

Semi-Discrete Optimal Transport: Hardness, Regularization and Numerical Solution



Soroosh Shafieezadeh Abadeh, Tepper School of Business, CMU

Optimal Transport: Old and New



Monge



Hitchcock



Kantorovich-Koopmans



Nobel '75



Vaserstein



Brenier



Villani

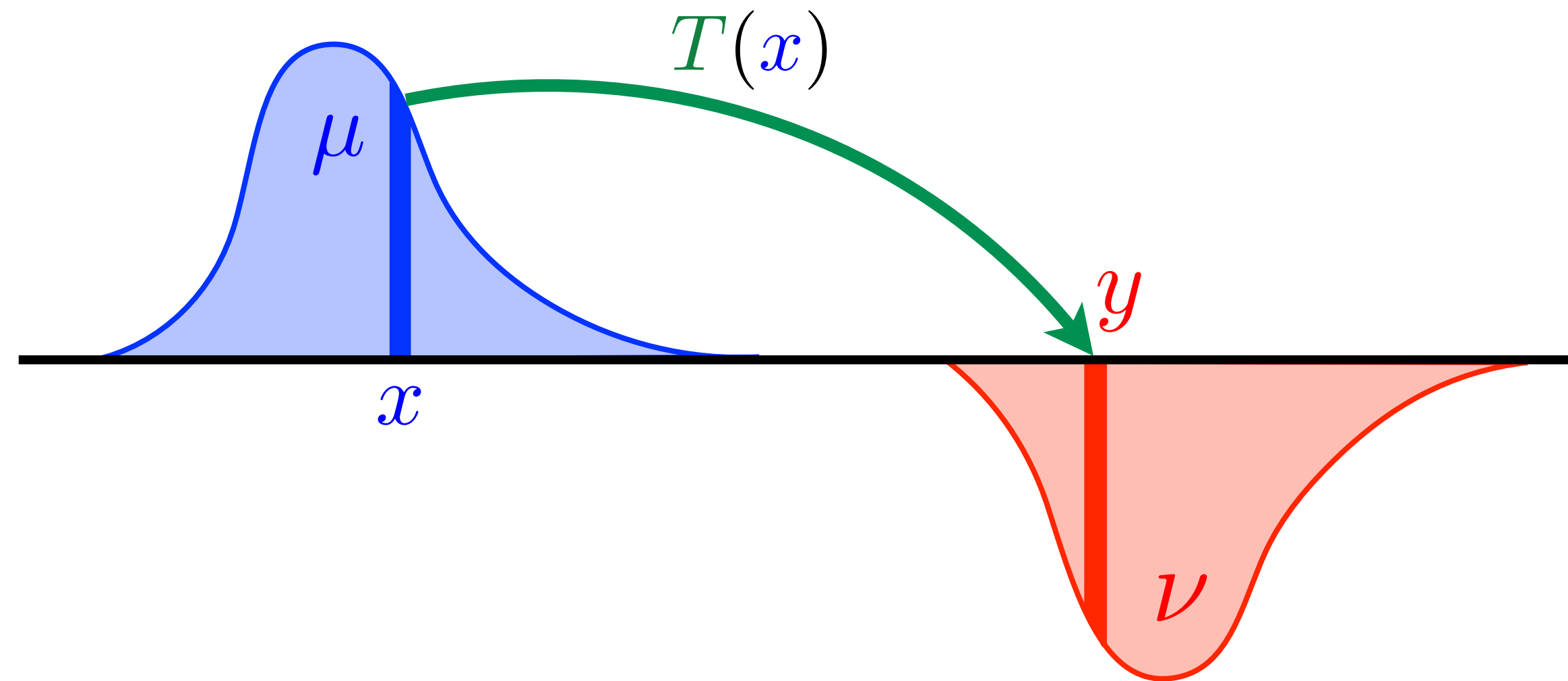
Fields '10



Figalli

Fields '18

What is Optimal Transport?



$$\begin{aligned} \inf_{T: X \rightarrow Y} & \quad \mathbb{E}_{\mu} [c(x, T(x))] \\ \text{s.t.} & \quad T\# \mu = \nu \end{aligned}$$

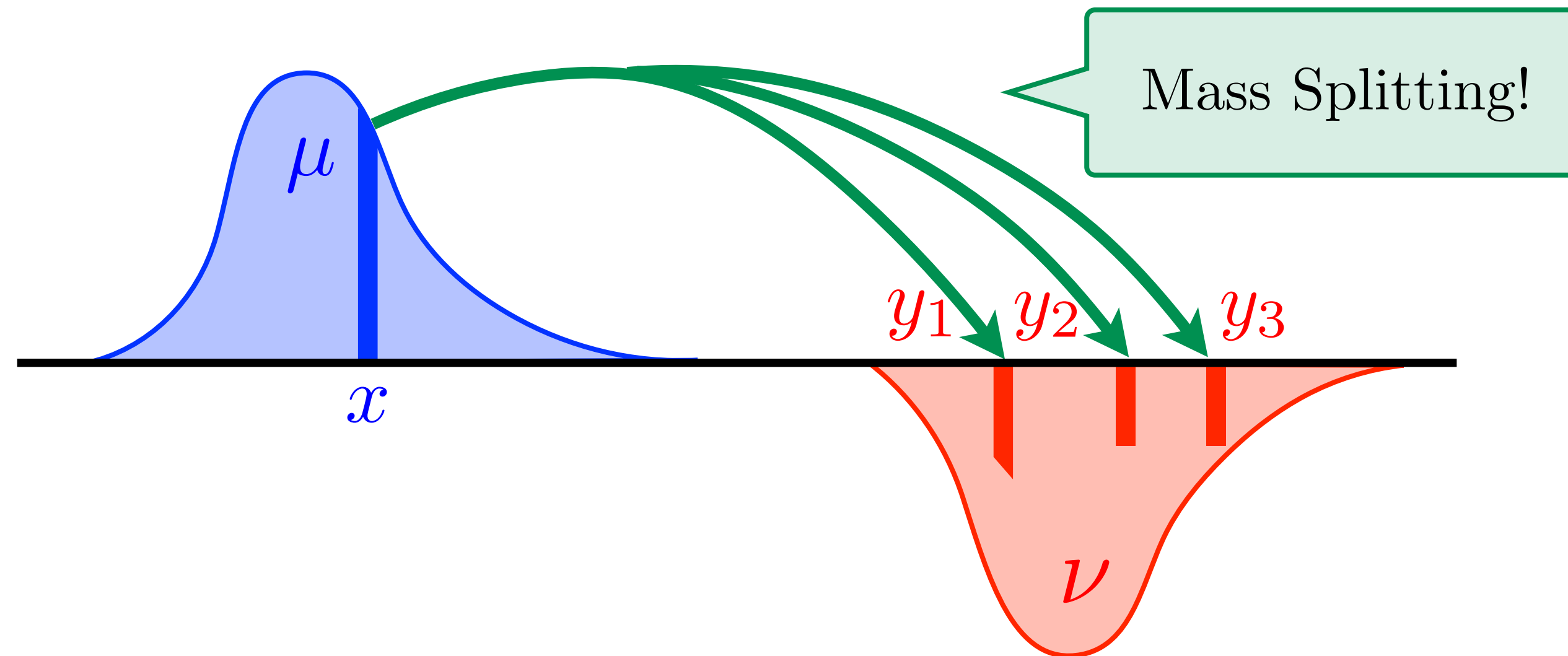
Nonconvex and infeasible!



M É M O I R E
SUR LA
THÉORIE DES DÉBLAIS
ET DES REMBLAIS.
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. Le prix du transport d'une molécule étant, toutes choses d'ailleurs égales, proportionnel à son poids & à l'espace qu'on lui fait parcourir, & par conséquent le prix du transport total devant être proportionnel à la somme des produits des molécules multipliées chacune par l'espace parcouru, il s'en suit que le déblai & le remblai étant donnés de figure & de position, il n'est pas indifférent que telle molécule du déblai soit transportée dans tel ou tel autre endroit du remblai, mais qu'il y a une certaine distribution à faire des molécules du premier dans le second, d'après laquelle la somme de ces produits fera la moindre possible, & le prix du transport total fera un *minimum*.

What is Optimal Transport?



$$\begin{aligned} \inf_{\pi \in \mathcal{M}(X, Y)} \quad & \mathbb{E}_{\pi} [c(x, y)] \\ \text{s.t.} \quad & \pi \in \Pi(\mu, \nu) \end{aligned}$$

Infinite Dimensional LP



THE DISTRIBUTION OF A PRODUCT FROM SEVERAL SOURCES TO NUMEROUS LOCALITIES
BY FRANK L. HITCHCOCK 1941

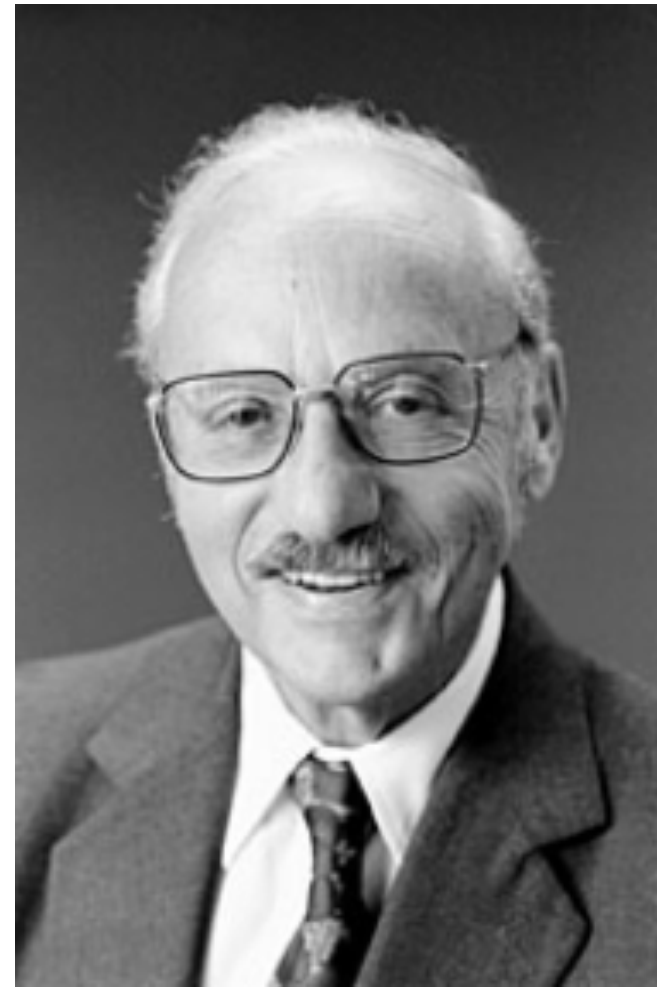
1. **Statement of the problem.** When several factories supply a product to a number of cities we desire the least costly manner of distribution. Due to freight rates and other matters the cost of a ton of product to a particular city will vary according to which factory supplies it, and will also vary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

OPTIMUM UTILIZATION OF THE TRANSPORTATION SYSTEM* 1949

by Tjalling C. Koopmans
Professor of Economics, The University of Chicago, and Research Associate, Cowles Commission for Research in Economics

The purpose of this paper is to give an application of the theory of optimum allocation of resources to one particular industry. I shall, therefore, not speak on that theory in general. I shall use one of its

Discrete OP Problems



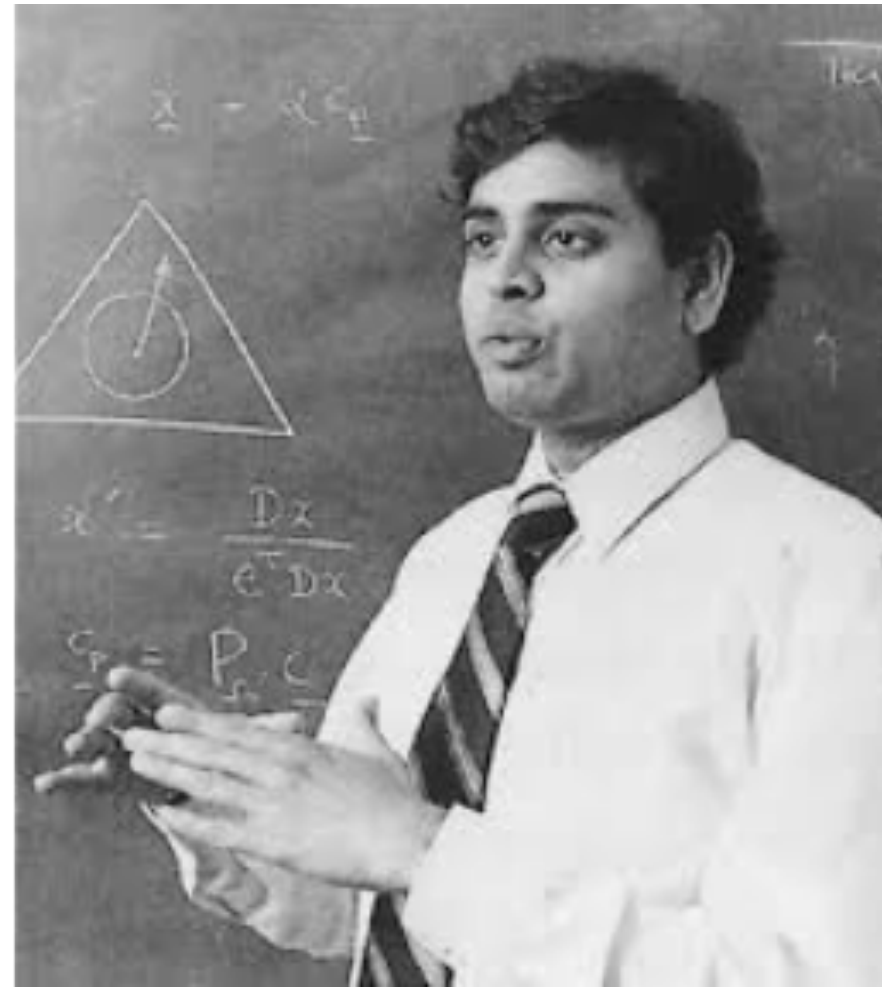
Dantzig

$$\mathcal{O}(2^n)$$



Khachiyan

$$\mathcal{O}(n^4)$$



Karmarkar

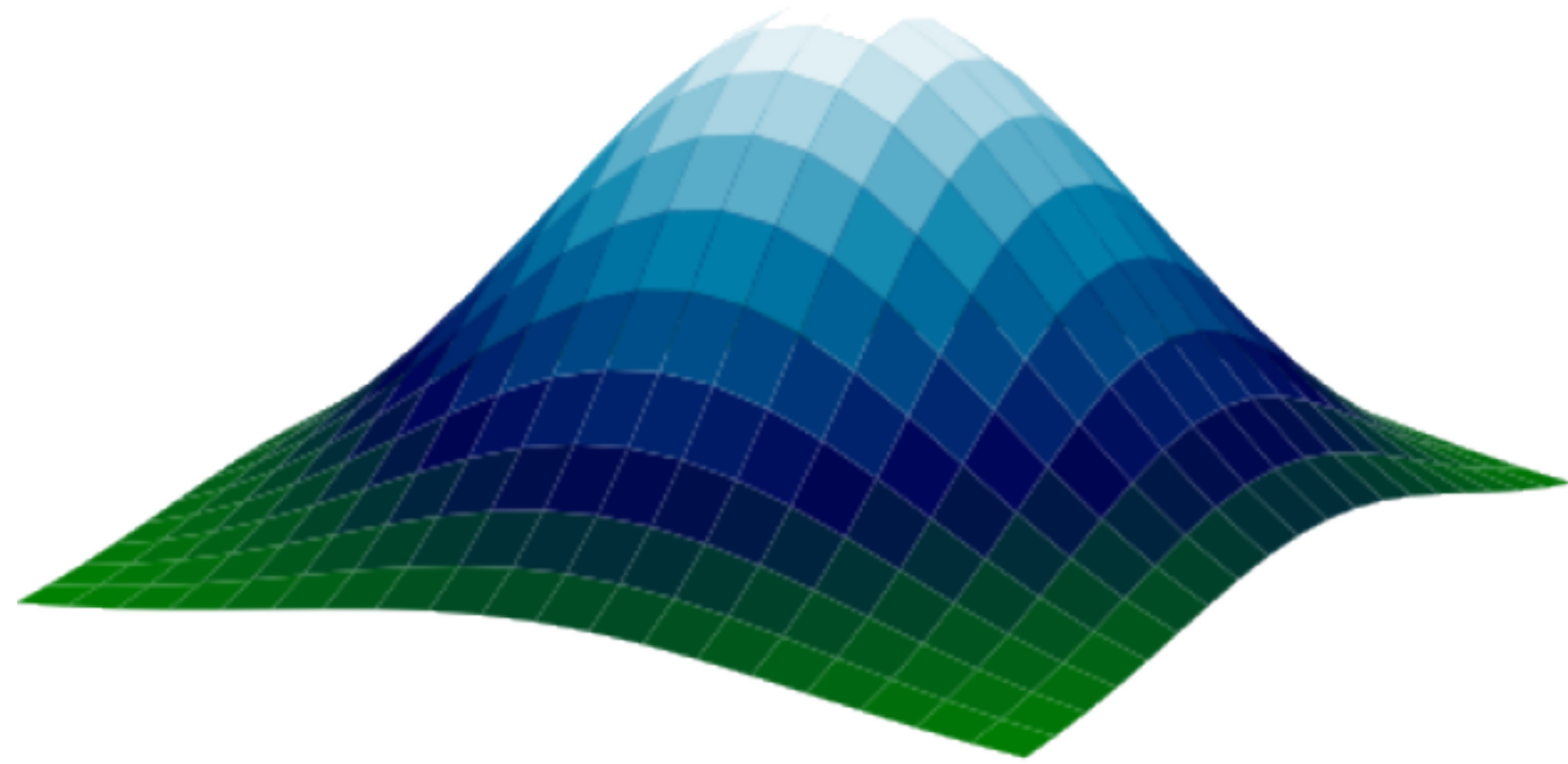
$$\mathcal{O}(n^{3.5})$$



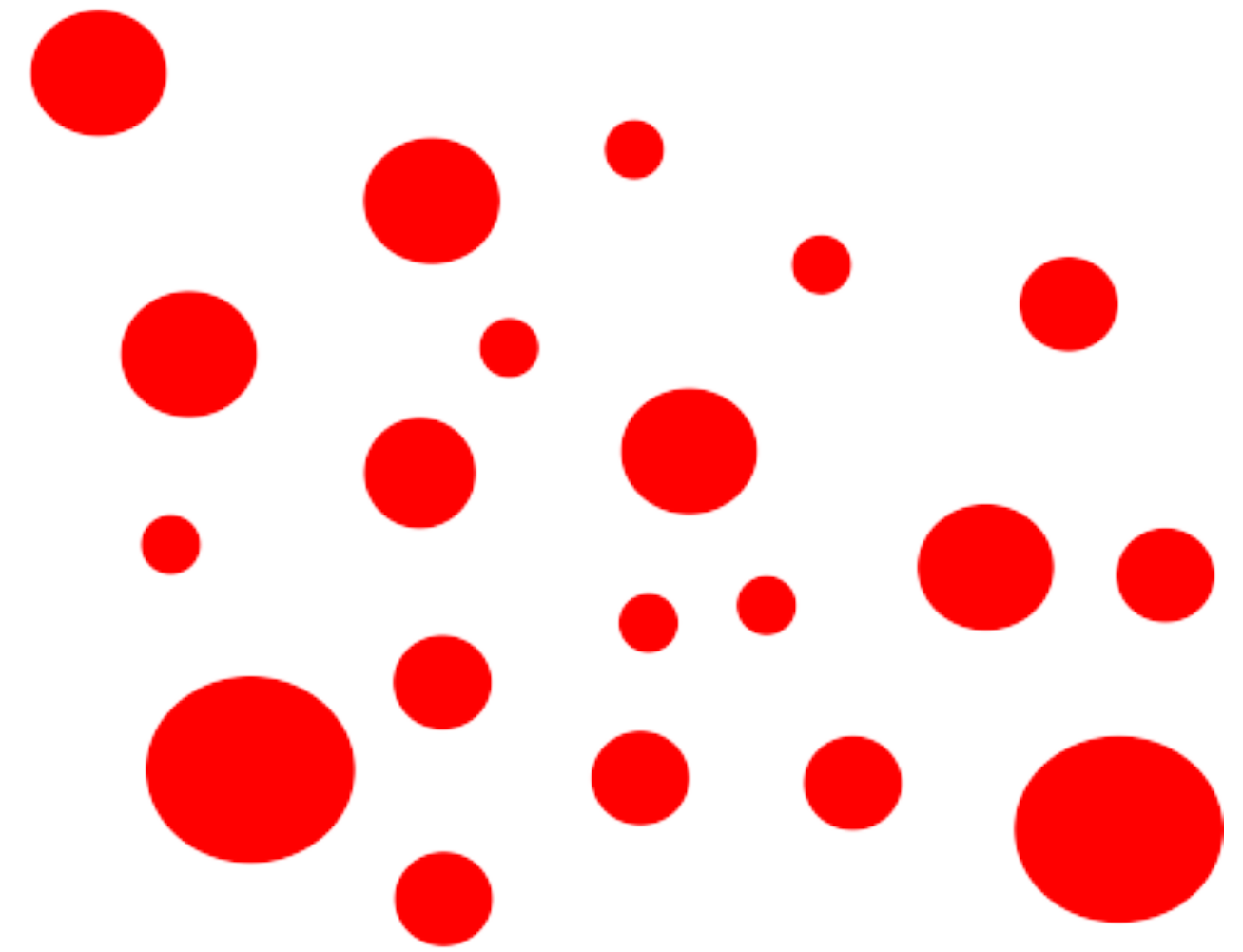
Sinkhorn

$$\mathcal{O}(n^2)$$

Semi-Discrete OT Problems



μ



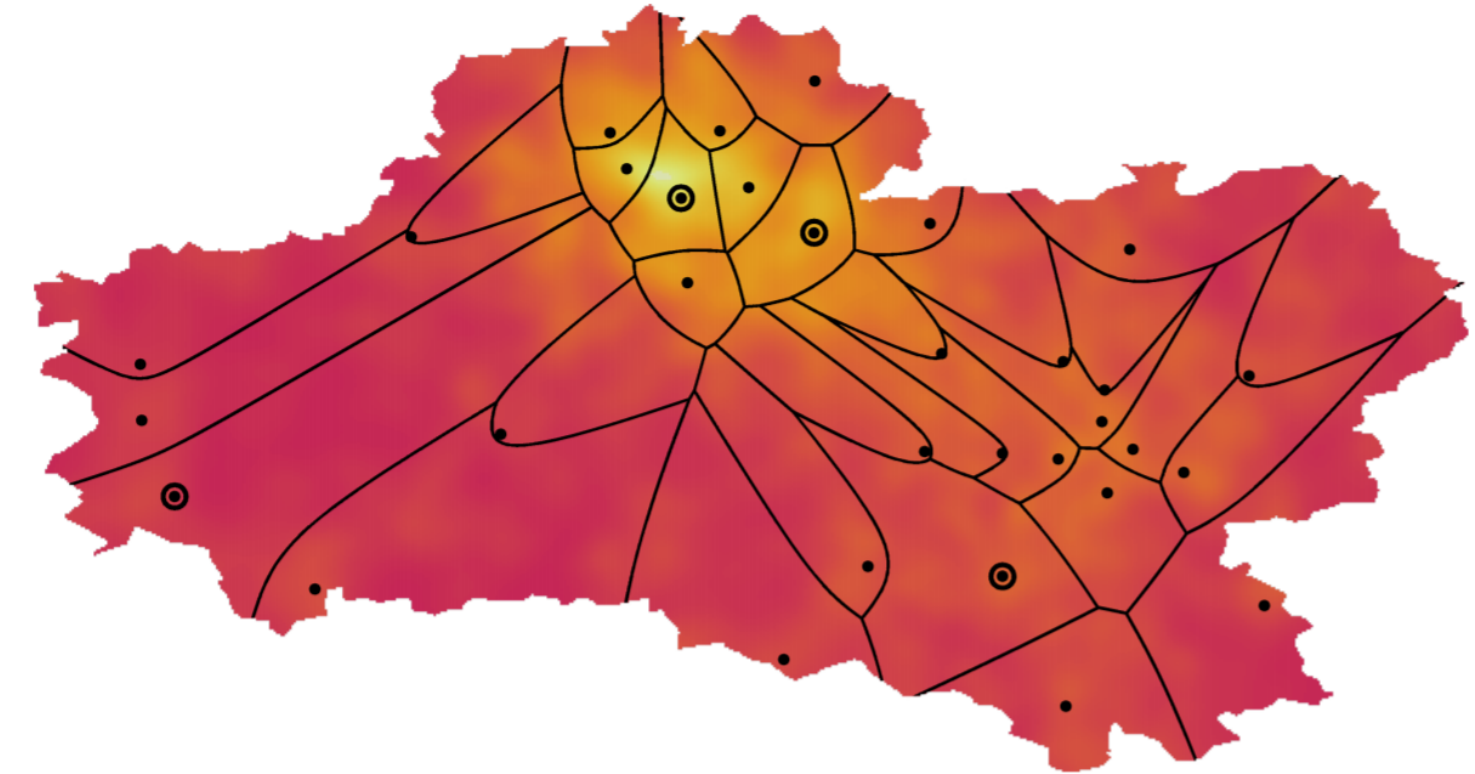
$$\nu = \sum_{i=1}^n \nu_i \delta_{y_i}$$

Applications

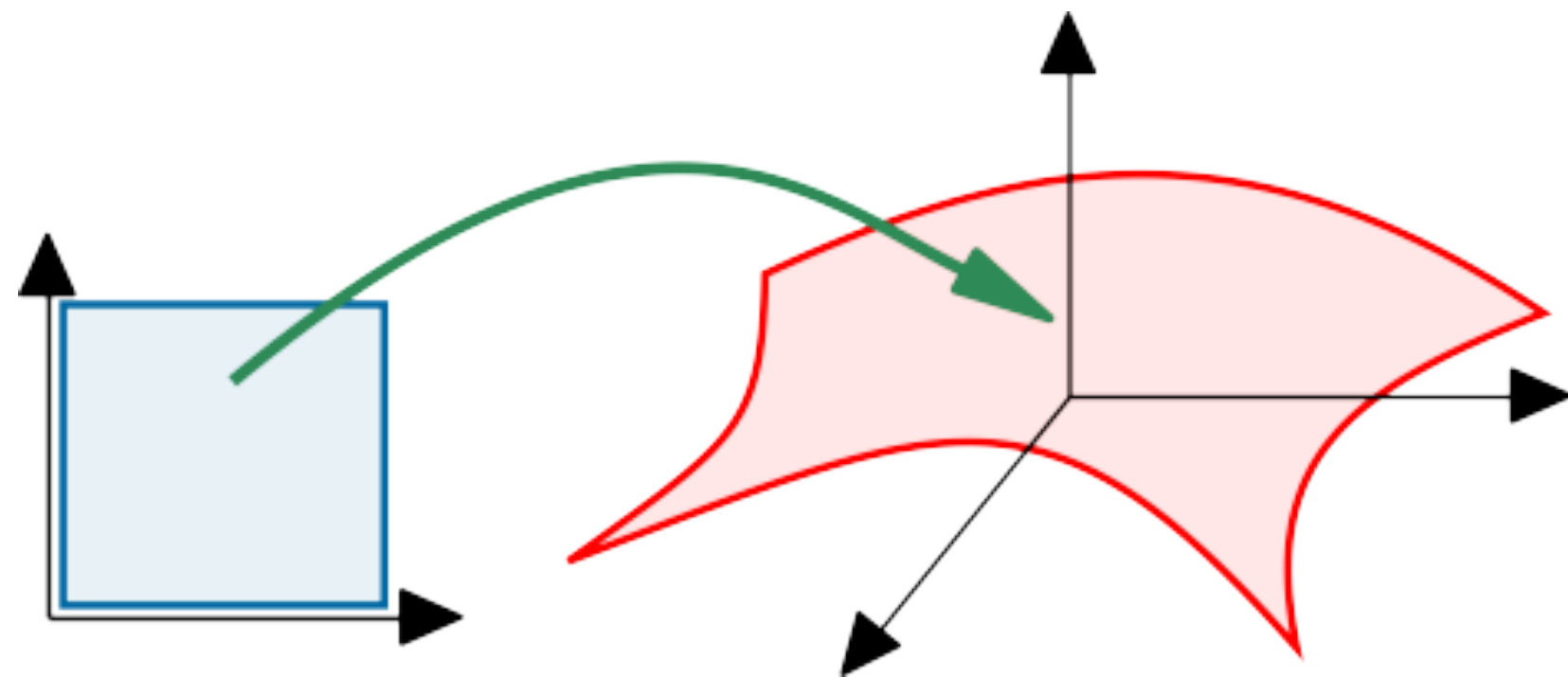
3D morphing [L14]



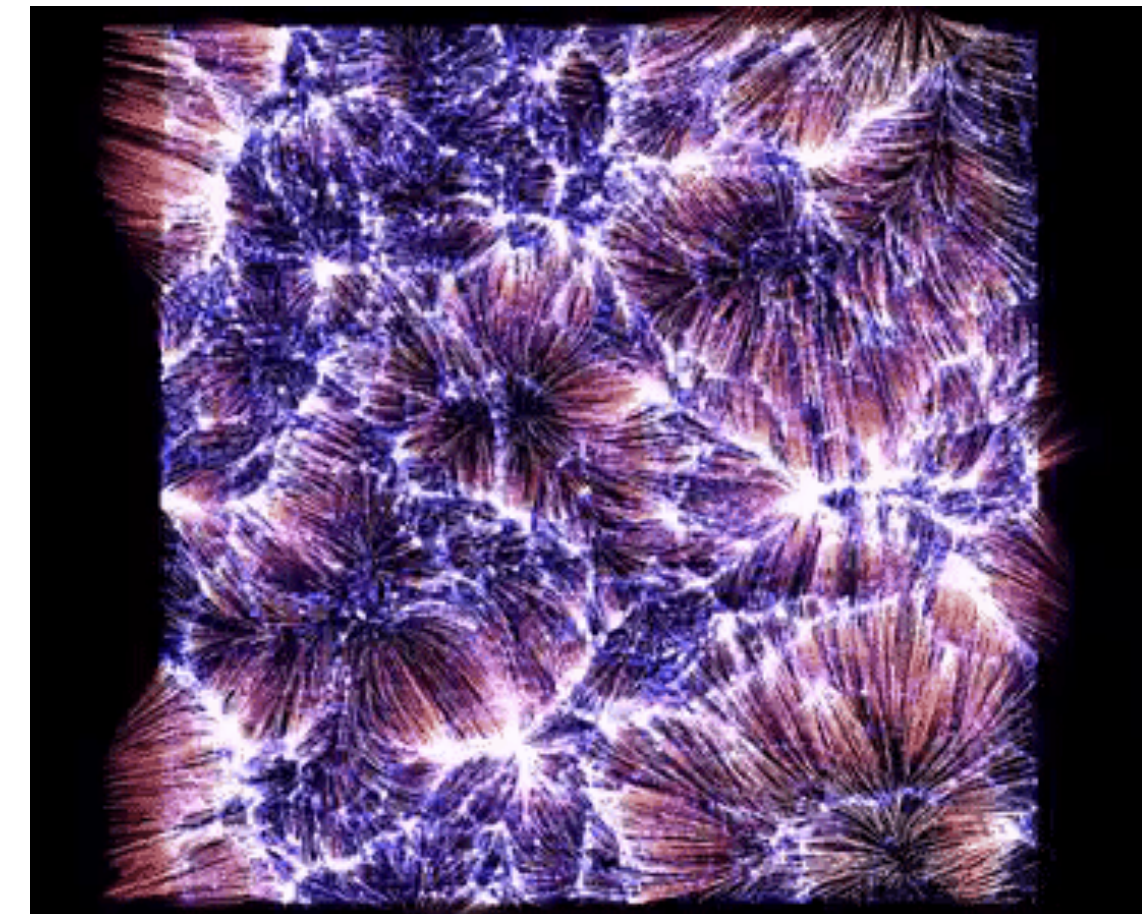
Resource allocation [HS20]



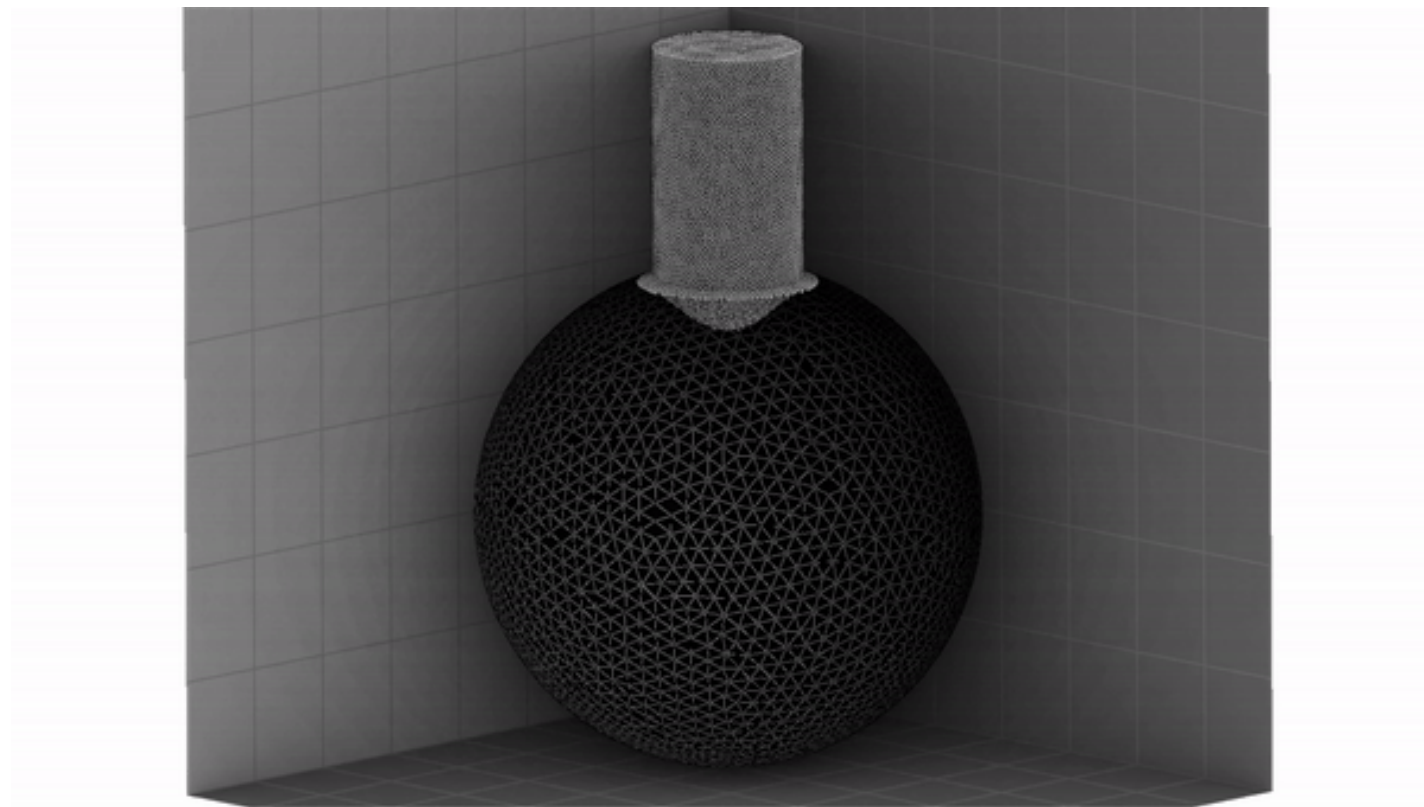
Generative models [HS20]



Reconstruction of early universe [LMHN21]

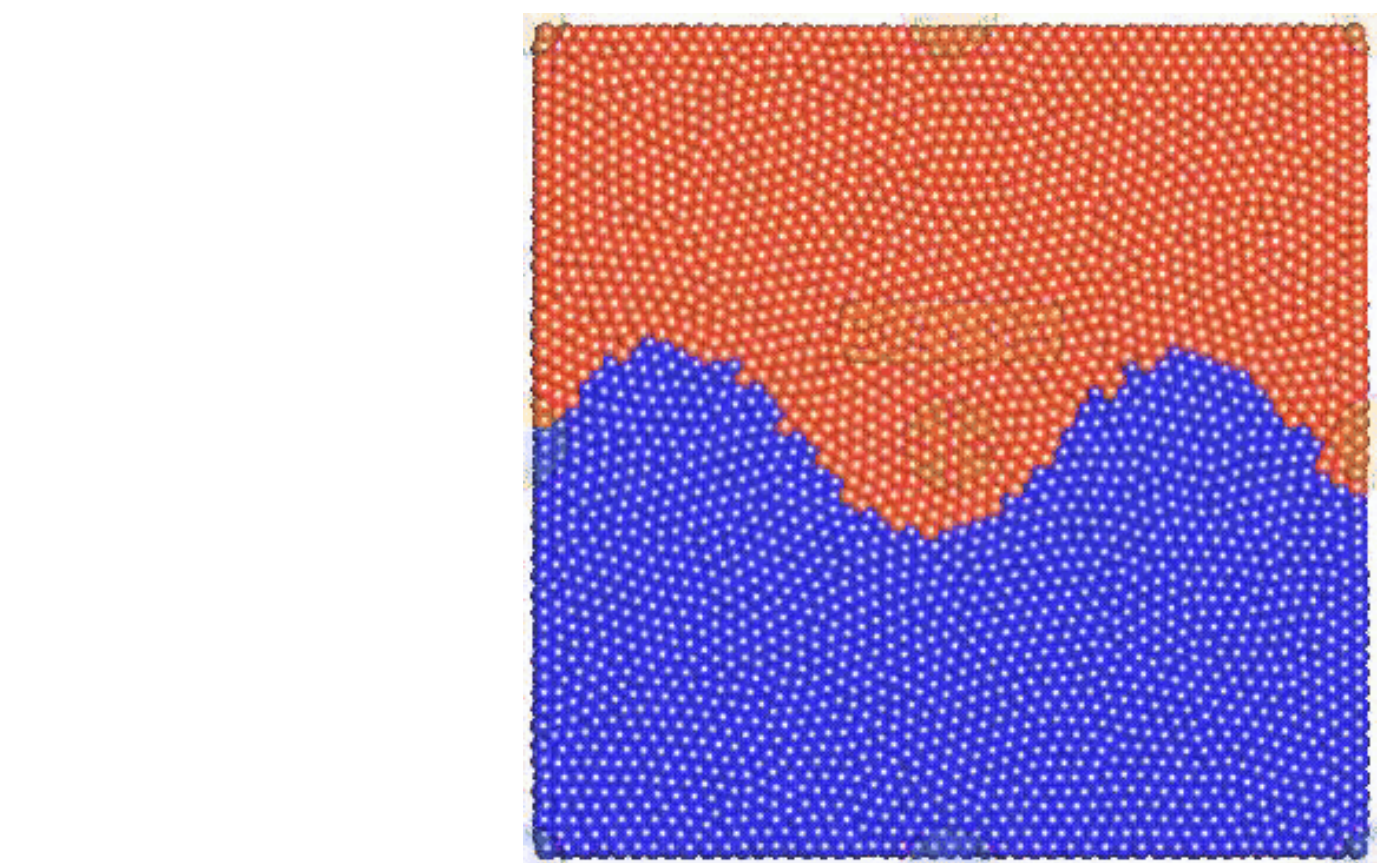
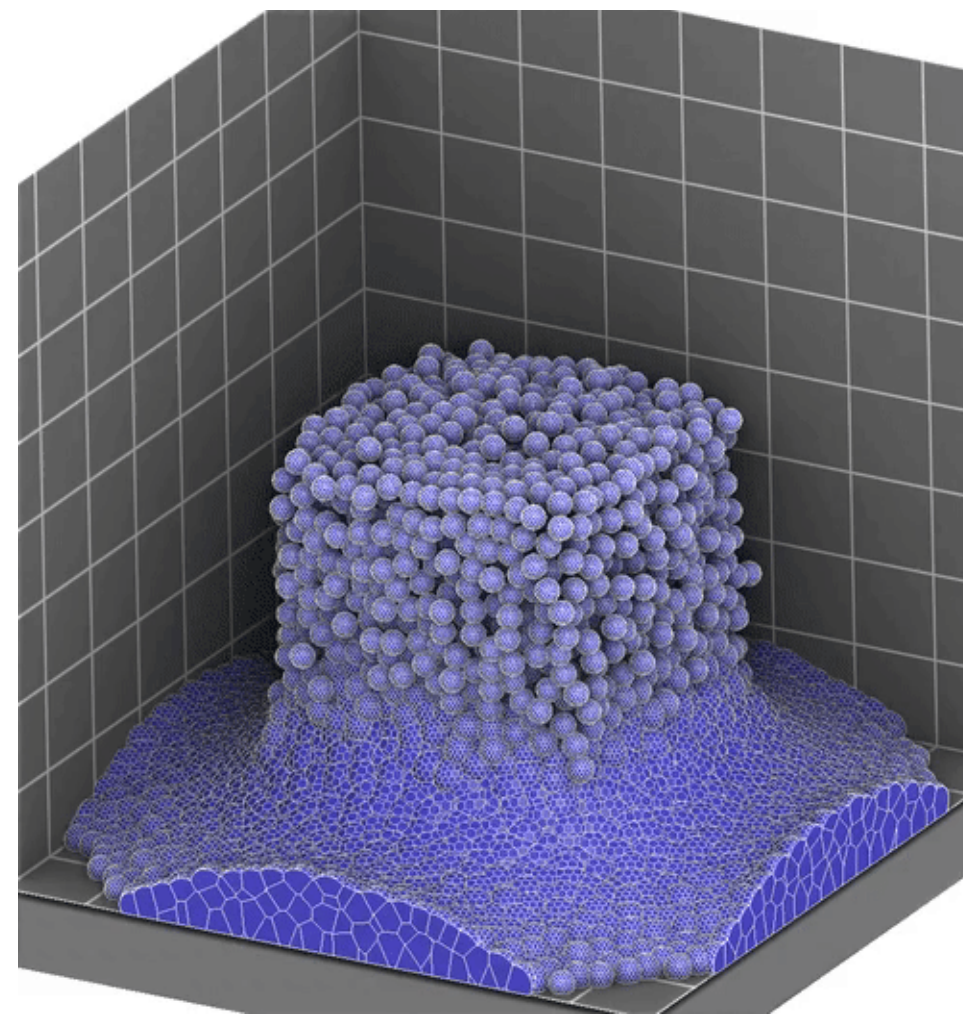


Applications



Simulation of an incompressible bi-phasic flow in a bottle [GWH18]

Free-surface fluid simulation with Gallouet-Merigot scheme [GWH18]



Taylor-Rayleigh instability using the Gallouet-Merigot scheme [GWH18]

Agenda

1- Complexity of OT?

2- Smooth OT?

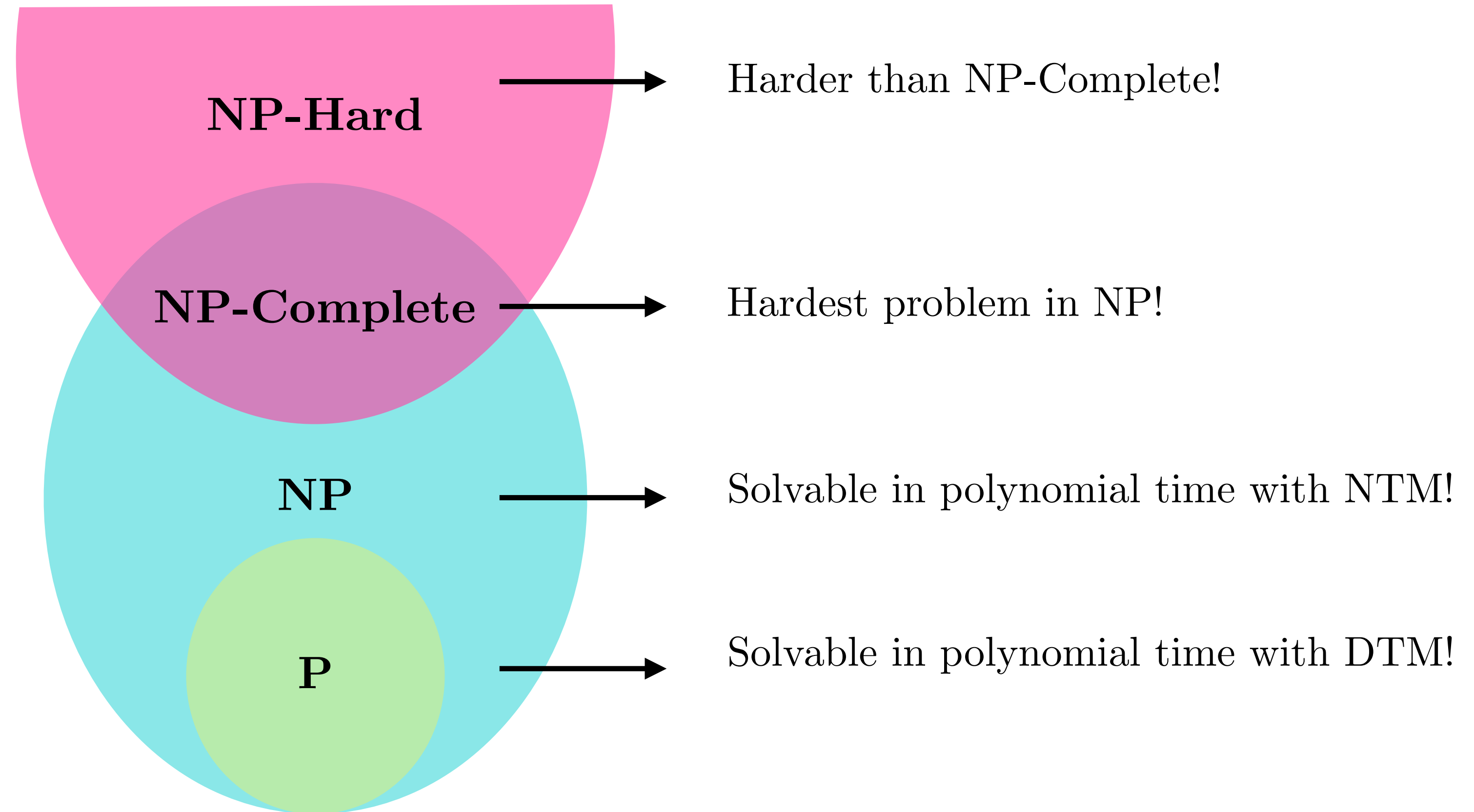
3- Algorithms for OT?

Computational Complexity



Yes

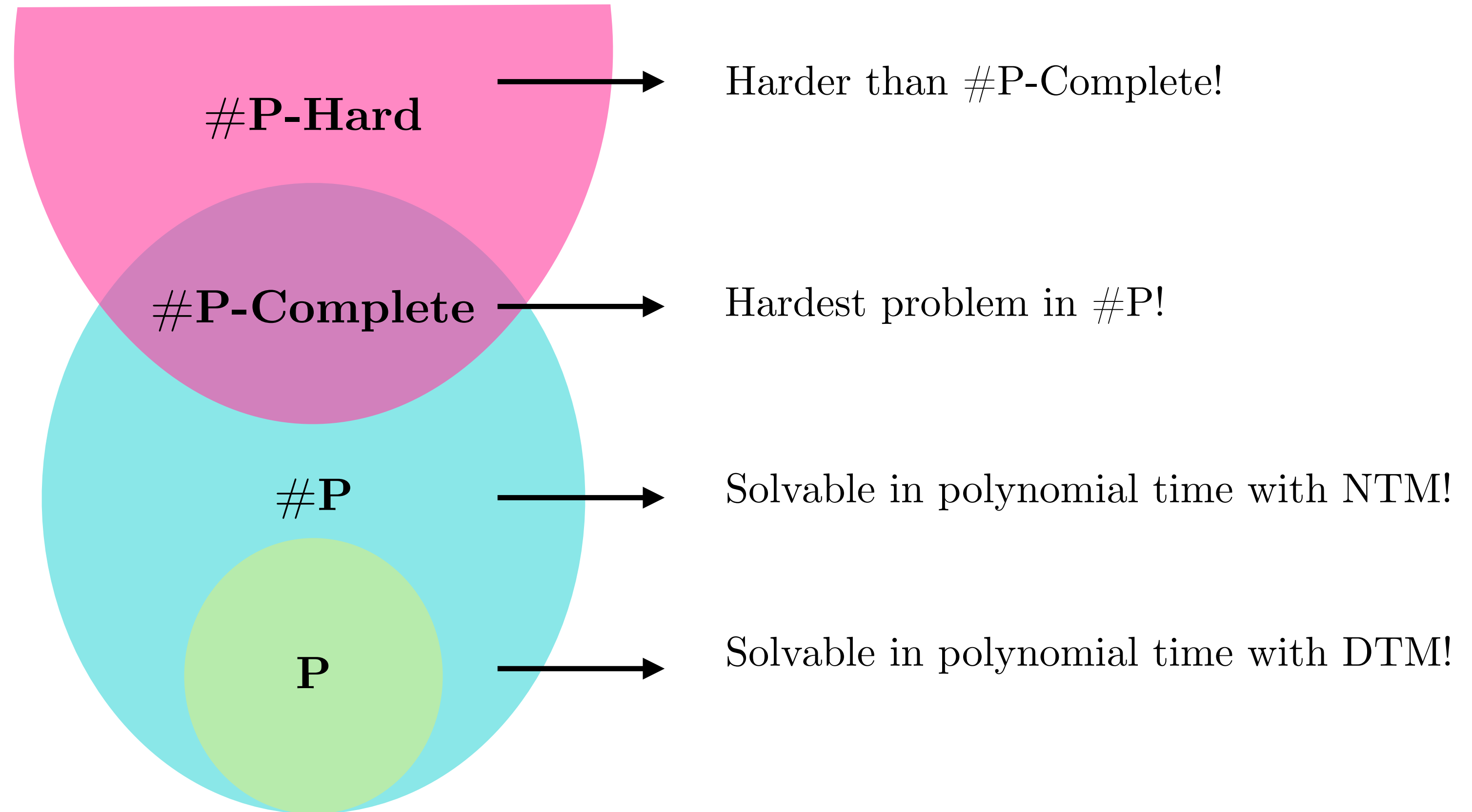
No



Computational Complexity



How many
yes?

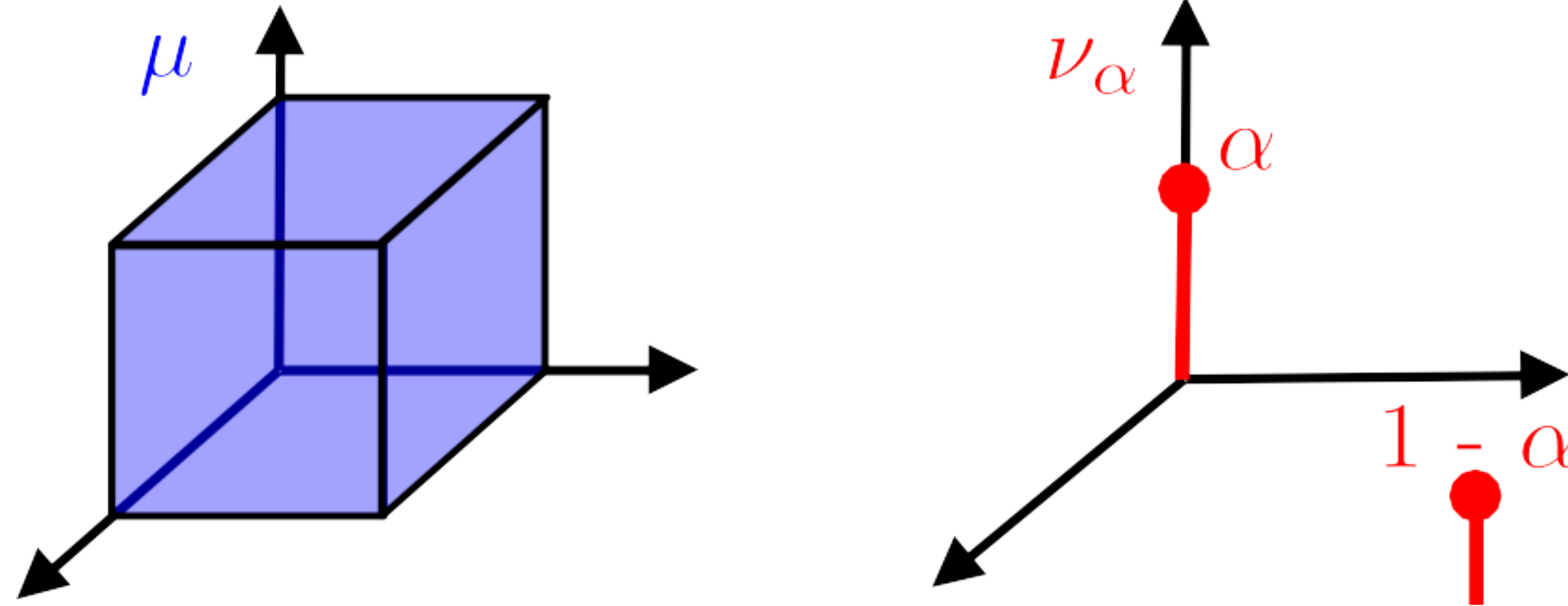


Computational Complexity of OT

Theorem 1. Computing the OT cost is $\#P$ -hard!

Computational Complexity of OT

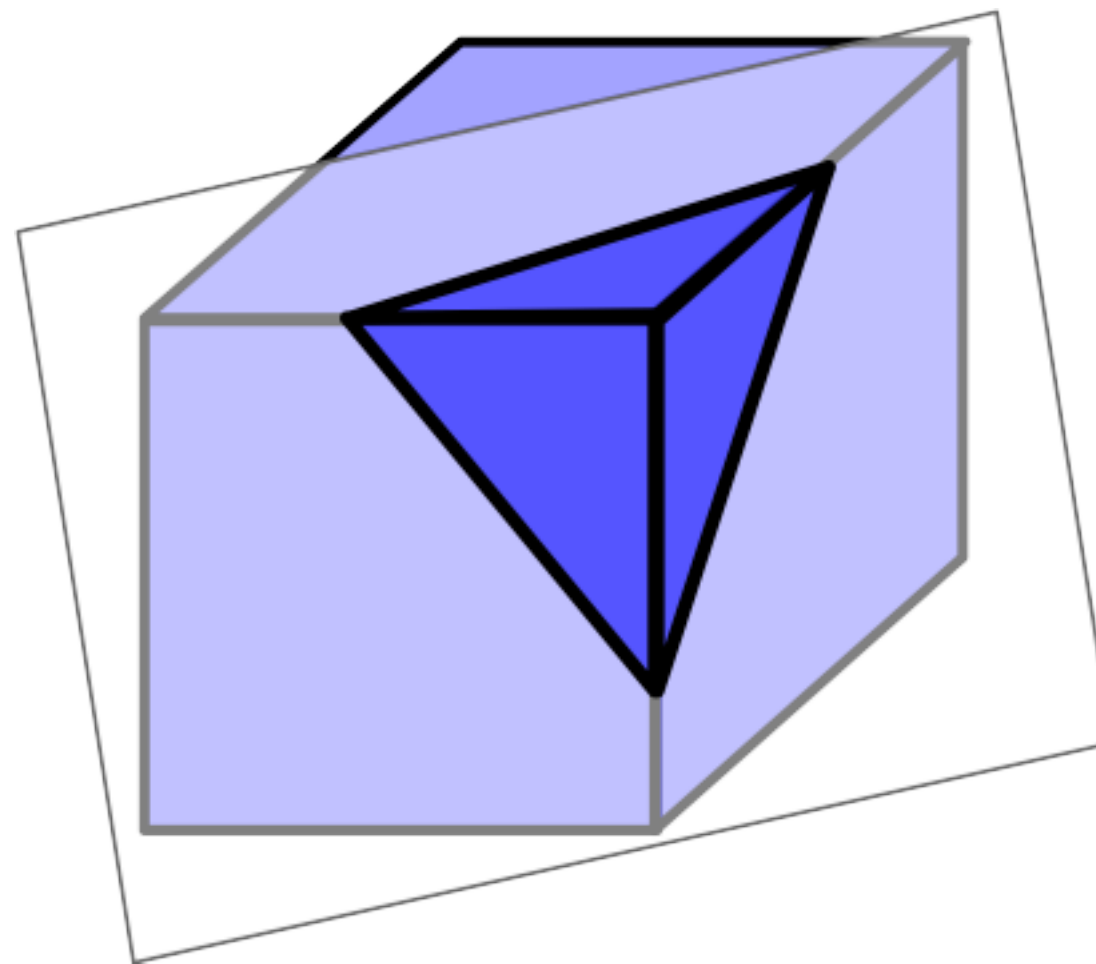
Theorem 1. Computing the OT cost is $\#P$ -hard!



$$\min_{\alpha \in [0,1]} W_c(\mu, \nu_\alpha)$$

Computational Complexity of OT

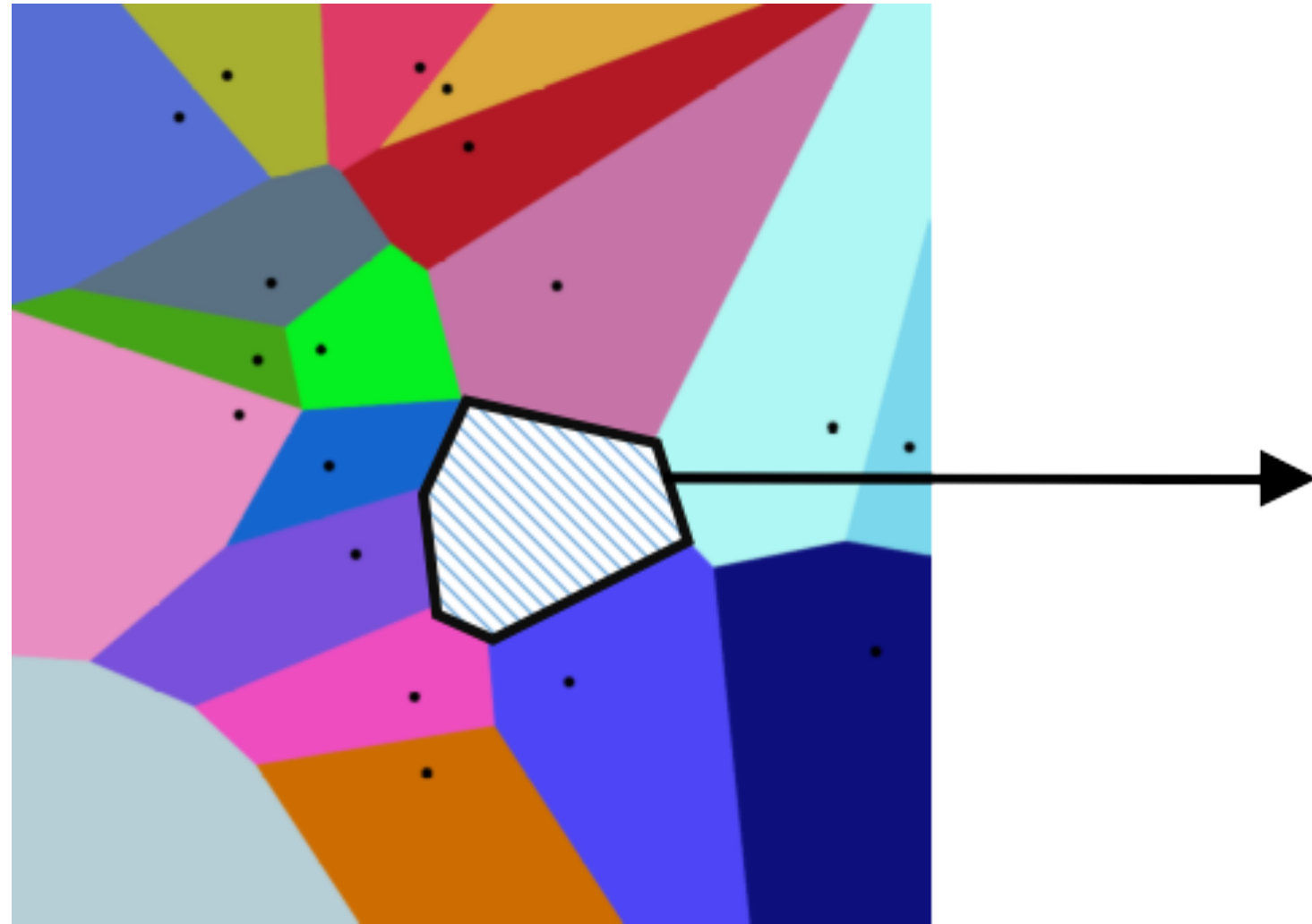
Theorem 1. Computing the OT cost is #P-hard!



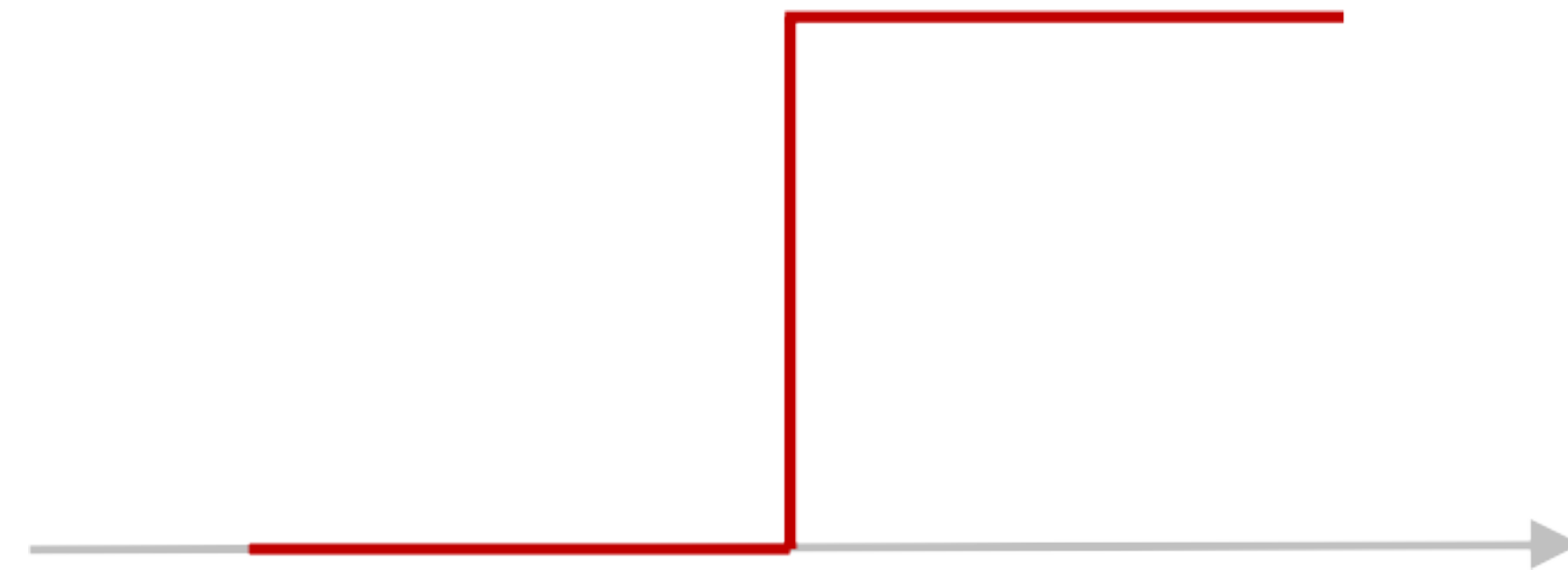
$\operatorname{argmin}_{\alpha \in [0,1]} W_c(\mu, \nu_\alpha) = \text{volume of Knapsack polytope}$

Why is #P-hard?

$$W_c(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{\pi} [c(x, y)]$$

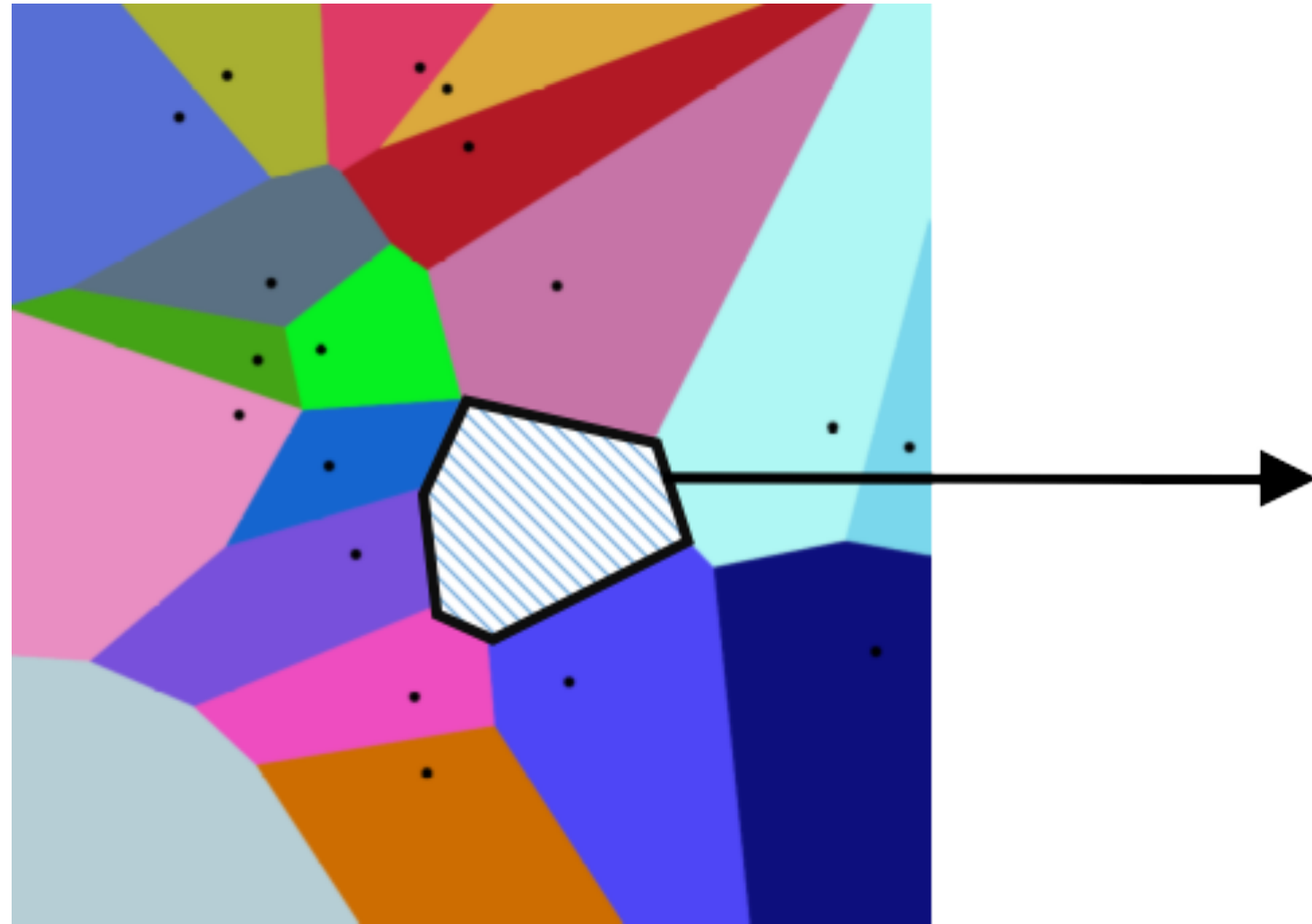


Non-smooth Integrand!

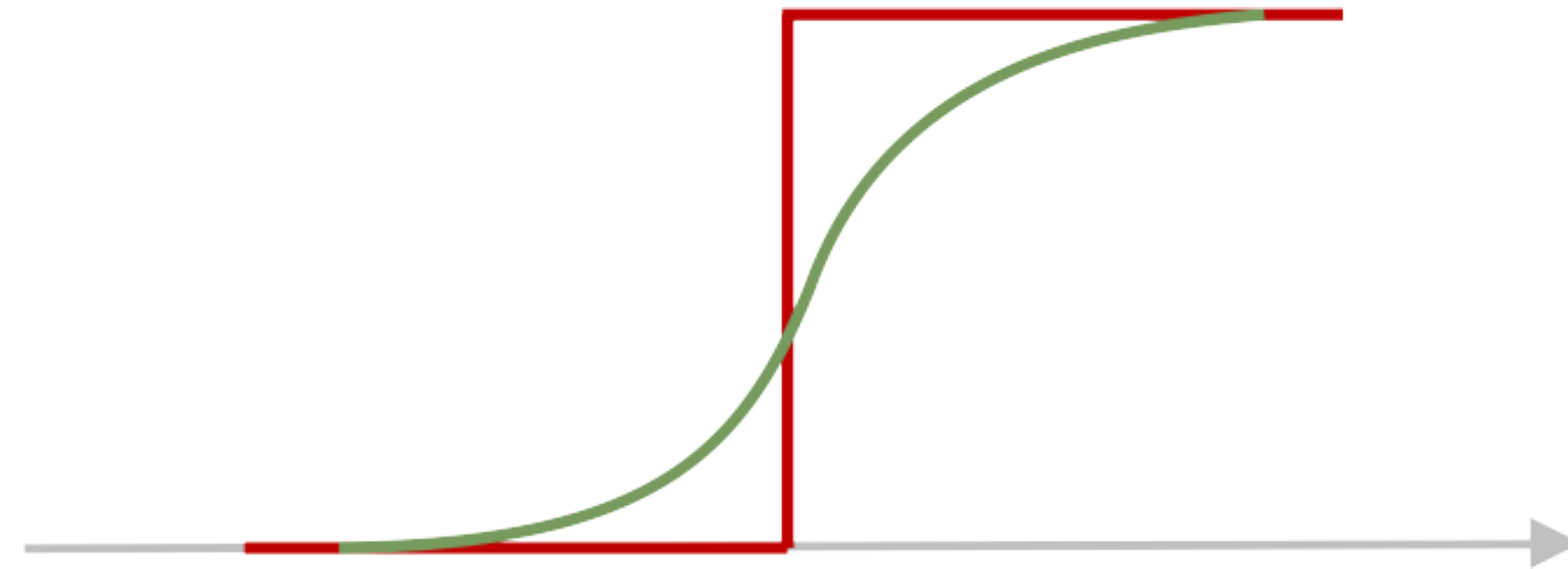


Why is #P-hard?

$$W_c(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_{\pi} [c(x, y)]$$



Smooth approximation!



What is #P-hard?



Input: 1- two distributions representable by n bits
2- Accuracy parameter ϵ

Output: OT distance between two distributions

No algorithm can return the OT cost to within accuracy ϵ in time $\mathcal{O}(\text{poly}(n) \log(1/\epsilon))$

What is #P-hard?



Input: 1- two distributions representable by n bits
2- Accuracy parameter ϵ

Output: OT distance between two distributions

~~No~~ algorithm can return the OT cost to within accuracy ϵ in time $\mathcal{O}(\text{poly}(n) \log(\text{poly}(1/\epsilon)))$

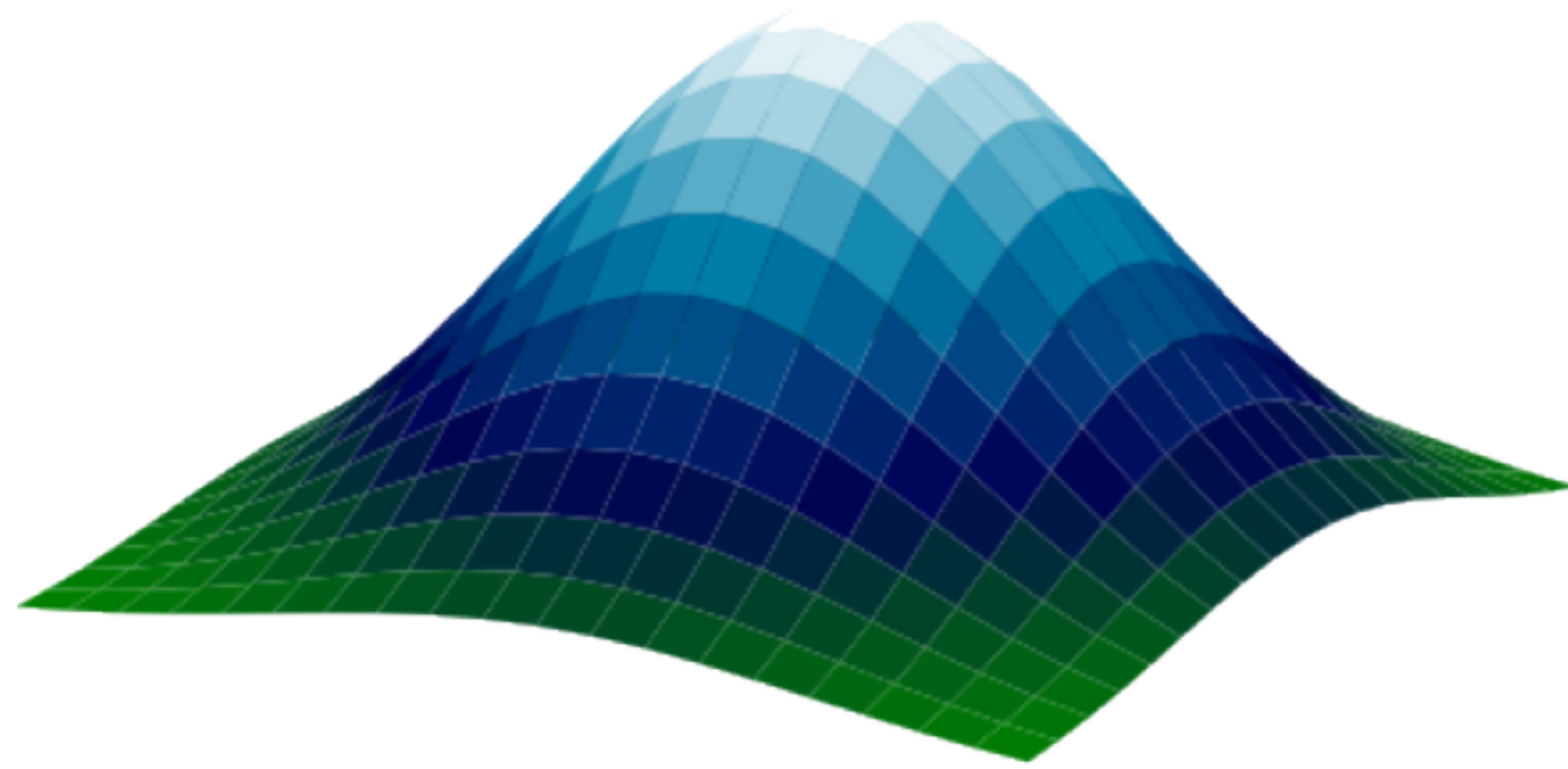
Agenda

1- Complexity of OT? #P-Hard

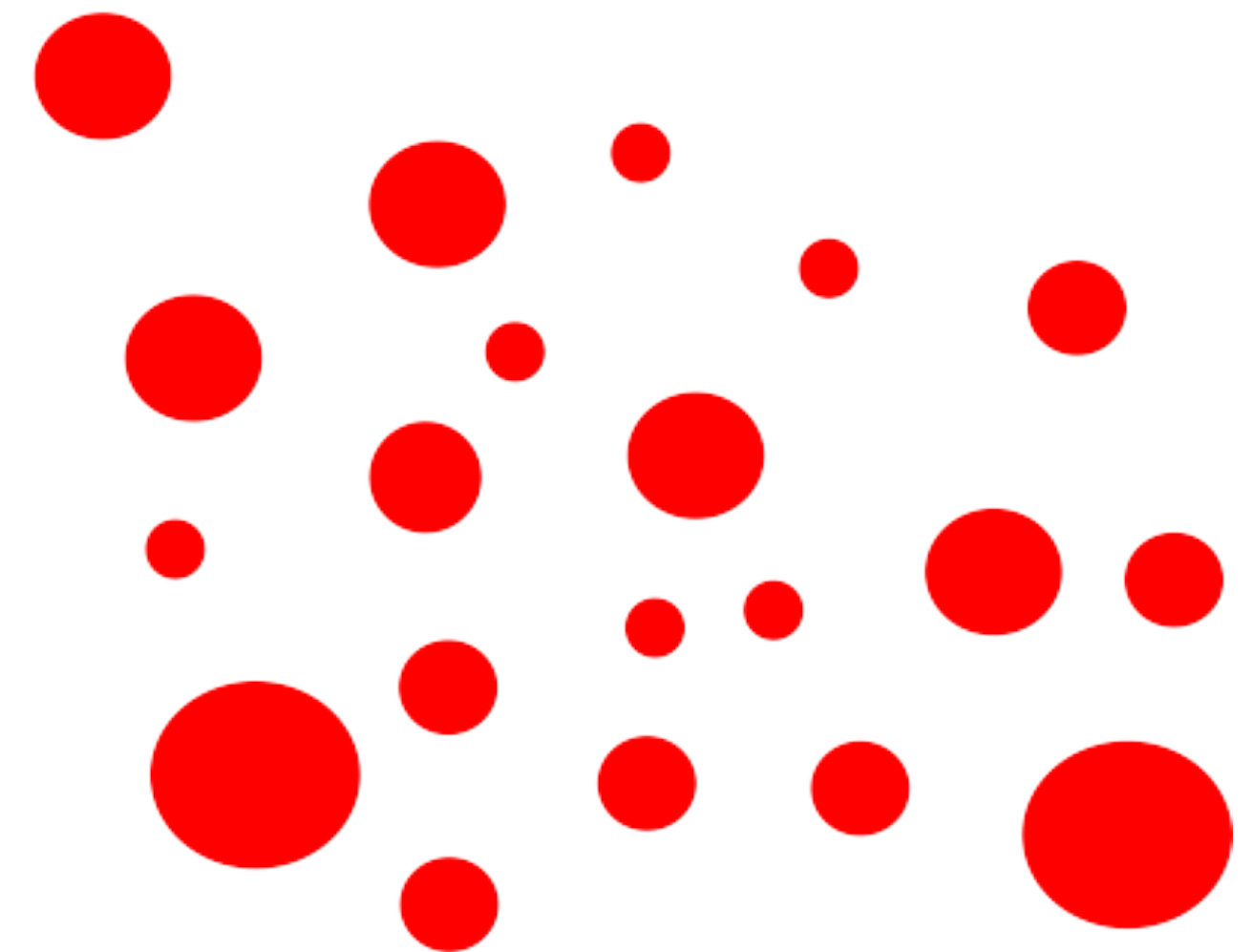
2- Smooth OT?

3- Algorithms for OT?

Semi-discrete case



μ

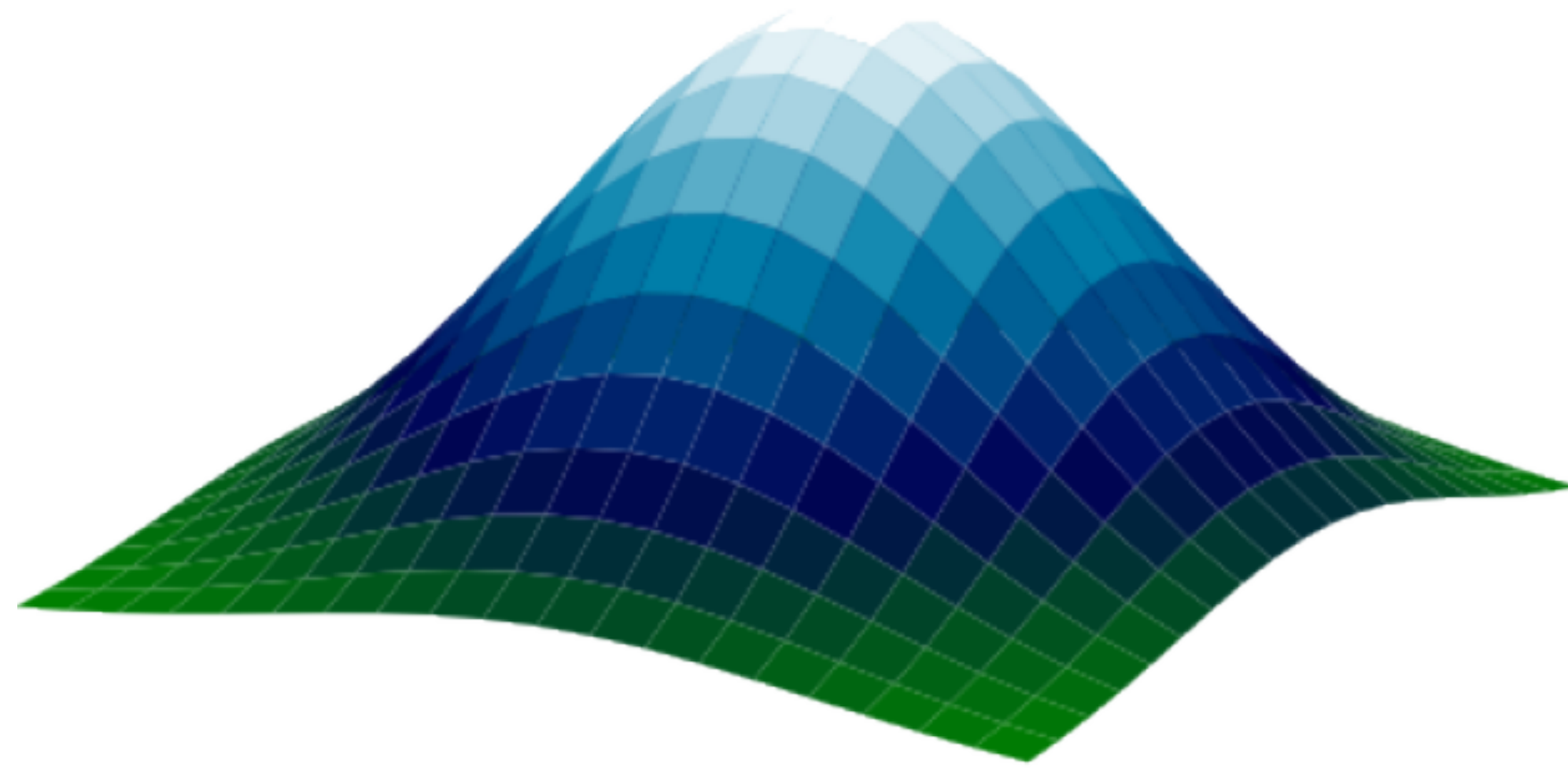


$$\nu = \sum_{i=1}^n \nu_i \delta_{y_i}$$

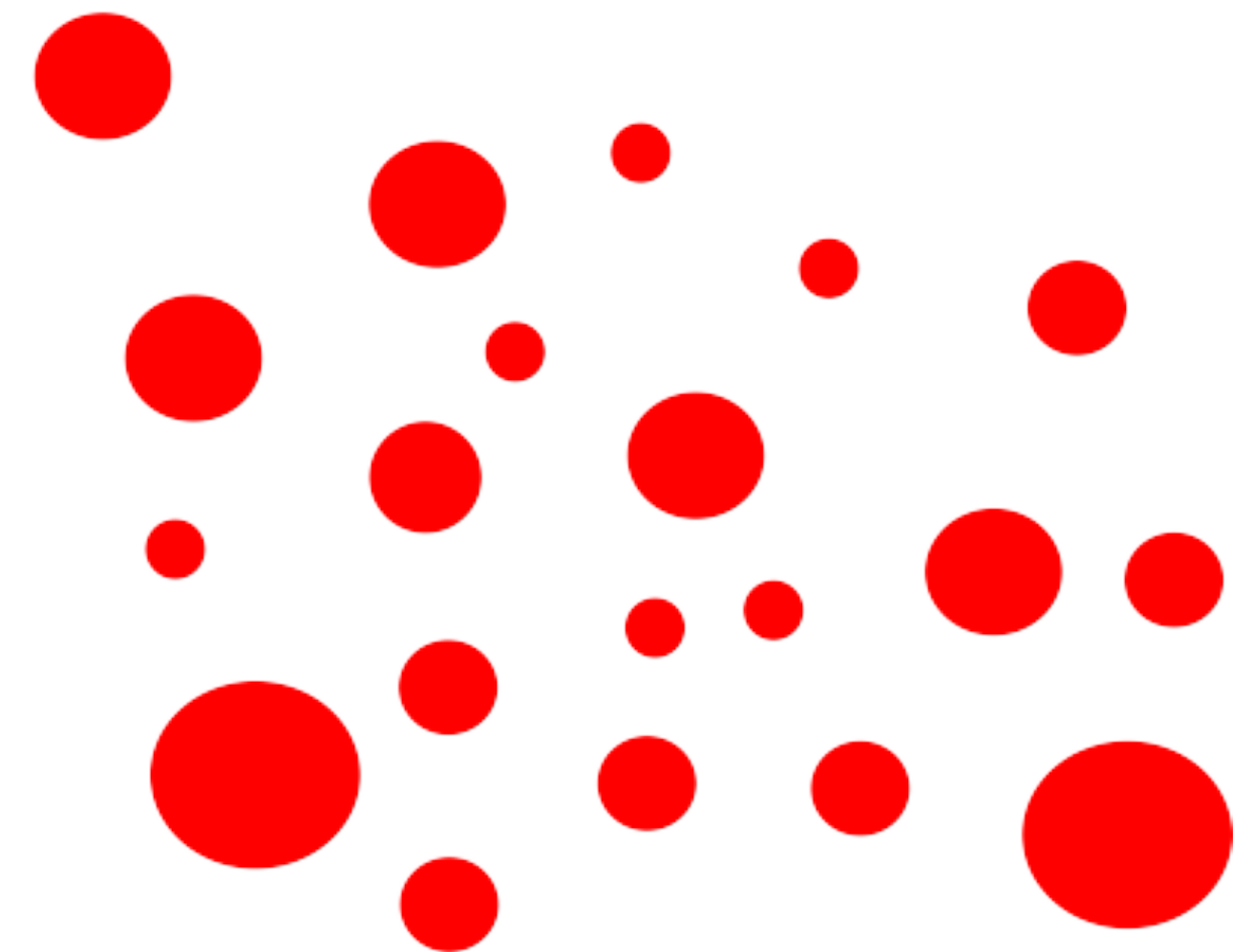
$$W_c(\mu, \nu) = \sup_{\phi \in \mathbb{R}^n} \sum_{i=1}^n \nu_i \phi_i - \mathbb{E}_\mu \left[\max_{i \in [n]} \phi_i - c(x, y_i) \right]$$

c-transform: $\psi_c(\phi, x)$

Semi-discrete case



μ

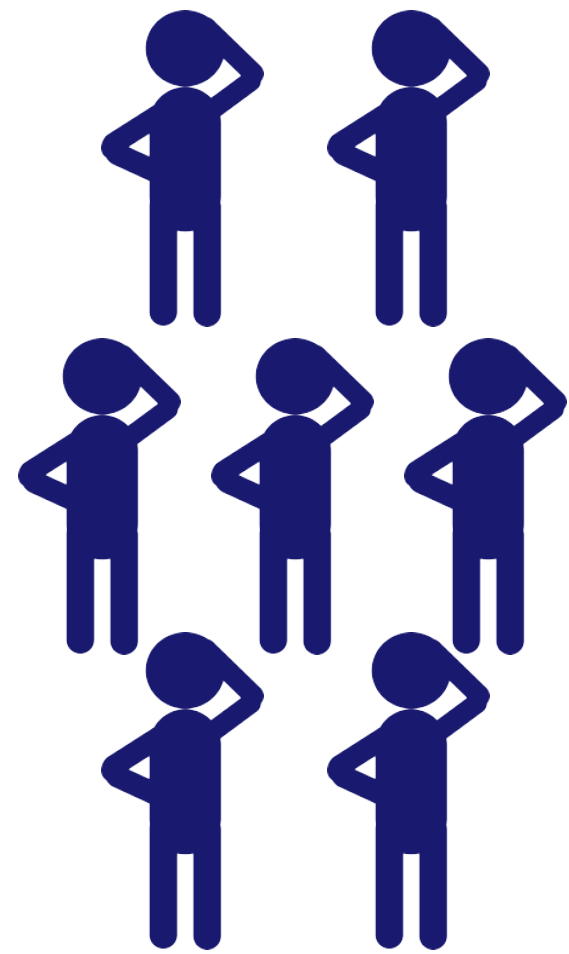


$$\nu = \sum_{i=1}^n \nu_i \delta_{y_i}$$

$$W_c(\mu, \nu) = \sup_{\phi \in \mathbb{R}^n} \sum_{i=1}^n \nu_i \phi_i - \mathbb{E}_{\mu} \left[\max_{i \in [n]} \phi_i - c(x, y_i) \right]$$

smooth c-transform: $\bar{\psi}_c(\phi, x) = \log\left(\sum_{i \in [n]} \exp(\phi_i - c(x, y_i))\right)$

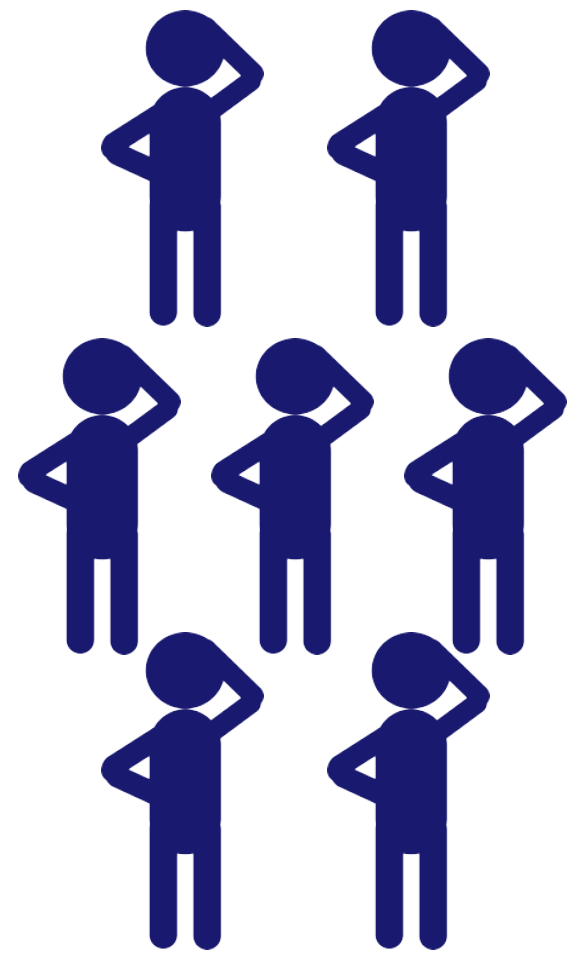
Multinomial Logit Models



$$\mathbb{E}_{z \sim \omega} \left[\max_{i \in [n]} u_i + z_i \right] = \log \left(\sum_{i=1}^n \exp(u_i) \right)$$

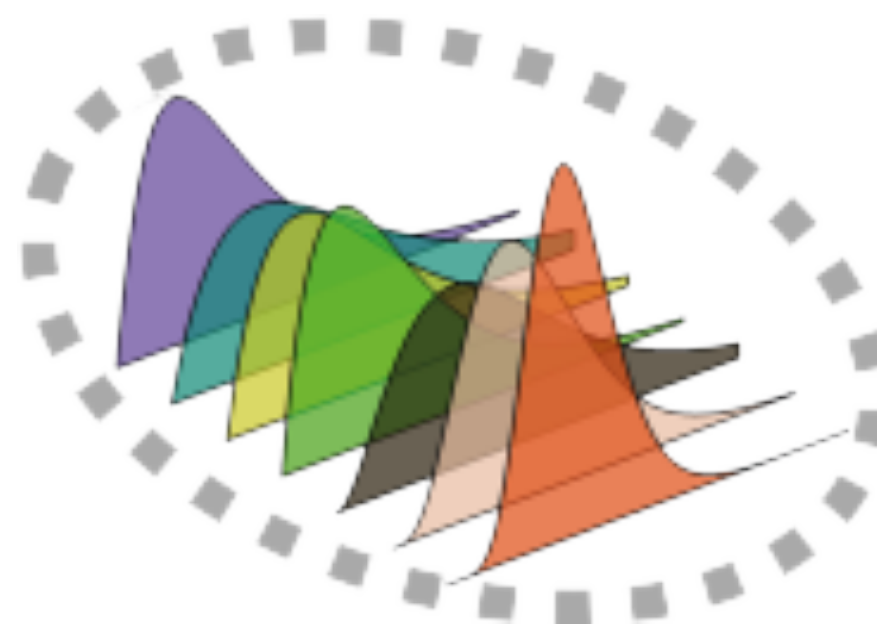
↓
Gumbel distribution

Ambiguous Discrete Choice Models



$$\sup_{\omega \in \Omega} \mathbb{E}_{z \sim \omega} \left[\max_{i \in [n]} u_i + z_i \right]$$

Fréchet Ambiguity Set



$$\Omega = \{ \omega \in \mathcal{M}(\mathbb{R}^n) : \omega(z_i \leq s) = F_i(s) \}$$

Ambiguous Discrete Choice Models



$$\begin{aligned} & \sup_{\omega \in \Omega} \mathbb{E}_{z \sim \omega} \left[\max_{i \in [n]} u_i + z_i \right] \\ &= \max_{p \in \Delta^n} \sum_{i \in [n]} u_i p_i + \int_{1-p_i}^1 F_i^{-1}(s) ds \end{aligned}$$

Smooth OT

- Smooth c-transform

$$\psi_c(\phi, x) = \max_{i \in [n]} \phi_i - c(x, y_i) \longrightarrow \bar{\psi}_c(\phi, x) = \sup_{\omega \in \Omega} \mathbb{E}_{z \sim \omega} \left[\max_{i \in [N]} \phi_i - c(x, y_i) + z_i \right]$$

- Smooth OT

$$\bar{W}_c(\mu, \nu) = \sup_{\phi \in \mathbb{R}^n} \sum_{i=1}^n \nu_i \phi_i - \mathbb{E}_{\mu} [\bar{\psi}_c(\phi, x)]$$

Regularized OT = Smooth OT

Theorem 2. If $\Omega = \{\omega \in \mathcal{M}(\mathbb{R}^n) : \omega(z_i \leq s) = F_i(s)\}$ with

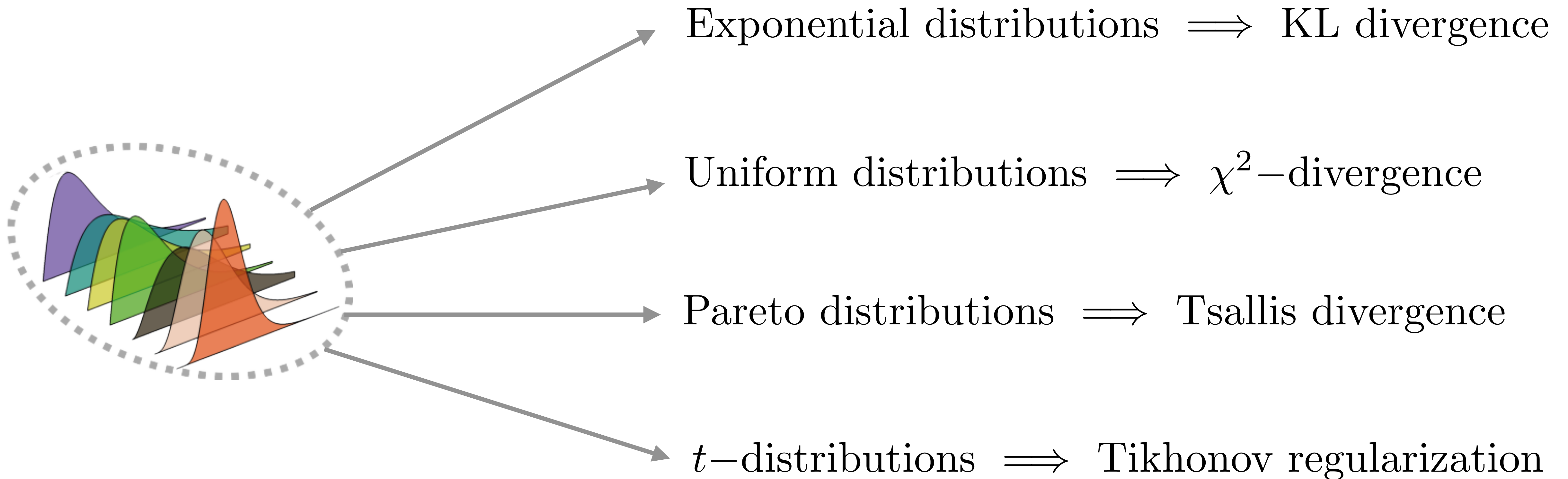
$$F_i(s) = \min\{1, \max\{0, 1 - \nu_i F(-s)\}\}$$

Let $f(s) = \int_0^s F^{-1}(t)dt$. Then we have

$$\overline{W}_c(\mu, \nu) = \inf_{\pi \in \Pi(\mu, \nu)} \mathbb{E}_\pi [c(x, y)] + D_f(\pi \parallel \mu \otimes \nu)$$

$$D_f(\tau \parallel \rho) = \int_{\mathcal{Z}} f(d\tau/d\rho(z))\rho(dz)$$

Examples



Agenda

1- Complexity of OT? #P-Hard

2- Smooth OT = Regularized OT

3- Algorithms for OT?

Averaged Stochastic Gradient Descent



$$\min_{\phi \in \mathbb{R}^n} \left\{ h(\phi) = \mathbb{E}_{\mu} [H(\phi, x)] \right\}$$

Input: $\phi_0 \in \mathbb{R}^n, \eta \in [0, 1]$

for $t = 0, \dots, T - 1$

$$\phi_{t+1} \leftarrow \phi_t - \eta \nabla_{\phi} H(\phi_t, x_t)$$

endfor

Output: average(ϕ_0, \dots, ϕ_T)

$$x \sim \mu \implies \nabla_{\phi} H(\phi, x)$$

Definitions

- **R-Lipschitz:** $|h(\phi) - h(\phi')| \leq R \|\phi - \phi'\|$

- **L-Smooth:** $\|\nabla h(\phi) - \nabla h(\phi')\| \leq L \|\phi - \phi'\|$

- **U-Strongly Convex:** $(\nabla h(\phi) - \nabla h(\phi'))^\top (\phi - \phi') \geq U \|\phi - \phi'\|^2$

- **M-(Generalized) Self-Concordant:** $v(s) = h(\phi + s(\phi' - \phi))$

$$\left| \frac{d^3 v(s)}{ds^3} \right| \leq M \|\phi - \phi'\| \frac{d^2 v(s)}{ds^2}$$

Averaged Stochastic Gradient Descent



$$\min_{\phi \in \mathbb{R}^n} \left\{ h(\phi) = \mathbb{E}_{\mu} [H(\phi, x)] \right\}$$

Input: $\phi_0 \in \mathbb{R}^n, \eta \in [0, 1]$

for $t = 0, \dots, T - 1$

$$\phi_{t+1} \leftarrow \phi_t - \eta \nabla_{\phi} H(\phi_t, x_t)$$

endfor

Output: $\text{average}(\phi_0, \dots, \phi_T)$

$$x \sim \mu \implies \nabla_{\phi} H(\phi, x)$$

Lipschitz: $\mathcal{O}(1/\sqrt{T})$

Smooth: $\mathcal{O}(1/\sqrt{T})$

Strongly Convex: $\mathcal{O}(1/T)$

Self-Concordant: $\mathcal{O}(1/T)$

Averaged Stochastic Gradient Descent for OT



$$\min_{\phi \in \mathbb{R}^n} \left\{ h(\phi) = \mathbb{E}_{\mu} [H(\phi, x)] \right\}$$

$$H(\phi, x) = \bar{\psi}_c(\phi, x) - \sum_{i \in [n]} \nu_i \phi_i$$

$$x \sim \mu \implies \nabla_{\phi} H(\phi, x) = \nabla_{\phi} \bar{\psi}_c(\phi, x) - \nu$$

Structural Results

Theorem 3. If $\Omega = \{\omega \in \mathcal{M}(\mathbb{R}^n) : \omega(z_i \leq s) = F_i(s)\}$, then $\bar{\psi}_c$ is

● **1-Lipschitz:** F_i is continuous

● **L-Smooth:** F_i is L-Lipschitz

● **M-(Generalized) Self-Concordant:**

$$\sup_{s \in F_i^{-1}(0,1)} \frac{|d^2 F_i(s)/ds^2|}{dF_i(s)/ds} \leq M$$

Inexact Averaged Stochastic Gradient Descent



$$\min_{\phi \in \mathbb{R}^n} \left\{ h(\phi) = \mathbb{E}_{\mu} [H(\phi, x)] \right\}$$

$$H(\phi, x) = \bar{\psi}_c(\phi, x) - \sum_{i \in [n]} \nu_i \phi_i$$

$$x \sim \mu \implies \nabla_{\phi} H(\phi, x) = \nabla_{\phi} \bar{\psi}_c(\phi, x) - \nu$$

computing gradients = solving regularized LPs

numerical errors

Inexact Averaged SGD

Theorem 4. Set $\varepsilon_t \leq \mathcal{O}(1/\sqrt{t})$ and $\eta = \mathcal{O}(1/\sqrt{T})$

● **Lipschitz:** $\mathbb{E} \left[h \left(\frac{1}{T} \sum_{t=1}^T \phi_{t-1} \right) \right] - h(\phi^*) \leq \mathcal{O}(1/\sqrt{T})$

● **Smooth:** $\mathbb{E} \left[h \left(\frac{1}{T} \sum_{t=1}^T \phi_t \right) \right] - h(\phi^*) \leq \mathcal{O}(1/T) + \mathcal{O}(1/\sqrt{T})$

● **Self-Concordant:** $\mathbb{E} \left[h \left(\frac{1}{T} \sum_{t=1}^T \phi_{t-1} \right) \right] - h(\phi^*) \leq \mathcal{O}(1/T)$

Convergence Behavior

