

Scalable Optimal Transport in High Dimensions for Graph Distances, Embedding Alignment, and More

Johannes Klicpera, Marten Lienen, Stephan Günnemann

www.daml.in.tum.de/lcn

Massively scalable optimal transport: Sparse and LCN-Sinkhorn

Entropy-regularized optimal transport

$$\bar{P} = \arg \min_P \langle P, C \rangle_F - \lambda H(P)$$

Solvable via Sinkhorn algorithm (alternating matrix normalization):

$$K = e^{-C/\lambda}$$

$$s^{(i)} = p \oslash (Kt^{(i-1)}), \quad t^{(i)} = q \oslash (K^T s^{(i)})$$

$$\bar{P} = \text{diag}(\bar{s}) K \text{diag}(\bar{t})$$

Parallelizable, runs in $\mathcal{O}(n^2)$ → K already has that many entries!

Sparse Sinkhorn

K exponential in C → only near neighbors important!

Sparsely approximate K :

$$K_{ij}^{\text{sp}} = \begin{cases} K_{ij} & \text{if } x_i, x_j \text{ are near,} \\ 0 & \text{otherwise} \end{cases}$$

Calculate P via Sinkhorn algorithm: $\bar{P} \approx \bar{P}^{\text{sp}} = \text{diag}(\bar{s}) K^{\text{sp}} \text{diag}(\bar{t})$
 Runs in $\mathcal{O}(n \log n)$!

Locally corrected Nyström (LCN)

How to incorporate short- and long-range interactions?

Correct low-rank Nyström approximation with sparse values:

$$K_{\text{Nys}} = UA^{-1}V$$

$$K_{\text{LCN}} = K_{\text{Nys}} - K_{\text{Nys}}^{\text{sp}} + K^{\text{sp}}$$

LCN is a general kernel approximation.

We can still use the Sinkhorn algorithm → LCN-Sinkhorn, runs in $\mathcal{O}(n \log n + nl^2)$

GTN: GNN with optimal transport -48% error for graph distance learning

Graph Transport Network (GTN)

Graph neural network (GNN) generates node embeddings
 Embeddings matched using Sinkhorn

→ Learns graph distances via backpropagation

Varying numbers of nodes → learnable unbalanced OT

Bipartite matching (BP) matrix for optimal transport
 (Scaled) norms as deletion cost

$$C_{BP} = \begin{bmatrix} C & C^{(p,\varepsilon)} \\ C^{(\varepsilon,q)} & C^{(\varepsilon,\varepsilon)} \end{bmatrix}, \quad C_{ij}^{(p,\varepsilon)} = \begin{cases} c_{i,\varepsilon} & i = j \\ 0 & i \neq j \end{cases}, \quad C_{ij}^{(\varepsilon,q)} = \begin{cases} c_{\varepsilon,j} & i = j \\ 0 & i \neq j \end{cases}, \quad C_{ij}^{(\varepsilon,\varepsilon)} = 0$$

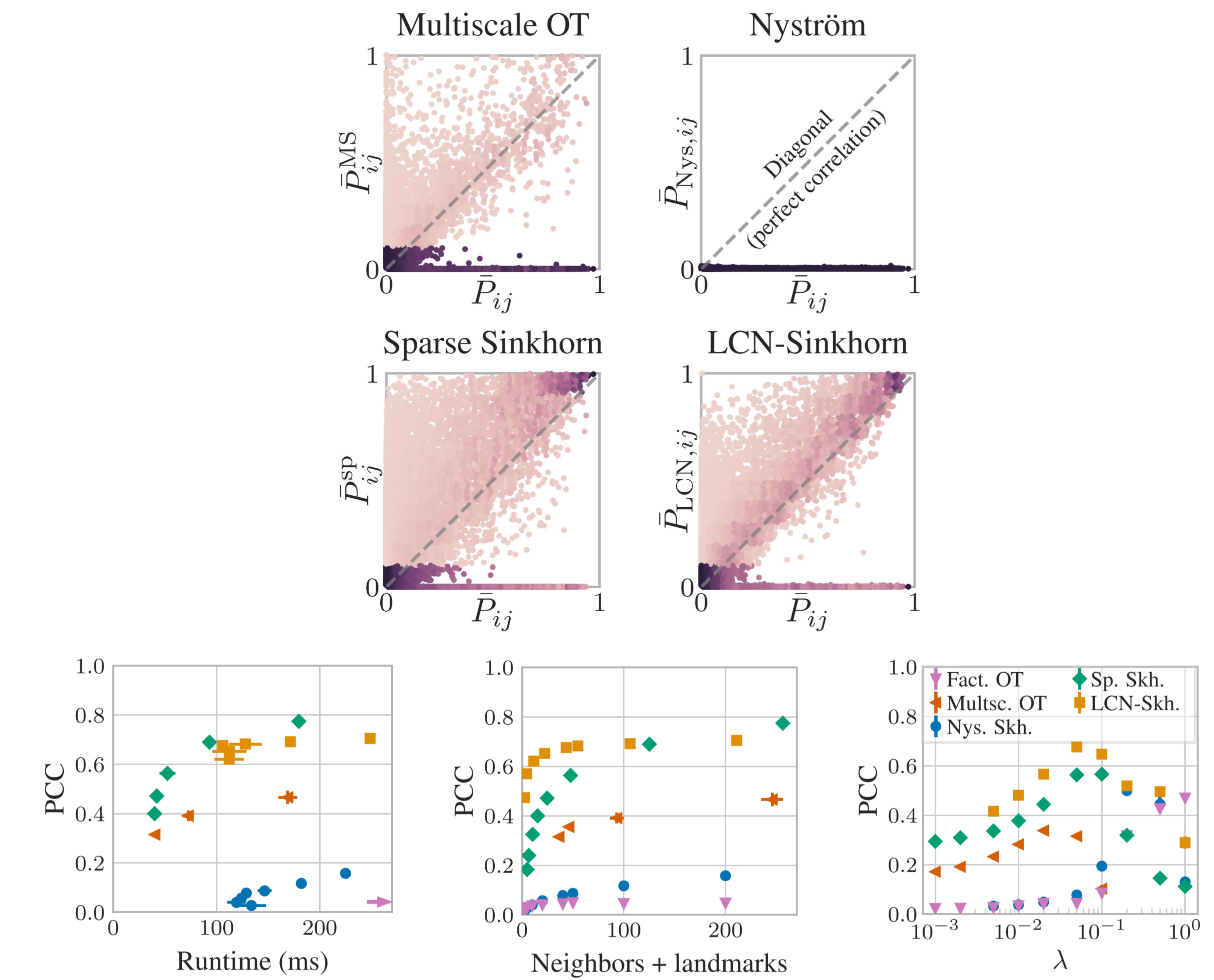
Multi-head OT: OT with multiple embeddings & varying regularization λ

Model	Linux	AIDS30	Pref. att.	Method	AIDS 30	Pref. att.	Pref. att. 200
SiamMPNN	0.09	13.8	12.1	Full Sinkhorn	3.7	4.5	1.3
SimGNN	0.039	4.5	8.3	Nyström Skh.	3.6	6.2	2.4
GMN	0.015	10.3	7.8	Multiscale OT	11.2	27.4	6.7
GTN, 1 head	0.022	3.7	4.5	Sparse Skh.	44	40.7	7.6
8 OT heads	0.012	3.2	3.5	LCN-Sinkhorn	4.0	5.1	1.4
Balanced OT	0.034	15.3	27.4				

LCN: Provably better approximation, same convergence

- Sparse corrections provably reduce kernel approximation error of Nyström for uniform and clustered data [Theorem 1&2].
- Better kernel approximation directly translates to better Sinkhorn approximation [Theorem 3].
- Sparse and LCN-Sinkhorn converge in $\mathcal{O}(\ln \min K / \varepsilon)$, like full Sinkhorn [Theorem 4].
- Sparse and LCN-Sinkhorn allow fast backpropagation via analytical gradients [Prop. 1].

Fast & accurate transport plan and embedding alignment



Embedding alignment

3x faster & 3.1pp more accurate

Method	Time (s)	Avg. accuracy
Original	268	67.9%
Full Sinkhorn	402	70.1%
Multiscale OT	88.2	26.8%
Nyström Sinkhorn	102	47.8%
Sparse Sinkhorn	49.2	68.4%
LCN-Sinkhorn	86.8	71.0%