047

048

049

050

051

052

053

054

055

056

057

# HierVision: Standardized and Reproducible Hierarchical Sources for Vision Datasets

Anonymous ICCVW submission

### Paper ID 10

#### Abstract

001 One-vs-rest training is a pervasive optimization regime in 002 deep learning, whether the problem is supervised, selfsupervised, or multi-modal in nature. The real world is 003 however not binary, but governed by hierarchies. Hier-004 archies provide key information about the semantic rela-005 006 tion between concepts, about which mistakes to avoid, and about the inherent organization of vision and language it-007 008 self. Hierarchical learning, therefore, has a long history in computer vision and has gained further traction with the 009 010 rise of hyperbolic deep learning. Currently, however, hierarchies are not standardized and centrally organized. In-011 stead, such knowledge is scattered around various reposito-012 ries, with inconsistent formatting, organizations, and avail-013 ability. The lack of a central hub for hierarchies in vision 014 015 datasets harms the utility and reproducibility of hierarchi-016 cal learning. This paper introduces HierVision, a central hub for hierarchical knowledge in vision datasets. This hub 017 contains 60+ hierarchical sources, spanning actions, con-018 cepts, fine-grained categories, vision-language, and more. 019 020 We outline a uniform coding of the hierarchies and procedures to embed them in existing pipelines. With this hub, 021 022 we hope to positively impact the broad use and re-use of hierarchies for deep learning in computer vision. 023

### **024 1. Introduction**

025 Hierarchies are ubiquitous data structures across all sci-026 ences; from lesion taxonomies in the medical domain to animal ontologies in biology and semantic trees in natural lan-027 guage [2, 82]. Such tree-like structures have been used for 028 centuries to organize our data and natural phenomena [1]. 029 Computer vision deals with categorizing concepts from the 030 real world, and datasets are therefore commonly organized 031 032 hierarchically. Consider for example the WordNet hierarchy behind ImageNet [26], the biological ontology of birds 033 in CUB [117], or the tree of verbs in Kinetics [58]. 034

035 Despite the widespread availability of hierarchical infor-

mation for computer vision, such knowledge is typically ig-036 nored when training deep networks. Instead, one-versus-037 rest optimization through cross-entropy and contrastive ob-038 jectives are default options [21, 50]. Such a setup presents a 039 binary view to categorization, where classes are either pos-040 itive or (equally) negative. As a result the standard deep 041 learning setup misses crucial hierarchical information about 042 class similarities. The lack of hierarchical usage negatively 043 impacts learning [53, 93, 126], generalization [93], error 044 severity [8], and more. 045

An important reason for the lack of hierarchical integration in modern deep learning for vision is a geometric mismatch. Deep learning is Euclidean by default. Hierarchies are, however, exponentially growing structures, which leads to distortion when embedding them in Euclidean space [103], as Euclidean volumes grow only polynomially with their radius [88]. Recently, hyperbolic learning has rapidly gained traction in computer vision [80], as a natural space for embedding hierarchies [36, 88, 101, 115] and therefore a natural solution for hierarchical computer vision [3, 5, 28, 57, 74]. As a result, there is a growing demand for hierarchical knowledge in vision datasets.

An important issue currently is that there is no central 058 hub for storing and sharing hierarchies. This does not align 059 with the best scientific practices and hampers research. Not 060 only are hierarchies arbitrarily hard to find depending on the 061 dataset, but they are also not standardized and can even be 062 altered. As such, it is unnecessarily hard to use hierarchies, 063 and reproducibility is low since it is unknown whether hi-064 erarchies are identical. This paper introduces HierVision, a 065 central hub for sharing hierarchies in vision datasets. Our 066 goal is simple: create a continuous effort to store hierar-067 chies for all vision datasets in a single place. Each hierarchy 068 is standardized in a single format for ease of use and repro-069 ducibility. The hub also contains pipelines for visualization, 070 analysis, and integration in deep learning and hyperbolic 071 embedding pipelines. With HierVision, we want to make 072 the community aware of the broad potential of hierarchi-073 cal knowledge and the need for a central hub to organized 074 computer vision hierarchically. 075

133

134

135

136

137

138

139

140

141

142

162

## 076 2. Related Work

#### 077 2.1. Hierarchical Datasets

For clarity, we categorize prominent dataset hierarchies into
two groups-those emphasizing semantic ontologies and
those following biological taxonomies.

081 Semantic hierarchies are typically derived from humandefined knowledge bases or lexical resources. A founda-082 tional example is the use of WordNet [82], a large lexical 083 database of English, to structure the ImageNet dataset [26]. 084 This provided a rich, structured ontology that has been 085 instrumental in the development of deep learning mod-086 087 els. Other prominent datasets include CIFAR-100 [61], PASCAL-VOC [31], and OpenImages [62]. The BREEDS 088 benchmark, derived from ImageNet, explicitly uses the 089 class hierarchy to study robustness [102]. Such semantic 090 structures are not limited to object recognition and extend 091 092 to domains like medical imaging with datasets like CheX-093 pert [109] and scene understanding with ADE20K [133].

The second category of datasets follows formal biolog-094 ical taxonomies, providing a scientifically grounded struc-095 096 ture for fine-grained visual categorization. These datasets 097 are critical for applications in biodiversity and conservation. For example, iNaturalist [114] organizes species observa-098 tions according to the taxonomic rank (kingdom, phylum, 099 class, etc.), ensuring that classes have a hierarchical rela-100 tionship. TreeOfLife-10M and Rare Species [110], Classic 101 fine-grained benchmarks CUB [117] and NABirds [113] are 102 103 built on the taxonomy of species and genera, and datasets 104 like AutoArborist [7] structure tree images by botanical tax-105 onomy.

#### **106 2.2. Hierarchies enhance vision tasks**

The use of hierarchies as a source of prior knowledge is 107 a long-standing concept in computer vision [26, 66, 78]. 108 Classical approaches from the pre-deep learning era explic-109 110 itly modeled the compositional nature of objects and scenes. Part-based approaches, such as pictorial structures [35] and 111 grammar-based models [134], organize objects and scenes 112 into parts to improve image recognition and interpretability. 113 [34, 35, 75, 107, 134]. 114

In the modern deep learning era, hierarchical knowledge 115 116 is integrated through various mechanisms such as class taxonomies, structured loss functions, and specialized archi-117 tectures [8, 9, 38, 85]. These approaches improve many 118 tasks such including image classification [19, 53, 93], ac-119 120 tion classification [43, 74], and robustness to distribution 121 shifts [102]. Hierarchical approaches were also used to 122 measure the severity of classification mistakes, where misclassifying an object as a close relative in the hierarchy is 123 penalized less severely [8, 9, 38, 39]. Furthermore, hierar-124 125 chical information has been successfully incorporated into 126 contrastive learning [12, 13, 46], and vision-language models [37, 90]. The utility of hierarchical methods also ex-<br/>tends to applied domains such as medical imaging [18] and<br/>autonomous driving [83]. A common theme of these works<br/>is their reliance on Euclidean geometry to model these hier-<br/>archical relationships.127<br/>128<br/>129

#### 2.3. Hyperbolic learning

Hyperbolic learning has emerged as a powerful paradigm for encoding and exploiting hierarchical relationships in visual data. Owing to the constant negative curvature of hyperbolic space, it can be thought of as a continuous version of a tree, making it a good choice to accommodates treelike structures while preserving distances [48, 112]. In recent years, hierarchical embeddings have been performed in hyperbolic space [88], leading to successfully embedding complex trees with low distortions [36, 63, 89, 101, 115, 129].

Many computer vision tasks inherently involve hierar-143 chies, for example, semantic grouping or biological tax-144 onomies (Sec. 2.1). A wide range of works have recently 145 shown the potential and effectiveness of using a hyper-146 bolic embedding space for both supervised and unsuper-147 vised learning [80]. Specifically, in supervised settings, the 148 hierarchical prior knowledge of the datasets can be embed-149 ded in hyperbolic space, after which the visual representa-150 tions can be mapped to the same space and optimized to 151 match this hierarchical organization [5, 57, 74]. Hyper-152 bolic embeddings have shown benefits in classification [41, 153 45, 116], segmentation [3, 17], out-of-distribution detec-154 tion [36, 57, 116], uncertainty quantification [3, 17], zero-155 shot learning [4, 51, 69], continual learning [5, 24, 111], 156 hierarchical representation learning [28, 29, 72, 74], con-157 trastive learning [40, 130], generative models [64, 95], and 158 vision-language models [27, 54, 92, 94]. Recently Ay-159 oughi et al. [6] discussed optimial tree structure for hyper-160 bolic embeddings. 161

#### **3. Hierarchies**

While the benefits of hierarchical information in computer 163 vision are well-established, its practical adoption has been 164 limited by fragmented and inconsistent hierarchy manage-165 ment across datasets. Hierarchies exist in various formats, 166 from simple text files and custom XML/JSON to folder-167 based organizations, which require custom parsing for each 168 data set. This fragmentation hinders reproducibility and the 169 development of hierarchical computer vision methods. To 170 address this, we introduce HierVision, a centralized hub that 171 standardizes hierarchy representation and enables seamless 172 integration with existing tools. In the following sections, 173 we first define a common hierarchy format that can be used 174 with graph processing library like NetworkX [47] (Sec-175 tion 3.1). Then we briefly discuss describe our standardiza-176 tion process (Section 3.2. We then discuss in detail the cur-177

242

243

244

245

246

247

248

249

250

251

252

rent datasets (Section 3.3) and finally present some datasethierarchy statistics (Section 3.4).

#### **180 3.1. Hierarchy format**

181 To enable standardized storage and use of hierarchical information across diverse vision datasets, we adopt a uniform 182 graph-based representation format. Each hierarchy is mod-183 eled as a directed, rooted tree encoded in JSON. Formally, 184 we represent a hierarchy as a graph, G = (V, E), where 185 V is the set of nodes (e.g. dataset classes or parents) and 186  $E \in V \times V$  is the set of directed edges, such that an edge 187  $(u, v) \in E$  denotes a parent-child relationship and node v 188 is a subclass of node u. In practice, each hierarchy is stored 189 in a JSON with the following components: 190

- "nodes": A list of nodes, each with an integer "id"
  and a human-readable "label". These define the concepts in the hierarchy.
- "links": A list of directed edges, each represented as a with "source" and "target" keys indicating parent and child node IDs, respectively.
- "directed": A boolean flag, always set to true.
- "multigraph": A boolean flag, always set to false,
  enforcing at most one edge between any pair of nodes.

200 An excerpt of a simplified hierarchy is shown below:

```
201
            ł
202
               "directed": true,
      2
203
               "multigraph": false,
204
      4
               "nodes": [
205
                 {"id": 3, "label": "root"},
      5
206
                 {"id": 2, "label": "animal"},
      6
207
                 {"id": 1, "label": "dog"},
                 {"id": 0, "label": "cat"}
208
      8
209
      9
              1,
210
               "links": [
     10
211
                 {"source": 3, "target": 2},
     11
                 {"source": 2, "target": 0},
212
     12
213
                 {"source": 2, "target": 1}
     13
214
     14
              ]
215
     15
            }
```

Listing 1. A sample JSON format of simple hierarchy. All hierarchies in *Hiervision* follow this structure, allowing easy visualization and use in downstream applications

216 In this example, "root" is the global ancestor of all nodes, "animal" is a direct child of "root", and 217 both "cat" and "dog" are subclasses of "animal". 218 This structure is fully compatible with standard graph-219 220 processing libraries, particularly NetworkX [47], which we 221 use throughout our implementation. Hierarchies can be di-222 rectly loaded as networkx.DiGraph objects, enabling efficient hierarchy traversal, visualization, validation, and 223 integration into hierarchical deep learning pipelines, includ-224 ing those that rely on hyperbolic embeddings or hierarchy-225 226 aware loss functions.

#### 3.2. Standardizing the dataset hierarchies

We collect hierarchies from over 61 vision datasets span-228 ning various domains such as object recognition, fine-229 grained classification, action understanding, scene interpre-230 tation, video analysis, and medical imaging. These hierar-231 chies originate from either the dataset creators themselves 232 or other papers that construct or refine hierarchies of exist-233 ing datasets for specific tasks. Across these sources, we ob-234 serve significant variation in format: some hierarchies are 235 in graph-structured files (e.g., JSON, XML trees). Others 236 use flat lists with indentation or prefix-based identifiers to 237 imply structure. Some are only documented visually or em-238 bedded in figures or tables in papers. These sources vary 239 significantly in format, structure, and accessibility, requir-240 ing a standardization process. 241

To unify these into a standardized format, we add each hierarchy into the JSON graph structure described in Sec. 3.1. We parse and extract the node and edge structure. We resolve any label inconsistences and add a single root node. We validate the resulting graph using NetworkX to ensure it forms a connect, directed tree.

In cases where multiple versions of a hierarchy exist (e.g., fine vs. coarse levels), we store each version explicitly. This ensures that downstream tasks can choose the level of abstraction best suited to their needs.

#### **3.3.** Dataset coverage

Our *HierVision* hub currently covers over 50 vision datasets 253 spanning a wide range of domains and recognition tasks. 254 The collection includes image classification datasets such as 255 CIFAR-100, ImageNet-100, ImageNet-1K, ImageNet-21K, 256 and iNaturalist; semantic segmentation datasets including 257 ADE20K and Cityscapes; action recognition datasets such 258 as Kinetics, ActivityNet, and FineSports; fine-grained cat-259 egorization with CUB, FGVC-Aircraft, and TreeOfLife-260 10M; and medical imaging datasets like CheXpert and 261 DeepLesion. Hierarchies in these datasets range from shal-262 low groupings of a few categories to deeply nested struc-263 tures with thousands of classes, reflecting both semantic and 264 biological taxonomies. Further details, including taxonomy 265 type and hierarchy source, are summarized in Table 1. 266

Figure 1 shows visualizations of the hierarchical struc-267 tures for six representative datasets: CIFAR-100 (semantic 268 object classification with a two-level hierarchy), Cityscapes 269 (urban scene segmentation with a class grouping tree), 270 RareSpecies (a deeply nested biological taxonomy), Ac-271 tivityNet (action recognition with a hierarchical verb on-272 tology), Moments in Time (an action dataset with a flat 273 class structure), and COCO-Stuff-10k (scene parsing with a 274 multi-level segmentation hierarchy). These examples high-275 light the diversity in both the structure and scale of the hi-276 erarchies, ranging from compact, balanced trees to large, 277 irregular graphs with hundreds or thousands of nodes. 278

#### ICCVW 2025 Submission #10. CONFIDENTIAL REVIEW COPY. DO NOT DISTRIBUTE.

303



Figure 1. Hierarchy visualizations for a selection of datasets form *HierVision*. The coverage in the repository hub ranges from small number of classes and shallow trees to very large number of classes and deep trees.

#### **3.4. Hierarchy Statistics** 279

281

285

286

287

288

280 We visualize statistics across all collected hierarchies to assess their structural diversity. Figure 2 summarizes four key 282 aspects: node count, maximum depth, the relationship between node count and depth, and average branching factor. 283 284

Most hierarchies contain fewer than 1,000 classes, though there is a long tail of large-scale taxonomies such as TreeOfLife-10M [110], ImageNet-21K [98], and Bamboo [132], that reach tens or hundreds of thousands of nodes (Figure 2a).

289 The majority of hierarchies are shallow, with depths of 2 or 3, as seen in datasets like CIFAR-100 and ADE20K 290 291 (Figure 2b). However, a subset—including TreeOfLife-292 10M [110] and Visual Genome [60]-features deeply nested structures, with depths exceeding 10. The scat-293 ter plot of node count versus maximum depth (Figure 2c) 294 highlights this diversity: some datasets combine high node 295 counts and depth, while others are both small and shallow. 296 297 Average branching factor calculates the average number

of children for each node, which also varies widely (Fig-298 ure 2d). Biological taxonomies tend to be deep and bal-299 anced, with low branching factors (1-3), while semantic hi-300 erarchies such as ADE20K [133] and Objects-365 [104] are 301 broad and flat, with extremely high branching at the root. 302

### 4. Integration into Deep Learning Pipelines

The standardized graph-based format adopted by HierVi-304 sion enables seamless integration of hierarchical informa-305 tion into modern deep learning workflows. In this section, 306 we outline typical approaches for incorporating hierarchies, 307 ranging from data loading and label preprocessing to ad-308 vanced hierarchical loss design and representation learning. 309 310

Hierarchies stored in our JSON format can be loaded 311 directly using widely adopted graph libraries such as Net-312 workX: 313 Table 1. All hierarchies available currently in *HierVision*. Datasets with multiple hierarchy versions (e.g., coarse/fine) are marked with \*. If the hierarchy was sourced from another paper, it is cited in the "Hierarchy Source" column.

	Dataset	Hierarchy Source	Original Format	Nodes	Edges	Depth	Classes
4	ActivityNet [14]	-	JSON	245	244	3	200
	DLD3V-10K [68]	-	JSON	81	80	2	64
	FineSports [124]	-	PKL	65	64	2	52
	HDM05 [86]	-	JSON	26	25	2	20
	HowTo100M [81]	-	JSON	142	141	2	129
Act	HumanAct12 [44]	-	NPY	47	46	2	34
ior	MAdverse [100] Matador [10]	-	JSON	82	000	4 5	5/8
/ <b>2</b> I	Mini-Kinetics-200 [58, 123]	-	ISON	240	239	3	200
Vi	MIntRec* [131]	-	TSV	240	237	2	200
de	Moments in Time [84]	-	JSON	486	485	4	339
0	Pseudo-Adverbs (ActivityNet) [30]	-	CSV	758	757	2	643
	Pseudo-Adverbs (MSRVTT) [30]	-	CSV	571	570	2	464
	Pseudo-Adverbs (VATEX) [30]	-	CSV	1686	1685	2	1550
	Something-Something V2* [42, 76]	[76]	JSON	225	224	2	174
	0CF101 [108]	-		12/	120	3	<u> </u>
	BioTrove-Balanced [127]	-	CSV	818	803 873	5 7	202
	BioTrove-LifeStages [127]	-	CSV	19	18	6	272 5
	BioTrove-Unseen [127]	-	ČŠV	4413	4703	7	1918
H	CUB-200-2011 [117]	-	JSON	251	250	3	200
Biological	iNaturalist [114]	-	CSV	4214	4213	2	4200
	MammalNet [20]	-	TSV	260	264	3	173
	Marine Tree [11]	-	CSV	79	78	5	62
	NABIRDS [115] Pare Species [110]	-	1X1 CSV	1011	1010	4	335 385
	Tree of Life [110]	-	CSV	635463	635462	22	537235
	VegFru-Fru92 [52]	-	JSON	103	102	2	92
	VegFru-Veg200 [52]	-	JSON	216	215	2	200
Ζ	CheXpert [55]	-	CSV	15	15	2	11
edi	DeepLesion [125]	-	CSV	172	2/3	2	117
cal	OpenCell [22]	-	CSV	13	15	3	11
	Domboo [122]		ISON	208207	245557	2	295240
	Caltech-101 [65]	-	Folder	113	113	4	285340
	CIFAR-100 [61]	-	PKL	121	120	2	100
	CMU MoCap [25]	-	WEB	306	305	3	280
	COCO-10K [15]	[3]	JSON	234	233	8	171
	COD10K [32]	-	JSON	75	74	2	69
	CORe50-Balanced [71]	-	TXT	70	69	4	50
	CORe50-Unbalanced [/1]	-	IXI	1665	00 1664	5 14	50
Q	EgoObjects [55] EGVC-Aircraft [77]	-	JSON	201	200	14	100
je	Fashionpedia-Attributes [56]	-	ISON	306	305	2	294
cts	Fashionpedia-Categories [56]	-	JSON	62	61	3	46
) (	IP102 [119]	-	TXT	113	112	3	102
en	ImageNet-100* [26]	[67]	Folder	121	120	2	100
ers	ImageNet-1K [26]	-	Folder	1763	1776	13	1010
<b>_</b>	ImageNet-21K* [98]	-	Folder	10891	11734	18	7414
	Objects 365 [104]	[98]	Folder	954	907 376	13	030 365
	OpenI ORIS [104]	-	ISON	98	97	27	69
	OpenImages [62]	-	JSON	603	648	5	525
	PASCAL VOC [31]	[3]	XML	36	35	6	21
	Portrait Mode 400 [49]	-	CSV	446	445	3	400
	Stanford Cars [59]	-	Folder	206	205	2	196
	Stanford Online Products [91]	-	Folder	22634	22633	2	22622
	ADE20K* [133] Cityscapes [23]	[3]	JSON WEB	1116 39	1115 38	2 2	1105 30
7.0	Grocery [79]	-	CSV	125	124	2	81
ce	Mapillary Vistas [87]	-	JSON	81	80	3	66
nes	Million-AID [73]	-	XML	74	73	3	51
\$/1	PACO-EGO4D [96]	-	JSON	515	514	2	441
Pla	PACO-LVIS [96] SUN260 [122]	-	JSON	532	531	2	458
ces	SUN300 [122] SUN307 [121]	-	WED CSV	570 /18	577 577	3	300
	SUN908 [121]	-	WEB	925	1088	3	904
	Visual Genome [60]	-	JSON	10503	10502	18	6114

336

337

338

339

340

341

342

343

344

345

346

347

348

349

350

351

352



Figure 2. Statistics and distributions of hierarchies in our HierVision collection.

```
314 1 import json, networkx as nx
315 2 with open('hierarchy.json') as f:
316 3 graph_dict = json.load(f)
317 4 G = nx.node_link_graph(graph_dict)
```

Listing 2. Loading a hierarchy into networkx.

318 Once loaded as a graph object, the hierarchy can be queried to obtain ancestor or descendant sets for each class, 319 compute semantic distances between nodes, or extract sub-320 hierarchies for specialized tasks. Below, we discuss the rel-321 322 evance of hierarchies to deep learning in euclidean space 323 (Section 4.1) and hyperbolic space (Section 4.2. Additionally in hyperbolic learning, we show how the few of the 324 hierarchies can be embedded into hyperbolic space. 325

#### **326 4.1. Relevance to Hierarchy-Aware Supervision**

#### 4.1.1. Hierarchy-Aware Label Representations and Multi-Level Outputs

HierVision enables augmenting dataset labels with hierar chical context by allowing straightforward retrieval of par ent and ancestor labels for any fine-grained class. This
 facilitates multi-task learning, where models predict class
 probabilities at multiple hierarchical levels (e.g., object
 category and super-category) simultaneously, typically via

level [118, 126]. Such coarse-to-fine supervision guides feature learning: high-level layers discern broad distinctions, while deeper layers refine for fine-grained classification.Alternatively, hierarchical classification [106] predicts

auxiliary output heads and backpropagating loss at each

sequentially, predicting a coarse category before specializing to specific subclasses [16, 70, 93, 118]. In both approaches, *HierVision*'s standardized graph simplifies retrieving relevant ancestor or child classes, ensuring hierarchy-consistent predictions and providing interpretable outputs at multiple levels of detail.

#### 4.1.2. Hierarchy-Aware Loss Functions

*HierVision* facilitates the design of loss functions that account for inter-class relationships, moving beyond standard cross-entropy's uniform error penalty. By leveraging the hierarchy, errors can be weighted by their distance in the tree or reflected in structured objectives.

A classic example is hierarchical softmax [85], which353factors the prediction over a tree of classes. Instead of a354flat N-class prediction, the model predicts a path from root355to leaf, decomposing the problem into a sequence of smaller356

382

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

binary classification tasks. For a tree with internal nodes  $\mathcal{N}$ , let y denote the true class and z the logits. The hierarchical softmax decomposes the probability of class y as:

360 
$$P(y \mid \mathbf{z}) = \prod_{n \in \text{path}(y)} P(n \mid \text{parent}(n), \mathbf{z})$$
(1)

where path(y) is the sequence of nodes from the root to y. The loss is then computed as the negative log-probability along this path:

$$\mathcal{L}_{\text{hier-softmax}} = -\log P(y \mid \mathbf{z}) \tag{2}$$

This approach not only reduces computational cost for large
 *N* but also implicitly aligns internal representations with the
 taxonomy.

Beyond hierarchical softmax, cost-sensitive losses penalize mistakes according to taxonomy distance, e.g., based on the length of the shortest path or the depth of the lowest common ancestor (LCA) between the predicted class  $\hat{y}$  and true class y [8]. Let  $d(y, \hat{y})$  denote the tree distance between classes. The loss can be defined as:

374 
$$\mathcal{L}_{\text{hier}-\text{LCA}} = \sum_{i=1}^{N} d(y_i, \hat{y}_i) \cdot \ell(y_i, \hat{y}_i)$$
(3)

where  $\ell$  is the standard loss (e.g., cross-entropy), and  $d(y, \hat{y})$ is typically the path length or a normalized form thereof.

377Similarly, hierarchy-based label smoothing replaces one-378hot targets with soft distributions, assigning higher proba-379bility mass to classes near the ground truth in the hierar-380chy [8, 97, 99]. For each target y, define the smoothed label381 $\tilde{y}$  as:

$$\tilde{y}_j = \frac{\exp(-\lambda \cdot d(y,j))}{\sum_{k=1}^N \exp(-\lambda \cdot d(y,k))}$$
(4)

where d(y, j) is the distance in the hierarchy between y and j, and  $\lambda > 0$  is a hyperparameter controlling the smoothing. The loss is then the standard cross-entropy between the prediction and  $\tilde{y}$ :

387 
$$\mathcal{L}_{\text{hier-smooth}} = -\sum_{j=1}^{N} \tilde{y}_j \log p_j$$
(5)

where  $p_j$  is the predicted probability for class j.

These hierarchy-aware objectives guide models to make
 semantically meaningful predictions, improve error robust ness, and support finer-grained evaluation of learning per formance.

#### **393 4.2. Hyperbolic Learning**

Hyperbolic learning uses the properties of hyperbolic
space to better represent hierarchical relationships in vision
datasets. Unlike Euclidean space, hyperbolic space expands

	Poincaré [88]		Entailment [36]		BHE [57]	
	Dist	mAP	Dist	mAP	Dist	mAP
CIFAR-100 [61]	0.713	0.162	0.18	0.623	0.026	0.885
ImageNet-100 [26]	0.450	0.119	0.20	65.91	0.095	0.746
CityScapes [23]	0.540	0.173	0.250	72.41	0.050	0.967
PASCAL-VOC [31]	0.477	0.122	0.182	0.692	0.05	0.837

Table 2. Hyperbolic embeddings of CIFAR-100, ImageNet-100, Cityscapes and PASCAL-VOC in using different hyperbolic embedding methods. Distortion and mAP [101] of the embeddings measure how well the tree distances are embedded in the hyperbolic space.

exponentially, making it ideal for embedding tree-like taxonomies with minimal distortion. In practice, both class labels and image features are embedded as points in a hyperbolic manifold (e.g., the Poincaré ball), so that distances between points correspond to semantic or taxonomic proximity in the hierarchy.

For hyperbolic learning with vision datasets, a typical process has two main steps. Step 1 is embedding the hierarchy itself: the class tree is mapped into hyperbolic space so that classes which are close in the hierarchy are also close together in the embedding. This is usually done using methods such as Poincaré embeddings [88] or entailment-based approaches [36], which are specifically designed to preserve the distances and relationships from the original tree. The effectiveness of this embedding is measured by two metrics following the hyperbolic learning literature [88, 101, 103]: distortion, which captures how well the hyperbolic distances match the true tree structure (lower is better), and mean average precision (mAP), which reflects how well nearest neighbors in the embedding correspond to true neighbors in the hierarchy (higher is better).

Step 2, as done in works like [5, 57, 74], involves learning: projecting image features produced by a neural network into hyperbolic space and training them to align with their respective class embeddings. This allows the model to make predictions that are consistent with the structure of the hierarchy.

In this section, we focus on Step 1 and analyze how well 424 different hyperbolic embedding methods can represent the 425 hierarchy itself in Table 2. Using standardized trees from 426 HierVision, we evaluate and compare several approaches 427 on CIFAR-100, ImageNet-1k, ActivityNet, and CUB. For 428 each dataset, we report distortion and mAP scores to as-429 sess how faithfully the class relationships are captured in 430 hyperbolic space. We use the hyperparameters defined in 431 Kasarla et al. [57] for generating the hyperbolic embeddings 432 of the methods, keeping hyperbolic dim = 64. In Table 2, 433 BHE [57] show better faithful tree embeddings. However, 434 for downstream tasks, any of these embeddings can be used 435 depending on the utility for the task. 436

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

437 The standardized JSON hierarchies in our hub make it easy to plug in such methods - for example, one can directly 438 439 use NetworkX to compute ancestor relations or to feed the graph into a Poincaré embedding algorithm to obtain initial 440 441 class vectors. The result is an integrated hyperbolic pipeline where both the data and the output of the model are easily 442 intergrated in the hyperbolic space. 443

#### **5.** Conclusion 444

We introduce *HierVision*, a hub for standardized hierarchi-445 cal knowledge across a broad spectrum of visual recognition 446 447 datasets. By consolidating and curating over 61 hierarchies 448 from diverse domains and encoding them in a unified graphbased format, we provide consistent, reproducible access to 449 450 structured label information. Our framework supports direct integration with graph libraries like NetworkX and en-451 ables hierarchical loss design, hyperbolic embeddings, and 452 large-scale benchmarking. 453

We believe that *HierVision* serves as a critical resource 454 455 for the vision community, promoting reproducibility, accelerating hierarchy-informed research, and enabling rigorous 456 benchmarking across a wide range of hierarchical struc-457 458 tures.

We will make the GitHub repo for HierVision public af-459 ter the review period. We shortlisted 40+ more hierarchies 460 in the pipeline to be added in the future. Recently, Ay-461 462 oughi et al. [6] discussed optimial tree structure for hyper-463 bolic embeddings, which can be further used to refine the existing hierarchies. We invite community contributions to 464 465 further expand and refine HierVision, and we hope this hub will catalyze advances in hierarchical representation learn-466 467 ing, benchmarking, and structured visual understanding at 468 scale.

#### References 469

470

471

473

474

475

476

477

478

479

480

481

486

487

- [1] Aristotle. Categories. Clarendon Press, Oxford, 350 BCE. Translated and edited by J.L. Ackrill, 1963. 1
- [2] Michael Ashburner, Catherine A Ball, Judith A Blake, 472 David Botstein, Heather Butler, J Michael Cherry, Allan P Davis, Kara Dolinski, Selina S Dwight, Janan T Eppig, et al. Gene ontology: tool for the unification of biology. Nature genetics, 25(1):25-29, 2000. 1
  - [3] Mina Ghadimi Atigh, Julian Schoep, Erman Acar, Nanne Van Noord, and Pascal Mettes. Hyperbolic image segmentation. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 4453-4462, 2022. 1, 2, 5
- 482 [4] Mina Ghadimi Atigh, Stephanie Nargang, Martin Keller-483 Ressel, and Pascal Mettes. Simzsl: Zero-shot learning be-484 yond a pre-defined semantic embedding space. Interna-485 tional Journal of Computer Vision, pages 1–17, 2025. 2
  - [5] Melika Ayoughi, Mina Ghadimi Atigh, Mohammad Mahdi Derakhshani, Cees GM Snoek, Pascal Mettes, and Paul

Groth. Continual hyperbolic learning of instances and 488 classes. arXiv preprint arXiv:2506.10710, 2025. 1, 2, 7 489

- [6] Melika Ayoughi, Max van Spengler, Pascal Mettes, and Paul Groth. Designing hierarchies for optimal hyperbolic embedding. In European Semantic Web Conference, pages 362-382. Springer, 2025. 2, 8
- [7] Sara Beery, Guanhang Wu, Trevor Edwards, Filip Pavetic, Bo Majewski, Shreyasee Mukherjee, Stanley Chan, John Morgan, Vivek Rathod, and Jonathan Huang. The auto arborist dataset: a large-scale benchmark for multiview urban forest monitoring under domain shift. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 21294-21307, 2022. 2
- [8] Luca Bertinetto, João F Henriques, and Philip HS Torr. Making the most of mistake-making: Hierarchical classification with partial labels. In European Conference on Computer Vision (ECCV), pages 706-722. Springer, 2020. 1, 2,
- [9] Luca Bertinetto, Romain Mueller, Konstantinos Tertikas, Sina Samangooei, and Nicholas A Lord. Making better mistakes: Leveraging class hierarchies with deep networks. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 12506-12515, 2020.
- [10] M. Beveridge and S. K. Nayar. (h)ierarchical (m)aterial (r)ecognition (f)rom (l)ocal (a)ppearance, 2025. 5
- [11] Thomas Boone et al. Marine-tree: Marine species classification. GitHub repository, 2021. 5
- [12] Guy Bukchin, Eli Schwartz, Kate Saenko, Ori Shahar, Rogerio Feris, Raja Giryes, and Leonid Karlinsky. Finegrained angular contrastive learning with coarse labels. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 8730-8740, 2021. 2
- [13] Hadas Bukchin, Yuval Alaluf, and Daniel Cohen-Or. Finegrained object recognition via deep contrastive learning. In Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops, pages 1534-1543, 2021. 2
- [14] Fabian Caba Heilbron, Victor Escorcia, Joao Carreira, Juan Carlos Niebles, and Bernard Ghanem. Activitynet: A large-scale video benchmark for human activity understanding. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 961–970, 2015. 5
- [15] Holger Caesar, Jasper Uijlings, and Vittorio Ferrari. Cocostuff: Thing and stuff classes in context. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 1209-1218, 2018. 5
- [16] Dongliang Chang, Kaiyue Pang, Yixiao Zheng, Zhanyu Ma, Yi-Zhe Song, and Jun Guo. Your" flamingo" is my" bird": Fine-grained, or not. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 11476-11485, 2021. 6
- [17] Bike Chen, Wei Peng, Xiaofeng Cao, and Juha Röning. Hyperbolic uncertainty aware semantic segmentation. IEEE Transactions on Intelligent Transportation Systems, 25(2): 1275-1290, 2023. 2

578

579

580

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

- [18] Haomin Chen, Shun Miao, Daguang Xu, Gregory D Hager,
  and Adam P Harrison. Deep hierarchical multi-label classification of chest x-ray images. In *International conference on medical imaging with deep learning*, pages 109–120. PMLR, 2019. 2
- [19] Jingzhou Chen, Peng Wang, Jian Liu, and Yuntao Qian.
  Label relation graphs enhanced hierarchical residual network for hierarchical multi-granularity classification. In *Proceedings of the IEEE/CVF Conference on Computer Vi-*sion and Pattern Recognition, pages 4858–4867, 2022. 2
- [20] Jun Chen, Ming Hu, Darren J Coker, Michael L Berumen,
  Blair Costelloe, Sara Beery, Anna Rohrbach, and Mohamed
  Elhoseiny. Mammalnet: A large-scale video benchmark for
  mammal recognition and behavior understanding. In *Proceedings of the IEEE/CVF conference on computer vision*and pattern recognition, pages 13052–13061, 2023. 5
- 561 [21] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Ge562 offrey Hinton. A simple framework for contrastive learning
  563 of visual representations. In *International Conference on*564 *Machine Learning*, pages 1597–1607. PMLR, 2020. 1
- [22] Nathan H Cho, Keith C Cheveralls, Andreas-David Brunner, Kibeom Kim, André C Michaelis, Preethi Raghavan, Hirofumi Kobayashi, Laura Savy, Jason Y Li, Hera Canaj, et al. Opencell: Endogenous tagging for the cartography of human cellular organization. *Science*, 375(6585):eabi6983, 2022. 5
- 571 [23] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo
  572 Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe
  573 Franke, Stefan Roth, and Bernt Schiele. The cityscapes
  574 dataset for semantic urban scene understanding. In *Pro-*575 *ceedings of the IEEE Conference on Computer Vision and*576 *Pattern Recognition*, pages 3213–3223, 2016. 5, 7
  - [24] Yawen Cui, Zitong Yu, Wei Peng, Qi Tian, and Li Liu. Rethinking few-shot class-incremental learning with open-set hypothesis in hyperbolic geometry. *IEEE Transactions on Multimedia*, 26:5897–5910, 2023. 2
- 581 [25] Fernando De la Torre, Jessica Hodgins, Adam Bargteil,
  582 Xavier Martin, Justin Macey, Alex Collado, and Pep Bel583 tran. Guide to the carnegie mellon university multimodal
  584 activity (cmu-mmac) database. 2009. 5
- 585 [26] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li,
  586 and Li Fei-Fei. Imagenet: A large-scale hierarchical image
  587 database. In 2009 IEEE Conference on Computer Vision
  588 and Pattern Recognition, pages 248–255. IEEE, 2009. 1, 2,
  589 5, 7
- [27] Karan Desai, Maximilian Nickel, Tanmay Rajpurohit, Justin Johnson, and Shanmukha Ramakrishna Vedantam. Hyperbolic image-text representations. In *International Conference on Machine Learning*, pages 7694–7731.
  PMLR, 2023. 2
- 595 [28] Ankit Dhall, Anastasia Makarova, Octavian Ganea, Dario
  596 Pavllo, Michael Greeff, and Andreas Krause. Hierarchical
  597 image classification using entailment cone embeddings. In
  598 *CVPR Workshop on Differential Geometry in Computer Vi-*599 *sion and Machine Learning*, 2020. 1, 2
- [29] Lars Doorenbos, Pablo Márquez-Neila, Raphael Sznitman,
   and Pascal Mettes. Hyperbolic random forests. arXiv
   preprint arXiv:2308.13279, 2023. 2

- [30] Hazel Doughty and Cees G. M. Snoek. How Do You Do It?
  Fine-Grained Action Understanding with Pseudo-Adverbs.
  In *The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 5
- [31] Mark Everingham, Luc Van Gool, Christopher KI
  Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. 2, 5, 7
- [32] Deng-Ping Fan, Ge-Peng Ji, Guolei Sun, Ming-Ming Cheng, Jianbing Shen, and Ling Shao. Camouflaged object detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2777–2787, 2020. 5
- [33] Jiyang Fan, Bo Xiong, Lu Jiang, Yin Cui, Chen Sun, Manmohan Chandraker Jain, and Kristen Grauman. Egoobjects: A large-scale egocentric dataset for object understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021. 5
- [34] Clement Farabet, Camille Couprie, Laurent Najman, and Yann LeCun. Learning hierarchical features for scene labeling. *IEEE transactions on pattern analysis and machine intelligence*, 35(8):1915–1929, 2012. 2
- [35] Pedro F Felzenszwalb, Ross B Girshick, and David McAllester. Cascade object detection with deformable part models. In 2010 IEEE Computer society conference on computer vision and pattern recognition, pages 2241–2248. Ieee, 2010. 2
- [36] Octavian Ganea, Gary Bécigneul, and Thomas Hofmann. Hyperbolic entailment cones for learning hierarchical embeddings. In *International conference on machine learning*, pages 1646–1655. PMLR, 2018. 1, 2, 7
- [37] Yuting Gao, Jinfeng Liu, Zihan Xu, Jun Zhang, Ke Li, Rongrong Ji, and Chunhua Shen. Pyramidclip: Hierarchical feature alignment for vision-language model pretraining. Advances in neural information processing systems, 35:35959–35970, 2022. 2
- [38] Ashima Garg, Depanshu Sani, and Saket Anand. Learning hierarchy aware features for reducing mistake severity. In *European Conference on Computer Vision*, pages 252–267. Springer, 2022. 2
- [39] Ashima Garg, Depanshu Sani, and Saket Anand. Learning hierarchy aware features for reducing mistake severity. In *European Conference on Computer Vision*, pages 252–267. Springer, 2022. 2
- [40] Songwei Ge, Shlok Mishra, Simon Kornblith, Chun-Liang Li, and David Jacobs. Hyperbolic contrastive learning for visual representations beyond objects. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pages 6840–6849, 2023. 2
- [41] Mina Ghadimi Atigh, Martin Keller-Ressel, and Pascal Mettes. Hyperbolic busemann learning with ideal prototypes. Advances in neural information processing systems, 34:103–115, 2021. 2
- [42] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video

661

662

690

691

692

693

694

700

701

725

726

727

728

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

database for learning and evaluating visual common sense. In Proceedings of the IEEE international conference on computer vision, pages 5842-5850, 2017. 5

- 663 [43] Sadaf Gulshad, Teng Long, and Nanne van Noord. Hier-664 archical explanations for video action recognition. In Pro-665 ceedings of the IEEE/CVF Conference on Computer Vision 666 and Pattern Recognition, pages 3703-3708, 2023. 2
- [44] Chuan Guo, Xinxin Zuo, Sen Wang, Shihao Zou, Qingyao 667 668 Sun, Annan Deng, Minglun Gong, and Li Cheng. Ac-669 tion2motion: Conditioned generation of 3d human motions. In Proceedings of the 28th ACM International Conference 670 on Multimedia, pages 2021-2029, 2020. 5 671
- 672 [45] Yunhui Guo, Xudong Wang, Yubei Chen, and Stella X Yu. Clipped hyperbolic classifiers are super-hyperbolic classi-673 fiers. In Proceedings of the IEEE/CVF Conference on Com-674 675 puter Vision and Pattern Recognition, pages 11-20, 2022. 676 2
- [46] Yuanfan Guo, Minghao Xu, Jiawen Li, Bingbing Ni, Xu-677 anyu Zhu, Zhenbang Sun, and Yi Xu. Hcsc: Hierarchi-678 679 cal contrastive selective coding. In Proceedings of the 680 IEEE/CVF Conference on Computer Vision and Pattern 681 Recognition, pages 9706–9715, 2022. 2
- 682 [47] Aric Hagberg, Pieter J Swart, and Daniel A Schult. Explor-683 ing network structure, dynamics, and function using networkx. Technical report, Los Alamos National Laboratory 684 685 (LANL), Los Alamos, NM (United States), 2008. 2, 3
- 686 [48] Matthias Hamann. On the tree-likeness of hyperbolic spaces. In Mathematical proceedings of the cambridge 687 philosophical society, pages 345-361. Cambridge Univer-688 sity Press, 2018. 2 689
  - [49] Mingfei Han, Linjie Yang, Xiaojie Jin, Jiashi Feng, Xiaojun Chang, and Heng Wang. Video recognition in portrait mode. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 21831-21841, 2024. 5
- 695 [50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 696 Deep residual learning for image recognition. Proceedings 697 of the IEEE Conference on Computer Vision and Pattern 698 Recognition, pages 770-778, 2016. 1
- 699 [51] Jie Hong, Zeeshan Hayder, Junlin Han, Pengfei Fang, Mehrtash Harandi, and Lars Petersson. Hyperbolic audiovisual zero-shot learning. In Proceedings of the IEEE/CVF 702 international conference on computer vision, pages 7873-703 7883, 2023. 2
- 704 [52] Qiwei Hou, Hongzhi Wu, Zilei Ye, Qian Qiu, Xiaokang 705 Yang, and Meng Wang Tang. Vegfru: A domain-specific 706 dataset for fine-grained visual categorization. In Proceed-707 ings of the IEEE Conference on Computer Vision and Pat-708 tern Recognition, pages 541-549, 2017. 5
- 709 [53] Thomas Hoyoux, Antonio J Rodríguez-Sánchez, and Jus-710 tus H Piater. Can computer vision problems benefit from 711 structured hierarchical classification? Machine Vision and 712 Applications, 27(8):1299-1312, 2016. 1, 2
- 713 [54] Sarah Ibrahimi, Mina Ghadimi Atigh, Nanne Van Noord, Pascal Mettes, and Marcel Worring. Intriguing properties of 714 715 hyperbolic embeddings in vision-language models. Trans-716 actions on Machine Learning Research, 2024. 2

- [55] Jeremy Irvin, Pranav Rajpurkar, Michael Ko, Yifan Yu, 717 Silviana Ciurea-Ilcus, Christopher Chute, Henrik Mark-718 lund, Behzad Haghgoo, Robyn Ball, Katie Shpanskaya, 719 Jake Seekins, David Mong, Safwan Halabi, Jacob Sand-720 berg, Russell Jones, David Larson, Curtis Langlotz, B Patel, 721 and Matthew Lungren. Chexpert: A large chest radiograph 722 dataset with uncertainty labels and expert comparison. Ra-723 diology: Artificial Intelligence, 1(1):e180057, 2019. 5 724
- [56] Menglin Jia, Mengyun Shi, Mikhail Sirotenko, Yin Cui, Claire Cardie, Bharath Hariharan, Hartwig Adam, and Serge Belongie. Fashionpedia: Ontology, segmentation, and an attribute localization dataset. In Computer vision-ECCV 2020: 16th European conference, glasgow, UK, August 23-28, 2020, proceedings, part i 16, pages 316-332. Springer, 2020. 5
- [57] Tejaswi Kasarla, Max van Spengler, and Pascal Mettes. Balanced hyperbolic embeddings are natural out-ofdistribution detectors. arXiv preprint arXiv:2506.10146, 2025. 1, 2, 7
- [58] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, Mustafa Suleyman, and Andrew Zisserman. The kinetics human action video dataset. arXiv preprint arXiv:1705.06950, 2017. 1, 5
- [59] Jonathan Krause, Michael Stark, Jia Deng, and Li Fei-Fei. 3d object representations for fine-grained categorization. In 4th International IEEE Workshop on 3D Representation and Recognition (3dRR-13), Sydney, Australia, 2013. 5
- [60] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David Shamma, Michael Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. International Journal of Computer Vision, 123(1):32–73, 2017. 4, 5
- [61] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009. 2, 5, 7
- [62] Alina Kuznetsova, Hassan Rom, Neil Alldrin, Jasper Uijlings, Ivan Krasin, Jordi Pont-Tuset, Shahab Kamali, Stefan Popov, Marco Malloci, Alexander Kolesnikov, Tom Duerig, et al. The open images dataset v4: Unified image classification, object detection, and visual relationship detection at scale. International Journal of Computer Vision, 128(7):1956–1981, 2020. 2, 5
- [63] Marc Law, Renjie Liao, Jake Snell, and Richard Zemel. Lorentzian distance learning for hyperbolic representations. In International Conference on Machine Learning, pages 3672-3681. PMLR, 2019. 2
- [64] Diego Lazcano, Nicolás Fredes Franco, and Werner Creixell. Hgan: Hyperbolic generative adversarial network. IEEE Access, 9:96309–96320, 2021. 2
- [65] Fei-Fei Li, Marco Andreeto, Marc'Aurelio Ranzato, and Pietro Perona. Caltech 101, 2022. 5
- Li-Jia Li, Chong Wang, Yongwhan Lim, David M Blei, [66] and Li Fei-Fei. Building and using a semantivisual image hierarchy. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pages 3336-3343. IEEE, 2010. 2

810

811

812

813

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

- [67] Randolph Linderman, Jingyang Zhang, Nathan Inkawhich,
  Hai Li, and Yiran Chen. Fine-grain inference on out-ofdistribution data with hierarchical classification. In *Confer- ence on Lifelong Learning Agents*, pages 162–183. PMLR,
  2023. 5
- [68] Lu Ling, Yichen Sheng, Zhi Tu, Wentian Zhao, Cheng Xin,
  Kun Wan, Lantao Yu, Qianyu Guo, Zixun Yu, Yawen Lu,
  et al. Dl3dv-10k: A large-scale scene dataset for deep
  learning-based 3d vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*,
  pages 22160–22169, 2024. 5
- [69] 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 *Proceed- ings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9273–9281, 2020. 2
- [70] Yang Liu, Lei Zhou, Pengcheng Zhang, Xiao Bai, Lin Gu,
  Yiaohan Yu, Jun Zhou, and Edwin R Hancock. Where to
  focus: Investigating hierarchical attention relationship for
  fine-grained visual classification. In *European Conference on Computer Vision*, pages 57–73. Springer, 2022. 6
- [71] Vincenzo Lomonaco and Davide Maltoni. Core50: a new dataset and benchmark for continuous object recognition.
  In *Conference on Robot Learning*, pages 17–26, 2017. 5
- 799 [72] Teng Long, Pascal Mettes, Heng Tao Shen, and Cees GM
  800 Snoek. Searching for actions on the hyperbole. In *Proceed-*801 *ings of the IEEE/CVF conference on computer vision and*802 *pattern recognition*, pages 1141–1150, 2020. 2
- [73] Yang Long, Gui-Song Xia, Shengyang Li, Wen Yang,
  Michael Ying Yang, Xiao Xiang Zhu, Liangpei Zhang, and
  Deren Li. On creating benchmark dataset for aerial image
  interpretation: Reviews, guidances, and million-aid. *IEEE Journal of selected topics in applied earth observations and remote sensing*, 14:4205–4230, 2021. 5
  - [74] Zhitao Long, Deng Cai, and Wenqian Gan. Searching for an effective and efficient hyperbolic representation for supervised learning. In Advances in Neural Information Processing Systems (NeurIPS), pages 18776–18788, 2020. 1, 2, 7
- [75] Chao Ma, Jia-Bin Huang, Xiaokang Yang, and Ming-Hsuan
  Yang. Hierarchical convolutional features for visual tracking. In *Proceedings of the IEEE international conference on computer vision*, pages 3074–3082, 2015. 2
- [76] Farzaneh Mahdisoltani, Guillaume Berger, Waseem Gharbieh, David Fleet, and Roland Memisevic. On the effective-ness of task granularity for transfer learning. *arXiv preprint arXiv:1804.09235*, 2018. 5
- 822 [77] Subhransu Maji, Raviv Raich, Greg Shakhnarovich, Rogerio Feris, Ludwig Schmidt, Ivan Laptev, and Josef Sivic.
  824 Fine-grained visual classification of aircraft. In *Proceed-*825 *ings of the IEEE Conference on Computer Vision and Pat-*826 *tern Recognition Workshops*, 2013. 5
- 827 [78] Marcin Marszalek and Cordelia Schmid. Semantic hierarchies for visual object recognition. In 2007 IEEE Conference on Computer Vision and Pattern Recognition, pages
  830 1–7. IEEE, 2007. 2
- [79] Patricia Martinez-Gonzalez, Hedvig Kjellstrom, and
   Javier J Romero. A grocery store image dataset with vi-

sual and semantic labels for object recognition. In 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), 2019. 5

- [80] 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. 1, 2
- [81] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. HowTo100M: Learning a Text-Video Embedding by Watching Hundred Million Narrated Video Clips. In *ICCV*, 2019. 5
- [82] George A Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995. 1, 2
- [83] Safaa Abdullahi Moallim Mohamud, Minjin Baek, and Dong Seog Han. Hierarchical question-answering for driving scene understanding using vision-language models. *arXiv preprint arXiv:2506.02615*, 2025. 2
- [84] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Bargal, Tian Yan, Yinan Lin, Lisa Brown, Qingqiu Fan, Daniel Gutfruend, Carl Vondrick, Aude Oliva, and Antonio Torralba. Moments in time dataset: one million videos for event understanding. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 739–748, 2019. 5
- [85] F Morin and Y Bengio. Hierarchical probabilistic neural network language model. In *Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics* (AISTATS), pages 246–252, 2005. 2, 6
- [86] M. Müller, T. Röder, M. Clausen, B. Eberhardt, B. Krüger, and A. Weber. Documentation mocap database hdm05. Technical Report CG-2007-2, Universität Bonn, 2007. 5
- [87] Gerhard Neuhold, Tobias Ollmann, Samuel Rota Bulo, and Peter Kontschieder. The mapillary vistas dataset for semantic urban scene understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4990– 4999, 2017. 5
- [88] Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. In Advances in Neural Information Processing Systems (NeurIPS), 2017. 1, 2, 7
- [89] Maximillian Nickel and Douwe Kiela. Learning continuous hierarchies in the lorentz model of hyperbolic geometry. In *International conference on machine learning*, pages 3779– 3788. PMLR, 2018. 2
- [90] Zachary Novack, Julian McAuley, Zachary Chase Lipton, and Saurabh Garg. Chils: Zero-shot image classification with hierarchical label sets. In *International Conference on Machine Learning*, pages 26342–26362. PMLR, 2023. 2
- [91] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4004–4012, 2016. 5
- [92] Avik Pal, Max van Spengler, Guido Maria D'Amely di Melendugno, Alessandro Flaborea, Fabio Galasso, and Pascal Mettes. Compositional entailment learning for hyperbolic
   889

891

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1001

1002

1003

1004

vision-language models. In *The Thirteenth International* Conference on Learning Representations, 2025. 2

- [93] Seulki Park, Youren Zhang, X Yu Stella, Sara Beery, and
  Jonathan Huang. Visually consistent hierarchical image
  classification. In *The Thirteenth International Conference on Learning Representations*, 2025. 1, 2, 6
- 896 [94] Tobia Poppi, Tejaswi Kasarla, Pascal Mettes, Lorenzo
  897 Baraldi, and Rita Cucchiara. Hyperbolic safety-aware
  898 vision-language models. In *Proceedings of the Computer*899 *Vision and Pattern Recognition Conference*, pages 4222–
  900 4232, 2025. 2
- 901 [95] Eric Qu and Dongmian Zou. Autoencoding hyperbolic
   902 representation for adversarial generation. *arXiv preprint* 903 *arXiv:2201.12825*, 2022. 2
- 904 [96] Vignesh Ramanathan, Anmol Kalia, Vladan Petrovic, Yi
  905 Wen, Baixue Zheng, Baishan Guo, Rui Wang, Aaron Mar906 quez, Rama Kovvuri, Abhishek Kadian, et al. Paco: Parts
  907 and attributes of common objects. In *Proceedings of the*908 *IEEE/CVF Conference on Computer Vision and Pattern*909 *Recognition*, pages 7141–7151, 2023. 5
- 910 [97] Sarah Rastegar, Yuki M Asano, Hazel Doughty, and Cees
  911 G. M. Snoek. Generalized category discovery with hierar-912 chical label smoothing, 2024. 7
- [98] Tal Ridnik, Ethan Ben-Baruch, Assaf Noy, and Lihi ZelnikManor. Imagenet-21k pretraining for the masses. arXiv
  preprint arXiv:2104.10972, 2021. 4, 5
- [99] Luuk Romeijn, Andrius Bernatavicius, and Duong Vu. Mycoai: Fast and accurate taxonomic classification for fungal its sequences. *Molecular Ecology Resources*, 24(8):
  e14006, 2024. 7
- 920[100] Amruth Sagar, Rishabh Srivastava, Venkata Kesav Venna,<br/>Ravi Kiran Sarvadevabhatla, et al. Madverse: A hierarchi-<br/>cal dataset of multi-lingual ads from diverse sources and<br/>categories. In *Proceedings of the IEEE/CVF Winter Con-*<br/><br/>924923*ference on Applications of Computer Vision*, pages 8087–<br/>8096, 2024. 5
- 926[101] Frederic Sala, Chris De Sa, Albert Gu, and Christopher Ré.927Representation tradeoffs for hyperbolic embeddings. In In-928ternational conference on machine learning, pages 4460-9294469. PMLR, 2018. 1, 2, 7
- [102] Shibani Santurkar, Dimitris Tsipras, and Aleksander
   Madry. Breeds: Benchmarks for subpopulation shift. In *International Conference on Learning Representations*, 2021.
   2
- [103] Rik Sarkar. Low distortion delaunay embedding of trees
  in hyperbolic plane. In *International symposium on graph drawing*, pages 355–366. Springer, 2011. 1, 7
- [104] Shuai Shao, Zeming Zhao, Bo Li, Tiancheng Xiao, Gang
  Yu, Xiangyu Zhang, and Jian Sun. Objects365: A largescale, high-quality dataset for object detection. In *Proceed- ings of the IEEE/CVF International Conference on Com- puter Vision*, pages 8430–8439, 2019. 4, 5
- 942 [105] Qi She, Liang Zhang, Yongchao Zhuang, Haotian Zhao, Jun
  943 Yang, Guohao Ma, Jiawei Liang, Jiaqi Dong, Mingcheng
  944 Cao, Bin Han, Minghao Wang, Chengdong Jiang, Yi
  945 Zhang, Yi Hu, Yongxiang Liu, and Jiaqi Li. Openloris946 object: A robotic vision dataset and benchmark for lifelong

learning. International Journal of Computer Vision, 128: 2413–2430, 2020. 5

- [106] Carlos N Silla and Alex A Freitas. A survey of hierarchical classification across different application domains. *Data mining and knowledge discovery*, 22:31–72, 2011. 6
- [107] Josef Sivic, Bryan C Russell, Andrew Zisserman, William T Freeman, and Alexei A Efros. Unsupervised discovery of visual object class hierarchies. In 2008 IEEE Conference on Computer Vision and Pattern Recognition, pages 1–8. IEEE, 2008. 2
- [108] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah.
   Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
   5
- [109] Satwik Srivastava and Deepak Mishra. Severity of error in hierarchical datasets. *Scientific Reports*, 13(1), 2023. 2
- [110] Olivia Stevens, David Gunning, Wenjie Miao, Wenbin Zeng, Ana Saldanha, Leslie Smith, Lavanya Baskaran, Pramod Mathai, Deep Joshi, Wei Peng, et al. Bioclip: A vision foundation model for the tree of life. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 2024. 2, 4, 5
- [111] Tanuj Sur, Samrat Mukherjee, Kaizer Rahaman, Subhasis Chaudhuri, Muhammad Haris Khan, and Biplab Banerjee. Hyperbolic uncertainty-aware few-shot incremental point cloud segmentation. In Proceedings of the Computer Vision and Pattern Recognition Conference, pages 11810–11821, 2025. 2
- [112] Abraham Ungar. A gyrovector space approach to hyperbolic geometry. Springer Nature, 2022. 2
- [113] Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In *Proceedings of the IEEE conference* on computer vision and pattern recognition, pages 595– 604, 2015. 2, 5
- [114] Grant Van Horn, Oisin Mac Aodha, Yang Song, Chenyi Cui, Chen Sun, Andrew Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8769–8778, 2018. 2, 5
- [115] Max van Spengler and Pascal Mettes. Low-distortion and gpu-compatible tree embeddings in hyperbolic space. *arXiv preprint arXiv:2502.17130*, 2025. 1, 2
- [116] Max Van Spengler, Erwin Berkhout, and Pascal Mettes. Poincaré resnet. In Proceedings of the IEEE/CVF International Conference on Computer Vision, pages 5419–5428, 2023. 2
- [117] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The caltech-ucsd birds-200-2011 dataset. Technical report, California Institute of Technology, 2011. 1, 2, 5
- [118] Jonatas Wehrmann, Ricardo Cerri, and Rodrigo Barros. Hierarchical multi-label classification networks. In *International conference on machine learning*, pages 5075–5084. PMLR, 2018. 6

- 1005[119] Xinpeng Wu, Peng Chen, Shuhan Wang, Bing Xu, Qiang1006Zhang, Yanan Zheng, Yihua Zhao, Yuliang Zhang, Jiawei1007Ren, Zhiyong Liu, and Guijun Yang. Ip102: A large-1008scale benchmark dataset for insect pest recognition. In Pro-1009ceedings of the AAAI Conference on Artificial Intelligence,1010pages 9177–9184, 2019. 5
- [1011 [120] Yongqin Xian, Christoph H. Lampert, Bernt Schiele, and
  [1012 Zeynep Akata. Zero-shot learning—a comprehensive eval1013 uation of the good, the bad and the ugly. *IEEE Transactions*1014 *on Pattern Analysis and Machine Intelligence*, 41(9):2251–
  1015 2265, 2019. 5
- [121] Jianxiong Xiao, James Hays, Krista A Ehinger, Aude Oliva, and Antonio Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In 2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pages 3485–3492, 2010. 5
- 1021 [122] Jianxiong Xiao, Krista A Ehinger, James Hays, Antonio
  1022 Torralba, and Aude Oliva. Recognizing scene viewpoint
  1023 using panoramic place representation. In 2012 IEEE Con1024 ference on Computer Vision and Pattern Recognition, pages
  1025 2695–2702, 2012. 5
- [123] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European conference on computer vision (ECCV)*, pages 305–321, 2018. 5
- 1031 [124] Jinglin Xu, Guohao Zhao, Sibo Yin, Wenhao Zhou, and
  1032 Yuxin Peng. Finesports: A multi-person hierarchical sports
  1033 video dataset for fine-grained action understanding. In *Pro-*1034 *ceedings of the IEEE/CVF Conference on Computer Vision*1035 *and Pattern Recognition*, pages 21773–21782, 2024. 5
- 1036 [125] Ke Yan, Xiaoshuang Wang, Le Lu, and Ronald M Summers. Deeplesion: automated mining of large-scale lesion annotations and universal lesion detection with deep learning. *Journal of Medical Imaging*, 5(3):036501, 2018. 5
- 1040[126] Zhicheng Yan, Hao Zhang, Robinson Piramuthu, Vignesh1041Jagadeesh, Dennis DeCoste, Wei Di, and Yizhou Yu. Hd-1042cnn: hierarchical deep convolutional neural networks for1043large scale visual recognition. In *Proceedings of the IEEE*1044*international conference on computer vision*, pages 2740–10452748, 2015. 1, 6
- [127] Chih-Hsuan Yang, Benjamin Feuer, Talukder Jubery, Zi
  Deng, Andre Nakkab, Md Zahid Hasan, Shivani Chiranjeevi, Kelly Marshall, Nirmal Baishnab, Asheesh Singh, et al. Biotrove: A large curated image dataset enabling ai
  for biodiversity. Advances in Neural Information Processing Systems, 37:102101–102120, 2024. 5
- 1052 [128] William Yang, Byron Zhang, and Olga Russakovsky.
  1053 Imagenet-ood: Deciphering modern out-of-distribution de1054 tection algorithms. In *The Twelfth International Conference*1055 *on Learning Representations*, 2024. 5
- 1056[129] Zhen Yu, Toan Nguyen, Yaniv Gal, Lie Ju, Shekhar S Chan-1057dra, Lei Zhang, Paul Bonnington, Victoria Mar, Zhiyong1058Wang, and Zongyuan Ge. Skin lesion recognition with1059class-hierarchy regularized hyperbolic embeddings. In In-1060ternational conference on medical image computing and1061computer-assisted intervention, pages 594–603. Springer,10622022. 2

- [130] Yun Yue, Fangzhou Lin, Kazunori D Yamada, and Ziming Zhang. Hyperbolic contrastive learning. arXiv preprint arXiv:2302.01409, 2023. 2
   1063
   1064
   1065
- [131] Hanlei Zhang, Hua Xu, Xin Wang, Qianrui Zhou, Shaojie
  Zhao, and Jiayan Teng. Mintrec: A new dataset for multimodal intent recognition. In *Proceedings of the 30th ACM International Conference on Multimedia*, page 1688–1697, 2022. 5
- [132] Yuanhan Zhang, Qinghong Sun, Yichun Zhou, Zexin He, Zhenfei Yin, Kun Wang, Lu Sheng, Yu Qiao, Jing Shao, and Ziwei Liu. Bamboo: Building mega-scale vision dataset continually with human-machine synergy, 2022. 4, 5
  1071
  1072
  1073
  1074
- [133] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 633–641, 2017. 2, 4, 5
  [133] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela 1075
  [133] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela 1076
  [133] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela 1076
  [133] Barriuso, and Antonio Torralba. Scene parsing through 1076
  [133] Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela 1075
  [133] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1076
  [133] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1076
  [133] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1076
  [133] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1076
  [133] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1076
  [134] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1077
  [135] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [136] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1077
  [136] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [136] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [136] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1077
  [136] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [137] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1077
  [137] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [138] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [137] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
  [138] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1077
  [138] Bolei Zhou, Hang Zhao, Yavier Puig, Sanja Fidler, Adela 1078
- [134] Song-Chun Zhu, David Mumford, et al. A stochastic grammar of images. *Foundations and Trends*® in Computer Graphics and Vision, 2(4):259–362, 2007. 2