# DeepHGCN: Toward Deeper Hyperbolic Graph Convolutional Networks

Jiaxu Liu, Xinping Yi, and Xiaowei Huang

arXiv:2310.02027v4 [cs.LG] 19 Jul 2024

*Abstract*—Hyperbolic graph convolutional networks (HGCNs) have demonstrated significant potential in extracting information from hierarchical graphs. However, existing HGCNs are limited to shallow architectures due to the computational expense of hyperbolic operations and the issue of over-smoothing as depth increases. Although treatments have been applied to alleviate over-smoothing in GCNs, developing a hyperbolic solution presents distinct challenges since operations must be carefully designed to fit the hyperbolic nature. Addressing these challenges, we propose DeepHGCN, the first deep multi-layer HGCN architecture with dramatically improved computational efficiency and substantially reduced over-smoothing. DeepHGCN features two key innovations: (1) a novel hyperbolic feature transformation layer that enables fast and accurate linear mappings, and (2) techniques such as hyperbolic residual connections and regularization for both weights and features, facilitated by an efficient hyperbolic midpoint method. Extensive experiments demonstrate that DeepHGCN achieves significant improvements in link prediction and node classification tasks compared to both Euclidean and shallow hyperbolic GCN variants.

*Impact Statement*—Graph-structured data presents unique challenges in machine learning due to its complex nature. Traditional Graph Convolutional Networks (GCNs) primarily use Euclidean space for node embeddings, which struggle to capture the intricacies of hierarchical graphs. While HGCNs offer a solution, they often face limitations due to computational challenges and over-smoothing issues, particularly in deeper architectures. In this article, we introduce DeepHGCN, a novel approach for deep HGCNs. This architecture employs a new hyperbolic feature transformation layer and several auxiliary techniques to address these challenges. Testing on datasets with various geometries shows that DeepHGCN outperforms both Euclidean-based and other hyperbolic GCN methods in tasks such as link prediction and node classification. This demonstrates the efficacy of our approach in enhancing hierarchical graph-structured learning across various domains.

*Index Terms*—Graph neural networks, Riemannian manifold, hyperbolic operations, deep model architecture.

## I. INTRODUCTION

Graph convolutional networks (GCN) [1], [2], [3] have emerged as a promising approach for analyzing graph-structured data [4], *e.g.,* social networks [5], protein interaction networks [6], human skeletons [7], drug molecules [8], cross-modal retrieval [9], to name a few. Conventional GCN methods embed node representations into Euclidean latent space for downstream tasks. However, Bourgain's theorem [10] indicates that the Euclidean space with arbitrary dimensions fails to embed hierarchical graphs with low distortion, suggesting the inadequacy of the Euclidean space to accommodate complex hierarchical data [5], [11], [12], [13].

Recently, the hyperbolic space *a.k.a.* Riemannian manifold of constant negative sectional curvature [14], [15], [16], has gained increasing attention in processing non-Euclidean data. Since the exponentially expanding capacity of hyperbolic space satisfies the demand for hierarchical data that requires an exponential amount of branching space, embedding graphs to such a manifold naturally promotes learning hierarchical information. Based on two prevalent isomorphic models for hyperbolic space (Fig. 1), [17] and [18] introduced the basic operations for constructing hyperbolic neural networks (HNN). Subsequently, the researchers generalized GCN operations to hyperbolic domains and derived a series of hyperbolic graph convolutional network (HGCN) variants [19], [20], [21], [22], [23], [24]. These hyperbolic models are more capable of generating high-quality representations with low embedding dimensions, making them particularly advantageous in low-memory circumstances.

Despite their popularity, most HGCNs only achieve competitive performance with a 2-layer model. This limitation hinders their ability to effectively gather information from higher-order neighbors. However, developing a deeper HGCN model faces two main challenges. First, the computational complexity involved in hyperbolic operations, particularly feature transformation, prevents HGCNs from going deeper. Second, when more layers are added, node representations within the same connected component become increasingly indistinguishable.

This phenomenon is inherited from their GCN counterparts, known as *over-smoothing* [25], which could severely degrade the performance of multi-layer GCNs. For Euclidean GCNs, the over-smoothing issue has been defined and extensively studied in [26], [27], [28], [29], [30], where the proposed techniques ensure that deep Euclidean GCNs outperform shallow counterparts as the depth increases. Given the hyperbolic representation, it is evidenced empirically that HGCNs still suffer from over-smoothing issues as network depth increases.

Aiming to address above challenges, in this paper, we propose a HGCN variant that is adaptive to depth variation, namely the DeepHGCN (Fig. 4(b)), by stacking a number of carefully constructed HGCN layers with computationally-efficient feature transformation, which can effectively prevent over-smoothing and deliver improved accuracy over state-of-the-art models with a deep architecture. Our contributions toward deep HGCNs are summarized as follows:

- **Scalable and Efficient Backbone.** Dealing with the

J. Liu and X. Huang are with the School of Electrical Engineering, Electronics and Computer Science, University of Liverpool, Liverpool, L69 3GJ UK e-mail: {jiaxu.liu, xiaowei.huang}@liverpool.ac.uk.

X. Yi is with the National Mobile Communications Research Laboratory, Southeast University, Nanjing, China. email: xyi@seu.edu.cn.

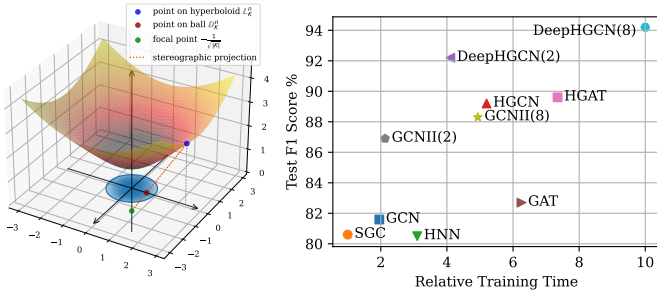This paragraph will include the Associate Editor who handled your paper.

Fig. 1. **Left**: Two prevalent hyperbolic models, isometric projection through the red line, where $\mathrm{P}_{\mathbb{D} \to \mathbb{L}}$: $\bullet \to \bullet$ and $\mathrm{P}_{\mathbb{L} \to \mathbb{D}}$: $\bullet \to \bullet$; **Right**: Performance over training time on *Airport* in 5k epochs. DeepHGCN(2-layer) outperforms existing hyperbolic models and is more efficient. Increasing depth to Deep-HGCN(8) bring further improvements.

computational complexity when increasing the depth of multi-layer HGCNs, we derive a novel hyperbolic fully connected layer that offers better efficiency and expressiveness for Poincaré feature transformation. In addition, Möbius gyromidpoint [31], [32] is also carefully incorporated as an accurate hyperbolic midpoint method that serves as the basis for not only message aggregation, but also all hyperbolic operations in our proposed techniques within the DeepHGCN architecture.

- **Extensive Techniques.** To address the over-smoothing issue, we generalize the concept of Dirichlet energy to the hyperbolic space, effectively tracking the smoothness of hyperbolic embeddings. Guided by the measure of hyperbolic Dirichlet energy, DeepHGCN is specifically powered by three techniques, namely the initial residual, weight alignment, and feature regularization. Evidently, DeepHGCN effectively alleviates the over-smoothing problem occurred in the Poincaré ball and can be naturally generalized to the Lorentz model.

- **Experiments and Ablation Studies.** Results on benchmark datasets (Tab. IV-VII) have validated the efficacy of our method for both node classification and link prediction under various layer settings. Additionally, ablation study in Sec. V-B4 demonstrate that all techniques are necessary for mitigating the over-smoothing issue.

The code is available at https://github.com/ljxw88/deephgcn

## II. BACKGROUND

### A. Brief Review of Riemannian Geometry

A *manifold* $\mathcal{M}$ is a topological space that is locally Euclidean. The *tangent space* $\mathcal{T}_{\mathbf{x}}\mathcal{M}$ is a real vector space of the same dimension as $\mathcal{M}$ attached to every point $\mathbf{x} \in \mathcal{M}$. All vectors in $\mathcal{T}_{\mathbf{x}}\mathcal{M}$ pass tangentially through $\mathbf{x}$. The *metric tensor* $g_{\mathbf{x}}$ at point $\mathbf{x}$ defines an inner product on the associated tangent space, *i.e.*, $g_{\mathbf{x}} : \mathcal{T}_{\mathbf{x}}\mathcal{M} \times \mathcal{T}_{\mathbf{x}}\mathcal{M} \to \mathbb{R}$. A *Riemannian manifold* $(\mathcal{M}, g)$ is defined as a manifold equipped with Riemannian metric $g$. The metric tensor provides local geometric properties such as angles and curve lengths. A *geodesic* is the shortest path between two points on the manifold. The *exponential map* $\exp_{\mathbf{x}} : \mathcal{T}_{\mathbf{x}}\mathcal{M} \to \mathcal{M}$ defines a mapping of a tangent vector to a point on the manifold, and the *logarithmic map* is the

inverse $\log_{\mathbf{x}} : \mathcal{M} \to \mathcal{T}_{\mathbf{x}}\mathcal{M}$. Extended reviews are provided in Appendix A.

### B. Hyperbolic Geometry

The hyperbolic space $\mathbb{H}^n_\kappa$ is a smooth Riemannian manifold with a constant sectional curvature $\kappa < 0$ [33]. It is usually defined via five isometric hyperbolic models, among which the $n$-dimensional Poincaré ball model $\mathbb{D}^n_\kappa = (\mathcal{D}^n_\kappa, g^{\mathbb{D}})$ and the Lorentz model (hyperboloid) $\mathbb{L}^n_\kappa = (\mathcal{L}^n_\kappa, g^{\mathbb{L}})$ are frequently used. The manifolds of $\mathbb{D}^n_\kappa$ are the projections of $\mathbb{L}^n_\kappa$ onto the $n$-dimensional space-like hyperplanes (Fig. 1 left). In this paper, we align with [17] and build our model upon the Poincaré ball model. As we still need the knowledge of the Lorentz model to complete our theories, we delegate the instruction and basic operations to Appendix A and Tab. VIII.

*Poincaré Ball Model.* The $n$-dimensional Poincaré ball is defined as the Riemannian manifold $\mathbb{D}^n_\kappa = (\mathcal{D}^n_\kappa, g^{\mathbb{D}})$, with point set $\mathcal{D}^n_\kappa = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < -\frac{1}{\kappa}\}$ and Riemannian metric $g^{\mathbb{D}}_\kappa(\mathbf{x}) = (\lambda^\kappa_{\mathbf{x}})^2 g^{\mathbb{E}}$, where $\lambda^\kappa_{\mathbf{x}} = \frac{2}{1 + \kappa\|\mathbf{x}\|^2}$ (the conformal factor) and $g^{\mathbb{E}} = \mathbf{I}_n$. The Poincaré metric tensor induces various geometric properties, *e.g.*, distances $d^\kappa_{\mathbb{D}}(\mathbf{x}, \mathbf{y})$, inner products $\langle \mathbf{u}, \mathbf{v} \rangle^\kappa_{\mathbf{x}}$, geodesics $\gamma_{\mathbf{x},\mathbf{v}}(t)$ [34], and more. The geodesics also induce the definition of exponential and logarithmic maps. which are denoted at $\mathbf{x} \in \mathbb{D}^n_\kappa$ as $\exp^\kappa_{\mathbf{x}}$ and $\log^\kappa_{\mathbf{x}}$, respectively. The Möbius gyrovector space [16] offers an algebraic framework to treat the Poincaré coordinates as vector-like mathematical objects (gyrovectors). The gyrovectors are equipped with series of operations, *e.g.* the vector addition $\oplus_\kappa$ and matrix-vector multiplication $\otimes_\kappa$. For brevity, we give instruction to Poincaré operations in Appendix A and Tab. VIII.

## III. AUGMENTED HGCN BACKBONE

### A. Efficient Feature Transformation

Training multi-layer GCNs requires a fast and accurate linear transformation $\mathcal{F} : \mathbb{R}^{d_1} \to \mathbb{R}^{d_2}$ as backbone. However, an obvious deficiency of traditional hyperbolic linear layer is the propagation efficiency. In this paper, we propose a better unified linear layer $\mathcal{F}^\kappa_{\mathbb{D}}$ for efficient feature transformation within the Poincaré ball. With synthetic hyperbolic dataset, we show that $\mathcal{F}^\kappa_{\mathbb{D}}$ is faster and more expressive than both the naive HNN [17] and PFC layer [31].

**Theorem 1.** *Given* $\mathbf{h} \in \mathbb{D}^{d_1}_\kappa$, *Euclidean weight and bias parameter* $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$ *and* $\mathbf{b}_1, \mathbf{b}_2 \in \mathbb{R}^{d_2}$. *A more computational-efficient and expressive feature transformation* $\mathcal{F}^\kappa_{\mathbb{D}} : \mathbb{D}^{d_1}_\kappa \to \mathbb{D}^{d_2}_\kappa$ *within high dimensional Poincaré ball can be formulated by*

$$\mathcal{F}^\kappa_{\mathbb{D}}(\mathbf{h}; \mathbf{W}, \mathbf{b}) := \frac{\phi(\mathbf{h}; \mathbf{W}, \mathbf{b})}{1 + \sqrt{|\kappa|\|\phi(\mathbf{h}; \mathbf{W}, \mathbf{b})\|^2 + 1}}, \quad (1)$$

*where* $\phi(\cdot)$ *is formulated as*

$$\phi(\mathbf{h}; \mathbf{W}, \mathbf{b}) = \frac{2\sqrt{|\kappa|}\mathbf{W}\mathbf{h} + \mathbf{b}_1(1 - \kappa\|\mathbf{h}\|^2)}{\sqrt{|\kappa|}(1 + \kappa\|\mathbf{h}\|^2)} + \mathbf{b}_2. \quad (2)$$

Given a point $\mathbf{x}$ on the hyperboloid, an arbitrary transformation matrix $\mathbf{W} \in \mathbb{R}^{d_2 \times (d_1+1)}$ can be multiplied to $\mathbf{x}$. The Lorentzian *re-normalization trick* ensures $\mathbf{x}$ still lies in

the hyperboloid. Specifically, given $\mathbf{x} \in \mathbb{L}_\kappa^{d_1}(\subset \mathbb{R}^{d_1+1})$ and $\mathbf{W} \in \mathbb{R}^{d_2 \times (d_1+1)}$, multiplying them gives $\mathbf{W}\mathbf{x} \in \mathbb{R}^{d_2}$ which is not essentially in the hyperboloid. To force the Lorentz constraint $\forall \mathbf{x} \in \mathbb{L}_\kappa : \langle \mathbf{x}, \mathbf{x} \rangle_\mathcal{L} = \frac{1}{\kappa}$, we re-normalize the time axis as $\sqrt{\|\mathbf{W}\mathbf{x}\|^2 - \frac{1}{\kappa}}$, such that the point set constraint is not violated. Thus a simple Lorentz transformation with re-normalization trick is expressed as

$$\mathbf{x}' \leftarrow \begin{bmatrix} \sqrt{\|\mathbf{W}\mathbf{x}\|^2 - \frac{1}{\kappa}} \\ \mathbf{W}\mathbf{x} \end{bmatrix}. \tag{3}$$

In a more general setting with bias parameter, $\mathbf{W}\mathbf{x}$ can be $\phi(\mathbf{x}, \mathbf{W}, \mathbf{b}) : \mathbb{L}^{d_1} \to \mathbb{R}^{d_2}$. It is easy to verify the Lorentzian constraint of Eq. (3). Next, we give the bijection between an arbitrary point on the hyperboloid $\mathbf{z} = \begin{bmatrix} z_t \\ \mathbf{z}_s \end{bmatrix} \in \mathbb{L}_\kappa^n$ and the corresponding point on the Poincaré ball $\mathbf{x} \in \mathbb{D}_\kappa^n$ (Fig. 1 left) as follows

$$\mathbb{L}_\kappa^n \to \mathbb{D}_\kappa^n : \mathrm{P}_{\mathbb{L} \to \mathbb{D}}(\mathbf{z}) = \frac{\mathbf{z}_s}{1 + \sqrt{|\kappa|} z_t}, \tag{4}$$

$$\mathbb{D}_\kappa^n \to \mathbb{L}_K^n : \mathrm{P}_{\mathbb{D} \to \mathbb{L}}(\mathbf{x}) = \begin{bmatrix} \frac{1 - \kappa \|\mathbf{x}\|^2}{\sqrt{|\kappa|} + \kappa \sqrt{|\kappa|}\|\mathbf{x}\|^2} \\ \frac{2\mathbf{x}}{1 + \kappa \|\mathbf{x}\|^2} \end{bmatrix}. \tag{5}$$

Therefore using Eq. (5), given a point $\mathbf{x} \in \mathbb{D}_\kappa^{d_1}$, we derive the corresponding point on the hyperboloid $\hat{\mathbf{x}} = \begin{bmatrix} \frac{1 - \kappa \|\mathbf{x}\|^2}{\sqrt{|\kappa|} + \kappa \sqrt{|\kappa|}\|\mathbf{x}\|^2} \\ \frac{2\mathbf{x}}{1 + \kappa \|\mathbf{x}\|^2} \end{bmatrix} \in \mathbb{L}_\kappa^{d_1}(\subset \mathbb{R}^{d_1+1})$. Employing Eq. (3), with a transformation matrix $\mathbf{W} \in \mathbb{R}^{d_2 \times (d_1+1)}$, we get $\mathbf{x}' \leftarrow \begin{bmatrix} \sqrt{\|\phi(\hat{\mathbf{x}})\|^2 - \frac{1}{\kappa}} \\ \phi(\hat{\mathbf{x}}) \end{bmatrix} \in \mathbb{L}_\kappa^{d_2}(\subset \mathbb{R}^{d_2+1})$, where $\phi(\hat{\mathbf{x}}) = \mathbf{W}\hat{\mathbf{x}} + \mathbf{b}$. Applying the reverse mapping Eq. (4) gives

$$\mathrm{P}_{\mathbb{L} \to \mathbb{D}}(\mathbf{x}') = \phi(\hat{\mathbf{x}}) \left( 1 + \sqrt{|\kappa| \|\phi(\hat{\mathbf{x}})\|^2 - \mathrm{sgn}(\kappa)} \right)^{-1}, \tag{6}$$

which gives us the form in Eq. (1) as $\kappa < 0$ in the Poincaré ball model. Since $\mathbf{W}$ is an arbitrary parameter in $\mathbb{R}$, we slice $\mathbf{W} \in \mathbb{R}^{d_2 \times (d_1+1)}$ as concat $[W_t \in \mathbb{R}^{d_2} \| \mathbf{W}_s \in \mathbb{R}^{d_2 \times d_1}]$, therefore term $\mathbf{W}\hat{\mathbf{x}}$ can be reformulated as

$$\frac{2\sqrt{|\kappa|}\mathbf{W}_s \mathbf{x} + W_t(1 - \kappa \|\mathbf{x}\|^2)}{\sqrt{|\kappa|} + \kappa \sqrt{|\kappa|}\|\mathbf{x}\|^2}, \tag{7}$$

which arrives at the form in Eq. (2). Eq. (6-7) together concludes Thm. 1.

**Proposition 2.** *Given the $i$-th node representation $\mathbf{h}_i \in \mathbb{D}_\kappa^{d_1}$, the hyperbolic feature transformation in Thm. 1 $\mathbf{h}_i \leftarrow \mathcal{F}_\mathbb{D}^\kappa(\mathbf{h}_i; \mathbf{W}, \mathbf{b})$ yields more expressive node embeddings.*

To verify Prop. 2, we seek the connection of our approach to the formulation of PFC layer [31], which is expressed as

$$\mathbf{x}' \leftarrow \mathcal{F}_{\mathrm{PFC}}^\kappa(\mathbf{x}; \mathbf{W}, \mathbf{b}) := \frac{\boldsymbol{\omega}}{1 + \sqrt{|\kappa| \|\boldsymbol{\omega}\|^2 + 1}}, \tag{8}$$

where $\boldsymbol{\omega} := \left( \frac{1}{\sqrt{|\kappa|}} \sinh\left(\sqrt{|\kappa|}\nu_j^\kappa(\mathbf{x})\right) \right)_{j=1}^{d_2}. \tag{9}$

where $\nu_i^\kappa(\mathbf{x})$ is the unidirectional re-generalization of hyperbolic multinomial logistic regression.

**Proposition 3.** *Given $\mathbf{x} \in \mathbb{D}_\kappa^{d_1}$ and $\mathbf{x}' = \mathcal{F}_{\mathrm{PFC}}^\kappa(\mathbf{x}) \in \mathbb{D}_\kappa^{d_2}$, the corresponding point of $\mathbf{x}'$ on the hyperboloid via stereographic projection is $\mathbf{h}' = \begin{bmatrix} h_t \\ \mathbf{h}_s \end{bmatrix} \in \mathbb{L}_\kappa^{d_2+1}$ where $\mathbf{h}_s = \boldsymbol{\omega}$ and $h_t = \sqrt{\|\boldsymbol{\omega}\|^2 - \frac{1}{\kappa}}$.*

*Proof:* We start with $\mathbf{h} \leftarrow \mathcal{F}_{\mathrm{PFC}}^\kappa(\mathbf{x})$ (where $\mathbf{x} \in \mathbb{D}_\kappa^{d_1}$ and $\mathbf{h} \in \mathbb{D}_\kappa^{d_2}$ are respectively the input and output within the Poincaré ball) in Eq. (8), applying stereographic projection in Eq. (5) we obtain the corresponding point on the hyperboloid $\mathbf{h}^\mathbb{L} = \begin{bmatrix} h_t \\ \mathbf{h}_s \end{bmatrix} \in \mathbb{L}_\kappa^{d_2+1}$, where

$$\mathbf{h}_s = 2 \frac{\mathbf{h}}{1 + \kappa \|\mathbf{h}\|^2} \tag{10}$$

$$= \frac{\frac{2\boldsymbol{\omega}}{\gamma}}{1 + \kappa \|\frac{\boldsymbol{\omega}}{\gamma}\|^2} \quad \text{where } \gamma = 1 + \sqrt{1 - \kappa \|\boldsymbol{\omega}\|^2} \tag{11}$$

$$= \frac{2\boldsymbol{\omega}\gamma}{\gamma^2 + \kappa \|\boldsymbol{\omega}\|^2} \tag{12}$$

$$= \frac{\boldsymbol{\omega}(2 + 2\sqrt{1 - \kappa \|\boldsymbol{\omega}\|^2})}{1 + 2\sqrt{1 - \kappa \|\boldsymbol{\omega}\|^2} + 1 - \kappa \|\boldsymbol{\omega}\|^2 + \kappa \|\boldsymbol{\omega}\|^2} \tag{13}$$

$$= \frac{\boldsymbol{\omega}(2 + 2\sqrt{1 - \kappa \|\boldsymbol{\omega}\|^2})}{2 + 2\sqrt{1 - \kappa \|\boldsymbol{\omega}\|^2}} = \boldsymbol{\omega}, \tag{14}$$

and similarly

$$h_s = \frac{1 - \kappa \|\mathbf{h}\|^2}{\sqrt{|\kappa|} + \kappa \sqrt{|\kappa|}\|\mathbf{h}\|^2} = \sqrt{\|\boldsymbol{\omega}\|^2 - \frac{1}{\kappa}}. \tag{15}$$

This concludes the proof. ■

According to Prop. 3, geometrically, the $\boldsymbol{\omega}$ in Eq. (8) can be interpreted as a special feature transformation of the spatial component $\mathbf{h}_s$. The time component $h_t$ is a re-normalization according to $\mathbf{h}_s$ which stabilizes the point on the corresponding hyperboloid. The formulation in Eq. (8) and, identically, Eq. (1) ensures that **any** definition of $\boldsymbol{\omega}$ will **not** violate the Lorentzian constraint.

In essence, PFC breaks down to three core stages: *1) Project Poincaré ball points to their hyperboloid equivalents via Eq. (5)*; *2) Apply transformation $\boldsymbol{\omega}$ to the spatial component and **re-normalize** the time segment*; *3) Reposition back to the Poincaré ball via Eq. (4)*. We hold the view that step-2, thanks to the re-normalization trick, with arbitrary $\boldsymbol{\omega}$, the transformed points will still adhere to the point set constraint (*i.e.* $\{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < -\frac{1}{\kappa}\}$) after stereographic projection. Consequently, $\boldsymbol{\omega}$ can be an arbitrary linear transformation. Therefore, we argue to let $\boldsymbol{\omega} = \mathrm{MLP}(\mathbf{x})$ be a Euclidean neural net instead of Eq. (9), leading us to the expression in Eq. (2).

**Performance Evaluation.** In Fig. III-A, we illustrate the decision hyperplane of different single transformation layers with randomly sampled two batches of points in $\mathbb{D}_\kappa^2$. We observe that the Euclidean hyperplane fail to adapt to non-linearity of the data samples due to the linear nature. The two hyperbolic baselines, limited by strict hyperbolic constraints on their formulations, also fails in fitting the data. Such a phenomenon is especially obvious on Fig. III-A.(b, f and g). In comparison, our reformulated FC layer offers the best performance on fitting the data samples.

(a) Euc (97.3%) (b) HNN (94.7%) (c) PFC (98.8%) (d) Ours (99.8%)

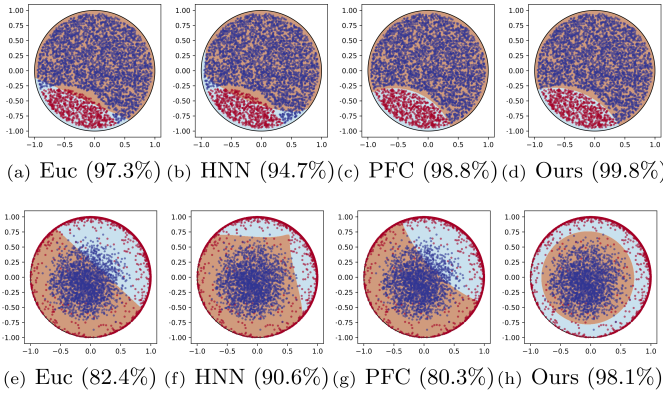(e) Euc (82.4%) (f) HNN (90.6%) (g) PFC (80.3%) (h) Ours (98.1%)

Fig. 2. Decision hyperplane of various feature transformations on synthetic binary classification tasks (Task #1: (a)-(d) and Task #2: (e)-(h)).

TABLE I
COMPARISON OF ACCURACY AND CALCULATION TIME (MS) OF VARIOUS HYPERBOLIC FC LAYERS. 2K+2K POINTS ARE SAMPLED WITHIN HIGH DIMENSIONAL POINCARÉ BALL. WE REPORT MEAN $\pm$ STD.

| Task | Euclidean | | HNN | | PFC | | Ours | |
|---|---|---|---|---|---|---|---|---|
| | Acc | Time | Acc | Time | Acc | Time | Acc | Time |
| # 1 | $97.3_{\pm 3.2\text{e-}04}$ | $81.8_{\pm 4.1\text{e}+01}$ | $95.6_{\pm 5.8\text{e-}04}$ | $465.6_{\pm 3.1\text{e}+02}$ | $99.8_{\pm 1.7\text{e-}02}$ | $291.3_{\pm 1.8\text{e}+01}$ | $100.0_{\pm 2.6\text{e-}04}$ | $197.6_{\pm 4.9\text{e}+00}$ |
| # 2 | $82.6_{\pm 4.3\text{e-}03}$ | $77.7_{\pm 1.2\text{e}+00}$ | $90.6_{\pm 2.2\text{e-}16}$ | $447.7_{\pm 6.4\text{e}+01}$ | $94.8_{\pm 5.7\text{e-}02}$ | $290.0_{\pm 1.5\text{e}+01}$ | $98.2_{\pm 5.3\text{e-}04}$ | $198.7_{\pm 6.0\text{e}+00}$ |

**Computation Cost Analysis.** Assume the feature $\mathbf{h}$ is in $d_1$-dimension, and the feature transformation matrix $\mathbf{W} \in \mathbb{R}^{d_2 \times d_1}$. Theoretically, the complexity of HNN and our method are both $\mathcal{O}(d_1 \times d_2 + d_1 + 2d_2)$, which is similar to that of the Euclidean linear transformation $\mathbf{Wh} + \mathbf{b}$ ($\mathcal{O}(d_1 \times d_2 + d_2)$ complexity). Below, we provide detailed complexity analysis, and give the reason why the computation cost of HNN blows up while our method costs similar to standard MLP.

As detailed in Appendix A-B, the HNN consists of exponential map, matrix-vector multiplication, and logarithmic map. We consider the simplest situation where the base is the north pole $\mathbf{0}$, consequently, the Möbius additions $\oplus_\kappa$ are canceled and the conformal factor $\lambda_\mathbf{0}^\kappa$ has $\mathcal{O}(1)$ complexity. Thus $\log_\mathbf{0}^\kappa(\mathbf{h})$ and $\exp_\mathbf{0}^\kappa(\mathbf{h})$ as in Eq. (29-28) both have $\mathcal{O}(d)$ complexity, meaning the transformation $\exp_\mathbf{0}^\kappa(\mathbf{W} \log_\mathbf{0}^\kappa(\mathbf{h}))$ is of $\mathcal{O}(d_1 \times d_2 + d_1 + d_2)$ complexity. However, when considering bias translation, the $\oplus_\kappa$ must be accounted. As shown in Eq. (27), $\oplus_\kappa$ basically consists inner products, norms and vector addition, all of them have $\mathcal{O}(d)$ complexity. The division has $\mathcal{O}(1)$ complexity, so overall $\oplus_\kappa$ is of $\mathcal{O}(4d)$ complexity. Therefore, the HNN $\exp_\mathbf{0}^\kappa(\mathbf{W} \log_\mathbf{0}^\kappa(\mathbf{h})) \oplus_\kappa \mathbf{b}$ has $\mathcal{O}(d_1 \times d_2 + d_1 + 4d_2)$ complexity. Moreover, as evidenced in [35], [36], the hyperbolic trigonometric functions substantially increase the GPU burden on parallelism, which result in $4 \sim 10$ times slower computation. With the repeat usage of $\tanh$ and $\tanh^{-1}$ in exp/log maps, the practical time consumption can be sufficiently higher than our theoretical analysis. Lastly, we analyze the complexity of our method. Our method belongs to the same family with PFC, hence similar complexities. In particular, the $\phi(\cdot)$ in Eq. (2) consists of matrix/scalar-vector multiplications, norm and addition, the complexity is $\mathcal{O}(d_1 \times d_2 + d_1 + 2d_2)$. With the result $\phi(\mathbf{h}; \mathbf{W}, \mathbf{b})$, the complexity of Eq. (1) (with only norm and



(a) Tangential (n=2) (b) Gyromid (n=2) (c) Fréchet (n=2)

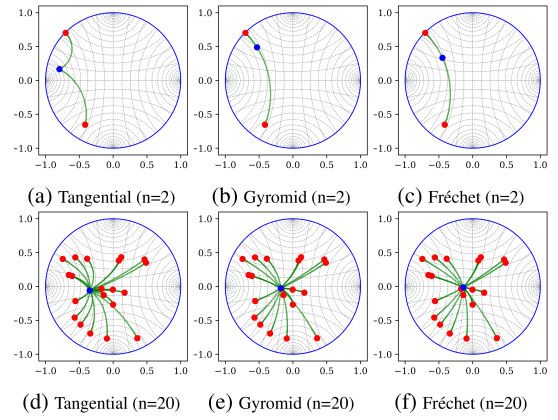(d) Tangential (n=20) (e) Gyromid (n=20) (f) Fréchet (n=20)

Fig. 3. Compare tangential midpoint [20], Möbius gyromidpoint [16] and differentiable Fréchet mean [37] in the 2-dimension Poincaré disk. For each method we illustrate the weighted midpoint (blue) for double and multiple randomly sampled points (red) with randomly initialized weights.

TABLE II
COMPARE VARIOUS HYPERBOLIC AVERAGING METHODS ON PRECISION AND CALCULATION TIME (MS) VS THE BASELINE (1000-ITER FRECHÉT MEAN). 4K POINTS ARE RANDOMLY SAMPLED WITHIN HIGH DIMENSIONAL POINCARÉ BALL. WE REPORT MEAN $\pm$ STD.

| Dim | Baseline | Tangential | | **Möbius Gyro** | | Frechét-Early | |
|---|---|---|---|---|---|---|---|
| | Time | MSE | Time | MSE | Time | MSE | Time |
| 8 | $152.4_{\pm 0.5}$ | $4.8\text{e-}5_{\pm 4.7\text{e-}6}$ | $3.4_{\pm 0.3}$ | $1.5\text{e-}6_{\pm 1.5\text{e-}7}$ | $2.3_{\pm 0.3}$ | $6.2\text{e-}29_{\pm 2.8\text{e-}29}$ | $2.3_{\pm 0.0}$ |
| 16 | $219.4_{\pm 1.1}$ | $6.5\text{e-}5_{\pm 4.7\text{e-}6}$ | $2.8_{\pm 0.3}$ | $9.1\text{e-}7_{\pm 6.8\text{e-}8}$ | $1.7_{\pm 0.2}$ | $6.9\text{e-}30_{\pm 1.9\text{e-}30}$ | $2.8_{\pm 0.1}$ |
| 64 | $371.1_{\pm 15.1}$ | $3.5\text{e-}5_{\pm 1.2\text{e-}6}$ | $3.3_{\pm 0.4}$ | $2.4\text{e-}10_{\pm 8.9\text{e-}12}$ | $2.2_{\pm 0.3}$ | $3.1\text{e-}31_{\pm 1.4\text{e-}31}$ | $3.9_{\pm 0.2}$ |

devision) is $\mathcal{O}(d_1 \times d_2 + d_1 + 4d_2)$.

From above we conclude that, all methods can achieve an approximately $\mathcal{O}(d_1 \times d_2)$ complexity. The HNN, however, requires heavy usage of trigonometric funcions, which are computationally costly as evidenced in the literature, thus hindering its scalability. Our approach, instead, requires no hyperbolic trigonometric functions, thus leads to similar cost as Euclidean MLP. In Tab. I, we illustrate the performance of transformation layers and observe that our approach achieves the best accuracy on two simple classification tasks. Notably, our approach requires approximately 2 times of the computation time of Euclidean linear layer, while other approaches require 3 times and more. This coincides with our analysis.

### B. Efficient Message Aggregation

**Definition 1** (Möbius gyromidpoint [16]). *Given the gyrovectors $\{\mathbf{x}_i \in \mathbb{D}_\kappa^d\}_{i=1}^N$ and the weights $\{w_i \in \mathbb{R}\}_{i=1}^N$, the weighted gyromidpoint in the Poincaré ball $f_{\mathrm{MG}}^\kappa : \mathbb{D}_\kappa^{N \times d} \times \mathbb{R}^N \to \mathbb{D}_\kappa^d$ is defined as*

$$f_{\mathrm{MG}}^\kappa(\{\mathbf{x}_i\}_{i=1}^N, \{w_i\}_{i=1}^N) = \frac{1}{2} \otimes_\kappa \left( \frac{\sum_{i=1}^N w_i \lambda_{\mathbf{x}_i}^\kappa \mathbf{x}_i}{\sum_{i=1}^N |w_i|(\lambda_{\mathbf{x}_i}^\kappa - 1)} \right).$$

The mean operator is an essential building block for neural networks. In non-Euclidean geometries, the mean computation cannot be performed simply by averaging the inputs, as the averaged vector may be out of manifold. Basically, there are three types of generalized weighted mean to hyperbolic space that can be used to guarantee the summed vectors on the

TABLE III
PEAK MEMORY USAGE COMPARISON (4K POINT, 64 DIM, MEAN±STD)

| Method | Frct-MaxIter | Tangential | **Möbius Gyro** | Frct-Early |
|---|---|---|---|---|
| Peak Mem (KiB) | $2727.8_{\pm 82.4}$ | $536.9_{\pm 24.6}$ | $\mathbf{411.1_{\pm 85.6}}$ | $2670.4_{\pm 110.8}$ |

manifold and to be differentiable, namely, the tangential aggregation [20], hyperbolic gyromidpoint [16] and differentiable Frechét mean [37]. In this work, we employ the hyperbolic gyromidpoint as the unified faster and accurate mean operator, defined in Def. 1. With Def. 1, we define the convolution for hyperbolic node feature $\mathbf{h}^{(l)} \in \mathbb{D}_\kappa^{|\mathcal{V}| \times d}$ upon Möbius gyromidpoint. Given the augmented normalized adjacency matrix $\tilde{\mathbf{P}} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$, we define our message aggregation as

$$f_{\mathrm{NA}}^\kappa(\tilde{\mathbf{P}}, \mathbf{h}) = \left( \frac{1}{2} \otimes_\kappa \left( \frac{\sum_{j=1}^d \tilde{\mathbf{P}}_i \lambda_{\mathbf{h}_j}^\kappa \mathbf{h}_j}{\sum_{j=1}^d |\tilde{\mathbf{P}}_i| (\lambda_{\mathbf{h}_j}^\kappa - 1)} \right) \right)_{i=1}^{|\mathcal{V}|}. \quad (16)$$

**Performance Evaluation.** We demonstrate the precision and calculation time of three methods in Tab. II, regarding the standard Fréchet mean [38], [39] as baseline. We observe: (1) the tangential aggregation is efficient but inaccurate (comparing Fig. 3.a to 3.b and 3.c), where the resulting midpoint is no longer on the geodesic between the inputs; (2) computing the differentiable Fréchet mean [37] requires an iterative solver, and obtaining an accurate result requires considerable computation; (3) Möbius gyromidpoint gives a close solution to the Fréchet mean while significantly reducing the complexity. Furthermore, we provide *memory evaluation* in Tab. III, where we report peak memory usage comparison on the same synthetic dataset under 1000 runs. The employed gyromidpoint approach exhibits the lowest memory cost, which highlights its advantage in computation efficiency and scalability, allowing us to further expand our model to a multi-layer scheme.

## IV. TOWARD DEEPER HYPERBOLIC GCN

### A. Over-Smoothness Analysis

First, we derive the hyperbolic Dirichlet energy $f_{\mathrm{DE}}^\mathbb{D}(\cdot)$ as a measure of smoothness for the Poincaré embeddings.

**Definition 2.** *Given the embedding* $\mathbf{h} = \{\mathbf{h}_i \in \mathbb{D}_\kappa^d\}_{i=1}^{|\mathcal{V}|}$, *the hyperbolic Dirichlet energy* $f_{\mathrm{DE}}^\kappa(\mathbf{h})$ *is defined as*

$$\frac{1}{2} \sum_{(i,j) \in \mathcal{E}} d_\mathbb{D}^\kappa \left( \exp_\mathbf{o}^\kappa \left( \frac{\log_\mathbf{o}^\kappa(\mathbf{h}_i)}{\sqrt{1+d_i}} \right), \exp_\mathbf{o}^\kappa \left( \frac{\log_\mathbf{o}^\kappa(\mathbf{h}_j)}{\sqrt{1+d_j}} \right) \right)^2,$$

*where* $d_{i/j}$ *denotes the node degree of node* $i/j$. *The distance* $d_\mathbb{D}^\kappa(\mathbf{x}, \mathbf{y})$ *between two points* $\mathbf{x}, \mathbf{y} \in \mathbb{D}$ *is the geodesic length, we detail the closed form expression in Appendix A.*

**Proposition 4.** *Hyperbolic message aggregation reduces the Dirichlet energy. i.e.* $f_{\mathrm{DE}}^\kappa(\tilde{\mathbf{P}} \otimes_\kappa \mathbf{h}^{(l)}) \leq f_{\mathrm{DE}}^\kappa(\mathbf{h}^{(l)})$.

The Dirichlet energy of node representation in each layer can be viewed as the weighted sum of the distances between normalized node pairs. Prop. 4 indicates that the energy of node representation will decay after the aggregation step. When multiple aggregations in HGCN are performed, the
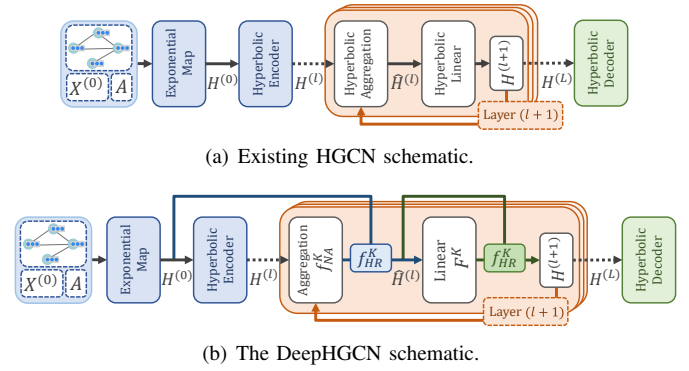


(a) Existing HGCN schematic.



(b) The DeepHGCN schematic.

Fig. 4. Comparison between the existing HGCN architecture and the proposed DeepHGCN. At the $l$-th layer, (a) performs linear transformation directly after the aggregation and regards the transformed feature as next layer's input, causing over-smoothing as $l$ increases; (b) performs hyperbolic residual connection after aggregation and linear layer to alleviate over-smoothing, such that the hyperbolic residual operator retains the feature on the manifold and the global hyperbolic geometry is preserved.

energy will converge to zero, indicating lower embedding expressiveness, which may lead to oversmoothing.

### B. Hyperbolic Graph Residual Connection

We define the hyperbolic graph residual connection $f_{\mathrm{HR}}^\kappa$ upon the notion of Möbius gyromidpoint. Given the node embedding matrices $\mathbf{h}^1, \mathbf{h}^2 \in \mathbb{D}_K^{|\mathcal{V}| \times d}$ and residual weight coefficients $w^1, w^2 \in \mathbb{R}$, define

$$f_{\mathrm{HR}}^\kappa(\mathbf{h}^1, \mathbf{h}^2; w^1, w^2)$$
$$= \left( \frac{1}{2} \otimes_\kappa \left( \frac{w^1 \lambda_{\mathbf{h}_i^1}^\kappa \mathbf{h}_i^1 + w^2 \lambda_{\mathbf{h}_i^2}^\kappa \mathbf{h}_i^2}{|w^1|(\lambda_{\mathbf{h}_i^1}^\kappa - 1) + |w^2|(\lambda_{\mathbf{h}_i^2}^\kappa - 1)} \right) \right)_{i=1}^{|\mathcal{V}|}. \quad (17)$$

This operation ensures the feature after residual connection still lies in the Poincaré ball. One also recovers the arithmetic mean as $\kappa \to 0$.

*1) Hyperbolic Initial Residual:* The graph residual connection is firstly introduced in the standard GCN [2], in which the current layer representation $\sigma(\tilde{\mathbf{P}}\mathbf{h}^{(l)}\mathbf{W}^{(l)})$ is connected to previous layer $\mathbf{h}^{(l)}$ to facilitate a deeper model. Some works [40], [26] empirically find the efficacy of adding residual connections to initial layer $\mathbf{h}^{(0)}$. It was also claimed in [29] that residual connections to both initial layer $\mathbf{h}^{(0)}$ and previous layer $\mathbf{h}^{(l)}$ can prevent the Dirichlet energy being below the lower energy limit that causes over-smoothing. Based on these studies, we formulate the *hyperbolic residual layer* by

$$\hat{\mathbf{h}}^{(l)} = f_{\mathrm{HR}}^\kappa \left( \mathbf{h}^{(t)}, f_{\mathrm{NA}}^\kappa(\tilde{\mathbf{P}}, \mathbf{h}^{(l)}); \alpha_l, 1 - \alpha_l \right), \quad (18)$$

where $t$ in this case can be 0 (initial layer) or $l-1$ (previous layer). Empirically, we find that adding previous residual does not contribute to the overall performance than adding only initial residual. Therefore, we set $t = 0$ in all case of our study. $f_{\mathrm{HR}}$ and $f_{\mathrm{NA}}$ are respectively defined in Eq. (17) and (16). The hyperparameter $\alpha_l$ for the $l$-th layer indicates the proportion of residual representation that the current layer retains. In practice, this value could be relatively small so it does not conceal the variation of embedding whilst alleviating over-smoothing.

*2) Hyperbolic Weight Alignment:* It is shown in [2], [26], however, that simply applying residual operation to initial/previous layer only partially relieves the over-smoothing issue and still degrades the performance when more layers are stacked. To fix such deficiency, [26] borrows the idea from ResNet [41] to align the weight matrix $\mathbf{W}^{(l)}$ in each layer to an identity matrix $\mathbf{I}$, which can be formulated as

$$\mathbf{h}^{(l+1)} = \hat{\mathbf{h}}^{(l)} \left( \beta_l \mathbf{W}^{(l)} + (1 - \beta_l)\mathbf{I}_{|\mathcal{V}|} \right) \tag{19}$$

$$= \beta_l(\hat{\mathbf{h}}^{(l)}\mathbf{W}^{(l)}) + (1 - \beta_l)\hat{\mathbf{h}}^{(l)}, \tag{20}$$

where $\hat{\mathbf{h}}^{(l)} \in \mathbb{D}_K^{|\mathcal{V}| \times d}$ is the aggregated feature. If $\beta_l$ is sufficiently small, the model ignores the weight matrix and simulates the behaviour of APPNP [42]. Further, it forces a small $\|\mathbf{W}\|$, which implies a small $s^L$ ($s$ is the maximum singular value of $\mathbf{W}^{(L)}$). According to [43], the loss of information in the $L$-layer GCN is relieved.

To generalize Eq. (19-20) to hyperbolic setting, one can leverage the power of Möbius gyromidpoint as Def. 1. Following the expanded Eq. (20), whose form is the weighted average between $\hat{\mathbf{h}}^{(l)}$ and $\hat{\mathbf{h}}^{(l)}\mathbf{W}^{(l)}$, thus midpoint can be employed to calculate the weighted mean between hyperbolic representation $\hat{\mathbf{h}}^{(l)}$ and the hyperbolic transformed representation $\mathbf{W} \otimes_K \hat{\mathbf{h}}^{(l)}$. Define the reformulation

$$\mathbf{h}^{\star(l+1)} = f_{\mathrm{HR}}^{\kappa}\left(\mathbf{W}^{(l)} \otimes_K \hat{\mathbf{h}}^{(l)}, \hat{\mathbf{h}}^{(l)}; \beta_l, 1 - \beta_l\right), \tag{21}$$

where $f_{\mathrm{HR}}^{\kappa}$ is defined in Eq. (17). Finally, we augment vanilla Möbius transformation $\mathbf{W} \otimes_{\kappa} \hat{\mathbf{h}}$ by our proposed $\mathcal{F}^K$ for the Poincaré ball-based model, which yields the formulation of proposed *hyperbolic weight alignment*

$$\mathbf{h}^{(l+1)} = f_{\mathrm{HR}}^{\kappa}\left(\mathcal{F}^{\kappa}(\hat{\mathbf{h}}^{(l)}; \mathbf{W}^{(l)}), \hat{\mathbf{h}}^{(l)}; \beta_l, 1 - \beta_l\right). \tag{22}$$

*3) Hyperbolic Feature Regularization:* Notably, [44] empirically finds that the learnt embeddings near the boundary of the hyperbolic ball are easier to classify. To push the nodes to the boundary, [45] proposed a regularization that first identifies a root node in hyperbolic space and then encourages the embeddings to move away from the root. We employ a similar approach to regularize the training of DeepHGCN. Firstly, the root node can be identified via gyromidpoint in Def. 1 as

$$\mathbf{h}_{\mathrm{root}}^{(l)} = \frac{1}{2} \otimes_{\kappa} \left( \frac{\sum_{i=1}^{|V|} \lambda_{\mathbf{h}_i^{(l)}}^{\kappa} \mathbf{h}_i^{(l)}}{\sum_{i=1}^{|V|} (\lambda_{\mathbf{h}_i^{(l)}}^{\kappa} - 1)} \right), \tag{23}$$

where $\mathbf{h}_i$ is the row-vector of embedding matrix $\mathbf{H}^{(l)}$. Each node is aligned with the root by $\bar{\mathbf{h}}_i^{(l)} = \mathbf{h}_i^{(l)} \ominus_{\kappa} \mathbf{h}_{\mathrm{root}}^{(l)}$. By the definition of the Poincaré ball we have $\|\bar{\mathbf{h}}_i^{(l)}\| < \frac{1}{|\kappa|}$. Pushing the aligned embeddings to the boundary deduces the norm of aligned embeddings $\|\bar{\mathbf{h}}_i^{(l)}\| \to \frac{1}{|\kappa|}$, which is equivalent to minimizing the inverse of the norm. The regularization term can therefore be defined as the inverse quadratic mean of the norms, formulated as

$$\mathcal{L}_{\mathrm{reg}}^{(l)} = \frac{1}{\sqrt{\frac{1}{|V|} \sum_{i=1}^{|V|} \|\bar{\mathbf{h}}_i^{(l)}\|^2}}. \tag{24}$$

In practice, computing $\mathcal{L}_{\mathrm{reg}}^{(l)}$ for a $L$-layer deep network is costly. Since the embedding is densely connected through

---

**Algorithm 1** DeepHGCN forward propagation

**Input**: graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$; node embeddings $\{\mathbf{x}_i\}_{i=1}^{|\mathcal{V}|}$; number of layers $L$; feature dim $d_f$, hidden dim $d_h$ and class dim $d_c$; activation $\sigma(\cdot)$; hyper-params $\{\alpha_l, \beta_l\}_{l=1}^L$ and $\gamma$
**Parameter**: $\{\mathbf{W}^{(l)}\}_{l=0}^L$ and $\{\mathbf{b}^{(l)}\}_{l=1}^L$
**Output**: Loss for back-propagation

1: generate hyperbolic representation $\{\mathbf{h}_i^{(0)}\}_{i=1}^{|\mathcal{V}|}$ via exponential map Eq. (28) transform $\mathbf{h}^{(0)} \to \mathbf{h}^{(1)}$ from $d_f$ to $d_h$ via Eq. (1)
2: **for** $l = 0$ to $L$ **do**
3:   neighborhood aggregation on $\mathbf{h}^{(l)}$ via Eq. (16)
4:   residual connect $\mathbf{h}^{(1)}$ and $\mathbf{h}^{(l)}$ via Eq (18)
5:   **if** $l == L$ **then**
6:     break the loop
7:   **end if**
8:   transform $\mathbf{h}^{(l)} \to \mathbf{h}^{(l+1)}$ within $\mathbb{D}^{d_h}$ via Eq. (1)
9:   weight alignment on $\mathbf{h}^{(l)}$ and $\mathbf{h}^{(l+1)}$ via Eq. (22)
10: **end for**
11: compute task-oriented loss $\mathcal{L}$ and feature regularization $\mathcal{L}_{\mathrm{reg}}$ via Eq. (24) with final embedding $\mathbf{H}^{(L)}$
12: **return** $\mathcal{L} + \gamma\mathcal{L}_{\mathrm{reg}}$

---

each layer in DeepHGCN, we can consider only the last layer regularization $\mathcal{L}_{\mathrm{reg}}^{(L)}$. Hence, the optimization target is formulated as $\mathcal{L} + \gamma\mathcal{L}_{\mathrm{reg}}^{(L)}$ where $\gamma$ is a hyperparameter.

*4) Rethinking DeepHGCN through Dirichlet Energy:* As per Def. 2, Dirichlet energy is essentially the weighted average of distances between normalized node pairs. To prevent over-smoothing, we want the energy of the last layer to be sufficiently large. On one hand, the initial residual and weight alignment ensure that the final representation contains at least a portion of the initial and previous layers. Since the energy of the starting layers is high, the residual connection mitigates the degradation of energy and retains the energy of the final layer at the same level as the previous layers. On the other hand, hyperbolic feature regularization encourages the node representation to move away from the center. This increases the distance between nodes and thus also alleviates the energy degradation.

### C. Training

The DeepHGCN forward pass is described in Alg. 1. For link prediction (LP) and node classification (NC) task, we employ the same downstream decoders and loss in [20], *i.e.* Fermi-Dirac decoder for LP and cross entropy for NC. Since all DeepHGCN parameters are resided in Euclidean space, a standard optimizer (*e.g.* Adam [46]) instead of Riemannian optimizers [47] can be leveraged.

## V. EMPIRICAL RESULTS

### A. Experiment Setup

*1) Datasets:* Four homophilic benchmark datasets are considered for node classification and link prediction: Airport [20], PubMed [49], CiteSeer [50] and Cora [51], with statistics in Tab. V. Airport dataset is a transductive dataset in which

TABLE IV
ROC AUC RESULTS (%) FOR LINK PREDICTION (LP) AND ACCURACY (%) FOR NODE CLASSIFICATION (NC). $\delta$ REFERS TO GROMOVS $\delta$-HYPERBOLICITY. A GRAPH IS MORE HYPERBOLIC AS $\delta \to 0$ AND IS A TREE WHEN $\delta = 0$. RESULTS PARTIALLY FROM [48].

| Dataset ($\delta$) | Airport ($\delta=1$) | | PubMed ($\delta=3.5$) | | CiteSeer ($\delta=5$) | | Cora ($\delta=11$) | |
|---|---|---|---|---|---|---|---|---|
| **Task** | LP | NC | LP | NC | LP | NC | LP | NC |
| MLP | 89.81±0.56 | 68.90±0.46 | 83.33±0.56 | 72.40±0.21 | 93.73±0.63 | 59.53±0.90 | 83.33±0.56 | 51.59±1.28 |
| HNN | 90.81±0.23 | 80.59±0.46 | 94.69±0.06 | 69.88±0.43 | 93.30±0.52 | 59.50±1.28 | 90.92±0.40 | 54.76±0.61 |
| GCN | 89.31±0.43 | 81.59±0.61 | 89.56±3.66 | 78.10±0.43 | 82.56±1.92 | 70.35±0.41 | 90.47±0.24 | 81.50±0.53 |
| GAT | 90.85±0.23 | 82.75±0.36 | 91.46±1.82 | 78.21±0.44 | 86.48±1.50 | 71.58±0.80 | 93.17±0.20 | 83.03±0.50 |
| GraphSAGE | 90.41±0.53 | 82.19±0.45 | 86.21±0.82 | 77.45±2.38 | 92.05±0.39 | 67.51±0.76 | 85.51±0.50 | 77.90±2.50 |
| SGC | 89.83±0.32 | 80.59±0.16 | 94.10±0.12 | 78.84±0.18 | 91.35±1.68 | 71.44±0.75 | 91.50±0.21 | 81.32±0.50 |
| HGNN | 96.42±0.44 | 84.71±0.98 | 92.75±0.26 | 77.13±0.82 | 93.58±0.33 | 69.99±1.00 | 91.67±0.41 | 78.26±1.19 |
| HGCN | 96.43±0.12 | 89.26±1.27 | 95.13±0.14 | 76.53±0.64 | 96.63±0.09 | 68.04±0.59 | 93.81±0.14 | 78.03±0.98 |
| HGAT | 97.86±0.09 | 89.62±1.03 | 94.18±0.18 | 77.42±0.66 | 95.84±0.37 | 68.64±0.30 | **94.02**±0.18 | 78.32±1.39 |
| HyboNet | 97.30±0.30 | 90.90±1.40 | 95.80±0.20 | 78.00±1.00 | 96.70±0.80 | 69.80±0.60 | 93.60±0.30 | 80.20±1.30 |
| DeepHGCN | **98.13**±0.33 | **94.70**±0.90 | **96.15**±0.17 | **79.43**±0.92 | **97.45**±0.44 | **73.31**±0.70 | 93.90±0.78 | **83.64**±0.40 |

nodes represent airports and edges indicate airline routes as derived from OpenFlights. The labels of nodes are determined by the population of the country to which the airport belongs. Cora, PubMed, and CiteSeer are benchmarks for citation networks, where nodes represent papers connected by citations. We report the hyperbolicity of each datatset (lower is more hyperbolic) as defined in [14].

TABLE V
STATISTICS OF THE GENERAL HGNN BENCHMARK DATASETS.

| Dataset | # Nodes | # Edges | Classes | Features | $\delta$ |
|---|---|---|---|---|---|
| Disease | 1,044 | 1,043 | 2 | 1,000 | 0 |
| Airport | 3,188 | 18,631 | 4 | 4 | 1 |
| PubMed | 19,717 | 44,338 | 3 | 500 | 3.5 |
| CiteSeer | 3,327 | 4,732 | 6 | 3,703 | 5 |
| Cora | 2,708 | 5,429 | 7 | 1,433 | 11 |

For heterophilic datasets, we evaluate node classification on three benchmarks, respectively, CORNELL, TEXAS and WISCONSIN [52] from the WebKB dataset (webpage networks). Detailed statistics are summarized in Tab. VI. We use the original fixed 10 split datasets. In Tab. VI, we report the homophily level $\mathcal{H}$ of each dataset, a sufficiently low $\mathcal{H} \leq 0.3$ means that the dataset is more heterophilic when most of neighbours are not in the same class.

TABLE VI
STATISTICS OF HETEROPHILIC BENCHMARK DATASETS.

| Dataset | # Nodes | # Edges | Classes | Features | $\mathcal{H}$ |
|---|---|---|---|---|---|
| Texas | 183 | 295 | 5 | 1,703 | 0.11 |
| Cornell | 182 | 295 | 5 | 1,703 | 0.21 |
| Wisconsin | 251 | 499 | 5 | 1,703 | 0.30 |

*2) Baselines:* Both (1) *Euclidean-hyperbolic comparison* and (2) *deep model comparison* are conducted in the experiment. For (1), we compare our model to 2 feature-based models, 4 Euclidean and 4 hyperbolic graph-based models. Feature-based models: without utilizing graph structure, we feed node feature into MLP and its hyperbolic variant HNN [17] to predict node labels or links. Graph-based models: we

consider GCN [53], GAT [54], GraphSAGE [3], and SGC [55] as Euclidean GNN methods. We consider HGNN [19], HGCN [20], HGAT [21] and HyboNet [24] as hyperbolic variants GNNs. For (2), we compare our model to the state-of-the-art deep GCN model GCNII [26], and also show the performance of vanilla GCN and HGCN under different layer settings.

*3) Parameter Settings:* Under homophilic setting, for node classification, the train/val/test percentage for Airport is 70/15/15% and standard splits [53], [56] on citation networks. For link prediction, we employ 85/5/10% edge splits on all datasets. For heterophilic datasets, the splits are in 60/20/20% and use all 10 random splits to compute the averaged results. All models are 16-dimension to ensure a fair comparison. We use Adam for training DeepHGCN, for other methods, we use the Adam and Riemannian Adam optimizer, respectively, for Euclidean and hyperbolic models. We set $\alpha$ to 0.1 and $\beta_l$ for the $l$-th layer to $\log(\frac{\lambda}{l} + 1)$ following [26]. The other hyperparameters are obtained via grid search, where $\lambda$: [0.4, 0.5, 0.6], $\gamma$: [1e-3, 1e-4, 1e-5], weight decay: [1e-3, 5e-4, 1e-4], and dropout: [0.0-0.6]. We conduct experiments averaging 10 sets of embeddings' quality using different random seeds. More details are listed in the Appendix.

### B. Experiment Results

*1) Comparison with Hyperbolic Models:* In Tab. IV, we report the averaged ROC AUC for link prediction and F1 score for node classification on various datasets. The dimensions for all models are (16) except for HyboNet (64) since it cannot be trained stably with low embedding dimension. We present our model with the hyperparameter settings that produce the best outcomes. Compared with baselines, our proposed DeepHGCN achieves the best NC and LP performance among all datasets with high hyperbolicity $\delta$. On the Cora dataset with low $\delta$, the DeepHGCN does not outperform the attention-based model on LP and only outperforms the Euclidean GCN on NC by a small margin, trading off the expense of training time. This suggests attention and Euclidean geometry are more suitable for non-hierarchical graph structures.

*2) Comparison with Deep Models:* In Tab. VII, we evaluate the deep models with different numbers of layers. Instead of

TABLE VII
TEST ACCURACY OF DIFFERENT 16-DIM MODELS WITH VARIOUS DEPTH.
NUMBERS IN *bold* DENOTE THE BEST OF EACH MODEL, AND IN *red bold*
HIGHLIGHT THE BEST MODELS FOR EACH DATASET.

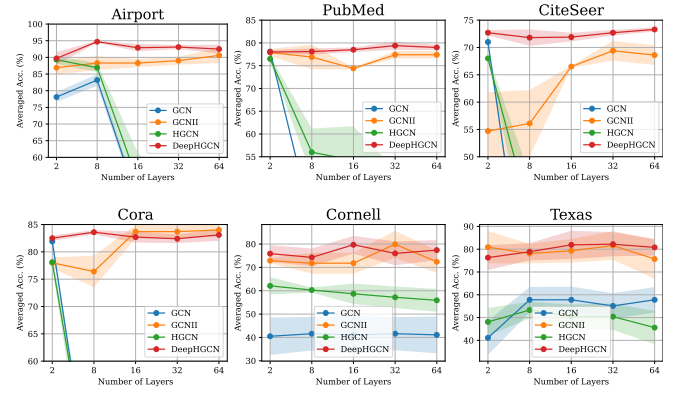| | Model | Layers | | | | |
|---|---|---|---|---|---|---|
| | | 2 | 8 | 16 | 32 | 64 |
| **Airport** | GCN | $78.1_{\pm1.5}$ | $\mathbf{83.2}_{\pm4.7}$ | $47.0_{\pm1.4}$ | $47.0_{\pm1.4}$ | $47.0_{\pm1.4}$ |
| | GCNII | $86.9_{\pm2.1}$ | $88.3_{\pm2.1}$ | $\mathbf{88.3}_{\pm1.2}$ | $86.6_{\pm1.2}$ | $86.6_{\pm2.3}$ |
| | HGCN | $\mathbf{89.3}_{\pm1.2}$ | $86.9_{\pm4.3}$ | $52.8_{\pm8.6}$ | $47.2_{\pm1.2}$ | $44.5_{\pm1.6}$ |
| | DeepHGCN | $89.7_{\pm1.9}$ | $\mathbf{94.7}_{\pm0.9}$ | $93.9_{\pm1.2}$ | $93.1_{\pm0.9}$ | $92.9_{\pm1.3}$ |
| **PubMed** | GCN | $\mathbf{78.1}_{\pm0.5}$ | $40.7_{\pm0.0}$ | $40.7_{\pm0.0}$ | $40.7_{\pm0.0}$ | $40.7_{\pm0.0}$ |
| | GCNII | $77.9_{\pm0.8}$ | $76.9_{\pm2.5}$ | $78.5_{\pm1.0}$ | $\mathbf{79.4}_{\pm0.5}$ | $79.1_{\pm0.8}$ |
| | HGCN | $\mathbf{76.5}_{\pm0.6}$ | $56.0_{\pm5.2}$ | $54.2_{\pm7.5}$ | $45.3_{\pm6.9}$ | $43.7_{\pm1.5}$ |
| | DeepHGCN | $78.0_{\pm0.4}$ | $78.1_{\pm0.7}$ | $78.5_{\pm0.5}$ | $\mathbf{79.4}_{\pm0.9}$ | $79.0_{\pm0.6}$ |
| **CiteSeer** | GCN | $\mathbf{71.0}_{\pm1.0}$ | $18.1_{\pm0.0}$ | $18.1_{\pm0.0}$ | $18.1_{\pm0.0}$ | $18.1_{\pm0.0}$ |
| | GCNII | $54.7_{\pm7.4}$ | $56.1_{\pm6.1}$ | $66.5_{\pm6.2}$ | $\mathbf{69.4}_{\pm1.8}$ | $68.6_{\pm1.8}$ |
| | HGCN | $\mathbf{68.0}_{\pm0.6}$ | $39.5_{\pm7.1}$ | $30.3_{\pm3.7}$ | $31.2_{\pm4.9}$ | $24.7_{\pm1.3}$ |
| | DeepHGCN | $72.7_{\pm0.5}$ | $71.8_{\pm1.5}$ | $71.9_{\pm0.8}$ | $72.7_{\pm0.5}$ | $\mathbf{73.3}_{\pm0.4}$ |
| **Cora** | GCN | $\mathbf{81.9}_{\pm1.1}$ | $31.9_{\pm0.0}$ | $31.9_{\pm0.0}$ | $31.9_{\pm0.0}$ | $31.9_{\pm0.0}$ |
| | GCNII | $78.0_{\pm3.0}$ | $76.4_{\pm2.2}$ | $77.2_{\pm3.5}$ | $83.7_{\pm0.9}$ | $\mathbf{84.0}_{\pm1.3}$ |
| | HGCN | $\mathbf{78.1}_{\pm0.9}$ | $33.5_{\pm5.5}$ | $30.8_{\pm3.9}$ | $24.1_{\pm7.3}$ | $21.7_{\pm7.8}$ |
| | DeepHGCN | $82.5_{\pm0.5}$ | $\mathbf{83.6}_{\pm0.4}$ | $82.7_{\pm1.0}$ | $82.4_{\pm0.6}$ | $83.1_{\pm0.7}$ |
| **Cornell** | GCN | $40.5_{\pm8.0}$ | $\mathbf{41.6}_{\pm7.2}$ | $41.1_{\pm8.1}$ | $41.6_{\pm7.1}$ | $41.1_{\pm7.9}$ |
| | GCNII | $\mathbf{72.9}_{\pm4.1}$ | $71.8_{\pm4.5}$ | $71.8_{\pm5.0}$ | $69.9_{\pm5.8}$ | $72.4_{\pm4.3}$ |
| | HGCN | $\mathbf{62.1}_{\pm3.7}$ | $60.3_{\pm4.1}$ | $58.7_{\pm4.4}$ | $57.2_{\pm4.6}$ | $55.9_{\pm4.9}$ |
| | DeepHGCN | $75.9_{\pm3.6}$ | $74.3_{\pm3.8}$ | $\mathbf{79.7}_{\pm3.8}$ | $76.0_{\pm5.1}$ | $77.4_{\pm4.3}$ |
| **Texas** | GCN | $41.1_{\pm7.5}$ | $57.8_{\pm5.6}$ | $\mathbf{57.8}_{\pm5.8}$ | $55.1_{\pm5.6}$ | $57.8_{\pm5.6}$ |
| | GCNII | $80.9_{\pm7.1}$ | $78.1_{\pm4.3}$ | $79.3_{\pm5.2}$ | $\mathbf{81.6}_{\pm6.3}$ | $75.7_{\pm8.7}$ |
| | HGCN | $48.1_{\pm6.1}$ | $\mathbf{53.3}_{\pm3.7}$ | $50.5_{\pm5.2}$ | $50.4_{\pm5.8}$ | $45.6_{\pm7.3}$ |
| | DeepHGCN | $76.3_{\pm5.4}$ | $78.9_{\pm3.8}$ | $81.9_{\pm6.2}$ | $\mathbf{82.2}_{\pm5.3}$ | $80.8_{\pm3.6}$ |
| **Wisconsin** | GCN | $43.1_{\pm3.7}$ | $47.8_{\pm8.2}$ | $\mathbf{49.0}_{\pm8.1}$ | $47.9_{\pm9.4}$ | $47.8_{\pm8.8}$ |
| | GCNII | $83.7_{\pm7.2}$ | $\mathbf{84.1}_{\pm3.9}$ | $83.1_{\pm1.9}$ | $83.5_{\pm3.5}$ | $80.4_{\pm4.3}$ |
| | HGCN | $\mathbf{57.9}_{\pm5.8}$ | $55.7_{\pm6.3}$ | $54.2_{\pm6.7}$ | $51.1_{\pm7.4}$ | $48.5_{\pm8.0}$ |
| | DeepHGCN | $79.2_{\pm4.2}$ | $81.3_{\pm5.5}$ | $\mathbf{84.0}_{\pm4.4}$ | $83.7_{\pm5.7}$ | $82.1_{\pm5.3}$ |



Fig. 5. Averaged performance of different models with various numbers of layers. We include hyperbolic, homophilic, and heterophilic datasets. The deeper models that overcome over-smoothing generally perform better.

using 64-dimensional hidden layers like [26], we apply 16-dimensions for all models. We observe that the performance of GCN and HGCN rapidly degrades when the number of layers surpasses 8, indicating that they are susceptible to over-smoothing. On the other hand, the GCNII does not outperform the shallow models with 2 layers, while consistently improving with more layers and achieving the best result on Cora and PubMed. The DeepHGCN can not only perform well with 2 layers, but outperforms all models on Airport and CiteSeer with $8 \sim 64$ layers with competitive results. Under the heterophilic datasets Cornell, Texas, and Wisconsin, the DeepHGCN also demonstrates strong performance with deeper architectures, which achieves the highest accuracy of 79.7% on Cornell with 16 layers, 82.2% on Texas with 32 layers, and 84.0% on Wisconsin with 16 layers. These results suggest our model can alleviate the over-smoothing issue of vanilla HGCN and capable of retrieving information from higher-order neighbors.

*Remark on benefit and limitation of depth.* In Fig. 5 illustrates the performance comparison through layers. It is worth noting that models with 2 layers are not sufficient for achieving the best performance. Our DeepHGCN achieves strong performance among groups with generally $8 \sim 32$ layers, similar to GCNII according to our validation. As 8-layer is already considered deep for graph models, we conclude that depth is beneficial for HGCN performance. By carefully tuning the depth, DeepHGCN could achieve

significantly better performance than its 2-layer baselines. However, on some datasets like Texas, the performance of all models declines with added depth beyond a certain point (*e.g* $32 \to 64$ layers). This suggests that while our approach helps, extremely deep models may have diminishing returns or even detriments on some datasets. This suggests that DeepHGCN is a remedy rather than a solution to hyperbolic model over-smoothing. Further research is needed to fully address the limitations of over-smoothness and realize the full potential in hyperbolic graph learning.
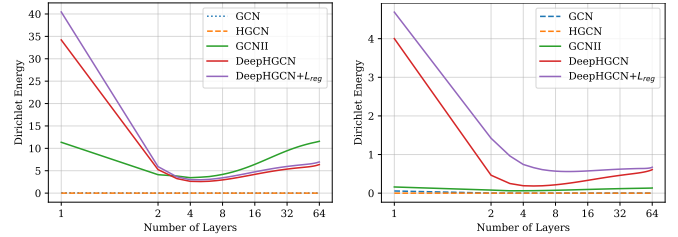


Fig. 6. Dirichlet energy variation through layers of different models on Cora (**Left**) and CiteSeer (**Right**). The x-axis is plotted in log-scale with base 2.

*3) Over-Smoothing Analysis:* In Fig. 6, we show the Dirichlet energy at each layer of a 64-layer DeepHGCN, comparing with GCN, HGCN and GCNII. Due to the over-smoothing issue, the energy of node embeddings in GCN and HGCN converges rapidly to zero. By fine-tuning the initial residual and weight alignment coefficients, the proposed DeepHGCN is able to obtain energy levels comparable to GCNII on Cora and CiteSeer. We also notice the feature regularized model tends to have higher energy through all layers compared with the unregularized ones. This validates our arguments in Sec. IV-B4.

*4) Ablation Study:* We study the contributions of our introduced techniques. In Fig. 7, we show the performance of DeepHGCN with different depths compared to HGCNs with one of the proposed techniques applied, further, we provide ablation studies on various component permutation on Airport and Cora datasets in Fig. 8. We find that hyperbolic initial residual is the most effective module for alleviating over-
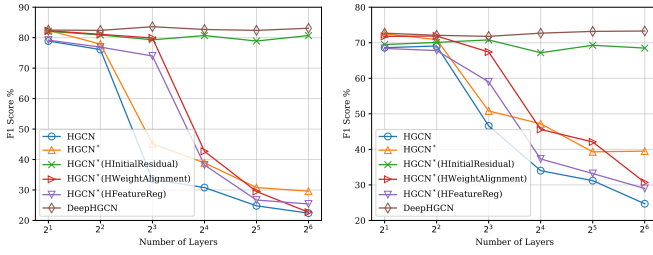
Fig. 7. Ablation study on hyperbolic initial residual, weight alignment and feature regularization on Cora (**Left**) and CiteSeer (**Right**) dataset. HGCN* denotes the HGCN with the proposed efficient backbone.
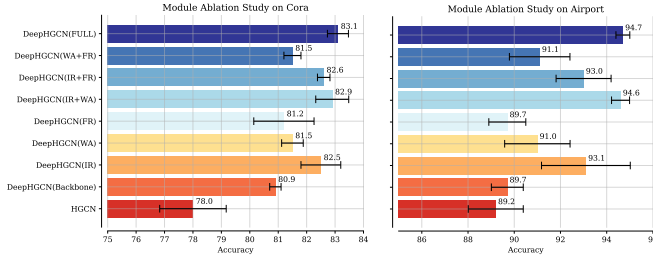


Fig. 8. Performance of permuted components on Cora (**Left**) and Airport (**Right**) in 10 evaluations, we show mean ± std for each model.



Fig. 9. Ablation study on the impact of different $\alpha_l$ in Eq. (18) on Airport (**Left**) and CiteSeer (**Right**). A higher $\alpha_l$ value indicates a stronger emphasis on initial embedding, vice versa.



Fig. 10. Visualize of node embeddings (after t-SNE) variation through layers on CiteSeer. Different colors represent different ground truth class labels.

smoothing, although performance of the 2-layer model still dominates. Directly applying weight alignment and feature regularization also relieves over-smoothing but the models still degrade after depth beyond $2^3$. This is because a small $\beta_l$ value in Eq. (22) forces the norm of the weight matrix $\|\mathbf{W}\|$ to be small. According to the theory in [43], the loss of information in an $L$-layer GCN is related to the maximum singular value s of $\mathbf{W}^{(L)}$. A smaller norm $\|\mathbf{W}\|$ implies a smaller $s^L$, thus relieving the loss of information and over-smoothing. However, when the number of layers grows very deep (*e.g.* beyond $2^3$), the multiplicative effect of $\beta$s over many layers ($\prod \beta_l$) may become too small. This could make the model overly reliant on the initial features and not have enough learning capacity, leading to degraded performance. As for feature regularization, applying the technique over a very deep model with many layers may force the node embeddings to be pushed too far apart. This could potentially distort the embed space and make it difficult to preserve the similarities between nodes that are needed for tasks like classification. As a result, performance degrades when the regularization is applied to an overly deep model. Therefore, the DeepHGCN simultaneously apply all three techniques, which assures that performance improves as network depth increases. This suggests that all of the techniques are indispensable for resolving the over-smoothing problem.

Additionally, as the hyperbolic initial residual in Eq. (18) contribute significantly to the overall performance, we study the impact of different choices of $\alpha_l$s. With various depth $(8, 16, 32)$ of DeepHGCN backbone, we select $\alpha_l$ from 0 to 1 with a step size 0.1. From Fig. 9 we can observe that with Airport, the performance is stably well across $\alpha_l = 0.1 \sim 0.9$. Under CiteSeer, the performance of DeepHGCN drastically
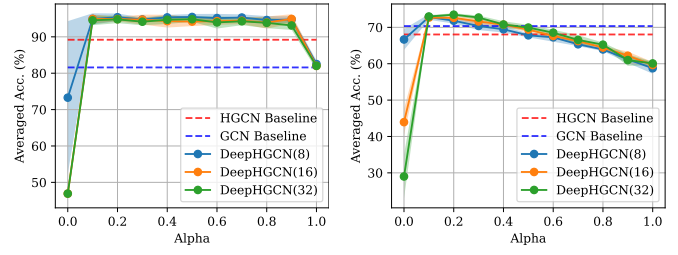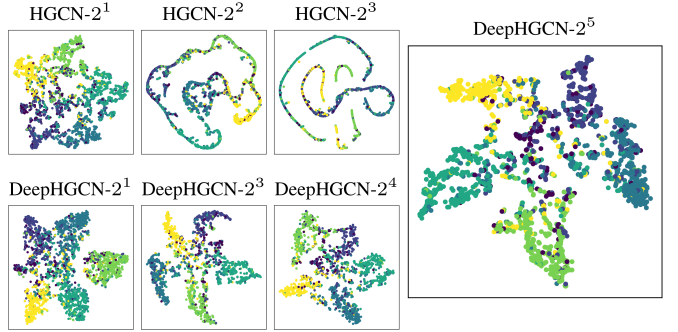
increases when $\alpha_l$ goes from 0 to 0.1, then gradually decreases as $\alpha_l$ increases. This suggests that DeepHGCN is adaptive to various $\alpha_l$s, while generally a small portion of initial residual is enough for achieving good performance.

*5) Embedding Visualization:* To better illustrate the effectiveness of DeepHGCN in maintaining distinguishable node representations in deeper layers, we visualize the 16-dimensional node embeddings of both HGCNs and Deep-HGCNs using t-SNE [57]. As shown in Fig. 10, as the depth increases, the embeddings of distinct classes in HGCNs tend to converge, resulting in indistinguishable representations. This phenomenon is a clear indication of the over-smoothing problem, which hinders the model's ability to capture meaningful node information in deeper layers. In contrast, DeepHGCN is able to obtain node embeddings that remain well-separated and clearly defined, even in higher layer settings. This evidenced the efficacy of our proposed techniques in alleviating the over-smoothing issue and preserving the discriminative power of node representations in deep hyperbolic networks.

## VI. CONCLUSION

In this paper, we introduce DeepHGCN, a novel deep hyperbolic graph neural network model that addresses the challenges of developing deeper HGCN architectures while mitigating the over-smoothing problem. Our model is powered by a computationally-efficient and expressive backbone, featuring a fast and accurate hyperbolic linear layer for feature transformation. Additionally, we propose a set of techniques, including hyperbolic initial residual connections, weight alignment, and feature regularization, which work together to effectively

prevent over-smoothing while preserving the manifold constraint. Extensive experiments on both hierarchical and non-hierarchical datasets demonstrate that DeepHGCN achieves competitive performance compared to existing Euclidean and shallow hyperbolic GCN models, highlighting the efficacy of our approach in capturing complex graph structures. Future research directions include extending our model to mixed-curvature manifolds and semi-Riemannian manifolds, as well as further investigating techniques to fully resolve the over-smoothing issue in extremely deep architectures.

## APPENDIX A
## HYPERBOLIC MODEL DETAILS

### A. Extended Review of Riemannian Geometry

A manifold $\mathcal{M}$ of $n$-dimension is a topological space that is locally-Euclidean, *i.e.* each point's neighborhood can be approximated by Euclidean space $\mathbb{R}^n$. The tangent space $\mathcal{T}_\mathbf{x}\mathcal{M}$ at $\mathbf{x} \in \mathcal{M}$ is the vector space of all tangent vectors at $\mathbf{x}$, the tangent space is isomorphic to $\mathbb{R}^n$. A Riemannian manifold $(\mathcal{M}, g)$ is a manifold $\mathcal{M}$ equipped with Riemannian metric $g = (g_\mathbf{x})_{\mathbf{x} \in \mathcal{M}}$, $g$ is a smooth collection of inner products on the tangent space of $\mathbf{x} \in \mathcal{M}$, *i.e.* $g_\mathbf{x} : \mathcal{T}_\mathbf{x}\mathcal{M} \times \mathcal{M}_\mathbf{x}\mathcal{M} \to \mathbb{R}$. It is natural to deduce Riemannian norm using metric $g$, *i.e.* for any vector $\mathbf{v} \in \mathcal{T}_\mathbf{x}\mathcal{M}$, $\|\mathbf{v}\|_{g_\mathbf{x}} = \sqrt{g_\mathbf{x}(\mathbf{v}, \mathbf{v})}$. The definition of inner-product and the induced norm can induce various geometric notions such as distances between points on $\mathcal{M}$, and angles between vectors on $\mathcal{T}_\mathbf{x}\mathcal{M}$.

*1) Length of Curve and Geodesic:* In the notion of differential geometry, a *curve* $\gamma$ is defined as a mapping from an interval to the manifold, *i.e.* $\gamma : [a, b] \to \mathcal{M}$. The length of curve is defined as $L(\gamma) = \int_{t_1}^{t_2} \sqrt{g_{\gamma(t)}(\gamma'(t), \gamma'(t))} dt$ where $t \in [a, b]$. Given the curve, let $\mathbf{x} = \gamma(a) \in \mathcal{M}$ and $\mathbf{y} = \gamma(b) \in \mathcal{M}$, the minimum distance between $\mathbf{x}$ and $\mathbf{y}$ on the manifold is called the *geodesic*

$$d(\mathbf{x}, \mathbf{y}) := \inf_\gamma L(\gamma) = \inf_\gamma \int_a^b \sqrt{g_{\gamma(t)}(\gamma'(t), \gamma'(t))} dt, \quad (25)$$

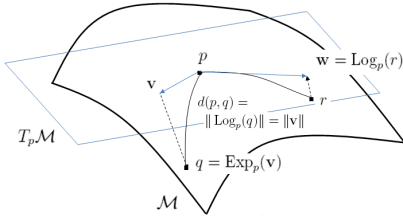it can be considered as a curve that minimizes the length.



Fig. 11. Exponential and logarithmic maps of non-Euclidean (Riemannian) manifolds. Illustration from [58].

*2) Exponential and Logarithmic Map:* The exponential map $\exp_\mathbf{x} : \mathcal{T}_\mathbf{x}\mathcal{M} \to \mathcal{M}$ at point $\mathbf{x}$ defines a way to project a vector $\mathbf{v}$ of tangent space $\mathcal{T}_\mathbf{x}\mathcal{M}$ at point $\mathbf{x}$, to a point $\mathbf{y} = \exp_\mathbf{x}(\mathbf{v}) \in \mathcal{M}$ on the manifold. The exponential map is generally used to parameterize a geodesic $\gamma$ uniquely defined by $\gamma(0) = \mathbf{x}$ and $\gamma'(0) = \mathbf{v}$. The logarithmic map $\log_\mathbf{x} : \mathcal{M} \to \mathcal{T}_\mathbf{x}\mathcal{M}$ is the inverse of exponential map, it defines a mapping of an arbitrary vector on $\mathcal{M}$ to the tangent space $\mathcal{T}_\mathbf{x}\mathcal{M}$. Different

manifolds have different ways to define exponential maps and logarithmic maps.

*3) Parallel Transport:* For two points $\mathbf{x}, \mathbf{y} \in \mathcal{M}$, the parallel transport $\mathcal{PT}_{\mathbf{x}\to\mathbf{y}} : \mathcal{T}_\mathbf{x}\mathcal{M} \to \mathcal{T}_\mathbf{y}\mathcal{M}$ defines a mapping from a vector $\mathbf{v}$ in $\mathcal{T}_\mathbf{x}\mathcal{M}$ to $\mathcal{T}_\mathbf{y}\mathcal{M}$ that moves $\mathbf{v}$ along the geodesic from $\mathbf{x}$ to $\mathbf{y}$. The parallel transport preserves the metric tensors.

### B. Poincaré Ball Model

The Poincaré ball is defined as the Riemannian manifold $\mathbb{D}_\kappa^n = (\mathcal{D}_\kappa^n, g^\mathbb{D})$, with point set $\mathcal{D}_\kappa^n = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < -\frac{1}{\kappa}\}$ and Riemannian metric

$$g_\mathbf{x}^\mathbb{D} = (\lambda_\mathbf{x}^\kappa)^2 g^\mathbb{E}, \quad (26)$$

where $\lambda_\mathbf{x}^\kappa = \frac{2}{1+\kappa\|\mathbf{x}\|^2}$ (the conformal factor) and $g^\mathbb{E} = \mathbf{I}_n$ (the Euclidean metric tensor). $\kappa < 0$ is a hyperparameter denoting the sectional curvature of the manifold.

*1) Gyrovector Addition:* Existing studies [17], [31], [59] adopt the gyrovector space framework [32], [16] as a non-associative algebraic formalism for hyperbolic geometry. The gyrovector operation $\oplus_\kappa$ is termed *Möbius addition*

$$\mathbf{x} \oplus_\kappa \mathbf{y} := \frac{(1 - 2\kappa\langle\mathbf{x}, \mathbf{y}\rangle - \kappa\|\mathbf{y}\|^2)\mathbf{x} + (1 + \kappa\|\mathbf{x}\|^2)\mathbf{y}}{1 - 2\kappa\langle\mathbf{x}, \mathbf{y}\rangle + \kappa^2\|\mathbf{x}\|^2\|\mathbf{y}\|^2}, \quad (27)$$

where $\mathbf{x}, \mathbf{y} \in \mathbb{D}_\kappa^n$. The Möbius addition defines the addition of two gyrovectors that preserves the summation on the manifold. The induced Möbius subtraction $\ominus_\kappa$ is defined as $\mathbf{x} \ominus_\kappa \mathbf{y} = \mathbf{x} \oplus_\kappa (-\mathbf{y})$.

*2) Tangent Space Operations:* As described in Sec. II-A, the *exponential map* projects a vector $\mathbf{v} \in \mathcal{T}_\mathbf{x}\mathbb{D}_\kappa$ in the tangent space at $\mathbf{x}$ to a point on $\mathbb{D}_\kappa$, while the *logarithmic map* projects the manifold vector back to the tangent space

$$\exp_\mathbf{x}^\kappa(\mathbf{v}) = \mathbf{x} \oplus_\kappa \left(\tanh(\sqrt{|\kappa|}\frac{\lambda_\mathbf{x}^\kappa\|\mathbf{v}\|}{2})\frac{\mathbf{v}}{\sqrt{|\kappa|}\|\mathbf{v}\|}\right), \quad (28)$$

$$\log_\mathbf{x}^\kappa(\mathbf{y}) = \frac{2\tanh^{-1}(\sqrt{|\kappa|}\|-\mathbf{x} \oplus_\kappa \mathbf{y}\|)(-\mathbf{x} \oplus_\kappa \mathbf{y})}{\sqrt{|\kappa|}\lambda_\mathbf{x}^\kappa\|-\mathbf{x} \oplus_\kappa \mathbf{y}\|}, \quad (29)$$

where $\lambda_\mathbf{x}^\kappa = \frac{2}{1+\kappa\|\mathbf{x}\|^2}$ is the *conformal factor* at point $\mathbf{x}$.

*3) Gyrovector Multiplication:* $\otimes_\kappa$ is the *Möbius scalar multiplication*. It defines the multiplication between scalar $r$ and a gyrovector. As provided in [17], [20], the scalar multiplication can be obtained by

$$r \otimes_\kappa \mathbf{x} = \exp_\mathbf{o}^\kappa(r \log_\mathbf{o}^\kappa(\mathbf{x})), \quad (30)$$

where $\mathbf{x} \in \mathbb{D}_\kappa^n/\mathbb{L}_\kappa^n$. One can further extend Eq. (30) to *matrix-vector multiplication*, formulated by

$$\mathbf{M} \otimes_\kappa \mathbf{x} = \exp_\mathbf{o}^\kappa(\mathbf{M} \log_\mathbf{o}^\kappa(\mathbf{x})), \quad (31)$$

where $\mathbf{M} \in \mathbb{R}^{m \times n}$. With broadcasting mechanism [60], we can derive the Möbius addition and matrix-vector multiplication manipulating on batched representations.

*4) Distance Metric:* Since the *geodesic* is the generalized straight line on Riemannian manifold, the distance between two points is essentially the *geodesic length*. For $\mathbf{x}, \mathbf{y} \in \mathbb{D}_\kappa^n$, the distance is given by:

$$d_\mathbb{D}^\kappa(\mathbf{x}, \mathbf{y}) = \frac{2}{\sqrt{|\kappa|}} \tanh^{-1}\left(\sqrt{|\kappa|}\|-\mathbf{x} \oplus_\kappa \mathbf{y}\|\right). \quad (32)$$

TABLE VIII
SUMMARY OF OPERATIONS IN THE POINCARÉ BALL MODEL AND THE LORENTZ (HYPERBOLOID) MODEL ($\kappa < 0$)

| | **Poincaré Ball Model** $\mathbb{D}$ | **Lorentz Model** $\mathbb{L}$ |
|---|---|---|
| **Point Set** | $\mathcal{D}_\kappa^n = \left\{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < -\frac{1}{\kappa} \right\}$ | $\mathcal{L}_\kappa^n = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} : \langle \mathbf{x}, \mathbf{x} \rangle_\mathcal{L} = \frac{1}{\kappa} \right\}$ |
| **Metric Tensor** | $g_\mathbf{x}^{\mathbb{D}_\kappa} = (\lambda_\mathbf{x}^\kappa)^2 g^\mathbb{E}$, where $\lambda_\mathbf{x}^\kappa = \frac{2}{1+\kappa\|\mathbf{x}\|^2}$ and $g^\mathbb{E} = \mathbf{I}$ | $g_\mathbf{x}^{\mathbb{L}_\kappa} = \eta$, where $\eta$ is $I$ except $\eta_{0,0} = -1$ |
| **Geodesic Length** | $d_\mathbb{D}^\kappa(\mathbf{x}, \mathbf{y}) = \frac{2}{\sqrt{|\kappa|}} \tanh^{-1}\left( \sqrt{|\kappa|}\| - \mathbf{x} \oplus_\kappa \mathbf{y}\| \right)$ | $d_\mathbb{L}^\kappa(\mathbf{x}, \mathbf{y}) = \frac{1}{\sqrt{|\kappa|}} \cosh^{-1}(\kappa\langle \mathbf{x}, \mathbf{y} \rangle_\mathcal{L})$ |
| **Exponential Map** | $\exp_\mathbf{x}^\kappa(\mathbf{v}) = \mathbf{x} \oplus_\kappa \left( \tanh(\sqrt{|\kappa|}\frac{\lambda_\mathbf{x}^\kappa\|\mathbf{v}\|}{2}\frac{\mathbf{v}}{\sqrt{|\kappa|}\|\mathbf{v}\|}) \right)$ | $\exp_\mathbf{x}^\kappa(\mathbf{v}) = \cosh\left( \sqrt{|\kappa|}\|\mathbf{v}\|_\mathcal{L} \right)\mathbf{x} + \mathbf{v}\frac{\sinh\left(\sqrt{|\kappa|}\|\mathbf{v}\|_\mathcal{L}\right)}{\sqrt{|\kappa|}\|\mathbf{v}\|_\mathcal{L}}$ |
| **Logarithmic Map** | $\log_\mathbf{x}^\kappa(\mathbf{y}) = \frac{2}{\sqrt{|\kappa|}\lambda_\mathbf{x}^\kappa} \tanh^{-1}(\sqrt{|\kappa|}\| - \mathbf{x} \oplus_\kappa \mathbf{y}\|)\frac{-\mathbf{x}\oplus_\kappa\mathbf{y}}{\|-\mathbf{x}\oplus_\kappa\mathbf{y}\|}$ | $\log_\mathbf{x}^\kappa(\mathbf{y}) = \frac{\cosh^{-1}(\kappa\langle\mathbf{x},\mathbf{y}\rangle_\mathcal{L})}{\sinh\left(\cosh^{-1}(\kappa\langle\mathbf{x},\mathbf{y}\rangle_\mathcal{L})\right)}(\mathbf{y} - \kappa\langle\mathbf{x},\mathbf{y}\rangle_\mathcal{L}\mathbf{x})$ |
| **Parallel Transport** | $\mathcal{PT}_{\mathbf{x}\to\mathbf{y}}^\kappa(\mathbf{v}) = \frac{\lambda_\mathbf{x}^\kappa}{\lambda_\mathbf{y}^\kappa}\, \mathrm{gyr}[\mathbf{y}, -\mathbf{x}]\mathbf{v}$ | $\mathcal{PT}_{\mathbf{x}\to\mathbf{y}}^\kappa(\mathbf{v}) = \mathbf{v} - \frac{\kappa\langle\mathbf{y},\mathbf{v}\rangle_\mathcal{L}}{1+\kappa\langle\mathbf{x},\mathbf{y}\rangle_\mathcal{L}}(\mathbf{x} + \mathbf{y})$ |

## C. Lorentz Model

The Lorentz model, *a.k.a.* the hyperboloid model is defined as the Riemannian manifold $\mathbb{L}_\kappa^n = (\mathcal{L}_\kappa^n, g^\mathbb{L})$, with point set $\mathcal{L}_\kappa^n = \left\{ \mathbf{x} \in \mathbb{R}^{n+1} : \langle \mathbf{x}, \mathbf{x} \rangle_\mathcal{L} = \frac{1}{\kappa} \right\}$ and Riemannian metric:

$$g^\mathbb{L} = \mathrm{diag}(-1, 1, \cdots, 1). \tag{33}$$

The point set $\mathcal{L}_\kappa^n$ is geometrically the upper sheet of hyperboloid in an $(n + 1)$-dimensional Minkowski space with the origin $(\sqrt{-\frac{1}{\kappa}}, 0, \cdots, 0)$. Each point in $\mathbb{L}_\kappa^n$ has the form $\mathbf{x} = \begin{bmatrix} x_t \\ \mathbf{x}_s \end{bmatrix}$, where $x_t \in \mathbb{R}$ is a scalar and $\mathbf{x}_s \in \mathbb{R}^n$. Given $\mathbf{x}, \mathbf{y} \in \mathbb{L}_\kappa^n$, the Lorentzian inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle_\mathcal{L} = -x_t y_t + \mathbf{x}_s^T \mathbf{y}_s \tag{34}$$
$$= \mathbf{x}^T \mathrm{diag}(-1, 1, \cdots, 1)\mathbf{y}. \tag{35}$$

## D. Summary of Operations

We summarize the hyperbolic operations of the Poincaré ball model and Lorentz model in Tab. VIII. We denote $\|\cdot\|$ and $\langle \mathbf{x}, \mathbf{y} \rangle_2$ as the Euclidean L2-norm and inner product, $\langle \mathbf{x}, \mathbf{y} \rangle_\mathcal{L}$ as the Lorentzian inner product $\mathbf{x}^T \mathrm{diag}(-1, 1, \cdots, 1)\mathbf{y}$ and $\|\cdot\|_\mathcal{L}$ as Lorentzian norm where $\|\mathbf{x}\|_\mathcal{L}^2 = \langle \mathbf{x}, \mathbf{x} \rangle_\mathcal{L}$. For systematic gyrovector space treatment, please refer to [16].

## APPENDIX B
### IMPLEMENTATION NOTES

### A. Fixing Numerical Instability via Clipping

In our implementation, 32-bit float tensors are used for all manipulations. In practice, computing the square of a float tensor $\mathbf{x}$ requires lower numeric limit, for instance, if $\mathbf{x} = 10^{-a}$, the precision required for $\mathbf{x}^2$ is at least $10^{-2a}$. The smallest value for a 32-bit float tensor is approximately $1.175494 \times 10^{-38}$, thus if $a > 19$, $\mathbf{x}^2$ will likely be out of memory and result in Nan value. In such cases, the tensors are out of hyperbolic space and could mislead the training. To avoid numerical instability, we employ feature clipping:

$$\mathrm{Clip}(\mathbf{x}; a, \kappa) = \begin{cases} \frac{1-\epsilon}{\sqrt{|\kappa|}\|\mathbf{x}\|}\mathbf{x}, & \|\mathbf{x}\| \geq \frac{1-\epsilon}{\sqrt{|\kappa|}} \\ \frac{a}{\|\mathbf{x}\|}\mathbf{x}, & \|\mathbf{x}\| < a \\ \mathbf{x}, & \text{otherwise} \end{cases} \tag{36}$$

where the $a$ is usually fixed to $10^{-15}$. The feature clipping is adopted in many parts of our implementation where we are likely to get Nan values.

### B. Dropout

Dropout is an essential technique for preventing over-fitting in hyperbolic models. In our implementation, we perform dropout on $\boldsymbol{\omega} = \phi(\cdot)$ in Eq. (1). In the following, we verify that the hyperbolic representation in the Poincaré ball model after dropout is manifold-preserving.

For arbitrary $\boldsymbol{\omega}$, following Eq. (1) we have

$$\|\mathcal{F}_\mathbb{D}^\kappa(\cdot)\| = \|\frac{\boldsymbol{\omega}}{1 + \sqrt{1 - \kappa\|\boldsymbol{\omega}\|^2}}\| \tag{37}$$

$$= \sqrt{\frac{\|\boldsymbol{\omega}\|^2}{(1 + \sqrt{1 - \kappa\|\boldsymbol{\omega}\|^2})^2}} \tag{38}$$

$$= \sqrt{\frac{\|\boldsymbol{\omega}\|^2}{1 + 2\sqrt{2 - \kappa\|\boldsymbol{\omega}\|^2} - \kappa\|V\|^2}} \tag{39}$$

$$= \sqrt{\frac{\|\boldsymbol{\omega}\|^2}{1 + 2\sqrt{2 + |\kappa|\|\boldsymbol{\omega}\|^2} + |\kappa|\|V\|^2}} \tag{40}$$

$$= \sqrt{\frac{1}{\frac{1}{\|\boldsymbol{\omega}\|^2} + 2\sqrt{\frac{2}{\|\boldsymbol{\omega}\|^4} + \frac{|\kappa|}{\|\boldsymbol{\omega}\|^2}} + |\kappa|}} \tag{41}$$

$$< \frac{1}{\sqrt{|\kappa|}}. \tag{42}$$

Notably, $\|\mathcal{F}_\mathbb{D}^\kappa(\cdot)\|$ reaches $\frac{1}{\sqrt{|\kappa|}}$ and $0$ respectively when $\boldsymbol{\omega}$ approximate infinity and when $\boldsymbol{\omega} = 0$. Thus the range of each component in $\boldsymbol{\omega}$ is $[-\infty, \infty]$, which is in coincidence with the Euclidean space representation. Hence, the dropout can be directly applied to $\boldsymbol{\omega}$ without further generalization.

### C. Non-linear Activation

The non-linear activation prevents multi-layer GCNs from collapsing into single-layer networks. Activation functions are typically applied after neighborhood aggregation and before linear transformation step for optimal performance [53]. In the Poincaré ball model, applying ReLU is manifold-preserving (*i.e.* $\forall \mathbf{x} \in \mathbb{D}_\kappa^n$ we have $\sigma(\mathbf{x}) \in \mathbb{D}_\kappa^n$) since ReLU only cut-off the negative half and remain the positive half unchanged.

## APPENDIX C
## PROOFS AND DERIVATIONS

### A. Proof of Proposition 4

*Proof:* We start the proof from the definition of hyperbolic Dirichlet energy on the Poincaré ball model. Given the closed form solution of the distance function $d_{\mathbb{D}}$, we have

$$f_{\text{DE}}^{\mathbb{D}}(\tilde{\mathbf{P}} \otimes_\kappa \mathbf{H}) = \tag{43}$$
$$\frac{1}{2}\sum_{i,j} a_{ij} d_{\mathbb{D}}^2\left(\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \tilde{\mathbf{P}} \otimes_\kappa \mathbf{h}_i, \frac{1}{\sqrt{1+d_j}} \otimes_\kappa \tilde{\mathbf{P}} \otimes_\kappa \mathbf{h}_j\right)$$
$$= \frac{1}{2}\sum_{i,j} a_{ij}\left(\frac{2}{\sqrt{|\kappa|}} \tanh^{-1}\left(\sqrt{|\kappa|}\|\rho(\mathbf{H})\|\right)\right)^2, \tag{44}$$

where function $\rho$ can be expanded as

$$\rho(\mathbf{H}) =$$
$$\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa (\tilde{\mathbf{P}} \otimes_\kappa \mathbf{h}_i)\right) \oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa (\tilde{\mathbf{P}} \otimes_\kappa \mathbf{h}_j)\right) \tag{45}$$
$$= \left(-\frac{\tilde{\mathbf{P}}}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right) \oplus_\kappa \left(\frac{\tilde{\mathbf{P}}}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right) \tag{46}$$
$$= \tilde{\mathbf{P}}\left(\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right) \oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right)\right). \tag{47}$$

One can easily prove that for all $\mathbf{A} \in \mathbb{R}^{n \times n}$ and $\mathbf{x} \in \mathbb{R}^n$, $\|\mathbf{A}\mathbf{x}\| \le \|\mathbf{A}\|_F\|\mathbf{x}\|$, thus the $\|\rho(\tilde{\mathbf{P}}, \mathbf{H}^{(l)})\|$ in Eq. (44) can be further written as

$$\|\rho(\mathbf{H})\|$$
$$= \left\|\tilde{\mathbf{P}}\left(\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right) \oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right)\right)\right\| \tag{48}$$
$$\le \|\tilde{\mathbf{P}}\|_F \left\|\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right) \oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right)\right\|. \tag{49}$$

Since $\tilde{\mathbf{P}}$ is normalized (can be row normalize or diagonal normalize), the Frobenius norm $\|\tilde{\mathbf{P}}\|_F \ge 1$ always establish, thus the norm of $\rho$:

$$\|\rho(\mathbf{H})\|$$
$$\le \left\|\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right) \oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right)\right\|. \tag{50}$$

Known that $\tanh^{-1}(\cdot)$ is a monotonically increasing function, Eq. (44) can be further derived as

$$\frac{1}{2}\sum_{i,j} a_{ij}\left(\frac{2}{\sqrt{|\kappa|}} \tanh^{-1}\left(\sqrt{|\kappa|}\|\rho(\tilde{\mathbf{P}}, \mathbf{H})\|\right)\right)^2 \tag{51}$$
$$\le \frac{1}{2}\sum_{i,j} a_{ij}\left(\frac{2}{\sqrt{|\kappa|}} \tanh^{-1}\left(\sqrt{|\kappa|}\left\|\left(-\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i\right)\right.\right.\right.$$
$$\left.\left.\left.\oplus_\kappa \left(\frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right)\right\|\right)\right)^2 \tag{52}$$
$$= \frac{1}{2}\sum_{i,j} a_{ij} d_{\mathbb{D}}^2\left(\frac{1}{\sqrt{1+d_i}} \otimes_\kappa \mathbf{h}_i, \frac{1}{\sqrt{1+d_j}} \otimes_\kappa \mathbf{h}_j\right) \tag{53}$$
$$= f_{\text{DE}}^{\mathbb{D}}(\mathbf{H}). \tag{54}$$

Eq. (43-54) concludes the proof. ∎

## APPENDIX D
## APPENDIX D: ADDITIONAL EVALUATIONS

### A. Additional Assessment on DISEASE Dataset

TABLE IX
MODEL EVALUATION IN DISEASE.

| Dataset ($\delta$) | Disease ($\delta = 0$) | |
|---|---|---|
| **Task** | LP | NC |
| GCN | 58.00±1.41 | 69.79±0.54 |
| GAT | 58.16±0.92 | 70.04±0.49 |
| GraphSAGE | 65.93±0.29 | 70.10±0.49 |
| SGC | 65.34±0.28 | 70.94±0.59 |
| GCNII (8) | / | 88.83±1.32 |
| GCNII (16) | / | **96.71**±2.78 |
| HGNN | 81.54±1.22 | 81.27±3.53 |
| HGCN | 90.80±0.31 | 88.16±0.76 |
| HGAT | 87.63±1.67 | 90.30±0.62 |
| HyboNet | **96.80**±0.40 | 96.00±1.00 |
| DeepHGCN (2) | 92.10±0.44 | 89.90±1.33 |
| DeepHGCN (8) | 95.70±0.32 | 92.51±2.10 |
| DeepHGCN (16) | 95.51±1.52 | 93.70±1.52 |

We provide the performance comparisons of models in Disease dataset in Tab. IX. We observed that, when increasing the depth of GCNII, the validation accuracy of node classification will reach almost 99%, suggesting Euclidean space is also capable for embedding DISEASE. Although the hyperbolic space is more natural for embedding tree-like data and therefore could yield improved performance on the DISEASE dataset, it looks the potential improvement we could expect over GCNII is marginal.

According to our experiments, GCNII, HyboNet and Deep-HGCN are all capable of fitting the data, that is, the training accuracy can reach 100% while other models are unable to fit the data even without dropout and weight regularization. We infer that the performance gap observed in the test set is attributed to the poor generalization ability of hyperbolic classifiers, thereby suggesting an intriguing direction for future research.

## REFERENCES

[1] M. Defferrard, X. Bresson, and P. Vandergheynst, "Convolutional neural networks on graphs with fast localized spectral filtering," *Advances in neural information processing systems*, vol. 29, 2016.

[2] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," 2017.

[3] W. Hamilton, Z. Ying, and J. Leskovec, "Inductive representation learning on large graphs," *Advances in neural information processing systems*, vol. 30, 2017.

[4] M. Li, A. Micheli, Y. G. Wang, S. Pan, P. Lió, G. S. Gnecco, and M. Sanguineti, "Guest editorial: Deep neural networks for graphs: Theory, models, algorithms, and applications," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 35, no. 4, pp. 4367–4372, 2024.

[5] A. Clauset, C. Moore, and M. E. Newman, "Hierarchical structure and the prediction of missing links in networks," *Nature*, vol. 453, no. 7191, pp. 98–101, 2008.

[6] M. Zitnik, R. Sosič, M. W. Feldman, and J. Leskovec, "Evolution of resilience in protein interactomes across the tree of life," *Proceedings of the National Academy of Sciences*, vol. 116, no. 10, pp. 4426–4433, 2019.

[7] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," in *Thirty-second AAAI conference on artificial intelligence*, 2018.

[8] D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams, "Convolutional networks on graphs for learning molecular fingerprints," *Advances in neural information processing systems*, vol. 28, 2015.

[9] M. Li, S. Zhou, Y. Chen, C. Huang, and Y. Jiang, "Educross: Dual adversarial bipartite hypergraph learning for cross-modal retrieval in multimodal educational slides," *Information Fusion*, p. 102428, 2024.

[10] N. Linial, E. London, and Y. Rabinovich, "The geometry of graphs and some of its algorithmic applications," *Combinatorica*, vol. 15, no. 2, pp. 215–245, 1995.

[11] D. Krioukov, F. Papadopoulos, M. Kitsak, A. Vahdat, and M. Boguná, "Hyperbolic geometry of complex networks," *Physical Review E*, vol. 82, no. 3, p. 036106, 2010.

[12] F. Papadopoulos, M. Kitsak, M. Serrano, M. Boguná, and D. Krioukov, "Popularity versus similarity in growing networks," *Nature*, vol. 489, no. 7417, pp. 537–540, 2012.

[13] L. Bai, L. Cui, Y. Wang, M. Li, J. Li, S. Y. Philip, and E. R. Hancock, "Haqjsk: Hierarchical-aligned quantum jensen-shannon kernels for graph classification," *IEEE Transactions on Knowledge and Data Engineering*, 2024.

[14] M. Gromov, "Hyperbolic groups," in *Essays in group theory*. Springer, 1987, pp. 75–263.

[15] M. Hamann, "On the tree-likeness of hyperbolic spaces," in *Mathematical proceedings of the cambridge philosophical society*, vol. 164, no. 2. Cambridge University Press, 2018, pp. 345–361.

[16] A. A. Ungar, "A gyrovector space approach to hyperbolic geometry," *Synthesis Lectures on Mathematics and Statistics*, vol. 1, no. 1, pp. 1–194, 2008.

[17] O. Ganea, G. Bécigneul, and T. Hofmann, "Hyperbolic neural networks," *Advances in neural information processing systems*, vol. 31, 2018.

[18] M. Nickel and D. Kiela, "Learning continuous hierarchies in the lorentz model of hyperbolic geometry," in *International Conference on Machine Learning*. PMLR, 2018, pp. 3779–3788.

[19] Q. Liu, M. Nickel, and D. Kiela, "Hyperbolic graph neural networks," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[20] I. Chami, Z. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," *Advances in neural information processing systems*, vol. 32, 2019.

[21] C. Gulcehre, M. Denil, M. Malinowski, A. Razavi, R. Pascanu, K. M. Hermann, P. Battaglia, V. Bapst, D. Raposo, A. Santoro *et al.*, "Hyperbolic attention networks," *arXiv preprint arXiv:1805.09786*, 2018.

[22] Y. Zhang, X. Wang, C. Shi, N. Liu, and G. Song, "Lorentzian graph convolutional networks," in *Proceedings of the Web Conference 2021*, 2021, pp. 1249–1261.

[23] J. Dai, Y. Wu, Z. Gao, and Y. Jia, "A hyperbolic-to-hyperbolic graph convolutional network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 154–163.

[24] W. Chen, X. Han, Y. Lin, H. Zhao, Z. Liu, P. Li, M. Sun, and J. Zhou, "Fully hyperbolic neural networks," *arXiv preprint arXiv:2105.14686*, 2021.

[25] Q. Li, Z. Han, and X.-M. Wu, "Deeper insights into graph convolutional networks for semi-supervised learning," in *Thirty-Second AAAI conference on artificial intelligence*, 2018.

[26] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," in *International conference on machine learning*. PMLR, 2020, pp. 1597–1607.

[27] K. Xu, C. Li, Y. Tian, T. Sonobe, K.-i. Kawarabayashi, and S. Jegelka, "Representation learning on graphs with jumping knowledge networks," in *International Conference on Machine Learning*. PMLR, 2018, pp. 5453–5462.

[28] Y. Rong, W. Huang, T. Xu, and J. Huang, "Dropedge: Towards deep graph convolutional networks on node classification," *arXiv preprint arXiv:1907.10903*, 2019.

[29] K. Zhou, X. Huang, D. Zha, R. Chen, L. Li, S.-H. Choi, and X. Hu, "Dirichlet energy constrained learning for deep graph neural networks," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[30] C. Huang, M. Li, F. Cao, H. Fujita, Z. Li, and X. Wu, "Are graph convolutional networks with random weights feasible?" *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 2751–2768, 2022.

[31] R. Shimizu, Y. Mukuta, and T. Harada, "Hyperbolic neural networks++," *arXiv preprint arXiv:2006.08210*, 2020.

[32] A. A. Ungar, *Analytic hyperbolic geometry: Mathematical foundations and applications*. World Scientific, 2005.

[33] R. Benedetti and C. Petronio, *Lectures on hyperbolic geometry*. Springer Science & Business Media, 1992.

[34] M. Nickel and D. Kiela, "Poincaré embeddings for learning hierarchical representations," *Advances in neural information processing systems*, vol. 30, 2017.

[35] L. Derczynski, "Power consumption variation over activation functions," *arXiv preprint arXiv:2006.07237*, 2020.

[36] N. Choudhary and C. K. Reddy, "Towards scalable hyperbolic neural networks using taylor series approximations," *arXiv preprint arXiv:2206.03610*, 2022.

[37] A. Lou, I. Katsman, Q. Jiang, S. J. Belongie, S.-N. Lim, and C. D. Sa, "Differentiating through the fréchet mean," *ArXiv*, vol. abs/2003.00335, 2020.

[38] H. Karcher, *Riemannian comparison constructions*. SFB 256, 1987.

[39] ——, "Riemannian center of mass and so called karcher mean," *arXiv preprint arXiv:1407.2087*, 2014.

[40] J. Klicpera, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank," *arXiv preprint arXiv:1810.05997*, 2018.

[41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.

[42] J. Gasteiger, A. Bojchevski, and S. Günnemann, "Predict then propagate: Graph neural networks meet personalized pagerank," *arXiv preprint arXiv:1810.05997*, 2018.

[43] K. Oono and T. Suzuki, "Graph neural networks exponentially lose expressive power for node classification," *arXiv preprint arXiv:1905.10947*, 2019.

[44] V. Khrulkov, L. Mirvakhabova, E. Ustinova, I. Oseledets, and V. Lempitsky, "Hyperbolic image embeddings," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6418–6428.

[45] M. Yang, M. Zhou, J. Liu, D. Lian, and I. King, "Hrcf: Enhancing collaborative filtering via hyperbolic geometric regularization," in *Proceedings of the ACM Web Conference 2022*, 2022, pp. 2462–2471.

[46] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[47] G. Bécigneul and O.-E. Ganea, "Riemannian adaptive optimization methods," *arXiv preprint arXiv:1810.00760*, 2018.

[48] S. Zhu, S. Pan, C. Zhou, J. Wu, Y. Cao, and B. Wang, "Graph geometry interaction learning," *Advances in Neural Information Processing Systems*, vol. 33, pp. 7548–7558, 2020.

[49] G. Namata, B. London, L. Getoor, B. Huang, and U. Edu, "Query-driven active surveying for collective classification," in *10th International Workshop on Mining and Learning with Graphs*, vol. 8, 2012, p. 1.

[50] C. L. Giles, K. D. Bollacker, and S. Lawrence, "Citeseer: An automatic citation indexing system," in *Proceedings of the third ACM conference on Digital libraries*, 1998, pp. 89–98.

[51] P. Sen, G. Namata, M. Bilgic, L. Getoor, B. Galligher, and T. Eliassi-Rad, "Collective classification in network data," *AI magazine*, vol. 29, no. 3, pp. 93–93, 2008.

[52] H. Pei, B. Wei, K. C.-C. Chang, Y. Lei, and B. Yang, "Geom-gcn: Geometric graph convolutional networks," *arXiv preprint arXiv:2002.05287*, 2020.

[53] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv preprint arXiv:1609.02907*, 2016.

[54] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[55] F. Wu, A. Souza, T. Zhang, C. Fifty, T. Yu, and K. Weinberger, "Simplifying graph convolutional networks," in *International conference on machine learning*.   PMLR, 2019, pp. 6861–6871.

[56] Z. Yang, W. Cohen, and R. Salakhudinov, "Revisiting semi-supervised learning with graph embeddings," in *International conference on machine learning*.   PMLR, 2016, pp. 40–48.

[57] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne." *Journal of machine learning research*, vol. 9, no. 11, 2008.

[58] J. Angulo, "Structure tensor image filtering using riemannian $l_1$ and $l_\infty$ center-of-mass," *Image Analysis & Stereology*, vol. 33, no. 2, pp. 95–105, 2014.

[59] E. Mathieu, C. Le Lan, C. J. Maddison, R. Tomioka, and Y. W. Teh, "Continuous hierarchical representations with poincaré variational auto-encoders," *Advances in neural information processing systems*, vol. 32, 2019.

[60] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.