

Vision-Language Understanding in Hyperbolic Space

Sarthak Srivastava, Kathy Wu
sarthasr@amazon.com, rhaow@amazon.com

ABSTRACT

State-of-the-art performance has been achieved in recent years on tasks such as search, recommendation and classification using Visuo-Lingual Multi-Modal models. While the pretrained Vision-Language models like Contrastive Language-Image Pre-training (CLIP) have achieved promising zero-shot performance on several generalized tasks by learning vision-language concepts in a common space, the natural hierarchical relationship between them remains unexplored. In this work we propose PoinCLIP: a hyperbolic Poincaré geometry based vision-language model that learns joint text-image representation considering the hierarchical relation between the two. We compare the performance of PoinCLIP with CLIP model for zero-shot image classification and retrieval tasks to demonstrate the efficacy of the proposed method.

CCS CONCEPTS

• **Computing methodologies** → **Semi-supervised learning settings.**

KEYWORDS

Vision Language Model, Contrastive Learning, Poincare Ball, Hyperbolic Geometry

ACM Reference Format:

Sarthak Srivastava, Kathy Wu. 2024. Vision-Language Understanding in Hyperbolic Space. In *Proceedings of Multimodal Representation and Retrieval Workshop (SIGIR '24)*. ACM, Washington D.C., USA, 7 pages.

1 INTRODUCTION

Vision Language Models Large vision-language models like CLIP [30] and ALIGN [14] learn visual concepts from their natural language description via multi-modal contrastive learning. In contrastive learning [16], an anchor item representation is compared with a similar and a dissimilar item with the aim of bringing similar item representation together and pushing different ones away. The effectiveness [34] of these models results from their pretraining over a diverse large-scale image-text dataset sources from the web, allowing them to learn diverse concept from real world resulting in their impressive generalizability over a variety of tasks in zero-shot setting like classification and retrieval. These models assume the geometry of the higher dimensional representation space as affine Euclidean [12, 25], making it harder to capture the visual-text hierarchical concepts. The entity containing more *general* concepts

should be located close to the root of the hierarchy tree than the entity encapsulating a more *specific* and complex information. Hyperbolic spaces [4, 31] are natural candidate for capturing this hierarchical information about data points as their volume grows exponentially away from the origin, against polynomial growth in case of Euclidean space. Hyperbolic space can be thought of as a continuous version of a tree with its root at the origin.

Vision Language Hierarchy The saying "A picture is worth a thousand words" conveys the information difference between an image and words describing them. For example, in Figure 1, the picture can be broken down into individual concepts consisting of "kitty" and "doggo", which might be transformed in different manner to generate caption, for e.g. 'my dog's innocence brings smile to my face', 'a dog and a cat having fun in field', etc. Following equivalence, a many words can be put together encapsulating complex concept to build an informative image. Injecting these inductive biases in the training of multi-modal models [30, 32] will allow them to learn a more generalizable and interpretable representation.

Hyperbolic Space Representation with PoinCLIP In this work we project the image-text concepts onto a Poincaré ball model of hyperbolic space while following the state of the art contrastive methodology, to help capture the hierarchical information about the image-text pair, in addition to their semantic similarity. The contribution of this work can be described as:

- We introduce PoinCLIP, a Poincaré ball based hyperbolic representation model trained using ViTs and Transformer encoder based contrastive loss using RedCap dataset containing 12M image-text pairs.
- We introduce an embedding entropy based entailment loss to enforce the hierarchy between image-text in the Poincaré space.

We compare the performance of the proposed method with strong baseline CLIP and MERU to demonstrate its competitiveness.

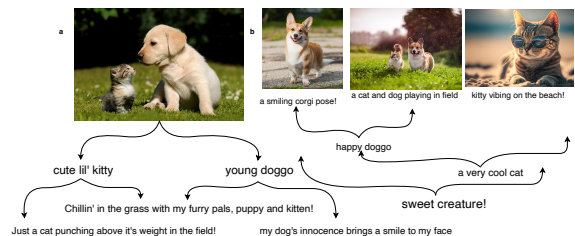


Figure 1: A picture is worth a thousand words. Left: Given an informative image it is possible to generate several textual concepts leveraging the visuo-lingual hierarchy. Right: Likewise, beginning from a simple text concept, it is possible to come up with complex visuo-lingual concepts by leveraging their hierarchical relation.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR '24, July 18, 2024, Washington D.C., USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.

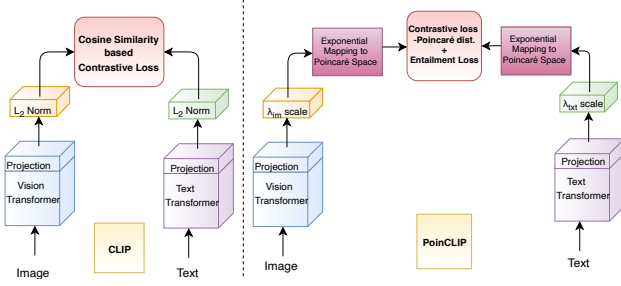


Figure 2: Overall Model Architecture. Left: Describes the baseline CLIP architecture based on which we have defined PoinCLIP. Image and text are encoded by Vision and Text Transformers respectively before being normalized and compared for contrastive loss calculation. Right: Describes PoinCLIP architecture. It differs from CLIP in aspect that encoder output is scaled and projected onto Poincaré space before computing contrastive loss and entailment loss for optimization.

2 RELATED WORK

The idea of hyperbolic space to better represent multi-modal entities is very recent and there are few related work in this field. MERU [9] attempts to capture the image-text semantics using Lorentz hyperboloid space. However, Lorentz manifold has less representation capacity compared to Poincaré ball as described in [28]. [10] discusses application of hyperbolic space based approach to learn hierarchical information just between the different image samples.

3 PRELIMINARIES

Hyperbolic geometry [2], also known as Lobachevskian geometry [29] is a non-Euclidean geometry where the Euclid’s fifth postulate of parallels don’t hold true and the space has a constant negative curvature. Hyperbolic spaces can be thought of as a continuous versions of tree data structure where the number of nodes until level h grow exponentially with the value l as $((b+1)b^h - 2)/(b - 1)$ where b is the branching factor. This tree grows from origin where h is 0 and it grows in terms of nodes exponentially away from origin. Such a structural arrangement is not possible in \mathbb{R}^2 Euclidean space [15] as the area and circumference of the hypercircle only grows quadratically and linearly respectively against an exponential growth in case of hyperbolic space. A manifold [3] is a topological space that locally resembles Euclidean space. Riemannian manifold [20] refers to a real and smooth manifold with Riemannian tensor, a metric tensor defined by a family a inner products as follow: Suppose p is a point on the curve of manifold M with $p \in M$ and denote the tangent space by $T_p(M) \in \mathbb{R}^n$, for any two tangent vectors $X(p)$ and $Y(p)$,

$$q : T_p M \times T_p M \rightarrow \mathbf{R}$$

defines a smooth function for the point $p \in M$

3.1 Curvature

In simple terms, curvature of a curve is its measure of deviation from a straight line and that of a surface is the measure of its

deviation from a plane. In terms of space, a curved space refers to spatial geometry which shows some finite curvature w.r.t a plane surface.

The curvature of Riemannian manifolds can be measured by Riemann curvature tensor, which assigns a tensor for each point at Riemannian manifold. The sectional curvature at a point P can be written as:

$$S(u, v) = \frac{\langle R(u, v)v, u \rangle}{\langle u, u \rangle \langle v, v \rangle - \langle u, v \rangle^2}$$

where $\langle \cdot, \cdot \rangle$ is the inner product, R is Riemannian curvature tensor, the u and v are vectors at tangent space of P that are linearly independent. There are special cases for Riemannian manifolds which have constant curvature at each point. For instance, Euclidean geometry has constant zero curvature, and Hyperbolic geometry has constant negative curvature.

Hyperbolic Space Hyperbolic n -space [18], denoted \mathbb{H}^n , is the unique simply connected, n -dimensional Riemannian manifold which has constantly negative sectional curvature.

3.2 Poincaré model of Hyperbolic Geometry

A Poincaré disk is a hyperbolic geometric model in which we represent a line as an arc of a circle whose ends are perpendicular to the disk’s diameter. It’s a useful model that uses hyperbolic geometry to discover continuous hierarchical relations among data pairs by embedding them into n dimensional Poincaré hypersphere. Mathematically, we can define an n -dimensional Poincaré ball in constant negative curvature value of $K = -1$ as:

$$\mathbb{P}_{K=-1}^n = \{x \in \mathbb{R}^n : \|x\|^2 < 1\} \quad (1)$$

Where $\|\cdot\|$ represents the Euclidean norm of a data point. The metric tensor for a Poincaré ball is represented as $g_x^{\mathbb{P}_{K=-1}^n} = (Y_x^{K=-1})^2 g_x^{\mathbb{E}}$ where $Y_x^{K=-1} = \frac{1}{1-\|x\|^2}$ is the conformity factor and $g_x^{\mathbb{E}}$ is the metric tensor for euclidean space represented as $g_x^{\mathbb{E}} = \text{diag}([1, 1, \dots, 1])$. The distance $d_h(p_1, p_2)$ between two samples p_1 and p_2 in the Poincaré space $\mathbb{P}_{K=-k}^n$ is calculated as:

$$d_h(p_1, p_2) = \frac{2}{\sqrt{k}} \tanh^{-1}(\sqrt{k} \|(-p_1) \oplus_k p_2\|_2) \quad (2)$$

Where $\|\cdot\|$ represents the Euclidean norm of a data point and \oplus_k is Möbius addition described shortly. We map Euclidean feature into hyperbolic Poincaré ball manifold via $h_i = \exp_0^{K=-1}(x_i^{Euc})$ where h_i represents the transformed x_i value in the hyperbolic space. The exponential map value \exp_x^K for a vector p in a space having curvature value K is calculated as:

$$\exp_x^K(p) = x \oplus_K \left(\tanh \left(\frac{\sqrt{-K} Y_x^K \|p\|}{2} \right) \frac{p}{\sqrt{-K} \|p\|} \right) \quad (3)$$

To reverse map a vector p from Hyperbolic space of curvature value K to Euclidean space, we apply logarithmic mapping as following:

$$\log_x^K(p) = \frac{2}{\sqrt{-K} Y_x^K} \text{arctanh} \left(\sqrt{-K} \|v\| \right) \frac{v}{\|v\|} \quad (4)$$

Where v is calculated as $-x \oplus_K p$ and \oplus_K represents the Möbius addition defined as follow:

$$x \oplus_K y = \frac{(1 - 2K \langle x, y \rangle - K \|y\|^2) x + (1 + K \|x\|^2) y}{1 - 2K \langle x, y \rangle + K^2 \|x\|^2 \|y\|^2} \quad (5)$$

Where $\langle x, y \rangle$ represents the inner product between x and y in hyperbolic space.

4 METHODOLOGY

In this section we discuss the learning objective and modelling details of PoinCLIP to learn the hierarchy aware representations for input text and images. PoinCLIP is based on CLIP methodology consisting of a vision transformer based image encoder and a text transformer based text encoder using byte pair encoding. Both encoders generate image and text representations for input image and text respectively, which are then passed into a projection layer to obtain embeddings of a fixed size n . In addition to these steps, we 1.) Transfer the embeddings from euclidean space to hyperbolic Poincaré space, and 2.) augment the loss function to enforce the partial order hierarchical relation between image and text.

Transfer of Embeddings onto the Poincaré Space While training, the image and text samples are passed to ViT and Text Transformer encoders respectively followed by a projection layer as shown in fig. 2. This is followed by transformation of the embeddings (v_{img}, v_{txt}) from Euclidean geometry to Hyperbolic Poincaré geometry as (h_{img}, h_{txt}) following the eq. 3 w.r.t the origin.

Numerical Overflow Prevention Since euclidean space to hyperbolic space to calculate (h_{img}, h_{txt}) requires an exponential operation, the norm of embeddings changes from order of \sqrt{n} to $e^{\sqrt{n}}$, potentially causing numerical overflow. To fix this, embedding scaling is applied before exponential mapping via two learnable parameters λ_{img} and λ_{txt} initialized to $1/\sqrt{n}$ to prevent the norm of the embedding from numerical overflow in the Poincaré space.

Training Objectives Our training objective is to enforce semantic similarity as well as structural partial order relation between given image-text pairs to improve the generalization capability of vision-language models. To this end, we optimize for image-text contrastive loss and entailment loss.

4.1 Contrastive Loss

We have implemented same multi-class N-pair version of the contrastive loss as used in CLIP [30] with an important difference that we calculate the similarity via distances in Poincaré space from eq.3 instead of cosine similarity. For a given batch size N we use the negative Poincaré space distance to compute contrastive loss between 1 positive and $N - 1$ negative pair per image and per text. The average of image wise and text wise loss is used as overall contrastive loss \mathcal{L}_{cont} to enforce image-text semantic similarity.

4.2 Entailment Loss

We apply an additional entailment loss from [9] with modification to enforce partial order relationship between image-text pairs. In [9], the assumption is that text always entails the image within the entailment cone. In contrast, we adopt an entropy based strategy to determine correct entailment order between text and image per instance. In Physics, the space-time structure is knitted together by the causal associations represented by the causal graph, the analog of entailment cone. An entailment cone is essentially a structure representing the ‘‘time evolution’’ from a particular initial condition [36]. Keeping this view in perspective and given that image-text embeddings from respective transformers are learned in same latent

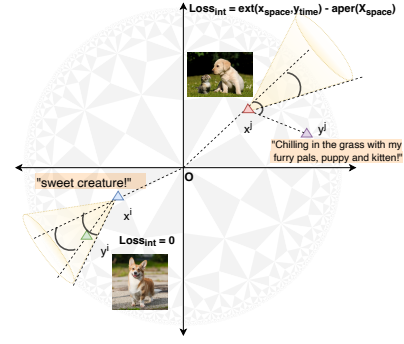


Figure 3: Entailment Cone (projection from Poincaré space on Euclidean Space). Loss pushes y_{time} embedding inside an entailment cone projected by embedding x and is defined as the difference between exterior angle $\angle Oxy$, and half aperture of the cone. Loss is zero if the y_{time} is already inside the cone. Indices i and j in superscripts represent two different instances of image-text pairs

space, we can determine the relative position in entailment cone comparing the entropy of embeddings with the assumption that entropy increases with evolution of time along the entailment cone. For a given image-text pair, the simpler concept with lower entropy should be entailing more complex concept with higher entropy with time. We calculate the information entropy [33] of embeddings as:

$$H(x_{emb}) = - \sum_{i=1}^n x_i \log_2 x_i \quad (6)$$

where H is the entropy of embedding x_{emb} and x_i represents the content of size n embedding for i^{th} dimension. We define $x = x_{img}$, the image embedding if $H(x_{img}) < H(x_{txt})$ else, $x = x_{txt}$. Similarly define $y = x_{txt}$, the image embedding if $H(x_{img}) < H(x_{txt})$ else, $x = x_{txt}$ Figure 3 gives an overview of the entailment loss as projected in euclidean space. Exterior angle $\angle Oxy$ is defined as:

$$ext(\angle Oxy) = \arccos \left(\frac{\langle x, y \rangle (1 + \|x\|^2) - \|x\|^2 ((1 + \|y\|^2))}{\|x\| \cdot \|x - y\| \sqrt{1 + \|x\|^2 \|y\|^2 - 2\langle x, y \rangle}} \right) \quad (7)$$

While the aperture of the entailment cone is defined as:

$$aper(x) = \arcsin \left(K \frac{1 - \|x\|^2}{\|x\|} \right) \quad (8)$$

We calculate the entailment loss as:

$$\mathcal{L}_{entail}(x, y) = \max(0, ext(\angle Oxy) - aper(x)) - \lambda_{reg} ext(\angle Oxy) \quad (9)$$

where λ_{reg} is the regularization coefficient. Hence, the overall loss to be optimized becomes $\mathcal{L} = \mathcal{L}_{cont} + \lambda \mathcal{L}_{entail}$ where λ is entailment regularization factor. The ablation study related to λ and λ_{reg} is discussed in the appendix.

5 EXPERIMENTS

To establish the competitiveness of Poincaré hyperbolic representations of PoinCLIP compared to Euclidean representations obtained from CLIP-style models, we compare zero shot classification and retrieval performances of PoinCLIP, MERU and CLIP. We train

		Food-101[1]	CIFAR-10[17]	CIFAR-100[17]	CUB[37]	SUN397[38]	Aircraft[23]	DTD[6]	Pets[27]	Caltech-101[11]	Flowers[26]	STL-10[7]	EuroSAT[13]	RESISC45[5]	Country211[30]	MNIST[19]	CLEVR[39]	PCAM[35]	SST2[30]
ViT-S/16	CLIP	74.5	60.1	24.4	33.8	27.5	1.4	15.0	73.7	63.9	47.0	88.2	18.6	31.4	5.2	10.0	19.4	50.2	50.1
	MERU	75.6	52.0	24.7	33.7	28.0	1.3	16.2	72.3	64.1	49.2	91.1	30.4	32.0	4.8	7.5	14.5	51.0	50.0
	PoinCLIP	75.1	53.6	27.7	35.1	27.6	1.6	17.6	71.9	62.1	47.9	90.9	30.8	32.1	5.1	10.4	14.8	53.8	50.8
ViT-B/16	CLIP	78.9	65.5	33.4	33.3	29.8	1.4	17.0	77.9	68.5	50.9	92.2	25.6	31.0	5.8	10.4	14.3	54.1	51.5
	MERU	78.8	67.7	32.7	34.8	30.9	1.7	17.2	79.3	68.5	52.1	92.5	30.2	34.5	5.6	13.0	13.5	49.8	49.9
	PoinCLIP	78.4	70.4	35.4	34.9	31.3	2.1	17.9	78.5	67.5	51.3	91.9	31.7	33.5	5.5	12.1	15.0	49.6	50.0
ViT-L/16	CLIP	80.3	72.0	36.4	36.3	32.0	1.1	16.5	78.8	68.3	48.6	93.7	26.7	35.4	6.1	14.8	13.6	51.2	51.1
	MERU	80.6	68.7	35.5	37.2	33.0	2.2	17.2	80.0	67.5	52.1	93.7	28.1	36.5	6.2	11.8	13.1	52.7	49.3
	PoinCLIP	80.6	74.3	38.8	37.5	33.3	2.6	18.5	80.1	66.0	51.3	93.8	27.9	37.2	6.5	12.0	13.4	55.7	50.0

Table 1: Comparison of Proposed Method PoinCLIP vs Baseline Methods on different datasets. The metrics in color represent best performance metric for particular dataset. We observe that PoinCLIP outperforms all methods in 13 out of 18 datasets.

		<i>text</i> \rightarrow <i>image</i>		<i>image</i> \rightarrow <i>text</i>	
		R5	R10	R5	R10
ViT-S/16	CLIP	29.9	40.1	37.5	48.1
	MERU	30.5	40.9	39.0	50.5
	PoinCLIP	30.5	40.2	40.4	50.7
ViT-B/16	CLIP	32.9	43.3	41.4	52.7
	MERU	33.2	44.0	41.8	52.9
	PoinCLIP	33.2	43.7	42.1	53.4
ViT-L/16	CLIP	31.7	42.2	40.6	51.3
	MERU	32.6	43.0	41.9	53.3
	PoinCLIP	32.6	42.7	43.2	53.8

Table 2: Zero Shot Image and Text Retrieval on COCO Dataset. Metric in color represent best performance for the task.

PoinCLIP on public RedCaps dataset [8] consisting of 12M image-text pairs for 120k iterations with 2048 batch size (20 epochs) on 8xV100 GPUs. **Model** We use different size versions of Vision Transformers (S/B/L) as vision encoder using patch size of 16, freezing the positional encoding layer of the model. Text encoder is same as that of CLIP with 12 layer 512 dimensional Transformer with 77 maximum length byte pair encoding. Poincaré ball of 512 dimensions and learnable curvature K is used for Poincaré space transformation post embedding scaling. **Optimizer** We use AdamW Optimizer [22] with weight decay of 0.2 and $(\beta_1, \beta_2) = (0.9, 0.98)$. Weight decay is disabled for all gains, biases, and learnable scalars. model is trained for 120K iterations with batch size 1024 (10 epochs). The maximum learning rate is 5×10^{-4} , which increases linearly for first 4K iterations, followed by cosine decay to 0 [21]. We evaluate the performance the PoinCLIP with CLIP and MERU on 18 datasets for zero shot classification and on COCO dataset for retrieval task.

5.1 Results

From table 1, we compare PoinCLIP’s performance for zero shot classification and observe it performing better than the Euclidean space CLIP for 14 out of 18 datasets and than Lorentz model based MERU for 15 out of 18 datasets and on 13 out of 18 datasets overall. Comparing Top N retrieval recall for COCO dataset in table 2 we see that PoinCLIP performs better than CLIP on 4 out of 4 tasks while it performs better than MERU on 3 out of 4 tasks. Overall, PoinCLIP performs better than all methods on 3 out of 4 tasks demonstrating the competitiveness of the proposed method. Ablation results for regularization terms λ_{reg} and λ will be shared in appendix.

6 DISCUSSION

We obtain better performance for PoinCLIP over Euclidean space CLIP owing to hyperbolic nature of Poincaré geometry which allows the capture of partial order relation between image and text, in addition to the semantic relation for learning representation. The incremental benefit over MERU can be attributed to 2 reasons: 1. Use of Poincaré space over Lorentz space: As per the work done in [24] Poincaré geometry has a relatively larger capacity than the Lorentz model for correctly representing points 2. Entailment loss based on entropy derived relative hierarchy between image and text at instance level (refer ablation study in appendix table 5).

7 CONCLUSION

In this work we discussed Poincaré geometry based large scale image-text model that learns image-text partial order hierarchical relation, in order to capturing their semantic similarity. The main contribution of this work can be summarised as: 1. Poincaré Hyperbolic Geometry based Image-Text model capturing image-text semantic information along with their hierarchical-relation. 2. Embedding entropy based method to decide the entailment order of image-text when enforcing partial order relation ship. We demonstrate the efficacy of the proposed method via experiments comparing accuracy for zero shot classification and recall for zero shot retrieval.

REFERENCES

- [1] Guillaumin M. Van Gool L. Bossard, L. 2014. Food-101 – Mining Discriminative Components with Random Forests. In: *Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds) Computer Vision – ECCV 2014. ECCV 2014. Lecture Notes in Computer Science, vol 8694. Springer, Cham. https://doi.org/10.1007/978-3-319-10599-4_29* (2014).
- [2] James W Cannon, William J Floyd, Richard Kenyon, Walter R Parry, et al. 1997. Hyperbolic geometry. *Flavors of geometry* 31, 59-115 (1997), 2.
- [3] Jack Carr. 2012. *Applications of centre manifold theory*. Vol. 35. Springer Science & Business Media.
- [4] Ines Chami, Albert Gu, Vaggos Chatziafratis, and Christopher Ré. 2020. From trees to continuous embeddings and back: Hyperbolic hierarchical clustering. *Advances in Neural Information Processing Systems* 33 (2020), 15065–15076.
- [5] Han J. Cheng, G. and X Lu. 2017. Remote Sensing Image Scene Classification: Benchmark and State of the Art. *Proc. IEEE* (2017).
- [6] Maji S. Kokkinos I. Mohamed S. Cimpoi, M. and A. Vedaldi. 2014. Describing Textures in the Wild. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014).
- [7] Ng A. Coates, A. and H Lee. 2011. An Analysis of Single Layer Networks in Unsupervised Feature Learning. In *Proceedings of the International Conference on Artificial Intelligence and Statistics (AISTATS)* (2011).
- [8] Karan Desai, Gaurav Kaul, Zubin Aysola, and Justin Johnson. 2021. Redcaps: Web-curated image-text data created by the people, for the people. *arXiv preprint arXiv:2111.11431* (2021).
- [9] Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin Johnson, and Shanmukha Ramakrishna Vedantam. 2023. Hyperbolic image-text representations. In *International Conference on Machine Learning*. PMLR, 7694–7731.
- [10] Aleksandr Ermolov, Leyla Mirvakhabova, Valentin Khruikov, Nicu Sebe, and Ivan Oseledets. 2022. Hyperbolic Vision Transformers: Combining Improvements in Metric Learning. *arXiv:2203.10833 [cs.CV]*
- [11] Fergus R. Fei-Fei, L. and Perona. 2004. Learning Generative Visual Models from Few Training Examples: An Incremental Bayesian Approach Tested on 101 Object Categories. *CVPR Workshop* (2004).
- [12] Nóra Frankl, Andrey Kupavskii, and Konrad J Swanepoel. 2020. Embedding graphs in Euclidean space. *Journal of Combinatorial Theory, Series A* 171 (2020), 105146.
- [13] Bischke B. Dengel A. R. Helber, P. and D Borth. 2019. EuroSAT: A Novel Dataset and Deep Learning Benchmark for Land Use and Land Cover Classification. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing* (2019).
- [14] Chao Jia, Yinfei Yang, Ye Xia, Yi-Ting Chen, Zarana Parekh, Hieu Pham, Quoc Le, Yun-Hsuan Sung, Zhen Li, and Tom Duerig. 2021. Scaling up visual and vision-language representation learning with noisy text supervision. In *International conference on machine learning*. PMLR, 4904–4916.
- [15] Roger A Johnson. 2013. *Advanced euclidean geometry*. Courier Corporation.
- [16] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2021. Supervised Contrastive Learning. *arXiv:2004.11362 [cs.LG]*
- [17] A Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. *URL https://www.cs.toronto.edu/kriz/learning-features-2009-TR.pdf*. (2009).
- [18] Serge Lang. 2013. *Introduction to complex hyperbolic spaces*. Springer Science & Business Media.
- [19] Cortes C. LeCun, Y. and C Burges. 2010. MNIST handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist> (2010).
- [20] John M Lee. 2018. *Introduction to Riemannian manifolds*. Vol. 2. Springer.
- [21] Ilya Loshchilov and Frank Hutter. 2017. SGDR: Stochastic Gradient Descent with Warm Restarts. *arXiv:1608.03983 [cs.LG]*
- [22] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. *arXiv:1711.05101 [cs.LG]*
- [23] Rahtu E. Kannala J. Blaschko M. B. Maji, S. and A. Vedaldi. 2013. Fine-Grained Visual Classification of Aircraft. *arXiv preprint arXiv:1306.5151* (2013).
- [24] Gal Mishne, Zhengchao Wan, Yusu Wang, and Sheng Yang. 2023. The Numerical Stability of Hyperbolic Representation Learning. *arXiv:2211.00181 [cs.LG]*
- [25] Kevin P Murphy. 2012. *Machine learning: a probabilistic perspective*. MIT press.
- [26] M.-E. Nilsback and A. Zisserman. 2008. Automated Flower Classification over a Large Number of Classes. *ICVGIP* (2008).
- [27] Vedaldi A. Zisserman A. Parkhi, O. and C. V. Jawahar. 2012. Cats and Dogs. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2012).
- [28] Wei Peng, Tuomas Varanka, Abdelrahman Mostafa, Henglin Shi, and Guoying Zhao. 2021. Hyperbolic deep neural networks: A survey. *IEEE Transactions on pattern analysis and machine intelligence* 44, 12 (2021), 10023–10044.
- [29] Andrey Popov and Andrei Iacob. 2014. *Lobachevsky geometry and modern non-linear problems*. Springer.
- [30] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.
- [31] Frederic Sala, Chris De Sa, Albert Gu, and Christopher Ré. 2018. Representation tradeoffs for hyperbolic embeddings. In *International conference on machine learning*. PMLR, 4460–4469.
- [32] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [33] C.E. Shannon. 1948. Claude Shannon introduced the concept of information entropy in his 1948 paper, "A Mathematical Theory of Communication. *A Mathematical Theory of Communication*. *Bell System Technical Journal*, 27, 379-423. <http://dx.doi.org/10.1002/j.1538-7305.1948.tb01338.x> (1948).
- [34] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [35] Linmans J. Winkens J. Cohen T. Veeling, B. S. and M. Welling. 2018. CNNs for Digital Pathology. *arXiv preprint arXiv:1806.03962* (2018).
- [36] Athanasios Vlontzos, Henrique Bergallo Rocha, Daniel Rueckert, and Bernhard Kainz. 2020. Causal future prediction in a minkowski space-time. *arXiv preprint arXiv:2008.09154* (2020).
- [37] Branson S. Welinder P. Perona P. Wah, C. and S. J. Belongie. [n. d.]. The Caltech-UCSD Birds-200-2011 Dataset. ([n. d.]).
- [38] Hays J. Ehinger K. A. Oliva A. Xiao, J. and A. Torralba. 2010. SUN database: Large-scale scene recognition from abbey to zoo. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2010).
- [39] Puigcerver J. Kolesnikov A. Ruyssen P. Riquelme C. Lucic M. Djolonga J. Pinto A. S. Neumann M. Dosovitskiy A. Beyer L. Bachem O. Tschannen M. Michalski M. Bousquet O. Gelly S. and Houlsby N. Zhai, X. 2019. A Large-scale Study of Representation Learning with the Visual Task Adaptation Benchmark. In *Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019).

A APPENDIX

In this section, we describe ablation studies and supplementary material related to our experiments. In table 3 and table 5, we observe the difference between the proposed Poincaré embedding with the entropy inferred text-image order entailment loss compared with Poincaré embedding without entropy inferred text-image order entailment where text always entails image in entailment loss. The experiment has been conducted for ViT S/16 model for 120000 iterations using same optimizer and learning rate as the proposed method. As can be observed, the addition of entailment leads to improvement in 14 out of 18 datasets in zero shot classification setting while improving performance in all 4 zero shot retrieval tasks for COCO dataset, proving the efficacy of the proposed entropy based image-text order entailment loss.

In our experiments, λ_{reg} component in the loss function acts as a regularization against exterior angle reduction indefinitely to accommodate the tailing component inside the entailment cone of the entailing component at the cost of generalizability. λ_{reg} component in the loss function acts as a regularization against overall entailment loss. In table 4 we run the ablation study for λ_{reg} and λ by training ViT-S/16 PoinCLIP model for 1 epoch (6k iterations) and compare the average zero shot retrieval accuracy for COCO dataset. We find that $\lambda_{reg} = 0.1$ and $\lambda = 0.1$ provides the best performance and was chosen as the value for our experiments. The entailment loss described in eq. 9 depends on λ_{reg} and λ for calculation of the overall entailment loss.

		<i>text</i> \rightarrow <i>image</i>		<i>image</i> \rightarrow <i>text</i>	
		R5	R10	R5	R10
ViT-S/16	Poincaré	30.1	40.2	39.0	50.2
	PoinCLIP	30.5	40.2	40.4	50.7

Table 3: Zero Shot Image and Text Retrieval on COCO Dataset. Metric in color represent best performance for the task. Row corresponding to Poincaré represents the case where no entropy derived entailment order is enforced in the entailment loss and we assume that text always entail the image as assumed in MERU. The row corresponding to PoinCLIP represent the case where embedding entropy derived image-text entailment order is applied in entailment loss.

		λ				
		0	0.01	0.1	0.5	1
λ_{reg}	0	20.2	17.5	18.6	16.5	18.7
	0.01	20.2	15.4	22.1	16.7	16.0
	0.1	20.2	21.1	22.3	18.9	19.0
	0.5	20.2	19.2	18.6	16.8	15.7
	1	20.2	18.6	18.1	15.4	19.5

Table 4: To select proper values of λ and λ_{reg} we run a grid search for different values and compare the average of average zero shot retrieval accuracy for different retrieval tasks for COCO dataset and zero shot classification accuracy for CIFAR 100 dataset by training ViT-S/16 model for 1 epoch (6K iterations). We find the best performance at $\lambda = 0.1$ and $\lambda_{reg} = 0.1$. The best performance metric in color

		Food-101[1]	CIFAR-10[17]	CIFAR-100[17]	CUB[37]	SUN397[38]	Aircraft[23]	DTD[6]	Pets[27]	Caltech-101[11]	Flowers[26]	STL-10[7]	EuroSAT[13]	RESISC45[5]	Country211[30]	MNIST[19]	CLEVR[39]	PCAM[35]	SST2[30]
ViT-S/16	Poincaré	74.9	55.3	27.5	34.1	28.3	1.5	16.4	72.9	60.0	48.4	90.7	28.3	30.6	4.9	8.3	14.4	48.9	50.2
	PoinCLIP	75.1	53.6	27.7	35.1	27.6	1.6	17.6	71.9	62.1	47.9	90.9	30.8	32.1	5.1	10.4	14.8	53.8	50.8

Table 5: Comparison of Proposed Method PoinCLIP implementing entropy inferred image-text entailment order, with PoinCLIP without entropy inferred image-text entailment order where text always entail image on different datasets. The metrics in color represent best performance metric for particular dataset. We observe that PoinCLIP outperforms the Poincaré method where we always assume text to be entailing image, in 14 out of 18 datasets.