039

040

041

042

043

HAPPI: Hyperbolic Hierarchical Prototypes for Image Recognition

Anonymous ICCV-BEW submission

Paper ID 11

Abstract

001 Prototypical neural networks have gained prominence in computer vision due to their inherent interpretability, en-002 003 abling decisions based on representative examples without the need for post-hoc explainability methods. However, 004 many prototype-based models overlook the hierarchical re-005 lationships within image features, treating them indepen-006 dently and often resulting in suboptimal performance for 007 008 tasks requiring complex structural understanding. To address this limitation, we propose HAPPI (Hierarchical And 009 **P**rototy**P**ical Image recognition), a framework that lever-010 ages hyperbolic geometry to organize prototypes hierarchi-011 cally within a Lorentzian manifold. By arranging generic 012 013 features near the origin of the hyperboloid and specific features farther away, HAPPI enables the learning of generic 014 prototypes for consistent and defining patterns and specific 015 prototypes for fine-grained and variable details, effectively 016 capturing hierarchical relationships in image data. Our ap-017 018 proach is model-agnostic and can be applied to various prototypical neural networks and backbones. We evalu-019 ate HAPPI across several prototypical model baselines and 020 datasets, demonstrating its versatility and showing that hy-021 perbolic prototypes consistently match or outperform their 022 023 Euclidean counterparts in quantitative accuracy while pro-024 viding additional interpretability. Qualitative visualizations reveal that generic prototypes capture consistent, semanti-025 cally important distinctions between classes. In contrast, 026 specific prototypes capture fine-grained variations within 027 028 each class, providing essential nuances for detailed clas-029 sification and enhancing the model's ability to differentiate 030 across multiple levels of abstraction. Our code can be found 031

032 1. Introduction

Prototypical neural networks have emerged as a prominent
field in computer vision due to their ability to make
interpretable and explainable decisions. They hold an
advantage over traditional black-box models as they can
provide the reasoning behind their results without reliance



Figure 1. Organization of hierarchical prototypes in hyperbolic space, where generic prototypes that capture defining consistent features, entail specific prototypes that capture intra-class variations such as color, texture, and orientation.

on post-hoc explainability methods, leading to more trustworthy decision making. Prototypical neural networks such as ProtoPNet [1], ProtoPNet variants [3, 27] and PIP-Net [22] have shown success in a variety of image classification tasks.

A common shortcoming of many prototype-based mod-044 els is their tendency to treat image features independently, 045 without considering the hierarchical relationships between 046 them. This limitation can result in suboptimal perfor-047 mance on tasks requiring a deeper understanding of im-048 age structure, such as object recognition, scene under-049 standing, and fine-grained classification. In real-world im-050 ages, features are often organized hierarchically, with more 051 generic features capturing consistent and defining charac-052 teristics, while more specific features reflect variable de-053 tails. Generic features are particularly effective for distin-054 guishing broad categories due to their relevance and distinc-055 tiveness across classes. However, when differentiating be-056 tween visually similar classes—such as distinguishing be-057 tween dog breeds-specific features become more valuable, 058

1

108

as they capture fine-grained variations necessary for precise
classification. Therefore, by neglecting this hierarchy, existing prototypical models may miss critical context needed
for accurate predictions.

As illustrated in Fig. 1, organizing prototypes hierarchi-063 cally allows models to capture both defining and variable 064 class features in a structured manner. Generic prototypes fo-065 cus on fundamental and consistent traits that reliably char-066 acterize a class-for instance, in animal images, they may 067 capture the typical shape of a nose, a paw, or an ear that 068 remains stable across variations. In contrast, specific proto-069 070 types highlight within-class diversity by attending to more variable aspects, such as different regions of the body (e.g., 071 a tail instead of a nose) or feature combinations that vary in 072 color, texture, or orientation. By structuring prototypes this 073 way, the model can differentiate classes based on both their 074 defining traits and the finer variations that help distinguish 075 076 similar-looking instances.

To address this limitation, we introduce HAPPI, a 077 Hierarchical And PrototyPical Image recognition method-078 ology, which learns prototypes in a hyperbolic space us-079 ing a Lorentzian manifold instead of Euclidean space. Our 080 method structures the embedding space so that generic 081 features-representing stable and defining characteristics 082 of a class-are positioned near the hyperboloid origin, 083 while specific features-capturing variations within the 084 class-are positioned farther away. This arrangement en-085 ables the model to learn generic prototypes that focus on 086 087 fundamental, class-defining traits, while specific prototypes capture within-class diversity, such as differences in body 088 regions, color patterns, or pose variations. By explicitly 089 modeling this hierarchy, our approach enhances the model's 090 091 ability to differentiate between classes while also capturing nuanced intra-class variations. While demonstrated on sev-092 093 eral backbones, our approach has the potential to generalize to other prototypical neural networks and image classifiers. 094 Our contributions are as follows: 095

We propose a model-agnostic framework to optimize prototypical neural networks on the Lorentzian manifold in hyperbolic space. We explicitly learn generic and specific prototypes while ensuring hierarchical consistency among prototypes of each class. These prototypes are structured and learned within the hyperboloid to preserve class relationships.

Extensive experiments across various prototypical architectures and datasets demonstrate that our hyperbolic approach consistently outperforms existing Euclidean-based optimization methods.

2. Related Works

2.1. Prototypical Neural Networks

Prototypical neural networks are inherently interpretable 109 models due to their structure. These models evaluate in-110 puts explicitly based on their similarity to learned discrim-111 inative features, or "prototypes", for each class. By visual-112 izing the appearance and location of prototypes, as well as 113 the coefficients relating prototype similarity to the class log-114 its, users can validate the model's decision-making process. 115 ProtoPNet [1] introduced the approach of classification us-116 ing representative prototypes from each of the classes in the 117 training set and providing visual validation by highlighting 118 feature proximity and localization within input images. 119

Subsequent work has focused on enhancing inter-120 pretability and prototype efficiency. ProtoPool [28] and 121 ProtoPShare [27] reduce the total number of prototypes 122 by sharing them across classes. PIP-Net [22] adopts a 123 self-supervised approach to produce sparse prototypes and 124 human-understandable features, while XProtoNet [13] ex-125 tends interpretability by allowing different prototype sizes, 126 applying prototypical neural networks successfully to radi-127 ology datasets. SPANet [32] further improves explainability 128 by combining part prototypes with semantic concepts, pro-129 viding clearer interpretations of what each prototype rep-130 resents. These approaches enhance interpretability by ei-131 ther structuring prototypes more efficiently or aligning them 132 with semantically meaningful concepts, supporting more 133 intuitive and transparent decision-making. 134

Other approaches have modified the learning formula-135 tion to improve performance. TesNet [35] employs a Grass-136 mann manifold to create distinct class subspaces, and ST-137 ProtoPNet [34] introduces an SVM-like method to learn 138 boundary-supporting prototypes. Prototypical networks 139 have also been integrated into transformer architectures [4], 140 as seen in ProtoPFormer [37], ProtoFormer [5], and most 141 recently, ProtoViT [16]. ProtoTree [21] uses a decision-142 tree structure to reduce the number of prototype compar-143 isons, enhancing classification efficiency. Recently, ProtoP-144 NeXt [36] demonstrated that cosine similarity and Bayesian 145 tuning could improve ProtoPNet's transparency and perfor-146 mance across various architectures and classification meth-147 ods. 148

Beyond interpretability and architectural improve-149 ments, recent works have explored hierarchical represen-150 tations within prototypical networks. MCPNet [33] and 151 HPDR [10] introduce hierarchical prototypes, albeit in dif-152 ferent ways. MCPNet learns hierarchical representations 153 through multi-scale prototypes, while HPDR refines fea-154 ture distributions using hierarchical prototypes in hyper-155 bolic space. These methods leverage hierarchy to better 156 capture structural relationships within data, enhancing rep-157 resentation learning and classification performance. 158

237

257

159 2.2. Hierarchical Representations

Image datasets inherently exhibit diverse types of hierar-160 161 chical relationships. One such hierarchy involves clearambiguous relationships, where clear images are associated 162 with specific classes, while ambiguous images (e.g., blurred 163 or occluded) are more generic and tend to exhibit features 164 that overlap with multiple classes. Another type of hier-165 archy is the whole-part relationship, where a global view 166 167 captures whole objects or scenes, while local views focus on finer details of objects or scene parts. Each global view 168 is associated with multiple local views, allowing local fea-169 tures to be contextualized within the broader scope of the 170 image (e.g. a leaf is part of a tree and a tree is part of a 171 park). Recognizing these hierarchical relationships is es-172 173 sential for understanding the intricate relationships within complex visual objects or scenes [12, 24]. 174

Prior works have proposed explicit methods to represent 175 hierarchical concepts within images. For example, Pyra-176 midCLIP [6] captured multi-granular image representations 177 178 by learning at global (whole image), intermediate (large image patch), and local levels (cropped object images). This 179 approach modeled relationships across these levels using 180 self-supervised learning guided by language, resulting in 181 182 more separable representation spaces and improved zeroshot classification. 183

184 2.3. Hyperbolic Learning

Hyperbolic learning has demonstrated strong potential for 185 capturing hierarchical structures within data [18]. Unlike 186 187 Euclidean space, which has a flat manifold and struggles to represent hierarchical relationships effectively, hyperbolic 188 189 manifolds offer a curved geometry that naturally preserves these relationships. It can accurately reflect the tree-like 190 structure of hierarchical data, where leaf nodes are more 191 specialized and represent local features and intermediate 192 193 nodes represent higher-level features. There are two main approaches to learning visual representations in hyperbolic 194 195 space: the Poincaré model [11, 12, 17] and the Lorentzian model. The Poincaré model represents hyperbolic space as 196 the interior of a disk (in 2D) or a ball (in higher dimensions) 197 within Euclidean space. Distances and angles in this model 198 are distorted to reflect the negative curvature of hyperbolic 199 200 geometry. However, the Poincaré approach is challenging to optimize due to its susceptibility to vanishing gradient 201 problems, which hinder the model's convergence [20]. 202

The Lorentzian model is a more recently proposed 203 204 hyperbolic representation approach that is less prone to vanishing gradient problems, and therefore much eas-205 ier to optimize [20]. Empirically, approaches using the 206 Lorentzian model have superior performance to those us-207 ing the Poincaré model [15, 23]. Prior works show that the 208 Lorentzian model leads to a better latent space with hierar-209 210 chical representations and improvement in image retrieval

and classification tasks [2, 26]. Desai et al. [2] used the 211 Lorentz model to enhance vision-language representation 212 learning. Their proposed approach, MERU, learned a con-213 trastive model (i.e. CLIP) in hyperbolic space where hierar-214 chy in visual and language concepts was preserved. MERU 215 adapted the contrastive loss to minimize the proximity be-216 tween the associated image and text using Lorentzian dis-217 tance instead of cosine similarity. Additionally, MERU pro-218 posed an entailment loss that enables the model to learn that 219 text represents more abstract, generic concepts than corre-220 sponding images (e.g., text of 'dog' encompasses many dog 221 images). The entailment loss pushes image embeddings re-222 lated to a specific text to exist within a cone-shaped space 223 emanating from the text embedding, indicating that the text 224 entails a set of associated images. Inspired by MERU, 225 our work, HAPPI, extends these ideas to prototypical neu-226 ral networks, tailoring the embedding space and loss func-227 tions for prototype-based learning. While prior work has 228 explored hyperbolic prototypes, these approaches typically 229 compute prototypes as the mean of class embeddings in hy-230 perbolic space, rather than explicitly learning class-specific 231 prototypes [7, 8, 12, 14]. In contrast, HAPPI learns dis-232 tinct, optimizable prototypes directly within the hyperbolic 233 space, ensuring they capture both generic and specific class 234 features beyond simple class averages. 235

3. Method

3.1. Prototypical Classification Framework

Prototypical methods classify inputs by measuring their 238 similarity to class-specific prototypes. Given an input x, 239 a feature extractor $f(\cdot)$ generates a feature map $f(x) \in$ 240 $\mathbb{R}^{H \times W \times D}$, where H and W are spatial dimensions and 241 D is feature depth. These features are then processed for 242 comparison with prototypes. Let $V_{\text{euc}} \in \mathbb{R}^{\vec{D}}$ denote each 243 extracted feature in Euclidean space to be compared with 244 prototypes. Prototypes are denoted as $p_{c,k} \in \mathbb{R}^D$, where 245 c and k represent the class and the index of the prototype 246 within that class, respectively, with typically K prototypes 247 per class. The comparison between V_{euc} and $p_{c,k}$ yields 248 similarity scores $s(V_{euc}, p_{c,k})$, usually computed using Eu-249 clidean distance or cosine similarity. These scores are then 250 fed into a fully connected layer $FC : \mathbb{R}^{K \times C} \to \mathbb{R}^C$, pro-251 ducing class prediction scores $\hat{y} \in \mathbb{R}^C$, where C is the num-252 ber of classes. This prototype-based approach enhances 253 interpretability by basing classifications on similarities to 254 known examples, facilitating clearer pattern identification 255 in the data. 256

3.2. Transition to Hyperbolic Space

Our method extends traditional prototypical networks by
projecting the prototypes and extracted features into hyper-
bolic space, specifically using the Lorentz model of the hy-258
259

290

291

293

294

295

296

306



Figure 2. An overview of our proposed approach. We use a given prototypical network's feature encoder to extract both generic and specific features, then lift those features to the hyperboloid using exponential mapping. The lifted features are then compared to their respective prototypes, and similarity scores are generated based on their distances to the prototypes. These similarity scores ultimately result in the activations that are averaged to create the class logits.

261 perboloid. In this model, the V_{euc} is mapped to a (D + 1)-262 dimensional hyperbolic space by adding an additional time dimension. This mapping results in a hyperbolic feature 263 $V_h = [V_{\text{space}}, V_{\text{time}}]$, where $V_{\text{space}} \in \mathbb{R}^D$ represents the spa-264 tial components and $V_{\text{time}} \in \mathbb{R}$ is the time dimension [19]. 265 This representation can capture hierarchical relationships 266 267 and varying levels of detail more effectively than Euclidean space. 268

To obtain the V_{space} , a mapping operation, shown as $\mathcal{H}(.)$ in Figure 2, projects Euclidean space features V_{euc} onto the hyperboloid. In our case, we use the simplifying assumption from [2] which considers only the exponential map centered at the origin. Under these conditions, the map simplifies to:

$$V_{\text{space}} = \frac{\sinh\left(\sqrt{c}\|V_{\text{euc}}\|\right)}{\sqrt{c}\|V_{\text{euc}}\|} V_{\text{euc}}$$
(1)

276 where *c* is the negative curvature of the hyperbolic space. 277 Since the hyperbolic features are constrained to reside on 278 the hyperboloid manifold, the time dimension V_{time} can be 279 derived based on V_{space} as follows:

280
$$V_{\text{time}} = \sqrt{\frac{1}{c} + \|V_{\text{space}}\|^2}$$
 (2)

281As a result, we can effectively map the V_{euc} from Euclidean282space into the hyperboloid by calculating both V_{space} and283 V_{time} . We can then calculate the Lorentzian distance of the284resulting features from the learnable prototypes to identify285the most prominent prototypes and classify the input image.

For example, let $a = [a_{\text{space}}, a_{\text{time}}]$ and $b = [b_{\text{space}}, b_{\text{time}}]$ 286 denote two points on the hyperboloid. The Lorentzian distance $d_L(a, b)$ is defined as: 288

$$d_L(a,b) = \frac{1}{\sqrt{c}} \cosh^{-1}(-c\langle a,b\rangle_L), \qquad (3) \qquad \text{289}$$

where $\langle \cdot, \cdot \rangle_L$ denotes the Lorentzian inner product. This inner product is defined as:

$$\langle a, b \rangle_L = \langle a_{\text{space}}, b_{\text{space}} \rangle - a_{\text{time}} b_{\text{time}},$$
 (4) 292

where $\langle \cdot, \cdot \rangle$ is the standard Euclidean inner product. This distance metric is foundational in aligning the prototypes with features in the embedding space.

3.3. Hyperbolic-Based Losses and Prototypes

In this section, we introduce the adaptation of prototypi-297 cal losses and prototype structures to hyperbolic space. We 298 describe the conversion of Euclidean-based clustering and 299 separation losses to their hyperbolic counterparts and intro-300 duce generic prototypes for capturing defining and consis-301 tent features. We also present an entailment loss inspired by 302 MERU [2] to encourage hierarchical structure in hyperbolic 303 space. Prototype optimization methods (such as clustering 304 and separation losses [1]) are based on hyperbolic distances. 305

3.3.1. Feature Extraction and Prototype Assignment

We introduce two types of prototypes per class: generic prototypes, which capture defining class features, and specific 308

364

374

382

390

391

392

prototypes, which capture fine-grained variations such as
 pose, color, or texture. Generic features are explicitly ex tracted based on the backbone architecture:

- Transformer-based: We use the [CLS] token to extract generic features, following ProtoPFormer [37].
- CNN-based: Attention maps, interpreted as occurrence
 maps, are applied to extract generic features, inspired by
 XProtoNet [13].

In contrast, specific features naturally emerge from the host
architecture's feature extraction pipeline and are used directly without additional processing.

320 For both feature types, we dedicate corresponding prototypes: generic prototypes for generic features and specific 321 prototypes for specific features. Once extracted, we utilize 322 323 Equations 1 and 2 to project both generic (V_{euc-g}) and spe-324 cific (V_{euc-s}) features to the hyperbolic space, forming V_{h-g} 325 and V_{h-s} , respectively. Their corresponding prototypes are 326 also projected, resulting in hyperbolic generic prototypes P_{h-g} and hyperbolic specific prototypes P_{h-s} . 327

328 3.3.2. Prototype-based Classification

Prototype-based classification leverages both generic and 329 specific prototypes to enhance predictive accuracy. 330 We compute the similarities of extracted features to both sets 331 332 of prototypes. To obtain the final class activation logits, we 333 calculate two separate sets of logits based on these similarities. Each set of similarities is passed through a distinct 334 335 final layer to produce corresponding logits: z_s for specific prototypes and z_g for generic prototypes. The final activa-336 tion is a joint decision that is computed as the average of 337 these two sets of logits. This approach allows the model to 338 339 integrate both generic and specific information for classification, improving its ability to perform effectively across 340 varying levels of granularity. 341

342 3.3.3. Clustering and Separation Losses

Clustering and separation losses are designed to encourage 343 the learning of effective prototypes that capture the struc-344 ture of the data. Clustering loss promotes intra-class com-345 346 pactness by encouraging features to be close to the proto-347 types of their corresponding classes, while separation loss 348 enhances inter-class separability by pushing features away from prototypes of different classes [1]. Similar to exist-349 350 ing prototypical approaches, we cluster extracted specific 351 and generic features that are mapped to the hyperboloid, V_h, with the prototypes of the corresponding classes. Like-352 353 wise, we separate the lifted features from the prototypes of different classes, using the following equations: 354

$$\mathcal{L}_{\text{clst}} = \min_{i} d_L(V_t^i, P_i^y), \quad \mathcal{L}_{\text{sep}} = -\mathbb{E}_{c \neq y, i} \left[d_L(V_t), P_i^c \right] \right],$$
(5)

where V_t represents generic (g) or specific (s) feature embeddings, and P_i^c denotes class-specific prototypes. Both clustering and separation criteria, denoted for specific prototypes as $L_{clst}^{specific}$ and $L_{sep}^{specific}$, and generically as $L_{cluster}^{generic}$ and $L_{sep}^{generic}$ respectively, utilize the Lorentzian distance, Equation 3, as $d_L(V_h, p)$ where p represents the prototype of interest. This adaptation maintains clustering and separation fidelity in hyperbolic space. 363

3.3.4. Entailment Loss

Entailment loss is designed to encourage a hierarchical 365 structure in the hyperbolic space by enforcing relationships 366 between specific and generic prototypes. It ensures that 367 each specific prototype is entailed by at least one generic 368 prototype of its class, creating a coherent hierarchy of rep-369 resentations. To encourage this hierarchical encoding, we 370 implement an entailment loss based on MERU's formula-371 tion [2]. This is achieved by minimizing the exterior angle 372 of $ext(V_{h-g}, V_{h-s})$ relative to the half-aperture $aper(V_{h-g})$: 373

$$L_{\text{entail}}(V_{\text{h-g}}, V_{\text{h-s}}) = \max(0, \text{ext}(V_{\text{h-g}}, V_{\text{h-s}}) - \operatorname{aper}(V_{\text{h-g}})).$$
(6)

The half-aperture of a cone for entailment relationships is calculated as: 376

$$\operatorname{aper}(a) = \sin^{-1}\left(\frac{2R}{\sqrt{c}\|a_{\operatorname{space}}\|}\right),\tag{7}$$

where R is a constant that controls the boundary conditions378near the origin. Furthermore, the exterior angle between379two points on the hyperboloid is given by:380

$$\operatorname{ext}(a,b) = \cos^{-1}\left(\frac{b_{\operatorname{time}} + a_{\operatorname{time}} \cdot c\langle a, b \rangle_L}{\|a_{\operatorname{space}}\|\sqrt{(c\langle a, b \rangle_L)^2 - 1}}\right).$$
 (8) 381

3.3.5. Combined Loss Function

The overall loss function \mathcal{L} uses cross entropy, \mathcal{L}_{CE} , for classification and integrates generic and specific clustering and separation losses, entailment loss as shown below: 385

$$\mathcal{L} = \mathcal{L}_{CE} + \lambda_{clst_g} \mathcal{L}_{clst}^{generic} + \lambda_{sep_g} \mathcal{L}_{sep}^{generic} + \lambda_{clst_s} \mathcal{L}_{clst}^{specific} + \lambda_{sep_s} \mathcal{L}_{sep}^{specific} + \lambda_{entail} \mathcal{L}_{entail}$$
(9) 386

where each λ denotes the coefficient corresponding to each 387 loss value. 388

4. Experiments 389

4.1. Datasets

We used two different datasets for our evaluations: CUB-200-2011 and Oxford-IIIT Pets.

CUB-200-2011 (Caltech-UCSD Birds 200) [31]: The CUB-200-2011 dataset is a widely-used benchmark for fine-grained classification tasks, focusing specifically on bird species recognition. It consists of 11,788 images across 200 bird species, each with a varying number of 397

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

images. The dataset is split into 5,994 images for trainingand 5,794 images for testing.

Oxford-IIIT Pets [25]: The Oxford-IIIT Pets dataset is used for pet breed classification and includes a balanced mix of cat and dog breeds. It contains 7,349 images spanning 37 breeds, with around 200 images per breed. The dataset is split into 3,680 images for training and 3,669 images for testing.

While each of these datasets has additional attribute annotations such as bounding boxes or segmentations, we
only used the labels to train our models. We selected
these datasets for their frequent use in prototypical neural
network research and their diversity: CUB includes 200
classes, while Pets includes 37 classes.

412 4.2. Implementation Details

413 To ensure a comprehensive evaluation, we selected a range of backbone architectures that cover diverse structural ap-414 415 proaches. Specifically, we included part-based prototypical neural networks like ProtoPNet and ST-ProtoPNet, 416 which use patch-based prototypes; XProtoNet, which uses 417 adaptable prototype sizes; and ProtoPFormer, which uses 418 transformer-based backbone. For CNN-based evaluations, 419 all models were tested with a ResNet-50 backbone [9], 420 421 while different transformer-based models were also evaluated to assess performance across architectural types. Ad-422 ditionally, a black-box version of the backbone classifica-423 tion model was evaluated for comparison, as shown in the **424** results section. We also evaluated MCPNet [33], a state-of-425 the-art hierarchical prototypical model in Euclidean space, 426 and PIPNet [22], a state-of-the-art prototypical model. All 427 evaluations were performed using 10 specific prototypes per 428 429 class, a standard choice in prototypical neural network literature, along with 2 generic prototypes per class, determined 430 empirically. The entailment loss coefficient, λ_{entail} , was set 431 432 to 0.1 and R was set to 0.1. The loss weights, λ , for the 433 remaining components were selected to match the values used in the baseline model. Input images were resized to 434 224x224, with shear and flip transformations for data aug-435 mentation, while test images were resized to 224x224 with-436 out additional cropping to maintain consistency. 437

438 For ProtoPFormer, we adopted their approach of using 439 the CLS token as a generic prototype and image tokens as 440 specific prototypes. Instead of ProtoPFormer's prototypical part concentration (PPC) loss, we implemented our cluster-441 ing and separation loss functions. While PPC loss focuses 442 443 on concentrating prototypes on distinct, centralized representative parts for each class, our clustering and separation 444 losses also enforce the learning of distinct and represen-445 tative prototypes, with a formulation that translates more 446 effectively to hyperbolic space. For ProtoPNet and XPro-447 toNet, we report two accuracy metrics: end-to-end (E2E) 448 449 training accuracy where the prototypes and features are optimized jointly, and multi-stage training accuracy. Multi-
stage training consists of initial prototype warm-up epochs,
followed by joint training, and then final layer optimiza-
tion for additional epochs. To ensure a fair comparison, we
retained the original backbone feature learning rates from
each model's source paper. Further implementation details
can be found in the supplementary materials.450
451
452

4.3. Quantitative Results

Our method consistently delivers strong quantitative performance across prototypical architectures, often surpassing or matching state-of-the-art baselines. This highlights the effectiveness of our hierarchical approach in enhancing feature representation.

As shown in Table 1, HAPPI significantly improves performance over simpler prototypical architectures like ProtoPNet. On CUB-200-2011, ProtoPNet achieves 45.24% accuracy, whereas integrating HAPPI boosts it to 61.44%, a notable 16.2% improvement. Similarly, on Oxford-IIIT Pets, hyperbolic ProtoPNet achieves 88.25%, up from the baseline's 54.65%, an increase of over 33.6%. These results demonstrate HAPPI's ability to mitigate the limitations of simpler methods by capturing complex hierarchical relationships.

For more advanced architectures such as XProtoNet 473 and ST-ProtoPNet, which already incorporate sophisticated 474 feature representations, performance gains are smaller. 475 On CUB-200-2011, XProtoNet improves from 73.82% to 476 75.46% with HAPPI, while ST-ProtoPNet sees a marginal 477 increase from 86.54% to 87.21%. A similar trend is ob-478 served for transformer-based models like ProtoPFormer 479 with DeiT-Ti and DeiT-S backbones, where improvements 480 range between 0.5% and 4.0%. This suggests that while 481 HAPPI enhances all architectures, its impact is most pro-482 nounced in simpler baselines, where hierarchical model-483 ing in hyperbolic space provides the greatest benefit. This 484 behavior likely reflects a saturation of feature representa-485 tion in advanced architectures, where additional hierarchi-486 cal structuring yields diminishing returns in accuracy. How-487 ever, even for these high-performing models, HAPPI still 488 enhances interpretability by structuring feature representa-489 tions hierarchically, ensuring that both generic and specific 490 prototypes contribute meaningfully to decision-making. 491

Despite its consistent advantages, HAPPI does not al-492 ways yield improvements. In some configurations, such as 493 ProtoPFormer with a DeiT-Ti backbone on cropped CUB-494 200-2011, accuracy remains unchanged at 79.65%. A likely 495 explanation is that DeiT-Ti, being the smallest model tested, 496 lacks the capacity to fully leverage hyperbolic represen-497 tations, limiting performance gains. This suggests that 498 HAPPI's effectiveness depends on model complexity and 499 dataset-specific optimization strategies. Future work could 500 explore adaptive hyperbolic scaling for smaller backbones 501

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

	Methods	CUB		CUB Cropped		Pets	
Backbone		Baseline (Euclidean)	HAPPI (Hyper- bolic)	Baseline (Euclidean)	HAPPI (Hyper- bolic)	Baseline (Euclidean)	HAPPI (Hyper- bolic)
ResNet50 [9]	Black-box	75.16	-	76.34	-	90.49	-
	PIP-Net [22]	73.52	-	82.00	-	88.50	-
	MCPNet [33]	74.28	-	70.78	-	90.32	-
	ProtoPNet [1]	45.24	61.44	62.30	80.51	54.65	88.25
	ProtoPNet - E2E	67.04	76.44	71.53	82.39	77.35	88.34
	XProtoNet [13]	73.82	75.46	80.84	81.03	89.48	89.40
	XProtoNet - E2E	75.87	77.91	82.09	82.90	90.00	91.50
	ST-ProtoPNet [34]	86.54	87.21	86.94	86.92	77.43	81.93
DeiT-Ti [29]	Black-box	75.27	-	78.13	-	88.85	-
	ProtoPFormer [37]	77.30	76.18	79.65	79.65	86.89	86.89
DeiT-S [29]	Black-box	79.86	-	82.53	-	92.23	-
	ProtoPFormer [37]	77.68	78.67	80.60	84.59	88.72	91.25
CaiT-XXS-24 [30]	Black-box	80.76	-	82.56	-	92.86	-
	ProtoPFormer [37]	80.86	80.81	82.24	83.91	91.47	91.88

Table 1. Classification accuracy (%) on full / cropped CUB-200-2011 and Oxford-IIIT Pets using various backbones and methods.

502 or refined tuning strategies to address these limitations.

Compared to black-box baselines, HAPPI achieves 503 504 competitive accuracy while significantly improving interpretability. On CUB-200-2011, the ResNet-50 black-box 505 model attains 75.16% accuracy, whereas HAPPI-enhanced 506 ProtoPNet reaches 76.44%, a modest improvement. How-507 ever, this gain comes with the added benefit of interpretabil-508 ity, as HAPPI's generic and specific prototypes provide 509 insight into the model's decision-making process. Simi-510 larly, on Oxford-IIIT Pets, hyperbolic XProtoNet achieves 511 512 91.50% accuracy, surpassing the ResNet-50 black-box model's 90.49%. These results highlight HAPPI's ability 513 to balance high performance with transparency, offering a 514 515 compelling alternative to opaque black-box models.

516 4.4. Qualitative Results

517 To illustrate the qualitative effectiveness of HAPPI, we provide visualizations of generic and specific prototypes in 518 Figure 3, learned using the XProtoNet-E2E architecture on 519 the Oxford-IIIT Pets [25] dataset. The figure includes four 520 521 classes: Abyssinian (top left), Bengal (bottom left), American Bulldog (top right), and Basset Hound (bottom right). 522 523 Each group has two rows—the bottom row shows a heatmap overlay highlighting attended regions, while the top row 524 masks less relevant areas by shadowing them. 525

526 Generic prototypes capture consistent features that 527 broadly distinguish classes. Across all classes, they focus 528 on the nose and mouth area, a key trait distinguishing not 529 only between cats and dogs but also among different breeds 530 within each group. The American Bulldog's generic proto-531 types highlight facial wrinkles and white fur, while the Bas-532 set Hound's emphasize long ears and a brown head. The Abyssinian cat's prototypes focus on fur texture, nose, and eyes, while the Bengal cat's emphasize its distinct tiger-like fur pattern.

In contrast, specific prototypes capture variations in body shapes, poses, and colors. The American Bulldog's specific prototypes highlight a different color combination—white and brown—unlike its generic prototypes, which focus only on white fur. In the last specific prototypes of both the American Bulldog and Bengal cat, a diffused heatmap or less pronounced shadowing suggests a broader focus, capturing some background details.

Overall, generic prototypes identify prominent classdefining features, while specific prototypes capture finegrained variations, ensuring a hierarchical, multi-level understanding of class distinctions. This layered approach enhances the model's capacity to differentiate classes across various levels of abstraction.

Specific prototypes, in contrast, provide complementary information that captures specific body shapes, poses, and angles, adding variable, fine-grained detail that further supports classification. This layered approach, where generic prototypes identify defining distinctions and specific prototypes capture detailed variations, enhances the model's capacity to differentiate classes across multiple levels of abstraction.

4.5. Ablation Study

We performed an ablation study using XProtoNet-E2E to assess the impact of hyperbolic space and hierarchical structuring. As shown in Table 2, transitioning to hyperbolic space and introducing generic and specific prototypes improved performance over the Euclidean base-563

ICCV-BEW 2025 Submission #11. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

605



Figure 3. Qualitative visualization of learned prototypes in the XProtoNet-E2E architecture for the Oxford-IIIT Pets dataset [25]. Generic prototypes capture defining class features, such as a 'white American Bulldog' (top right), while specific prototypes capture intra-class variations like color differences. As prototype distance from the hyperboloid root increases, it covers broader areas or more variations.

line, even without entailment loss ($\lambda_{entail} = 0$). In-564 creasing λ_{entail} slightly reduced accuracy, highlighting a 565 tradeoff between classification performance and hierarchi-566 567 cal structuring, consistent with observations in [2]. Larger 568 λ_{entail} values enhanced structure but posed optimization challenges. With $\lambda_{entail} = 0.1$, the model maintained high 569 accuracy while preserving a meaningful hierarchy, making 570 it the preferred setting. 571

572 To further analyze the effect of prototype allocation, we 573 varied the number of generic (K_g) and specific (K_s) prototypes per class, as shown in Table 3. Increasing the number 574 of prototypes generally improved performance, but gains 575 plateaued around $K_g = 2$ and $K_s = 5$. Notably, models 576 with more specific prototypes performed better than those 577 with more generic ones, reinforcing the importance of cap-578 579 turing intra-class variability. The setup with $K_g = 5$ and $K_s = 1$ performed worse than $K_g = 1$ and $K_s = 5$, in-580 dicating that generic prototypes alone are not sufficient for 581 fine-grained classification. When specific prototypes were 582 583 sufficient, reducing the number of generic prototypes from 5 to 2 further improved optimization stability. These re-584 585 sults suggest that the ideal prototype configuration depends on the complexity of the dataset. A careful balance be-586 tween generic and specific prototypes is crucial, as different 587 588 datasets may require different levels of hierarchical struc-589 turing to achieve optimal performance.

590 5. Conclusion

In this paper, we introduced HAPPI, a model-agnostic approach for adapting prototypical networks to hyperbolic
space to effectively learn hierarchical image representations. Our method showed quantitative improvements over
single-scale Euclidean baselines across diverse datasets,

Mathad	λ_{entail}					
Method	0	0.1	0.2	0.5		
Euclidean	90.00	-	-	-		
Hyperbolic	91.50	91.50	91.47	91.20		

Table 2. Ablation study, using XProtoNet-E2E, on the use of hyperbolic space and entailment in terms of Accuracy (%) on the Pets dataset. The first row refers to the Euclidean model which has 10 prototypes per class, while the other rows are hyperbolic versions with 2 generic and 10 specific prototypes.

		K_s			
		1	2	5	10
K_g	1	88.83	90.27	91.06	91.06
	2	90.19	90.90	91.91	91.50
	5	90.79	88.53	91.77	91.52

Table 3. Accuracy (%) by number of generic and specific prototypes on the Pets dataset.

highlighting its ability to capture complex, structured rela-596 tionships within the data. Furthermore, qualitative analyses 597 reveal that HAPPI provides an interpretable understanding **598** of image hierarchies, allowing insights into the features the 599 model attends to when making predictions. Our findings 600 suggest that hyperbolic prototypical networks, as exempli-601 fied by HAPPI, hold significant potential for applications 602 requiring both hierarchical understanding and interpretabil-603 ity in visual tasks. 604

References

 Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. This looks like that: deep learn 606
 607

609

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

716

ing for interpretable image recognition. Advances in neural information processing systems, 32, 2019. 1, 2, 4, 5, 7

- [2] Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin
 Johnson, and Ramakrishna Vedantam. Hyperbolic imagetext representations. *arXiv (Cornell University)*, 2023. 3, 4,
 5, 8
- [3] Jon Donnelly, Alina Jade Barnett, and Chaofan Chen. Deformable protopnet: An interpretable image classifier using
 deformable prototypes. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
 pages 10265–10275, 2022. 1
- [4] Alexey Dosovitskiy. An image is worth 16x16 words:
 Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 2
- 622 [5] Ashkan Farhangi, Ning Sui, Nan Hua, Haiyan Bai, Arthur
 623 Huang, and Zhishan Guo. Protoformer: Embedding proto624 types for transformers. In *Pacific-Asia Conference on Knowl-*625 *edge Discovery and Data Mining*, pages 447–458. Springer,
 626 2022. 2
- 627 [6] Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, Ron628 grong Ji, and Chunhua Shen. Pyramidclip: Hierarchical fea629 ture alignment for vision-language model pretraining. *arXiv*630 (*Cornell University*), 2022. 3
- [7] Mina Ghadimi Atigh et al. Hyperbolic busemann learning
 with ideal prototypes. *Advances in neural information pro- cessing systems*, 34:103–115, 2021. 3
- [8] Yunhui Guo, Xudong Wang, Yubei Chen, and Stella X Yu.
 Clipped hyperbolic classifiers are super-hyperbolic classifiers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11–20, 2022. 3
- 638 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun.
 639 Deep residual learning for image recognition. In *Proceed-*640 *ings of the IEEE conference on computer vision and pattern*641 *recognition*, pages 770–778, 2016. 6, 7
- [10] Chengyang Hu, Ke-Yue Zhang, Taiping Yao, Shouhong
 Ding, and Lizhuang Ma. Rethinking generalizable face
 anti-spoofing via hierarchical prototype-guided distribution
 refinement in hyperbolic space. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1032–1041, 2024. 2
- [11] Chengyang Hu, Ke-Yue Zhang, Taiping Yao, Shouhong Ding, and Lizhuang Ma. Rethinking generalizable face anti-spoofing via hierarchical prototype-guided distribution refinement in hyperbolic space. In 2024 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 1032–1041. IEEE, 2024. 3
- [12] Valentin Khrulkov, Leyla Mirvakhabova, Evgeniya Ustinova, Ivan Oseledets, and Victor Lempitsky. Hyperbolic
 image embeddings. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages
 658 6418–6428, 2020. 3
- [13] Eunji Kim, Siwon Kim, Minji Seo, and Sungroh Yoon. Xprotonet: diagnosis in chest radiography with global and local explanations. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15719– 15728, 2021. 2, 5, 7

- [14] Sungyeon Kim, Boseung Jeong, and Suha Kwak. Hier: Metric learning beyond class labels via hierarchical regularization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 19903–19912, 2023. 3
- [15] Marc Law, Renjie Liao, Jake Snell, and Richard Zemel. Lorentzian distance learning for hyperbolic representations. In *Proceedings of the 36th International Conference on Machine Learning*, pages 3672–3681. PMLR, 2019. 3
- [16] Chiyu Ma, Jon Donnelly, Wenjun Liu, Soroush Vosoughi, Cynthia Rudin, and Chaofan Chen. Interpretable image classification with adaptive prototype-based vision transformers. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. 2
- [17] Paolo Mandica, Luca Franco, Konstantinos Kallidromitis, Suzanne Petryk, and Fabio Galasso. Hyperbolic learning with multimodal large language models. *arXiv preprint arXiv:2408.05097*, 2024. 3
- [18] Pascal Mettes, Mina Ghadimi Atigh, Martin Keller-Ressel, Jeffrey Gu, and Serena Yeung. Hyperbolic deep learning in computer vision: A survey. *International Journal of Computer Vision*, 132(9):3484–3508, 2024. 3
- [19] Hermann Minkowski. Raum und Zeit. *Physikalische Zeitschrift*, 1908. 4
- [20] Gal Mishne, Zhengchao Wan, Yusu Wang, and Sheng Yang. The numerical stability of hyperbolic representation learning. arXiv (Cornell University), 2022. 3
- [21] Meike Nauta, Ron Van Bree, and Christin Seifert. Neural prototype trees for interpretable fine-grained image recognition. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 14933–14943, 2021. 2
- [22] Meike Nauta, Jörg Schlötterer, Maurice Van Keulen, and Christin Seifert. Pip-net: Patch-based intuitive prototypes for interpretable image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2744–2753, 2023. 1, 2, 6, 7
- [23] Maximilian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. *arXiv (Cornell University)*, 2018. 3
- [24] Avik Pal, Max van Spengler, Guido Maria D'Amely di Melendugno, Alessandro Flaborea, Fabio Galasso, and Pascal Mettes. Compositional entailment learning for hyperbolic vision-language models. arXiv preprint arXiv:2410.06912, 2024. 3
- [25] Omkar M Parkhi, Andrea Vedaldi, Andrew Zisserman, and CV Jawahar. Cats and dogs. In 2012 IEEE conference on computer vision and pattern recognition, pages 3498–3505. IEEE, 2012. 6, 7, 8
- [26] Sameera Ramasinghe, Violetta Shevchenko, Gil Avraham, and Ajanthan Thalaiyasingam. Accept the modality gap: An exploration in the hyperbolic space. pages 27253–27262. IEEE, 2024. 3
- [27] Dawid Rymarczyk, Łukasz Struski, Jacek Tabor, and Bartosz Zieliński. Protopshare: Prototypical parts sharing for similarity discovery in interpretable image classification. In *Proceedings of the 27th ACM SIGKDD Conference on Knowl* 720

- redge Discovery & Data Mining, pages 1420–1430, 2021. 1,
 2
- [28] Dawid Rymarczyk, Łukasz Struski, Michał Górszczak, Koryna Lewandowska, Jacek Tabor, and Bartosz Zieliński. Interpretable image classification with differentiable prototypes assignment. In *European Conference on Computer Vision*, pages 351–368. Springer, 2022. 2
- [29] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco
 Massa, Alexandre Sablayrolles, and Hervé Jégou. Training
 data-efficient image transformers & distillation through attention. In *International conference on machine learning*,
 pages 10347–10357. PMLR, 2021. 7
- [30] Hugo Touvron, Matthieu Cord, Alexandre Sablayrolles,
 Gabriel Synnaeve, and Hervé Jégou. Going deeper with image transformers. In *Proceedings of the IEEE/CVF interna- tional conference on computer vision*, pages 32–42, 2021. 7
- [31] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011
 dataset. 2011. 5
- [32] Qiyang Wan, Ruiping Wang, and Xilin Chen. Interpretable
 object recognition by semantic prototype analysis. In *Proceedings of the IEEE/CVF Winter Conference on Applica-*tions of Computer Vision, pages 800–809, 2024. 2
- [33] Bor-Shiun Wang, Chien-Yi Wang, and Wei-Chen Chiu.
 Mcpnet: An interpretable classifier via multi-level concept prototypes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10885– 10894, 2024. 2, 6, 7
- [34] Chong Wang, Yuyuan Liu, Yuanhong Chen, Fengbei Liu, Yu
 Tian, Davis McCarthy, Helen Frazer, and Gustavo Carneiro.
 Learning support and trivial prototypes for interpretable image classification. In *Proceedings of the IEEE/CVF Inter- national Conference on Computer Vision*, pages 2062–2072,
 2023. 2, 7
- [35] Jiaqi Wang, Huafeng Liu, Xinyue Wang, and Liping Jing.
 Interpretable image recognition by constructing transparent embedding space. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 895–904, 2021.
 2
- [36] Frank Willard, Luke Moffett, Emmanuel Mokel, Jon Donnelly, Stark Guo, Julia Yang, Giyoung Kim, Alina Jade Barnett, and Cynthia Rudin. This looks better than that: Better interpretable models with protopnext. *arXiv preprint arXiv:2406.14675*, 2024. 2
- [37] Mengqi Xue, Qihan Huang, Haofei Zhang, Lechao Cheng, Jie Song, Minghui Wu, and Mingli Song. Protopformer: Concentrating on prototypical parts in vision transformers for interpretable image recognition. *arXiv preprint arXiv:2208.10431*, 2022. 2, 5, 7