

Bridging Arbitrary and Tree Metrics via Differentiable Gromov Hyperbolicity



Pierre Houedry



Nicolas Courty



Florestan Martin-Baillon



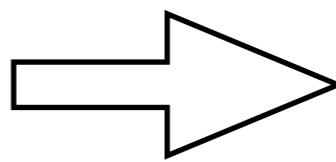
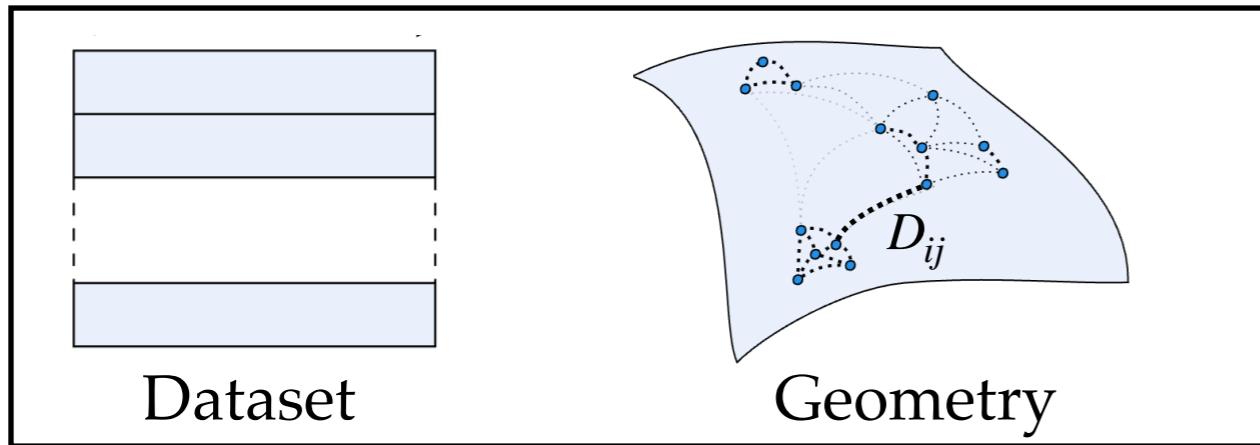
Laetitia Chapel



Titouan Vayer

Motivations

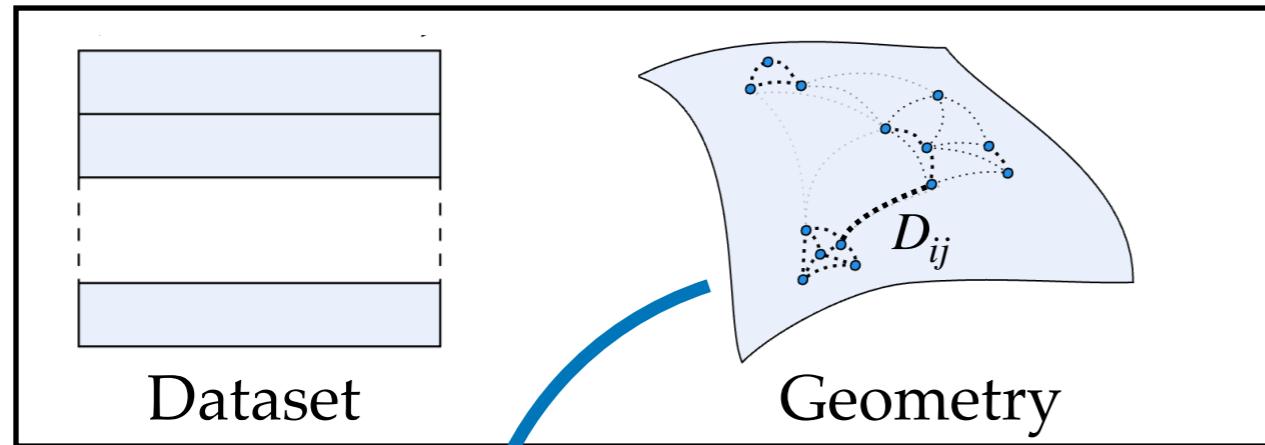
◆ Representation learning



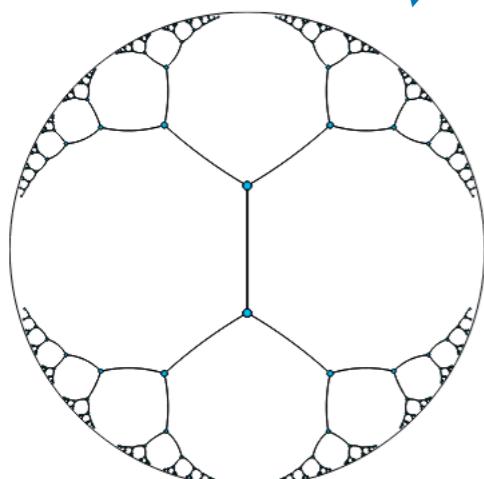
« Interesting » embedding space

Motivations

◆ Representation learning



« Interesting » embedding space

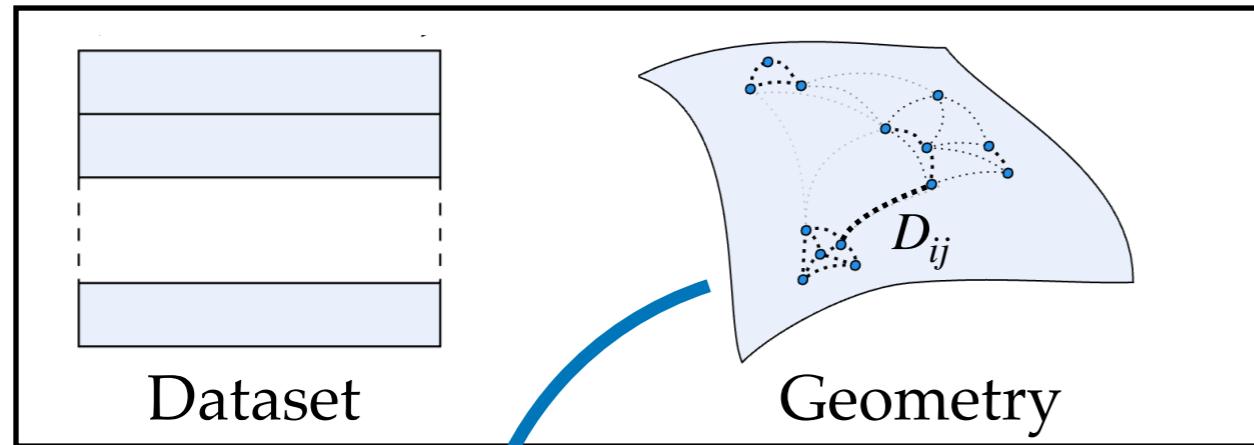


Trees embedding

◆ Hierarchical representations

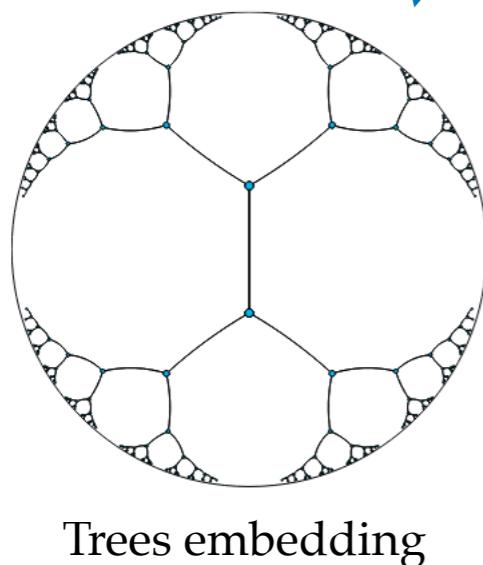
Motivations

◆ Representation learning

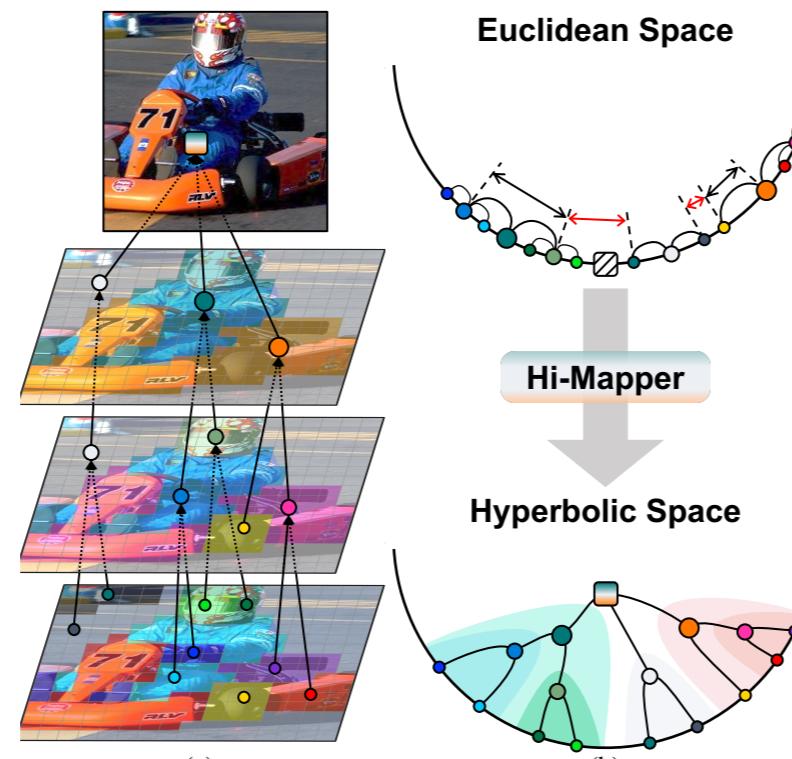


« Interesting » embedding space

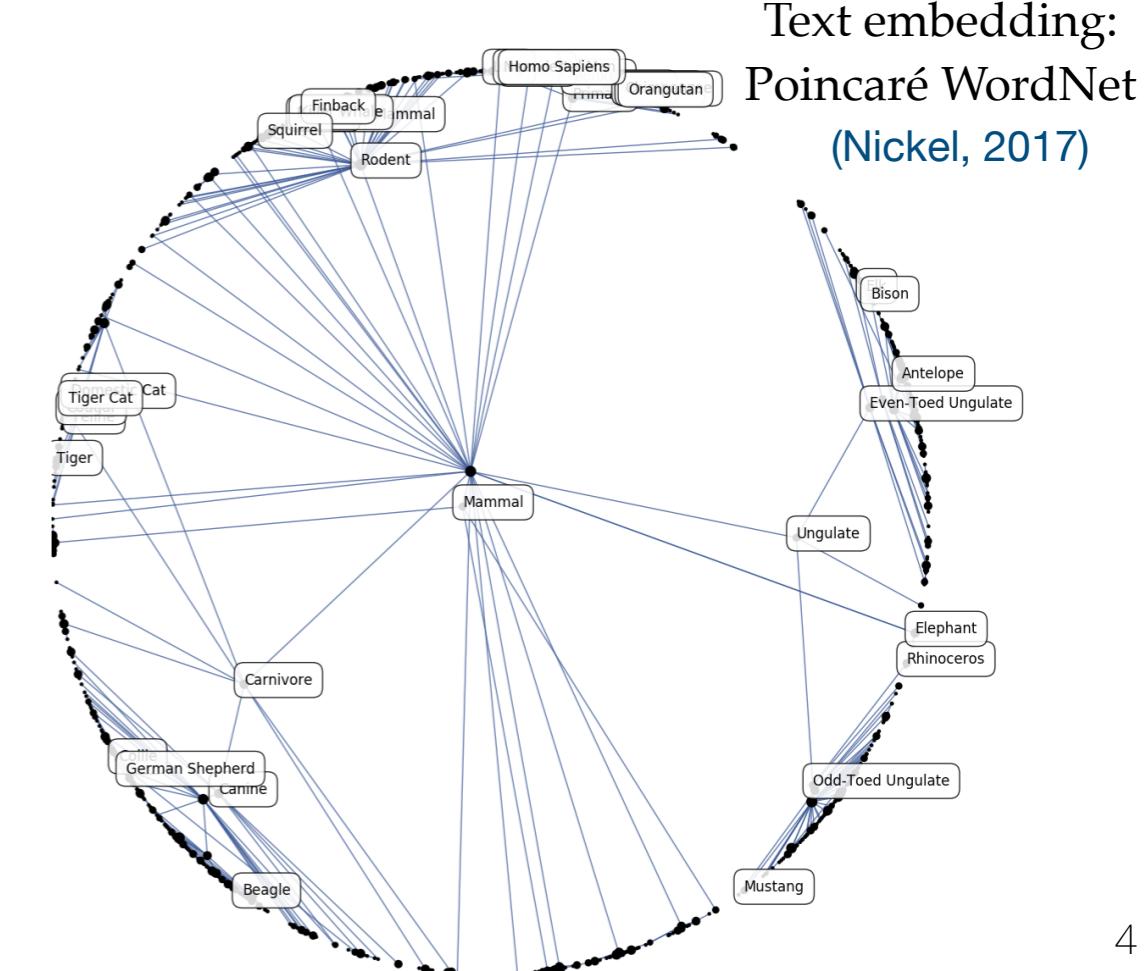
◆ Hierarchical representations



Trees embedding

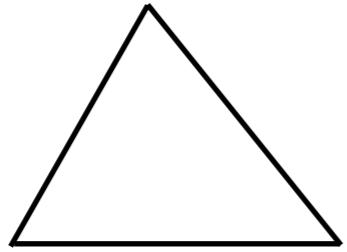


Computer vision (Kwon, 2024)

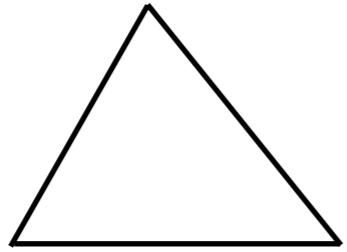
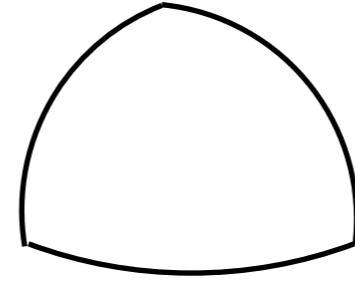


Text embedding:
Poincaré WordNet
(Nickel, 2017)

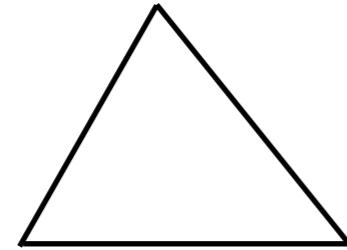
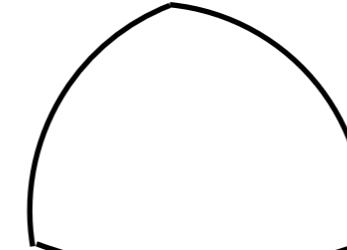
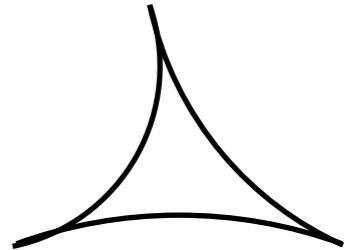
Geometries

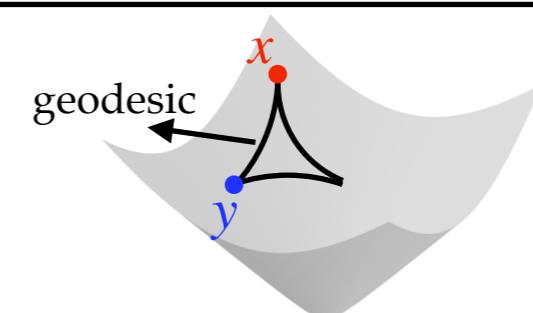
	Euclidean
Curvature	0
Parallel lines	1
Triangles are	normal
Shape	

Geometries

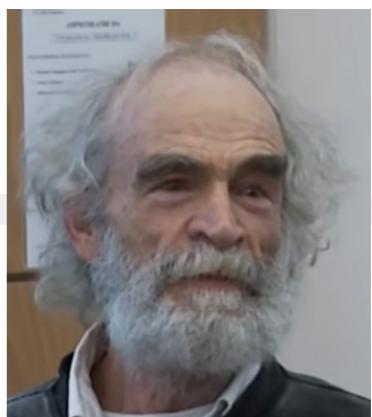
	Euclidean	Spherical
Curvature	0	> 0
Parallel lines	1	0
Triangles are	normal	thick
Shape		

Geometries

	Euclidean	Spherical	Hyperbolic
Curvature	0	> 0	< 0
Parallel lines	1	0	∞
Triangles are	normal	thick	thin
Shape			

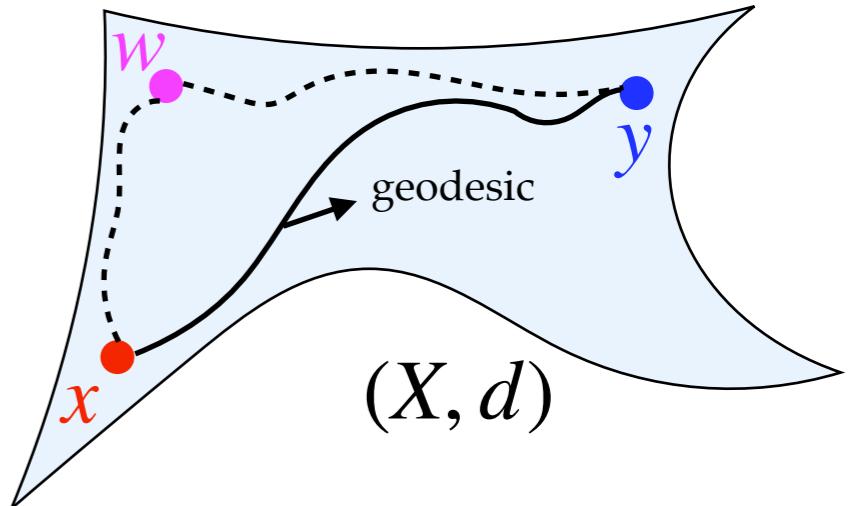


Gromov hyperbolicity



♦ Gromov product

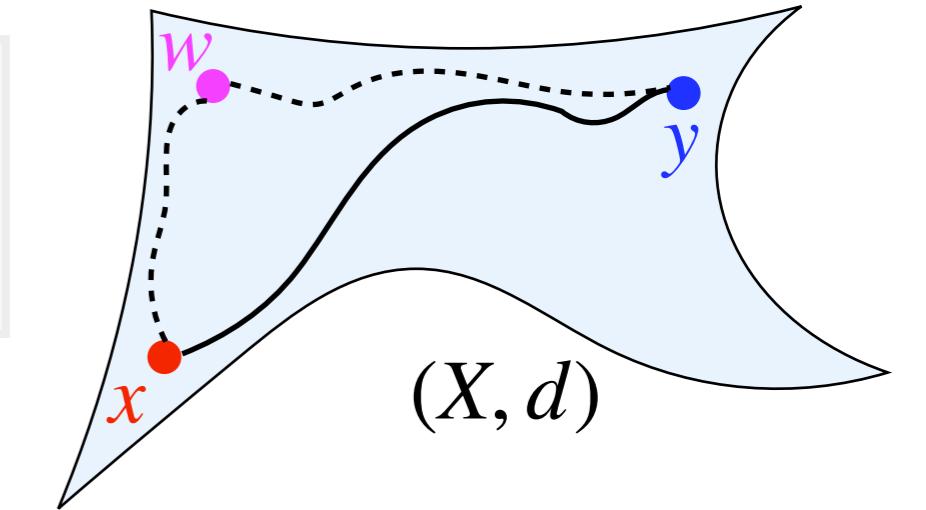
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



Gromov hyperbolicity

◆ Gromov product

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



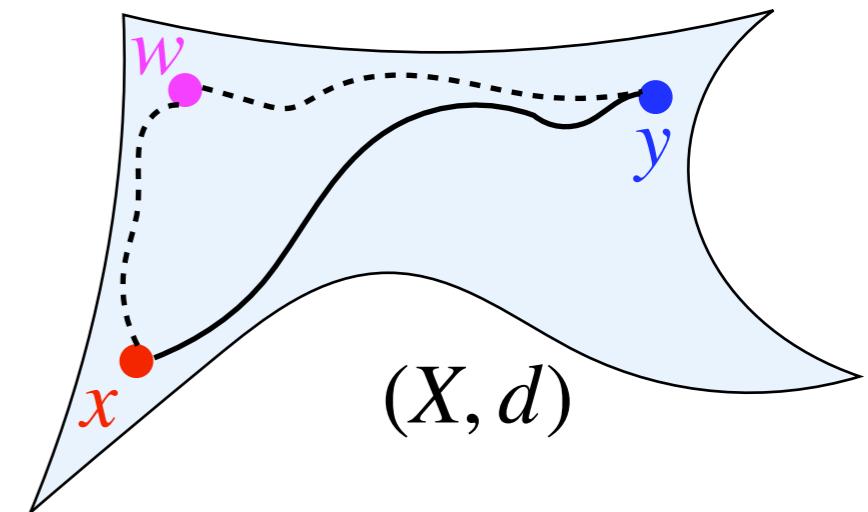
$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

$$\text{car } (x|y)_w \leq \frac{1}{2} (d(x,y) + d(y,w) + d(w,x) - d(x,y))$$

Gromov hyperbolicity

◆ Gromov product

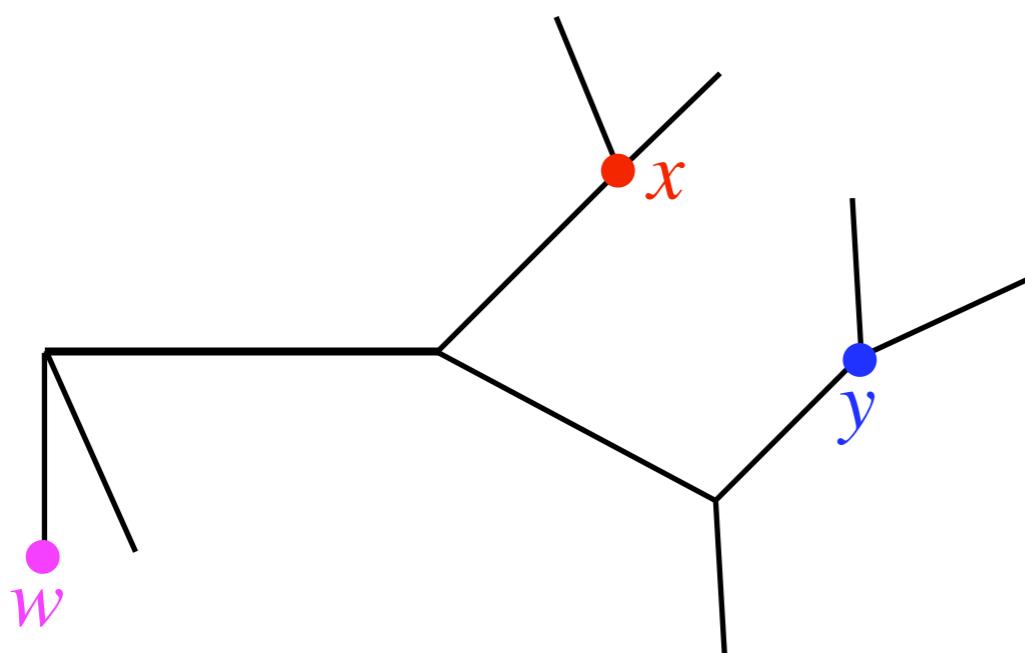
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Interpretation

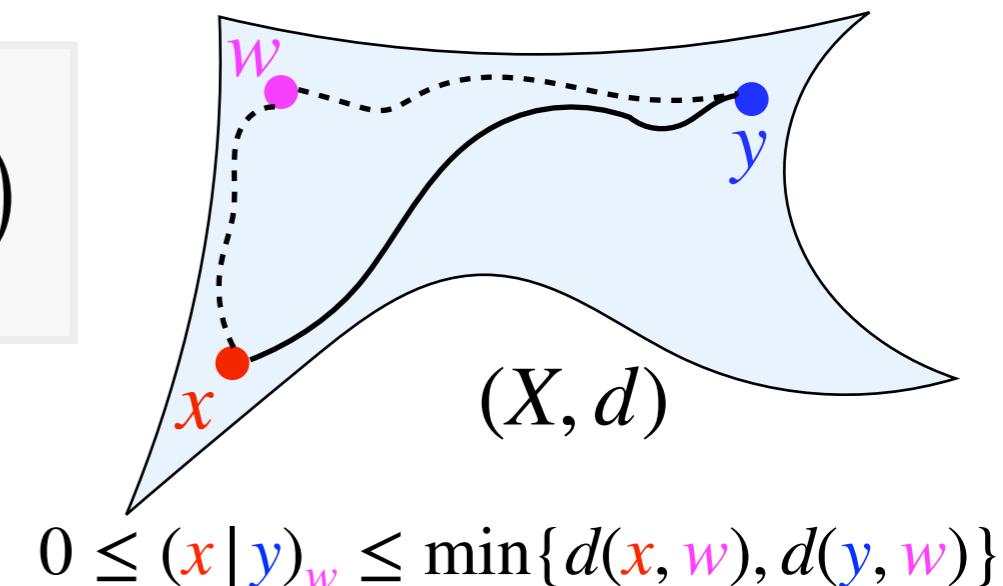
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

◆ Gromov product

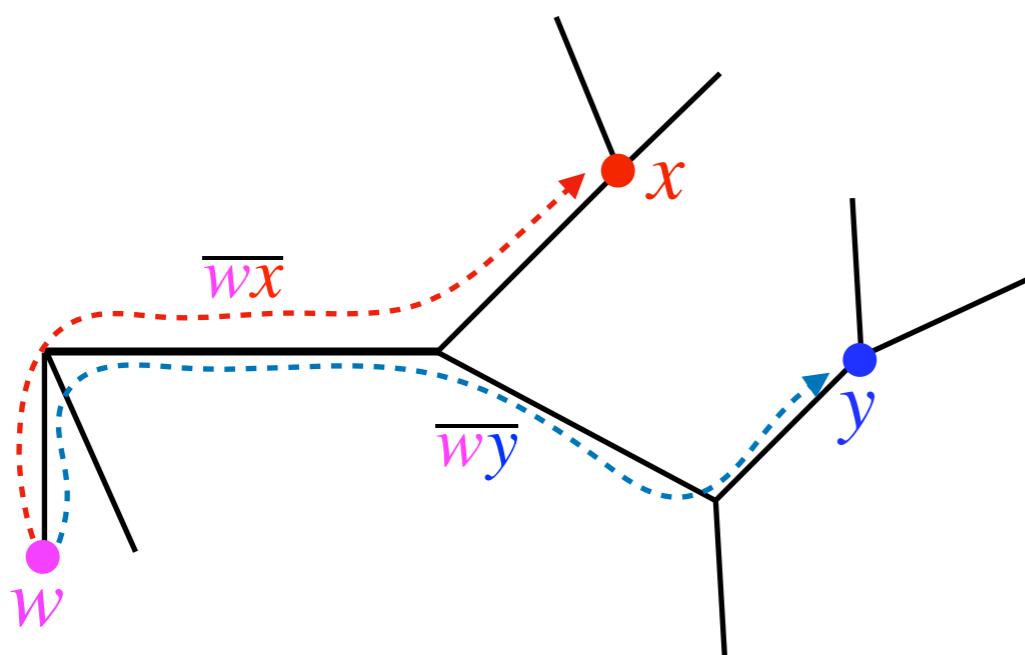
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



◆ Interpretation

$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

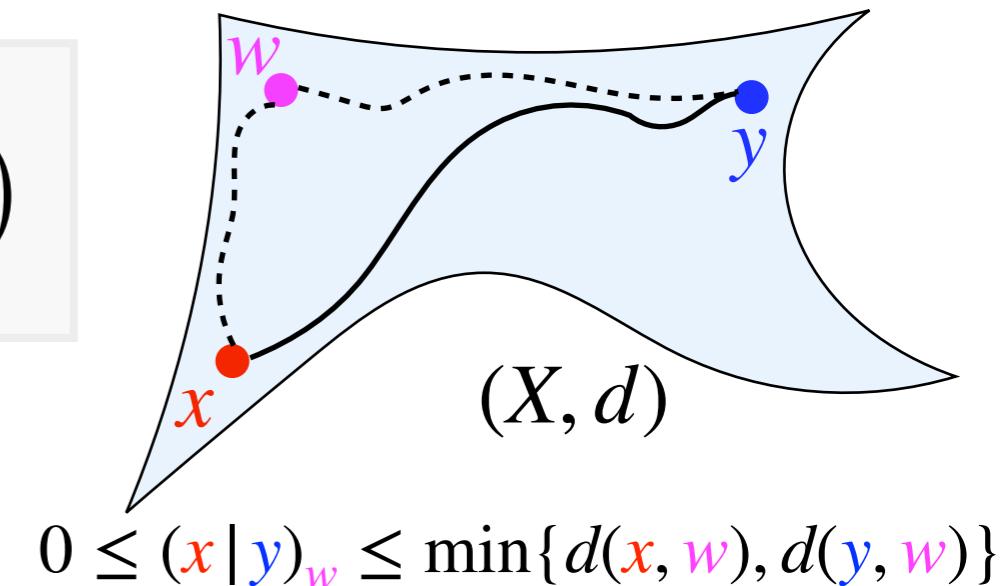
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

◆ Gromov product

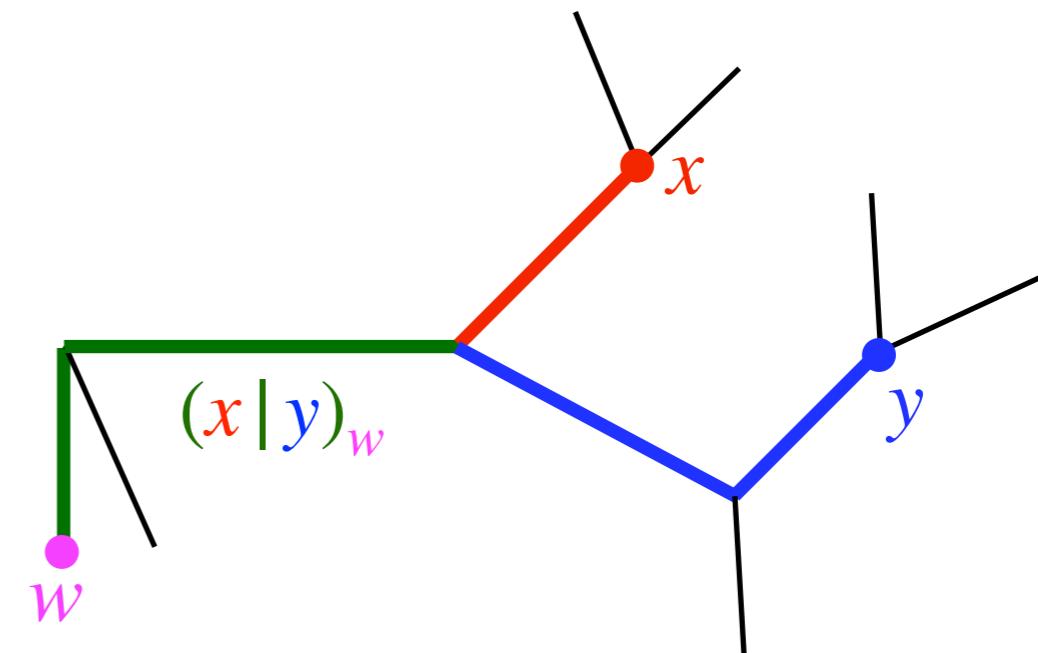
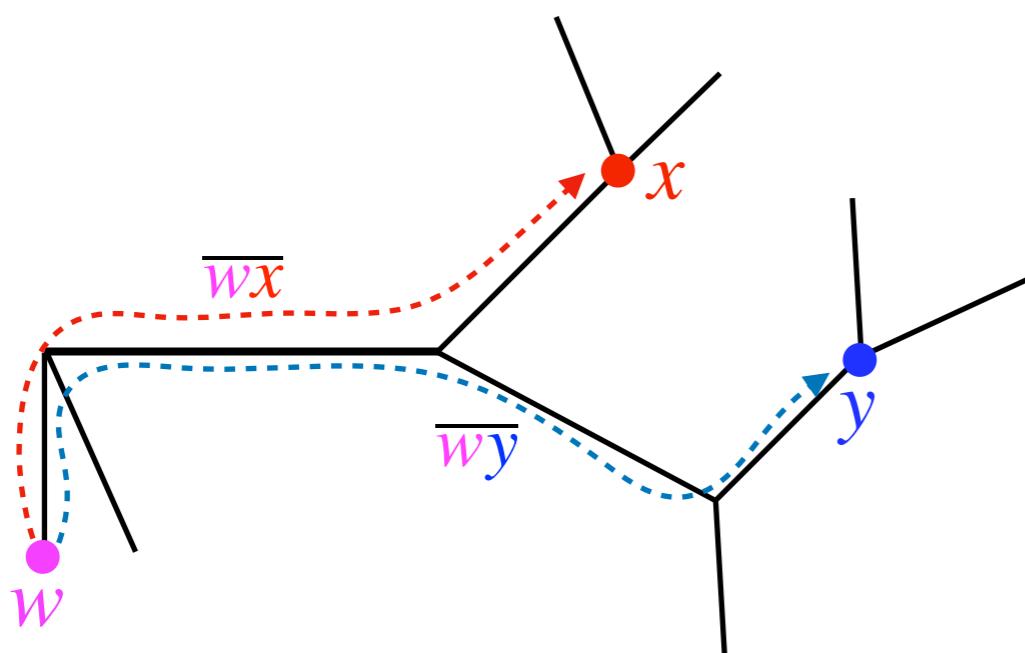
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Interpretation

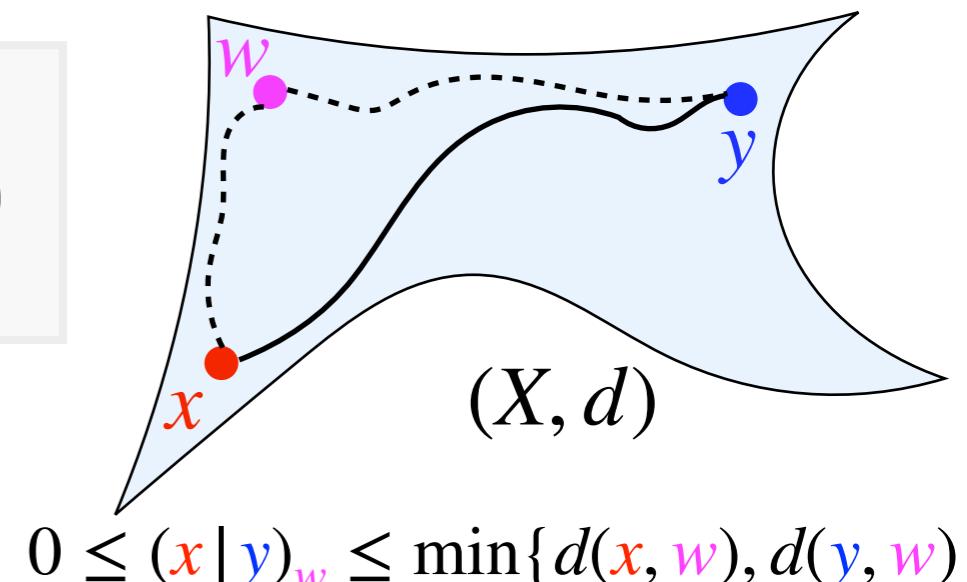
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

◆ Gromov product

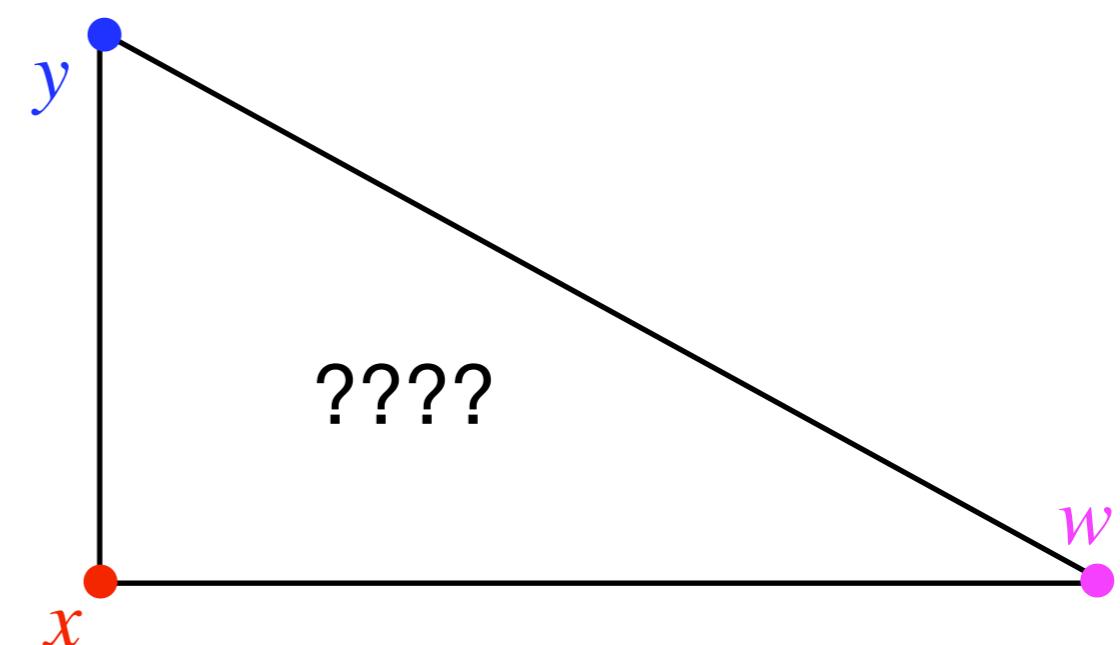
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Interpretation

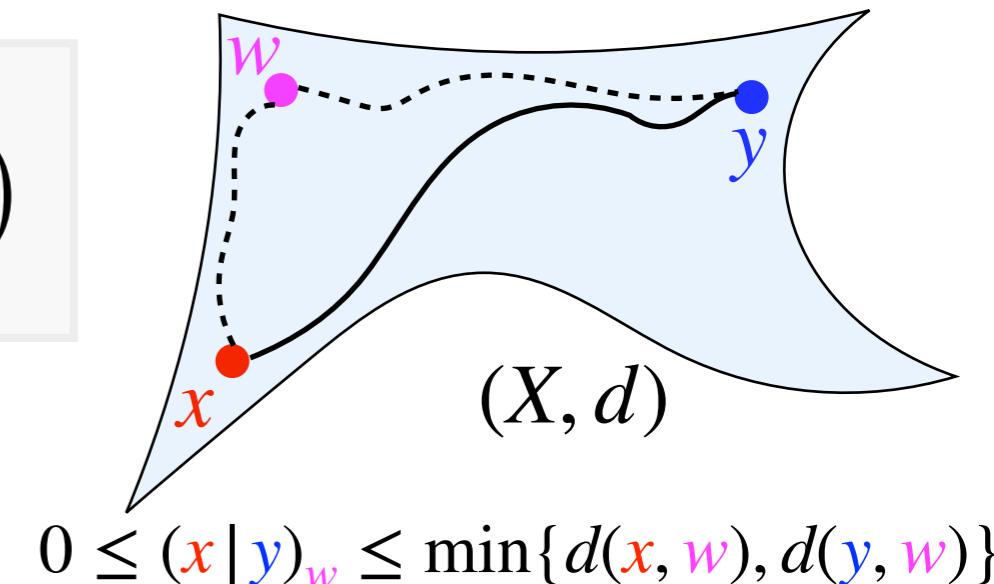
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

◆ Gromov product

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



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Interpretation

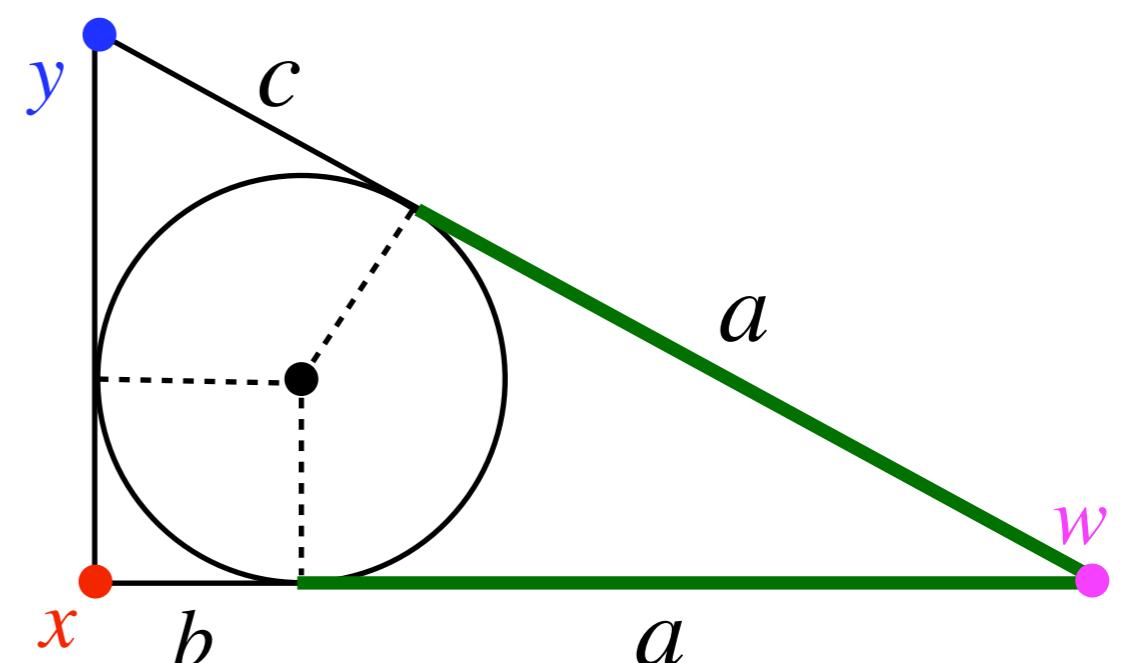
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging

$$\exists !(a, b, c) \geq 0$$

$$d(x,w) = a + b$$

$$d(x,y) = b + c$$

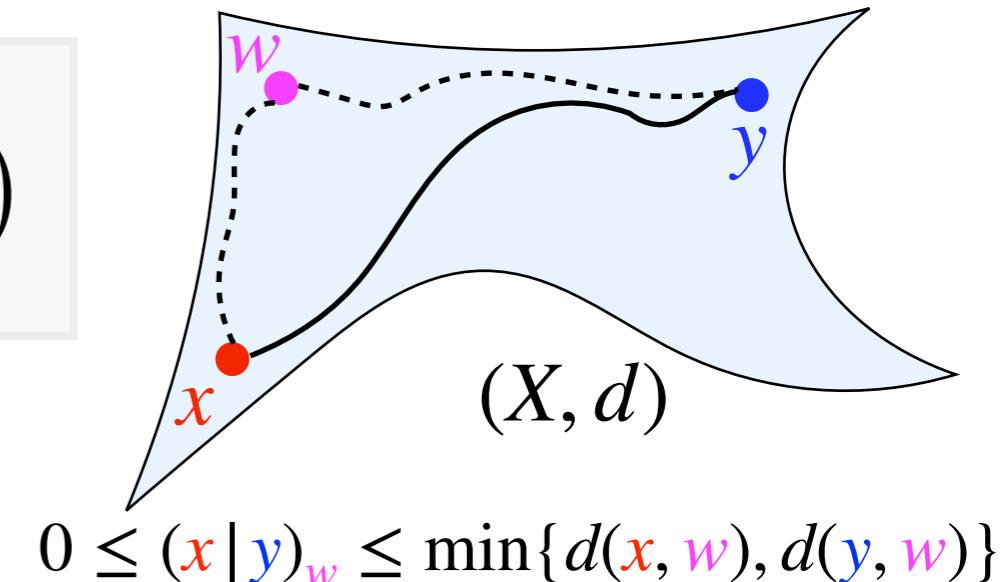
$$d(y,w) = a + c$$



Gromov hyperbolicity

◆ Gromov product

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



◆ Interpretation

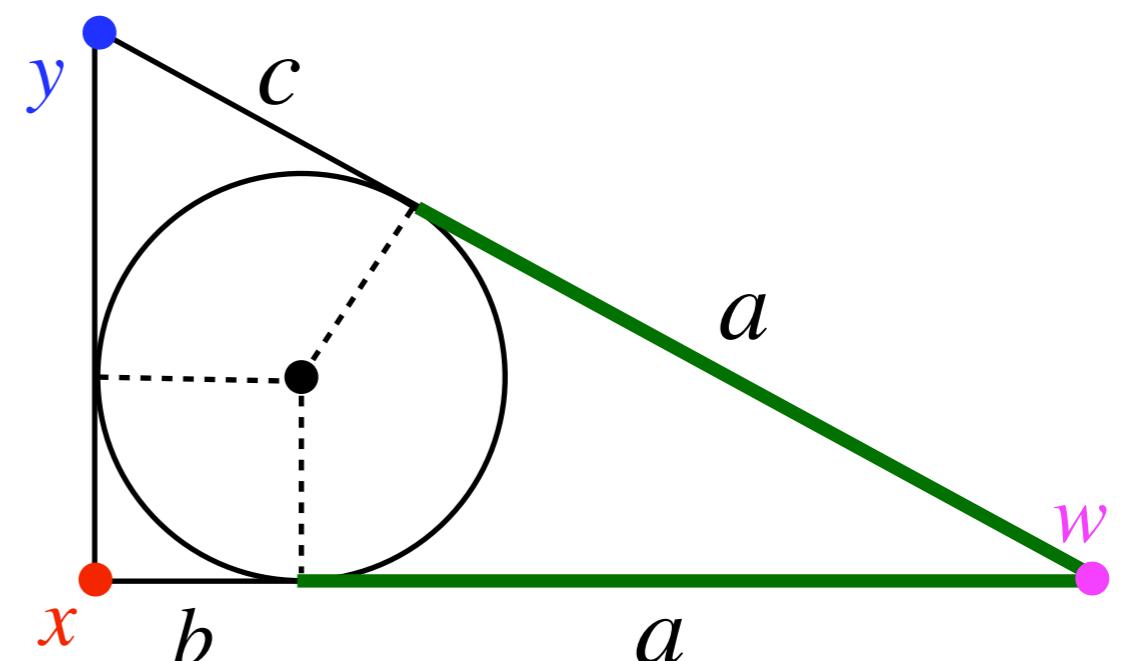
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging

$$\exists ! (a, b, c) \geq 0$$

$$d(x, w) = a + b$$

$$d(x, y) = b + c \implies a = (x|y)_w$$

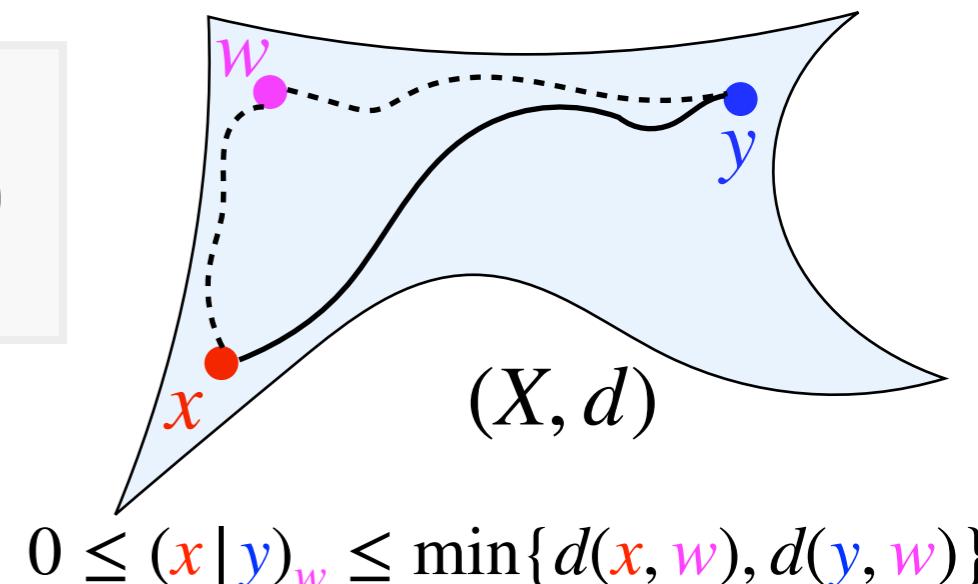
$$d(y, w) = a + c$$



Gromov hyperbolicity

◆ Gromov product

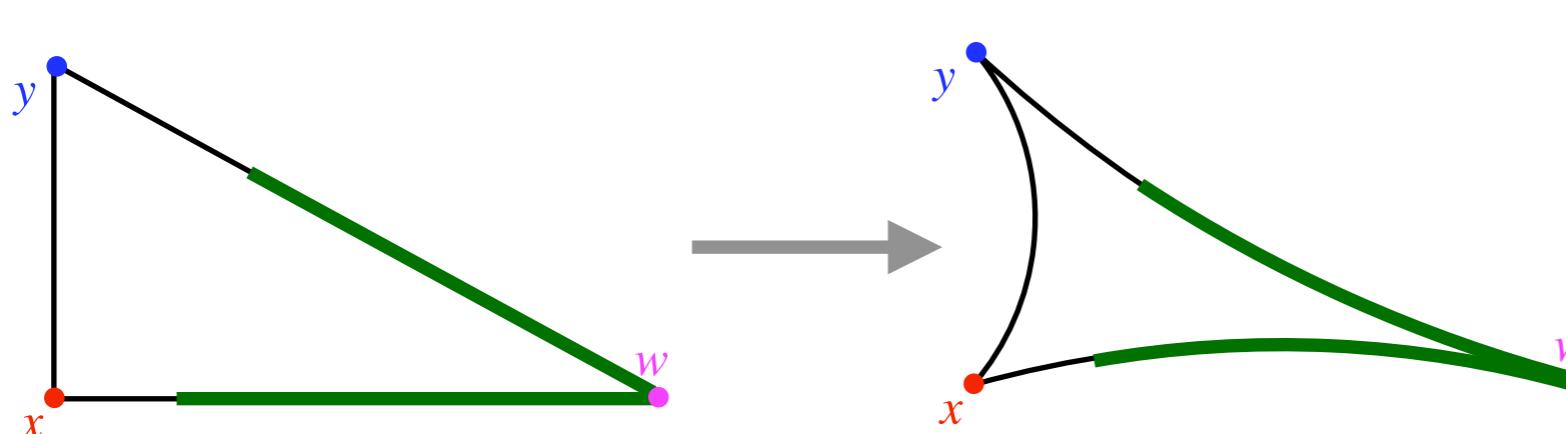
$$(x|y)_w = \frac{1}{2} (d(x,w) + d(y,w) - d(x,y))$$



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Interpretation

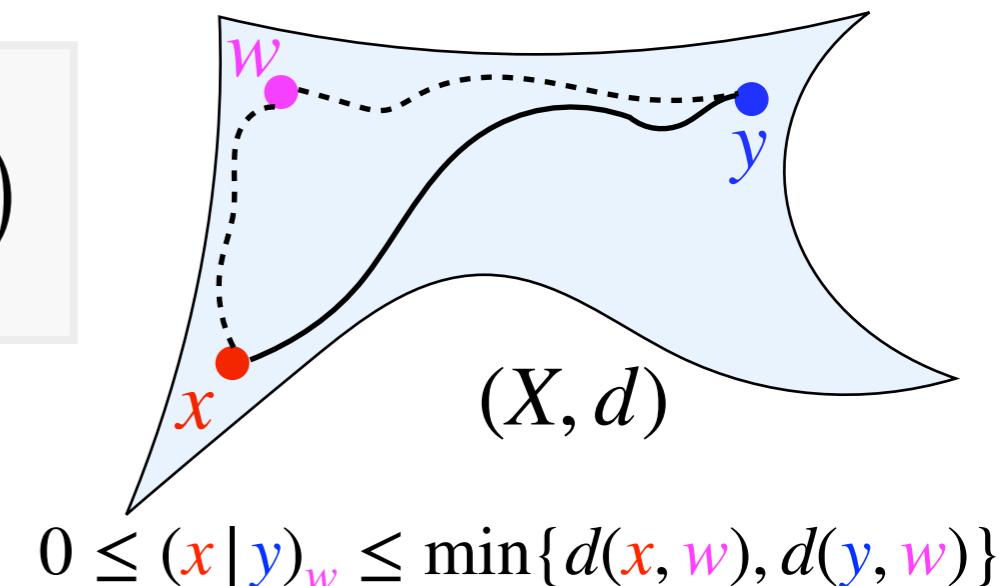
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

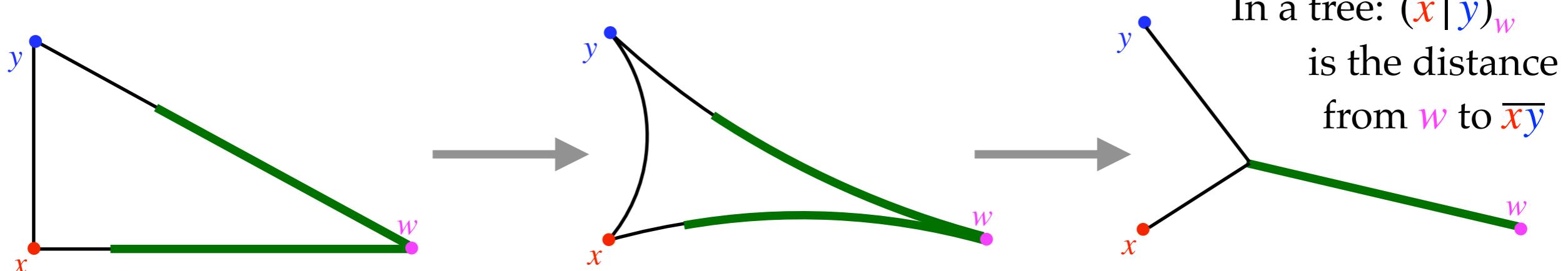
◆ Gromov product

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



◆ Interpretation

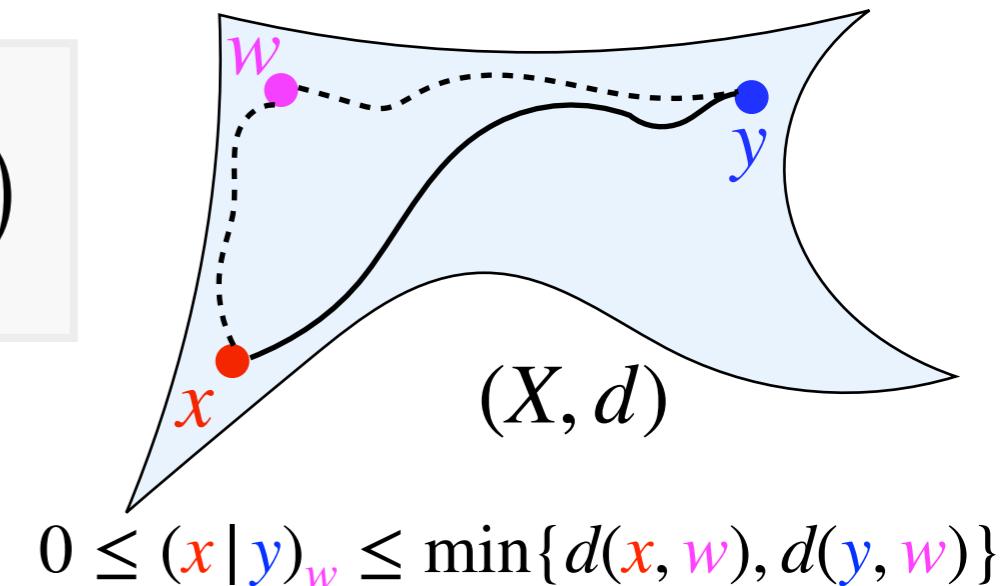
$(x|y)_w$ measures how long geodesics \overline{wx} and \overline{wy} travel the same distance before diverging



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

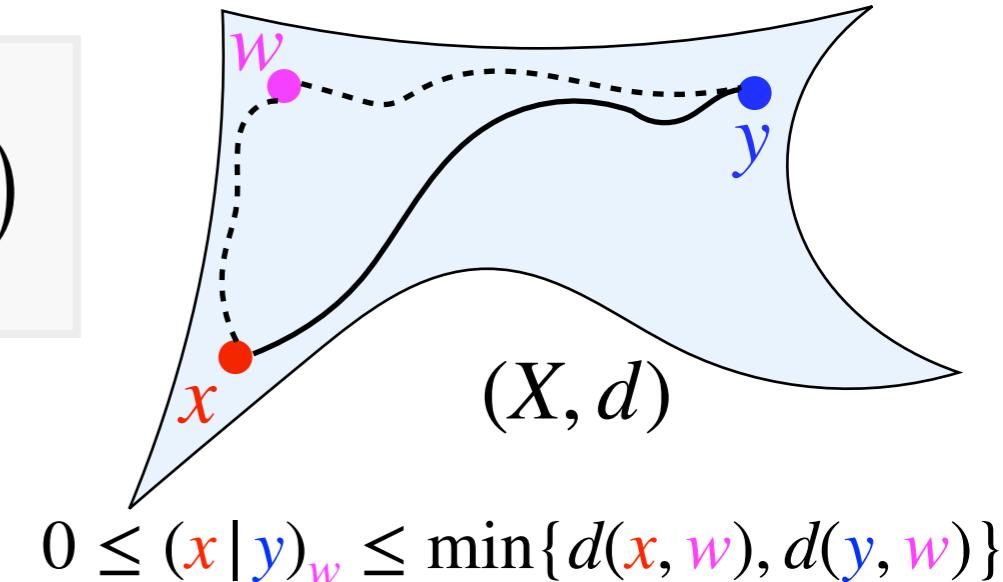
(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

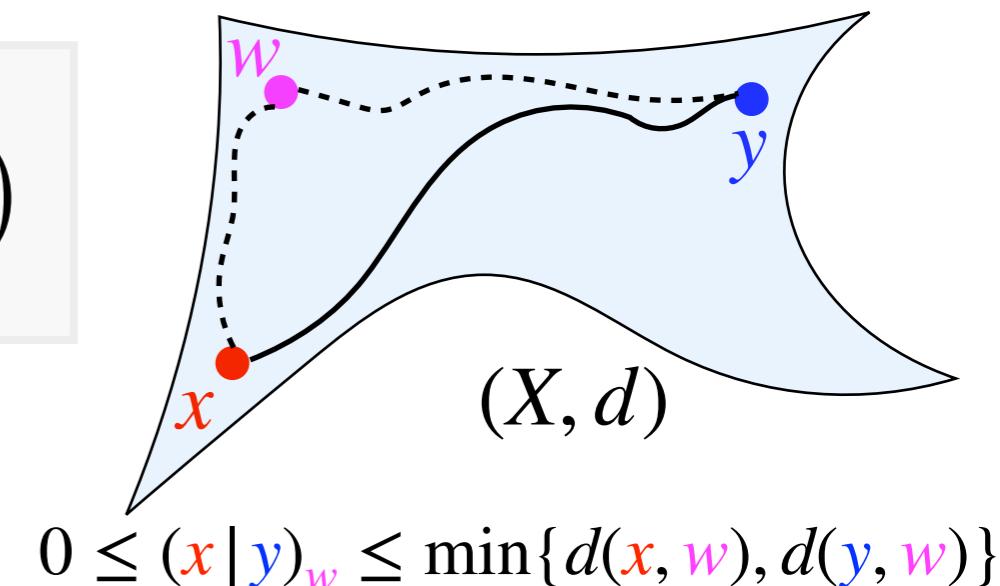
♦ Interpretation: δ -slim triangles

(X, d) is δ -hyperbolic iff each side of any Δ is contained within the δ -neighborhood of the union of the two other sides

Gromov hyperbolicity

♦ Gromov product

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



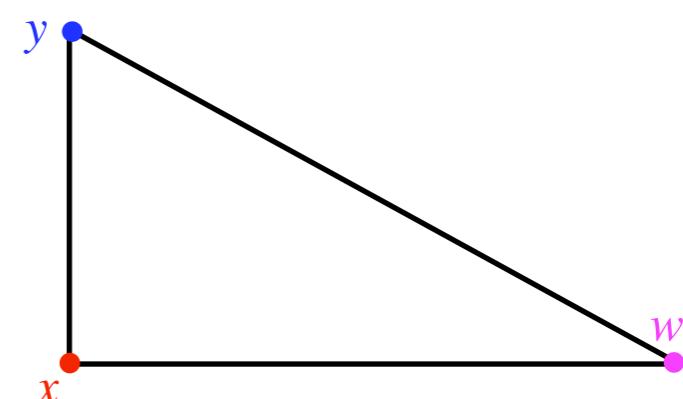
♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Interpretation: δ -slim triangles

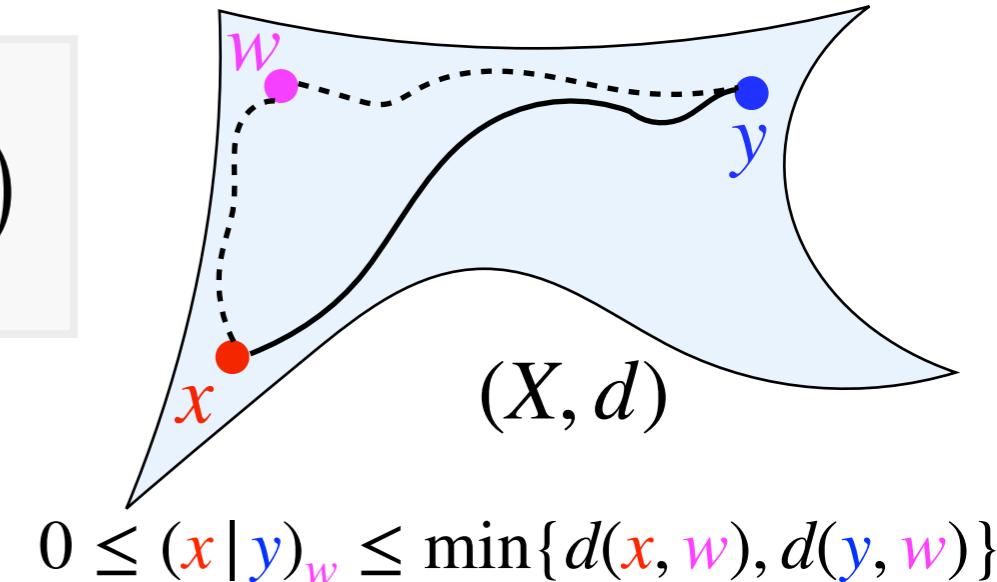
(X, d) is δ -hyperbolic iff each side of any Δ is contained within the δ -neighborhood of the union of the two other sides



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

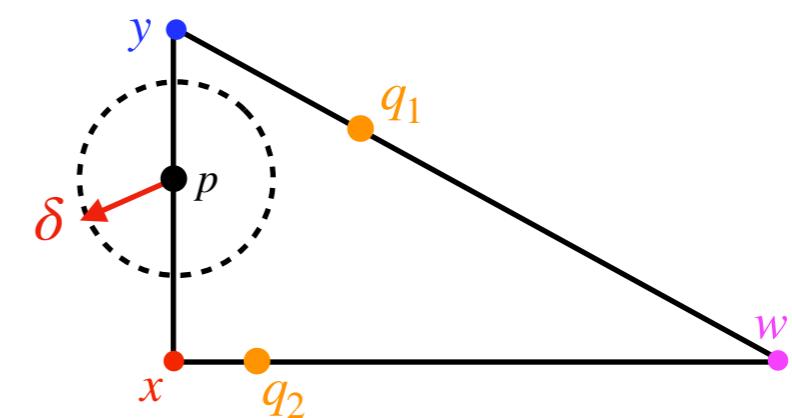
(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Interpretation: δ -slim triangles

(X, d) is δ -hyperbolic iff each side of any Δ is contained within the δ -neighborhood of the union of the two other sides

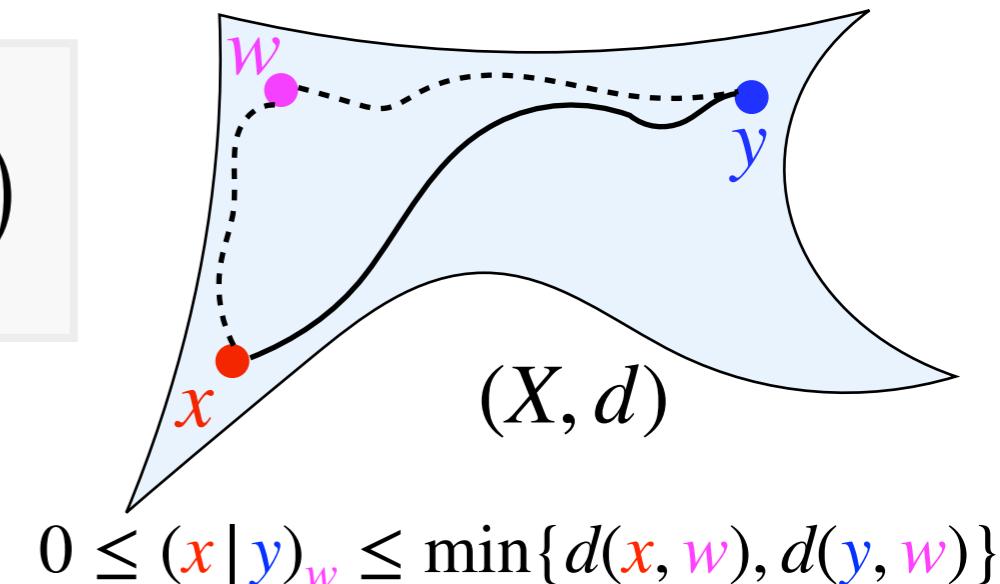
$$\forall p \in [x, y], \exists q \in [y, w] \cup [x, w], d(p, q) \leq \delta$$



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

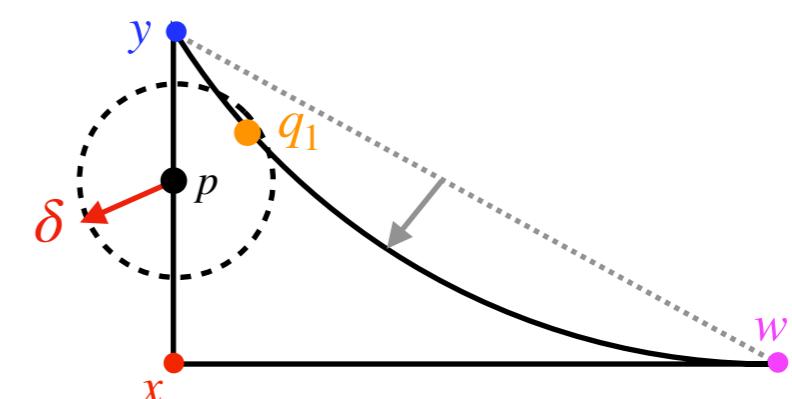
(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Interpretation: δ -slim triangles

(X, d) is δ -hyperbolic iff each side of any Δ is contained within the δ -neighborhood of the union of the two other sides

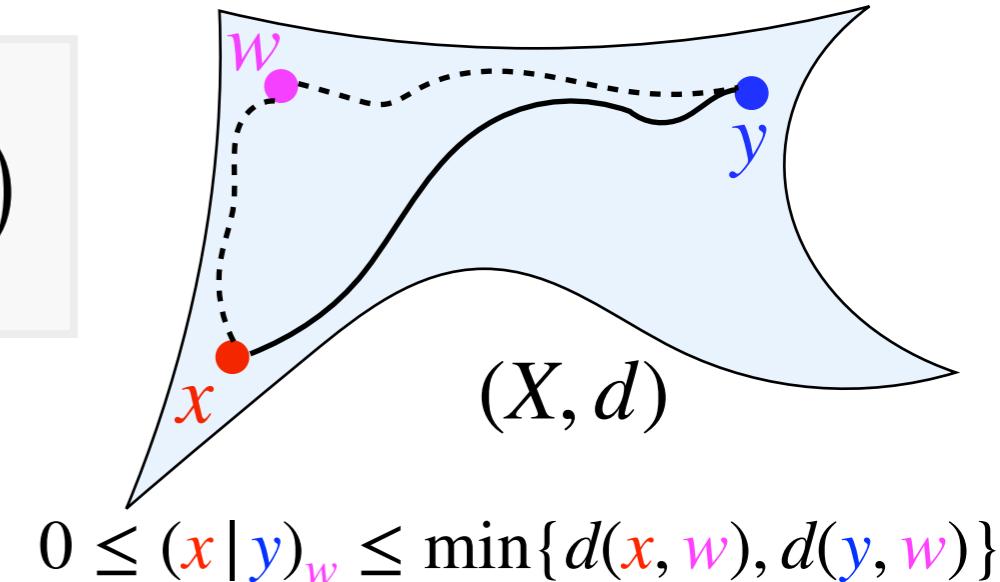
$\forall p \in [x, y], \exists q \in [y, w] \cup [x, w], d(p, q) \leq \delta$



Gromov hyperbolicity

◆ Gromov product

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



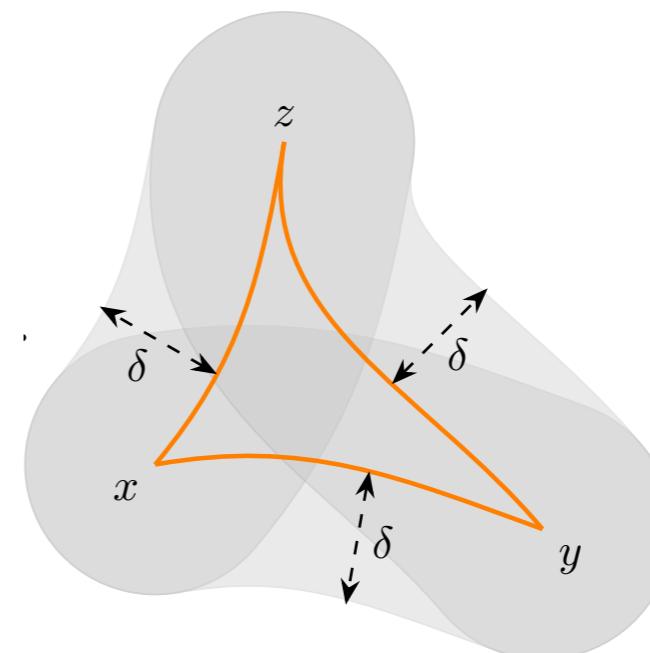
◆ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

◆ Interpretation: δ -slim triangles

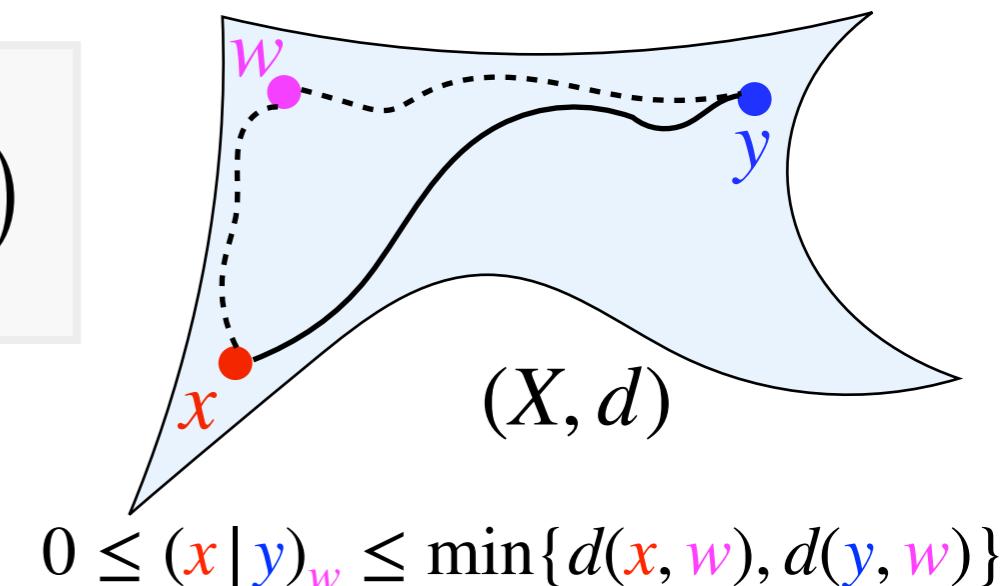
(X, d) is δ -hyperbolic iff each side of any Δ is contained within the δ -neighborhood of the union of the two other sides



Gromov hyperbolicity

♦ Gromov product

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



$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

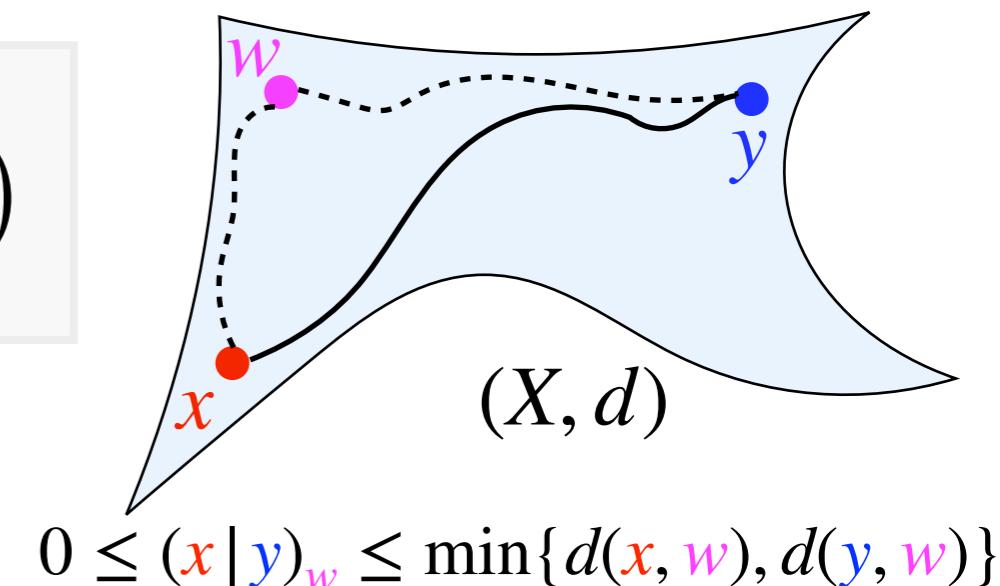
♦ Interpretation: δ -slim triangles

- ♦ Triangles in δ -hyperbolic spaces are « thin »
- ♦ The smaller δ the thinner

Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

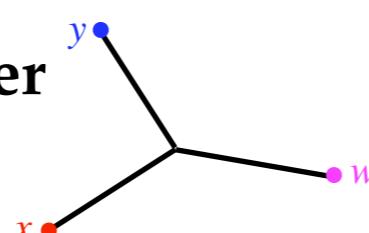
$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Interpretation: δ -slim triangles

♦ Triangles in δ -hyperbolic spaces are « thin »

♦ The smaller δ the thinner

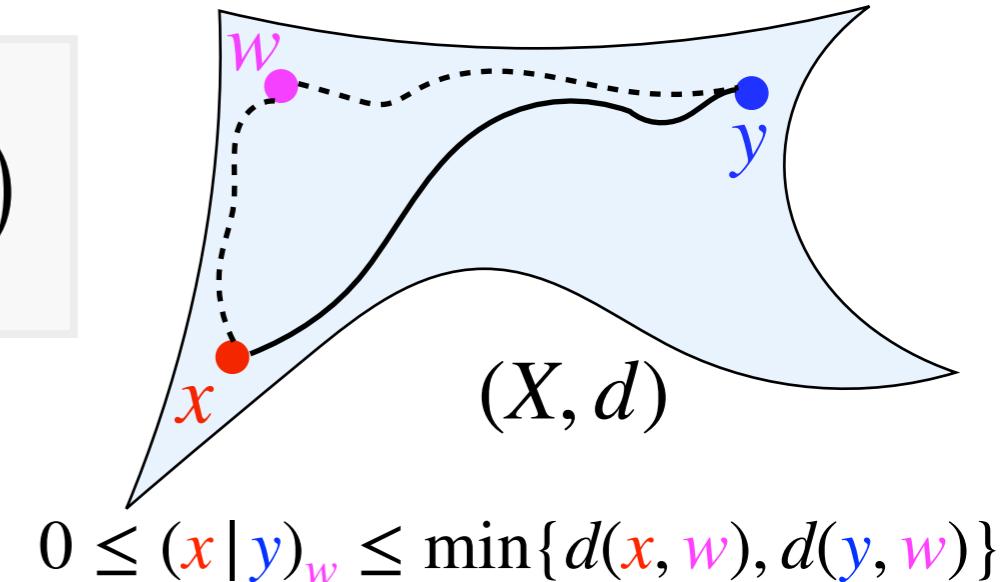
♦ Trees are 0-hyperbolic



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

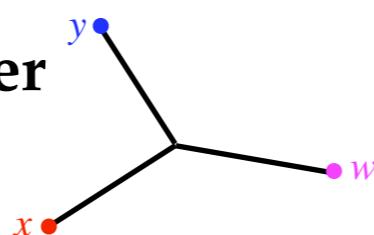
$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Interpretation: δ -slim triangles

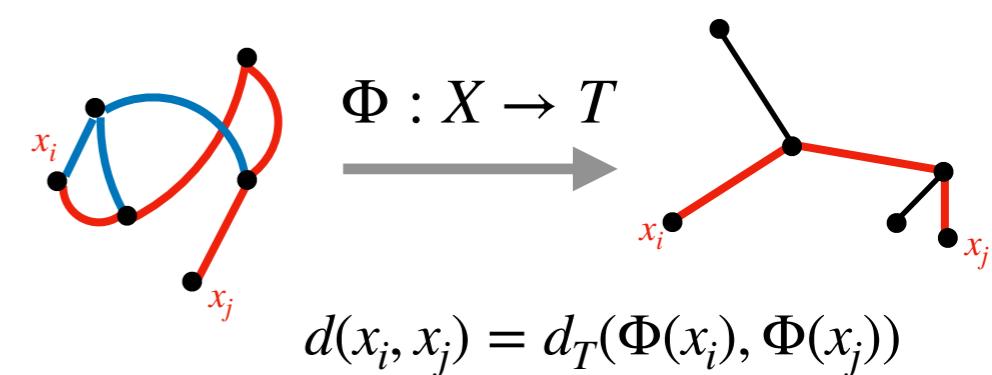
♦ Triangles in δ -hyperbolic spaces are « thin »

♦ The smaller δ the thinner

♦ Trees are 0-hyperbolic



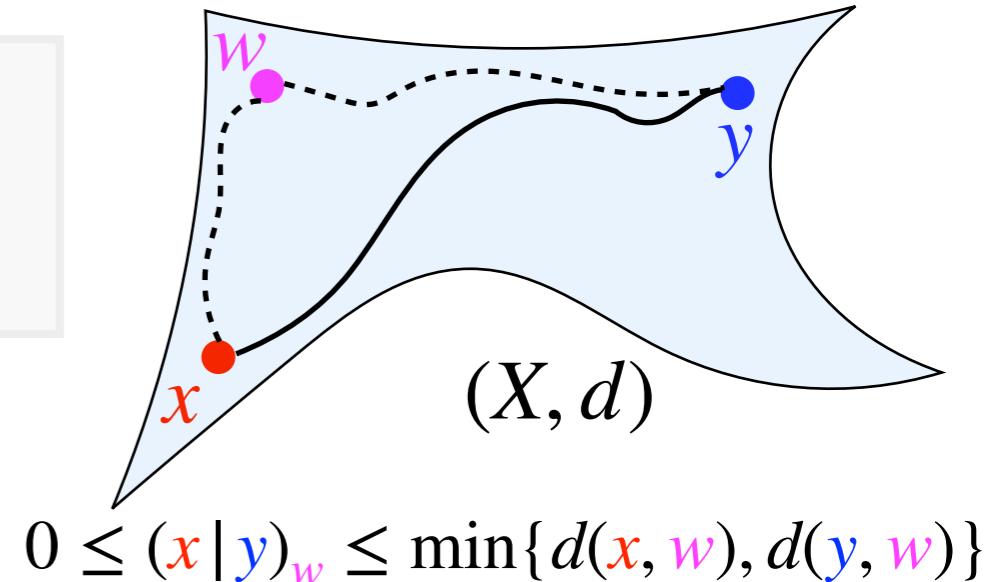
♦ If (X, d) is 0-hyperbolic it is isometric to a tree



Gromov hyperbolicity

◆ Gromov product

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



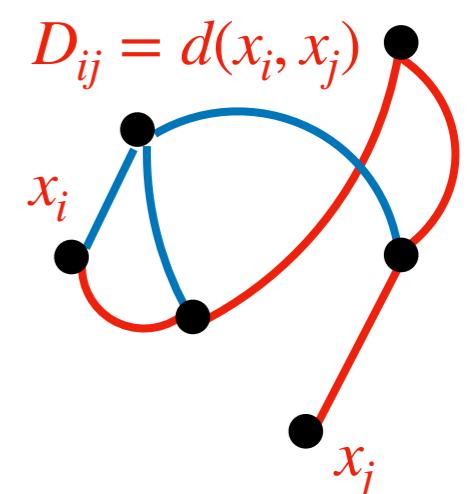
$$0 \leq (x|y)_w \leq \min\{d(x,w), d(y,w)\}$$

◆ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

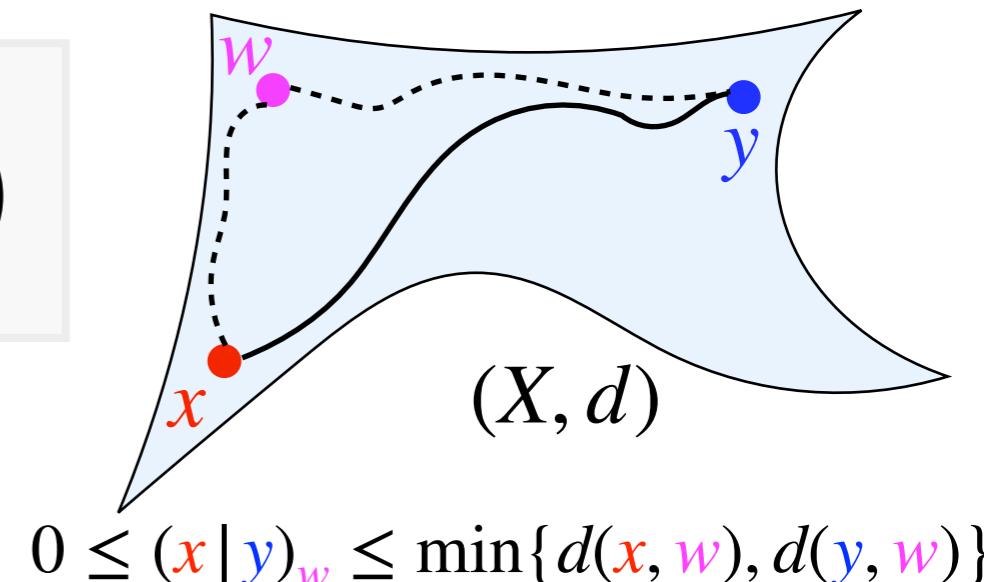
◆ Finite metric space $X = \{x_1, \dots, x_n\}$



Gromov hyperbolicity

◆ Gromov product

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



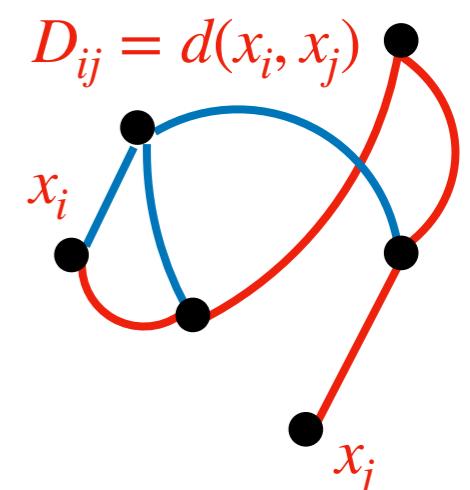
◆ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

◆ Finite metric space $X = \{x_1, \dots, x_n\}$

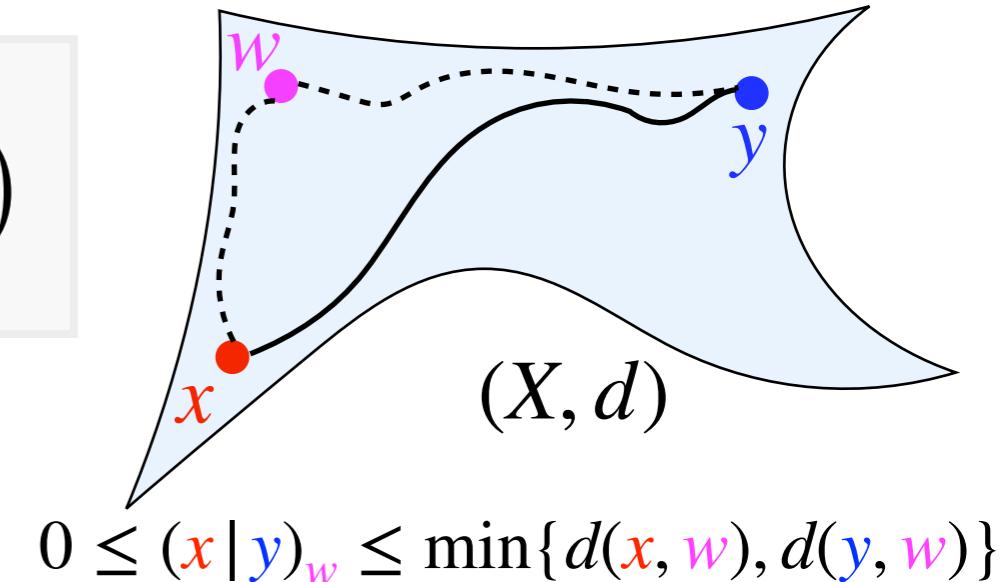
$$\delta_X = \max_{x,y,z,w} (\min\{(y|z)_w, (x|z)_w\} - (x|y)_w)$$



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

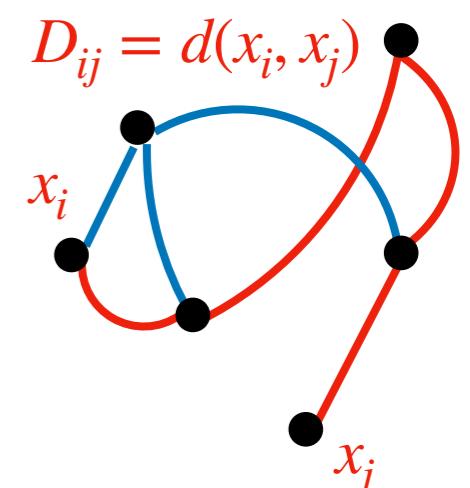
♦ Finite metric space $X = \{x_1, \dots, x_n\}$

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

♦ Take $z = w = y$ this implies $\delta_X \geq 0$

♦ Computation in $O(n^4)$

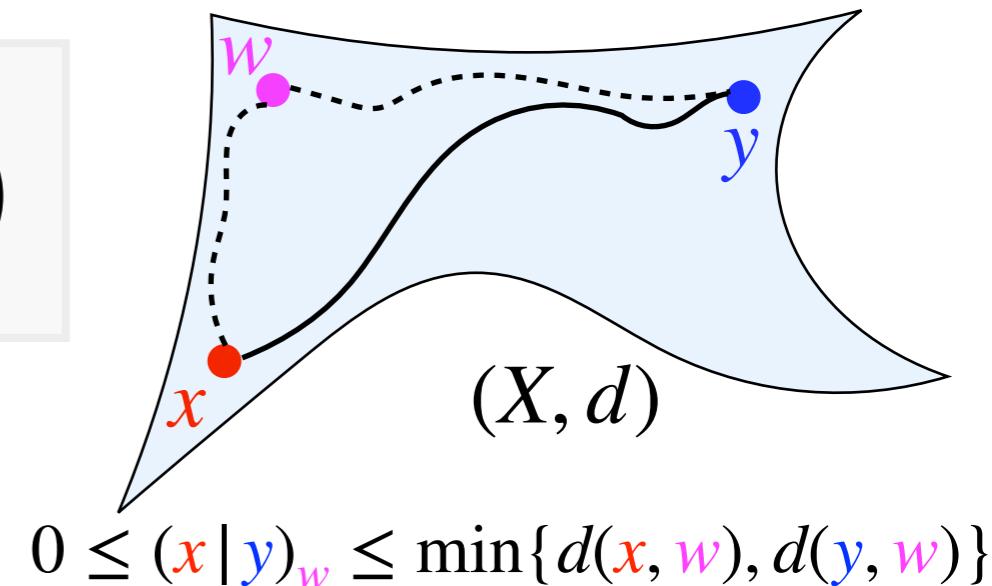
♦ Depends only on d , we note $\delta_d, \delta_D \dots$



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

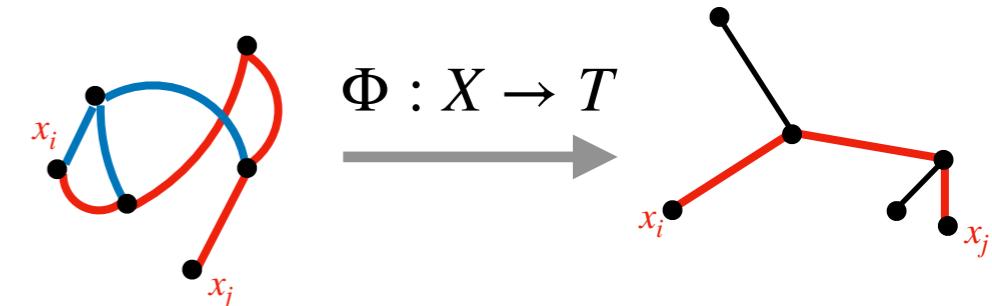
(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Finite metric space $X = \{x_1, \dots, x_n\}$

$\exists T$ a tree, $\Phi : X \rightarrow T$

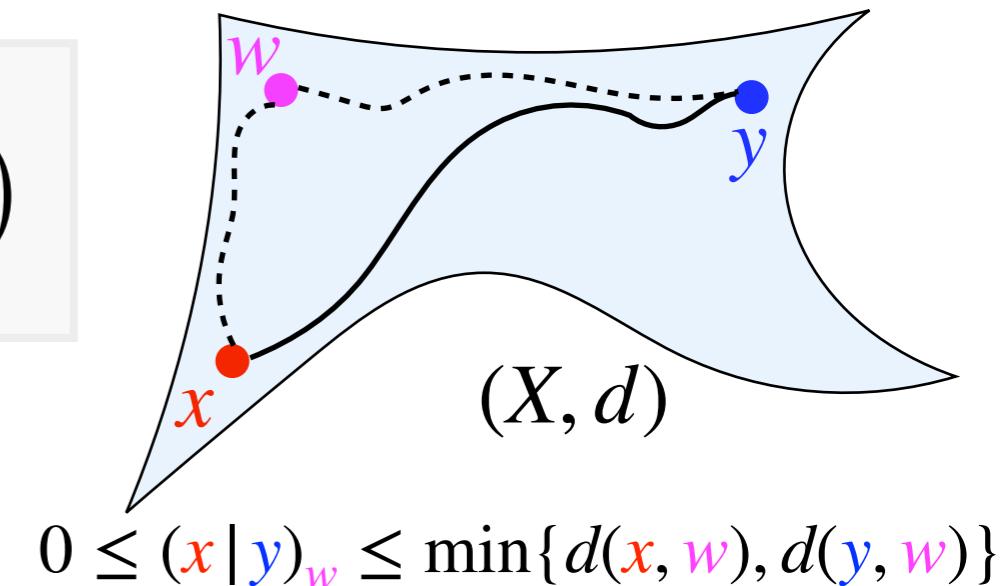
$$d(x_i, x_j) - 2\delta_X \log(n-2) \leq d_T(\Phi(x_i), \Phi(x_j)) \leq d(x_i, x_j)$$



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

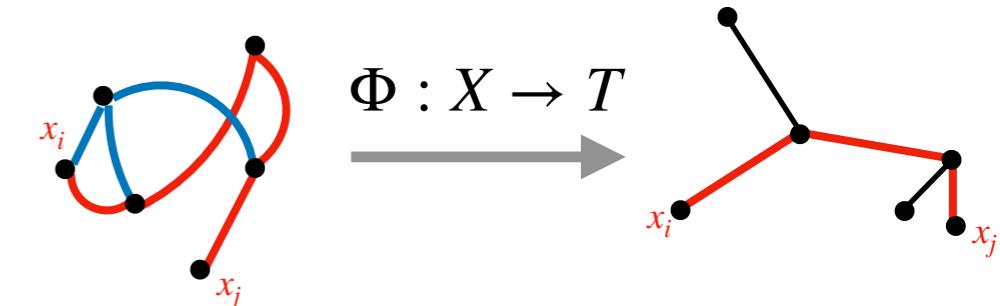
$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Finite metric space $X = \{x_1, \dots, x_n\}$

$\exists T$ a tree, $\Phi : X \rightarrow T$

$$d(x_i, x_j) - 2\delta_X \log(n-2) \leq d_T(\Phi(x_i), \Phi(x_j)) \leq d(x_i, x_j)$$

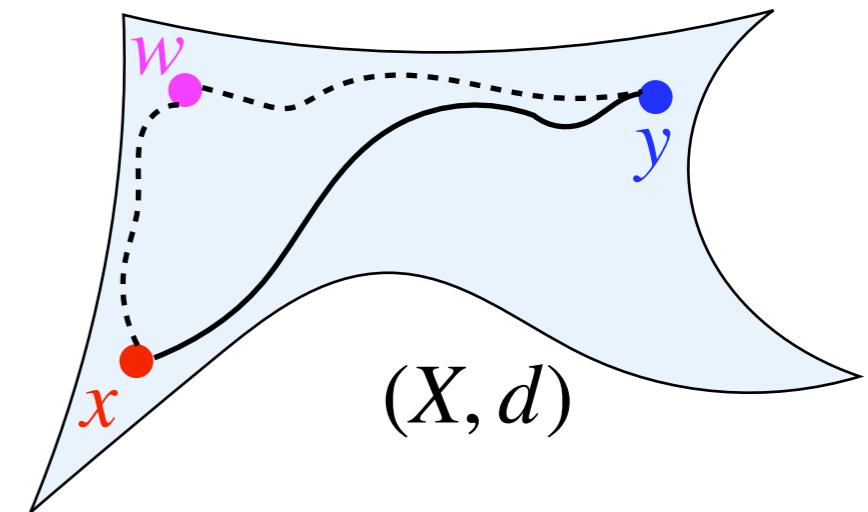
- ♦ It can be computed in $O(n^2)$
- ♦ Single Linkage Hierarchical Clustering algorithm
- ♦ This is the **Gromov embedding** $\Phi, T = \text{Gromov}(X, d)$



Gromov hyperbolicity

♦ Gromov product

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



♦ Gromov δ -hyperbolicity

(X, d) is δ -hyperbolic (for $\delta \geq 0$) if

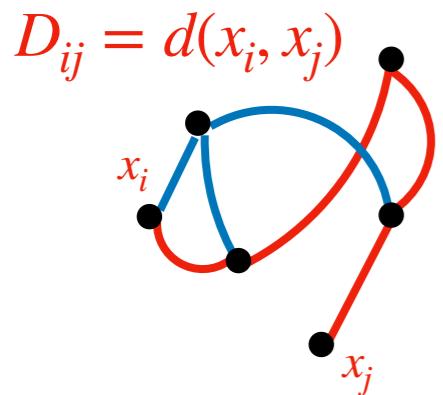
$$\forall x, y, z, w \in X, (x|y)_w \geq \min\{(y|z)_w, (x|z)_w\} - \delta$$

♦ Conclusion

- ◆ δ -hyperbolicity quantifies to which extent (X, d) is hyperbolic
- ◆ The smaller δ the more (X, d) has a « tree structure »

To DeltaZero

♦ Motivation

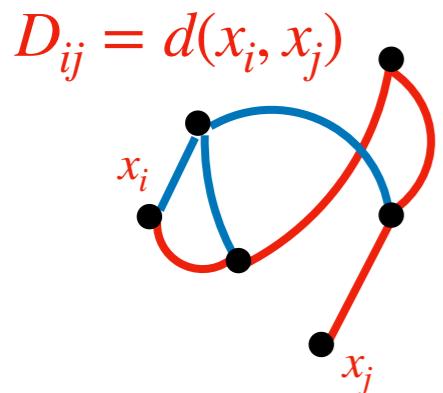


Embedding with a hierarchical structure

- ◆ Control how much « tree structure » we want
- ◆ Faithful to the original metric
- ◆ Reasonable in terms of computation

To DeltaZero

♦ Motivation



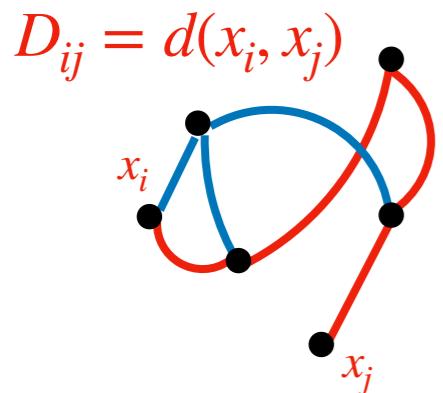
Embedding with a hierarchical structure

- ◆ Control how much « tree structure » we want
- ◆ Faithful to the original metric
- ◆ Reasonable in terms of computation

◆ Observation: $| d^* \in \operatorname{argmin}_{d' \in M_n} \mu \|d - d'\|_\infty + \delta_{d'}$
 $\forall i, j \ d'(x_i, x_j) \leq d(x_i, x_j)$

To DeltaZero

♦ Motivation



Embedding with a hierarchical structure

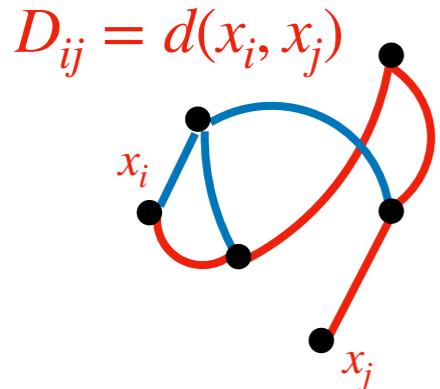
- ◆ Control how much « tree structure » we want
- ◆ Faithful to the original metric
- ◆ Reasonable in terms of computation

◆ Observation: $| d^\star \in \operatorname{argmin}_{d' \in M_n} \mu \|d - d'\|_\infty + \delta_{d'}$
 $\forall i, j \ d'(x_i, x_j) \leq d(x_i, x_j)$

$|\Phi, T = \text{Gromov}(X, d^\star)$

To DeltaZero

♦ Motivation



Embedding with a hierarchical structure

- ◆ Control how much « tree structure » we want
- ◆ Faithful to the original metric
- ◆ Reasonable in terms of computation

◆ Observation: $| d^\star \in \underset{d' \in M_n}{\operatorname{argmin}} \mu \|d - d'\|_\infty + \delta_{d'}$
 $\forall i, j \ d'(x_i, x_j) \leq d(x_i, x_j)$

$|\Phi, T = \text{Gromov}(X, d^\star)$

$$d(x_i, x_j) - 2\delta_X \log(n-2) + (2 \log(n-2)\mu - 1) \|d - d^\star\|_\infty \leq d_T(\Phi(x_i), \Phi(x_j)) \leq d(x_i, x_j)$$

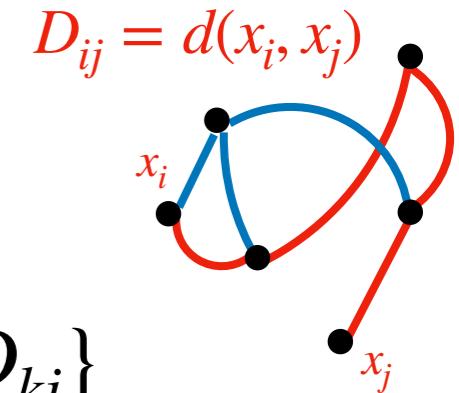
- ◆ When $\mu \geq 1/(2 \log(n-2))$ we improve the lower bound

To DeltaZero

♦ Optimization problem

- ♦ Space of metrics on n points

$$\mathcal{D}_n = \{D : \text{diag}(D) = 0, D = D^\top, D_{ij} \leq D_{ik} + D_{kj}\}$$



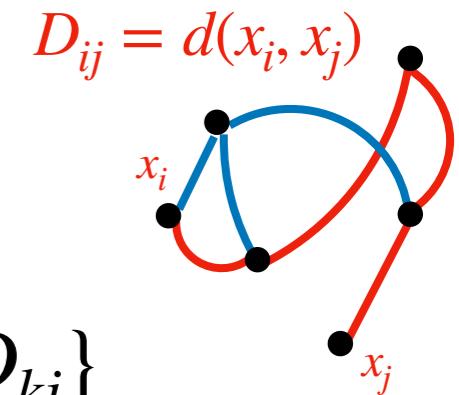
$$\min_{D' \in \mathcal{D}_n} L(D) := \mu \|D - D'\|_F^2 + \delta_{D'}$$

To DeltaZero

♦ Optimization problem

- ♦ Space of metrics on n points

$$\mathcal{D}_n = \{D : \text{diag}(D) = 0, D = D^\top, D_{ij} \leq D_{ik} + D_{kj}\}$$



$$\min_{D' \in \mathcal{D}_n} L(D) := \mu \|D - D'\|_F^2 + \delta_{D'}$$

- ♦ Trade-off between fidelity to D and small δ -hyperbolicity
- ♦ $D \rightarrow \delta_D$ is **piecewise affine** and **not convex**
- ♦ Complexity in $O(n^4)$ + not everywhere differentiable
- ♦ How to handle the constraints ?

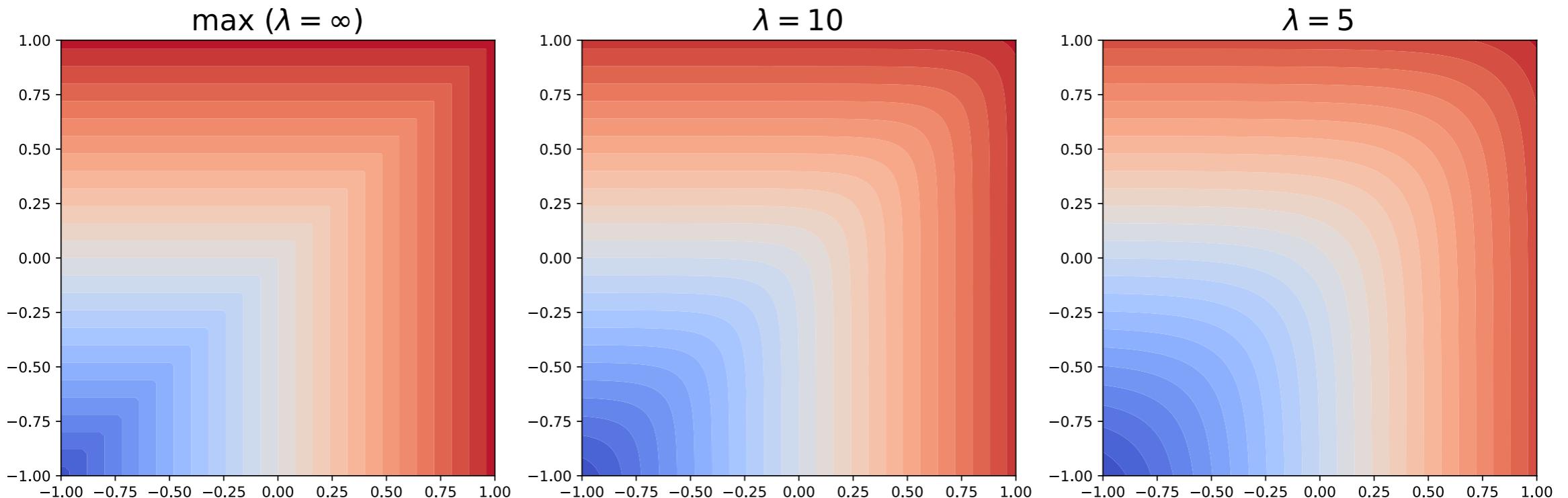
Smooth and batched δ -hyperbolicity

◆ Smoothing δ -hyperbolicity

- ◆ We replace the \min, \max in δ -hyperbolicity by a smooth surrogate

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

$$\max(x) \leq \text{LSE}_\lambda(x) \leq \max(x) + \frac{\log(n)}{\lambda}$$



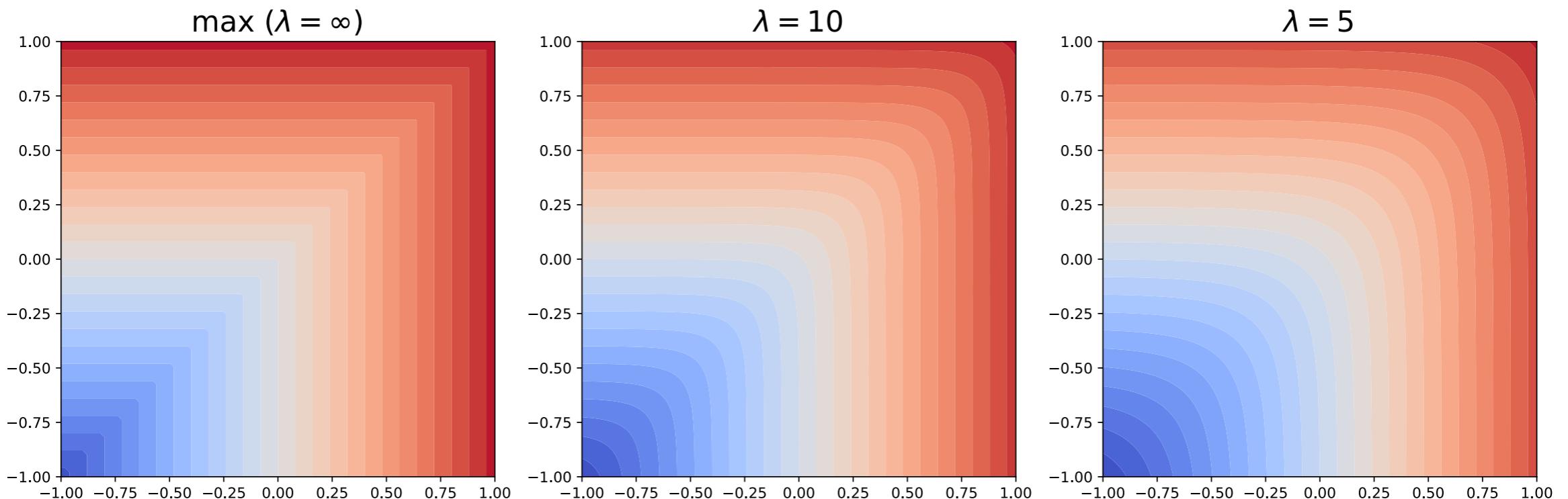
Smooth and batched δ -hyperbolicity

◆ Smoothing δ -hyperbolicity

- ◆ We replace the \min, \max in δ -hyperbolicity by a smooth surrogate

$$\max(x) \leq \text{LSE}_\lambda(x) \leq \max(x) + \frac{\log(n)}{\lambda}$$

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



- ◆ Corresponds to regularizing a linear problem

$$\max(x) = \max_{q \in \Sigma_n} \langle x, q \rangle \longrightarrow \text{LSE}_\lambda(x) = \operatorname{argmax}_{q \in \Sigma_n} \langle x, q \rangle - \frac{1}{\lambda} \text{entropy}(q)$$

Smooth and batched δ -hyperbolicity

◆ Smoothing δ -hyperbolicity

- ◆ We replace the min , max in δ -hyperbolicity by a smooth surrogate

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

- ◆ Smoothed δ -hyperbolicity: differentiable but still $O(n^4)$

$$\delta^{(\lambda)} = \text{LSE}_\lambda \left(\left\{ \text{LSE}_{-\lambda} \{ (\textcolor{blue}{y} | z)_{\textcolor{violet}{w}}, (\textcolor{red}{x} | z)_{\textcolor{violet}{w}} \} - (\textcolor{red}{x} | \textcolor{blue}{y})_{\textcolor{violet}{w}} \right\}_{x,y,z,w} \right)$$

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

Smooth and batched δ -hyperbolicity

◆ Smoothing δ -hyperbolicity

- ◆ We replace the min , max in δ -hyperbolicity by a smooth surrogate

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

- ◆ Smoothed δ -hyperbolicity: differentiable but still $O(n^4)$

$$\delta^{(\lambda)} = \text{LSE}_\lambda \left(\left\{ \text{LSE}_{-\lambda} \{ (\textcolor{blue}{y} | z)_{\textcolor{violet}{w}}, (\textcolor{red}{x} | z)_{\textcolor{violet}{w}} \} - (\textcolor{red}{x} | \textcolor{blue}{y})_{\textcolor{violet}{w}} \right\}_{x,y,z,w} \right)$$

◆ Batched δ -hyperbolicity

- ◆ Sample m points among the n
- ◆ Do that K times: gives $X_1, \dots, X_K \subset X$
- ◆ Compute $\delta_{X_1}^{(\lambda)}, \dots, \delta_{X_K}^{(\lambda)}$
- ◆ Complexity $O(K \cdot m^4)$

Smooth and batched δ -hyperbolicity

◆ Smoothing δ -hyperbolicity

- ◆ We replace the min, max in δ -hyperbolicity by a smooth surrogate

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

- ◆ Smoothed δ -hyperbolicity: differentiable but still $O(n^4)$

$$\delta^{(\lambda)} = \text{LSE}_\lambda \left(\left\{ \text{LSE}_{-\lambda} \{ (\textcolor{blue}{y} | z)_{\textcolor{violet}{w}}, (\textcolor{red}{x} | z)_{\textcolor{violet}{w}} \} - (\textcolor{red}{x} | \textcolor{blue}{y})_{\textcolor{violet}{w}} \right\}_{x,y,z,w} \right)$$

◆ Batched δ -hyperbolicity

- ◆ Sample m points among the n

- ◆ Do that K times: gives $X_1, \dots, X_K \subset X$

- ◆ Compute $\delta_{X_1}^{(\lambda)}, \dots, \delta_{X_K}^{(\lambda)}$

- ◆ Complexity $O(K \cdot m^4)$
- ◆ Under some hypothesis, close to δ with high prob.

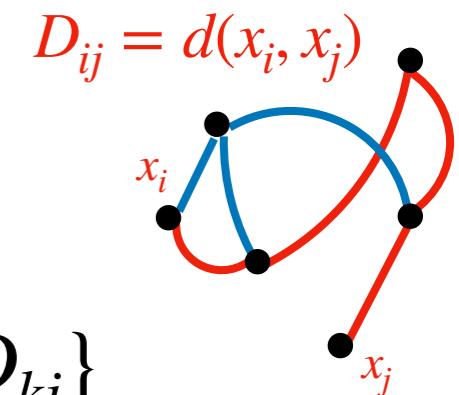
$$\delta_{K,m}^{(\lambda)} = \text{LSE}_\lambda(\delta_{X_1}^{(\lambda)}, \dots, \delta_{X_K}^{(\lambda)})$$

DeltaZero

♦ Optimization problem

- ♦ Space of metrics on n points

$$\mathcal{D}_n = \{D : \text{diag}(D) = 0, D = D^\top, D_{ij} \leq D_{ik} + D_{kj}\}$$



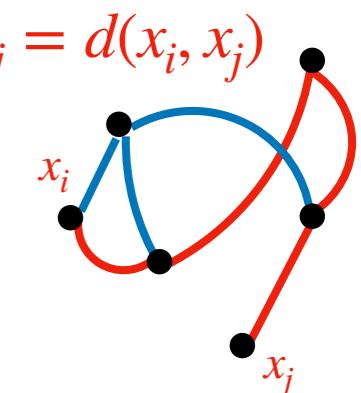
$$\min_{D' \in \mathcal{D}_n} L(D) := \mu \|D - D'\|_F^2 + \delta_{D', K, m}^{(\lambda)}$$

DeltaZero

♦ Optimization problem

- ♦ Space of metrics on n points

$$\mathcal{D}_n = \{D : \text{diag}(D) = 0, D = D^\top, D_{ij} \leq D_{ik} + D_{kj}\}$$



$$\min_{D' \in \mathcal{D}_n} L(D) := \mu \|D - D'\|_F^2 + \delta_{D', K, m}^{(\lambda)}$$

- ♦ Inspired by projected gradient descent

Algorithm

$$G_t = \nabla L(D_t)$$

$$D_{t+\frac{1}{2}} = \text{Adam}(G_t, D_t)$$

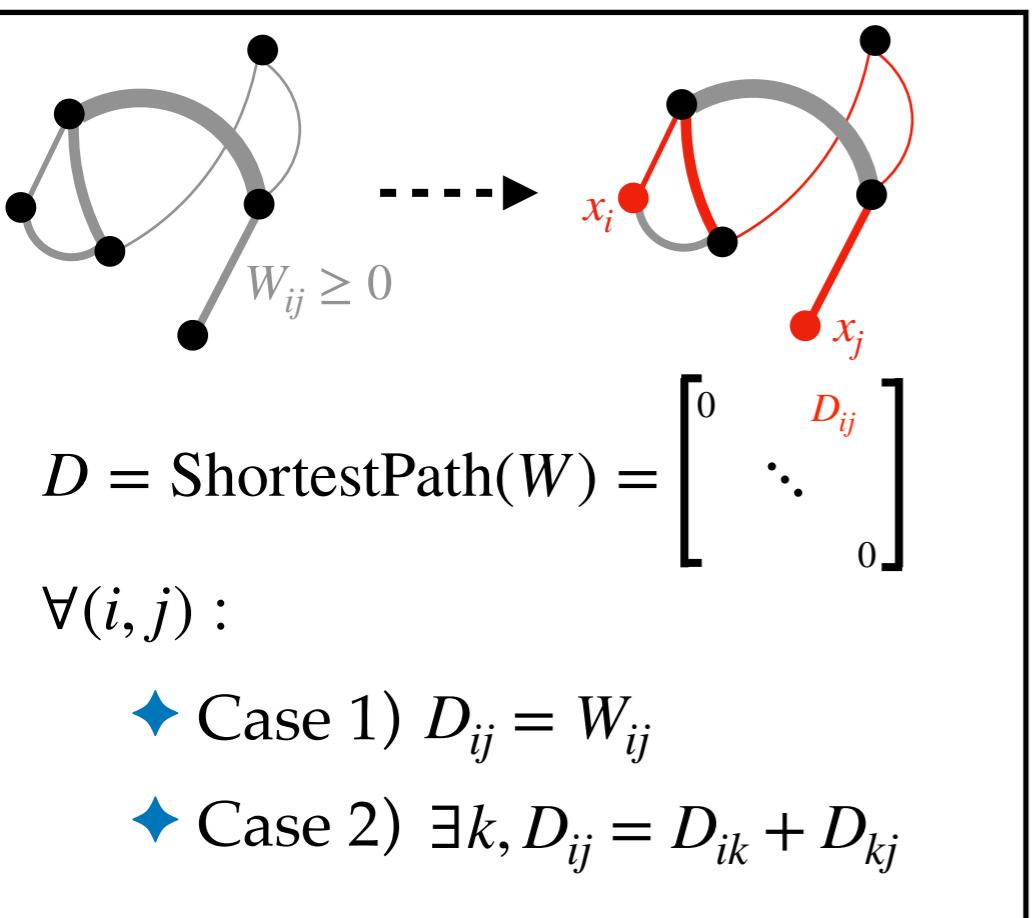
$$D_t = \Pi(D_{t+\frac{1}{2}})$$

$$\Pi(W) = \underset{D \in \mathcal{D}_n : D \leq W}{\operatorname{argmin}} \|D - W\|_F^2$$

- ♦ How to compute this projection ?
- ♦ Answer: this is a **shortest path problem**

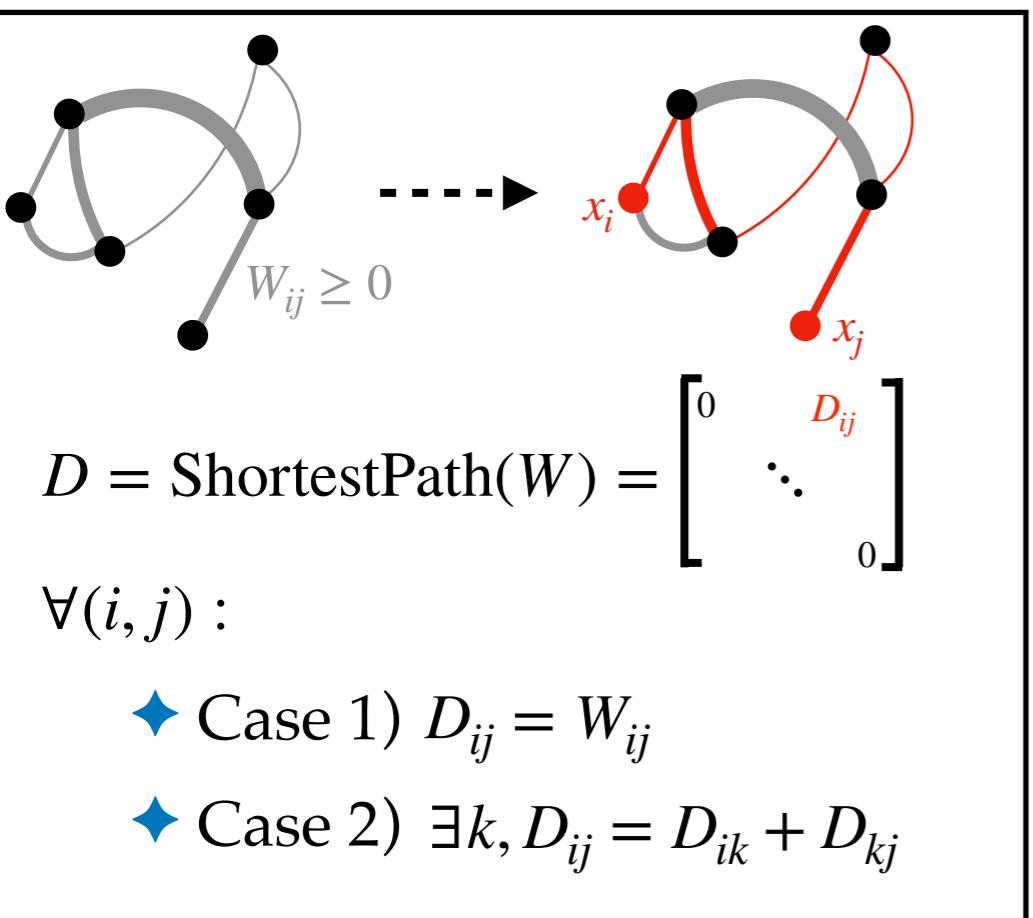
DeltaZero

◆ The metric nearest problem



DeltaZero

◆ The metric nearest problem

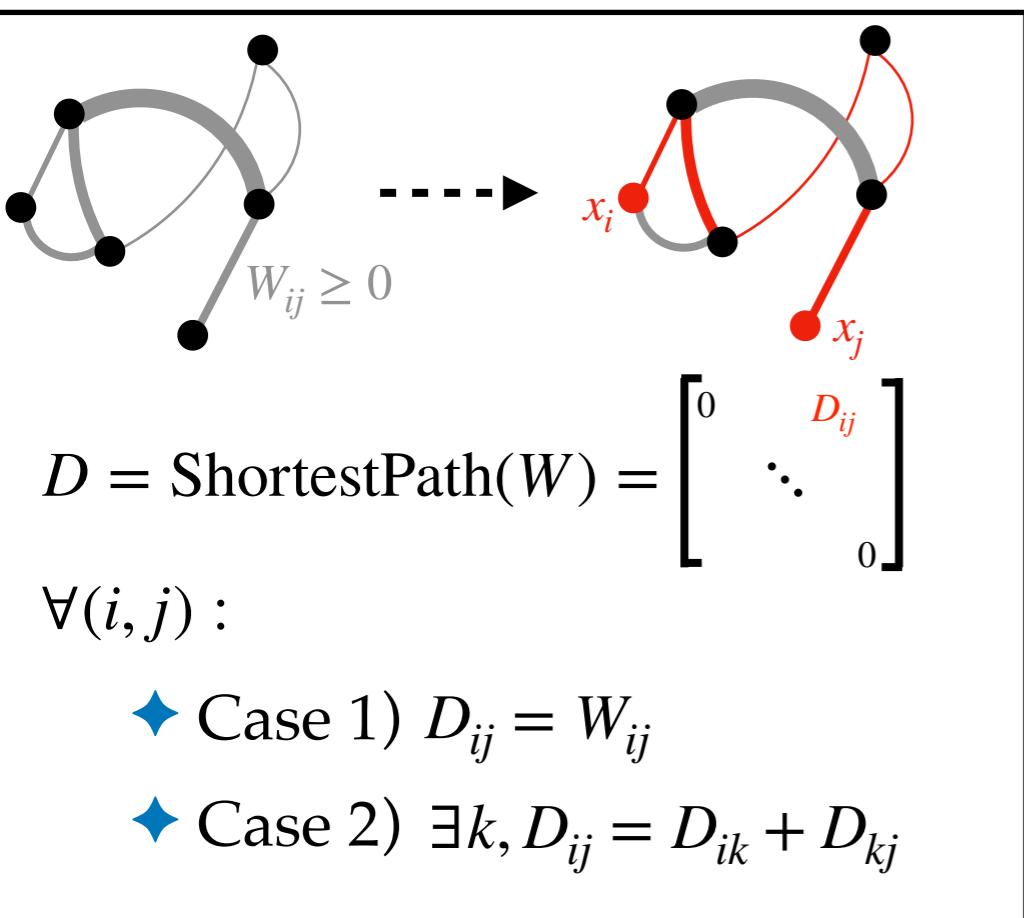


(Brickell, 2008)

$$\forall p, \underset{D \in \mathcal{D}_n : D \leq W}{\operatorname{argmin}} \|D - W\|_{\ell_p} = \text{ShortestPath}(W)$$

DeltaZero

◆ The metric nearest problem



(Brickell, 2008)

$$\forall p, \underset{D \in \mathcal{D}_n : D \leq W}{\operatorname{argmin}} \|D - W\|_{\ell_p} = \text{ShortestPath}(W)$$

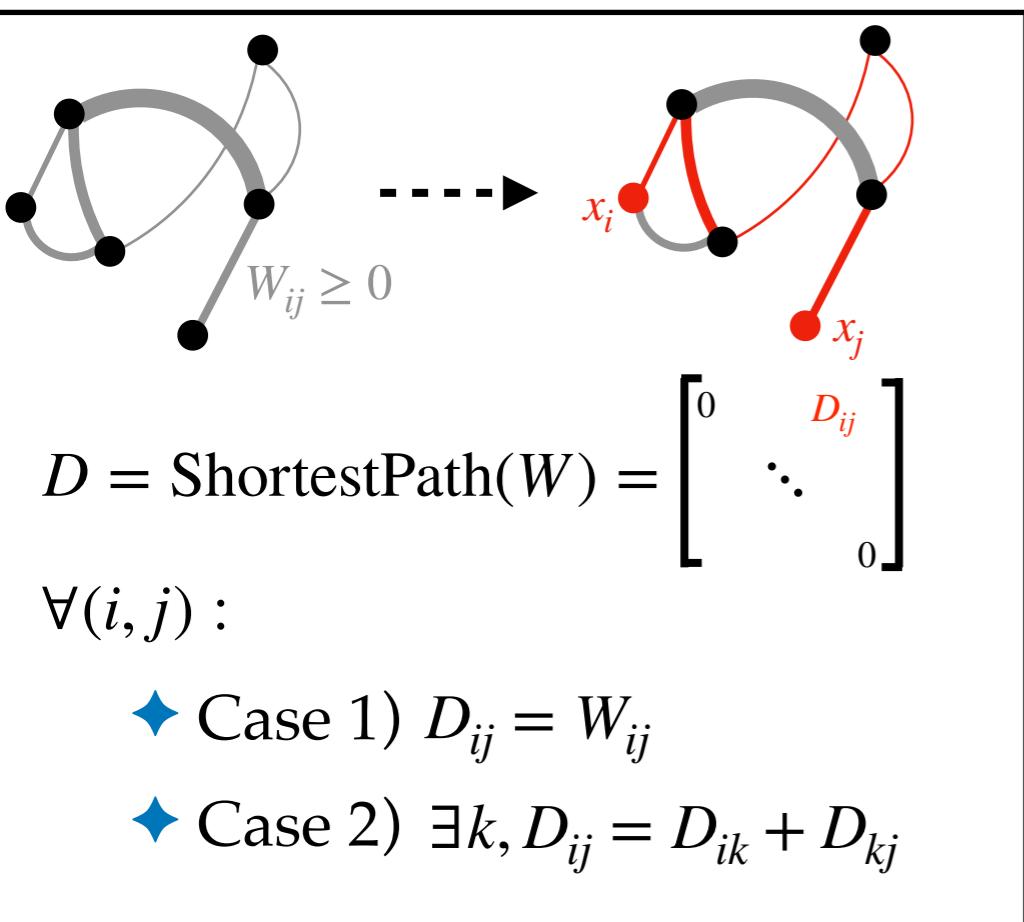
◆ Floyd-Warshall

```
 $D = W$ 
for  $k \in \{1, \dots, n\}$ 
    for  $i \in \{1, \dots, n\}$ 
        for  $j \in \{1, \dots, n\}$ 
             $D_{ij} = \min\{D_{ij}, D_{ik} + D_{kj}\}$ 
```

- ◆ Finds $\text{ShortestPath}(W)$
- ◆ Runs in $O(n^3)$

DeltaZero

◆ The metric nearest problem



(Brickell, 2008)

$$\forall p, \underset{D \in \mathcal{D}_n : D \leq W}{\operatorname{argmin}} \|D - W\|_{\ell_p} = \text{ShortestPath}(W)$$

◆ Floyd-Warshall

```
 $D = W$ 
for  $k \in \{1, \dots, n\}$ 
    for  $i \in \{1, \dots, n\}$ 
        for  $j \in \{1, \dots, n\}$ 
             $D_{ij} = \min\{D_{ij}, D_{ik} + D_{kj}\}$ 
```

- ◆ Finds $\text{ShortestPath}(W)$
- ◆ Runs in $O(n^3)$

◆ DeltaZero

$$G_t = \nabla L(D_t) \quad // O(K \cdot m^4 + n^2)$$

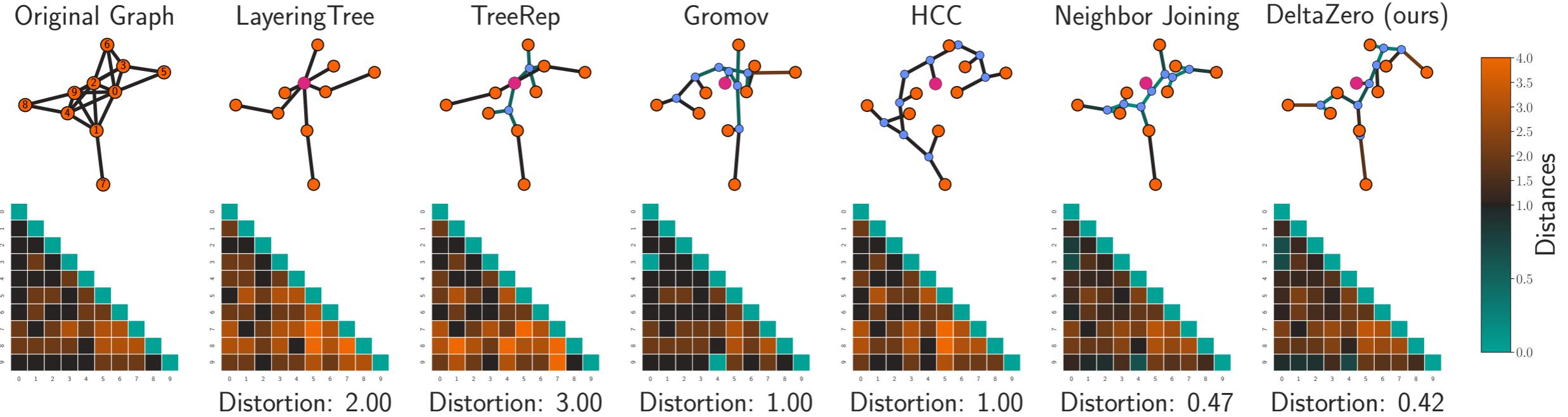
$$D_{t+\frac{1}{2}} = \text{Adam}(G_t, D_t)$$

$$D_t = \text{FloydWarshall}(D_{t+\frac{1}{2}}) \quad // O(n^3)$$

if output tree: $\Phi, T = \text{Gromov}(D_\infty) \quad // O(n^2)$

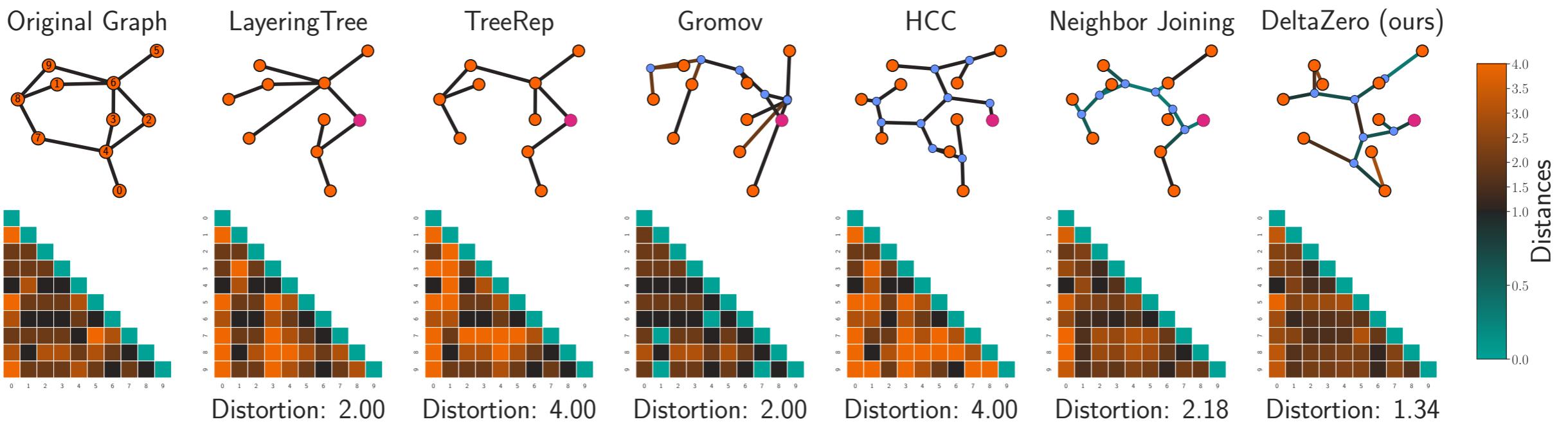
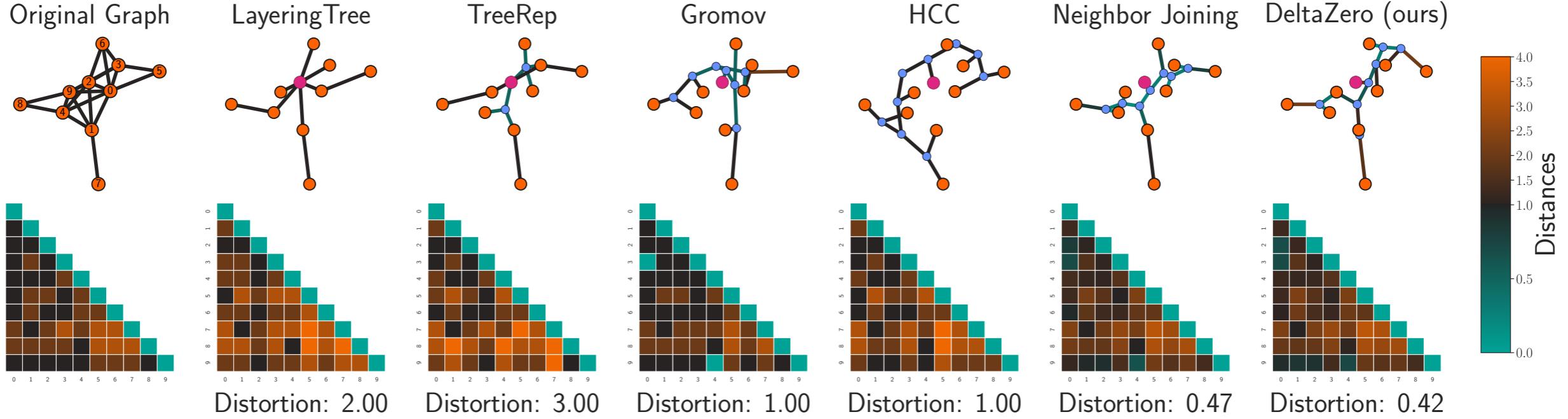
DeltaZero

❖ Illustrations



DeltaZero

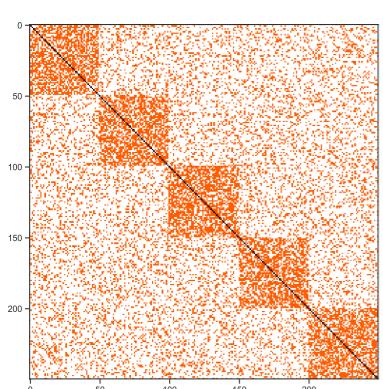
❖ Illustrations



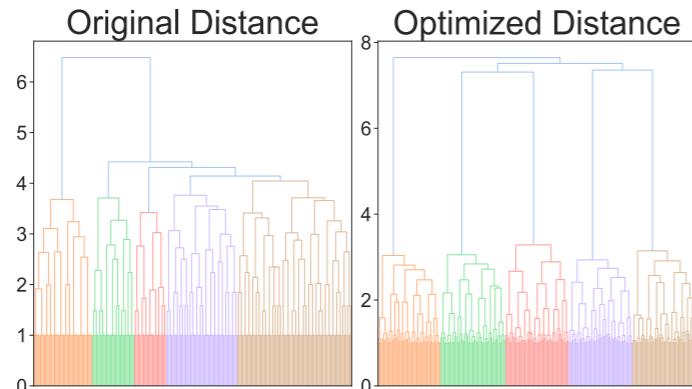
DeltaZero

♦ On stochastic block model

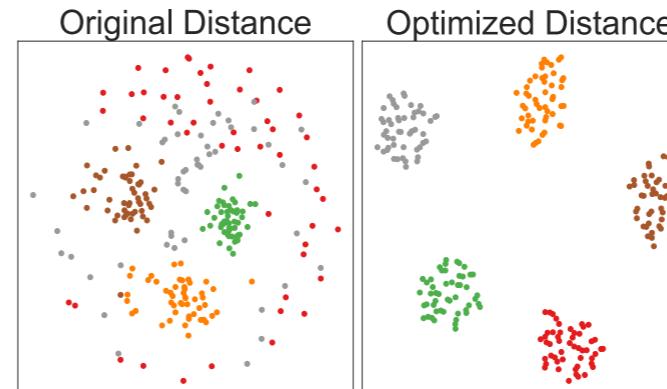
- ♦ SBM with 5 communities. Gives a shortest path matrix D
- ♦ Objective: clustering the nodes of the graph given D
- ♦ We compute $D' = \text{DeltaZero}(D)$
- ♦ We compare clustering (single linkage) with D vs D'



(a) All pairs Shortest-Paths
distance matrix D



(b) Dendrograms from original and opti-
mized distance matrices

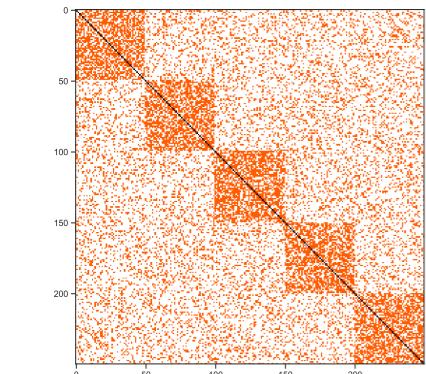


(c) t-SNE plots from original and opti-
mized distance matrices

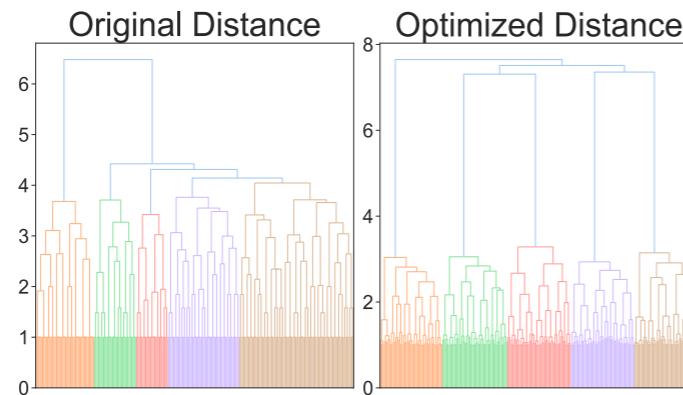
DeltaZero

◆ On stochastic block model

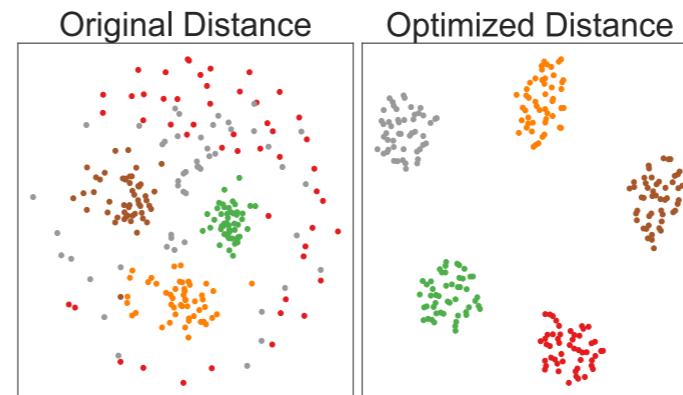
- ◆ SBM with 5 communities. Gives a shortest path matrix D
- ◆ Objective: clustering the nodes of the graph given D
- ◆ We compute $D' = \text{DeltaZero}(D)$
- ◆ We compare clustering (single linkage) with D vs D'



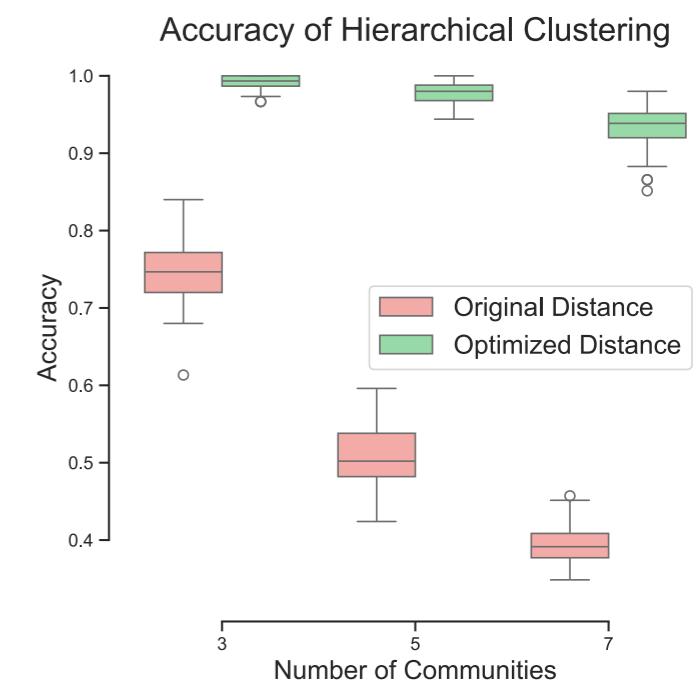
(a) All pairs Shortest-Paths
distance matrix D



(b) Dendograms from original and opti-
mized distance matrices



(c) t-SNE plots from original and opti-
mized distance matrices



DeltaZero

◆ Distortion on real datasets

- ◆ We compute the tree metric $D_T = \text{DeltaZero}(D) + \text{Gromov}$
- ◆ We evaluate $\|D_T - D\|_\infty$

Datasets	Unweighted graphs					Non-graph metrics	
	C-ELEGAN	CS PHD	CORA	AIRPORT	WIKI	ZEISEL	IBD
	n	452	1025	2485	3158	2357	3005
Diameter	7	28	19	12	9	0.87	0.99
NJ	<u>2.97</u>	16.81	13.42	4.18	6.32	0.51	<u>0.90</u>
TR	5.90 ± 0.72	21.01 ± 3.34	16.86 ± 2.11	10.00 ± 1.02	9.97 ± 0.93	0.66 ± 0.10	1.60 ± 0.22
HCC	4.31 ± 0.46	23.35 ± 2.07	12.28 ± 0.96	7.71 ± 0.72	7.20 ± 0.60	0.53 ± 0.07	1.25 ± 0.11
LayeringTree	5.07 ± 0.25	25.48 ± 0.60	7.76 ± 0.54	2.97 ± 0.26	4.08 ± 0.27	–	–
Gromov	3.33 ± 0.45	<u>13.28 ± 0.61</u>	9.34 ± 0.53	4.08 ± 0.27	5.54 ± 0.49	<u>0.43 ± 0.02</u>	1.01 ± 0.04
DELTAZERO	1.87 ± 0.08	10.31 ± 0.62	7.59 ± 0.38	2.79 ± 0.15	3.56 ± 0.20	0.24 ± 0.00	0.70 ± 0.03
Improvement (%)	43.8%	22.3%	2.3%	6.0%	12.7%	44.1 %	22.2%

DeltaZero

◆ Distortion on real datasets

- ◆ We compute the tree metric $D_T = \text{DeltaZero}(D) + \text{Gromov}$
- ◆ We evaluate $\|D_T - D\|_\infty$

Datasets	Unweighted graphs					Non-graph metrics	
	C-ELEGAN	CS PHD	CORA	AIRPORT	WIKI	ZEISEL	IBD
	n	452	1025	2485	3158	2357	3005
Diameter	7	28	19	12	9	0.87	0.99
NJ	<u>2.97</u>	16.81	13.42	4.18	6.32	0.51	<u>0.90</u>
TR	5.90 ± 0.72	21.01 ± 3.34	16.86 ± 2.11	10.00 ± 1.02	9.97 ± 0.93	0.66 ± 0.10	1.60 ± 0.22
HCC	4.31 ± 0.46	23.35 ± 2.07	12.28 ± 0.96	7.71 ± 0.72	7.20 ± 0.60	0.53 ± 0.07	1.25 ± 0.11
LayeringTree	5.07 ± 0.25	25.48 ± 0.60	<u>7.76 ± 0.54</u>	<u>2.97 ± 0.26</u>	<u>4.08 ± 0.27</u>	–	–
Gromov	3.33 ± 0.45	<u>13.28 ± 0.61</u>	9.34 ± 0.53	4.08 ± 0.27	5.54 ± 0.49	<u>0.43 ± 0.02</u>	1.01 ± 0.04
DELTAZERO	1.87 ± 0.08	10.31 ± 0.62	7.59 ± 0.38	2.79 ± 0.15	3.56 ± 0.20	0.24 ± 0.00	0.70 ± 0.03
Improvement (%)	43.8%	22.3%	2.3%	6.0%	12.7%	44.1 %	22.2%

◆ Sensitivity analysis

