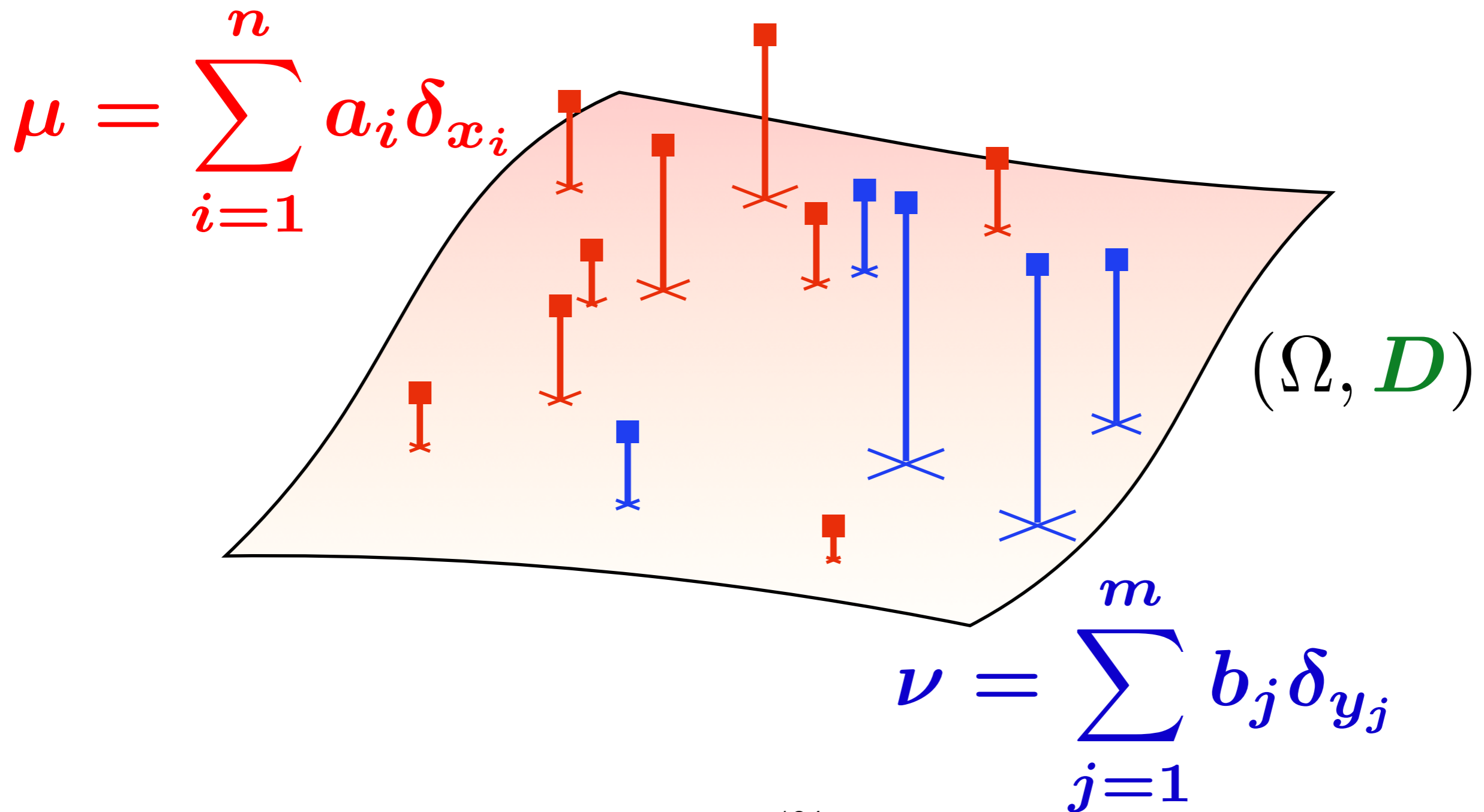


$$a \leftarrow a + \Delta a$$

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

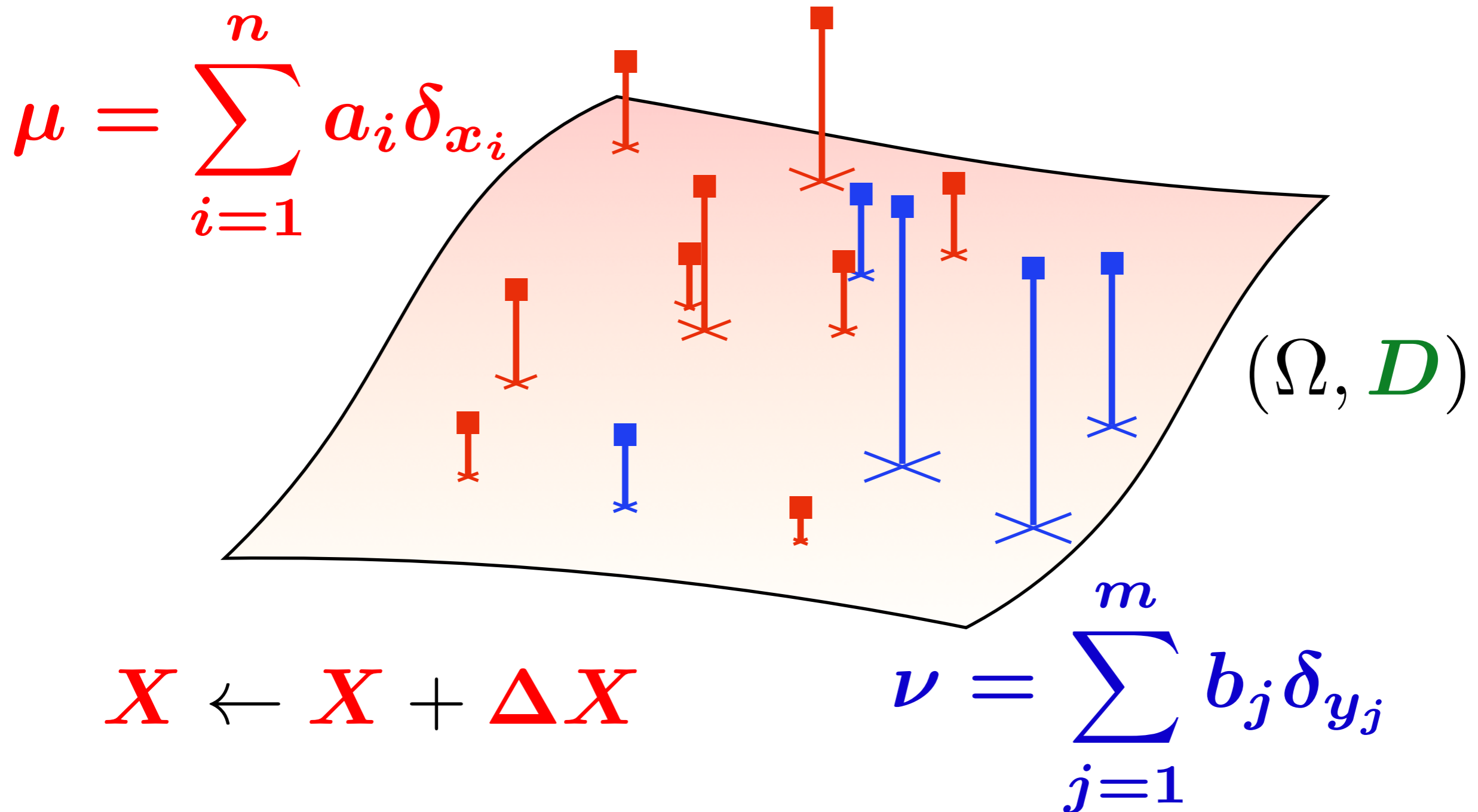
Sinkhorn \rightsquigarrow *Differentiability*

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



Sinkhorn \rightsquigarrow *Differentiability*

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



How to decrease W ? change weights

$$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu}) = \max_{\substack{\boldsymbol{\alpha} \in \mathbb{R}^n, \boldsymbol{\beta} \in \mathbb{R}^m \\ \boldsymbol{\alpha} \oplus \boldsymbol{\beta} \leq M_{\mathbf{X}\mathbf{Y}}}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}.$$

DUAL

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

$$\partial_{\mathbf{a}} W = \arg_{\boldsymbol{\alpha}} \max_{\boldsymbol{\alpha} \oplus \boldsymbol{\beta} \leq M_{\mathbf{X}\mathbf{Y}}} \boldsymbol{\alpha}^T \mathbf{a} + \boldsymbol{\beta}^T \mathbf{b}.$$

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

How to decrease W ? change locations

$$W_2^2(\boldsymbol{\mu}, \boldsymbol{\nu}) = \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P\mathbf{1}_m = \mathbf{a}, P^T\mathbf{1}_n = \mathbf{b}}} \langle P, \mathbf{1}_n \mathbf{1}_d^T X^2 + Y^{2T} \mathbf{1}_d \mathbf{1}_m - 2X^T Y \rangle$$

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

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

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

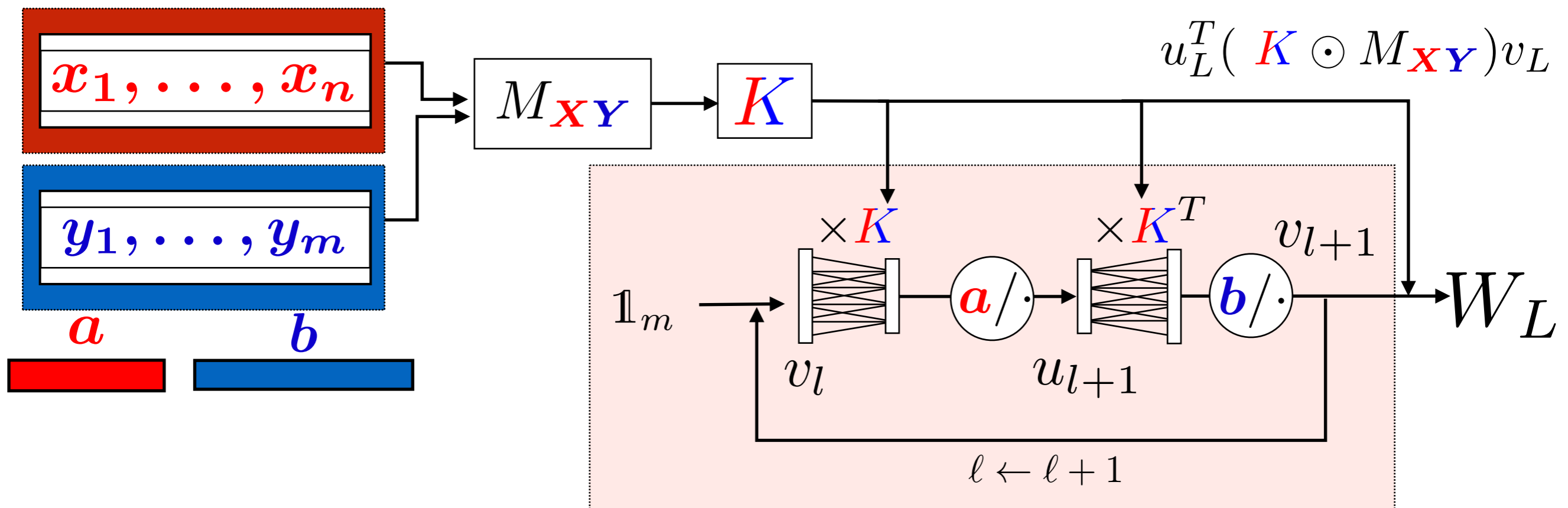
$$\boldsymbol{v}_0 = \mathbf{1}_m; l \geq 0, \boldsymbol{u}_l \stackrel{\text{def}}{=} \boldsymbol{a} / \boldsymbol{K} \boldsymbol{v}_l, \boldsymbol{v}_{l+1} \stackrel{\text{def}}{=} \boldsymbol{b} / \boldsymbol{K}^T \boldsymbol{u}_l.$$

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.

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

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

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.

[Hashimoto'16] [Bonnel'16][Shalit'16]

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



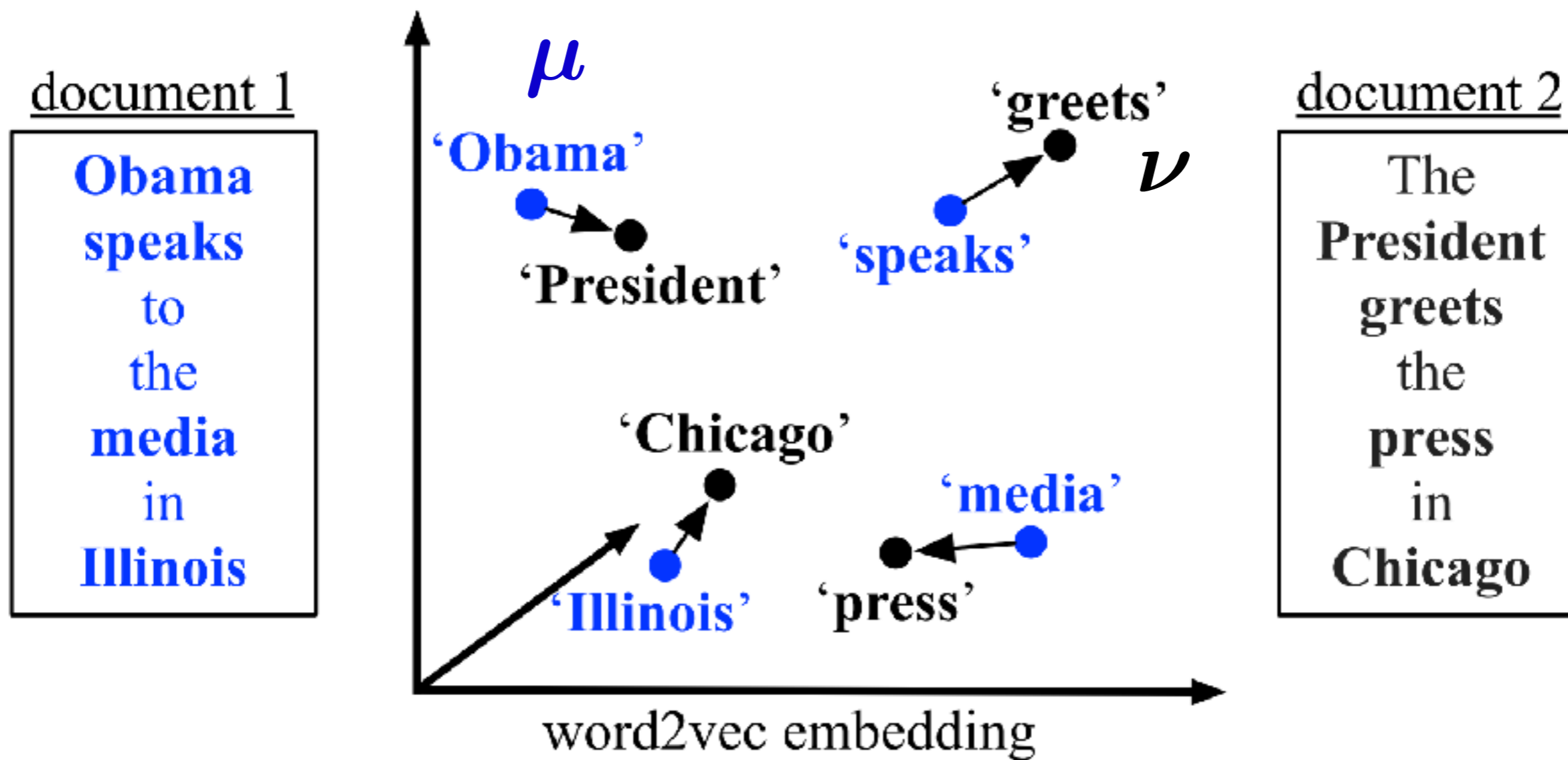
The Earth Mover's Distance



[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

“Wasserstein + Data” Problems

- Quantization, k -means problem [Lloyd’82]

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

- [McCann’95] Interpolant

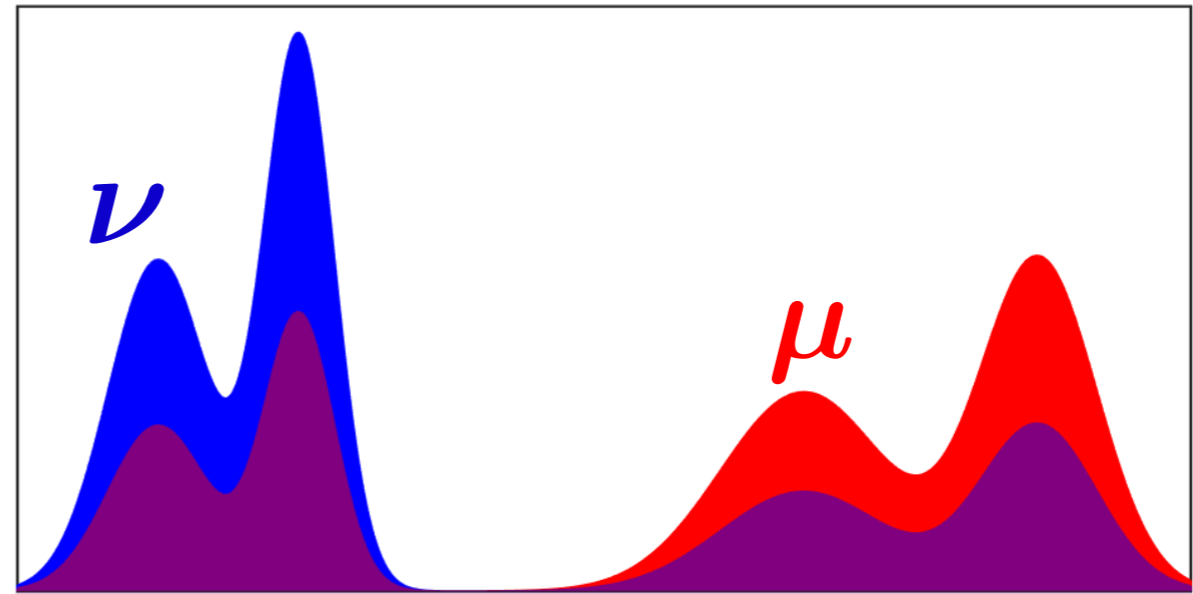
$$\min_{\mu \in \mathcal{P}(\Omega)} (1 - t)W_2^2(\mu, \nu_1) + tW_2^2(\mu, \nu_2)$$

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

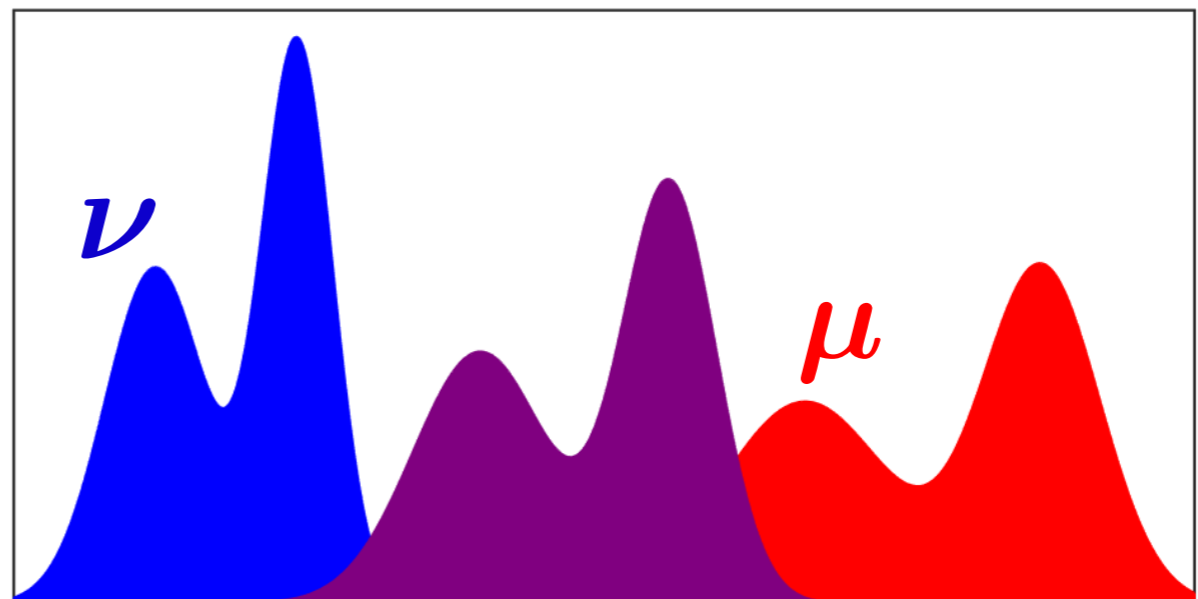
$$\mu_{t+1} = \operatorname{argmin}_{\mu \in \mathcal{P}(\Omega)} J(\mu) + \lambda_t W_p^p(\mu, \mu_t)$$

Averaging Measures

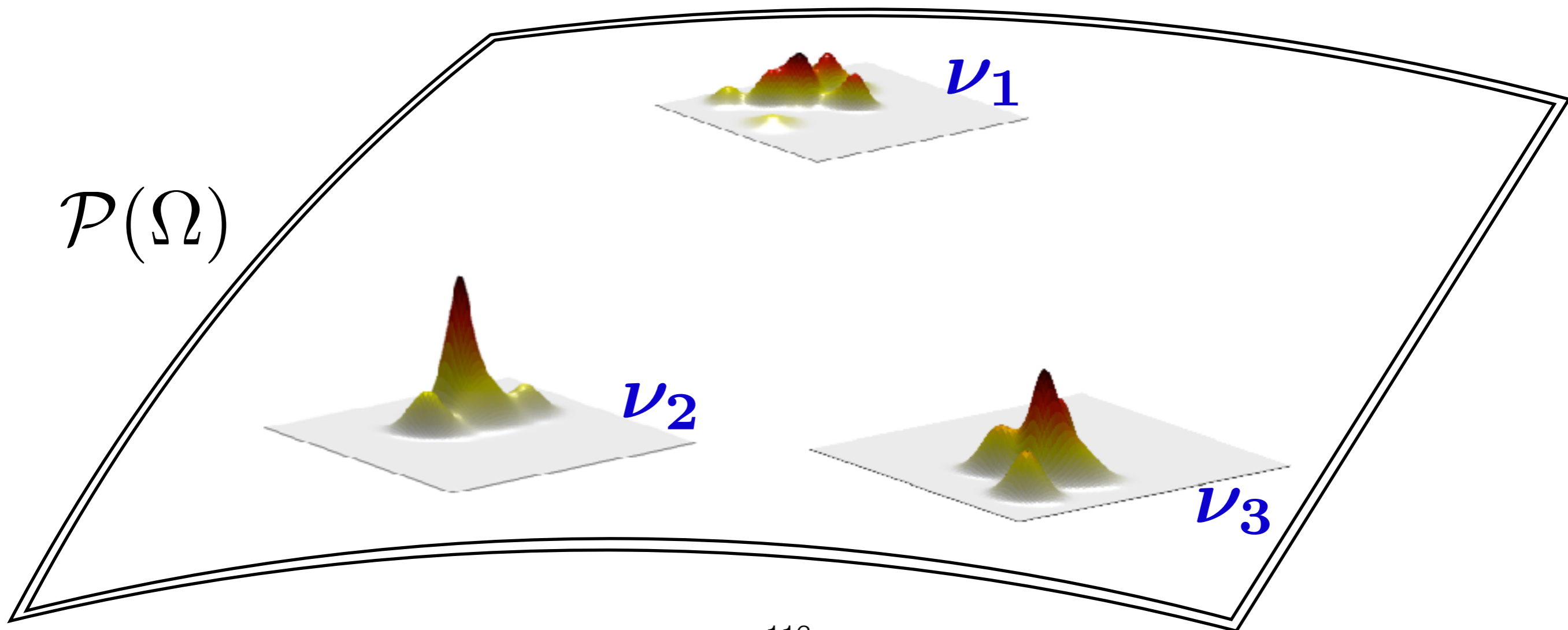
L_2 average



W average



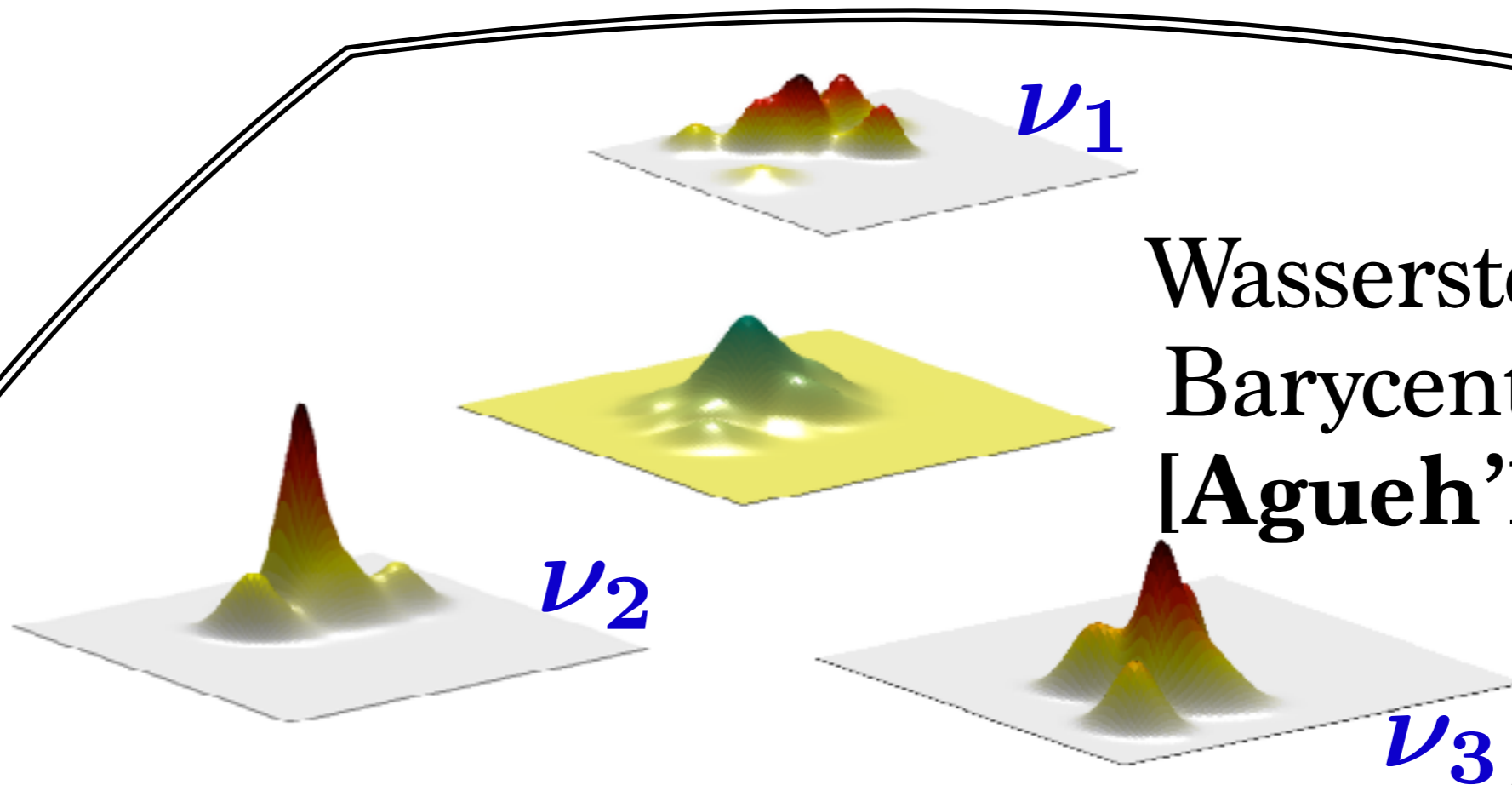
Barycenter for Measures?



Barycenter for Measures?

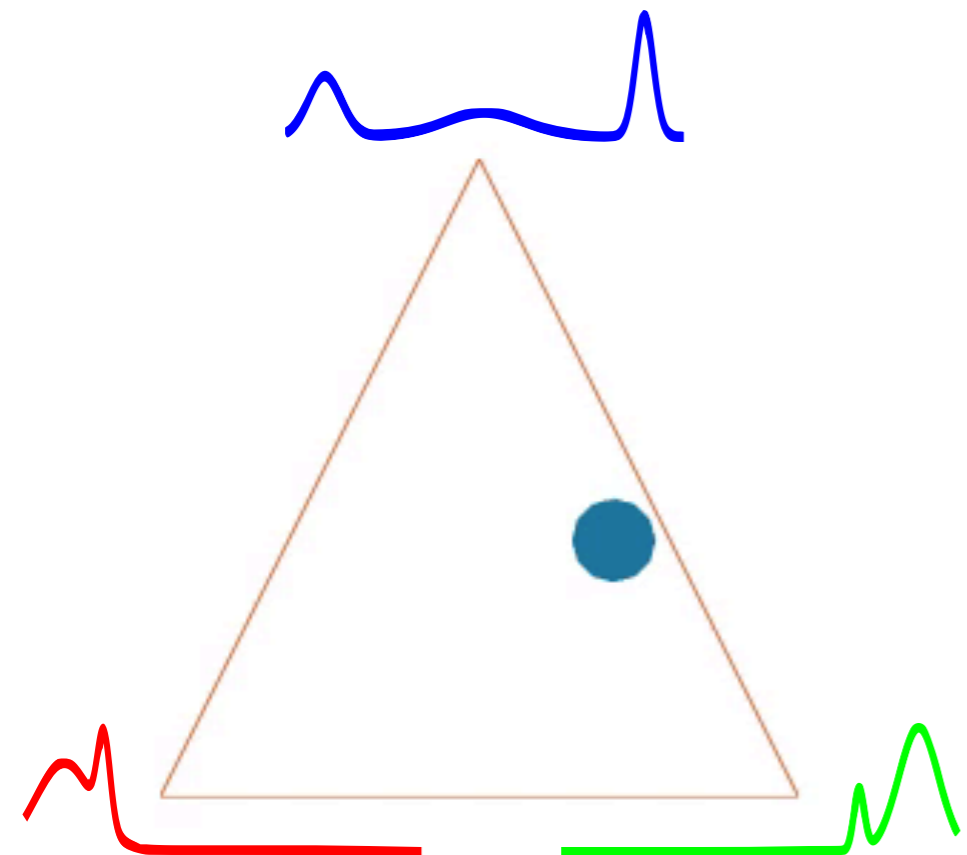
$$\min_{\mu \in \mathcal{P}(\Omega)} \sum_{i=1}^N \lambda_i W_p^p(\mu, \nu_i)$$

$\mathcal{P}(\Omega)$

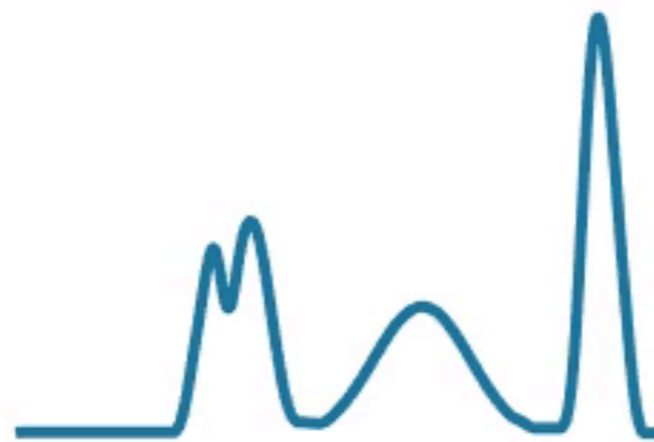


Wasserstein
Barycenter
[Agueh'11]

Barycenter for Measures?



$$\lambda \in \Sigma_3$$

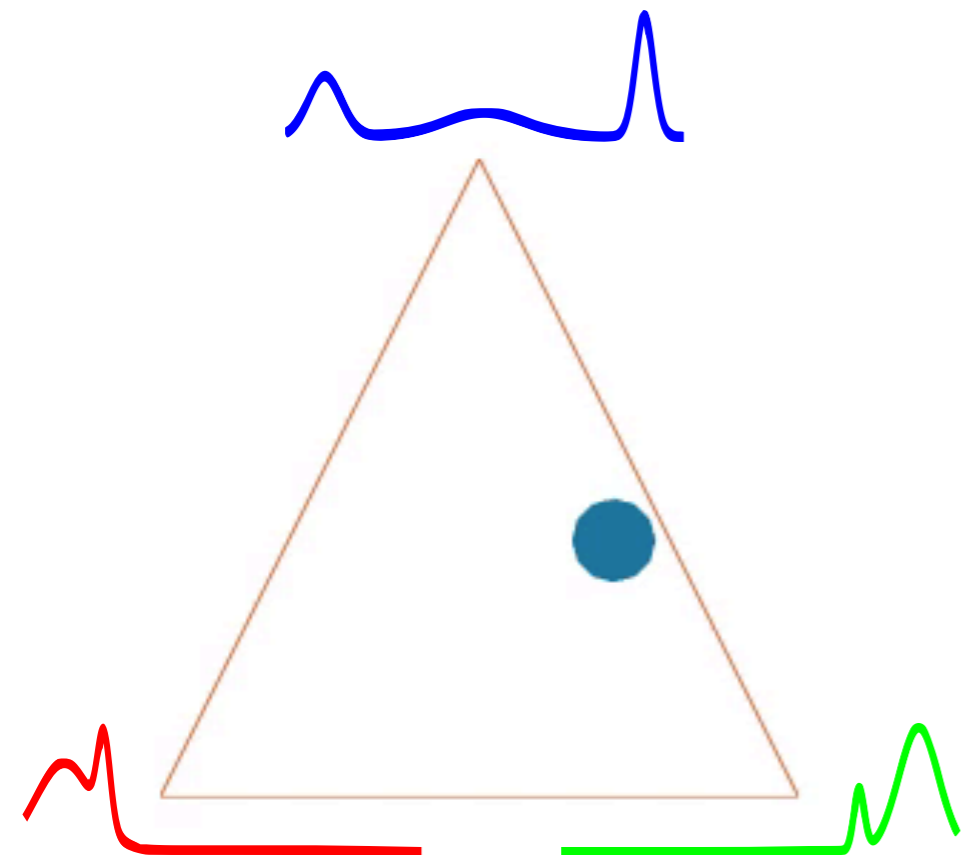


Wasserstein mean

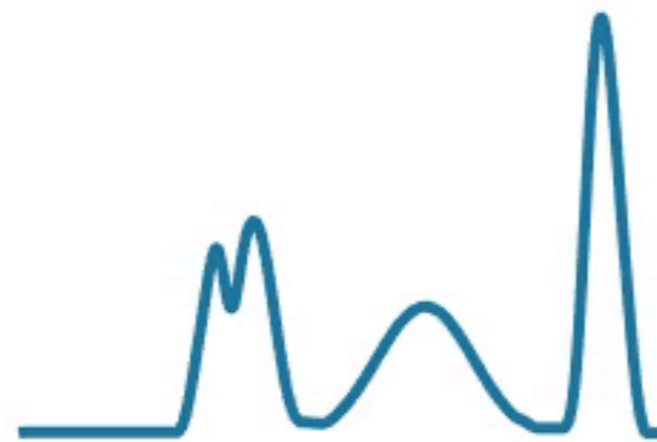


L_2 mean

Barycenter for Measures?



$$\lambda \in \Sigma_3$$



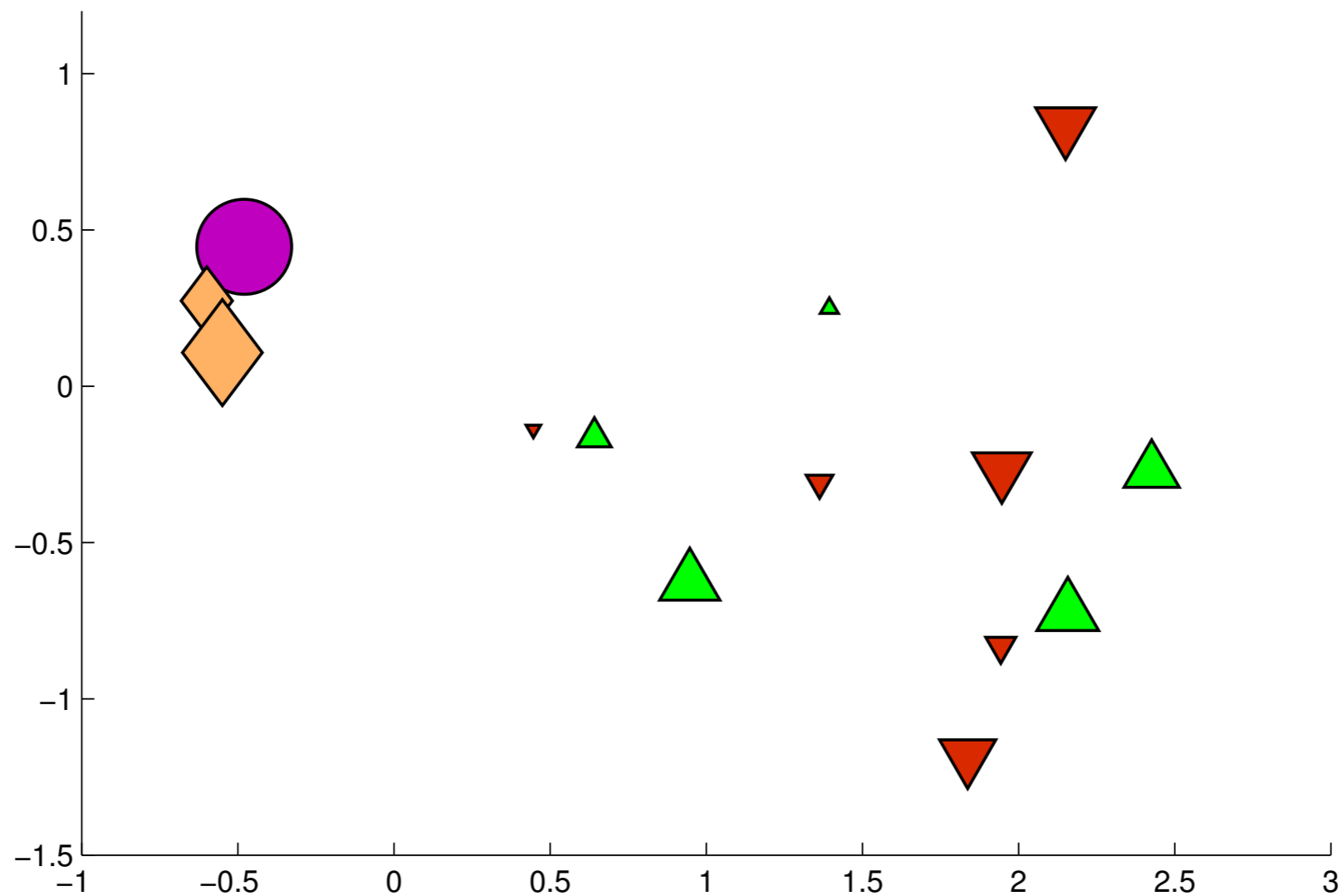
Wasserstein mean



L_2 mean

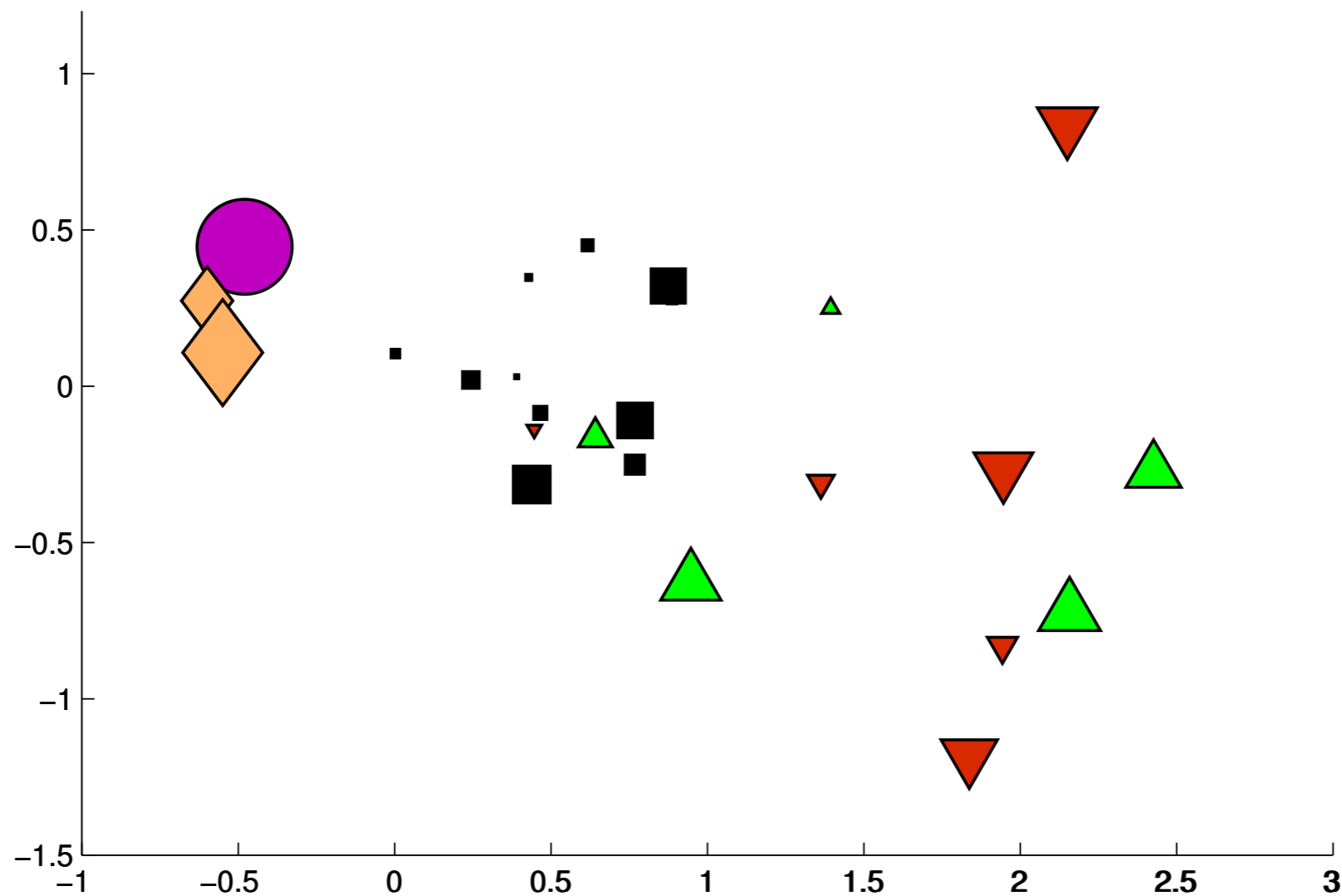
Multimarginal Formulation

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

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

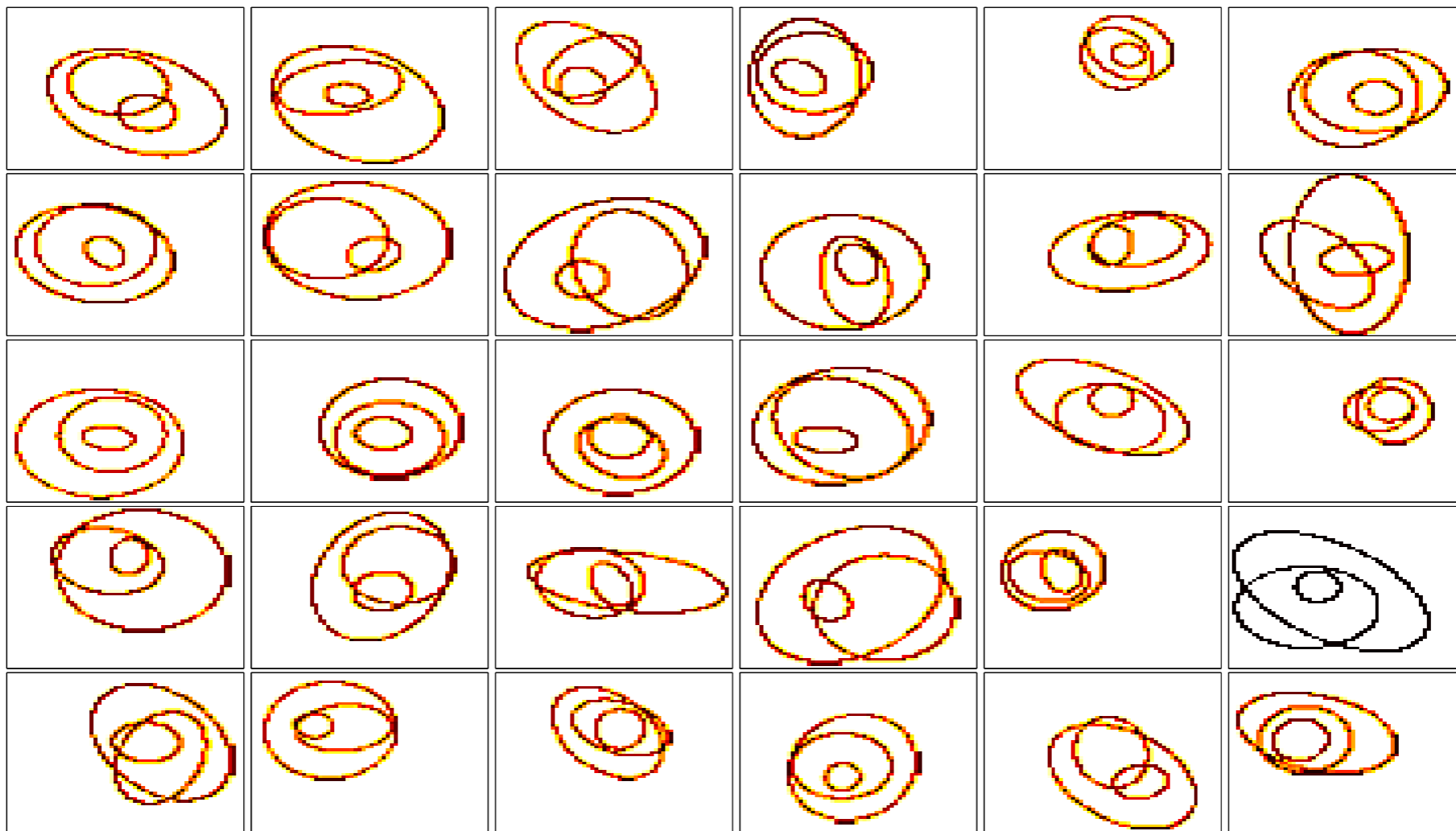
When Ω is a **finite metric space** defined by M .

$$\begin{aligned} \min_{P_1, \dots, P_N, \mathbf{a}} \quad & \sum_{i=1}^N \lambda_i \langle P_i, M \rangle \\ \text{s.t.} \quad & P_i^T \mathbf{1}_n = \mathbf{b}_i, \forall i \leq N, \\ & P_1 \mathbf{1}_n = \dots = P_N \mathbf{1}_d = \mathbf{a}. \end{aligned}$$

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

Primal Descent on Regularized W

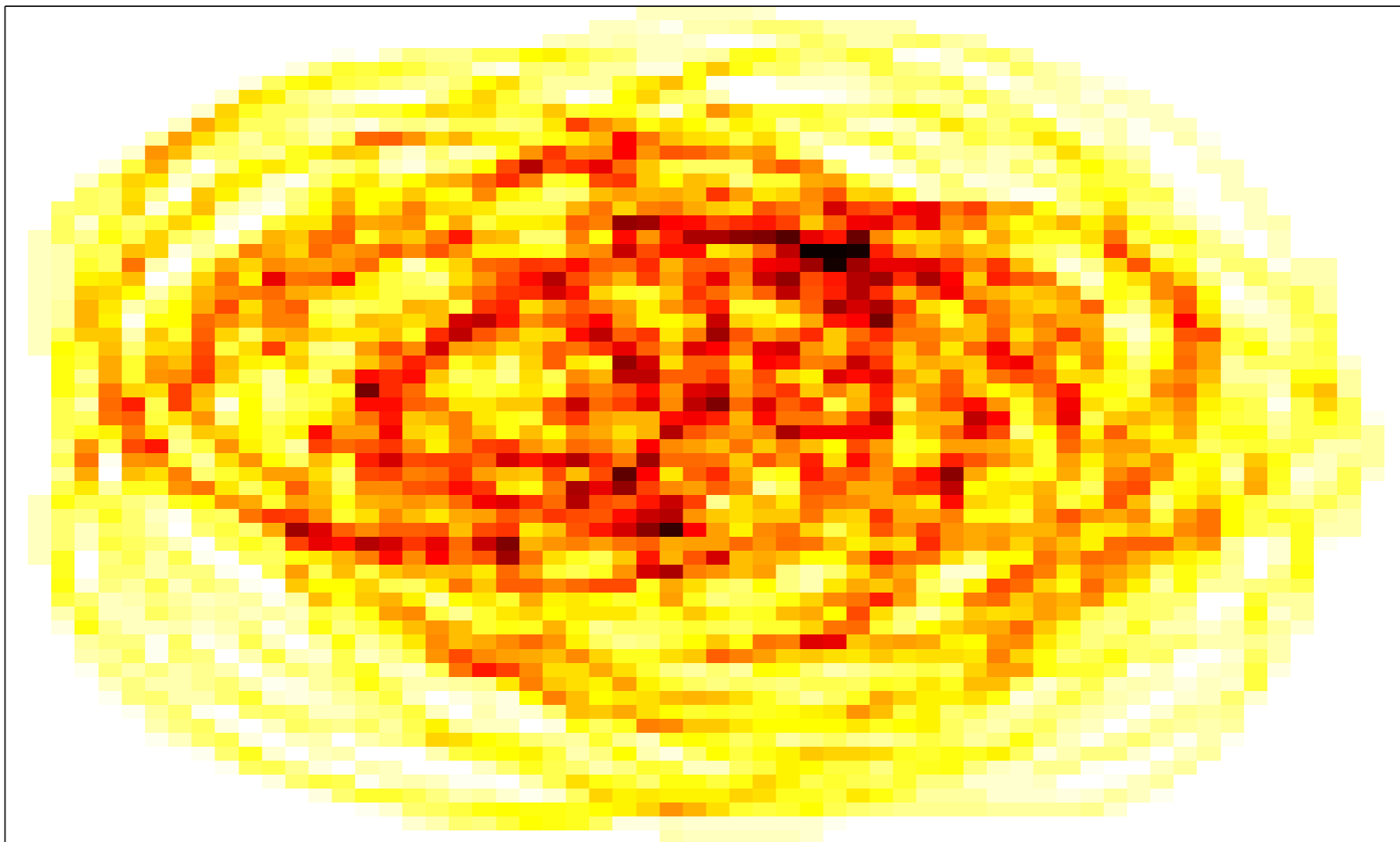
$$\min_{\mathbf{a} \in \Sigma_{h \times h}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$



[Cuturi'14]

Primal Descent on Regularized W

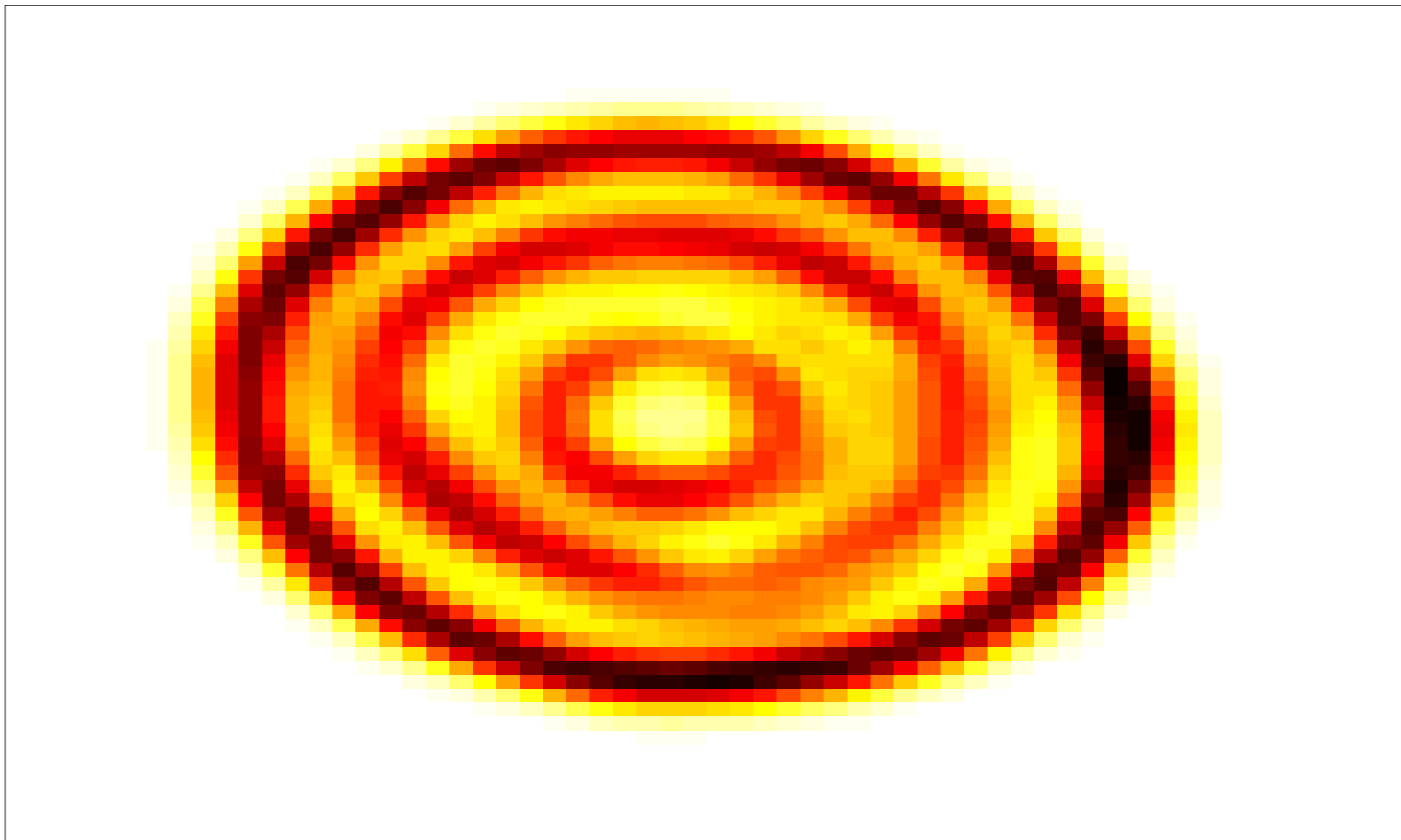
$$\min_{\mathbf{a} \in \Sigma_{h \times h}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i)$$



[Cuturi'14]

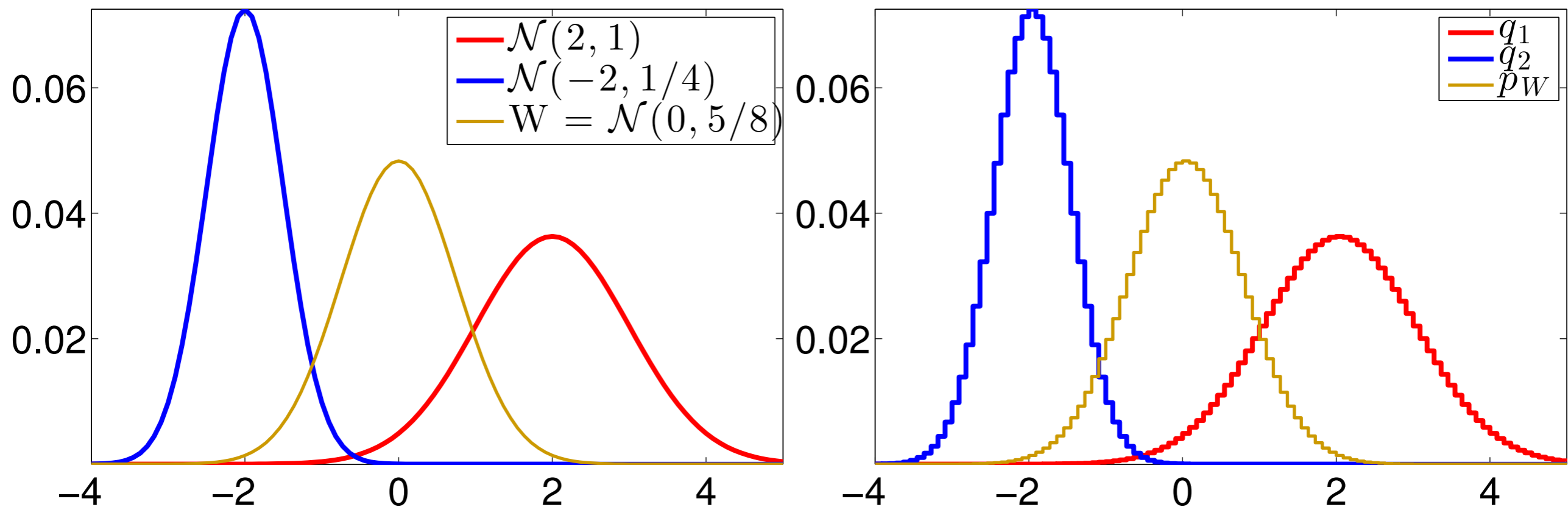
Primal Descent on Regularized W

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

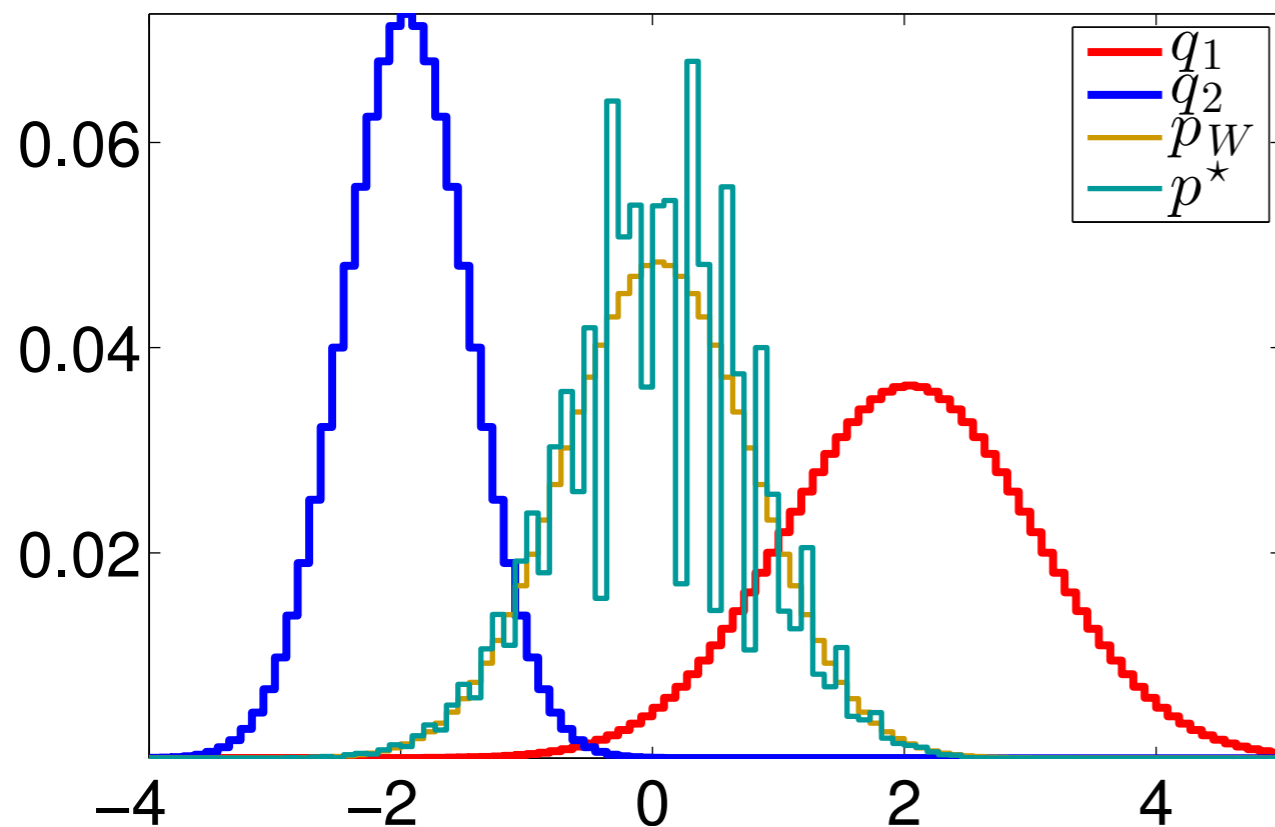


[Cuturi'14]

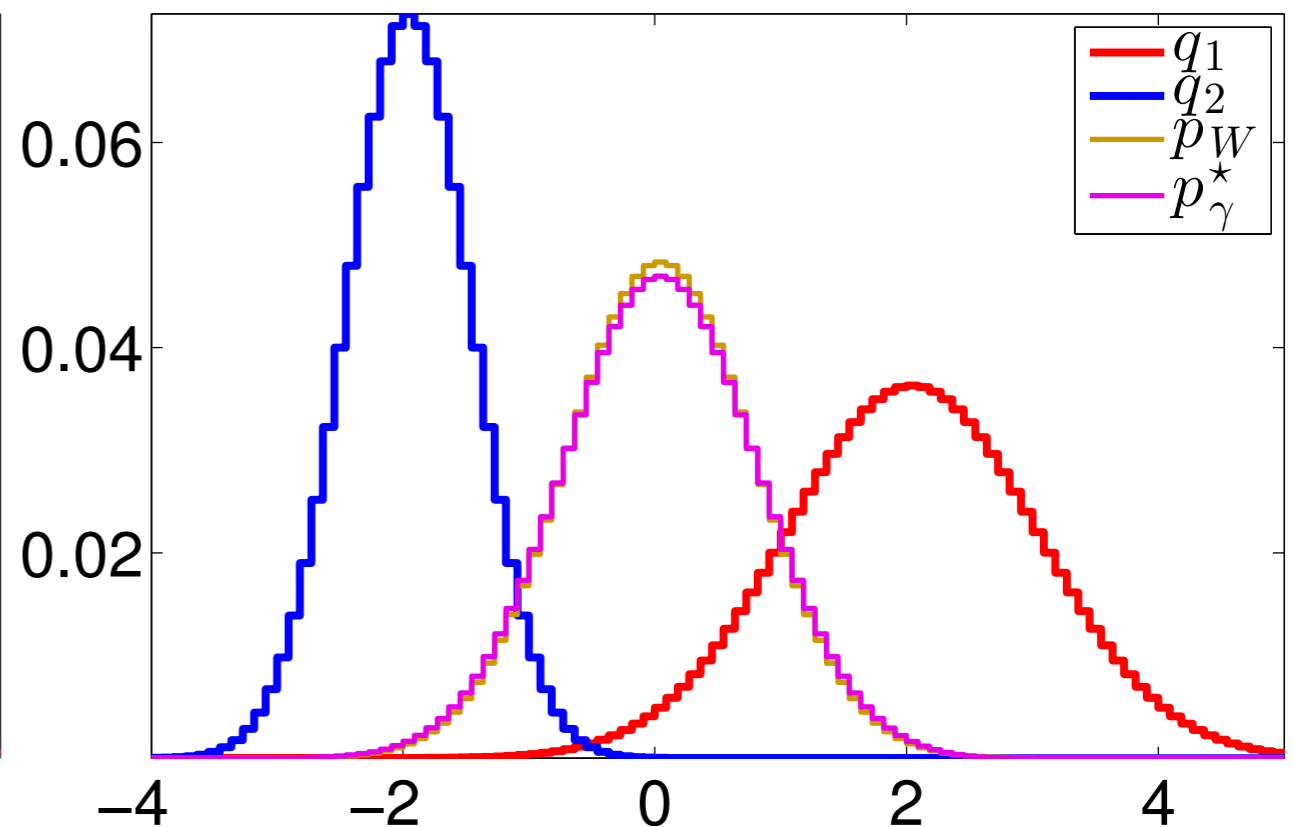
On Regularizing or Not



On Regularizing or Not

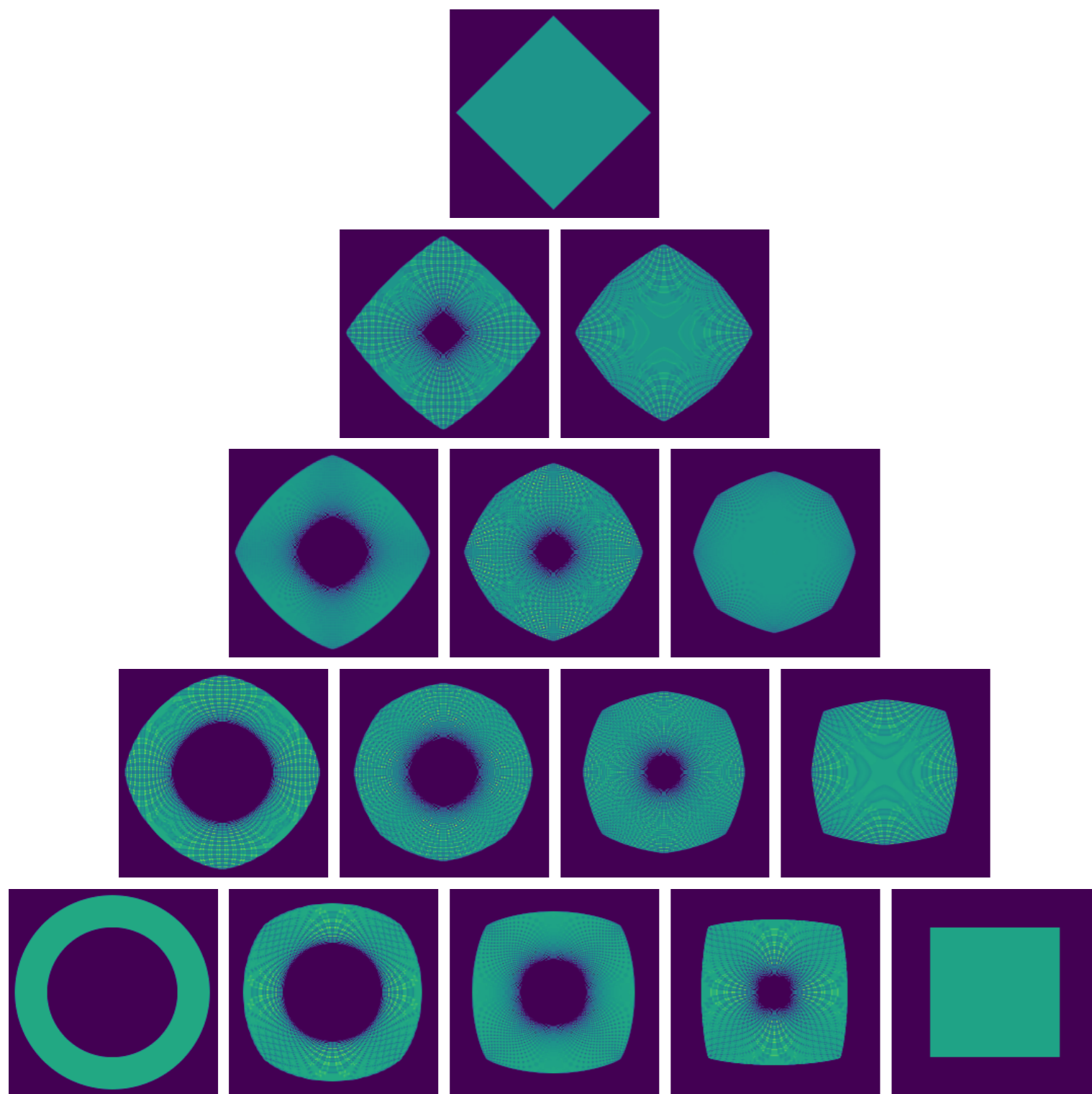


True barycenter



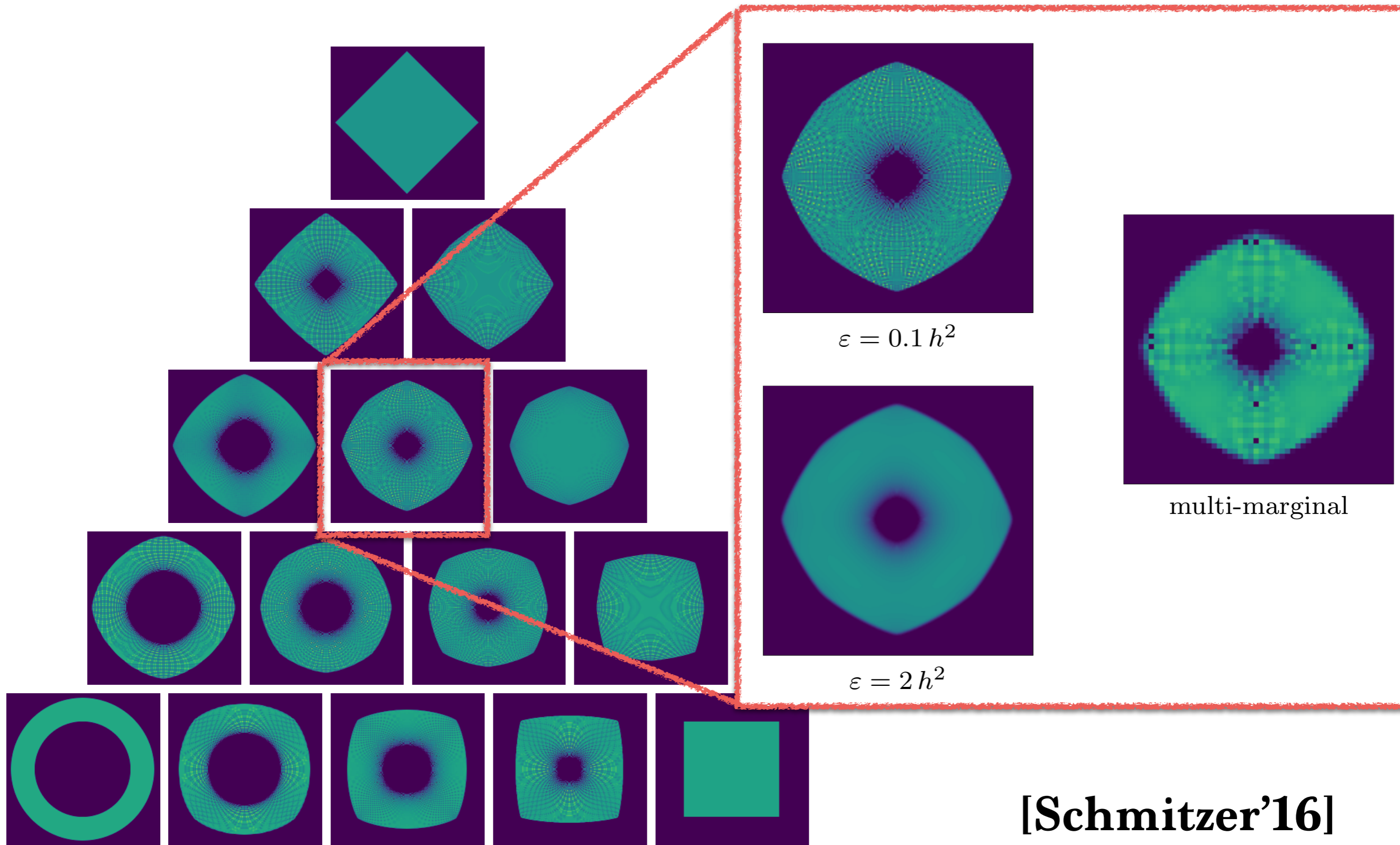
*Barycenter using
regularized OT*

On Regularizing or Not



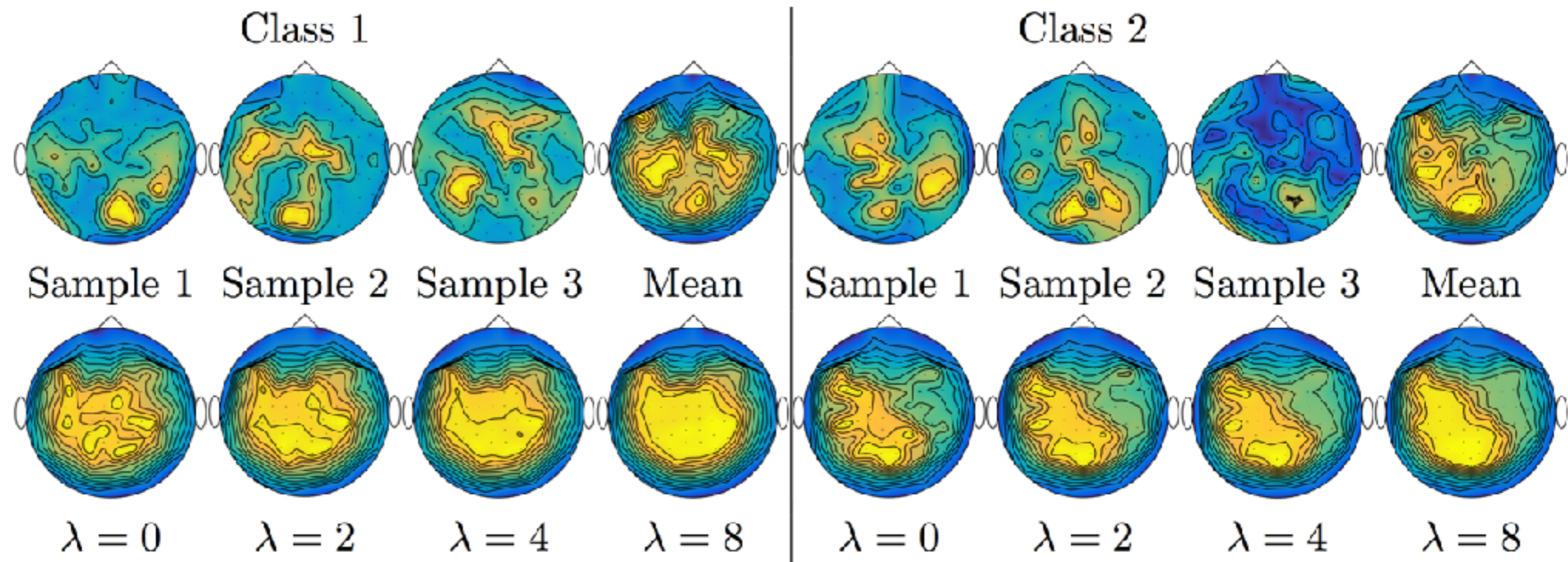
[Schmitzer'16]

On Regularizing or Not



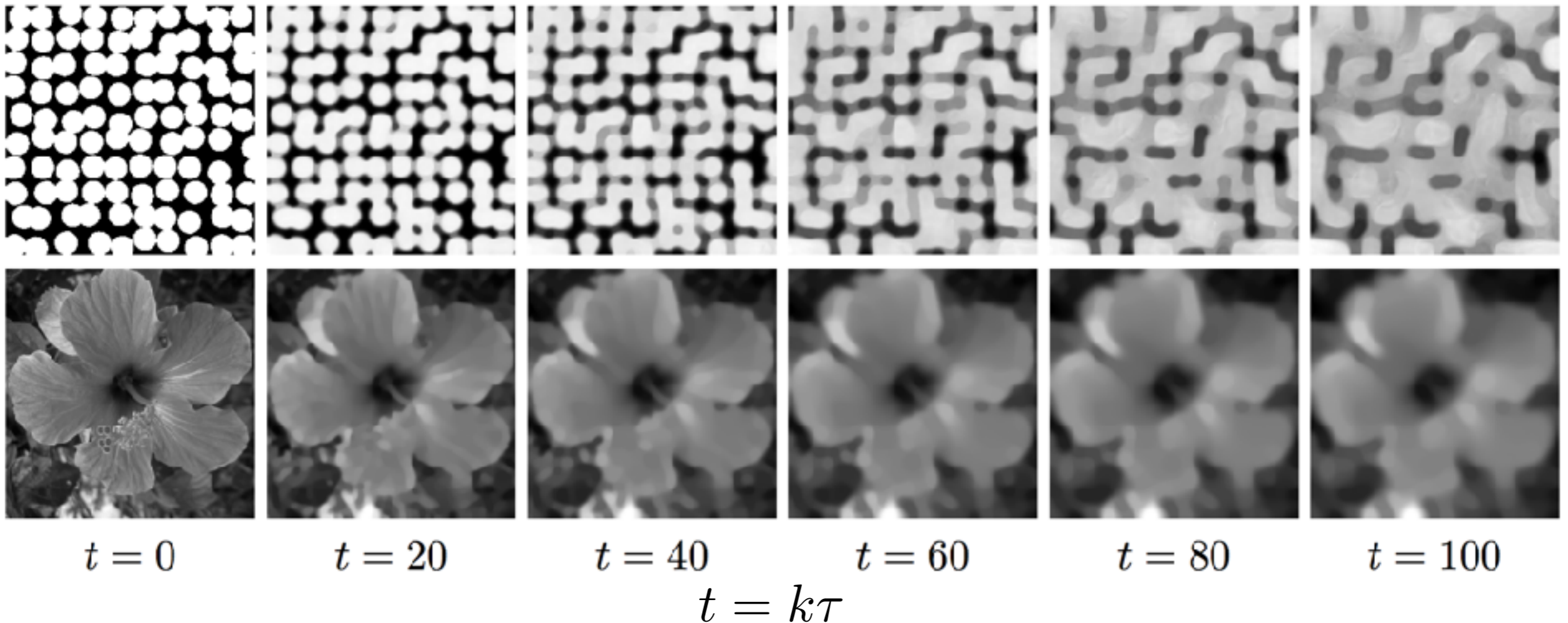
Duality: Regularized Barycenters

$$\min_{\mu \in \mathcal{P}(\Omega)} \sum_{i=1}^N W_{\gamma}(\mu, \nu_i) + \lambda \text{TV}(\mu)$$



Duality: TV Gradient Flow

$$\mu_{k+1} = \operatorname{argmin}_{\mu \in \mathcal{P}(\Omega)} W_{\gamma}(\mu, \mu_k) + \tau \operatorname{TV}(\mu)$$



Regularized OT as KL Projection

$$\mathbf{KL}(P \mid K) = \sum_{ij} P_{ij} \log (P_{ij} / K_{ij})$$
$$\langle P, M_{\mathbf{X}\mathbf{Y}} \rangle - \gamma E(P) = \gamma \mathbf{KL}(P \mid K)$$

Prop. $P_\gamma = \text{Proj}_{C_{\mathbf{a}} \cap C'_{\mathbf{b}}}(K)$

$$C_{\mathbf{a}} = \{P \mid P \mathbf{1}_m = \mathbf{a}\}, \quad C'_{\mathbf{b}} = \{P \mid P^T \mathbf{1}_n = \mathbf{b}\}$$

Regularized OT as KL Projection

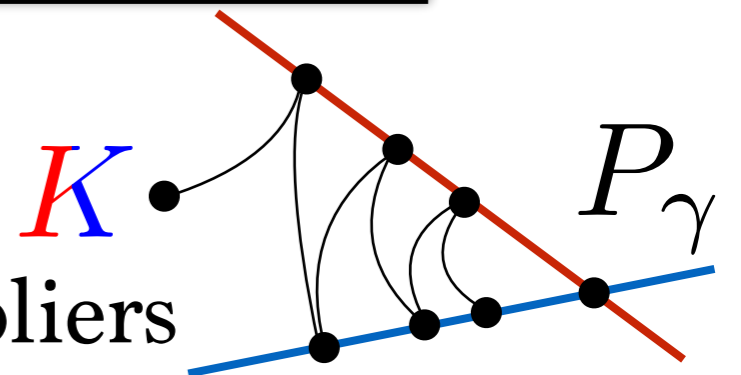
$$\text{Prop. } P_\gamma = \text{Proj}_{C_a \cap C'_b}(K)$$

$$C_a = \{P \mid P\mathbf{1}_m = \mathbf{a}\}, \quad C'_b = \{P \mid P^T\mathbf{1}_n = \mathbf{b}\}$$

$$\text{Proj}_{C_a}(P) = \mathbf{D} \left(\frac{\mathbf{a}}{P\mathbf{1}_m} \right) P,$$

$$\text{Proj}_{C'_b}(P) = P \mathbf{D} \left(\frac{\mathbf{b}}{P^T\mathbf{1}_n} \right).$$

1. Sinkhorn = Dykstra's alternate projection
2. Only need to store & update diagonal multipliers



Wasserstein Barycenter = KL Projections

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

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} = [P_1, \dots, P_N] \\ \mathbf{P} \in \mathbf{C}_1 \cap \mathbf{C}_2}} \sum_{i=1}^N \lambda_i \mathbf{KL}(P_i | K)$$

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

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

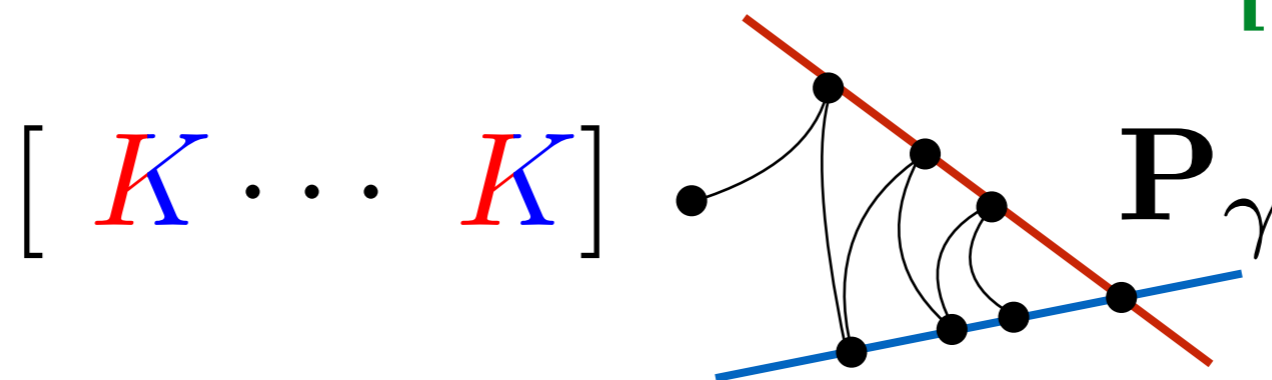
Wasserstein Barycenter = KL Projections

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P} = [P_1, \dots, P_N] \\ \mathbf{P} \in \mathcal{C}_1 \cap \mathcal{C}_2}} \sum_{i=1}^N \lambda_i \text{KL}(P_i | K)$$

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

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

[BCCNP'15]



Wasserstein Barycenter = KL Projections

$$\min_{\mathbf{a}} \sum_{i=1}^N \lambda_i W_{\gamma}(\mathbf{a}, \mathbf{b}_i) = \min_{\substack{\mathbf{P}=[P_1, \dots, P_N] \\ \mathbf{P} \in \mathcal{C}_1 \cap \mathcal{C}_2}} \sum_{i=1}^N \lambda_i \text{KL}(P_i | K)$$

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

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

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

```
while not converged
```

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

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

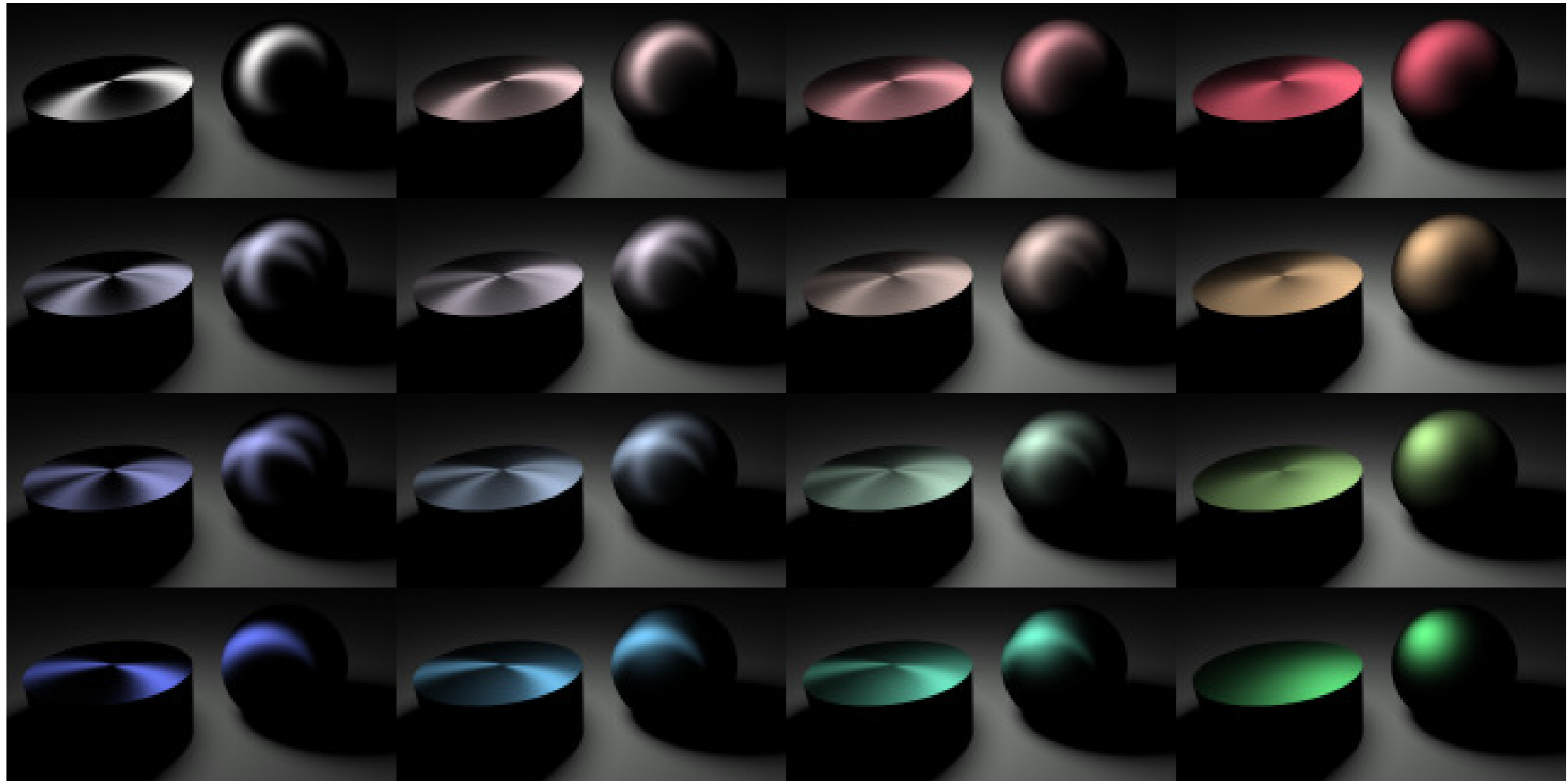
```
end
```

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

[BCCNP'15]

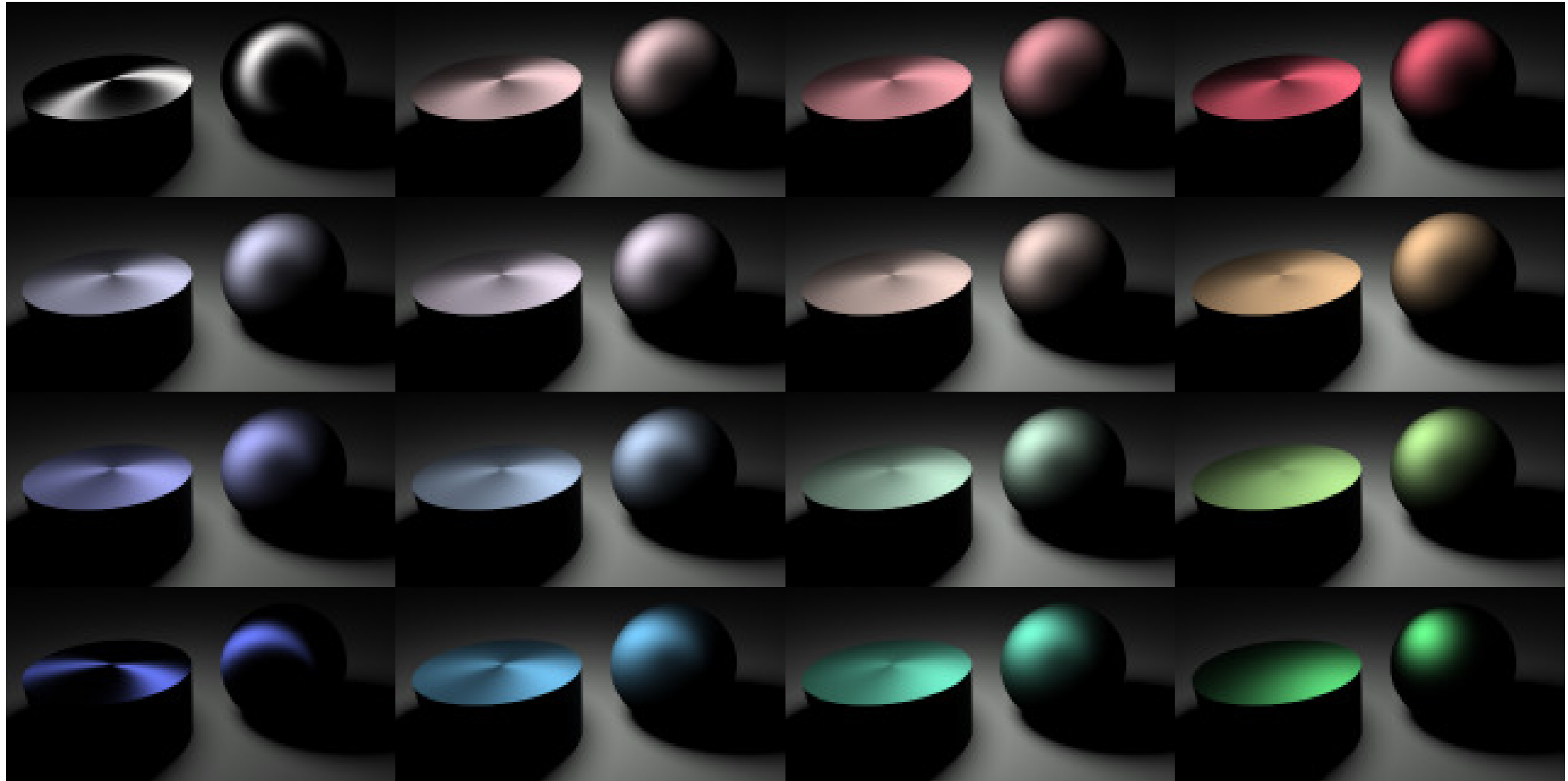
*Iterative Bregman Projections for
Regularized Transportation Problems*
SIAM J. on Sci. Comp. 2015

Applications in Imaging



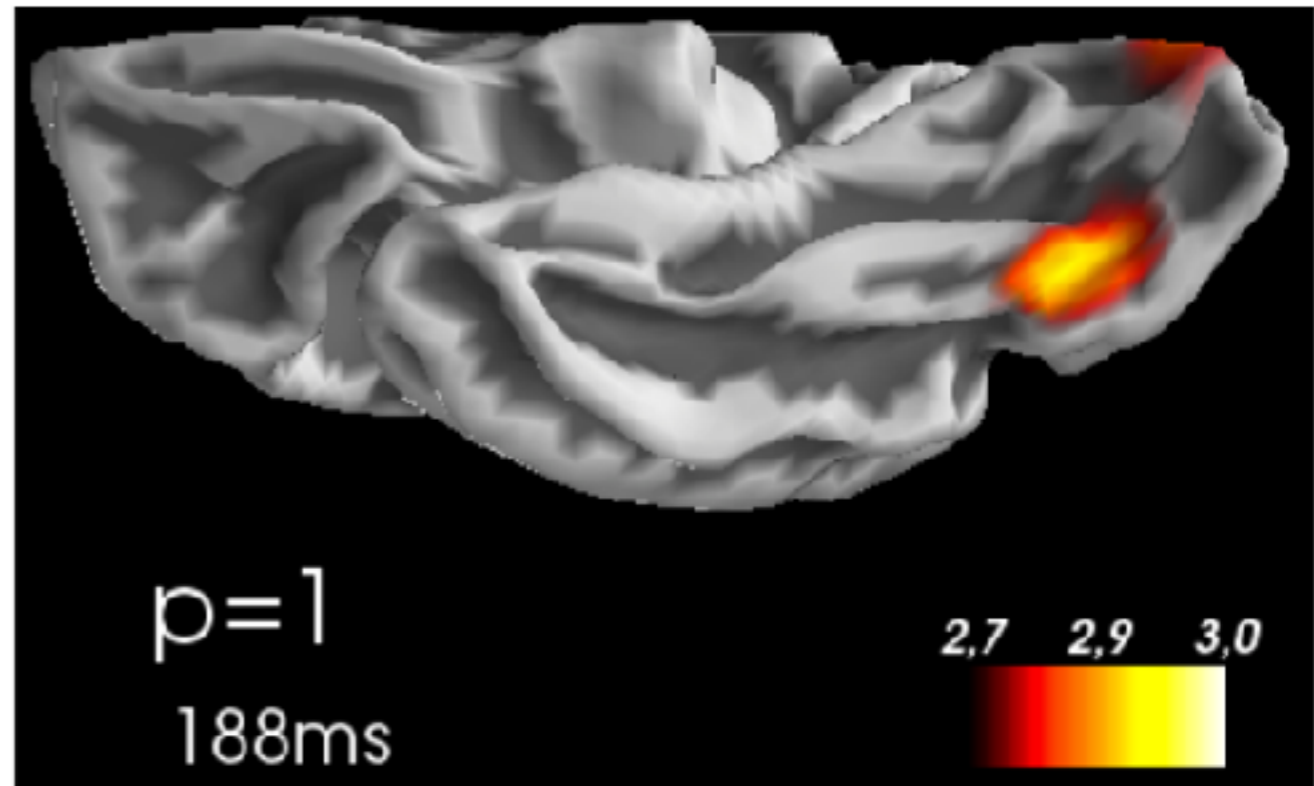
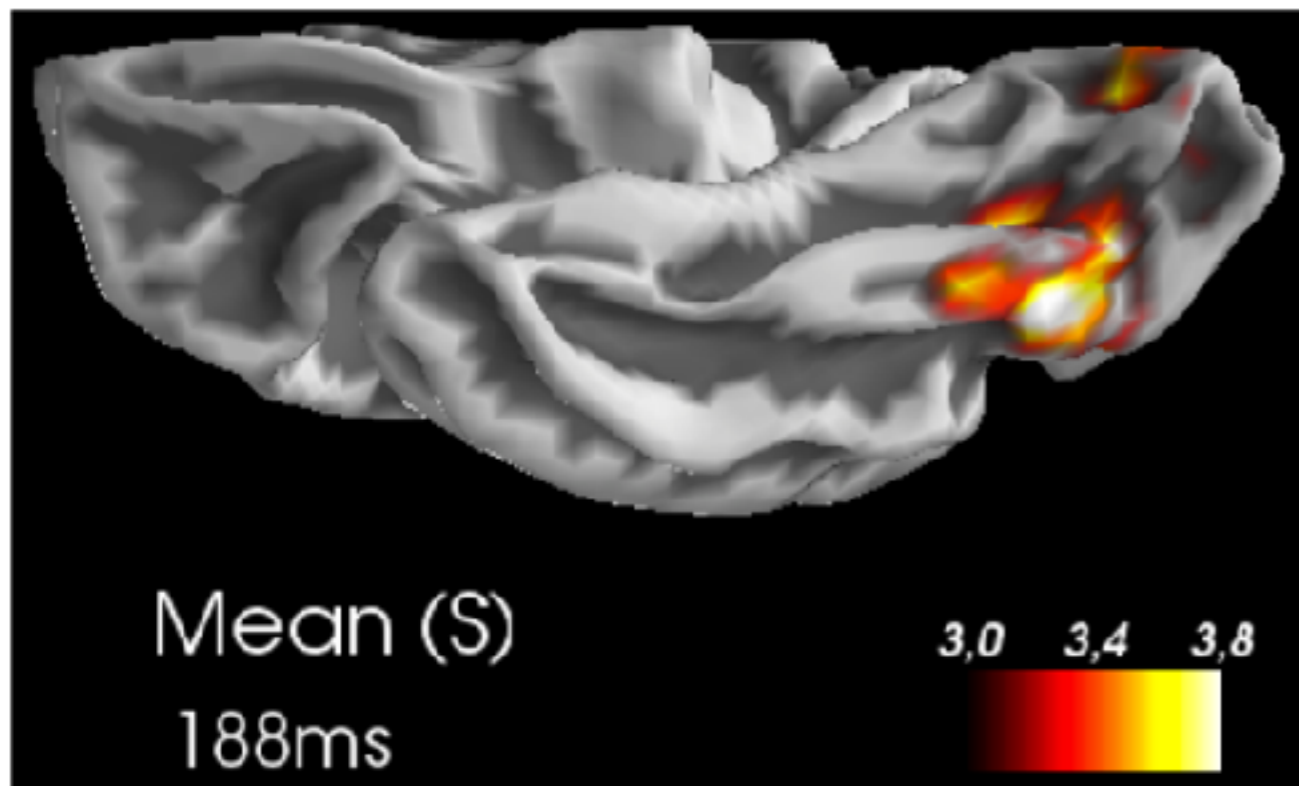
[Solomon'15]

Applications in Imaging



[Solomon'15]

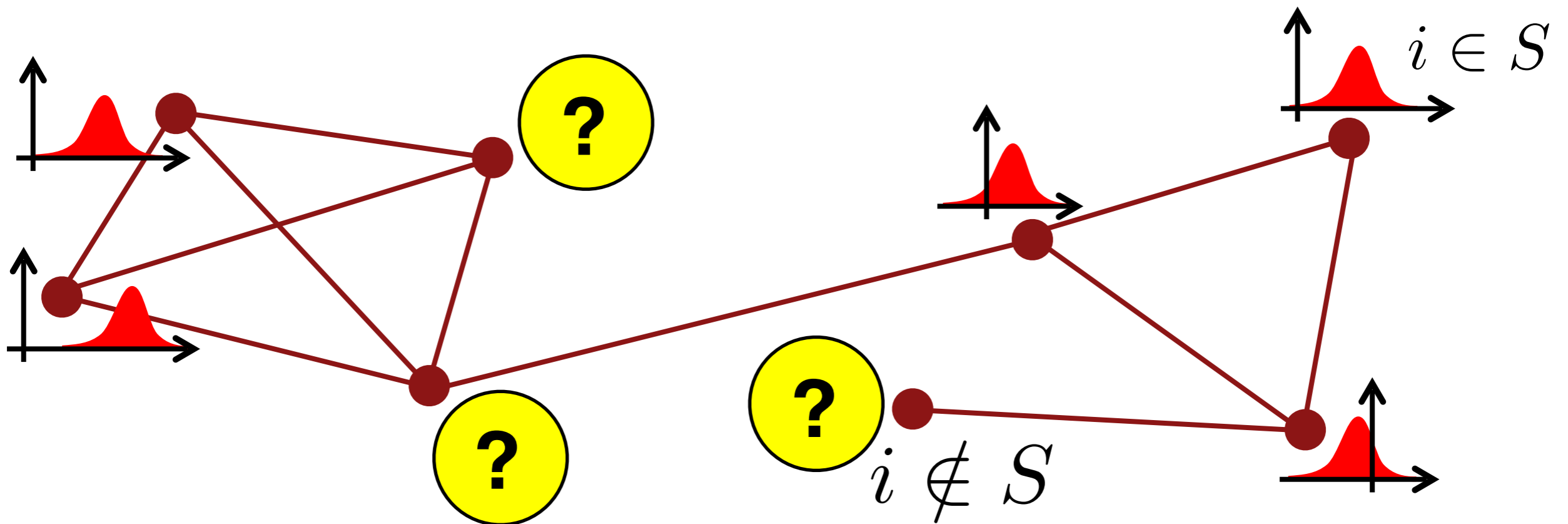
Applications: Brain Imaging



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

[Gramfort'16]

Wasserstein Propagation



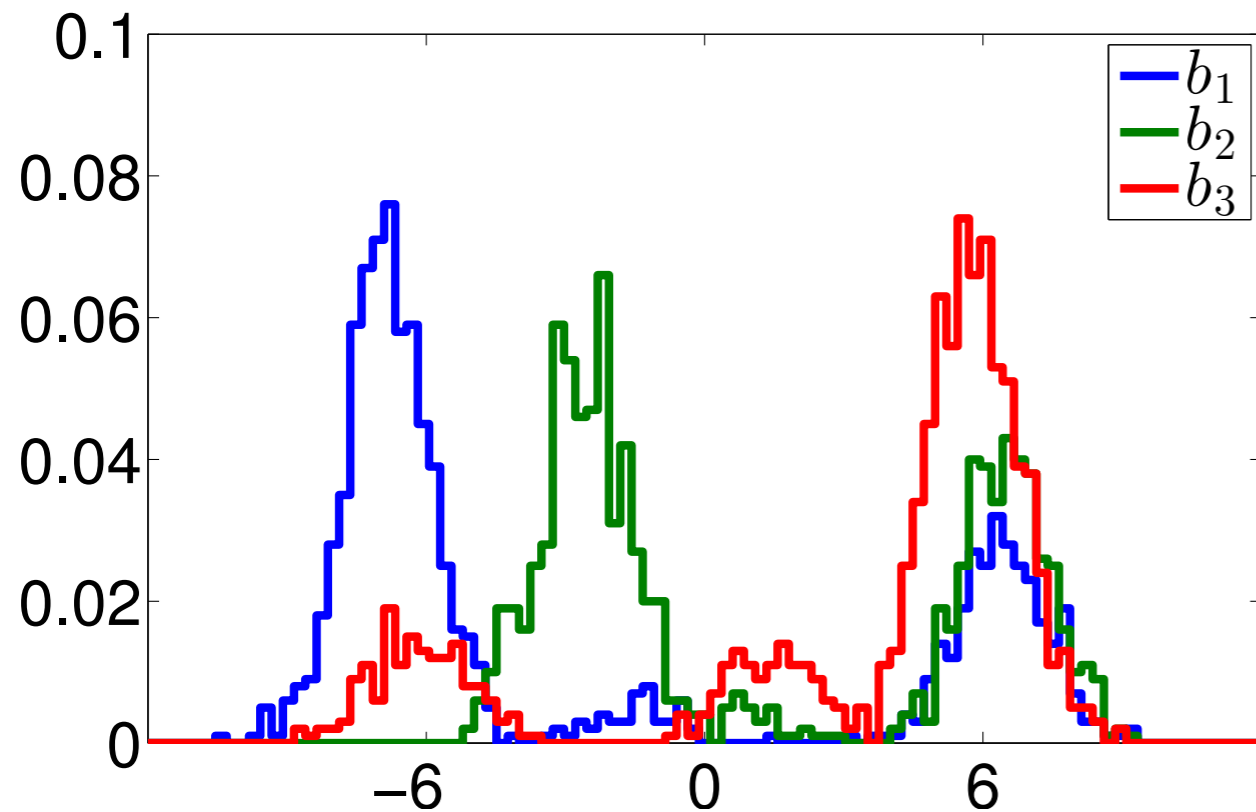
$$\min_{\substack{\mu_i \in \mathcal{P}(\Omega) \\ \mu_i \text{ fixed for } i \in S}} \sum_{(e_1, e_2) \in E} W_2^2(\mu_{e_1}, \mu_{e_2})$$

[Solomon'14]

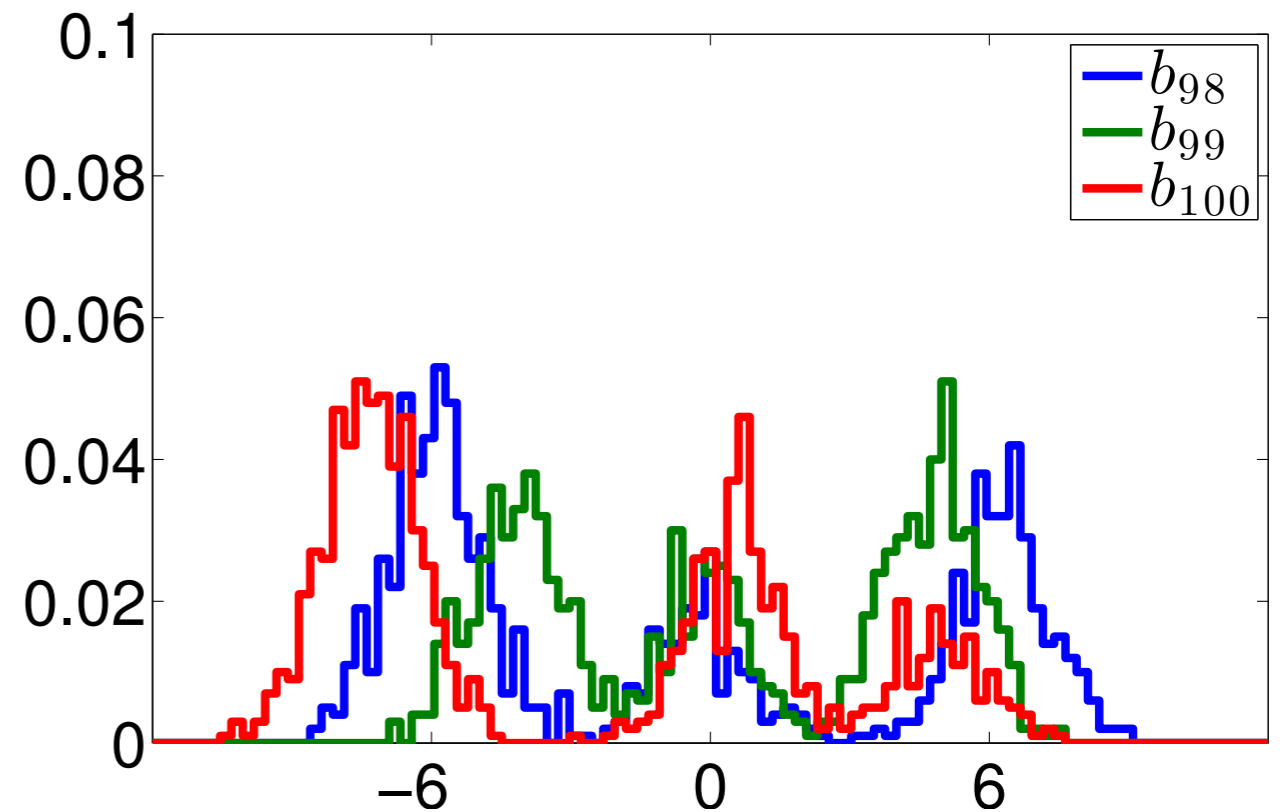
Dictionary Learning

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

Data samples



Data samples

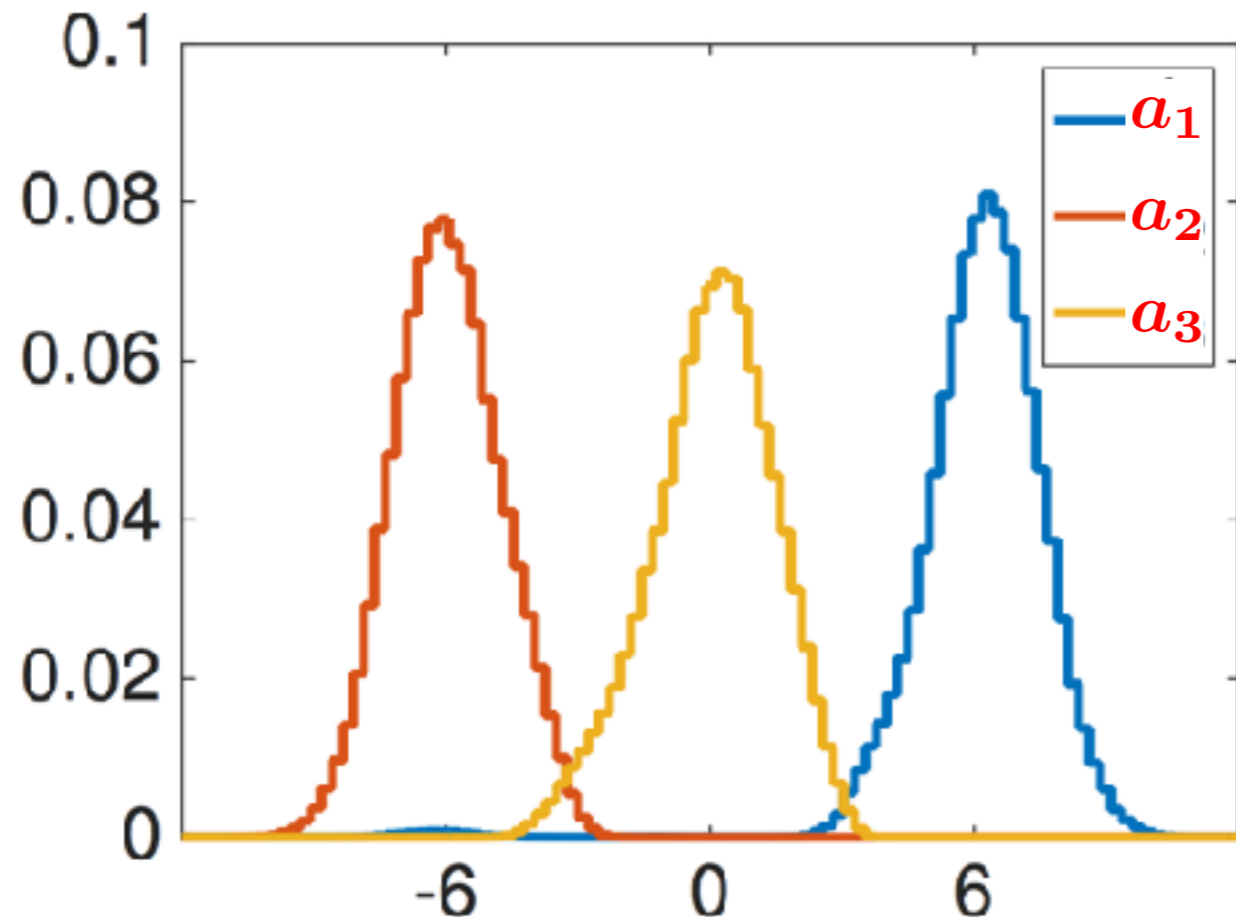


[Sandler'11] [Zen'14] [Rolet'16]

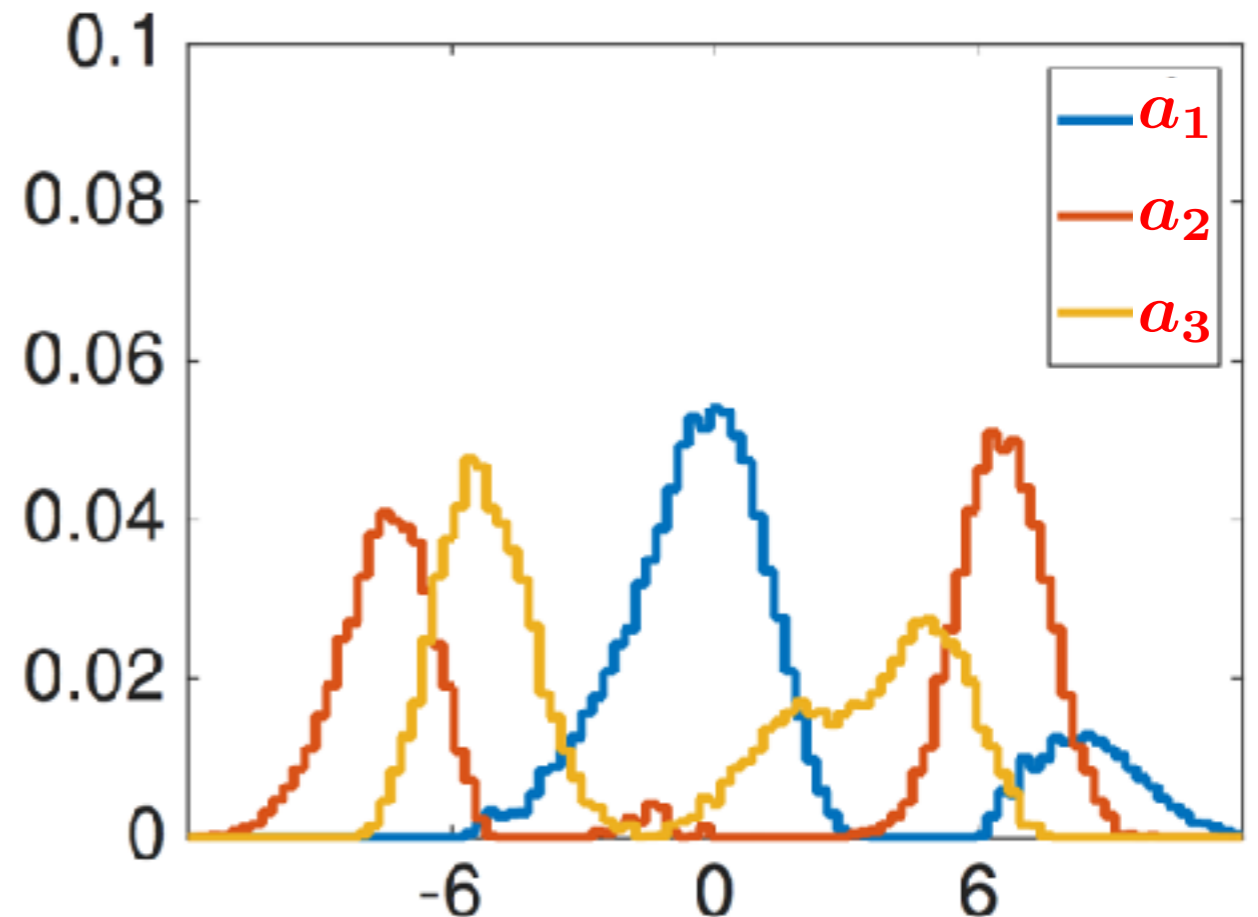
Dictionary Learning

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

Wasserstein NMF



KL NMF



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

OT Dictionary Learning

- **[Hoffman'98]** proposed to learn dictionaries (topics) for text, seen as histograms-of-words.

$$\Omega = \{\text{words}\}, \quad |\Omega| \approx 13,000$$

- Vector embeddings for words **[Mikolov'13]**
[Pennington'14] defines geometry:

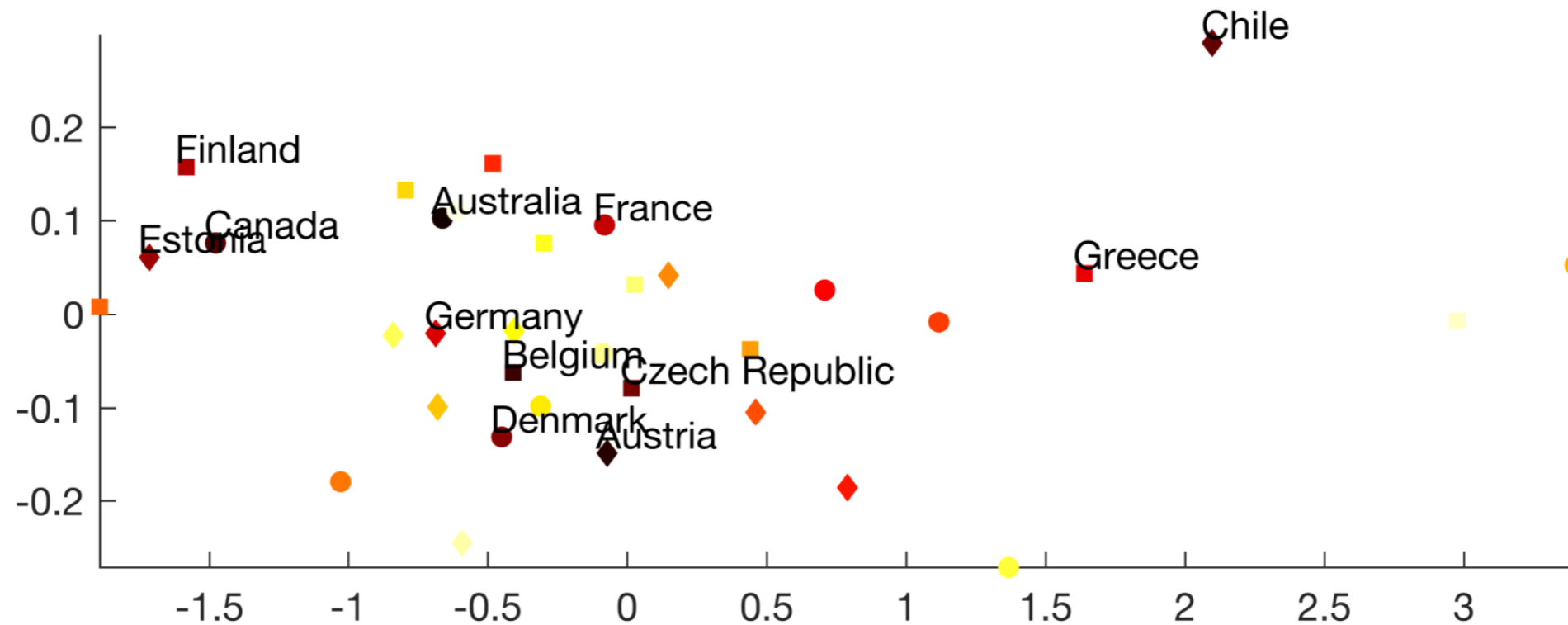
$$D(\text{public}, \text{car}) = \|x_{\text{public}} - x_{\text{car}}\|^2$$

- Data: 7,034 Reuters, 737 BBC sports news articles

Elliptical Embeddings

Multidimensional Scaling [MDS]

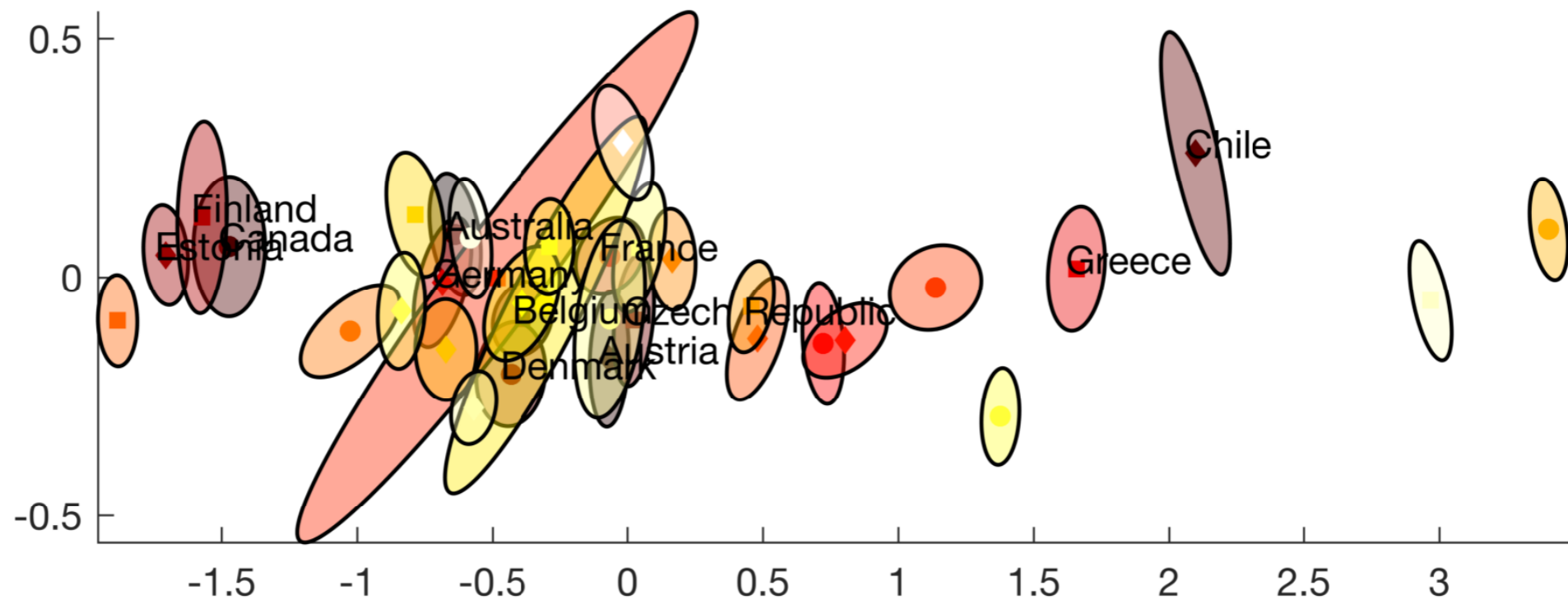
embed a metric space in R^2



Elliptical Embeddings

Multidimensional Scaling [MDS]

embed a metric space in elliptical distributions in $P(\mathbb{R}^2)$, W_2

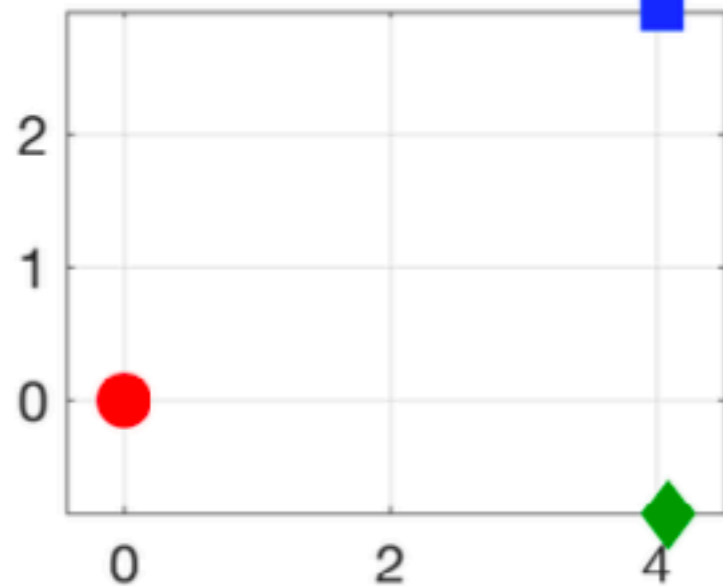


Elliptical Embeddings

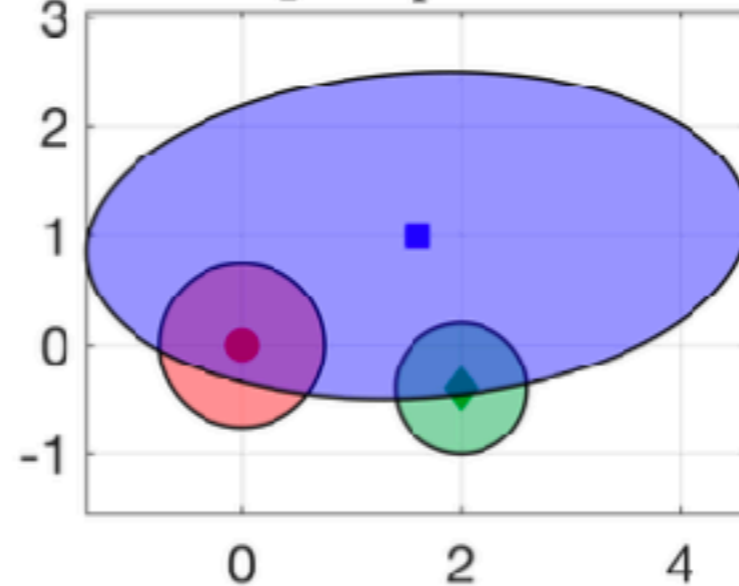
Visualization issue

need to shift to precision matrix to recover intuition

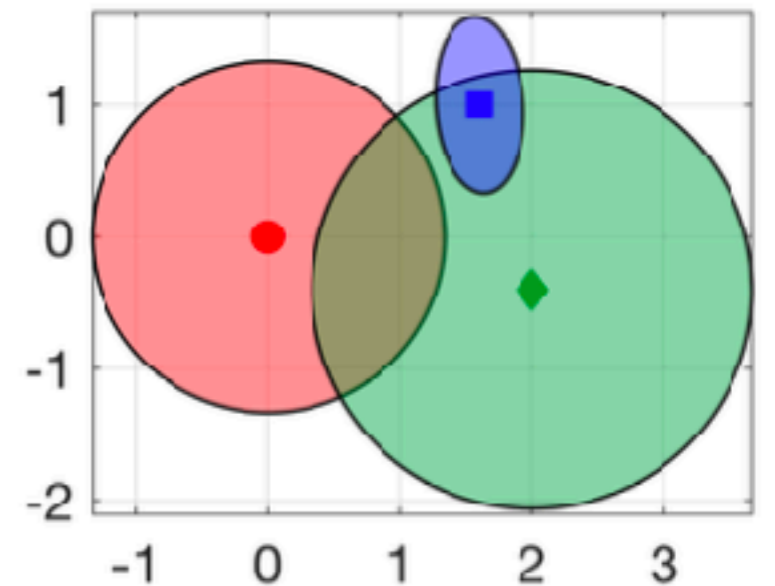
Points in 2D



Isometric W_2 Elliptical Embedding



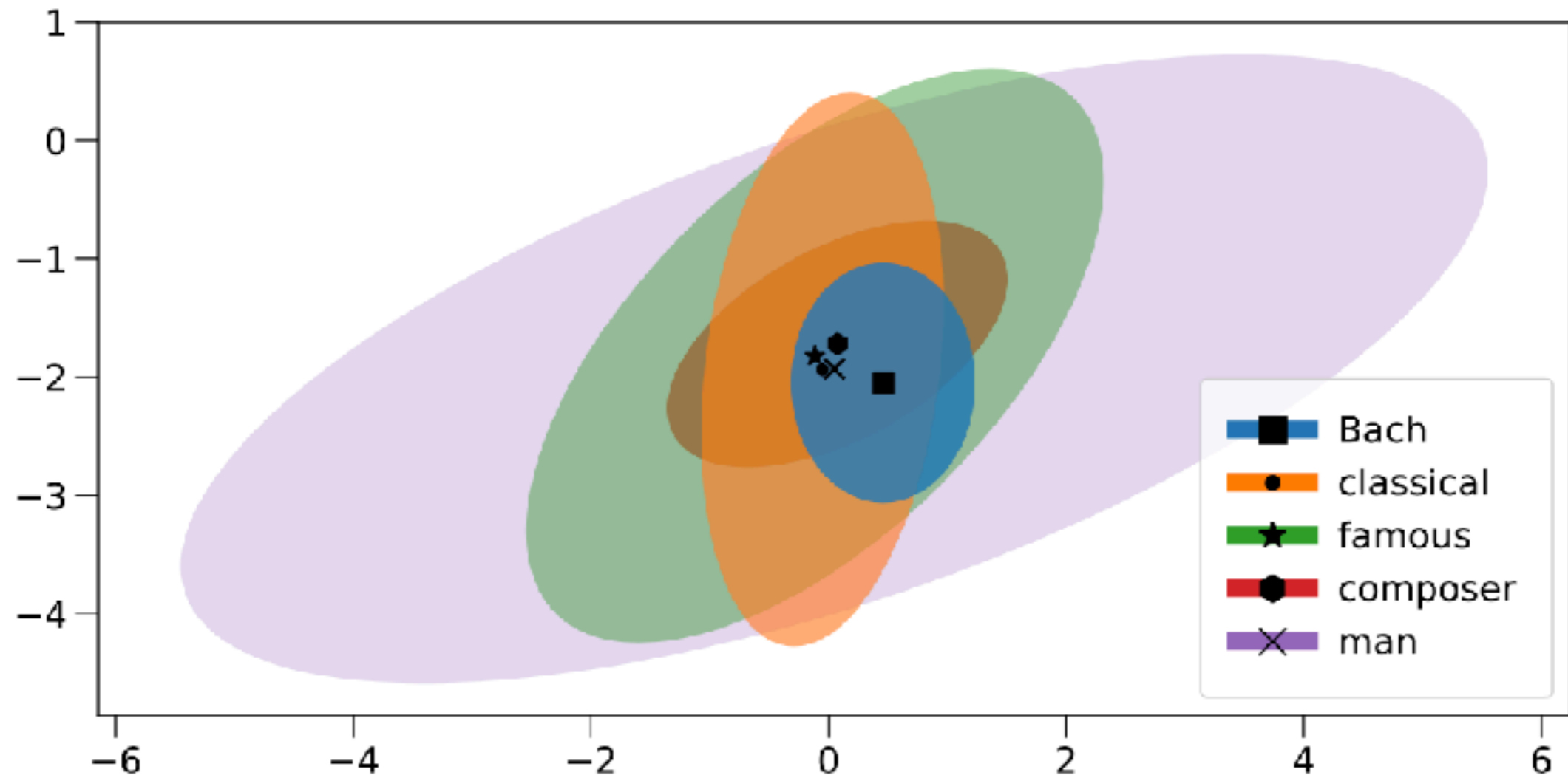
Precision Matrix Visualization



Elliptical Embeddings

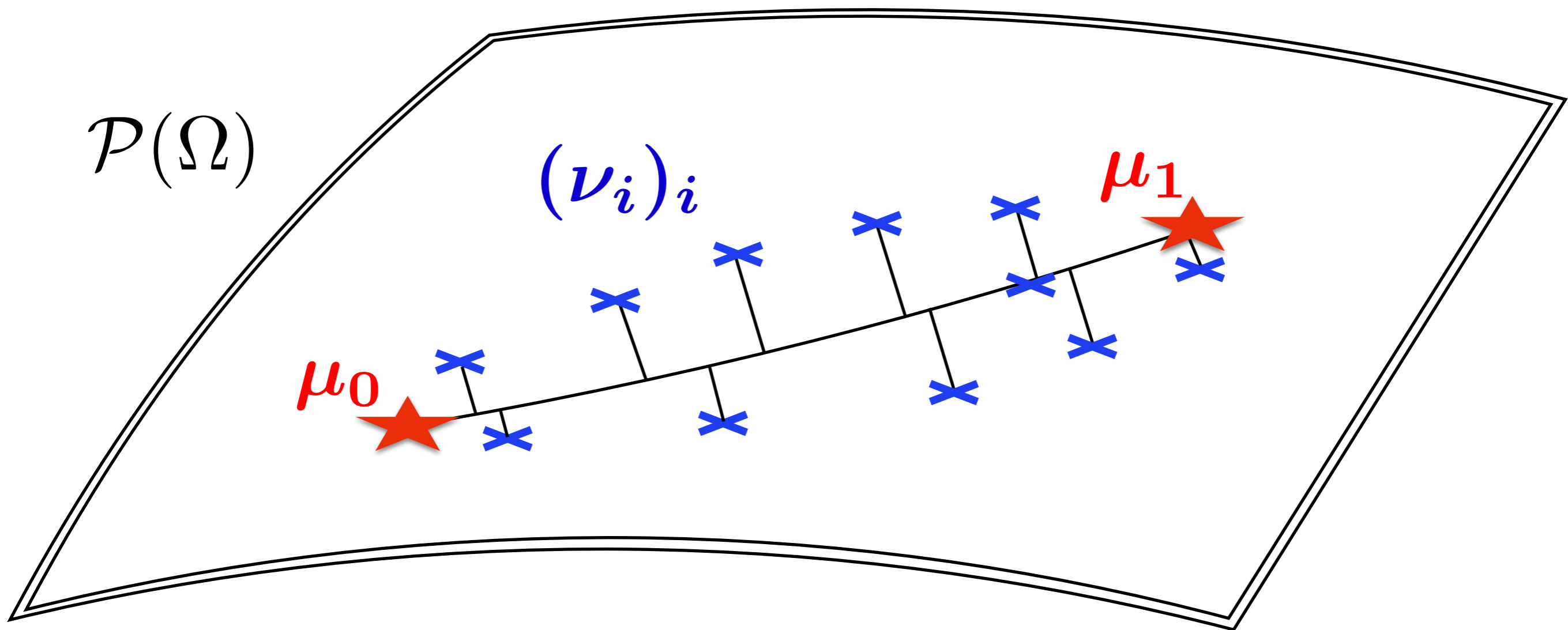
Word Embeddings

Compute elliptical distribution representations for Words

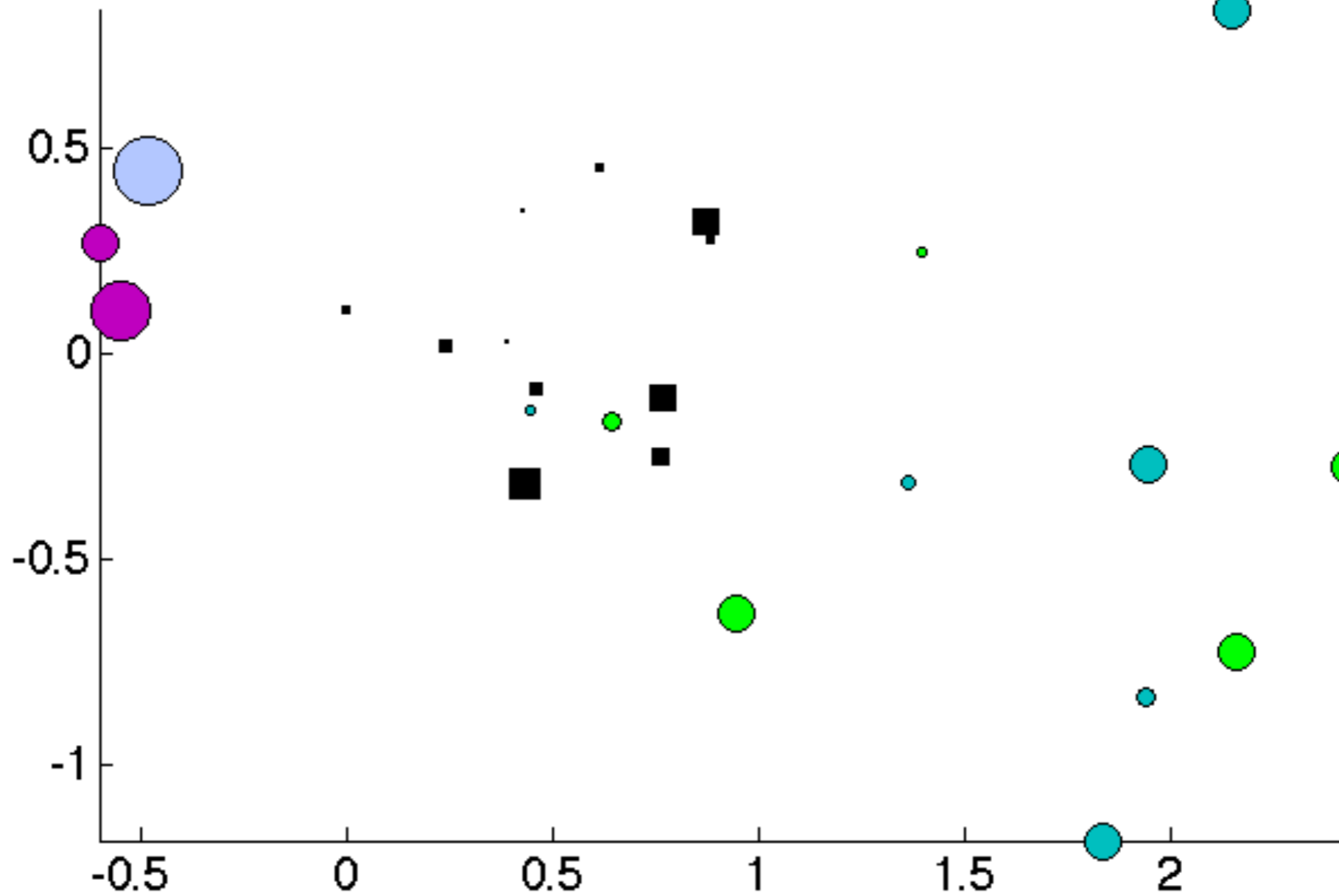


Wasserstein PCA

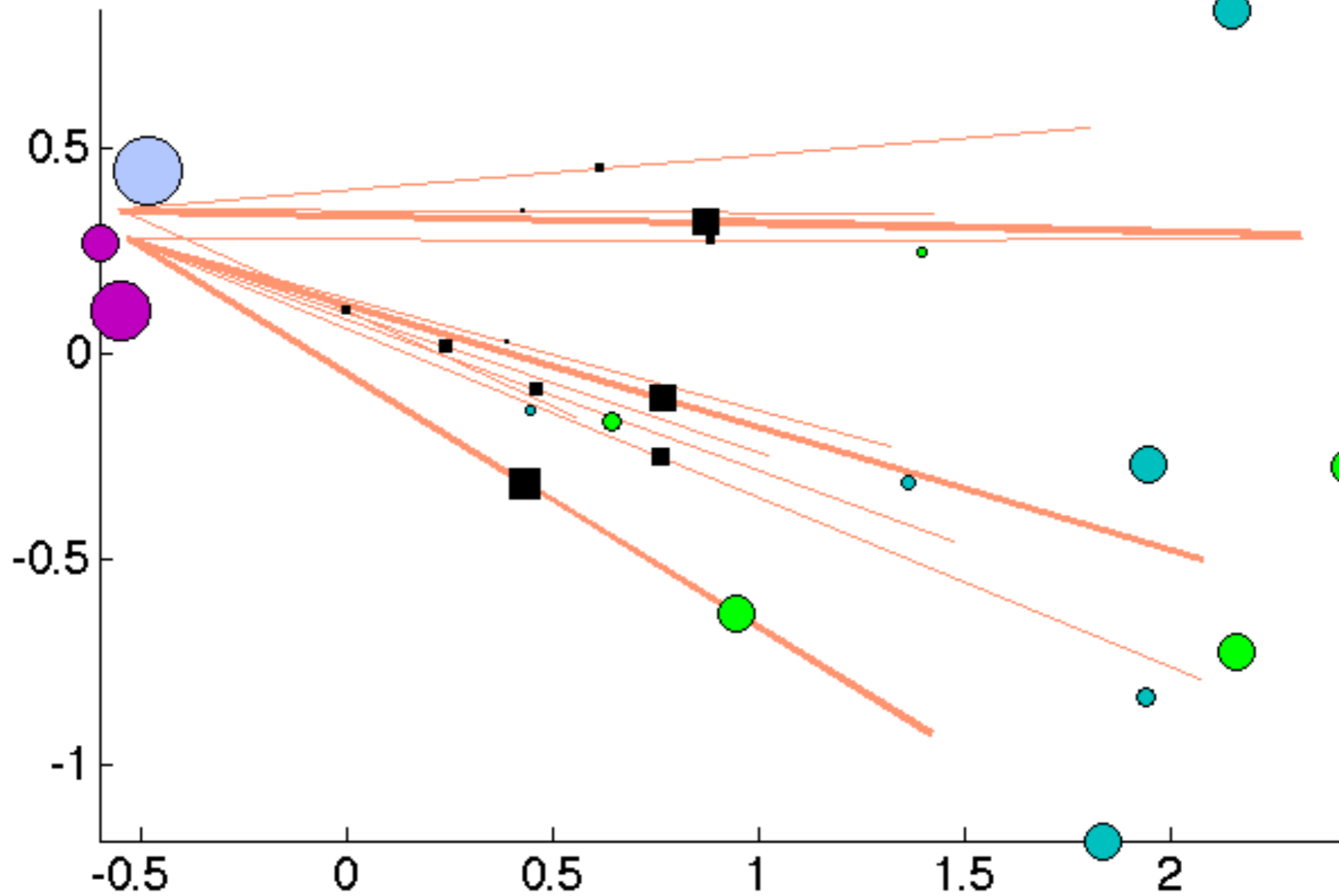
$$\min_{\mu_0, \mu_1} \sum_{i=1}^N \min_t W_2^2(\rho_{\mu_0 \rightarrow \mu_1}^t, \nu_i)$$



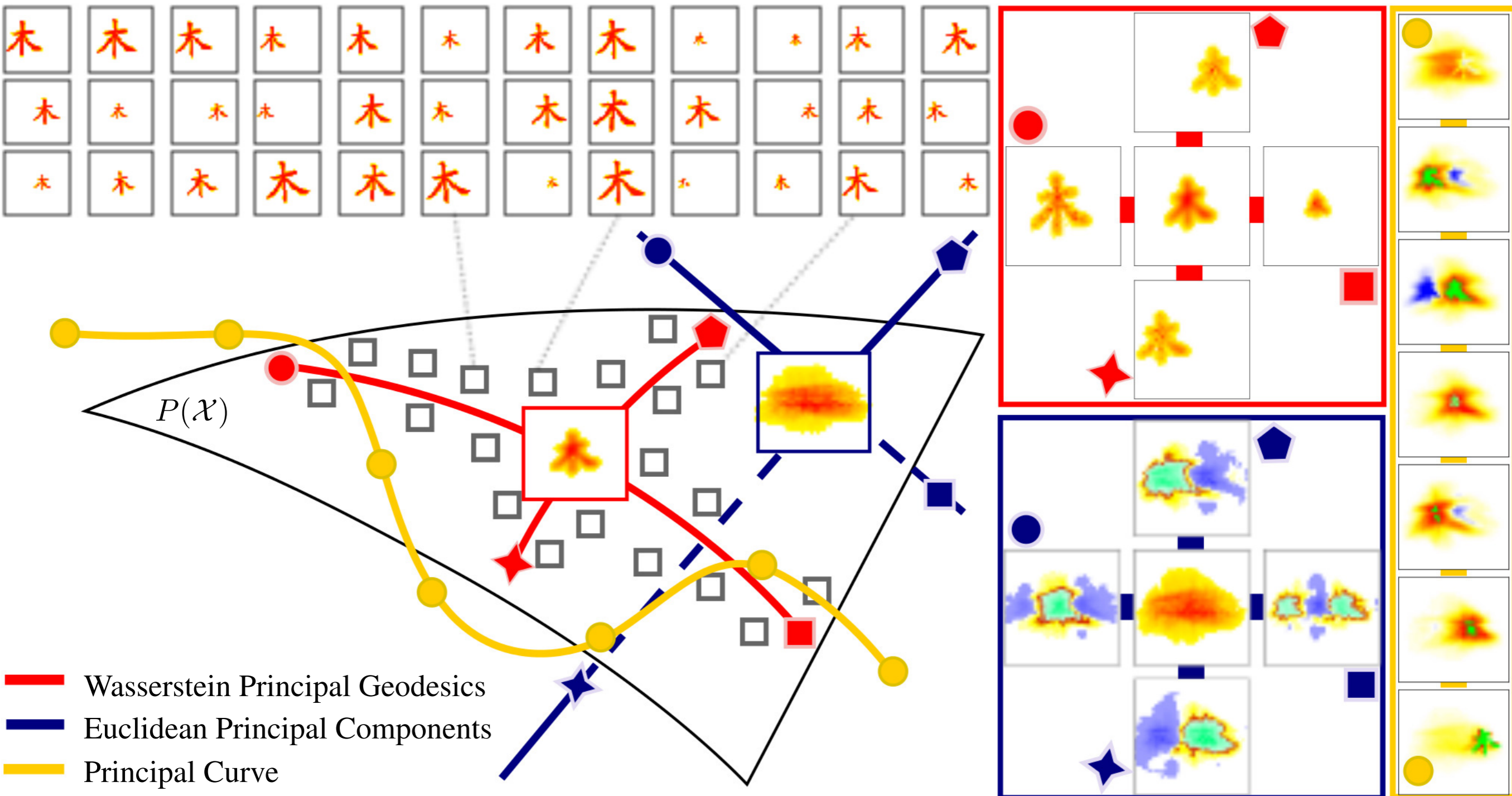
Wasserstein PCA



On Empirical Measures



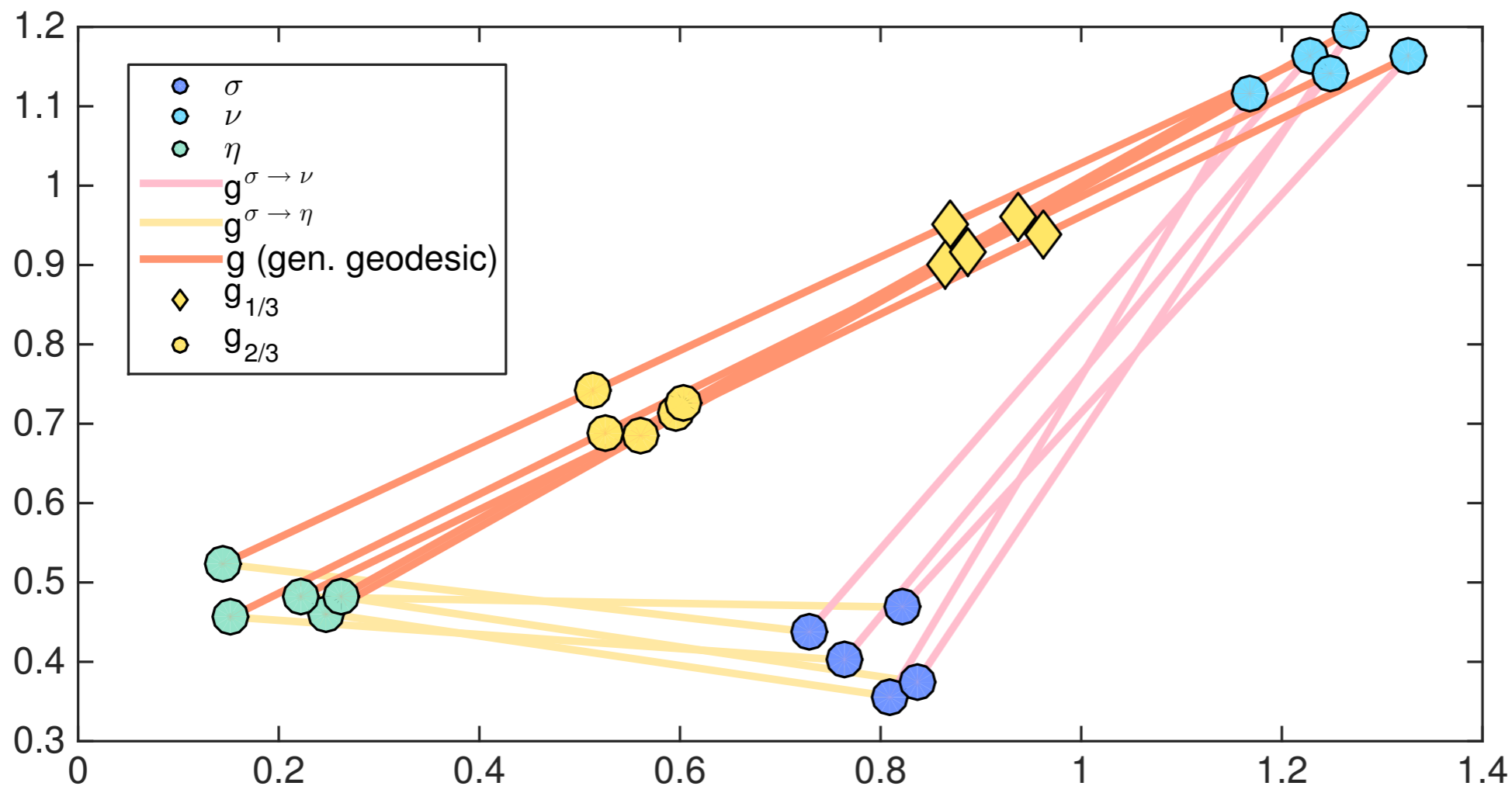
Wasserstein PCA vs. Euclidean PCA



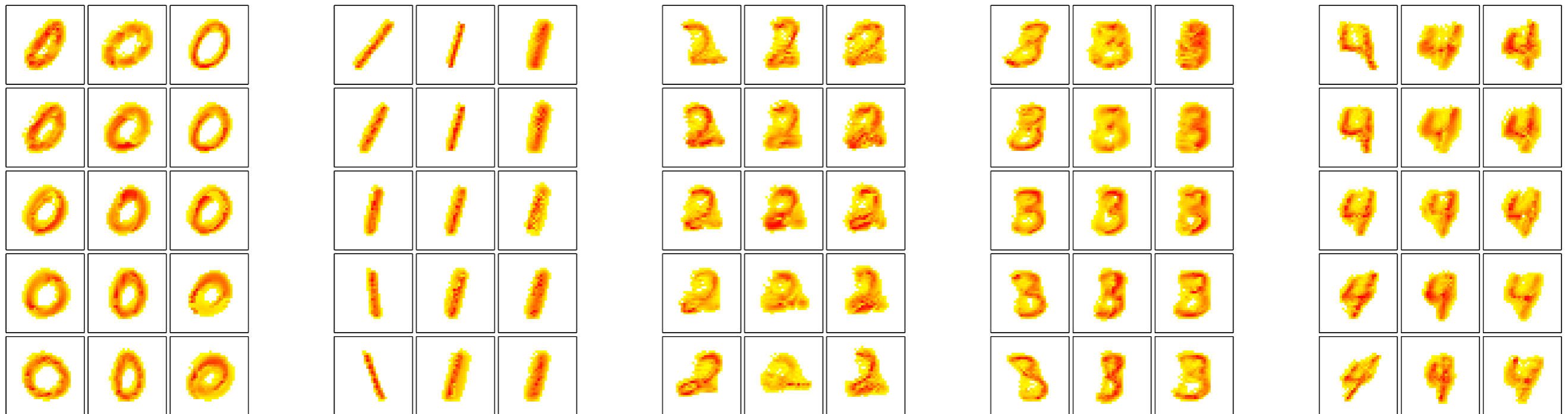
[Ambrosio'06] Generalized Geodesics

$$\min_{\mathbf{v}_1, \mathbf{v}_2 \in L^2(\bar{\nu}, \Omega)} \sum_{i=1}^N \min_{t \in [0, 1]} W_2^2(g_t(\mathbf{v}_1, \mathbf{v}_2), \nu_i) + \lambda R(\mathbf{v}_1, \mathbf{v}_2),$$

subject to $\begin{cases} g_t(\mathbf{v}_1, \mathbf{v}_2) = (\text{Id} - \mathbf{v}_1 + t(\mathbf{v}_1 + \mathbf{v}_2)) \# \bar{\nu} \\ \text{Id} - \mathbf{v}_1 \text{ and } \text{Id} + \mathbf{v}_2 \text{ are Monge maps from } \bar{\nu} \end{cases}$



Generalized Principal Geodesics



For each digit, 1,000 MNIST images

[Seguy'15]

Inverse Wasserstein Problems

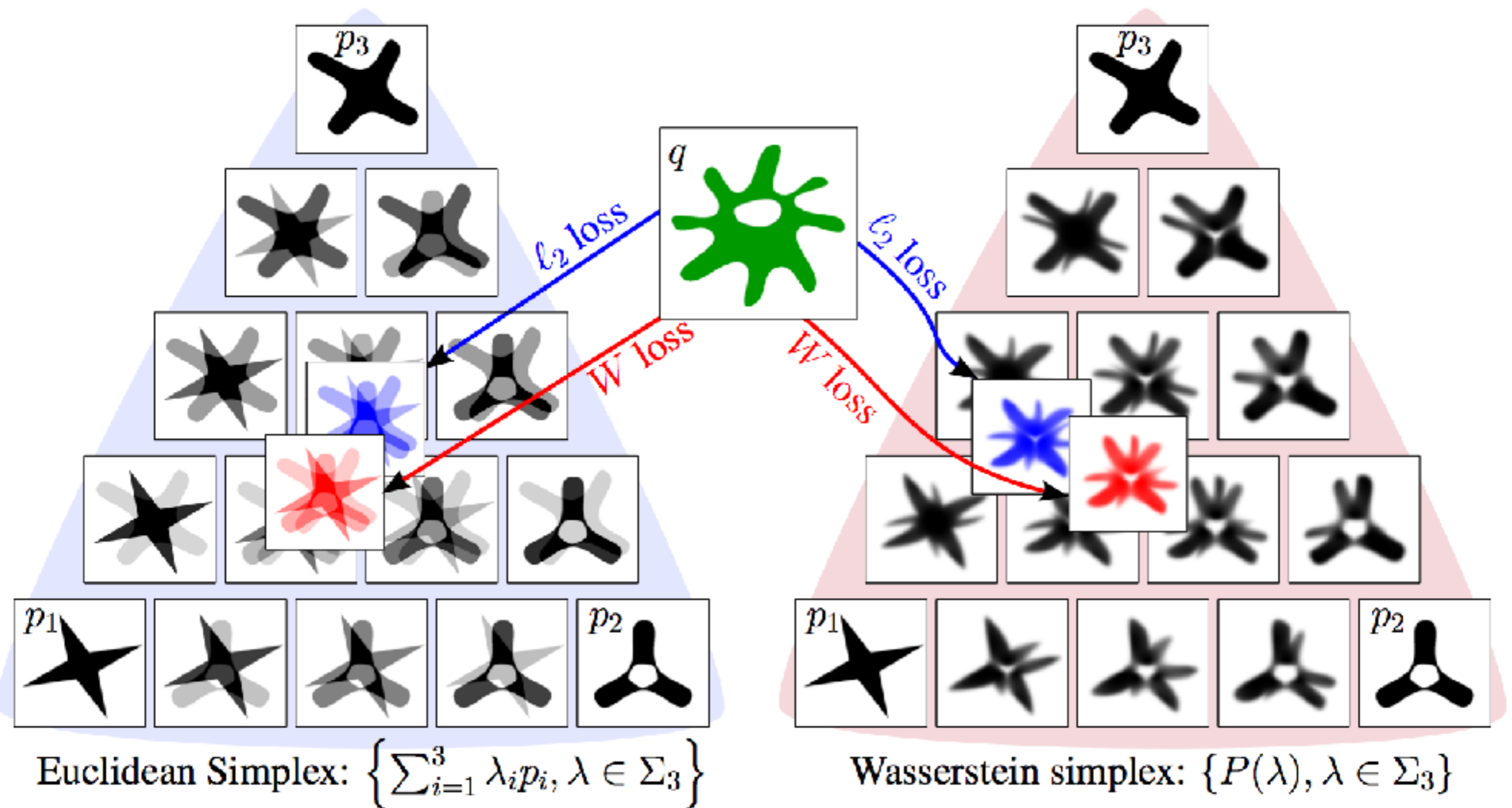
- consider Barycenter operator:

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

- address now **Wasserstein inverse problems**:

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

Wasserstein Inverse Problems



Barycenters = Fixed Points

Prop. [BCCNP'15] Consider $B \in \Sigma_d^N$ and let $U_0 = \mathbf{1}_{d \times N}$, and then for $l \geq 0$:

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

Using Truncated Barycenters

- instead of using the exact barycenter

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

- use instead the L-iterate barycenter

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

- Differentiate using **the chain rule.**

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

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

[Bonneel'16]

Application: Color Grading



Application: Color Grading



$$\lambda_0 = 0.03$$



$$\lambda_1 = 0.12$$



$$\lambda_2 = 0.40$$

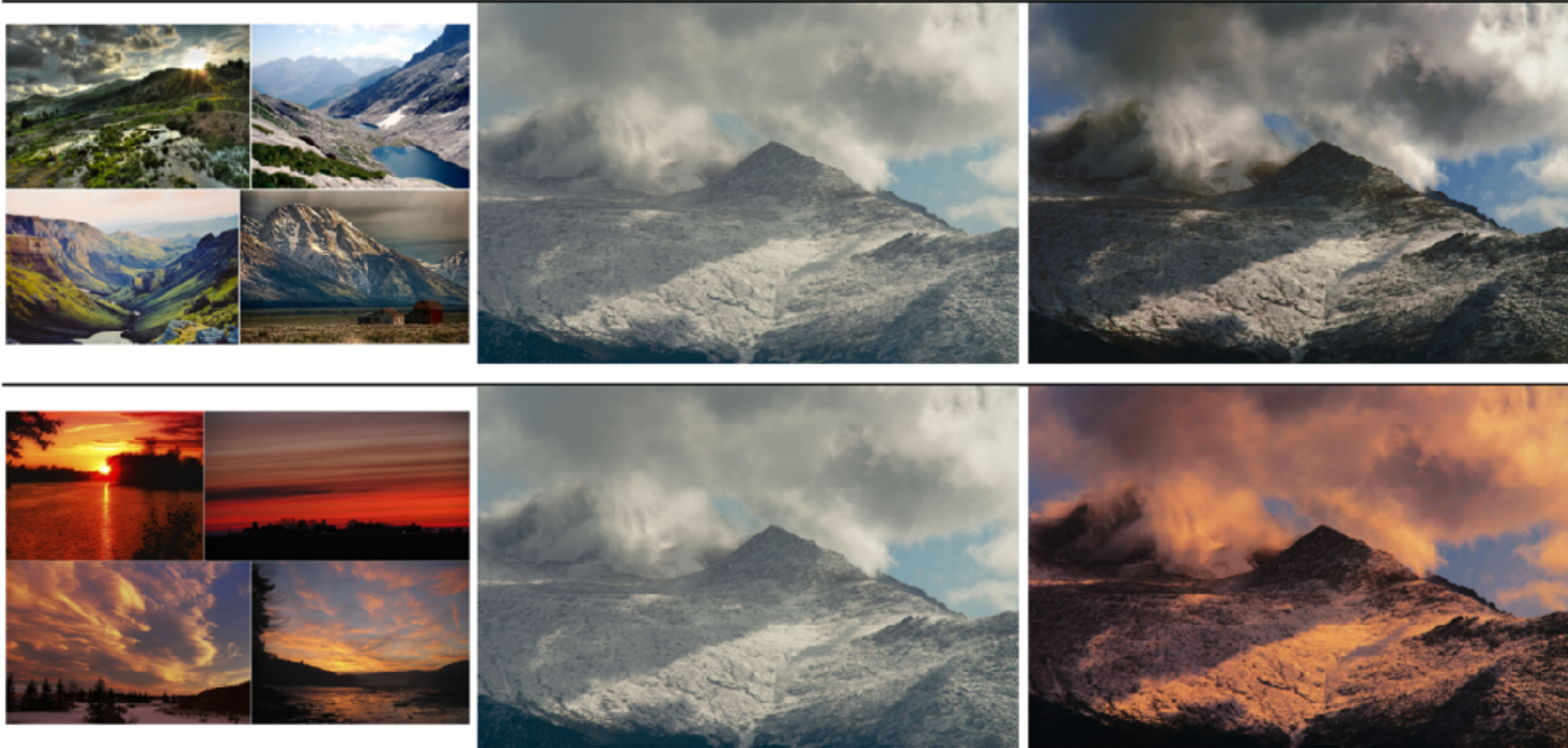


$$\lambda_3 = 0.43$$

Application: Color Grading



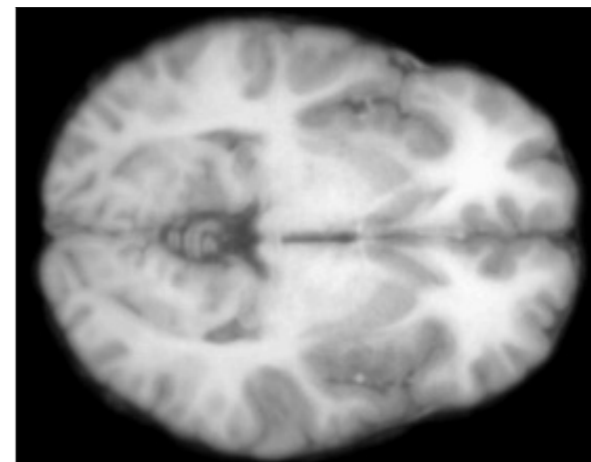
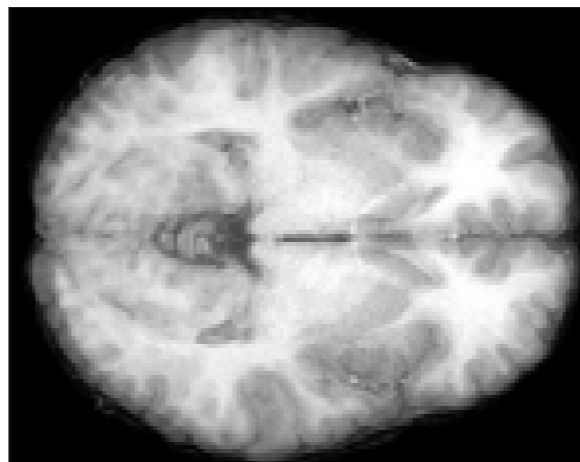
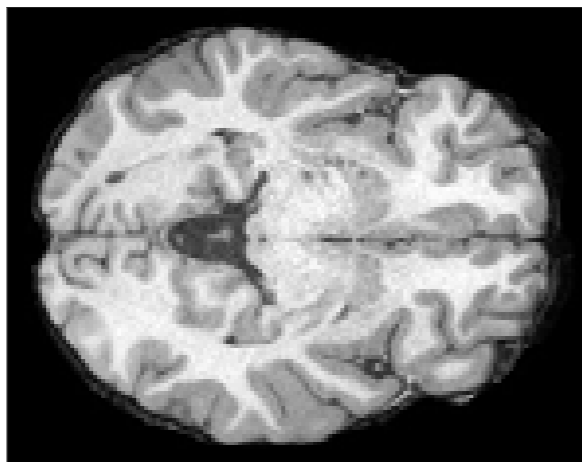
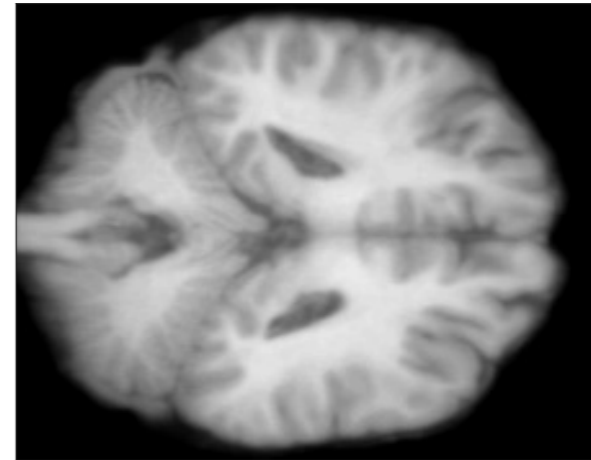
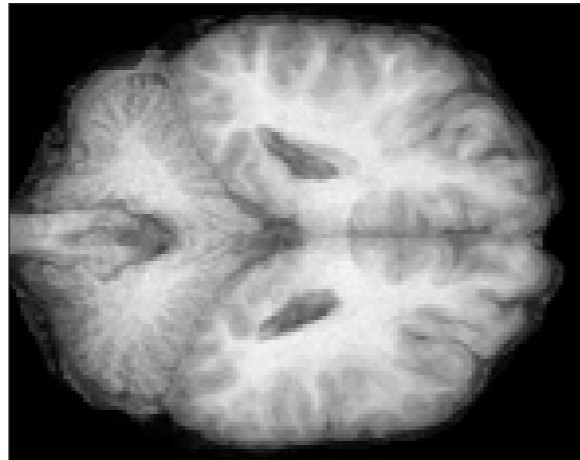
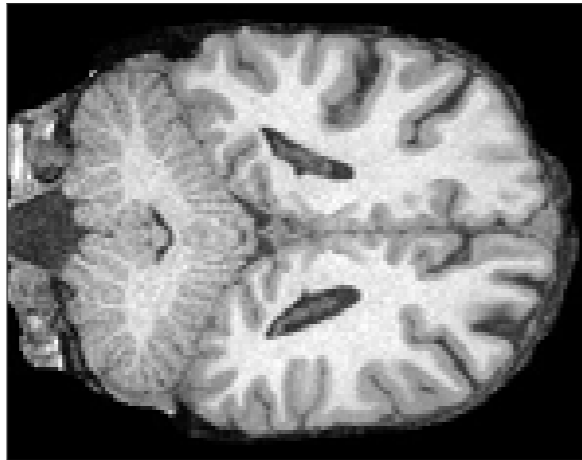
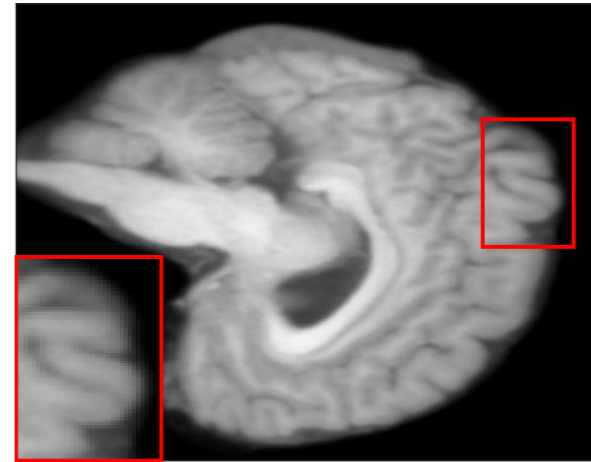
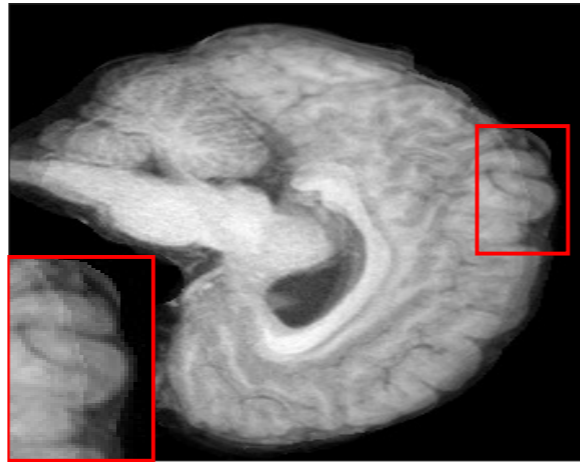
Application: Color Grading



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

[BPC'16]

Application: Brain Mapping

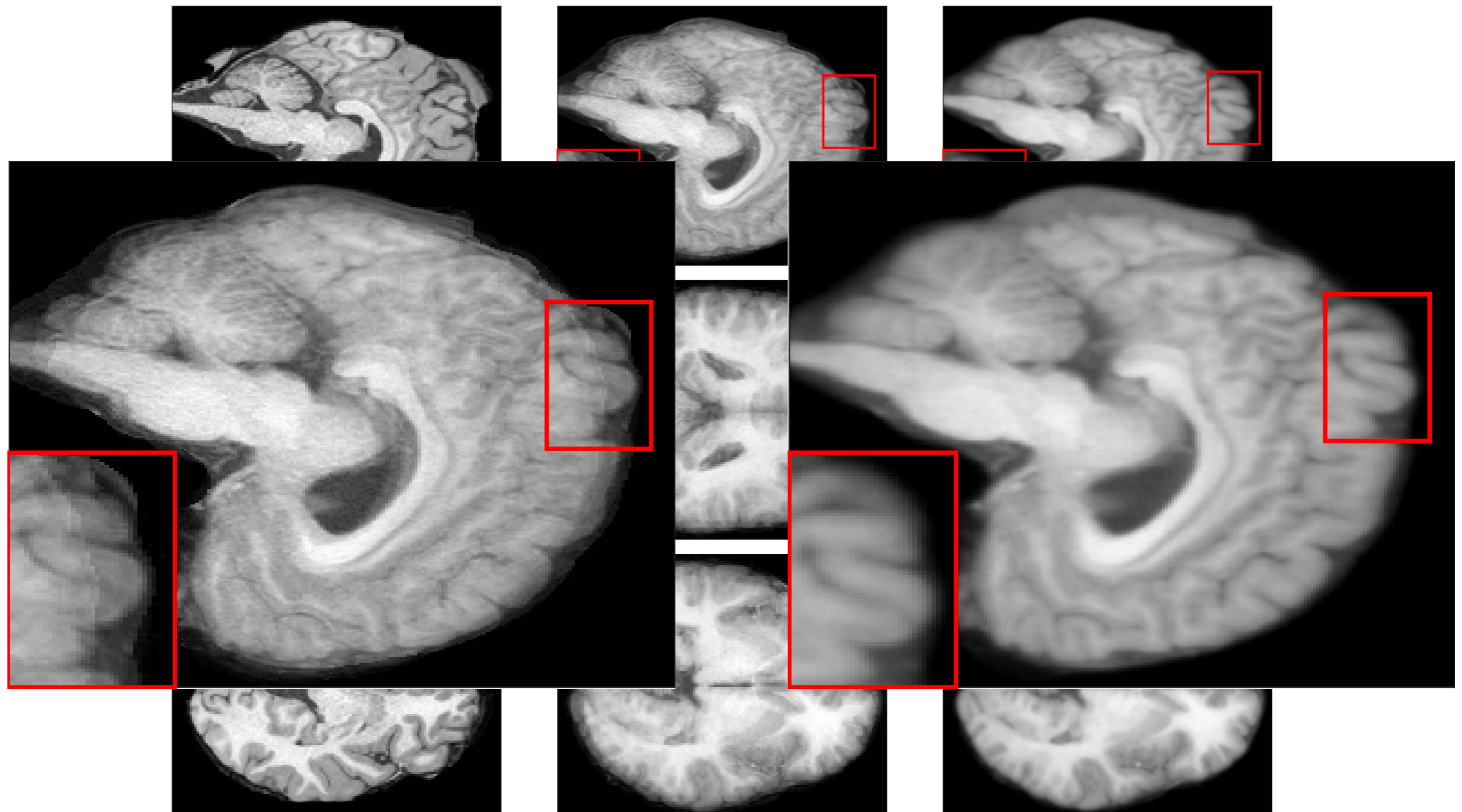


Original

Euclidean
projection

Wasserstein
projection

Application: Brain Mapping



Original

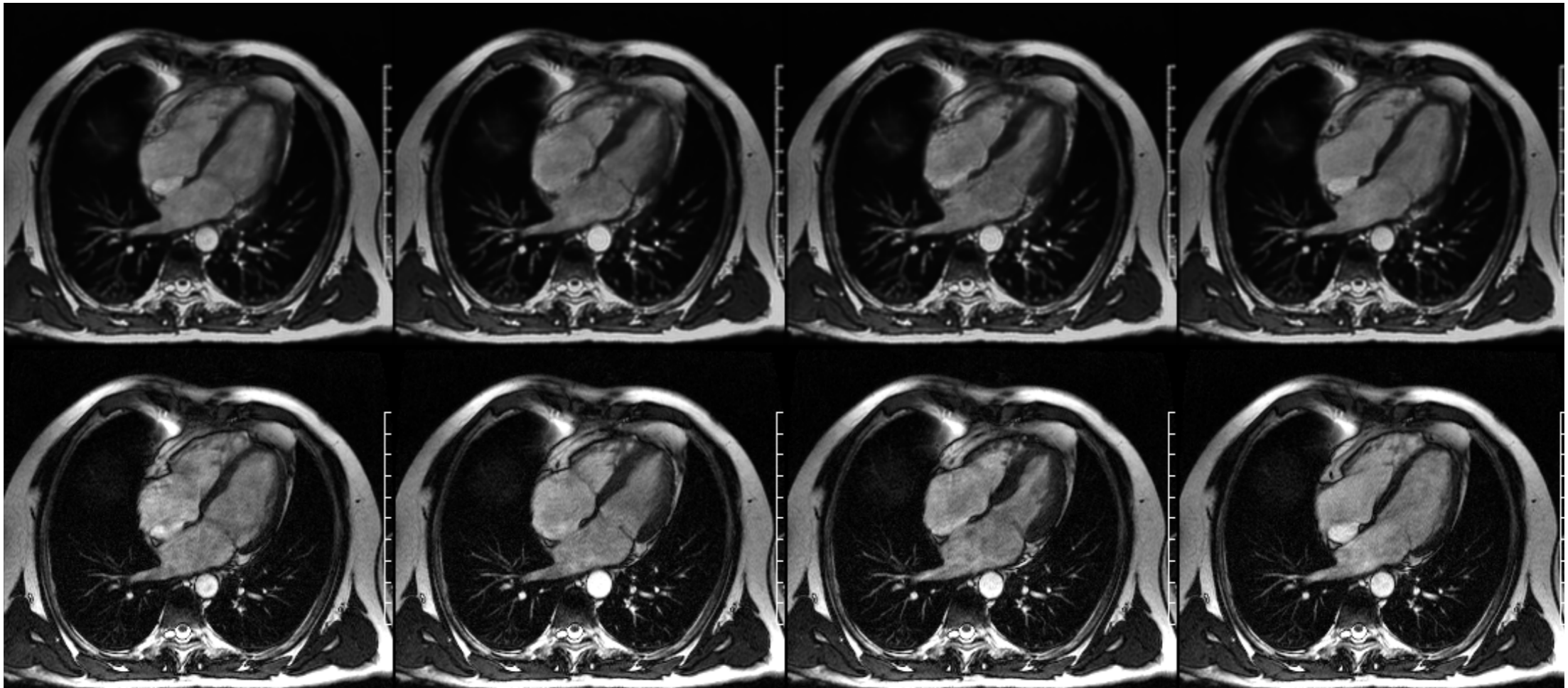
Euclidean
projection

Wasserstein
projection

end-to-end W Dictionary Learning

$$\min_{\mathbf{A} \in (\Sigma_n)^K, \mathbf{\Lambda} \in (\Sigma_K)^N} \sum_{i=1}^N \mathcal{L}(b_i, a(\lambda_i))$$

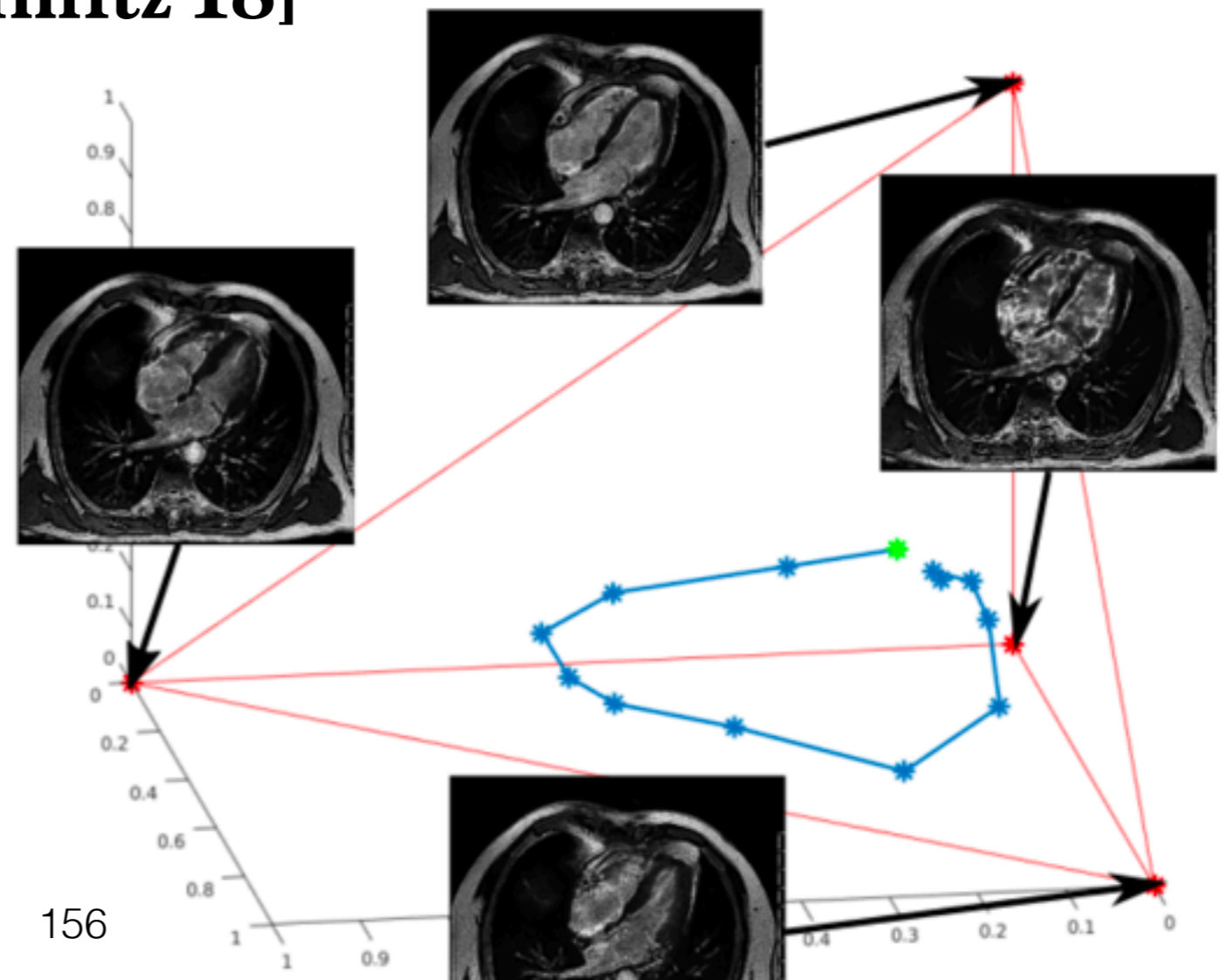
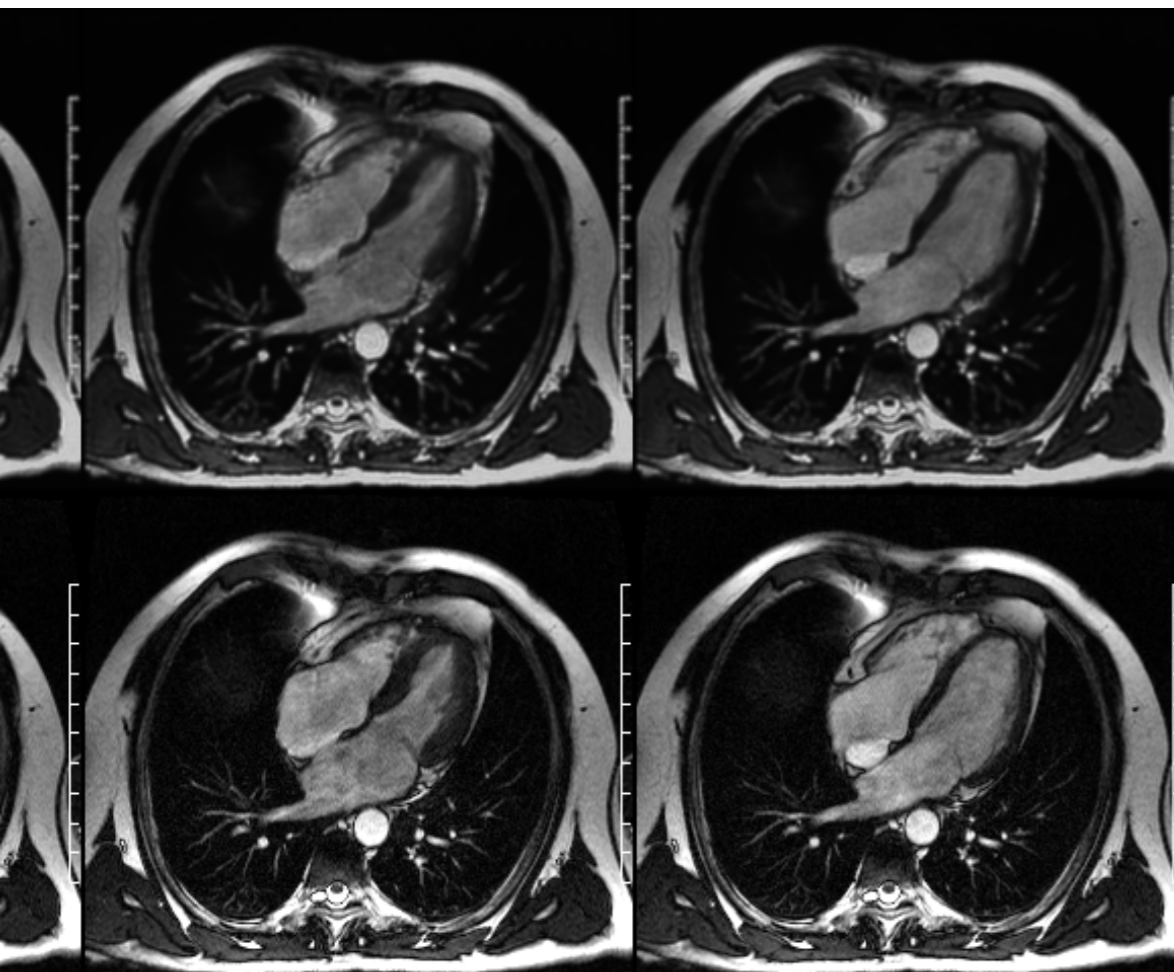
[Schmitz'18]



end-to-end W Dictionary Learning

$$\min_{\mathbf{A} \in (\Sigma_n)^K, \mathbf{\Lambda} \in (\Sigma_K)^N} \sum_{i=1}^N \mathcal{L}(b_i, a(\lambda_i))$$

[Schmitz'18]



Distributionally Robust Optimization

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

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

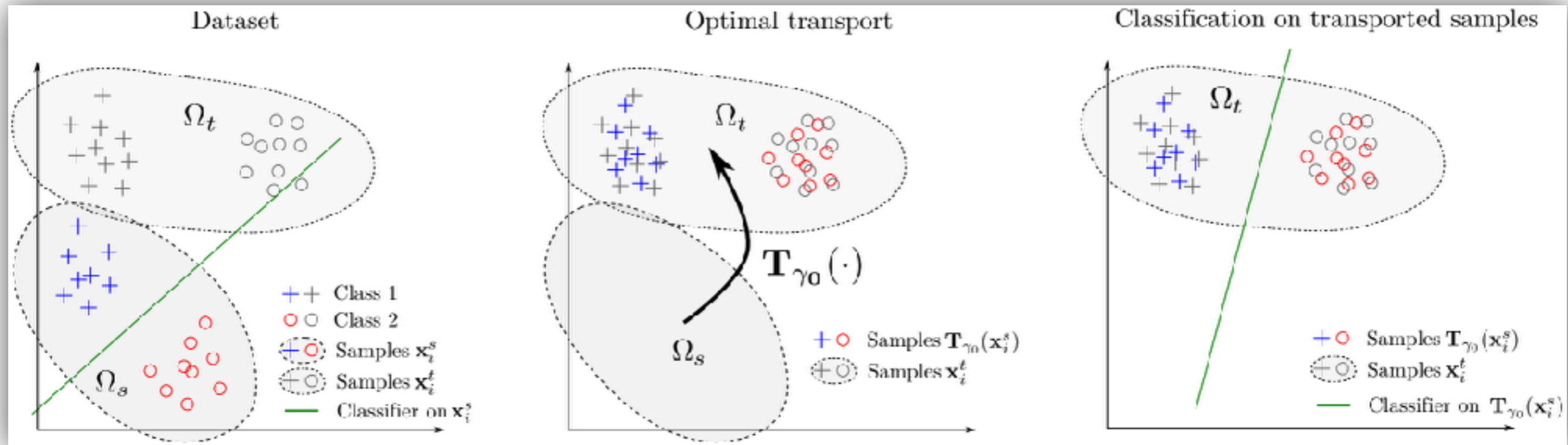
Learning with Wasserstein Ambiguity

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

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

[Courty'16]

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]

Minimum Kantorovich Estimation



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.

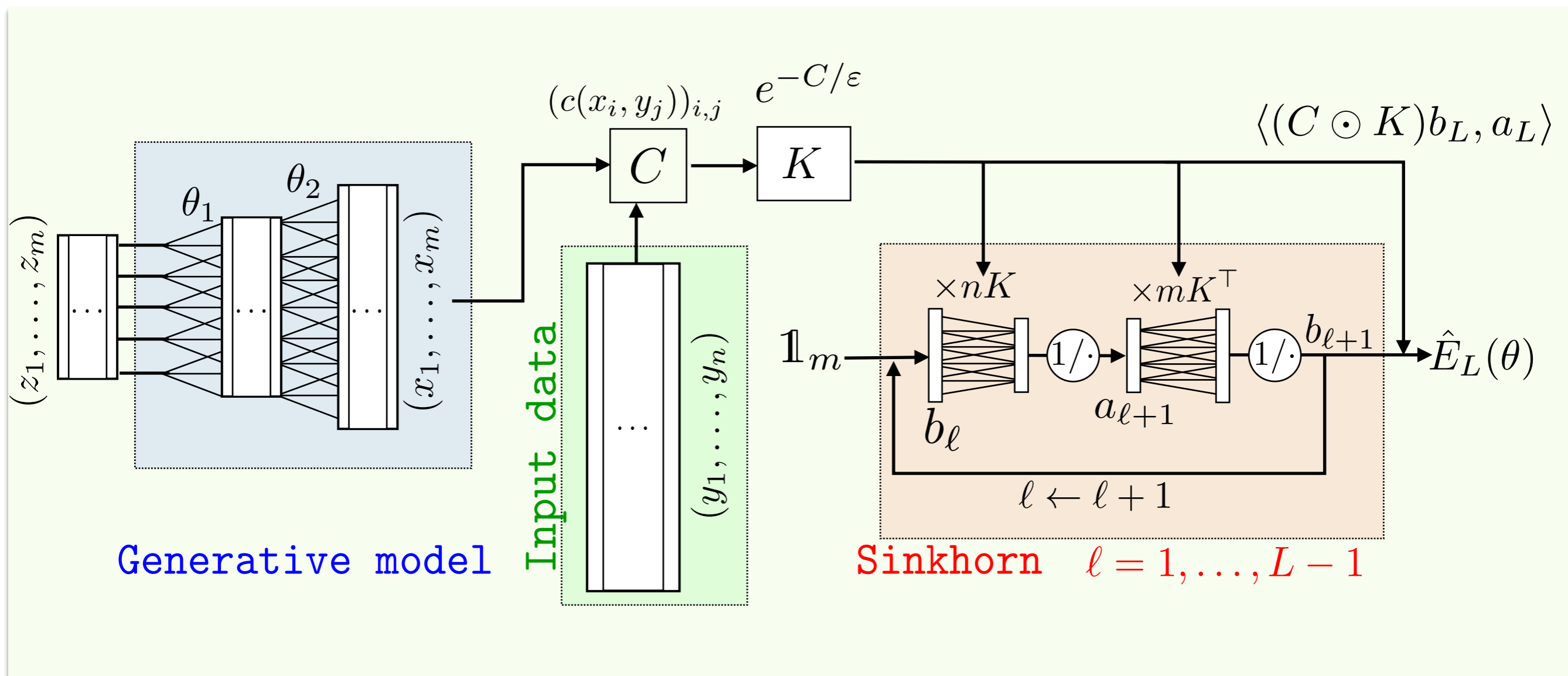
[Arjovsky'17] (WGAN) uses a NN to approximate dual solutions and recover gradient w.r.t. parameter

[Bernton'17] (*Wasserstein ABC*)

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

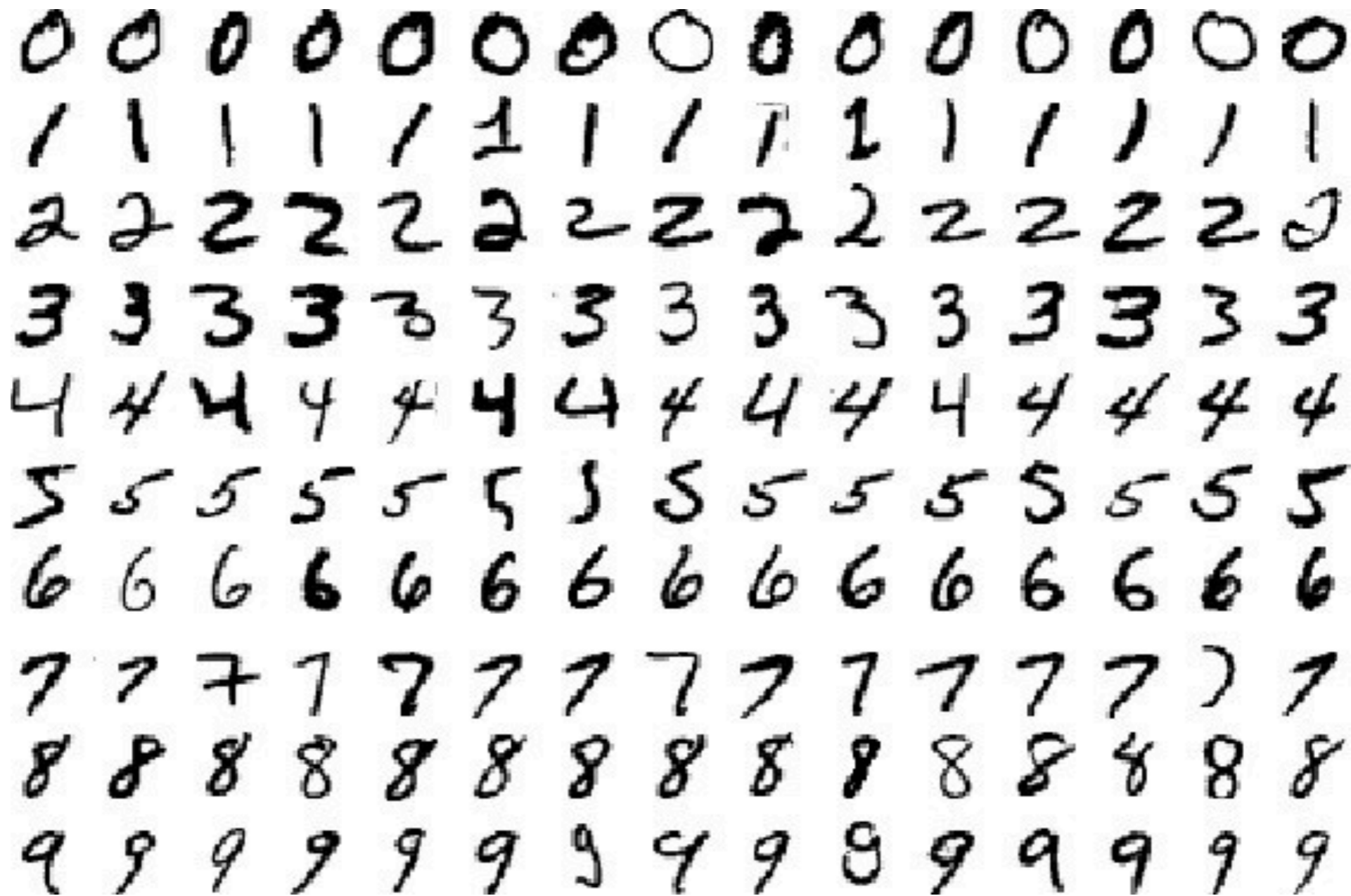
Proposal: Autodiff OT using Sinkhorn

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



[GPC'17]

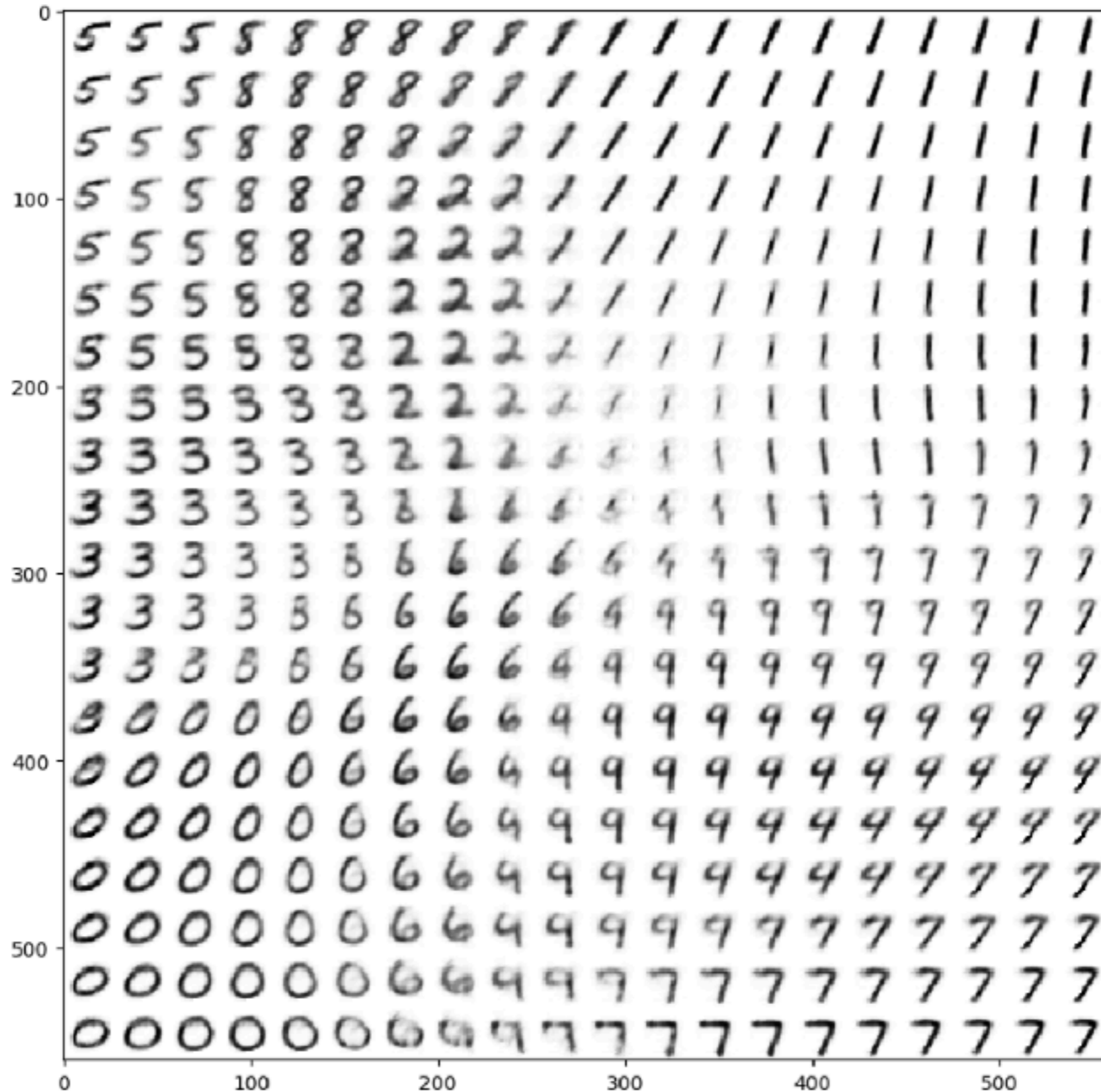
Example: MNIST, Learning f_{θ}



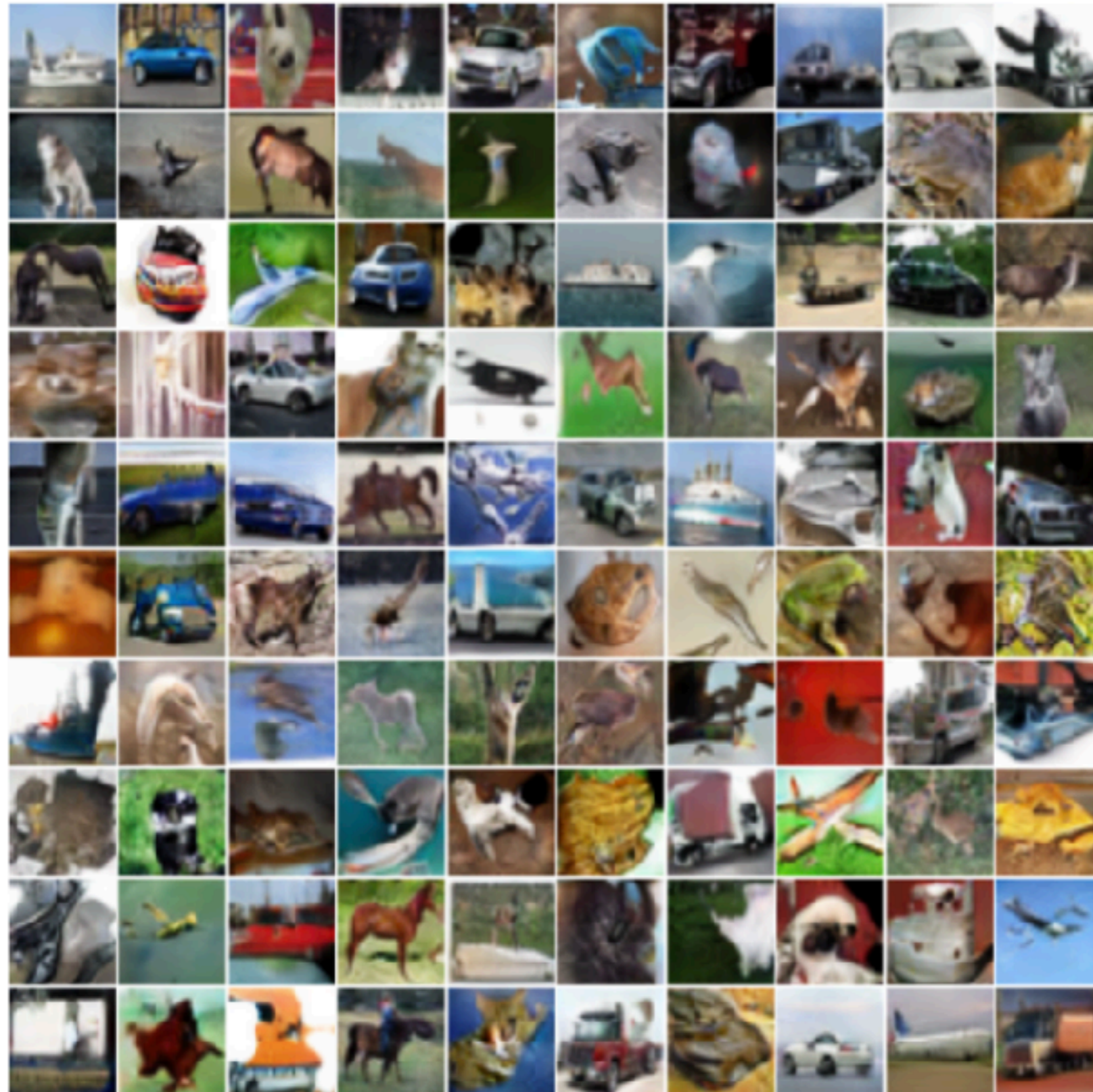
Example: MNIST, Learning f_{θ}

Latent space

$[0, 1]^2$



Example: Generation of Images



Example: Generation of Images



Concluding Remarks

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

What I could not talk about...

- Very large 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
- Applications to sampling.

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