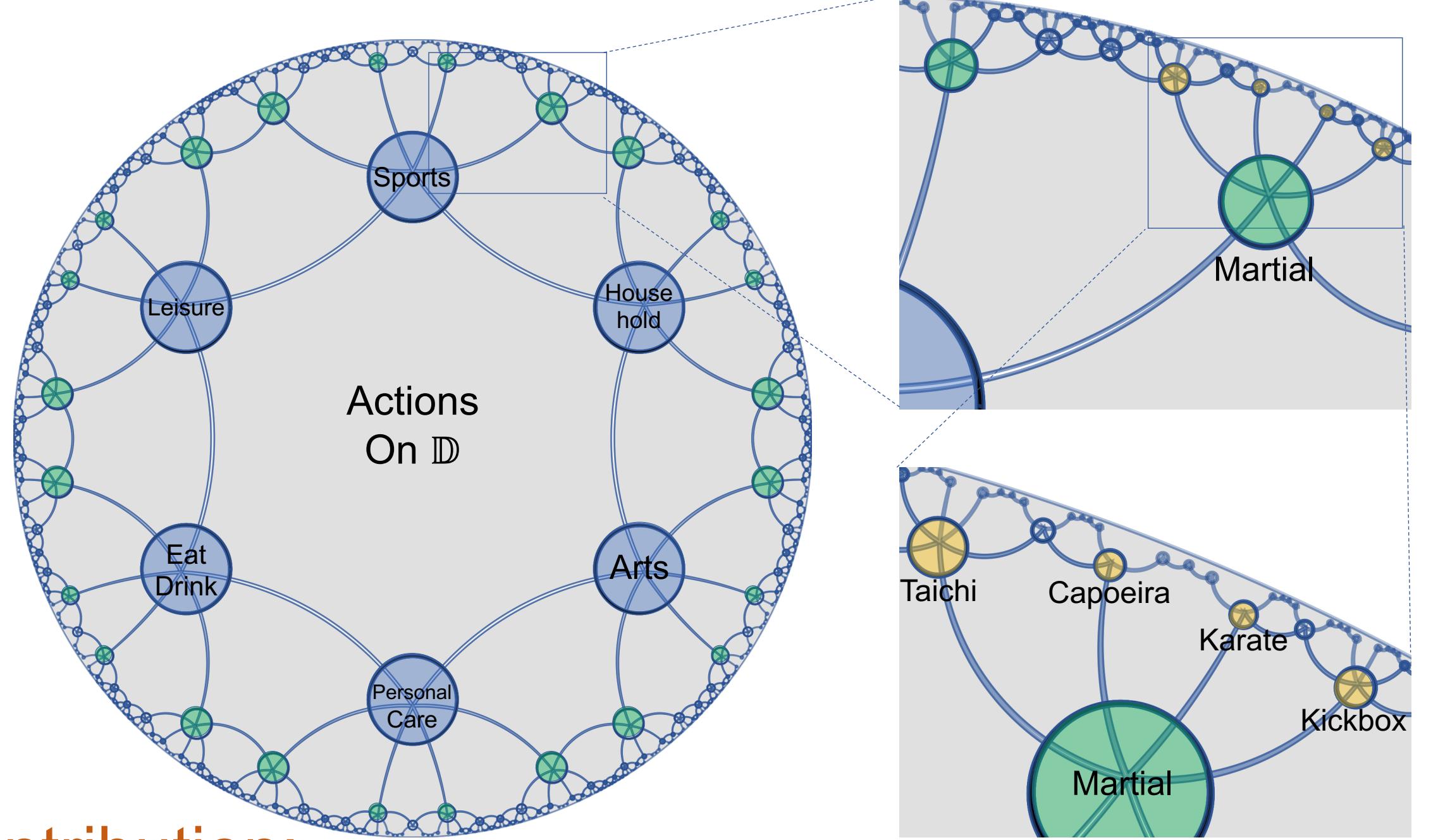


Searching for Actions on the Hyperbole

Teng Long, Pascal Mettes, Heng Tao Shen, Cees Snoek

Problem

Can hierarchical knowledge benefit action retrieval?



Contribution:

- We position action hierarchy with discriminative hyperbolic embeddings
- We propose hyperbolic matching between actions and videos on the hyperbole \mathbb{D}
- We perform hierarchical action search by name/video, on supervised/zero-shot setting.

Preliminary of \mathbb{D}

Mobius Addition $a \oplus_c b$:

The vector addition in \mathbb{D} with curvature c

$$\frac{(1 + 2c\langle a, b \rangle + c\|b\|^2) a + (1 - c\|a\|^2) b}{1 + 2c\langle a, b \rangle + c^2\|a\|^2\|b\|^2}$$

Exponential Map $\exp_x(v)$:

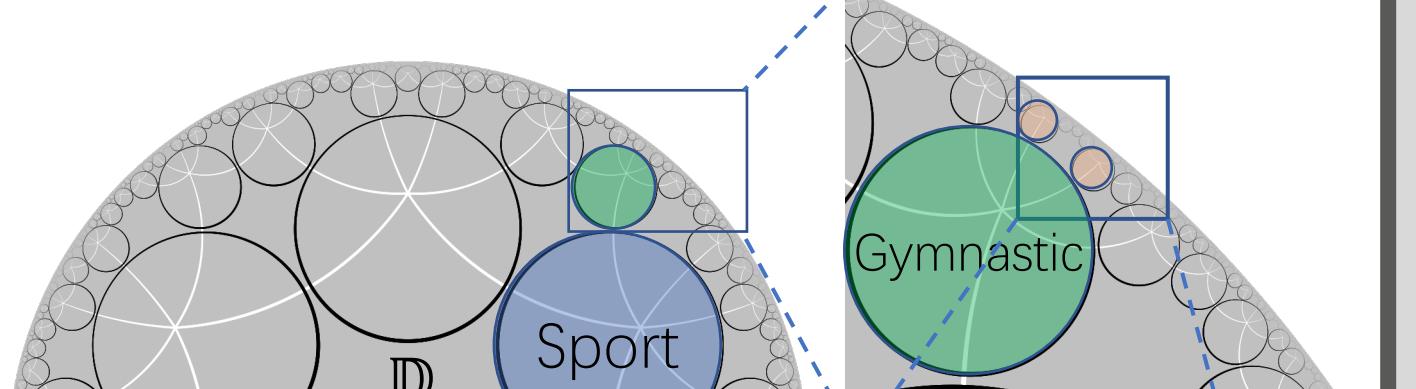
The transform from \mathbb{R} to \mathbb{D}

$$x \oplus_c \left(\tanh\left(\frac{\sqrt{c}\lambda_x\|v\|}{2}\right) \frac{v}{\sqrt{c}\|v\|} \right)$$

Conformality: angular invariance for \mathbb{D}/\mathbb{R}

$$\frac{g_x^\mathbb{D}(u, v)}{\sqrt{g_x^\mathbb{D}(u, u)}\sqrt{g_x^\mathbb{D}(v, v)}} = \frac{\langle u, v \rangle}{\|u\|\|v\|}$$

Why \mathbb{D} suits hierarchy?



Riemannian optimization:

Gradient Descent based on Riemannian gradient ∇_R

$$\mathbf{P}_{t+1} = \mathbf{P}_t - \eta_t \nabla_R \mathcal{L}(\mathbf{P}_t)$$

Search

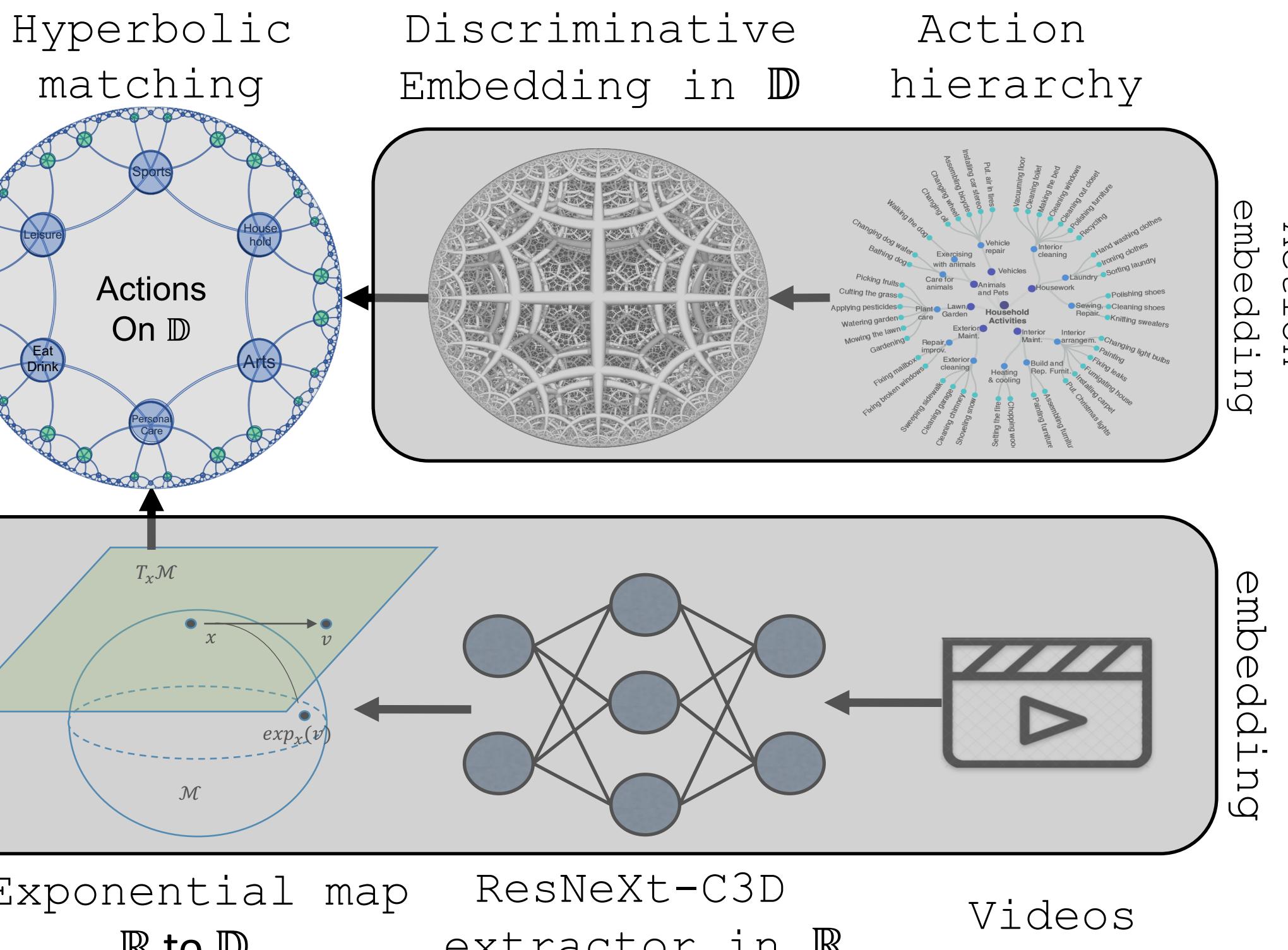
Cosine distance is a fast approximation for distance ranking in \mathbb{D}

Code:



$$d_q(\mathbf{x}_i) = 1 - \cos(\mathbf{q}, \mathbf{x}_i)$$

Model



Exponential map \mathbb{R} to \mathbb{D} ResNeXt-C3D extractor in \mathbb{R} Videos

Action hierarchy embedding:

$$\mathcal{L}_1(\mathcal{P}, \mathcal{N}, \mathbf{P}) = \mathcal{L}_H(\mathcal{P}, \mathcal{N}) + \lambda \cdot \mathcal{L}_S(\mathbf{P})$$

\mathcal{L}_H preserve the parent-child relationship

\mathcal{L}_S separate leaf node in \mathbb{D}

\mathcal{L}_2 further construct entailment cone

$$\mathcal{L}_2 = \sum_h E(\mathbf{u}, \mathbf{v}) + \sum_{\neg h} (\gamma - E(\mathbf{u}', \mathbf{v}')^+$$

Matching videos to actions:

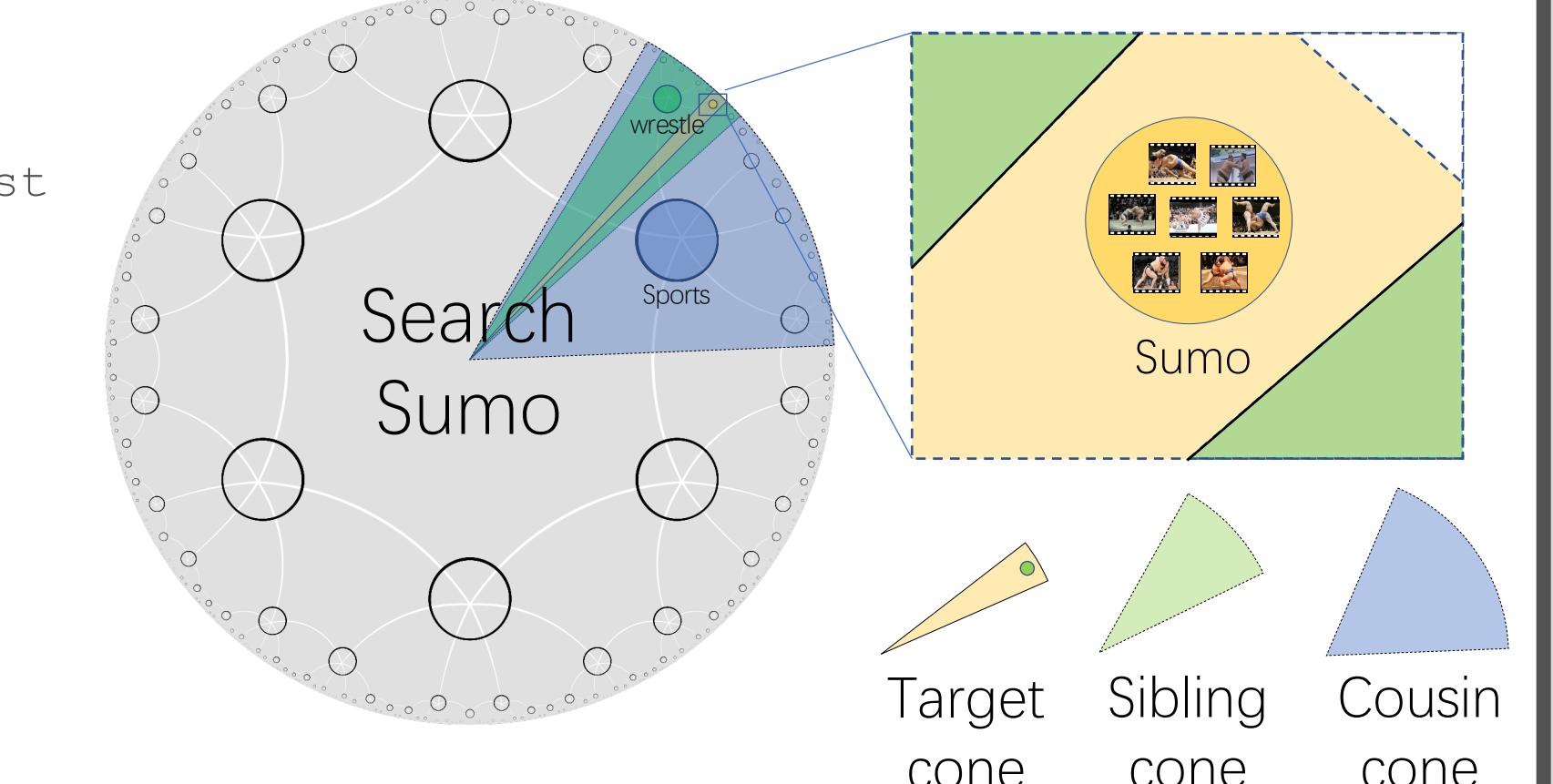
Softmax classifier based on hyperbolic distance

$$d_c(\mathbf{a}, \mathbf{b}) := \frac{2}{\sqrt{c}} \operatorname{arctanh}(\sqrt{c}\|\mathbf{b} \oplus_c \mathbf{a}\|)$$

$$p = \frac{\exp(-d_c(\Psi_e(v; \theta), \phi_c(k)))}{\sum_{k'} \exp(-d_c(\Psi_e(v; \theta), \phi_c(k')))}$$

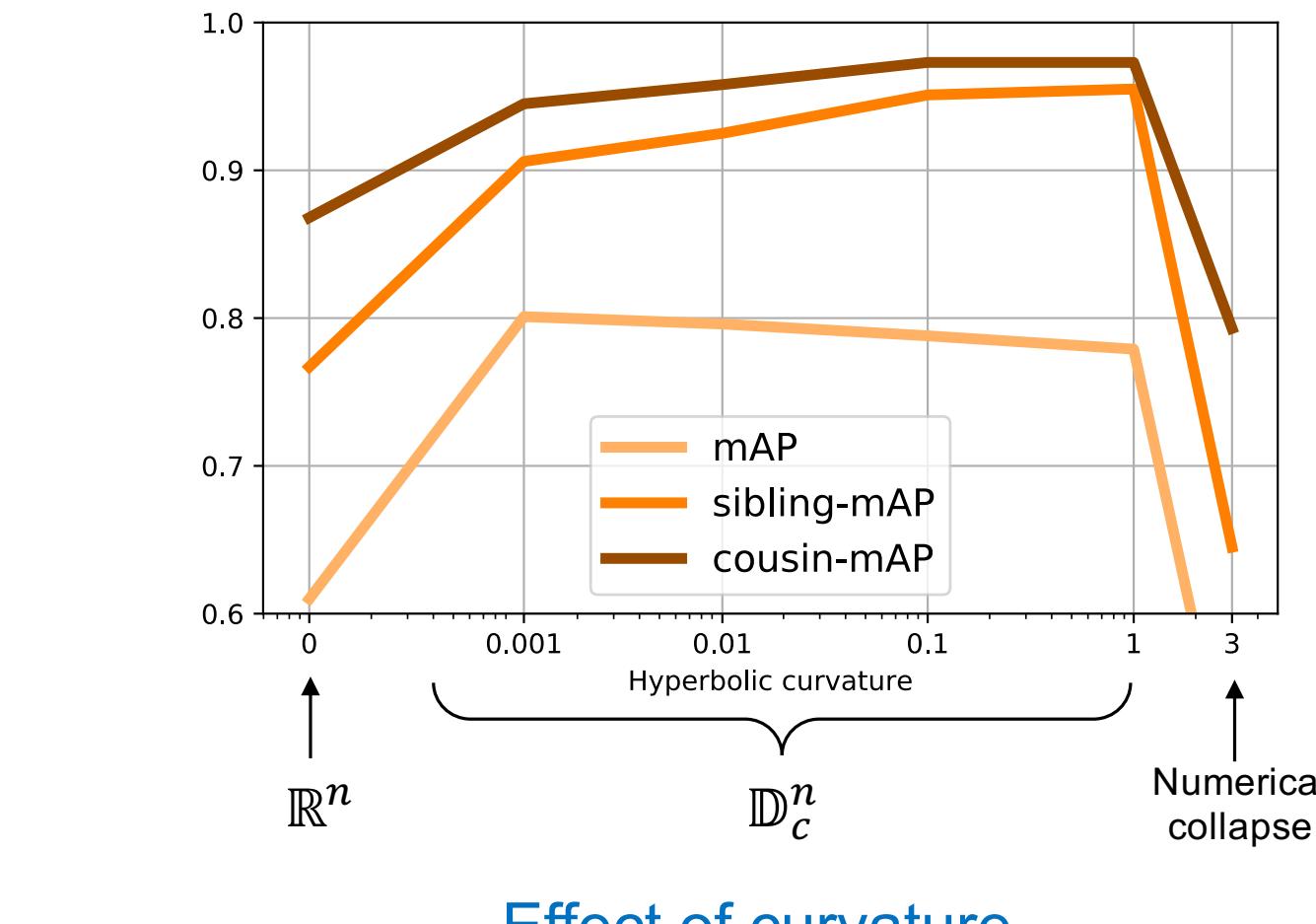
Riemannian optimization:

Gradient Descent based on Riemannian gradient ∇_R



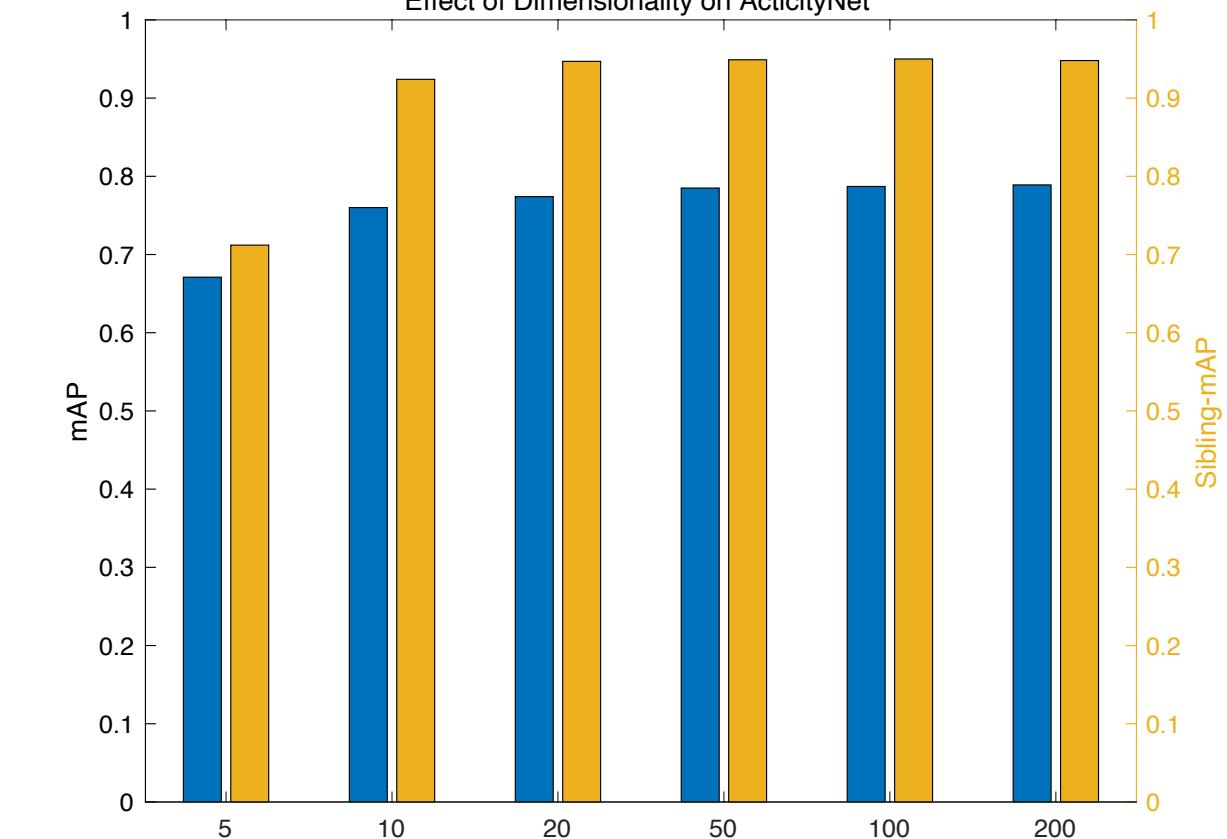
$$d_q(\mathbf{x}_i) = 1 - \cos(\mathbf{q}, \mathbf{x}_i)$$

Results



Effect of curvature

$c \rightarrow 0$: degrade into \mathbb{R}
 $c \rightarrow \infty$: extremely curved \mathbb{D} , numerical collapse



Effect of dimensionality

$d \rightarrow 20$: performance saturated at low dimensionality

| | H-ActivityNet | | H-Kinetics | | H-Moments | |
|----------------------|--------------------|--------------|--------------|--------------|--------------|--------------|
| space | mAP | S-mAP | mAP | S-mAP | mAP | S-mAP |
| ResNextC3D [18] | Δ^n | 0.592 | 0.761 | 0.532 | 0.733 | 0.145 |
| DeViSE [14] | \mathbb{R}^n | 0.609 | 0.761 | 0.553 | 0.715 | 0.134 |
| Li et al. [22] | Δ^n | 0.583 | 0.760 | 0.552 | 0.753 | - |
| Mettes et al. [25] | \mathbb{S}^{n-1} | 0.587 | 0.760 | 0.551 | 0.754 | 0.142 |
| Barz and Denzler [4] | \mathbb{S}^{n-1} | 0.583 | 0.747 | 0.547 | 0.725 | 0.143 |
| This Paper | \mathbb{R}^n | 0.610 | 0.767 | 0.565 | 0.738 | 0.172 |
| This Paper | \mathbb{D}_c^n | 0.678 | 0.843 | 0.593 | 0.824 | 0.163 |
| | | | | | | 0.201 |

\mathbb{D} v.s. Δ^n , \mathbb{R}^n , \mathbb{S}^{n-1} :

Hyperbolic space suits hierarchical search while Benefit plain search too. Other spaces show little.

Qualitative results:



Zero-shot generalization:

| | H-ActivityNet | | H-Moments | |
|----------------------|---------------|--------------|--------------|--------------|
| | mAP | S-mAP | mAP | S-mAP |
| Zhang et al. [42] | 0.397 | 0.449 | 0.026 | 0.027 |
| Li et al. [22] | 0.389 | 0.461 | - | - |
| Mettes et al. [25] | 0.235 | 0.281 | 0.169 | 0.171 |
| Barz and Denzler [4] | 0.453 | 0.527 | 0.216 | 0.219 |
| This Paper | 0.543 | 0.627 | 0.222 | 0.225 |

\mathbb{D} not only benefit supervised search, but also promote zero-shot search.