

Hyperbolic Deep Neural Networks: A Survey

Wei Peng, Tuomas Varanka, Abdelrahman Mostafa, Henglin Shi, Guoying Zhao* *Senior Member, IEEE*

Abstract—Recently, hyperbolic deep neural networks (HDNNs) have been gaining momentum as the deep representations in the hyperbolic space provide high fidelity embeddings with few dimensions, especially for data possessing hierarchical structure. Such a hyperbolic neural architecture is quickly extended to different scientific fields, including natural language processing, single-cell RNA-sequence analysis, graph embedding, financial analysis, and computer vision. The promising results demonstrate its superior capability, significant compactness of the model, and a substantially better physical interpretability than its counterpart in the Euclidean space. To stimulate future research, this paper presents a comprehensive review of the literature around the neural components in the construction of HDNN, as well as the generalization of the leading deep approaches to the hyperbolic space. It also presents current applications of various tasks, together with insightful observations and identifying open questions and promising future directions.

Index Terms—Neural Networks on Riemannian Manifold, Hyperbolic Neural Networks, Poincaré Model, Lorentz Model.

1 INTRODUCTION

Data with tree structure and hierarchy are ubiquitous in natural scenarios and the real-world [1], [2], [3], [4]. Lots of such data, *e.g.*, Internet [5], brain networks [6], the world trade network [7], and financial networks [8], exhibit a highly non-Euclidean latent anatomy and negatively curved properties [9], [10], [11]. Similar properties can be observed from tasks such as recommendation system [12], social media analysis [13], knowledge graph embedding [14], single-cell RNA-sequence analysis [15] and image retrieval [16]. In parallel, current machine learning, especially deep learning [17] has provided a powerful data-driven way to analyze and understand data. At its core lies the expectation of learning low-dimensional and semantically rich representations of data. Deep neural networks, constructed with a multi-layered structure, parameterized with millions of parameters, and boosted with comprehensive and highly-optimized deep learning libraries [18], [19], [20], theoretically have the potential to fit any complex functions, leading to the domination of many research fields, such as image classification [21], [22], machine translation tasks [23], and even video games playing [24].

Nevertheless, in such successful applications, regular grid data in the Euclidean space, *e.g.*, text and images, is the main focus. Besides, the learning process is conducted in the intuition-friendly Euclidean space, which is a flat space with zero curvature. However, neural architectures operating in the Euclidean space rely heavily on the locality and are originally designed for grid data, thus not necessarily providing the most powerful or meaningful geometrical representations for structured data in a non-Euclidean space. Typically, the non-Euclidean spaces including the elliptic space (a sphere) with a constant positive curvature [25] and the hyperbolic space with constant negative (sectional) curvature [26] should be considered. Neural networks in the elliptic space [27], [28], [29], [30], using spherical harmonic transform or Laplacian-based graph convolution, have been successfully applied to spherical signals. Analogous to this, there is a strong expectation

to construct neural networks in the hyperbolic space for data possessing hierarchies, as hierarchies can be represented in such space with low distortion [31].

Apart from the structured data, the underlying relationships between regular samples and important developmental progresses can be also modeled by hierarchy in the hyperbolic space. From the perspective of cognitive science, it is widely accepted that human beings use hierarchy to organize actions [32] and object categories [33], [34], [35]. An interesting study [36] finds that both natural odors and human perceptual descriptions of smells can be described using a three-dimensional hyperbolic space. In single-cell analysis, the cell developmental processes show strong hierarchical relationships and can be described by a hyperbolic space [15]. Zhou *et al.* [37] also found that genetic variation and their expression can be modeled by a low-dimensional hyperbolic geometry. Interestingly, current study [16] even finds that the features of miniImageNet [38] and CIFAR [39], which are learned in the Euclidean space using neural architectures like VGG [40], Inception [41], and Resnet [21], have apparent hyperbolic properties. Therefore, it is necessary and advantageous to construct hyperbolic deep neural networks (HDNNs) to efficiently deal with far more complex irregular data, and interpret and reason more complex relations beyond Euclidean space, thus producing success in the hyperbolic space [31], [42].

Recently, numerous HDNN architectures [4], [43], [44], [45], [46], [47], [48], [49], [50] are developed to solve a variety of different machine learning tasks. Hyperbolic spaces have been proposed as an alternative continuous approach to learn hierarchical representations for data from textual [51], [52], [53] and graph-structured data [46], [48], to biology [15], [37] and images [16], [54]. The reason for this interest is that the hyperbolic metric approximates the exponential expansion of possible states of the system described by a hierarchical tree-like process. The negative-curvature of the hyperbolic space results in drastically different geometric properties, which makes the circle circumference ($2\sinh r$) and disc area ($2\pi(\cosh r - 1)$) grow exponentially with radius r , as opposed to the Euclidean spaces where they only grow linearly and quadratically. Therefore, hyperbolic spaces have recently gained momentum to model data in the space that exhibits certain desirable geometric hierarchical characteristics.

• W. Peng, T. Varanka, A. Mostafa, H. Shi, and G. Zhao are with the Center for Machine Vision and Signal Analysis, University of Oulu, Oulu, Finland.
E-mail: firstname.lastname@oulu.fi

Manuscript received April 19, 2020; revised August 26, 2020.

To summarize, there are several potential advantages of utilizing hyperbolic deep neural networks to represent data:

- A better generalization capability of the model, with less overfitting, computational complexity, and requirement of training data.
- Reduction in the number of model parameters and embedding dimensions.
- A low distortion embedding, which preserves the local and global geometric information.
- A better model understanding and interpretation, which can provide a conformal mapping to Euclidean space such that it is friendly to down-stream tasks.

Although constructing neural networks in hyperbolic spaces has gained considerable attention recently, this is an extremely challenging task as Euclidean neural operators will not consistently work in the space. Generalizing Euclidean operations, from basic arithmetic operations, *e.g.*, additions and multiplications, to neural operators like convolutions and poolings, to the hyperbolic space is also arduous. To the best of our knowledge, a survey of hyperbolic deep neural networks does not exist in this field. This article makes the first attempt and aims to provide a comprehensive review of the literature around hyperbolic deep neural networks for machine learning tasks. Our goals are to 1) provide a concise context and explanation to enable the reader to become familiar with the basics of hyperbolic geometry, 2) review the current literature related to algorithms and applications about hyperbolic deep learning, and 3) identify open questions and promising future directions.

The article is organized as follows. In Section 2, we introduce the fundamental concepts about hyperbolic geometry, making the paper self-contained. Section 3 introduces the generalization of important Euclidean neural components to the hyperbolic space. We then review the constructions for hyperbolic deep neural networks, including building networks on two commonly used hyperbolic models, Lorentz Model and Poincaré Model in Section 4. In Section 5, we describe applications for testing hyperbolic deep neural networks and discuss the performance of different approaches under different settings. Finally, in Section 6 we identify open problems and possible future research directions.

2 HYPERBOLIC GEOMETRY

2.1 Mathematical preliminaries

Manifold: A manifold \mathcal{M} of dimension n is a topological space of which each point's neighborhood can be locally approximated by the Euclidean space \mathbb{R}^n .

Tangent space: For each point $x \in \mathcal{M}$, the tangent space $\mathcal{T}_x\mathcal{M}$ of \mathcal{M} at x is defined as an n -dimensional vector-space approximating \mathcal{M} around x at a first order.

Riemannian metric: The metric tensor gives a local notion of angle, length of curves, surface area, and volume. For a manifold \mathcal{M} , a Riemannian metric g is a smooth family of inner products on the associated tangent space: $\langle \cdot, \cdot \rangle_x : \mathcal{T}_x\mathcal{M} \times \mathcal{T}_x\mathcal{M} \rightarrow \mathbb{R}$. A given smooth manifold can be equipped with many different Riemannian metrics.

Riemannian manifold: A Riemannian manifold [55] is then defined as a manifold equipped with a group of Riemannian metrics g , which is formulated as a tuple (\mathcal{M}, g) [56].

Geodesics: Geodesics is the generalization of a straight line in the Euclidean space. It is the constant speed curve giving the shortest (straightest) path between pairs of points.

Exponential map: The exponential map takes a vector $v \in \mathcal{T}_x\mathcal{M}$ of a point $x \in \mathcal{M}$ to a point on the manifold \mathcal{M} , *i.e.*, $\text{Exp}_x : \mathcal{T}_x\mathcal{M} \rightarrow \mathcal{M}$ by moving a unit length along the geodesic uniquely defined by $\gamma(0) = x$ with direction $\gamma'(0) = v$. Different manifolds have their own way to define the exponential maps. Generally, this is very useful when computing the gradient, which provides update that the parameter moves along the geodesic emanating from the current parameter position.

Logarithmic map: As the inverse of the aforementioned exponential map, the logarithmic map projects a point $z \in \mathcal{M}$ on the manifold to the tangent space of another point $x \in \mathcal{M}$, which is $\text{Log}_x : \mathcal{M} \rightarrow \mathcal{T}_x\mathcal{M}$. Like the exponential map, there are different logarithmic maps for different manifolds.

Parallel Transport: Parallel Transport $\mathcal{PT}_{u \rightarrow v}$ from vector $u \in \mathcal{M}$ to $v \in \mathcal{M}$ defines the transporting of the local geometry along smooth curves that preserves the metric tensors. It is a map from tangent space $\mathcal{T}_u\mathcal{M}$ to $\mathcal{T}_v\mathcal{M}$ that carries a vector in $\mathcal{T}_u\mathcal{M}$ along the geodesic from u to v .

Gromov δ -hyperbolicity: Gromov δ -hyperbolicity [57], [58] is used to evaluate the hyperbolicity of a dataset/space. Normally, it is defined under four-point condition, say points a, b, c, v . A metric space (X, d) is δ -hyperbolic if there exists a $\delta > 0$ such that these four points in X : $\langle a, b \rangle_v \geq \min\{\langle a, c \rangle_v, \langle b, c \rangle_v\} - \delta$, where the $\langle \cdot, \cdot \rangle_v$ with respect to a third point v is the Gromov product [31] of two points and it is defined as $\langle a, b \rangle_v = \frac{1}{2}(d(a, v) + d(b, v) - d(a, b))$ with $d(\cdot, \cdot)$ as the distance function. For instance, Euclidean space \mathbb{R}^n is not δ -hyperbolic, Poincaré disc (\mathbb{B}^2) is $(\log 3)$ -hyperbolic.

2.2 Isometric Models in the Hyperbolic Space

The five isometric models [59], [60], *i.e.*, the Lorentz (Hyperboloid) model, the Poincaré ball model, the Poincaré half space model, the Klein model, and the Hemishpere model, are the well-known models of hyperbolic space¹. They are embedded sub-manifolds of ambient real vector spaces. We detail the most commonly used three models, *i.e.*, the Lorentz, Klein, and Poincaré models, as illustrated in Fig. 1, for constructing hyperbolic deep neural networks. Please refer to the Appendix for more details about the other models.

2.2.1 Lorentz Model

The Lorentz model \mathbb{L}^n of an n dimensional hyperbolic space is a manifold embedded in the $n+1$ dimensional Minkowski space. The Lorentz model is defined as the upper sheet of a two-sheeted n -dimensional hyperbola with the metric $g^{\mathbb{L}}$, which is

$$\mathbb{L}^n = \{x = (x^0, \dots, x^n) \in \mathbb{R}^{n+1} : \langle x, x \rangle_{\mathbb{L}} = -1, x^0 > 0\}, \quad (1)$$

in which the $\langle \cdot, \cdot \rangle_{\mathbb{L}}$ represents the Lorentzian inner product:

$$\langle x, y \rangle_{\mathbb{L}} = x^T g^{\mathbb{L}} y = -x^0 y^0 + \sum_{i=1}^n x^i y^i, \quad x \text{ and } y \in \mathbb{R}^{n+1}, \quad (2)$$

where $g^{\mathbb{L}}$ is a diagonal matrix with entries of 1s, except for the first element being -1. For any $x \in \mathbb{L}^n$, we can get that $x^0 = \sqrt{1 + \sum_{i=1}^n (x^i)^2}$. The distance in the Lorentz Model is defined as

$$d(x, y) = \text{arcosh}(-\langle x, y \rangle_{\mathbb{L}}). \quad (3)$$

The main advantage of this parameterization model is that it provides an efficient space for Riemannian optimization. An

1. We introduce models with constant sectional curvature of -1. This can be easily generalized to other (negative) curvatures.

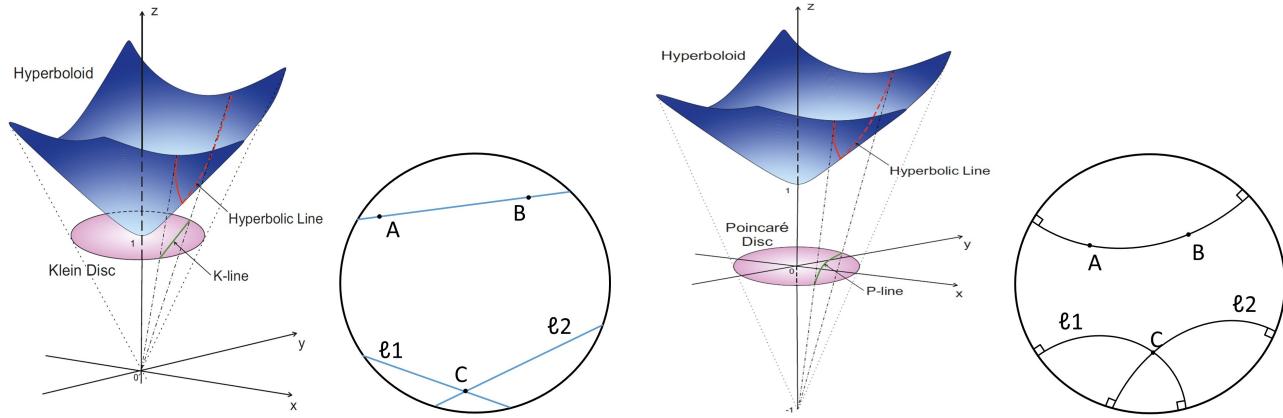


Fig. 1. Illustration of Klein model (left two) and Poincaré model (Right two) in the hyperbolic space. Leftmost: the relationships between Lorentz model and Klein model. We provide the examples of 'straight line' in Klein model (Second from the left). Rightmost: the Poincaré model and its relationship with Lorentz model is provided in the second from the right.

additional advantage is that its distance function avoids numerical instability, when compared to Poincaré model, in which the instability arises from the fraction.

2.2.2 Klein Model

Klein model is also known as the Beltrami–Klein model, named after the Italian mathematician Eugenio Beltrami and German mathematician Felix Klein. The Klein model of hyperbolic space is a subset of \mathbb{R}^n . As illustrated in Fig. 1. It is the isometric image of the Lorentz model under the stereographic projection [60]. The Klein model is obtained by mapping $x \in \mathbb{L}^{n+1}$ to the hyperplane $x^0 = 1$, using rays emanating from the origin. Formally, the Klein model is defined as

$$\mathbb{K}^n = \{x \in \mathbb{R}^n : \|x\| < 1\}. \quad (4)$$

The distance is

$$d(x, y) = \text{arcosh} \left(1 + \frac{1 - \langle x, y \rangle}{\sqrt{(1 - \|x\|^2)(1 - \|y\|^2)}} \right). \quad (5)$$

A straight line, e.g., line \overline{AB} in the second figure from the left of Fig. 1, in Klein model is an intersection of a plane with the disk, thus it is still straight like in Euclidean space. Therefore, Klein model is commonly used to compute the middle point. This model is not conformal to the Euclidean model, which means that angles and circles are distorted,

2.2.3 Poincaré Model

The Poincaré model, as shown in Fig. 1, is given by projecting each point of \mathbb{L}^n onto the hyperplane $x^0 = 0$, using the rays emanating from $(-1, 0, \dots, 0)$. The Poincaré model \mathbb{B} is a manifold equipped with a Riemannian metric $g^{\mathbb{B}}$. This metric is conformal to the Euclidean metric $g^E = I^n$ with the conformal factor $\lambda_x = \frac{2}{1 - \|x\|^2}$, and $g^{\mathbb{B}} = \lambda_x^2 g^E$. Formally, an n dimensional Poincaré unit ball (manifold) is defined as

$$\mathbb{B}^n = \{x \in \mathbb{R}^n : \|x\| < 1\}, \quad (6)$$

where $\|\cdot\|$ denotes the Euclidean norm. The distance between $x, y \in \mathbb{B}^n$ is defined as:

$$d(x, y) = \text{arcosh} \left(1 + 2 \frac{\|x - y\|^2}{(1 - \|x\|^2)(1 - \|y\|^2)} \right). \quad (7)$$

As illustrated in Fig. 2, the left one compares the hyperbolic distance between two points (curves in blue) with that in the

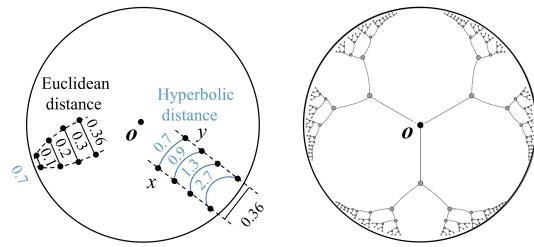


Fig. 2. Illustration of Poincaré Disk (2D Poincaré model), Left: the distances comparison between Euclidean space (in black) and hyperbolic space (in blue). Right: an example of modeling a tree using hyperbolic model.

Euclidean space (lines in black). Like in the left line group, given a constant hyperbolic distance 0.7, the corresponding Euclidean distances decrease dramatically when the points approach the unit boundary. On the contrary, as shown in the right line group, when fix the Euclidean distance to 0.36, the hyperbolic distances grow significantly when the points are closed to the border. The hyperbolic distance $d(x, y) \approx 2\|x - y\|$ when both x and y are closed to the origin with zero norms, which means the model resembles Euclidean geometry near the center. As the points move away from the origin and approach the border, the norms are close to one and the distance grows exponentially [61]. As a comparison, for a regular tree with branching factor b , there are $(b + 1)b^{l+1}$ nodes at level l . The number of nodes in this tree grows exponentially when the tree goes deeper. Therefore, this exponential growing property of hyperbolic space fits very well with the depth increasing in the tree. Besides, when both x and y are approaching the boundary, the distance between x and y can be $d(x, y) \approx d(O, x) + d(O, y)$, which means the shortest path between x and y approaches the path through the origin O . From the definition of Gromov-product [61], we know that such data has a very small δ value of the Gromov δ -hyperbolicity. This is analogous to the property of trees in which the shortest path between two sibling nodes is the path through their parent. Therefore, the Poincaré disk is very suitable for modeling a tree [62], as shown in the rightmost of Fig. 2.

Since they are isometric models in the hyperbolic space, they can be transferred between each other by mapping functions. Fig. 1 provides the visualizations of their relationships. For

Lorentz and Poincaré models, the map can be described as

$$x = (x^0, \dots, x^n) \in \mathbb{L}^n \Leftrightarrow \left(\frac{x^1}{1+x^0}, \dots, \frac{x^n}{1+x^0} \right) \in \mathbb{B}^n. \quad (8)$$

Likewise, the Klein model can be transferred from Lorentz model by the projection

$$x = (x^0, \dots, x^n) \in \mathbb{L}^n \Leftrightarrow \left(\frac{x^1}{x^0}, \dots, \frac{x^n}{x^0} \right) \in \mathbb{K}^n. \quad (9)$$

3 GENERALIZING EUCLIDEAN OPERATIONS TO THE HYPERBOLIC SPACE

Although the use of hyperbolic embeddings (first proposed by Kleinberg *et al.* [63]) in machine learning was already introduced early in 2007, only recently there are methods being extended to deep neural networks. Traditional Euclidean neural networks heavily depend on locality, thus cannot directly be applied to hyperbolic space. Fundamental neural operations like linear projection, average pooling, and feature concatenation, will not work, as the outputs would not necessarily lie in the manifold. Therefore, generalizing Euclidean operations to the hyperbolic space plays a key role in constructing hyperbolic neural networks.

Constructing deep neural networks in the hyperbolic space is not as easy as it is on the Euclidean space. One of the most crucial reasons is that it is the non-trivial or impossible principled generalizations of basic operations, *e.g.*, vector addition, matrix-vector multiplication. Work [43] provides a pioneer study of how classical Euclidean deep learning tools can be generalized to the hyperbolic space. Fueled by this, many current works generalize various deep learning operations as it is the key step towards to hyperbolic deep neural networks. In this section, we review the research literature which is trying to generalize operations, *e.g.*, basic addition, mean and neural network layers, to the hyperbolic space.

3.1 Basic Arithmetic Operations

Basic mathematical operations, like addition and multiplication, are fundamental components of neural networks. They are everywhere in the neural network components, like in convolutional filters, fully connected layers, and activation functions.

One simple way to perform these computations is to approximate them by employing the tangent space. However, as observed in some works [47], [64], the approximation in the tangent space can have a negative impact on the learning process. Specifically, the procedures follow a manifold-tangent manifold scheme, basically by transforming features between hyperbolic spaces and tangent spaces via the logarithmic and exponential maps, in which the logarithmic and exponential maps require a series of hyperbolic and inverse hyperbolic functions. However, the mapping between the manifold and the tangent space is only locally diffeomorphic, which may distort the global structure of the hyperbolic manifold [65], [66]. Besides, the compositions of these functions are complicated and usually range to infinity, significantly weakening the stability of models.

Another good choice is the Gyrovector space [67], which is a generalization of Euclidean vector spaces to models of hyperbolic space based on Möbius transformations. Specifically, for a model \mathbb{B} , the Gyrovector space provides a non-associative algebraic formulation for studying hyperbolic geometry, in analogy to the

way vector spaces are used in Euclidean geometry. In the Gyrovector space, the **Möbius addition** \oplus for x and y in model \mathbb{B} is defined as

$$x \oplus y = \frac{(1 + 2 \langle x, y \rangle + \|y\|^2)x + (1 - \|x\|^2)y}{1 + 2 \langle x, y \rangle + \|x\|^2\|y\|^2}. \quad (10)$$

This is a generalization of the addition in Euclidean space. $x \oplus y$ will recover to $x + y$ when the curvature goes to zero. In addition, the Möbius subtraction \ominus is simply defined as: $x \ominus y = x \oplus (-y)$.

Then the **Möbius scalar multiplication** \otimes is defined as

$$r \otimes x = \begin{cases} \tanh(r \operatorname{artanh}(\|x\|)) \frac{x}{\|x\|}, & x \in \mathbb{B}^n \\ 0, & x = 0, \end{cases} \quad (11)$$

where r is a scalar factor. In fact, all above-mentioned operations can also be conducted in the tangent space by using the exponential and logarithmic maps. As provided by [44], the **Möbius scalar multiplication** can be obtained by projecting x in the tangent space at 0, multiplying this projection by the scalar r in the tangent space. Then projecting it back on the manifold with the exponential map, which means

$$r \otimes x = \operatorname{Exp}_0(r \operatorname{Log}_0(x)). \quad (12)$$

With a similar strategy, the authors derived the **Möbius vector multiplication** $M^\otimes(x)$ between the matrix M and input x , which is defined as

$$M^\otimes(x) = \tanh\left(\frac{\|Mx\|}{\|x\|}\operatorname{actanh}(\|x\|)\right) \frac{Mx}{\|Mx\|}. \quad (13)$$

Based on the Möbius transformations, the authors [43] also derived a closed-form expression of Möbius exponential and logarithmic maps for the Poincaré model. For a vector $v \in T_x \mathcal{M}$ in the tangent space, the exponential map is defined as

$$\operatorname{Exp}_x(v) = x \oplus \left(\tanh\left(\frac{\lambda_x\|v\|}{2}\right) \frac{v}{\|v\|} \right), \quad (14)$$

and as the inverse operation of the exponential map, for a point $y \in \mathbb{B}^n$ on the manifold, the logarithmic map is defined as

$$\operatorname{Log}_x(y) = \frac{2}{\lambda_x} \operatorname{artanh}(\|-x \oplus y\|) \frac{-x \oplus y}{\|-x \oplus y\|}, \quad (15)$$

where λ_x is the conformal factor, as mentioned in Sec. 2.2.3.

3.2 Mean in the Hyperbolic Space

The simple but valuable mean computation is one of the most fundamental operations for machine learning approaches. For example, it can be used to build the batch normalization [68] and average pooling [41] in deep neural networks, as well it can be employed to learn the latent distribution, *e.g.*, in variational Auto-Encoder [69], [70], of data (feature). The weighted average counts much for information aggregation in graph convolutional networks (GCNs) [71]. However, unlike in the Euclidean space, the mean computation cannot be conducted simply by averaging the inputs, which may lead to a result out of the manifold. Basically, the primary approaches to generalize the mean to hyperbolic space are tangent space aggregation [47], Einstein midpoint method [44], and the Fréchet mean method [64].

Tangential aggregations is one of the most straightforward ways to compute the mean in the hyperbolic space. Generally, the mean aggregation μ in Euclidean is defined as a weighted average on the involved neighbors, $\mathcal{N}(i)$, of a center i , which is

$$\mu = \sum_{j \in \mathcal{N}(i)} w_j x_j. \quad (16)$$

However, directly utilizing Eq. (16) in the hyperbolic space is not reasonable since the resulting average may be out of the manifold. Work [47] turns to the tangent space and an attention based aggregation is proposed to compute the aggregated information. Specifically, given the corresponding hyperbolic feature representation, one can compute the attention weights w_j first, then the mean (aggregated information) μ is

$$\mu = \text{Exp}_x \left(\sum_{j \in N(i)} w_j \text{Log}_x(x_j) \right). \quad (17)$$

Work [44] proposes to compute it with **Einstein midpoint**. Einstein midpoint is an extension of the mean operation to hyperbolic spaces, which has the most concise form in the Klein coordinates. The Einstein midpoint of N samples is defined as

$$\mu = \frac{\sum_{i=1}^N \gamma_i x_i}{\sum_{i=1}^N \gamma_i}, \quad (18)$$

in which the x_i is the i -th sample represented using coordinates in Klein model. The $\gamma_i = \frac{1}{\sqrt{1 - ||x_i||^2}}$ are the Lorentz factors. One can easily execute midpoint computations by simply projecting to the Klein model from various models of hyperbolic space since all of them are isomorphic. For instance, from Eqs. (8) and (9), the transition between the Poincaré ($x_{\mathbb{B}}$) and Klein ($x_{\mathbb{K}}$) models can be derived as $x_{\mathbb{K}} = \frac{2x_{\mathbb{B}}}{1 + ||x_{\mathbb{B}}||^2}$. Therefore, based on the Einstein midpoint in Eq. (18), we can get the mean in Poincaré model as

$$\mu_{\mathbb{B}} = \frac{\mu}{1 + \sqrt{1 - ||\mu||^2}} \quad (19)$$

There is also a closed-form expression for Poincaré model to compute the average (midpoint) in the Gyrovector spaces. Work [72] defines a gyromidpoint, which is

$$m(x_{(1)}, \dots, x_{(N)}; \alpha) = \frac{1}{2} \oplus \left(\sum_{i=1}^N \frac{\alpha_i \gamma_i}{\sum_{j=1}^N \alpha_j (\gamma_j - 1)} x_{(i)} \right) \quad (20)$$

with $\alpha = (\alpha_1, \dots, \alpha_N)$ as the weights for each sample $x_{(i)}$ and $\gamma_i = \frac{2}{||x_{(i)}||^2}$.

Recently, using the **Fréchet mean** [73], Luo *et al.* derives a closed-form gradient expression for the mean on Riemannian manifolds [64]. There are considerable years for generalizing Euclidean mean in non-Euclidean geometries [74], using Fréchet mean. However, the Fréchet mean does not have a closed-form solution, and its computation involves the argmin operation that cannot be easily differentiated. Besides, as mentioned by work [64], computing the Fréchet mean relies on some iterative solvers, which is computationally inefficient, numerical unstable thus not friendly to deep neural networks. Therefore, the authors derived an optimization objective for mean (and variance) computation in the hyperbolic space, which is

$$\mu_{fr} = \underset{\mu \in \mathcal{M}}{\text{argmin}} \frac{1}{N} \sum_{i=1}^N d(x_i, \mu)^2, \quad (21)$$

$$\delta_{fr}^2 = \min_{\mu \in \mathcal{M}} \frac{1}{N} \sum_{i=1}^N d(x_i, \mu)^2. \quad (22)$$

However, the general formulation of Fréchet mean requires an argmin operation and offers no closed-form solution, thus both computation and differentiation are problematic. Inspired by work [75], the authors provided its generalization which allows

to differentiate the argmin operation on the manifold. Therefore, they provided their closed-form of the Fréchet mean.

Based on the aforementioned various ways to compute the mean, fundamental neural components, *e.g.*, average pooling layer, batch normalization layer, can be constructed.

3.3 Concatenation and Split Operations

Concatenation and split are commonly used in current deep neural networks, *e.g.*, feature fusion in multimodal learning, operations like GCN filters, and attention mechanism [23]. Simple as they are, they cannot directly be applied in the hyperbolic space as the operations are not manifold preserving ones. For example, points $(0.99, 0)$ and $(0, 0.99)$ are both on the Poincaré disk. However, the resultant point $(0.99, 0, 0, 0.99)$ from concatenation are surely not on a Poincaré model as its norm is larger than one. Therefore, new approaches should be provided for concatenation and split in the hyperbolic space.

As previously described operations, the hyperbolic concatenation and split can be obtained by using the tangent space. Specifically, for an n -dimensional feature embedding $x \in \mathbb{B}^n$ (Note this can be easily generalized to other hyperbolic models) in the hyperbolic space, it can be split into N feature representations V ,

$$V = \{v_1 \in \mathbb{R}^{n_1}, \dots, v_N \in \mathbb{R}^{n_N}\} = \text{Log}_0(x) \quad (23)$$

subject to $\sum_{i=1}^N n_i = n$. Then, the tangent vector can be mapped to the hyperbolic space using the exponential map. Likewise, for N parts feature representation V in the hyperbolic space, the tangent space can also be used to perform concatenation, which is

$$x = \text{Exp}_0(\text{Log}_0(v_1) | \text{Log}_0(v_2), \dots, | \text{Log}_0(v_N)) \quad (24)$$

where $|$ denotes the concatenation operation in the tangent space and v_i represents one feature in the hyperbolic space.

However, as pointed out by work [4], merely splitting the coordinates will lower the norm of the output Gyrovector, which will limit the representational power. Therefore, work [4] proposes a β -split and β -concatenation, as an analogy to the generalization criterion in Euclidean neural networks [76].

The β -split and β -concatenation provided by [4] introduce a scalar coefficient $\beta_n = B(\frac{n}{2}, \frac{1}{n})$, where B is the Beta distribution. With this scalar coefficient, the tangent vectors are scaled before being projected back to the hyperbolic space. Therefore β -split is

$$V = \{\text{Exp}_0(\beta_{n_1} \beta_n^{-1} v_1), \dots, \text{Exp}_0(\beta_{n_N} \beta_n^{-1} v_N)\}. \quad (25)$$

For the β -concatenation,

$$x = \text{Exp}_0(\beta_n \beta_{n_1}^{-1} v_1 | \beta_n \beta_{n_2}^{-1} v_2, \dots, | \beta_n \beta_{n_N}^{-1} v_N). \quad (26)$$

Work [43] presents another way to perform vector concatenation in the hyperbolic space, by introducing linear projection functions based on the Möbius transformations. Specifically, for a set of hyperbolic representations $\{x_1 \in \mathbb{B}^{n_1}, \dots, x_N \in \mathbb{B}^{n_N}\}$, a group of projection functions $\{M_1 \in \mathbb{B}^{n_1 \times n_1}, \dots, M_N \in \mathbb{B}^{n_N \times n_N}\}$ is introduced. Then the concatenated result is :

$$x = M_1 \otimes x_1 \oplus, \dots, \oplus M_N \otimes x_N. \quad (27)$$

However, compared to the previous β methods, this one applies the Möbius transformations (addition and multiplication) many times, as mentioned by [4], which incurs a heavy computational cost and an unbalanced priority in each input sub-Gyrovector.

3.4 Fully-Connected Layers

Fully-connected layer (FC), or linear transform layer, defined as $y = Ax + b$, is also one important component of deep neural networks, in which all inputs from one layer are connected to every activation unit of the next layer. With analogy to Euclidean FC layers, works [4], [43] generalized it to the hyperbolic space. In work [43], fully-connected layers are constructed by Matrix-vector multiplication

$$y = A \otimes x \oplus b = \text{Exp}_0(A \text{Log}_0(x)) \oplus b. \quad (28)$$

Here, the bias translations can be further conducted by Möbius translation, which first maps the bias to the tangent space of origin and then parallel transports it to the tangent space of the addend, finally maps back the result to the manifold, which means

$$x \oplus b = \text{Exp}_0(\mathcal{PT}_{0 \rightarrow x}(\text{Log}_0(b))) = \text{Exp}_x\left(\frac{\lambda_0}{\lambda_x} \text{Log}_0(b)\right), \quad (29)$$

in which λ_x is the conformal factors defined in Sec. 2.2.3. However, work [4] points out that with the Möbius translation, such a surface is no longer a hyperbolic hyperplane. Besides, the shape of the contour surfaces is determined since the norm of each row vector a_k and bias b contribute to the total scale and shift. To deal with this problem, work [4] provides a Poincaré FC layer based on a linear transformation. They argued that the circular reference between a_k and q_{a_k, r_k} can be unraveled by considering the tangent vector at the origin, $z_k \in \mathcal{T}_0 \mathbb{B}^n$, from which a_k is parallel transported. In this way, a new linear transformation is formulated as

$$v_k(x) = 2||z_k|| \text{arsinh}\left(\lambda_x \langle x, \frac{z_k}{||z_k||} \rangle \cosh(2r_k) - (\lambda_x - 1) \sinh(2r_k)\right) \quad (30)$$

where z_k is the generalization of the parameter a_k in A. Based on this, they provided their FC layer, which is

$$Y = \frac{w}{1 + \sqrt{1 + ||w||^2}}, \quad (31)$$

where $w = (\sinh(v_k(x)))$. In this way, we avoid using the tangent space to approximate the hyperplane such that can more properly make use of the hyperbolic nature.

3.5 Convolutional Neural Network Operations

There is limited research about convolutional layers in this space. Work [16] proposes to address the common computer vision tasks, e.g., image classification and person re-identification, using hyperbolic geometry. However, only the decision hyperplanes are established in the hyperbolic space. Thus, the authors did not generalize CNN to hyperbolic space.

Basically, the generalization of a CNN can also be simply conducted by using the tangent space. However, as pointed out by [46], stacking multiple CNNs in the tangent space may collapse to a vanilla Euclidean CNN. This is because the exponential map at k -th layer would have been cancelled by the logarithmic map at next layer. This can be avoided by either applying activation function after the exponential map or adding bias b using the parallel transport.

However, the advantages of hyperbolic geometry may not be well adapted if only using the tangent space. Work [4] provides a novel method to bridge this gap. By using the β -concatenation and the Poincaré fully-connected (FC) layer, the authors presented

a method to build the convolutional layer. In particular, given a C-channel input tensor $x \in \mathbb{B}^{C \times W \times H}$ on the Poincaré ball, for each of the $W \times H$ feature pixels, the representations in the reception field of a convolutional filter with size K are concatenated into a single vector $x \in \mathbb{B}^{nK}$, using the β -concatenation. Then naturally, a Poincaré FC layer, which will be detailed in Sec. 3.4, can be employed to transfer the feature on the manifold. Let C' be the output channels of the CNN layer, then there will be C' groups of such transformations.

3.6 Recurrent Neural Network Operations

Recurrent neural networks (RNNs) [77], [78] are commonly used in sequence learning tasks. Formally, a RNN can be defined by

$$h_{t+1} = \delta(Wh_t + Ux_t + b), \quad (32)$$

where h_{t+1} is the hidden state of the next step, which is updated using the current hidden state h_t and input x_t . δ is a non-linear function. W and U are learnable parameters, and b is the corresponding bias. Work [43] generalizes the RNN to the hyperbolic space, leveraging the Möbius operations in Gyrovector space. The RNN in the hyperbolic space can be defined by

$$h_{t+1} = \delta^\otimes(W \otimes h_t + U \otimes x_t \oplus b), \quad (33)$$

where \oplus and \otimes are the generalization of original $+$ and \times in Gyrovector space, as defined in Section 2. The authors also extended the same idea into the gated recurrent unit (GRU) architecture [79], with the same strategy.

Existing works are limited to the Poincaré model with the corresponding operations defined in the Gyrovector space. However, these kinds of operations are always costly when compared to the Euclidean counterparts. Future works can explore more efficient ways and extend to other hyperbolic models.

3.7 Activation function

Activation function provides a non-linear projection of the feed-in features such that much richer semantic representation can be learned. The expected activation function for hyperbolic model should be the non-linear function which well preserve the manifold. Fortunately, for some hyperbolic models, e.g., Poincaré and Klein models, the manifold is defined only by the norm constraints. Therefore, the activation function is manifold preserving once it is a norm decreasing operation. Hence, any pointwise Euclidean activation functions which do not increase the norm can be directly applied to these models. For example, Liu *et al.* [46] directly applied the norm decreasing ReLU [80] and leaky ReLU [81] activation functions in the Poincaré model. However, manifold preserving activation functions are different for different manifolds. For Lorentz model, as the origin in Poincaré model is the pole vector in Lorentz model, these activation functions cannot be applied as they will depart the point away from the manifold. As Lorentz model is isometric to Poincaré model, work [46] maps the point from Lorentz model to Poincaré model, conducts the above mentioned Euclidean activation functions, and then maps feature back to the original model.

Ganea *et al.* proposed a Möbius version of projection function [43], which can also be employed to activation functions. In particular, for a function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$, its Möbius version f^\otimes is

$$f^\otimes(x) = \text{Exp}_0(f(\text{Log}_0(x))). \quad (34)$$

The authors utilized the tangent space of the origin to perform the function f . A hyperbolic activation function can be also realized

by this way, in which case, the input dimension n is equal to the output dimension m . Following the same idea, work [47] provides a similar activation function for graph convolutional networks. The only difference is that they considered the curvatures of different layers. Thus, the logarithmic map and exponential map are defined at the point origins in the manifold with different curvatures.

Interestingly, work [4] removes the activation functions, since the authors thought that the operation on the manifold itself is non-linear, which obviates the need for activation functions.

3.8 Batch Normalization

Batch Normalization (BN) [68] limits the internal covariate shift by normalizing the activation of each layer. It is commonly used to speed up the training procedure of neural networks, as well as to make the training process more stable. The basic idea behind is normalizing of the feature representations by re-centering and re-scaling. Specifically, given a batch of m data-points, The BN algorithm will first compute the mean μ of this batch. Based on μ , the mini-batch variance σ is also computed. Then, two learnable parameters are introduced, which are the scale parameter γ and the shift parameter β . The input activation x is then re-centered and re-scaled, which is $y = \gamma x + \beta$. Theoretically, this BN operation can be easily generalized to the manifold via transferring to the tangent space. Work [64] provides an alternative based on Fréchet mean [73]. In particular, the authors formulated the Riemannian extension of the standard Euclidean Batch Normalization by a differentiable Fréchet mean, as described in Sec. 3.2.

3.9 Classifiers and Multiclass Logistic Regression

Classifier is an essential component for classification tasks. In the context of deep learning, multiclass logistic regression (MLR) or softmax regression is commonly used to perform multi-class classification in Euclidean space. Formally, given K classes, MLR is introduced to predict the probabilities of each class $k \in \{1, 2, 3, \dots, K\}$ based on the input representation x ,

$$p(y = k, x) \propto \exp(\langle a_k, x \rangle - b) \quad (35)$$

where a_k denotes the normal vector and $b \in \mathbb{R}$ is the scalar shift. Then, the decision hyperplane determined by $a \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ and b is defined by $H_{a,b} = \{x \in \mathbb{R}^n : \langle a, x \rangle - b = 0\}$. Note that \exp is the exponential function, not the manifold map function Exp . According to [82], the MLR can be reformulated as

$$p(y = k, x) \propto \exp(\text{sign}(\langle a_k, x \rangle - b) \|a_k\| d(x, H_{a_k, b_k})) \quad (36)$$

where $d(x, H_{a_k, b_k})$ is the Euclidean distance between x and the hyperplane H_{a_k, b_k} . To further generalize it to the hyperbolic space, work [43] re-parameterizes the scalar term $b \in \mathbb{R}$ with a new set of parameters $p_k \in \mathbb{R}^n$, by which they reformulated the hyperplane: $\hat{H}_{a,p} = \{x \in \mathbb{R}^n : \langle a, x - p \rangle = 0\}$, and $\hat{H}_{a,p} = \hat{H}_{a, \langle a, p \rangle}$. In this way, the MLR is rewritten as

$$p(y = k, x) \propto \exp(\text{sign}(\langle a_k, x - p_k \rangle) \|a_k\| d(x, \hat{H}_{a_k, p_k})). \quad (37)$$

Then, the definition of the hyperbolic setting is simply achieved by replacing the addition $+$ with Möbius addition \oplus .

However, this causes an undesirable increase in the parameters from $n+1$ to $2n$ in each class k . As pointed out by [4], there is no need to introduce countless choices of p_k to determine the same discriminative hyperplane. Instead, they introduced another scalar parameter $r_k \in \mathbb{R}$ such that the bias vector $q_{a_k, r_k} = r_k \frac{a_k}{\|a_k\|}$ parallels to the normal vector a_k . That is

$$\hat{H}_{a_k, r_k} = \{x \in \mathbb{R}^n | \langle a_k, -q_{a_k, r_k} + x \rangle = 0\} = H_{a_k, r_k \|a_k\|}. \quad (38)$$

Based on this, the MLR is reformulated as

$$p(y = k, x) \propto \exp(\text{sign}(\langle a_k, -q_{a_k, r_k} + x \rangle) \|a_k\| d(x, \hat{H}_{a_k, r_k})). \quad (39)$$

In addition to the MLRs in the hyperbolic space, there are also works constructing classifiers in this space. Work [83] introduces a hyperbolic formulation of support vector machine classifier (HSVM). Work [84] provides theoretical guarantees for learning a SVM in the hyperbolic space. Besides, an efficient strategy is introduced to learn a large-margin hyperplane, by the injection of adversarial examples. Please refer to Appendix for more details.

3.10 Optimization

Optimizer plays a crucial role in training deep neural networks. It largely influences the convergence of the training process, the training speed, and the final predictive performance. Therefore, generalizing optimizers to hyperbolic space is as important as constructing hyperbolic neural architectures [85]. In terms of stochastic optimizers on Riemannian manifolds, one pioneer work should be the Riemannian stochastic gradient descent (RSGD), provided by Bonnabel *et al.* [86]. They pointed out that for the standard stochastic gradient descent in \mathbb{R}^n , seeking the matrix with certain rank, which best approximates the updated matrix, can be numerically costly, especially for very large parameter matrix. To enforce the rank constraint, a more natural way is to endow the parameter space with a Riemannian metric and perform a gradient step within the manifold of fixed-rank matrices. To address this issue, they proposed to replace the usual update in SGD using an exponential map (Exp) with the following update

$$W_{t+1} = \text{Exp}_{W_t}(-\alpha g_t) \quad (40)$$

where $g_t \in \mathcal{T}_{W_t} \mathcal{M}$ denotes the Riemannian gradient of the objective at point W_t . The provided algorithm is completely intrinsic, which does not depend on a specific embedding of the manifold or on the choice of local coordinates. For manifolds where Exp function is not known in a closed form, it is common to replace it with a first order approximation [87].

Current deep learning optimization approaches prefer adaptive strategies, *e.g.*, Adagrad [88] and Adam [89], which dynamically incorporates knowledge of the geometry of the data observed in earlier iterations to perform more informative gradient-based learning. Although the input data is with very high dimension, useful features have very different frequencies. When the gradient vectors are sparse, the update of the neural network can often be performed in time proportional to the support of the gradient. One of the current challenges of generalizing the adaptivity of these optimization methods to hyperbolic space is that the manifold does not provide an intrinsic coordinate system, while Riemannian manifold only allows to work in a certain local coordinate systems. Therefore, it is non-trivial to extend the optimizers to hyperbolic space in an intrinsic manner (coordinate-free). Suggested by [90], one solution can be fixing a canonical coordinate system in the tangent space and then parallelly transporting it along the optimization trajectory. However, general Riemannian manifold depends on both the chosen path and the curvature, which will give the parallel transport a rotational component (holonomy). This will break the gradient sparsity and thus harm the benefit of adaptivity. For instance, imagining the vector x is initially sparse, *e.g.*, $x = (2.5, 0, 0)$, with the parallel transport, the vector may be rotated and thus has other components $x = (1.2, 0.4, 0.1)$, which definitely brakes the coordinate-wise updates and sparsity.

To avoid such problem, work [90] represent a n -dimensional manifold \mathcal{M} by a Cartesian product of n manifolds, which means $\mathcal{M} = \mathcal{M}_1 \times \dots \times \mathcal{M}_n$. Based on this, the authors provided the Riemannian Adagrad, which is defined by

$$W_{t+1} = \text{Exp}_{W_t}(-\alpha g_t) \quad (41)$$

$$W_{t+1}^i = \text{Exp}_{W_t^i}\left(-\frac{\alpha g_t^i}{\sqrt{\sum_{k=1}^t \|g_k^i\|_{x_k^i}^2}}\right), \quad (42)$$

where the $\|g_k^i\|_{x_k^i}^2 = \mathbf{g}(g_k^i, g_k^i)$ is a Riemannian norm. They further extended it to Riemannian Adam by introducing the momentum term and an adaptive term.

4 HYPERBOLIC DEEP NEURAL NETWORKS

Overwhelming number of studies using deep neural networks [17] ranging from convolutional neural networks (CNN) [91], recurrent neural networks (RNN) [78], [92], graph neural networks (GNN) [71], [93] to generative models like variational auto-encoder (VAE) [69] have showcased the superiority of various deep neural networks in different research tasks. With the success of all these neural architectures in Euclidean space, we have the reason to expect their generalization and performance in the hyperbolic space. Fortunately, we do observe a considerable number of research extending such Euclidean architectures to the hyperbolic space, mainly using Poincaré and Lorentz models. This section describes various hyperbolic deep neural architectures. We try our best to collect and summarize all advanced hyperbolic machine learning methods, as illustrated in Table 1. One can find that most of the works are from the prominent conferences/journals, e.g., Nature, Cell, NeurIPS, ICML, and ICLR. Besides, increasing institutions are getting into this potential research field. From the method perspective, the Poincaré model and Lorentz model are prominent hyperbolic models for generalizing neural networks. At the same time, the former one (Poincaré model) is dominated in the hyperbolic deep neural networks. In the following part, we detail different kinds of hyperbolic architectures, including hyperbolic embedding, hyperbolic clustering, hyperbolic attention networks, hyperbolic graph neural networks, hyperbolic normalizing flow, hyperbolic variational auto-encoder, and hyperbolic neural networks with mixed geometries.

4.1 Hyperbolic Embeddings

Embedding [94] data into a standard space makes it possible to use properties of the target space as additional structure in the original dataset, and brings to front information that is hard to detect in the raw input. Embedding has been found relevant to a variety of subjects such as data visualization, network analysis, routing, localization, machine learning, statistics, biology and many others.

There have been many works [63], [95], [96], [97], [98], [99], [100], [101] considering hyperbolic space as an alternative for various embedding tasks. Walter [95] provided a construction of a distance preserving embedding of high-dimensional data into the hyperbolic space for interactive visualization. In Internet routing, Kleinberg [63] presented a constructive proof that every finite, connected, undirected graph has a greedy embedding in two-dimensional hyperbolic space, thus introducing hyperbolic geometry for greedy routing in geographic communication networks. Later, Cvetkovski *et al.* [100] extended it to dynamic graphs, i.e., communication networks whose topology changes over time. Similarly, Shavitt *et al.* [96] embedded the Internet distance metric

in a hyperbolic space and by carefully selecting the curvature, they improved the accuracy of Internet distance embedding. Boguñá *et al.* [5] resolved the serious scaling limitations of Internet by the embeddings of the AS (autonomous system) Internet topology in the hyperbolic space to perform greedy shortest path routing. For complex networks, Krioukov [99] showed that heterogeneous degree distributions and strong clustering can emerge by assuming an underlying hyperbolic geometry, thus developed a geometric framework to model complex networks using hyperbolic space. Bläsius *et al.* [101] constructed and implemented a new maximum likelihood estimation algorithm that embeds scale-free graphs in the hyperbolic space. Coalescent embedding [102] was the first model-free unsupervised kernel learning based solution for inferring the graph embedding in the Poincaré disk.

Currently, Nickel *et al.* [103] proposed to learn a Poincaré embedding for symbolic data, while considering the latent hierarchical structures. This work proves that Poincaré embeddings can outperform Euclidean embeddings significantly on data with latent hierarchies, in terms of both representation capacity and generalization ability. Following this study, many new embedding methods in the hyperbolic space are proposed. They can be roughly summarized into four categories, *i.e.*, tangent optimized, fully Riemannian optimization, numerical stable embedding, and combinatorially tree embedding (in Appendix).

Specifically, the Poincaré embedding [103] is trying to find embeddings $\Theta = \{\theta_i\}_{i=1}^n$, where $\theta \in \mathbb{B}^d, \|\theta_i\| < 1$ in the unit d -dimensional Poincaré ball for a set of symbols with size of n . Therefore, the optimization problem can be framed as $\Theta^* = \operatorname{argmin}_{\Theta} \mathcal{L}(\Theta)$, where $\mathcal{L}(\cdot)$ is a task-related loss function. For instance, in the hierarchy embedding task, the loss function over the entire dataset \mathcal{D} can be represented as

$$\mathcal{L}(\Theta) = \sum_{(u,v) \in \mathcal{D}} \log \frac{e^{-d_{\mathbb{B}}(u,v)}}{\sum_{v' \in \mathcal{N}(u)} e^{-d_{\mathbb{B}}(u,v')}}, \quad (43)$$

where $\mathcal{N}(u)$ denotes a set of negative examples for u . This loss function encourages related objects to be closer to each other than objects without an obvious relationship. The embedding is further optimized utilizing the RSGD with the exponential map, which is a scaled version of the Euclidean gradient. Specifically, the retraction operation $\mathcal{R}_w(v) = w + v$ is utilized as the approximation of the exponential map, Exp . Then, a projection, which normalizes the embeddings with a norm bigger than one, is utilized to ensure the embeddings remain within the Poincaré model. Thus,

$$W_{t+1} = \text{proj}\left(W_t - \eta \frac{(1 - \|W_t\|^2)^2}{4} \nabla_E\right), \quad (44)$$

where ∇_E is the Euclidean gradient, and η is the learning rate. The normalization function is $\text{proj}(x) = x/\|x\|$ if $\|x\| \geq 1$, otherwise $\text{proj}(x) = x$. Similarly, work [52] proposes text and sentence embedding with Poincaré model. The authors proposed to re-parametrize the Poincaré embeddings such that the projection step is not required. On top of the existing encoder architectures, *e.g.*, LSTMs, a reparameterization technique is introduced to map the output of the encoder to the Poincaré ball, which is defined as

$$\Theta = \delta(\phi_{norm}(h)) \frac{\phi_{dir}(h)}{\|\phi_{dir}(h)\|}, \quad (45)$$

where h represents the hidden representation of encoder, and δ denotes the sigmoid function. The function $\phi_{dir} : \mathbb{R}^d \rightarrow \mathbb{R}^d$ is used to compute a direction vector. Function $\phi_{norm} : \mathbb{R}^d \rightarrow \mathbb{R}$ is a norm function. In this way, the authors mapped the encoder

embeddings to the Poincaré ball and the Adam is introduced to optimize the parameters of the encoder.

Work [43] presents hyperbolic entailment cones to especially deal with more complicated connections like asymmetric relations in directed acyclic graphs, thus further improving the performance of hyperbolic embeddings. The authors pointed out that most of the embedding points of the method [103] collapse on the boundary of Poincaré ball. To address this issue, they generalized the idea of order embeddings [104], which views hierarchical relations as partial orders defined on a family of nested geodesically convex cones, into the hyperbolic space.

Specifically, in a vector space, a convex cone S is a set that is closed under non-negative linear combinations, which means for vectors $v_1, v_2 \in S$, then $\alpha v_1 + \beta v_2 \in S, \forall \alpha, \beta > 0$. Since the cones are defined in the vector space, the authors proposed to build the hyperbolic cones using the exponential map, which leads to the definition of S -cone at point x , which is

$$\mathfrak{S}_x = \text{Exp}_x(S), S \in \mathcal{T}_x \mathcal{M}. \quad (46)$$

To avoid heavy cone intersections and scale exponentially with the space dimension, the authors further constructed the angular cones in the Poincaré ball. To achieve this, they introduced the so-called cone aperture functions $\phi(x)$ such that the angular cones $\mathfrak{S}_x^{\phi(x)}$ follow four intuitive properties, including axial symmetry, rotation invariant, continuous of cone aperture function, and the transitivity of nested angular cones. With all of this the authors provided a closed-form of the angular cones, which are

$$\mathfrak{S}_x^{\phi(x)} = \left\{ (\pi - \angle Oxy) \leq \arcsin \left(K \frac{1 - \|x\|^2}{\|x\|} \right) \right\}, \quad (47)$$

where O in angle $\angle Oxy$ is the origin, K is a constant. They optimize the objective using **fully Riemannian optimization**, instead of using the first-order approximation as work [103] did.

In addition, Tifrea *et al.* [51] extends an unsupervised word embedding algorithm, Glove [105], to the Riemannian manifolds. They proposed to embed words in a Cartesian product of hyperbolic spaces which they theoretically connected to the Gaussian word embeddings and their Fisher geometry.

Instead of providing new embedding methods in the hyperbolic space, work [106] addresses the **numerical instability** issue of the networks. Like mentioned in work [106], the difficulty is caused by floating-point computation and amplified by the ill-conditioned Riemannian metrics. As points move far away from the origin, the error caused by using floating-point numbers to represent them will be unbounded. For the Poincaré model, the distance changes rapidly when the points are close to the ball boundary such that it is not well conditioned. While for Lorentz model, it is not bounded such that it will experience large numerical error when the points are far away from the origin. Therefore, for the representation in the hyperbolic space, it is desirable to find a method that can represent any point with small fixed bounded error in an efficient way. To this end, work [106] presents a tiling-based model to utilize the integer-lattice square tiling (or tessellation) [107] in the hyperbolic space to construct a constant-error representation. They proved that the representation error, the error of computing distances, and the error of computing the gradient are bounded by a fixed value that is independent of distance to the origin.

4.2 Hyperbolic Cluster Learning

Hierarchical Clustering (HC) [136], which generally constructs a hierarchy over clusters with the form of a multi-layered tree whose

leaves correspond to samples and internal nodes correspond to clusters, is a fundamental problem in data analysis, visualization, and mining the underlying relationships. Mainstream methods including bottom-up linkage methods [137], and recently, cost function based methods [138]. However, these methods are either not amenable for stochastic gradient methods or computationally expensive. Gradient-based hyperbolic hierarchical clustering, gHHC [121], a geometric heuristic to provide an approximate distribution over lowest common ancestor (LCA), is proposed over continuous representations of tree in the hyperbolic space (Poincaré model), based on the observation that child-parent relationships can be modelled by the **distances** and **norms** of the embedded node representations. The authors used the norm of vectors to model depth in the tree, requiring child nodes to have a larger norm than their parents. The root is near the origin of the space and leaves near the edge of the ball. Formally, let $Z = \{z_1, z_2, \dots, z_k\}, z_i \in \mathbb{B}^d$ represent the node representation in the d -dimensional Poincaré ball. Then a child-parent dissimilarity function is used to encourage the children have a smaller norm than the parent, which is

$$d_{cp}(\mathcal{T}_c, \mathcal{T}_p) = d_{\mathbb{B}}(z_c, z_p)(1 + \max(\|z_p\| - \|z_c\|, 0)), \quad (48)$$

where z_c, z_p are the hyperbolic embedding of nodes $\mathcal{T}_c, \mathcal{T}_p$. If the norm of the parent node is smaller than the child, then the dissimilarity will just be the distance in the hyperbolic space. Otherwise, this dissimilarity will be bigger than the distance. Then this dissimilarity function is used to model a distribution over the tree structure to encode the uncertainty, which is $P_{par}(\mathcal{T}_p|\mathcal{T}_c, Z) \propto \exp(-d_{cp}(\mathcal{T}_c, \mathcal{T}_p))$, thus the tree distribution over embedding will be

$$P_{par}(\mathcal{T}|Z) \propto \prod_{\mathcal{T}_p} \prod_{\mathcal{T}_c \in \text{children}(\mathcal{T}_p)} P_{par}(\mathcal{T}_p|\mathcal{T}_c, Z). \quad (49)$$

In contrast with previous gradient-based approaches, this approach has theoretical guarantees in terms of clustering quality and empirically outperforms agglomerative heuristics.

4.3 Hyperbolic Attention Network

Currently, attention mechanism [23] for various neural networks becomes one of the most attractive research topics. Outstanding architectures include the neural Transformer [23], BERT [139], and even the graph attention networks [122], [140], [141]. While attention mechanisms have become the de-facto standard for NLP tasks, their momentum has continuously been extended to computer vision applications [142]. At its core lies the strategy of focusing on the most relevant parts of the input to make decisions. Different from Euclidean space, the distances defined in the hyperbolic space highly depend on their locations (*e.g.*, close to the origin or to the boundary) thus they have their own ways to measure the similarity/dissimilarity, when computing the attention. Therefore, it is important to design hyperbolic attention network such that the correlations can be captured reasonably in terms of the hyperbolic topology and semantic representations.

Work [44] extends it into the hyperbolic space (Lorentz model), utilizing the Lorentzian distance and the Einstein midpoint method to conduct the score matching and aggregation. First, the data representation is organized by a pseudo-polar coordinate, in which an activation $x \in \mathbb{R}^{n+1} = (\mathbf{d}, r)$ is constructed by a n dimensional normalized angle $\mathbf{d}, \|\mathbf{d}\| = 1$ and a scalar radius r . Then, a well-developed map function is introduced to project the activation to the hyperbolic space, which is $\pi((\mathbf{d}, r)) =$

TABLE 1

Summary of the advanced machine learning methods in the hyperbolic space. Here, G. means the type of Geometry. Its value 'Mixed' means the methods combine more than one different geometries, i.e., Euclidean, elliptic, and hyperbolic geometries, to build the model.

Method	Year	Architecture	Tasks	G.	Institution	Source
PEmbedding [103]	2017	Embedding	NLP and Graph	B	Facebook	NeurIPS
Coalescent [102]	2017	Embedding	Graph	B	TU-Dresden	Nature
HyperGraph [108]	2017	Embedding	Network Vertex Classification	B	ICL	MLG W
TextHyper [52]	2018	AE(GRU)	Text Embedding	B	Google	NAACL W
h-MDS [109]	2018	Embedding	Tree and tree-like data modeling	B	Stanford	PMLR
HyperQA [110]	2018	Encoder-Decoder	Neural Question Answering	B	NTU	WSDM
HyperConv [111]	2018	Embedding	NLP and Graph	B	ETH Zürich	ICML
HNN [43]	2018	RNN	NLP(textual entailment and noisy prefixes)	B	ETH Zürich	NeurIPS
Lorentz [45]	2018	Embedding	Taxonomies Embedding, Graph and Historical Linguistics	L	Facebook	ICML
HyperBPR [112]	2018	BPR [113]	Recommender Systems	B	NTU	AAAI
ProductM [114]	2018	Embedding	Tree(with Cycle)	Mixed	Stanford	ICLR
HAN [44]	2019	Attention Module	Graph, Machine translation and Relational Modeling	L(K)	Deepmind	ICLR
HGCN [47]	2019	GNN	Graph	L	Stanford	NeurIPS
PGlove [51]	2019	Glove [105]	Word Embedding	Mixed	ETH Zürich	ICLR
HGN [46]	2019	GNN	Graph	B,L	Facebook	NeurIPS
Tiling [106]	2019	Tiling	NLP and Compressing	L	Cornell	NeurIPS
PTaxo [115]	2019	Embedding	Taxonomy Induction	B	U. Hamburg	ACL
H-SVM [83]	2019	SVM	NLP(Word E) Graph (Node C)	B	MIT	AISTATS
H-Recom [12]	2019	BPR	Recommender Systems	L	ASOS	CoRR
MuRP [116]	2019	Bilinear	knowledge Graph	B	U. Edinburgh	NeurIPS
RAO [90]	2019	Optimizer	NLP	B	ETH Zürich	ICLR
WrapN [117]	2019	VAE	NLP,MNIST and Atari trajectories	B	U.Tokyo	ICML
LDistance [118]	2019	Embedding	NLP	B	U. Toronto	ICML
PVAE [119]	2019	VAE	NLP, and MNIST	B	Oxford	NeurIPS
CCM-AAE [120]	2019	AE	MNIST C, Graph	Mixed	U.Lugano	ASOC
HWAE [54]	2019	VAE	G MNIST, Graph and Tree	B	ETH Zürich	-
gHHC [121]	2019	Clustering	Clustering ImagNet, Multi-Task	B	U. Mass	KDD
HGAT [122]	2020	GNN	Graph	B	BUPT	AAAI
H-STGCN [48]	2020	GNN	Skeleton Action	B	U. Oulu	ACM MM
HyperKG [14]	2020	Translational	Knowledge Graph	B	EPFL	ESWC
HyperML [123]	2020	Metric Learning	Recommender Systems	B	NTU	WSDM
LorentzFM [124]	2020	Triangle Inequality	Recommender Systems	B	eBay	AAAI
APo-VAE [53]	2020	VAE	NLP (dialog-response generation)	B	Duke	-
H-Image [16]	2020	Embedding	Image c, few-shot	B	SIST	CVPR
HyperText [125]	2020	RNN	NLP(text classification)	B	Huawei	EMNLP
k-GCN [49]	2020	GNN	Graph	Mixed	ETH Zürich	ICML
FMean [64]	2020	GNN	Graph	B	Cornell	ICML
H-NormF [126]	2020	Norm. Flow	Graph	B	McGill	ICML
k-Stereographic [127]	2020	GNN	Graph	Mixed	SIST	ICML
L-Group [128]	2020	Group	jet physics	L	U. Chicago	ICML
MVAE [129]	2020	VAE	Image reconstruction and Tree systhesis	Mixed	ETH Zürich	ICLR
HypHC [130]	2020	Clustering	Clustering(e.g., CIFAR-100)	B	Stanford	NeurIPS
R-NormF [131]	2020	Norm. Flow	Earth sciences	B	Oxford	NeurIPS
RH-SVM [84]	2020	SVM	ImageNet(Pick 2classes)	L	Princeton	NeurIPS
TREEREP [132]	2020	Embedding	Tree	B	U. Michigan	NeurIPS
UltraH [133]	2020	Embedding	Graph	Mixed	NVIDIA	NeurIPS
GIL [134]	2020	GNN	Graph	B	CAS	NeurIPS
P-Maps [15]	2020	Embedding	Biology	B	Facebook	Nature
HNN++ [4]	2021	CNN,Transformer	NLP, Clustering, Machine Translation	B	U. Tokyo	ICLR
Hyper-gene [37]	2021	Embedding	Biology	B	Salk	Cell
scPhere [135]	2021	VAE	Biology	L	MIT and Harvard	Nature

$(\sinh(r)\mathbf{d}, \cosh(r))$. It is easy to see that the projected point lies in the Lorentz model. Then for hyperbolic matching, the authors took $\alpha(q_i, k_j) = f(-\beta d_{\mathcal{L}(q_i, k_j)} - c)$, in which the negative Lorentzian distance (scaled by $-\beta$ and shifted by c) is utilized to measure the correlation (matching score α). Since there is no natural definition of mean on the manifold, they turn to Einstein midpoint to conduct hyperbolic aggregation. Specifically, the aggregated message m_i can be represented as

$$m_i(\{\alpha_{ij}\}_j, \{v_{ij}\}_j) = \sum_j \left[\frac{\alpha_{ij} \gamma(v_{ij})}{\sum_l \alpha_{il} \gamma(v_{il})} \right] v_{ij}, \quad (50)$$

where $\gamma(v_{ij})$ is the Lorentz factor at point v_{ij} , and v_{ij} is defined on the Klein model. On the top of the proposed hyperbolic attention network, the authors further formulated the Hyperbolic Transformer model, which is proved to have the superiority when compared to the Euclidean Transformer.

Based on the Euclidean graph attention network (GAT) [140], work [122] generalizes it to a hyperbolic GAT. The idea is very simple, just replacing the Euclidean operation with Möbius operations, which means the matching score

$$\alpha_{ij} = f(W \otimes h_i, W \otimes h_j). \quad (51)$$

They further defined the f function based on the hyperbolic distance. Just as work [44], the negative of the distance between nodes is utilized as the matching score. The scores are finally normalized using softmax function, otherwise all the scores are negative. The hyperbolic aggregation is simply conducted on the tangent space, as it is done in Euclidean space.

Shimizu *et al.* [4] proves that three different kinds of hyperbolic centroids, including the Möbius gyromidpoint [72], Einstein midpoint [72] and the centroid of the squared Lorentzian distance [118], are the same midpoint operations projected on each manifold and exactly matches each other. Based on this observation, they explored on Möbius gyromidpoint and generalized

it by extending to the entire real value weights (previously, it is defined under the condition of non-negative weights) by regarding a negative weight as an additive inverse operation. So the centroid with real weights $\{v_i \in \mathbb{R}\}_{i=1}^N$ is

$$\mu = [x_i, v_i] = \frac{1}{2} \oplus \left(\frac{\sum_{i=1}^N v_i \lambda x_i}{\sum_{i=1}^N \|v_i\| \lambda x_i} \right). \quad (52)$$

With the above weights, the authors computed the attention in the hyperbolic space. Given the source and target as sequences of Gyrovector, first, the proposed Poincaré FC layers (in Sec. 3.4) are utilized to construct the queries, keys, and values. Then, in order to build the multi-head attention, they are broken down into several parts. Like previous methods, the negative distances are also employed to measure the matching scores. Finally, the message from multi-head is aggregated using the proposed Poincaré weighted centroid. The authors built a hyperbolic set transformer model and compared to its Euclidean counterpart [143]. The result shows that the hyperbolic one can at least get equivalent performance, at the same time showing a remarkable stability and consistently converges.

4.4 Hyperbolic Graph Neural Network

Recently, there has been a growing passion of modeling graphs in the hyperbolic space [46], [48], [49]. A core reason for that is learning hierarchical representations of graphs is easier in the hyperbolic space due to the curvature and the geometrical properties of the hyperbolic space. Such spaces were shown by Gromov to be well suited to represent tree-like structures [31] with low distortion as objects requiring an exponential number of dimensions in Euclidean space can be represented in a polynomial number of dimensions in the hyperbolic space. As an alternative, graph neural networks (GNNs) [71], [93], [144] are powerful tools for data with non-Euclidean graph structures. However, the operations are still in Euclidean space, which does not make full use of the geometric property. Current studies of GNN in the hyperbolic space [46], [47], [48] show a superiority, in terms of both model compactness and the predictive performance, when compared to their counterparts in the Euclidean space. This suggests the essential of generalizing graph neural networks to hyperbolic space.

GNN can be interpreted as performing message passing between nodes, which can be formulated as

$$h_i^{k+1} = \sigma \left(\sum_{j \in N(i)} A_{ij} W^k h_j^k \right), \quad (53)$$

where h_i^{k+1} represents the hidden representation of the i -th node at the $(k+1)$ -layer, W^k denotes the weight of the network at k layer. The A_{ij} is the entry of the normalized adjacency matrix A . Eq. (53) performs the information aggregation around the neighbor nodes $N(i)$ of node i to update the representation of this node.

Works [46], [48] provide a straightforward way to extend the graph neural network to hyperbolic space, using tangent space. Work [46] utilizes the logarithmic map $\text{Log}_{x'}$ at a chosen point x' , such that the functions with trainable parameters are executed there. Thus, the graph neural operation in a hyperbolic space is

$$h_i^{k+1} = \sigma \left(\text{Exp}_{x'} \left(\sum_{j \in N(i)} A_{ij} W^k \text{Log}_{x'}(h_j^k) \right) \right), \quad (54)$$

where an exponential map $\text{Exp}_{x'}$ is applied afterwards to map the learned feature back to the manifold. The authors moved the

activation function σ to the tangent space, since they suggested that otherwise the hyperbolic operation will collapse to the vanilla Euclidean GCN as the exponential map will be canceled by the logarithmic map at next layer.

Work [48] shares the same idea to construct spatial temporal graph convolutional networks [145] in the hyperbolic space and applied this for dynamic graph sequences. They further explored the projection dimension in the tangent space, using neural architecture search (NAS) [146]. Chami *et al.* [47] decoupled the message passing procedure of GCN before generalization. The operation of GCN is divided into three parts, including feature transform, neighborhood aggregation, and activating by the activation function. Then, the authors provided operations for corresponding parts in the hyperbolic space.

Bachmann *et al.* [49] presented a novel κ -GCN in the hyperbolic space, which is a mathematically grounded generalization of GCN to constant curvature spaces. Specifically, they extended the operations in Gyrovector space to the space with constant positive curvature, which is the stereographic spherical projection models in their study. By this way, they provided a uniform GCN for spaces with different kinds of curvature (0, negative, and positive), which is called the κ -GCN. Work [127] further improves this method by providing a more reasonable definition of the gradient of curvature at zero, since the original one is incomplete.

4.5 Hyperbolic Normalizing Flows

Normalizing flows [147] involve learning a series of invertible transformations, which are used to transform a sample from a simple base distribution to a sample from a richer distribution. The models produce tractable distributions where both sampling and density evaluation can be efficient and exact. However, for the current Euclidean normalizing flows, data with hierarchies embedded in the Euclidean space will suffer high embedding distortion [94]. Besides, sampling from densities defined on Euclidean space cannot guarantee the generated points still lie on the underlying hyperbolic space. Therefore, it is fundamental to construct normalizing flows in the hyperbolic space.

Literally, normalizing flows have already been extended to Riemannian manifolds (spherical spaces) [50], [148], [149]. Work [126] is the pioneer one to present a new normalizing flow in the hyperbolic space. They proposed first elevated normalizing flows to hyperbolic space (leveraging Lorentz model) using coupling transforms defined on the tangent bundles. Then, they explicitly utilized the geometric structure of hyperbolic spaces and further introduced Wrapped Hyperboloid Coupling (WHC), which is a fully invertible and learnable transformation.

Based on the tangent space, work [126] provides a method called Tangent Coupling, which builds upon the real-valued non-volume preserving transformations (RealNVP flow) [150] and introduces the efficient affine coupling layers. Specifically, this work follows normalizing flows designed with partially-ordered dependencies [150]. They defined a class of invertible parametric hyperbolic functions $f_i^{\text{TC}} : \mathbb{L}^k \rightarrow \mathbb{L}^k$. The coupling layer is implemented using a binary mask, and partitions the input x into two sets $x_1 = x^{1:d}, x_2 = x^{d+1:n}$. For the first set x_1 , its elements are transformed elementwise independently of other dimensions, while the transform of second set x_2 is based on the first one. Thus, the overall transformation of one layer is

$$f^{\text{TC}}(\hat{x}) = \begin{cases} \hat{z}_1 = \hat{x}_1 \\ \hat{z}_2 = \hat{x}_2 \odot \delta(s(\hat{x}_1)) + t(\hat{x}_1) \end{cases}$$

$$f^{\mathcal{T}C}(x) = \text{Exp}_0(\hat{f}^{\mathcal{T}C}(\text{Log}_0(x))) \quad (55)$$

where \odot and δ are pointwise multiplication and pointwise non-linearity, respectively. s and t are map functions, which are implemented as linear neural layers and conduct the projection from $\mathcal{T}_0\mathbb{L}^d \rightarrow \mathcal{T}_0\mathbb{L}^{n-d}$. So the transformed result will be \hat{z} , which is the concatenation of resulted and mapped \hat{z}_1 and \hat{z}_2 on the manifold. They also provided an efficient expression for the Jacobian determinant by using the chain rule and identity.

To make full use of the expression power of the manifold, the authors conducted operations using parallel transport, instead of only operating in the tangent space of origin. Likewise, they provided an efficient expression for the Jacobian determinant. In this way, they built two different kinds of normalizing flows in the Lorentz model and improved the performance.

4.6 Hyperbolic Variational Auto-Encoders

Variational Auto-Encoder (VAE) [69], [70] is a popular probabilistic generative model, composed of an auxiliary encoder that draws samples of latent code from the approximate posterior (conditional density), and a decoder generating observations $x \in \mathcal{X}$ from latent variables $z \in \mathcal{Z}$. However, the vanilla VAE posterior is parameterized as a unimodal distribution such that it is not able to allow the structure assumption for data distributed in the hyperbolic space. Unfortunately, such a normal prior for the low-dimensional latent variables will encourage the low-dimensional representations of different samples to the center of the latent space, even for data consisting of distinct structures. Besides, embedding non-Euclidean data to a Euclidean space introduces significant distortion for commonly used dimensionality reduction tools, which is not good for visualization. Thus, it is meaningful to construct hyperbolic VAEs.

One of the main challenges of generalizing VAE to the hyperbolic space is the generalization of the latent distribution learning. As mentioned in [50], [119], there are mainly three ways to model the normal distributions in the hyperbolic space.

First, the Riemannian Normal [151] with Fréchet mean [73] μ and dispersion parameter σ . Sometimes, it is also referred as maximum entropy normal. In particular, the Riemannian normal distribution is defined as

$$\mathcal{N}_{\mathcal{M}}(z|\mu, \sigma^2) = \frac{1}{Z(\sigma)} \exp\left(\frac{-d_{\mathcal{M}}(\mu, z)}{2\sigma^2}\right) \quad (56)$$

where $d_{\mathcal{M}}$ is the induced distance [119], [152] and $Z(\sigma)$ is a normalization constant.

Second, the Restricted Normal. Restricting the sampled points from the normal distribution to sub-manifolds. It has also been treated as the maximum entropy distribution with respect to the ambient euclidean metric [153]. For instance, in the work Hyperpherical variational auto-encoders [154], the authors presented a novel VAE model, called S -VAE, via von Mises-Fisher (vMF) distribution, in which the encoder is a homeomorphism and can provide an invertible and globally continuous mapping.

Third, the Wrapped Normal [117], [119], [155]. This distribution is constructed by utilizing the exponential map of a Gaussian distribution on the tangent space centered at the mean value. Specifically, there are four steps to get a wrapped normal distribution. First, sample one point from the Euclidean normal distribution $\mathcal{N}(0, \sigma)$. Second, concatenate 0 as the zeroth coordinate of this point and transfer it to the tangent space of the origin. Third, parallel transport the sample from the current tangent space

to the tangent space at the point μ . Finally, map the point from the tangent space to the manifold. In this way, a latent sample on the manifold is obtained. As mentioned in [53], the Wrapped Normal has the following reparametrizable form

$$z = \text{Exp}_{\mu}(\mathcal{P}\mathcal{T}_{0 \rightarrow \mu}(v)), v \in \mathcal{N}(0, \sigma). \quad (57)$$

However, the Riemannian Normal distribution could be computationally expensive for sampling if it is based on rejection sampling. The Restricted Normal like vMF has only a single scalar covariance parameter, while other approaches can parametrize the covariance in different dimensions separately. On the contrary, sampling with Wrapped Normal distributions are very computationally efficient. Based on the aforementioned generalization of the normal distribution in the non-Euclidean space, there are numerous works [53], [54], [117], [119], [120] constructing VAE in the hyperbolic space, aiming at imposing structure information on the latent space.

Ovinnikov *et al.* [54] provided a closed-form definition of Gaussian distribution in the hyperbolic space, as well as the sampling rules for the prior and posterior distribution, to endow a VAE latent space with the ability to model underlying structure via the Poincaré ball model. Specifically, based on the maximum entropy generalization of Guassian distribution [151], they derived the normalization constant in Eq. (56) by decomposing it into radial and angular components. Based on the Wasserstein Autoencoder [156] framework, which is introduced to circumvent the high variance associated with the Monte-Carlo approximation, they built each layer using the hyperbolic feedforward layer provided by [43]. They also provided a generalization of the reparameterization trick by using the Möbius transformations. They further relaxed the constraint to the posterior by using the maximum mean discrepancy [157] and the network is optimized by RSGD [86]. However, as pointed out by [119], the authors had to choose a Wasserstein Auto-Encoder framework since they could not derive a closed-form solution of the ELBO's entropy term. Besides, work [117] mentions that the approximation of the likelihood and its gradient can be avoided.

Nagano *et al.* [117] provided a new normal distribution function in the hyperbolic space (Lorentz model), which is called Pseudo-Hyperbolic Gaussian, and it can be utilized to construct and learn a probabilistic model like VAE in this non-Euclidean space. The authors emphasized that this distribution is computed analytically and could be sampled efficiently. Pseudo-Hyperbolic Gaussian can be constructed with four steps, as mentioned above in the wrapped normal. The author further highlighted their contributions by deriving the density of Pseudo-Hyperbolic Gaussian distribution $\mathcal{G}(\mu, \Sigma)$ due to the exponential map and the parallel transport in the wrapped normal are differentiable. Since they provided a closed-form of the density function, they could evaluate the ELBO exactly and no need to introduce the reparameterization trick in this hyperbolic VAE.

However, as pointed out by [119], the neural layers are still the Euclidean ones, which do not take into account the hyperbolic geometry. Therefore, work [119] introduces a VAE that respects the geometry of the hyperbolic latent space. This is done by adding a generalization of the decision hyperplane in Euclidean space. Normally, the Euclidean linear affine transformation is $f(z) = \text{sign}(\langle a, z - p \rangle) \|a\| d(z, H_{a,p})$, where a is the coefficient, p is the intercept (offset). $H_{a,p}$ denotes a hyperplane passing through p with a as the normal direction, thus $d(z, H_{a,b})$ means

the Euclidean distance of z to the hyperplane. Analogue to the Euclidean linear function $f(z)$, they generalized it like

$$f_{a,p}(z) = \text{sign}(\langle a, \text{Log}_p(z) \rangle) \|a\|_p d^{\mathbb{B}}(z, H_{a,p}^{\mathbb{B}}) \quad (58)$$

Inspired by the MLR in work [43], the first layer of the decoder, which is called the gyroplane layer and is chosen to be a concatenation with a Poincaré operator f . Then it is then composed with a standard feed-forward neural network.

The gyroplane layer is then composed with a standard feed-forward neural network. For the encoder part, the author also changed the last layer by adding an exponential map for the Fréchet mean, and a softplus function for the positive defined Σ .

Then, the ELBO of VAE is optimized via an unbiased Monte Carlo (MC) estimator with two main Gaussian generalisations, which are wrapped normal and Riemannian normal generalization. Through a hyperbolic polar change of coordinates, they provided efficient and reparameterizable sampling schemes to calculate the probability density functions.

Compared to the above-mentioned methods, the work [53] highlights an implicit posterior and data-driven prior. They proposed an adversarial Poincaré variational autoencoder (APo-VAE), using a wrapped normal distribution as the prior and the variational posterior for a more expressive generalization. However, they replaced the tangent space sampling step in Eq. (57) with a more flexible implicit distribution from $\mathcal{N}(0, I)$, inspired by work [158]. Then, a geometry-aware Poincaré decoder is constructed, which shares the same idea as it for the decoder in work [119].

ApoVAE further optimized the variational bound by adversarially training this model by exploiting the primal-dual formulation of Kullback-Leibler (KL) divergence based on the Fenchel duality [159]. The training procedure is following the training scheme of coupled variational Bayes (CVB) from work [160] and implicit VAE [158]. Meanwhile, inspired by [161], they replaced the prior with a data-driven alternative to reduce the induced bias.

5 APPLICATIONS AND PERFORMANCE

Various applications from different research fields can benefit from hyperbolic deep neural networks, since the latent hierarchical structure is a generic property of real-world data. Thus, we will introduce the applications for processing data which contains hierarchical structures, such as graph embedding learning, natural language processing (NLP), and the analysis of data with tree-like properties. Moreover, we also notice there exists an increasing potential for the application of hyperbolic neural networks on data that has no obvious hierarchical structures, such as images. As a result, how hyperbolic networks can be adapted to computer vision tasks is also introduced in this section. Besides, another interesting application is about hyperbolic models for biology. Current studies demonstrate attractive results for measuring cells and their activities. Thus, we also introduce advanced approaches for computational biology with hyperbolic geometry.

5.1 Hyperbolic Models for Graph Applications

Hierarchies are ubiquitous in graph data. There are numerous works leveraging deep hyperbolic neural networks dealing with the graph tasks, including node classification [47], graph classification [46], [48], link prediction [47], and graph embedding [49]. Here, we only concentrate on the hyperbolic models on top of the GNN architecture, although there are many other existing hyperbolic embedding methods, like PVAE [119], HAT [44],

which also consider the modeling of graph (mainly for natural language or networks).

Works [102], [108] are pioneering works to introduce the concept of graph embeddings in the hyperbolic space. Recently, HGNN [46] and HGCN [47] are almost proposed at the same time to build GNN in the hyperbolic space. The HGCN is built on the Lorentz model. The authors applied HGCN to both node classification and link prediction tasks. Experiments show that for a dataset with low δ -hyperbolicity, the HGCN can get a performance much better than models built in the Euclidean space. The HGNN provides models on both the Poincaré model and Lorentz model. The authors dealt with graph classification tasks, in which the graphs are synthetically generated with three distinct graph generation algorithms. The results show that when the embedding dimension is smaller than 20, the HGNN has a significant superiority. We can also see that for most cases, HGNN on the Lorentz model performs better than that on the Poincaré model. While increasing the dimensions, this advantages disappeared and when the dimension is larger than 256, the HGNN even performs worse than its Euclidean counterpart. Currently, work [49] derives a general version of GCN with constant curvature, κ -GCN, which significantly minimizes the graph embedding distortion and gets a superior performance on the node classification tasks.

Most of previous mentioned tasks are focused on static graphs, while work [48] considered the graph classification task with a dynamic graph input. They constructed and searched for the optimal ST-GCN in the Poincaré model. Another interesting work [127] conducted extensive experiments on four different kinds of graph tasks, including node classification, link prediction, graph classification, and graph embedding, to provide a profound analysis of when the hyperbolic space can provide a superior performance. The experimental results suggest that the non-Euclidean space are not always a better choice than the Euclidean counterpart. As summarized by work [127], in the task that labels depend only on the local neighborhood of the node, hyperbolic models may be inferior to their Euclidean counterparts. However, many other factors, like the way to build the manifold and the corresponding optimization strategy, may also lead to this result. Therefore, more explorations are needed to draw this conclusion.

5.2 Knowledge Graph Completion

Knowledge graph is a multi-relational graph representation of a collection of facts \mathcal{F} , formed in a set of triplet (e_s, r, e_o) , where e_s is the subject entity, e_o is the object entity and r is a binary relation (typed directed edges) between them. $(e_s, r, e_o) \in \mathcal{F}$ denotes subject entity e_s is related to object entity e_o by the relation r . As mentioned by work [116], knowledge graphs often exhibit multiple hierarchies simultaneously. For instance, nodes near the root of the tree under one relation may be leaf nodes under another. The challenge for representing multi-relational data lies in the difficulty to represent entities shared across relations, such that different hierarchies are formed under different relations.

Most of the previous works [162], [163] model in the Euclidean space, relying on the inner product or Euclidean distance as a similarity measure, which can be categorised as translational models and bilinear models, respectively. Work [164] proposes to embed the entities in a Riemannian manifold, where each relation is modeled as a move to a point and they also defined specific novel distance dissimilarities for the relations. However, as pointed out by [116], this model defined in the hyperbolic space does not outperform Euclidean models. Work [116] presents a

new bilinear model called MuRP to embed hierarchical multi-relational data in the Poincaré ball model of the hyperbolic space. MuRP defines a basis score function for multi-relational graph embedding and generalizes it to the hyperbolic space. Experiments show that MuRP can get superior performances on the link prediction task. Besides, it requires far fewer dimensions than Euclidean embeddings to achieve comparable performance. Work [165] points out that MuRP is not able to encode some logical properties of relationships. Therefore, the authors leveraged the hyperbolic isometry to simultaneously exhibit logical patterns and hierarchies, and achieved the current best performance.

On the contrary, work [14] proposes a new translational model in Poincaré model, where both the entities and the relations are embedded. Compared to the translational models in Euclidean space, this method almost doubles the performance in terms of the mean reciprocal rank (MRR) metric. However, as mentioned by the authors, the HyperKG excludes from the comparison with many recent works that explores advanced techniques, thus this method is not comparable to the state-of-the-art methods, as listed in results in [116].

5.3 Natural language processing (NLP)

As summarized in Table 1, there are more than one third hyperbolic methods being presented to deal with NLP tasks, of which specific tasks include text classification [125], taxonomy induction [115], taxonomies embedding [45], word embeddings [51], [52], Lexical Entailment [103] and text generation [53]. Natural language often conceives a latent hierarchical structure, *e.g.*, linguistic ontology. It is natural to turn to the hyperbolic space. Another advantage of modeling in the hyperbolic space, as mentioned by work [53], is that the latent representation allows more control of the sentences we want to generate. In the following part, we will detail different NLP tasks using hyperbolic geometry, from tasks of embedding, generation, to classification.

In the task of **word embeddings**, work [117] proposes a Gaussian-like distribution in the hyperbolic space, which is called pseudo-hyperbolic Gaussian. Based on this, a hyperbolic VAE is presented to deal with the word embeddings. Work [51] adapts the Glove [105] algorithm to learn unsupervised word embeddings in this type of Riemannian manifolds. To this end, they proposed to embed words in a Cartesian product of hyperbolic spaces which they theoretically connected to the Gaussian word embedding and the Fisher geometry. Some notable finds are, based on their method, the fully unsupervised model can almost outperform all supervised Euclidean counterparts. Once trained with a small amount of weakly supervision for the hypernymy [166] score, they can obtain significant improvements and this result is much better than the models in Euclidean space.

In the **taxonomy embedding** task, work [103] is one pioneer research piece that can learn embeddings in the hyperbolic space. To evaluate its ability to infer hierarchical relationships **without supervision**, they trained on data where the hierarchy of objects is not explicitly encoded. A significant improvement was witnessed in the taxonomy embedding task. Dhingra *et al.* [52] proposed a re-parametrization of Poincaré embeddings that removes the need for the projection step and allows the use of any of the popular optimization techniques in deep learning, such as Adam [89]. After this several works are presented to deal with the stability of the hyperbolic embedding. Nickel *et al.* [45] pointed out that the Lorentz model is substantially more efficient thus proposed to learn a continuous representation using the model and

further improves the performance. Marc *et al.* [118] mentioned that Poincaré distance is numerically unstable and suggested the squared Lorentzian distance as a better choice. Based on this, they learned a closed-form squared Lorentzian distance and thus improved the performance on the this task. Besides, Yu *et al.* [106] constructed the hyperbolic model for dealing with the numerical instability of the previous hyperbolic networks. Based on the Lorentz model, they provide a very efficient model to learn the embedding. For instance, they can even compress the embeddings down to 2% of a Poincaré embedding on the WordNet Nouns. Very recently, Shimizu *et al.* [4] shows the superior parameter efficiency of their methods compared to conventional hyperbolic components, and the stability and outperform over their Euclidean counterparts.

Then, we introduce the hyperbolic models for generation tasks. In the **text generation** task, HyperQA [110] is the first to model QA pairs in the hyperbolic space. HyperQA is an extremely fast and parameter efficient model that achieves very competitive results on multiple QA benchmarks, especially when compared to Euclidean methods at that time. ApoVAE [53] also deals with the dialog-response generation problem. It optimized the variational bound by adversarially training and exploited the primal-dual formulation of KL divergence based on the Fenchel duality [159]. In **neural machine translation**, based on the hyperbolic attention mechanism, work [44] provides a hyperbolic Transformer, which shows the improvement over the original one, especially when the model capacity is restricted.

Finally, hyperbolic models are also commonly used for text-related classification tasks. In the task of **sentence entailment classification**, works [43] and [4] proposed the hyperbolic MLRs. The results confirm the tendency of the hyperbolic MLRs to outperform the Euclidean version in all settings. At the same time, the hyperbolic MLR from work [4] shows more stable training, relatively narrower confidence intervals, and at least comparable performance with only half of the parameters compared to the MLR in [43]. In the task of **text classification**, HyperText [125] performs significantly better than the state-of-the-art text classifier, FastText [167], in Euclidean space. Besides, HyperText with 50-dimension achieves better performance to FastText with 300-dimension, which proves the hyperbolic space is a better choice for this task. Also, works [125] and [168] benefits from the hyperbolic methods in the Chinese text analysis tasks. **Textual entailment**, which is also called natural language inference, is a binary classification task to predict whether the second sentence (hypothesis) can be inferred from the first one (premise). HNN [43] embedded two sentences using two distinct hyperbolic RNNs. With the corresponding distances, the sentence embeddings are then fed into a feedforward network and predicted with an MLR. Interestingly, the results shows that the fully Euclidean baselines might even have an advantage over hyperbolic models. On top of pre-trained Poincaré embeddings [103], they conducted experiments on the WordNet noun hierarchy to evaluate the hyperbolic MLR. On this subtree classification task, hyperbolic MLR displays a clear advantage to the Euclidean counterpart. Nevertheless, in what case hyperbolic is more suitable is not clear enough.

5.4 Hyperbolic Space for Recommender Systems

One of the most important factors to the success of a recommender system is the accurate representation of user preferences and item characteristics, modelled by complex networks. As mentioned in [99], hyperbolic geometry naturally emerges from network

heterogeneity in the same way that network heterogeneity emerges from hyperbolic geometry. Therefore, given the complex nature of these networks, the hyperbolic space is more suitable to embed them than its Euclidean counterpart. Based on the above observations, work [12] embeds bipartite user-item graphs in the hyperbolic space. This recommendation algorithm learns to rank loss that represents user-item correlations, using of hyperbolic representations through an analogy with complex networks. This algorithm shows a clear advantage when compared with the system in the Euclidean space. The recommender system based on this algorithm is scaled to millions of users.

Although this system shows obvious superiority, work [123] points out that it does not learn the embeddings in a metric learning manner. Therefore, the authors explored the connections between metric learning and collaborative filtering, thus proposed a highly effective model for recommender systems. They constructed an input triplet tuple with the user, an item liked by the user, and the item unliked by the user. Then they proposed to learn the user-item joint metric in the hyperbolic space. At the same time, they introduced so-called local and global factor to better embed user-item pairs to the hyperbolic space and preserve good structure quality for metric learning.

As pointed out by work [124], in the factorization machine model [169], the naive inner product is not expressive enough for spurious or implicit feature interactions. Therefore, higher-order factorization machine [170] is proposed to learn higher-order feature interactions efficiently. As suggested by collaborative metric learning [171], learning distance instead of inner product has advantages to provide a fine-grained embedding space which could capture the representation for item-user interactions, item-item and user-user distances at the same time. Thus, the triangle inequality is more preferred than the inner product. Inspired by this, work [124] proposes a model named Lorentzian Factorization Machine (LorentzFM), which learns feature interactions with a score function measuring the validity of triangle inequalities. The authors argued that the feature interaction between two points can be learned by the **sign** of the triangle inequality for Lorentz distance, rather than using the distance itself. Based on this, they presented the model in the hyperbolic space.

5.5 Hyperbolic Models for Computer Vision

The passion of solving computer vision tasks using hyperbolic models is inspired from the observation that similar hierarchical relations between images are also common in computer vision tasks [16]. Besides, hierarchies investigated in NLP can be also transcended to the visual domain, like the knowledge graph.

Work [16] is one of the pioneering methods to model images in the hyperbolic space. To prove it is reasonable to utilize the hyperbolic space for image-based tasks, the δ -Hyperbolicity is used to measure the property of “negatively curved” of the features extracted from the embedding network. The authors concluded that the feature embeddings from current famous architectures like ResNet [21], VGG19 [40], and InceptionV3 [41] are with a small δ thus suggested the learned features process strong hierarchical relationships. Based on these observations, the authors constructed the analogues of layers in the hyperbolic spaces. They evaluated their models in computer vision tasks, including person re-identification [172], and few-shot classification [173], and results proved its superiority.

Another study in this field is trying to generalize generative models like VAE to the hyperbolic space and deal with

image reconstruction or generation tasks. Work [54] provides a Wasserstein Autoencoder for the Poincaré model and applied it to the task of generating binarized MNIST digits in order to get an intuition for the properties of the latent hyperbolic geometry. However, they did not provide a better reconstruction results when compared to the Euclidean counterpart. As mentioned by the authors, both models meet with a dimension mismatch problem [174] such that reconstructed samples present a deteriorating quality as the dimensionality increases despite the lower reconstruction error. Compared to work [119], the authors derived a closed-form solution of the ELBO with two different kinds of normal distributions in the hyperbolic space. Their VAE model outperforms its Euclidean counterpart, especially for low latent dimensions. However, as the latent dimension increases, the embeddings quality decreases, hence the gain from the hyperbolic geometry is reduced, just as also observed by work [103]. The same situation is also found in the work [129], where a mixed geometry space is introduced. On the MINIST reconstruction task, they displayed clear advantage when setting the latent dimension to six. However, this advantage to Euclidean spaces immediately disappears when they double the dimension of latent space. This can be also caused by the property of the datasets. More studies are needed to answer whether the hyperbolic space has its advantages to address such computer vision tasks. Besides, it needs to extend to more complicated settings in larger-scale cases.

5.6 Computational Biology with Hyperbolic Geometry

Computational Biology [175] is an interdisciplinary area using biological data-driven computational models to understand biological systems and relationships. Meanwhile, hierarchical representations, such as phylogenetic trees and clustering clades have long been applied to characterize differences between cells, proteins, the activity within cells [37]. It is natural to consider hyperbolic metric as an alternative when modeling biological data in computational biology. In line with this, multiple advances have been made in computational biology towards the goal of discovering and analyzing hierarchical structures from single-cell measurements. This section presents advanced results [15], [37], [135] in single-cell RNA-seq analysis, cells developmental processes, and gene expression analysis in the hyperbolic space.

For single-cell analysis, one major difficulty stems from how to reveal the progression of cells along continuous trajectories with multiple tree-like branches. Especially for complex hierarchies, modeling with efficient low-dimensional embeddings in Euclidean space will substantially distort distances between measurements, which definitely is an unwanted issue for modeling the progression. Based on the hyperbolic embedding [103], Klimovskaya *et al.* [15] provided Poincaré maps for the discovery and analyzing of complex hierarchies in single-cell data. In addition, the approach deals with all these different tasks such as clustering, lineage detection, and pseudotime inference using a single embedding in an unsupervised manner, which is impossible for previous works like t-SNE [176], PCA [177], and UMAP [178].

There are three steps in this method, of which the first two are used to approximate an unknown manifold and the last is to learn the hyperbolic embedding. First, a connected k-nearest-neighbor graph (kNNG) [179] is constructed to embed each cell and measure their Euclidean distances. This step codes the local geometries of an underlying manifold. Second, based on kNNG built before, to estimate the intrinsic geometry of the underlying manifold, the global geodesic distances are computed. Then, the

third step learns a two-dimensional Poincaré embedding for each cell, which preserves the topology. This can place nodes with small distances (like cells in early developmental stage) close to the center of the Poincaré disk and nodes with large distances close to the boundary. Poincaré maps produce state-of-the-art two-dimensional representations of cell trajectories on multiple scRNAseq datasets.

Following this work, Ding *et al.* [135] focused on eliminating the batch-correction and addressing visual crowding issues of conventional generative modeling approaches for Single-cell RNA-seq (scRNA-seq). The motivations are: First, normally, scRNA-seq is very high-dimensional data with typically low intrinsic dimensionality due to the co-expressed property of genes. Second, the crowding issues caused by multidimensional normal prior assumption in Euclidean space lead to unreasonable gathering at the center of the latent space, even for cells from distinct cell types. Third, datasets typically have multiple technical and biological factors which cannot be handled by current VAE or batch-correction method. To address above issues, Ding *et al.* provided to represent and infer of branched developmental trajectories in the hyperbolic spaces via deep generative embedding model, of which a wrapped normal [117], [119] distribution is used as the prior for the latent variable. In this way, latent representation is learned in the hyperbolic space and the resulted latent structure accounts for the multiple batch effects. Visualization on large datasets with multiple cell types and hierarchical structures shows a much better results without cell crowding problems, which proves the effectiveness of the learned representations.

Another interesting and valuable work is applying hyperbolic geometry to gene expressions [37] analysis, which forms an important part to understand how the genotype of an organism impacts its phenotype, like disease. The challenge of this problem lies at the complexity of their relationship as thousands of genes can affect a single phenotype of interest [180]. Fortunately, the widespread correlations between genes suggests that a low-dimensional geometry can be applied to model the genetic variation and their expressions [181]. Therefore, the authors in work [37] developed a quantitative test for distinguishing the curvature of the underlying low-dimensional geometry. Besides, by incorporating hyperbolic metric into the t-SNE method, they provided visualization tools for data that exhibit a low-dimensional hyperbolic geometry. Results on several gene expression datasets from mouse and human prove that gene expression can be effectively described using low-dimensional hyperbolic metric.

6 DISCUSSION AND OPEN PROBLEMS

6.1 When to expect benefits from Hyperbolic networks

One observation is that hyperbolic embeddings or hyperbolic neural networks cannot consistently work better than the Euclidean counterparts [43]. Many times the Euclidean counterpart can recover its superiority when adding the embedding dimensions. It is not clear whether the hyperbolic model can only provide a more compact model or it can provide a more efficient model with significant performance improvement. A better understanding of when and why the use of hyperbolic geometry is justified is crucial. Based on current study results, we summarize when to expect such potential benefits.

First, when there is hierarchical data. This is quite clear as hyperbolic space is naturally a better place for such data and could provide a low distortion embedding. Therefore, for data like realistic complex networks [182], data with tree-structure, graph

with power-law distribution (this can be organized with latent hierarchy), generally we can expect a better performance using the hyperbolic space. Second, when the data or feature has a low Gromov δ -hyperbolicity [57]. Basically, for hierarchical data, its δ will be very small. For data without a clear observed hierarchy, Gromov δ -hyperbolicity is a better way to measure this underlying structure. From work [16], we know for Poincaré disk, they get experimental value $\delta = 0.18$, and for the feature representations for miniImageNet learned from VGG, $\delta = 0.17$. Thus, data or feature with similar value can be treated as data possessing such hierarchy thus can expect an improvement by applying the hyperbolic space. Third, when the process is expected to have a hierarchical development. Although the data itself does not possess hierarchy, the evolutionary process of samples may have strong hierarchical relationships, *e.g.*, the developmental progress of single cell [15]. Therefore, we can expect success in such relationship modeling task or other reason induction cases. Fourth, when dealing with huge dataset with extremely limited resources. Many high-dimensional data exhibit an underlying low-dimensional hyperbolic geometry [37]. As well, hyperbolic embeddings could provide an extremely low-dimensional coding with low distortion. In such case, the hyperbolic space can provide richer information and a more compact model. Fifth, when model interpretation is much more important than the performance, the hyperbolic space is also a better choice. One of the biggest issues of current Euclidean neural network is the lack of interpretability. The black-box property of deep learning keeps causing concern in real-life applications. For example, a great deep model can perform very well on the seen-dataset, while it could provide extremely wrong predicts with very high confidence on the unseen dataset. On the contrary, hyperbolic model has much better interpretation as it would give an uncertain prediction with low confidence [15]. Last, hyperbolic neural models can also be introduced as a supplemental information branch for helping traditional Euclidean neural model. For instance, work [134] develops more advanced architectures by the interaction of Euclidean and hyperbolic spaces. The learned representations show a superiority.

6.2 Which hyperbolic model?

As shown in Table 1, majority of the research in the literature is based on the Poincaré model. However, on one hand, current research has found that the Lorentz model has better numerical stability properties [45], due to its large variance when close to the boundary. On the top of the Lorentz model, Yu *et al.* provided a more stable model [106] by introducing tiling method. On the other hand, Lorentz model is un-bounded from the definition, which is not friendly to modern neural networks. Therefore, more studies are needed to choose the right hyperbolic model.

6.3 Neural Networks with Mixed Geometries

Currently, we also find a trend to construct neural networks utilizing mixed geometries². As mentioned by work [114], the quality of the learned representations is determined by how well the geometry of the embedding space matches the underlying structure of the data. Therefore, for the real-world data possessing multiple complicated structures, like in the rightmost of Fig. 3, learning representations via hybrid geometry may be a better choice. Nonetheless, one major issue is how to construct the mixed manifold and at the same time determine its signature. Besides,

² Note, the 'Mixed' methods listed in Table 1 combine geometries with different types of curvature, not different type of hyperbolic models.

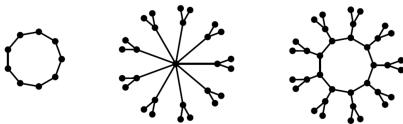


Fig. 3. Illustration of data with different structures. Leftmost shows a data with cycle. Middle shows a tree-structured data and the rightmost is a combination.

current optimizations require a costly grid search to tune hyperparameters. Efficient Optimization is also highly expected for the neural network with mixed geometries. Fortunately, there are some representative attempts for mixed geometries, *i.e.*, First, product spaces [114], which utilizes a Riemannian product manifold of model spaces [183], including hyperbolic spaces , spherical space, and the flat Euclidean space. Second, mixed-curvature method [49], [129], which partitions the space into multiple component spaces and learn a curvature for each part. Third, pseudo-Riemannian manifolds of constant nonzero curvature [133], [184], which include the hyperbolic and spherical geometries and whose non-degenerate metric tensor is not constrained to be positive definite.

6.4 Advanced Hyperbolic Networks

One potential direction can be the combination of Riemannian neural network with advanced deep learning technology. For instance, exploring new Riemannian neural architectures with advanced automatic machine learning methods, like NAS [146]. Work [48] provides to search the best projection dimension in the Poincaré model, utilizing the NAS method. However, there is much room to improve by automatically designing the neural modules, instead of only searching for optimal projection dimensions. Another important research topic can be the generalization of more sophisticated Euclidean optimization algorithms. In many cases, as mentioned in [43], fully Euclidean baseline models might have an advantage over hyperbolic baselines. One possible reason is that Euclidean space is equipped with much more advanced optimization tools. Once the hyperbolic neural networks are also equipped with such tools, we can expect more from the powerful hyperbolic networks.

ACKNOWLEDGMENTS

The authors would like to thank Octavian-Eugen Ganea, from MIT for the useful discussion. We also want to thank Emile Mathieu, from University of Oxford, for the explanation regarding the Gyroplane layer in their Poincaré Variational Auto-Encoder. As well, we thank professor Stan Z. Li from Westlake University for the detailed comments of the Riemannian geometry.

This work is supported by the Academy of Finland for ICT 2023 project (grant 328115), Academy Professor project EmotionAI (grants 336116, 345122), project MiGA (grant 316765) and Infotech Oulu.

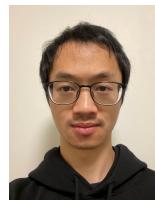
REFERENCES

- [1] M. E. Newman, “Power laws, pareto distributions and zipf’s law,” *Contemporary physics*, 2005.
- [2] H. W. Lin and M. Tegmark, “Critical behavior in physics and probabilistic formal languages,” *Entropy*, 2017.
- [3] K. Katayama and E. W. Maina, “Indexing method for hierarchical graphs based on relation among interlacing sequences of eigenvalues,” *Journal of information processing*, 2015.
- [4] R. Shimizu, Y. Mukuta, and T. Harada, “Hyperbolic neural networks++,” 2021.
- [5] M. Boguñá, F. Papadopoulos, and D. Krioukov, “Sustaining the internet with hyperbolic mapping,” *Nature communications*, 2010.
- [6] B. Tadić, M. Andjelović, and M. Šuvakov, “Origin of hyperbolicity in brain-to-brain coordination networks,” *Frontiers in Physics*, 2018.
- [7] G. García-Pérez, M. Boguñá, A. Allard, and M. Á. Serrano, “The hidden hyperbolic geometry of international trade: World trade atlas 1870–2013,” *Scientific reports*, 2016.
- [8] M. Keller-Ressel and S. Nargang, “The hyperbolic geometry of financial networks,” *Scientific reports*, 2021.
- [9] R. Krauthgamer and J. R. Lee, “Algorithms on negatively curved spaces,” in *2006 47th Annual IEEE Symposium on Foundations of Computer Science (FOCS’06)*. IEEE, 2006.
- [10] E. Begelfor and M. Werman, “The world is not always flat or learning curved manifolds,” *School of Engineering and Computer Science, Hebrew University of Jerusalem, Tech. Rep*, 2005.
- [11] M. M. Bronstein, J. Bruna, Y. LeCun, A. Szlam, and P. Vandergheynst, “Geometric deep learning: going beyond euclidean data,” *IEEE Signal Processing Magazine*, 2017.
- [12] B. P. Chamberlain, S. R. Hardwick, D. R. Wardope, F. Dzogang, F. Daolio, and S. Vargas, “Scalable hyperbolic recommender systems,” *CoRR*, 2019.
- [13] K. Verbeek and S. Suri, “Metric embedding, hyperbolic space, and social networks,” in *Proceedings of the thirtieth annual symposium on Computational geometry*, 2014.
- [14] P. Kolyvakis, A. Kalousis, and D. Kiritsis, “Hyperkg: hyperbolic knowledge graph embeddings for knowledge base completion,” *arXiv*, 2019.
- [15] A. Klimovskaya, D. Lopez-Paz, L. Bottou, and M. Nickel, “Poincaré maps for analyzing complex hierarchies in single-cell data,” *Nature communications*, 2020.
- [16] 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.
- [17] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, 2015.
- [18] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard *et al.*, “Tensorflow: A system for large-scale machine learning,” in *12th {USENIX} symposium on operating systems design and implementation ({OSDI} 16)*, 2016.
- [19] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer, “Automatic differentiation in pytorch,” 2017.
- [20] 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,” in *Advances in neural information processing systems*, 2019.
- [21] 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.
- [22] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *IEEE conference on computer vision and pattern recognition*. IEEE, 2009.
- [23] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, 2017.
- [24] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” *arXiv*, 2013.
- [25] J. L. Coolidge, *The elements of non-Euclidean geometry*. At the Clarendon press, 1909.
- [26] M. R. Bridson and A. Haefliger, *Metric spaces of non-positive curvature*. Springer Science & Business Media, 2013.
- [27] C. Esteves, C. Allen-Blanchette, A. Makadia, and K. Daniilidis, “Learning so (3) equivariant representations with spherical cnns,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [28] T. S. Cohen, M. Geiger, J. Köhler, and M. Welling, “Spherical cnns,” *International Conference on Learning Representations (ICLR)*, 2018.
- [29] M. Defferrard, M. Milani, F. Gusset, and N. Perraudin, “Deepsphere: a graph-based spherical cnn,” *International Conference on Learning Representations (ICLR)*, 2020.
- [30] N. Perraudin, M. Defferrard, T. Kacprzak, and R. Sgier, “Deepsphere: Efficient spherical convolutional neural network with healpix sampling for cosmological applications,” *Astronomy and Computing*, 2019.
- [31] M. Gromov, “Hyperbolic groups,” in *Essays in group theory*. Springer, 1987.
- [32] R. G. Barker and H. F. Wright, “Midwest and its children: The psychological ecology of an american town.” 1955.

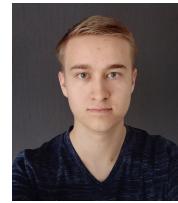
- [33] A. M. Collins, M. R. Quillian *et al.*, "Retrieval time from semantic memory," *Journal of verbal learning and verbal behavior*, 1969.
- [34] F. C. Keil, *Semantic and conceptual development: An ontological perspective*. Harvard University Press, 2013.
- [35] D. M. Roy, C. Kemp, V. K. Mansinghka, and J. B. Tenenbaum, "Learning annotated hierarchies from relational data," in *Advances in neural information processing systems*, 2007.
- [36] Y. Zhou, B. H. Smith, and T. O. Sharpee, "Hyperbolic geometry of the olfactory space," *Science advances*, 2018.
- [37] Y. Zhou and T. O. Sharpee, "Hyperbolic geometry of gene expression," *Iscience*, 2021.
- [38] O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra *et al.*, "Matching networks for one shot learning," *Advances in neural information processing systems*, 2016.
- [39] A. Krizhevsky *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [40] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv*, 2014.
- [41] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015.
- [42] Wiki - Non-euclidean space. [Online]. Available: https://en.wikipedia.org/wiki/Non-Euclidean_geometry
- [43] O. Ganea, G. Béćigneul, and T. Hofmann, "Hyperbolic neural networks," *Advances in neural information processing systems*, 2018.
- [44] 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*, 2018.
- [45] M. Nickel and D. Kiela, "Learning continuous hierarchies in the lorentz model of hyperbolic geometry," *Proceedings of the 35-th International Conference on Machine Learning, PMLR*, 2018.
- [46] Q. Liu, M. Nickel, and D. Kiela, "Hyperbolic graph neural networks," in *Advances in Neural Information Processing Systems*, 2019.
- [47] I. Chami, Z. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," in *Advances in neural information processing systems*, 2019.
- [48] W. Peng, J. Shi, Z. Xia, and G. Zhao, "Mix dimension in poincaré geometry for 3d skeleton-based action recognition," in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020.
- [49] G. Bachmann, G. Béćigneul, and O. Ganea, "Constant curvature graph convolutional networks," in *International Conference on Machine Learning*. PMLR, 2020.
- [50] D. J. Rezende, G. Papamakarios, S. Racanière, M. S. Albergo, G. Kanwar, P. E. Shanahan, and K. Cranmer, "Normalizing flows on tori and spheres," *Proceedings of the 37th International Conference on Machine Learning*, 2020.
- [51] A. Tifrea, G. Béćigneul, and O.-E. Ganea, "Poincaré glove: Hyperbolic word embeddings," *International Conference on Learning Representations (ICLR)*, 2019.
- [52] B. Dhingra, C. J. Shallue, M. Norouzi, A. M. Dai, and G. E. Dahl, "Embedding text in hyperbolic spaces," *arXiv*, 2018.
- [53] S. Dai, Z. Gan, Y. Cheng, C. Tao, L. Carin, and J. Liu, "Apo-vae: Text generation in hyperbolic space," *arXiv*, 2020.
- [54] I. Ovinnikov, "Poincaré's wasserstein autoencoder," *arXiv*, 2019.
- [55] M. P. d. Carmo, *Riemannian geometry*. Birkhäuser, 1992.
- [56] S. Gallot, D. Hulin, and J. Lafontaine, *Riemannian geometry*. Springer, 1990.
- [57] A. B. Adcock, B. D. Sullivan, and M. W. Mahoney, "Tree-like structure in large social and information networks," in *2013 IEEE 13th International Conference on Data Mining*. IEEE, 2013.
- [58] M. Bonk and O. Schramm, "Embeddings of gromov hyperbolic spaces," in *Selected Works of Oded Schramm*. Springer, 2011.
- [59] E. Beltrami, "Teoria fondamentale degli spazi di curvatura costante," *Annali di Matematica Pura ed Applicata (1867-1897)*, 1868.
- [60] J. W. Cannon, W. J. Floyd, R. Kenyon, W. R. Parry *et al.*, "Hyperbolic geometry," *Flavors of geometry*, 1997.
- [61] M. Hamann, "On the tree-likeness of hyperbolic spaces," in *Mathematical Proceedings of the Cambridge Philosophical Society*. Cambridge University Press, 2018.
- [62] A. Dyubina and I. Polterovich, "Explicit constructions of universal - trees and asymptotic geometry of hyperbolic spaces," *Bulletin of the London Mathematical Society*, 2001.
- [63] R. Kleinberg, "Geographic routing using hyperbolic space," in *IEEE INFOCOM 2007-26th IEEE International Conference on Computer Communications*. IEEE, 2007.
- [64] A. Lou, I. Katsman, Q. Jiang, S. Belongie, S.-N. Lim, and C. De Sa, "Differentiating through the fr'echet mean," *Thirty-seventh International Conference on Machine Learning*, 2020.
- [65] Z. Huang, R. Wang, S. Shan, L. Van Gool, and X. Chen, "Cross euclidean-to-riemannian metric learning with application to face recognition from video," *IEEE transactions on pattern analysis and machine intelligence*, 2017.
- [66] O. Tuzel, F. Porikli, and P. Meer, "Pedestrian detection via classification on riemannian manifolds," *IEEE transactions on pattern analysis and machine intelligence*, 2008.
- [67] A. A. Ungar, "Hyperbolic trigonometry and its application in the poincaré ball model of hyperbolic geometry," *Computers & Mathematics with Applications*, 2001.
- [68] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," *arXiv*, 2015.
- [69] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," 2014.
- [70] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," *International Conference on Machine Learning (ICML)*, 2014.
- [71] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," *arXiv*, 2016.
- [72] A. A. Ungar, "A gyrovector space approach to hyperbolic geometry," *Synthesis Lectures on Mathematics and Statistics*, 2008.
- [73] M. R. Fréchet, "Les éléments aléatoires de nature quelconque dans un espace distancié," *Annales de l'institut Henri Poincaré*, 1948.
- [74] M. Fréchet, "Les éléments aléatoires de nature quelconque dans un espace distancié," in *Annales de l'institut Henri Poincaré*, 1948.
- [75] S. Gould, B. Fernando, A. Cherian, P. Anderson, R. S. Cruz, and E. Guo, "On differentiating parameterized argmin and argmax problems with application to bi-level optimization," *arXiv*, 2016.
- [76] K. He, X. Zhang, S. Ren, and J. Sun, "Delving deep into rectifiers: Surpassing human-level performance on imagenet classification," in *Proceedings of the IEEE international conference on computer vision*, 2015.
- [77] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, 1997.
- [78] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013.
- [79] K. Cho, B. v. M. C. Gulcehre, D. Bahdanau, F. B. H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation."
- [80] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *ICML*, 2010.
- [81] B. Xu, N. Wang, T. Chen, and M. Li, "Empirical evaluation of rectified activations in convolutional network," *arXiv*, 2015.
- [82] G. Lebanon and J. Lafferty, "Hyperplane margin classifiers on the multinomial manifold," in *Proceedings of the twenty-first international conference on Machine learning*, 2004.
- [83] H. Cho, B. DeMeo, J. Peng, and B. Berger, "Large-margin classification in hyperbolic space," in *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019.
- [84] M. Weber, M. Zaheer, A. S. Rawat, A. Menon, and S. Kumar, "Robust large-margin learning in hyperbolic space," *Advances in Neural Information Processing Systems*, 2020.
- [85] M. Kochurov, R. Karimov, and S. Kozlukov, "Geoopt: Riemannian optimization in pytorch," 2020.
- [86] S. Bonnabel, "Stochastic gradient descent on riemannian manifolds," *IEEE Transactions on Automatic Control*, 2013.
- [87] H. Zhang and S. Sra, "First-order methods for geodesically convex optimization," in *Conference on Learning Theory*, 2016.
- [88] J. Duchi, E. Hazan, and Y. Singer, "Adaptive subgradient methods for online learning and stochastic optimization." *Journal of machine learning research*, 2011.
- [89] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *International Conference on Learning Representations (ICLR)*, 2015.
- [90] G. Béćigneul and O.-E. Ganea, "Riemannian adaptive optimization methods," 2019.
- [91] Y. LeCun, Y. Bengio *et al.*, "Convolutional networks for images, speech, and time series," *The handbook of brain theory and neural networks*, 1995.
- [92] M. Schuster and K. K. Paliwal, "Bidirectional recurrent neural networks," *IEEE transactions on Signal Processing*, 1997.
- [93] J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun, "Spectral networks and locally connected networks on graphs," *arXiv*, 2013.

- [94] R. Sarkar, "Low distortion delaunay embedding of trees in hyperbolic plane," in *International Symposium on Graph Drawing*. Springer, 2011.
- [95] J. A. Walter, "H-mds: a new approach for interactive visualization with multidimensional scaling in the hyperbolic space," *Information systems*, 2004.
- [96] Y. Shavitt and T. Tanel, "Hyperbolic embedding of internet graph for distance estimation and overlay construction," *IEEE/ACM Transactions on Networking*, 2008.
- [97] D. Krioukov, F. Papadopoulos, A. Vahdat, and M. Boguná, "Curvature and temperature of complex networks," *Physical Review E*, 2009.
- [98] A. Cvetkovski and M. Crovella, "Hyperbolic embedding and routing for dynamic graphs," in *IEEE INFOCOM*. IEEE, 2009.
- [99] D. Krioukov, F. Papadopoulos, M. Kitsak, A. Vahdat, and M. Boguná, "Hyperbolic geometry of complex networks," *Physical Review E*, 2010.
- [100] A. Cvetkovski and M. Crovella, "Multidimensional scaling in the poincaré disk," *arXiv preprint arXiv:1105.5332*, 2011.
- [101] T. Bläsius, T. Friedrich, A. Krohmer, and S. Laue, "Efficient embedding of scale-free graphs in the hyperbolic plane," *IEEE/ACM transactions on Networking*, 2018.
- [102] A. Muscoloni, J. M. Thomas, S. Ciucci, G. Bianconi, and C. V. Cannistraci, "Machine learning meets complex networks via coalescent embedding in the hyperbolic space," *Nature communications*, 2017.
- [103] M. Nickel and D. Kiela, "Poincaré embeddings for learning hierarchical representations," in *Advances in neural information processing systems*, 2017.
- [104] I. Vendrov, R. Kiros, S. Fidler, and R. Urtasun, "Order-embeddings of images and language," *International Conference on Learning Representations (ICLR)*, 2016.
- [105] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, 2014.
- [106] T. Yu and C. M. De Sa, "Numerically accurate hyperbolic embeddings using tiling-based models," in *Advances in Neural Information Processing Systems*, 2019.
- [107] J. H. Conway, H. Burgiel, and C. Goodman-Strauss, *The symmetries of things*. CRC Press, 2016.
- [108] B. P. Chamberlain, J. Clough, and M. P. Deisenroth, "Neural embeddings of graphs in hyperbolic space," *arXiv*, 2017.
- [109] C. De Sa, A. Gu, C. Ré, and F. Sala, "Representation tradeoffs for hyperbolic embeddings," *Proceedings of machine learning research*, 2018.
- [110] Y. Tay, L. A. Tuan, and S. C. Hui, "Hyperbolic representation learning for fast and efficient neural question answering," in *Proceedings of the Eleventh ACM International Conference on Web Search and Data Mining*, 2018.
- [111] O.-E. Ganea, G. Bécigneul, and T. Hofmann, "Hyperbolic entailment cones for learning hierarchical embeddings," *International Conference on Machine Learning (ICML)*, 2018.
- [112] T. D. Q. Vinh, Y. Tay, S. Zhang, G. Cong, and X.-L. Li, "Hyperbolic recommender systems," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.
- [113] S. Rendle, C. Freudenthaler, Z. Gantner, and L. Schmidt-Thieme, "Bpr: Bayesian personalized ranking from implicit feedback," *arXiv*, 2012.
- [114] A. Gu, F. Sala, B. Gunel, and C. Ré, "Learning mixed-curvature representations in product spaces," in *International Conference on Learning Representations*, 2018.
- [115] R. Aly, S. Acharya, A. Ossa, A. Köhn, C. Biemann, and A. Panchenko, "Every child should have parents: A taxonomy refinement algorithm based on hyperbolic term embeddings," in *Proceedings of the 57th Annual Meeting of the ACL*, 2019.
- [116] I. Balazevic, C. Allen, and T. Hospedales, "Multi-relational poincaré graph embeddings," in *Advances in Neural Information Processing Systems*, 2019.
- [117] Y. Nagano, S. Yamaguchi, Y. Fujita, and M. Koyama, "A wrapped normal distribution on hyperbolic space for gradient-based learning," *Proceedings of the 36-th International Conference on Machine Learning*, 2019.
- [118] M. Law, R. Liao, J. Snell, and R. Zemel, "Lorentzian distance learning for hyperbolic representations," in *International Conference on Machine Learning*, 2019.
- [119] E. Mathieu, C. Le Lan, C. J. Maddison, R. Tomioka, and Y. W. Teh, "Continuous hierarchical representations with poincaré variational auto-encoders," in *Advances in neural information processing systems*, 2019.
- [120] D. Grattarola, L. Livi, and C. Alippi, "Adversarial autoencoders with constant-curvature latent manifolds," *Applied Soft Computing*, 2019.
- [121] N. Monath, M. Zaheer, D. Silva, A. McCallum, and A. Ahmed, "Gradient-based hierarchical clustering using continuous representations of trees in hyperbolic space," in *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2019.
- [122] Y. Zhang, X. Wang, X. Jiang, C. Shi, and Y. Ye, "Hyperbolic graph attention network," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [123] L. V. Tran, Y. Tay, S. Zhang, G. Cong, and X. Li, "Hyperml: A boosting metric learning approach in hyperbolic space for recommender systems," in *WSDM*, 2020.
- [124] C. Xu and M. Wu, "Learning feature interactions with lorentzian factorization machine," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [125] Y. Zhu, D. Zhou, J. Xiao, X. Jiang, X. Chen, and Q. Liu, "Hypertext: Endowing fasttext with hyperbolic geometry," in *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing*, 2020.
- [126] A. J. Bose, A. Smofsky, R. Liao, P. Panangaden, and W. L. Hamilton, "Latent variable modelling with hyperbolic normalizing flows," *arXiv*, 2020.
- [127] M. Kochurov, S. Ivanov, and E. Burnaev, "Are hyperbolic representations in graphs created equal?" 2020.
- [128] A. Bogatskiy, B. Anderson, J. T. Oeffermann, M. Roussi, D. W. Miller, and R. Kondor, "Lorentz group equivariant neural network for particle physics," *International Conference on Machine Learning (ICML)*, 2020.
- [129] O. Skopek, O.-E. Ganea, and G. Bécigneul, "Mixed-curvature variational autoencoders," *International Conference on Learning Representations (ICLR)*, 2020.
- [130] I. Chami, A. Gu, V. Chatziafratis, and C. Ré, "From trees to continuous embeddings and back: Hyperbolic hierarchical clustering," *Advances in Neural Information Processing Systems*, 2020.
- [131] E. Mathieu and M. Nickel, "Riemannian continuous normalizing flows," *Advances in Neural Information Processing Systems*, 2020.
- [132] R. Sonthalia and A. C. Gilbert, "Tree! i am no tree! i am a low dimensional hyperbolic embedding," *Advances in neural information processing systems*, 2020.
- [133] M. T. Law and J. Stam, "Ultrahyperbolic representation learning," *Advances in Neural Information Processing Systems*, 2020.
- [134] S. Zhu, S. Pan, C. Zhou, J. Wu, Y. Cao, and B. Wang, "Graph geometry interaction learning," *Advances in Neural Information Processing Systems*, 2020.
- [135] J. Ding and A. Regev, "Deep generative model embedding of single-cell rna-seq profiles on hyperspheres and hyperbolic spaces," *Nature communications*, 2021.
- [136] V. Cohen-Addad, V. Kanade, F. Mallmann-Trenn, and C. Mathieu, "Hierarchical clustering: Objective functions and algorithms," *Journal of the ACM (JACM)*.
- [137] O. Yim and K. T. Ramdeen, "Hierarchical cluster analysis: comparison of three linkage measures and application to psychological data," *The quantitative methods for psychology*, 2015.
- [138] S. Dasgupta, "A cost function for similarity-based hierarchical clustering," in *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, 2016.
- [139] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *ACL*, 2019.
- [140] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv*, 2017.
- [141] W. Peng, X. Hong, H. Chen, and G. Zhao, "Learning graph convolutional network for skeleton-based human action recognition by neural searching," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [142] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," in *International Conference on Learning Representations*, 2020.
- [143] J. Lee, Y. Lee, J. Kim, A. Kosirok, S. Choi, and Y. W. Teh, "Set transformer: A framework for attention-based permutation-invariant neural networks," in *International Conference on Machine Learning*. PMLR, 2019.
- [144] M. Henaff, J. Bruna, and Y. LeCun, "Deep convolutional networks on graph-structured data," *arXiv*, 2015.
- [145] S. Yan, Y. Xiong, and D. Lin, "Spatial temporal graph convolutional networks for skeleton-based action recognition," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2018.

- [146] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," 2019.
- [147] D. J. Rezende and S. Mohamed, "Variational inference with normalizing flows," *Proceedings of the 32nd International Conference on Machine Learning*, 2015.
- [148] M. C. Gemici, D. Rezende, and S. Mohamed, "Normalizing flows on riemannian manifolds," *arXiv*, 2016.
- [149] J. Brehmer and K. Cranmer, "Flows for simultaneous manifold learning and density estimation," *Advances in neural information processing systems*, 2020.
- [150] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," *arXiv*, 2016.
- [151] X. Pennec, "Intrinsic statistics on riemannian manifolds: Basic tools for geometric measurements," *Journal of Mathematical Imaging and Vision*, 2006.
- [152] S. Said, L. Bombrun, and Y. Berthoumieu, "New riemannian priors on the univariate normal model," *Entropy*, 2014.
- [153] K. V. Mardia, *Statistics of directional data*. Academic press, 2014.
- [154] T. R. Davidson, L. Falorsi, N. De Cao, T. Kipf, and J. M. Tomczak, "Hyperspherical variational auto-encoders," *34th Conference on Uncertainty in Artificial Intelligence (UAI-18)*, 2018.
- [155] L. Falorsi, P. de Haan, T. Davidson, and P. Forré, "Reparameterizing distributions on lie groups," 2019.
- [156] I. Tolstikhin, O. Bousquet, S. Gelly, and B. Schoelkopf, "Wasserstein auto-encoders," *International Conference on Learning Representations (ICLR)*, 2018.
- [157] K. M. Borgwardt, A. Gretton, M. J. Rasch, H.-P. Kriegel, B. Schölkopf, and A. J. Smola, "Integrating structured biological data by kernel maximum mean discrepancy," *Bioinformatics*, 2006.
- [158] L. Fang, C. Li, J. Gao, W. Dong, and C. Chen, "Implicit deep latent variable models for text generation," *EMNLP*, 2019.
- [159] R. T. Rockafellar *et al.*, "Extension of fenchel'duality theorem for convex functions," *Duke mathematical journal*, 1966.
- [160] B. Dai, H. Dai, N. He, W. Liu, Z. Liu, J. Chen, L. Xiao, and L. Song, "Coupled variational bayes via optimization embedding," in *Advances in Neural Information Processing Systems*, 2018.
- [161] J. Tomczak and M. Welling, "Vae with a vampprior," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2018.
- [162] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," *Advances in neural information processing systems*, 2013.
- [163] B. Yang, W.-t. Yih, X. He, J. Gao, and L. Deng, "Embedding entities and relations for learning and inference in knowledge bases," *International Conference on Learning Representations*, 2015.
- [164] A. Suzuki, Y. Enokida, and K. Yamanishi, "Riemannian transe: Multi-relational graph embedding in non-euclidean space," 2018.
- [165] I. Chami, A. Wolf, D.-C. Juan, F. Sala, S. Ravi, and C. Ré, "Low-dimensional hyperbolic knowledge graph embeddings," in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020.
- [166] I. Vulic, D. Gerz, D. Kiela, F. Hill, and A. Korhonen, "Hyperlex: A large-scale evaluation of graded lexical entailment," *Computational Linguistics*, 2017.
- [167] A. Joulin, É. Grave, P. Bojanowski, and T. Mikolov, "Bag of tricks for efficient text classification," in *Proceedings of the 15th Conference of the European Chapter of the ACL*, 2017.
- [168] M. V. Micic and H. Chu, "Hyperbolic deep learning for chinese natural language understanding," *arXiv*, 2018.
- [169] S. Rendle, "Factorization machines," in *2010 IEEE International Conference on Data Mining*. IEEE, 2010.
- [170] M. Blondel, A. Fujino, N. Ueda, and M. Ishihata, "Higher-order factorization machines," in *Advances in Neural Information Processing Systems*, 2016.
- [171] C.-K. Hsieh, L. Yang, Y. Cui, T.-Y. Lin, S. Belongie, and D. Estrin, "Collaborative metric learning," in *Proceedings of the 26th international conference on world wide web*, 2017.
- [172] L. Zheng, Y. Yang, and A. G. Hauptmann, "Person re-identification: Past, present and future," *arXiv*, 2016.
- [173] J. Snell, K. Swersky, and R. Zemel, "Prototypical networks for few-shot learning," *Advances in neural information processing systems*, 2017.
- [174] P. K. Rubenstein, B. Schoelkopf, and I. Tolstikhin, "On the latent space of wasserstein auto-encoders," *arXiv*, 2018.
- [175] M. Huerta, G. Downing, F. Haseltine, B. Seto, and Y. Liu, "Nih working definition of bioinformatics and computational biology," *US National Institute of Health*, 2000.
- [176] L. Van der Maaten and G. Hinton, "Visualizing data using t-sne," *Journal of machine learning research*, 2008.
- [177] M. Ringnér, "What is principal component analysis?" *Nature biotechnology*, 2008.
- [178] L. McInnes, J. Healy, and J. Melville, "Umap: Uniform manifold approximation and projection for dimension reduction," *arXiv preprint arXiv:1802.03426*, 2018.
- [179] U. Von Luxburg, "A tutorial on spectral clustering," *Statistics and computing*, 2007.
- [180] T. A. Manolio, F. S. Collins, N. J. Cox, D. B. Goldstein, L. A. Hindorff, D. J. Hunter, M. I. McCarthy, E. M. Ramos, L. R. Cardon, A. Chakravarti *et al.*, "Finding the missing heritability of complex diseases," *Nature*, 2009.
- [181] J. Novembre, T. Johnson, K. Bryc, Z. Kutalik, A. R. Boyko, A. Auton, A. Indap, K. S. King, S. Bergmann, M. R. Nelson *et al.*, "Genes mirror geography within europe," *Nature*, 2008.
- [182] G. Bianconi and C. Rahmede, "Emergent hyperbolic network geometry," *Scientific reports*, 2017.
- [183] J. M. Lee, *Riemannian manifolds: an introduction to curvature*. Springer Science & Business Media, 2006.
- [184] B. O'neill, *Semi-Riemannian geometry with applications to relativity*. Academic press, 1983.



Wei Peng is currently a Ph.D. candidate with the Center for Machine Vision and Signal Analysis, University of Oulu, Oulu, Finland. He received the M.S. degree in computer science from the Xiamen University, Xiamen, China, in 2016. His articles have published in mainstream conferences and journals, such as AAAI, ICCV, ACM Multimedia, Transactions on Image Processing. His current research interests include machine learning, affective computing, medical imaging, and human action analysis.



Tuomas Varanka is currently a Ph.D. candidate with the Center for Machine Vision and Signal Analysis, University of Oulu, Oulu, Finland. He received his B.S. and M.S. degree in computer science and engineering from the University of Oulu in 2019 and 2020, respectively. His work has focused on micro-expression recognition. His current research interests include machine learning, and affective computing.



Abdelrahman Mostafa is currently a researcher and a Ph.D. candidate with the Center for Machine Vision and Signal Analysis, University of Oulu, Oulu, Finland. He received the M.S. degree in Artificial Intelligence, University of Oulu, Finland, in 2020. His research interests are Machine Learning, Deep Learning and Computer Vision.



Henglin Shi received the M.S. degree in Artificial Intelligence, University of Oulu, Finland. He is currently a researcher and PhD student, University of Oulu, Finland. His articles have published in mainstream conferences and journals, such as BMVC, TMM and TIP. His research interests are Machine Learning, Deep Learning and Computer Vision.



Guoying Zhao received the Ph.D. degree (2005) in computer science from the Chinese Academy of Sciences, China. She is currently an Academy Professor with the Center for Machine Vision and Signal Analysis, University of Oulu, Finland. Her work has focused on affective computing and machine learning. She is IEEE fellow, IAPR fellow and ELLIS member, and associate editor for Pattern Recognition, IEEE Transactions on Circuits and Systems for Video Technology, and Image and Vision Computing Journals.