





Continuous Geometry-Aware Graph Diffusion via Hyperbolic Neural PDE

Jiaxu Liu¹ , Xinpeng Yi² , Sihao Wu¹, Xiangyu Yin¹, Tianle Zhang¹, Xiaowei Huang¹, and Shi Jin²

¹ Department of Computer Science, University of Liverpool, Liverpool, UK
{jiaxu.liu, sihao.wu, xiangyu.yin, tianle.zhang, xiaowei.huang}@liverpool.ac.uk

² National Mobile Communications Research Laboratory, Southeast University, Nanjing, China
{xyi, jinshi}@seu.edu.cn

Abstract. While Hyperbolic Graph Neural Network (HGNN) has recently emerged as a powerful tool dealing with hierarchical graph data, the limitations of scalability and efficiency hinder itself from generalizing to deep models. In this paper, by envisioning depth as a continuous-time embedding evolution, we decouple the HGNN and reframe the information propagation as a partial differential equation, letting node-wise attention undertake the role of diffusivity within the Hyperbolic Neural PDE (HPDE). By introducing theoretical principles *e.g.*, field and flow, gradient, divergence, and diffusivity on a non-Euclidean manifold for HPDE integration, we discuss both implicit and explicit discretization schemes to formulate numerical HPDE solvers. Further, we propose the Hyperbolic Graph Diffusion Equation (HGDE) – a flexible vector flow function that can be integrated to obtain expressive hyperbolic node embeddings. By analyzing potential energy decay of embeddings, we demonstrate that HGDE is capable of modeling both low- and high-order proximity with the benefit of local-global diffusivity functions. Experiments on node classification and link prediction and image-text classification tasks verify the superiority of the proposed method, which consistently outperforms various competitive models by a significant margin.

Keywords: Continuous GNN · Hyperbolic Space · Neural ODE

1 Introduction

Graphs play a vital role in various disciplines, including social network analysis [12], bioinformatics [48], and computer vision [37]. The advent of Graph Neural Networks (GNNs, [23]) has significantly enhanced the analysis of these structures to capture complex relationships between nodes in a graph. However, traditional GNNs operate within the borders of Euclidean space, which may not be sufficiently expressive for data with inherent hierarchical or complex structures. To improve, this paper delves into the realm of hyperbolic geometry, a Riemannian

manifold demonstrated to be particularly effective for embedding hierarchical data [16, 18]. We focus on the development of HGNNs [6], which leverage the unique properties of hyperbolic space to enhance the embedding of GNNs.

The principal challenge confronted by HGNNs is their architectural design, which primarily consists of combinations of aggregation and transformation within layers. This fusion presents a unique problem, particularly the difficulty of training attention weights and manifold parameters (*e.g.*, curvature of the hyperbolic manifold) layer-wise in a **deeply** layered scheme. With such challenge, we pose our initial questions: **Q1**: *Considering hyperbolic space slows down layer-wise attention and propagation [11, 25], how to develop a deeply-layered attentive HGNN?* **Q2**: *How to incorporate high-order info to benefit a deeply layered scheme?* **Q3**: *Deep GNNs suffer from embedding smoothing, how should the node smoothness be measured when there is no defined metric for hyperbolic smoothness. And how to tackle over-smoothing within hyperbolic manifold constraints?*

Motivated by above questions, in this paper, we propose to decouple the functions within layers of HGNNs so as to deal with each of them separately. Unlike traditional decoupling-GNN approaches [15, 35] that aggregate all information from the neighbors, we view information propagation as a distillation process, such that unimportant information is filtered out and significant information is weighted and contributes to the continuous variation of embeddings. More explicitly, by letting the transformation layer manifest as an encoder-decoder scheme, the aggregation layer is re-envisioned to solve the partial differential equation (Neural ODE/PDE, [8]) - essentially, the graph diffusion equation [4] in hyperbolic space, which essentially simulates an infinitely deep HGNN with single layer parameters. In specific, in response to **Q1**, we consider the PDE reformulation and developed Hyperbolic-PDE (HPDE) solvers, which only leverage single-layer parameters. To answer **Q2**, we formulate the Hyperbolic Graph Diffusion Equation (HGDE), a low-high order vector flow function that can be integrated by HPDE. Tackling **Q3**, we firstly introduce the hyperbolic adaptation of Dirichlet energy and augmented HGDE with a hyperbolic residual, powered by Poincaré midpoint. Deconstructions above introduce extensive mathematical principles, including for instance: manifold vector field, flow, gradient, divergence, diffusivity, numerical HPDE solvers and hyperbolic residuals for bounding embedding energy decay. Through these concepts, we open new pathways to fully exploit the unique potential of hyperbolic space in the contextual analysis of graph-based data. In summary, the contributions of this paper are listed as follows.

(I) We present the geometric intuition for designing projective numerical integration methods that solve hyperbolic ODE/PDE, and examine the connection to Riemannian gradient descent methods. Focusing on fixed-grid solvers, we derive both hyperbolic generalizations of explicit schemes (Euler, Runge-Kutta) and implicit schemes (Adams-Moulton).

(II) We formulate the HGDE, which acts as the *vector flow* of the HPDE, and thereby induces concepts such as gradient, divergence and diffusivity within HGDE. The proposed framework is flexible and efficient for generating expressive (endowed by the depth) hyperbolic graph embeddings.

(III) We instantiate the diffusivity function as a mixed-order multi-head attention to account for both homophilic (local) and heterophilic (global) relations. Besides, we introduce hyperbolic residual technique to benefit the optimization and prevent over-smoothing.

Through extensive experiments and comparison with the state-of-the-art on multiple real-world datasets, we show that HGDE framework can not only learn comparably high-quality node embeddings as Euclidean models on non-hierarchical datasets, but outperform all compared hyperbolic models variants on highly-hierarchical datasets with improved efficiency and accuracy. [The code and appendix can be found in https://github.com/ljxw88/HyperbolicGDE.](https://github.com/ljxw88/HyperbolicGDE)

2 Preliminaries

Riemannian Geometry and Hyperbolic Space. A Riemannian manifold \mathcal{M} of n -dimension is a topological space associated with a metric tensor g , denoted as (\mathcal{M}, g) , which extends curved surfaces to higher dimensions and can be locally approximated by \mathbb{R}^n . At any point $\mathbf{x} \in \mathcal{M}$, the tangent space $\mathcal{T}_{\mathbf{x}}\mathcal{M} \cong \mathbb{R}^n$ represents the first-order approximation of a small perturbation around \mathbf{x} , isomorphic to Euclidean space. The Riemannian metric g on the manifold determines a smoothly varying positive definite inner product on the tangent space, enabling the definition of diverse properties *e.g.* geodesic length, angles, and curvature.

The hyperbolic space \mathbb{H}^n is a smooth Riemannian manifold with a constant negative sectional curvature $\kappa < 0$. Its coordinates can be represented via various isometric models. [3] established the equivalence of hyperbolic and Euclidean geometry through the utilization of the n -dimensional *Poincaré ball model*, which equips an open ball $\mathbb{D}_{\kappa}^n = (\mathcal{D}_{\kappa}^n, g^{\mathbb{D}})$, with point set $\mathcal{D}_{\kappa}^n = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < -\frac{1}{\kappa}\}$ and Riemannian metric $g_{\mathbf{x}}^{\mathbb{D}} = (\lambda_{\mathbf{x}}^{\kappa})^2 \mathbf{I}_n$, where the conformal factor $\lambda_{\mathbf{x}}^{\kappa} = \frac{2}{1+\kappa\|\mathbf{x}\|^2}$. The Poincaré metric tensor induces various geometric properties *e.g.* distances $d_{\mathbb{D}}^{\kappa}(\mathbf{x}, \mathbf{y})$, inner products $\langle \mathbf{u}, \mathbf{v} \rangle_{\mathbf{x}}^{\kappa}$, geodesics $\gamma_{\mathbf{x} \rightarrow \mathbf{y}}(t)$ and more [26]. Geodesics also induce the definition of *exponential* and *logarithmic* maps [13]. At point $\mathbf{x} \in \mathbb{D}_{\kappa}^n$, the exponential map $\exp_{\mathbf{x}}^{\kappa} : \mathcal{T}_{\mathbf{x}}\mathbb{D}_{\kappa}^n \rightarrow \mathbb{D}_{\kappa}^n$ essentially maps a small perturbation of \mathbf{x} by $\mathbf{v} \in \mathcal{T}_{\mathbf{x}}\mathbb{D}_{\kappa}^n$ to $\exp_{\mathbf{x}}^{\kappa}(\mathbf{v}) \in \mathbb{D}_{\kappa}^n$, so that $t \in [0, 1] : \exp_{\mathbf{x}}^{\kappa}(t\mathbf{v})$ is the geodesic from \mathbf{x} to $\exp_{\mathbf{x}}^{\kappa}(\mathbf{v})$. The logarithmic map $\log_{\mathbf{x}}^{\kappa} : \mathbb{D}_{\kappa}^n \rightarrow \mathcal{T}_{\mathbf{x}}\mathbb{D}_{\kappa}^n$ is defined as the inverse of $\exp_{\mathbf{x}}^{\kappa}$. Finally, the parallel transport $\mathcal{PT}_{\mathbf{x} \rightarrow \mathbf{y}} : \mathcal{T}_{\mathbf{x}}\mathbb{D}_{\kappa}^n \rightarrow \mathcal{T}_{\mathbf{y}}\mathbb{D}_{\kappa}^n$ moves a tangent vector $\mathbf{v} \in \mathcal{T}_{\mathbf{x}}\mathbb{D}_{\kappa}^n$ along the geodesic to $\mathcal{T}_{\mathbf{y}}\mathbb{D}_{\kappa}^n$ while preserving the metric tensor. For closed-form expression of above operations, please refer to Appendix B.

Diffusion Equations. The process of generating representations of individual data points through information flows can be characterized by an an-isotropic *diffusion* process, a concept borrowed from physics used to describe heat diffusion on Riemannian manifold. Denote the manifold as \mathcal{M} , and let $z(t)$ denote a family of functions on $\mathcal{M} \times [0, \infty)$ and $z(u, t)$ be the density at location $u \in \mathcal{M}$ and times t . The general framework of diffusion equations is expressed as a PDE

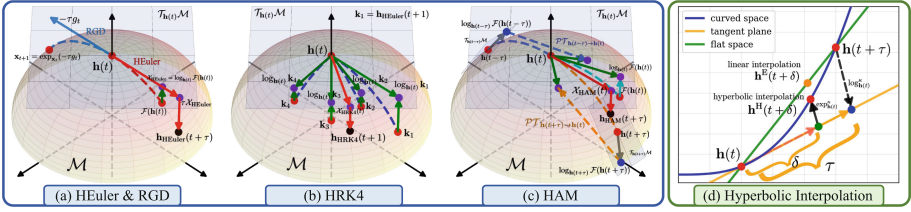


Fig. 1. (a–c) Illustration of various numerical integration methods with comparison to RGD. In each time-step, an explicit scheme calibrates the vector field within only the tangent space of time t , while an implicit scheme requires multiple tangent spaces to estimate future slopes, thus requiring parallel transport for aligning the directions of vectors in different spaces. (d) Illustration of hyperbolic interpolation method.

$$\partial z(u, t) / \partial t = \operatorname{div}(a(z(u, t)) \nabla z(u, t)), \quad t > 0 \quad (1)$$

where $a(\cdot)$ defines the *diffusivity* function controlling the diffusion strength between any location pair at time t . The gradient operator $\nabla : \mathcal{M} \rightarrow \mathcal{TM}$ describes the steepest change at point $u \in \mathcal{M}$. $\operatorname{div}(\cdot) : \mathcal{TM} \rightarrow \mathcal{M}$ is the divergence operator that summarizes the flow of the diffusivity-scaled vector field $(a(\cdot)\nabla)$. Equation (1) can be physically viewed as a variation of heat based on time at the location i , identical to the heat that flows through that point from the surrounding areas.

Graph Diffusion Equation. Let $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ denote an undirected graph with the node set \mathcal{V} and the edge set \mathcal{E} . Let $\mathbf{x} = \{\mathbf{x}_i \in \mathbb{R}^d\}_{i=1}^{|\mathcal{V}|}$ be the node features and $\mathbf{z}(t)$ be node embeddings at time t . Process Eq. (1) can be re-written as

$$\partial \mathbf{z}_i(t) / \partial t = \operatorname{div}(\mathbf{A}(\mathbf{z}(t)) \nabla \mathbf{z}_i(t)), \quad (2)$$

where \mathbf{A} is generally realised by a time-independent $n \times n$ attention matrix [4, 5], consistent with the flow of *heat flux* in/out node i . The formulation of Eq. (2) as a PDE allows leveraging vast existing numerical integration methods to solve the continuous dynamics.

3 Hyperbolic Numerical Integrators

Consider the continuous form of ODE/PDE specified by a neural network parameterized by θ , expressed as

$$d\mathbf{h}(t) / dt = f_\theta(\mathbf{h}(t), t), \quad \mathbf{h}(0) = \mathbf{h}_0 \quad (3)$$

where the time step $t = [0, T]$. Equation (3) essentially tells that the *rate of change* of $\mathbf{h}(t) \in \mathbb{R}^n$ at each time step is given by the vector field $f_\theta : \mathbb{R}^n \times \mathbb{R} \rightarrow \mathbb{R}^n$. Equation (3) is integrated to obtain $\mathbf{h}(T)$. In our context, we are interested in formulating a PDE recipe that is aware of hyperbolic geometry.

Definition 1. A time-dependent **manifold vector field** is a mapping $\mathcal{X} : \mathcal{M} \times \mathbb{R} \rightarrow \mathcal{TM}$, which assigns each point in \mathcal{M} at t a tangent vector. The particle's time-evolution according to \mathcal{X} is then given by the following PDE

$$d\mathbf{h}(t)/dt = \mathcal{X}_\theta(\mathbf{h}(t), t). \quad (4)$$

Definition 2. A **vector flow** is a mapping generated by vector field, i.e. $\mathcal{F} \equiv \pi(\mathcal{X})$, where $\pi : \mathcal{M} \rightarrow \mathcal{M}$ is a smooth projection of vector field to manifold of their local coordinates. Vice versa, if π is a diffeomorphism, then $\mathcal{X} \equiv \pi^{-1}(\mathcal{F})$.

In hyperbolic geometry, where π and π^{-1} are properly defined exp and log maps, our concern lies in the particle's location on the manifold subsequent to integration, i.e. integrate through the path defined by flow \mathcal{F} . This can be achieved via the spirit of projective method [17]. In the following, we derive numerical solvers for estimating the integral of field \mathcal{X} or flow \mathcal{F} w.r.t. time t using, respectively, the explicit and implicit schemes.

3.1 Hyperbolic Projective Explicit Scheme

In an explicit scheme, the state at the next time step is computed directly from the current state and its derivatives. In this part, we derive the hyperbolic generalization of the explicit scheme. To illustrate high-level ideas, we introduce both *single step method* and *multi-step method*. We also discuss the geometric intuition and strong analogy between one-step explicit scheme and Riemannian gradient descent (RGD).

H-Explicit Euler (HEuler). Consider a small time step τ . Iteratively, we seek an approximation for $\mathbf{h}(t + \tau)$ based on $\mathbf{h}(t)$ and vector field $f(\cdot)$. In Euclidean space, the explicit Euler method is written as

$$\mathbf{h}(t + \tau) \approx \mathbf{h}(t) + \tau f_\theta(\mathbf{h}(t), t), \quad (5)$$

which is a discrete version of Eq. (3). Similarly in hyperbolic space, we discretize Eq. (4), and have the stepping function formulated as

$$\mathbf{h}_{\text{HEuler}}(t + \tau) = \exp_{\mathbf{h}(t)}^\kappa(\tau \mathcal{X}_{\text{HEuler}}(t)), \quad (6)$$

where the vector field \mathcal{X} gives the direction at time t according to flow $\mathcal{F}_\theta^\kappa$

$$\mathcal{X}_{\text{HEuler}}(t) = \log_{\mathbf{h}(t)}^\kappa(\mathcal{F}_\theta^\kappa(\mathbf{h}(t), t)) \in \mathcal{T}_{\mathbf{h}(t)}\mathbb{D}_\kappa^n. \quad (7)$$

Geometric Intuition. The equation in Eq. (5) signifies a transition from $\mathbf{h}(t)$ in the direction of f by a distance proportional to τ . In hyperbolic space, where $\mathbf{h}(t) \in \mathbb{D}_\kappa^n$ and we presume $\mathcal{X}^\kappa : \mathbb{D}_\kappa^n \rightarrow \mathcal{T}\mathbb{D}_\kappa^n$, the transition follows the geodesic dictated by the direction of \mathcal{X}^κ . Recall the definition of exponential map: given $x \in \mathbb{D}_\kappa$, $\exp_x^\kappa(v)$ takes $v \in \mathcal{T}_x\mathbb{D}_\kappa$ and returns a point in \mathbb{D}_κ reached by moving from x along the geodesic determined by the tangent vector v . Thus Eqs. (6–7)

can be essentially viewed as a geometric transportation of points on manifold along the curve defined by \mathcal{F} .

Connection to RGD. As visualized in Fig. 1(a), the explicit Euler can be viewed as reversed RGD, where the direction $\mathcal{X}_{\text{HEuler}}(t)$ plays similar role as the Riemannian gradient g_t at $\mathbf{h}(t)$. Similar to RGD, when (\mathcal{M}, ρ) is Euclidean space $(\mathbb{R}^n, \mathbf{I}_n)$, then Eq. (6) converges to Eq. (5) since we have $\exp_h^\kappa(v) \xrightarrow{\kappa \rightarrow 0} h + v$. This property is useful on developing higher-order integrators.

H-Runge-Kutta (HRK). With a similar geometric intuition, we derive the hyperbolic extension of the Runge-Kutta method. Define the s -order HRK stepping function

$$\mathbf{h}_{\text{HRK}}(t + \tau) = \exp_{\mathbf{h}(t)}^\kappa(\tau \mathcal{X}_{\text{HRK}}(t)), \tag{8}$$

where the vector field is estimated by

$$\mathcal{X}_{\text{HRK}}(t) = \left(\sum_{i=1}^s \phi_i \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_i) \right) / \sum_{i=1}^s \phi_i. \tag{9}$$

In Eq. (9), \mathbf{k} denotes the vector flow functions, $\{\phi_i\}$ are coefficients determined by the order. Specifically for 4th order Runge-Kutta (HRK4), we have $\{\phi_{1\dots 4}\} = \{1, 3, 3, 1\}$ derived from Taylor series expansion as in [8]. The vector flows $\mathbf{k}_{1\dots 4}$ are respectively formulated by

$$\mathbf{k}_1 = \mathbf{h}_{\text{HEuler}}(t + \tau), \tag{Eq. (6)} \tag{10}$$

$$\mathbf{k}_2 = \mathcal{F}_\theta^\kappa(\exp_{\mathbf{h}(t)}^\kappa(\tau \mathcal{X}_{\mathbf{k}_2}), t + \tau/3), \text{ where } \mathcal{X}_{\mathbf{k}_2} = \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_1)/3.$$

$$\mathbf{k}_3 = \mathcal{F}_\theta^\kappa(\exp_{\mathbf{h}(t)}^\kappa(\tau \mathcal{X}_{\mathbf{k}_3}), t + 2\tau/3), \text{ where } \mathcal{X}_{\mathbf{k}_3} = \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_2) - \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_1)/3.$$

$$\mathbf{k}_4 = \mathcal{F}_\theta^\kappa(\exp_{\mathbf{h}(t)}^\kappa(\tau \mathcal{X}_{\mathbf{k}_4}), t + \tau), \text{ where } \mathcal{X}_{\mathbf{k}_4} = \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_1) - \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_2) + \log_{\mathbf{h}(t)}^\kappa(\mathbf{k}_3).$$

As illustrated in Fig. 1(b), this method approximates the solution to the PDE within a small interval, considering not only the derivative at the initial time (as in Eq. (5)), but also at intermediate points and the end of the interval.

3.2 Hyperbolic Projective Implicit Scheme

In an implicit scheme, the state of the next iteration is computed by incorporating its own value. This requires solving a linear system to obtain $\mathbf{h}(t + \tau)$ based on $\mathbf{h}(t)$. In below, we illustrate a hyperbolic generalization of the implicit solver.

H-Implicit Adams-Moulton (HAM). Adams numerical integration methods are introduced as families of multi-step methods. With order $s = 0$, Adams methods are identical to the Euler’s method. Principally, there are two types of Adams methods, namely, Adams-Bashforth (explicit) and Adams-Moulton (implicit). Our emphasis is on the latter.

The implicit nature of AM requires the initialization of first several steps with a different method. We use the hyperbolic Runge-Kutta (Eq. (8)) for initialization. With the input $\mathbf{h}(t) \in \mathbb{D}_\kappa^n$ and flow \mathcal{F}^κ , define the warm up

$$\mathbf{h}_{\text{HAM}}(i\tau) = \mathbf{h}_{\text{HRK4}}(i\tau), \quad 0 \leq i < s_{\min} \quad (11)$$

where s_{\min} is the min order. During the whole warm up process, we maintain a queue \mathbf{q} of tangent vectors and the points spanning the tangent space. In each time step of Eq. (11), we push $[q_0 = \mathcal{X}_{\text{RK4}}(i\tau), q_1 = \mathbf{h}(i\tau)]$ to the head of \mathbf{q} . When $\text{len}(\mathbf{q}) \geq s_{\min}$, we start the time-stepping

$$\mathbf{h}_{\text{HAM}}(t + \tau) = \exp_{\mathbf{h}(t)}^{\kappa_t}(\tau \mathcal{X}_{\text{HAM}}(t)), \quad (12)$$

where the vector field is expressed as

$$\begin{aligned} \mathcal{X}_{\text{HAM}}(t) = & \phi_0 \mathcal{P}\mathcal{T}_{\mathbf{h}(t+\tau) \rightarrow \mathbf{h}(t)}(\log_{\mathbf{h}(t+\tau)}^{\kappa_t}(\mathcal{F}_{\theta}^{\kappa_t}(\mathbf{h}(t + \tau), t + \tau))) \\ & + \sum_{i=1}^s \phi_i \mathcal{P}\mathcal{T}_{\mathbf{q}_{i,1} \rightarrow \mathbf{h}(t)}(\mathbf{q}_{i,0}). \end{aligned} \quad (13)$$

The order $s = \min(\text{len}(\mathbf{q}), s_{\max})$, $\{\phi_i\}$ are coefficients determined by the order, which are typically within a pre-defined look-up table. As illustrated in Fig. 1(c), since the reference point $\mathbf{h}(t)$'s stored in \mathbf{q} are different, the parallel transport $\mathcal{P}\mathcal{T}$ is leveraged for aligning tangent spaces for different slopes. When $\mathbf{h}_{\text{HAM}}(t + \tau)$ is accepted as converged, $[\log_{\mathbf{h}(t)}^{\kappa_t}(\mathbf{h}_{\text{HAM}}(t + \tau)), \mathbf{h}(t)]$ is pushed to \mathbf{q} for the next iteration and the last element is popped if $\text{len}(\mathbf{q})$ reaches s_{\max} . We refer readers to Appendix C for detailed explanation of the algorithms.

3.3 Interpolation on Curved Space

Fixed grid PDE solvers typically use their own internal step sizes τ to advance the solution of the PDE. For certain time step t , given $\mathbf{h}(t)$ and $\mathbf{h}(t + \tau)$, we may want to obtain the solution at time point $t + \delta$ where $0 < \delta < \tau$. Since δ does not lie on the grid defined by $\{0, \tau\}$, interpolation methods are invoked to estimate $\mathbf{h}(t + \delta)$. For hyperbolic geometry that $\mathbf{h} \in \mathbb{D}_{\kappa}^n$, define the interpolation

$$\mathbf{h}(t + \delta) = \exp_{\mathbf{h}(t)}^{\kappa_t} \left(\delta \log_{\mathbf{h}(t)}^{\kappa_t}(\mathbf{h}(t + \tau)) / \tau \right). \quad (14)$$

Proposition 1 (Proved in Appendix D). *For any step size $0 < \delta < \tau$, the interpolation $\mathbf{h}(t + \delta)$ via Eq. (14) is on the geodesic between $\mathbf{h}(t)$ and $\mathbf{h}(t + \tau)$ on the manifold, and $\frac{d_{\mathbb{D}}^{\kappa_t}(\mathbf{h}(t), \mathbf{h}(t + \delta))}{d_{\mathbb{D}}^{\kappa_t}(\mathbf{h}(t), \mathbf{h}(t + \tau))} = \frac{\delta}{\tau}$ where $d_{\mathbb{D}}^{\kappa}$ is the geodesic length.*

4 Diffusing Graphs in Hyperbolic Space

4.1 Hyperbolic Graph Diffusion Equation

We study the diffusion process of graphs with node representation residing in hyperbolic geometry. Given the diffusion time $t \in [0, T]$, embedding space $\mathbb{D}_{\kappa_t}^d$ with learnable curvature κ_t at time t , node embedding $\mathbf{z}_*(t) \in \mathbb{D}_{\kappa_t}^d$ and $\mathcal{C}(\cdot)$ being the correlated coordinates of certain node, we formulate the vector flow $\mathcal{F}_{\theta}^{\kappa}$ of the i th representation as

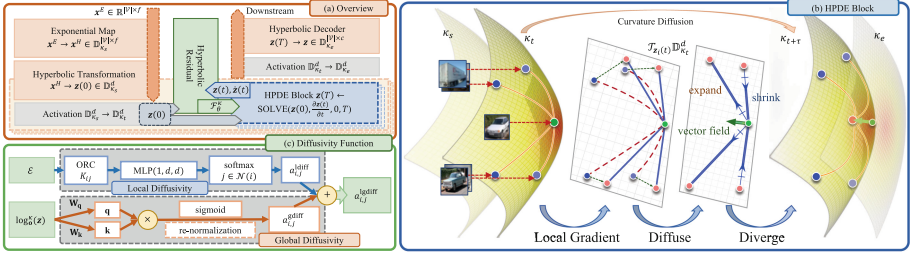


Fig. 2. Schematic of HGDE. (a) The pipeline of our method includes hyperbolic projection, feature transformation, and HPDE block that integrates the GDE. After that, a decoder is applied to the embeddings for specific downstream tasks. (b) The visualization of the diffusion process within the HPDE block: first, map local gradients of \mathbf{z}_i onto the tangent space, calculate the diffusivity, and diverge to obtain the vector flow (green arrow), then perform one-step integration on the manifold with the guidance of continuous curvature diffusion. (c) The details of attention-powered local-global diffusivity function. (Color figure online)

$$\underbrace{\exp^{\kappa_t}_{\mathbf{z}_i(t)}}_{\text{divergence}} \left(\sigma \left[\underbrace{\sum_{j \in \mathcal{C}(i)} a(\mathbf{z}_i(t), \mathbf{z}_j(t))}_{\text{diffusivity}} \underbrace{\log^{\kappa_t}_{\mathbf{z}_i(t)}(\mathbf{z}_j(t))}_{\text{gradient}} \right] \right), \quad (15)$$

where σ can be either identity/non-linear activation. With initial state encoded by learnable feature transformation ψ , *i.e.* $\mathbf{z}(0) = \psi(\mathbf{x}) \in \mathbb{D}_{\kappa}^d$, the final state can be numerically estimated by our proposed HPDE integrators, *i.e.* $\mathbf{z}_i(T) = \text{HPDESolve}(\mathbf{z}_i(0), \frac{\partial \mathbf{z}_i(t)}{\partial t}, 0, T)$. In matrix form, the vector flow is expressed as

$$\mathcal{F}_{\theta}^{\kappa}(\mathbf{z}(t), t) = \exp^{\kappa_t}_{\mathbf{z}(t)} \left(\sigma[\mathbf{S}(\mathbf{z}(t))\nabla\mathbf{z}(t)] \right), \quad (16)$$

where $\mathbf{S}(\mathbf{z}(t)) = (a(\mathbf{z}_i(t), \mathbf{z}_j(t)))$ is a normalized $|\mathcal{V}| \times |\mathcal{V}|$ similarity matrix, and $\nabla\mathbf{z}(t)_{i,j} := \log^{\kappa_t}_{\mathbf{z}_i(t)}(\mathbf{z}_j(t))$. In below, we discuss the key ingredients of Eq. (15, 16).

Gradient. The gradient of a function $z(u, t)$ at location u in a discrete space can be approximated as the difference between the function values at neighboring points. In graph space, let \mathbf{z}_i and $\{\mathbf{z}_j\}_{j \in \mathcal{C}(i)}$, respectively, denote the target node and the correlated positions of i that can be modeled by edge connectivity or self-attention. The graph diffusion process [4, 5] treats nodes as Euclidean representations, such that the analogy of gradient operator $(\nabla\mathbf{z}(t))_{i,j} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is expressed as $\mathbf{z}_j(t) - \mathbf{z}_i(t)$. However, when nodes are embedded in Riemannian manifolds, the gradient of a node is no longer the difference between itself and neighboring points. Instead, we take vectors in the tangent space at \mathbf{z}_i that are obtained by taking the derivative of \mathbf{z} in all possible directions, *i.e.* $(\nabla\mathbf{z}(t))_{i,j} : \mathbb{D}_{\kappa}^d \rightarrow T\mathbb{D}_{\kappa}^d$ that can be formulated as $\log^{\kappa_t}_{\mathbf{z}_i(t)}(\mathbf{z}_j(t))$. One recovers the discrete Euclidean gradient as the curvature $\kappa \rightarrow 0$.

Diffusivity. The diffusivity scales the gradient, with either isotropic or anisotropic behavior. For graph diffusion, the **isotropic** formula is presented

by the normalized adjacency matrix [23], where $a_{i,j} = \frac{1}{\sqrt{d_i d_j}}$ iff. $(i, j) \in \mathcal{E}$ and d is the degree.

Alternatively, the **anisotropic** approach incorporates the attention mechanism [33] to account for the asymmetric relationship between pairs of nodes. This paper considers *local*, *global* and *local-global* schemes based on structure information. Define the schemes

$$\begin{cases} a_{i,j}^{\text{ldiff}} = \text{normalize}_{j \in \mathcal{N}(i)} (f_\theta(\mathbf{z}_i(t), \mathbf{z}_j(t))) & \text{(local scheme)} \\ a_{i,j}^{\text{gdiff}} = \beta \text{normalize}_{j \in \mathcal{V}} (g_\phi(\mathbf{z}_i(t), \mathbf{z}_j(t))) + \frac{1-\beta}{\sqrt{d_i d_j}} & \text{(global scheme)} \\ a_{i,j}^{\text{lgdiff}} = \beta \text{normalize}_{j \in \mathcal{V}} (g_\phi(\mathbf{z}_i(t), \mathbf{z}_j(t))) + (1-\beta) a_{i,j}^{\text{ldiff}} & \text{(local-global scheme)} \end{cases}$$

where f_θ/g_ϕ are learnable functions that compute the diffusivity weight between node pair $(i, j) \in \mathcal{E}$. β can be constant or trainable parameters that adjust the emphasis on homophilic (local attention) and heterophilic (high-order, global attention) relations. In comparison, the local attention scheme implicitly incorporates the graph information since only neighboring elements are considered based on $\mathcal{N}(i)$. Whereas for global attention, it neglects the graph topology and hence requires manual incorporation.

Low-Order Local Diffusivity. A straightforward approach is to leverage the formula of graph attention [34], which is extended to the hyperbolic space by [6], where the weights are calculated tacitly in the tangent space. An alternative method to consider is the Oliver-Ricci Curvature (ORC) [27] attention, introduced in [38, 40] to drive message propagation. This approach is not limited by the non-Euclidean property of node feature, as it computes attention weight via the ORC value derived from the graph topology, thus allowing adoption without leveraging tangent space.

High-Order Global Diffusivity. Propagation of high-order node pairs results in exponentially increasing complexity compared to f_θ . [36, 42] introduced a series of scalable and efficient node-level transformers. With a similar notion in the hyperbolic space, we first project the embeddings onto the tangent space of the origin. Subsequently, the weights can be obtained using existing graph transformer architectures. We adopt energy-constrained transformers [36] with a sigmoid kernel, which performs well in most scenarios.

* Figure 2(c) presents the high-level schematic of *diffuse*. The implementation and algorithmic details are delegated to Appendix C.

Divergence. For simplicity, we assume any $\mathbf{x}_i \in \mathbb{R}^d$ to be scalar-valued. The divergence at a point \mathbf{z}_i is a measure of how much the vector field $\mathcal{X} = \{\nabla \mathbf{z}(t)\}_{i,j} \in \mathcal{C}(i)$ is expanding or contracting at \mathbf{z}_i . In a Euclidean space, the divergence would indeed be the sum of the components of the gradient, *i.e.*, $\text{div}_i = \sum_j (\nabla \mathbf{z}(t))_{i,j}$, producing a scalar (with dimensionality $\mathcal{T}_{\mathbf{z}_i} \mathbb{D}^d \cong \mathbb{R}^d$). In our context, we are interested in how \mathbf{z}_i is varied in the manifold rather than in the tangent space; thus an exponential map is applied to the sum of gradients on $\mathcal{T}_{\mathbf{z}_i} \mathbb{D}^d$, giving $\text{div}_i = \exp_{\mathbf{z}_i(t)}(\sum_j a_{i,j} (\nabla \mathbf{z}(t))_{i,j})$. This also satisfies the form of $\mathcal{F}_\theta^\kappa$ in Definition 2, and thus can be numerically integrated through HPDESolve.

Continuous Curvature Diffusion. Equation (16) implicitly guides the manifold towards its optimal geometry for embedding $\mathbf{z}(t)$ as the manifold parameter κ_t also accumulates and is updated during backpropagation. Similar to the attention parameters θ , we let κ be time independent based on the assumption that $\lim_{\tau \rightarrow 0} \frac{\kappa_{t+\tau} - \kappa_t}{\tau} = 0$.

4.2 Convergence of Dirichlet Energy

Definition 3 Given the node embedding $\{\mathbf{z}_i \in \mathbb{D}_\kappa^d\}_{i=1}^{|\mathcal{V}|}$, the hyperbolic Dirichlet energy is defined as

$$f_{\text{DE}}^\kappa(\mathbf{z}) = \frac{1}{2} \sum_{(i,j) \in \mathcal{E}} d_{\mathbb{D}}^\kappa \left(\exp_{\mathfrak{o}}^\kappa \left(\frac{\log_{\mathfrak{o}}^\kappa(\mathbf{z}_i)}{\sqrt{1+d_i}} \right), \exp_{\mathfrak{o}}^\kappa \left(\frac{\log_{\mathfrak{o}}^\kappa(\mathbf{z}_j)}{\sqrt{1+d_j}} \right) \right)^2, \quad (17)$$

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 B.

Definition 3 introduces a node-similarity measure to quantify over-smoothness in hyperbolic space. f_{DE}^κ of node representation can be viewed as the weighted sum of distance between normalized node pairs. [25, Prop. 4] proved that hyperbolic energy f_{DE}^κ diminishes after message passing, and multiple aggregations result in converging towards zero energy, indicating reduced embedding expressiveness that could potentially cause over-smoothing. Also as proved in [43, Prop. 2] that over-smoothing is an intrinsic property of first-order continuous GNN. In a continuous diffusion process, where each iteration can be viewed as a layer in HGNNs, as supported by Fig. 3, we also observe a convergence of hyperbolic Dirichlet energy of $\mathbf{z}(t)$ w.r.t. time t .

Residual-Empowered Flow. Empirically, studies in multi-layer GNNs [15, 24] demonstrated the efficacy of adding residual connections to the initial layer. It is also claimed in [45] that using residual connections for both initial and previous layers can prevent the Dirichlet energy from reaching a lower energy limit, thus avoiding over-smoothing. Building upon these studies, we define the hyperbolic residual empowered vector flow

$$\mathcal{F}_\theta^\kappa(\mathbf{z}(t), t) = \mu_{\mathbb{D}}^\kappa(\{\dot{\mathbf{z}}(t), \mathbf{z}(t), \mathbf{z}(0)\}; \{\eta\}_{j=1}^J), \quad (18)$$

where $\dot{\mathbf{z}}(t) = \exp_{\mathbf{z}(t)}^{\kappa_t}(\sigma[\mathbf{S}(\mathbf{z}(t))\nabla\mathbf{z}(t)])$ is the manifold dynamic as in Eq. (16). $\{\eta\}_{j=1}^J$ are the weight coefficients. $\mu_{\mathbb{D}}^\kappa$ is the node-wise hyperbolic averaging. We instantiate it via *Möbius Gyromidpoint* [32] for its trade-off between computational cost and precision. Define

$$\mu_{\mathbb{D}}^\kappa(\{\mathbf{z}\}_{j=1}^J; \{\eta\}_{j=1}^J) = \left(\frac{1}{2} \otimes_{\kappa} \left(\frac{\sum_j \eta_j \lambda_{\mathbf{z}_i}^{\kappa(j)} \mathbf{z}_i^{(j)}}{\sum_j |\eta_j| (\lambda_{\mathbf{z}_i}^{\kappa(j)} - 1)} \right) \right)_{i=1}^{|\mathcal{V}|}. \quad (19)$$

This operation ensures the point set constraint of \mathbb{D} for the residual flow. We recover the arithmetic mean as $\kappa \rightarrow 0$. During diffusion, Eq. (18) retains at least a portion of the initial and prior embeddings. Since the initial embedding possesses high energy, the residual connection mitigates energy degradation and retains the energy of the final iteration at the same level as the preceding iterations.

5 Empirical Results

5.1 Experiment Setup

Datasets. Under homophilic setting, we consider 5 datasets for node classification and link prediction: DISEASE, AIRPORT (transductive datasets, provided in [6] to investigate the tree-likeness modeling), PUBMED, CITESEER and CORA ([39] widely used citation networks), which are summarized in the table in Appendix A. Additionally, we report the Gromov’s hyperbolicity δ given by [16] for each dataset. A graph is more hyperbolic as $\delta \rightarrow 0$ and is a tree when $\delta = 0$.

For heterophilic datasets, we evaluate node classification on three heterophilic graphs, respectively, CORNELL, TEXAS and WISCONSIN [29] from the WebKB dataset (webpage networks). Detailed statistics are summarized in Appendix A. We use the original fixed 10 split datasets. In addition, 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.

Baselines. We compare our models to (1) *Euclidean-hyperbolic* baselines, (2) *discrete-continuous depth* baselines and (3) *heterophilic relationship* baselines. For (1), we compare against feature-based models, Euclidean, and hyperbolic graph-based models. Feature-based models: without using graph structure, we feed node feature directly to MLP and HNN [14]; Euclidean graph-based models: GCN [23], GAT [34], GraphSAGE [19], and SGC [35]; Hyperbolic graph-based models: HGCN [6], κ GCN [1], LGCN [44] and HyboNet [10]. For (2), we compare our models on citation networks with the discrete-continuous depth models. Discrete depth: GCNII [7], C-DropEdge [20]; Discrete-decouple: HyLa-SGC [41]; Continuous depth: GDE [30], GRAND and BLEND [4, 5]. For (3), we compare to the prevalent GNNs: GCN, GAT, HGCN, HyboNet, and those optimized for heterophilic relationships: H2GCN [46], GCNII, GraphSAGE and GraphCON [31]. The test results are partially derived from the above works. For fairness, we compare to models with no more than 16 layers/iterations. Please refer to Appendix A for more details regarding the compared baselines. We detail the parameter settings for model and evaluation metric in Appendix C.

5.2 Experiment Results

Euclidean-Hyperbolic Baselines. We investigate our methods with different solvers with $\tau = 1$, *i.e.* HGDE-E (multi-step explicit integrator, HRK4) and HGDE-I (multi-step implicit integrator, HAM). The experimental results are summarized in Tables 1 and 2. (1) Our proposed models outperform previous

Table 1. Test accuracy (%) for node classification task.

Dataset	DISEASE	AIRPORT	PUBMED	CITeseer	CORA
δ	0	1	3.5	5	11
MLP	32.5 \pm 1.1	60.9 \pm 3.4	72.4 \pm 0.2	59.5 \pm 0.9	51.6 \pm 1.3
HNN	45.5 \pm 3.3	80.6 \pm 0.5	69.9 \pm 0.4	59.5 \pm 1.2	54.7 \pm 0.6
GCN	69.7 \pm 0.4	81.6 \pm 0.6	78.1 \pm 0.4	70.3 \pm 0.4	81.5 \pm 0.5
GAT	70.4 \pm 0.4	82.7 \pm 0.4	78.2 \pm 0.4	71.6\pm0.8	83.0 \pm 0.5
SAGE	69.1 \pm 0.6	82.2 \pm 0.5	77.5 \pm 2.4	67.5 \pm 0.7	79.9 \pm 2.5
SGC	69.5 \pm 0.2	80.6 \pm 0.2	78.8 \pm 0.2	71.4 \pm 0.8	81.3 \pm 0.5
HGCN	82.8 \pm 0.8	89.2 \pm 1.3	80.3\pm0.3	68.0 \pm 0.6	79.9 \pm 0.2
κ GCN	82.1 \pm 1.1	84.4 \pm 0.4	78.3 \pm 0.6	71.1 \pm 0.6	80.8 \pm 0.6
LGCN	84.4 \pm 0.8	90.9\pm1.0	78.8 \pm 0.5	71.1 \pm 0.3	83.3\pm0.5
HyboNet	96.0\pm1.0	90.9 \pm 1.4	78.0 \pm 1.0	69.8 \pm 0.6	80.2 \pm 1.3
HGDE-E	92.1\pm1.6	95.1\pm0.4	81.2\pm0.5	74.1\pm0.5	84.4\pm0.7
HGDE-I	90.9\pm2.5	93.9\pm0.8	81.0\pm0.3	73.5\pm0.7	84.0\pm0.4

Table 2. Test ROC AUC (%) results for link prediction task.

Dataset	DISEASE	AIRPORT	PUBMED	CITeseer	CORA
δ	0	1	3.5	5	11
MLP	69.9 \pm 3.4	68.9 \pm 0.5	83.3 \pm 0.6	93.7 \pm 0.6	83.3 \pm 0.6
HNN	70.2 \pm 0.1	80.6 \pm 0.5	94.7 \pm 0.1	93.3 \pm 0.5	90.9 \pm 0.4
GCN	64.7 \pm 0.5	89.3 \pm 0.4	89.6 \pm 3.7	82.6 \pm 1.9	90.5 \pm 0.2
GAT	69.8 \pm 0.3	90.9 \pm 0.2	91.5 \pm 1.8	86.5 \pm 1.5	93.2 \pm 0.2
SAGE	65.9 \pm 0.3	90.4 \pm 0.5	86.2 \pm 0.8	92.1 \pm 0.4	85.5 \pm 0.5
SGC	65.1 \pm 0.2	89.8 \pm 0.3	94.1 \pm 0.1	91.4 \pm 1.7	91.5 \pm 0.2
HGCN	91.2 \pm 0.6	96.4 \pm 0.1	95.1 \pm 0.1	96.6\pm0.1	93.8\pm0.1
κ GCN	92.0 \pm 0.5	92.5 \pm 0.5	94.9 \pm 0.3	95.1 \pm 0.6	92.6 \pm 0.4
LGCN	96.6\pm0.6	96.0 \pm 0.6	96.6\pm0.1	95.8 \pm 0.4	93.6 \pm 0.4
HyboNet	96.8\pm0.4	97.3\pm0.3	95.8 \pm 0.2	96.7\pm0.8	93.6 \pm 0.3
HGDE-E	96.2\pm0.5	98.2\pm0.2	96.6\pm0.2	96.7\pm0.7	94.1\pm0.4
HGDE-I	95.6 \pm 0.5	97.6\pm0.5	96.2\pm0.7	96.4 \pm 0.7	94.5\pm0.8

Euclidean and hyperbolic models in four out of five datasets, suggesting that graph learning in hyperbolic space through topological diffusion is beneficial. (2) Hyperbolic models typically exhibit poor performance on datasets that are less hyperbolic (*e.g.*, CORA), while our method surprisingly exceeds Euclidean

Table 3. Discrete-continuous depth GNN comparison.

Type	Model	CORA	CITeseer	PUBMED
Discrete	GCNII	84.6\pm0.8	72.9 \pm 0.5	80.2 \pm 0.4
	C-DropEdge	82.6 \pm 0.9	71.0 \pm 1.0	77.8 \pm 1.0
Decouple (Hyp PosEnc)	HyLa-SGC	82.5 \pm 0.5	72.6 \pm 1.0	80.3 \pm 0.9
Continuous	GDE	83.8 \pm 0.5	72.5 \pm 0.5	79.9 \pm 0.3
	GRAND	82.9 \pm 0.7	73.6 \pm 0.3	81.0\pm0.4
Continuous (Hyp PosEnc)	BLEND	84.2\pm0.6	74.4\pm0.7	80.7 \pm 0.7
Continuous (Hyp Embed)	HGDE(4)	83.4 \pm 0.5	73.0 \pm 0.3	80.2 \pm 0.6
	HGDE(8)	83.7 \pm 0.6	73.5 \pm 0.7	80.8 \pm 0.4
	HGDE(12)	84.2\pm0.6	74.1\pm0.5	81.2\pm0.5
	HGDE(16)	84.4\pm0.7	73.8\pm0.7	80.9\pm0.3

Table 4. Heterophilic relationship GNN comparison.

Type	\mathcal{H}	TEXAS 0.11	WISCONSIN 0.21	CORNELL 0.30
Euclidean	GCN	55.1 \pm 5.2	51.8 \pm 3.1	60.5 \pm 5.3
	GAT	52.2 \pm 6.6	49.4 \pm 4.1	61.9 \pm 5.1
Hyperbolic	HGCN	55.7 \pm 6.3	48.1 \pm 6.1	62.1 \pm 3.7
	HyboNet	60.0 \pm 4.1	51.2 \pm 3.3	62.3 \pm 3.5
High-Order GNNs	H2GCN	84.9\pm7.2	87.7\pm5.0	82.7\pm5.3
	GCNII	77.6 \pm 3.8	80.4 \pm 3.4	77.9 \pm 3.8
	SAGE	82.4 \pm 6.1	81.2 \pm 5.6	76.0 \pm 5.0
	GraphCON	85.4\pm4.2	87.8\pm3.3	84.3\pm4.8
Ours	HGDE	85.9\pm2.8	86.2\pm2.4	85.0\pm5.3

GAT on datasets with lower δ , indicating the necessity of curvature diffusion in adapting to datasets with scarce hierarchical structures and modeling long-term dependency via the local-global diffusivity function. (3) HGDE and other hyperbolic models achieve superior performance compared to Euclidean counterparts in link prediction due to the larger embedding space in hyperbolic geometry, which better preserves structural dependencies and allows for improved node arrangement. (4) HGDE-E generally outperforms HGDE-I with lower memory

Table 5. Evaluation on image (CIFAR/STL) and text (20News) classification (Left) and Memory & Runtime comparison (Right). \star indicate OOM.

Dataset		MLP	LabelProp	ManiReg	GCN-kNN	GAT-kNN	DenseGAT	GLCN	HGDE
CIFAR	100 labels	65.9 \pm 1.3	66.2	67.0 \pm 1.9	66.7 \pm 1.5	66.0 \pm 2.1	\star	66.6 \pm 1.4	68.9\pm2.1
	500 labels	73.2 \pm 0.4	70.6	72.6 \pm 1.2	72.9 \pm 0.4	72.4 \pm 0.5	\star	72.7 \pm 0.5	74.0\pm1.8
	1000 labels	75.4 \pm 0.6	71.9	74.3 \pm 0.4	74.7 \pm 0.5	74.1 \pm 0.5	\star	74.7 \pm 0.3	76.3\pm0.9
STL	100 labels	66.2 \pm 1.4	65.2	66.5 \pm 1.9	66.9 \pm 0.5	66.5 \pm 0.8	\star	66.4 \pm 0.8	66.9\pm1.3
	500 labels	73.0\pm0.8	71.8	72.5 \pm 0.5	72.1 \pm 0.8	72.0 \pm 0.8	\star	72.4 \pm 1.3	72.5 \pm 0.2
	1000 labels	75.0 \pm 0.8	72.7	74.2 \pm 0.5	73.7 \pm 0.4	73.9 \pm 0.6	\star	74.3 \pm 0.7	75.1\pm0.6
20News	1000 labels	54.1 \pm 0.9	55.9	56.3 \pm 1.2	56.1 \pm 0.6	55.2 \pm 0.8	54.6 \pm 0.2	56.2 \pm 0.8	56.3\pm0.9
	2000 labels	57.8 \pm 0.9	57.6	60.0 \pm 0.8	60.6 \pm 1.3	59.1 \pm 2.2	59.3 \pm 1.4	60.2 \pm 0.7	61.0\pm1.0
	4000 labels	62.4 \pm 0.6	59.5	63.6 \pm 0.7	64.3 \pm 1.0	62.9 \pm 0.7	62.4 \pm 1.0	64.1 \pm 0.8	64.1\pm0.8
T/τ		Model (with Att)		Memory ($\times 10^6$)		Runtime (ms)			
2		HGDCN		4045		9.25			
		HGDCN (LocalAtt)		4246		1310.31			
		LGCN		4630		16.64			
		HyboNet		4368		14.90			
		HGDE		62		13.66			
4		HGDCN		10255		29.77			
		HGDCN (LocalAtt)		10578		4086.08			
		LGCN		12675		40.88			
		HyboNet		11931		35.23			
		HGDE		73		20.28			
8		HGDCN		22674		67.24			
		HGDCN (LocalAtt)		OOM		\star			
		LGCN		23712		160.75			
		HyboNet		23403		122.55			
		HGDE		112		34.11			
16		All Baselines		OOM		\star			
		HGDE		188		61.95			

consumption and better precision, indicating that a larger τ may be necessary for implicit solvers. To align with multi-layer GNN schema (step-size is analogous to depth), we employ HGDE with HRK4 ($\tau = 1$) for further evaluation.

Discrete-Continuous Depth Baselines. In Table 3, we compare our models with discrete and continuous-depth baselines. We observe that our method with $T \in [12, 16]$ achieves competitive results with the state-of-the-art models. Notably, HGDE models consistently outperform discrete models and continuous models with Euclidean embeddings, highlighting the benefits of utilizing hyperbolic embeddings in a continuous-depth framework. Compared to position encoding approaches (*e.g.*, HyLa, BLEND), HGDE exhibits superior performance, indicating the feasibility of using hyperbolic space embeddings directly over initial position encoding. Interestingly, we find HGDE models performs better when increasing T up to 12, but slightly worse at $T = 16$. This may due to the capacity of Poincaré ball or potential over-smoothing. Overall, the results underscore the effectiveness of the proposed HGDE models in harnessing the power of hyperbolic space for graph data modeling.

Heterophilic Relationship Baselines. We show that HGDE is also capable in managing heterophilic relationship. In Table 4, HGDE achieves the highest scores on the TEXAS and CORNELL and a competitive score on WISCONSIN. This shows that hyperbolic space is beneficial in learning hierarchical heterophilic relationships. It also reflects the flexibility of HGDE as a hyperbolic vector flow for embedding high-order structures, with our model, powered by HPDE, outperforming other baselines on average.

Image and Text Classification. We follow the experiment setup in [36] and conduct additional experiments on the CIFAR, STL, and 20NEWS datasets to evaluate HGDE in multiple scenarios with limited label rates. We employ the SimCLR [9] extracted embedding as provided in [36] for image classification. For the pre-processed 20NEWS [28] for text classification, we take 10 topics and regard words with TF-IDF > 5 as features. For graph-based models, we use k NN to construct a graph over input features. For HGDE (hyperbolic), we map the initial feature to \mathbb{D}_κ via $\exp_\circ^\kappa(\cdot)$ before the embedding process. As depicted in Table 5(Left), HGDE consistently surpasses its opponents, including MLP, LabelProp [47], ManiReg [2], GCN- k NN, GAT- k NN, DenseGAT, and GLCN [21]. Across all datasets, HGDE outperforms the Euclidean models, underscoring its proficiency in understanding the potential hierarchical structure of image embeddings [22] and text embeddings. Furthermore, HGDE exhibits good performance compared to static graph-based baselines *e.g.*, GAT- k NN and GLCN, which underlines the distinct advantage of the evolving diffusivity mechanism in understanding the potential hierarchical structure of image/text embeddings.

5.3 Ablation Study

Efficacy of Hyperbolic Residual. Figure 3 visualizes the convergence of hyperbolic energy through iterations. We observe that, without residuals, the averaged energy rapidly decreases to near-zero values, supporting the hypothesis that, without residual connections, the embedding can evolve to an overly smoothed state that is potentially low in expressiveness. However, with hyperbolic residuals, for all three integrators, the average energy decreases over the first few iterations and then appears to stabilize around a certain value above zero. This behavior is consistent across both datasets, suggesting that the system is able to converge to a stable state with non-zero energy (Fig. 4).

Efficacy of Diffusivity Function. Figure 5 visualizes sampled node embeddings and their edge diffusivity on CORA. The blue edges are inherently determined by the graph structure. Red ones are determined by global attention, showing that a^{lgdiff} accounts for high-order relations. The bar graph shows the average accuracy on various datasets produced by HGDE with different diffusivity functions. We find that anisotropic approaches generally outperform the isotropic approach, suggesting the necessity of directional information in the diffusion process. Although the performance degrades on CITESEER when using a^{lgdiff} , there are significant improvements on other graphs, certifying the benefit of higher-order proximity induced by local-global diffusivity.

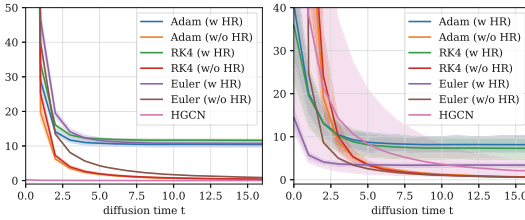


Fig. 3. Hyperbolic Dirichlet energy $f_{DE}^K(\cdot)$ variation through t on CORA (left) and CITESEER (right). Models are compared with different integrators w or w/o hyperbolic residual.

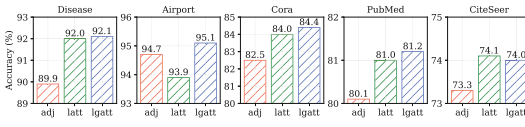


Fig. 4. Averaged node classification performance comparison of models with different diffusivity functions on various datasets.

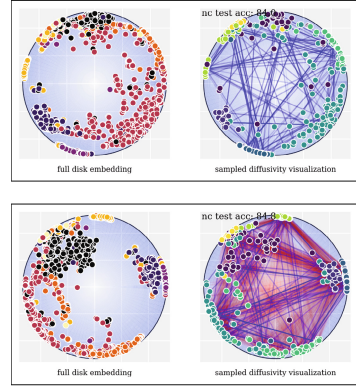


Fig. 5. CORA diffusivity (400 node sampled from \mathbb{D}_K^2 embeddings) produced by a^{ldiff} (left) and a^{lgdiff} (right), blue and red lines denote local and global attention; bolder lines indicate more attentiveness. (Color figure online)

Parameter Efficiency. In Table 5 (Right), we provide an additional comparison of peak GPU memory usage and per-epoch running time on the CORA. We tested HGDE-E where all models have a 16 hidden dim. Our model significantly outperforms the other baselines in both training time (for ≥ 4 layers) and memory consumption. The memory reduction is primarily due to the utilization of sparse attention, and the advantages of a weight-tied network (requiring only single-layer parameters) as a nature of HPDE. The training time efficiency is achieved by eliminating layer-wise feature transformation, implementing weight-tying, and applying scattered-agg for hyperbolic representation.

6 Conclusion

We developed multiple numerical integrators for HPDE, and proposed the first hyperbolic continuous-time embedding diffusion framework – HGDE. Being capable of capturing both low and high order proximity, HGDE outperforms both Euclidean and hyperbolic baselines on various datasets. The effectiveness of HGDE was further validated by the ablation studies on hyperbolic energy and diffusivity functions. The superiority of HGDE underscores the potential of developing PDE-based non-Euclidean models.

Limitation. While HGDE presents strong performance in modeling graph data, hyperbolic spaces may not always be optimal, particularly for data without clear

hierarchical structures. For instance, HGDE is difficult to beat natural Euclidean deep models (*e.g.* GCNII) on the non-hierarchical CORA. Moreover, a higher memory complexity and lower training time only tells the efficiency rather than scalability of HGDE, since our models are evaluated with fixed number of parameters (which is natural for ODE-based models), increasing T is not necessarily *scaling up*. Future work include addressing these limitations and exploring the scalability and generalizability of HGDE in diverse real-world settings.

Acknowledgments. XH is financially supported by the U.K. EPSRC through End-to-End Conceptual Guarding of Neural Architectures [EP/T026995/1]. JL is supported by Liverpool-CSC scholarship [202208890034].

References

1. Bachmann, G., Bécigneul, G., Ganea, O.: Constant curvature graph convolutional networks. In: International Conference on Machine Learning, pp. 486–496. PMLR (2020)
2. Belkin, M., Niyogi, P., Sindhvani, V.: Manifold regularization: a geometric framework for learning from labeled and unlabeled examples. *J. Mach. Learn. Res.* **7**(11) (2006)
3. Beltrami, E.: Teoria fondamentale degli spazii di curvatura costante: memoria. F. Zanetti (1868)
4. Chamberlain, B., Rowbottom, J., Gorinova, M.I., Bronstein, M., Webb, S., Rossi, E.: Grand: graph neural diffusion. In: International Conference on Machine Learning, pp. 1407–1418. PMLR (2021)
5. Chamberlain, B., Rowbottom, J., Eynard, D., Di Giovanni, F., Dong, X., Bronstein, M.: Beltrami flow and neural diffusion on graphs. In: Advances in Neural Information Processing Systems, vol. 34, pp. 1594–1609 (2021)
6. Chami, I., Ying, Z., Ré, C., Leskovec, J.: Hyperbolic graph convolutional neural networks. In: Advances in Neural Information Processing Systems, vol. 32 (2019)
7. Chen, M., Wei, Z., Huang, Z., Ding, B., Li, Y.: Simple and deep graph convolutional networks. In: International Conference on Machine Learning, pp. 1725–1735. PMLR (2020)
8. Chen, R.T., Rubanova, Y., Bettencourt, J., Duvenaud, D.K.: Neural ordinary differential equations. In: Advances in Neural Information Processing Systems, vol. 31 (2018)
9. Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International Conference on Machine Learning, pp. 1597–1607. PMLR (2020)
10. Chen, W., et al.: Fully hyperbolic neural networks. arXiv preprint [arXiv:2105.14686](https://arxiv.org/abs/2105.14686) (2021)
11. Choudhary, N., Reddy, C.K.: Towards scalable hyperbolic neural networks using Taylor series approximations. arXiv preprint [arXiv:2206.03610](https://arxiv.org/abs/2206.03610) (2022)
12. Clauset, A., Moore, C., Newman, M.E.: Hierarchical structure and the prediction of missing links in networks. *Nature* **453**(7191), 98–101 (2008)
13. Ganea, O., Bécigneul, G., Hofmann, T.: Hyperbolic entailment cones for learning hierarchical embeddings. In: International Conference on Machine Learning, pp. 1646–1655. PMLR (2018)

14. Ganea, O., Bécigneul, G., Hofmann, T.: Hyperbolic neural networks. In: *Advances in Neural Information Processing Systems*, vol. 31 (2018)
15. Gasteiger, J., Bojchevski, A., Günnemann, S.: Predict then propagate: graph neural networks meet personalized pagerank. arXiv preprint [arXiv:1810.05997](https://arxiv.org/abs/1810.05997) (2018)
16. Gromov, M.: Hyperbolic groups. In: Gersten, S.M. (ed.) *Essays in Group Theory*. Mathematical Sciences Research Institute Publications, vol. 8, pp. 75–263. Springer, New York (1987). https://doi.org/10.1007/978-1-4613-9586-7_3
17. Hairer, E.: Solving differential equations on manifolds, **8** (2011). <https://www.unige.ch/~hairer/poly-sde-mani.pdf>
18. Hamann, M.: On the tree-likeness of hyperbolic spaces. In: *Mathematical Proceedings of the Cambridge Philosophical Society*, vol. 164, pp. 345–361. Cambridge University Press (2018)
19. Hamilton, W., Ying, Z., Leskovec, J.: Inductive representation learning on large graphs. In: *Advances in Neural Information Processing Systems*, vol. 30 (2017)
20. Huang, W., et al.: Towards deepening graph neural networks: a GNTK-based optimization perspective. arXiv preprint [arXiv:2103.03113](https://arxiv.org/abs/2103.03113) (2021)
21. Jiang, B., Zhang, Z., Lin, D., Tang, J., Luo, B.: Semi-supervised learning with graph learning-convolutional networks. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11313–11320 (2019)
22. Khrulkov, V., Mirvakhabova, L., Ustinova, E., Oseledets, I., Lempitsky, V.: Hyperbolic image embeddings. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6418–6428 (2020)
23. Kipf, T.N., Welling, M.: Semi-supervised classification with graph convolutional networks. arXiv preprint [arXiv:1609.02907](https://arxiv.org/abs/1609.02907) (2016)
24. Li, G., Muller, M., Thabet, A., Ghanem, B.: DeepGCNs: can GCNs go as deep as CNNs? In: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9267–9276 (2019)
25. Liu, J., Yi, X., Huang, X.: DeephGCN: recipe for efficient and scalable deep hyperbolic graph convolutional networks (2024)
26. Nickel, M., Kiela, D.: Poincaré embeddings for learning hierarchical representations. In: *Advances in Neural Information Processing Systems*, vol. 30 (2017)
27. Ollivier, Y.: Ricci curvature of Markov chains on metric spaces. *J. Funct. Anal.* **256**(3), 810–864 (2009)
28. Pedregosa, F., et al.: Scikit-learn: machine learning in Python. *J. Mach. Learn. Res.* **12**, 2825–2830 (2011)
29. Pei, H., Wei, B., Chang, K.C.C., Lei, Y., Yang, B.: Geom-GCN: geometric graph convolutional networks. arXiv preprint [arXiv:2002.05287](https://arxiv.org/abs/2002.05287) (2020)
30. Poli, M., Massaroli, S., Park, J., Yamashita, A., Asama, H., Park, J.: Graph neural ordinary differential equations. arXiv preprint [arXiv:1911.07532](https://arxiv.org/abs/1911.07532) (2019)
31. Rusch, T.K., Chamberlain, B., Rowbottom, J., Mishra, S., Bronstein, M.: Graph-coupled oscillator networks. In: *International Conference on Machine Learning*, pp. 18888–18909. PMLR (2022)
32. Ungar, A.A.: A gyrovector space approach to hyperbolic geometry. *Synth. Lect. Math. Stat.* **1**(1), 1–194 (2008)
33. Vaswani, A., et al.: Attention is all you need. In: *Advances in Neural Information Processing Systems*, vol. 30 (2017)
34. Veličković, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y.: Graph attention networks. arXiv preprint [arXiv:1710.10903](https://arxiv.org/abs/1710.10903) (2017)
35. Wu, F., Souza, A., Zhang, T., Fifty, C., Yu, T., Weinberger, K.: Simplifying graph convolutional networks. In: *International Conference on Machine Learning*, pp. 6861–6871. PMLR (2019)

36. Wu, Q., Yang, C., Zhao, W., He, Y., Wipf, D., Yan, J.: Difformer: scalable (graph) transformers induced by energy constrained diffusion. arXiv preprint [arXiv:2301.09474](https://arxiv.org/abs/2301.09474) (2023)
37. Yan, S., Xiong, Y., Lin, D.: Spatial temporal graph convolutional networks for skeleton-based action recognition. In: Thirty-Second AAAI Conference on Artificial Intelligence (2018)
38. Yang, M., Zhou, M., Pan, L., King, I.: κ HGCN: tree-likeness modeling via continuous and discrete curvature learning (2023)
39. Yang, Z., Cohen, W., Salakhudinov, R.: Revisiting semi-supervised learning with graph embeddings. In: International Conference on Machine Learning, pp. 40–48. PMLR (2016)
40. Ye, Z., Liu, K.S., Ma, T., Gao, J., Chen, C.: Curvature graph network. In: International Conference on Learning Representations (2019)
41. Yu, T., De Sa, C.: Random Laplacian features for learning with hyperbolic space. arXiv preprint [arXiv:2202.06854](https://arxiv.org/abs/2202.06854) (2022)
42. Yun, S., Jeong, M., Kim, R., Kang, J., Kim, H.J.: Graph transformer networks. In: Advances in Neural Information Processing Systems, vol. 32 (2019)
43. Zhang, Y., Gao, S., Pei, J., Huang, H.: Improving social network embedding via new second-order continuous graph neural networks. In: Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, pp. 2515–2523 (2022)
44. Zhang, Y., Wang, X., Shi, C., Liu, N., Song, G.: Lorentzian graph convolutional networks. In: Proceedings of the Web Conference 2021, pp. 1249–1261 (2021)
45. Zhou, K., et al.: Dirichlet energy constrained learning for deep graph neural networks. In: Advances in Neural Information Processing Systems, vol. 34, pp. 21834–21846 (2021)
46. Zhu, J., Yan, Y., Zhao, L., Heimann, M., Akoglu, L., Koutra, D.: Beyond homophily in graph neural networks: current limitations and effective designs. In: Advances in Neural Information Processing Systems, vol. 33, pp. 7793–7804 (2020)
47. Zhu, X., Ghahramani, Z., Lafferty, J.D.: Semi-supervised learning using Gaussian fields and harmonic functions. In: Proceedings of the 20th International Conference on Machine Learning (ICML-03), pp. 912–919 (2003)
48. Zitnik, M., Sosič, R., Feldman, M.W., Leskovec, J.: Evolution of resilience in protein interactomes across the tree of life. *Proc. Natl. Acad. Sci.* **116**(10), 4426–4433 (2019)