

# Hyperbolic Learning: Theory and Applications

---

Presenters: Pengxiang Li<sup>1</sup>, Peilin Yu<sup>1</sup>, Yangkai Xue<sup>1</sup>, Yuwei Wu<sup>1</sup>, Zhi Gao<sup>1</sup>

1 Hyperbolic geometry

2 Hyperbolic learning

# Hyperbolic geometry

- Why hyperbolic?
- Riemannian manifold
- Hyperbolic Models
- Basic Operation
- $\delta$ -Hyperbolicity

# Hyperbolic geometry

Why hyperbolic?

Riemannian manifold

Hyperbolic Models

Basic Operation

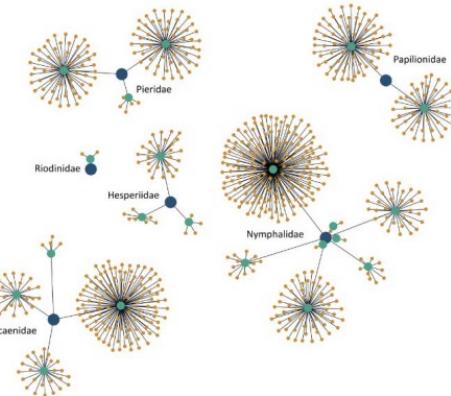
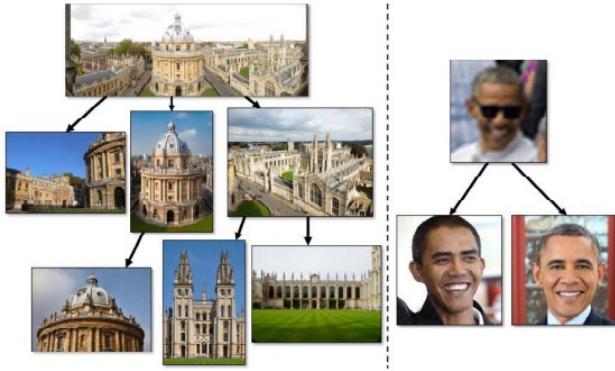
$\delta$ -Hyperbolicity

# Why hyperbolic?

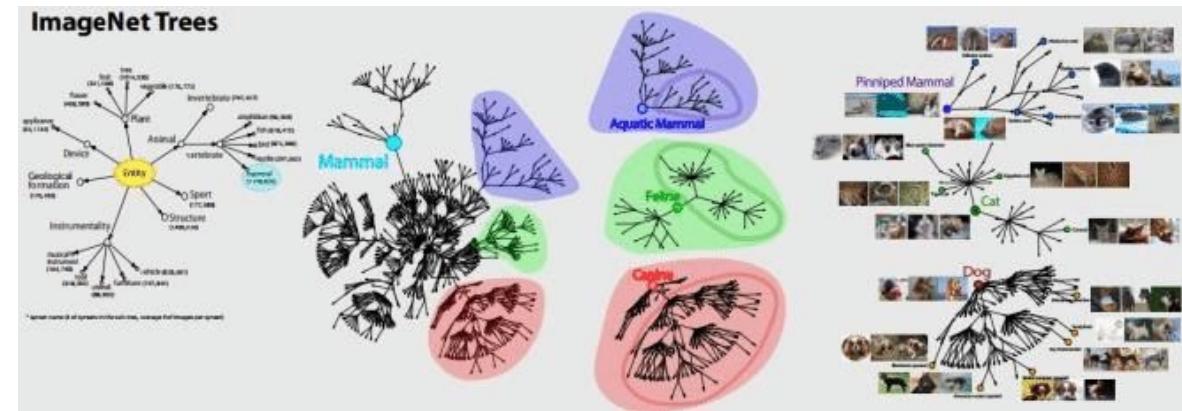
## 1. High modeling capacity

- Hierarchical/tree-like structures

Hierarchical data



Hierarchical knowledge [3]



Visual hierarchies [1] Semantic hierarchies [2]

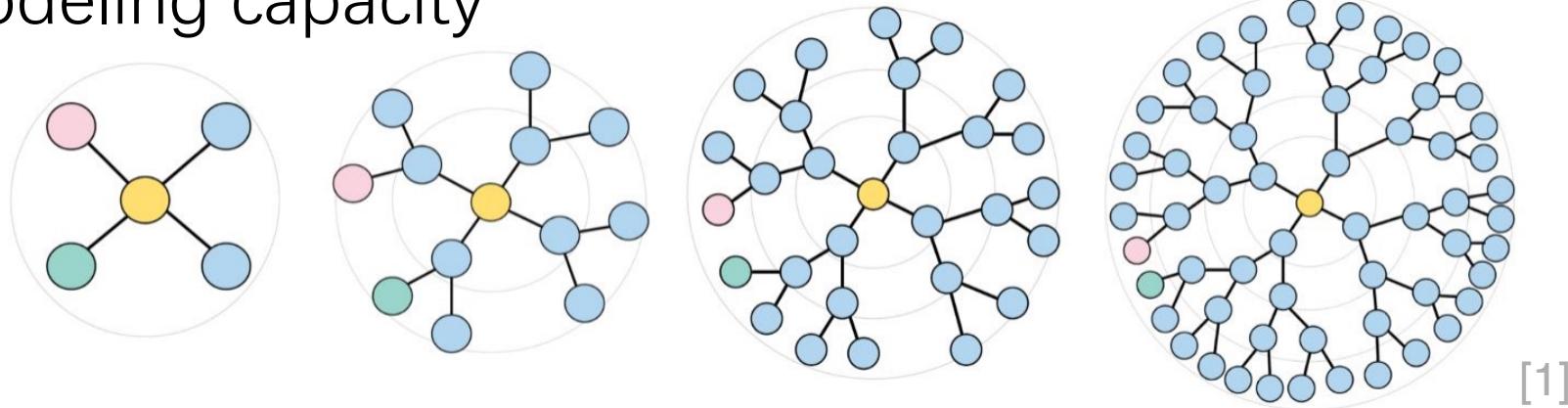
[1] Hyperbolic Image Embeddings, Mirvakhabova et al. CVPR 2020.

[2] Hierarchical Image Classification using Entailment Cone Embeddings. Dhall et al. CVPRW 2020.

[3] Hyperbolic Representation Learning for Computer Vision. Mettes et al. ECCV 2022

# Why hyperbolic?

## 1. High modeling capacity



[1]

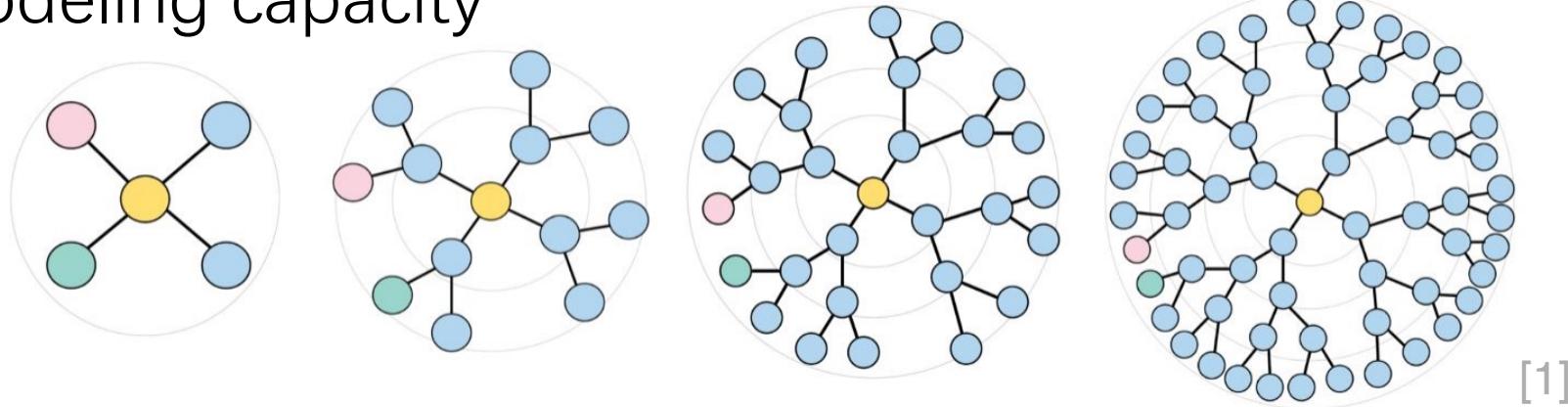
Hierarchies grow **exponentially** in depth.

Euclidean space grows linearly with norm.

Distances in hyperbolic space grow **exponentially**.

# Why hyperbolic?

## 1. High modeling capacity



Hierarchies grow **exponentially** in depth.

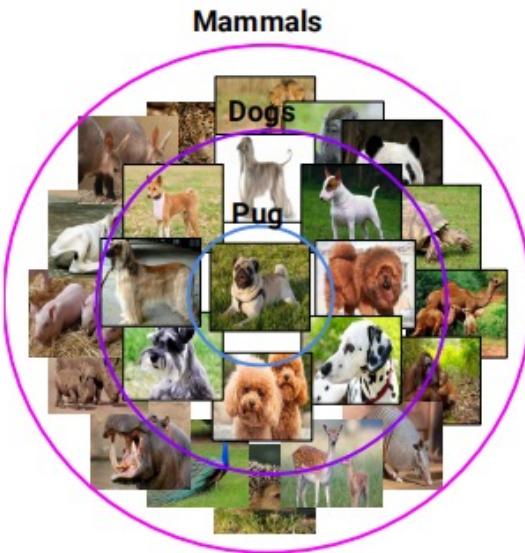
Euclidean space grows linearly with norm.

Distances in hyperbolic space grow **exponentially**.

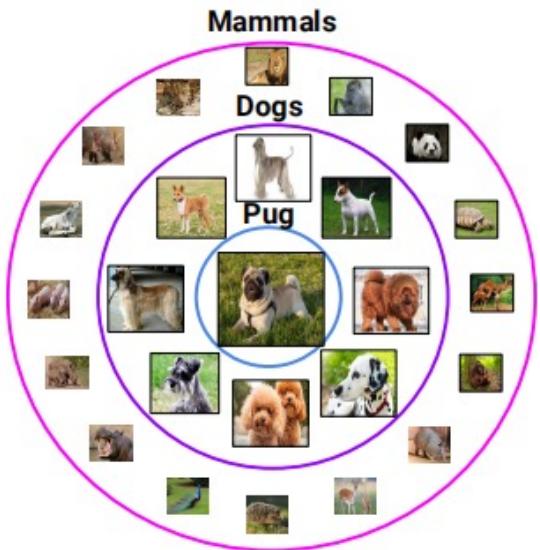
Hyperbolic spaces **more closely resemble hierarchical structures** than Euclidean space.

# Why hyperbolic?

## 2. More space than Euclidean geometry

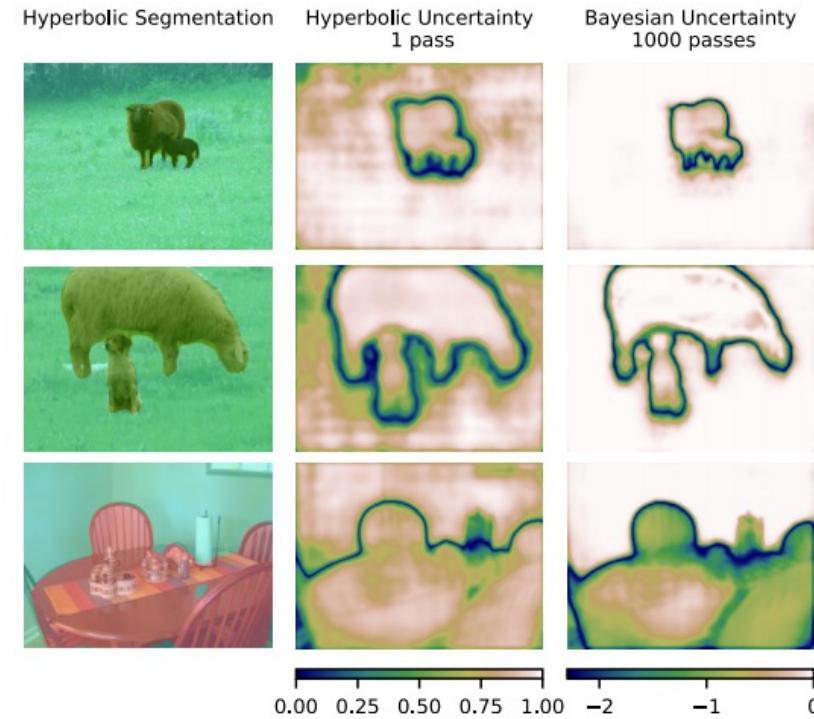


Euclidean Space [1]



Hyperbolic Space [1]

## 3. Natural measure of uncertainty



Hyperbolic vs Euclidean uncertainty [2]

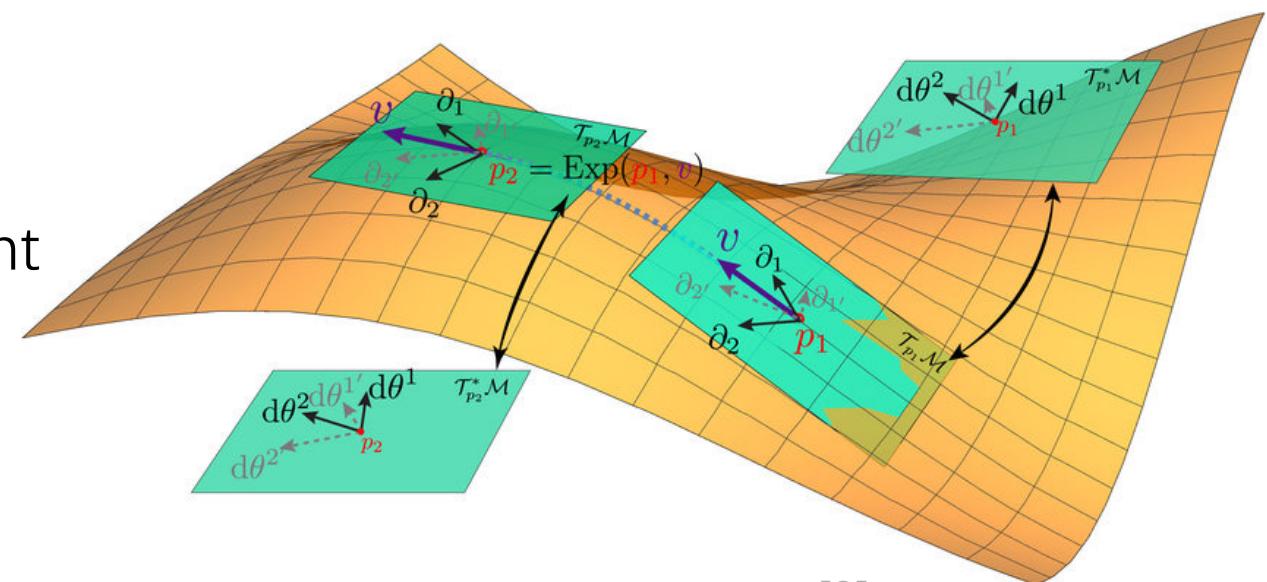
# Hyperbolic geometry

- { Why hyperbolic?
- Riemannian manifold
- Hyperbolic Models
- Basic Operation
- $\delta$ -Hyperbolicity

# Riemannian manifold

In differential geometry, a Riemannian manifold or Riemannian space  $(M, g)$ , is a real, **smooth** manifold  $M$  equipped with a positive-definite inner product  $g_p$  on the tangent space  $T_p M$  at each point  $p$ .<sup>[1]</sup>

The tangent bundle of a smooth manifold  $M$  assigns to each point  $p$  of  $M$  a vector space  $T_p M$  called the tangent space of  $M$  at  $p$ .<sup>[1]</sup>



Riemannian manifold [2]

[1] [https://en.wikipedia.org/wiki/Riemannian\\_manifold](https://en.wikipedia.org/wiki/Riemannian_manifold)

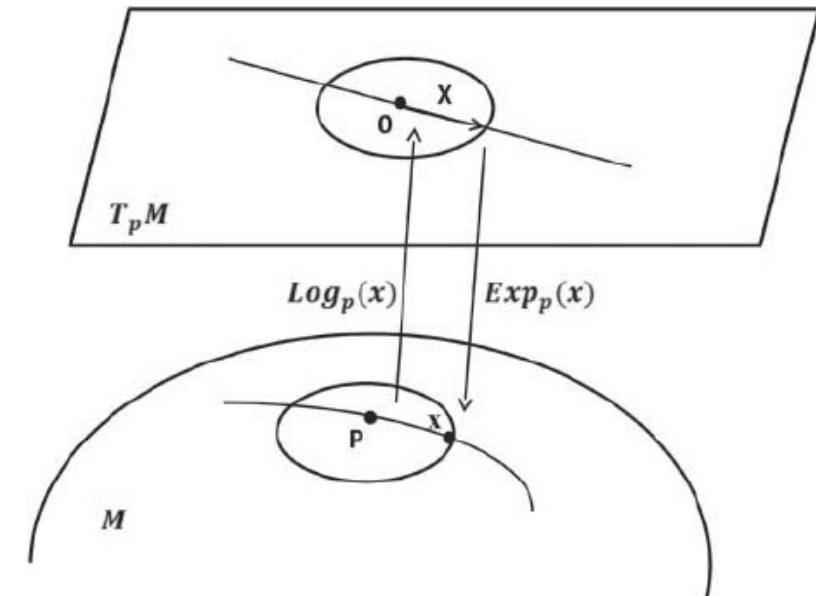
# Riemannian manifold

A Riemannian metric (by its definition) assigns to each  $p$  a positive-definite inner product  $g_p: T_p M \times T_p M \rightarrow \mathbb{R}$ . The smooth manifold  $M$  endowed with this metric  $g$  is a Riemannian manifold, denoted  $(M, g)$ . [1]

$$Exp_p : T_p M \rightarrow M$$

$$Log_p : M \rightarrow T_p M$$

As a Riemannian manifold, a **hyperbolic space** is a complete contractible **smooth manifold** of constant **negative Riemann curvature** (equivalently, Gaussian curvature).

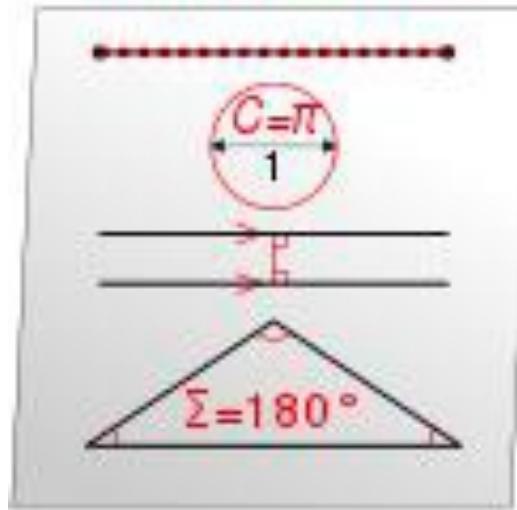


Log and exp mapping [2]

[1] [https://en.wikipedia.org/wiki/Riemannian\\_manifold](https://en.wikipedia.org/wiki/Riemannian_manifold)

[2] Nodehi, Anahita & Golalizadeh, Mousa & Heydari, Abbas. (2015). Dihedral angles principal geodesic analysis using nonlinear statistics. Journal of Applied Statistics.

# Euclidean Space vs Non-Euclidean Space



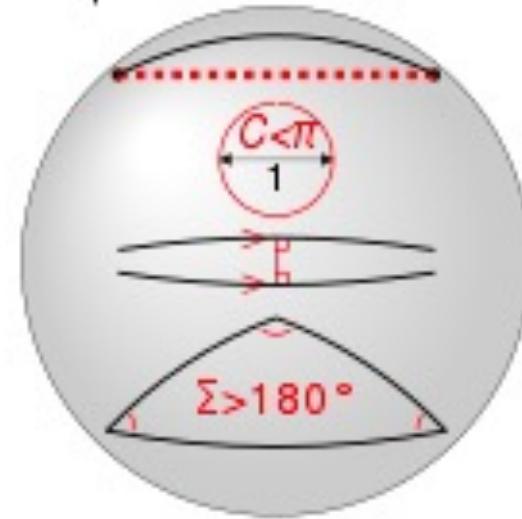
Euclidean

Curvature = 0

Circumference/Diameter =  $\pi$

One parallel line

Sum of angles in a  $\Delta$  =  $180^\circ$



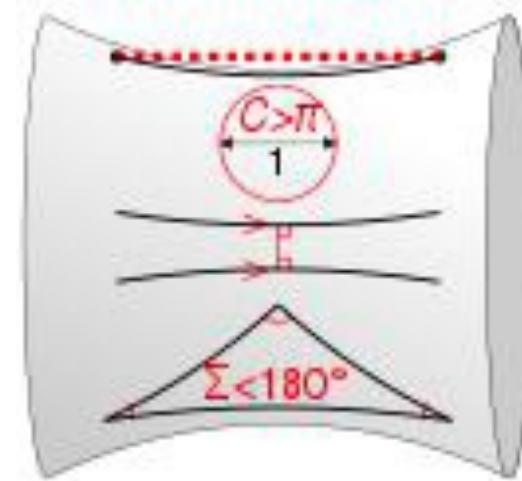
Spherical

Curvature > 0

Circumference/Diameter >  $\pi$

No parallel lines

Sum of angles in a  $\Delta$  >  $180^\circ$



Hyperbolic

Curvature < 0

Circumference/Diameter <  $\pi$

Infinitely many parallel lines

Sum of angles in a  $\Delta$  <  $180^\circ$

# Geodesic in Different Geometries



**Geodesic:** Some kind of shortest distance between two points on the manifold. Applies to graphs too.

Hyperbolic  
geodesic

$$\lambda(t) = \cosh(t\sqrt{\kappa})x + \sinh(t\sqrt{\kappa})v \quad \lambda''(t) - \kappa\lambda(t) = 0$$

Euclidean  
geodesic

$$\lambda(t) = x + tv \quad \lambda''(t) = 0$$

# Euclidean Space vs Hyperbolic Space

Five postulates for plane geometry.

## Euclidean

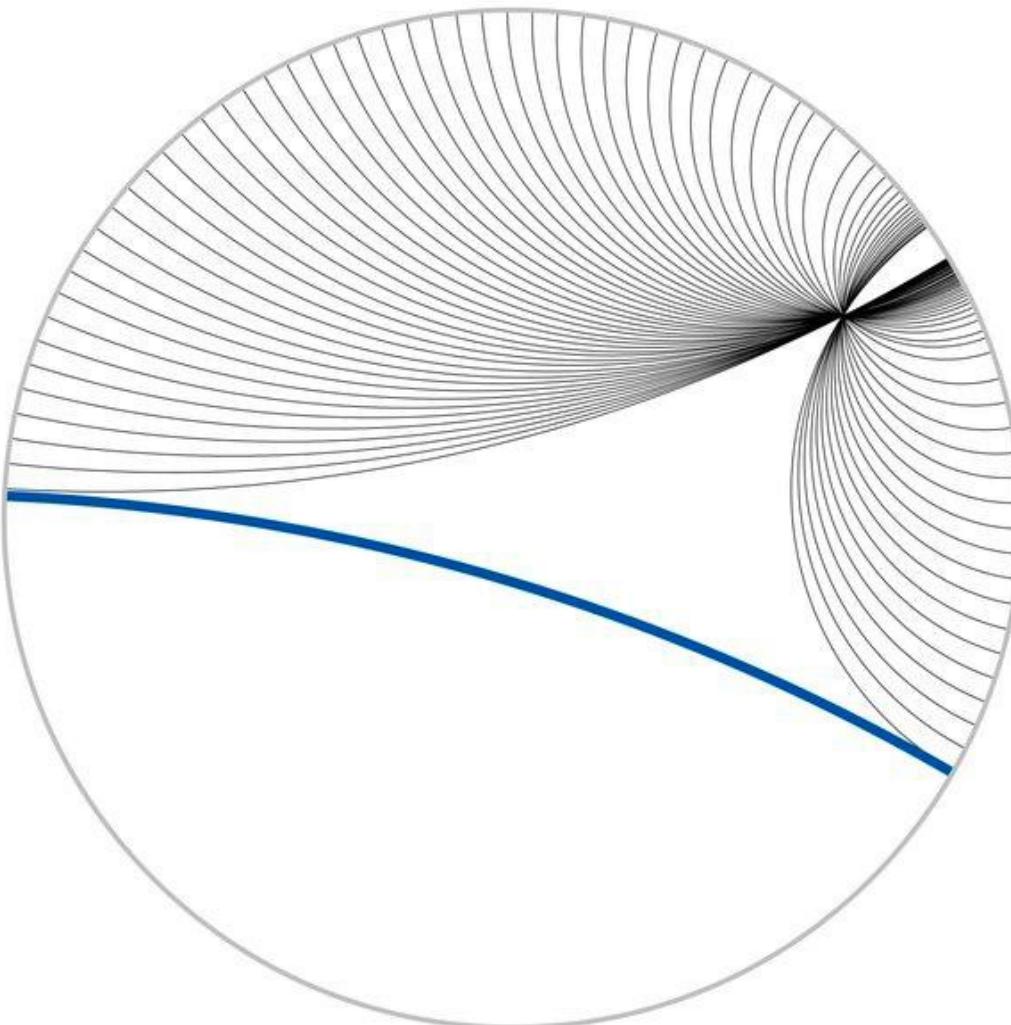
1. Each pair of points can be joined by one and only one straight line segment.
2. Any straight line segment can be indefinitely extended in either direction.
3. There is exactly one circle of any given radius with any given center.
4. All right angles are congruent to one another.
5. **(Parallel Postulate)** In a plane, through a point not on a given straight line, at most one line can be drawn that never meets the given line.

## Hyperbolic

1. Each pair of points can be joined by one and only one straight line segment.
2. Any straight line segment can be indefinitely extended in either direction.
3. There is exactly one circle of any given radius with any given center.
4. All right angles are congruent to one another.
5. **(Hyperbolic Alternative)** Given a line and a point not on it, there is **more than** one line going through the given point that does not intersect the given line.

# Euclidean Space vs Hyperbolic Space

Five postulates for plane geometry.



## Hyperbolic

1. Each pair of points can be joined by one and only one straight line segment.
2. Any straight line segment can be indefinitely extended in either direction.
3. There is exactly one circle of any given radius with any given center.
4. All right angles are congruent to one another.
5. (**Hyperbolic Alternative**) Given a line and a point not on it, there is **more than** one line going through the given point that does not intersect the given line.

# Hyperbolic geometry

- Why hyperbolic?
- Riemannian manifold
- Hyperbolic Models
- Basic Operation
- $\delta$ -Hyperbolicity

# Hyperbolic function



$$\sinh x = \frac{e^x - e^{-x}}{2}$$

$$\operatorname{arsinh} x = \ln(x + \sqrt{x^2 + 1})$$

$$\cosh x = \frac{e^x + e^{-x}}{2}$$

$$\operatorname{arcosh} x = \ln(x \pm \sqrt{x^2 - 1})$$

$$\tanh x = \frac{\sinh x}{\cosh x} = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1}$$

$$\operatorname{artanh} x = \frac{1}{2} \ln \left( \frac{1+x}{1-x} \right)$$

# The five isometric models

$\mathbb{L}$ : Lorentz Model

$\mathbb{B}$ : Poincaré Model

$\mathbb{K}$ : Klein Model

$\mathbb{H}$ : Poincaré Half Model

$\mathbb{J}$ : Hemisphere model

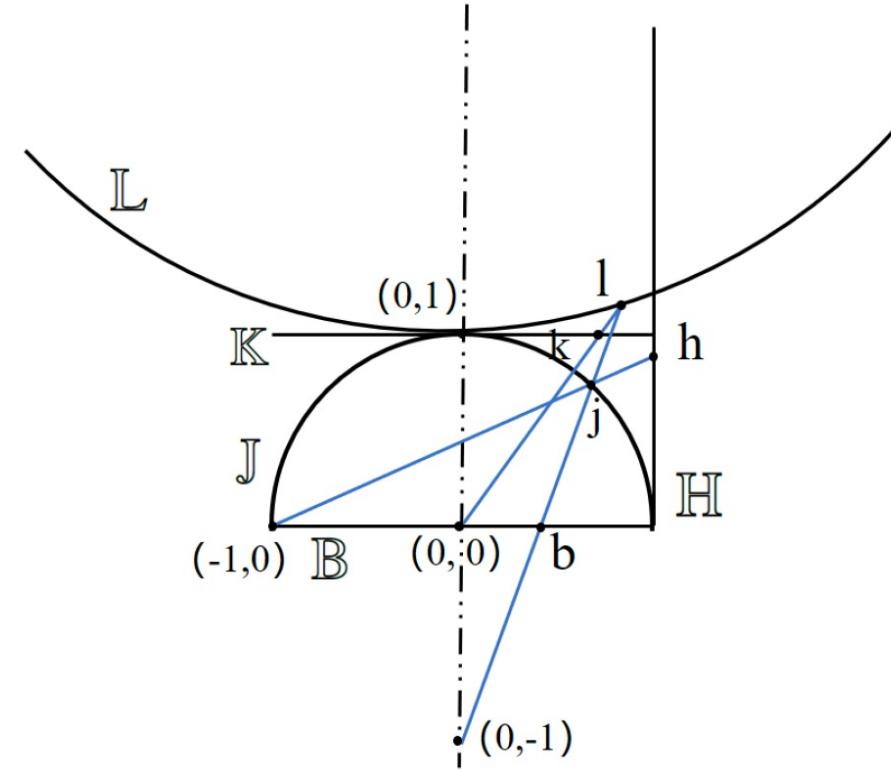
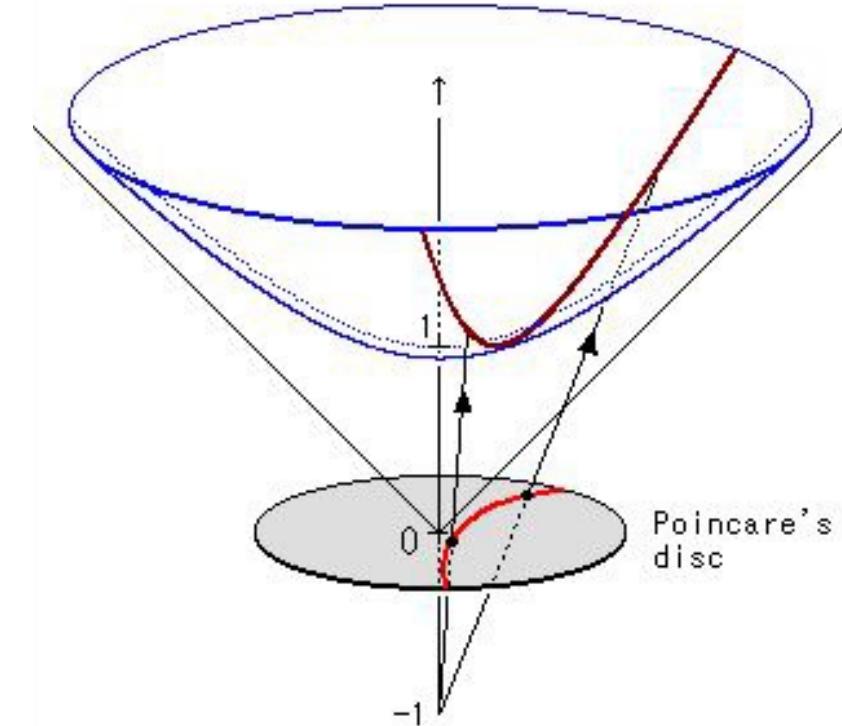


Illustration of five isometric model[1]

# Isometric Models

Lorentz Model  $\Leftrightarrow$  Poincaré Model

$$x = (x_0, \dots, x_n) \in \mathbb{L}^n \Leftrightarrow \left( \frac{x_1}{1+x_0}, \dots, \frac{x_n}{1+x_0} \right) \in \mathbb{B}^n$$



# Isometric Models

Poincaré Model  $\Leftrightarrow$  Poincaré Half Model

$$x = (x_0, \dots, x_{n-1}) \in \mathbb{B}^n \Leftrightarrow \left( \frac{1 - \|x\|^2}{1 + 2x_0 + \|x\|^2}, 2x_1, \dots, 2x_{n-1} \right) \in \mathbb{H}^n$$

Lorentz Model  $\Leftrightarrow$  Klein Model

$$x = (x_0, \dots, x_n) \in \mathbb{L}^n \Leftrightarrow \left( \frac{x_1}{x_0}, \dots, \frac{x_n}{x_0} \right) \in \mathbb{K}^n$$

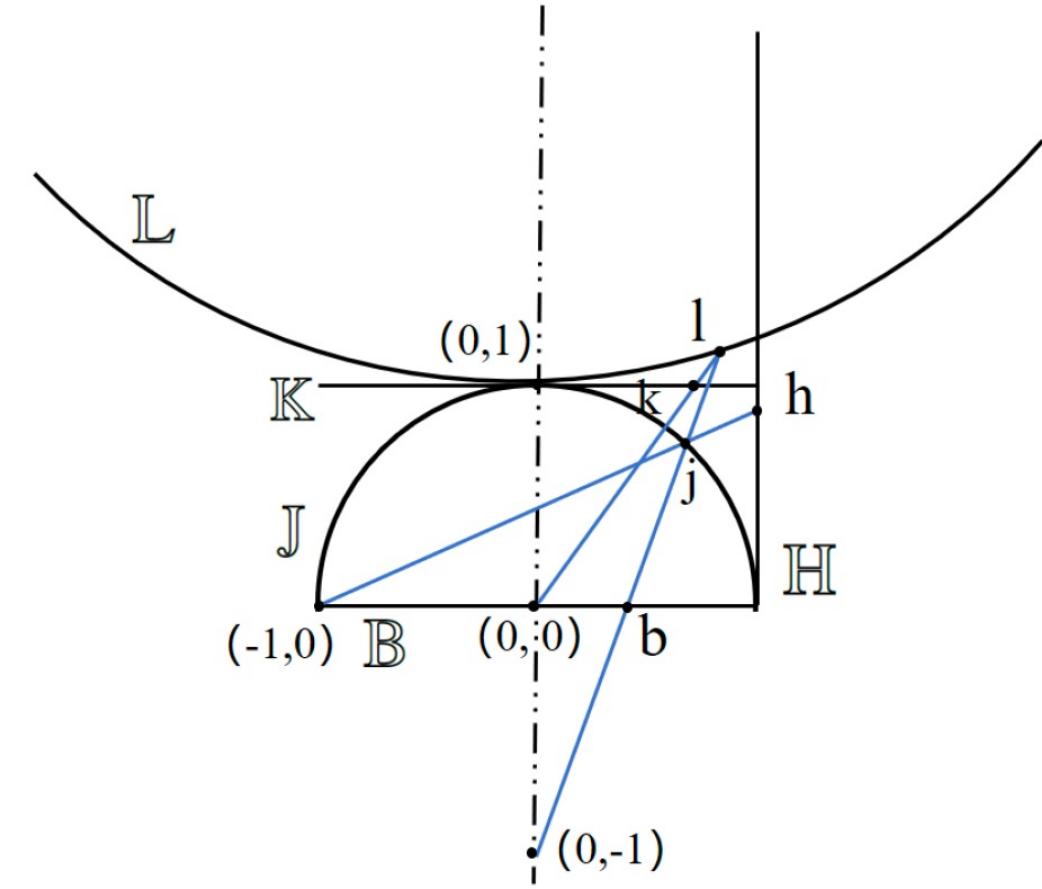
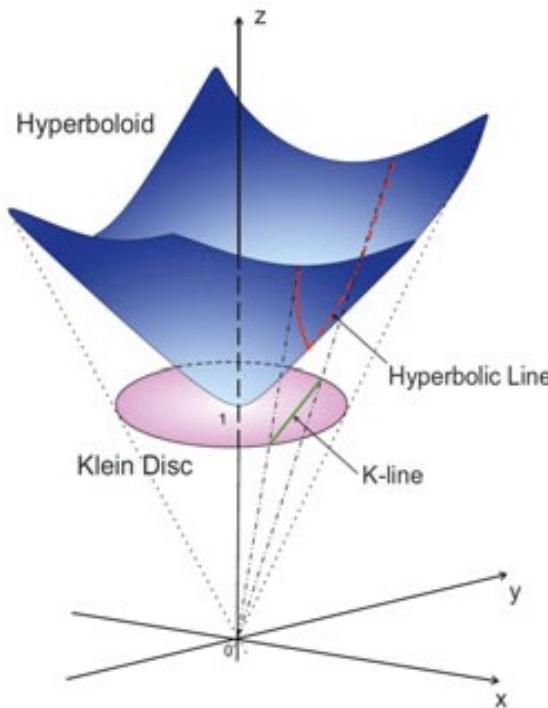


Illustration of five isometric model [1]

# Isometric Models

Poincaré Model  $\Leftrightarrow$  Klein Model

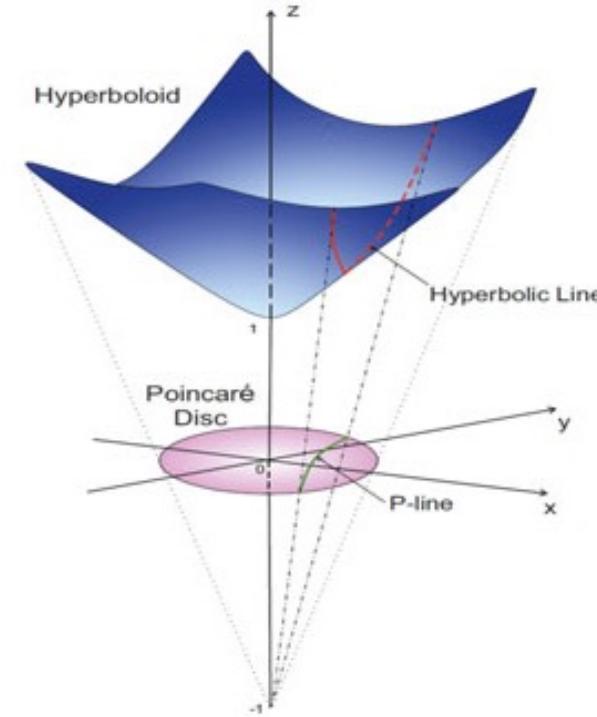
Lorentz Model



[1]

Klein Model

Lorentz Model



[1]

Poincaré Model

# The five isometric models

$\mathbb{L}$ : Lorentz Model

$\mathbb{B}$ : Poincaré Model

$\mathbb{K}$ : Klein Model

$\mathbb{H}$ : Poincaré Half Model

$\mathbb{J}$ : Hemisphere model

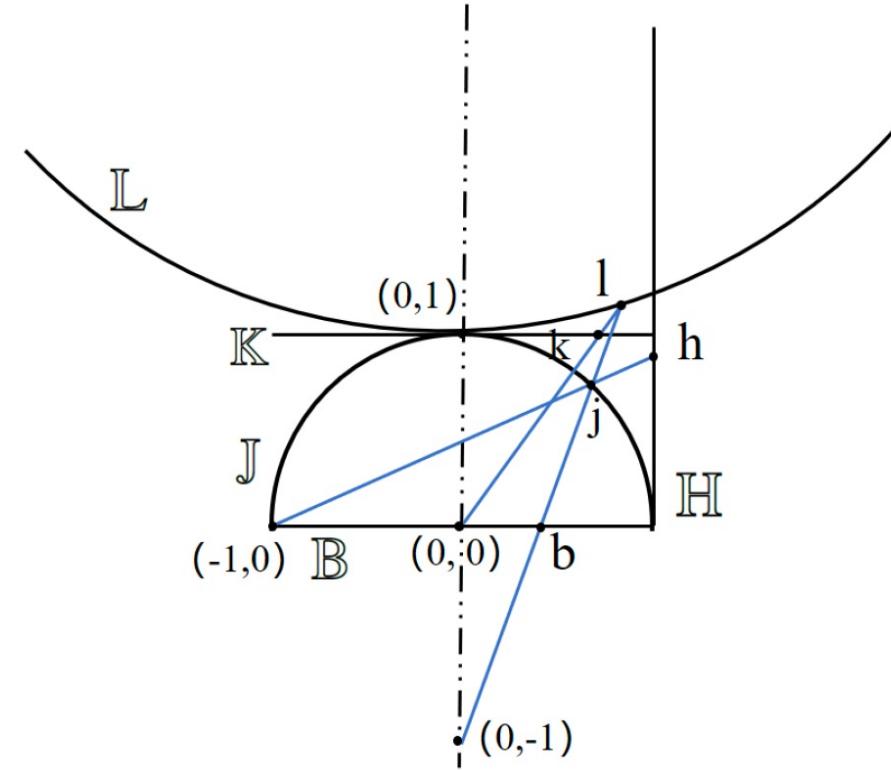


Illustration of five isometric model[1]

# Lorentz Model

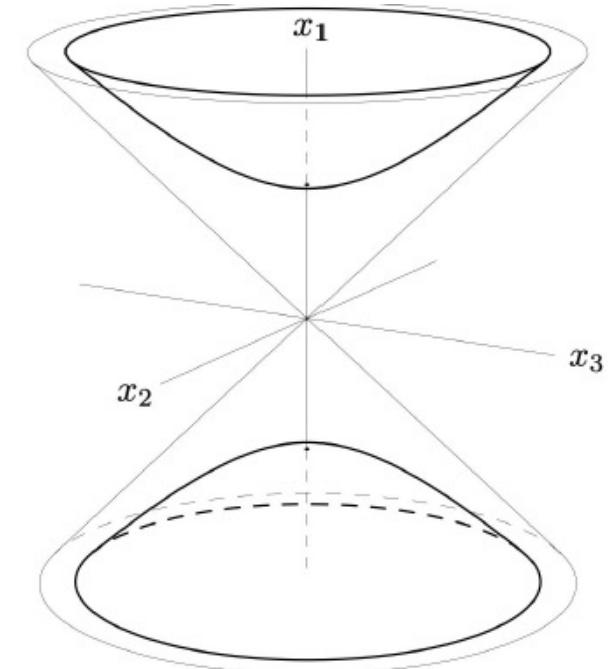
The Lorentz model  $\mathbb{L}^n$  of an  $n$  dimensional hyperbolic space is a manifold embedded in the  $n + 1$  dimensional Minkowski space.

$$\mathbb{L}^n = \{x = (x_0, \dots, x_n) \in \mathbb{R}^{n+1} : \langle x, x \rangle_{\mathbb{L}} = -1, x_0 > 0\},$$

in which  $\langle x, x \rangle_{\mathbb{L}}$  is the Lorentzian inner product

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

where  $\mathbf{g}^{\mathbb{L}}$  is a diagonal matrix with entries of 1s, except for the first element being -1.



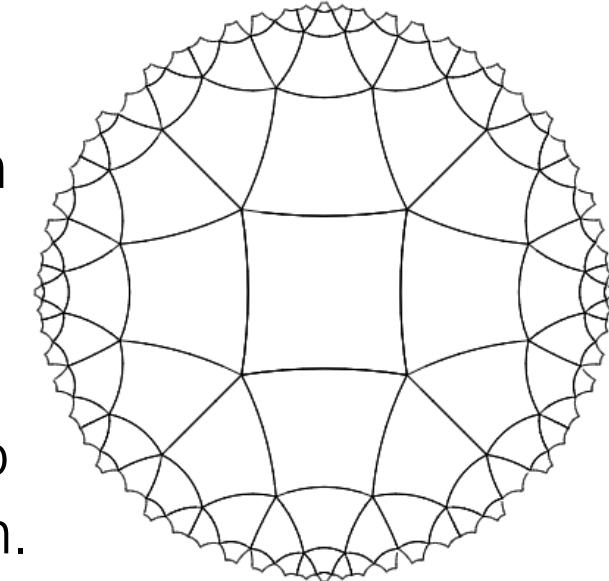
$$d(x, y) = \operatorname{arcosh}(-\langle x, y \rangle_{\mathbb{L}}).$$

# Poincaré Model

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

A Poincaré model is a manifold equipped with a Riemannian metric  $\mathbf{g}^B = \lambda_x^2 \mathbf{g}^B$ , where  $\lambda_x = \frac{2}{1 - \|x\|^2}$ .

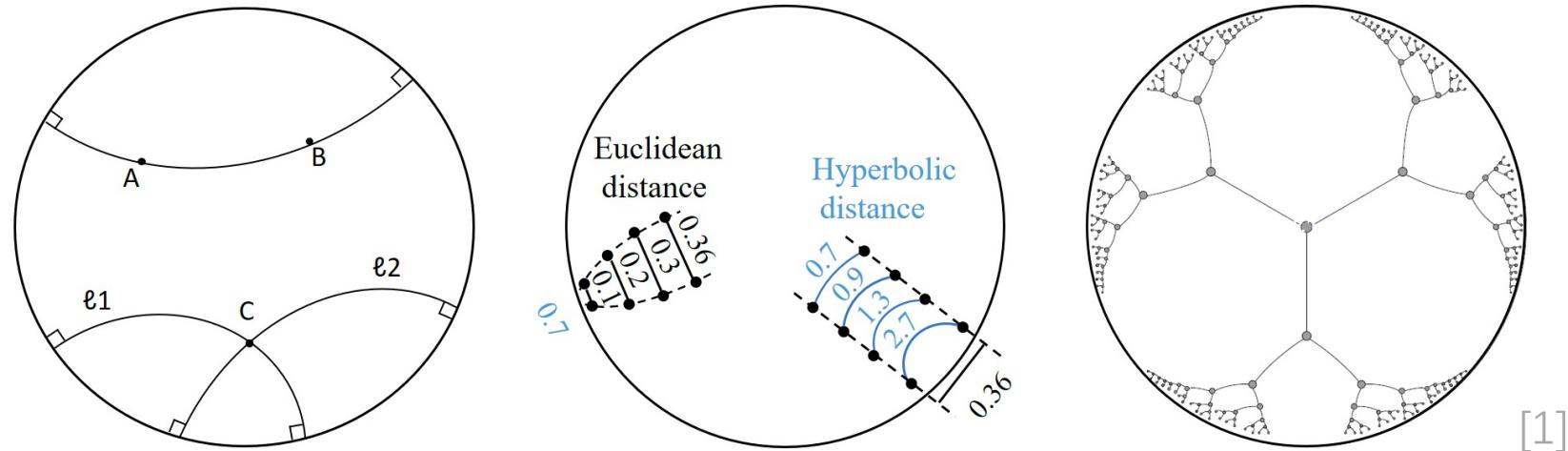
The Poincaré model  $\mathbb{B}^n$  is obtained by mapping  $x \in \mathbb{L}^{n+1}$  to the hyperplane  $x_0 = 0$ , using rays emanating from the origin.



# Poincaré Model

The distance between  $p, q \in \mathbb{B}^n$  is defined as

$$d(p, q) = \frac{1}{\sqrt{\kappa}} \operatorname{arcosh} \left( 1 + \frac{2|p - q|^2}{(1 - |p|^2)(1 - |q|^2)} \right)$$



# Most commonly used two models.



## Lorentz Model

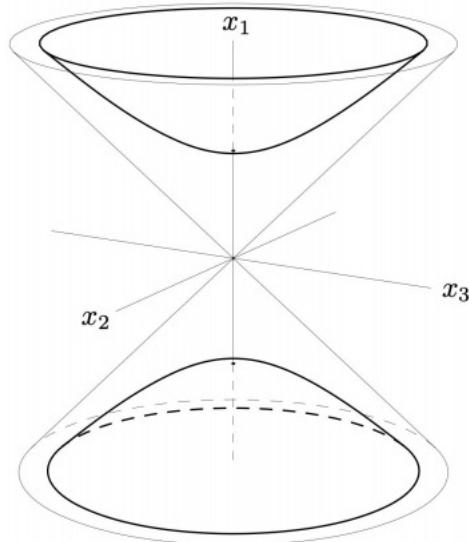
- + has linear structure (Translations and other isometries are linear maps)
- + simpler formulas for geodesics, distance, etc.
- needs one redundant dimension
- not as easy to visualize

## Poincaré Model

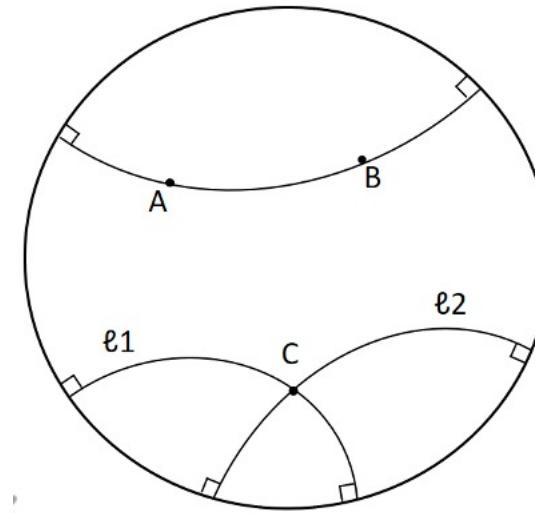
- + visually attractive
- + conformal model (Euclidean angle = hyperbolic angle)
- + geodesic lines are Euclidean circles or line segments
- no linear structure
- formulas for geodesics, distance, etc. not as nice, but: gyrovectorspace calculus

# Most commonly used two models.

Lorentz Model



Poincaré Model



The Klein model  $\mathbb{K}^n$  is obtained by mapping  $x \in \mathbb{L}^{n+1}$  to the hyperplane  $x_0 = 1$ , using rays emanating from the origin.

$$\mathbb{K}^n = \{x \in \mathbb{R}^n : ||x|| < 1\},$$

The distance between  $x, y \in \mathbb{K}^n$  is

$$d(x, y) = \text{arcosh} \left( 1 + \frac{(y_0 - x_0)^2 + (y_1 - x_1)^2}{2x_1 y_1} \right).$$

A Poincaré half model is a manifold equipped with a Riemannian metric  $(\mathbb{H}^n, g^H)$ , where

$$\mathbb{H}^n = \{x \in \mathbb{R}^n : x_n > 0\} \quad g^H = \frac{g^E}{x_n^2}.$$

The Poincaré half model can be obtained by taking the inverse of the Poincaré model and the distance for  $x, y \in \mathbb{H}^n$  is

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

# Hemisphere Model



The hemisphere model is generally used to visualize the transformation between various models.

$$\mathbb{J}^n = \{x = (x_0, \dots, x_n) \in \mathbb{R}^{n+1} : \|x\| = 1, x_0 > 0\},$$

# Hyperbolic geometry

- Why hyperbolic?
- Riemannian manifold
- Hyperbolic Models
- Basic Operation
- $\delta$ -Hyperbolicity

# Basic Operation in Hyperbolic Space



- Exponential map
- Logarithmic map
- Möbius addition
- Möbius scalar multiplication
- Möbius vector multiplication
- Einstein midpoint

To keep the hyperbolicity of the output.

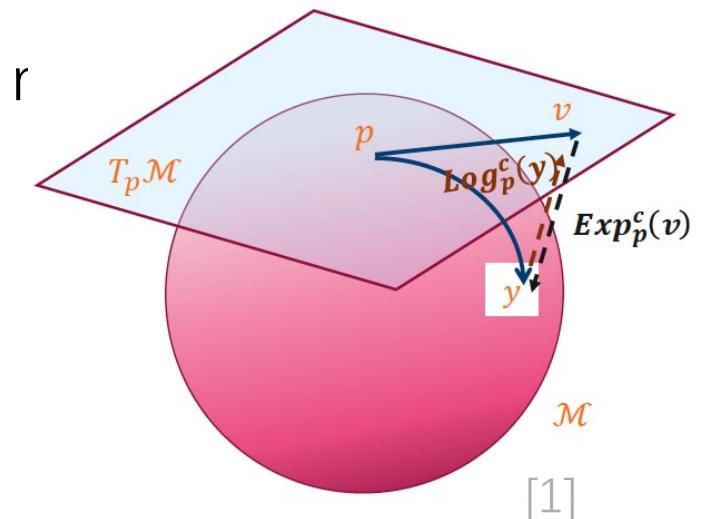
# Exponential map

The exponential map takes a vector  $v \in T_x \mathcal{M}$  of a point  $x \in \mathcal{M}$  to a point on the manifold  $\mathcal{M}$ .

$$\text{Exp}_x : T_x \mathcal{M} \rightarrow \mathcal{M}$$

For a vector  $v \in T_x \mathcal{M}$  in the tangent space, the exponential  $r$  is defined as

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



# Logarithmic map



The logarithmic map is the inverse of the aforementioned exponential map, which projects a point  $y \in \mathcal{M}$  on the manifold to the tangent space of another point  $x \in \mathcal{M}$ .

$$\text{Log}_x : \mathcal{M} \rightarrow \mathcal{T}_x \mathcal{M}$$

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

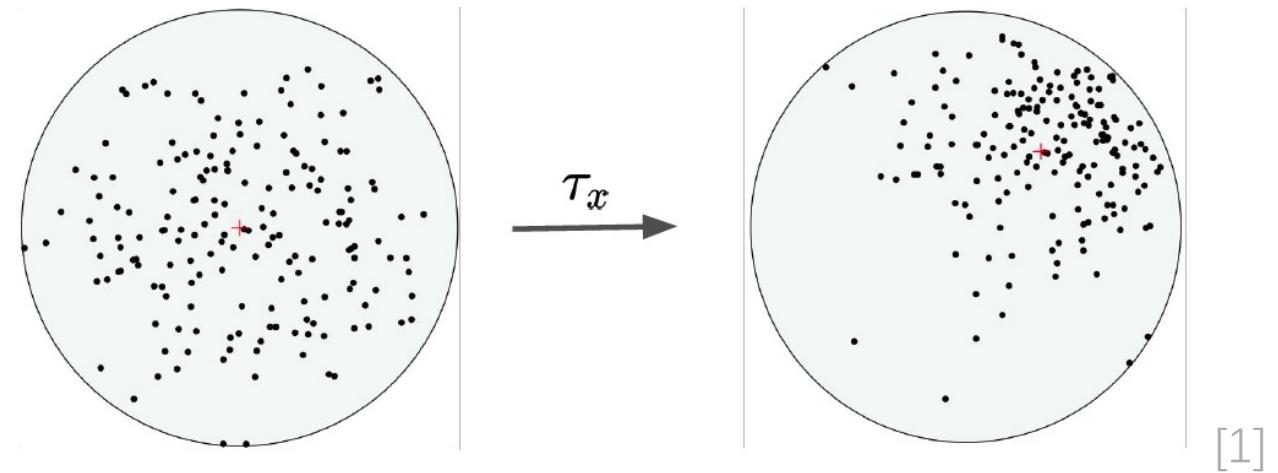
$$\text{Log}_x(y) = \frac{2}{\lambda_x} \operatorname{artanh} (|| -x \oplus y ||) \frac{-x \oplus y}{|| -x \oplus y ||}.$$

# Möbius addition

In the Gyrovector space, the Möbius addition  $\oplus$  for  $x$  and  $y$  in model  $\mathbb{B}$  is defined as

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

$x \oplus y$  will recover to  $x + y$  when the curvature goes to zero. the Möbius subtraction  $\ominus$  is simply defined as:  $x \ominus y = x \oplus (-y)$ .



# Möbius scalar multiplication



In the Gyrovector space, the Möbius scalar multiplication  $\otimes$  for  $x$  in model  $\mathbb{B}^n$  is defined as

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

where  $r$  is a scalar factor.

In addition, the Möbius scalar multiplication can also be conducted in the tangent space by using the exponential and logarithmic maps.

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

# Möbius vector multiplication



In the Gyrovector space, the Möbius vector multiplication  $M^\otimes(x)$  for metric  $M$  and  $x$  in model  $\mathbb{B}^n$  is defined as

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

# Mean in hyperbolic space



Tangential aggregation. Map  $x$  to tangent space, calculate the mean and back to hyperbolic space.

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

Lou et al. NIPS 2019

# Einstein midpoint

Einstein midpoint is an extension of the mean operation to hyperbolic spaces, which has the most concise form in the Klein coordinates.

$$\mu = \frac{\sum_{i=1}^N \gamma_i x^{(i)}}{\sum_{i=1}^N \gamma_i}, \quad \gamma_i = \frac{1}{\|x^{(i)}\|^2}$$

in which the  $x_i$  is the i-th sample represented using coordinates in Klein model. The mean in Poincaré model can be conducted from the results in Klein model by

$$\mu_{\mathbb{B}} = \frac{\mu}{1 + \sqrt{1 - \|\mu\|^2}}.$$

# Einstein midpoint

A closed-form expression for Poincaré model to compute the average (midpoint) in the Gyrovector spaces.

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

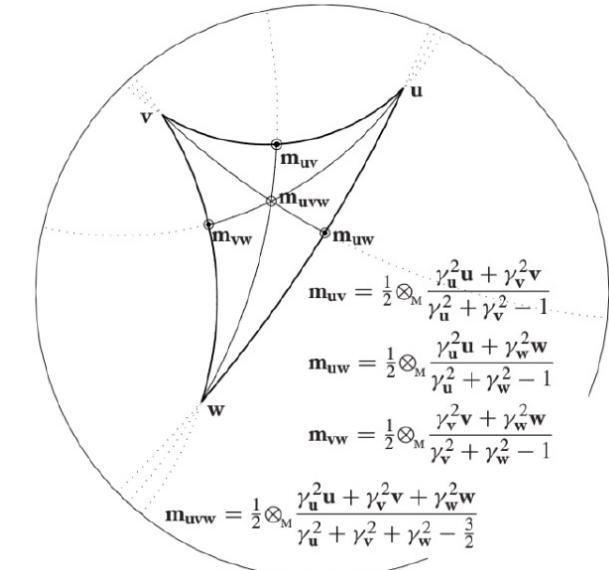


Illustration of the Möbius gyromidpoint. [1]

with  $\alpha = (\alpha_1, \dots, \alpha_N)$  as the weights, and Lorentz factor as  $\gamma_i = \frac{2}{||x^{(i)}||^2}$

# Hyperbolic geometry

- Why hyperbolic?
- Riemannian manifold
- Hyperbolic Models
- Basic Operation
- $\delta$ -Hyperbolicity

# $\delta$ -Hyperbolicity

Gromov  $\delta$ -hyperbolicity is used to evaluate the hyperbolicity of a dataset/space  $X$ .

Gromov product

$$(y, z)_x = \frac{1}{2}(d(x, y) + d(x, z) - d(y, z))$$

$\delta$  is defined as the minimal value such that the following four-point condition holds for all points  $x, y, z, w \in X$ :

$$(x, z)_w \geq \min((x, y)_w, (y, z)_w) - \delta$$

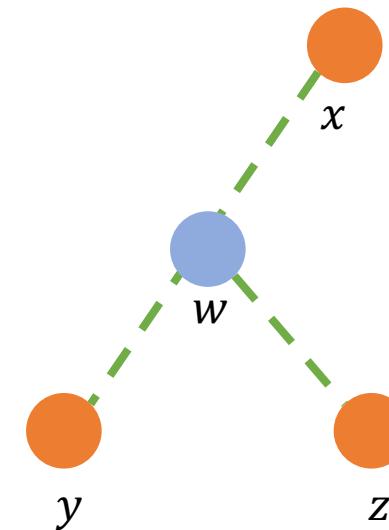
# $\delta$ -Hyperbolicity

Gromov  $\delta$ -hyperbolicity is used to evaluate the **hyperbolicity** of a dataset/space  $X$ .

Gromov product  $(y, z)_x = \frac{1}{2}(d(x, y) + d(x, z) - d(y, z))$

$\delta$  is defined as the minimal value such that the following four-point condition holds for all points  $x, y, z, w \in X$ :

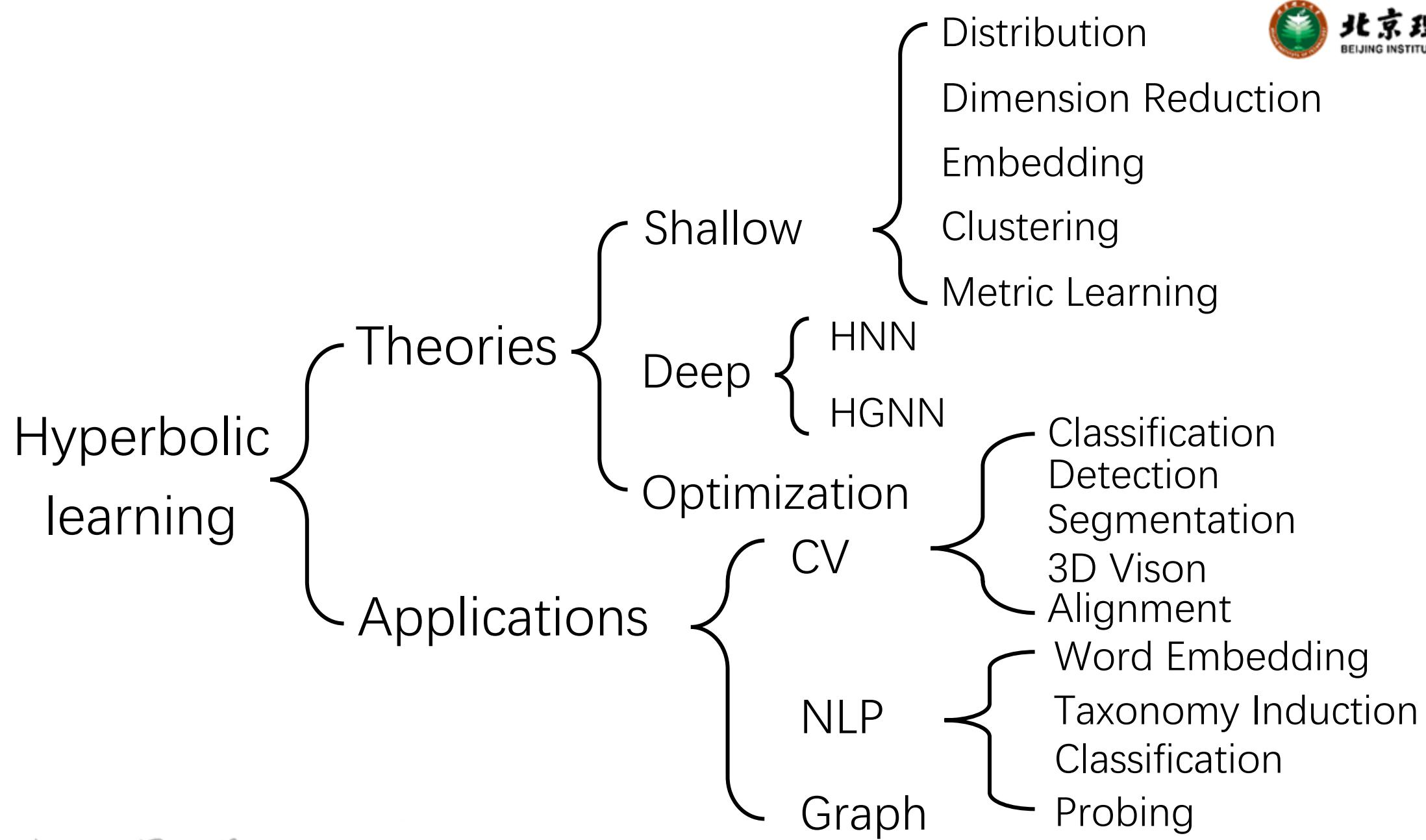
$$(x, z)_w \geq \min((x, y)_w, (y, z)_w) - \delta$$

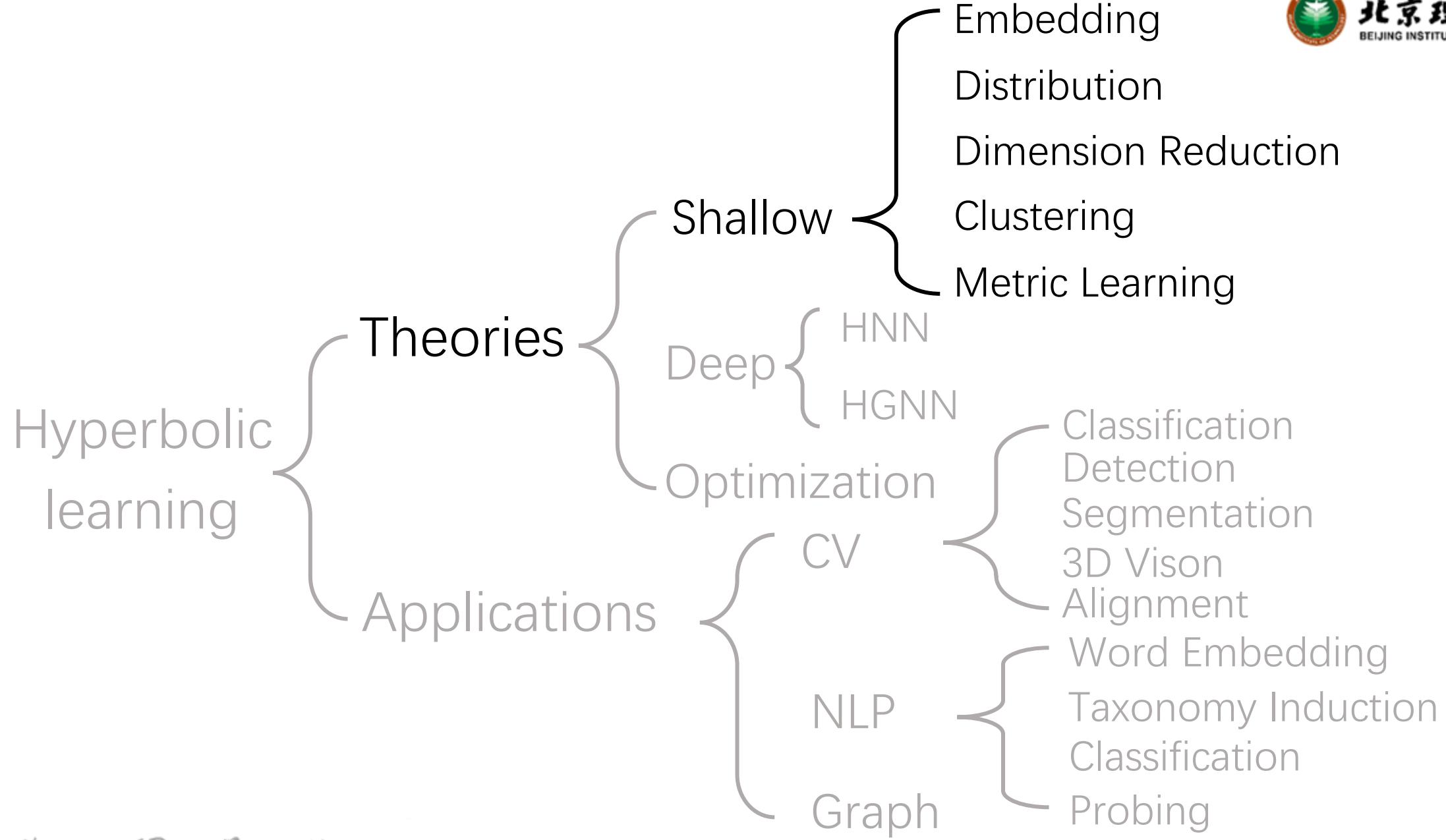


Encoder	Dataset			
	CIFAR10	CIFAR100	CUB	<i>MiniImageNet</i>
Inception v3 [49]	0.25	0.23	0.23	0.21
ResNet34 [14]	0.26	0.25	0.25	0.21
VGG19 [42]	0.23	0.22	0.23	0.17

The relative delta  $\delta_{rel}$  values calculated for different datasets. [1]

# Hyperbolic learning





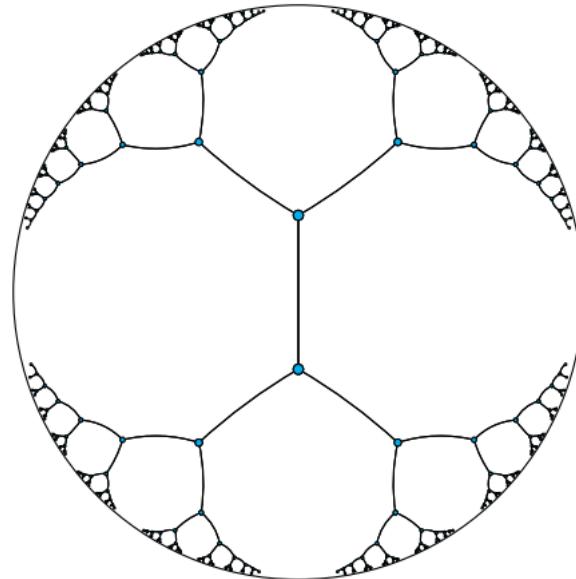
# Theories

Shallow Learning

## Why do we model embedding hierarchies in hyperbolic space?

- Hierarchies grow exponentially in depth, Euclidean space grows linearly with norm.
- The exponential growth of distances in hyperbolic space is a direct fit with hierarchies.

## Poincaré embedding<sup>[1]</sup>



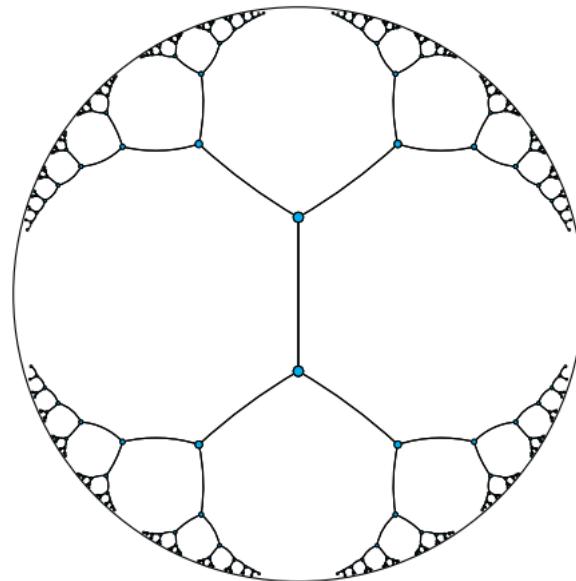
Embedding of a tree

The earliest paper that studied hyperbolic embeddings.

Learning hierarchical representations of data in a low-dimensional space while preserving the inherent geometric structure of the data.

# Embedding

## Poincaré embedding<sup>[1]</sup>



Embedding of a tree

INPUT      Nodes      Positive Samples

$$\mathcal{S} = \{x_i\}_{i=1}^n \quad \mathcal{D} = \{(u, v)\}$$

Negative Samples

$$\mathcal{N}(u) = \{v' \mid (u, v') \notin \mathcal{D}\} \cup \{v\}$$

OPTIMI  
ZATION

$$\Theta' \leftarrow \arg \min_{\Theta} \mathcal{L}(\Theta) \quad \text{s.t. } \forall \theta_i \in \Theta : \|\theta_i\| < 1$$

$$\mathcal{L}(\Theta) = \sum_{(u, v) \in \mathcal{D}} \log \frac{e^{-d(u, v)}}{\sum_{v' \in \mathcal{N}(u)} e^{-d(u, v')}}$$

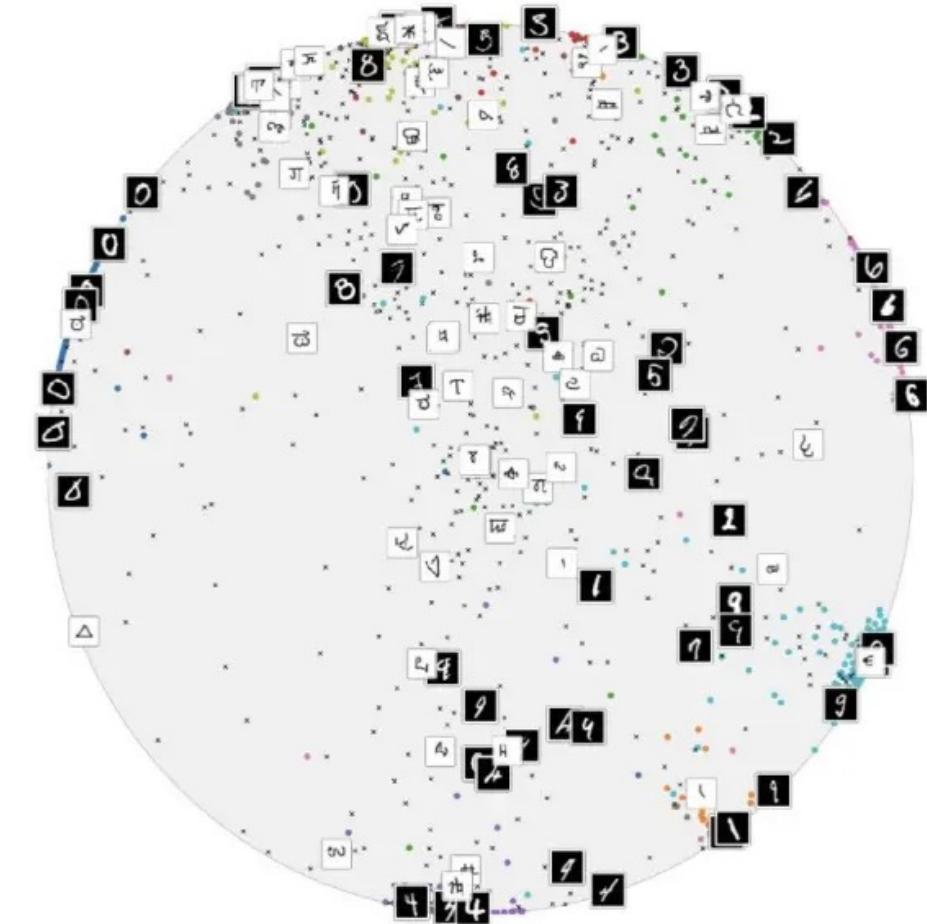
OUTPUT      Embeddings in poincaré model.

$$\Theta = \{\theta_i\}_{i=1}^n$$

Hyperbolic embeddings for computer vision.

## Hyperbolic Image embeddings<sup>[1]</sup>

- Develop a new method for embedding high-dimensional data, such as images, in a **low-dimensional** space that **retains the important relationships** between the data points.
- The motivation for using hyperbolic space is that it better captures the hierarchical structure present in many real-world data sets, such as image categories.



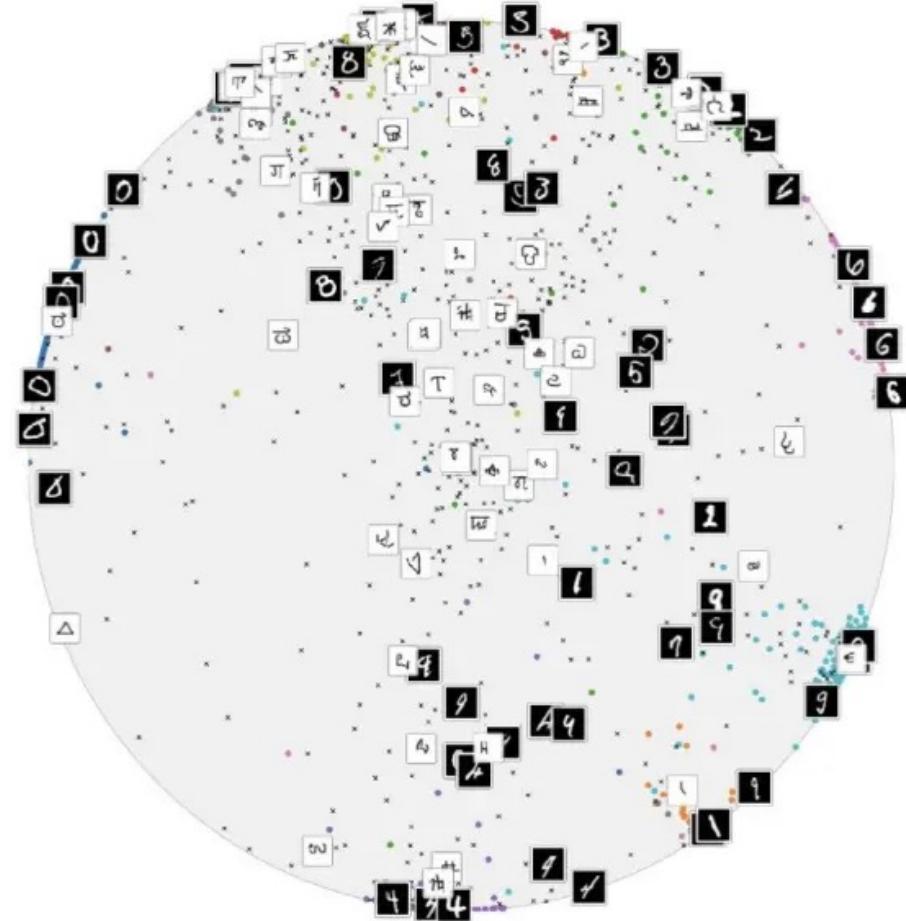
Hyperbolic embeddings for computer vision.

## Hyperbolic Image embeddings<sup>[1]</sup>

- Input the images to CNNs.
- Map the features to hyperbolic space.
- Generate prototypes in hyperbolic spaces.

$$\text{HypAve}(\mathbf{x}_1, \dots, \mathbf{x}_N) = \sum_{i=1}^N \gamma_i \mathbf{x}_i / \sum_{i=1}^N \gamma_i,$$

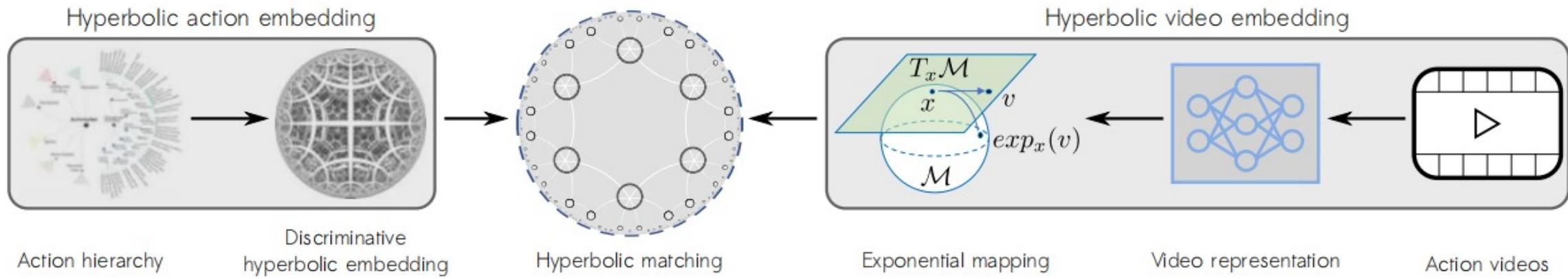
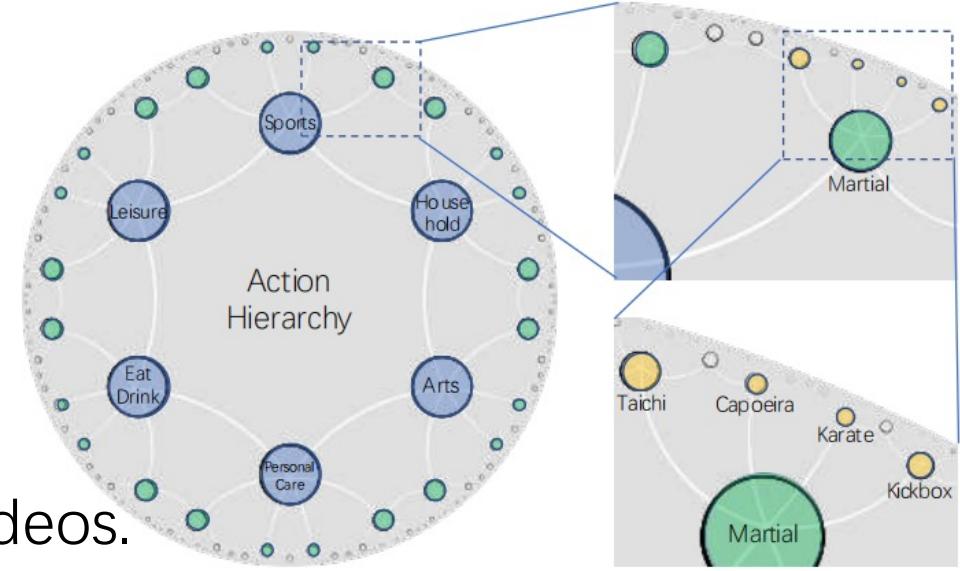
$$\gamma_i = \frac{1}{\sqrt{1 - c \|\mathbf{x}_i\|^2}}$$



# Embedding

## Searching for Actions on the Hyperbole<sup>[1]</sup>

1. Extract poincaré embedding of actions.
2. Map the representation of videos in Euclidean space to hyperbolic space.
3. Minimize the distance between actions and videos.



## Wrapped normal distribution<sup>[1]</sup>

A normal distribution in hyperbolic space

$$\mathcal{N}_{\mathbb{B}_c^d}^W(z|\boldsymbol{\mu}, \Sigma) = \mathcal{N}(\lambda_{\boldsymbol{\mu}}^c \log_{\boldsymbol{\mu}}(z)|\mathbf{0}, \Sigma) \left( \frac{\sqrt{c} d_p^c(\boldsymbol{\mu}, z)}{\sinh(\sqrt{c} d_p^c(\boldsymbol{\mu}, z))} \right)^{d-1}$$

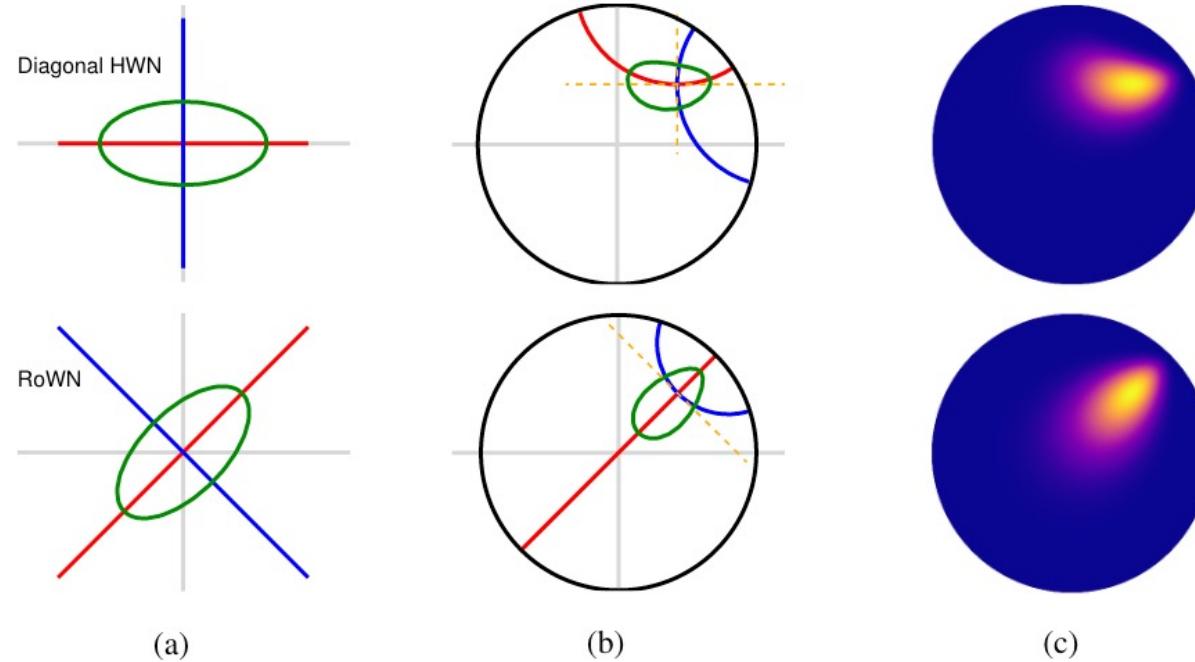
## Riemannian (max-entropy) normal distribution<sup>[2]</sup>

$$\mathcal{N}_{\mathbb{B}_c^d}^R(z|\boldsymbol{\mu}, \sigma^2) = \frac{d\nu^R(z|\boldsymbol{\mu}, \sigma^2)}{d\mathcal{M}(z)} = \frac{1}{Z^R} \exp\left(-\frac{d_p^c(\boldsymbol{\mu}, z)^2}{2\sigma^2}\right)$$

[1] A Wrapped Normal Distribution on Hyperbolic Space for Gradient-Based Learning. Yoshihiro Nagano et al. ICML 2019.

[2] Continuous Hierarchical Representations with Poincaré Variational Auto-Encoders. Mathieu et al. NIPS 2019.

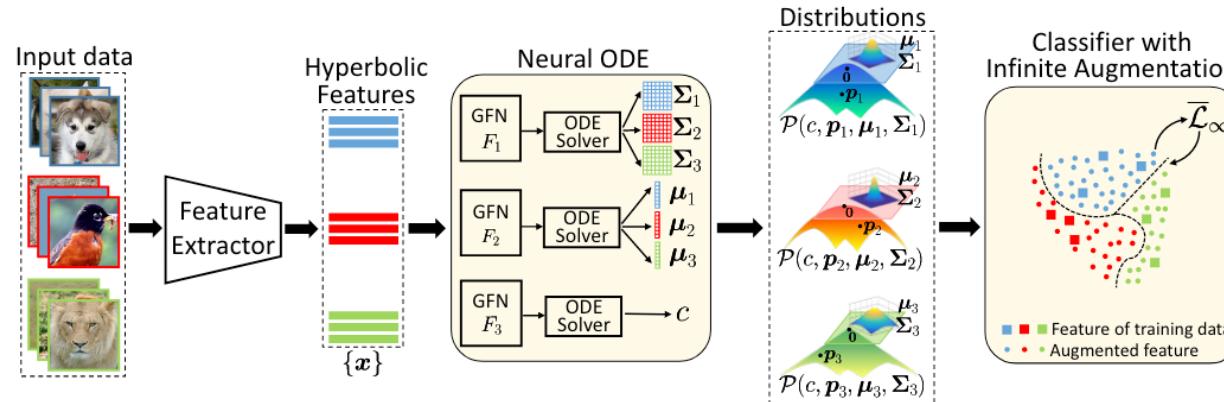
## RoHWN (Rotated Hyperbolic Wrapped Normal Distribution)



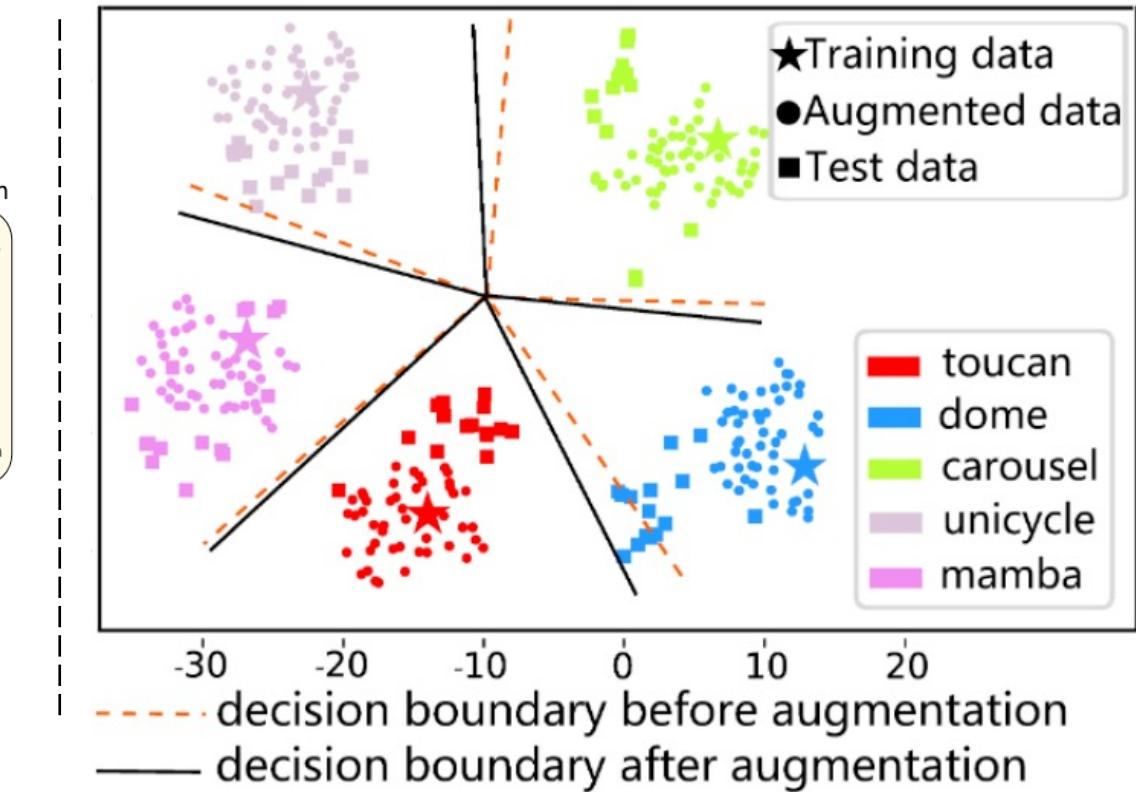
Analyze HWN's geometry and limits in representation learning.

[1] A Rotated Hyperbolic Wrapped Normal Distribution for Hierarchical Representation Learning. Seunghyuk Cho et al. CSED POSTECH

## Hyperbolic Feature Augmentation via Distribution Estimation<sup>[1]</sup>



Sample an infinite number of data augmentation points from HWN.

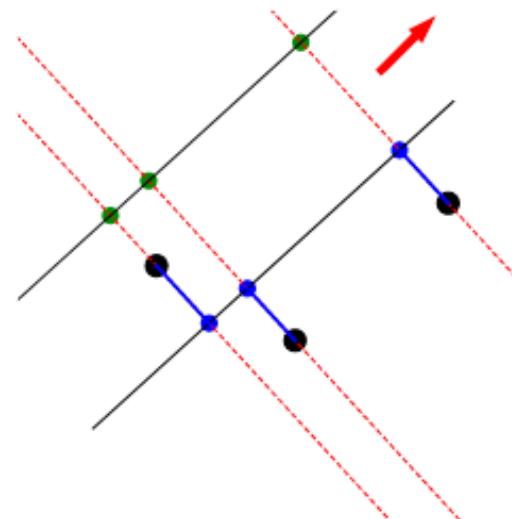


[1] Hyperbolic Feature Augmentation via Distribution Estimation and Infinite Sampling on Manifolds. Gao Zhi et al. NIPS 2022

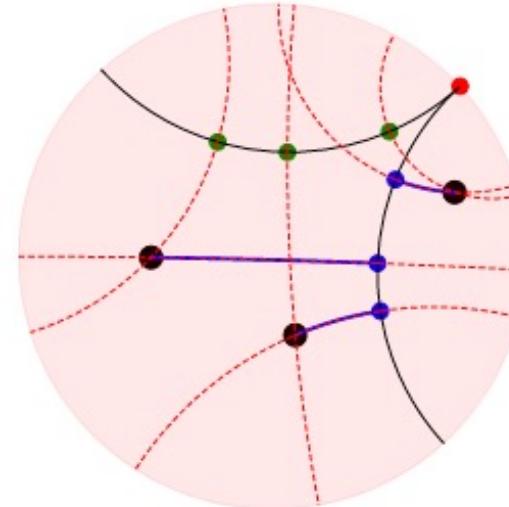
# Dimension Reduction

## HoroPCA

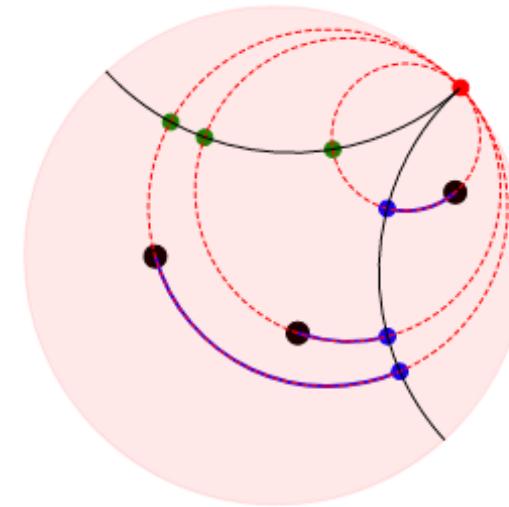
Look for directions that explain the data in hyperbolic space.



(a) Euclidean projections.



(b) Hyperbolic geodesic projections.

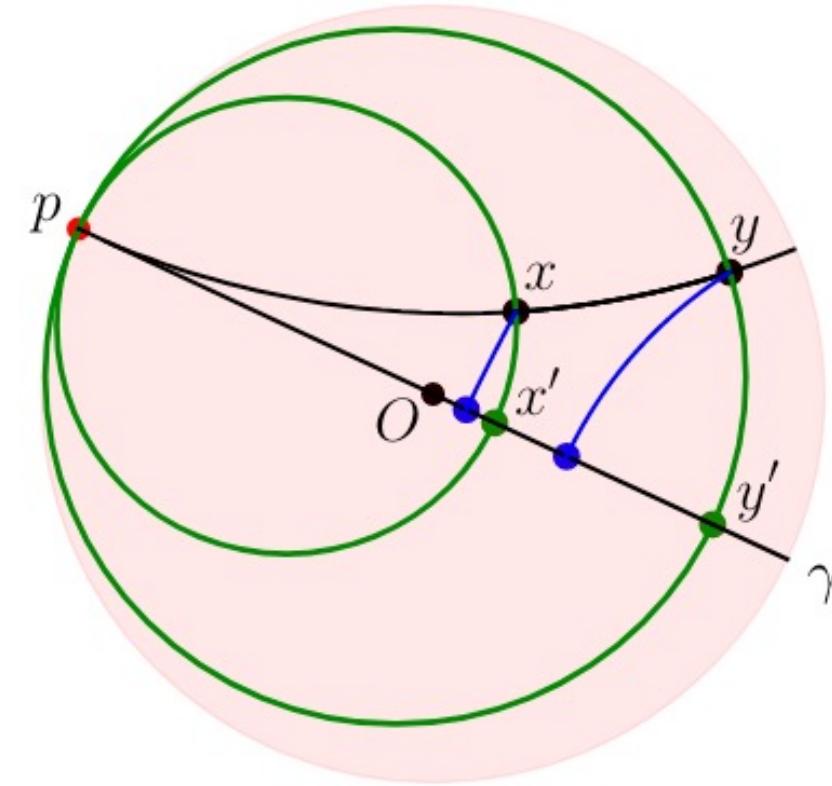
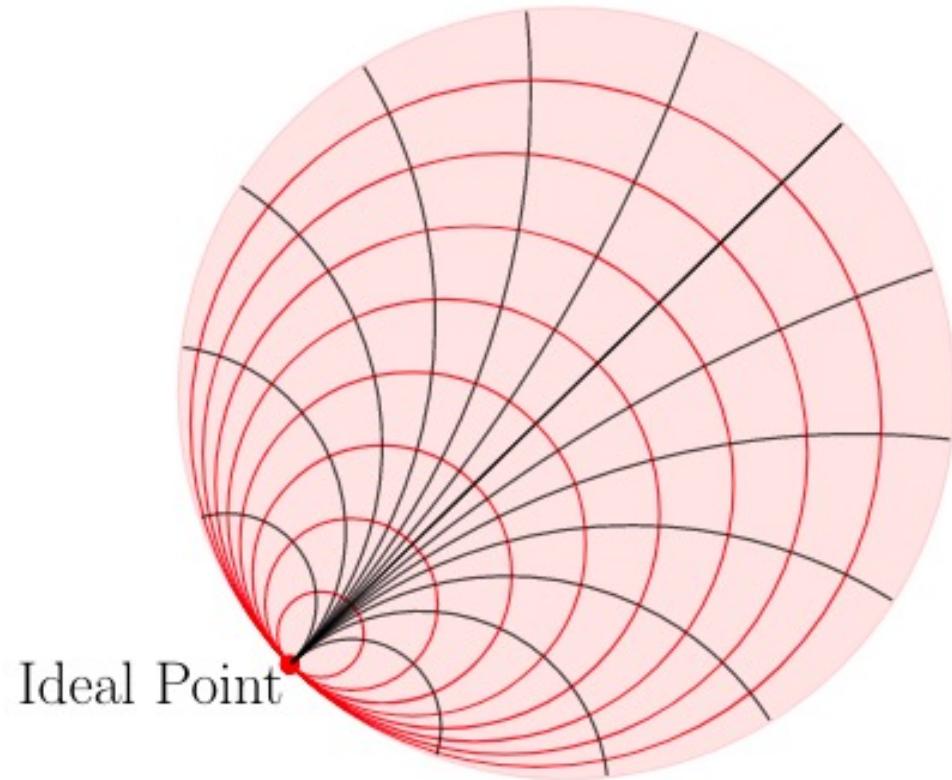


(c) Hyperbolic horospherical projections.

[1]

[1] HoroPCA: Hyperbolic Dimensionality Reduction via Horospherical Projections. Ines Chami et al. ICML 2021.

# Dimension Reduction

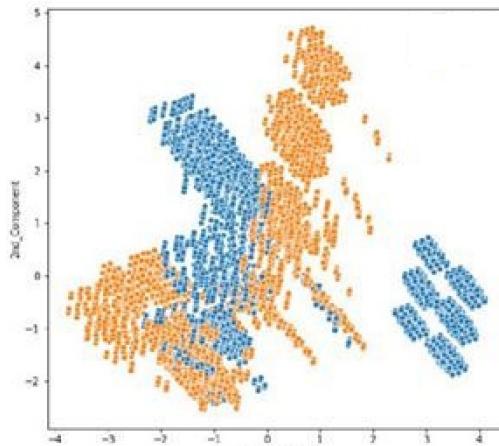


PGA vs HoroPCA [1]

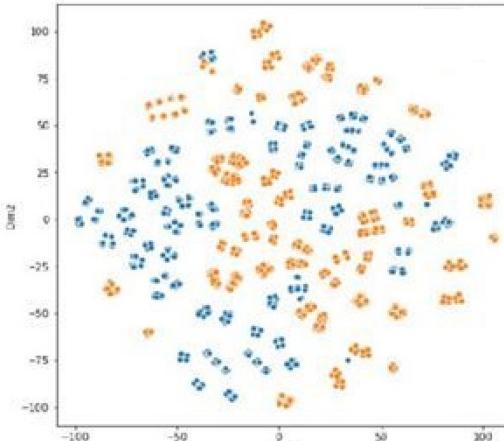
[1] HoroPCA: Hyperbolic Dimensionality Reduction via Horospherical Projections. Ines Chami et al. ICML 2021.

# Dimension Reduction

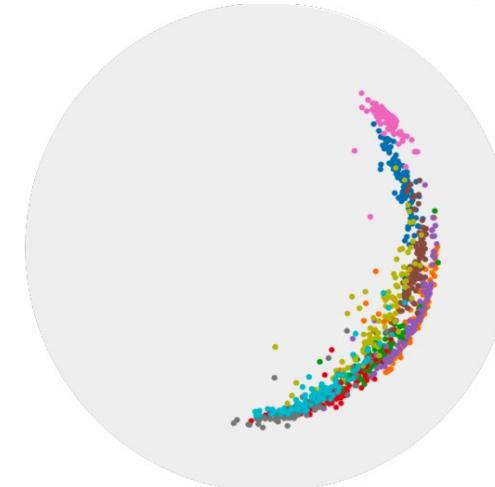
## CO-SNE [1]



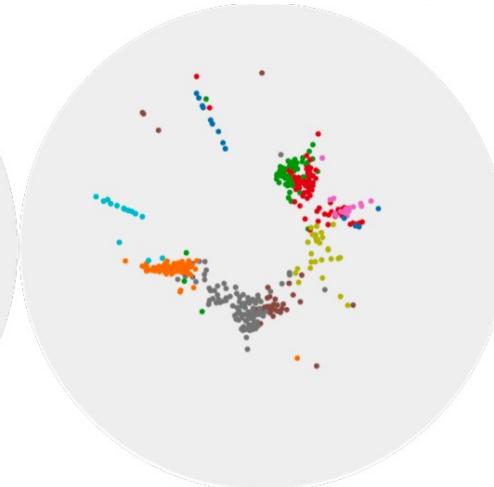
PCA



t-SNE



b) HoroPCA



c) CO-SNE

[1]

HoroPCA cannot preserve the **local** similarity, T-SNE cannot preserve the **global** hierarchy.

CO-SNE can preserve **both the global hierarchy and local similarity** of high-dimensional hyperbolic embeddings in a low-dimensional hyperbolic space.

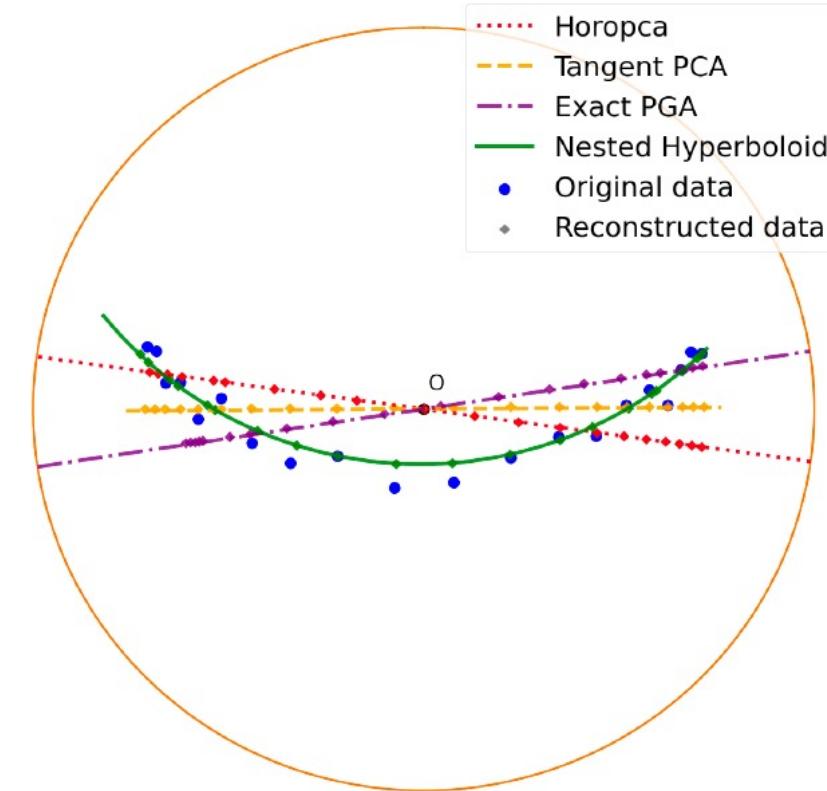
[1] CO-SNE: Dimensionality Reduction and Visualization for Hyperbolic Data. Yunhui Guo et al. CVPR 2022



# Dimension Reduction

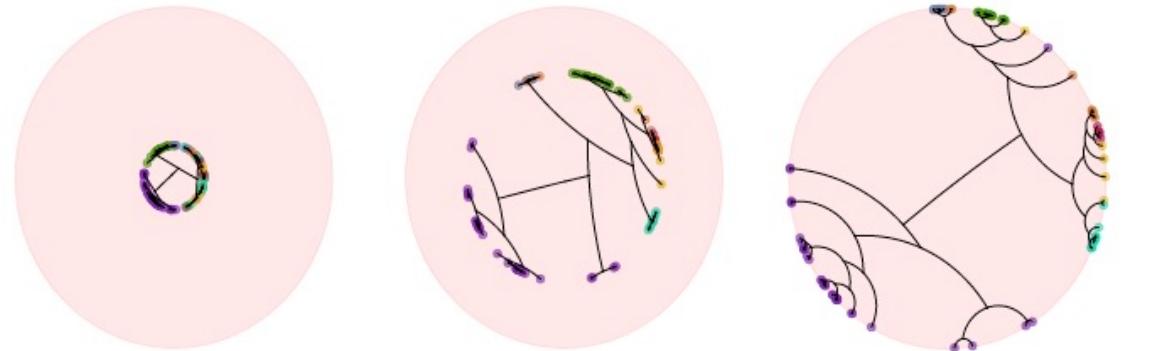
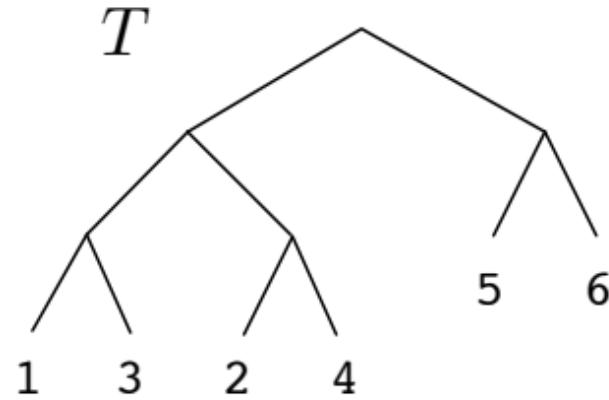
## Nested Hyperboloid [1]

Capture the main trend of the data.  
Proposed curved principal component axes,  
better suited to the data



[1] Nested Hyperbolic Spaces for Dimensionality Reduction and Hyperbolic NN Design. Xiran Fan et al. CVPR 2022.

## HypHC (Hyperbolic Hierarchical Clustering)<sup>[1]</sup>



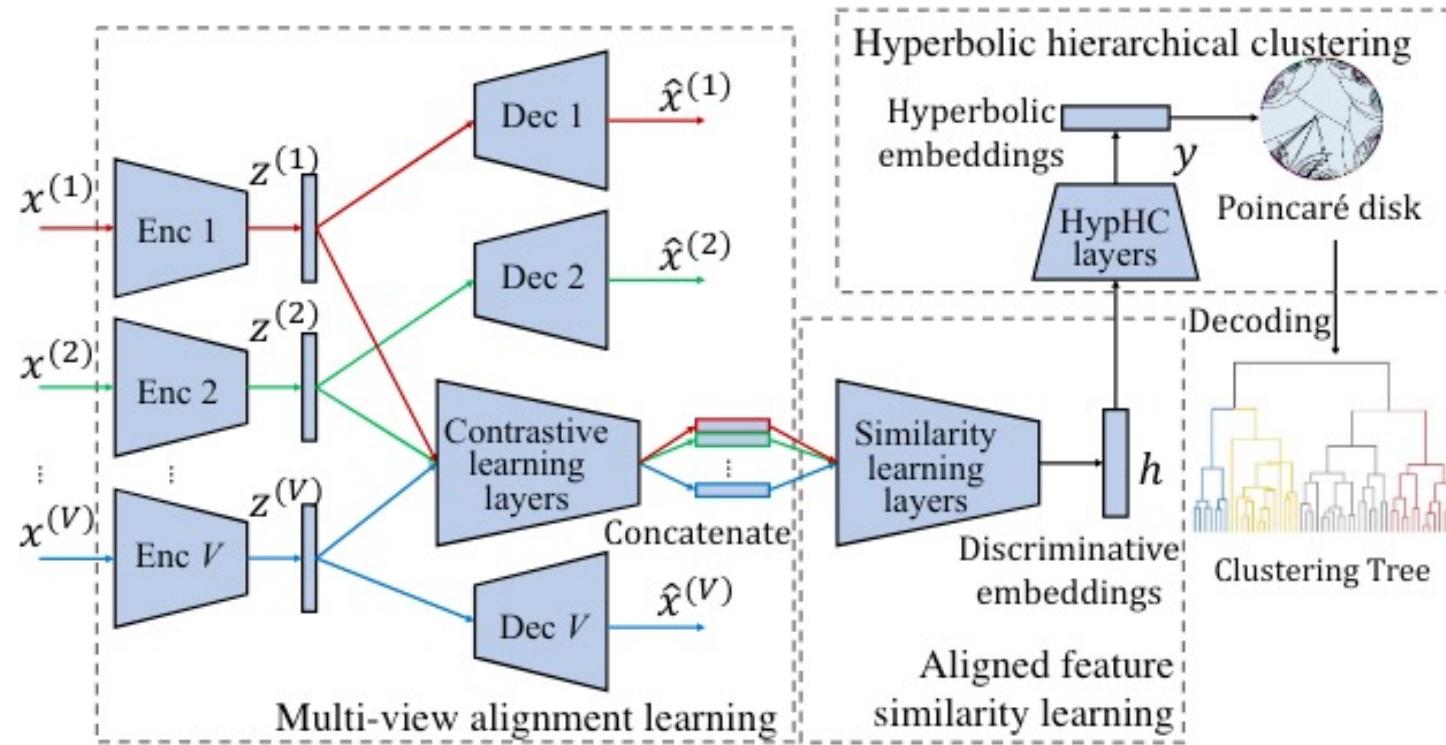
[1]

Show a direct correspondence from discrete trees to **continuous** representations and back, allowing us to search the space of discrete binary trees with continuous optimization.

- [1] From Trees to Continuous Embeddings and Back: Hyperbolic Hierarchical Clustering. Ines Chami et al. NIPS 2020.
- [2] A cost function for similarity-based hierarchical clustering. Sanjoy Dasgupta.

# Clustering

## CMHHC (Contrastive Multi-view Hyperbolic Hierarchical Clustering)<sup>[1]</sup>



Align data of different modalities in feature space to obtain hierarchical clustering in hyperbolic space

[1] Contrastive Multi-view Hyperbolic Hierarchical Clustering. Fangfei Lin et al. IJCAI 2022.

# Metric Learning

## Kernel Methods

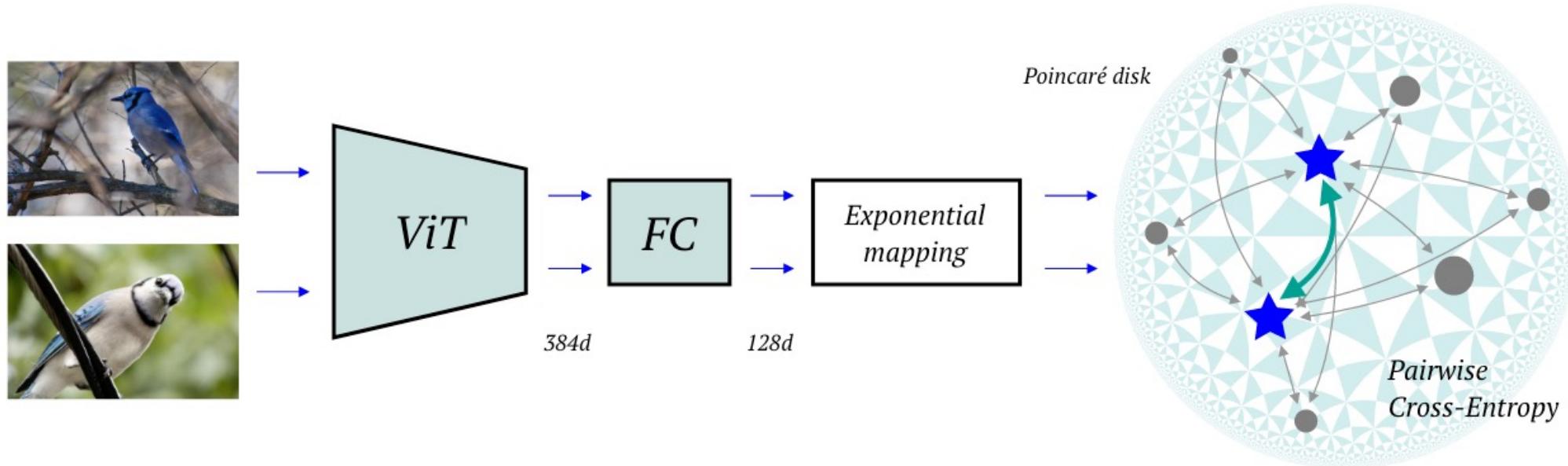
Table 1. Summary of the proposed positive definite kernels in hyperbolic spaces and their properties.

Kernel	Formulation: $k(z_i, z_j)$	Condition	Properties
	$f_{\mathbb{D}}(z) = \tanh^{-1}(\sqrt{c}\ z\ ) \frac{z}{\sqrt{c}\ z\ }, c > 0 \text{ and } z \in \mathbb{D}_c^n$		
Hyperbolic tangent kernel	$k^{\tan}(z_i, z_j) = \langle f_{\mathbb{D}}(z_i), f_{\mathbb{D}}(z_j) \rangle$	-	pd
Hyperbolic RBF kernel	$k^{\text{rbf}}(z_i, z_j) = \exp(-\xi \ f_{\mathbb{D}}(z_i), f_{\mathbb{D}}(z_j)\ ^2)$	$\xi > 0$	pd, universal
Hyperbolic Laplace kernel	$k^{\text{lap}}(z_i, z_j) = \exp(-\xi \ f_{\mathbb{D}}(z_i), f_{\mathbb{D}}(z_j)\ )$	$\xi > 0$	pd, universal
Generalized Hyperbolic Laplace kernel	$k^{\text{glap}}(z_i, z_j) = \exp(-\xi \ f_{\mathbb{D}}(z_i), f_{\mathbb{D}}(z_j)\ ^{2\alpha})$	$\xi > 0, 0 < \alpha < 1$	pd, universal
Hyperbolic binomial kernel	$k^{\text{bin}}(z_i, z_j) = (1 - \langle f_{\mathbb{D}}(z_i), f_{\mathbb{D}}(z_j) \rangle)^{-\alpha}$	$\alpha > 0$	pd, universal

- Benefit from kernel machines and hyperbolic embeddings. [1]
- Simplify various operations involving hyperbolic data because of the rich structure of the Hilbert spaces.

# Metric Learning

Hyperbolic metric works well on vision transformer.



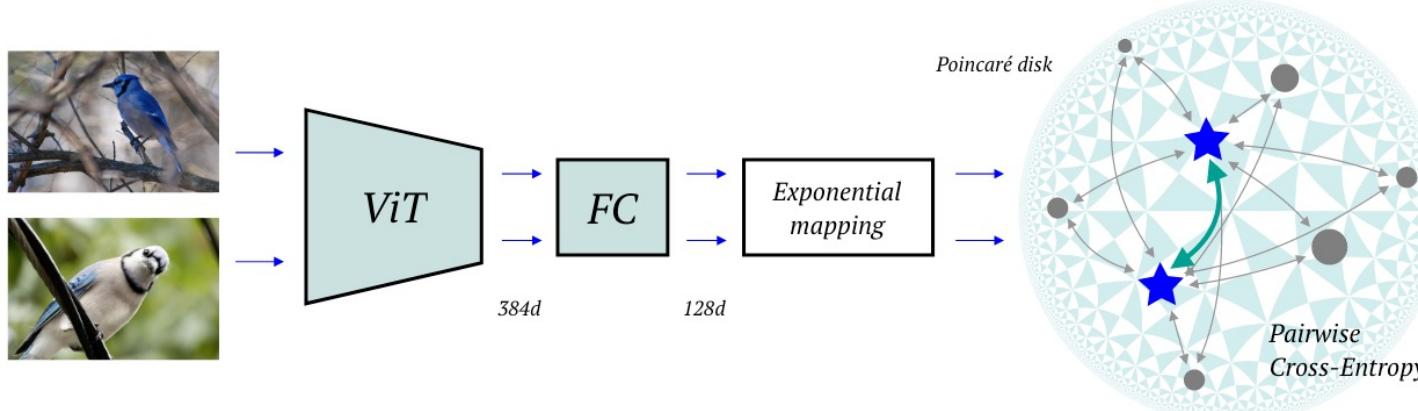
New state-of-the-art on well-known fine-grained datasets.

[1]

[1] Hyperbolic Vision Transformers: Combining Improvements in Metric Learning. Aleksandr Ermolov et al. CVPR2022.

# Metric Learning

Hyperbolic metric works well on vision transformer.



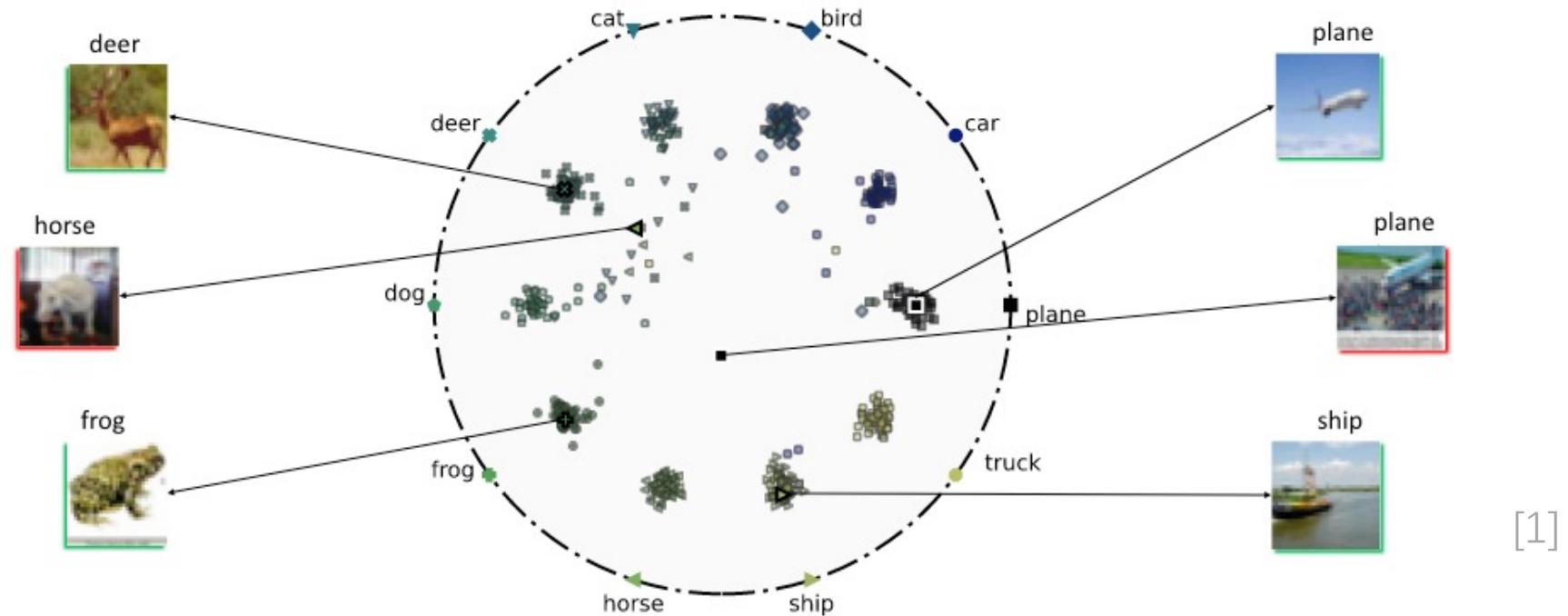
$$D_{cos}(\mathbf{z}_i, \mathbf{z}_j) = \left\| \frac{\mathbf{z}_i}{\|\mathbf{z}_i\|_2} - \frac{\mathbf{z}_j}{\|\mathbf{z}_j\|_2} \right\|_2 = 2 - 2 \frac{\langle \mathbf{z}_i, \mathbf{z}_j \rangle}{\|\mathbf{z}_i\|_2 \cdot \|\mathbf{z}_j\|_2}$$

$$l_{i,j} = -\log \frac{\exp(-D(\mathbf{z}_i, \mathbf{z}_j)/\tau)}{\sum_{k=1, k \neq i}^K \exp(-D(\mathbf{z}_i, \mathbf{z}_k)/\tau)}, \quad [1]$$

[1] Hyperbolic Vision Transformers: Combining Improvements in Metric Learning. Aleksandr Ermolov et al. CVPR2022.

## Hyperbolic Busemann Learning<sup>[1]</sup>

Samples closer to the boundary are more robust.



Embedding hyperbolic prototypes on ideal boundaries of Poincaré without prior label knowledge.

## Hyperbolic Busemann Learning<sup>[1]</sup>

Ideal points are at infinite geodesic distance from all other points in  $\mathbb{B}_d$ .

The Busemann function with respect to  $\mathbf{p}$  is defined for  $\mathbf{z} \in \mathbb{B}_d$  as

$$b_{\mathbf{p}}(\mathbf{z}) = \lim_{t \rightarrow \infty} (d_{\mathbb{B}}(\gamma_{\mathbf{p}}(t), \mathbf{z}) - t).$$

$\mathbf{p}$ : Ideal point.

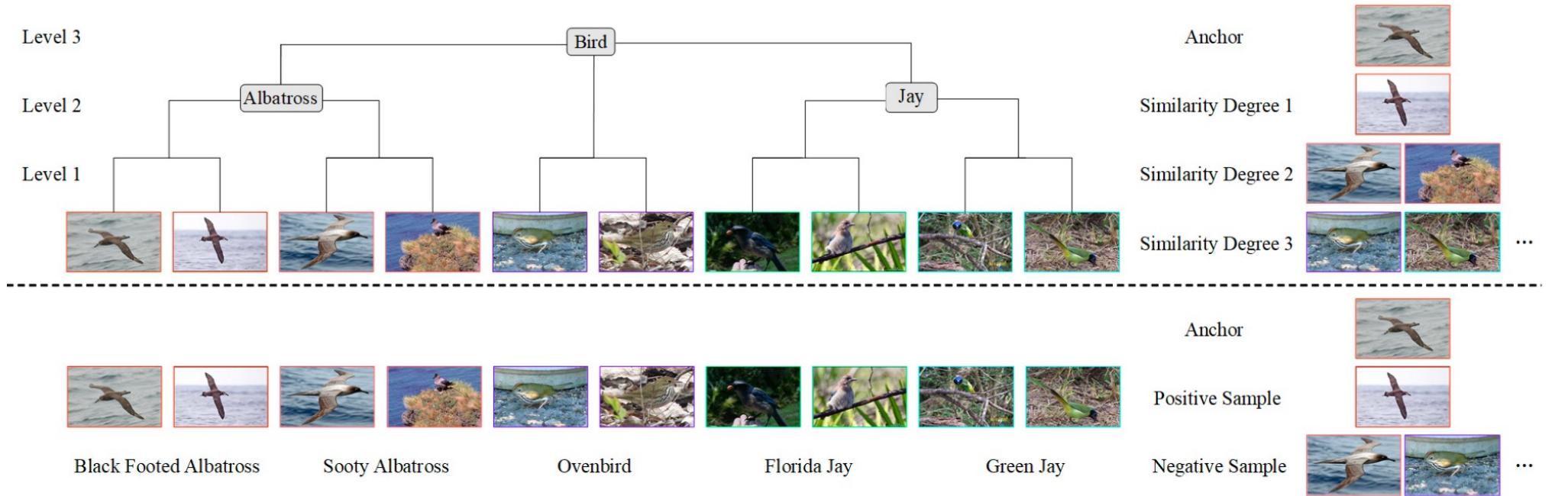
$\gamma_p$ : Geodesic ray, parametrized by arc length, tending to  $p$ .

Busemann function in Poincaré Model

$$b_{\mathbf{p}}(\mathbf{z}) = \log \frac{\|\mathbf{p} - \mathbf{z}\|^2}{(1 - \|\mathbf{z}\|^2)}.$$

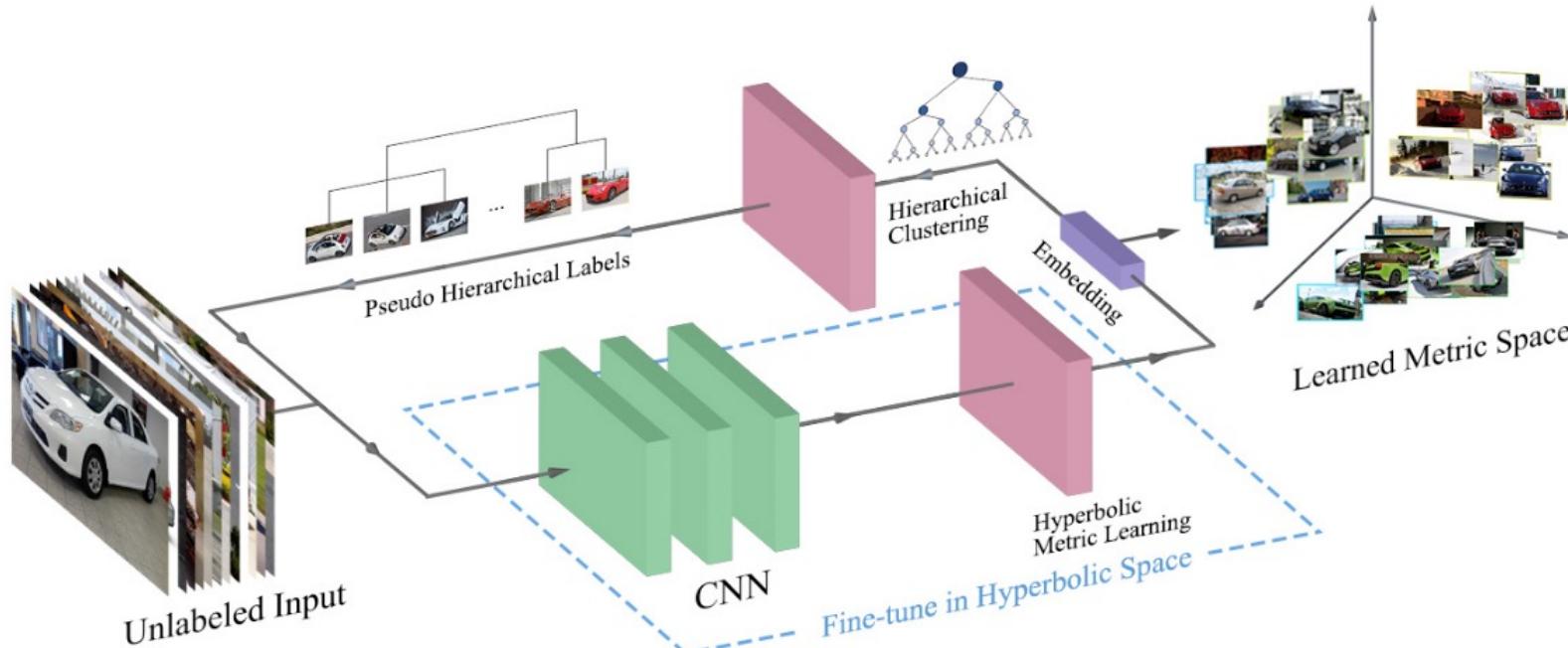
# Metric Learning

## Unsupervised Hyperbolic Metric Learning<sup>[1]</sup>



First hyperbolic unsupervised deep metric learning framework.

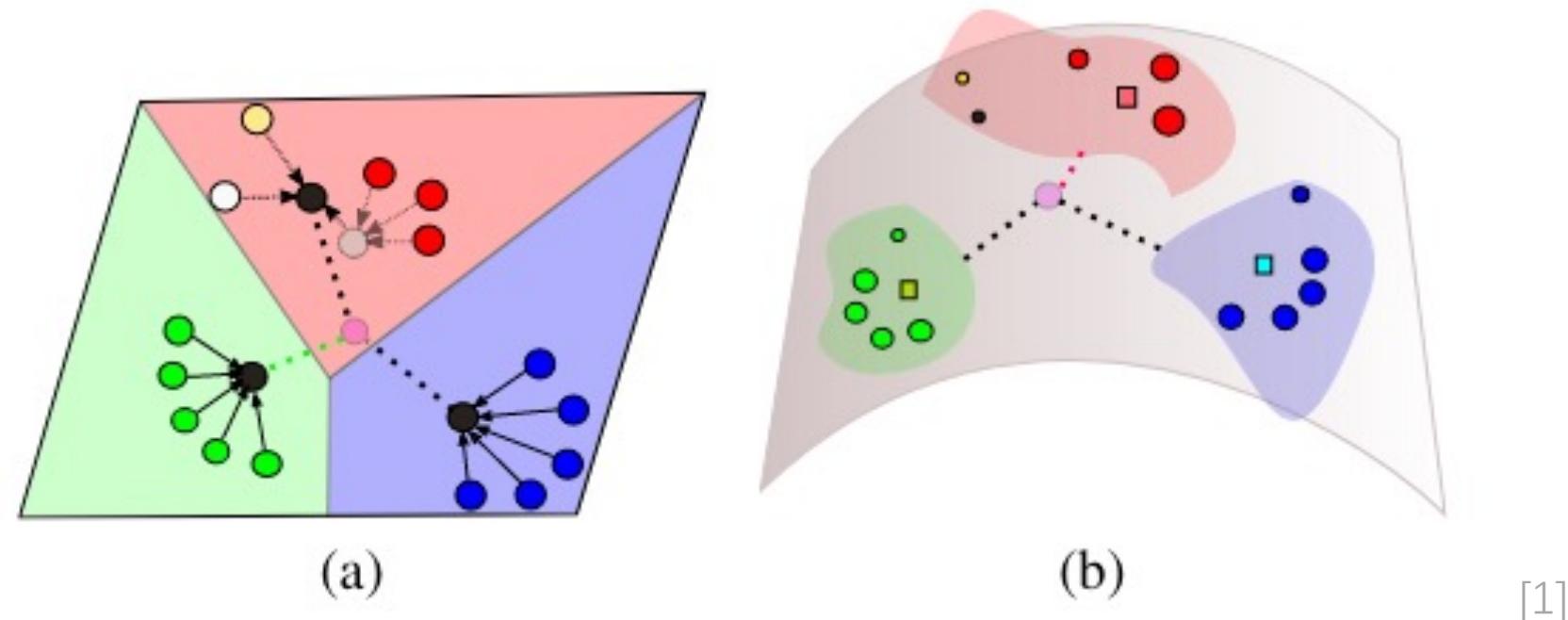
## Unsupervised Hyperbolic Metric Learning<sup>[1]</sup>



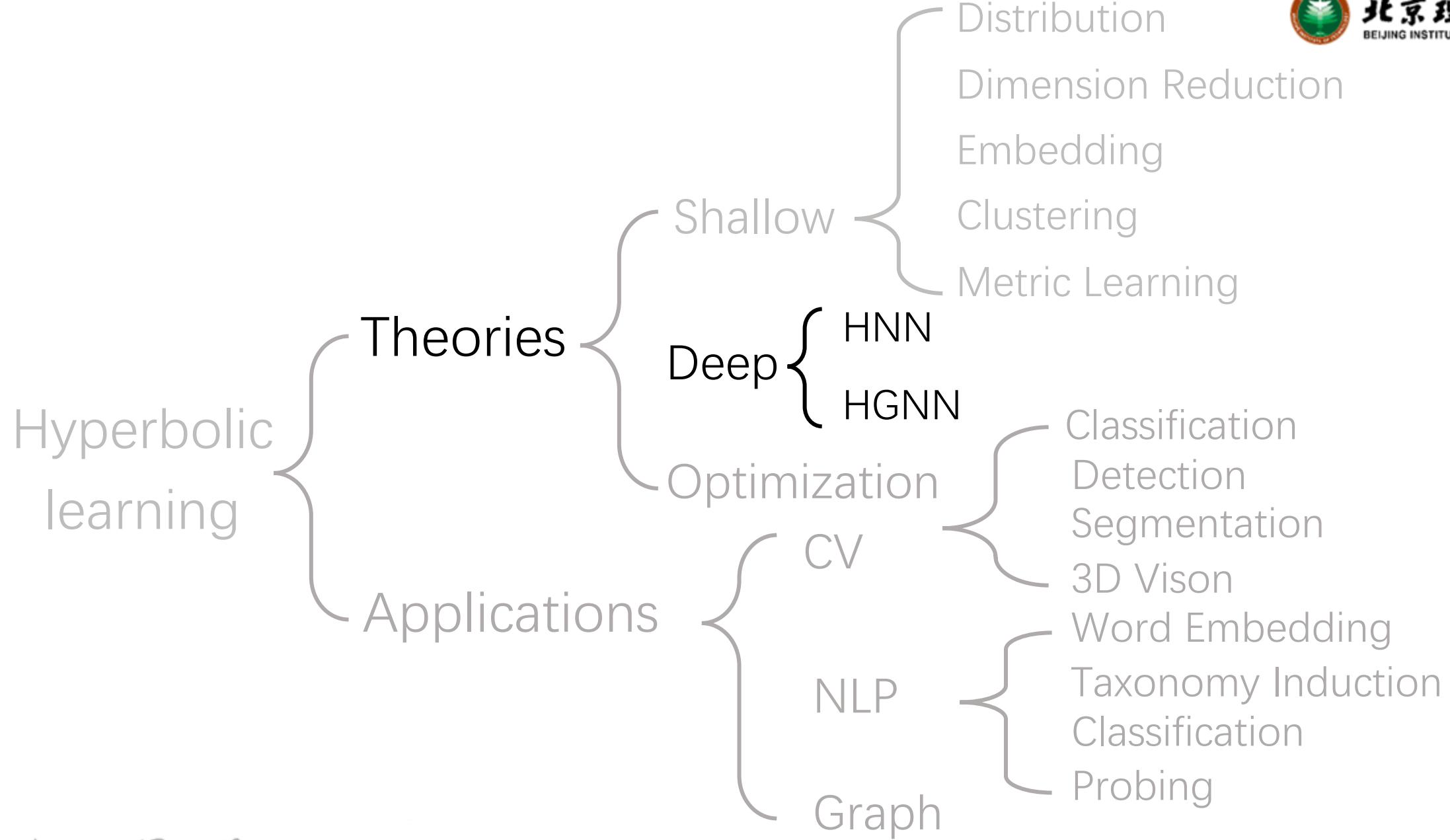
[1]

First hyperbolic unsupervised deep metric learning framework.

## Adaptive Poincaré Point to Set Distance<sup>[1]</sup>



Robust metric against **the noises and outliers** in few-shot learning.



# Theories

Deep Learning

## Hyperbolic Neural Networks<sup>[1]</sup>

Multinomial logistic regression, Feed-forward, Recurrent layers in hyperbolic space.<sup>[1]</sup>

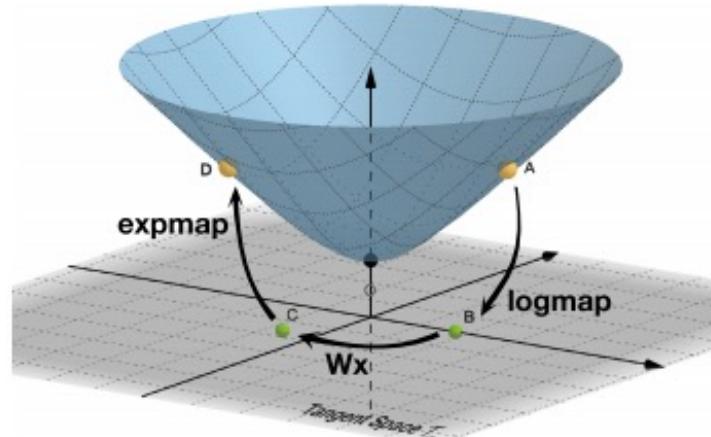
## Hyperbolic Neural Networks++<sup>[2]</sup>

Multinomial logistic regression, Fully-connected layers, Convolutional layers, Attention mechanisms.<sup>[2]</sup>

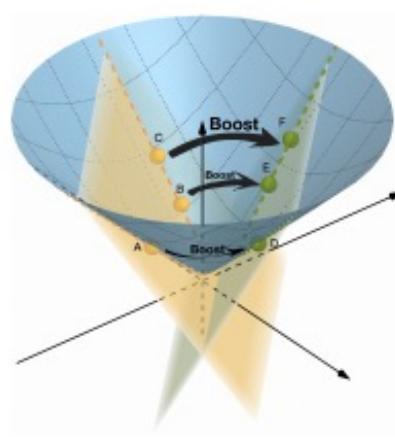
[1] Hyperbolic Neural Networks. Octavian-Eugen Ganea et al. NIPS 2018.

[2] Hyperbolic Neural Networks++. Ryohei Shimizu et al. ICLR 2021.

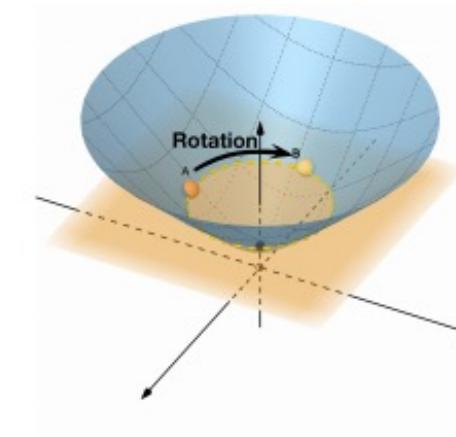
# Fully Hyperbolic Neural Networks<sup>[1]</sup>



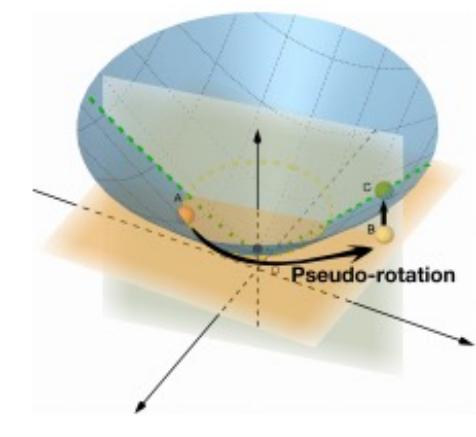
(a) Linear layer formalized in tangent space



(b) Lorentz boost



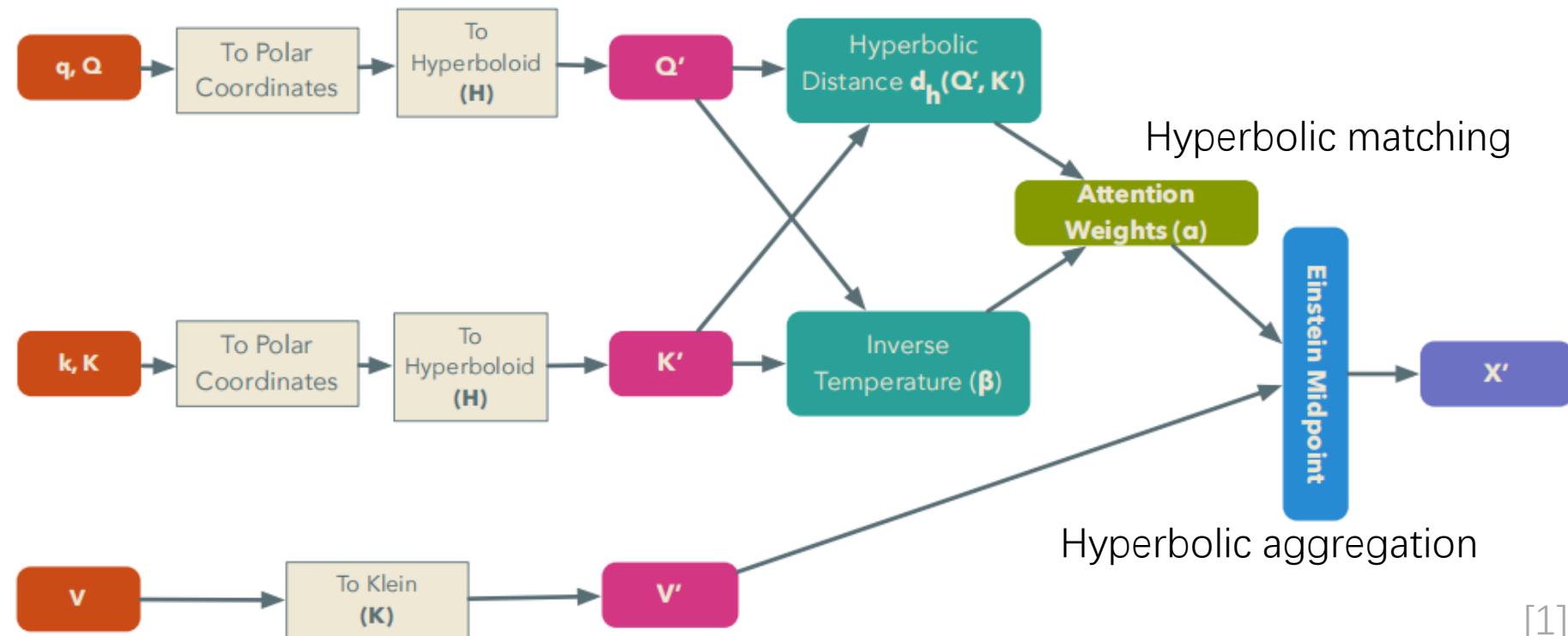
(c) Lorentz rotation



(d) Pseudo-rotation

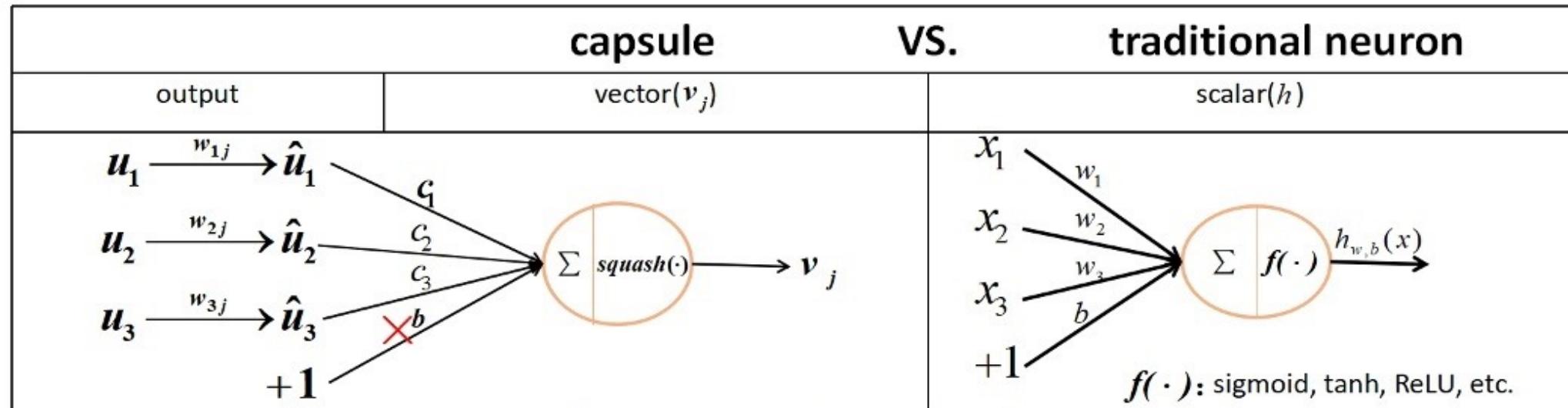
[1]

# Hyperbolic Attention Network<sup>[1]</sup>

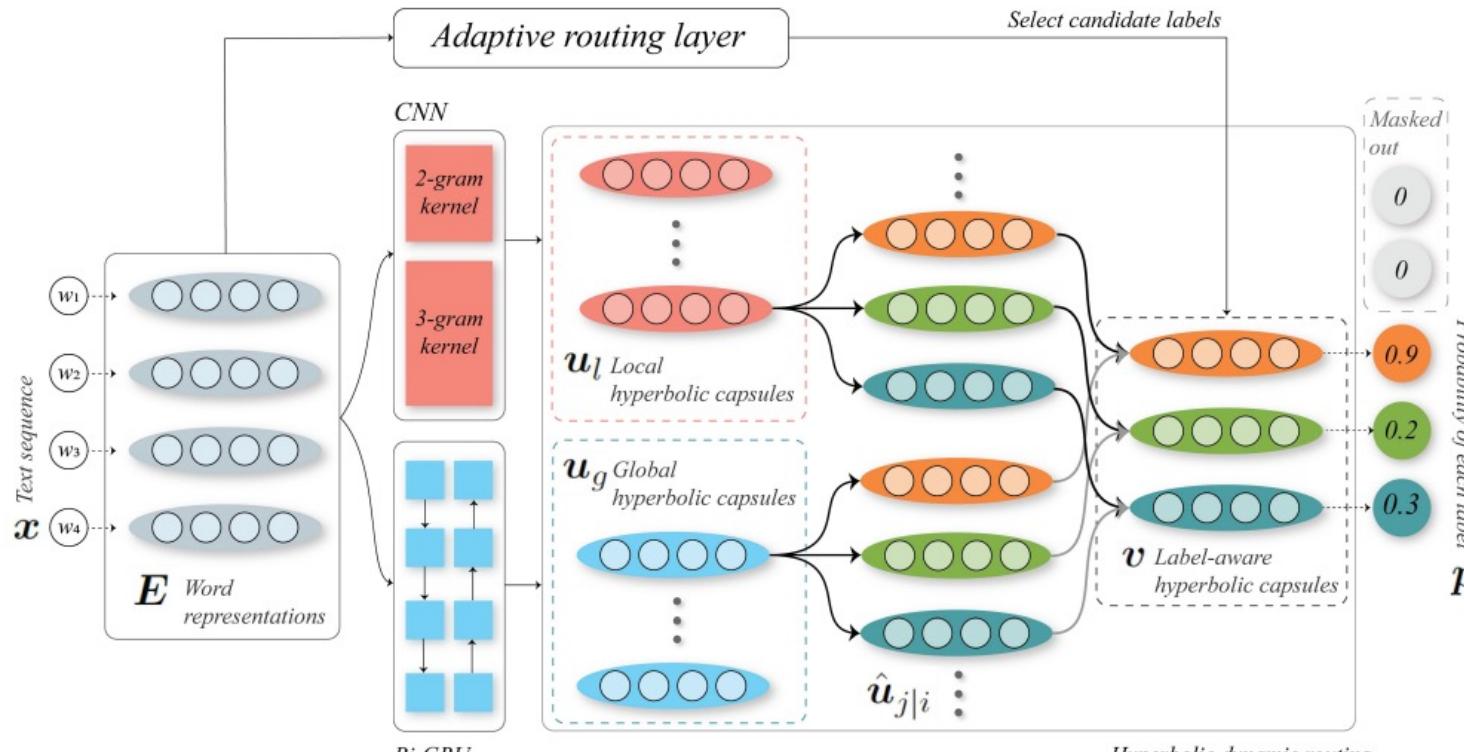


Make the algorithm match more complicated data.

## Hyperbolic Capsule Networks[1]



# Hyperbolic Capsule Networks<sup>[1]</sup>



Capture complicated structures among labels and documents

[1]

HGNN<sub>[1]</sub>

Propagation rule in GNN:

$$\mathbf{h}_u^{k+1} = \sigma \left( \sum_{v \in \mathcal{I}(u)} \tilde{\mathbf{A}}_{uv} \mathbf{W}^k \mathbf{h}_v^k \right)$$

Propagation rule in HGNN:

$$\mathbf{h}_u^{k+1} = \sigma \left( \exp_{\mathbf{x}'} \left( \sum_{v \in \mathcal{I}(u)} \tilde{\mathbf{A}}_{uv} \mathbf{W}^k \log_{\mathbf{x}'}(\mathbf{h}_v^k) \right) \right)$$

HGCN<sub>[1]</sub>

## GCN

$$\mathbf{h}_i^{\ell, E} = W^\ell \mathbf{x}_i^{\ell-1, E} + \mathbf{b}^\ell$$

$$\mathbf{x}_i^{\ell, E} = \sigma(\mathbf{h}_i^{\ell, E} + \sum_{j \in \mathcal{N}(i)} w_{ij} \mathbf{h}_j^{\ell, E})$$

## HGCN

$$\mathbf{h}_i^{\ell, H} = (W^\ell \otimes^{K_{\ell-1}} \mathbf{x}_i^{\ell-1, H}) \oplus^{K_{\ell-1}} \mathbf{b}^\ell$$

$$\mathbf{y}_i^{\ell, H} = \text{AGG}^{K_{\ell-1}}(\mathbf{h}^{\ell, H})_i$$

$$\mathbf{x}_i^{\ell, H} = \sigma^{\otimes^{K_{\ell-1}, K_\ell}}(\mathbf{y}_i^{\ell, H})$$

## Attention

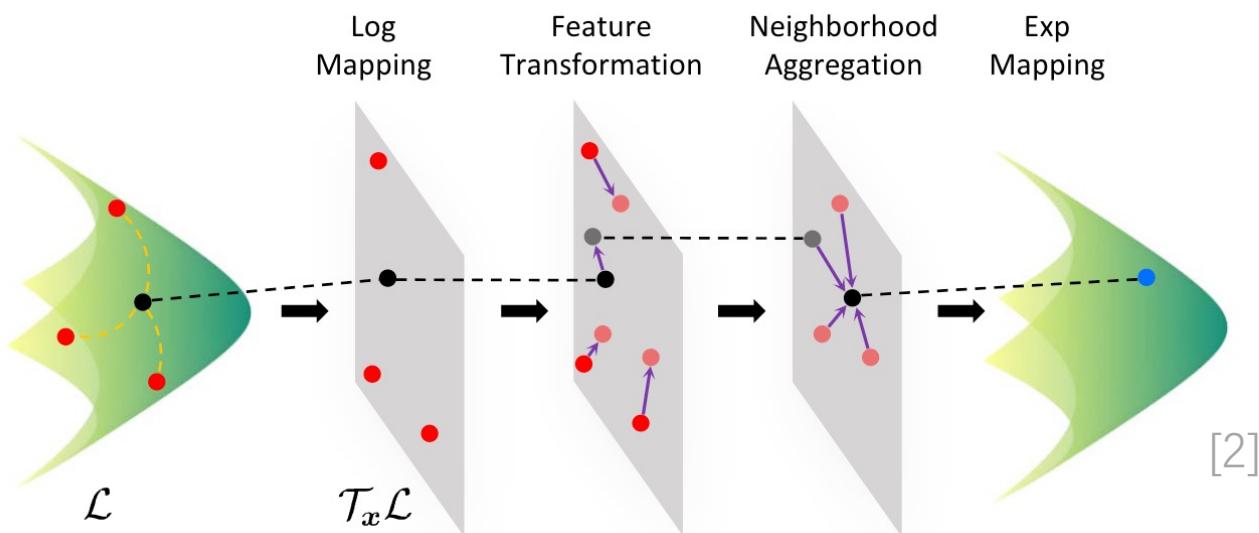
$$w_{ij} = \text{SOFTMAX}_{j \in \mathcal{N}(i)}(\text{MLP}(\log^K_{\mathbf{o}}(\mathbf{x}_i^H) || \log^K_{\mathbf{o}}(\mathbf{x}_j^H)))$$

$$\text{AGG}^K(\mathbf{x}^H)_i = \exp^K_{\mathbf{x}_i^H} \left( \sum_{j \in \mathcal{N}(i)} w_{ij} \log^K_{\mathbf{x}_i^H}(\mathbf{x}_j^H) \right).$$

[1] Hyperbolic Graph Convolutional Neural Networks. Chami et al. NIPS 2019.

## HGCN [1] / HGNN [2]

Map the nodes to tangent space for aggregation.



Transformations – Tangent space  
Neighborhood Aggregation - Tangent space  
Non-linear Activations - Tangent space

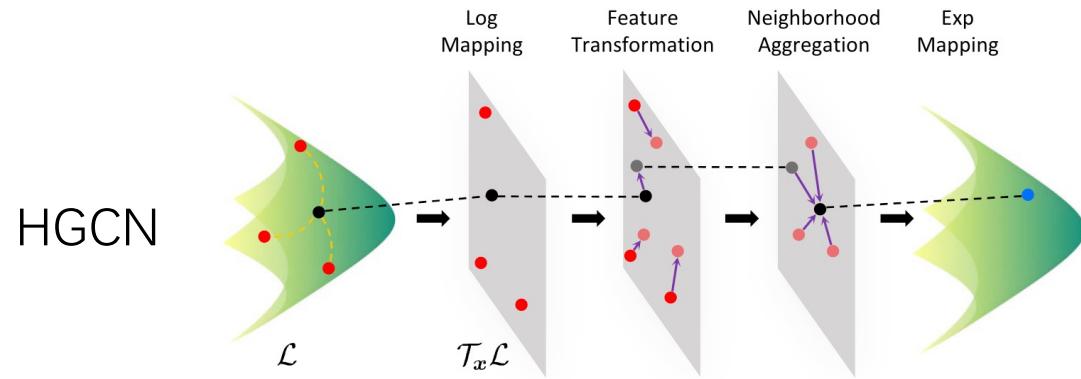
[1] Hyperbolic Graph Convolutional Neural Networks. Chami et al. NIPS 2019.

[2] Hyperbolic Graph Neural Networks. Liu et al. NIPS 2019..

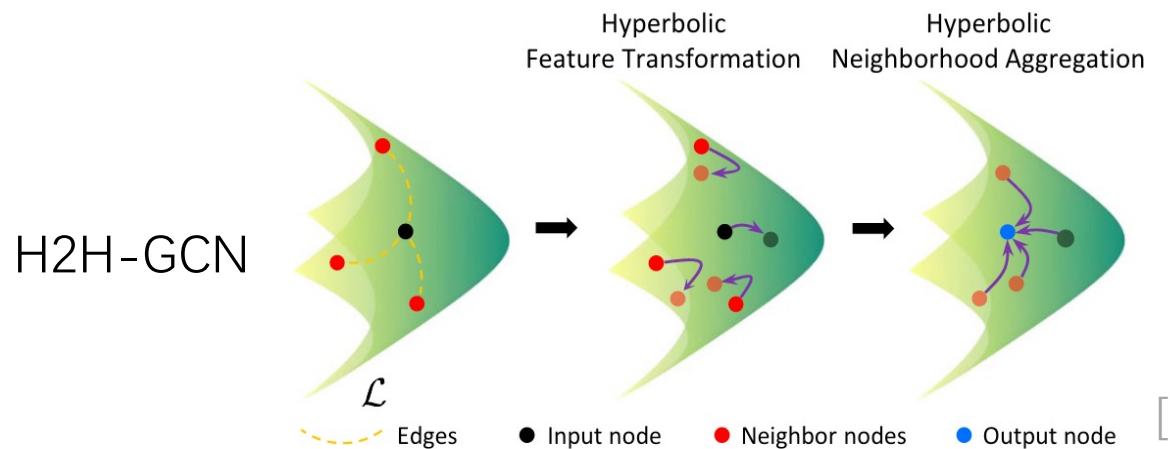
[3] A Hyperbolic-to-Hyperbolic Graph Convolutional Network. Jindou Dai. CVPR 2021.

H2H-GCN<sup>[1]</sup>

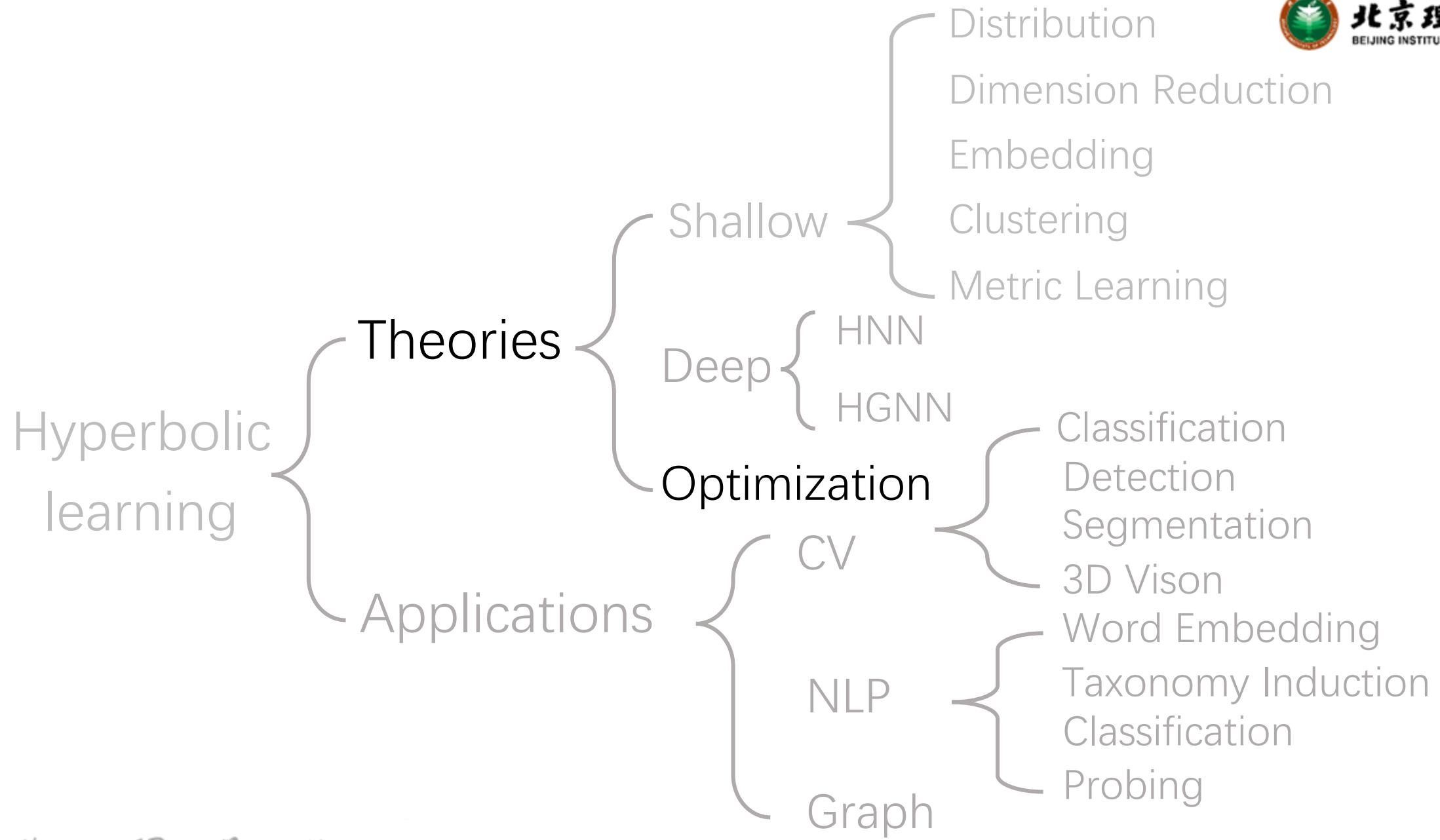
Fully hyperbolic GNN, avoiding approximations.



Transformations – Tangent space  
 Neighborhood Aggregation - Tangent space  
 Non-linear Activations - Tangent space



Transformations - Lorentz model  
 Neighborhood Aggregation - Klein model  
 Non-linear Activations - Poincaré model



# Theories

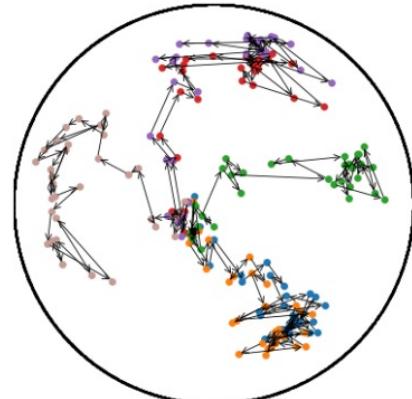
## Optimization

- Non-convex problem in the Euclidean view becomes geodesically convex with the right geometry.
- Hyperbolic form of Nesterov-like accelerated gradient
- No analogue of accelerated gradient descent for geodesically convex functions
- The hyperbolic plane needs to search over a much larger area than any fixed dimension of Euclidean space within a ball of radius  $r$ .

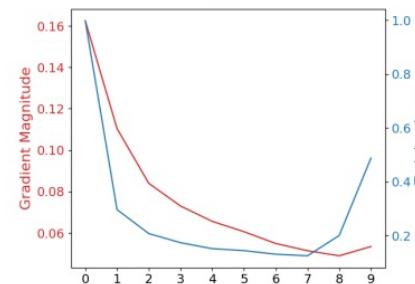
[1] A No-go Theorem for Robust Acceleration in the Hyperbolic Plane, Hamilton et al. NIPS2021

德以明理 学以精工

## Clipped Hyperbolic Classifiers



a)



## Vanishing Gradient Problem

$$p(y = k | \mathbf{x}) \propto \exp(\langle -\mathbf{p}_k \oplus_c \mathbf{x}, \mathbf{a}_k \rangle) \sqrt{\mathbf{g}_{\mathbf{p}_k}^c(\mathbf{a}_k, \mathbf{a}_k)} d_c(\mathbf{x}, \tilde{H}_{\mathbf{a}_k, \mathbf{p}}^c)$$

$\|\mathbf{x}^H\|^2$  approaches 1

$$\frac{\partial \ell}{\partial \mathbf{x}^H} = \frac{(1 - \|\mathbf{x}^H\|^2)^2}{4} \nabla \ell(\mathbf{x}^H)$$

$$\frac{\partial \ell}{\partial \mathbf{w}^E} = \left( \frac{\partial \mathbf{x}^H}{\partial \mathbf{w}^E} \right)^T \frac{\partial \ell}{\partial \mathbf{x}^H}$$

[1]

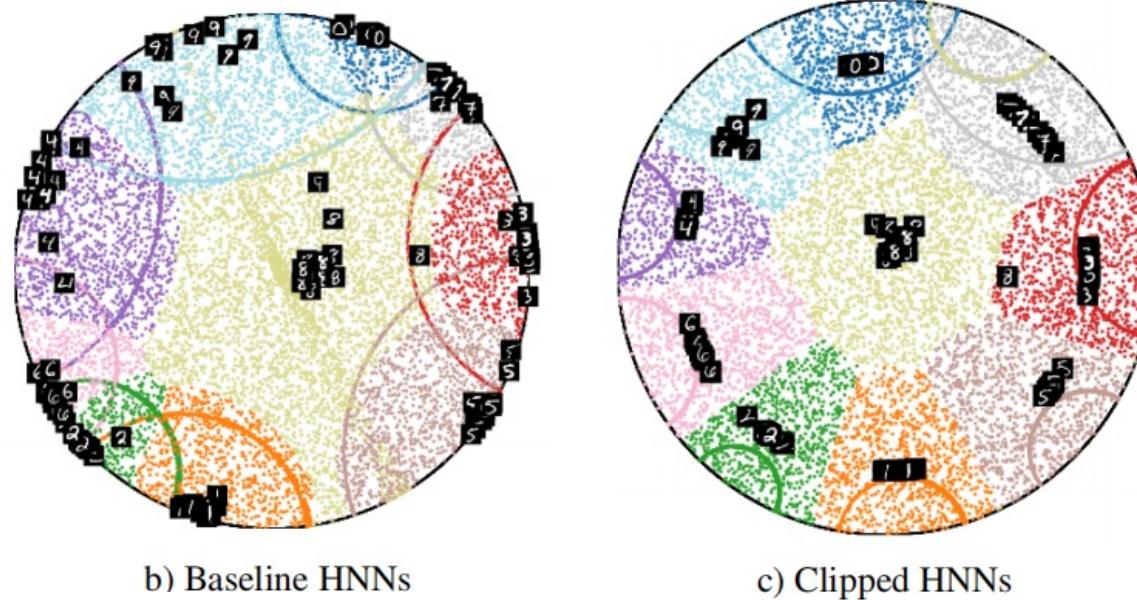
b)

Vanishing gradient problem  
during backpropagation

## Clipped Hyperbolic Classifiers

Euclidean Feature Clipping

$$\text{CLIP}(\mathbf{x}^E; r) = \min\left\{1, \frac{r}{\|\mathbf{x}^E\|}\right\} \cdot \mathbf{x}^E$$



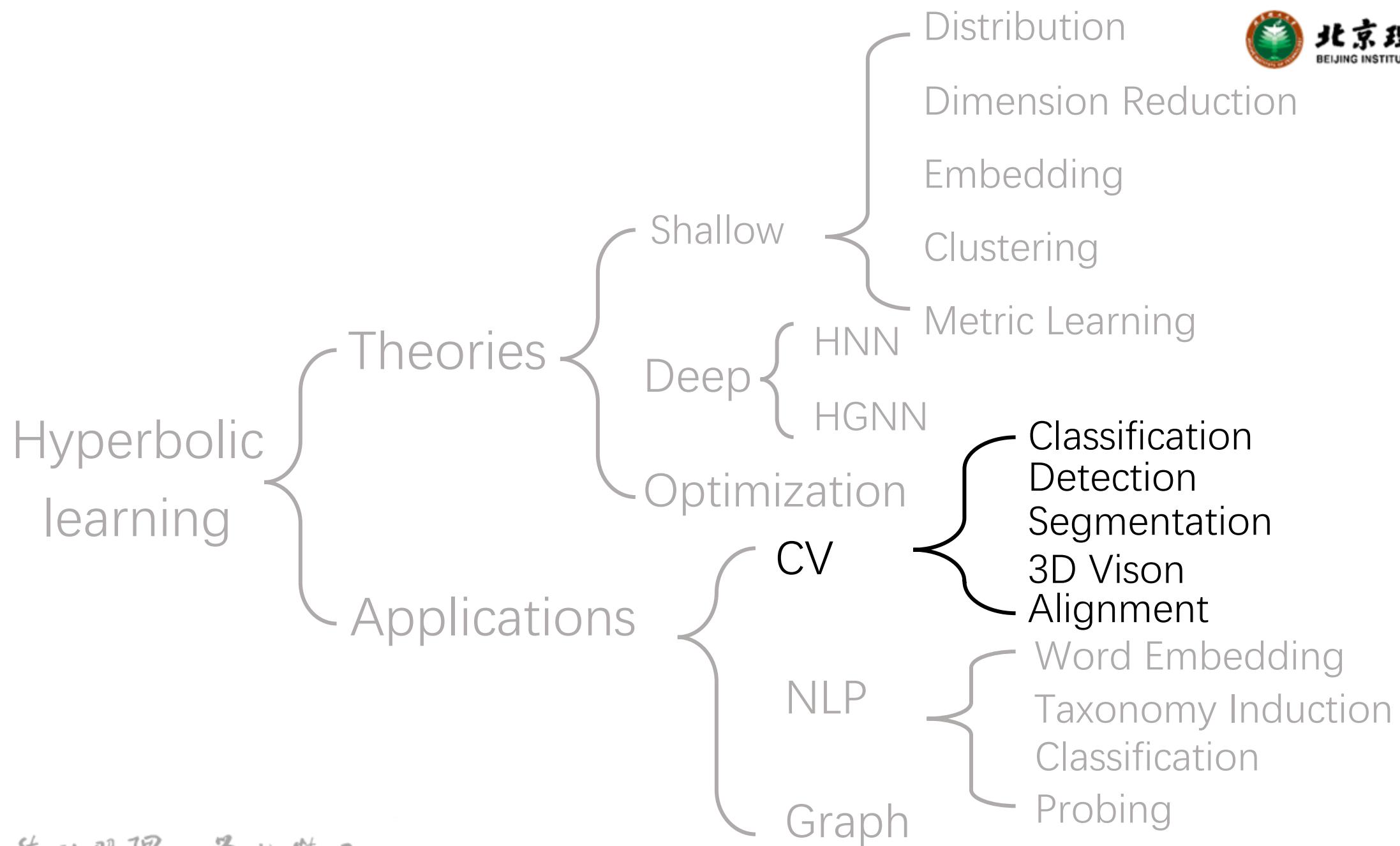
[1]

Alleviating the vanishing gradient problem without compromising accuracy

# Optimization

- Linear graph embedding (LGE): graph embedding using linear space equipped with an inner product function.
- Hyperbolic graph embedding (HGE): graph embedding in hyperbolic space.
- HGE's performance is limited to ideal noiseless. LGE and HGE's performance with noise has not been discussed.
- LGE and HGE cause a polynomial and exponential error with respect to the radius, and imbalanced data distribution can worsen the error.

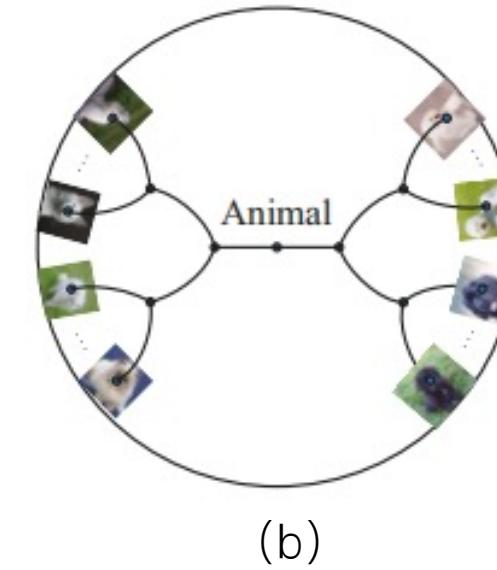
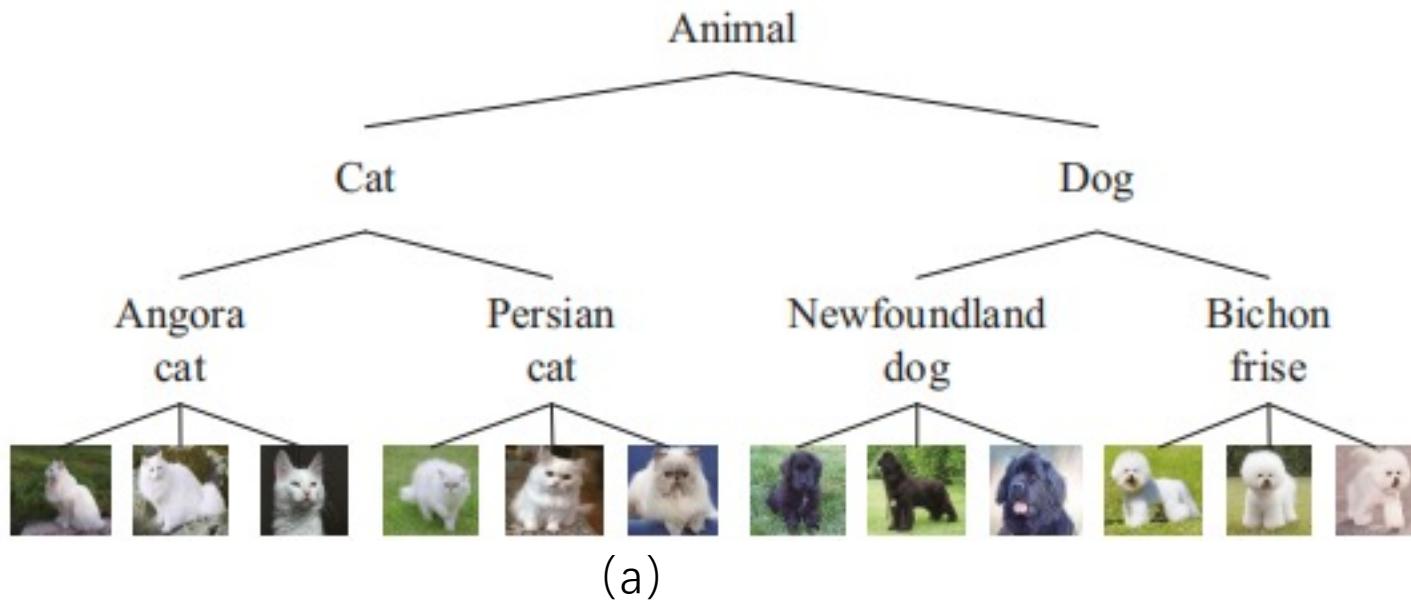
[1] Generalization Error Bounds for Graph Embedding Using Negative Sampling: Linear vs Hyperbolic, Suzuki et al. NIPS2021  
[2] Generalization Error Bound for Hyperbolic Ordinal Embedding, Suzuki et al. ICML2021



# Application

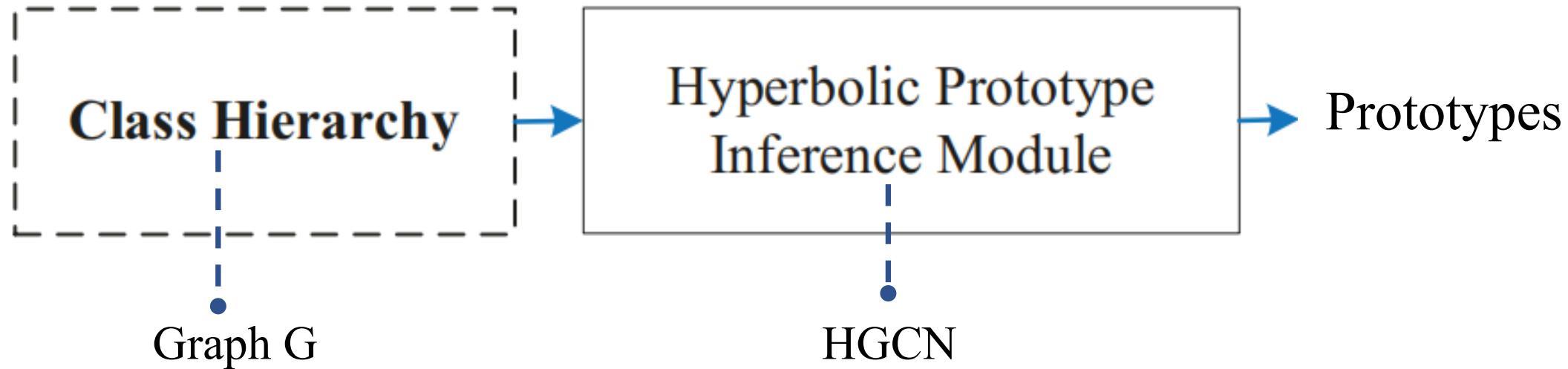
## Computer Vision

# Classification



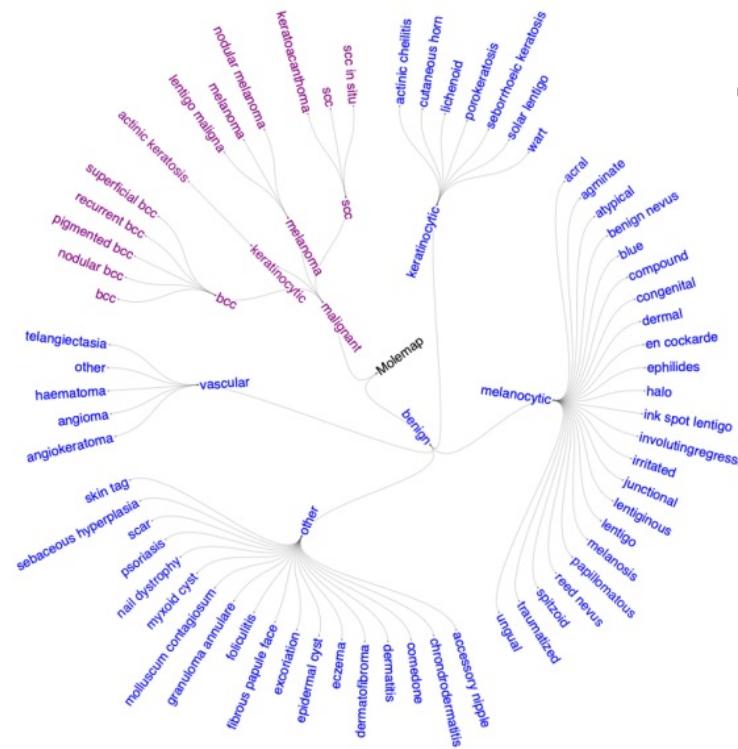
Class hierarchy [1]

# Classification

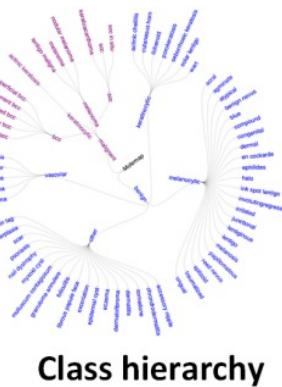


Use class hierarchy in hyperbolic space as prior for Few-shot Learning.

# Classification



Class distance guided hyperbolic prototype learning. [1]

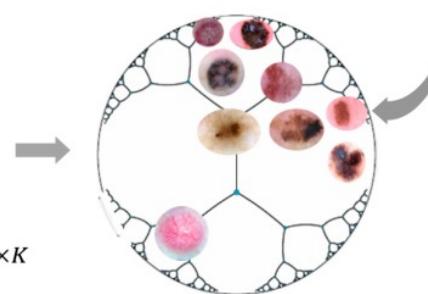


**Relation pairs**

$$\begin{aligned}\mathcal{P} &= \{(u, v) | u = h(v)\} \\ \mathcal{N}(u) &= \{v | (u, v) \notin \mathcal{P}\} \cup \{u\}\end{aligned}$$

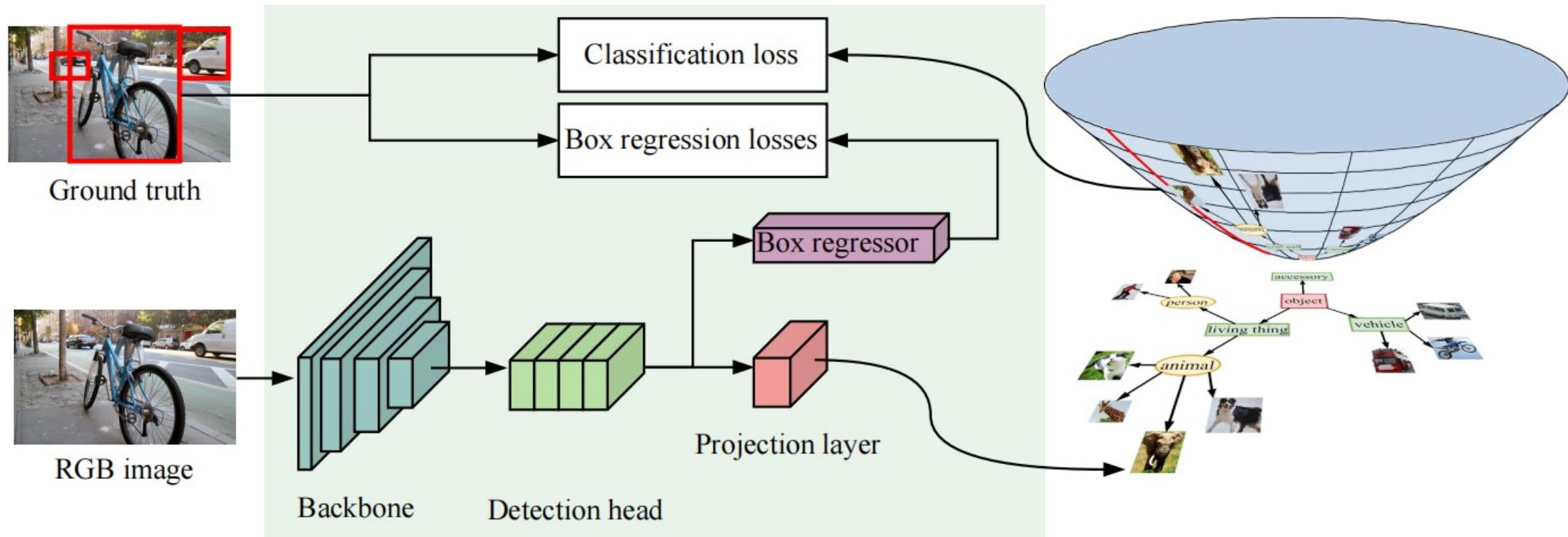
**Class relation encoding**

Class distance matrix:  $\mathbf{D} \in \mathbb{R}_+^{K \times K}$



# Detection

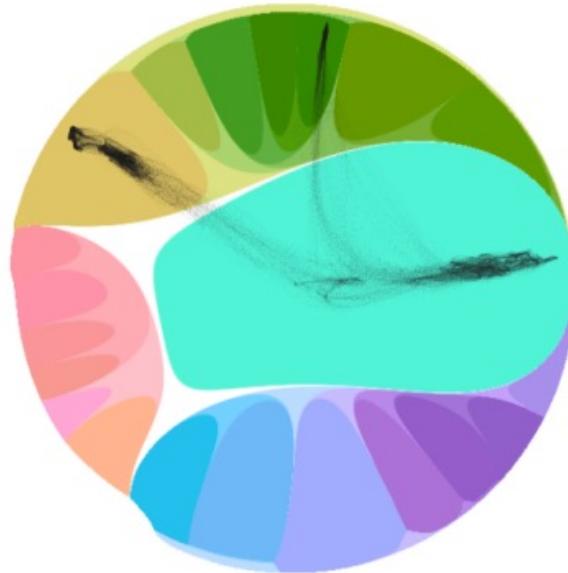
Use hyperbolic embeddings to learn class prototypes for object detector.



# Segmentation

Per-pixel classification in hyperbolic space.

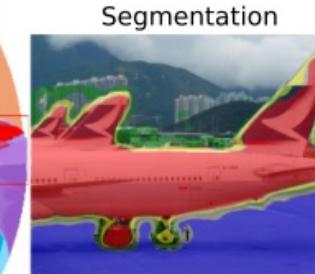
TABLE I [1]  
THEORETICAL  $\delta$ -HYPERBOLICITY AND REAL  $\delta$ -HYPERBOLICITY VALUES  
IN CITYSCAPES AND UAVID DATASETS.  $\mathbb{B}^n$  IS THE POINCARÉ BALL  
MODEL OF HYPERBOLIC SPACE.



(a) Prediction uncertainty for free



(b) Boundary information for free

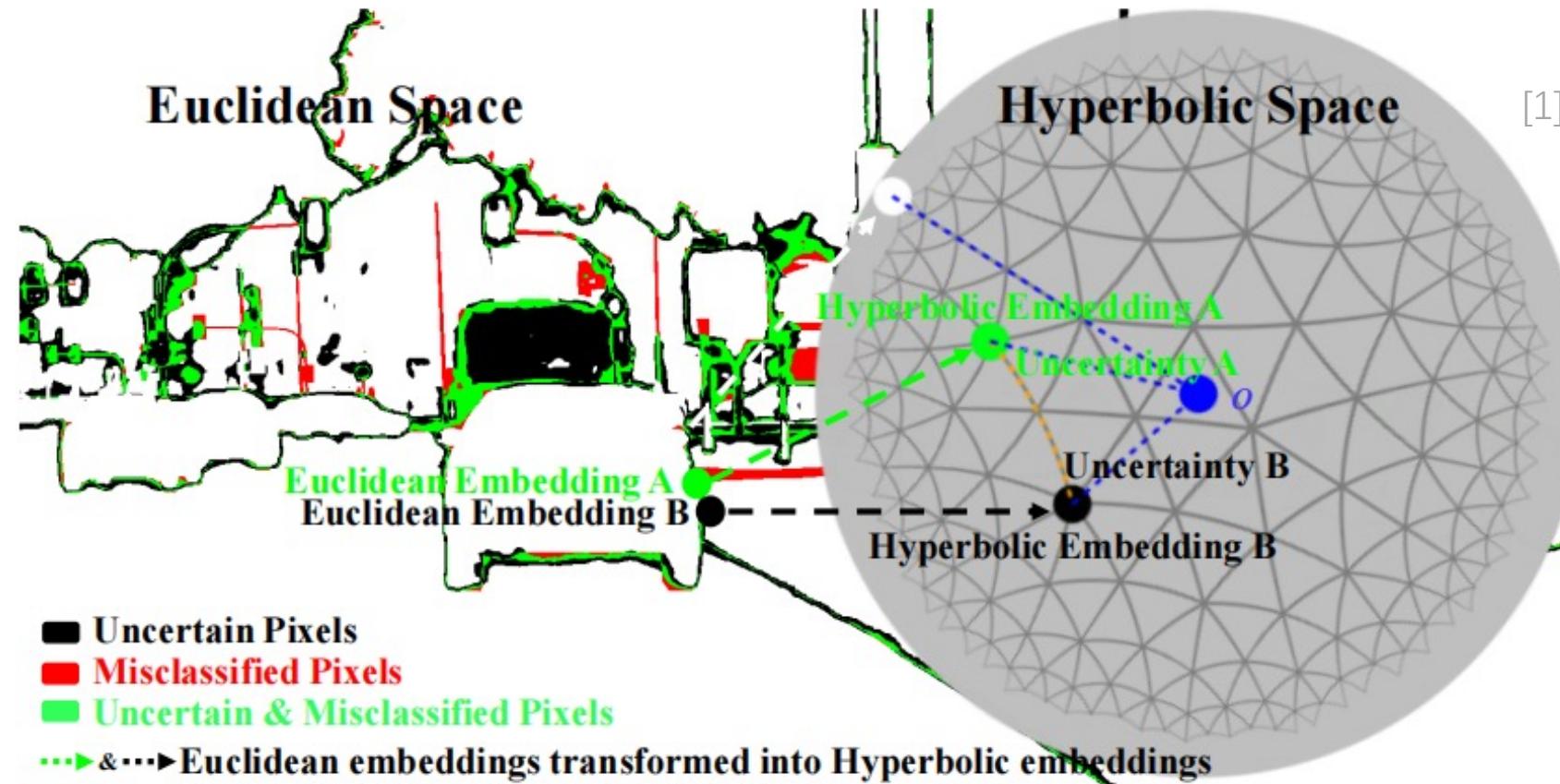


[2]

Spaces	$\mathbb{B}^n$
Theory	0
Cityscapes dataset ( $\delta$ )	0.16
UAVid dataset ( $\delta$ )	0.16

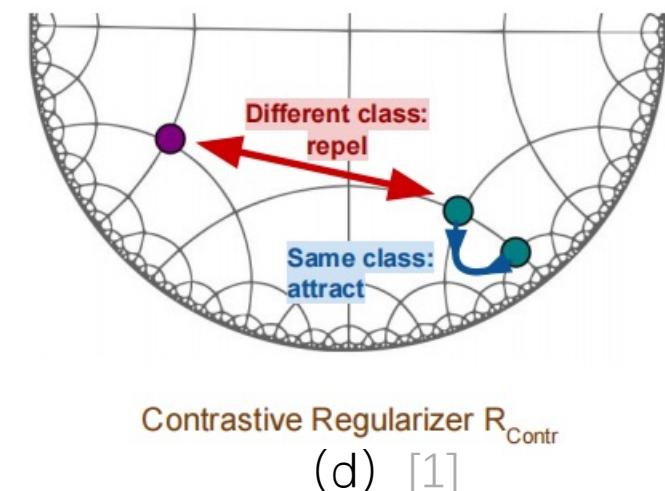
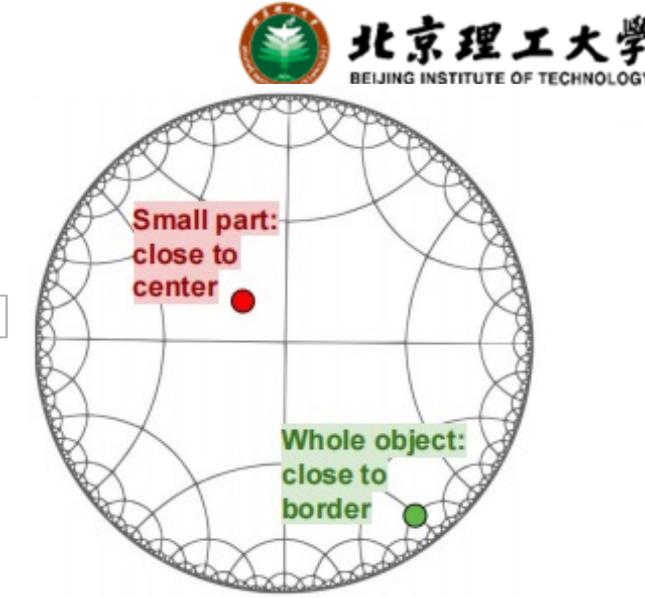
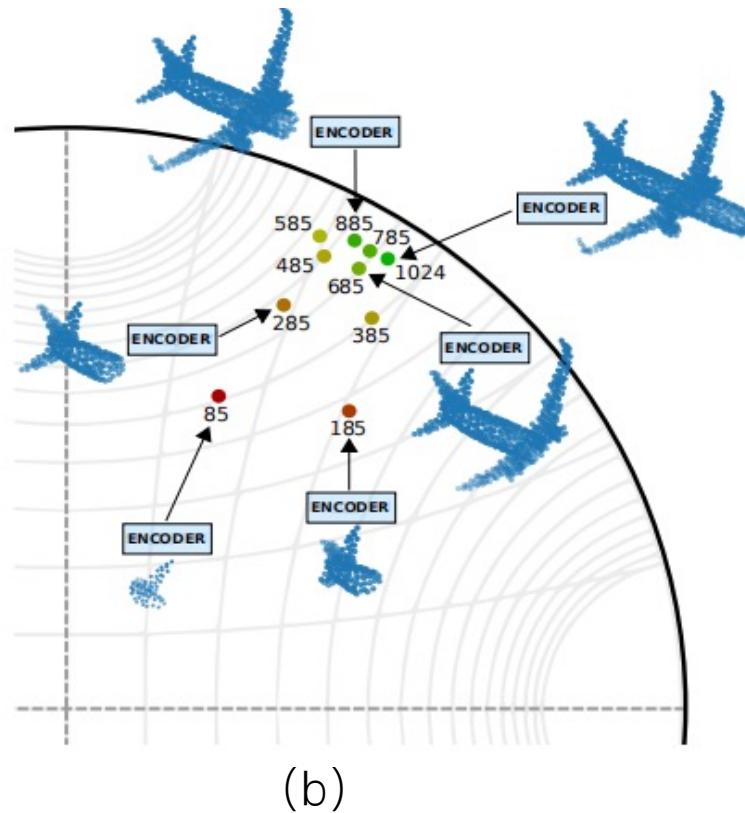
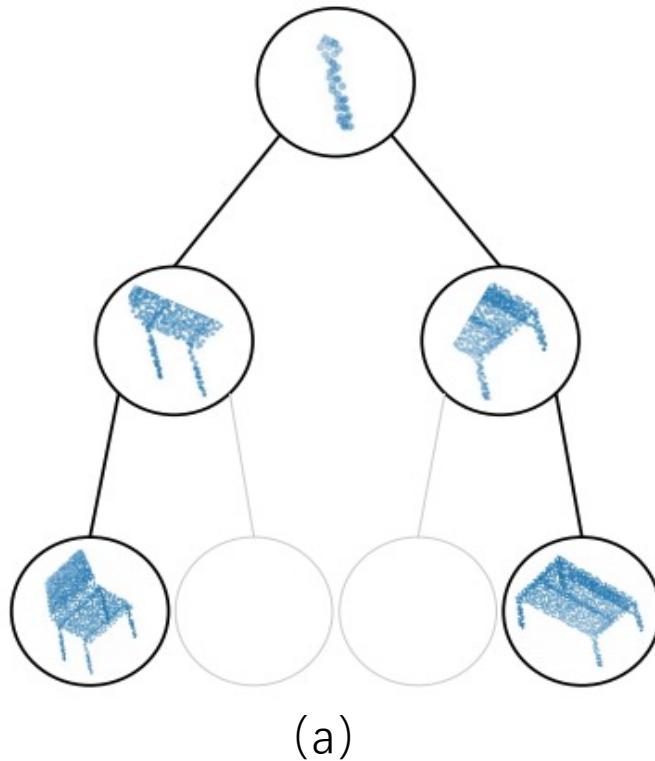
# Segmentation

Improve segmentation through hyperbolic embeddings with uncertainty.

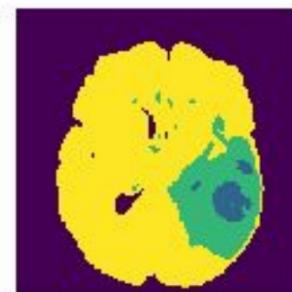
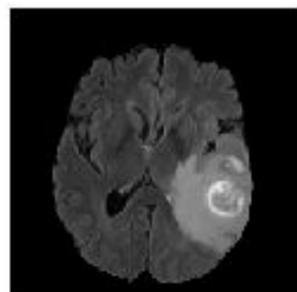


# 3D Vision

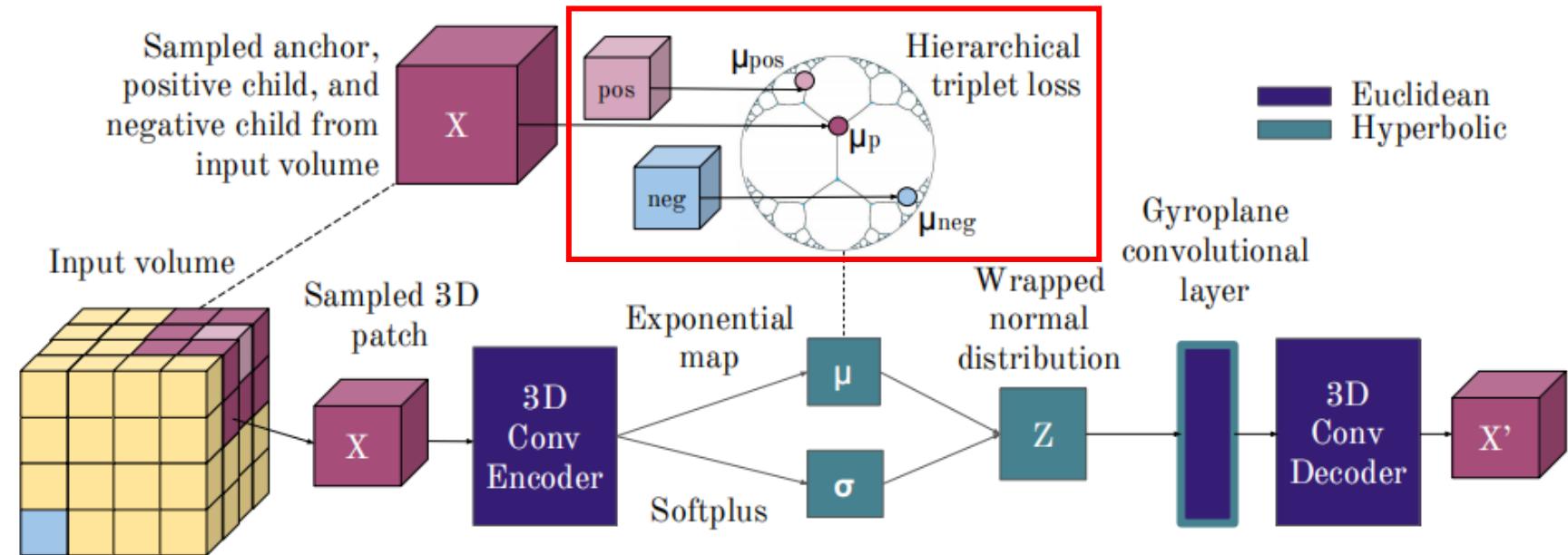
Compositional part-whole hierarchy in 3D objects. [1]



Capturing implicit hierarchical structure in 3D biomedical images.



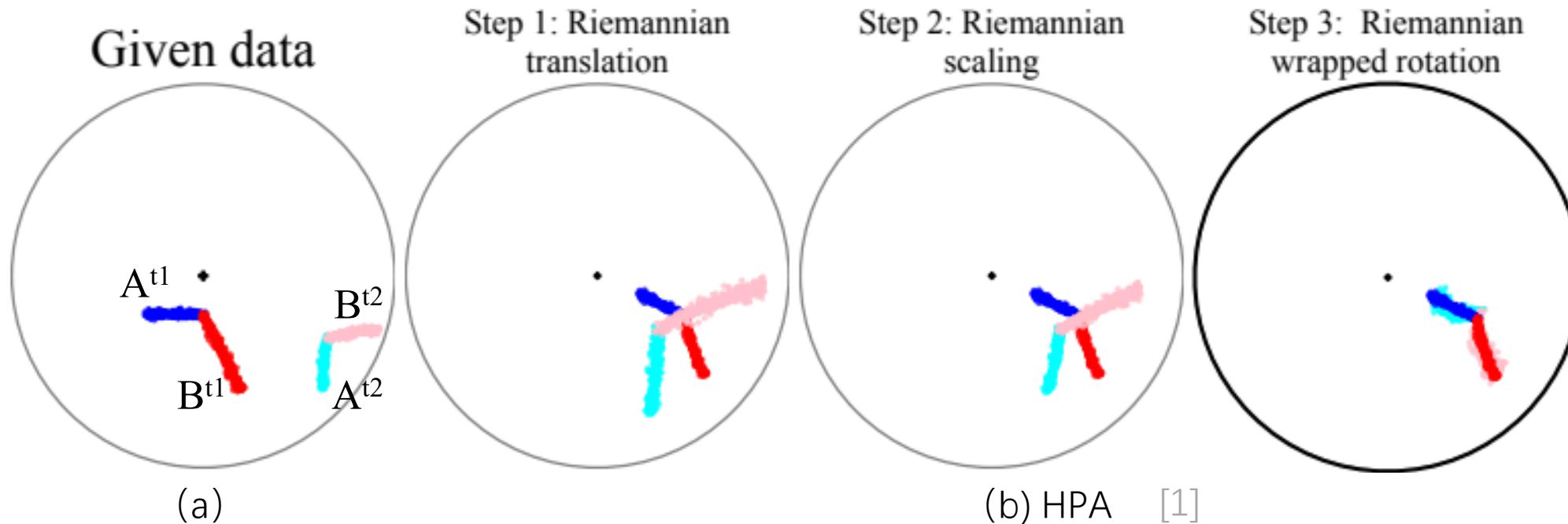
(a) BraTS 2019 [1]

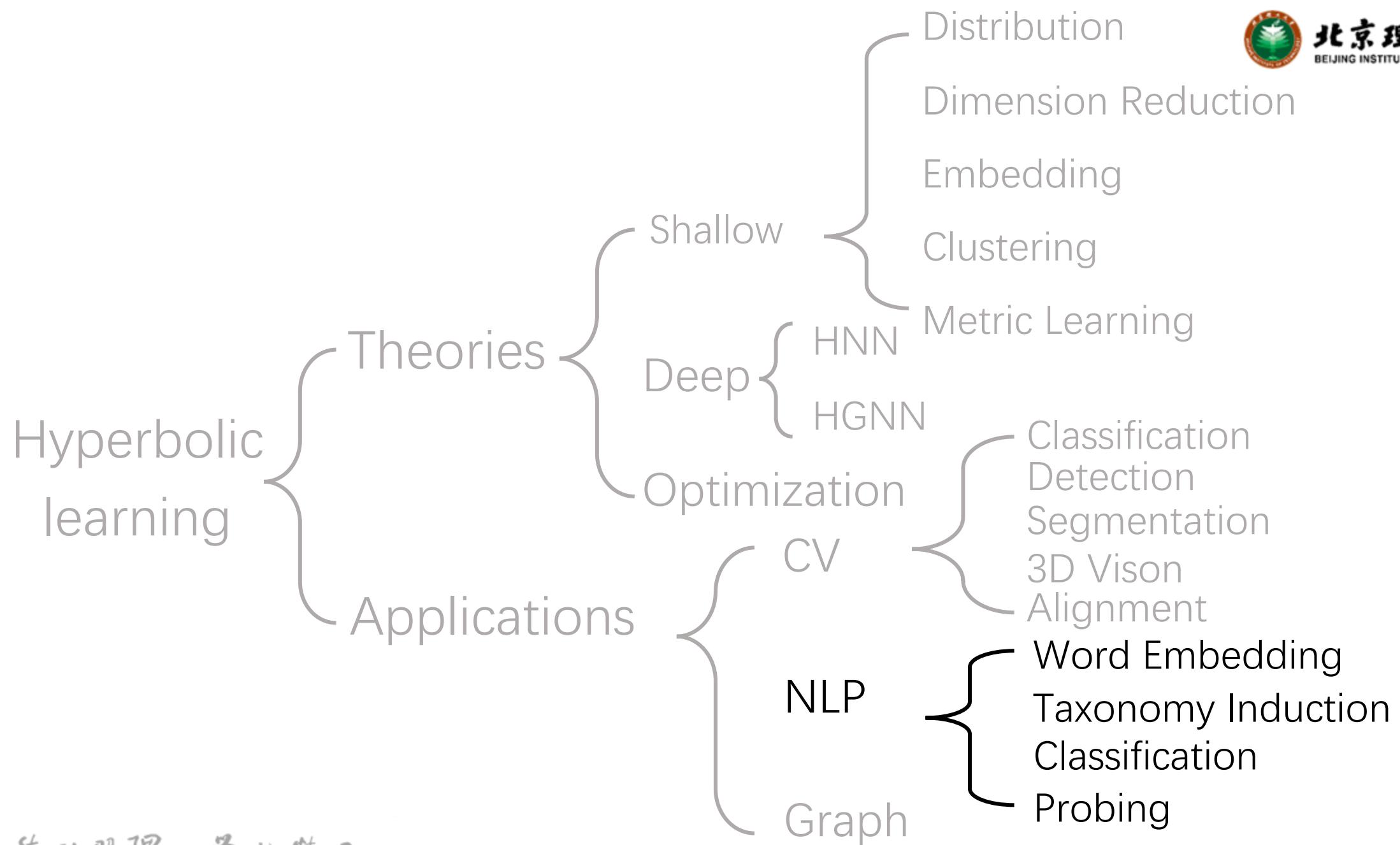


(b) 3D hyperbolic VAE [1]

# Alignment

Hyperbolic Procrustes analysis (HPA) for label-free alignment of hierarchical datasets.



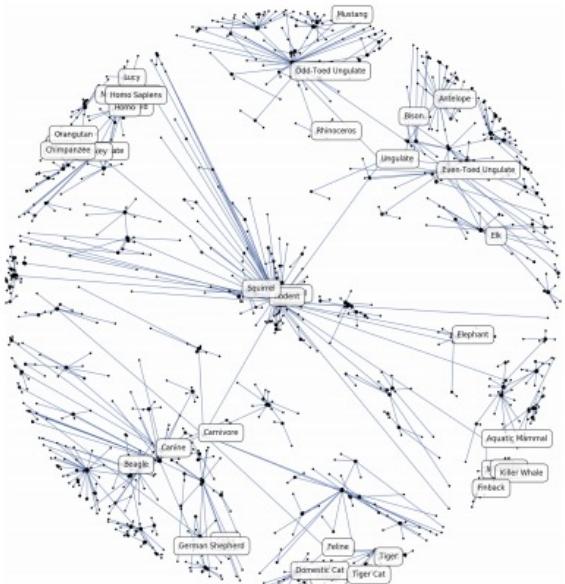


# Application

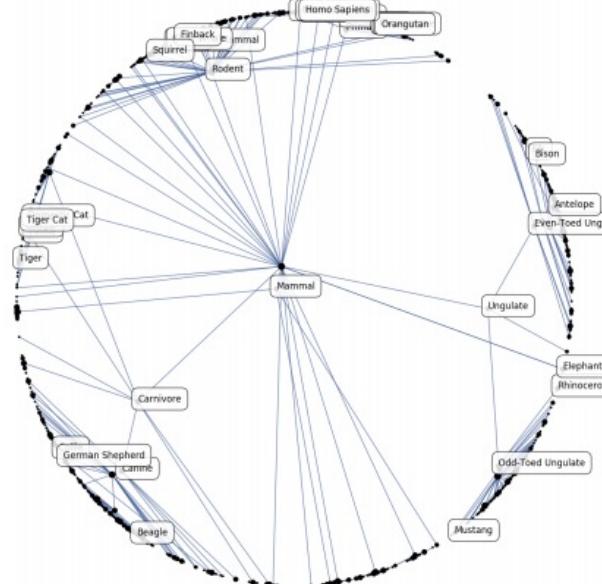
## NLP

# Word Embedding

Embed the WordNet Noun hierarchy into the Poincaré ball.<sup>[1]</sup>



(a) Intermediate embedding after 20 epochs

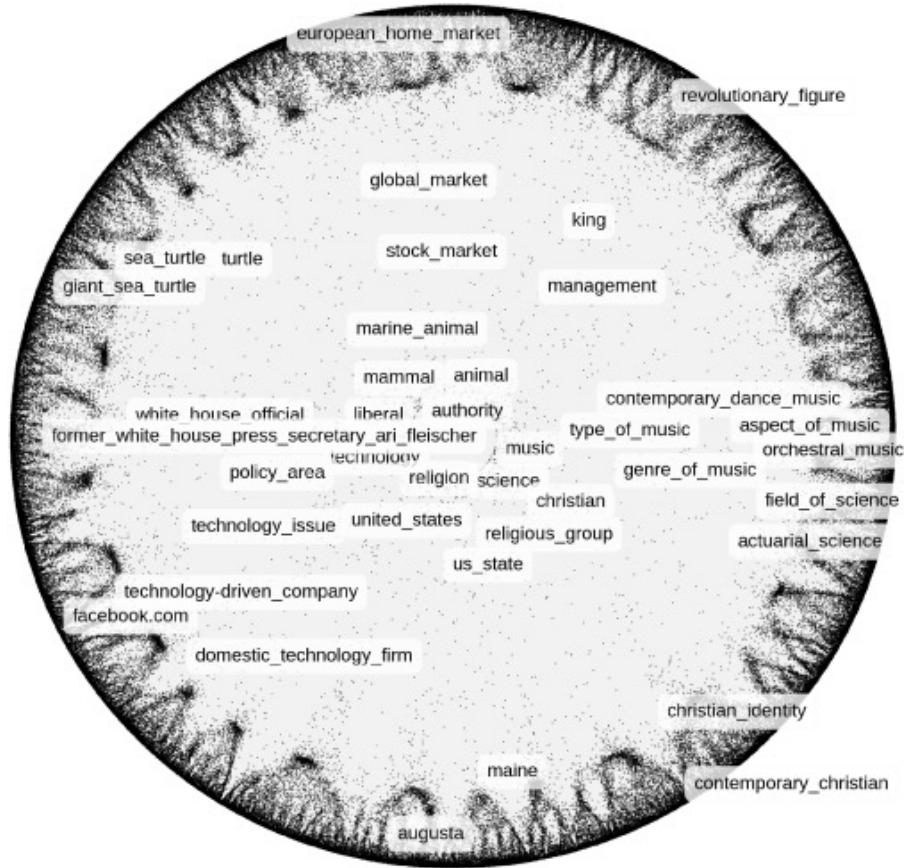


(b) Embedding after convergence

$$\mathcal{L}(\Theta) = \sum_{(u,v) \in \mathcal{D}} \log \frac{e^{-d(u,v)}}{\sum_{v' \in \mathcal{N}(u)} e^{-d(u,v')}}$$

# Word Embedding

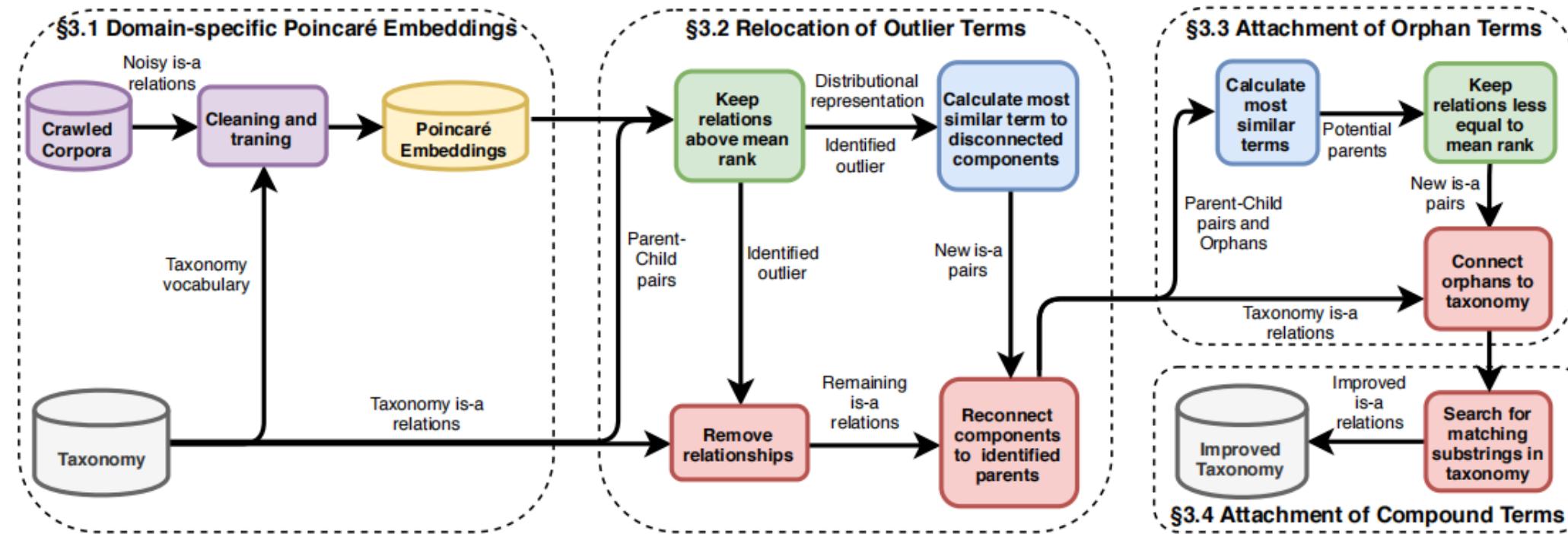
Hyperbolic embedding of Hearst Graph.<sup>[1]</sup>



- extract potential is-a relationships using Hearst patterns
- build Hearst Graph
- embed Hearst Graph in hyperbolic space

# Taxonomy Induction

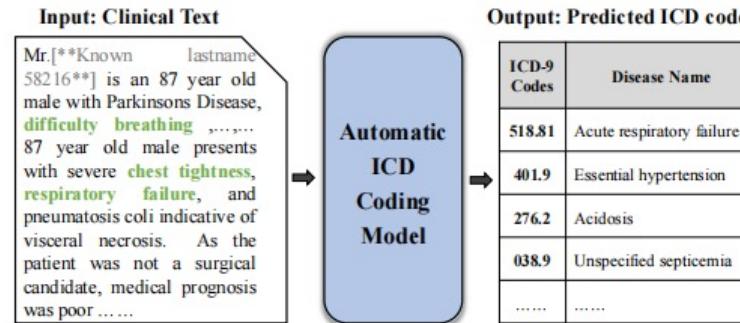
Taxonomy Refinement using Hyperbolic Word Embeddings.<sup>[1]</sup>



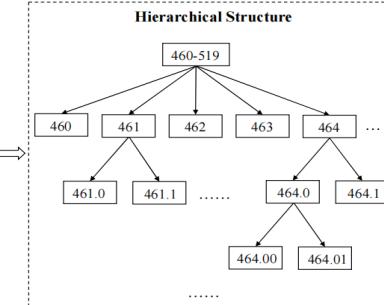
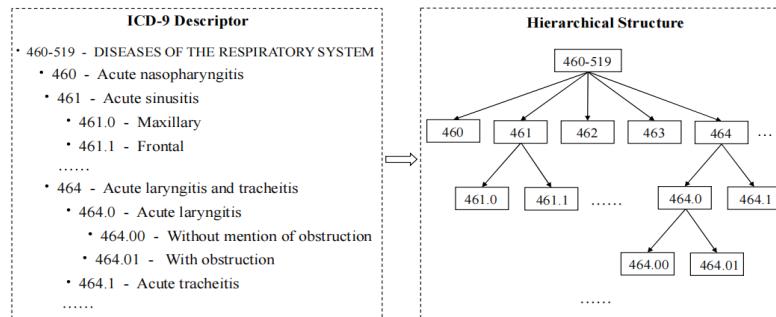
Better capture hierarchical lexical-semantic relationships

# Classification

Automatic ICD Coding——a multi-label classification task.



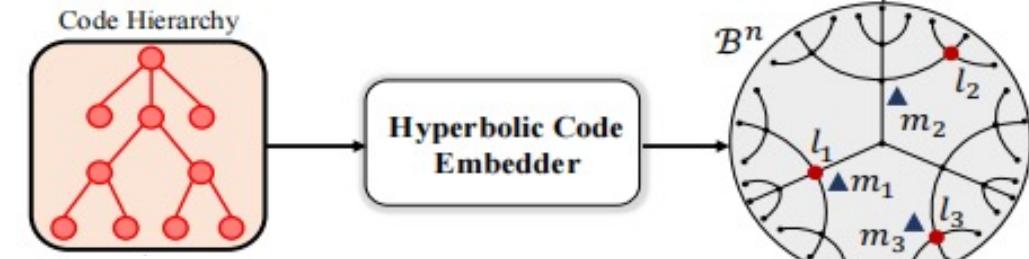
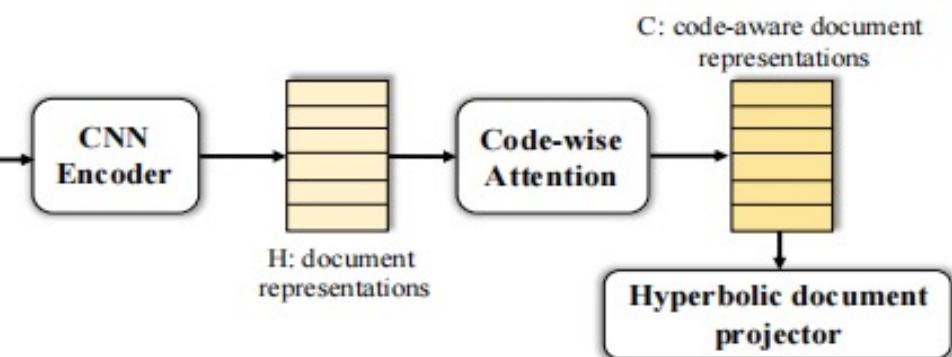
(a) ICD coding[1]



(b) Code Hierarchy[1]

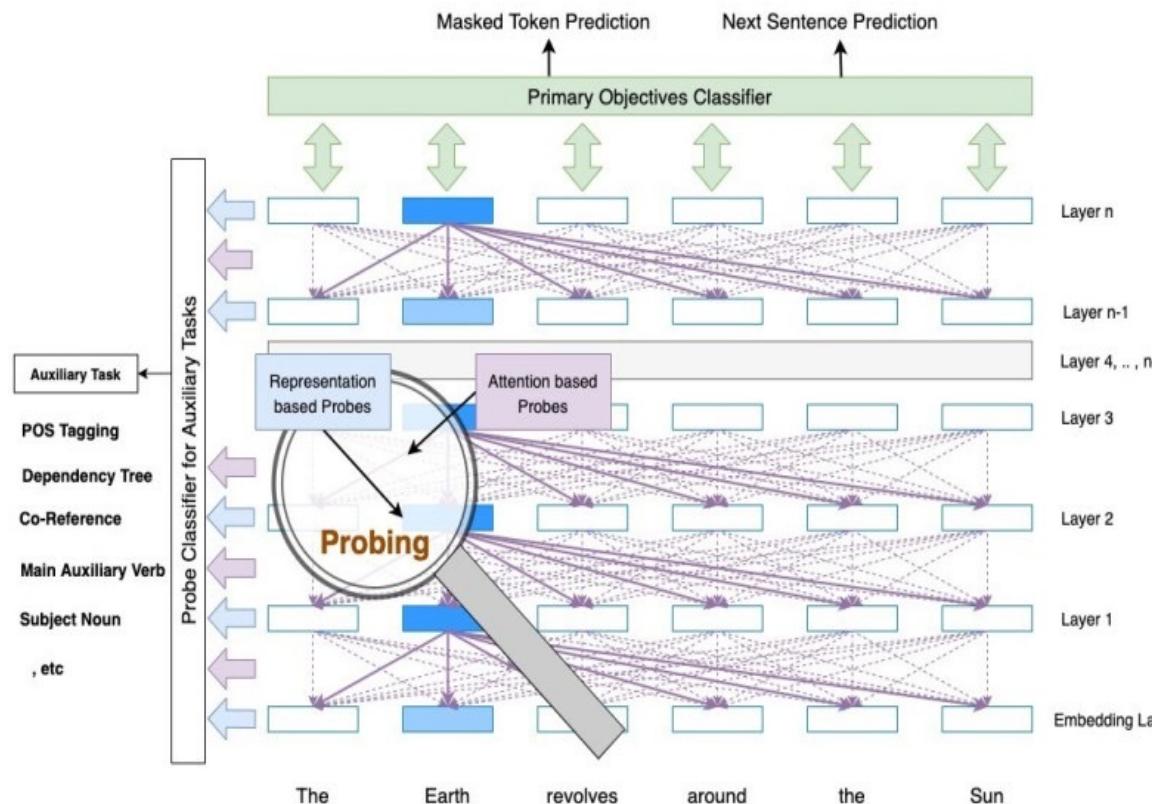
**A Clinical Text**

This was a 51 year old woman who entered via the emergency room after a fall. She was transferred from an outside hospital ...



(c) Hyperbolic module [1]

# Probing

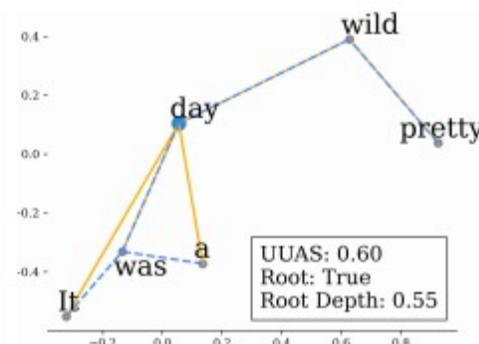


(a) Probing BERT[1]

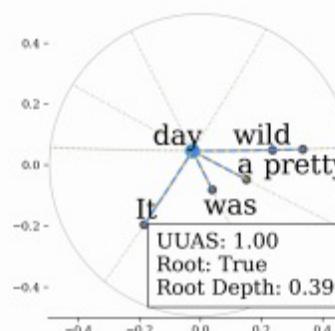
Poincaré probe discovers more syntactic information.

$$p_i = \exp_0(Ph_i)$$

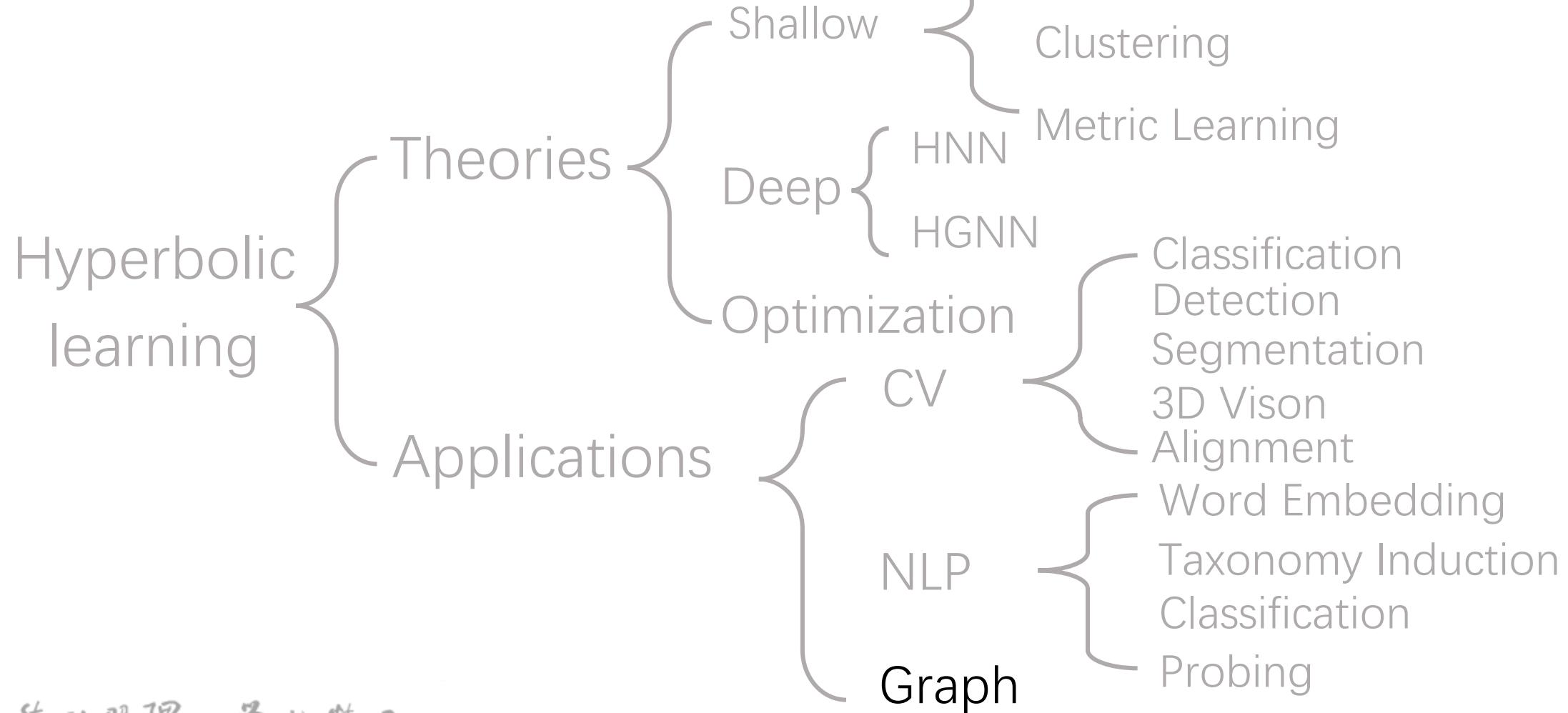
$$q_i = Q \otimes_c p_i$$



(C) BERTBASE7



(D) BERTBASE7 [2]

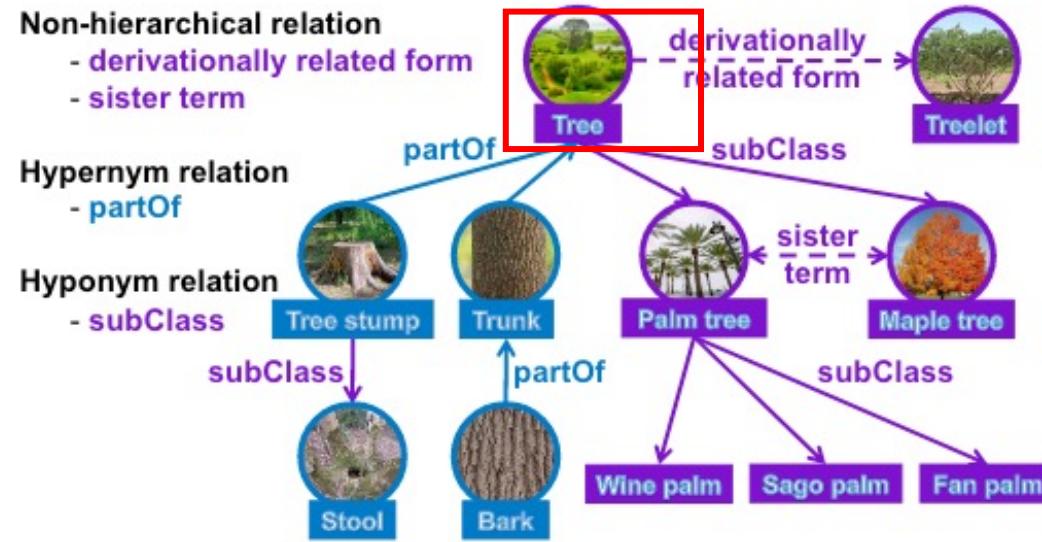


# Application

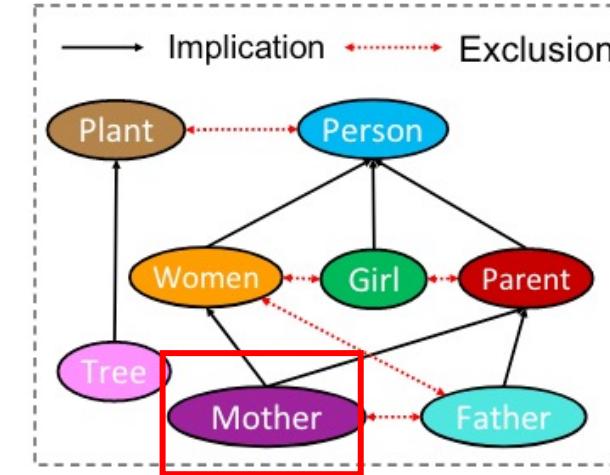
## Graph

# Knowledge Graph Embedding

## Multi-relation or Multi-label



Modeling Heterogeneous Hierarchies<sup>[1]</sup>



Hyperbolic structured multi-label prediction<sup>[1]</sup>

[1] Modeling Heterogeneous Hierarchies with Relation-specific Hyperbolic Cones, Yushi Bai et al. NIPS 2021.

[2] Hyperbolic Embedding Inference for Structured Multi-Label Prediction. Bo Xiong et al. NIPS 2022.

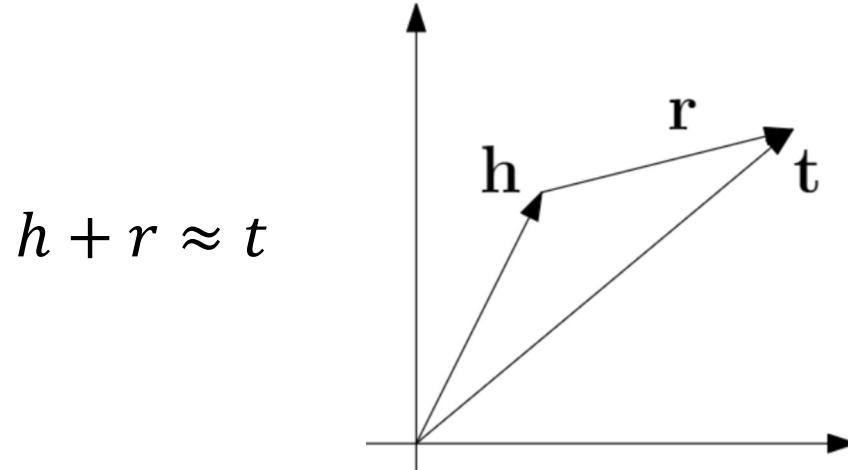
# Knowledge Graph Embedding

Solve multi-relation problem in hyperbolic KGE.

**Bilinear models:** regard relations as linear transformations.

For a triple  $(h, r, t)$ ,  $hW_r \approx t$

**Translational models:** regard relations as translations.



# Knowledge Graph Embedding

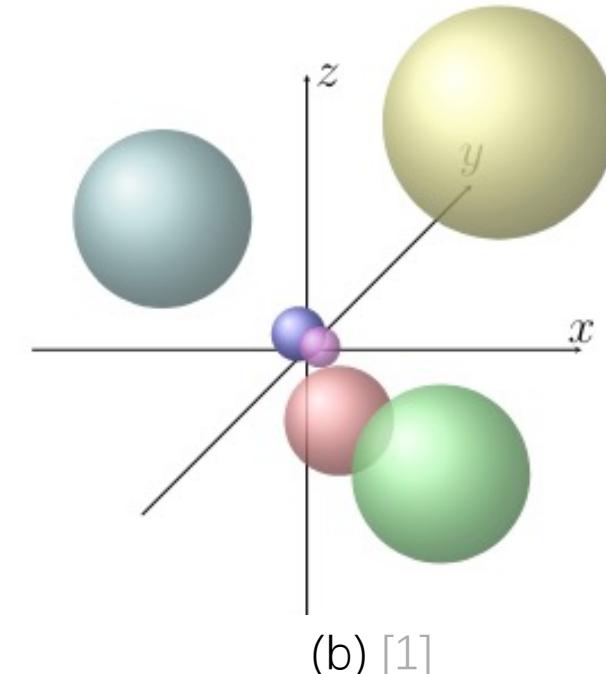
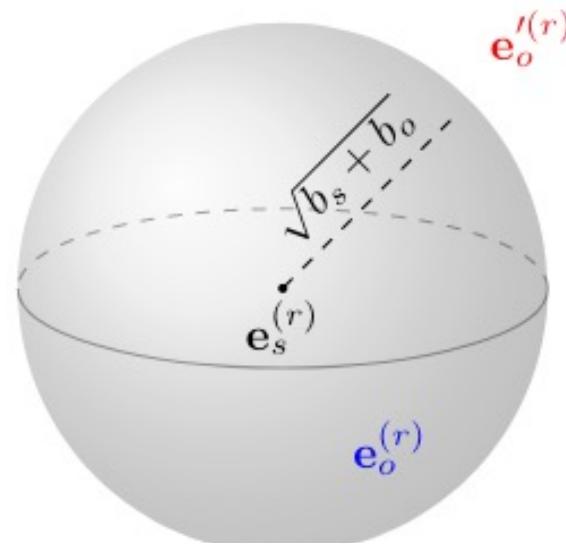
Different “balls” for different hierarchies.

Basic score function

$$\begin{aligned}\phi(e_s, r, e_o) &= -d(\mathbf{e}_s^{(r)}, \mathbf{e}_o^{(r)})^2 + b_s + b_o \\ &= -d(\mathbf{R}\mathbf{e}_s, \mathbf{e}_o + \mathbf{r})^2 + b_s + b_o\end{aligned}\quad (1)$$

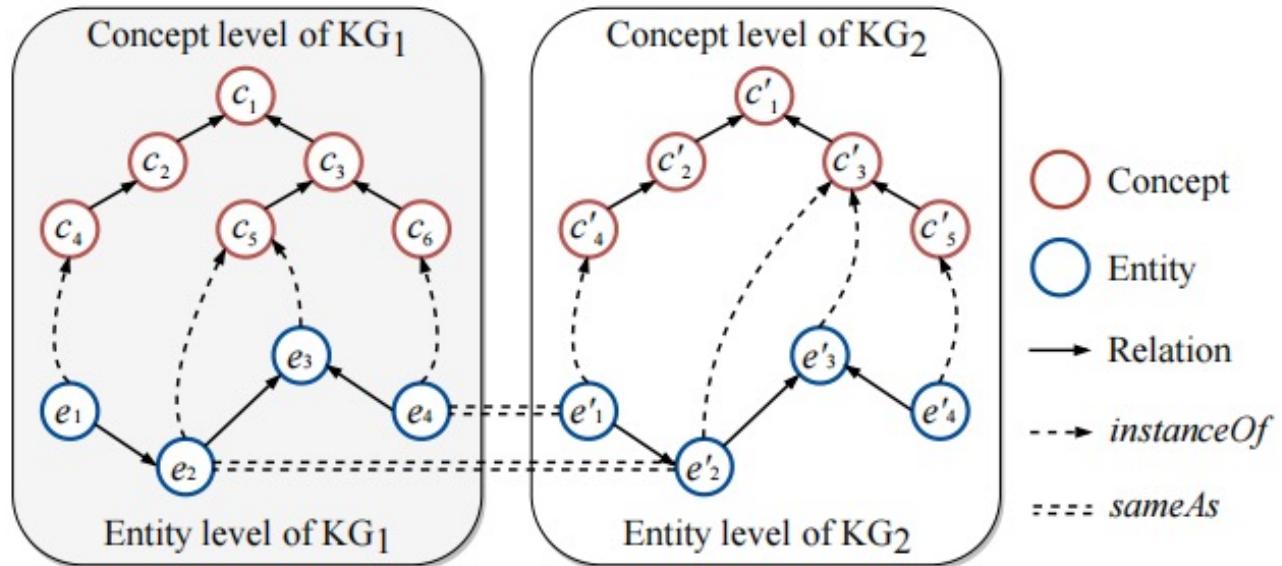
Hyperbolic score function

$$\begin{aligned}\phi_{\text{MuRP}}(e_s, r, e_o) &= -d_{\mathbb{B}}(\mathbf{h}_s^{(r)}, \mathbf{h}_o^{(r)})^2 + b_s + b_o \\ &= -d_{\mathbb{B}}(\exp_0^c(\mathbf{R}\log_0^c(\mathbf{h}_s)), \mathbf{h}_o \oplus_c \mathbf{r}_h)^2 + b_s + b_o\end{aligned}\quad (2)$$



# Knowledge Graph Embedding

Hyperbolic knowledge association with translation in the hyperbolic space.



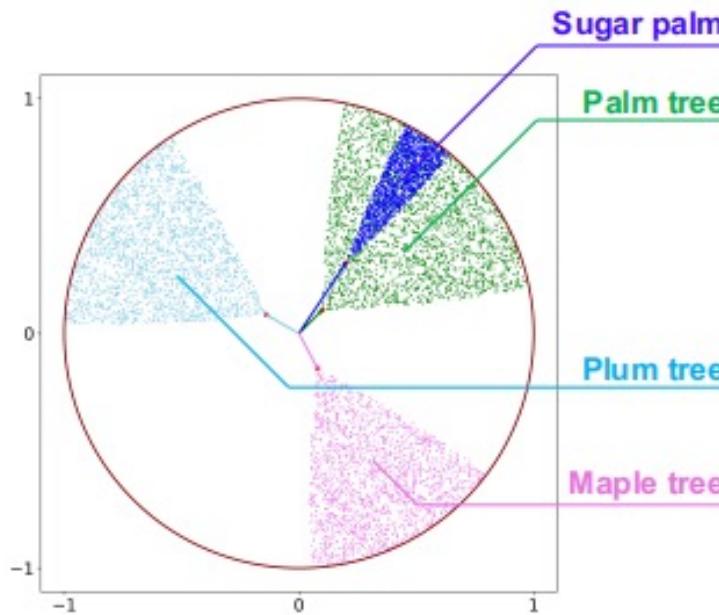
$$f(\tau) = d_{\mathbb{D}}(\mathbf{u}_h^{(0)} \oplus \mathbf{u}_r^{(0)}, \mathbf{u}_t^{(0)}) \quad (1)$$

Figure 1: Illustration of two kinds of knowledge associations (i.e., *sameAs* and *instanceOf*) in KGs. [1]

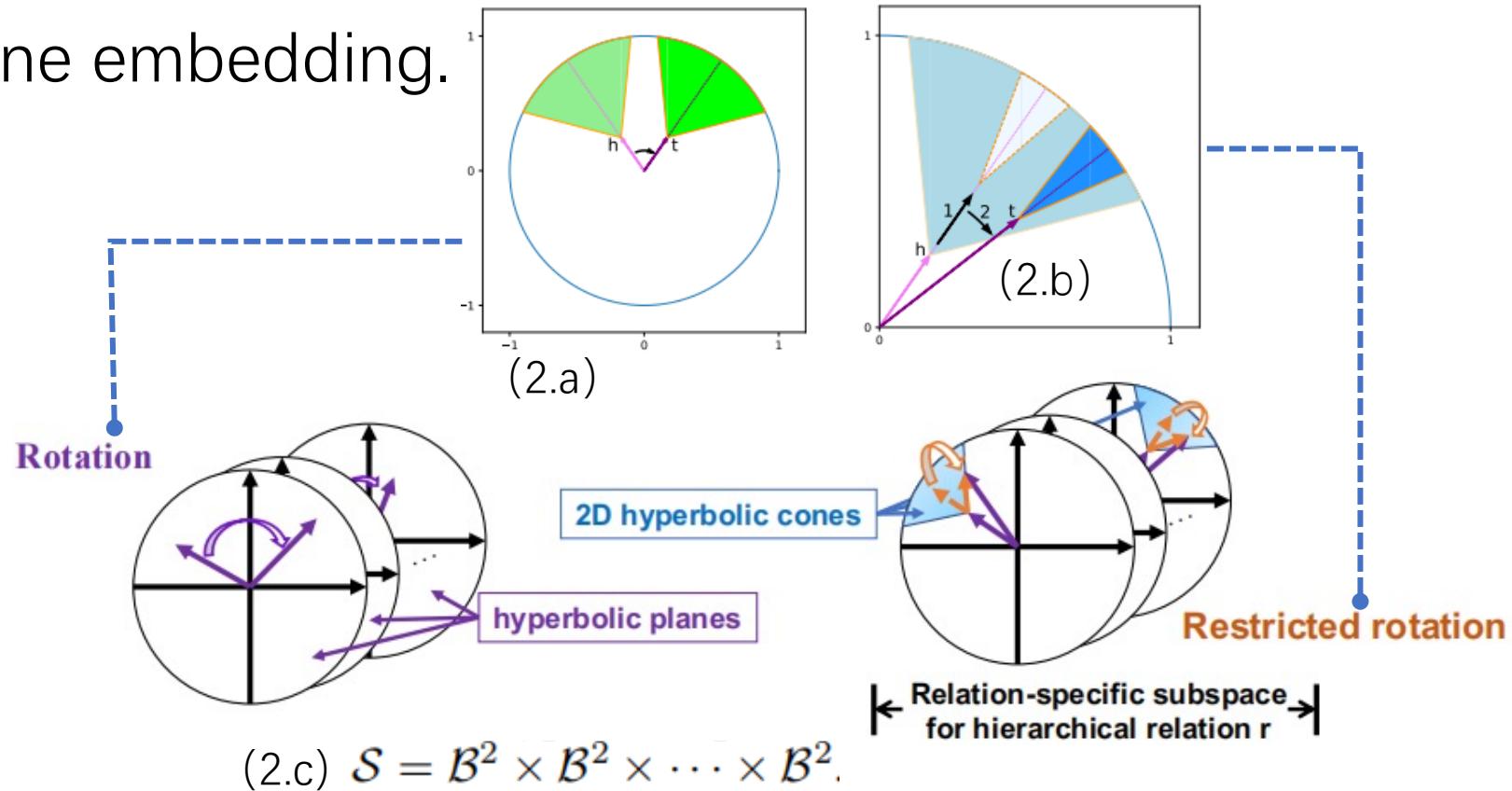
# Knowledge Graph Embedding



Model hierarchies by cone embedding.



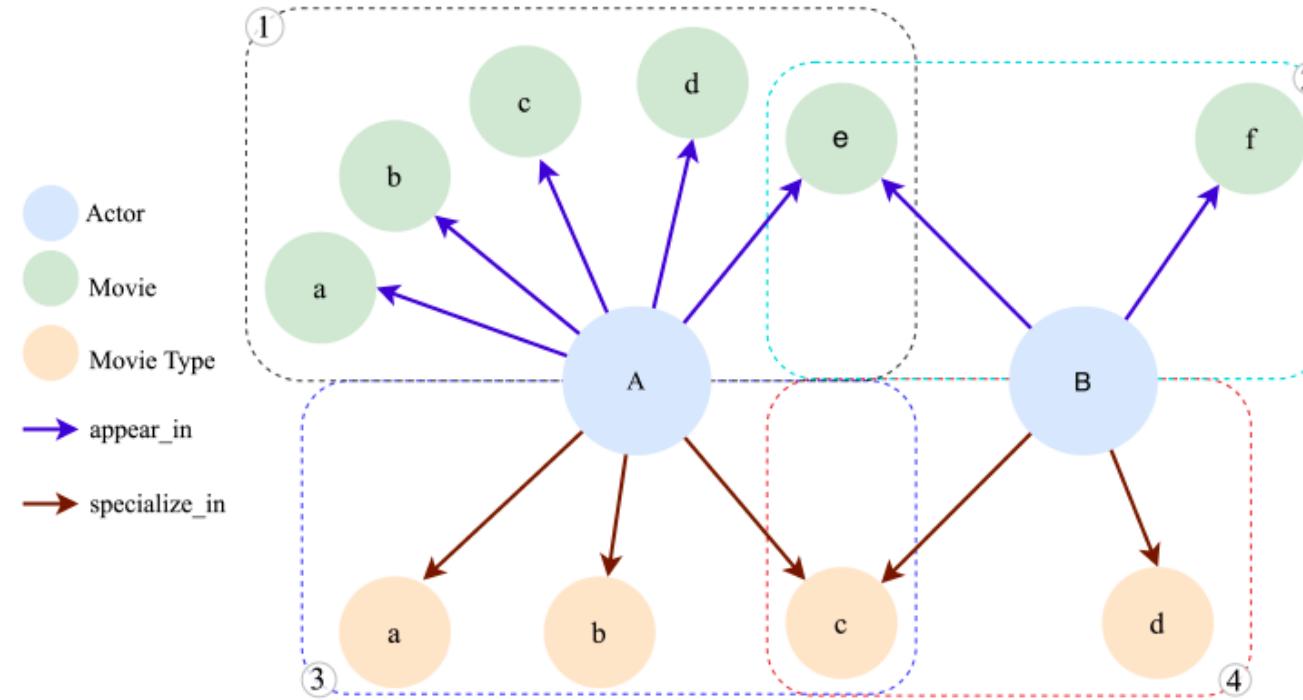
(1) Hyperbolic entailment cones [1]



(2) ConE Model[1]

# Knowledge Graph Embedding

Choose suitable curvatures for various hierarchies. [1]



$$s(h, r, t) = -d_{\mathbb{B}}^{c_{h,r}} (\mathbf{h}_{e,rot}^{\mathbf{H}} \oplus^{c_{h,r}} \boldsymbol{\varepsilon}_r^{\mathbf{H}}, \mathbf{t}^{\mathbf{H}})^2 + b_h + b_t$$

[1] Hyperbolic Hierarchical Knowledge Graph Embeddings for Link Prediction in Low Dimensions. Wenjie Zheng, et al. EMNLP 2021.

# Challenges

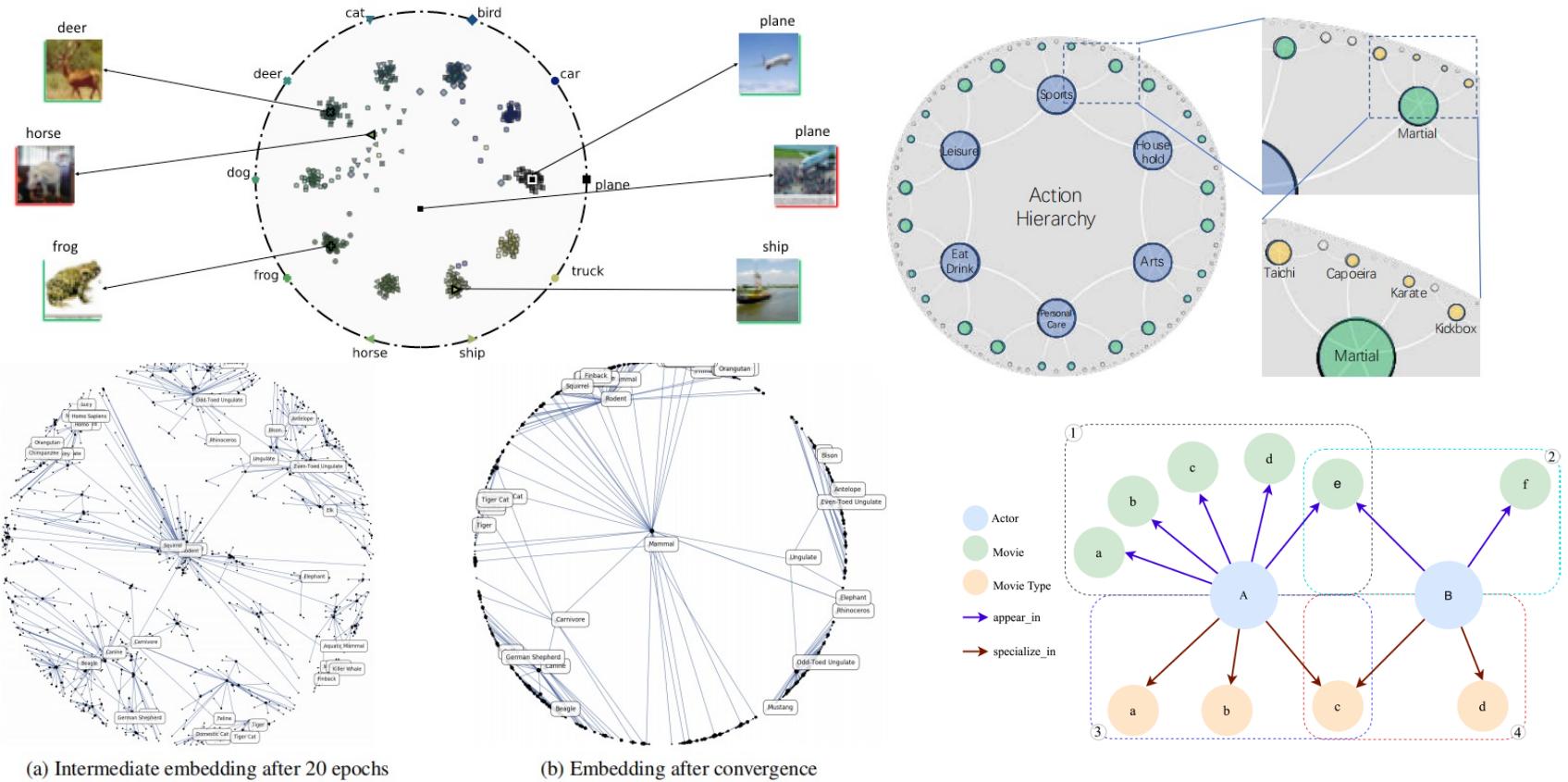
# Challenges in Hyperbolic learning



- Currently only limited to representation learning.
- Objective functions remains limited to hyperbolic distance, which is equivalent to L1 norm. There is scope for development of more complex classification and regression objectives.
- **Unstable** training of hyperbolic networks..
- Learn to **adjust curvature** to data and problem.
- **Large-scale** hyperbolic learning.

# Challenges in Hyperbolic learning

➤ Currently only **limited to** representation learning.



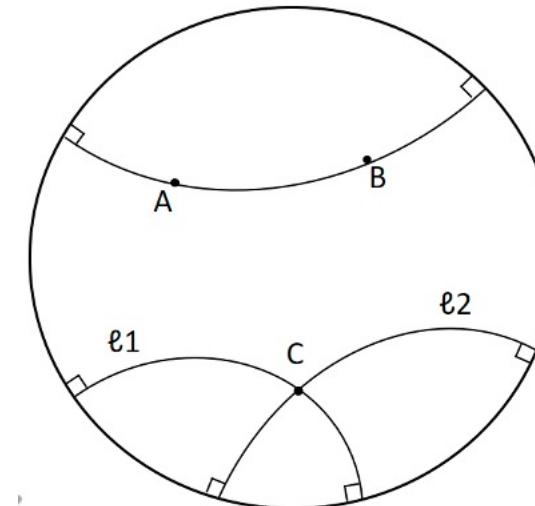
德以明理 学以精工

[1] Hyperbolic Networks: Theory, Architectures and Applications. Department of Computer Science, Virginia Tech, Arlington, VA.

# Challenges in Hyperbolic learning

- Currently only **limited to** representation learning
- Objective functions remains limited to hyperbolic distance, which is equivalent to L1 norm. There is scope for development of **more suitable** classification and regression objectives.

$$d(p, q) = \frac{1}{\sqrt{\kappa}} \text{arcosh} \left( 1 + \frac{2|p - q|^2}{(1 - |p|^2)(1 - |q|^2)} \right)$$



# Challenges in Hyperbolic learning



- Currently only **limited to** representation learning
- Objective functions remains limited to hyperbolic distance, which is equivalent to L1 norm. There is scope for development of **more suitable** classification and regression objectives.
- **Unstable** training of hyperbolic networks.

# Challenges in Hyperbolic learning



- Currently only **limited to** representation learning
- Objective functions remains limited to hyperbolic distance, which is equivalent to L1 norm. There is scope for development of **more suitable** classification and regression objectives.
- Further study into the **hyperbolic gradient descent** is needed to provide new techniques for stable training.
- Learn to **adjust curvature** to data and problem.

# Challenges in Hyperbolic learning



- Currently only **limited to** representation learning
- Objective functions remains limited to hyperbolic distance, which is equivalent to L1 norm. There is scope for development of **more suitable** classification and regression objectives.
- Further study into the **hyperbolic gradient descent** is needed to provide new techniques for stable training.
- Learn to **adjust curvature** to data and problem.
- **Large-scale** hyperbolic learning.

# Q & A

Thanks for your listening!