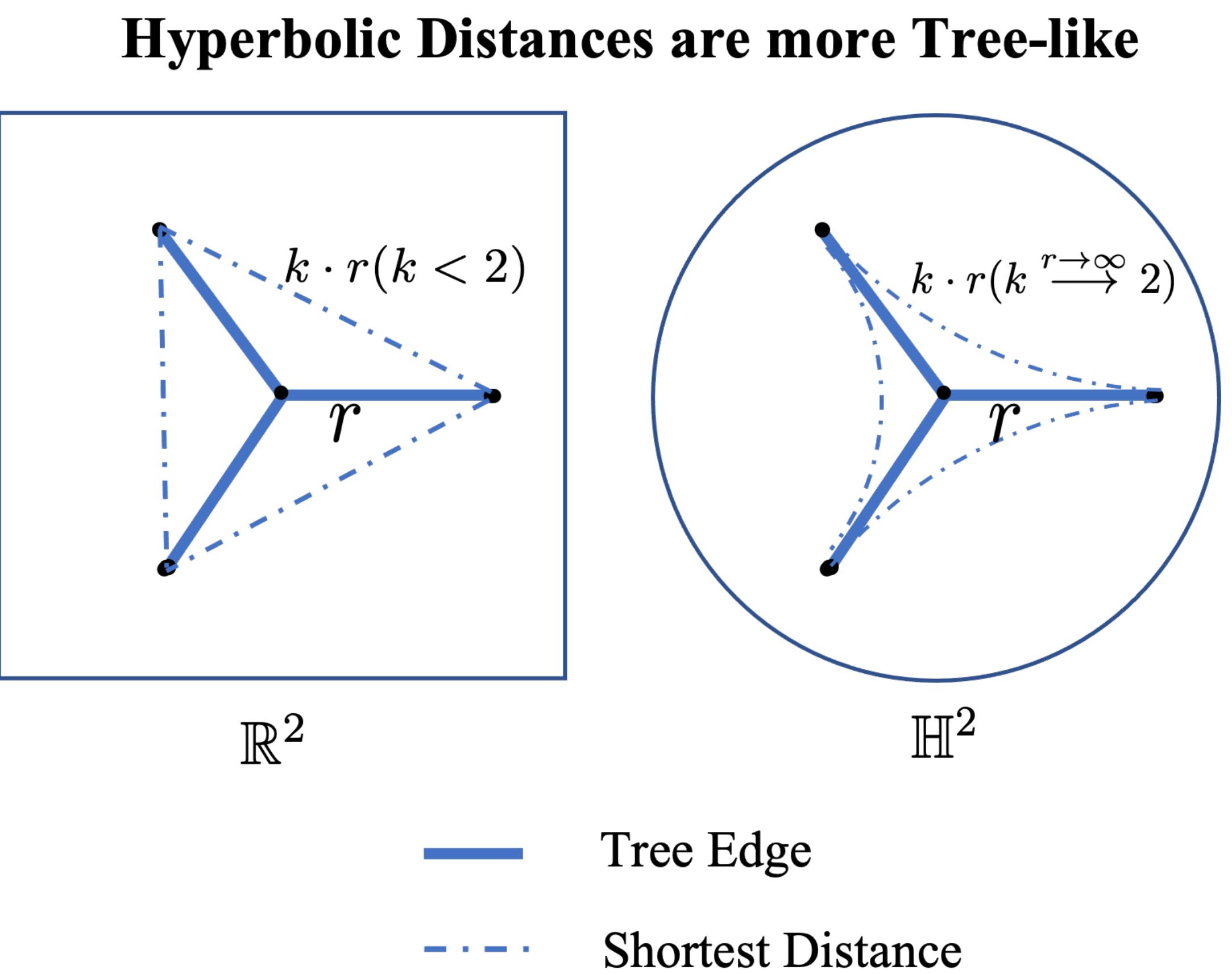


Gal Mishne<sup>1</sup>, Zhengchao Wan<sup>1</sup>, Yusu Wang<sup>1</sup>, Sheng Yang<sup>2</sup>  
<sup>1</sup>University of California - San Diego, <sup>2</sup>Harvard University

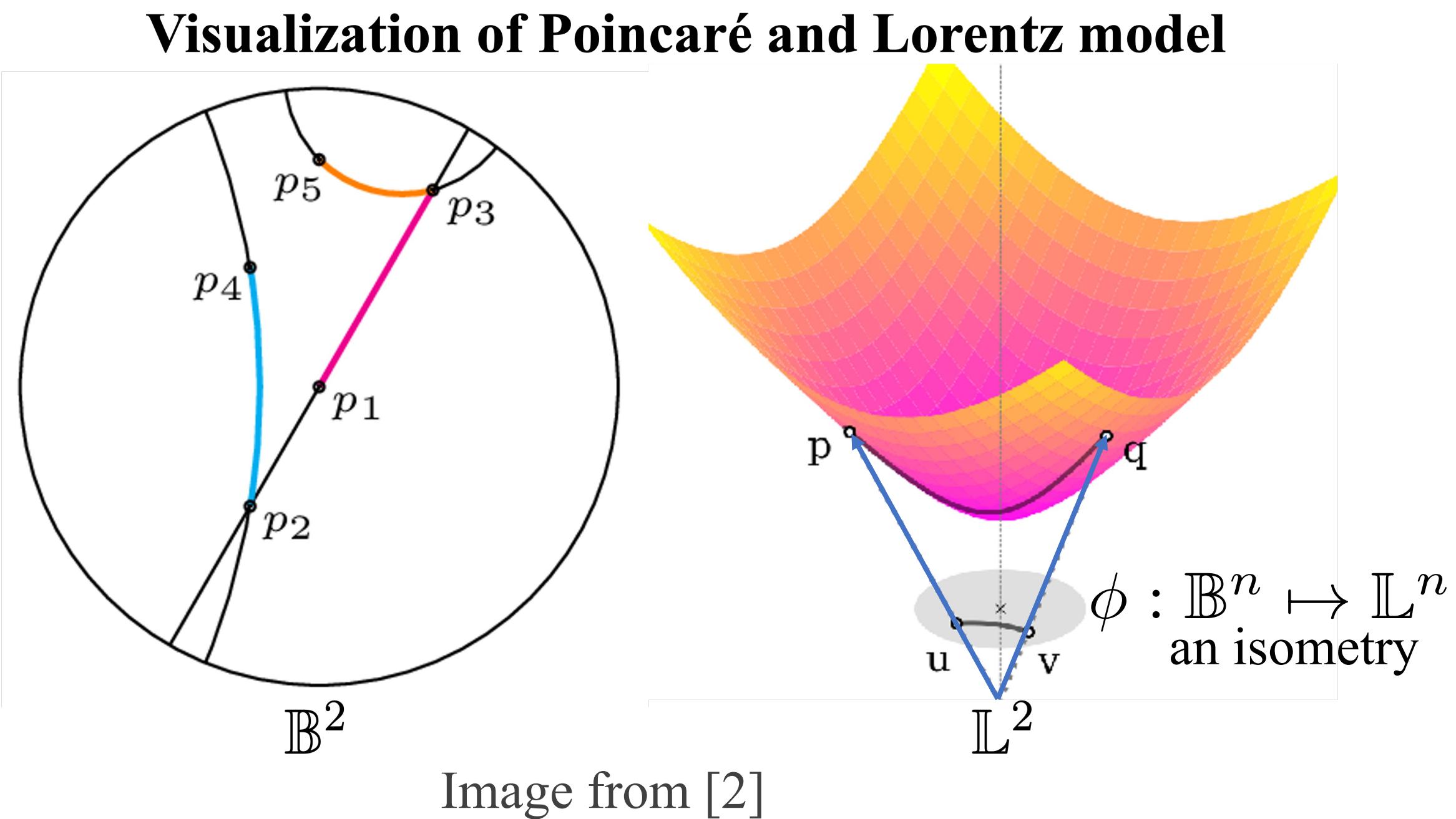
## Introduction

- Distances grow exponentially in hyperbolic spaces  $\mathbb{H}^n$
- Finite trees can be embedded into  $\mathbb{H}^n$  with arbitrarily low distortion<sup>1</sup>
- Motivates hyperbolic embeddings for hierarchical datasets such as words and images
- **Concern:** highly numerically unstable; mysteriously better performance of Lorentz than Poincaré despite being isometric.



## Main Contribution

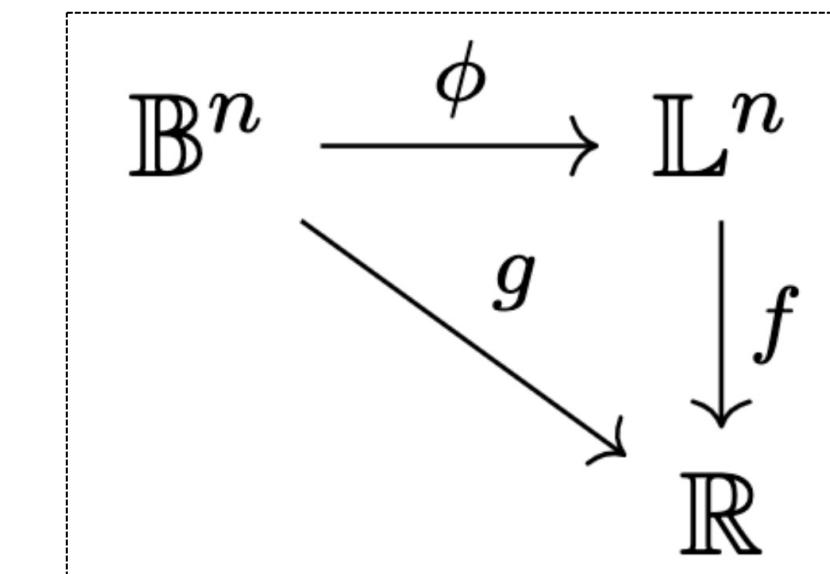
- Identify the source of numerical instability in two popular isometric Hyperbolic models, Lorentz  $\mathbb{L}^n$  and Poincaré  $\mathbb{B}^n$
- Clarify comparative advantage of  $\mathbb{L}^n$  and  $\mathbb{B}^n$ 
  - In **representation**,  $\mathbb{B}^n$  has a larger diameter than  $\mathbb{L}^n$
  - In **optimization**,  $\mathbb{L}^n$  has a larger correct update region
- Demonstrate empirical **superiority of  $\mathbb{L}^n$  in tree-embedding related tasks**



## Key Perspective

- By IEEE float64 standard, the addition of two numbers that differ by a magnitude of  $10^{16}$  gives just the large number.
- **Thm 1:** For a point  $\tilde{x} \in \mathbb{B}^n$  and  $x = \phi(\tilde{x}) \in \mathbb{L}^n$ , then  $\|x\| = \Omega(\frac{1}{1-\|\tilde{x}\|})$
- **Thm 2:** Consider a smooth function  $f : \mathbb{L}^n \mapsto \mathbb{R}$  and  $g := f \circ \phi$ , then for any  $x \in \mathbb{B}^n$ :  $\|x\| = 1 - \delta$  for small  $\delta$ ,

$$\begin{aligned} \|\nabla_{\mathbb{B}^n} f(x)\| &= \Omega(\delta^2 \|\nabla f(x)\|), \\ \|\nabla_{\mathbb{L}^n} g(\tilde{x})\| &= O(\|\nabla f(x)\|). \end{aligned}$$

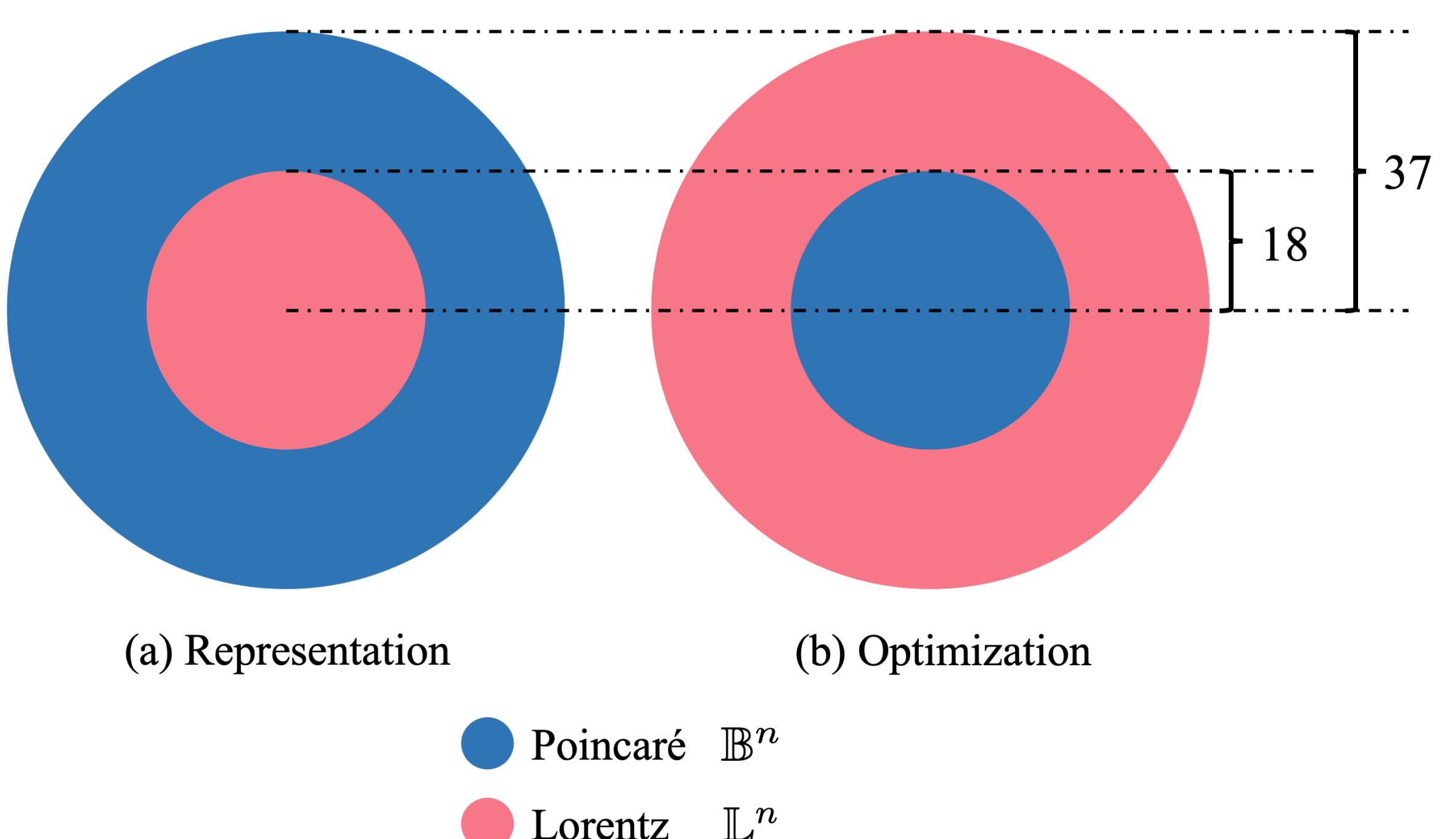


Relationship of definitions

## Interpretation

- **Thm 1** implies that  $\mathbb{B}^n$  and  $\mathbb{L}^n$  use similar number of bits to represent points
- A more refined analysis shows that  $\mathbb{B}^n$  has twice the diameter than  $\mathbb{L}^n$  under the IEEE constraint.
- **Thm 2** suggests that  $\mathbb{B}^n$  suffers more from a gradient vanishing problem when the magnitude is large due to its small Euclidean magnitude.

## Regions of Accurate Representation/Optimization



## Demonstration: Tree Embedding Optimization

- Create synthetic trees in  $\mathbb{R}^2$  as embeddings
  - Optimize using a proxy measure of distortion:  $d_E$  is the embedding distance and  $d_R$  the tree distances
- $$\mathcal{L}(\theta) = \frac{1}{|S|(|S|-1)} \sum_{x \neq y \in S} (\alpha \frac{d_E(x,y)}{d_R(x,y)} - 1)^2$$
- We use Riemannian Adam with fixed learning rate and number of epochs
  - Evaluate embedding using the average distortion

## Results<sup>3</sup>

- A good embedding comes with a large diameter  $d$  and small distortion  $\delta$
- In all tested trees,  $\mathbb{L}^n$  has larger diameter (more spread out) and smaller distortion
- The hierarchical structure is clearer in  $\mathbb{L}^n$  embeddings

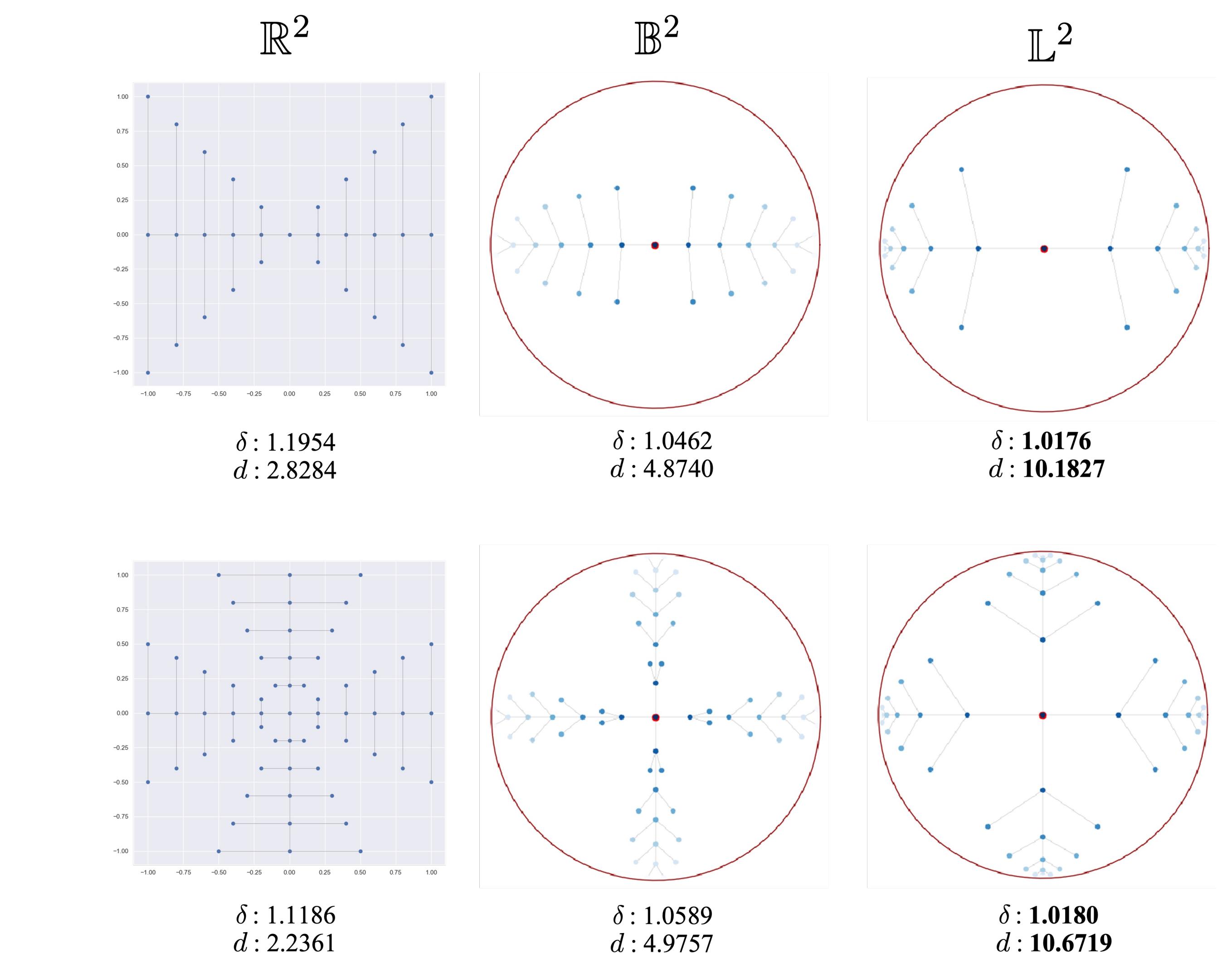


Table 1: Average distortion  $\delta$  and diameter  $d$  by Dataset and Model

| tree | manifold       | $\delta$      | $d$            | tree | manifold       | $\delta$      | $d$            |
|------|----------------|---------------|----------------|------|----------------|---------------|----------------|
| 1    | raw            | 1.1954        | 2.8284         | 5    | raw            | 1.3349        | 1.4142         |
|      | $\mathbb{D}^2$ | 1.0462        | 4.8740         |      | $\mathbb{D}^2$ | 1.0176        | 4.2523         |
|      | $\mathbb{L}^2$ | <b>1.0176</b> | <b>10.1827</b> |      | $\mathbb{L}^2$ | <b>1.0029</b> | <b>7.3850</b>  |
| 2    | raw            | 1.1186        | 2.2361         | 6    | raw            | 1.5080        | 1.4142         |
|      | $\mathbb{D}^2$ | 1.0589        | 4.9757         |      | $\mathbb{D}^2$ | 1.1243        | 4.3168         |
|      | $\mathbb{L}^2$ | <b>1.0180</b> | <b>10.6719</b> |      | $\mathbb{L}^2$ | <b>1.0392</b> | <b>8.9204</b>  |
| 3    | raw            | 1.0511        | 2.0396         | 7    | raw            | 1.2108        | 0.8321         |
|      | $\mathbb{D}^2$ | 1.0321        | 6.3564         |      | $\mathbb{D}^2$ | 1.2145        | 5.0657         |
|      | $\mathbb{L}^2$ | <b>1.0190</b> | <b>9.6830</b>  |      | $\mathbb{L}^2$ | <b>1.0166</b> | <b>19.1703</b> |
| 4    | raw            | 1.1539        | 2.8284         | 8    | raw            | 1.1421        | 0.9220         |
|      | $\mathbb{D}^2$ | 1.0754        | 5.0453         |      | $\mathbb{D}^2$ | 1.0324        | 4.5670         |
|      | $\mathbb{L}^2$ | <b>1.0421</b> | <b>10.4028</b> |      | $\mathbb{L}^2$ | <b>1.0157</b> | <b>8.8324</b>  |

## Acknowledgement

This work is partially supported by NSF under grants CCF-2112665, CCF-2217058, and by NIH under grant RF1MH125317.

## References

1. Sarkar, Rik. "Low distortion delaunay embedding of trees in hyperbolic plane." *International Symposium on Graph Drawing*. Springer, Berlin, Heidelberg, 2011.
2. Nickel, Maximilian, and Douwe Kiela. "Learning continuous hierarchies in the lorentz model of hyperbolic geometry." *ICML*. PMLR, 2018.
3. Mishne, Gal, et al. "The Numerical Stability of Hyperbolic Representation Learning." *arXiv preprint arXiv:2211.00181* (2022).