

Doc-Start

Doc-Start

Doc-Start

Doc-Start

Decentralized Computation of Wasserstein Barycenters

César A. Uribe

Alexey Kroshnin Darina Dvinskikh Pavel Dvurechensky
Alexander Gasnikov Nazarii Tupitsa Roman Krawtschenko
Ivan Lau Shiqian Ma Dmitry Pasechnyuk Angelia Nedić

May 19, 2025



RICE ENGINEERING

Electrical and Computer Engineering

Outline

- ▶ Distributed Optimization 101
- ▶ Primal-Dual and Decentralized Computations
- ▶ Discrete and Semi-Discrete Barycenters
- ▶ Semi-Discrete Barycenters and Quantization
- ▶ Equitable OT

Distributed Optimization 101

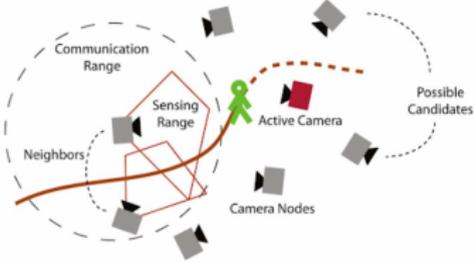
Motivation



(a) Sensor Networks in Agriculture



(b) (Mis)information Spread



(c) Camera Networks for Security



(d) Huge-scale ML

some intersecting points

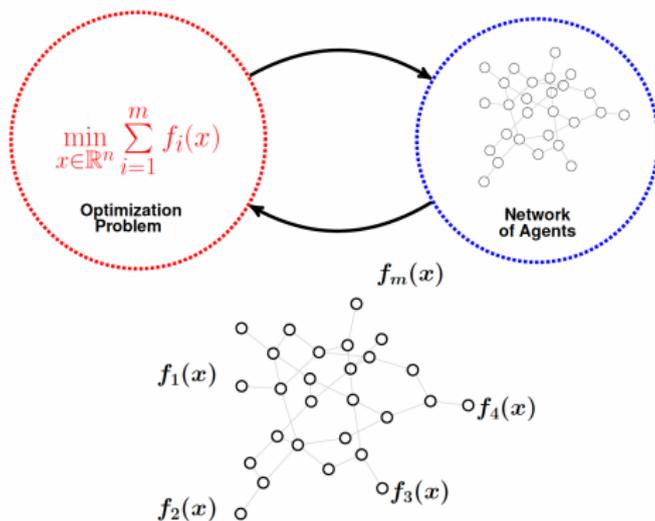
Characteristics

- ▶ Many components/units (we call them agents).
- ▶ Connected over networks.
- ▶ Cyber and Physical interactions.
- ▶ Distributed Storage.

Challenges

- ▶ Decentralization: distributed computations.
- ▶ Scalability: Price of decentralization.
- ▶ Optimality: Efficiency & Performance
- ▶ Robustness & Resiliency: Performance under failures

Distributed Optimization 101: Object of Study



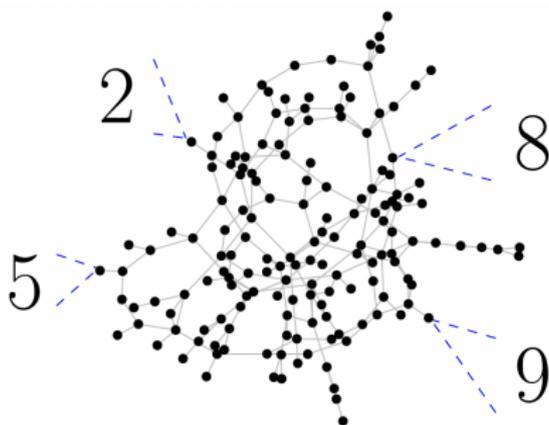
DECENTRALIZED

SCALABLE

OPTIMAL

An example

- ▶ There is a network of m agents, i.e. a graph $G = \{V, E\}$.
- ▶ Agent i holds an initial value $x_0^i \in \mathbb{R}$.
- ▶ Each agent needs to distributedly compute $\frac{1}{m} \sum_{i=1}^m x_0^i$.



Equivalently, solve $\min_{x \in \mathbb{R}} \frac{1}{2} \sum_{i=1}^m \|x - x_i\|_2^2$.

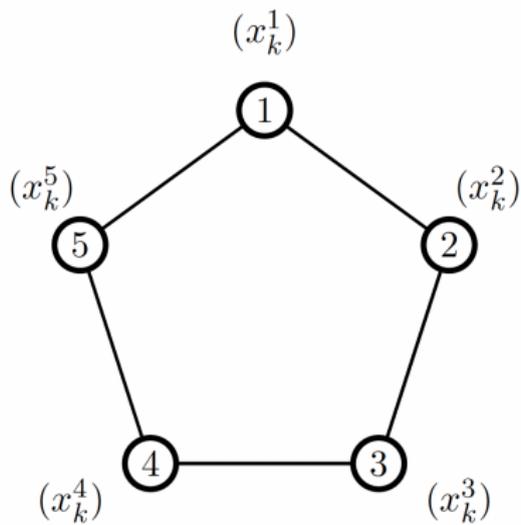
A fundamental result

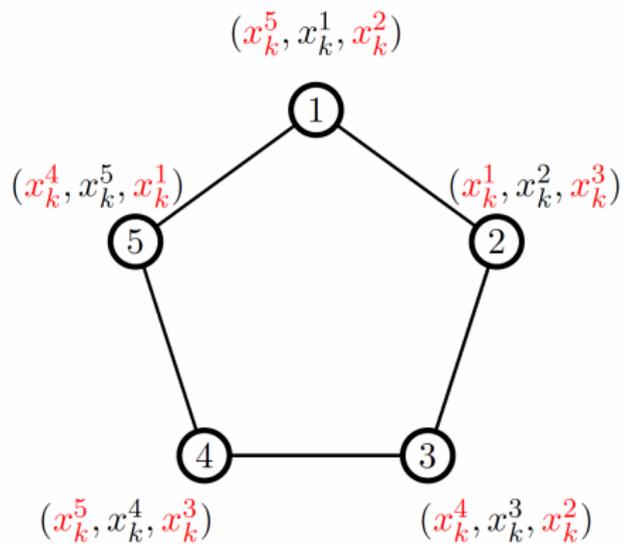
$$x_{k+1}^i = \sum_{j=1}^m [A]_{ij} x_k^j \quad (1)$$

FUNDAMENTAL RESULT:

If \mathcal{G} is connected, undirected, and static, and A is *doubly stochastic* (i.e. $[A]_{ij} > 0$ iff $(j, i) \in E$), then the iterates generated by **(1)** satisfy

$$\lim_{k \rightarrow \infty} x_k^i = \frac{1}{m} \sum_{j=1}^m x_0^j, \quad \forall i \in V.$$

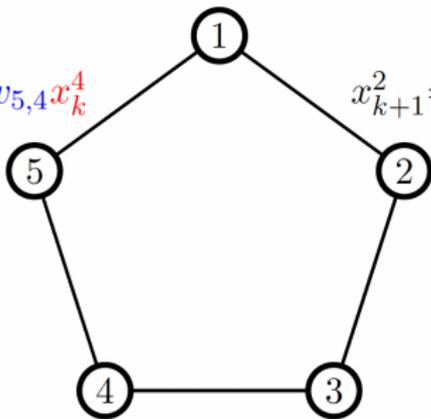




$$x_{k+1}^1 = w_{1,5}x_k^5 + w_{1,1}x_k^1 + w_{1,2}x_k^2$$

$$x_{k+1}^5 = w_{5,1}x_k^1 + w_{5,5}x_k^5 + w_{5,4}x_k^4$$

$$x_{k+1}^2 = w_{2,1}x_k^1 + w_{2,2}x_k^2 + w_{2,3}x_k^3$$



$$x_{k+1}^4 = w_{4,5}x_k^5 + w_{4,4}x_k^4 + w_{4,3}x_k^3$$

$$x_{k+1}^3 = w_{3,4}x_k^4 + w_{3,3}x_k^3 + w_{3,2}x_k^2$$

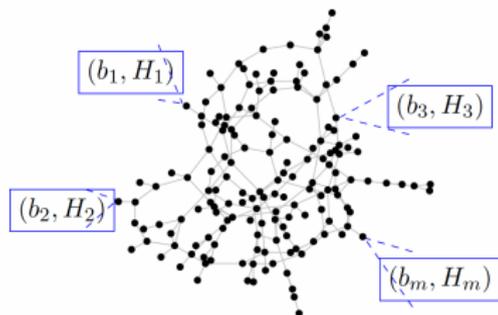
An Example: Distributed Ridge Regression

We want to estimate x assuming

$$b_i = H_i x + \text{noise},$$

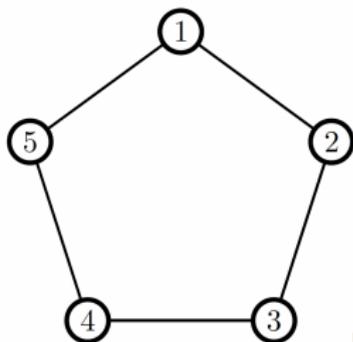
where

- ▶ $H_i \in \mathbb{R}^{d_i \times n}$: d_i data points of dimension n .
- ▶ $b_i \in \mathbb{R}^{d_i}$: d_i outputs.



$$\min_x \frac{1}{2} \frac{1}{m} \sum_{i=1}^m \|b_i - H_i x\|_2^2.$$

A little bit of analysis



$$\bar{W} = \begin{bmatrix} 2 & -1 & 0 & 0 & -1 \\ -1 & 2 & -1 & 0 & 0 \\ 0 & -1 & 2 & -1 & 0 \\ 0 & 0 & -1 & 2 & -1 \\ -1 & 0 & 0 & -1 & 2 \end{bmatrix}.$$

Note that:

- $Wx = 0$ if and only if $x_1 = \dots = x_m$.

$$x = \begin{bmatrix} x_1 \in \mathbb{R}^n \\ x_2 \in \mathbb{R}^n \\ \vdots \\ x_m \in \mathbb{R}^n \end{bmatrix}$$

Rewrite problem (2) in an equivalent form as follows:

$$\min_{x \in \mathbb{R}^n} \sum_{i=1}^m f_i(x) \quad \text{equivalent to} \quad \min_{Wx=0} \sum_{i=1}^m f_i(x_i), \quad (6)$$

$$\text{where } W = \bar{W} \otimes I_n.$$

Initially, consider the general problem

$$\min_{Ax=0} f(x). \quad (7)$$

We assume that the problem has optimal solutions. Later, we will derive the specific results when

$$A = \sqrt{\bar{W}} \quad \text{and} \quad f(x) = \sum_{i=1}^m f_i(x_i).$$

Approximate Solution Definition A point $x \in \mathbb{R}^{mn}$ is said to be an $(\varepsilon, \tilde{\varepsilon})$ -solution of (7) if

$$f(x) - f^* \leq \varepsilon \quad \text{and} \quad \|Ax\|_2 \leq \tilde{\varepsilon},$$

where f^* denotes the optimal value of (7).

The Lagrangian dual for the problem in **(7)** is given by

$$\min_{Ax=0} f(x) = \max_y \left\{ \min_x \{ f(x) - \langle A^T y, x \rangle \} \right\},$$

or equivalently

$$\min_y \varphi(y) \quad \text{where} \quad \varphi(y) \triangleq \max_x \{ \langle A^T y, x \rangle - f(x) \},$$

$$\text{with} \quad \nabla \varphi(y) = Ax^*(A^T y) \quad (\text{Demyanov–Danskin})$$

$$x^*(A^T y) = \arg \max_x \{ \langle A^T y, x \rangle - f(x) \}.$$

We say that f is **dual friendly** when we can determine a solution of the preceding problem efficiently (in a closed form ideally).

- ▶ $f(x)$ is μ -strongly convex $\iff \varphi(y)$ is L_φ -smooth with $L_\varphi = \lambda_{\max}(A^T A) / \mu$.
- ▶ $f(x)$ is L -smooth $\iff \varphi(y)$ is μ_φ -strongly convex on $\text{range}(A)$ with $\mu_\varphi = \lambda_{\min}^+(A^T A) / L$.

The dual problem $\min_y \varphi(y)$ may have multiple solutions of the form $y^* + \ker(A^T)$.

Informally: If $f(x)$ has condition number L/μ , then $\varphi(y)$ has condition number

$$\frac{\lambda_{\max}(A^T A)}{\lambda_{\min}^+(A^T A)} \frac{L}{\mu}.$$

Let's recall a version of Nesterov's fast gradient method for

$$\min_y \varphi(y)$$

$$y_{k+1} = \tilde{y}_k - \frac{1}{L_\varphi} \nabla \varphi(\tilde{y}_k),$$

$$\tilde{y}_{k+1} = y_{k+1} + \frac{\sqrt{L_\varphi} - \sqrt{\mu_\varphi}}{\sqrt{L_\varphi} + \sqrt{\mu_\varphi}} (y_{k+1} - y_k).$$

and

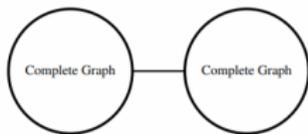
$$\varphi(y_k) - \varphi^* \leq L_\varphi \left(1 - \sqrt{\frac{\mu_\varphi}{L_\varphi}}\right)^k \|y_0 - y^*\|_2^2,$$

Set $A = \sqrt{\bar{W}}$, $z_k = \sqrt{\bar{W}} y_k$, and $\tilde{z}_k = \sqrt{\bar{W}} \tilde{y}_k$

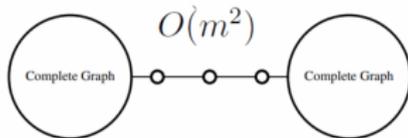
$$x_i^*(\tilde{z}_k^i) = \arg \max_{x_i} \{ \langle \tilde{z}_k^i, x_i \rangle - f_i(x_i) \}$$

$$z_{k+1}^i = \tilde{z}_k^i - \frac{\mu}{\lambda_{\max}(W)} \sum_{j=1}^m W_{ij} x_j^*(\tilde{z}_k^j)$$

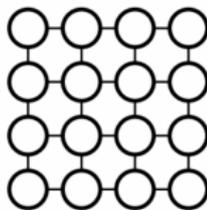
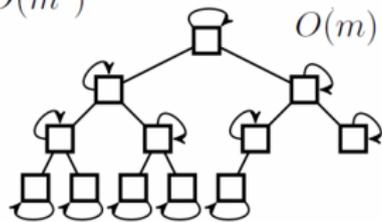
$$\tilde{z}_{k+1}^i = z_{k+1}^i + \frac{\sqrt{\lambda_{\max}(W)/\mu} - \sqrt{\lambda_{\min}^+(W)/L}}{\sqrt{\lambda_{\max}(W)/\mu} + \sqrt{\lambda_{\min}^+(W)/L}} (z_{k+1}^i - z_k^i).$$



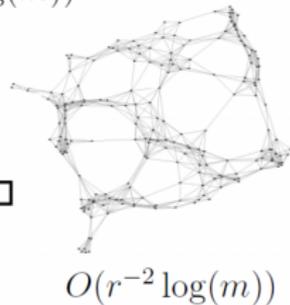
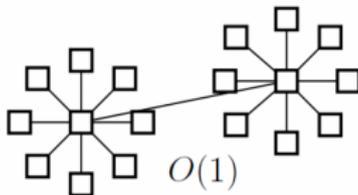
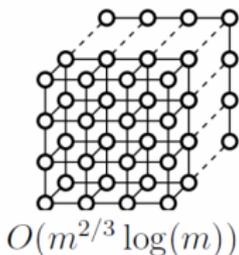
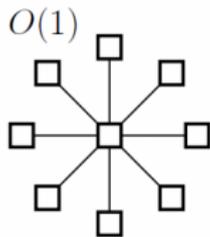
$O(m^2)$



$O(\log^2(m))$



$O(m \log(m))$



Decentralized Computations of Discrete Wasserstein Barycenters

Why Decentralized Wasserstein Barycenters?

- ▶ **Distributed-data reality**

Probability measures reside on edge devices (IoT sensors, mobile phones, robots, medical sites); raw data aggregation is infeasible or prohibited.

- ▶ **Geometry-preserving averaging**

Wasserstein barycenters provide the correct *geometric mean* of distributions, enabling domain adaptation, sensor-fusion maps, and manifold-aware clustering.

- ▶ **Privacy & compliance**

Local measures may contain sensitive content (e.g. medical images); decentralized protocols keep data in place, sharing only aggregated or quantized statistics.

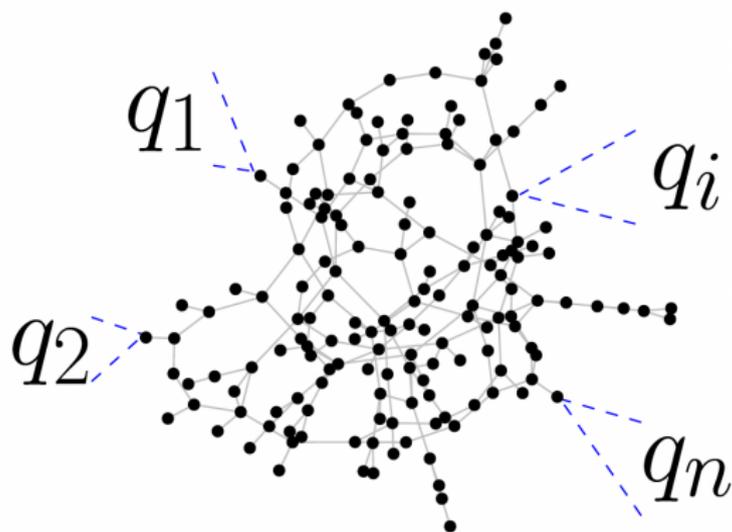
- ▶ **Scalability & fault-tolerance**

Centralized Sinkhorn/IBP requires $O(mn^2)$ memory; distributing both storage *and* computation removes single points of failure and scales to thousands of agents.

- ▶ **Long-term vision**

Provably-optimal decentralized barycenters underpin federated OT learning, real-time cooperative perception, and privacy-preserving generative modelling.

What if the data are probability distributions



Equivalently, solve $\min_{q \in \mathbb{P}} \frac{1}{2} \sum_{i=1}^m d^2(p, q_i)$.

Entropy Regularized Discrete Optimal Transport (OT)

Given $p, q \in \Delta_n$ and cost $C \in \mathbb{R}_{\geq 0}^{n \times n}$,

$$W(p, q) = \min_{\pi \in \Pi(p, q)} \langle \pi, C \rangle, \quad \Pi(p, q) = \{\pi \geq 0 : \pi \mathbf{1} = p, \pi^T \mathbf{1} = q\}.$$

Cuturi (2013) solves large OT by adding entropy:

$$W_\gamma(p, q) = \min_{\pi \in \Pi(p, q)} \langle \pi, C \rangle + \gamma \text{KL}(\pi \| \mathbf{1} \mathbf{1}^T).$$

Gains:

- ▶ Strong convexity \Rightarrow unique optimal plan π^γ .
- ▶ Sinkhorn scaling: $O(n^2)$ per iteration with linear rate $1 - \Theta(\gamma)$.

Wasserstein Barycenter

For $\{p_\ell\}_{\ell=1}^m \subset \Delta_n$ and weights $w \in \Delta_m$,

$$q^* = \arg \min_{q \in \Delta_n} \sum_{\ell=1}^m w_\ell W_\gamma(p_\ell, q).$$

Our objective: compute q^* in a network where p_ℓ reside on distinct agents.

Dual Representation — Theorem 1

Theorem 1: (Theorem 2.4 in (Cuturi et al., 2016))

For $\gamma > 0$, the Fenchel–Legendre dual function $W_{\gamma,q}^*(y)$ is differentiable and its gradient is $\frac{1}{\gamma}$ -Lipschitz in the 2-norm with

$$W_{\gamma,q}^*(y) = \gamma(E(q) + \langle q, \log K\alpha \rangle), \nabla W_{\gamma,q}^*(y) = \alpha \circ \frac{Kq}{K\alpha} \in S_1(n),$$

where $y \in \mathbb{R}^n$, $\alpha = \exp(y/\gamma)$, and $K = \exp(-M/\gamma)$.

Centralized vs. Decentralized Formulation

Centralized

Full data at one server; complexity dominated by OT sub-problems $O(mn^2)$ each iteration.

Decentralized

$$\min_{\substack{q_1 = \dots = q_m \\ q_i \in \Delta_n}} \sum_{i=1}^m W_\gamma(p_i, q_i).$$

Consensus enforced by graph constraints; complexity coupled with spectral gap $\lambda_{\min}^+(W)$.

Convergence Guarantee

Theorem 2: Distributed Computation of WB

Let $\varepsilon > 0$ and $\|\nabla \mathcal{W}_{\gamma,q}^*(y)\|_2 \leq G$ for $y \in B_R(0)$ with $R = \|y^*\|_2$. Then, for

$$N \geq \sqrt{\frac{16G^2 \lambda_{\max}(W)}{\gamma \cdot \varepsilon \lambda_{\min}^+(W)}},$$

$p_N = [(p_N)_1^T, \dots, (p_N)_m^T]^T$ and $y_N = [(y_N)_1^T, \dots, (y_N)_m^T]^T$ have the following properties:

$$\mathcal{W}_{\gamma,q}(p_N) + \mathcal{W}_{\gamma,q}^*(y_N) \leq \varepsilon \quad \text{and} \quad \|\sqrt{W}p_N\|_2 \leq \varepsilon/R.$$

Lemma 8 — Radius of the Optimal Dual

Lemma 1

Let q_γ^* be the optimal solution of problem (4) with minimal 2-norm. Then there exists an optimal dual solution

$$\lambda^* = [\lambda_1^*, \dots, \lambda_m^*]$$

for problem (29) satisfying

$$\|\lambda^*\|_2 \leq R, \quad R^2 = \frac{2n \sum_{l=1}^m w_l^2 \|C_l\|_\infty^2}{\lambda_{\min}^+(W)}.$$

Here $\lambda_{\min}^+(W)$ is the minimal positive eigenvalue of the matrix W .

AGD Complexity for Barycenter

Theorem 3: Distributed Computation of WB

The Primal-Dual Accelerated Gradient Based method after

$$N = \frac{1}{\varepsilon} \sqrt{64 \chi(\bar{W}) mn \ln n \sum_{l=1}^m w_l^2 \|C_l\|_\infty^2}$$

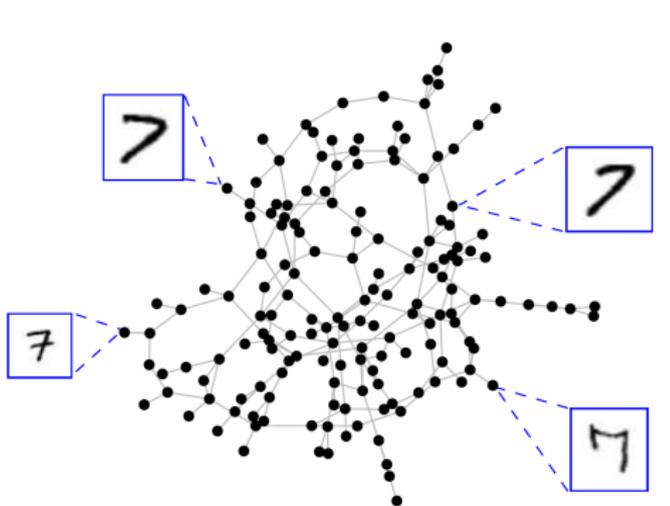
iterations generates an ε -solution of problem (2), i.e. finds a vector $q_N = [q_1^T, \dots, q_m^T]^T$ such that

$$\sum_{l=1}^m w_l W(p_l, q_{N,l}) - \sum_{l=1}^m w_l W(p_l, q^*) \leq \varepsilon, \quad \|\sqrt{W} q_N\|_2 \leq \varepsilon/(2R),$$

where q^* is an unregularized barycenter, and R is a bound on the solution to the dual problem. Moreover, the number of arithmetic operations is

$$O(N n (mn + \text{nnz}(\bar{W})) / \varepsilon).$$

Proof by Picture



Video 1

Decentralized Computation of Semi-Discrete Wasserstein Barycenters

Semi-Discrete Regularized OT

$$W_\gamma(\mu, \nu) = \min_{\pi \in \Pi(\mu, \nu)} \left\{ \sum_{i=1}^n \int_Y c_i(y) \pi_i(y) dy + \gamma \sum_{i=1}^n \int_Y \pi_i(y) \log\left(\frac{\pi_i(y)}{\xi}\right) dy \right\}$$

$$\Pi(\mu, \nu) = \left\{ \pi \in \mathcal{M}_+^1(Y) \times \mathcal{S}_1(n) : \sum_{i=1}^n \pi_i(y) = q(y), \int_Y \pi_i(y) dy = p_i \right\}$$

$$\min_{p \in \mathcal{S}_1(n)} \sum_{i=1}^m W_{\gamma, \mu_i}(p)$$

Dual Formulation & Gradient

$$\min_{\lambda \in \mathbb{R}^{mn}} \sum_{i=1}^m W_{\gamma, \mu_i}^*([\sqrt{W} \lambda]_i)$$

$$W_{\gamma, \mu}^*(\bar{\lambda}) = \mathbb{E}_{Y \sim \mu} \gamma \log\left(\frac{1}{q(Y)} \sum_{\ell=1}^n \exp\left(\frac{\bar{\lambda}_\ell - c_\ell(Y)}{\gamma}\right)\right)$$

$$(\nabla W_{\gamma, \mu}^*(\bar{\lambda}))_l = \mathbb{E}_{Y \sim \mu} \frac{\exp((\bar{\lambda}_l - c_l(Y))/\gamma)}{\sum_{\ell=1}^n \exp((\bar{\lambda}_\ell - c_\ell(Y))/\gamma)}, \quad l = 1, \dots, n$$

Smoothness & Stochastic Gradient Bounds

Lemma 1 (smoothness). For every μ , $W_{\gamma,\mu}^*$ is $\frac{1}{\gamma}$ -smooth.

Lemma 2. The gradient of the dual objective satisfies

$$\|\nabla W_{\gamma}^*(\lambda_1) - \nabla W_{\gamma}^*(\lambda_2)\|_2 \leq \frac{\lambda_{\max}(W)}{\gamma} \|\lambda_1 - \lambda_2\|_2,$$

and the mini-batch estimator with size M obeys

$$\mathbb{E} \|\tilde{\nabla} W_{\gamma}^*(\lambda) - \nabla W_{\gamma}^*(\lambda)\|_2^2 \leq \frac{\lambda_{\max}(W) m}{M}.$$

Convergence of Algorithm 4

Theorem 4: Distributed Computation of the Semi-Discrete WB

Let $\|\lambda^*\|_2 \leq R$. After

$$N = \sqrt{\frac{32 \lambda_{\max}(W) R^2}{\varepsilon \gamma}}$$

iterations, Algorithm 4 returns $\hat{\rho}_N$ such that

$$\sum_{i=1}^m W_{\gamma, \mu_i}(\mathbb{E}[\hat{\rho}_N]_i) - \sum_{i=1}^m W_{\gamma, \mu_i}(p_i^*) \leq \varepsilon, \quad \|\sqrt{W} \mathbb{E} \hat{\rho}_N\|_2 \leq \varepsilon/R,$$

with total arithmetic cost

$$O\left(mn \max\left\{\sqrt{\frac{\lambda_{\max}(W) R^2}{\varepsilon \gamma}}, \frac{\lambda_{\max}(W) m R^2}{\varepsilon^2}\right\}\right).$$

Proof by Picture 2

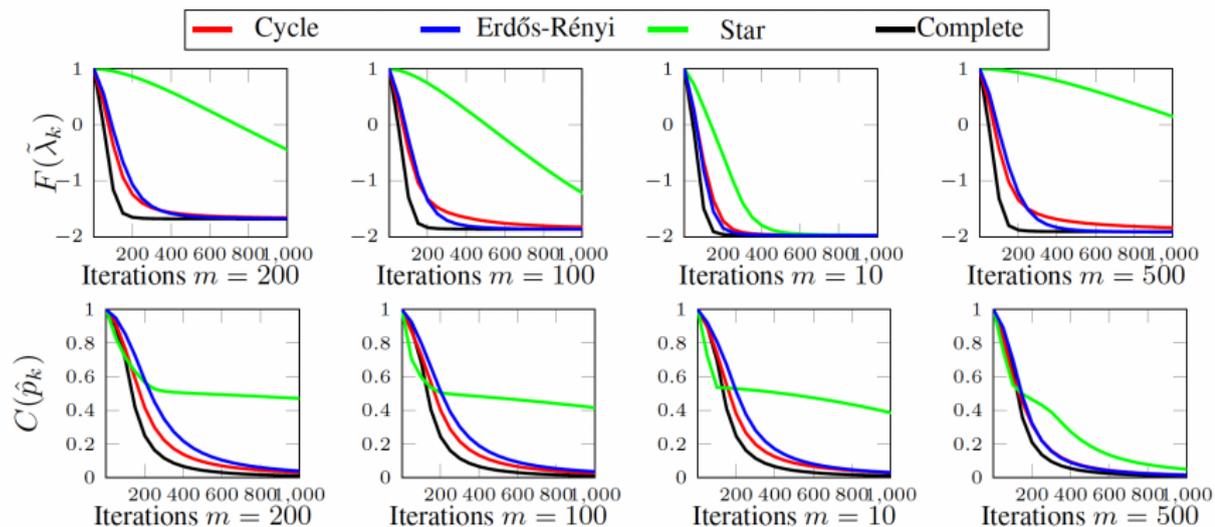
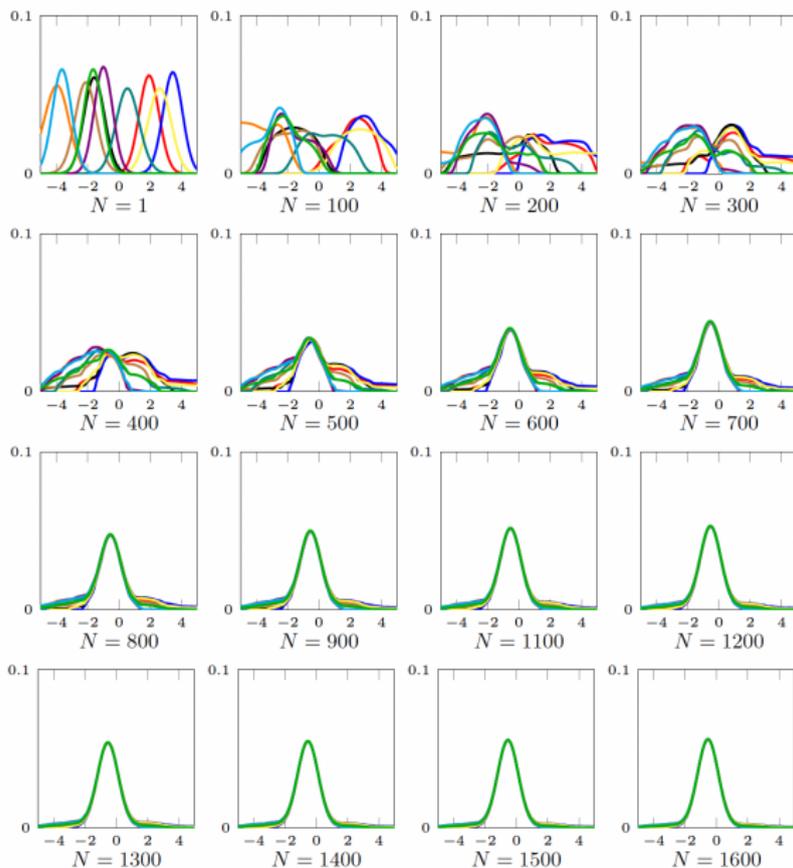


Figure 1: Dual function value and distance to consensus for 200, 100, 10, 500 agents, $M_k = 100$ and $\gamma = 0.1$.

Discrete Bounded Support Gaussian-Like



IBP: High-Accuracy Baseline

Core idea: Alternating KL projections on marginal constraints leveraging Sinkhorn kernel $K_\ell = \exp(-C_\ell/\gamma)$.

- ▶ Each outer iteration performs two Sinkhorn updates per distribution.
- ▶ Communication: vector of size n per neighbor (scalable).

Iterative Bregman Projections

$$\min_{\substack{\pi_l \mathbb{1} = p_l, \pi_l^T \mathbb{1} = \pi_{l+1}^T \mathbb{1} \\ \pi_l \in \mathbf{R}_+^{n \times n}, l=1, \dots, m}} \frac{1}{m} \sum_{l=1}^m \{ \pi_l, C_l + \gamma H(\pi_l) \}$$

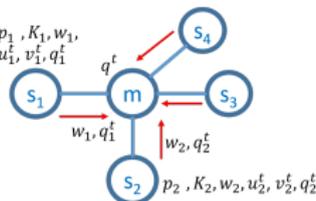
Dual problem:

$$\min_{\substack{\mathbf{u}, \mathbf{v} \\ \frac{1}{m} \sum_{l=1}^m v_l = 0}} f(\mathbf{u}, \mathbf{v}) := \frac{1}{m} \sum_{l=1}^m \{ \mathbb{1}, B_l(u_l, v_l) \mathbb{1} - u_l, p_l \},$$

$$\mathbf{u} = [u_1, \dots, u_m], \mathbf{v} = [v_1, \dots, v_m], u_l, v_l \in \mathbf{R}^n, \\ B_l(u_l, v_l) := \text{diag}(e^{u_l}) \exp(-C_l/\gamma) \text{diag}(e^{v_l}).$$

IBP is equivalent to alternating minimization for the dual problem.

- ▶ $u_l^{t+1} := \ln p_l - \ln K_l e^{v_l^t}, \mathbf{v}^{t+1} := \mathbf{v}^t$
- ▶ $v_l^{t+1} := \frac{1}{m} \sum_{k=1}^m \ln K_k^T e^{u_k^t} - \ln K_l^T e^{u_l^t},$
 $\mathbf{u}^{t+1} := \mathbf{u}^t$



Complexity of IBP

Theorem 5

Let $c := \max_{l=1, \dots, m} \|C_l\|_\infty$ and fix an accuracy parameter $\varepsilon' > 0$. Then Algorithm 1 (Iterative Bregman Projections) terminates after

$$N \leq 4 + \frac{44c}{\gamma \varepsilon'}$$

iterations.

The couplings produced at termination, $\pi_1^{(N)}, \dots, \pi_m^{(N)}$, together with $q^{(N)} = \sum_{l=1}^m w_l \pi_l^{(N)} \mathbf{1}$, satisfy

$$\sum_{l=1}^m w_l \|\pi_l^{(N)} \mathbf{1} - q^{(N)}\|_1 \leq \varepsilon',$$

and

$$\sum_{l=1}^m w_l \left(\langle C_l, \pi_l^{(N)} \rangle + \gamma H(\pi_l^{(N)}) \right) - \min_{q \in \mathcal{S}_n(1)} \sum_{l=1}^m w_l W_\gamma(p_l, q) \leq c \varepsilon'.$$

Barycenters and Quantization

Why Probability-Proportional-to-Size (PPS) Quantization?

- ▶ In decentralised optimisation the *communication* of dense gradient vectors dominates runtime.
- ▶ PPS converts an n -dimensional vector into an *unbiased* random sample of $M \ll n$ indices, with probabilities proportional to component magnitudes.
- ▶ Each iteration transmits only $\approx 2M \log_2 n$ bits yet preserves high-quality gradient information, matching the best known compression operators in bit-complexity while enjoying a smaller second moment when gradients lie in a simplex.
- ▶ Leveraging PPS, we derive *accelerated* primal and primal-dual methods with high-probability guarantees and near-optimal oracle and communication complexity.

Duality and Quantization

Lemma 2

Given a positive Radon probability measure $\mu \in \mathcal{M}_+^1(\mathcal{Y})$ with density $q(y)$ on a metric space \mathcal{Y} , the Fenchel-Legendre dual function of $\mathbf{W}_{\gamma,\mu}(p)$ can be written as

$$\mathbf{W}_{\gamma,\mu}^*(\bar{\lambda}) = \mathbf{E}_{Y \sim \mu} \left[\gamma \log \left(\frac{1}{q(Y)} \sum_{\ell=1}^n \exp \left(\frac{[\bar{\lambda}]_{\ell} - c_{\ell}(Y)}{\gamma} \right) \right) \right].$$

Moreover, $\mathbf{W}_{\gamma,\mu}^*(\bar{\lambda})$ has m/γ -Lipschitz gradient w.r.t. 2-norm, and its l -th coordinate, for $l = 1, \dots, n$, is

$$[\nabla \mathbf{W}_{\gamma,\mu}^*(\bar{\lambda})]_l = \mathbf{E}_{Y \sim \mu} \left[\frac{\exp(([\bar{\lambda}]_l - c_l(Y))/\gamma)}{\sum_{\ell=1}^n \exp(([\bar{\lambda}]_{\ell} - c_{\ell}(Y))/\gamma)} \right]. \quad (1)$$

Sampling

Given that at iteration k each agent i can obtain $M_{i,1}^k$ independent realizations of the random variable $Y^i \sim \mu_i$:

$$\widehat{\mathbf{V}} \mathbf{W}_{\gamma, \mu_i}^*(\bar{\lambda}_i) = \frac{1}{M_{i,1}} \sum_{r=1}^{M_{i,1}} p_i(\bar{\lambda}_i, Y_r^i), \quad (2)$$

where, for all $l = 1, \dots, n$,

$$[p_i(\bar{\lambda}_i, Y_r^i)]_l = \frac{\exp(([\bar{\lambda}_i]_l - c_l(Y_r^i))/\gamma)}{\sum_{k=1}^n \exp(([\bar{\lambda}_i]_k - c_k(Y_r^i))/\gamma)}, \quad (3)$$

IMPORTANTLY: $\widehat{\mathbf{V}} \mathbf{W}_{\gamma, \mu}^*(\bar{\lambda}) \in S_1(n)$.

Stochastic Gradient PPS Quantisation (per agent i at iter. k)

1. **Monte-Carlo OT gradient.** Draw M_1 i.i.d. samples $Y_r^i \sim \mu_i$ and form

$$\widehat{\nabla} W_{\gamma, \mu_i}^*(\lambda_i) = \frac{1}{M_1} \sum_{r=1}^{M_1} g(Y_r^i; \lambda_i), \quad \text{see (2.8).}$$

2. **PPS sampling.** Let Z^i be an independent categorical r.v. with $\Pr\{Z^i = l\} = [\widehat{\nabla} W_{\gamma, \mu_i}^*(\lambda_i)]_l$, $l = 1, \dots, n$. Take M_2 i.i.d. draws $\{Z_r^i\}_{r=1}^{M_2}$.
3. **Quantised gradient.** Construct the histogram

$$\widetilde{\nabla} W_{\gamma, \mu_i}^*(\lambda_i) = \frac{1}{M_2} \sum_{r=1}^{M_2} e_{Z_r^i},$$

where e_l is the l -th canonical basis vector.

This yields an *unbiased* estimator using only $M_2 \log_2 n$ bits of communication per agent at iteration k .

Unbiased estimator

Lemma 3

The full quantized stochastic gradient is unbiased, i.e.,

$\mathbf{E} \tilde{\nabla} \mathbf{W}_\gamma^*(\lambda) = \nabla \mathbf{W}_\gamma^*(\lambda)$. Furthermore, its variance is bounded as

$$\mathbf{E} \|\tilde{\nabla} \mathbf{W}_\gamma^*(\lambda) - \nabla \mathbf{W}_\gamma^*(\lambda)\|_2^2 \leq 2\lambda_{\max}(W) \sum_{i=1}^m \left(\frac{1}{M_{i,1}} + \frac{1}{M_{i,2}} \right).$$

Let's focus on the generic case:

$$f(x_*) = \min_{\substack{x \in \mathbb{R}^n \\ Ax=b}} f(x), \quad (4)$$

$$M_1 = r, \quad M_2 = M.$$

Lemma 1 — Sub-Gaussianity of PPS

Lemma 4

Let $\sigma_{r,M}^2 = 50 \left(\frac{2(1-1/n)B^2}{eM} + \frac{\sigma^2}{r} \right)$. Then, it holds that

$$\forall x \in \mathbb{R}^n, \mathbb{E}_{\xi, \mathbf{k}, I} [\exp(\|PPS(x, \xi, \mathbf{k}, I) - \nabla f(x)\|^2 / \sigma_{r,M}^2)] \leq e.$$

Theorem 1 — High-Probability Primal Error Bound

Theorem 6

For $\delta \in (0, 1)$ Algorithm 1 guarantees

$$\mathbb{P}(f(x_T) - f(x^*) \leq \varepsilon(T, \delta, \alpha, \beta, r, M)) > 1 - \delta,$$

where

$$\varepsilon = \frac{\beta_T R^2}{2A_T} + \frac{C'_1 R}{A_T} \sqrt{\sum_{t=0}^T \alpha_t^2 \sigma_{r_t, M_t}^2} + \frac{C'_2}{A_T} \sum_{t=0}^T \frac{A_t \sigma_{r_t, M_t}^2}{\beta_t - L},$$

$$C'_2 = 1 + \ln \frac{4}{\delta} = \mathcal{O}(\ln \frac{1}{\delta}) \text{ and } C'_1 = (2J(T) + \sqrt{2} - 1)(\sqrt{2} + (\sqrt{2} + 1)\sqrt{3 \ln \frac{4}{\delta}}) + \sqrt{2} - 2 = \mathcal{O}(\text{poly}(\ln T))\sqrt{\ln \frac{1}{\delta}}.$$

Theorem 2 — Primal–Dual Large-Deviation Bounds

Theorem 7

For $\delta \in (0, 1)$ Algorithm 2 (primal-dual) ensures

$$\mathbb{P}(f(x_T) - f(x^*) \leq \varepsilon) > 1 - \delta, \quad \mathbb{P}\left(\|Ax_T - b\| \leq \frac{\varepsilon}{R}\right) > 1 - \delta,$$

with

$$\varepsilon = \frac{\beta_T R^2}{2A_T} + \frac{C_3 R + C_4 L / \|A\|_2}{A_T} \sqrt{\sum_{t=0}^T \alpha_t^2 \sigma_{r_t, M_t}^2} + \frac{C_2}{A_T} \sum_{t=0}^T \frac{A_t \sigma_{r_t, M_t}^2}{\beta_t - L},$$

$$C_4 = \sqrt{2}(1 + \sqrt{3 \ln(5/\delta)}), \quad C_3 = C_1 + 2\sqrt{2}J(T)(1 + \sqrt{3 \ln(5/\delta)}).$$

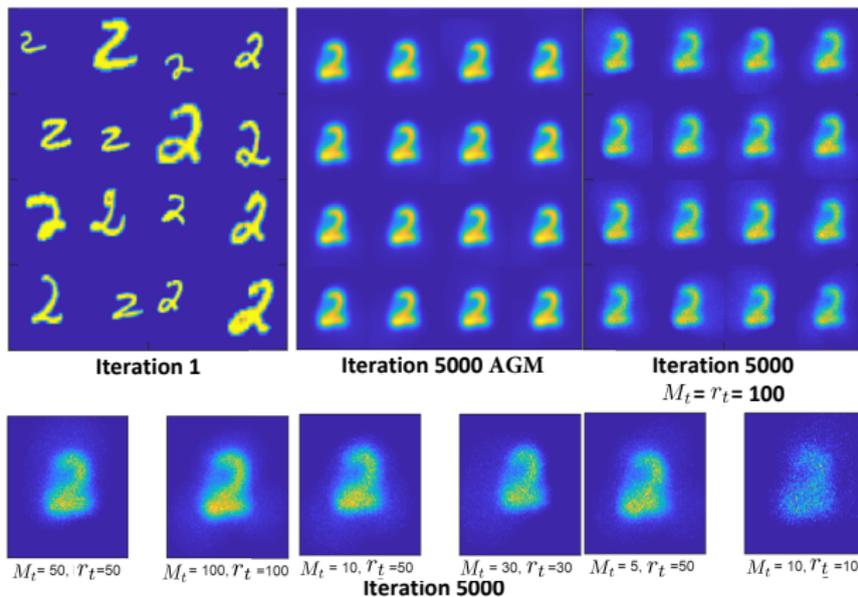
Communication Complexity

Theorem 8

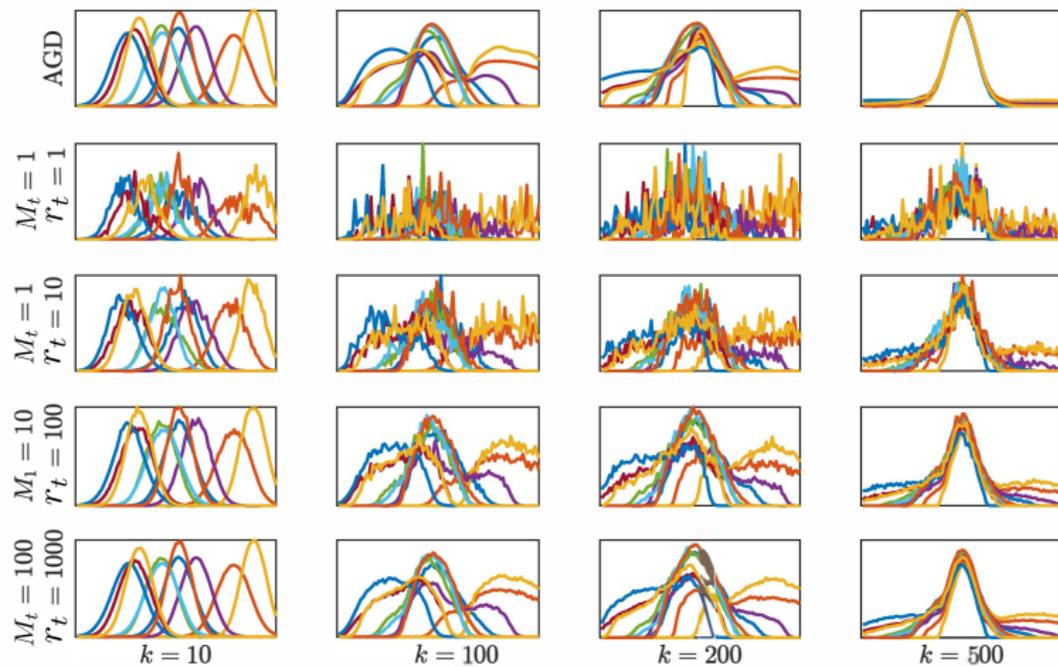
Given $\varepsilon > 0$, $0 < \delta < 1$, \mathbf{r} and \mathbf{M} selected appropriately, one can ensure $P(f(\mathbf{x}) - f(\mathbf{x}_*) \leq \varepsilon) > 1 - \delta$ and $P\left(\|\sqrt{\mathbf{W}}\mathbf{x}\| \leq \frac{\varepsilon\sqrt{2m}}{B_*}\right) > 1 - \delta$, if any node sends

$$\tilde{O}\left(B^2 d \log n \max\left\{\frac{1}{\sigma_i^2} \sqrt{\frac{DB_*^2}{\gamma \varepsilon m d}}, \frac{DB_*^2}{\varepsilon^2}, \frac{m^2}{\gamma^2 \varepsilon^2}\right\}\right) \text{ bits.}$$

Proof by Picture



Proof by Picture



Decentralized and Equitable OT

Decentralized Optimal Transport

There is an undirected and connected network of n , where each agent has access to a column of the matrix C .

$$\min_{x_i \in \mathbf{R}_+^n} \sum_{i=1}^n c_i^T x_i \text{ s.t. } \sum_{i=1}^n x_i = p \text{ and } x_i^T \mathbb{1}_n = q_i \quad (5)$$

- ▶ Problem setup involves multiple agents, each with access to part of the cost information.
- ▶ Agents collaborate to design a transportation plan that minimizes the total transportation cost.
- ▶ Communication is limited to immediate neighbors in a network.

Decentralized and Equitable Optimal Transport (DE-OT)

Each agent k has its own private cost function C^k .

- ▶ Combines D-OT and EOT, focusing on decentralized networks and equitable cost sharing.
- ▶ Problem formulation:

$$\begin{aligned} & \min_{X^k \in \mathbf{R}_+^{n \times n}} \sum_{k=1}^N \langle C^k, X^k \rangle \\ \text{s.t. } & \langle C^k, X^k \rangle = \langle C^l, X^l \rangle \text{ for all } k, l \in [N], \\ & \left(\sum_{k=1}^N X^k \right) \mathbf{1}_n = p \text{ and } \left(\sum_{k=1}^N X^k \right)^\top \mathbf{1}_n = q. \end{aligned} \quad (6)$$

- ▶ Ensures fairness among agents.

Reformulating D-OT as a DCCO problem

Proposition 5

The DOT problem (5) is equivalent to

$$\min_{x_i} \sum_{i=1}^n c_i^T x_i + \iota_{\geq 0}(x_i) \text{ s.t. } \sum_{i=1}^n M_i x_i = \begin{bmatrix} p \\ q \end{bmatrix}, \quad (7)$$

where the indicator function $\iota_{\geq 0}(y)$ is defined by

$$\iota_{\geq 0}(y) \begin{cases} 0 & \text{if } y \geq 0 \\ \infty & \text{otherwise,} \end{cases} \quad (8)$$

and the matrix $M_i \in \mathbf{R}^{2n \times n}$ is the i^{th} column-block of

$$[M_1, \dots, M_n] = M \begin{bmatrix} I_n & I_n & \cdots & I_n \\ \mathbb{1}_n^T & \mathbf{0}_n^T & \cdots & \mathbf{0}_n^T \\ \mathbf{0}_n^T & \mathbb{1}_n^T & \cdots & \mathbf{0}_n^T \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}_n^T & \mathbf{0}_n^T & \cdots & \mathbb{1}_n^T \end{bmatrix}. \quad (9)$$

Reformulating DE-OT as a DCCO problem

Proposition 6

Problem (6) is equivalent to

$$\min_{x_i} \sum_{i=1}^N c_i^\top x_i + \iota_{\geq 0}(x_i) \quad \text{s.t.} \quad \sum_{i=1}^N \begin{bmatrix} M \\ E_i \end{bmatrix} x_i = \begin{bmatrix} p \\ q \\ \mathbf{0}_{N-1} \end{bmatrix}, \quad (10)$$

where the indicator function $\iota_{\geq 0}$, and matrices E_i are defined as follows:

$$E_i \left(\tilde{L}_{*,i} \right) c_i^\top = \tilde{L} \begin{bmatrix} \mathbf{0}_{(i-1) \times n^2} \\ c_i^\top \\ \mathbf{0}_{(N-i) \times n^2} \end{bmatrix} \in \mathbf{R}^{(N-1) \times n^2}. \quad (11)$$

Common formulation

$$\begin{aligned} \min_{\mathbf{x}=(x_i) \in \mathbf{R}^{Nd}} f(\mathbf{x}) &\triangleq \mathbf{c}^T \mathbf{x} + \iota_{\geq \mathbf{0}}(\mathbf{x}) = \sum_{i=1}^N c_i^T x_i + \iota_{\geq \mathbf{0}}(x_i) \\ \text{s.t. } A\mathbf{x} &= \sum_{i=1}^N A_i x_i = b. \end{aligned} \tag{12}$$

Table 1: Algorithms for Problem (12).

Paper	Convergence rate	Singe-loop
Falsone et al., 2020	Asymptotic	No
Chang et al. 2014,2016	$O(1/k)$, ergodic	No
Su et al. 2021	$O(1/k)$, non-ergodic	No
Alghunaim et al. 2019 and Li et al. 2023	Asymptotic	Yes

Euclidean Regularization

$$\begin{aligned} \min_{\substack{\mathbf{x}=(x_i) \in \mathbf{R}^{Nd}, \\ \mathbf{y}=(y_i) \geq 0}} \quad f(\mathbf{x}) &\triangleq \sum_{i=1}^N c_i^T x_i + \frac{\eta}{2} \|x_i\|^2 \\ \text{s.t.} \quad A\mathbf{x} = \sum_{i=1}^N A_i x_i &= b \text{ and } \mathbf{y} - \mathbf{x} = \mathbf{0}. \end{aligned} \tag{13}$$

Algorithm: PDC-ADMM

Algorithm 1 PDC-ADMM for Decentralized Optimal Transport

Penalty parameters $\rho, \beta_i, \tau_i > 0$ satisfying conditions Optimal solution x^* Initialize $k = 0$ and $x^0, y^0, z^0, \lambda^0, \rho^0$

while stopping criterion not met **do** Exchange λ_i^k with neighbors N_i Update y_i^{k+1} :

$$y_i^{k+1} := \left(1 - \frac{1}{\beta_i \tau_i}\right) y_i^k + \frac{1}{\beta_i \tau_i} (x_i^k - \tau_i z_i^k)$$

Update x_i^{k+1} :

$$x_i^{k+1} := x_i^k - \frac{1}{\beta_i} \left(c_i + \eta x_i^k + \frac{1}{\tau_i} (x_i^k - y_i^{k+1} - \tau_i z_i^k) + \frac{A_i^T}{2\rho |N_i|} \left(A_i x_i^{k+1} - \frac{1}{N} b - \rho^k + \rho \sum_{j \in N_i} (\lambda_j^k + \lambda_i^k) \right) \right)$$

Update λ_i^{k+1} :

$$\lambda_i^{k+1} := \frac{1}{2|N_i|} \left(\sum_{j \in N_i} (\lambda_j^k + \lambda_i^k) + \frac{1}{\rho} \rho_i^{k+1} + \frac{1}{\rho} \left(A_i x_i^{k+1} - \frac{1}{N} b \right) \right)$$

Update z_i^{k+1} :

$$z_i^{k+1} := z_i^k + \frac{1}{\tau_i} (y_i^{k+1} - x_i^{k+1})$$

Update ρ_i^{k+1} :

$$\rho_i^{k+1} := \rho_i^k + \rho \sum_{j \in N_i} (\lambda_j^{k+1} - \lambda_j^k)$$

Increment k Output x_i^k : $x^* := \frac{1}{k} \sum_{t=1}^k x_i^t$

Theorem 9:

Let the graph be connected and undirected. Then the sequence $(\mathbf{x}^k, \mathbf{y}^k)$ generated by PDC-ADMM with parameters $\rho, \beta_i, \tau_i > 0$ satisfying

$$\beta_i \tau_i \geq 1 \text{ and } \left(\beta_i - \frac{\beta_i}{\beta_i \tau_i - 1} - 1 \right) I_d - \frac{1}{2\rho |\mathcal{N}_i|} A_i^\top A_i \succ \mathbf{0}, \quad (14)$$

converges to the optimal solution $(\mathbf{x}^*, \mathbf{y}^* = \mathbf{x}^*)$ of Problem (13). Furthermore,

$$|f(\bar{\mathbf{x}}^k) - f(\mathbf{x}^*)| + \left\| A \bar{\mathbf{x}}^k - b \right\| + \left\| \bar{\mathbf{y}}^k - \bar{\mathbf{x}}^k \right\| = O(1/k), \quad (15)$$

where $(\bar{\mathbf{x}}^k, \bar{\mathbf{y}}^k) = \frac{1}{k} \sum_{t=1}^k (\mathbf{x}^t, \mathbf{y}^t)$.

Numerical Results

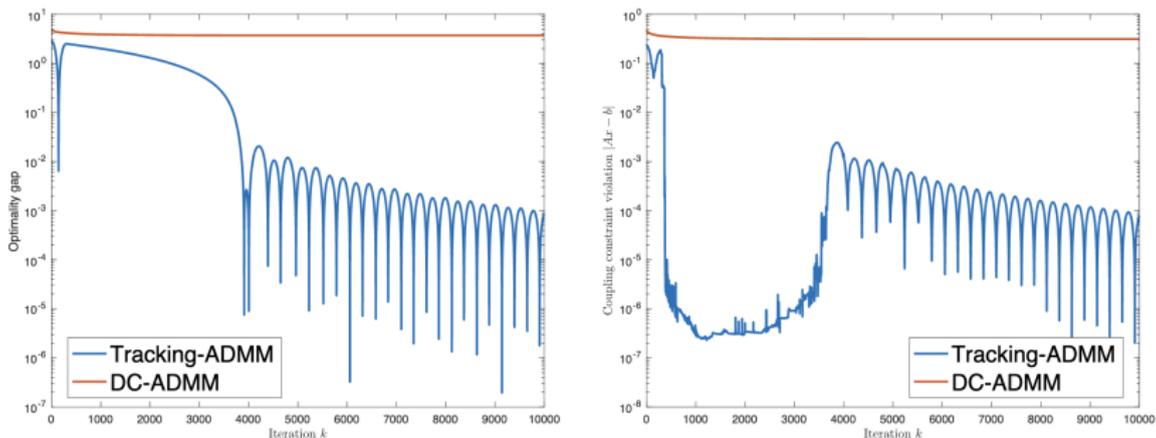


Figure 1: Optimality gap $f(\mathbf{x}^k) - f^*$ and feasibility violation $\|A^k - b\|$ for D-OT and DE-OT under Tracking-ADMM and DC-ADMM.

Decentralized Computation of Wasserstein Barycenters

César A. Uribe

Alexey Kroshnin Darina Dvinskikh Pavel Dvurechensky
Alexander Gasnikov Nazarii Tupitsa Roman Krawtschenko
Ivan Lau Shiqian Ma Dmitry Pasechnyuk Angelia Nedić

May 19, 2025



RICE ENGINEERING

Electrical and Computer Engineering

References I

- [1] A. Kroshnin, D. Dvinskikh, P. Dvurechensky, A. Gasnikov, N. Tupitsa, C. Uribe. On the complexity of approximating Wasserstein barycenter. 2020.
- [2] P. Dvurechensky, D. Dvinskikh, A. Gasnikov, C. Uribe, A. Nedić. Decentralize and randomize: Faster algorithm for Wasserstein barycenters. In *NeurIPS*, 2018.
- [3] C. Uribe, D. Dvinskikh, P. Dvurechensky, A. Gasnikov, A. Nedić. Distributed computation of Wasserstein barycenters over networks. In *IEEE CDC*, 2018.
- [4] R. Krawtschenko, C. Uribe, A. Gasnikov, P. Dvurechensky. Distributed optimization with quantization for computing Wasserstein barycenters. 2020.
- [5] D. Pasechnyuk, P. Dvurechensky, C. Uribe, A. Gasnikov. Decentralised convex optimisation with probability-proportional-to-size quantization. 2025.
- [6] I. Lau, S. Ma, C. Uribe. Decentralized and equitable optimal transport. 2024.