

# Low-Rank Sinkhorn Factorization

M. Scetbon



M. Cuturi



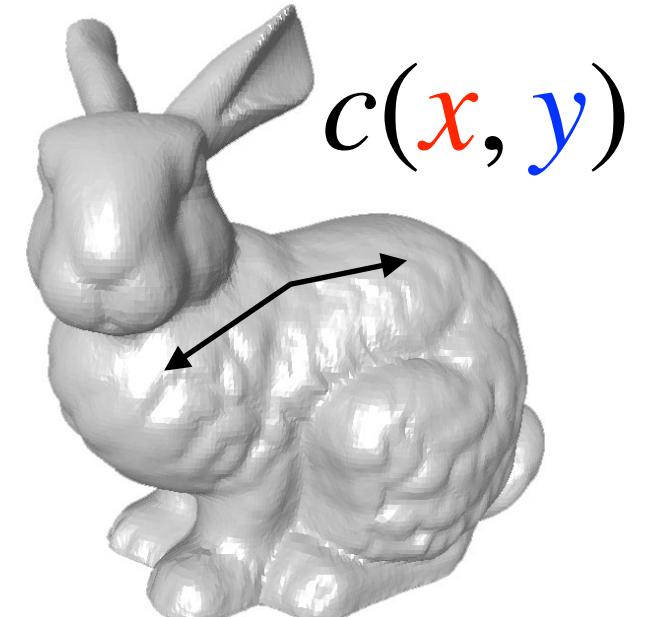
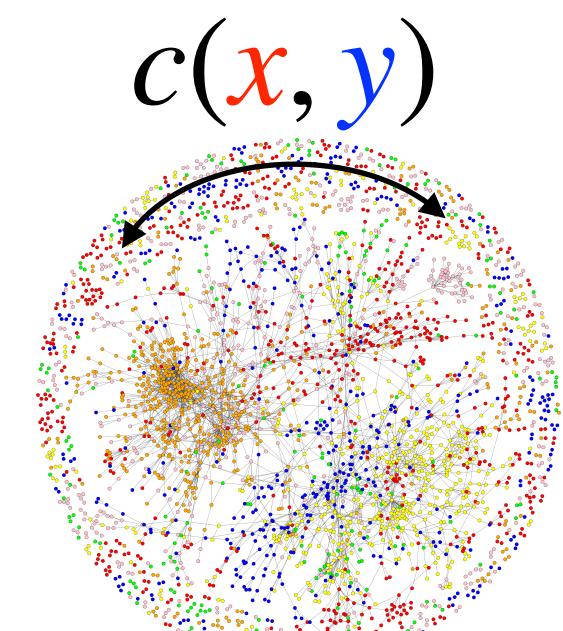
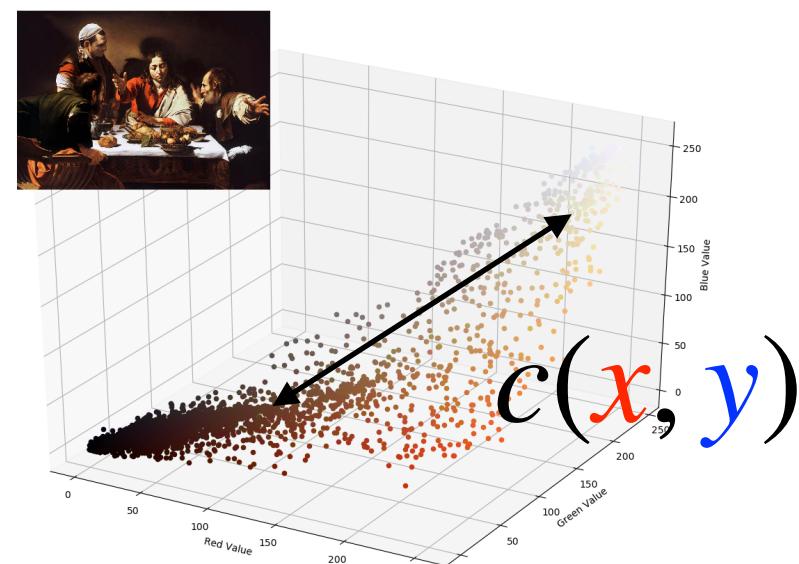
G. Peyré



Thirty-eighth International Conference on Machine Learning

# Optimal Transport: Comparing Distributions

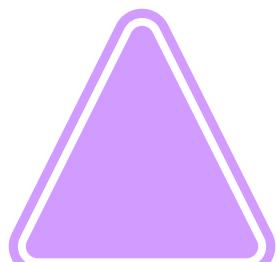
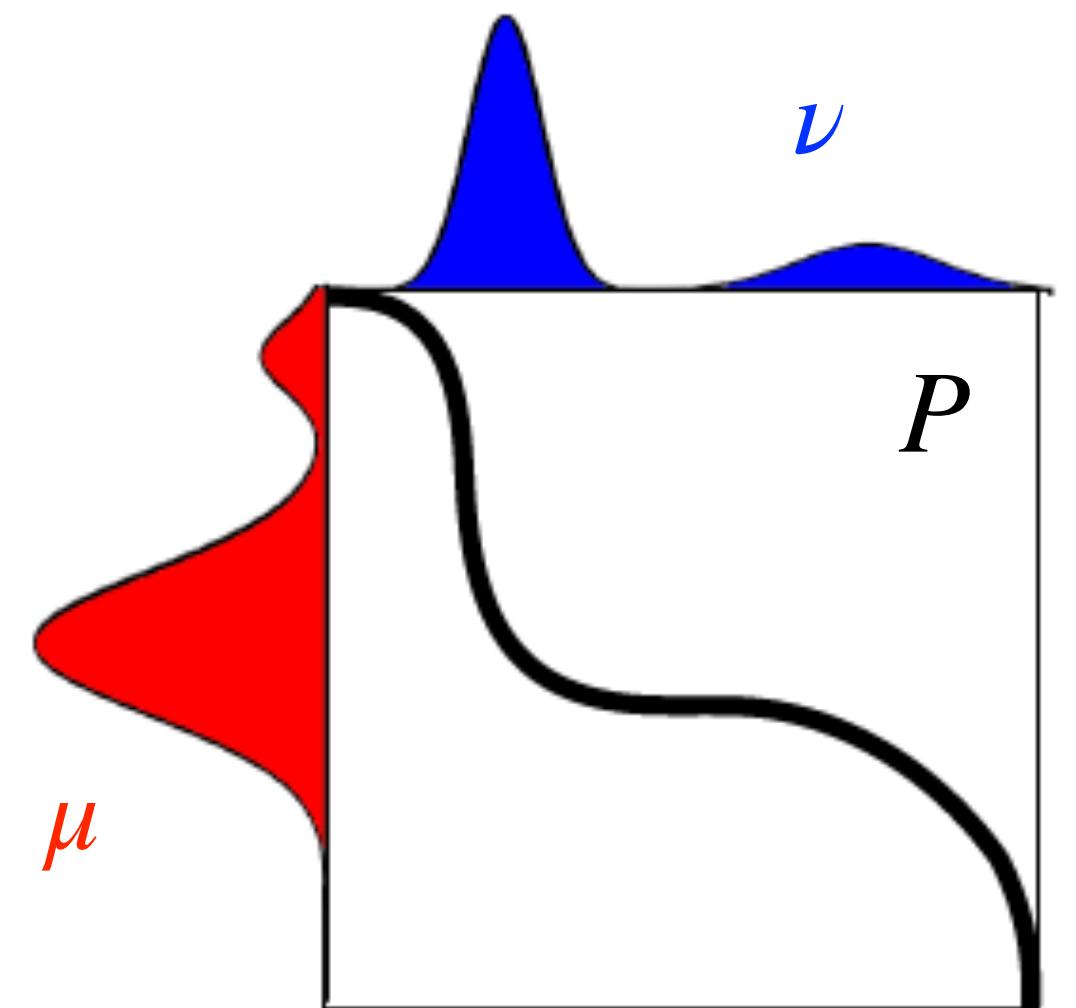
Some Examples of Probability Measures:



Discrete Optimal Transport:

- Discrete Distributions:  $\mu = \sum_{i=1}^n a_i \delta_{x_i} \in \mathcal{M}_1^+(\mathcal{X})$  ,  $\nu = \sum_{j=1}^m b_j \delta_{y_j} \in \mathcal{M}_1^+(\mathcal{Y})$
- Set of Couplings:  $\Pi_{\mu, \nu} = \{P \in \mathbb{R}_+^{n \times m} \text{ s.t. } P\mathbf{1}_m = \mu, P^T\mathbf{1}_n = \nu\}$
- Cost Matrix:  $\forall i, j \ C_{i,j} = c(x_i, y_j)$  where  $c : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$

$$W_c(\mu, \nu) = \min_{P \in \Pi_{\mu, \nu}} \langle P, C \rangle$$



Costly to compute  $\longrightarrow$  Linear Program:  $\mathcal{O}(n^3 \log(n))$  complexity

# Entropic Regularization

Shannon Entropy:  $H(P) = - \sum_{i,j} P_{i,j}(\log(P_{i,j}) - 1)$

Regularized Optimal Transport:

$$W_{c,\varepsilon}(\mu, \nu) := \min_{\substack{P \in \mathbb{R}_+^{n \times m} \\ P\mathbf{1}_m = \textcolor{red}{a}, P^T\mathbf{1}_n = \textcolor{blue}{b}}} \langle P, C \rangle - \varepsilon H(P)$$

Sinkhorn Algorithm:

- Kernel Matrix:  $\mathbf{K} := \exp(-C/\varepsilon)$
- Until convergence, at each iteration compute :  $v \leftarrow \frac{\textcolor{blue}{b}}{\mathbf{K}^T u}$ ,  $u \leftarrow \frac{\textcolor{red}{a}}{\mathbf{K} v}$
- Output:  $P_\varepsilon^* = \text{Diag}(u)\mathbf{K}\text{Diag}(v)$

Computing  $\mathbf{K}^T u$  and  $\mathbf{K} v$  requires  $\mathcal{O}(\textcolor{red}{nm})$  operations

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

Quadratic time algorithm

Low-Rank Approximation of the Kernel:

Replace  $\mathbf{K}$  in the Sinkhorn iterations by  $\widetilde{\mathbf{K}} = AB^T$  where  $(A, B) \in (\mathbb{R}_+^*)^{n \times \textcolor{red}{r}} \times (\mathbb{R}_+^*)^{m \times \textcolor{red}{r}}$

→ Computing  $BA^T u$  and  $AB^T v$  requires  $\mathcal{O}(\textcolor{green}{nr})$  algebraic operations

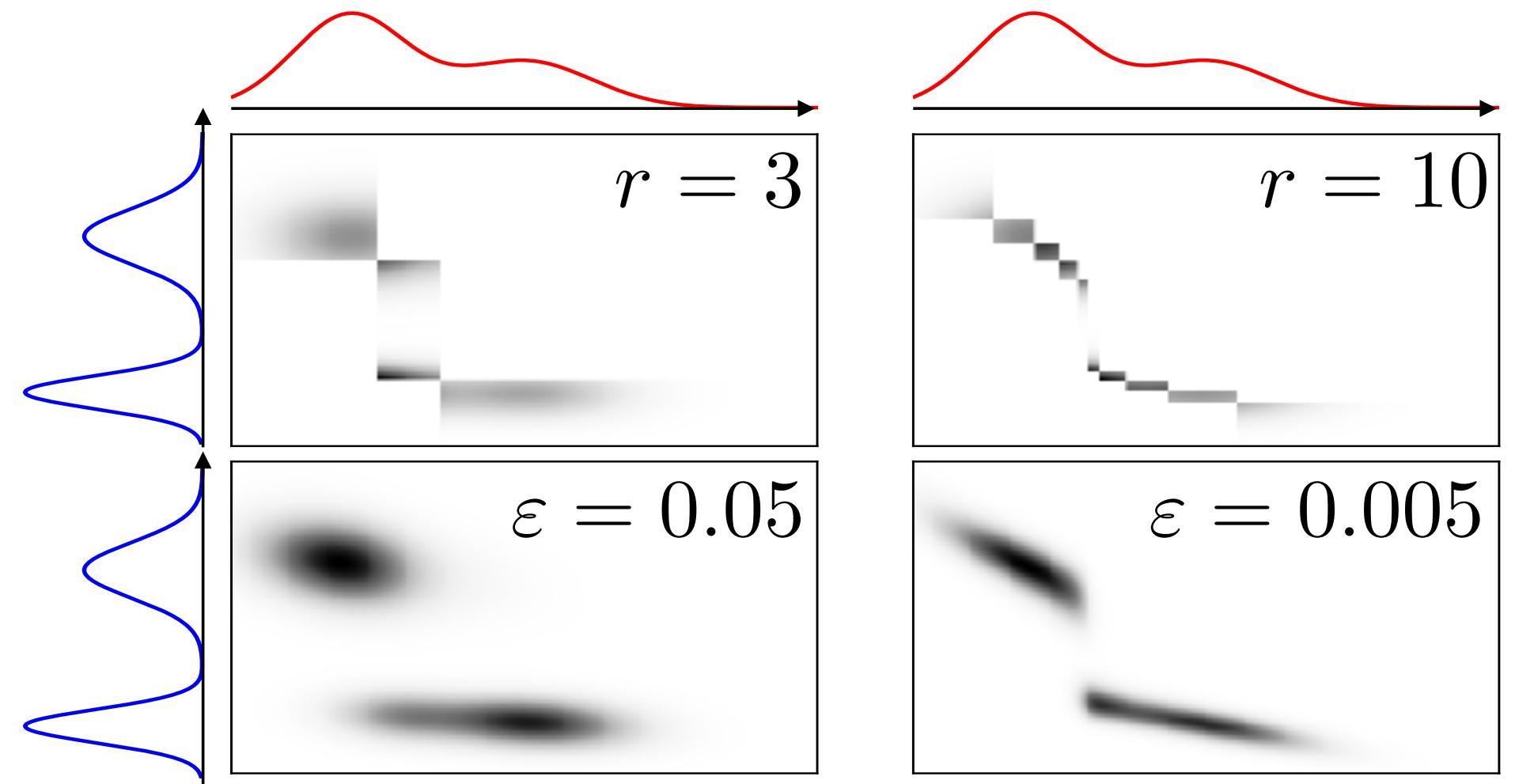
Linear time algorithm

In this work, we propose *instead* to directly constraint the coupling to admit a low-NN rank

# Low-Rank Optimal Transport

NN rank :  $\text{rk}_+(M) := \min \left\{ q \mid M = \sum_{i=1}^q R_i, \forall i, \text{rk}(R_i) = 1, R_i \geq 0 \right\}$

Low-NN Rank Couplings:  $\Pi_{\color{red}a,\color{blue}b}(\color{red}r) := \{P \in \Pi_{\color{red}a,\color{blue}b} \text{ s.t. } \text{rk}_+(P) \leq \color{red}r\}$



Definition of Low-rank Optimal Transport

$$\text{LOT}_{\color{red}r}(\color{red}\mu, \color{blue}\nu) := \min_{P \in \Pi_{\color{red}a,\color{blue}b}(\color{red}r)} \langle P, C \rangle$$

Characterization of Low-NN Rank Couplings:

$$\Pi_{\color{red}a,\color{blue}b}(\color{red}r) = \{P \in \mathbb{R}_+^{n \times m} \mid P = Q \text{Diag}(1/\color{purple}g) R^T, Q \in \Pi_{\color{red}a,g}, R \in \Pi_{\color{blue}b,g}, g > 0 \text{ and } \color{purple}g \in \Delta_{\color{red}r}\}$$

## Reparametrization of LOT

$$\text{LOT}_r(\mu, \nu) = \min_{(Q, R, g) \in \mathcal{C}_1(a, b, r) \cap \mathcal{C}_2(r)} \langle C, Q \text{Diag}(1/g) R^T \rangle$$

where  $\left\{ \begin{array}{l} \mathcal{C}_1(a, b, r) := \{(Q, R, g) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times (\mathbb{R}_+^*)^r \text{ s.t. } Q\mathbf{1}_r = a, R\mathbf{1}_r = b\} \\ \mathcal{C}_2(r) := \{(Q, R, g) \in \mathbb{R}_+^{n \times r} \times \mathbb{R}_+^{m \times r} \times (\mathbb{R}_+)^r \text{ s.t. } Q^T \mathbf{1}_n = R^T \mathbf{1}_m = g\} \end{array} \right.$

## Entropic Regularization of LOT

$$\text{LOT}_{r,\varepsilon}(\mu, \nu) := \inf_{(Q, R, g) \in \mathcal{C}_1(a, b, r) \cap \mathcal{C}_2(r)} \langle C, Q \text{Diag}(1/g) R^T \rangle - \varepsilon H((Q, R, g))$$

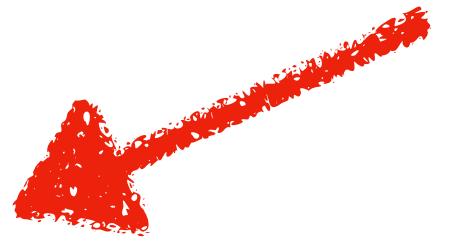
### Proposition

If  $\varepsilon = 0$  then the infimum is attained. If  $\varepsilon > 0$ , then if  $r = 1$ , the infimum is attained and for  $r \geq 2$ , the problem admits a minimum if  $\text{LOT}_{r,\varepsilon}(\mu, \nu) < \text{LOT}_{r-1,\varepsilon}(\mu, \nu)$ .

## Mirror Descent Scheme

$$(Q_{k+1}, R_{k+1}, g_{k+1}) := \min_{(\zeta_1, \zeta_2, \zeta_3) \in \mathcal{C}_1(\textcolor{red}{a}, \textcolor{blue}{b}, \textcolor{brown}{r}) \cap \mathcal{C}_2(\textcolor{brown}{r})} \text{KL}((\zeta_1, \zeta_2, \zeta_3), (\textcolor{green}{K}_k^{(1)}, \textcolor{green}{K}_k^{(2)}, \textcolor{green}{K}_k^{(3)}))$$

Quadratic time algorithm



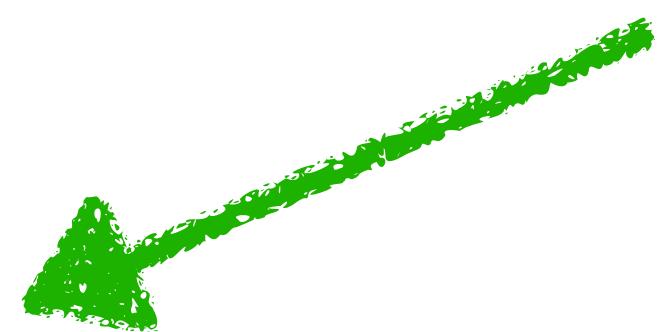
where  $(Q_0, R_0, \textcolor{violet}{g}_0) \in \mathcal{C}_1(\textcolor{red}{a}, \textcolor{blue}{b}, \textcolor{brown}{r}) \cap \mathcal{C}_2(\textcolor{brown}{r})$  ,  $\textcolor{green}{K}_k^{(1)} := \exp(-\gamma_k C R_k \text{Diag}(1/\textcolor{violet}{g}_k) - (\gamma_k \varepsilon - 1) \log(Q_k))$  ,

$\textcolor{green}{K}_k^{(2)} := \exp(-\gamma_k C^T Q_k \text{Diag}(1/\textcolor{violet}{g}_k) - (\gamma_k \varepsilon - 1) \log(R_k))$  ,  $\textcolor{green}{K}_k^{(3)} := \exp(\gamma_k \omega_k / \textcolor{violet}{g}_k^2 - (\gamma_k \varepsilon - 1) \log(\textcolor{violet}{g}_k))$  ,  $[\omega_k]_i := [Q_k^T C R_k]_{i,i}$  ,  $\gamma_k > 0$

**Remarks:** • Each iteration of the MD is a convex problem which can be solved efficiently using the IBP algorithm

- Given  $(\textcolor{green}{K}_k^{(i)})_{i=1}^3$ , each iteration of the IBP algorithm requires  $\mathcal{O}((n + m)r)$  algebraic operations
- Computing  $(\textcolor{green}{K}_k^{(i)})_{i=1}^2$  requires  $\mathcal{O}(nmr)$  algebraic operations

Linear time algorithm

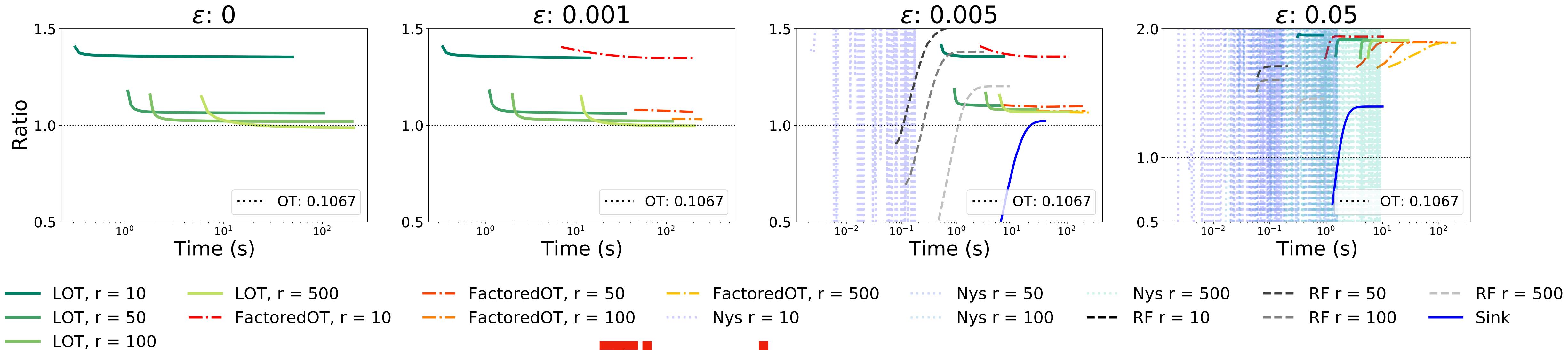


## Linear Time Approximation

If  $C \simeq AB^T$  where  $(A, B) \in \mathbb{R}^{n \times \textcolor{violet}{d}} \times \mathbb{R}^{m \times \textcolor{violet}{d}}$  and  $\textcolor{violet}{d} \ll \min(n, m)$   $\longrightarrow$   $(\textcolor{green}{K}_k^{(i)})_{i=1}^2$  can be performed in  $\mathcal{O}((n + m)dr)$

**Example:** The squared Euclidean distance, or more generally any distance matrix.

## Experiments



# Thank you

## Advantages of the proposed method

- Faster to compute than Sinkhorn.
- Its parametrization, which is the rank  $r$ , encodes directly a property of the resulting coupling.