



# COMPOSITIONAL ENTAILMENT LEARNING FOR HYPERBOLIC VISION-LANGUAGE MODELS

**Avik Pal<sup>1\*</sup>**    **Max van Spengler<sup>1</sup>**    **Guido Maria D'Amely di Melendugno<sup>2</sup>**

**Alessandro Flaborea<sup>2</sup>**    **Fabio Galasso<sup>2</sup>**    **Pascal Mettes<sup>1</sup>**

<sup>1</sup>University of Amsterdam

<sup>2</sup>Sapienza University of Rome

2025.1.13 Minghao Liu

# 目录

CONTENTS

- 01 Author
- 02 Background
- 03 Method
- 04 Experiments





北京大学  
PEKING UNIVERSITY

## PART 02

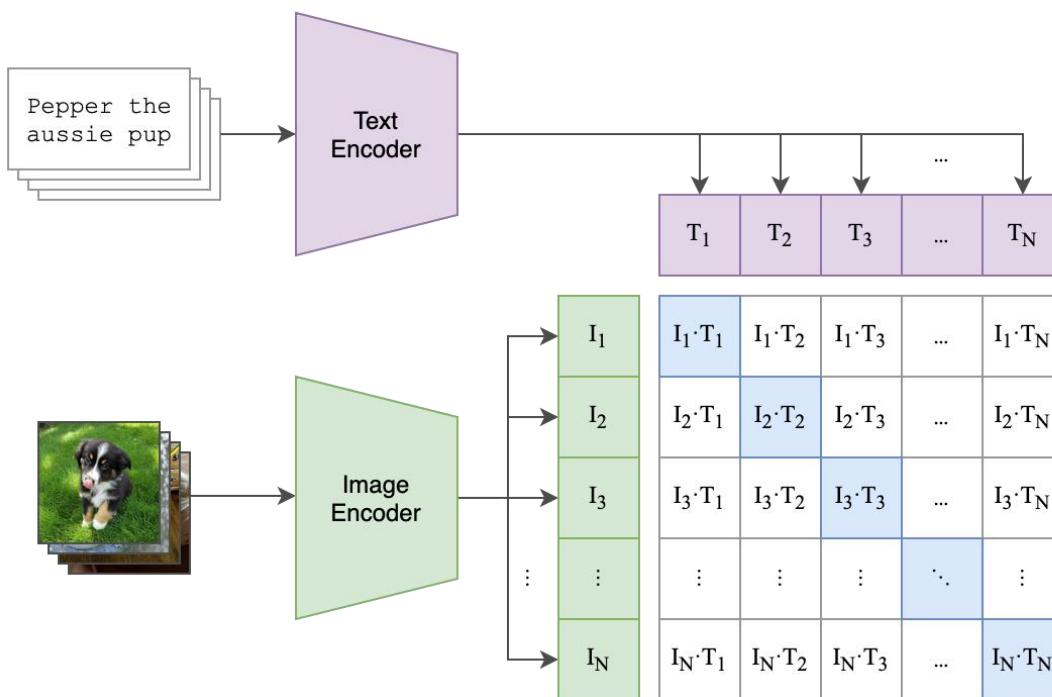
---

# Background

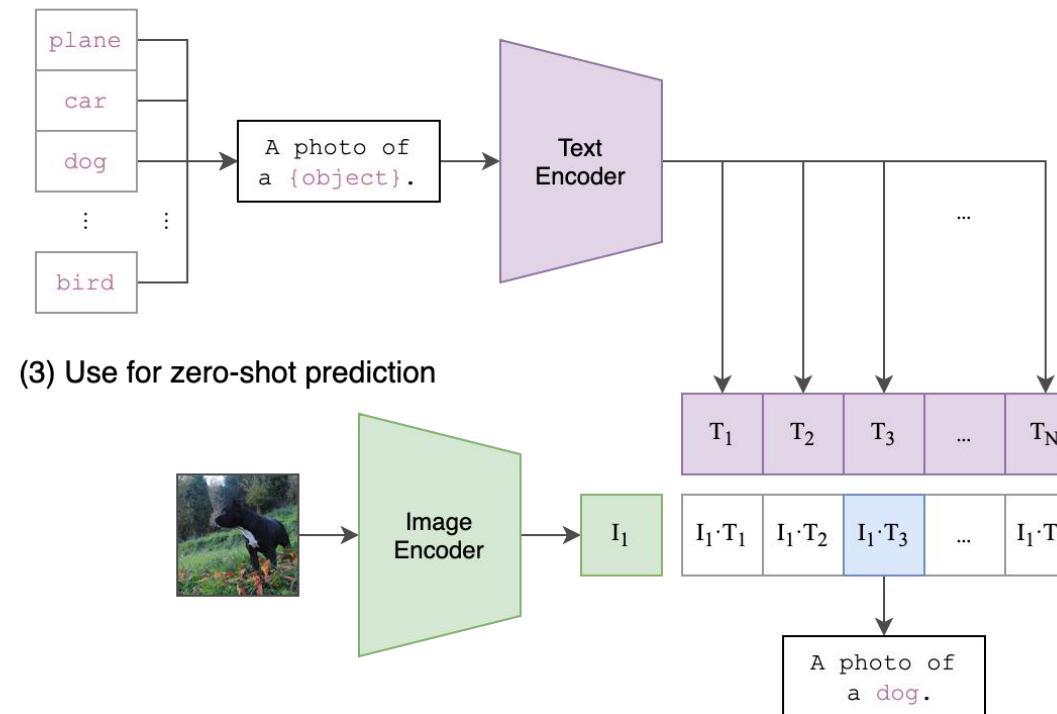


## CLIP (Contrastive Language-Image Pre-Training)

(1) Contrastive pre-training

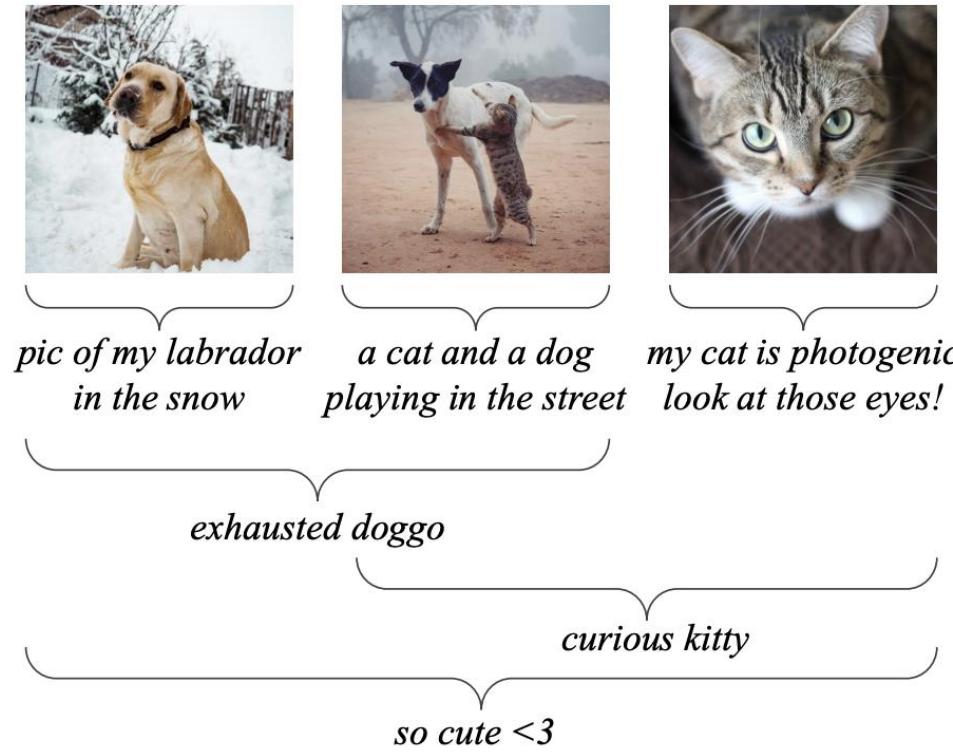


(2) Create dataset classifier from label text

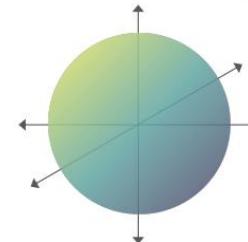


Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International conference on machine learning. PMLR, 2021.

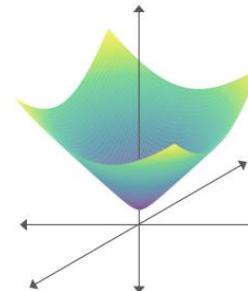
## MERU (Hyperbolic Image-Text Representations)



*CLIP: embed images and text in a Euclidean space*



*MERU: embed images and text in a hyperbolic space*

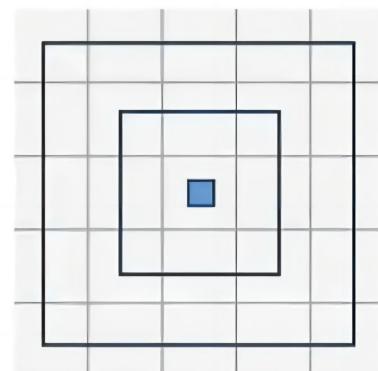


## Why Hyperbolic:

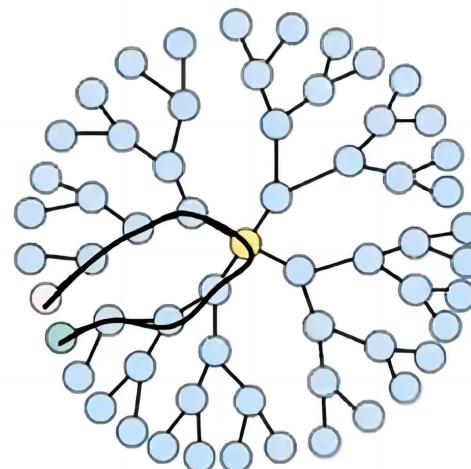
Better capture hierarchy informations

Geometric properties to embed tree-like data

Better representation space



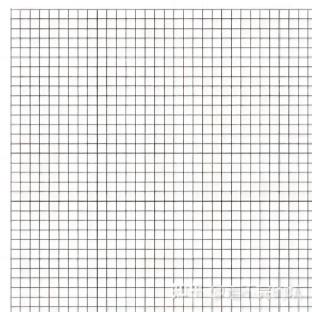
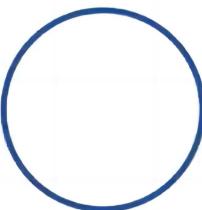
Grid-like  
Polynomial level



Tree-like  
Exponential level

## What is Hyperbolic:

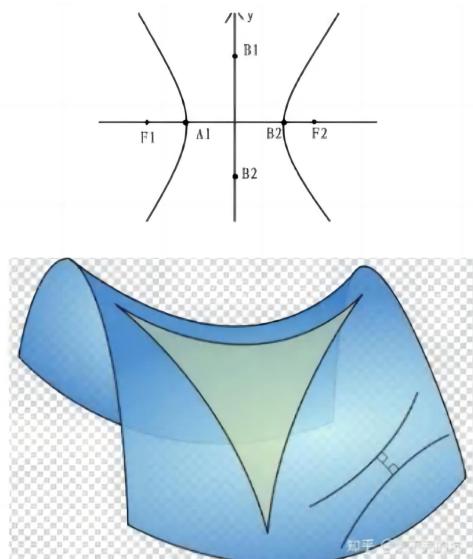
### Negative constant curvature space



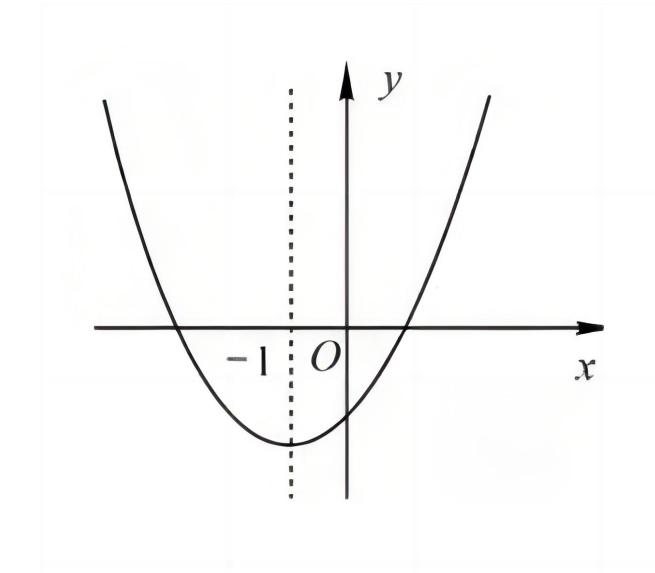
Euclidean  
 $k = 0$



Spherical  
 $k > 0$



Hyperbolic  
 $k < 0$

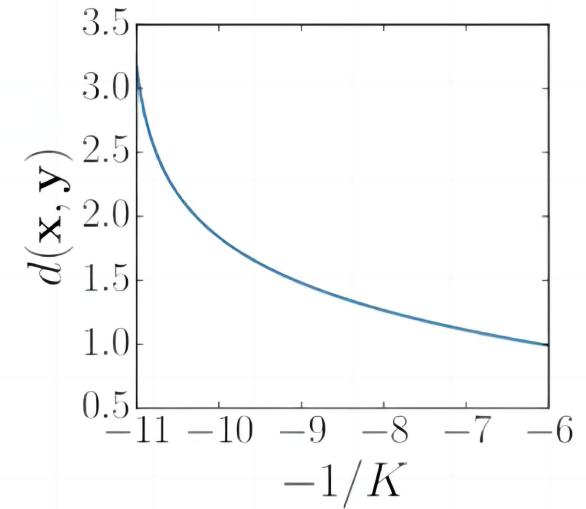
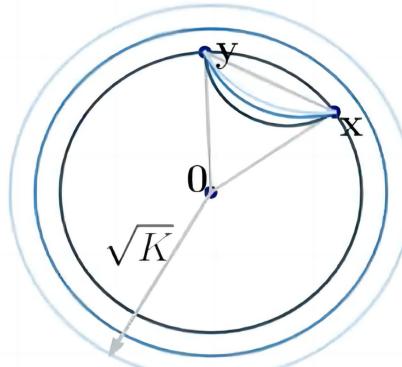
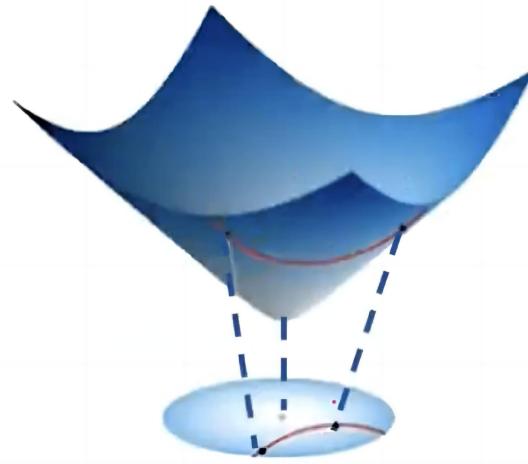


$k$  is not  
constant

# Background

---

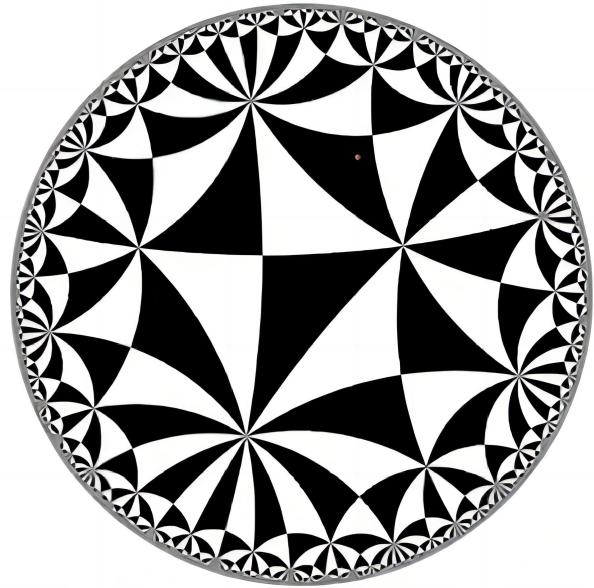
How to Represent Hyperbolic:  
e.g. Poincaré disk model



# Background

---

How to Represent Hyperbolic:  
e.g. Poincaré disk model



Two-dimensional  
Poincare disk model

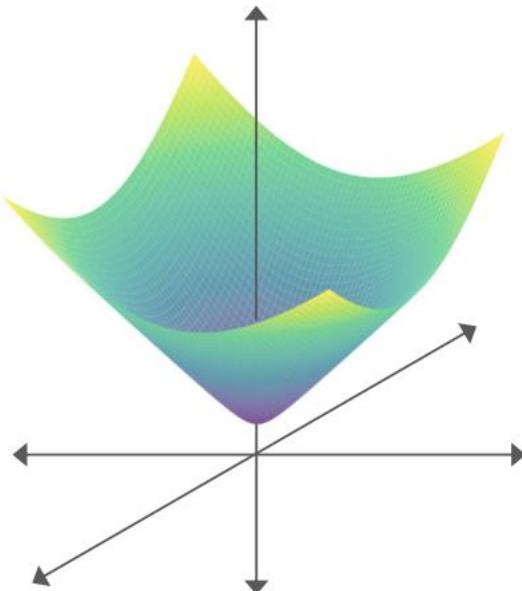
$$d_{\mathbb{D}}^1(\mathbf{x}, \mathbf{y}) = \operatorname{arcosh} \left( 1 + 2 \frac{\|\mathbf{x} - \mathbf{y}\|_2^2}{(1 - \|\mathbf{x}\|_2^2)(1 - \|\mathbf{y}\|_2^2)} \right).$$

# Background

---

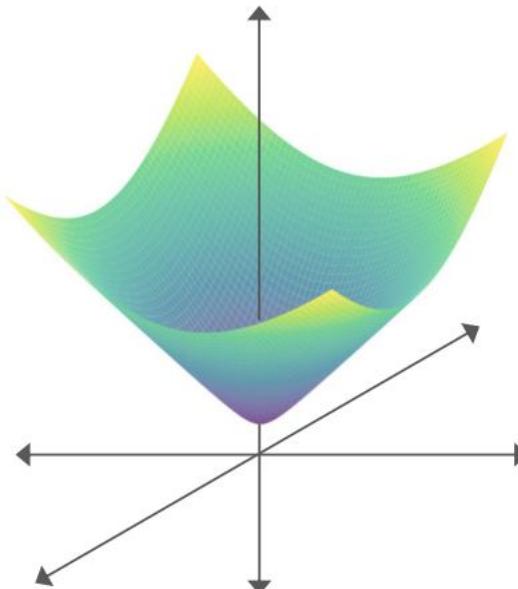
MERU uses Lorentz model

use  $R^{n+1}$  to represent  $L^n$



**MERU uses Lorentz model  
use  $R^{n+1}$  to represent  $L^n$**

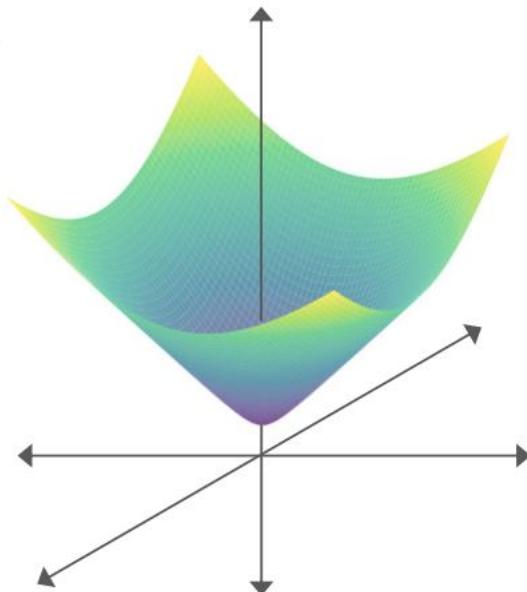
Every vector  $\mathbf{x} \in \mathbb{R}^{n+1}$  can be written as  $[\mathbf{x}_{space}, x_{time}]$ ,  
where  $\mathbf{x}_{space} \in \mathbb{R}^n$  and  $x_{time} \in \mathbb{R}$ .



Inner Product:  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}} = \langle \mathbf{x}_{space}, \mathbf{y}_{space} \rangle - x_{time} y_{time}$

MERU uses Lorentz model  
use  $\mathbb{R}^{n+1}$  to represent  $L^n$

Every vector  $\mathbf{x} \in \mathbb{R}^{n+1}$  can be written as  $[\mathbf{x}_{space}, x_{time}]$ ,  
where  $\mathbf{x}_{space} \in \mathbb{R}^n$  and  $x_{time} \in \mathbb{R}$ .



Inner Product:  $\langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}} = \langle \mathbf{x}_{space}, \mathbf{y}_{space} \rangle - x_{time} y_{time}$

Norm:  $\|\mathbf{x}\|_{\mathcal{L}} = \sqrt{|\langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}}|}$

Lorentz model with constant curvature  $-c$ :

$$\mathcal{L}^n = \{\mathbf{x} \in \mathbb{R}^{n+1} : \langle \mathbf{x}, \mathbf{x} \rangle_{\mathcal{L}} = -1/c, c > 0\}$$

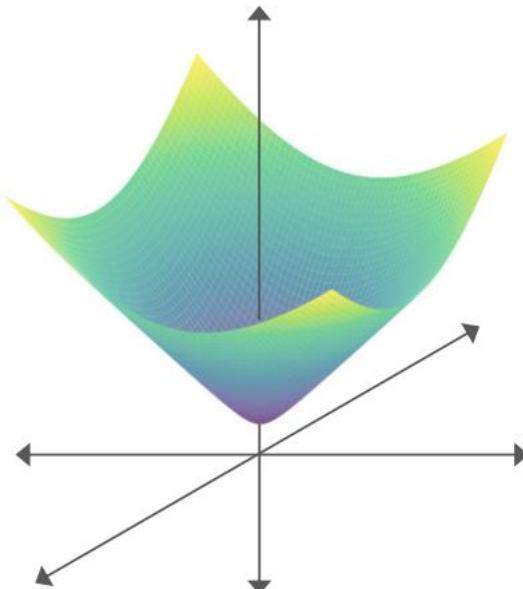
Satisfy:

$$x_{time} = \sqrt{1/c + \|\mathbf{x}_{space}\|^2}$$

# Background

---

MERU uses Lorentz model  
use  $\mathbb{R}^{n+1}$  to represent  $L^n$



Geodesics:  $d_{\mathcal{L}}(\mathbf{x}, \mathbf{y}) = \sqrt{1/c} \cdot \cosh^{-1}(-c \langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}})$

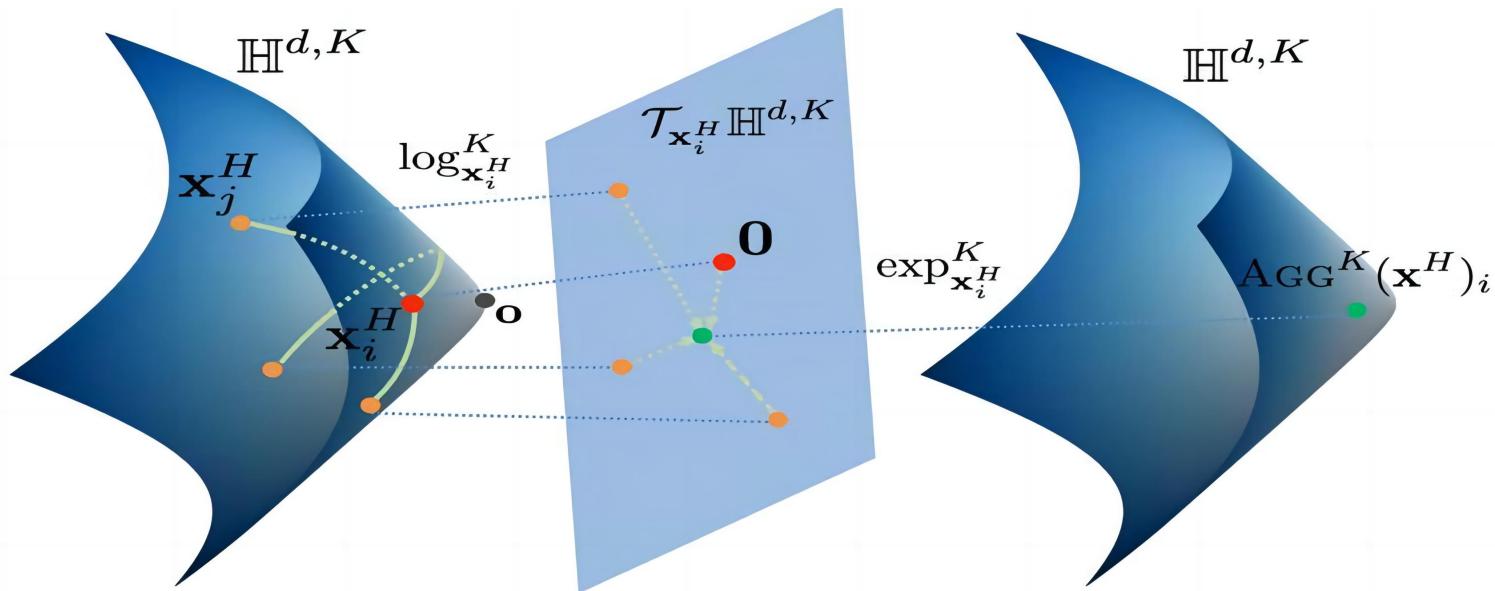
Tangent space for  $\mathbf{z}$ :  $\mathcal{T}_{\mathbf{z}}\mathcal{L}^n = \{\mathbf{v} \in \mathbb{R}^{n+1} : \langle \mathbf{z}, \mathbf{v} \rangle_{\mathcal{L}} = 0\}$

Project to Tangent space:

$$\mathbf{v} = \text{proj}_{\mathbf{z}}(\mathbf{u}) = \mathbf{u} + c \mathbf{z} \langle \mathbf{z}, \mathbf{u} \rangle_{\mathcal{L}}$$

# Background

---

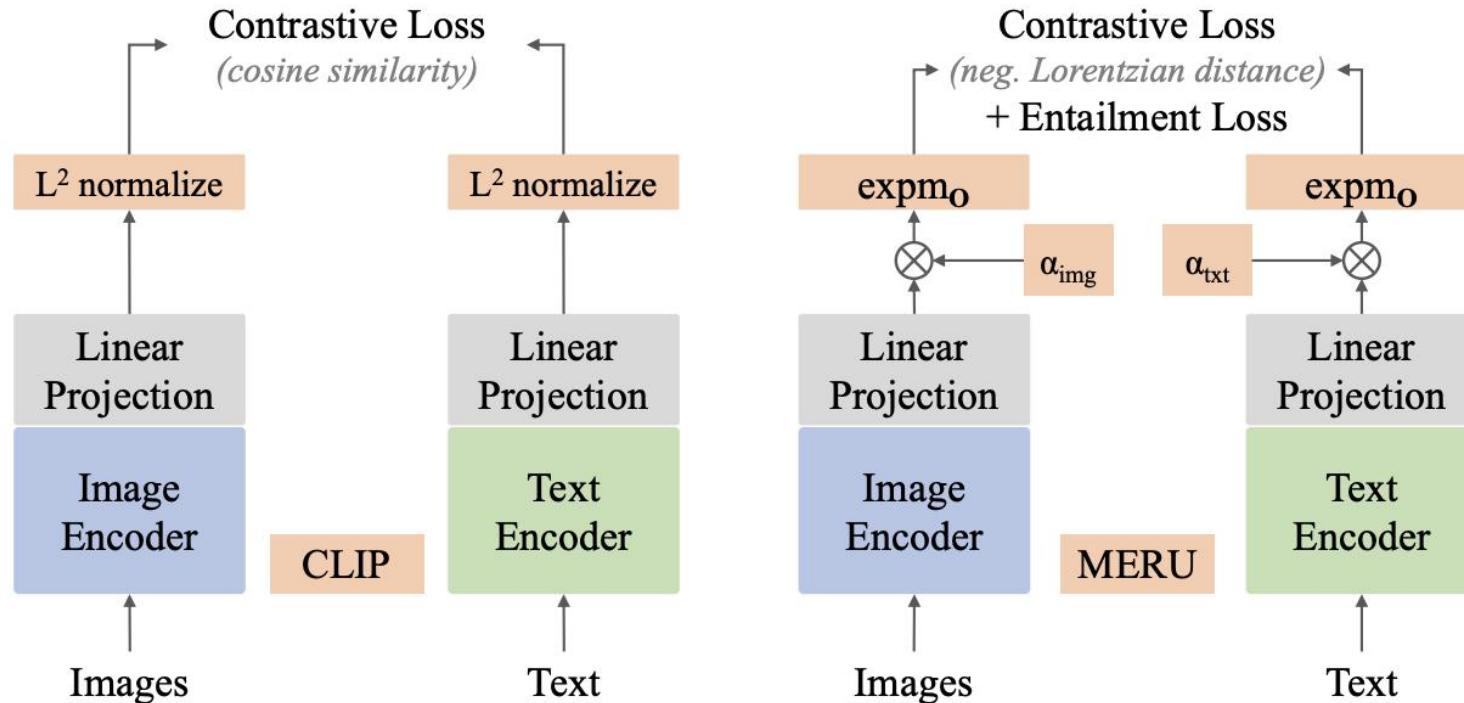


Exponential map:  $\mathbf{x} = \text{expm}_{\mathbf{z}}(\mathbf{v}) = \cosh(\sqrt{c} \|\mathbf{v}\|_{\mathcal{L}}) \mathbf{z} + \frac{\sinh(\sqrt{c} \|\mathbf{v}\|_{\mathcal{L}})}{\sqrt{c} \|\mathbf{v}\|_{\mathcal{L}}} \mathbf{v}$

Logarithmic map:  $\mathbf{v} = \text{logm}_{\mathbf{z}}(\mathbf{x}) = \frac{\cosh^{-1}(-c \langle \mathbf{z}, \mathbf{x} \rangle_{\mathcal{L}})}{\sqrt{(c \langle \mathbf{z}, \mathbf{x} \rangle_{\mathcal{L}})^2 - 1}} \text{proj}_{\mathbf{z}}(\mathbf{x})$

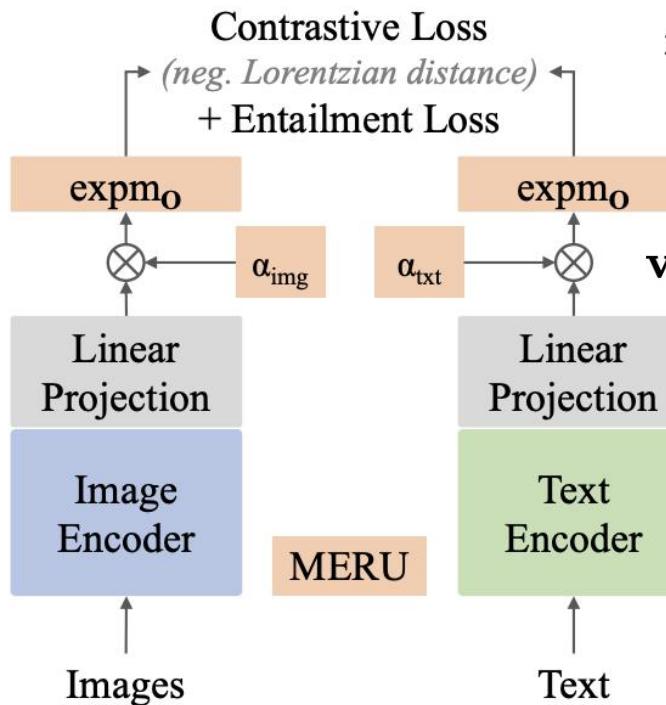
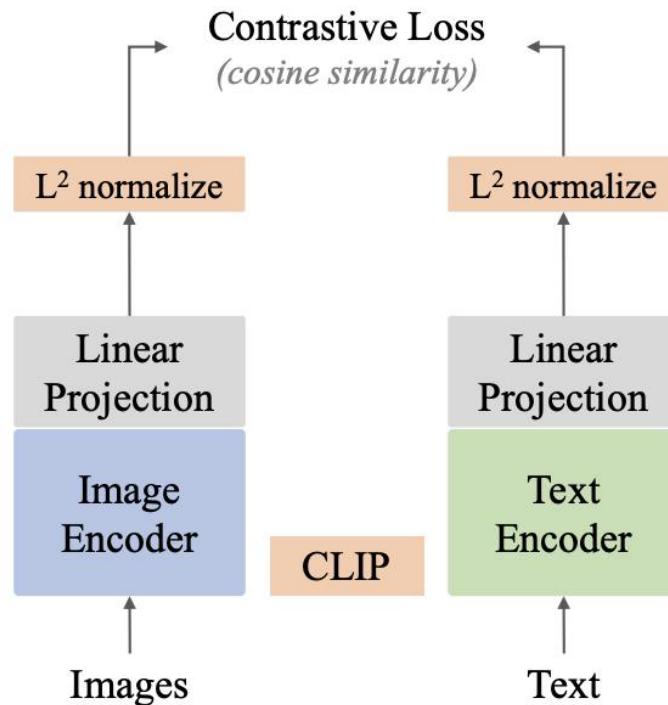
# Background

- MERU Model Design



# Background

- MERU Model Design



$$\mathbf{x}_{space} = \frac{\sinh(\sqrt{c} \|\mathbf{v}_{space}\|)}{\sqrt{c} \|\mathbf{v}_{space}\|} \mathbf{v}_{space}$$

$$\mathbf{v}_{enc} \in \mathbb{R}^n \quad \mathbf{v} = [\mathbf{v}_{enc}, 0] \in \mathbb{R}^{n+1}$$

# Background

- MERU Model Design

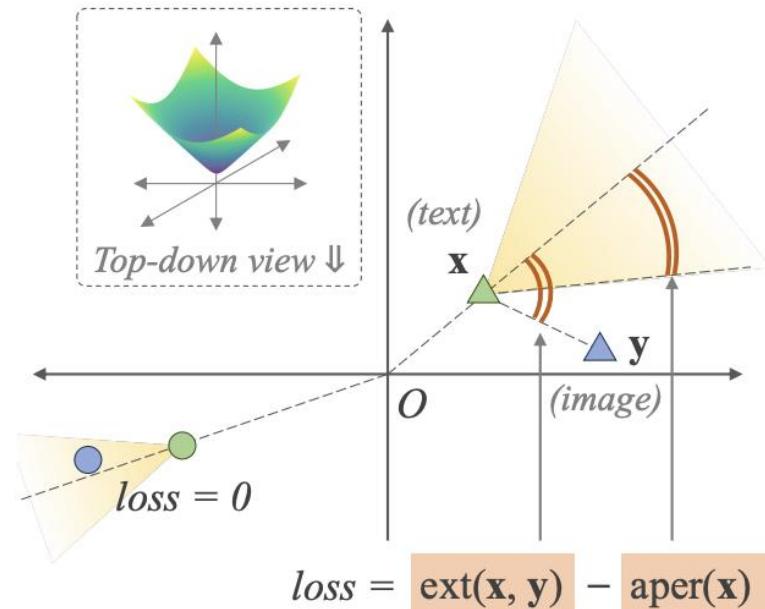
Contrastive Loss

$$L = -\frac{1}{N} \sum_{i=1}^N \log \frac{\exp(\cos(\mathbf{z}_i^I, \mathbf{z}_i^T)/\tau)}{\sum_{j=1}^N \exp(\cos(\mathbf{z}_i^I, \mathbf{z}_j^T)/\tau)}$$

Not cos similarity,  
but geodesics

Entailment loss

$$\mathcal{L}_{entail}(\mathbf{x}, \mathbf{y}) = \max(0, \text{ext}(\mathbf{x}, \mathbf{y}) - \text{aper}(\mathbf{x}))$$



$$\text{ext}(\mathbf{x}, \mathbf{y}) = \cos^{-1} \left( \frac{y_{time} + x_{time} c \langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}}}{\|\mathbf{x}_{space}\| \sqrt{(c \langle \mathbf{x}, \mathbf{y} \rangle_{\mathcal{L}})^2 - 1}} \right)$$

$$\text{aper}(\mathbf{x}) = \sin^{-1} \left( \frac{2K}{\sqrt{c} \|\mathbf{x}_{space}\|} \right)$$

Le, M., Roller, S., Papaxanthos, L., Kiela, D., and Nickel, M. Inferring Concept Hierarchies from Text Corpora via Hyperbolic Embeddings. In Proceedings of the Annual Meeting on Association for Computational Linguistics (ACL), 2019. 3, 4, 9



北京大学  
PEKING UNIVERSITY

## PART 03

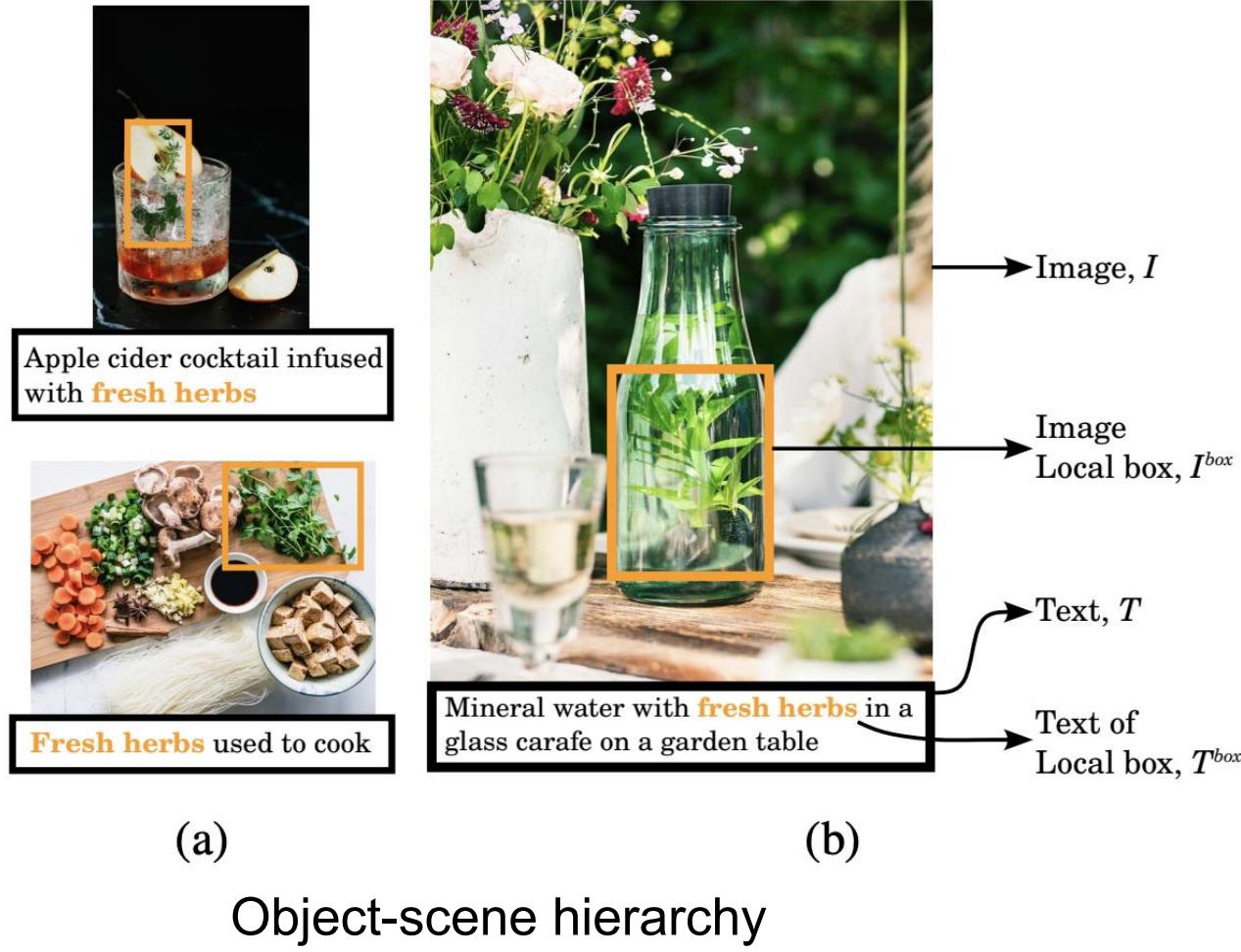
---

### Method



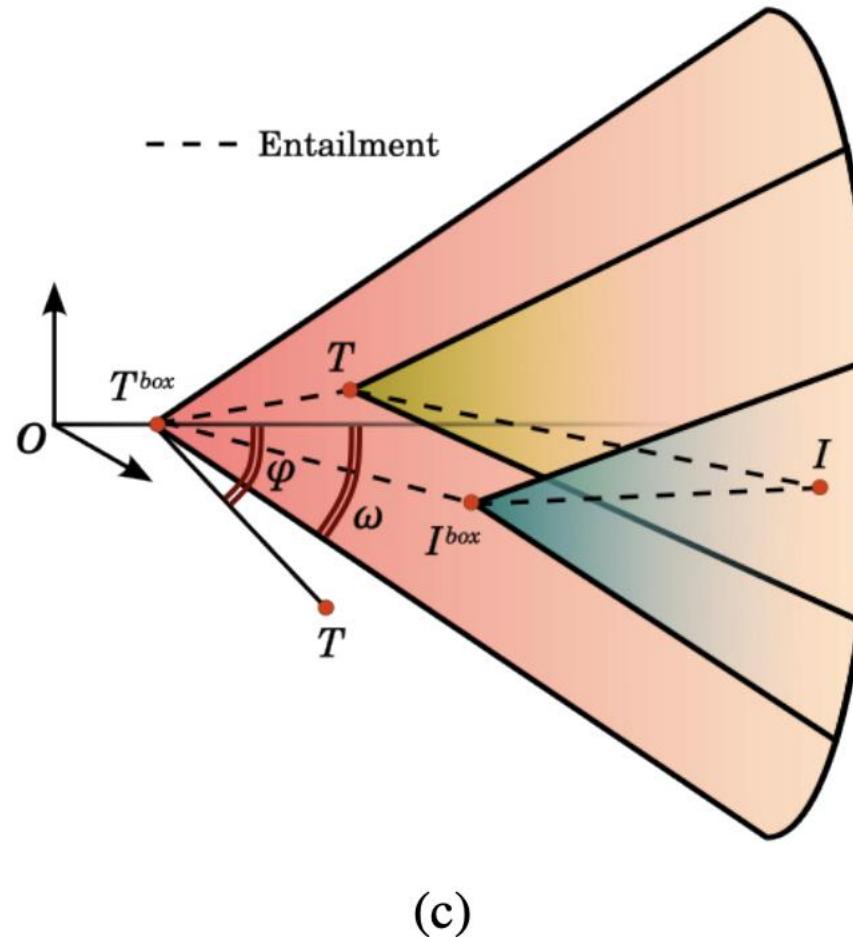
# Method

---



# Method

---



# Method

---

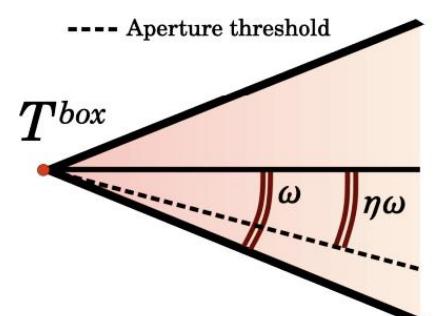
Contrastive:  $L_{cont}^*(I, T) = - \sum_{i \in B} \log \frac{\exp(d_{\mathbb{L}}(g_I(I_i), g_T(T_i))/\tau)}{\sum_{k=1, k \neq i}^B \exp(d_{\mathbb{L}}(g_I(I_i), g_T(T_k))/\tau)}$

$$hCC(I, T, I^{box}, T^{box}) = \frac{1}{4} \left( \underbrace{L_{cont}^*(I, T) + L_{cont}^*(T, I)}_{\text{specific-info contrast}} + \underbrace{L_{cont}^*(I^{box}, T) + L_{cont}^*(T^{box}, I)}_{\text{general-info contrast}} \right)$$

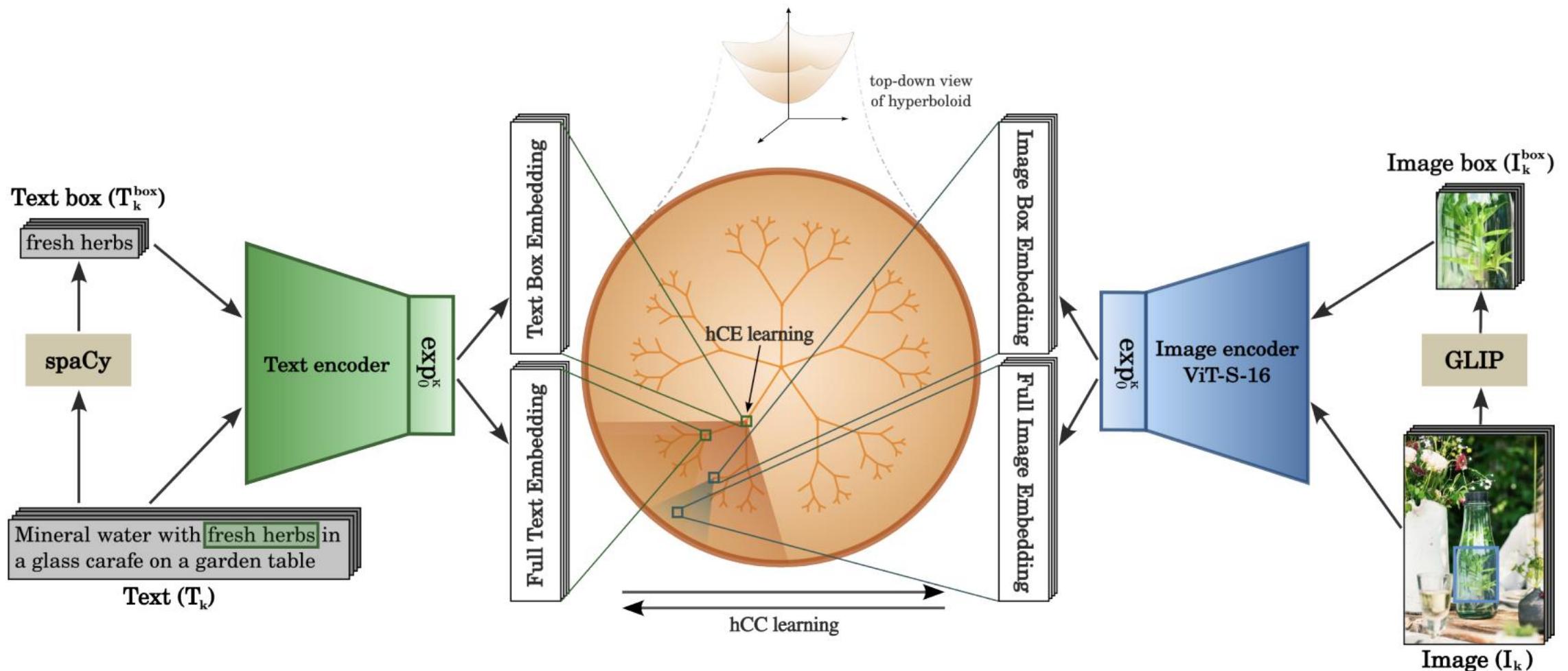
Entailment:  $L_{ent}^*(\mathbf{p}, \mathbf{q}) = \max(0, \phi(\mathbf{p}, \mathbf{q}) - \eta\omega(\mathbf{q}))$

$$hCE(I, T, I^{box}, T^{box}) = \underbrace{L_{ent}^*(I^{box}, T^{box}) + L_{ent}^*(I, T)}_{\text{inter-modality entailment}} + \underbrace{L_{ent}^*(I, I^{box}) + L_{ent}^*(T, T^{box})}_{\text{intra-modality entailment}}$$

Total Loss:  $hC = hCC + \gamma hCE$



# Method





北京大学  
PEKING UNIVERSITY

## PART 04

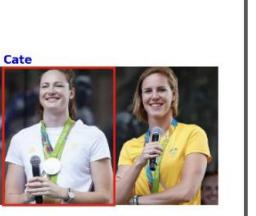
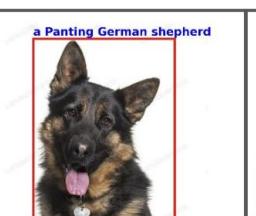
---

# Experiments



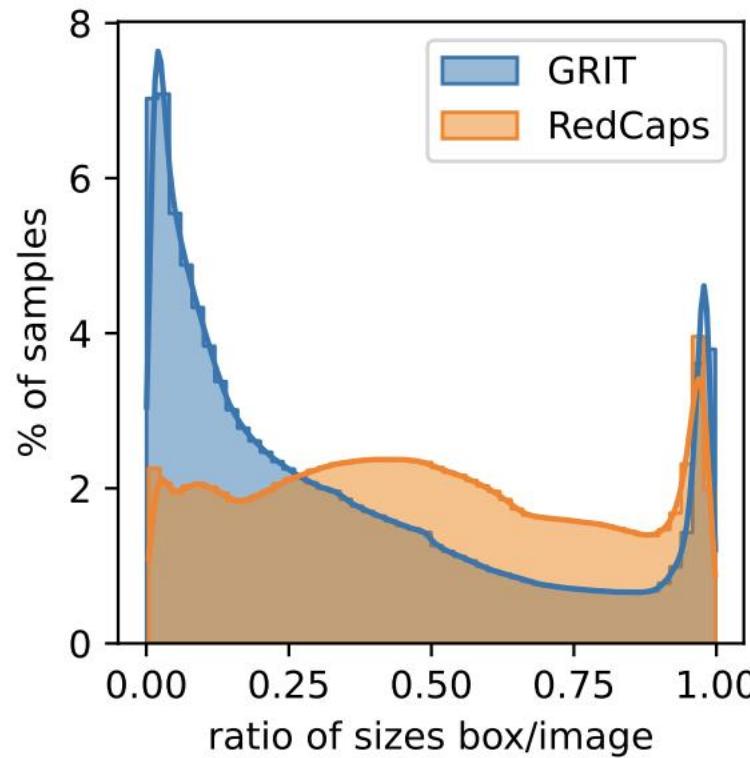
# Experiments

- Datasets
- Grounded vision-language pairs
  - Grounded Image-Text Pairs (GRIT) dataset: 20.5 million pairs
  - RedCaps
  - Conceptual Captions 3M (CC3M)

				
People sitting by <b>the pool</b> for a Virginia wedding reception	Tela Golden retriever against <b>decorated Chrostmas tree</b> in the room	Cucumber in <b>a glass</b> of water	For <b>Cate</b> (left), her experience at Rio left her "heartbroken."	<b>Young doctor</b> with <b>patient</b> preparing an <b>x-ray</b>
				
banquette : <b>Waitress</b> put dishes at <b>the table</b>	KHON2 pets and <b>their people</b>	<b>Mating sandhill cranes</b> dance in the air	Close-up of a <b>Panting German shepherd</b> sitting, isolated	<b>Flowers</b> on a field against an isolated background

# Experiments

---



# Experiments

		w/ boxes	samples (M)	General datasets						Fine-grained datasets						Misc. datasets			
				ImageNet	CIFAR-10	CIFAR-100	SUN397	Caltech-101	STL-10	Food-101	CUB	Cars	Aircraft	Pets	Flowers	DTD	EuroSAT	RESISC45	Country211
<b>ViT S/16</b>	<b>RedCaps</b>																		
	CLIP <sup>†</sup>	✗	11.4	32.5	66.7	35.8	26.7	60.8	89.8	72.5	29.8	11.1	1.3	72.5	44.9	16.4	30.1	27.7	5.0
	CLIP	✓	11.4 [6.3]	30.2	76.5	42.4	25.8	62.3	89.5	69.6	25.7	8.5	2.2	65.3	38.6	13.6	36.6	28.5	4.6
	MERU <sup>†</sup>	✗	11.4	31.4	65.9	35.2	26.8	58.1	89.3	71.4	29.0	8.3	1.6	71.0	40.9	17.0	29.9	29.3	4.7
	MERU	✓	11.4 [6.3]	29.9	76.4	39.9	26.6	62.3	89.5	68.4	25.4	8.9	1.2	67.2	37.6	13.0	30.5	27.6	4.3
<b>ViT S/16</b>	<b>HyCoCLIP</b>	✓	5.8 [6.3]	31.9	77.4	37.7	27.6	64.5	90.9	71.1	28.8	9.7	1.1	70.5	41.4	13.4	22.7	30.7	4.4
	<b>GRIT</b>																		
	CLIP	✗	20.5	36.7	70.2	42.6	49.5	73.6	89.7	44.7	9.8	6.9	2.0	44.6	14.8	22.3	40.7	40.1	5.1
	CLIP	✓	20.5 [35.9]	36.2	84.2	54.8	46.1	74.1	91.6	43.2	11.9	6.0	2.5	45.9	18.1	24.0	32.4	35.5	4.7
	MERU	✗	20.5	35.4	71.2	42.0	48.6	73.0	89.8	48.8	10.9	6.5	2.3	42.7	17.3	18.6	39.1	38.9	5.3
<b>ViT B/16</b>	MERU	✓	20.5 [35.9]	35.0	85.0	54.0	44.6	73.9	91.6	41.1	10.1	5.6	2.2	43.9	15.9	24.5	39.3	33.5	4.8
	HyCoCLIP	✓	20.5 [35.9]	41.7	85.0	53.6	52.5	75.7	92.5	50.2	14.7	8.1	4.2	52.0	20.5	22.3	33.8	45.7	5.2
	CLIP	✗	20.5	40.6	78.9	48.3	53.0	76.7	92.4	48.6	10.0	9.0	3.4	45.9	21.3	23.4	37.1	42.7	5.7
<b>ViT B/16</b>	MERU	✗	20.5	40.1	78.6	49.3	53.0	72.8	93.2	51.5	11.9	8.6	3.7	48.5	21.2	22.2	31.7	44.2	5.6
	HyCoCLIP	✓	20.5 [35.9]	45.8	88.8	60.1	57.2	81.3	95.0	59.2	16.4	11.6	3.7	56.8	23.9	29.4	35.8	45.6	6.5

Zero-shot classification

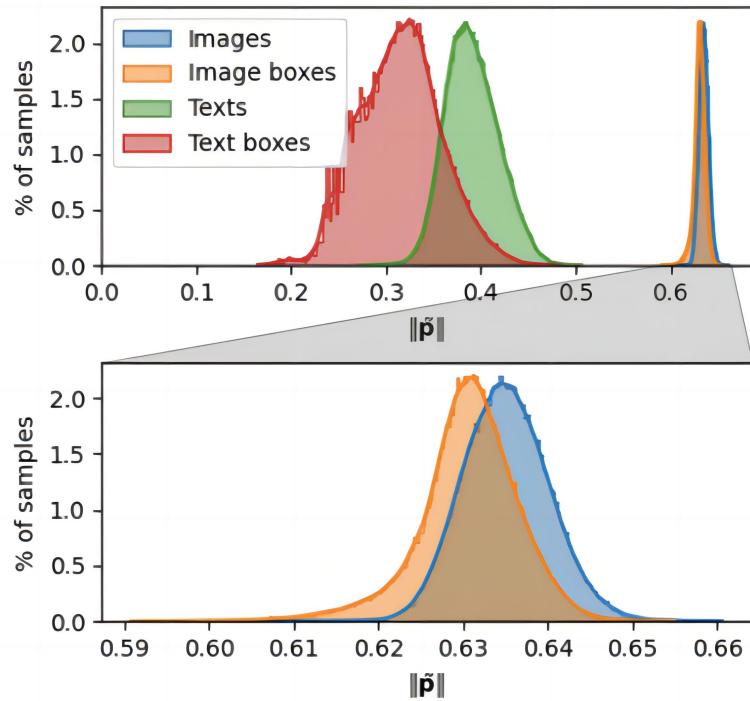
# Experiments

Vision encoder	Model	w/ boxes	Text retrieval				Image retrieval				Hierarchical metrics				
			COCO		Flickr		COCO		Flickr		WordNet				
			R@5	R@10	R@5	R@10	R@5	R@10	R@5	R@10	TIE( $\downarrow$ )	LCA( $\downarrow$ )	J( $\uparrow$ )	$P_H(\uparrow)$	$R_H(\uparrow)$
ViT S/16	CLIP	$\times$	69.3	79.1	<b>90.2</b>	<b>95.2</b>	53.7	65.2	81.1	87.9	4.02	2.39	0.76	0.83	0.84
	CLIP	$\checkmark$	60.7	71.8	84.2	91.3	47.1	58.6	73.1	82.1	4.03	2.38	0.76	0.83	0.83
	MERU	$\times$	68.8	78.8	89.4	94.8	53.6	65.3	80.4	87.5	4.08	2.39	0.76	0.83	0.83
	MERU	$\checkmark$	<b>72.7</b>	<b>81.9</b>	83.5	90.1	46.6	58.3	60.0	71.7	4.08	2.39	0.75	0.83	0.83
	HyCoCLIP	$\checkmark$	69.5	79.5	89.1	93.9	<b>55.2</b>	<b>66.6</b>	<b>81.5</b>	<b>88.1</b>	<b>3.55</b>	<b>2.17</b>	<b>0.79</b>	<b>0.86</b>	<b>0.85</b>
ViT B/16	CLIP	$\times$	71.4	81.5	<b>93.6</b>	<b>96.9</b>	57.4	68.5	83.5	89.9	3.60	2.21	0.79	0.85	0.85
	MERU	$\times$	<b>72.3</b>	82.0	93.5	96.2	57.4	68.6	84.0	90.0	3.63	2.22	0.78	0.85	0.85
	HyCoCLIP	$\checkmark$	72.0	<b>82.0</b>	92.6	95.4	<b>58.4</b>	<b>69.3</b>	<b>84.9</b>	<b>90.3</b>	<b>3.17</b>	<b>2.05</b>	<b>0.81</b>	<b>0.87</b>	<b>0.87</b>

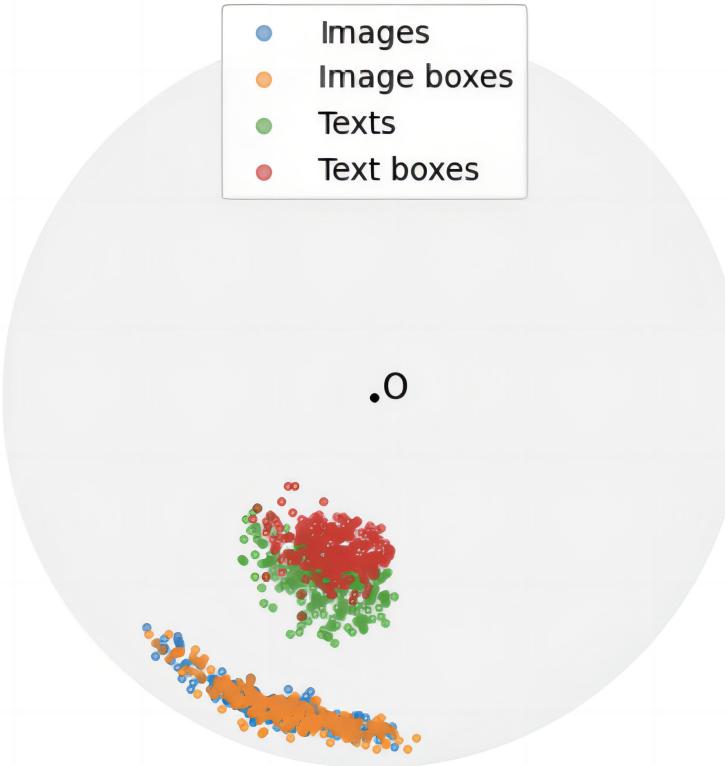
Zero-shot retrieval, detection, and hierarchical classification

Model	AP
CLIP	51.2
MERU	55.8
RegionCLIP	65.2
HyCoCLIP	<b>68.5</b>

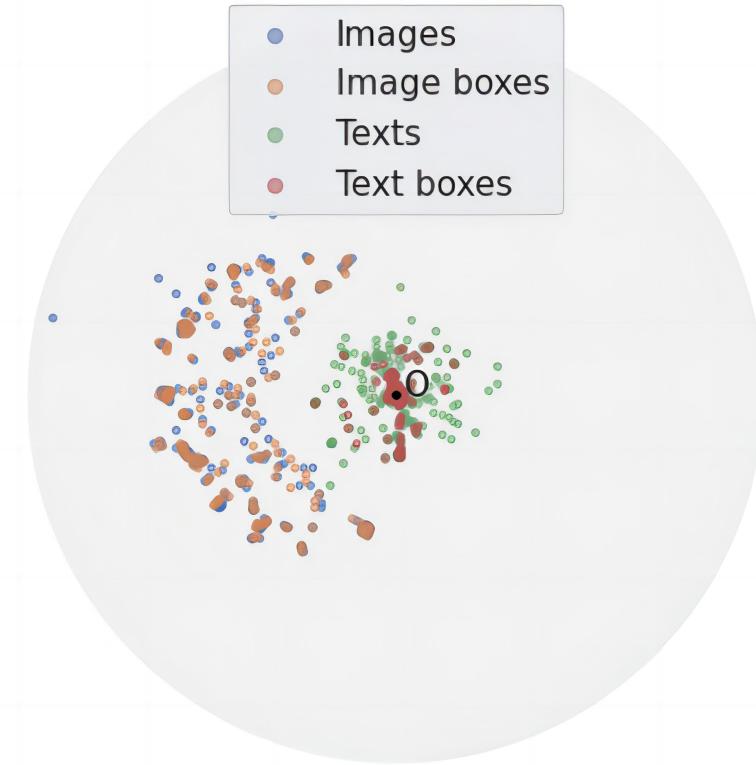
# Experiments



(a)



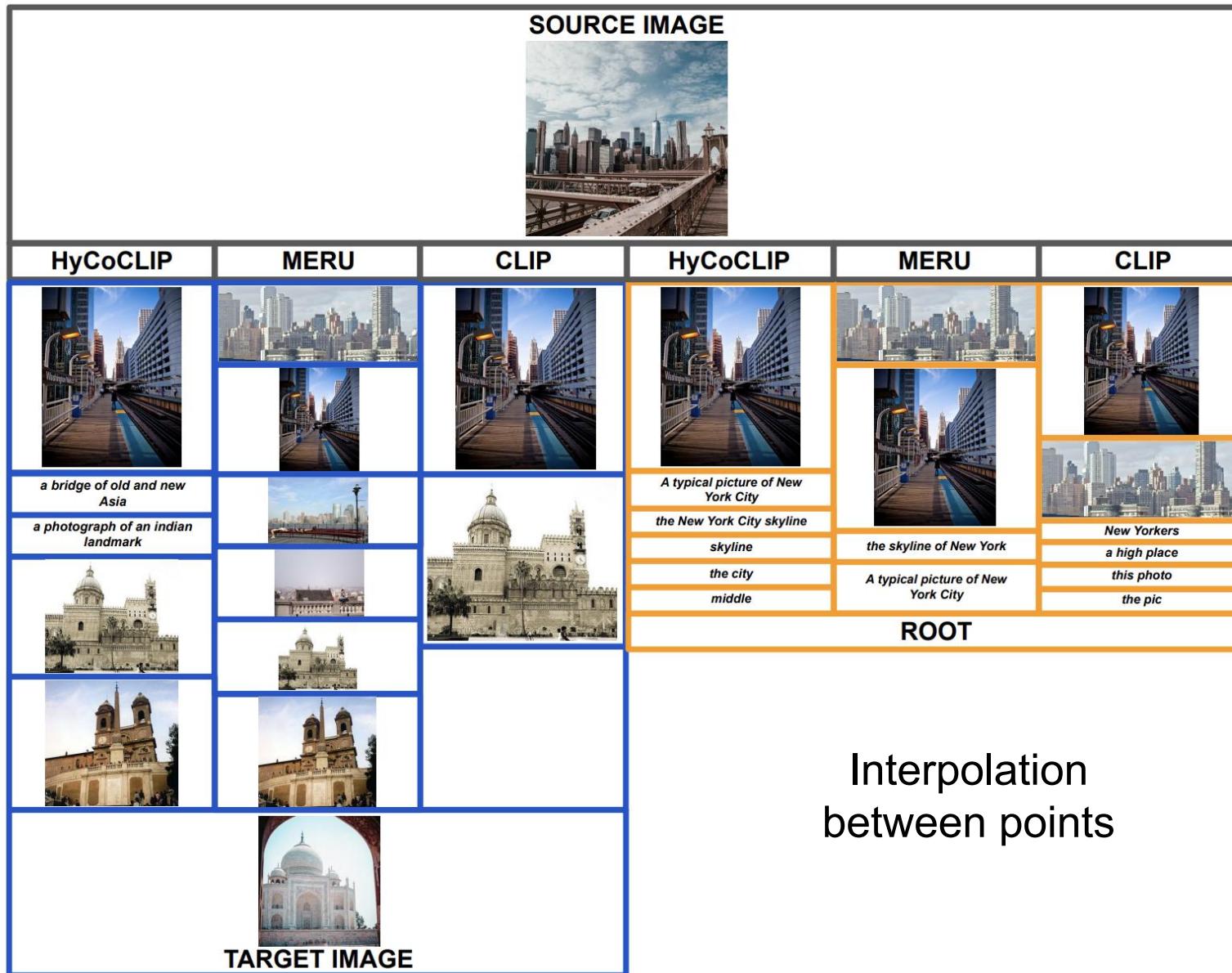
(b)



(c)

Visualizing the learned hyperbolic space of HyCoCLIP in lower dimensions

# Experiments





Thanks for  
Listening!