

A Space Efficient Approximation of Optimal Transport via Hierarchically Separated Trees

Ian Wu^{1,3}, Charles Zhang^{2,3}, Alina Ene³

Del Norte High School, 16601 Nighthawk Lane, San Diego, CA, 92127¹; Great Oak High School, 32555 Deer Hollow Way, Temecula, CA 92592²; Boston University, Boston, MA 02215³



Introduction

- **Optimal Transport (OT)** is a well-studied optimization problem with applications through machine learning (domain adaptation, resource allocation, etc.)
 - OT seeks to minimize the cost of transporting mass between two probability distributions. In other words, OT seeks to minimize $\min_{\mathbf{P} \in \mathbb{R}^{n \times n}} \langle \mathbf{C}, \mathbf{P} \rangle$ where $\mathbf{C} \in \mathbb{R}^{n \times n}$ and \mathbf{P} is a matrix representing mass transfer [1]
- **Sinkhorn**, the most widely used OT algorithm, utilizes entropic regularization to efficiently compute an approximate solution in $O(n^2)$ time [2]
 - Sinkhorn must store the entire transport matrix in memory, however, using $O(n^2)$ space
- **Hierarchical Refinement (HiRef)** attempts to solve this by grouping points in bijective couplings, but relies on several restrictive assumptions, limiting its applicability despite using $O(n)$ space [3]
 - Assumes same number of source and target points that all have equal weight, does not permit mass splitting
- We introduce a new algorithm to approximate OT in $O(n \log \Delta)$ space, an improvement upon Sinkhorn, by grouping points using a Hierarchically Well-Separated Tree Embedding (HST), while accomodating aforementioned limitations.
 - We test the algorithm's efficiency and accuracy on several synthetic datasets, in addition to comparing it with Sinkhorn

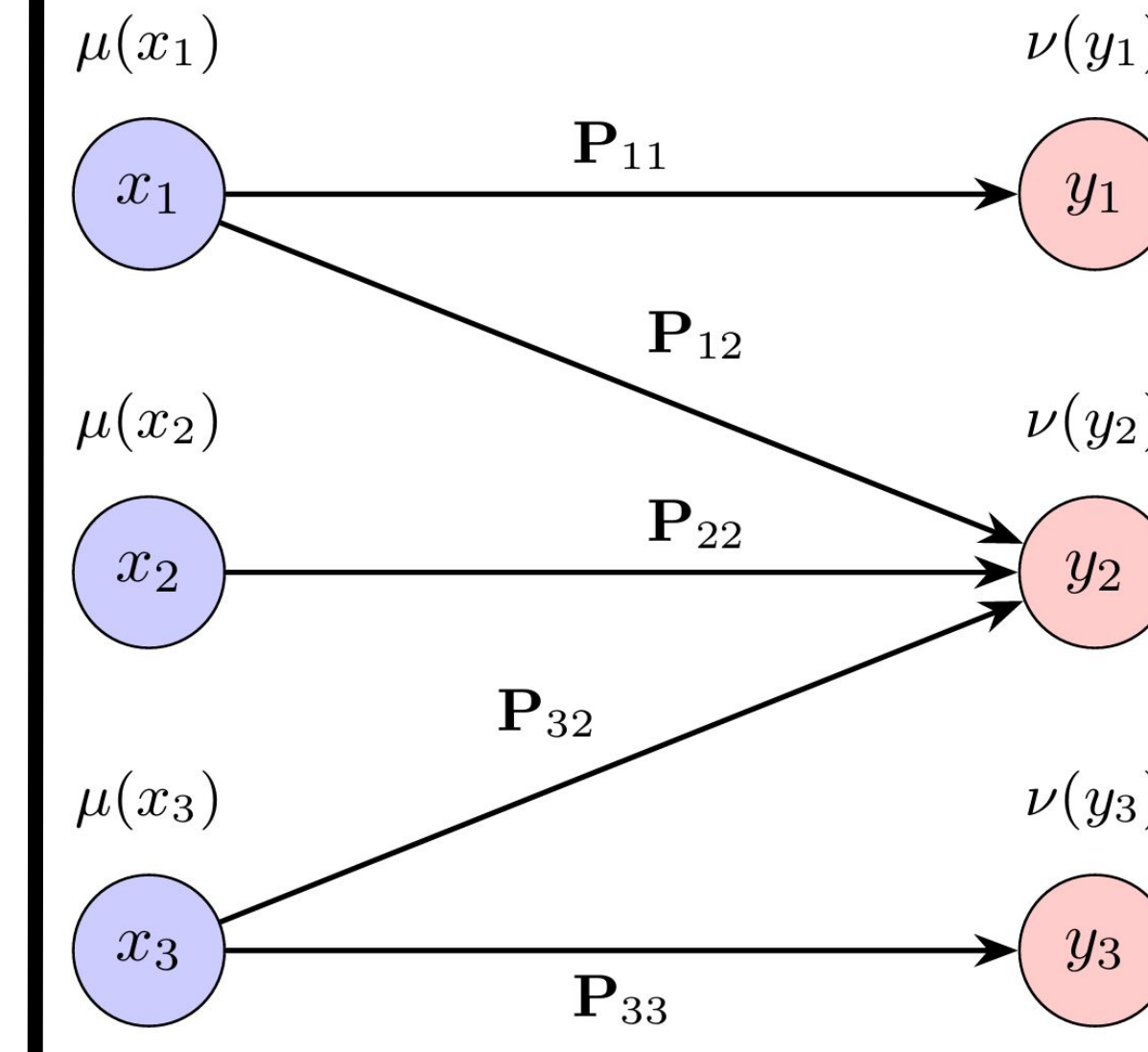


Fig 1: Possible optimal transport on system with three source points and three target points

Results

Theoretical Results

- **Space:** Our algorithm avoids instantiating a cost matrix and instead utilize a tree, which contains every point in each of $O(\log \Delta)$ layers, resulting in $O(n \log \Delta)$ space complexity
- **Time:** The time bottleneck of the algorithm is generating the FRT tree, which runs in $O(n^2)$ time.
- **Accuracy:** By using an FRT tree approximation, we incur a distortion of $O(\log n)$ such that $\mathbb{E}[c_T(u, v)] \leq O(\log n)c(u, v)$ however, we demonstrate through tests that the average error within the cost approximation is significantly less than $O(\log n)$

Method	Time	Space	Loss	Caveats
EMD	$O(n^3 \log n)$	$O(n^2)$	≈ 0	Computationally expensive.
Sinkhorn	$O(n^2/\epsilon^2)$	$O(n^2)$	$\leq +\epsilon$	$\epsilon \rightarrow 0$ incurs instability/slowdown.
HiRef	$O(n \log n)$	$O(n)$	≈ 0	Only solves for Monge maps.
HSTOT	$O(n^2)$	$O(n \log \Delta)$	$\leq \times O(\log n)$	Inaccurate for certain datasets.

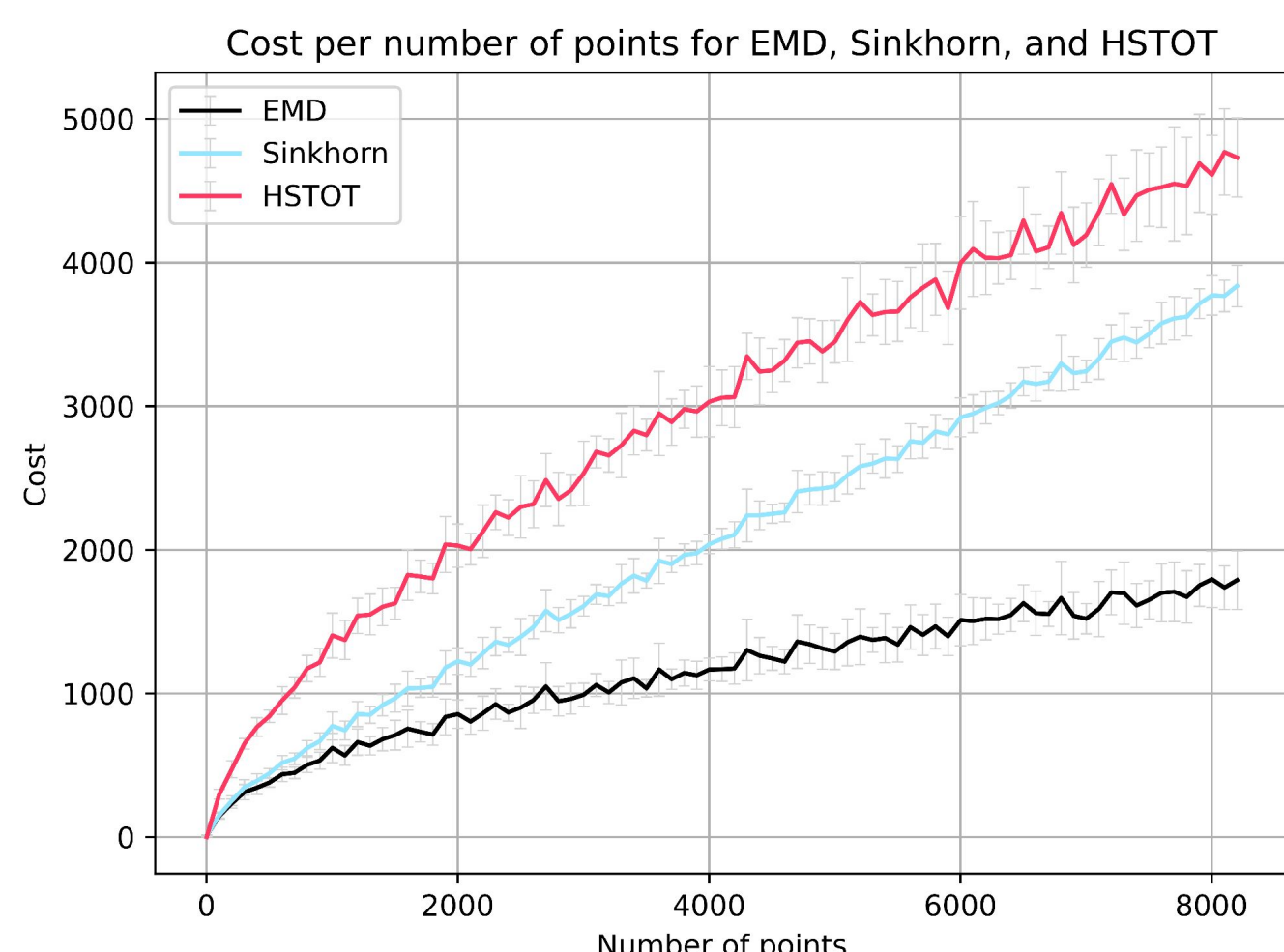
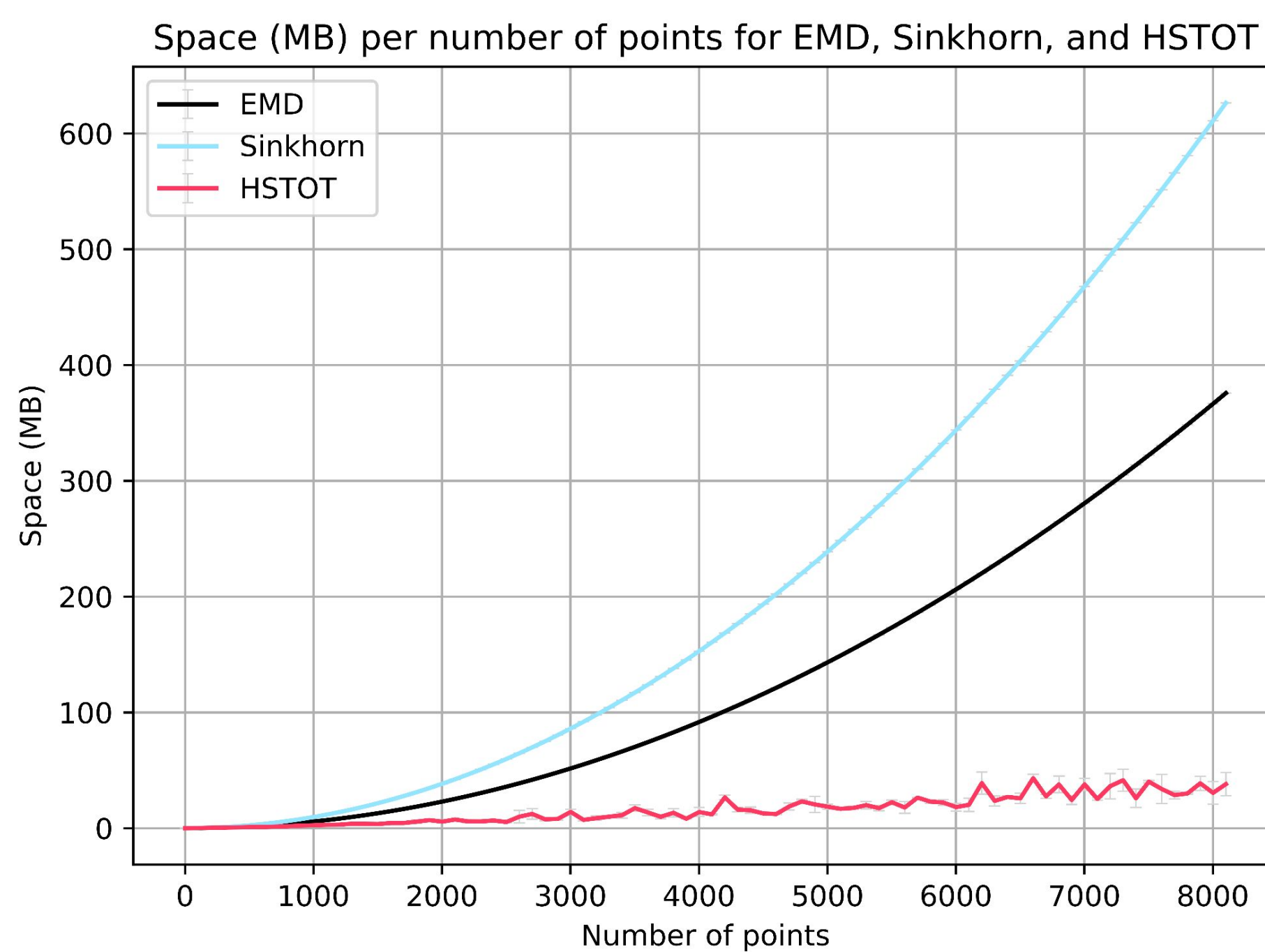
Table 1: Comparison of OT algorithms.

Empirical Results

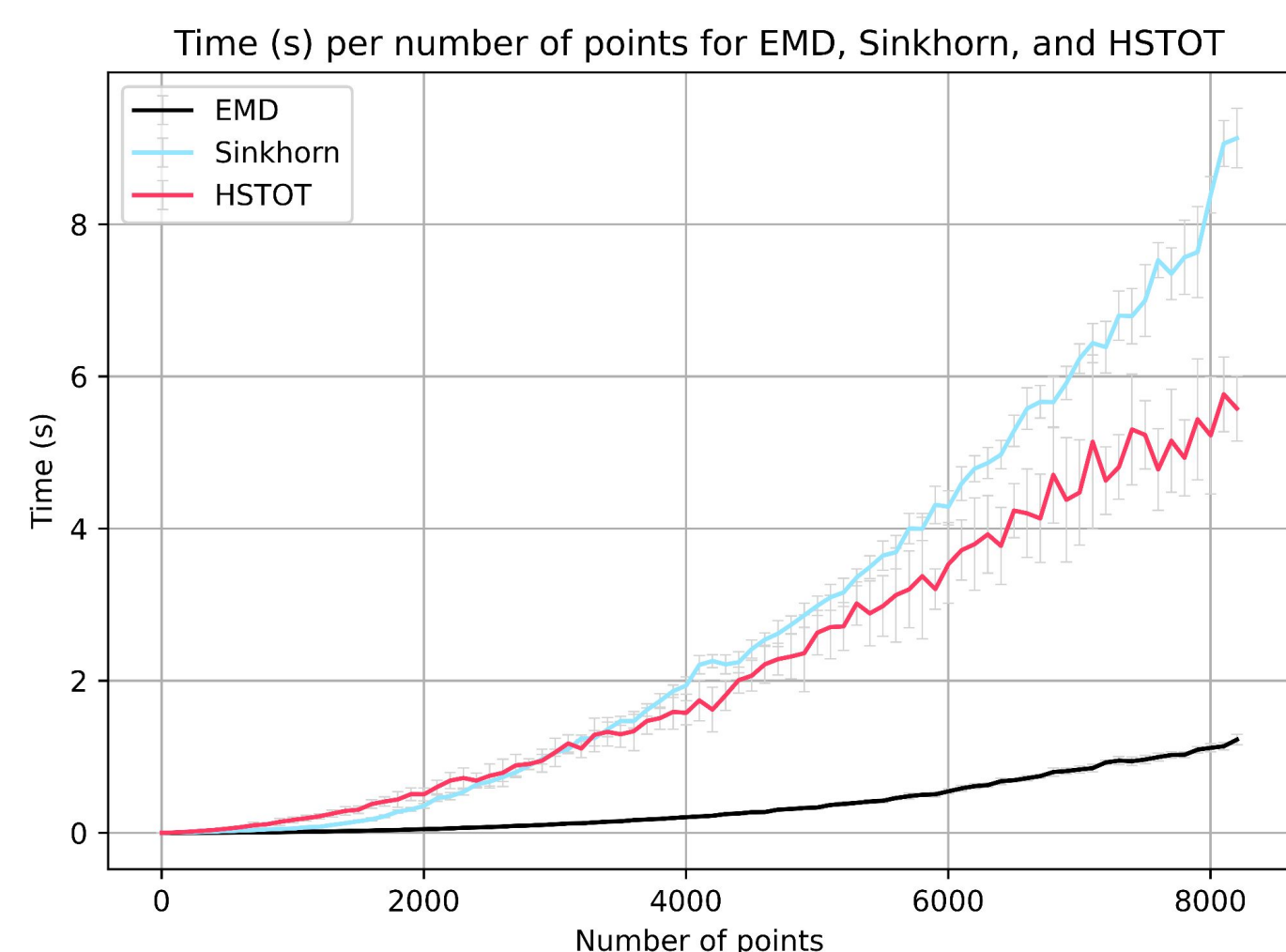
Fig. 3 (right): Space comparison between EMD, Sinkhorn, and our algorithm. EMD and Sinkhorn both use quadratic space, while our algorithm runs with significantly less speed.

Fig. 4 (bottom-left): Accuracy comparison between EMD, Sinkhorn, and our algorithm. Our algorithm provides a less accurate approximate than Sinkhorn, but when compared to actual costs (EMD), both algorithms' errors increase at the same rate.

Fig. 5 (bottom-right): Time comparison between EMD, Sinkhorn, and our algorithm. Our algorithm and Sinkhorn have comparable runtime. EMD has a cubic time complexity and is typically significantly slower in comparison, except in this case POT's EMD algorithm is implemented in C++ while Sinkhorn and ours are written in python, resulting in faster runtime.



Note: In Fig 3, 4, and 5, each algorithm was ran on 10 different randomly generated point clouds and weights, and the resulting measures were subsequently averaged to obtain means and standard deviations on the figures.



Acknowledgements

I would like to thank Dr. Alina Ene for her guidance and support throughout this project, who without this project would not have been possible. I would also like to thank the RISE program and BU for making this opportunity possible.

Algorithm/Methods

Algorithm

Our algorithm begins by generating an **FRT Tree Embedding** [4] for the data points, with the following modifications:

- If a node is only associated with one point, we stop generating its children and assign a difference value representing its weight
- After generating a node's children, remove the set of nodes it represents from memory
- Iterate over only node-represented points when generating its children, instead of every point

Then, we calculate the cost of optimal transport by pushing mass towards the root node on the tree and maximizing lower-cost matchings farther from the root node.

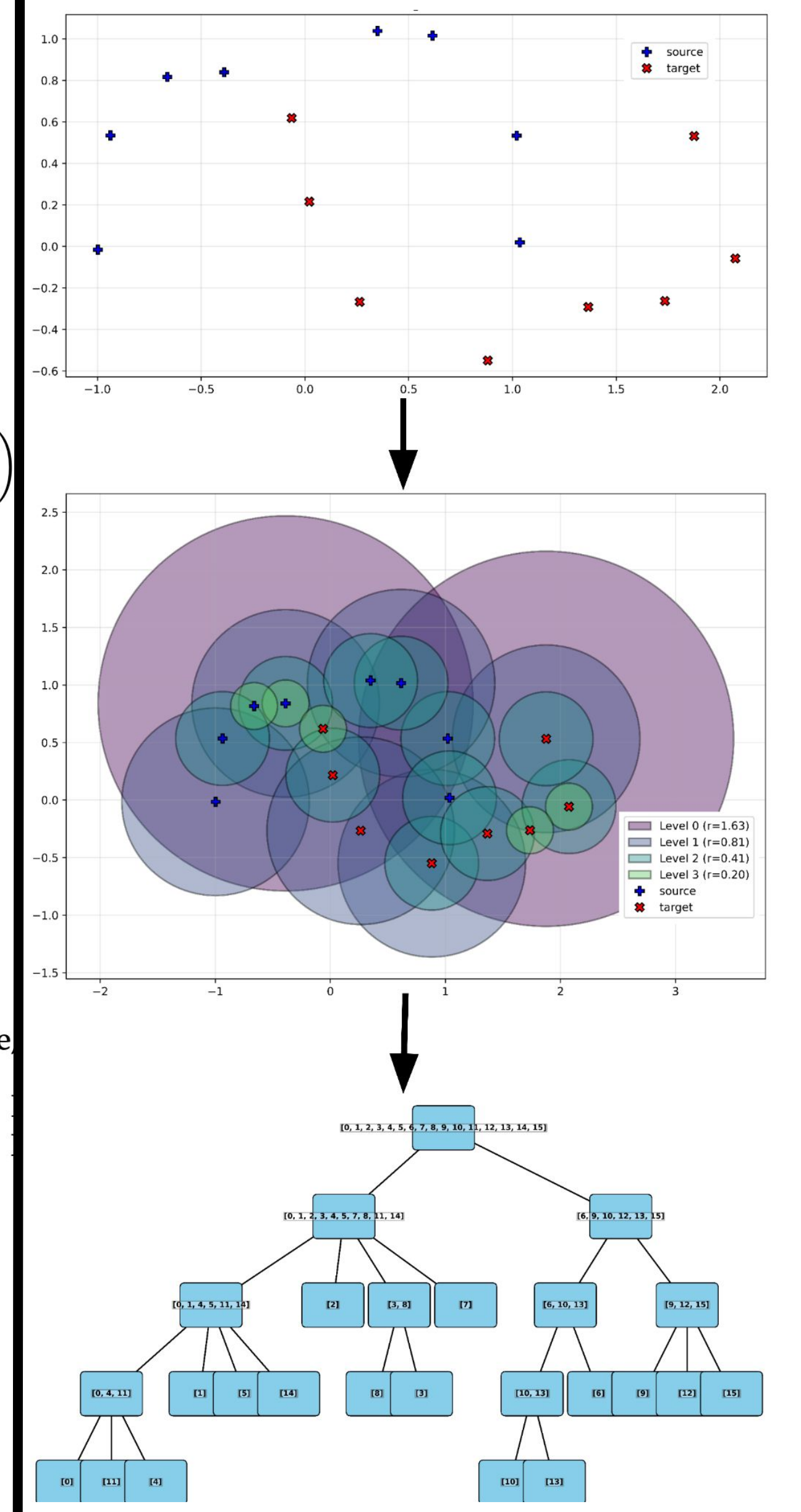
Algorithm 2 Mass-Pushing Cost Calculation

```

1: procedure CALCULATE_COST(node, π)
2:   if node.branches ≠ ∅ then
3:     sourceSet ← { b ∈ node.branches | b.difference > 0 }
4:     targetSet ← { b ∈ node.branches | b.difference < 0 }
5:     totalMass ← min(Σ s.difference, Σ (-t.difference))
6:     deferred ← []
7:     for all p1 ∈ sourceSet do
8:       for all p2 ∈ targetSet do
9:         mass' ← totalMass × (p1.difference / Σ sourceSet) × (-p2.difference / Σ targetSet)
10:        deferred.append(p1, p2, mass')
11:   for all (p1, p2, m) ∈ deferred do
12:     if p1.leaf ∧ p2.leaf then
13:       node.cost ← node.cost + dist(π[p1], π[p2]) × m
14:       p1.difference ← p1.difference - m
15:       p2.difference ← p2.difference + m
16:     else if p1.branch ∧ p2.leaf then
17:       p2.difference ← p2.difference + m
18:       PUSH_MASS(p1, m, true, p2, π)
19:     else if p2.branch ∧ p1.leaf then
20:       p1.difference ← p1.difference - m
21:       PUSH_MASS(p2, m, false, p1, π)
22:   else
23:     node.cost ← node.cost + PUSH_MASS_BETWEEN(node, p1, p2, m, π)
24:   for all b ∈ node.branches do
25:     CALCULATE_COST(b, π)

```

Fig. 2 (right): FRT Tree Embedding algorithm ran on a small two moons dataset. FRT greedily groups points using the ball function, as demonstrated by the middle image. The resulting tree (bottom image) groups nearby points into clusters with low inter-transport costs.



Comparison

- We compare our algorithm to Sinkhorn using the Python OT library with $\epsilon = 0.04 \times \text{mean}(\mathbf{C})$
- We test on a point cloud of n uniformly sampled source and target points with random Gaussian weights.
- We use the Earth Mover's Distance (EMD) as a baseline for approximation measurement

Discussion

Our algorithm attains a significantly lower space complexity but sacrifices accuracy when compared to Sinkhorn. Both algorithms have comparable speed.

Extensions

Unlike HiRef, our algorithm does not make restrictive assumptions about the data points and therefore can be applied to

- Partial Optimal Transport
 - If the source and target weights are not equivalent, leftover mass will be represented at the root node
- Non-Euclidean distance metrics
 - Does not rely on Euclidean distance, algorithm can be applied to other metrics as well

Future Work

- On some datasets (i.e. moon and blob), FRT may generate a tree with source and target points completely grouped together
 - Modifications could be made to FRT to avoid non-ideal groupings
- Testing on non-synthetic datasets

References

- [1] L. Kantorovitch. On the translocation of masses. *Management Science*, 5(1):1–4, 1958.
- [2] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transportation distances. 2013.
- [3] Peter Halmos, Julian Gold, Xinhao Liu, and Benjamin Raphael. Hierarchical refinement: Optimal transport to infinity and beyond. *In Forty-second International Conference on Machine Learning*, 2025.
- [4] Jittat Fakcharoenphol, Satish Rao, and Kunal Talwar. A tight bound on approximating arbitrary metrics by tree metrics. *Journal of Computer and System Sciences*, 69(3):485–497, 2004. Special Issue on STOC 2003.