

# Continual Hyperbolic Learning of Instances and Classes

Anonymous authors

Paper under double-blind review

## Abstract

Instance-level continual learning addresses the challenging problem of recognizing and remembering specific instances of object classes in an incremental setup, where new instances appear over time. Continual learning of instances forms a more fine-grained challenge than conventional continual learning, which is only concerned with incremental discrimination at the class level. In this paper, we argue that for real-world continual understanding, we need to recognize samples both at the instance and class level. We find that classes and instances form a hierarchical structure and propose HyperCLIC, a hyperbolic continual learning algorithm for visual instances and classes, to enable us to learn from this structure. We introduce continual hyperbolic classification and distillation objectives, allowing us to embed the hierarchical relations between classes and from classes to instances. Empirical evaluations show that HyperCLIC can operate effectively at both levels of granularity and with better hierarchical generalization, outperforming well-known continual learning algorithms. The code is included with this submission and will be made publicly available.

## 1 Introduction

Continual learning addresses a long-standing challenge in machine learning: learning from new classes leads to the catastrophic forgetting of old classes (Kirkpatrick et al., 2017; Wu et al., 2019). As a result of catastrophic forgetting, we cannot simply fine-tune models on new data. Many solutions have been proposed to enable incremental learning from new classes while retaining recognition performance from old classes. Well-known solutions include data replay (Bang et al., 2021; Wang et al., 2021), model regularization (Yin et al., 2021; Lee et al., 2020), and knowledge distillation (Kang et al., 2022; Dong et al., 2022).

Where continual learning typically performs class-level discrimination, recent works have broadened the scope to instance-level continual learning. Instance-level continual understanding is essential in many real-world domains such as robotics, where classifying the specific instance of objects helps the robot to decide where to place or how to use them (Ammirato, 2019; Singh et al., 2014; Held et al., 2016). The EgoObjects dataset (Zhu et al., 2023) and its latest instance-level continual challenge (Pellegrini et al., 2022) highlight the difficult and open-ended nature of this problem.

In this work, we strive to perform continual learning simultaneously at the instance and class level. In many real-world domains, recognition only at one granularity level is insufficient. On one hand, class-level continual learning is invariant to class instances by design. On the other hand, instance-level continual learning without class-level awareness is prone to making significant mistakes. For domains such as robotics, self-driving cars, and more, it is crucial that in the case of instance-level mistakes, the class-level prediction is still correct to avoid accidents and to be able to generalize to new instances quickly. Classes and instances are hierarchically related, as classes have their coarse-to-fine hierarchy (Fellbaum, 2010; Deng et al., 2009), and instances add an additional layer to the hierarchy (Yan et al., 2024). Figure 1 visualizes the hierarchical organization of classes and their instances.

We propose HyperCLIC, a hyperbolic continual learner that leverages the class-instance hierarchy for joint instance- and class-level recognition. We first show how the class-instance hierarchy can be embedded as prototypes in hyperbolic space. We then outline hyperbolic classification and distillation losses to enable incremental learning and hierarchical knowledge retention at both levels of granularity. Experimental evalua-

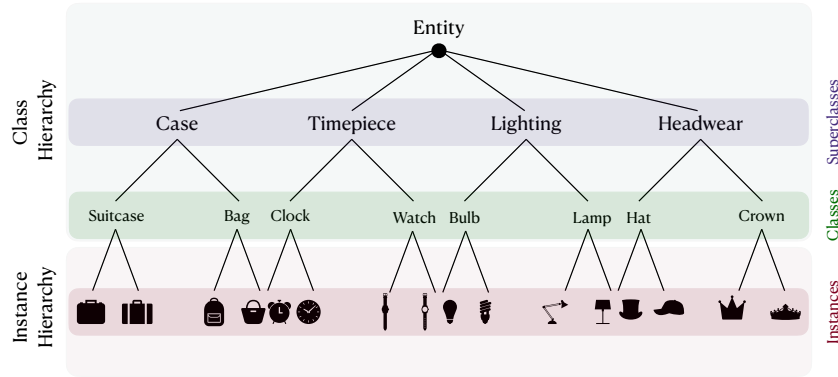


Figure 1: **Recognizing classes and other instances simultaneously** is important in many real-world applications. By adding instances as an additional layer to the object hierarchy and learning representations that capture the joint hierarchy, we can classify samples at multiple levels of granularity.

tion on EgoObjects and Split CIFAR-100 highlight the potential of our approach, with the ability to perform high-quality recognition at both instance- and class-level.

## 2 Literature Review

### 2.1 Continual Learning

Continual learning focuses on the ability of a model to learn and retain knowledge incrementally over time. This field is studied under three main scenarios (Van de Ven et al., 2022; Chen & Liu, 2018): Task-incremental, class-incremental, and domain-incremental learning. In task-incremental learning (Van de Ven et al., 2022; Kirkpatrick et al., 2017; Li & Hoiem, 2017), a model is trained on incremental tasks with clear boundaries. The task ID is known during test time. In contrast, in the class-incremental learning (Guo et al., 2022a; Kim et al., 2022; 2023; Wu et al., 2019; Rebuffi et al., 2017), the task ID is not provided, making task-incremental a particular case of class-incremental. Domain-incremental (Shi & Wang, 2024; Garg et al., 2022b; Kalb et al., 2021; Mirza et al., 2022; Wang et al., 2022b) focuses on scenarios where the data distribution’s incremental shift is explicitly modeled. In this work, we will concentrate on class-incremental learning.

Regardless of continual scenarios, the main challenge of continual learning is to alleviate catastrophic forgetting with only limited access to the previous data (Evron et al., 2022; Shi et al., 2021). Catastrophic forgetting means that performance on previously learned tasks degrades significantly when learning new tasks. To address this, various strategies have been developed. Zhou et al. (2023) groups class-incremental methods into data-centric, model-centric, and algorithm-centric. Data-centric methods concentrate on solving class-incremental learning with exemplars by using data replay (Prabhu et al., 2020; Bang et al., 2021; Chaudhry et al., 2018a; Aljundi et al., 2019; Rolnick et al., 2019; Shin et al., 2017; Wang et al., 2021) or data regularization techniques (Lopez-Paz & Ranzato, 2017; Chaudhry et al., 2018b; Zeng et al., 2019). Model-centric methods either regularize the model parameters from drifting away (Kirkpatrick et al., 2017; Zenke et al., 2017; Yin et al., 2021; Lee et al., 2020) or dynamically expand the network structure for stronger representation ability (Yan et al., 2021; Zhou et al., 2022b; Wang et al., 2022a; Douillard et al., 2022). Algorithm-centric methods either utilize knowledge distillation to resist forgetting (Rebuffi et al., 2017; Kang et al., 2022; Li & Hoiem, 2017; Hou et al., 2019; Buzzega et al., 2020; Dong et al., 2022; Douillard et al., 2020) or rectify the bias in the model (Wu et al., 2019; Belouadah & Popescu, 2019; Zhou et al., 2022a; Zhao et al., 2020). Unlike existing continual methods that do not consider data geometry and use Euclidean space for all data, we advocate using hyperbolic embeddings for continual learning to capture class-class and class-instance hierarchies.

While class-level continual learning has received significant attention, the problem of joint instance- and class-level continual learning remains relatively unexplored. The recent EgoObjects benchmark (Zhu et al., 2023), based on earlier benchmarks (Lomonaco & Maltoni, 2017; She et al., 2020), makes it possible to investigate continual learning at the instance level. Incremental learning of instances has many applications, such as in (i) robotics, *i.e.*, recognizing and interacting with specific objects in the environment; (ii) in NLP, *i.e.*, identifying and classifying new instances of entities (*e.g.*, people, locations, organizations); (iii) in medical imaging, *i.e.*, identifying new forms of lesions and tumors; and (iv) in product recognition, *i.e.*, product recommendations, inventory management, and visual search. Parshotam & Kilickaya (2020) realize the need to incrementally recognize different instances of the same object class using a metric learning approach under an object re-identification setting. In contrast, we propose continually learning both instance-level and class-level representations. According to the challenge winners (Pellegrini, 2022), iCaRL (Rebuffi et al., 2017) is the best-performing model for instance-level classification on the EgoObjects dataset. Therefore, our proposed hyperbolic continual learner takes inspiration from the core design choices of iCaRL.

## 2.2 Hyperbolic Learning

Hyperbolic learning has gained considerable attention in deep learning in embedding taxonomies and tree-like structures (Nickel & Kiela, 2017; Ganea et al., 2018a; Law et al., 2019; Nickel & Kiela, 2018), graphs (Liu et al., 2019; Chami et al., 2019; Bachmann et al., 2020; Dai et al., 2021), and text (Tifrea et al., 2018; Zhu et al., 2020; Dhingra et al., 2018; Leimeister & Wilson, 2018). Hyperbolic space is a space with constant negative curvature that can be thought of as a continuous version of a tree, making it a good choice to embed any finite tree while preserving the distances (Ungar, 2008; Hamann, 2018). Based on the tree-like behavior, Nickel & Kiela (2017) introduce a new approach to embed symbolic data in the Poincaré ball model, a particular hyperbolic space model. Ganea et al. (2018a) take a step forward and improve the Poincaré embeddings (Nickel & Kiela, 2017) using a model based on the geodesically convex entailment cones, showing the effectiveness when embedding data with a hierarchical structure. Furthermore, hyperbolic space has also been used to develop intermediate layers (Ganea et al., 2018a; Cho et al., 2019), and deep neural networks (Ganea et al., 2018b; Shimizu et al., 2020).

Following the initial success, hyperbolic space has shown success and gained attention in computer vision tasks in supervised and unsupervised learning (Mettes et al., 2024). Hyperbolic embeddings have shown benefits in classification and few-shot learning (Khrulkov et al., 2020; Gao et al., 2021; Guo et al., 2022b; Ghadimi Atigh et al., 2021b), zero-shot learning (Liu et al., 2020; Hong et al., 2023), segmentation (Atigh et al., 2022; Chen et al., 2023), out-of-distribution generalization (Ganea et al., 2018b; van Spengler et al., 2023), uncertainty quantification (Atigh et al., 2022; Chen et al., 2023), contrastive learning (Yue et al., 2023; Ge et al., 2023), hierarchical representation learning (Long et al., 2020; Dhall et al., 2020; Doorenbos et al., 2024), generative learning (Cho et al., 2024; Dai et al., 2020; Mathieu & Nickel, 2020), and vision-language representation learning (Ibrahimi et al., 2024; Desai et al., 2023). To our knowledge, Gao et al. (2023) is the only study that utilizes non-Euclidean geometries for class-level continual learning. Their approach introduces an expanding geometry of mixed-curvature space, particularly targeting low-memory scenarios. In contrast, we propose hyperbolic continual learning tailored for joint learning of instances and classes in a default memory regime.

## 3 HyperCLIC

In this section, we present our method for continual learning of instances and classes using hyperbolic geometry. Our goal is to continually learn hierarchy-aware representations, enabling classification at different levels of granularity: instance-, class-, and superclass-level. We follow a class-incremental learning setup, where new instance classes are introduced at each new task. Consider a sequence of  $T$  training tasks  $\mathcal{D}^1, \mathcal{D}^2, \dots, \mathcal{D}^T$  with non-overlapping instances, where  $\mathcal{D}^t = (x_i^t, y_i^t)_{i=1}^{n_t}$  is the  $t$ -th incremental step with  $n_t$  training samples. Each  $x_i^t \in \mathbb{R}^D$  is an example of instance  $y_i \in Y_t$ , where  $Y_t$  is the label space of task  $t$ . We formalize the hierarchy of instances, classes, superclasses, and other nodes as a tree  $\mathcal{T} = (V, E)$ . Each label in  $Y_t$  corresponds to a leaf node label of tree  $\mathcal{T}$ , as detailed in Section 3.1. We only have access to samples of instances in  $\mathcal{D}^t$  when training task  $t$ . Figure 2 illustrates the two main components of

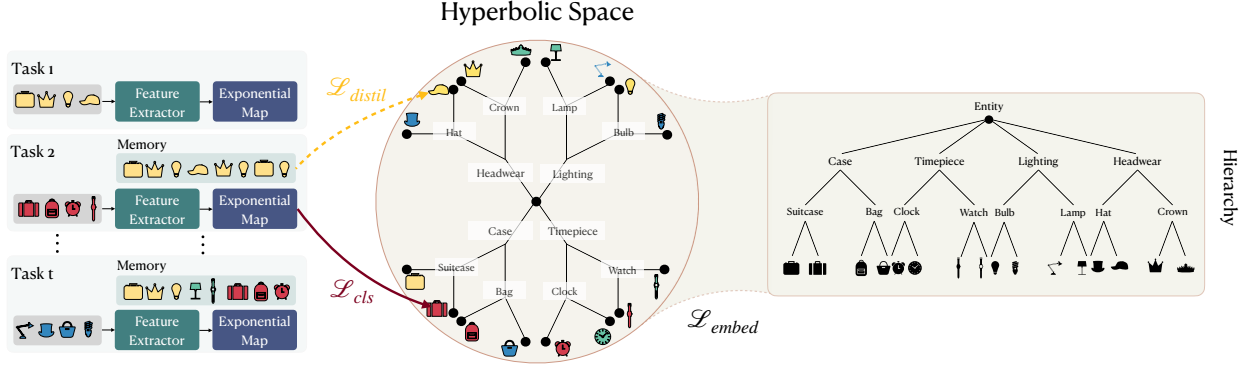


Figure 2: **Overview of HyperCLIC** . The class-instance hierarchy is projected into a fixed shared hyperbolic space. When learning an instance (e.g., instance A), samples of A go through the feature extractor and are mapped into the shared hyperbolic space using the exponential map. These samples are then pushed toward their hyperbolic instance prototype via classification loss and are encouraged to maintain hyperbolic logits from previous classes through distillation loss.

**HyperCLIC** : First, we embed the class-instance hierarchy into hyperbolic space using Poincaré embeddings and entailment cones, leveraging hyperbolic space’s ability to model hierarchical relationships. Second, we perform continual hyperbolic alignment between visual inputs and the embedded hierarchy. We apply a hyperbolic prototype-based loss for classifying new instances and a hyperbolic distillation loss to maintain the consistency of embeddings for previously seen instances, ensuring that the model respects the hierarchical structure during continual learning.

### 3.1 Embedding class-instance hierarchies in hyperbolic space

The joint class-instance tree hierarchy  $\mathcal{T} = (V, E)$  (Figure 1) consists of four types of vertices  $V = V_I \cup V_C \cup V_S \cup V_O$  and four types of directed edges  $E = E_{IC} \cup E_{CS} \cup E_{SO} \cup E_{OO}$ . Vertices  $V_I = \{i_1, i_2, \dots, i_{n_I}\}$  are the set of instances,  $V_C = \{c_1, c_2, \dots, c_{n_C}\}$  the set of classes that are parents of instances,  $V_S = \{s_1, s_2, \dots, s_{n_S}\}$  the set of superclasses that are parents of classes, and  $V_O = \{o_1, o_2, \dots, o_{n_O}\}$  are the set of other remaining nodes. Each node  $v \in V$  corresponds to a distinct label. The number of nodes  $|V|$  equals the number of all labels  $|C|$  in the hierarchy.  $E_{IC} = \{(i_j, c_k) \mid i_j \in V_I, c_k \in V_C\}$  represent edges between instances and classes,  $E_{CS} = \{(c_k, s_m) \mid c_k \in V_C, s_m \in V_S\}$  edges between classes and superclasses,  $E_{SO} = \{(s_m, o_p) \mid s_m \in V_S, o_p \in V_O\}$  edges between superclasses and other nodes, and  $E_{OO} = \{(o_p, o_q) \mid o_p \in V_O, o_q \in V_O\}$  represents the remaining 1-hop hypernymy relations. This formalization defines the hierarchical relationships in a tree structure where each node  $V$  (instance, class, superclass, and other) except the root has exactly one parent. Each instance  $V_I$  is a leaf node in the tree structure  $\mathcal{T}$ .

Given the class-instance tree hierarchy  $\mathcal{T}$ , our goal is to embed the symbolic tree representation *a priori* into a hyperbolic embedding space that incorporates the hierarchical relations between classes and from classes to instances. We define the hyperbolic manifold using the Poincaré ball model (Nickel & Kiela, 2017), as it is well-suited for gradient-based optimization. More formally, let  $\mathbb{B}_c^d = \{z \in \mathbb{R}^d \mid \|z\| < 1\}$  be the open d-dimensional unit ball, where  $\|\cdot\|$  denotes the Euclidean norm, and  $c$  denotes the curvature. The Poincaré ball model corresponds to the Riemannian manifold  $(\mathbb{B}_c^d, g_z^{\mathbb{B}})$  with the Riemannian metric tensor  $g_z^{\mathbb{B}} = 4(1 - c\|z\|^2)^{-2}\mathbf{I}_d$ .

Let  $d_{\mathcal{T}}(v_i, v_j)$  denote the graph distance between two nodes  $v_i, v_j \in V$  based on their hierarchical relations. Also, let  $d_{\mathbb{B}}(p_i, p_j)$  denote the hyperbolic distance between two points  $p_i, p_j \in \mathbb{B}^d$  in the Poincaré ball. We seek to obtain a set of prototypes  $P = \{p_i\}_{i=1}^{|V|}$  corresponding to  $V = \{v_i\}_{i=1}^{|V|}$  such that  $d_{\mathbb{B}}(p_i, p_j) \propto d_{\mathcal{T}}(v_i, v_j)$ , i.e., the hyperbolic distance  $d_{\mathbb{B}}$  is proportional to the edge distances in graph  $\mathcal{T}$ . The hyperbolic distance

$d_{\mathbb{B}}$  between two points  $p_1, p_2 \in \mathbb{B}^d$  is given by the following equation:

$$d_{\mathbb{B}}(p_1, p_2) = \frac{2}{\sqrt{c}} \operatorname{arctanh}(\sqrt{c} \|\cdot - p_1 \oplus_c p_2\|). \quad (1)$$

where  $\langle p_1, p_2 \rangle$  denotes the inner product of two vectors  $p_1$  and  $p_2$  and  $\oplus$  denotes the Möbius addition given by:

$$p_1 \oplus_c p_2 = \frac{(1 + 2c\langle p_1, p_2 \rangle + c\|p_2\|^2)p_1 + (1 - c\|p_1\|^2)p_2}{1 + 2c\langle p_1, p_2 \rangle + c^2\|p_1\|^2\|p_2\|^2} \quad (2)$$

We minimize the loss function  $\mathcal{L}_{\text{Poincaré}}$  defined by Nickel & Kiela (2017). Let  $\mathcal{R} = \{(u, v)\}$  denote the set of transitive closures in graph  $\mathcal{T}$ , meaning there is a path from node  $u$  to  $v$ . Ganea et al. (2018a) define  $\mathcal{R}$  as entailment relations, where  $v$  entails  $u$ , or equivalently, that  $u$  is a subconcept of  $v$ .  $\mathcal{L}_{\text{Poincaré}}$  encourages semantically similar objects to be close in the embedding space according to their Poincaré distance:

$$\mathcal{L}_{\text{Poincaré}} = \sum_{(u,v) \in \mathcal{R}} \log \frac{e^{-d_B(u,v)}}{\sum_{v' \in \mathcal{N}(u)} e^{-d_B(u,v')}}. \quad (3)$$

$\mathcal{N}(u) = \{v | (u, v) \notin \mathcal{R}\} \cup \{u\}$  denotes the set of negative examples for  $u$ . Negative examples  $\mathcal{N}(u)$  include all nodes  $v$  that do not entail  $u$ . While the Poincaré loss results in hyperbolic prototypes  $P$  in tree-shaped regions on  $\mathbb{B}^d$ , there are no guarantees of entailment, *i.e.*, of a partial order relationship that requires the region of each subtree to be fully covered by their parent tree. Therefore, following Ganea et al. (2018a), we apply a max-margin entailment loss  $\mathcal{L}_{\text{entail}}$  on the extracted Poincaré embeddings  $P$ , to enforce entailment regions:

$$\mathcal{L}_{\text{entailment}} = \sum_{(u,v) \in \mathcal{R}} E(u, v) + \sum_{(u',v') \in \mathcal{N}} \max(0, \gamma - E(u', v')). \quad (4)$$

The energy function  $E(u, v) := \max(0, \Xi(u, v) - \psi(v))$  measures how far point  $u$  from belonging to the entailment cone  $\psi(v)$  is. The first term encourages  $u$  to be part of the entailment cone  $\psi(v)$  for  $(u, v) \in \mathcal{R}$ . The second term pushes negative samples  $(u', v') \in \mathcal{N}$  angularly away for a minimum margin  $\gamma > 0$  if they don't share an entailment cone. For full details of Equation 4, we refer to Ganea et al. (2018a).

We observe that the entailment loss tends to bring the prototypes, especially those of the instances  $V_I$ , too close to each other. To counteract this, we apply a separation loss  $\mathcal{L}_S$ , similar to the approach in Long et al. (2020), which ensures that all prototypes  $P$  are adequately separated.

$$\mathcal{L}_S(P) = \bar{\mathbf{I}}^T (\bar{P} \bar{P}^T - I) \bar{\mathbf{I}}, \quad (5)$$

where  $\bar{P}$  denotes the vector-wise  $l_2$ -normalization of  $P$ . The proposed loss function minimizes the cosine similarity between any two prototypes. After obtaining the prototypes for all nodes in the tree  $\mathcal{T}$ , we only use the prototypes  $P = \{P_y\}_{y=1}^{|V_I|}$  corresponding to the instance (leaf) nodes  $V = \{v_y\}_{y=1}^{|V_I|}$  in the next section.

### 3.2 Continual hyperbolic learning

We use a backbone  $\varphi(\cdot; \theta_t)$  to extract features from a training sample  $x_i$ , where  $\theta_t$  denotes the model parameters at timestep  $t$ . The extracted features form a Euclidean representation  $\varphi(x_i; \theta_t)$ . We assume that the backbone output is in the tangent space  $T_x \mathcal{M}$ , while the extracted prototypes from Section 3.1 are in the hyperbolic space  $\mathcal{M}$ . Thus, we need to project the Euclidean representation to the hyperbolic space  $\mathcal{M}$  through the exponential map function.

$$z_i^t = \exp_0(\varphi(x_i; \theta^t)). \quad (6)$$

Here,  $z_i^t$  is the hyperbolic representation of  $\varphi(x_i; \theta_t)$ . The exponential map function is given as (Ganea et al., 2018b):

$$\exp_0^c(x) = \tanh(\sqrt{c}\|x\|) \cdot \frac{x}{\sqrt{c}\|x\|}, \quad (7)$$

where  $c$  denotes the curvature of the hyperbolic space. A higher  $c$  indicates a more curved space. The exponential map embeds visual representation into the hyperbolic space where the hierarchy of classes and instances are embedded as prototypes. Our goal is to minimize the hyperbolic distance of representation  $z$  with its instance prototype  $P_y$ . Thus, we define the hyperbolic logit  $h(z, y)$  as the negative hyperbolic distance of vector representation  $z$  and all prototypes  $y$ :

$$h(z, y) = -\frac{d_{\mathbb{H}}(z, P_y)}{\tau}. \quad (8)$$

The temperature parameter  $\tau$  controls the entropy of the probability distribution while preserving the relative ranks of each class. We can make the probability distribution more peaked or smoother by adjusting this parameter. The conditional probability  $p(y_i|z_i^t)$  can be derived by the softmax of the hyperbolic logit:

$$p(y_i|z_i^t) = \frac{e^{h(z_i^t, y_i)}}{\sum_{j=1}^{|C|} e^{h(z_i^t, y_j)}}. \quad (9)$$

Once the hyperbolic logits are calculated, we can compute the classification loss as the cross entropy between the hyperbolic logits and the targets:

$$\mathcal{L}_{cls} = -\frac{1}{|\mathcal{D}^t|} \sum_{(x_i, y_i) \in \mathcal{D}^t} \sum_{i=1}^{|C|} y_i \cdot \log\left(\frac{e^{h(z_i^t, y_i)}}{\sum_{j=1}^{|C|} e^{h(z_i^t, y_j)}}\right). \quad (10)$$

While the classification loss ensures that the new classes learn their representation, the ultimate goal of class-incremental learning is to continually learn a model that works both for the old and the new classes. Formally, the model should not only acquire the knowledge from the current task  $\mathcal{D}^t$  but also preserve the knowledge from former tasks  $\mathcal{D}^{b< t}$ . We use a hyperbolic distillation loss to ensure that the model predicts the same hyperbolic logits for the old classes  $\mathcal{D}^{b< t}$  at time  $t$  as time  $t-1$ . Following Rebuffi et al. (2017), we keep a copy of the model parameters at time  $t-1$  as  $\theta^{t-1}$  and define the probability distribution at time  $t-1$  as  $p(y_i|z_i^{t-1})$  where  $(x_i, y_i) \in \mathcal{D}^{b< t}$ :

$$p(y_i|z_i^{t-1}) = \frac{e^{h(z_i^{t-1}, y_i)}}{\sum_{j=1}^{|C|} e^{h(z_i^{t-1}, y_j)}}. \quad (11)$$

The hyperbolic distillation loss is calculated as the cross-entropy loss between the probability distribution at time  $t$  and  $t-1$ :

$$\mathcal{L}_{distil} = -\frac{1}{|\mathcal{D}^{b< t}|} \sum_{(x_i, y_i) \in \mathcal{D}^{b< t}} \sum_{i=1}^{|C^{b< t}|} \left( \frac{e^{h(z_i^{t-1}, y_i)}}{\sum_{j=1}^{|C^{b< t}|} e^{h(z_i^{t-1}, y_j)}} \right) \cdot \log\left(\frac{e^{h(z_i^t, y_i)}}{\sum_{j=1}^{|C^{b< t}|} e^{h(z_i^t, y_j)}}\right). \quad (12)$$

Samples from old classes  $\mathcal{D}^{b< t}$  are selected at the end of each task by the herding strategy (Welling, 2009), which is a commonly used strategy aiming to select the most representative samples of each class. These exemplars are saved in memory and are representative data points of each known instance class. In addition to distillation loss, the exemplars are also used during inference. During inference, the closeness to exemplars determines the final prediction for a given data point. We perform the nearest mean of exemplar classification following Rebuffi et al. (2017):

$$y_i^* = \arg \min_{y=1, \dots, n_I} \|\varphi(x_i) - \mu_y\| \quad (13)$$

Where  $\mu_y$  denotes the mean of exemplars for the instance class  $y$  and  $n_I$  denotes the number of instance classes. The exemplars and their means could be calculated in hyperbolic space, however, we observe that the Euclidean representations are already aligned with the hierarchy and conclude that there is no need for an extra hyperbolic computation. Overall, our method aligns classification and distillation losses with hyperbolic representations of a fixed hierarchy, leveraging the class-instance hierarchy with minor modifications to existing continual learning methods. This alignment allows continual classification at multiple levels of granularity such as instance and class-level.

## 4 Experimental Setup

**Datasets** We conduct experiments using the class-incremental instance-level continual learning benchmark EgoObjects (Zhu et al., 2023), first introduced at the 3rd CLVISION Challenge (Pellegrini et al., 2022). This dataset comprises a stream of 15 tasks, each created by cropping the main object from egocentric videos. In total, there are 1110 unique instances from 277 classes. Since the test set labels are not provided in the challenge, we split the training set into training (80%), testing (10%), and validation (10%) subsets. Following the challenge convention, we divide the dataset into non-overlapping sets of `gaia_ids`, which are identifiers for each unique video clip. Additionally, we experiment with iCIFAR-100 (Krizhevsky et al., 2009; Rebuffi et al., 2017), which contains  $i = 10$  tasks, each with 10 classes. Although CIFAR-100 is designed for class-level classification and lacks instances within classes, we treat the classes as instances to facilitate comparison with other class-incremental continual learning methods.

**Hierarchies** For the EgoObjects (Zhu et al., 2023) dataset, we construct a hierarchy using WordNet (Fellbaum, 2010), a comprehensive lexical database of the English language. We match each class in the dataset with the most frequent noun synset in WordNet. To find the path from the leaf nodes to the root node, we iteratively choose the hypernym with the shortest path to the root. If this path does not include the `physical_entity.n.01` node, we select the next shortest path that does. We manually add classes that do not match any synsets to the hierarchy. Categories without any instances (91 out of 277) are removed. We also remove nodes with only one child, except for parents of leaf nodes. The resulting unbalanced hierarchy has a maximum of 7 instances per category, a minimum of 1, and a depth of 12. Our automated hierarchy extraction method can be adapted to other datasets. In CIFAR-100 (Krizhevsky et al., 2009), classes are grouped into 20 superclasses. Each image has a *fine* label (the class) and a *coarse* label (the superclass). Following Yan (2021), we group the 20 superclasses into 7 ultra classes that are children of the root node.

**Implementation details** For HyperCLIC EgoObjects experiments, we train a ResNet34 backbone for 3 epochs following Pellegrini (2022), with a `batch_size` of 90, and a learning rate of 0.01 for 15 tasks. The memory size is 3500 according to the challenge. The pretrained models are pretrained with `ResNet34_Weights.IMAGENET1K_V1`, `Wide_ResNet50_2_Weights.IMAGENET1K_V2` and `RegNet_X_16GF_Weights.IMAGENET1K_V2`. For the hyperbolic parameters, the curvature of the Poincaré model (PoincaréBallExact Kochurov et al. (2020)) is 1, and the temperature  $\tau$  is 10. The hyperbolic prototypes are 64-dimensional using Ganea et al. (2018a) with 150 Poincaré epochs and 50 entailment epochs, plus separation following Long et al. (2020) with a learning rate of 1 for 500 epochs.

**Evaluation metrics** We perform the evaluation on both standard and continual hierarchical metrics. The standard metrics include average forgetting, average accuracy, and per-task accuracies. In our work, the accuracy metric reflects only the instance-level accuracy, so we refer to it as the continual instance-level accuracy. The instance-level accuracy on task  $t$  after incremental learning of task  $t'$  is defined as:

$$\text{Acc}_{instance}^{tt'} = \frac{1}{n_t} \sum_{i=1}^{n_t} 1\{\hat{y}_i^{tt'} = y_i^t\} \quad (14)$$

where  $\hat{y}_i^{tt'}$  is the predicted class for instance label  $y_i^t$  evaluated on the test set of  $t$ -th task after incremental learning of the  $t'$ -th task. To evaluate continual hierarchical consistency and robustness, we report continual class-level and superclass-level accuracies, inspired by sibling and cousin accuracy in Ghadimi Atigh et al. (2021a). For each sample  $x_i^t$  and its ground-truth instance label  $y_i^t$ , let  $p(y_i^t)$  be the parent class  $y_i^t$  and  $gp(y_i^t)$  be the grandparent class of  $y_i^t$ . In class-level accuracy, a prediction is also correct if it shares a parent with the target class. The class-level accuracy on task  $t$  after incremental learning of task  $t'$  is defined as:

$$\text{Acc}_{class}^{tt'} = \frac{1}{n_t} \sum_{i=1}^{n_t} 1\{p(\hat{y}_i^{tt'}) = p(y_i^t)\} \quad (15)$$

Table 1: **Experimental results on EgoObjects** with a ResNet34 backbone. Both with and without pre-training, HyperCLIC performs best.

	Pretrain	Accuracy $\uparrow$			LCA $\downarrow$	Forgetting $\downarrow$
		Instance	Class	Superclass		
Naive (fine-tune)		0.44	1.07	1.58	5.74	10.38
EWC (Kirkpatrick et al., 2017)		0.22	0.66	1.42	5.26	10.38
iCaRL (Rebuffi et al., 2017)		20.05	21.39	22.24	5.45	10.57
CoPE (De Lange & Tuytelaars, 2021)		3.14	4.45	4.91	5.63	25.38
GDumb (Prabhu et al., 2020)		0.50	1.27	1.88	5.50	5.87
DER (Yan et al., 2021)		19.59	20.59	21.14	5.61	35.39
HyperCLIC		<b>41.76</b>	<b>45.91</b>	<b>48.04</b>	<b>4.93</b>	<b>4.17</b>
Naive (fine-tune)	✓	6.06	15.55	18.24	4.70	93.03
EWC (Kirkpatrick et al., 2017)	✓	6.07	15.77	18.37	4.74	92.83
iCaRL (Rebuffi et al., 2017)	✓	81.63	87.02	87.81	3.89	<b>3.77</b>
CoPE (De Lange & Tuytelaars, 2021)	✓	37.15	49.89	52.51	4.22	32.40
GDumb (Prabhu et al., 2020)	✓	61.51	71.00	72.82	4.02	18.22
DER (Yan et al., 2021)	✓	80.00	84.40	85.20	4.28	10.59
HyperCLIC	✓	<b>84.81</b>	<b>91.94</b>	<b>92.67</b>	<b>2.99</b>	7.08

In the superclass-level accuracy, the predicted labels must share a grandparent with the target class to count as correct. The superclass-level accuracy on task  $t$  after incremental learning of task  $t'$  is defined as:

$$\text{Acc}_{superclass}^{tt'} = \frac{1}{n_t} \sum_{i=1}^{n_t} 1\{gp(\hat{y}_i^{tt'}) = gp(y_i^t)\} \quad (16)$$

Similar to Bertinetto et al. (2020); Garg et al. (2022a), we also report the continual distance to the Lowest Common Ancestor (LCA) that captures the mistake’s hierarchical severity. For the wrong predictions, LCA distance reports the average number of edges between the predicted node and the LCA of the predicted and ground-truth nodes. LCA distance reveals how hierarchically far off the inaccurate predictions are, considering the joint class-instance hierarchy as the ground truth.

Following the 3rd CLVision challenge (Verwimp et al., 2023) and the SSLAD competition (Pellegrini et al., 2022)), all continual hierarchical metrics are calculated as the average mean over all the tasks. Here, we take the continual instance-level accuracy as an example. Let  $\text{Acc}_{instance}^{ij} \in [0, 1]$  denote the instance-level classification accuracy evaluated on the test set of the  $i$ -th task after incremental learning of the  $j$ -th task. The average mean instance-level accuracy is defined as:

$$\text{Accuracy}_{instance-level} = \frac{1}{T^2} \sum_{j=1}^T \sum_{i=1}^T \text{Acc}_{instance}^{ij} \quad (17)$$

where  $T$  denotes the number of all tasks. We use the same formulation for average mean class-level accuracy and superclass-level accuracies. We also report the per-task accuracies at the end of the last task.

## 5 Results

### 5.1 Comparison to existing methods

**Hierarchical and continual metrics** Table 1 reports the results of different baselines on the EgoObjects dataset. If trained from scratch, HyperCLIC achieves a significantly higher instance-level accuracy (41.76%), class-level accuracy (45.91%), and superclass-level accuracy (48.04%) compared to all baselines. The closest competitor is iCaRL (Rebuffi et al., 2017) with an instance-, class-, and superclass-level accuracy of 20.05%,



Table 2: **Results on the Split CIFAR-100 dataset** with a ResNet34 backbone. HyperCLIC consistently outperforms iCaRL in hierarchical metrics. At the instance level, HyperCLIC is comparable with iCaRL without pretraining and better with pretraining.

	Pretrain	Accuracy $\uparrow$			LCA $\downarrow$	Forgetting $\downarrow$
		Instance	Class	Superclass		
iCaRL		47.55	62.32	73.38	2.22	14.36
HyperCLIC		45.52	62.31	74.76	2.15	24.76
iCaRL	$\checkmark$	55.86	70.40	79.82	2.12	11.50
HyperCLIC	$\checkmark$	56.46	74.11	84.52	1.95	15.39

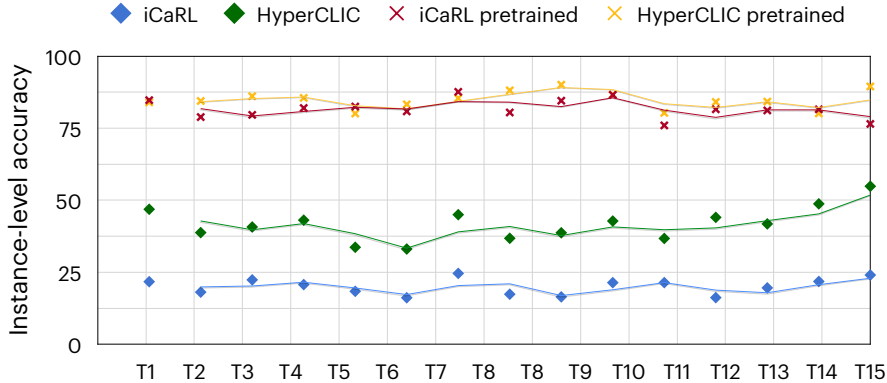


Figure 3: **The final per-task results and their moving average on EgoObjects** with a ResNet34 backbone. HyperCLIC outperforms iCaRL by maintaining prior knowledge while delivering higher accuracy on new tasks.

21.39%, and 22.24%, respectively, which is less than half of HyperCLIC’s performance. This demonstrates HyperCLIC’s strength in maintaining performance across different levels of granularity. HyperCLIC shows a lower LCA value (4.93) indicating that the inaccurate predictions are less hierarchically severe compared to other baselines. The second-best is EWC (Kirkpatrick et al., 2017) with 5.26. HyperCLIC also shows the lowest forgetting rate (4.17%), indicating better retention of learned knowledge. For the pretrained scenario, HyperCLIC again leads with the instance-level accuracy of 84.81%, followed by iCaRL (Rebuffi et al., 2017) with 81.63%, and DER (Yan et al., 2021) with 80%. DER is a strong non-transformer-based continual baseline. In both trained-from-scratch and pretrained scenarios, HyperCLIC consistently outperforms all baseline methods across all metrics. It excels in instance-, class-, and superclass-level accuracies while showing lower LCA values and minimal forgetting. This indicates that HyperCLIC not only learns a hierarchically-aware representation but also retains knowledge from previous tasks, proving its superiority and robustness over other existing methods.

Given that iCaRL (Rebuffi et al., 2017) has the highest performance after HyperCLIC, we use it as the primary comparison for evaluating our method in subsequent experiments. Figure 3 illustrates the instance-level accuracy for each task (T1 to T15) after the completion of all tasks. In both non-pretrained and pretrained models, iCaRL (Rebuffi et al., 2017) demonstrates higher stability but at the cost of not learning effective representations for new tasks, thereby compromising its performance on recent tasks to retain knowledge from previous ones. We can also observe this behavior in Table 1. In contrast, our method not only preserves knowledge from prior tasks but also achieves higher accuracy on the latest tasks compared to the baseline.

**Split CIFAR-100** In addition to evaluating our results on the EgoObjects dataset (Zhu et al., 2023), which serves as an instance-level continual learning benchmark, we also benchmark our method on iCIFAR-

Table 3: **EgoObjects results with different distillation losses** with a ResNet34 backbone. Cross-entropy loss consistently outperforms other distillation losses in both scenarios.

	From scratch					Pretrained				
	Instance	Class	Superclass	LCA	Forgetting	Instance	Class	Superclass	LCA	Forgetting
KL divergence	32.60	35.30	36.95	5.13	-1.40	49.76	62.45	65.46	3.75	-6.00
MSE	36.09	40.33	42.53	4.94	8.18	84.38	91.21	92.33	2.92	7.84
Cross-entropy	41.76	45.91	48.04	4.93	4.17	84.11	91.61	92.54	2.96	7.08

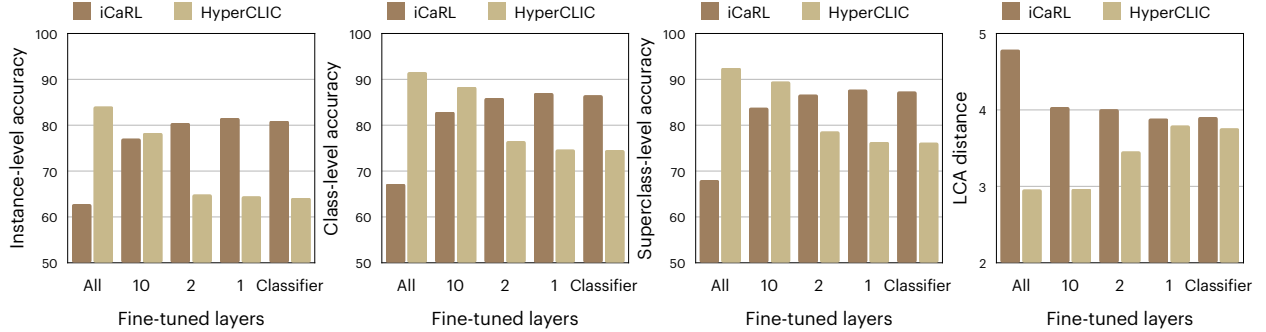


Figure 4: **Performance of HyperCLIC and iCaRL across four metrics** with varying numbers of fine-tuned layers. HyperCLIC excels when all layers are fine-tuned, while iCaRL performs better with more fixed layers. The Table 1 comparison in the pretrained scenario considers the optimal setup for each model.

100 (Krizhevsky et al., 2009) to enable comparisons with other methods. Table 2 presents the performance of HyperCLIC and iCaRL (Rebuffi et al., 2017) on the iCIFAR-100 dataset under both pretrained and training-from-scratch conditions. When training from scratch, HyperCLIC demonstrates performance comparable to iCaRL in terms of instance- and class-level accuracy while surpassing iCaRL in LCA distance and super-class accuracy, indicating that the model effectively learns hierarchical representations. A lower LCA distance signifies that the model’s errors adhere more closely to the underlying hierarchy. We also observe a higher degree of forgetting compared to the EgoObjects dataset. When pretrained, HyperCLIC consistently achieves superior hierarchical metrics than the baseline, although it still exhibits higher forgetting rates. Additionally, there is a noticeable difference between pretrained and from-scratch performance in EgoObjects versus iCIFAR-100, likely due to pretrained backbones being trained with class-level supervision that reduces intra-class variance.

## 5.2 Ablation study

**Distillation loss** A key component of our method is the distillation loss, which ensures that the model retains hyperbolic knowledge from previous tasks. To evaluate this, we tested three different loss functions: Kullback-Leibner (KL) divergence, Mean Squared Error (MSE), and cross-entropy loss. Table 3 presents the results for both pretrained and from-scratch scenarios. KL divergence shows the lowest forgetting rates, at  $-1.40\%$  and  $-6\%$  respectively, for from-scratch and pretrained scenarios. However, it performs the worst in terms of hierarchical metrics and instance accuracy, suggesting it disrupts the balance between plasticity and stability by favoring stability at the expense of plasticity. Cross-entropy loss performs the best across both scenarios, which is why our method is based on this loss function. There is also a more notable difference between MSE and cross-entropy when training from scratch compared to the pretrained scenario.

**Number of fine-tuned layers** We observed that the results of our method and iCaRL differ based on the number of layers fine-tuned in the pretrained scenario. Figure 4 illustrates the impact on four hierarchical metrics as the number of fine-tuned layers varies for both methods. For iCaRL, increasing the number of fine-tuned layers results in worse performance, with lower instance-, class-, and superclass-level accuracies,

Table 4: **Comparison of iCaRL and HyperCLIC on EgoObjects with various pretrained backbones.** HyperCLIC consistently outperforms iCaRL regardless of the architecture, as shown with ResNet34 (highlighted in gray). This demonstrates that HyperCLIC is adaptable to any backbone and, when paired with advanced models, enhances joint continual learning of instances and classes.

	iCaRL					HyperCLIC				
	Instance	Class	Superclass	LCA	Forgetting	Instance	Class	Superclass	LCA	Forgetting
WideResNet50	82.17	87.46	88.08	4.02	3.42	86.50	94.17	94.61	2.58	6.01
RegNetX_16GF	78.63	86.64	87.58	3.55	5.10	87.22	94.11	95.01	2.59	5.99
ResNet34	81.63	87.02	87.81	3.89	3.77	84.81	91.94	92.67	2.99	7.08

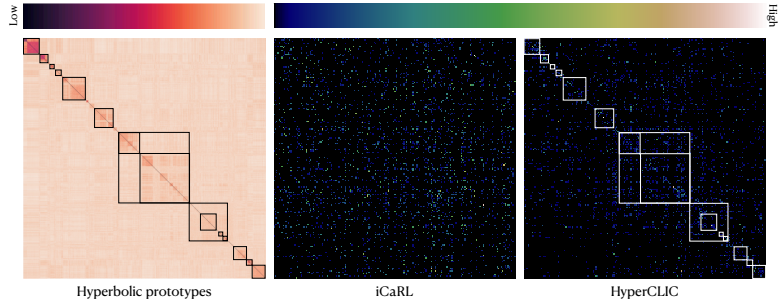


Figure 5: **Left:** the pair-wise distances of the instance-level hyperbolic prototypes. **Middle & right:** the class-level predictions only for the wrong instance-level predictions. The squares highlight the hierarchical structure of HyperCLIC mistakes compared with iCaRL and its similarity to the original hierarchy.

and higher LCA distance. Conversely, for HyperCLIC, more fine-tuned layers lead to improved results. We hypothesize that this is because HyperCLIC enforces a hierarchical structure on the representations, and since the pretrained model was trained in Euclidean space, fine-tuning more layers allows HyperCLIC to leverage the model’s capacity to learn hierarchically-aware representations from the first layer. Given these trends, we use the best scores for each method to compare HyperCLIC with the baseline in the pretrained scenario.

**Backbones** All results reported thus far utilize the ResNet34 backbone. Table 4 compares our method with the leading baseline, iCaRL (Rebuffi et al., 2017), using various pretrained backbones. The findings indicate that HyperCLIC consistently surpasses the baseline across the hierarchical metrics when using WideResNet50, RegNetX, and ResNet34 as the backbone architectures. Notably, HyperCLIC achieves the highest performance with the RegNet backbone, attaining instance-, class-, and superclass-level accuracies of 87.22%, 94.11%, and 95.01%, respectively, compared to the best iCaRL results for WideResNet, which are 82.17%, 87.46%, and 88.08%. Although HyperCLIC improves the LCA distance by nearly one edge, it shows a higher forgetting rate. These results suggest that our method, particularly when combined with advanced models such as transformers, has the potential to significantly enhance the joint continual learning of instances and classes.

**Class-level errors** Figure 5 compares the class-level confusion matrix for incorrectly predicted instances between iCaRL and HyperCLIC, alongside the distances of hyperbolic prototypes. In this matrix, the squares represent children of the same parent node, and larger squares encompass multiple subtrees. The aim of this experiment is to demonstrate that even when instances are misclassified, HyperCLIC’s class-level predictions still follow the hierarchical structure. As shown, iCaRL’s class-level predictions are evenly distributed, while HyperCLIC’s predictions are primarily concentrated along the diagonal. Additionally, a similar hierarchical organization is observed in HyperCLIC when compared with the pairwise hyperbolic prototype distances,














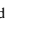
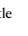
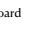














														
	Instance	Class	Instance	Class	Instance	Class	Instance	Class	Instance	Class	Instance	Class	Instance	Class
GT		Backpack		Christmas bush		Footwear		Jug		Picture frame		Oven		Cookie
iCaRL		Backpack		Christmas bush		Jug		Jug		Picture frame		Oven		Cookie
Ours		Backpack		Houseplant		Footwear		Skateboard		Mug		Earphone		Pineapple

Figure 6: **Qualitative examples from EgoObjects** comparing instance and class-level predictions for ground truth (1<sup>st</sup> row), iCaRL (2<sup>nd</sup> row), and our method (3<sup>rd</sup> row). **Green** indicates correct predictions, while **red** and **orange** denote severe and less severe hierarchical errors, respectively. HyperCLIC makes less severe (orange) errors compared to iCaRL.

represented by matching squares. This indicates that HyperCLIC maintains the hierarchical structure, even when making incorrect predictions, resulting in less severe errors.

**Qualitative examples: success and failure cases** Figure 6 presents several images from the test set and compares the instance and class-level predictions of our method against the baseline. In the far-left example, both methods easily make correct instance-level predictions. Whereas, the far-right example is challenging for both methods, as it depicts food items in a drawing, leading both to predict unrelated instances and classes. The middle examples reveal that when HyperCLIC makes an incorrect instance-level prediction, the class-level prediction remains hierarchically close. For example, it might mistake a “Christmas bush” for a “house plant” or a “bottle” for a “mug”. This suggests that by incorporating hierarchical information during training, HyperCLIC reduces the severity of hierarchical errors.

## 6 Conclusion and Discussion

Instance-level continual learning addresses the challenging task of recognizing and remembering specific instances of object classes in an incremental setup, where new instances appear over time. This approach forms a more fine-grained challenge than conventional continual learning, which typically focuses on incremental discrimination at the class level. In this paper, we introduced a method for continually learning at both instance and class levels, arguing that real-world continual understanding requires recognizing samples at multiple layers of granularity. We observed that classes and instances form a hierarchical structure that can be leveraged to enhance learning at both levels. To this end, we proposed HyperCLIC, a hyperbolic continual algorithm designed for jointly learning instances and classes. We introduced continual hyperbolic classification and hyperbolic distillation, which embed the hierarchical relationships between classes and from classes to instances. Our experiments demonstrated that HyperCLIC operates effectively at multiple levels of granularity and achieves superior hierarchical generalization, consistently outperforming strong continual learning baselines. We conducted ablations on different distillation losses, backbone architectures, and pretrained models. Additionally, we presented qualitative analysis on class-level errors and provided evidence that our method makes less hierarchically severe mistakes. HyperCLIC enables real-world continual understanding, where recognizing and remembering both instances and classes over time is crucial.

## References

- Rahaf Aljundi, Min Lin, Baptiste Goujaud, and Yoshua Bengio. Gradient based sample selection for online continual learning. *NeurIPS*, 2019.
- Phil Ammirato. *Recognizing Fine-Grained Object Instances for Robotics Applications*. PhD thesis, The University of North Carolina at Chapel Hill, 2019.
- Mina Ghadimi Atigh, Julian Schoep, Erman Acar, Nanne Van Noord, and Pascal Mettes. Hyperbolic image segmentation. In *CVPR*, 2022.
- Gregor Bachmann, Gary Bécigneul, and Octavian Ganea. Constant curvature graph convolutional networks. In *ICML*, 2020.
- Jihwan Bang, Heesu Kim, YoungJoon Yoo, Jung-Woo Ha, and Jonghyun Choi. Rainbow memory: Continual learning with a memory of diverse samples. In *CVPR*, 2021.
- Eden Belouadah and Adrian Popescu. Il2m: Class incremental learning with dual memory. In *ICCV*, 2019.
- Luca Bertinetto, Romain Mueller, Konstantinos Tertikas, Sina Samangooei, and Nicholas A Lord. Making better mistakes: Leveraging class hierarchies with deep networks. In *CVPR*, 2020.
- Pietro Buzzega, Matteo Boschini, Angelo Porrello, Davide Abati, and Simone Calderara. Dark experience for general continual learning: a strong, simple baseline. *NeurIPS*, 2020.
- Ines Chami, Zhitao Ying, Christopher Ré, and Jure Leskovec. Hyperbolic graph convolutional neural networks. *NeurIPS*, 2019.
- Arslan Chaudhry, Puneet K Dokania, Thalaiyasingam Ajanthan, and Philip HS Torr. Riemannian walk for incremental learning: Understanding forgetting and intransigence. In *ECCV*, 2018a.
- Arslan Chaudhry, Marc’Aurelio Ranzato, Marcus Rohrbach, and Mohamed Elhoseiny. Efficient lifelong learning with a-gem. *arXiv preprint arXiv:1812.00420*, 2018b.
- Bike Chen, Wei Peng, Xiaofeng Cao, and Juha Röning. Hyperbolic uncertainty aware semantic segmentation. *T-ITS*, 2023.
- Zhiyuan Chen and Bing Liu. *Lifelong machine learning*. Springer, 2018.
- Hyunghoon Cho, Benjamin DeMeo, Jian Peng, and Bonnie Berger. Large-margin classification in hyperbolic space. In *AISTat*, 2019.
- Seunghyuk Cho, Juyong Lee, and Dongwoo Kim. Hyperbolic vae via latent gaussian distributions. *NeurIPS*, 2024.
- Jindou Dai, Yuwei Wu, Zhi Gao, and Yunde Jia. A hyperbolic-to-hyperbolic graph convolutional network. In *CVPR*, 2021.
- Shuyang Dai, Zhe Gan, Yu Cheng, Chenyang Tao, Lawrence Carin, and Jingjing Liu. Apo-vae: Text generation in hyperbolic space. *arXiv preprint arXiv:2005.00054*, 2020.
- Matthias De Lange and Tinne Tuytelaars. Continual prototype evolution: Learning online from non-stationary data streams. In *ICCV*, 2021.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin Johnson, and Shanmukha Ramakrishna Vedantam. Hyperbolic image-text representations. In *ICML*, 2023.
- Ankit Dhall, Anastasia Makarova, Octavian Ganea, Dario Pavllo, Michael Greeff, and Andreas Krause. Hierarchical image classification using entailment cone embeddings. In *CVPR*, 2020.

- Bhuwan Dhingra, Christopher J Shallue, Mohammad Norouzi, Andrew M Dai, and George E Dahl. Embedding text in hyperbolic spaces. *arXiv preprint arXiv:1806.04313*, 2018.
- Jiahua Dong, Lixu Wang, Zhen Fang, Gan Sun, Shichao Xu, Xiao Wang, and Qi Zhu. Federated class-incremental learning. In *CVPR*, 2022.
- Lars Jelte Doorenbos, Pablo Márquez Neila, Raphael Sznitman, and Pascal Mettes. Hyperbolic random forests. *TMLR*, 2024.
- Arthur Douillard, Matthieu Cord, Charles Ollion, Thomas Robert, and Eduardo Valle. Podnet: Pooled outputs distillation for small-tasks incremental learning. In *ECCV*, 2020.
- Arthur Douillard, Alexandre Ramé, Guillaume Couairon, and Matthieu Cord. Dytox: Transformers for continual learning with dynamic token expansion. In *CVPR*, 2022.
- Itay Evron, Edward Moroshko, Rachel Ward, Nathan Srebro, and Daniel Soudry. How catastrophic can catastrophic forgetting be in linear regression? In *Conference on Learning Theory*, 2022.
- C Fellbaum. About wordnet. wordnet. princeton university, 2010.
- Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *ICML*, 2018a.
- Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic neural networks. *NeurIPS*, 2018b.
- Zhi Gao, Yuwei Wu, Yunde Jia, and Mehrtash Harandi. Curvature generation in curved spaces for few-shot learning. In *ICCV*, 2021.
- Zhi Gao, Chen Xu, Feng Li, Yunde Jia, Mehrtash Harandi, and Yuwei Wu. Exploring data geometry for continual learning. In *CVPR*, 2023.
- Ashima Garg, Depanshu Sani, and Saket Anand. Learning hierarchy aware features for reducing mistake severity. In *ECCV*, 2022a.
- Prachi Garg, Rohit Saluja, Vineeth N Balasubramanian, Chetan Arora, Anbumani Subramanian, and CV Jawahar. Multi-domain incremental learning for semantic segmentation. In *WACV*, 2022b.
- Songwei Ge, Shlok Mishra, Simon Kornblith, Chun-Liang Li, and David Jacobs. Hyperbolic contrastive learning for visual representations beyond objects. In *CVPR*, 2023.
- Mina Ghadimi Atigh, Martin Keller-Ressel, and Pascal Mettes. Hyperbolic busermann learning with ideal prototypes. *NeurIPS*, 2021a.
- Mina Ghadimi Atigh, Martin Keller-Ressel, and Pascal Mettes. Hyperbolic busermann learning with ideal prototypes. *NeurIPS*, 2021b.
- Yiduo Guo, Bing Liu, and Dongyan Zhao. Online continual learning through mutual information maximization. In *ICML*, 2022a.
- Yunhui Guo, Xudong Wang, Yubei Chen, and Stella X Yu. Clipped hyperbolic classifiers are super-hyperbolic classifiers. In *CVPR*, 2022b.
- Matthias Hamann. On the tree-likeness of hyperbolic spaces. In *Mathematical proceedings of the cambridge philosophical society*. Cambridge University Press, 2018.
- David Held, Sebastian Thrun, and Silvio Savarese. Robust single-view instance recognition. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2016.
- Jie Hong, Zeeshan Hayder, Junlin Han, Pengfei Fang, Mehrtash Harandi, and Lars Petersson. Hyperbolic audio-visual zero-shot learning. In *ICCV*, 2023.

- Saihui Hou, Xinyu Pan, Chen Change Loy, Zilei Wang, and Dahua Lin. Learning a unified classifier incrementally via rebalancing. In *CVPR*, 2019.
- Sarah Ibrahimi, Mina Ghadimi Atigh, Nanne Van Noord, Pascal Mettes, and Marcel Worring. Intriguing properties of hyperbolic embeddings in vision-language models. *TMLR*, 2024.
- Tobias Kalb, Masoud Roschani, Miriam Ruf, and Jürgen Beyerer. Continual learning for class-and domain-incremental semantic segmentation. In *IEEE IV*, 2021.
- Minsoo Kang, Jaeyoo Park, and Bohyung Han. Class-incremental learning by knowledge distillation with adaptive feature consolidation. In *CVPR*, 2022.
- Valentin Khrulkov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. Hyperbolic image embeddings. In *CVPR*, 2020.
- Gyuhak Kim, Changnan Xiao, Tatsuya Konishi, Zixuan Ke, and Bing Liu. A theoretical study on solving continual learning. *NeurIPS*, 2022.
- Gyuhak Kim, Changnan Xiao, Tatsuya Konishi, and Bing Liu. Learnability and algorithm for continual learning. In *ICML*, 2023.
- James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 2017.
- Max Kochurov, Rasul Karimov, and Serge Kozlukov. Geoopt: Riemannian optimization in pytorch. *arXiv preprint arXiv:2005.02819*, 2020.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
- Marc Law, Renjie Liao, Jake Snell, and Richard Zemel. Lorentzian distance learning for hyperbolic representations. In *ICML*, 2019.
- Janghyeon Lee, Hyeong Gwon Hong, Donggyu Joo, and Junmo Kim. Continual learning with extended kronecker-factored approximate curvature. In *CVPR*, 2020.
- Matthias Leimeister and Benjamin J Wilson. Skip-gram word embeddings in hyperbolic space. *arXiv preprint arXiv:1809.01498*, 2018.
- Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE PAMI*, 2017.
- Qi Liu, Maximilian Nickel, and Douwe Kiela. Hyperbolic graph neural networks. *NeurIPS*, 2019.
- Shaoteng Liu, Jingjing Chen, Liangming Pan, Chong-Wah Ngo, Tat-Seng Chua, and Yu-Gang Jiang. Hyperbolic visual embedding learning for zero-shot recognition. In *CVPR*, 2020.
- Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition. In *Conference on robot learning*. PMLR, 2017.
- Teng Long, Pascal Mettes, Heng Tao Shen, and Cees GM Snoek. Searching for actions on the hyperbole. In *CVPR*, 2020.
- David Lopez-Paz and Marc’Aurelio Ranzato. Gradient episodic memory for continual learning. *NeurIPS*, 2017.
- Emile Mathieu and Maximilian Nickel. Riemannian continuous normalizing flows. *NeurIPS*, 2020.
- Pascal Mettes, Mina Ghadimi Atigh, Martin Keller-Ressel, Jeffrey Gu, and Serena Yeung. Hyperbolic deep learning in computer vision: A survey. *IJCV*, 2024.
- M Jehanzeb Mirza, Marc Masana, Horst Possegger, and Horst Bischof. An efficient domain-incremental learning approach to drive in all weather conditions. In *CVPR*, 2022.

- Maximillian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. *NeurIPS*, 2017.
- Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *ICML*, 2018.
- Kishan Parshotam and Mert Kilickaya. Continual learning of object instances. In *CVPR*, 2020.
- Lorenzo Pellegrini. clvision-challenge-2022 reports. <https://github.com/ContinualAI/clvision-challenge-2022/tree/main/reports>, 2022.
- Lorenzo Pellegrini, Chenchen Zhu, Fanyi Xiao, Zhicheng Yan, Antonio Carta, Matthias De Lange, Vincenzo Lomonaco, Roshan Sumbaly, Pau Rodriguez, and David Vazquez. 3rd continual learning workshop challenge on egocentric category and instance level object understanding. *arXiv preprint arXiv:2212.06833*, 2022. URL <https://sites.google.com/view/clvision2022/challenge>.
- Ameya Prabhu, Philip HS Torr, and Puneet K Dokania. Gdumb: A simple approach that questions our progress in continual learning. In *ECCV*, 2020.
- Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, 2017.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. *NeurIPS*, 2019.
- Qi She, Fan Feng, Xinyue Hao, Qihan Yang, Chuanlin Lan, Vincenzo Lomonaco, Xuesong Shi, Zhengwei Wang, Yao Guo, Yimin Zhang, et al. Openloris-object: A robotic vision dataset and benchmark for lifelong deep learning. In *ICRA*, 2020.
- Guangyuan Shi, Jiabin Chen, Wenlong Zhang, Li-Ming Zhan, and Xiao-Ming Wu. Overcoming catastrophic forgetting in incremental few-shot learning by finding flat minima. *NeurIPS*, 2021.
- Haizhou Shi and Hao Wang. A unified approach to domain incremental learning with memory: Theory and algorithm. *NeurIPS*, 2024.
- Ryohei Shimizu, Yusuke Mukuta, and Tatsuya Harada. Hyperbolic neural networks++. *arXiv preprint arXiv:2006.08210*, 2020.
- Hanul Shin, Jung Kwon Lee, Jaehong Kim, and Jiwon Kim. Continual learning with deep generative replay. *NeurIPS*, 2017.
- Arjun Singh, James Sha, Karthik S Narayan, Tudor Achim, and Pieter Abbeel. Bigbird: A large-scale 3d database of object instances. In *IEEE international conference on robotics and automation (ICRA)*, 2014.
- Alexandru Tifrea, Gary Bécigneul, and Octavian-Eugen Ganea. Poincaré glove: Hyperbolic word embeddings. *arXiv preprint arXiv:1810.06546*, 2018.
- Abraham Albert Ungar. A gyrovector space approach to hyperbolic geometry. *Synthesis Lectures on Mathematics and Statistics*, 2008.
- Gido M Van de Ven, Tinne Tuytelaars, and Andreas S Tolias. Three types of incremental learning. *Nature Machine Intelligence*, 2022.
- Max van Spengler, Erwin Berkhout, and Pascal Mettes. Poincaré resnet. In *ICCV*, 2023.
- Eli Verwimp, Kuo Yang, Sarah Parisot, Lanqing Hong, Steven McDonagh, Eduardo Pérez-Pellitero, Matthias De Lange, and Tinne Tuytelaars. Clad: A realistic continual learning benchmark for autonomous driving. *Neural Networks*, 2023.
- Fu-Yun Wang, Da-Wei Zhou, Han-Jia Ye, and De-Chuan Zhan. Foster: Feature boosting and compression for class-incremental learning. In *ECCV*, 2022a.



- Liyuan Wang, Kuo Yang, Chongxuan Li, Lanqing Hong, Zhenguo Li, and Jun Zhu. Ordisco: Effective and efficient usage of incremental unlabeled data for semi-supervised continual learning. In *CVPR*, 2021.
- Yabin Wang, Zhiwu Huang, and Xiaopeng Hong. S-prompts learning with pre-trained transformers: An occam’s razor for domain incremental learning. *NeurIPS*, 2022b.
- Max Welling. Herding dynamical weights to learn. In *ICML*, 2009.
- Yue Wu, Yinpeng Chen, Lijuan Wang, Yuancheng Ye, Zicheng Liu, Yandong Guo, and Yun Fu. Large scale incremental learning. In *CVPR*, 2019.
- Cheng Yan, Feng Zhao, and Yudong Zhang. From concept to instance: Hierarchical reinforced knowledge graph reasoning. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2024.
- Edwin S Yan. Cifar-100 hierarchy, 2021. URL [https://github.com/edwin-yan/hierarchical\\_classification/blob/master/CIFAR100\\_helper.py](https://github.com/edwin-yan/hierarchical_classification/blob/master/CIFAR100_helper.py).
- Shipeng Yan, Jiangwei Xie, and Xuming He. Der: Dynamically expandable representation for class incremental learning. In *CVPR*, 2021.
- Haiyan Yin, Ping Li, et al. Mitigating forgetting in online continual learning with neuron calibration. *NeurIPS*, 2021.
- Yun Yue, Fangzhou Lin, Kazunori D Yamada, and Ziming Zhang. Hyperbolic contrastive learning. *arXiv preprint arXiv:2302.01409*, 2023.
- Guanxiong Zeng, Yang Chen, Bo Cui, and Shan Yu. Continual learning of context-dependent processing in neural networks. *Nature Machine Intelligence*, 2019.
- Friedemann Zenke, Ben Poole, and Surya Ganguli. Continual learning through synaptic intelligence. In *ICML*, 2017.
- Bowen Zhao, Xi Xiao, Guojun Gan, Bin Zhang, and Shu-Tao Xia. Maintaining discrimination and fairness in class incremental learning. In *CVPR*, 2020.
- Da-Wei Zhou, Fu-Yun Wang, Han-Jia Ye, Liang Ma, Shiliang Pu, and De-Chuan Zhan. Forward compatible few-shot class-incremental learning. In *CVPR*, 2022a.
- Da-Wei Zhou, Qi-Wei Wang, Han-Jia Ye, and De-Chuan Zhan. A model or 603 exemplars: Towards memory-efficient class-incremental learning. *arXiv preprint arXiv:2205.13218*, 2022b.
- Da-Wei Zhou, Qi-Wei Wang, Zhi-Hong Qi, Han-Jia Ye, De-Chuan Zhan, and Ziwei Liu. Deep class-incremental learning: A survey. *arXiv preprint arXiv:2302.03648*, 2023.
- Chenchen Zhu, Fanyi Xiao, Andrés Alvarado, Yasmine Babaei, Jiabo Hu, Hichem El-Mohri, Sean Culatana, Roshan Sumbaly, and Zhicheng Yan. Egoobjects: A large-scale egocentric dataset for fine-grained object understanding. In *ICCV*, 2023.
- Yudong Zhu, Di Zhou, Jinghui Xiao, Xin Jiang, Xiao Chen, and Qun Liu. Hypertext: Endowing fasttext with hyperbolic geometry. In *EMNLP*, 2020.