# Optimal Transport in Data Sciences: Why and How?

## Marco Cuturi

*Joint work with*
G. Peyré, F. Bach, A. Genevay *(ENS)*
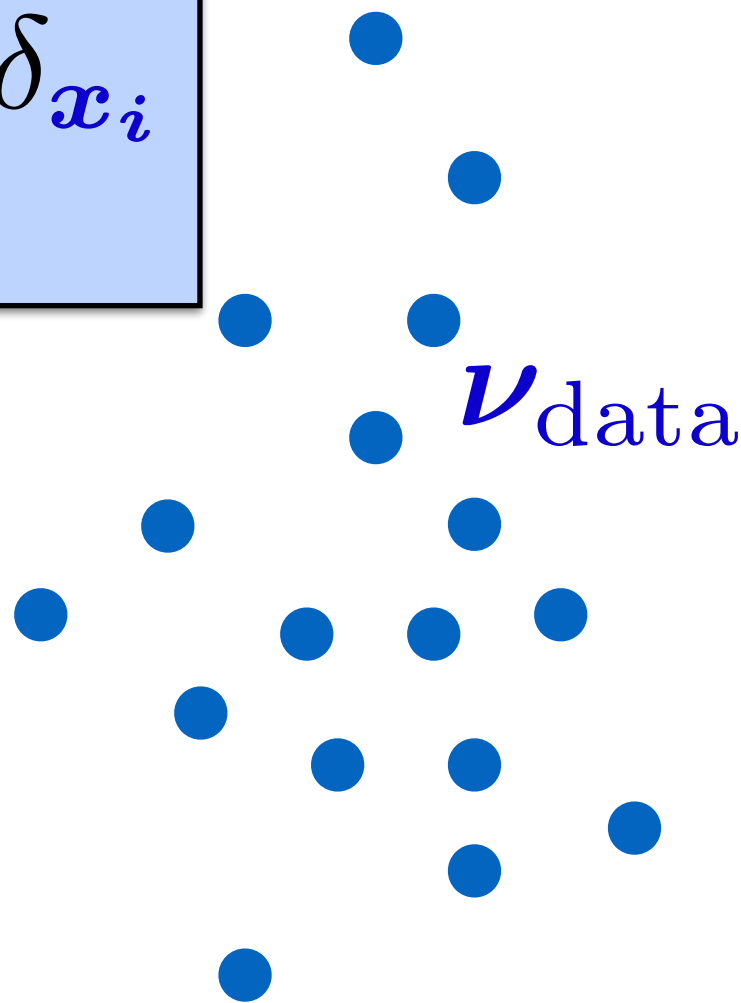N. Bonneel *(INRIA)* A. Rolet *(Kyoto)* J. Solomon *(MIT)*

**https://optimaltransport.github.io/**

We collect data

$$\boldsymbol{\nu}_{\text{data}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{\boldsymbol{x_i}}$$

$\boldsymbol{\nu}_{\text{data}}$

# Statistics 0.1 : Density Fitting

We collect data

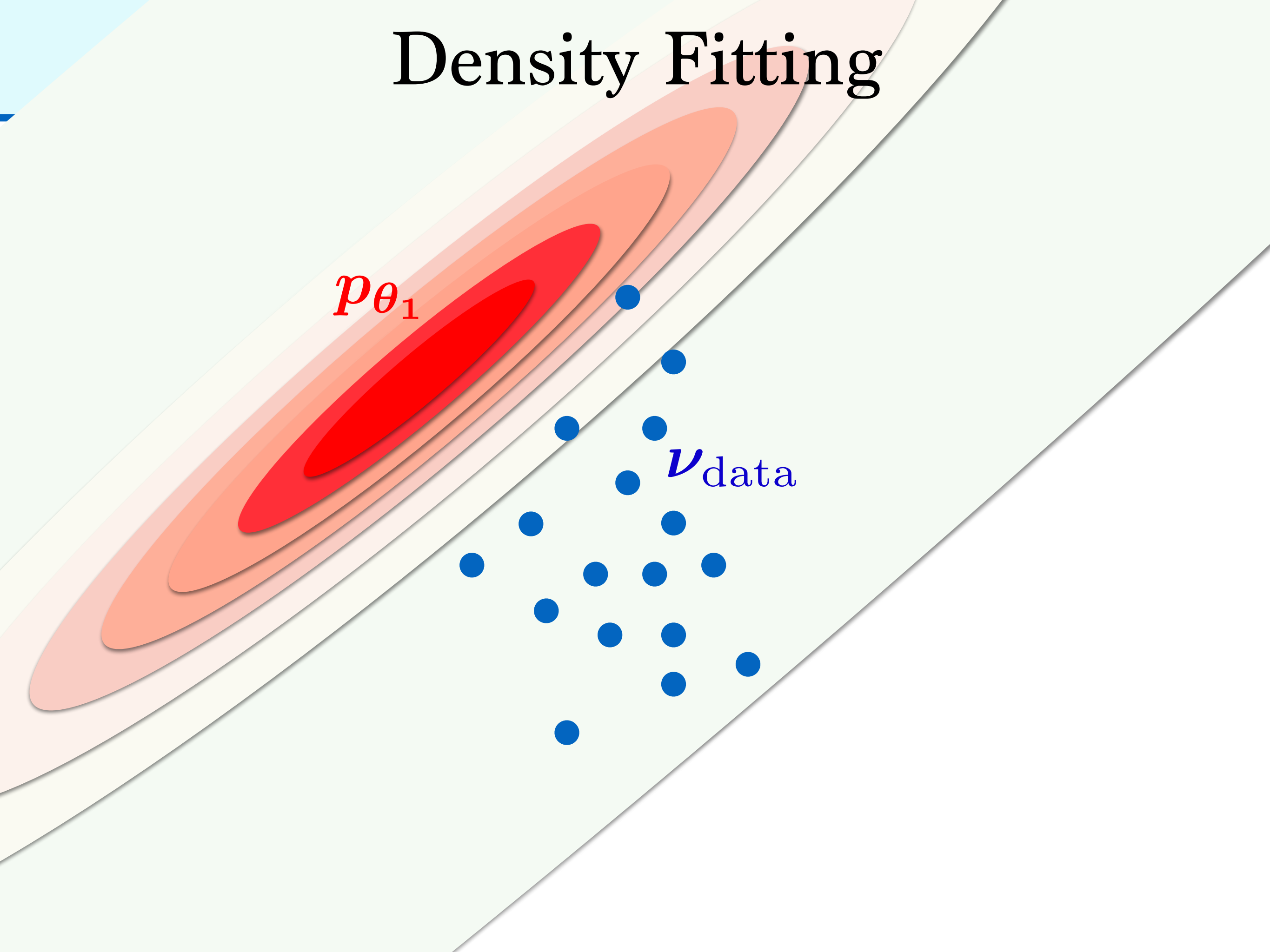$$\nu_{\text{data}} = \frac{1}{N} \sum_{i=1}^{N} \delta_{\boldsymbol{x_i}}$$

$p_{\boldsymbol{\theta}_0}$

$\nu_{\text{data}}$

We fit a parametric family of densities

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

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

Density Fitting

$p_{\theta_1}$

$\nu_{\text{data}}$
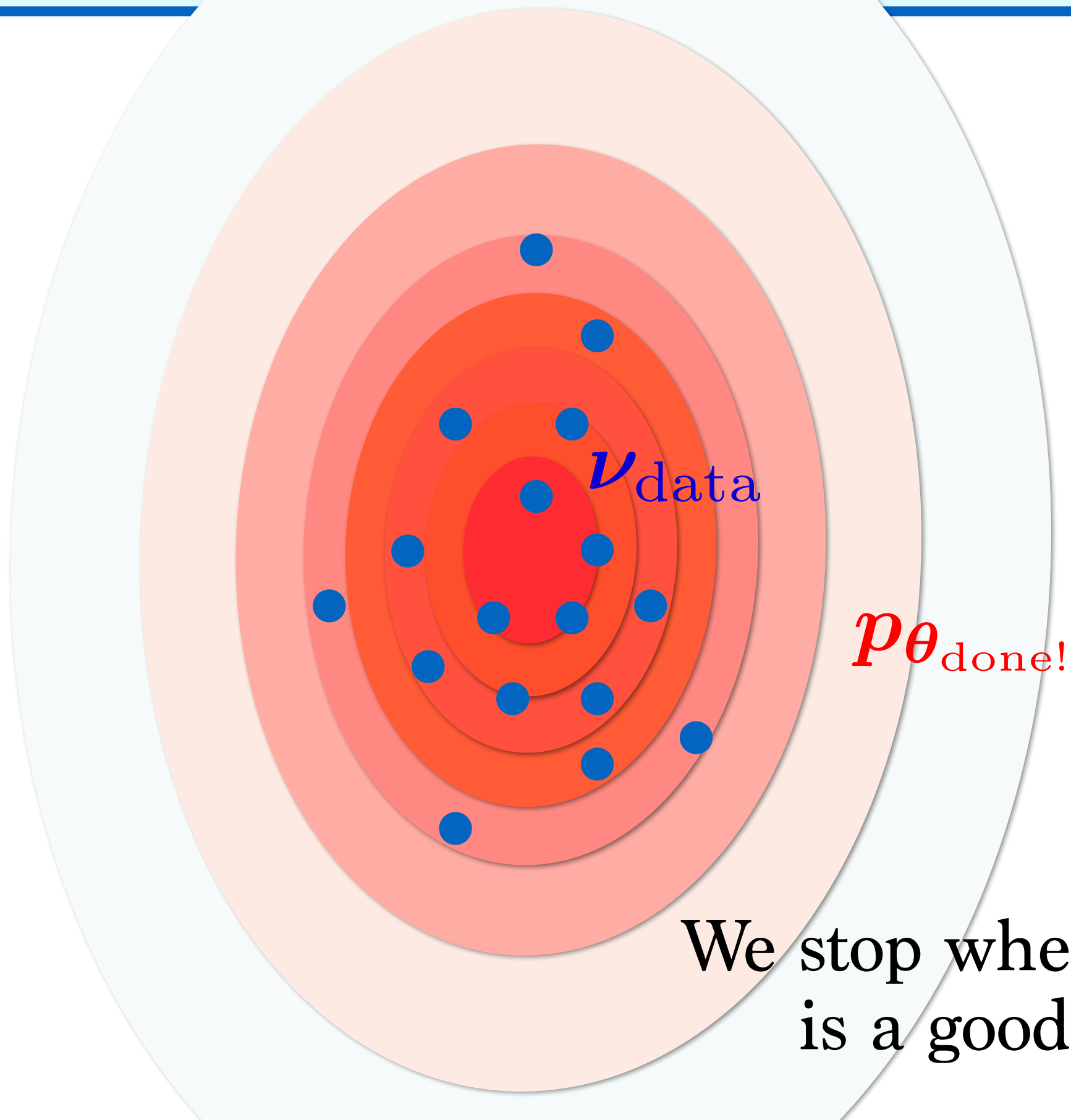
# Density Fitting



$\nu_{\text{data}}$

$p_{\boldsymbol{\theta}\,\text{done!}}$

We stop when there is a good fit.

# 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 a which appears to invalidate any results we ma

$\nu_{\mathrm{data}}$

$p_{\boldsymbol{\theta}\mathrm{done!}}$

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

# 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 a which appears to invalidate any results we ma
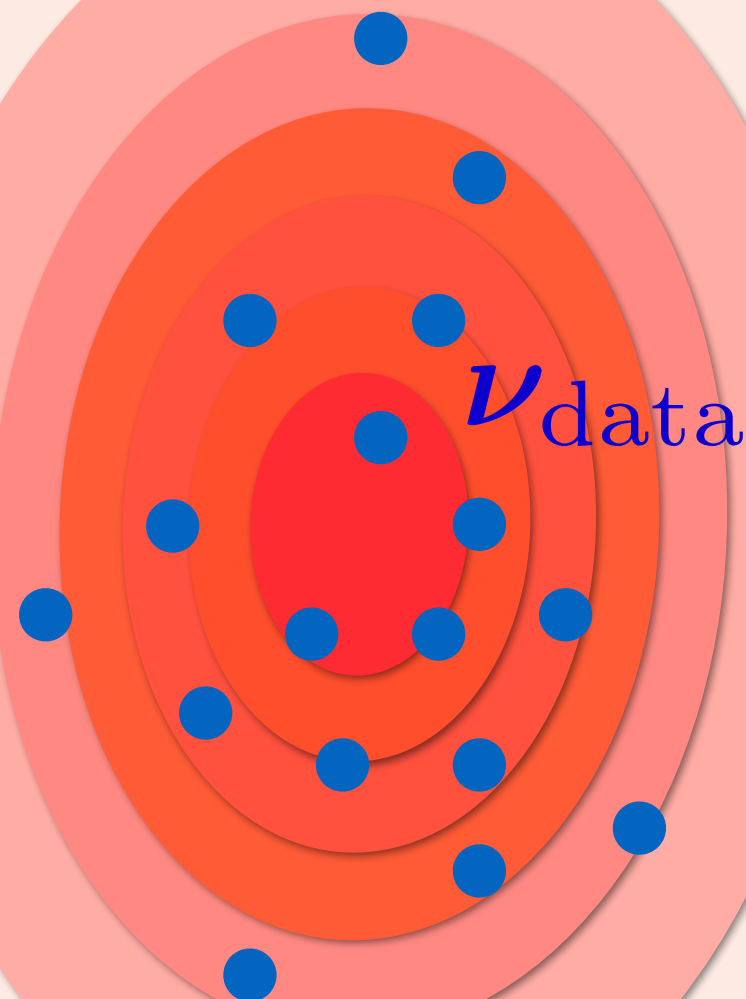
$\nu_{\text{data}}$

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

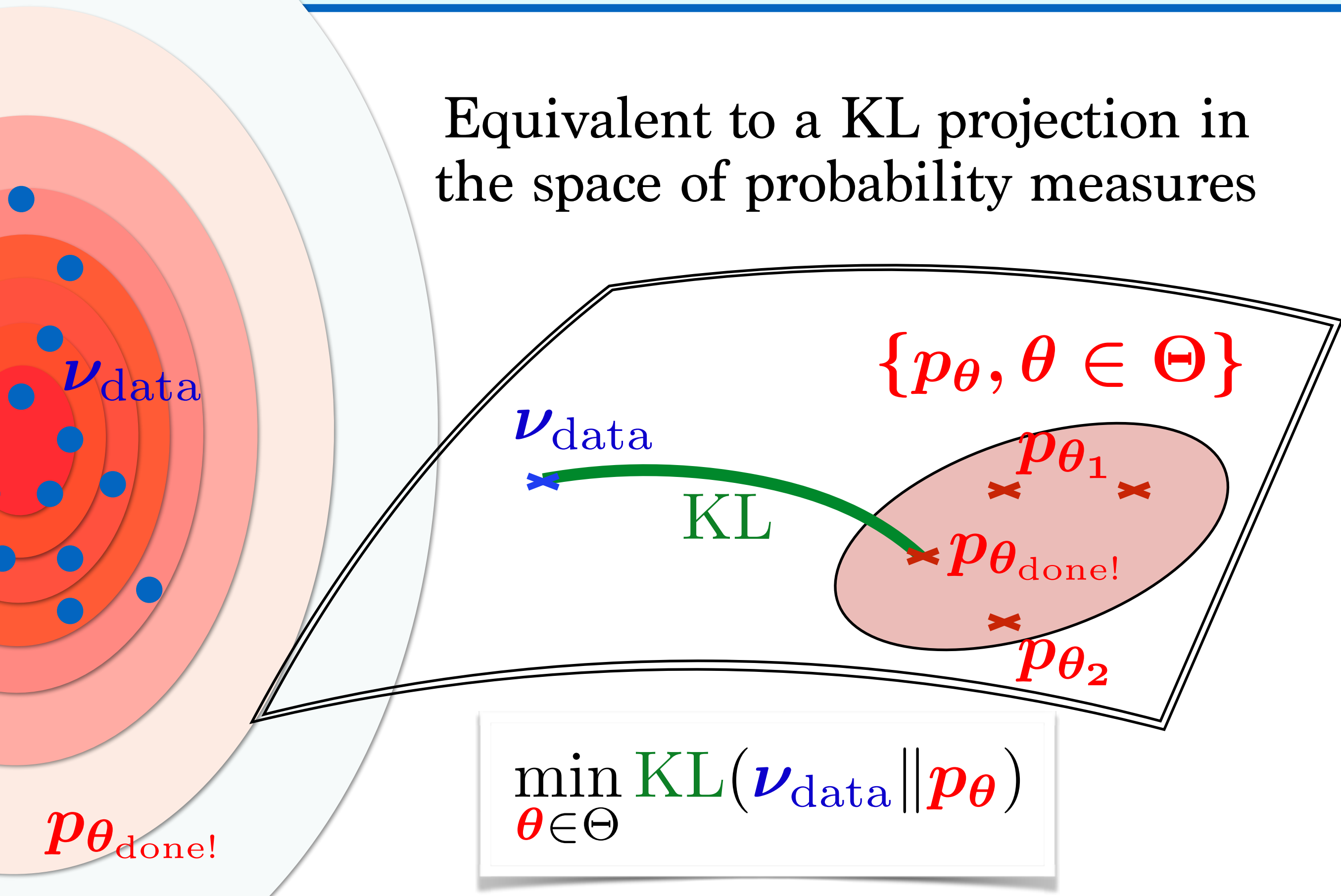$p_{\boldsymbol{\theta}}{}_{\text{done!}}$

$$\log 0 = -\infty$$

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

# Maximum Likelihood Estimation



Equivalent to a KL projection in the space of probability measures

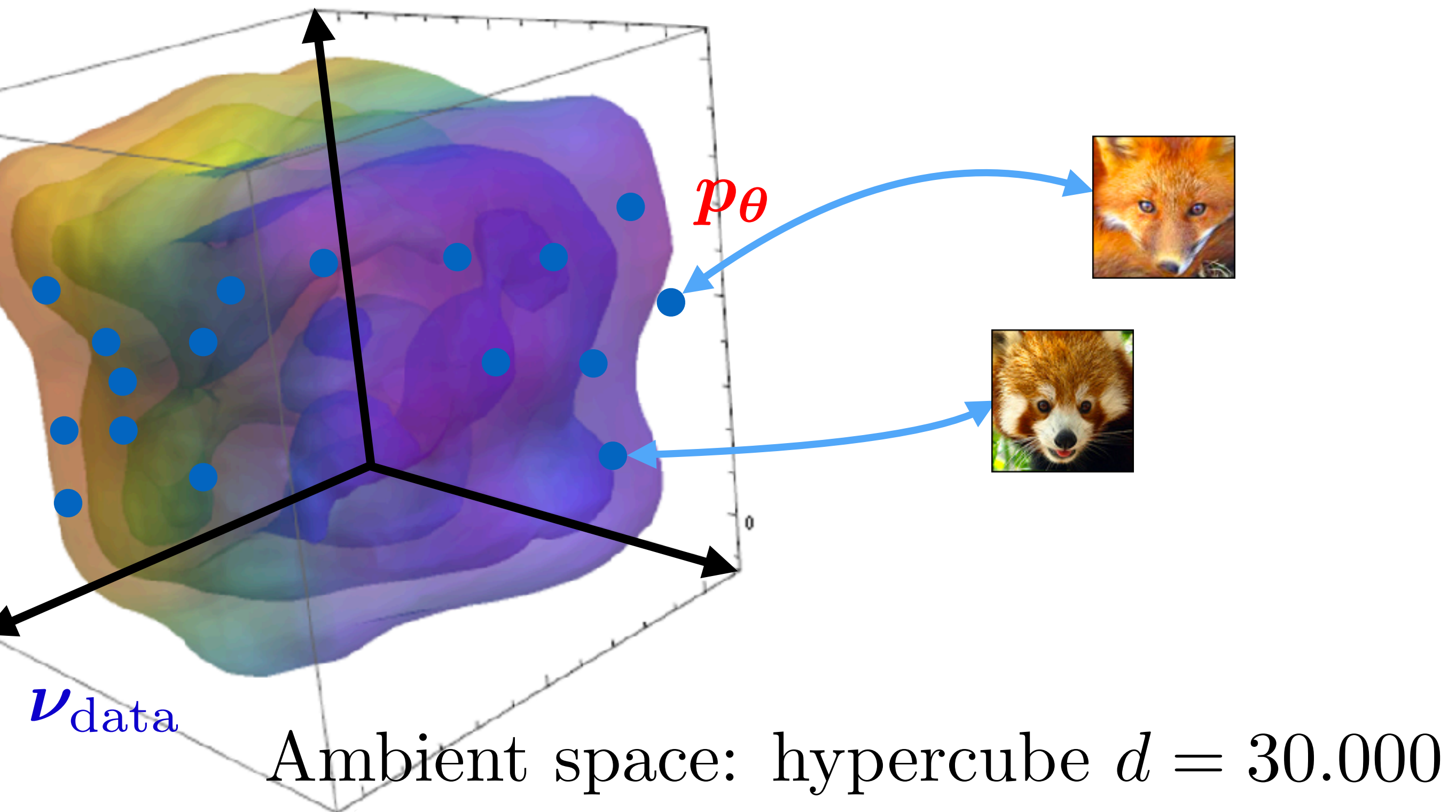$\boldsymbol{\nu}_{\text{data}}$

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

$p_{\boldsymbol{\theta}_1}$

KL

$p_{\boldsymbol{\theta}_{\text{done!}}}$

$p_{\boldsymbol{\theta}_2}$

$p_{\boldsymbol{\theta}_{\text{done!}}}$

$$\min_{\boldsymbol{\theta} \in \Theta} \text{KL}(\boldsymbol{\nu}_{\text{data}} \| p_{\boldsymbol{\theta}})$$

# In higher dimensional spaces...



$p_\theta$

$\nu_{\text{data}}$

Ambient space: hypercube $d = 30.000$

# Generative Models
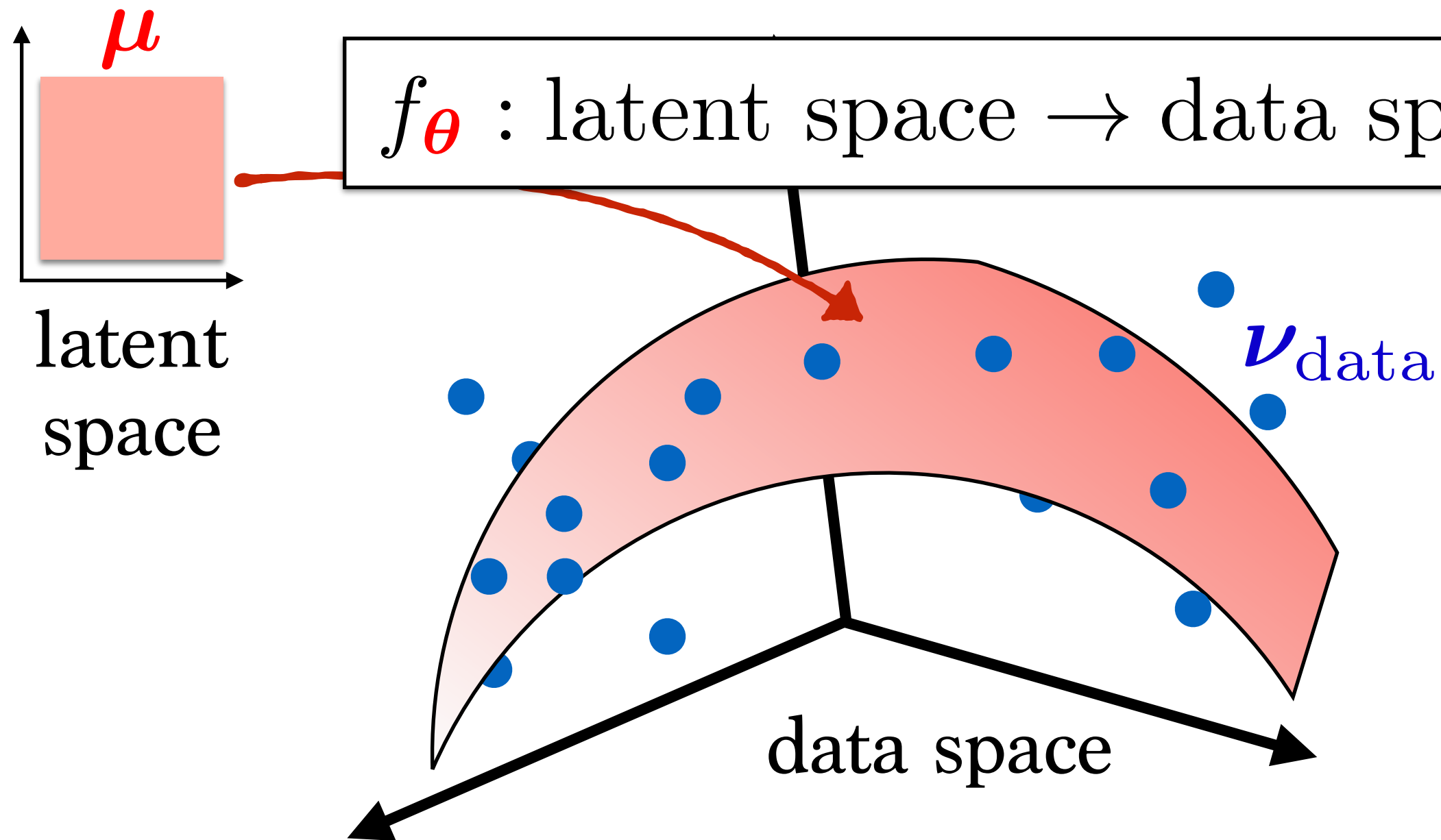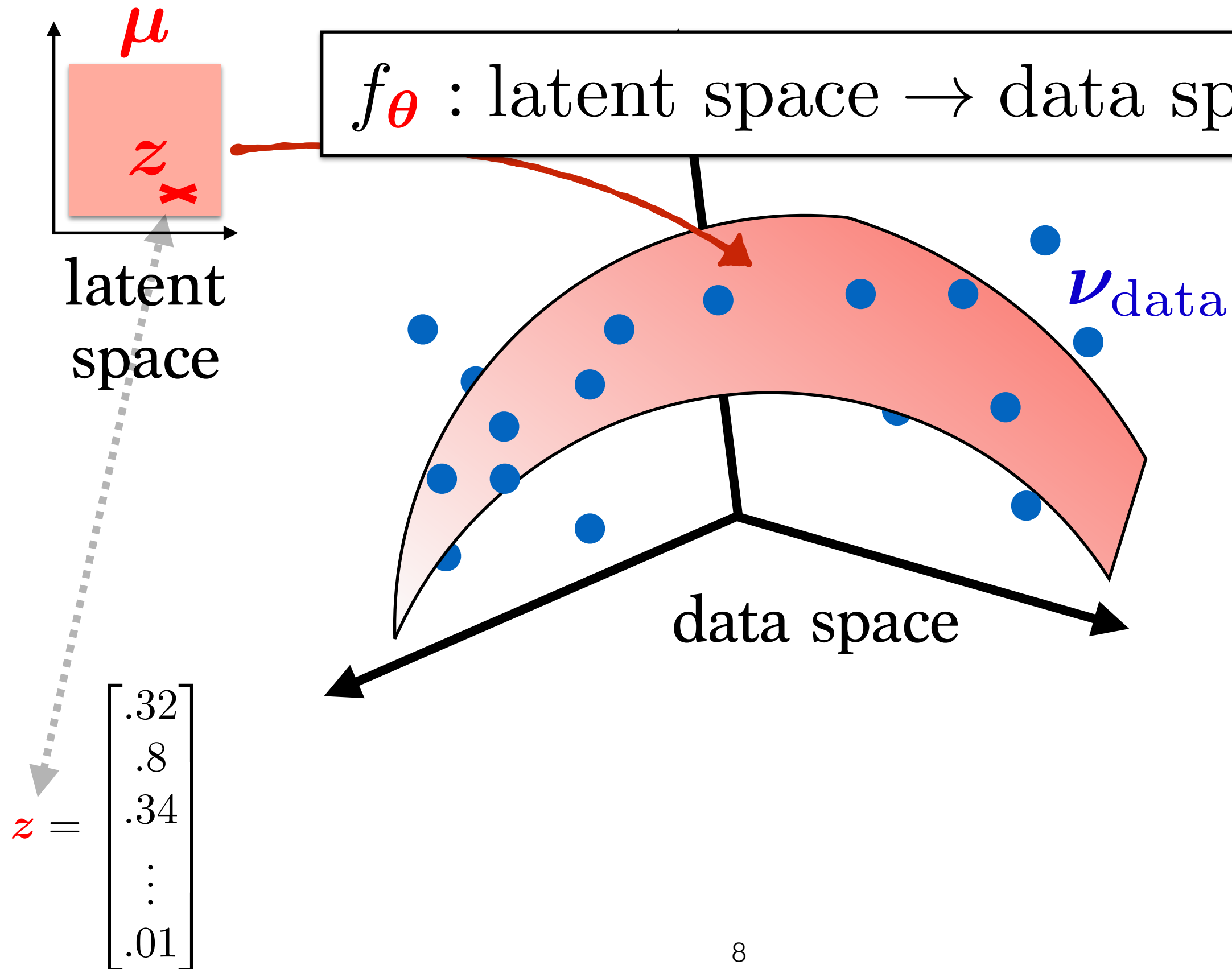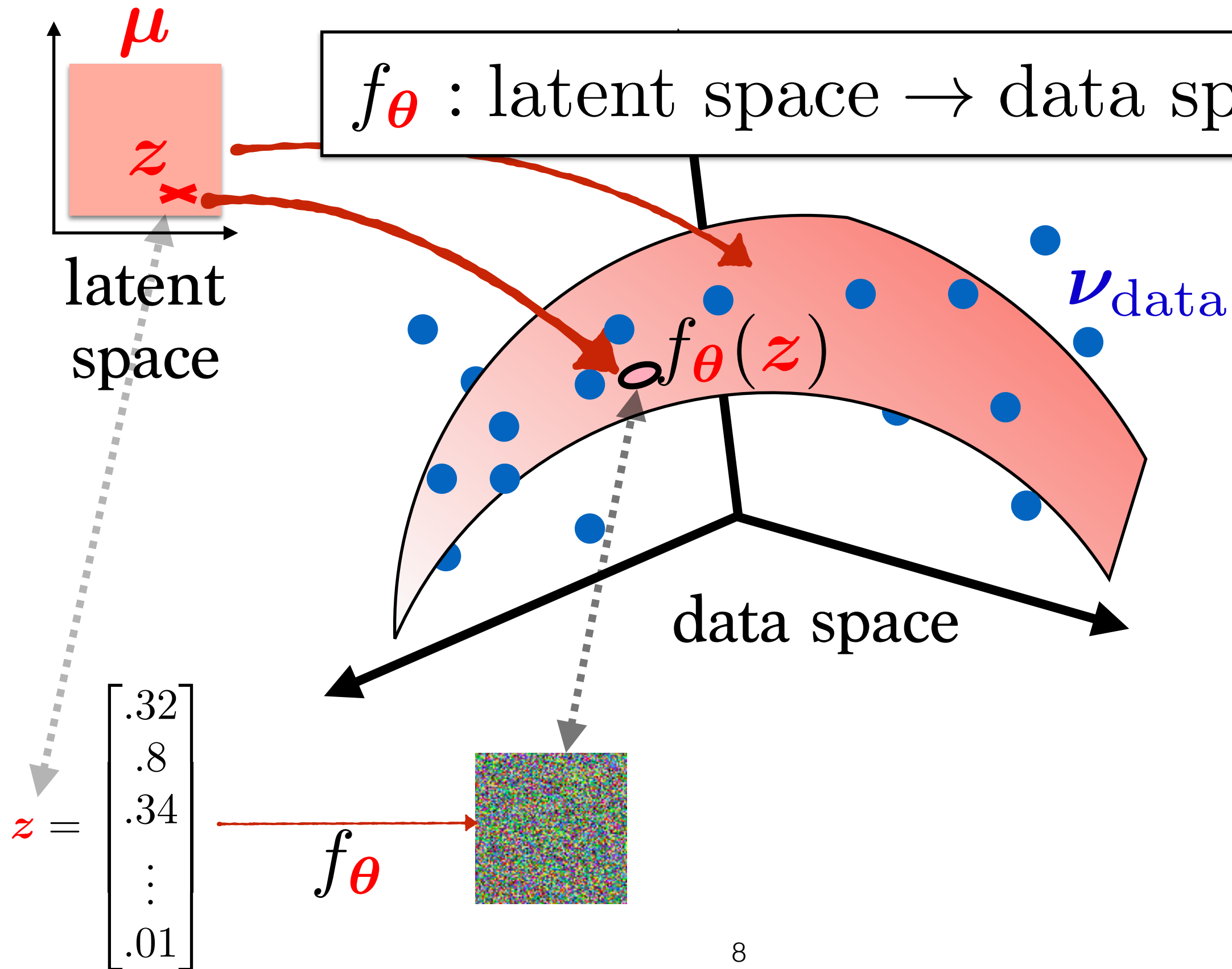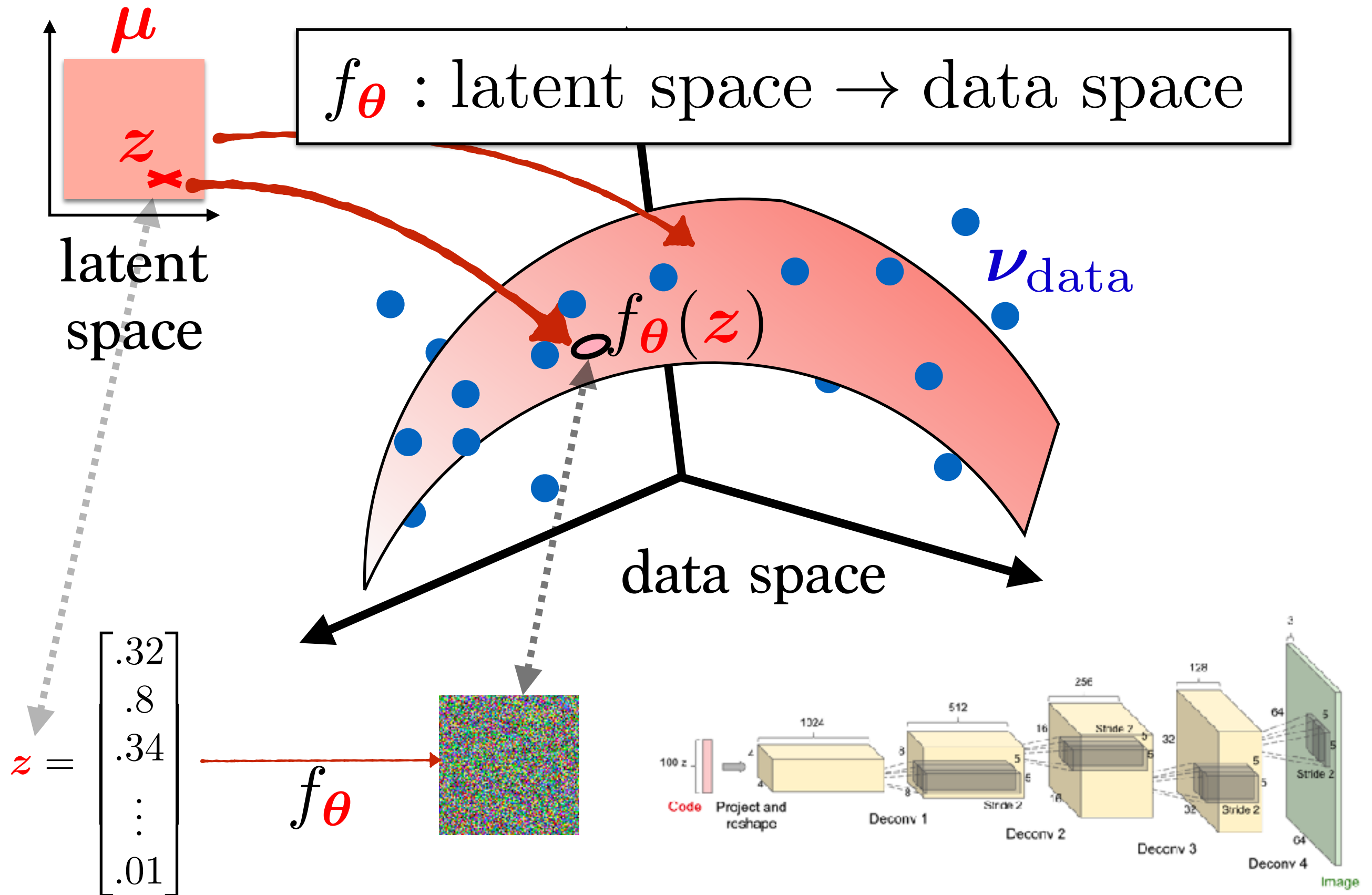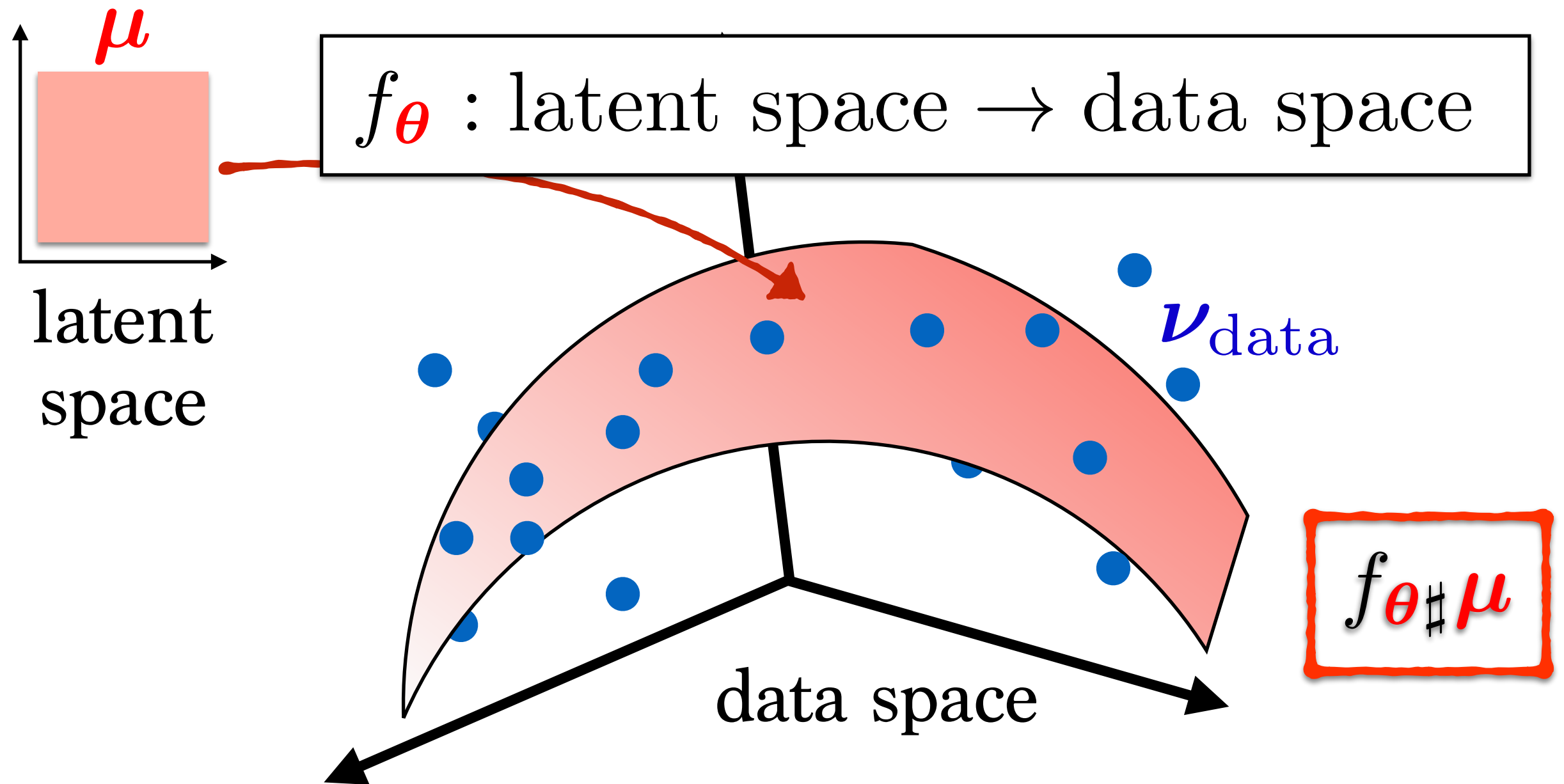
# Generative Models

# Generative Models

# Generative Models

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\mu}$

$\boldsymbol{z}$

$\boldsymbol{x}$

latent space

$\boldsymbol{\nu}_{\mathrm{data}}$

data space

$$\boldsymbol{z} = \begin{bmatrix} .32 \\ .8 \\ .34 \\ \vdots \\ .01 \end{bmatrix}$$

# Generative Models



$\boldsymbol{\mu}$

$\boldsymbol{z}$

latent space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}}(\boldsymbol{z})$

data space

$\boldsymbol{z} = \begin{bmatrix} .32 \\ .8 \\ .34 \\ \vdots \\ .01 \end{bmatrix}$

$f_{\boldsymbol{\theta}}$

8

# Generative Models



$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\mu}$

$\boldsymbol{z}$

latent space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}}(\boldsymbol{z})$

data space

$\boldsymbol{z} = \begin{bmatrix} .32 \\ .8 \\ .34 \\ \vdots \\ .01 \end{bmatrix}$

$f_{\boldsymbol{\theta}}$

# Generative Models



$\boldsymbol{\mu}$

latent space

$f_{\boldsymbol{\theta}} : \text{latent space} \rightarrow \text{data space}$

$\boldsymbol{\nu}_{\text{data}}$

data space

$f_{\boldsymbol{\theta}} \sharp \boldsymbol{\mu}$

# Generative Models



$\boldsymbol{\mu}$

latent space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

data space

Goal: find $\boldsymbol{\theta}$ such that $f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$ fits $\boldsymbol{\nu}_{\text{data}}$

# Generative Models



$\boldsymbol{\mu}$

latent space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\nu}_{\mathrm{data}}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

data space

Goal: find $\boldsymbol{\theta}$ such that $f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$ fits $\boldsymbol{\nu}_{\mathrm{data}}$

9

# Generative Models



$\boldsymbol{\mu}$

latent space

$f_{\boldsymbol{\theta}}$ : latent space $\rightarrow$ data space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

data space

Difference between fitting $f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$ *vs.* a density $p_{\boldsymbol{\theta}}$?

# Generative Models



$$\boxed{f_{\boldsymbol{\theta}} : \text{latent space} \to \text{data space}}$$

$\boldsymbol{\mu}$

latent
space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

data space

MLE

$$\max_{\boldsymbol{\theta} \in \Theta} \frac{1}{N} \sum_{i=1}^{N} \log \boldsymbol{p}_{\boldsymbol{\theta}}(\boldsymbol{x_i}) \quad = \quad \min_{\boldsymbol{\theta} \in \Theta} \text{KL}(\boldsymbol{\nu}_{\text{data}} \| \boldsymbol{p}_{\boldsymbol{\theta}})$$

11

# Generative Models

$\boldsymbol{\mu}$

$f_{\boldsymbol{\theta}}$ : latent space $\to$ data space

latent space

$\boldsymbol{\nu}_{\text{data}}$

$f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}$

data space

MLE

$$\max_{\boldsymbol{\theta}\in\Theta}\frac{1}{N}\sum_{i=1}^{N}\log f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}(\boldsymbol{x_i}) \qquad \min_{\boldsymbol{\theta}\in\Theta}\text{KL}(\boldsymbol{\nu}_{\text{data}}\|f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

# Generative Models

# Workarounds?



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

$$\min_{\boldsymbol{\theta} \in \Theta} \max_{\text{classifiers } \boldsymbol{g}} \text{Accuracy}_{\boldsymbol{g}} \left( (\boldsymbol{f}_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu}, +1), (\boldsymbol{\nu}_{\text{data}}, -1) \right)$$

- Use a **metric** $\triangle$ for probability measures, that can handle measures with non-overlapping supports:

$$\min_{\boldsymbol{\theta} \in \Theta} \triangle(\boldsymbol{\nu}_{\text{data}}, \boldsymbol{p}_{\boldsymbol{\theta}}), \quad \textbf{not} \min_{\boldsymbol{\theta} \in \Theta} \text{KL}(\boldsymbol{\nu}_{\text{data}} \| \boldsymbol{p}_{\boldsymbol{\theta}})$$

14

# Minimum △ Estimation

## MINIMUM 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

Springer Series in Statistics

Luc Devroye
Gábor Lugosi

Combinatorial
Methods in
Density
Estimation

$l_1$

## 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*

Statistics & Probability Letters 76 (2006) 1298–1302

STATISTICS &
PROBABILITY
LETTERS

www.elsevier.com/locate/stapro

ELSEVIER

## On minimum Kantorovich distance estimators

Federico Bassetti[a], Antonella Bodini[b], Eugenio Regazzini[a,*]

## Generative Moment Matching Networks

Yujia Li[1]                                    YUJIALI@CS.TORONTO.EDU
Kevin Swersky[1]                              KSWERSKY@CS.TORONTO.EDU
Richard Zemel[1,2]                            ZEMEL@CS.TORONTO.EDU

[1]Department of Computer Science, University of Toronto, Toronto, ON, CANADA
[2]Canadian 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 Cheng[2]   Yiming Yang[1]   Barnabás Póczos[1]
[1] Carnegie Mellon University,    [2]IBM Research
{chunlial,wchang2,yiming,bapoczos}@cs.cmu.edu   chengyu@us.ibm.com

## Training generative neural networks via Maximum Mean Discrepancy optimization

Gintare Karolina Dziugaite          Daniel M. Roy           Zoubin Ghahramani
University of Cambridge          University of Toronto       University of Cambridge

**Generative Moment Matching Networks**

Yujia Li[1]          YUJIALI@CS.TORONTO.EDU
Kevin Swersky[1]      KSWERSKY@CS.TORONTO.EDU
Richard Zemel[1,2]    ZEMEL@CS.TORONTO.EDU

[1]Department of Computer Science, University of Toronto, Toronto, ON, CANADA
[2]Canadian 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 Cheng[2]   Yiming Yang[1]   Barnabás Póczos[1]
[1] Carnegie Mellon University,    [2]IBM 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

**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

Wasserstein GAN

Martin Arjovsky[1], Soumith Chintala[2], and Léon Bottou[1,2]

[1]Courant Institute of Mathematical Sciences
[2]Facebook AI Research

# △ Generative Model Estimation

### Generative Moment Matching Networks

Yujia Li[1]                          YUJIALI@CS.TORONTO.EDU
Kevin Swersky[1]                     KSWERSKY@CS.TORONTO.EDU
Richard Zemel[1,2]                   ZEMEL@CS.TORONTO.EDU
[1]Department of Computer Science, University of Toronto, Toronto, ON, CANADA
[2]Canadian 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 Cheng[2]   Yiming Yang[1]   Barnabás Póczos[1]
[1] Carnegie Mellon University,   [2]IBM 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

## Training generative neural networks via Maximum Mean Discrepancy optimization

Gintare Karolina Dziugaite        Daniel M. Roy         Zoubin Ghahramani
University of Cambridge          University of Toronto      University of Cambridge

## Wasserstein Training of Restricted Boltzmann Machines

Grégoire Montavon                        Klaus-Robert Müller[*]
Technische Universität Berlin            Technische Universität Berlin
gregoire.montavon@tu-berlin.de    klaus-robert.mueller@tu-berlin.de

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

Wasserstein GAN

Martin Arjovsky[1], Soumith Chintala[2], and Léon Bottou[1,2]

[1]Courant Institute of Mathematical Sciences
[2]Facebook AI Research

Learning Generative Models with Sinkhorn Divergences

Aude Genevay              Gabriel Peyré              Marco Cuturi
CEREMADE,                 CNRS and DMA,              ENSAE CREST
Université Paris-Dauphine  École Normale Supérieure   Université Paris-Saclay

## Improving GANs Using Optimal Transport

Tim Salimans[*]            Han Zhang[‡]              Alec Radford
OpenAI                    Rutgers University          OpenAI
tim@openai.com          han.zhang@cs.rutgers.edu   alec@openai.com

Dimitris Metaxas
Rutgers University
dnm@cs.rutgers.edu

16

# Minimum Kantorovich Estimation

- Use optimal transport theory, namely *Wasserstein distances* to define discrepancy $\color{green}\triangle$.

$$\min_{\boldsymbol{\theta} \in \Theta} W(\boldsymbol{\nu}_{\mathrm{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

- Optimal transport? fertile field in mathematics.



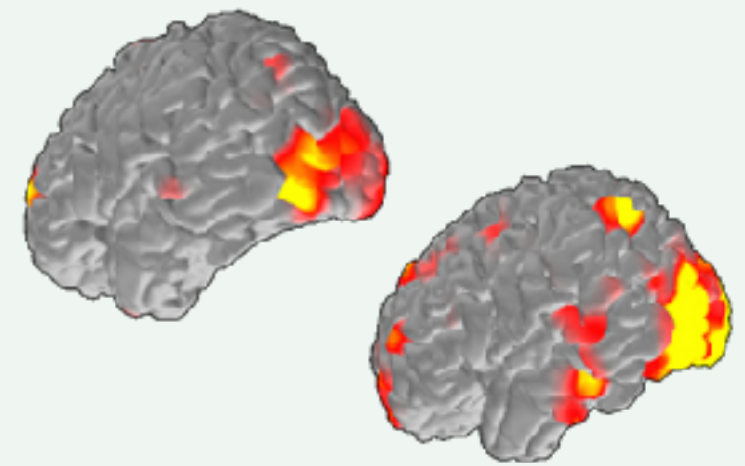| Monge | Kantorovich | Koopmans | Dantzig | Brenier | Otto | McCann | Villani |

Nobel '75

Fields '10

# What is Optimal Transport?

The natural geometry for probability measures



$p_\theta$

$p_{\theta'}$

*Statistical Models*



$h_1$

$h_2$

*Bags of features*



*Brain Activation Maps*



*Generative Models vs. data*

$\mu$

$\nu_{\text{data}}$

latent space



*Color Histograms*

18

# What is Optimal Transport?

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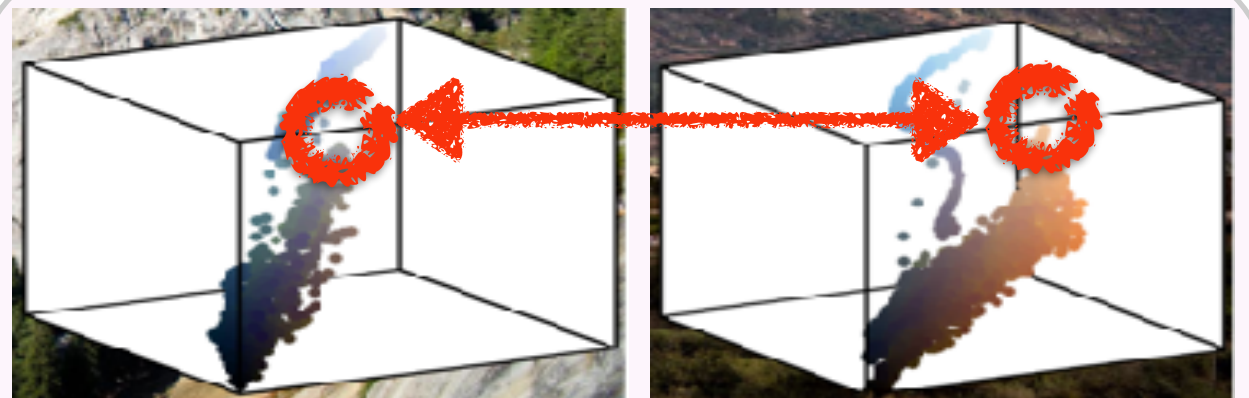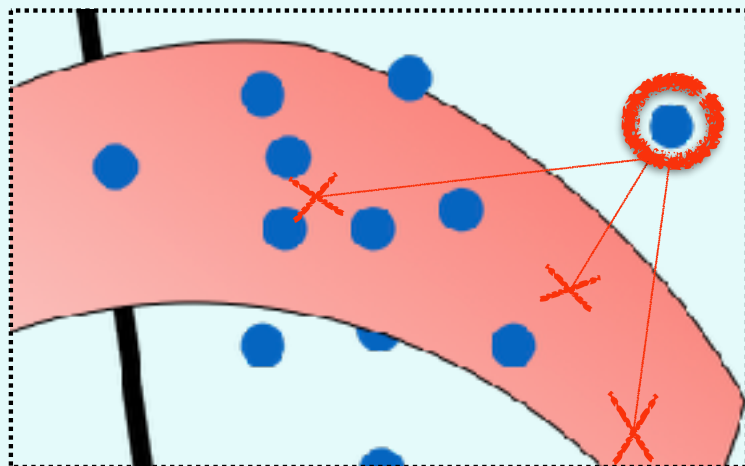


*Statistical Models*

*Bags of Features*

*Brain Activation Maps*

*Generative Models vs. Data*

$\mu$

$\nu_{\text{data}}$

latent space

*Color Histograms*

# What is Optimal Transport?

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


*Statistical Models*


*Bags of Features*


*Brain Activation Maps*


*Generative Models vs. Data*


*Color Histograms*

19

# What is Optimal Transport?

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



*Statistical Models*



*Bags of Features*



*Brain Activation Maps*



*Generative Models vs. Data*



*Color Histograms*

66 MÉMOIRES DE L'ACADÉMIE ROYALE

# MÉMOIRE

### SUR LA

# THÉORIE DES DÉBLAIS

## ET DES REMBLAIS.

### Par M. MONGE.

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.

*In the 21st Century...*

$\mu$

$\nu$

# Origins: Monge Problem

*In the 21st Century...*

*In 1781 however...*

$\mu$

$\nu$

*In 1781 however...*



$\mu$

$\mu(x)$

$x$

$\nu$

*In 1781 however...*

$\boldsymbol{\mu}$

$\boldsymbol{\mu}(x)$

$y = T(x)$

$\times$

$\times$
$x$

$\nu$

*In 1781 however…*

$\mu$

$\mu(x)$

$y = T(x)$

$x$

$\nu$

# Origins: Monge's Problem

*In 1781 however...*

$\mu$

$\mu(x)$

$y = T(x)$

$\mathrm{D}(x, T(x))$

$\nu$

# Origins: Monge's Problem

*In 1781 however...*

$\mu$

$\mu(x)$

$y = T(x)$

$x$

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

$\nu$

*T must map red to blue.*

*T must map red to blue.*

$\mu$

$B$

$\nu$

*T must map red to blue.*

$\mu$

$B$

$\nu$

*T must map red to blue.*

$\mu$

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

$\boldsymbol{B}$

23

$\nu$

# Origins: Monge's Problem

*T must map red to blue.*

$\mu$

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

$B$

$A_1 \; A_2 \; A_3$

$\nu$

*T must map red to blue.*

$\mu$

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

$A_1$ $A_2$ $A_3$

$B$

$\nu$

# Origins: Monge's Problem

*T must map red to blue.*

$\mu$

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

$A_1 \; A_2 \; A_3$

$B$

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

$\nu$

*T must map red to blue.*

$\mu$

$B$

$\nu$

*T must map red to blue.*

$\mu$

$B$

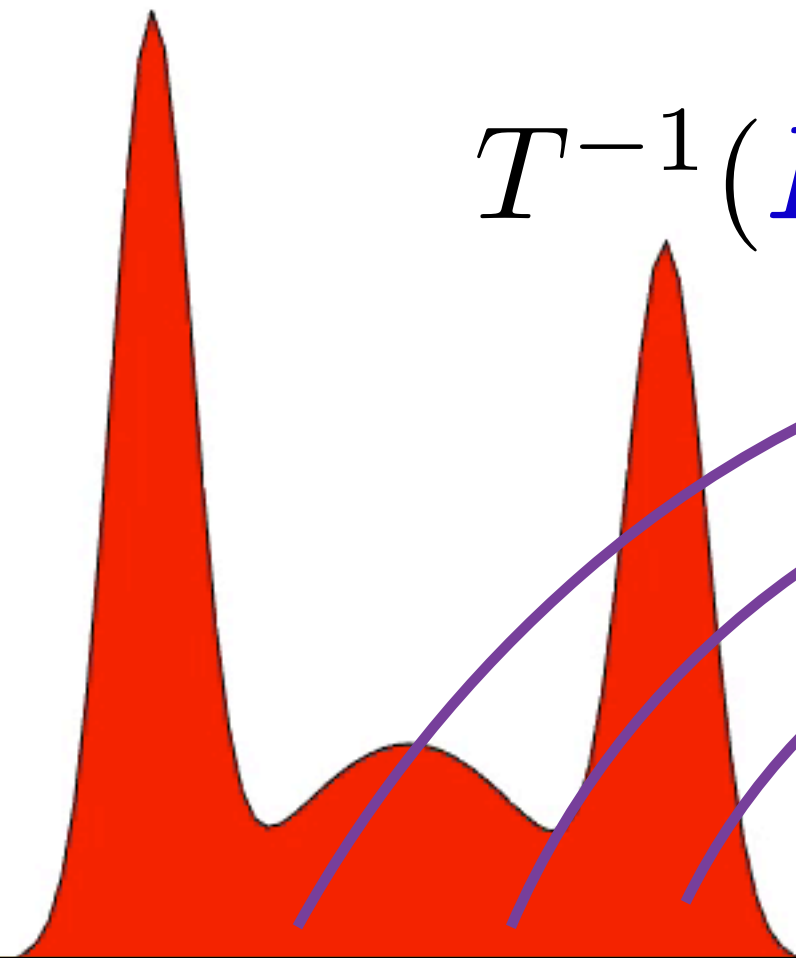$\nu$

*T must map red to blue.*

$\mu$

$B$

$\nu$

# Origins: Monge's Problem

*T must map red to blue.*



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

# Origins: Monge's Problem

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



$$T_\sharp \mu = \nu$$

# Origins: Monge's Problem

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

$\boldsymbol{\mu}$

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

$\boldsymbol{\nu}$

# Monge Problem

$\Omega$ a probability space, $\color{green}c$ $: \Omega \times \Omega \to \mathbb{R}$.
$\color{red}\mu\color{black}, \color{blue}\nu$ two probability measures in $\mathcal{P}(\Omega)$.

[**Monge'81**] problem: find a map $\color{brown}T$ $: \Omega \to \Omega$

$$\inf_{\color{brown}T\color{black}_\sharp \color{red}\mu\color{black}=\color{blue}\nu} \int_\Omega \color{green}c\color{black}(x, \color{brown}T\color{black}(x))\color{red}\mu\color{black}(dx)$$



$x$ $\times$ $\times$ $\color{brown}T\color{black}(x)$

# Monge Problem

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

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

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



$x \times \qquad \times T(x)$

# Monge Problem

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

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

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



$x \times$ $\times T(x)$

25

# Monge Problem

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

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

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

$\delta_x$

# [Kantorovich'42] Relaxation

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

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

# [Kantorovich'42] Relaxation

$$\Pi(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \stackrel{\text{def}}{=} \{\textcolor{darkred}{\boldsymbol{P}} \in \mathcal{P}(\Omega \times \Omega) | \forall \textcolor{red}{\boldsymbol{A}}, \textcolor{blue}{\boldsymbol{B}} \subset \Omega,$$
$$\textcolor{darkred}{\boldsymbol{P}}(\textcolor{red}{\boldsymbol{A}} \times \Omega) = \textcolor{red}{\boldsymbol{\mu}}(\textcolor{red}{\boldsymbol{A}}), \textcolor{darkred}{\boldsymbol{P}}(\Omega \times \textcolor{blue}{\boldsymbol{B}}) = \textcolor{blue}{\boldsymbol{\nu}}(\textcolor{blue}{\boldsymbol{B}})\}$$

# [Kantorovich'42] Relaxation

$$\Pi(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \stackrel{\mathrm{def}}{=} \{\textcolor{red}{\boldsymbol{P}} \in \mathcal{P}(\Omega \times \Omega) | \forall \textcolor{red}{\boldsymbol{A}}, \textcolor{blue}{\boldsymbol{B}} \subset \Omega,$$
$$\textcolor{red}{\boldsymbol{P}}(\textcolor{red}{\boldsymbol{A}} \times \Omega) = \textcolor{red}{\boldsymbol{\mu}}(\textcolor{red}{\boldsymbol{A}}), \textcolor{red}{\boldsymbol{P}}(\Omega \times \textcolor{blue}{\boldsymbol{B}}) = \textcolor{blue}{\boldsymbol{\nu}}(\textcolor{blue}{\boldsymbol{B}})\}$$



27

# Kantorovich Problem

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

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

PRIMAL

# Kantorovich Problem

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

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

PRIMAL

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

DUAL

# (Kantorovich) Wasserstein Distances

Let $p \geq 1$.     Let $c := D$, a metric.

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

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

# (Kantorovich) Wasserstein Distances

Let $p \geq 1$.          Let $c := D$, a metric.

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

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

# Optimal Transport Geometry

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

# Optimal Transport Geometry

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



$\mathcal{P}(\Omega)$

$W(p_\theta, p_{\theta'})$

$\nu$

$\mu$

Wasserstein Distance

# Optimal Transport Geometry

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



$\mathcal{P}(\Omega)$

$\boldsymbol{\nu}$
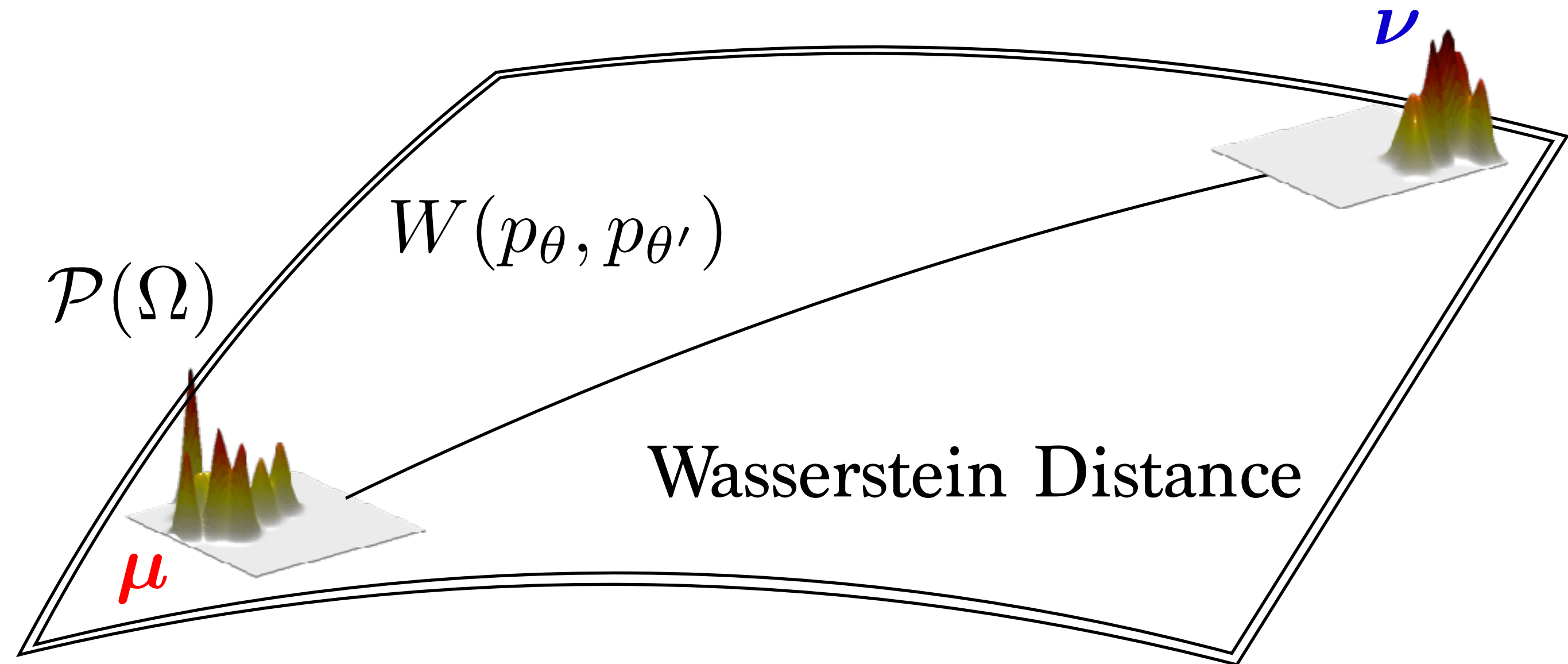
$\boldsymbol{\mu}$

**[McCann'95]** Interpolant

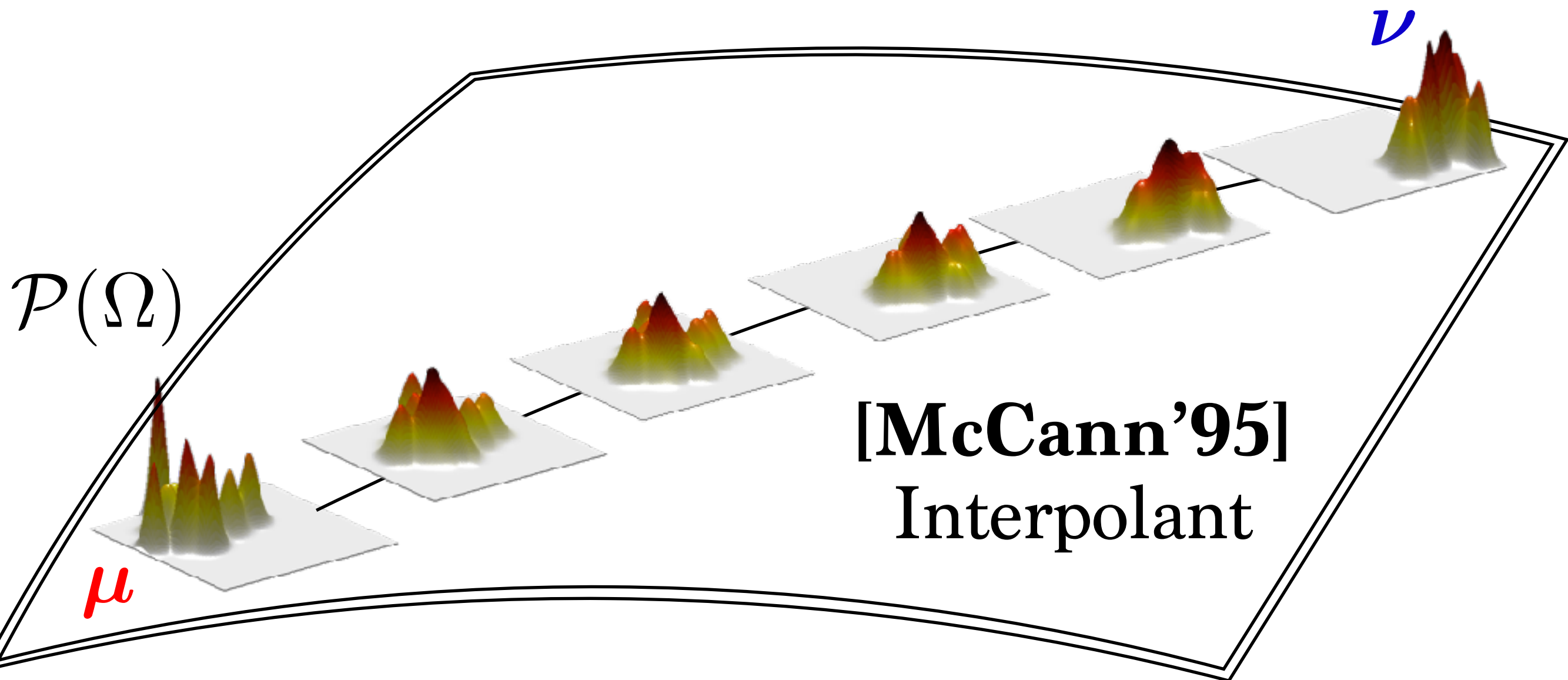# Optimal Transport Geometry

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

# Optimal Transport Geometry

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



$P(\Omega)$

$p_{\theta'}$

$p_\theta$

$p_{\theta''}$

# Optimal Transport Geometry

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



$p_{\theta'}$

Wasserstein Barycenter [**Agueh'11**]
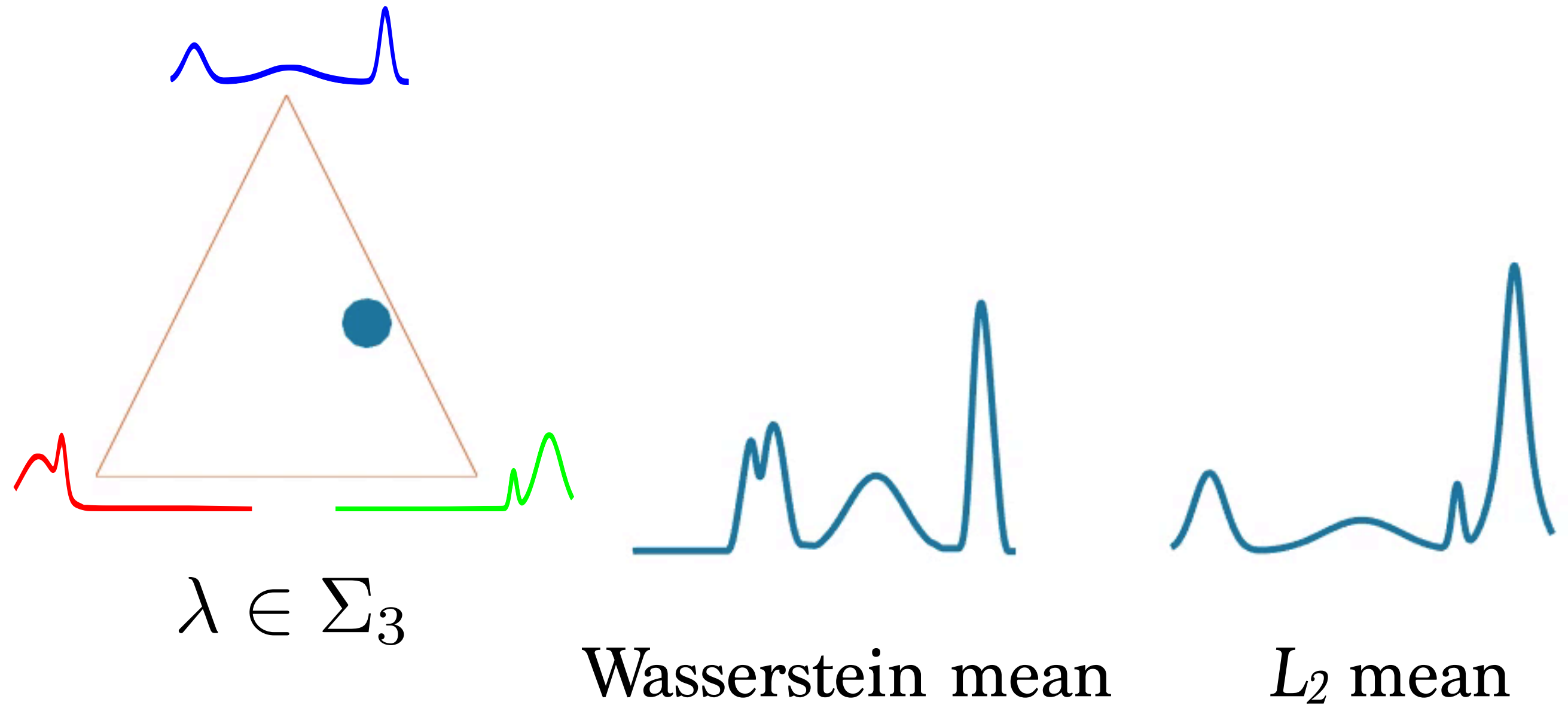
$P(\Omega)$

$p_\theta$

$p_{\theta''}$

# Optimal Transport Geometry



$\lambda \in \Sigma_3$

Wasserstein mean

$L_2$ mean

# Optimal Transport Geometry



$\lambda \in \Sigma_3$

Wasserstein mean

$L_2$ mean

# Computational OT

Up to 2010: OT solvers $\qquad W_p(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) = ?$

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

$$\underset{\textcolor{red}{\boldsymbol{\mu}} \in \mathcal{P}(\Omega)}{\mathrm{argmin}} \ F(W_p(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu_1}}), W_p(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu_2}}), \ldots, \textcolor{red}{\boldsymbol{\mu}}) = ?$$

$$\nabla_{\textcolor{red}{\boldsymbol{\mu}}} W_p(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu_1}}) = ?$$

# How can we compute OT?

Discrete - Discrete

Discrete - Continuous

Continuous - Continuous

34

# 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]**

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

$$(\Omega, D)$$

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

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



$$(\Omega, D)$$

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

# Wasserstein on Discrete Measures

**Consider** $\boldsymbol{\mu} = \sum_{i=1}^{n} a_i \boldsymbol{\delta}_{\boldsymbol{x}_i}$ and $\boldsymbol{\nu} = \sum_{j=1}^{m} b_j \boldsymbol{\delta}_{\boldsymbol{y}_j}$.

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

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

**Def.** Optimal Transport Problem

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

$M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}}$

$U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})$

$M_{\boldsymbol{X}\boldsymbol{Y}}$

*min cost flow solver used in practice.*

$\boldsymbol{O(n^3 \log(n))}$

$U(\boldsymbol{a}, \boldsymbol{b})$

$P^\star$

$M_{\boldsymbol{X}\boldsymbol{Y}}$

$U(\boldsymbol{a}, \boldsymbol{b})$

$P^\star$

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

*Solution $P^\star$ unstable and not always unique.*

$M_{\textcolor{red}{X}\textcolor{blue}{Y}}$

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

$U(\textcolor{red}{a}, \textcolor{blue}{b})$

$\{P^\star\}$

*Solution $P^\star$ unstable and not always unique.*

$M_{XY}$

$U(a, b)$

$\{P^\star\}$

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

*Solution $P^\star$ unstable and not always unique.*

$M_{\textbf{X}\textbf{Y}}$

$U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})$

$P^{\star}$

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

*Solution $P^{\star}$ unstable and not always unique.*

$M_{\mathbf{X}\mathbf{Y}}$

*min cost flow solver used in practice.*

$$O(n^3 \log(n))$$

$U(\mathbf{a}, \mathbf{b})$

*Solution $P^\star$ unstable and not always unique.*

$P^\star$

$W_p^p(\boldsymbol{\mu}, \boldsymbol{\nu})$ not differentiable.

# Solution: Regularization



$M_{\textbf{XY}}$

$U(\textbf{a}, \textbf{b})$

$P^\star$

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

# Entropic Regularization [Wilson'62]

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

$$W_\gamma(\textcolor{red}{\boldsymbol{\mu}}, \textcolor{blue}{\boldsymbol{\nu}}) \overset{\text{def}}{=} \min_{\textcolor{red}{\boldsymbol{P}} \in U(\textcolor{red}{\boldsymbol{a}}, \textcolor{blue}{\boldsymbol{b}})} \langle \textcolor{red}{\boldsymbol{P}}, M_{\textcolor{red}{\boldsymbol{X}}\textcolor{blue}{\boldsymbol{Y}}} \rangle - \gamma E(\textcolor{red}{\boldsymbol{P}})$$
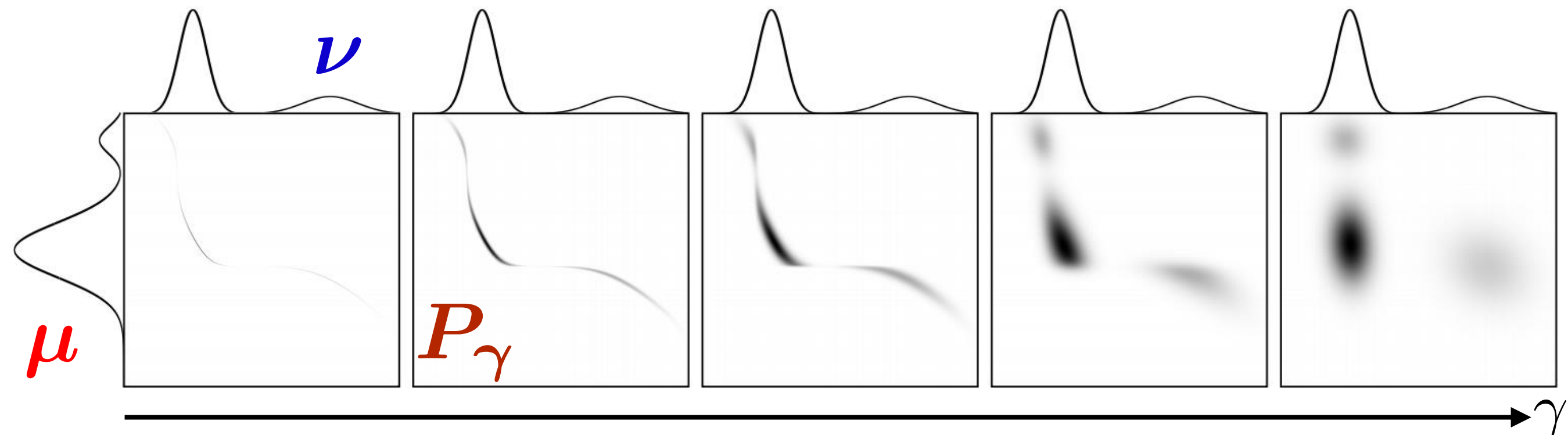
$$E(P) \overset{\text{def}}{=} - \sum_{i,j=1}^{nm} P_{ij}(\log P_{ij} - 1)$$

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

# Entropic Regularization [**Wilson'62**]

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

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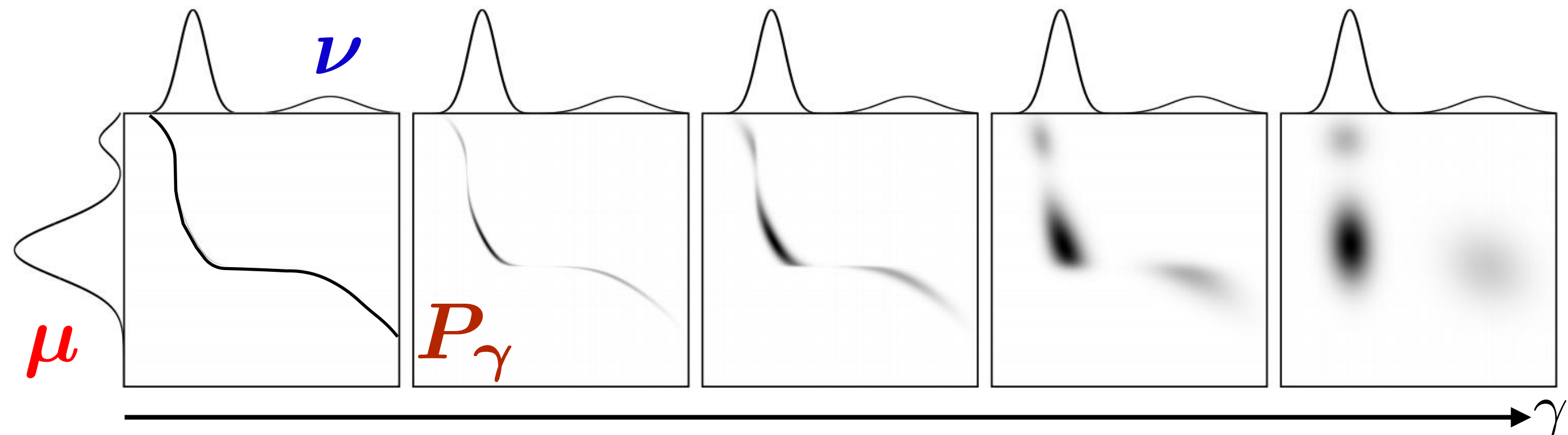


**Note: Unique** optimal solution because of strong concavity of entropy

# Entropic Regularization [**Wilson'62**]

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

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



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

# Fast & Scalable Algorithm

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

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

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

# Fast & Scalable Algorithm

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

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

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

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

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

$$(\partial L/\partial P_{ij} = 0) \Rightarrow P_{ij} = e^{\frac{\alpha_i}{\gamma}}\; e^{-\frac{M_{ij}}{\gamma}}\; e^{\frac{\beta_j}{\gamma}} = u_i\; K_{ij} v_j$$

# Fast & Scalable Algorithm

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

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

$$P_\gamma = \mathbf{diag}(\boldsymbol{u}) K \mathbf{diag}(\boldsymbol{v}), \quad K \overset{\text{def}}{=} e^{-M_{\boldsymbol{XY}}/\gamma}$$

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

# Fast & Scalable Algorithm

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

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

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

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

# Fast & Scalable Algorithm

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

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

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

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

43

# Fast & Scalable Algorithm

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

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

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

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

# Fast & Scalable Algorithm

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

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

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

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

43

# Fast & Scalable Algorithm

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

then $\exists! u \in \mathbb{R}^n_+, v \in \mathbb{R}^m_+$, such that

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

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

# Fast & Scalable Algorithm

Sinkhorn's Algorithm : Repeat

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

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

# Fast & Scalable Algorithm

Sinkhorn's Algorithm : Repeat

$$1. \quad u = a/Kv$$

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

- **[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]**
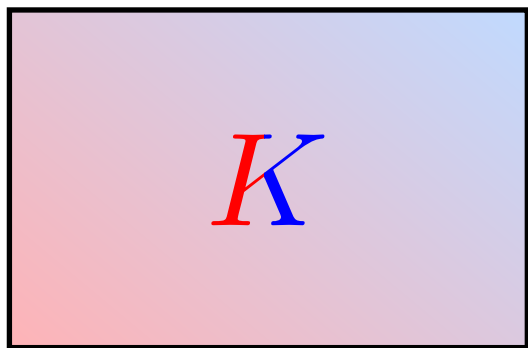
# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

$\textcolor{red}{a}$ $\textcolor{blue}{b}$

$\textcolor{blue}{v_0}$

$\textcolor{red}{K}$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a} / \textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b} / \textcolor{red}{K}^{T}\textcolor{red}{u}$$

$\textcolor{blue}{v_0}$

$\textcolor{red}{a}$ $\textcolor{blue}{b}$

$\textcolor{blue}{K}$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

$\textcolor{red}{a}$  $\textcolor{blue}{b}$

$K$

$\textcolor{blue}{v_0}$

$\textcolor{red}{a}$

$:= \quad \Big/ \quad K$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$





$$\textcolor{red}{u_1} := \textcolor{red}{a} / \textcolor{red}{K} \quad \textcolor{blue}{v_0}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a} / \textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b} / \textcolor{red}{K}^T \textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^{T}\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$



$$\textcolor{red}{u_2} := \textcolor{red}{a} / \textcolor{red}{K}\, \textcolor{blue}{v_1}$$

$$\textcolor{blue}{v} := \textcolor{blue}{b} / \textcolor{red}{K}^T\, \textcolor{red}{u_1}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

$\textcolor{red}{a}$ $\textcolor{blue}{b}$

$K$

$\textcolor{red}{a} := \textcolor{red}{a} / K$ $\textcolor{blue}{v_1}$

$\textcolor{blue}{b} := \textcolor{blue}{b} / K^T$ $\textcolor{red}{u_2}$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations for $(\textcolor{red}{u}, \textcolor{blue}{v})$

$$\textcolor{red}{u} \leftarrow \textcolor{red}{a}/\textcolor{red}{K}\textcolor{blue}{v}, \quad \textcolor{blue}{v} \leftarrow \textcolor{blue}{b}/\textcolor{red}{K}^T\textcolor{red}{u}$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations.

$$\mathbf{diag}(v_L)$$

$$\mathbf{diag}(u_L) \qquad K$$

$$a \quad b$$

$$P_L :=$$

$$K$$

# Fast & Scalable Algorithm

- **[Sinkhorn'64]** fixed-point iterations.

$a$ $b$

$K$

$$\mathbf{diag}(u_L) \quad K \quad \mathbf{diag}(v_L)$$

$$P_L := $$

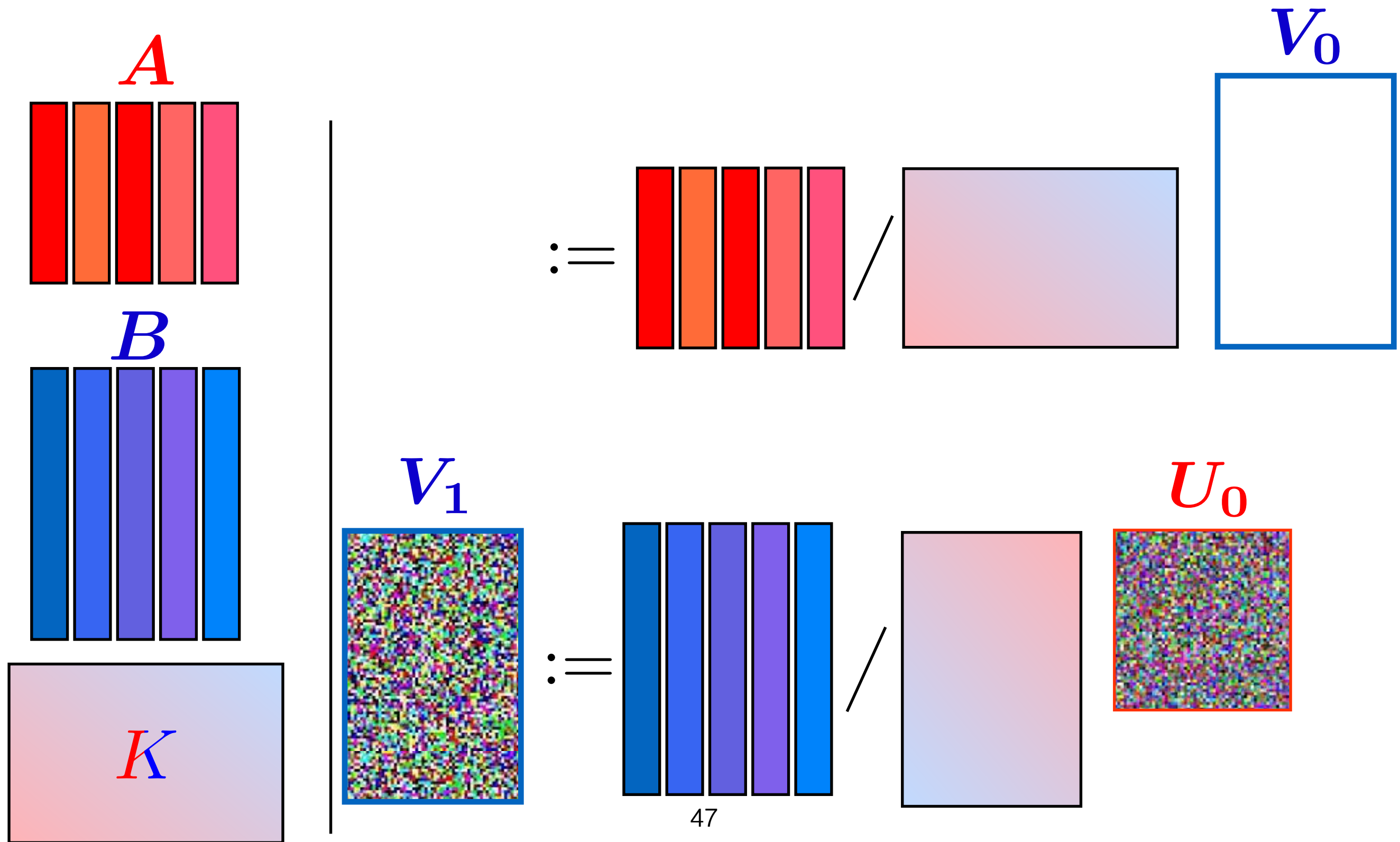$$\langle P_L, M_{XY} \rangle = u_L{}^T (K \odot M_{XY}) v_L$$

# Also embarrassingly parallel

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

# Also embarrassingly parallel

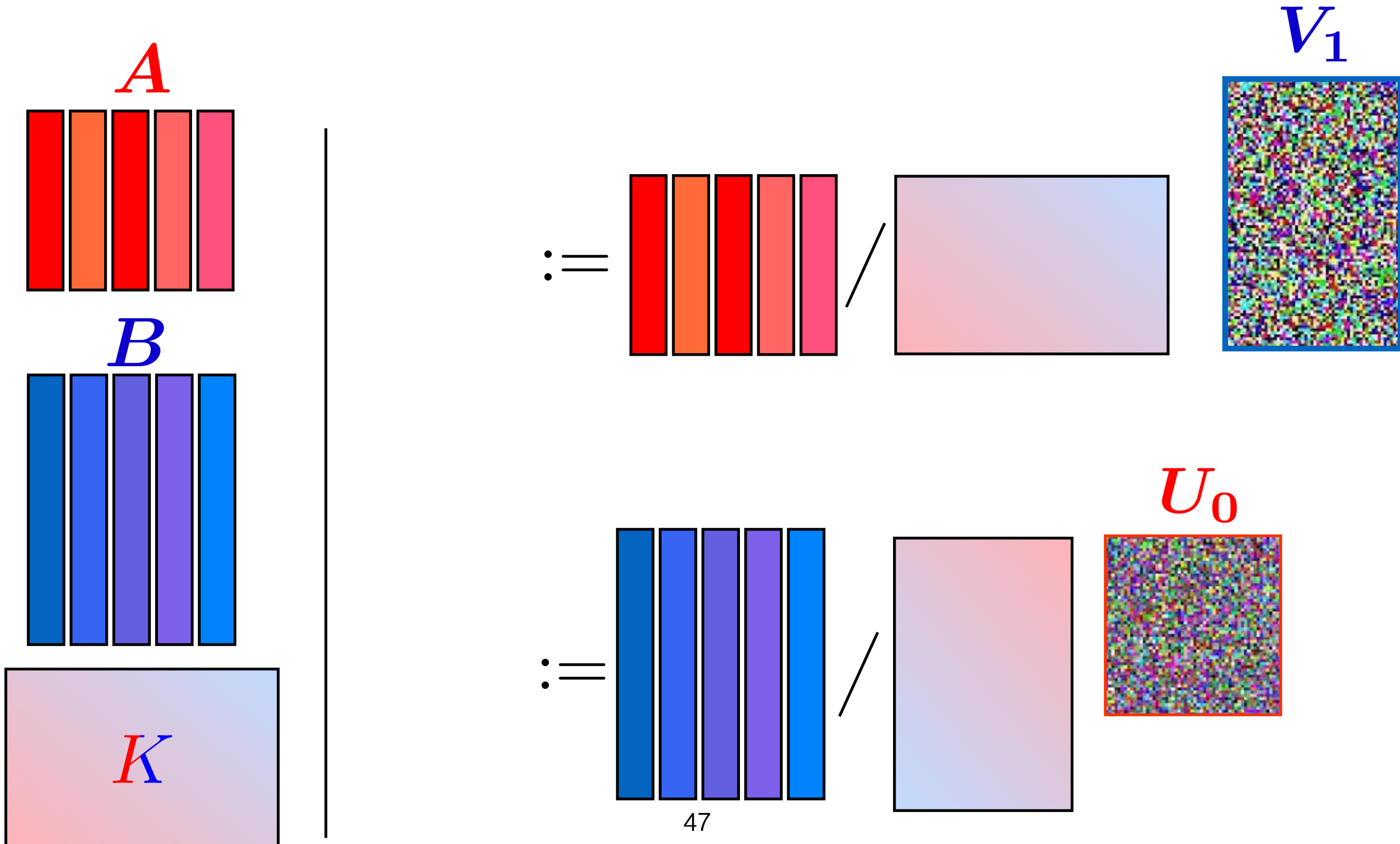- **[Sinkhorn'64]** with *matrix* fixed-point iterations
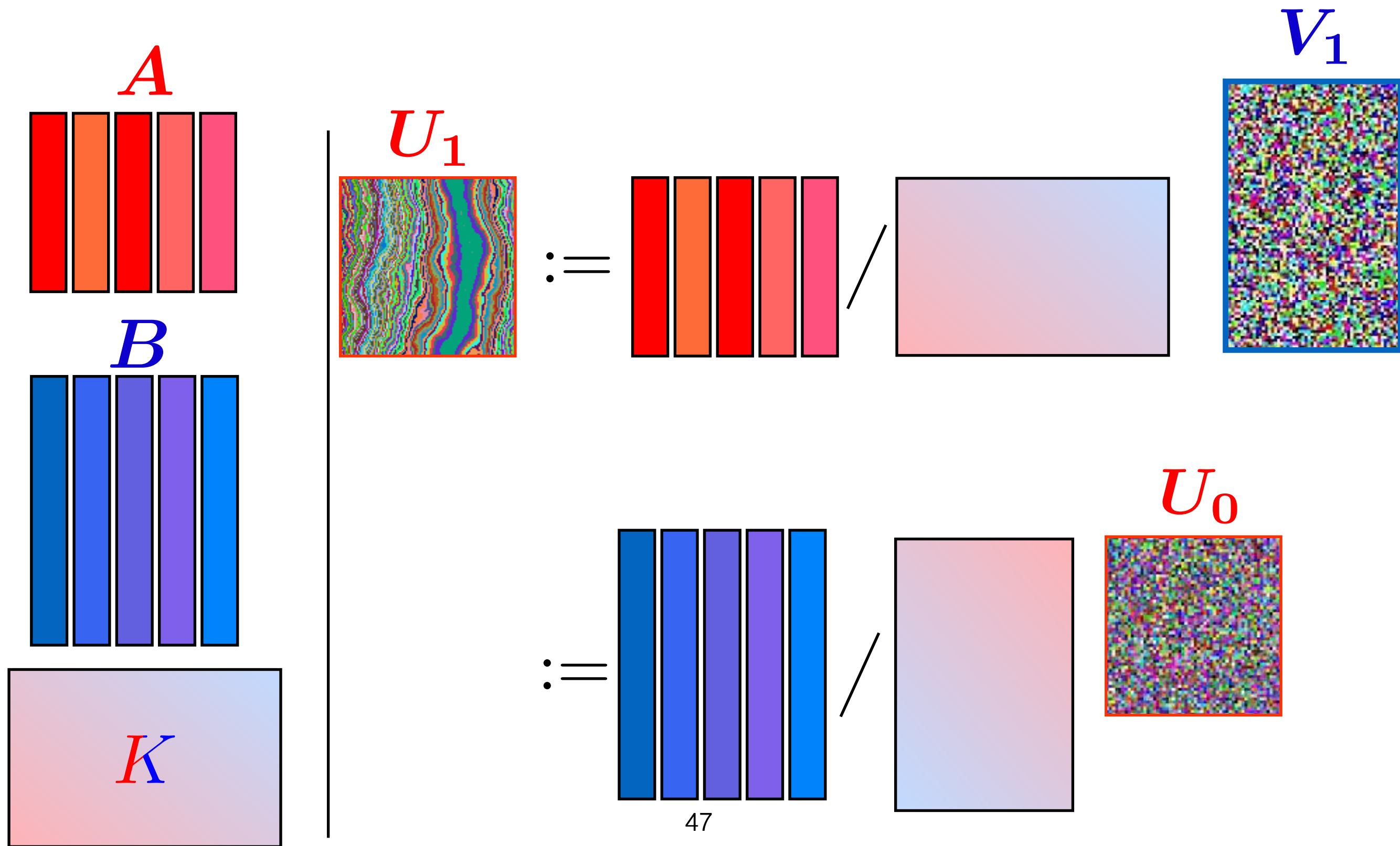
# Also embarrassingly parallel

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

# Also embarrassingly parallel

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

# Also embarrassingly parallel

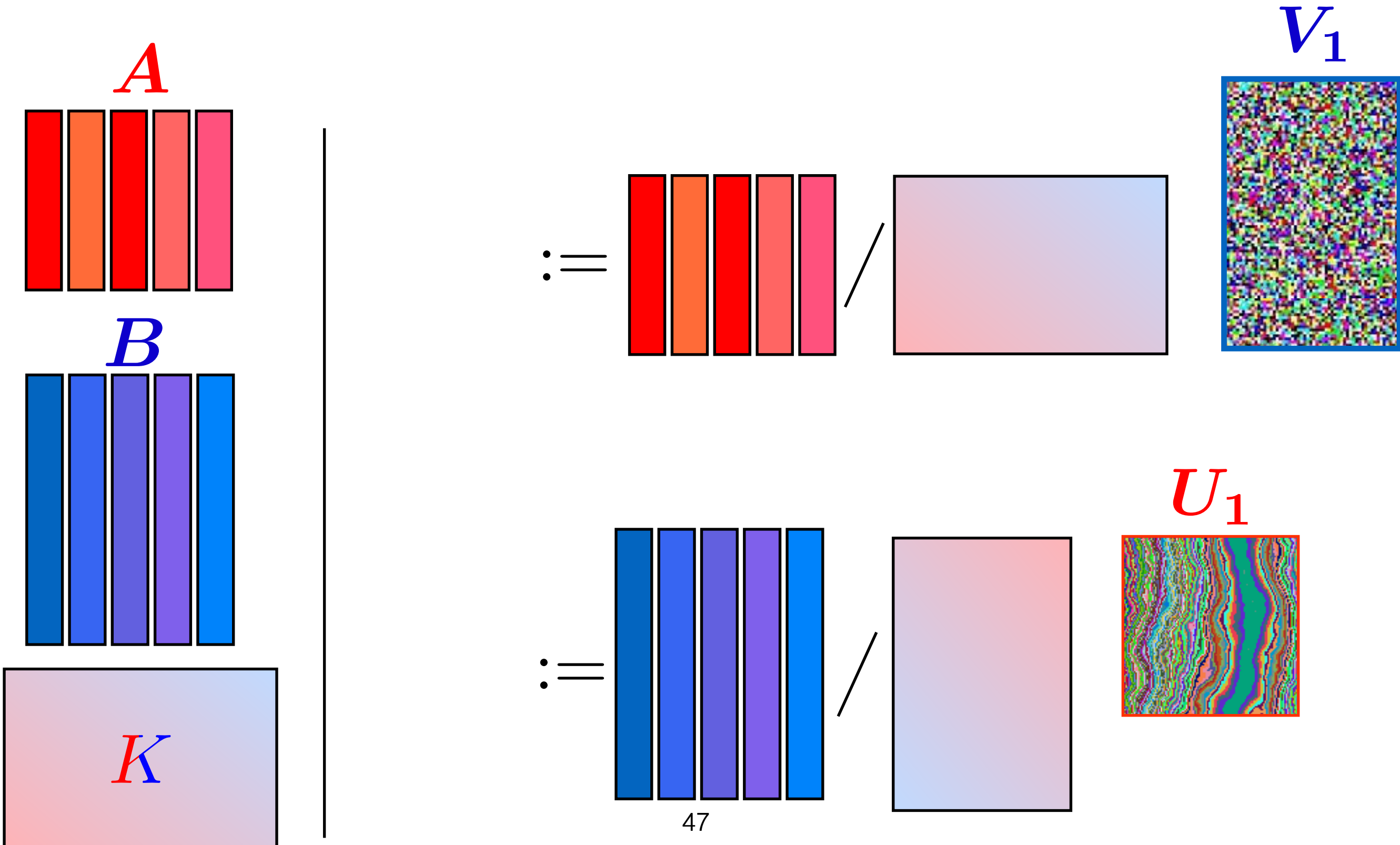- **[Sinkhorn'64]** with *matrix* fixed-point iterations

# Also embarrassingly parallel

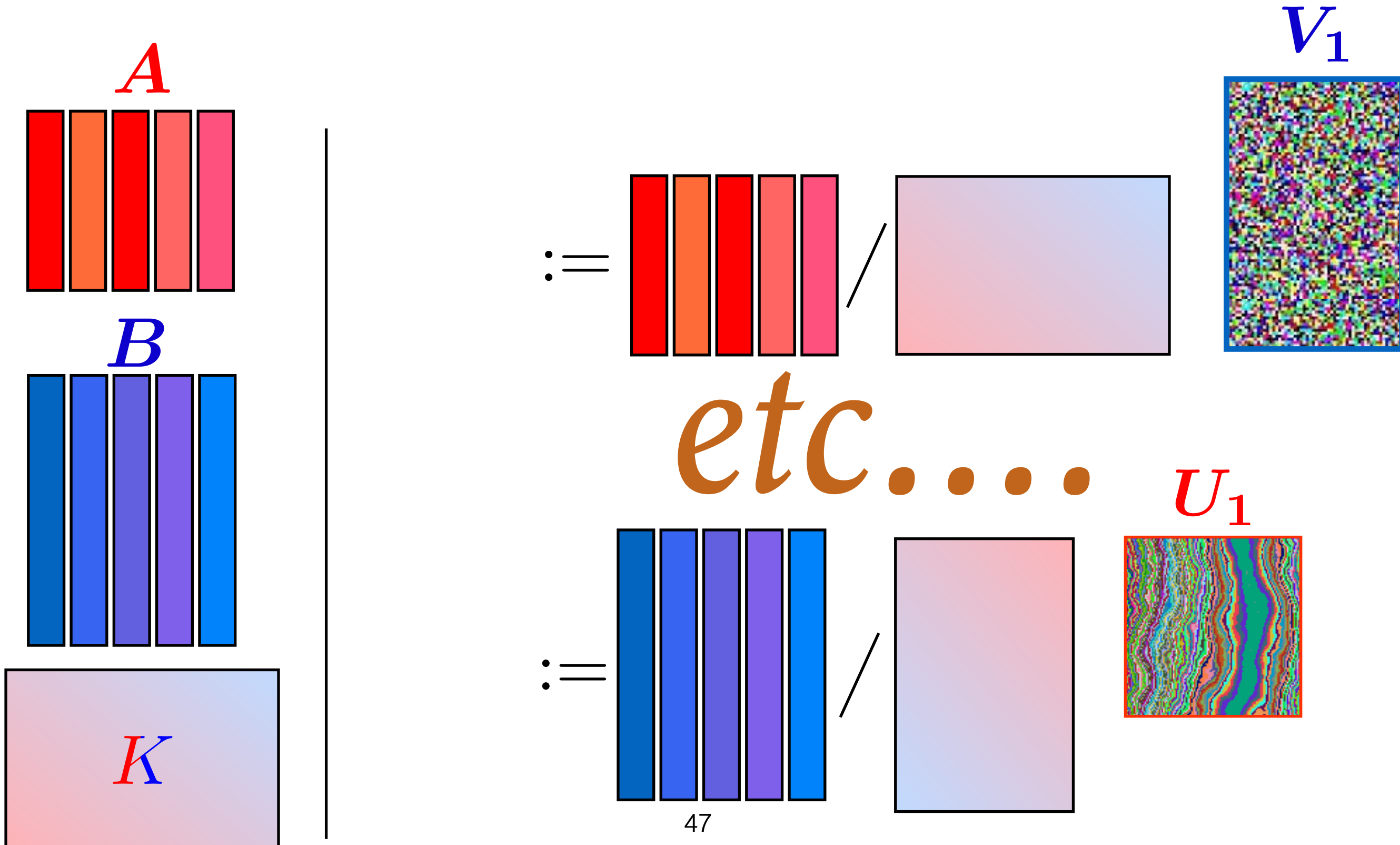- **[Sinkhorn'64]** with *matrix* fixed-point iterations

# 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



47

# Also embarrassingly parallel

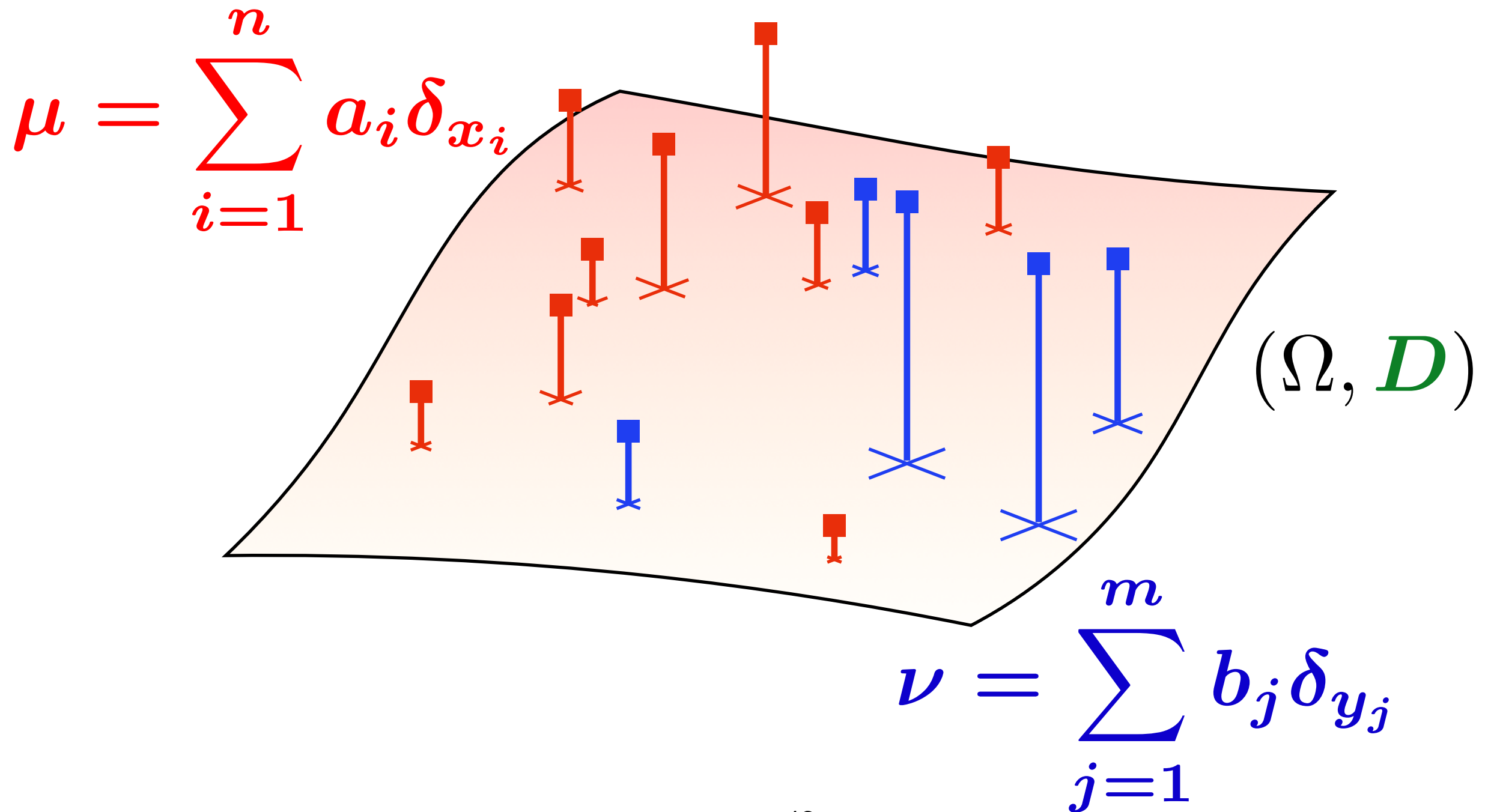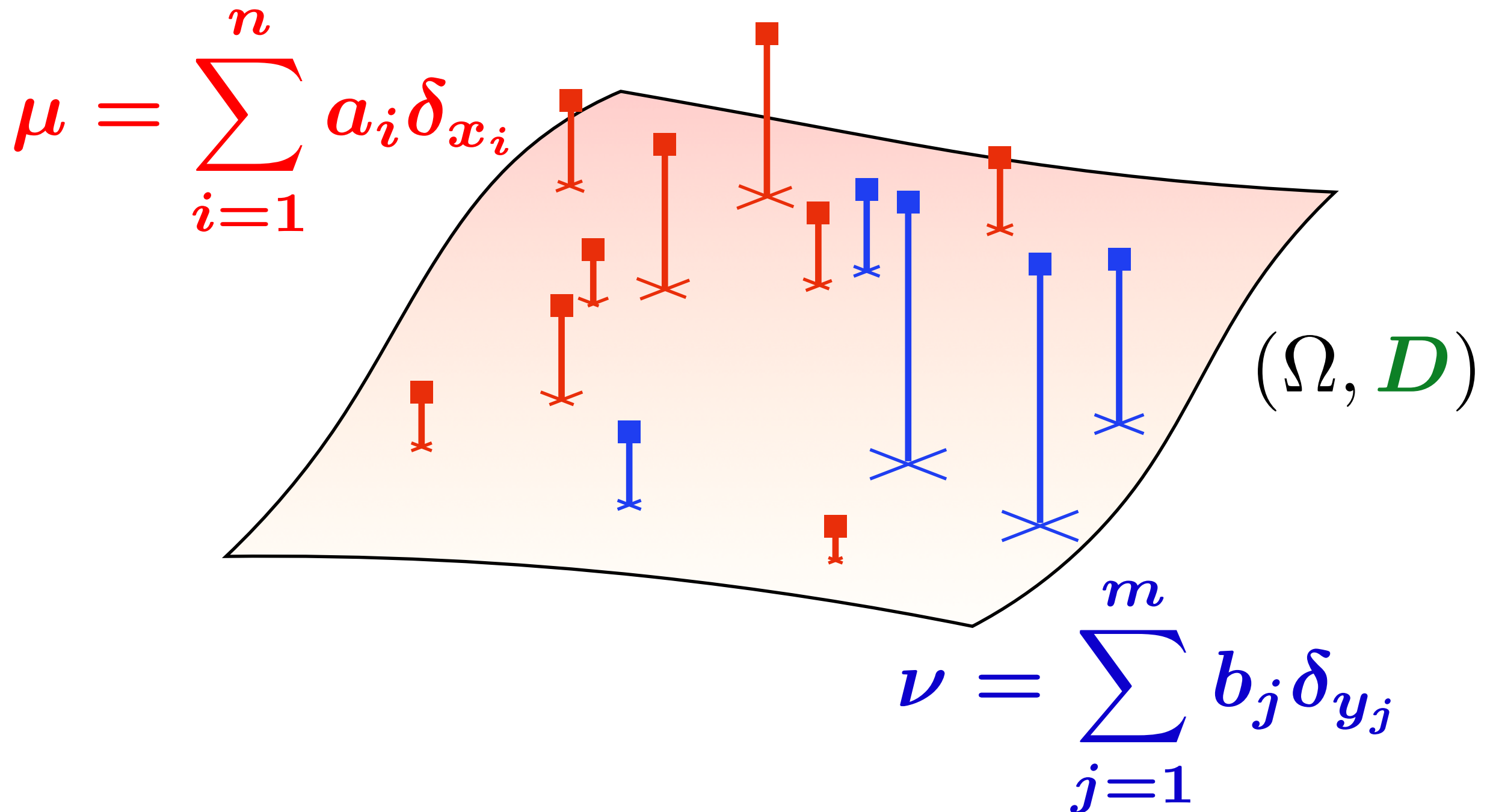- **[Sinkhorn'64]** with *matrix* fixed-point iterations



$V_1$

$A$

$B$

$K$

$:=$

*etc....*

$U_1$

$:=$

# Differentiability of $W$

$$W\big((\textcolor{red}{a, X}), (\textcolor{blue}{b, Y})\big)$$



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

$$(\Omega, \textcolor{green}{D})$$

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

48

# Differentiability of $W$

$$W((\textcolor{red}{a + \Delta a, X}), \textcolor{blue}{(b, Y)}) = W(\textcolor{red}{(a, X)}, \textcolor{blue}{(b, Y)}) + ??$$

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

$(\Omega, \textcolor{green}{D})$

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

48

# Differentiability of $W$

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

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



$(\Omega, D)$

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

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

$$W\big((\textcolor{red}{a, X + \Delta X}), (\textcolor{blue}{b, Y})\big) = W\big((\textcolor{red}{a, X}), (\textcolor{blue}{b, Y})\big) + ??$$

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



$$(\Omega, \textcolor{green}{D})$$

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

49

# Sinkhorn ⤳ *Differentiability*

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

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

$(\Omega, D)$

$$X \leftarrow X + \Delta X$$

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

49

# Sinkhorn: A Programmer View

**Def.** For $L \geq 1$, define
$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \stackrel{\text{def}}{=} \langle \boldsymbol{P_L}, M_{\boldsymbol{XY}} \rangle,$$



$$u_L^T(\, K \odot M_{\boldsymbol{XY}})v_L$$

Sinkhorn $\ell = 1, \ldots, L-1$

# Sinkhorn: A Programmer View

**Def.** For $L \geq 1$, define
$$W_L(\boldsymbol{\mu}, \boldsymbol{\nu}) \overset{\text{def}}{=} \langle \boldsymbol{P_L}, M_{\boldsymbol{XY}} \rangle,$$

**Prop.** $\frac{\partial W_L}{\partial \boldsymbol{X}}, \frac{\partial W_L}{\partial \boldsymbol{a}}$ can be computed recursively, in $O(L)$ kernel $K \times$ vector products.
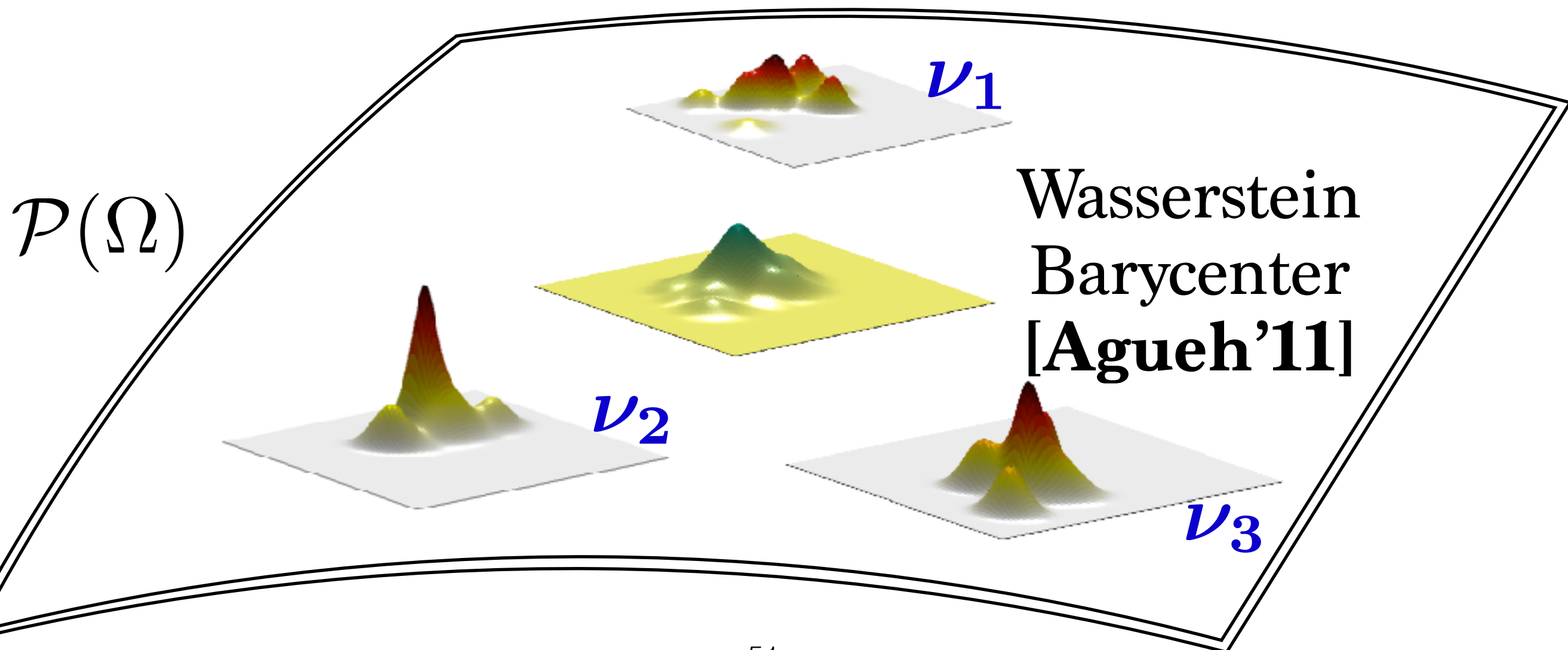
**[Hashimoto'16][Bonneel'16][Shalit'16][Flammary'16]**

$\mathcal{P}(\Omega)$

$\nu_1$

$\nu_2$

$\nu_3$

# OT: Barycenters

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



$\mathcal{P}(\Omega)$

$\boldsymbol{\nu_1}$

$\boldsymbol{\nu_2}$
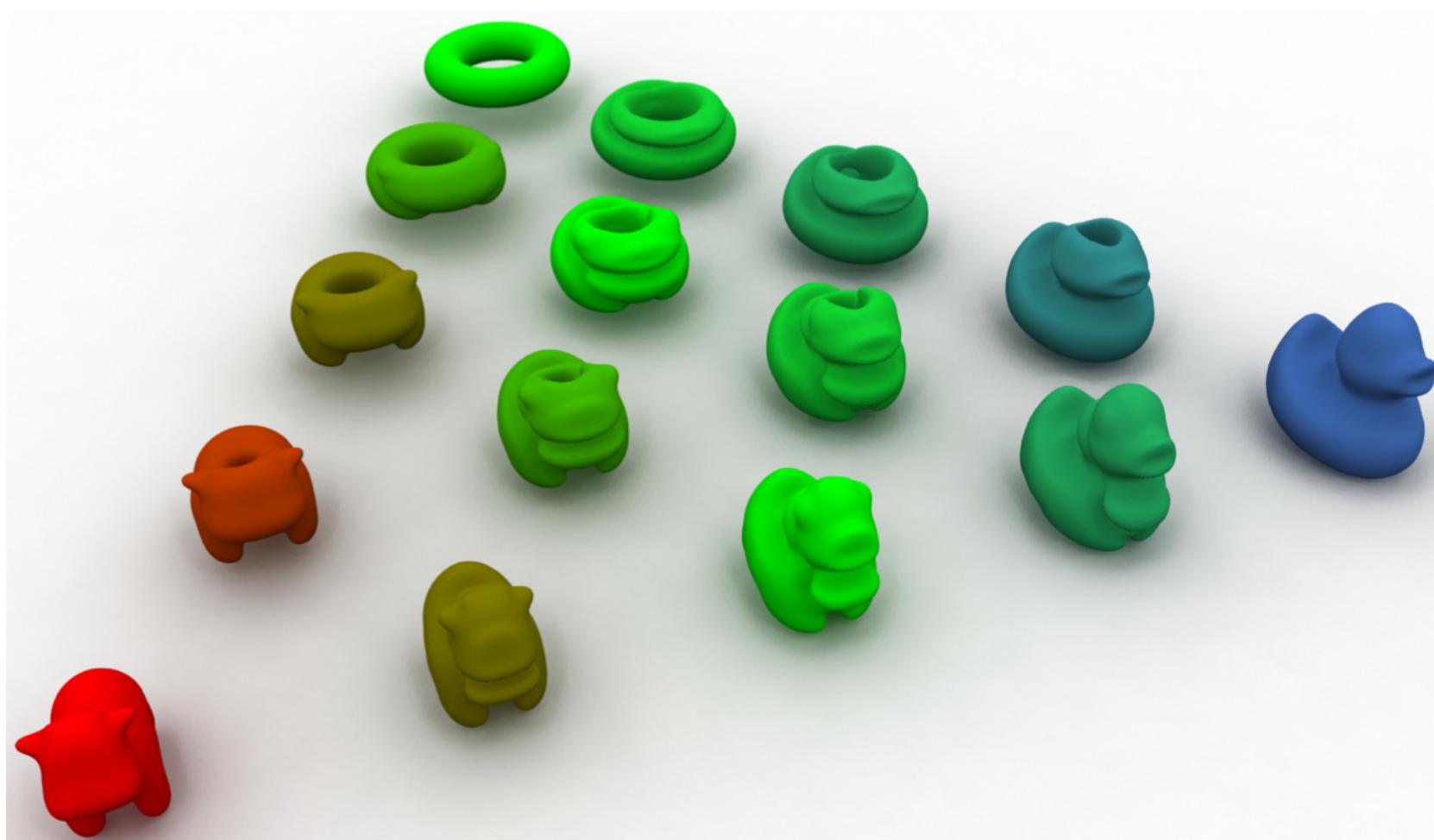
$\boldsymbol{\nu_3}$

Wasserstein Barycenter [**Agueh'11**]

# OT: Barycenters

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



**[Solomon'15]**

# OT: Barycenters

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



**[Solomon'15]**

# OT: Dictionary Learning

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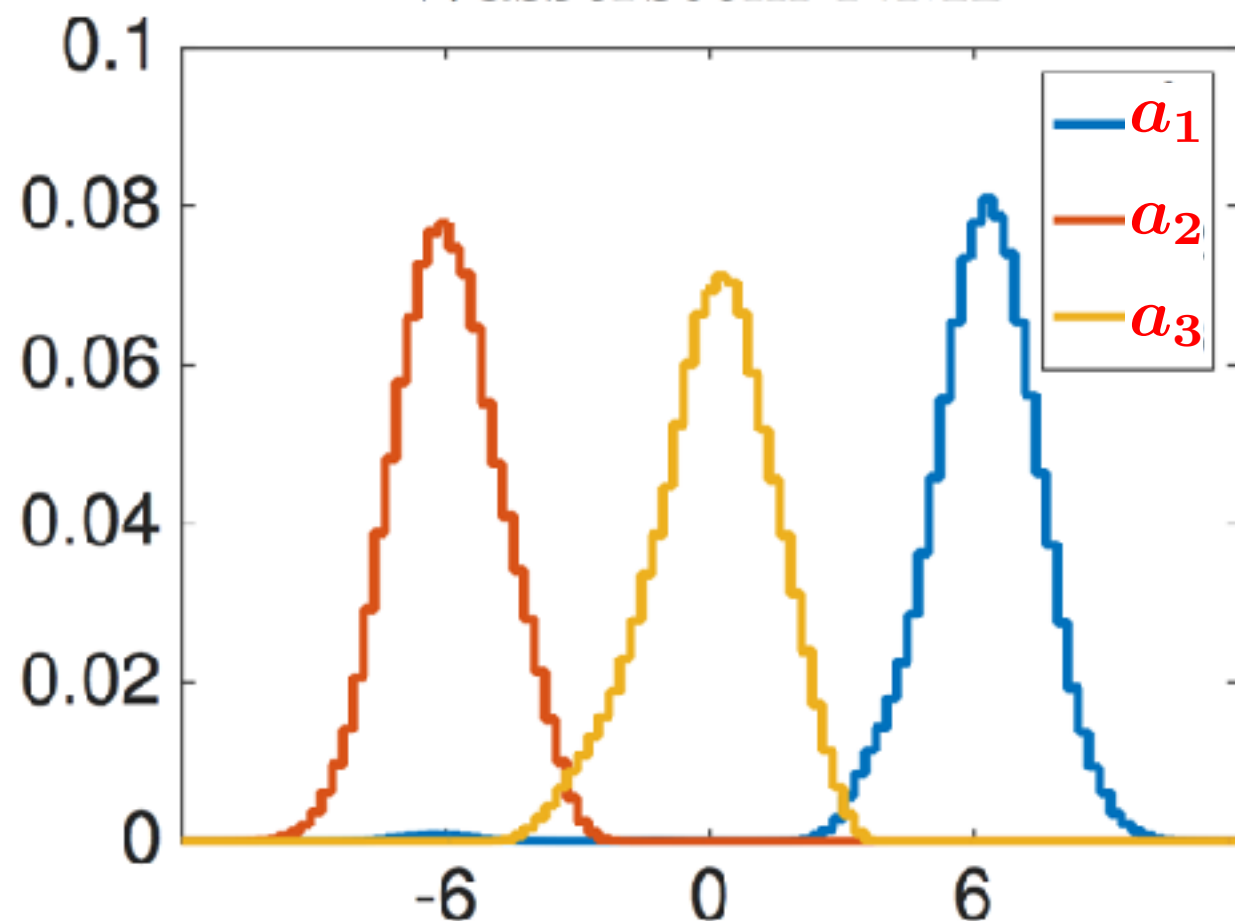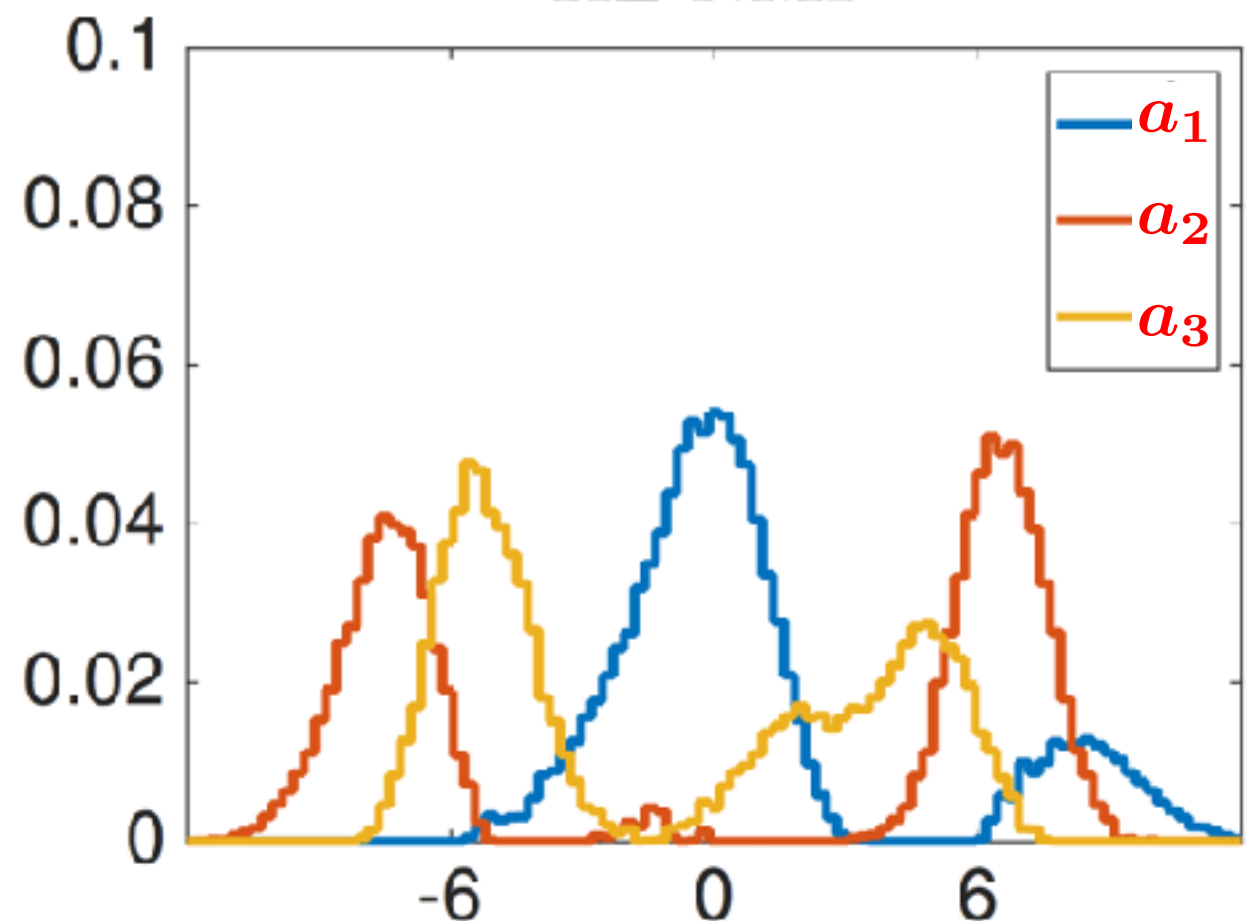


Data samples

Data samples

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

# OT: Dictionary Learning

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



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

53

# (Word Mover's Distance)



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

# Topic Models



[Rolet'16]

Euclidean Simplex: $\left\{ \sum_{i=1}^{3} \lambda_i p_i, \lambda \in \Sigma_3 \right\}$

Wasserstein simplex: $\left\{ P(\lambda), \lambda \in \Sigma_3 \right\}$

Shape database $(p_1, \ldots, p_5)$

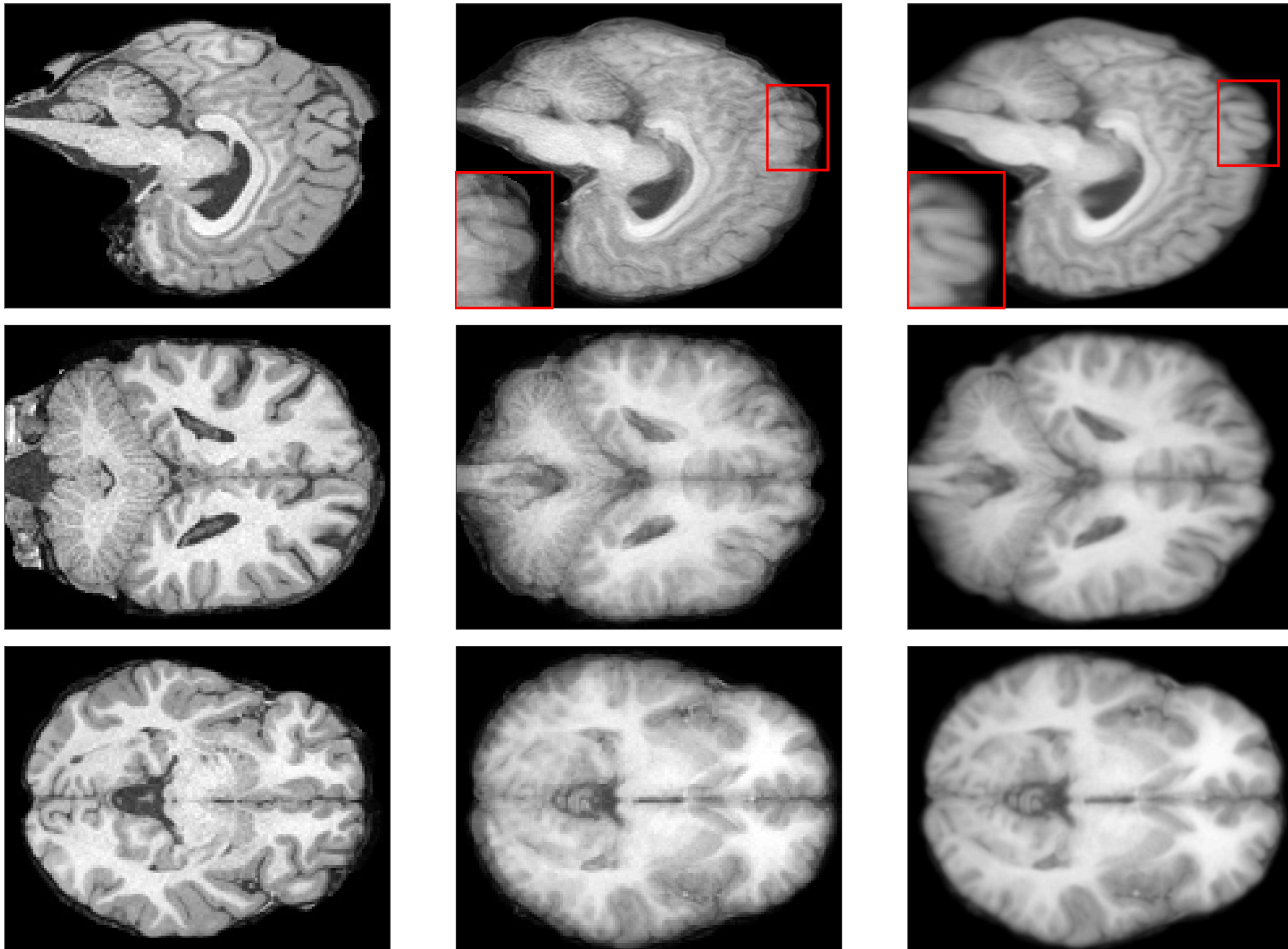Input shape $q$

Projection $P(\lambda)$

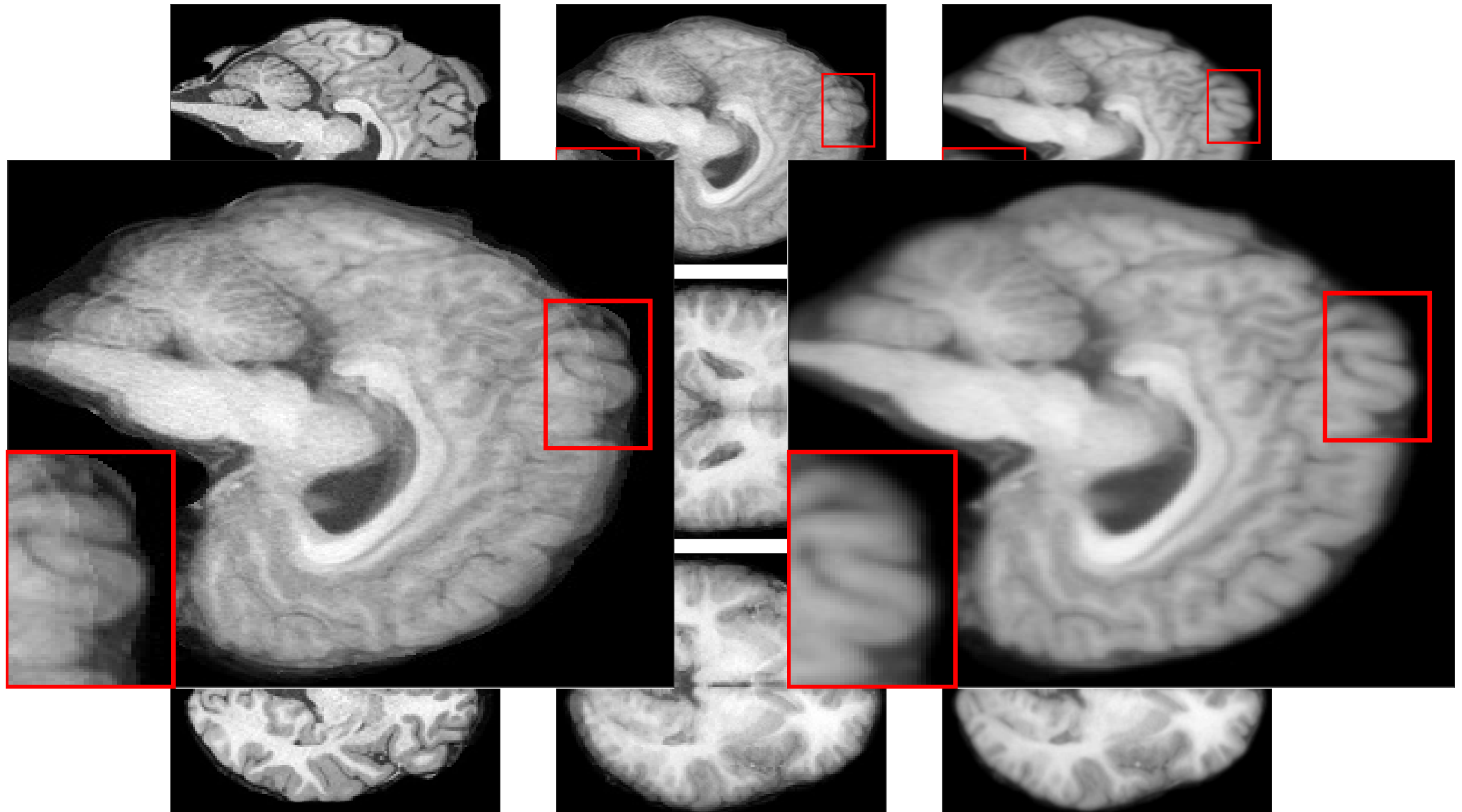Iso-surface

**[Bonneel'16]**

57

# Application: Brain Mapping



Original       Euclidean      Wasserstein
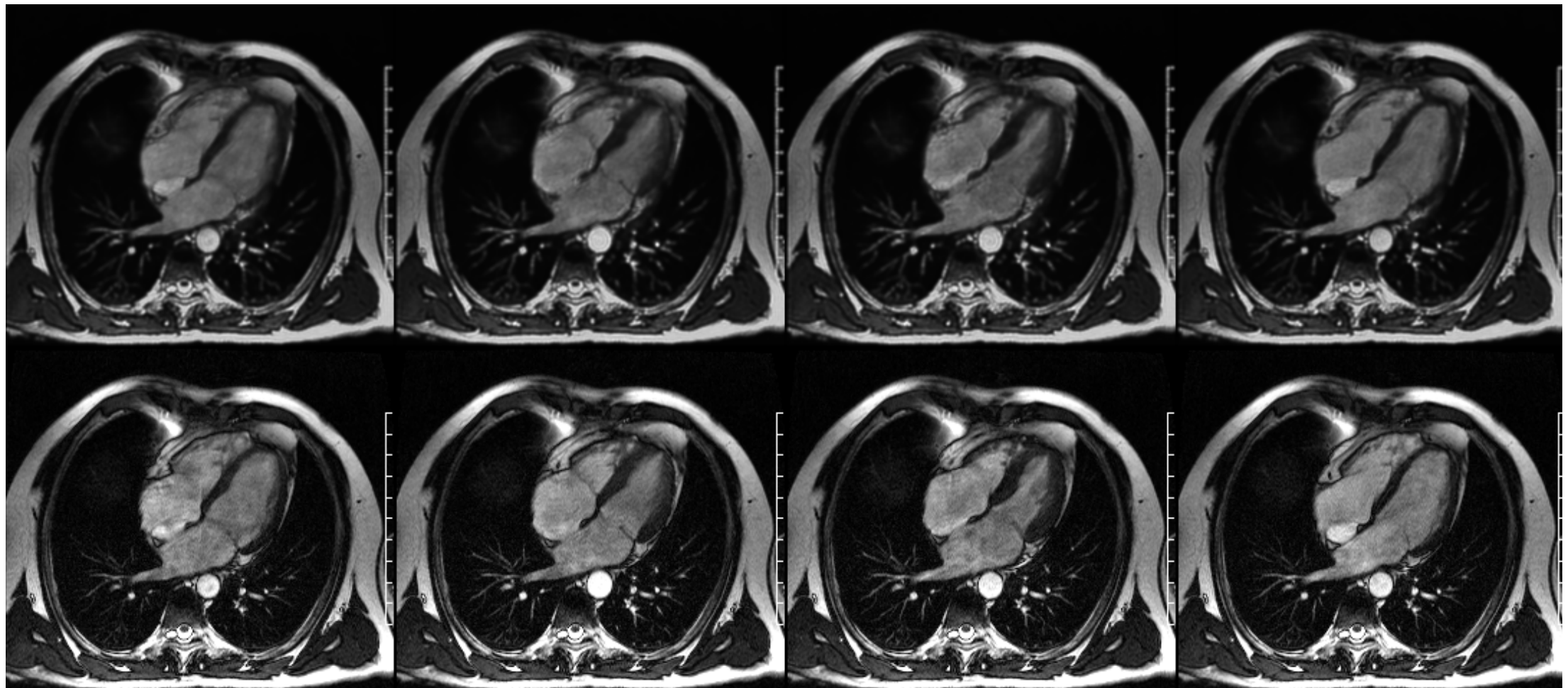projection     projection

# Application: Brain Mapping

$$\min_{\substack{\textcolor{red}{\boldsymbol{A}} \in (\Sigma_n)^K \; \textcolor{green}{\boldsymbol{\Lambda}} \in (\Sigma_K)^N}} \sum_{i=1}^{N} \mathcal{L}\left(\textcolor{blue}{\boldsymbol{b_i}}, \mathcal{B}_{\textcolor{green}{\boldsymbol{\Lambda}_k^i}}\left(\textcolor{red}{\boldsymbol{a_k}}\right)\right)$$

# Application: $W$ Dictionary Learning

$$\min_{\substack{\textcolor{red}{A} \in (\Sigma_n)^K \ \textcolor{green}{\Lambda} \in (\Sigma_K)^N}} \sum_{i=1}^{N} \mathcal{L}\left(\textcolor{blue}{b_i}, \mathcal{B}_{\textcolor{green}{\Lambda_k^i}}(\textcolor{red}{a_k})\right)$$



59

# 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 $f_{\boldsymbol{\theta}}$ : Images $\mapsto$ Labels

# Learning with a Wasserstein Loss

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



$x_i$

husky
snow
sled
slope
men

$y_i$

Which loss $\mathcal{L}$ could we use?

# Learning with a Wasserstein Loss

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

dog
driver
winter
ice

$f_{\boldsymbol{\theta}}(x_i)$

husky
snow
sled
slope
men

$y_i$

Which loss $\mathcal{L}$ could we use?

# Learning with a Wasserstein Loss

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

$$\mathcal{L}(\boldsymbol{a}, \boldsymbol{b}) = \min_{\boldsymbol{P} \in \mathbb{R}^{nm}} \langle \boldsymbol{P}, M \rangle + \varepsilon \text{KL}(\boldsymbol{P}\mathbf{1}, \boldsymbol{a})$$

$$+ \varepsilon \text{KL}(\boldsymbol{P}^T \mathbf{1}, \boldsymbol{b}) - \gamma E(\boldsymbol{P})$$

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

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

# Minimum Kantorovich Estimators

$$\min_{\boldsymbol{\theta} \in \Theta} W(\boldsymbol{\nu}_{\text{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})$$

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

$$\text{Challenge: } \nabla_{\boldsymbol{\theta}} W(\boldsymbol{\nu}_{\text{data}}, f_{\boldsymbol{\theta}\sharp}\boldsymbol{\mu})?$$
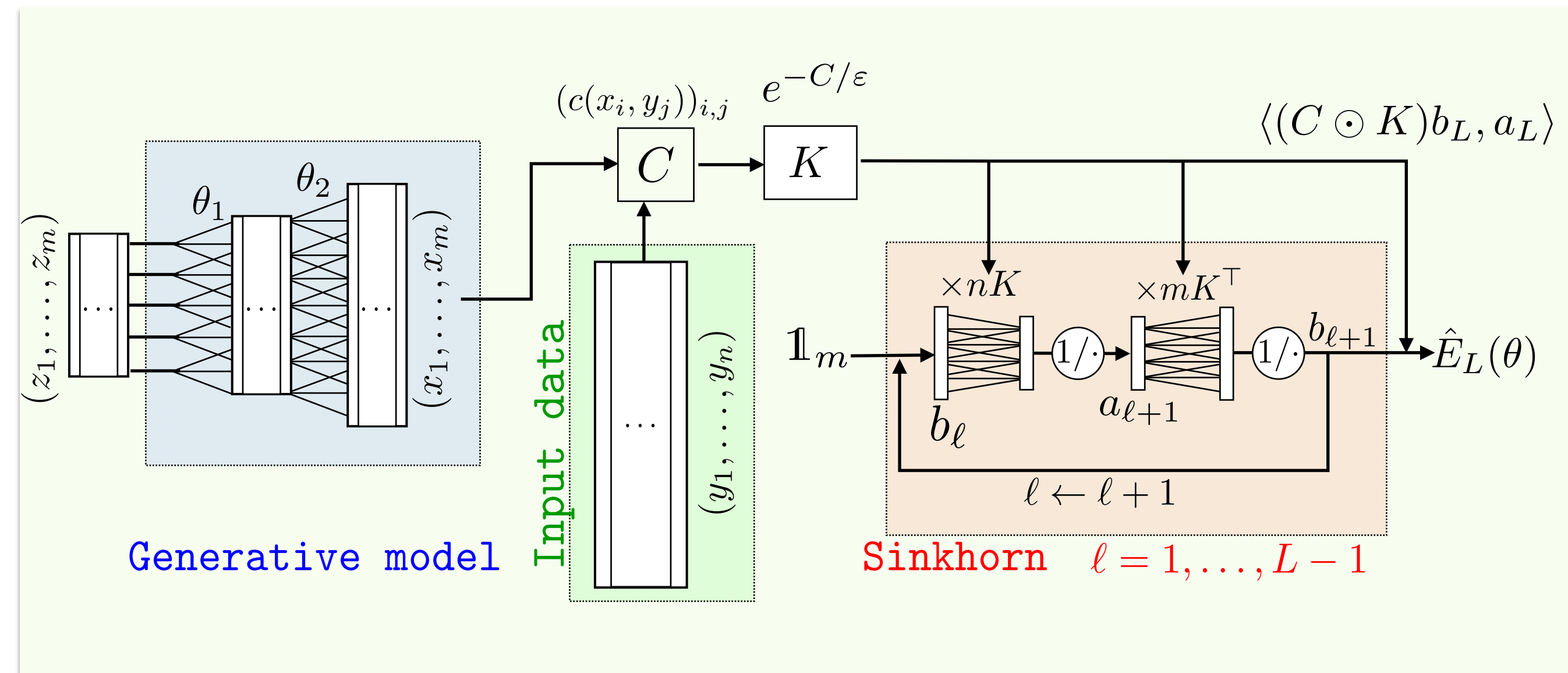
[**Montavon'16**] use regularized OT in a finite setting.

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

[**Bernton'17**] reject mechanism *W*(sample, data)

[**Genevay'17, Salimans'17**] *(Sinkhorn approach)*
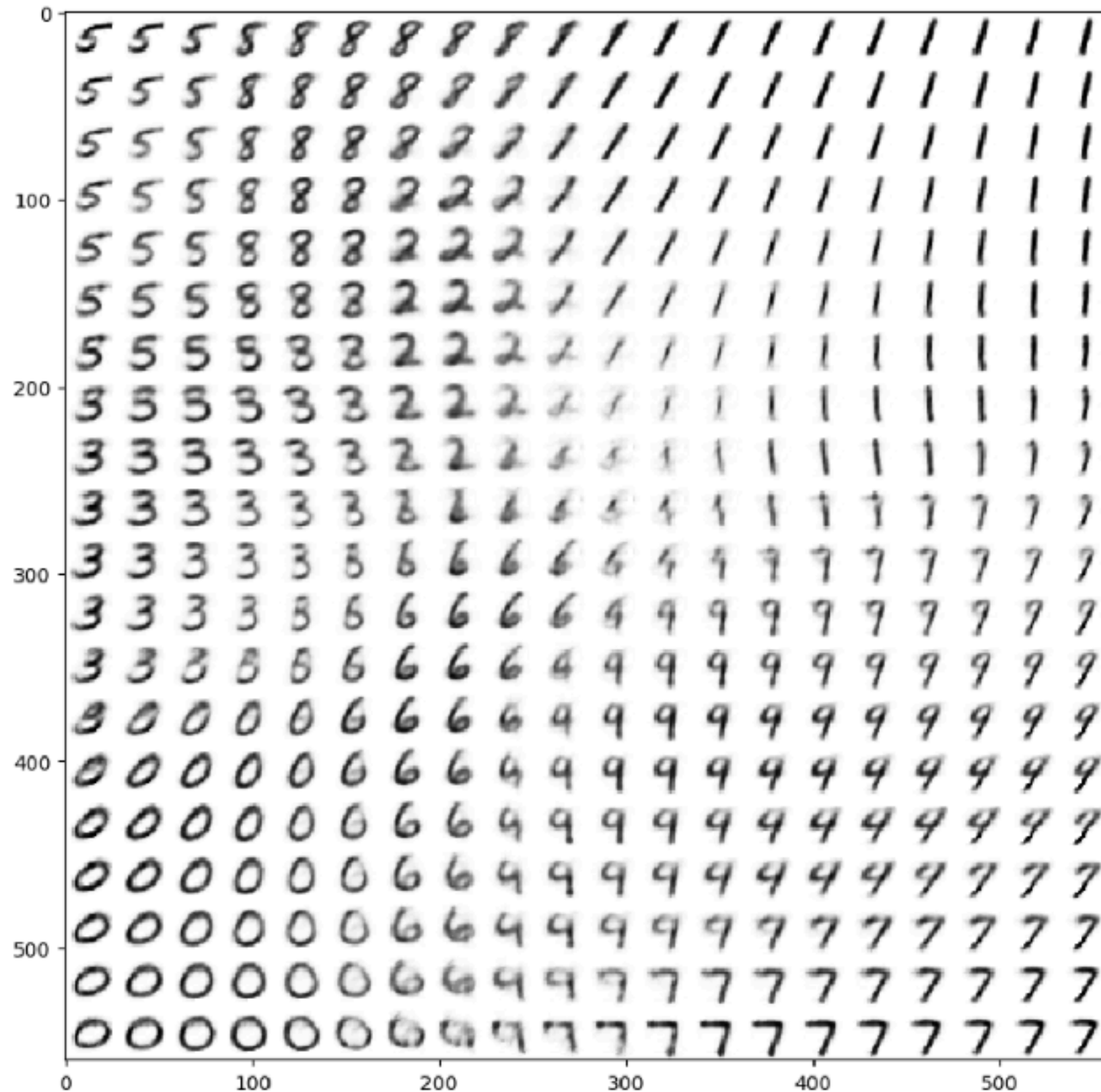
# Proposal: Autodiff OT using Sinkhorn

Approximate $W$ loss by the transport cost $\bar{W}_L$ after $L$ Sinkhorn iterations.
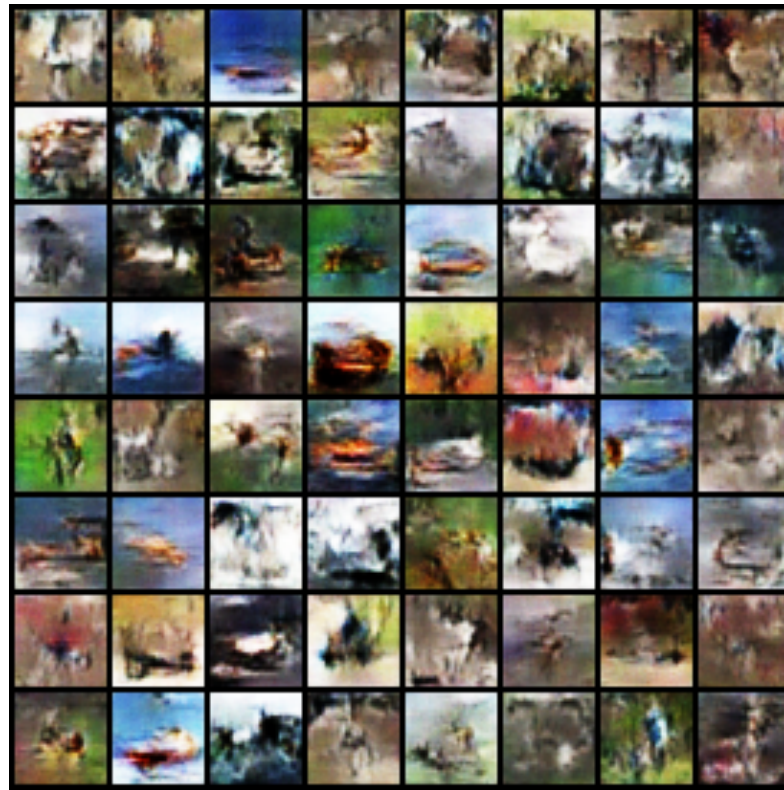
**[GPC'17]**

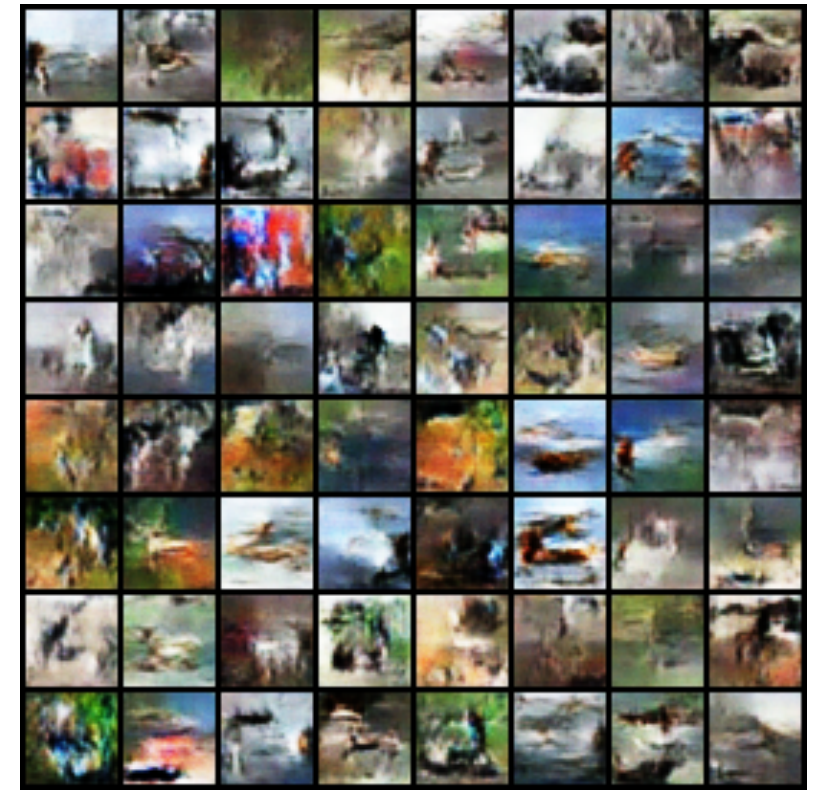# Example: MNIST, Learning $f_\theta$

# Example: Generation of Images
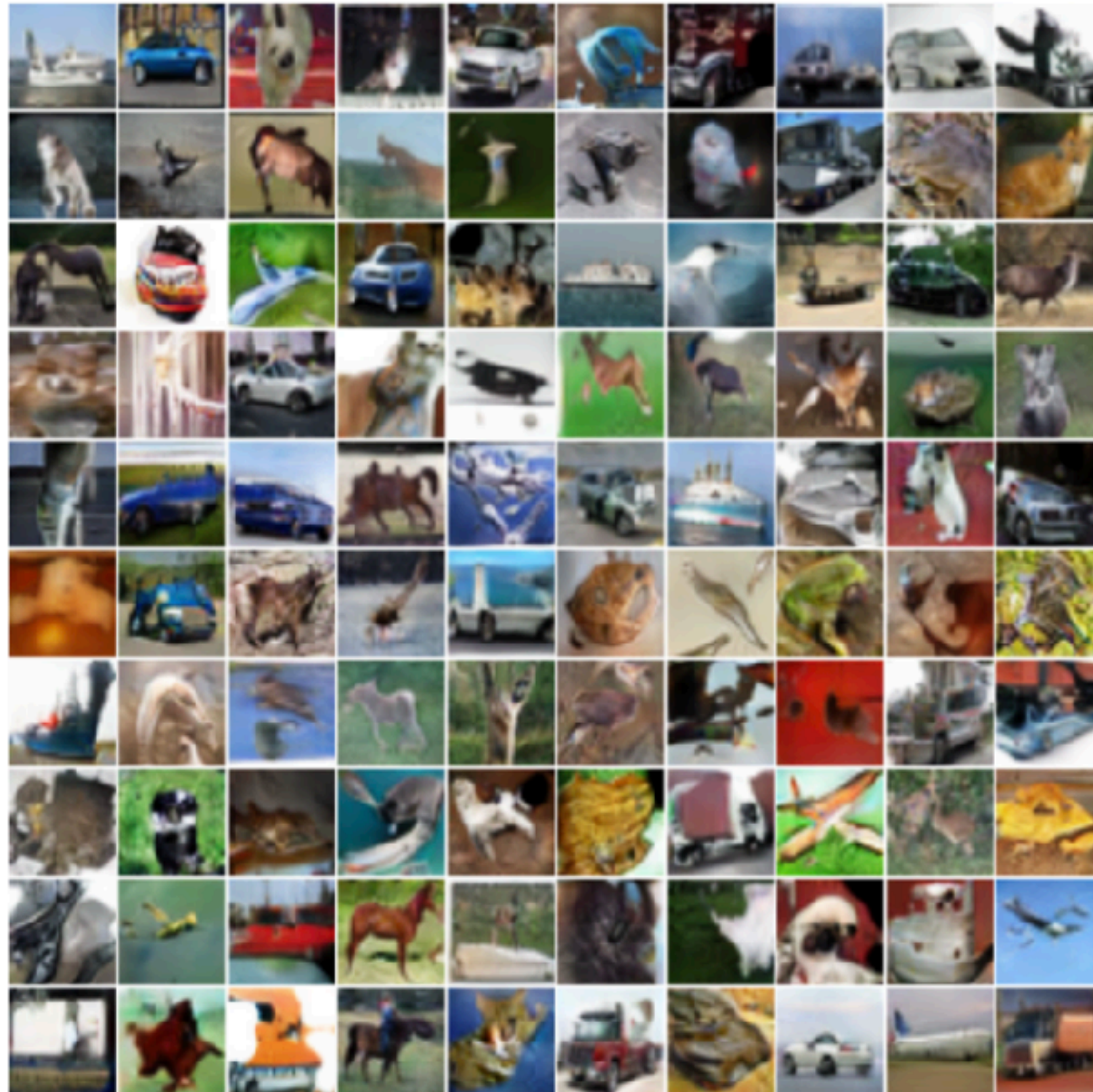


MMD-GAN                 τ = 1000                 τ=10

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

# 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* metrics.

- The solution of *regularized OT* is *"auto-differentiable"*.

- **Many open problems remain!**