

# Bridging Arbitrary and Tree Metrics via Differentiable Gromov Hyperbolicity

---

Pierre Houedry<sup>1</sup>   Nicolas Courty<sup>1</sup>   Florestan Martin-Baillon<sup>2</sup>  
Laetitia Chapel<sup>1</sup>   Titouan Vayer<sup>3</sup>

October 27, 2025

<sup>1</sup>IRISA, UMR 6074, CNRS, Université de Bretagne Sud

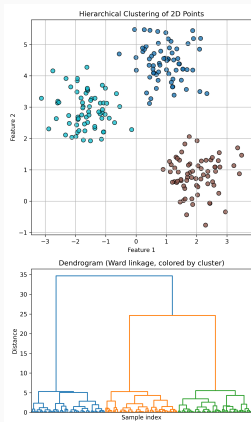
<sup>2</sup>IRMAR, UMR 6625, CNRS, Université de Rennes

<sup>3</sup>LIP, UMR 5668, CNRS, Univ. Lyon, ENS de Lyon, UCBL, CNRS, Inria

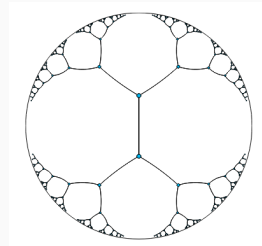
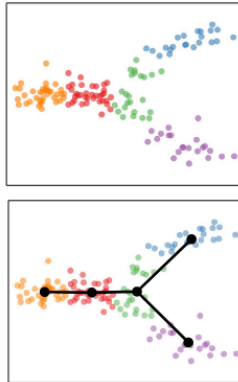
# Introduction

---

# Motivation



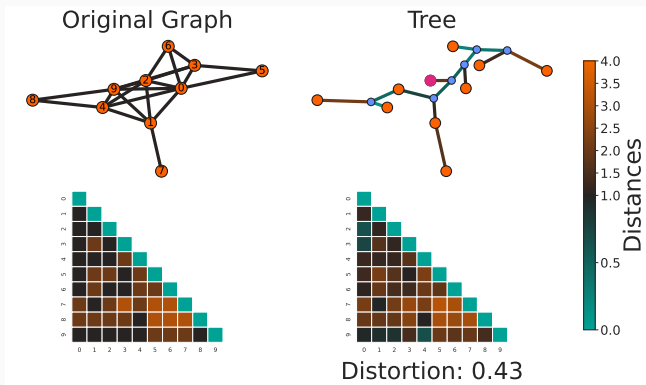
Hierarchical Clustering



Tree Embedding (Nickel and Kiela, 2017)

Single-cell trajectory inference (Street et al., 2018)

# Problem



A graph and a possible embedding of this graph into a tree.

## Problem Statement

- *Given:* A finite metric space  $(X, d_X)$  (e.g. pairwise distances between  $n$  items)
- *Find:*
  1. A weighted tree  $T$  with vertex set  $V(T)$
  2. An embedding  $\Phi : X \hookrightarrow V(T)$
- *Objective:* Minimize the  $\ell_\infty$  norm

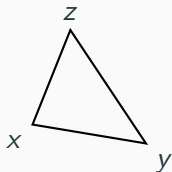
$$\max_{x,y \in X} |d_X(x,y) - d_T(\Phi(x), \Phi(y))|,$$

## Background on $\delta$ -hyperbolicity

---

# Geometric intuition

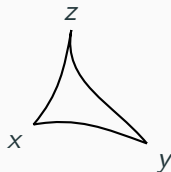
## Triangles



Euclidean

$$K = 0$$

$$\delta = \infty$$



Hyperbolic

$$K < 0$$

$$\delta > 0$$



Tree

$$K = -\infty$$

$$\delta = 0$$

Geometry

Curvature

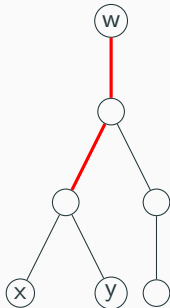
$\delta$ -Hyperbolicity

## Formal definition

### Gromov Product (Gromov, 1987)

Let  $(X, d_X)$  be a metric space and let  $x, y, w \in X$ . The *Gromov product* of  $x$  and  $y$  with respect to the basepoint  $w$  is defined as

$$(x|y)_w = \frac{1}{2} (d_X(x, w) + d_X(y, w) - d_X(x, y)).$$





### $\delta$ -hyperbolicity and Gromov hyperbolicity (Gromov, 1987)

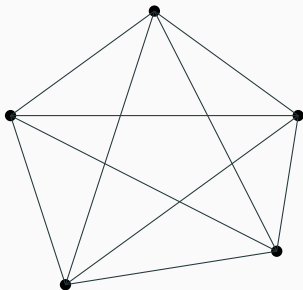
A metric space  $(X, d_X)$  is said to be  $\delta$ -hyperbolic if there exists  $\delta \geq 0$  such that for all  $x, y, z, w \in X$ , the Gromov product satisfies

$$(x|z)_w \geq \min \{(x|y)_w, (y|z)_w\} - \delta.$$

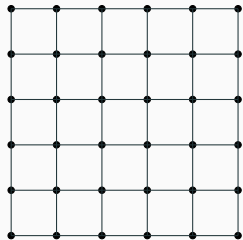
The *Gromov hyperbolicity*, denoted by  $\delta_X$ , is the smallest value of  $\delta$  that satisfies the above property. Consequently, every finite metric space  $(X, d_X)$  has a Gromov hyperbolicity equal to

$$\delta_X = \max_{x,y,z,w \in X} (\min \{(x|y)_w, (y|z)_w\} - (x|z)_w).$$

## Some examples



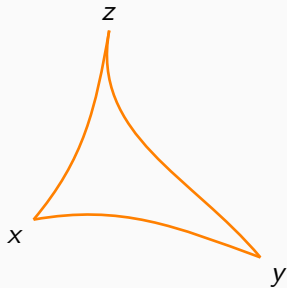
Complete graph  $K_5$ ,  $\delta_X = 0$



6x6 grid graph,  $\delta_X = 5$

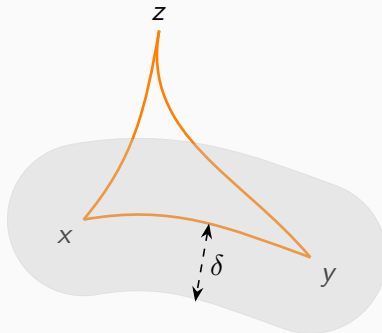
Two example graphs

## Geometrical view of the definition



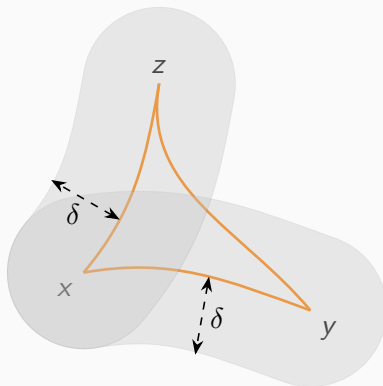
$\delta$ -thin triangle

## Geometrical view of the definition



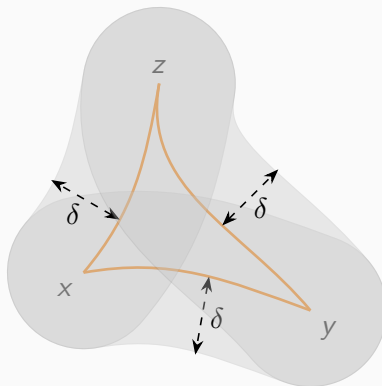
$\delta$ -thin triangle

## Geometrical view of the definition



$\delta$ -thin triangle

## Geometrical view of the definition



$\delta$ -thin triangle

## Geometrical view of the definition

## Tree Approximation of $\delta$ -Hyperbolic Spaces

### **Gromov's Theorem (Ghys et al., 1990, Ch. 2, §2, Thm. 12)**

Let  $(X, d_X)$  be a finite  $\delta$ -hyperbolic metric space over  $n$  points. For every  $w \in X$ , there exists a finite metric tree  $(T, d_T)$ , and an embedding  $\Phi_w : X \rightarrow T$  such that:



# Tree Approximation of $\delta$ -Hyperbolic Spaces

## **Gromov's Theorem (Ghys et al., 1990, Ch. 2, §2, Thm. 12)**

Let  $(X, d_X)$  be a finite  $\delta$ -hyperbolic metric space over  $n$  points. For every  $w \in X$ , there exists a finite metric tree  $(T, d_T)$ , and an embedding  $\Phi_w : X \rightarrow T$  such that:

1. Distance to the basepoint is preserved:

$$d_T(\Phi_w(x), \Phi_w(w)) = d_X(x, w).$$

# Tree Approximation of $\delta$ -Hyperbolic Spaces

## Gromov's Theorem (Ghys et al., 1990, Ch. 2, §2, Thm. 12)

Let  $(X, d_X)$  be a finite  $\delta$ -hyperbolic metric space over  $n$  points. For every  $w \in X$ , there exists a finite metric tree  $(T, d_T)$ , and an embedding  $\Phi_w : X \rightarrow T$  such that:

1. Distance to the basepoint is preserved:

$$d_T(\Phi_w(x), \Phi_w(w)) = d_X(x, w).$$

2. For every pair of points  $x, y \in X$ , the tree-embedding  $\Phi_w$  does not stretch distances and only contracts them by at most an additive factor of  $2\delta \log_2(n-2)$ .

$$d_X(x, y) - 2\delta \log_2(n-2) \leq d_T(\Phi_w(x), \Phi_w(y)) \leq d_X(x, y).$$

# Tree Approximation of $\delta$ -Hyperbolic Spaces

## Gromov's Theorem (Ghys et al., 1990, Ch. 2, §2, Thm. 12)

Let  $(X, d_X)$  be a finite  $\delta$ -hyperbolic metric space over  $n$  points. For every  $w \in X$ , there exists a finite metric tree  $(T, d_T)$ , and an embedding  $\Phi_w : X \rightarrow T$  such that:

1. Distance to the basepoint is preserved:

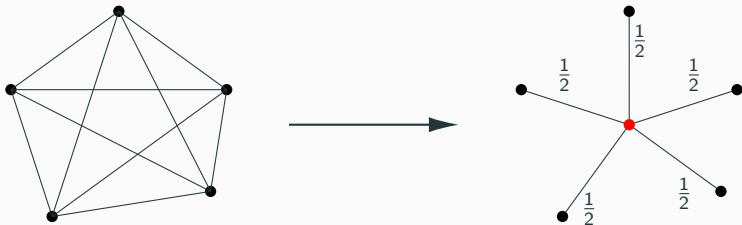
$$d_T(\Phi_w(x), \Phi_w(w)) = d_X(x, w).$$

2. For every pair of points  $x, y \in X$ , the tree-embedding  $\Phi_w$  does not stretch distances and only contracts them by at most an additive factor of  $2\delta \log_2(n-2)$ .

$$d_X(x, y) - 2\delta \log_2(n-2) \leq d_T(\Phi_w(x), \Phi_w(y)) \leq d_X(x, y).$$

$X$  can be embedded into a tree iff  $\delta_X = 0$  !!

## Example



- Gromov hyperbolicity requires checking **all point quadruples**  
 $\Rightarrow$  naive cost is  $O(n^4)$ .
- (Fournier et al., 2015) reduce base-point hyperbolicity to a (max, min) matrix product  $\Rightarrow$  Overall hyperbolicity in  $O(n^{3.69})$

# Structural Pruning via Far-apart Pairs

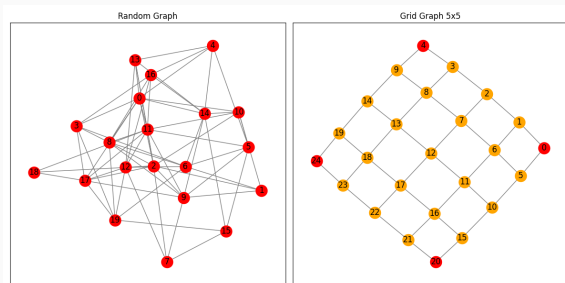


Illustration of far-apart pairs pruning

- Use far-apart vertex pairs to witness values.
- Prune by exploring only key quadruples.
- Delivers major speedups (Cohen et al., 2015; Coudert et al., 2022).

# Restatement of the problem

## Problem Statement

- *Given:* A finite metric space  $(X, d_X)$  (e.g. pairwise distances between  $n$  items)
- *Find:*
  1. A weighted tree  $T$  with vertex set  $V(T)$
  2. An embedding  $\Phi : X \hookrightarrow V(T)$
- *Objective:* Minimize the  $\ell_\infty$  norm

$$\max_{x,y \in X} |d_X(x,y) - d_T(\Phi(x), \Phi(y))|,$$

# Existing Methods Comparison

Method	HCC <sup>1</sup>	Gromov <sup>2</sup>	TR <sup>3</sup>	NJ <sup>4</sup>	LT <sup>5</sup>
Complexity	$O(n^2 \log n)$	$O(n^2)$	$O(n^2)$	$O(n^3)$	$O( E )$
Differentiable	✗	✗	✗	✗	✗
Requires a root	✓	✓	✗	✗	✓
Deterministic	✓	✓	✗	✓	✓
Non graph metrics	✓	✓	✓	✓	✗

Comparison of existing methods.

---

<sup>1</sup>(Yim and Gilbert, 2023)

<sup>2</sup>(Gromov, 1987)

<sup>3</sup>(Sonthalia and Gilbert, 2020)

<sup>4</sup>(Saitou and Nei, 1987)

<sup>5</sup>(Chepoi et al., 2012)



# Our Approach: Differentiable Hyperbolicity

- We introduce a **smooth, differentiable surrogate** of hyperbolicity.
- Enables **gradient-based optimization**.
- We propose a **batched approximation scheme**:
  - **Scalable** computation
  - **Independent of graph size**

**DeltaZero**

---

# Our approach

The set of  $n$ -points metrics:

$$\mathcal{D}_n := \left\{ \mathbf{D} \in \mathbb{R}_+^{n \times n} : \mathbf{D} = \mathbf{D}^\top, D_{ii} = 0, D_{ij} \leq D_{ik} + D_{kj} \right\}$$

is a closed, convex, polyhedral cone. Each metric  $d_X$  is represented by  $\mathbf{D}_X \in \mathcal{D}_n$ .

**Objective:** Find a metric  $\mathbf{D}_{X'}$  that is:

- **Close** to  $\mathbf{D}_X$  (via  $\ell_\infty$  distortion),
- Does not stretch  $\mathbf{D}_X$ ,
- **Tree-like** (via low hyperbolicity  $\delta_{X'}$ ).

# Our approach

The set of  $n$ -points metrics:

$$\mathcal{D}_n := \left\{ \mathbf{D} \in \mathbb{R}_+^{n \times n} : \mathbf{D} = \mathbf{D}^\top, D_{ii} = 0, D_{ij} \leq D_{ik} + D_{kj} \right\}$$

is a closed, convex, polyhedral cone. Each metric  $d_X$  is represented by  $\mathbf{D}_X \in \mathcal{D}_n$ .

**Objective:** Find a metric  $\mathbf{D}_{X'}$  that is:

- **Close** to  $\mathbf{D}_X$  (via  $\ell_\infty$  distortion),
- Does not stretch  $\mathbf{D}_X$ ,
- **Tree-like** (via low hyperbolicity  $\delta_{X'}$ ).

## Lagrangian type problem

$$\min_{\substack{\mathbf{D}_{X'} \in \mathcal{D}_n \\ \mathbf{D}_{X'} \leq \mathbf{D}_X}} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$$

## Minimizer + Gromov Embedding

$$\mathbf{D}_{X^*} = \arg \min_{\substack{\mathbf{D}_{X'} \in \mathcal{D}_n \\ \mathbf{D}_{X'} \leq \mathbf{D}_X}} \mathcal{L}_X(\mathbf{D}_{X'}, \mu),$$

$$(T^*, d_{T^*}) = \Phi_w(\mathbf{D}_{X^*}).$$

$$d_{T^*}(\Phi_w(x), \Phi_w(y)) \leq d_X(x, y)$$

$$d_X(x, y) - C_\mu \leq d_{T^*}(\Phi_w(x), \Phi_w(y))$$

## Gromov Embedding

$$(T, d_T) = \Phi_w(\mathbf{D}_X)$$

does not stretch distances

only contracts them

by at most  $2\delta \log_2(n-2)$

In particular if  $\mu \geq \frac{1}{2 \log_2(n-2)}$  then we have  $C_\mu < 2\delta \log_2(n-2)$ .

## How to make the objective function differentiable?

**Objective:**  $\min_{\substack{\mathbf{D}_{X'} \in \mathcal{D}_n \\ \mathbf{D}_{X'} \leq \mathbf{D}_X}} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$

# How to make the objective function differentiable?

**Objective:**  $\min_{\substack{\mathbf{D}_{X'} \in \mathcal{D}_n \\ \mathbf{D}_{X'} \leq \mathbf{D}_X}} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$

**Non-differentiable**

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty$$

**Differentiable**

$$\rightsquigarrow \|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2$$

# How to make the objective function differentiable?

**Objective:**  $\min_{\substack{\mathbf{D}_{X'} \in \mathcal{D}_n \\ \mathbf{D}_{X'} \leq \mathbf{D}_X}} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$

**Non-differentiable**

$$\begin{aligned} &\|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty \\ &\mathbf{D}_{X'} \leq \mathbf{D}_X \end{aligned}$$

**Differentiable**

$$\begin{aligned} &\rightsquigarrow \|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2 \\ &\rightsquigarrow \text{X} \end{aligned}$$



# How to make the objective function differentiable?

**Objective:**  $\min_{\mathbf{D}_{X'} \in \mathcal{D}_n} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$

**Non-differentiable**

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty$$

$\rightsquigarrow$

**Differentiable**

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2$$

# How to make the objective function differentiable?

**Objective:**  $\min_{\mathbf{D}_{X'} \in \mathcal{D}_n} \mathcal{L}_X(\mathbf{D}_{X'}, \mu) := \mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty + \delta_{X'}.$

**Non-differentiable**

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty$$

$$\delta_{X'}$$

**Differentiable**

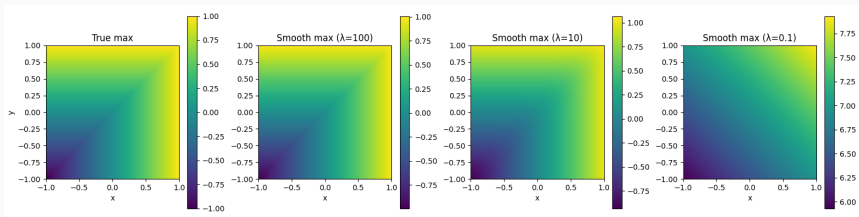
$$\rightsquigarrow \|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2$$

$$\rightsquigarrow ?$$

# Smooth min-max

We introduce a differentiable surrogate using log-sum-exp:

$$\text{LSE}_\lambda(\mathbf{x}) = \frac{1}{\lambda} \log \left( \sum_i e^{\lambda x_i} \right)$$



Smooth max visualization

$$\begin{array}{ll} \text{True } \delta_{\mathbf{X}} & \max(\min\{(x|y)_w, (y|z)_w\} - (x|z)_w) \\ \text{Smooth } \delta_{\mathbf{X}}^{(\lambda)} & \text{LSE}_{\lambda} \{ \text{LSE}_{-\lambda}((x|y)_w, (y|z)_w) - (x|z)_w \} . \end{array}$$

$$\begin{array}{ll} \text{True } \delta_X & \max(\min\{(x|y)_w, (y|z)_w\} - (x|z)_w) \\ \text{Smooth } \delta_X^{(\lambda)} & \text{LSE}_\lambda \{ \text{LSE}_{-\lambda}((x|y)_w, (y|z)_w) - (x|z)_w \} . \end{array}$$

## Bounds

$$\delta_X - \frac{\log 2}{\lambda} \leq \delta_X^{(\lambda)} \leq \delta_X + \frac{4 \log n}{\lambda}.$$

$\delta_X^{(\lambda)}$  still requires  $O(n^4)$  operations !

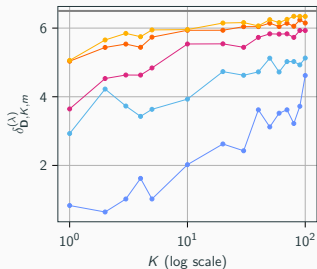
# Batched Approximation for Scalability

- We sample  $K$  subsets  $X_m^1, \dots, X_m^K \subset X$ , each of size  $m$ .
- Compute local estimates  $\delta_{X_m^i}^{(\lambda)}$  and aggregate:

$$\delta_{X,K,m}^{(\lambda)} = \text{LSE}_\lambda \left( \delta_{X_m^1}^{(\lambda)}, \dots, \delta_{X_m^K}^{(\lambda)} \right)$$

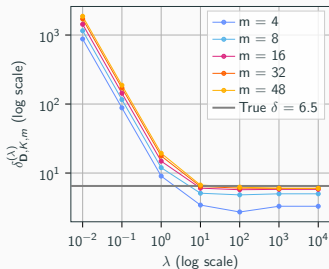
- Reduces complexity to  $O(K \cdot m^4)$ .
- Caveat: Rare, high-curvature configurations may be missed in random batches  $\Rightarrow$  potential underestimation.

# In practice



Evolution of the mean  $\delta_{D,K,m}^{(\lambda)}$  as a function of the number of batches  $K$ , with  $\lambda = 1000$ .

Mean of  $\delta_{D,K,m}^{(\lambda)}$  over 5 runs, computed on the CS-PHD dataset for different batch sizes  $m$ .



Estimation of  $\delta_{D,K,m}^{(\lambda)}$  as a function of  $\lambda$ , with  $K = 50$ .

# Final objective

Non-differentiable

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_\infty$$
$$\delta_{X'}$$

Differentiable

$$\|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2$$
$$\delta_{X,K,m}^{(\lambda)}$$

Cost

$$O(n^2)$$
$$O(Km^4)$$

$$\mathcal{L}_X(\mathbf{D}_{X'}, \mu)$$

$\rightsquigarrow$

$$\mu \|\mathbf{D}_X - \mathbf{D}_{X'}\|_2^2 + \delta_{X,K,m}^{(\lambda)}$$

$$O(Km^4 + n^2)$$



# Final Algorithm

---

## DELTAZERO

---

**Require:** A metric space  $(X, d_X)$ , root  $w \in X$ , learning rate  $\epsilon$ , batches  $K$ , size  $m$ , scale  $\lambda$ , regularization  $\mu$ , steps  $T$

- 1: Initialize  $\mathbf{D}_0 = \mathbf{D}_X$
  - 2: **for**  $t = 0$  to  $T - 1$  **do**
  - 3:     Sample  $K$  batches of  $m$  points
  - 4:     Compute  $\delta_{\mathbf{D}_t, K, m}^{(\lambda)}$
  - 5:      $\mathbf{G}_t = \nabla L_X(\mathbf{D}_t)$
  - 6:      $\tilde{\mathbf{D}}_{t+1} = \text{ADAMSTEP}(\mathbf{D}_t, \mathbf{G}_t, \epsilon)$
  - 7:      $\mathbf{D}_{t+1} = \text{FLOYDWARSHALL}(\tilde{\mathbf{D}}_{t+1}) \Rightarrow O(n^3)$
  - 8: **end for**
  - 9: **return**  $\text{GROMOVBEMBED}(\mathbf{D}_T, w)$
- 

Total cost of  $\mathbf{O}(\mathbf{T}(\mathbf{K}m^4 + n^3))$

**Setup:** 100 random roots for pivot-based methods, mean and std. of distortion reported. DELTAZERO hyperparameters selected via grid search:

- $\epsilon \in \{0.1, 0.01, 0.001\}$ ,
- $\mu \in \{0.1, 0.01, 1.0\}$ ,
- $\lambda \in \{0.01, 0.1, 1.0, 10.0\}$ ,
- $K \in \{100, 500, 1000, 3000, 5000\}$ ,
- $T = 1000$ ,
- $m = 32$ .

# Experiments: Unweighted Graphs

**Table 1:**  $\ell_\infty$  error on unweighted graphs (lower is better). Best result in bold, second-best underlined.

Dataset	C-ELEGAN	CS PhD	CORA	AIRPORT	WIKI
$n$	452	1025	2485	3158	2357
Diameter	7	28	19	12	9
NJ	<u>2.97</u>	16.81	13.42	4.18	6.32
TR	$5.90 \pm 0.72$	$21.01 \pm 3.34$	$16.86 \pm 2.11$	$10.00 \pm 1.02$	$9.97 \pm 0.93$
HCC	$4.31 \pm 0.46$	$23.35 \pm 2.07$	$12.28 \pm 0.96$	$7.71 \pm 0.72$	$7.20 \pm 0.60$
LT	$5.07 \pm 0.25$	$25.48 \pm 0.60$	<u><math>7.76 \pm 0.54</math></u>	<u><math>2.97 \pm 0.26</math></u>	<u><math>4.08 \pm 0.27</math></u>
Gromov	$3.33 \pm 0.45$	<u><math>13.28 \pm 0.61</math></u>	$9.34 \pm 0.53$	$4.08 \pm 0.27$	$5.54 \pm 0.49$
DELTAZERO	<b><math>1.87 \pm 0.08</math></b>	<b><math>10.31 \pm 0.62</math></b>	<b><math>7.59 \pm 0.38</math></b>	<b><math>2.79 \pm 0.15</math></b>	<b><math>3.56 \pm 0.20</math></b>
Improvement (%)	43.8%	22.3%	2.3%	6.0%	12.7%

## Experiments: General Metrics

**Table 2:**  $\ell_\infty$  error on general metrics (lower is better). Best result in bold, second-best underlined.

Dataset	ZEISEL	IBD
$n$	3005	396
Diameter	0.87	0.99
NJ	0.51	<u>0.90</u>
TR	$0.66 \pm 0.10$	$1.60 \pm 0.22$
HCC	$0.53 \pm 0.07$	$1.25 \pm 0.11$
LT	—	—
Gromov	<u><math>0.43 \pm 0.02</math></u>	$1.01 \pm 0.04$
DELTAZERO	<b><math>0.24 \pm 0.00</math></b>	<b><math>0.70 \pm 0.03</math></b>
Improvement (%)	44.1%	22.2%

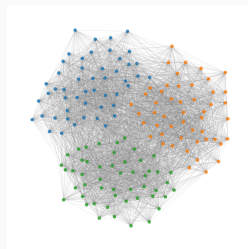
# Toy Application: Hierarchical Clustering

## Setup:

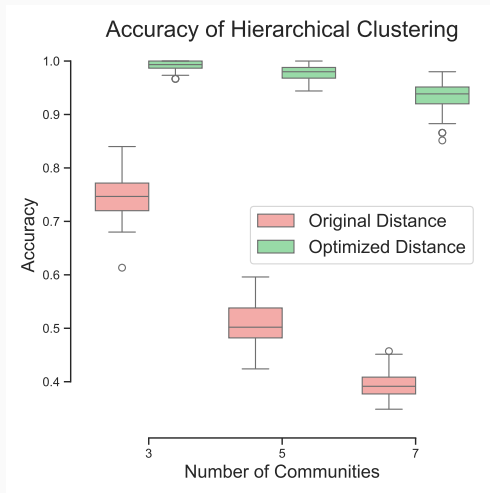
- Stochastic Block Model with  $N = 3, 5, 7$  communities, 50 nodes each
- $p_{in} = 0.6$ ,  $p_{out} = 0.2$

## Procedure:

1. Compute shortest-path distances and optimize them using DELTAZERO ( $\mu = 1$ ,  $\lambda = 100$ ).
2. Apply Ward's linkage on both original and optimized distance matrices.



# Toy Application: Hierarchical Clustering



Accuracy of hierarchical clustering with varying number of communities (3, 5, 7). Each setting is repeated 30 times.

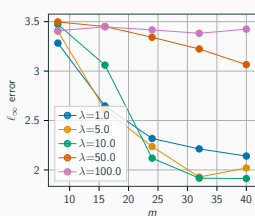
- Further theoretical investigations to:
  - Reduce the computational cost of  $\delta_{\mathbf{D},K,m}^{(\lambda)}$ .
  - Reduce the number of hyperparameters.
- Explore applications, such as:
  - Hierarchical clustering on real-world datasets,
  - Phylogenetic tree reconstruction,
  - Single-cell trajectory inference,
  - Hyperbolicity-aware learning ?

Thanks for listening!

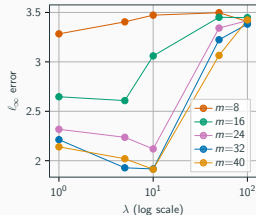
Any questions?



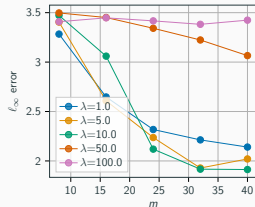
# Sensitivity Analysis



Effect of the distance regularization coefficient  $\mu$  on  $\ell_\infty$  distortion.



Impact of the log-sum-exp scale  $\lambda$  on  $\ell_\infty$  error across various batch sizes  $m$ .



Impact of the batch size  $m$  on  $\ell_\infty$  error for multiple  $\lambda$  values

Sensitivity analysis of optimization hyperparameters on the C-ELEGAN dataset. In each plot, non-varied hyperparameters are set to their optimal values from a prior grid search. Results are averaged over 5 runs, and distortion values averaged over 100 root samples per run.

## References

---

- Chepoi, V., Dragan, F. F., Estellon, B., Habib, M., Vaxès, Y., and Xiang, Y. (2012). Additive Spanners and Distance and Routing Labeling Schemes for Hyperbolic Graphs. *Algorithmica*, 62(3-4):713–732.
- Cohen, N., Coudert, D., and Lancin, A. (2015). On computing the gromov hyperbolicity. *ACM J. Exp. Algorithmics*, 20.
- Coudert, D., Nusser, A., and Viennot, L. (2022). Enumeration of far-apart pairs by decreasing distance for faster hyperbolicity computation. *ACM J. Exp. Algorithmics*, 27.
- Fournier, H., Ismail, A., and Vigneron, A. (2015). Computing the Gromov hyperbolicity of a discrete metric space. *Information Processing Letters*, 115(6):576–579.
- Ghys, F., De la Harpe, P., Oesterlé, J., and Weinstein, A. (1990).