

BIOCLIP 2: Emergent Properties from Scaling Hierarchical Contrastive Learning

Jianyang Gu^{1†}, Samuel Stevens¹, Elizabeth G Campolongo¹, Matthew J Thompson¹, Net Zhang¹, Jiaman Wu¹, Andrei Kopanev¹, Zheda Mai¹, Alexander E. White², James Balhoff³, Wasila Dahdul⁴, Daniel Rubenstein⁵, Hilmar Lapp⁶, Tanya Berger-Wolf¹, Wei-Lun Chao¹, Yu Su^{1†}

¹The Ohio State University, ²Smithsonian Institution, ³UNC Chapel Hill,
⁴University of California, Irvine, ⁵Princeton University, ⁶Duke University

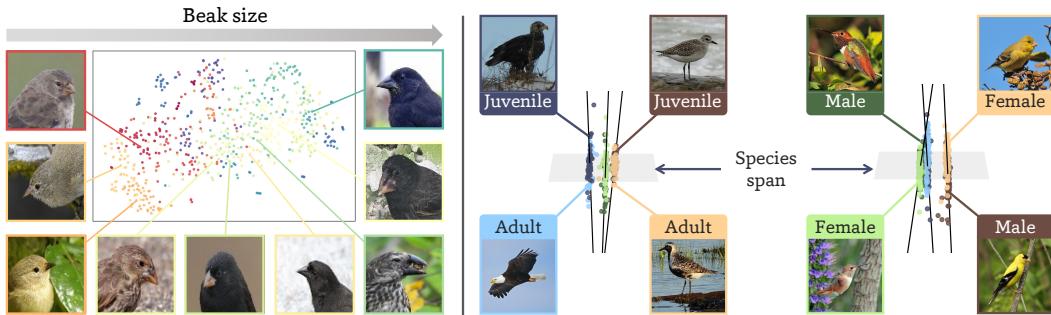


Figure 1: While BIOCLIP 2 is trained to distinguish species, it demonstrates emergent properties beyond the initial training objective. **Left:** At the *inter-species* level, the embedding distribution of different species aligns with ecological relationships; the embeddings of Darwin’s finches arrange themselves by beak size from left to right. **Right:** Instead of collapsing, the *intra-species* variations are preserved in subspaces orthogonal to the inter-species variation (the black lines point from the mean embedding of one variant to that of the other variant). Orthogonality increases with scale (see Figure 3c).

Abstract

Foundation models trained at scale exhibit remarkable emergent behaviors, learning new capabilities beyond their initial training objectives. We find such emergent behaviors in biological vision models via large-scale contrastive vision-language training. To achieve this, we first curate TREEOFLIFE-200M, comprising 214 million images of living organisms, the largest and most diverse biological organism image dataset to date. We then train BIOCLIP 2 on TREEOFLIFE-200M to distinguish different species. Despite the narrow training objective, BIOCLIP 2 yields extraordinary accuracy when applied to various biological visual tasks such as habitat classification and trait prediction. We identify emergent properties in the learned embedding space of BIOCLIP 2. At the inter-species level, the embedding distribution of different species aligns closely with functional and ecological meanings (*e.g.*, beak sizes and habitats). At the intra-species level, instead of being diminished, the intra-species variations (*e.g.*, life stages and sexes) are preserved and better separated in subspaces orthogonal to inter-species distinctions. We provide formal proof and analyses to explain why hierarchical supervision and contrastive objectives encourage these emergent properties. Crucially, our results reveal that these properties become increasingly significant with larger-scale training data, leading to a biologically meaningful embedding space.

Models, data, and code available at imageomics.github.io/bioclip-2. [†]{gu.1220, su.809}@osu.edu

1 Introduction

Recent advances in artificial intelligence (AI) are transforming core scientific workflows to become more efficient and automated [1, 2]. Tasks that once demanded overwhelming time and labor, like inferring atomic protein folds, designing functional materials, or producing global weather forecasts, can now be efficiently accomplished by large-scale predictive models [3, 4, 5, 6, 7, 8]. Particularly, a growing class of *domain-specific foundation models* demonstrate capabilities that arise without explicit definition during training [9, 10]. For example, language models trained purely on large-scale amino-acid strings unexpectedly develop an understanding of 3D chemistry to predict atomic folds with near-experimental accuracy [4, 11]. These scale-driven emergent abilities are reshaping scientific inference and opening new avenues for data-centric discovery.

In ecology and evolutionary biology, previous efforts leveraged hierarchical taxonomic labels and CLIP-style contrastive training [12] to achieve pronounced species classification accuracy across the tree of life [13, 14, 15]. This work asks a simple but intriguing question: *what properties emerge if we scale up hierarchical contrastive training?* To answer this question, we curate **TREEOFLIFE-200M**, comprising 214M organism images spanning 952K taxonomic classes, making it the largest and most diverse visual catalog of life to date. Through training at scale, our model **BIOCLIP 2** improves species classification accuracy by 18.0% over **BIOCLIP** [13]. More importantly, we explore whether representations learned solely through species-level supervision can generalize to diverse biological questions **beyond species classification**.

To probe these capabilities, we evaluate **BIOCLIP 2** on a variety of existing biological visual tasks, including habitat classification [16], trait prediction [17, 18], new-species identification [19], and agricultural disease detection [20]. These applications push beyond simple species recognition and apply directly to biodiversity conservation, trait organization, and agricultural health. Despite being trained primarily with species-level supervision, **BIOCLIP 2** outperforms both vision-language (*e.g.*, **SigLIP** [21]) and vision-only baselines (*e.g.*, **DINOv3** [22]) by an average margin of 10.3% on these tasks. We then look deeper into the embedding space of **BIOCLIP 2** and identify two *emergent properties* as the training scales up.

At the **inter-species** level, the embedding distribution of different species aligns with their ecological relationships. As shown on the left side of [Figure 1](#), **BIOCLIP 2** embeddings of Darwin’s finches demonstrate an increasing beak size from left to right, which is not observed in the original CLIP embedding space. We attribute the property to the adopted hierarchical taxonomic labels, which inherently encode functional and ecological information [23]. The hierarchical supervision at scale pushes related species to co-locate in functionally coherent “macro-clusters.” In such a way, **BIOCLIP 2** acquires functional trait knowledge without using explicitly labeled traits.

At the **intra-species** level, contrary to the intuition that fine-grained differences collapse after extensive training [24, 25], **BIOCLIP 2** keeps the intra-species variations (*e.g.*, life stages and sexes) distinct. On the right side of [Figure 1](#), three species form tight clusters when projected onto the “species plane,” while their intra-species variations fan out along axes orthogonal to the plane. Such variation cues are not encoded in taxonomic labels. We theoretically prove that when species prototypes are nearly orthogonal (*i.e.*, species are well separated), the contrastive objective prioritizes orthogonality between intra-species variations and inter-species differences over raw magnitude ([Theorem 5.1](#)). Furthermore, these variations are observed to be increasingly separable as training data scales up. This microstructure preserves the intra-species representational diversity without interfering with inter-species distinctions, enabling various attribute recognition applications ([§5.1](#)).

We show in [§5.2](#) through quantitative and qualitative analyses that larger-scale training improves both inter-species ecological alignment and separation of intra-species variants. These scale-amplified patterns make the embedding space more interpretable and biologically meaningful. **BIOCLIP 2** evidences that combining domain-specific scaling with structured supervision can unlock qualitatively new emergent behaviors in scientific vision models.

2 Related Work

Emergent properties in foundation models. Emergent properties refer to the capabilities implicitly acquired from the training process and generalized beyond the initial training objective. Large language models (LLMs) illustrate a variety of in-context learning skills after next-token-prediction

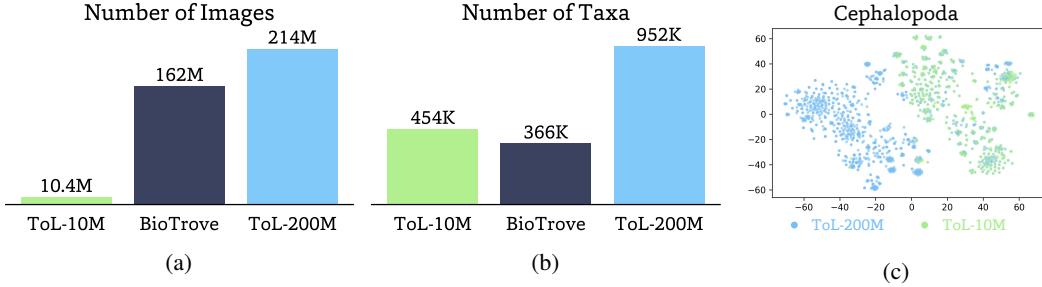


Figure 2: (a) Number of images across organismal biology datasets. (b) Biodiversity comparison across datasets (measured unique 7-tuples for TREEOFLIFE (ToL) datasets, species count provided by BioTrove). (c) The taxa distributional difference in the *Cephalopoda* class (octopuses, squids, etc.) between ToL-200M and ToL-10M.

pre-training [9, 26]. Such emergence also arises in computer vision [27]. DINO learns semantic segmentation through purely visual self-supervision [28], whereas GroupViT learns semantic segmentation solely through text supervision [29]. The studies closest to this work are [30, 31]. Alper and Averbuch-Elor explore the visual-semantic hierarchies in CLIP models [30]. Although hierarchical semantic structures are never explicitly presented as the training supervision, CLIP models acquire the capability of matching images with varying levels of descriptions. Abbasi *et al.* discover that CLIP models possess disentangled representational sub-spaces for different factors of variations [31]. It allows CLIP models to generalize across compositional out-of-distribution concepts. This work leverages CLIP to distinguish different species with hierarchical labels, which is a different scenario from the above work. We investigate the emergent properties under this setting.

Computer vision for ecology & evolutionary biology. Ecology and evolutionary biology are naturally challenging for computer vision systems due to long-tail distributions, extremely fine-grained classifications, and a wide variety of image distributions. Existing work formalizes these challenges into specific visual tasks such as attribute prediction [NeWT, 18], trait prediction [FishNet, 16], and plant disease detection [PlantDoc, 20].

Recent advancements in computer vision have led to the development of foundation models for biological applications. BIOCLIP incorporates taxonomic labels in the vision-language contrastive training, yielding promising species classification accuracy [13]. Follow-up work scaled data to 162M images [BioTrove, 14], specialized the data to camera traps [CATALOG and WildCLIP, 32, 33], and added additional model modalities [TaxaBind, 34]. We investigate both data and model scaling, with a focus on both broad biological applications and any emergent properties after extensive training.

3 TREEOFLIFE-200M

Large-scale, clean, diverse data drives progress in machine learning. There have been efforts such as TREEOFLIFE-10M [35] and BioTrove [14] to create large-scale biological organismal datasets for machine learning (ML). As shown in Figure 2, TREEOFLIFE-10M [35] improves on prior work such as iNat21 [36] and BIOSCAN-1M [37] by increasing taxa diversity by a factor of 45. BioTrove [14] increases data scale to 162M but fails to match the biodiversity of TREEOFLIFE-10M. In this work, we combine the vast breadth of Global Biodiversity Information Facility (GBIF) [38] images with those of the Encyclopedia of Life project (EOL) [39], BIOSCAN-5M [40], and FathomNet Database [41]. With nearly 214 million images representing 952,257 taxa, **TREEOFLIFE-200M** is the largest *and* most diverse public ML-ready dataset for computer vision models in biology.

Unlike BioTrove [14], which relies solely on iNaturalist and contains 162M images but only 366K unique species, our use of museum, camera-trap, and citizen-science contributions expands the taxonomic coverage to 2.6× more taxa. Our curation efforts ensure this breadth does not come at the cost of data quality. We quantify image type diversity from GBIF: 51.8M museum specimen, 617.8K camera trap, and 151M citizen science images. Beyond the taxa-wise diversity, these images also provide more observing perspectives for the focal species. In such a way, the robustness of models trained on **TREEOFLIFE-200M** is significantly enhanced against a variety of use cases. Specifically, we demonstrate that BIOCLIP 2 yields a 22.8% performance gap compared with BIOCLIP on camera trap images (See Table 1). We provide detailed statistics of the image number and taxa diversity from each data provider in Figure 11.

3.1 Images

TREEOFLIFE-200M consists of images curated from four core data providers: GBIF, EOL, BIOSCAN-5M, and FathomNet. GBIF and EOL aggregate biological data from various sources, such as iNaturalist [42], the Smithsonian Institution [43], and Flickr [44]. BIOSCAN-5M and FathomNet are curated collections of expert-annotated images designed to improve cataloging and identification of species from highly diverse and under-represented branches of the tree of life. BIOSCAN-5M is part of the ongoing DNA Barcoding project to improve insect identification (one of the most diverse classes, *Insecta*). FathomNet focuses on a habitat rather than a clade: all animals that live in the ocean. Together, these comprise TREEOFLIFE-200M.

3.2 Data Curation & Filtering

We retrieve data from data providers using [distributed-downloader](#) [45]. The initial retrieval gives us 222,065,140 images and 1,359,405 unique taxonomic hierarchies. We then cleaned the data, focusing on (1) aligning taxonomic labels, (2) image quality, and (3) eliminating duplication and data leakage. We summarize these efforts here, with more details provided in §I.

Taxonomic alignment. Curating large biological image collections from distributed data providers requires the alignment of noisy and inconsistent taxonomic labels both between and within providers. We develop a taxonomic alignment package [TaxonoPy](#) in consultation with taxonomists that resolves entries to both a seven-rank Linnaean hierarchy and a common name. TaxonoPy queries GNVerifier [46] against the GBIF Backbone, Catalogue of Life, and OpenTree hierarchies (in that order). From the original 1.36M taxa, TaxonoPy filters 407K taxonomic hierarchies (such as synonyms, provisional names, etc.) and yields 952K unique taxa.

Image-quality screens. Digital archives contain herbarium labels, empty camera-trap frames, and occasional people. None add biological signal, and faces raise privacy concerns, so we drop them to keep learning focused on organisms via three neural-network-based filters: (i) **Museum non-organism removal.** A pre-trained CLIP-L/14 [12] is used for a nearest-centroid classifier spanning 25 fine-grained subtypes (10 collection areas, each split into categories such as specimen, fossil, drawer-label, etc.). The classifier is fit to 8.5K manually-curated examples and predicts a subtype for all museum images; we drop all non-organismal images. (ii) **Camera-trap trimming.** We apply a pre-trained camera-trap model MegaDetector [47, 48] to filter for frames with visible animals. (iii) **Face removal.** We apply a pre-trained face-detection model MTCNN [49] to discard images containing human faces. We release the code used in processing the images in [TreeOfLife-toolbox](#).

Duplicate and leakage control. Exact duplicates in the training set are removed with MD5 hashes. GBIF includes images from iNaturalist, a popular source for computer vision ecology benchmarks. To prevent train-test leakage and inflated downstream scores, we compute both MD5 and perceptual PDQ [50] hashes for every test image and purge any near or exact duplicates from training.

3.3 Taxa Coverage

The International Union for Conservation of Nature (IUCN) estimates 2.14M species have been described [51]. Following curation, there are nearly 868K unique taxa labeled to the level of species in TREEOFLIFE-200M.¹ Based on the most recent IUCN Red List assessment [52], TREEOFLIFE-200M demonstrates a particularly strong representation of threatened species, with 77.1% coverage (36,370 species). This coverage establishes that the approach to integrating diverse data sources used in TREEOFLIFE-200M is a valuable resource for conservation research, providing representation for a substantial majority of species prioritized for global conservation action.

Notably, TREEOFLIFE-200M adds diverse clades that are extremely under-represented in prior work like TREEOFLIFE-10M and BioTrove. [Figure 2c](#) compares the distribution of taxonomic names in the class of *Cephalopoda* between TREEOFLIFE-10M and TREEOFLIFE-200M. While they share overlaps, there are clades almost completely absent in TREEOFLIFE-10M. These under-represented clades receive a substantial influx of samples in TREEOFLIFE-200M (1,102 new taxa). Another

¹Not all images contain full 7-rank Linnaean taxa; for instance, 93% of BIOSCAN-5M images are not labeled to the species level. Thus, the unique taxa with non-null species is a more appropriate comparison.

Table 1: Zero-, one-, and five-shot species classification top-1 accuracy across 10 tasks for different models. **Bold** and underlined entries indicate the **best** and second best accuracies, respectively. BIOCLIP 2 outperforms both strong general- (CLIP, SigLIP, DINOv3) and domain-specific- (BIOCLIP, BioTrove-CLIP) baselines. “Camera Trap” is mean performance across 5 camera-trap datasets; Appendix F.4 contains more details.

Model	Animals					Plants & Fungi					Rare Species	Mean
	NABirds	Plankton	Insects	Insects 2	Camera Trap	PlantNet	Fungi	PlantVillage	Med. Leaf			
Random Guessing	0.2	1.2	1.0	1.0	3.5	4.0	4.0	2.6	4.0	0.3	2.2	
<i>Zero-Shot Classification</i>												
CLIP (ViT-L/14)	66.5	1.3	9.0	11.7	29.5	61.7	7.6	6.5	25.6	35.2	25.5	
SigLIP	61.7	2.4	27.3	<u>20.7</u>	<u>33.7</u>	81.8	36.9	28.5	<u>54.5</u>	47.6	39.5	
BioTrove-CLIP	39.4	1.0	20.5	15.7	10.7	64.4	38.2	15.7	31.6	24.6	26.2	
BIOCLIP	58.8	6.1	<u>34.9</u>	20.5	31.7	<u>88.2</u>	<u>40.9</u>	19.0	38.5	37.1	37.6	
BIOCLIP 2	74.9	<u>3.9</u>	55.3	27.7	53.9	96.8	83.8	<u>25.1</u>	57.8	76.8	55.6	
<i>One-Shot Classification</i>												
CLIP (ViT-L/14)	42.7 \pm 0.8	28.9 \pm 0.6	29.0 \pm 0.4	17.0 \pm 0.8	36.0 \pm 2.8	58.7 \pm 2.8	20.7 \pm 2.2	56.7 \pm 2.1	74.4 \pm 1.9	34.3 \pm 0.8	39.8	
SigLIP	39.9 \pm 0.9	28.4 \pm 0.5	32.3 \pm 0.6	20.6 \pm 1.3	37.8 \pm 2.6	66.3 \pm 3.3	28.7 \pm 0.9	64.1 \pm 3.0	81.7 \pm 2.3	38.8 \pm 0.7	43.9	
Supervised-IN21K	43.8 \pm 0.8	23.5 \pm 1.2	15.2 \pm 1.1	18.2 \pm 1.5	30.6 \pm 2.4	63.8 \pm 4.4	26.4 \pm 1.4	52.8 \pm 3.5	75.2 \pm 4.4	31.6 \pm 0.7	38.1	
DINOv3	48.3 \pm 1.1	36.5 \pm 0.8	8.8 \pm 0.7	18.8 \pm 1.6	<u>42.6</u> \pm 2.4	66.8 \pm 4.7	27.5 \pm 1.8	<u>64.8</u> \pm 1.2	92.1 \pm 2.2	41.5 \pm 0.2	44.8	
BioTrove-CLIP	61.9 \pm 0.6	26.4 \pm 0.5	<u>57.1</u> \pm 1.4	<u>20.9</u> \pm 0.7	31.2 \pm 2.3	<u>69.7</u> \pm 3.4	<u>47.3</u> \pm 2.1	55.8 \pm 3.4	83.5 \pm 1.1	34.9 \pm 0.4	48.9	
BIOCLIP	57.4 \pm 1.2	29.7 \pm 1.1	<u>57.1</u> \pm 1.0	20.4 \pm 0.9	35.0 \pm 2.8	67.7 \pm 3.9	44.6 \pm 2.0	59.5 \pm 2.5	83.7 \pm 1.8	<u>44.9</u> \pm 0.7	50.0	
BIOCLIP 2	82.4 \pm 1.1	<u>32.0</u> \pm 0.4	74.6 \pm 0.4	<u>28.4</u> \pm 0.7	<u>48.1</u> \pm 2.2	85.8 \pm 4.5	70.3 \pm 2.6	67.6 \pm 1.1	<u>92.0</u> \pm 1.9	59.5 \pm 0.9	64.1	
<i>Five-Shot Classification</i>												
CLIP (ViT-L/14)	68.2 \pm 0.3	48.2 \pm 1.5	50.6 \pm 0.7	30.1 \pm 0.7	53.9 \pm 2.2	75.9 \pm 1.2	31.4 \pm 2.5	78.3 \pm 1.4	92.6 \pm 0.7	53.3 \pm 0.4	58.3	
SigLIP	64.2 \pm 0.3	47.4 \pm 1.1	54.9 \pm 0.7	<u>35.2</u> \pm 0.5	56.9 \pm 2.0	81.6 \pm 1.4	45.5 \pm 1.6	81.1 \pm 0.7	94.1 \pm 0.7	57.8 \pm 0.6	61.9	
Supervised-IN21K	57.5 \pm 0.4	40.1 \pm 0.6	30.1 \pm 0.7	30.3 \pm 0.2	48.0 \pm 2.2	77.2 \pm 1.4	39.6 \pm 1.9	78.0 \pm 1.1	92.8 \pm 0.9	48.8 \pm 0.4	54.2	
DINOv3	72.3 \pm 0.4	57.8 \pm 1.6	19.7 \pm 0.7	34.2 \pm 0.6	<u>63.3</u> \pm 2.8	83.3 \pm 1.4	44.1 \pm 1.6	<u>82.4</u> \pm 1.1	98.8 \pm 0.5	62.6 \pm 0.5	61.8	
BioTrove-CLIP	<u>78.5</u> \pm 0.2	44.6 \pm 0.6	77.0 \pm 0.8	34.2 \pm 0.6	47.9 \pm 2.0	<u>86.0</u> \pm 1.0	<u>65.2</u> \pm 0.8	75.1 \pm 0.8	96.2 \pm 0.7	51.3 \pm 0.2	65.6	
BIOCLIP	78.2 \pm 0.3	49.2 \pm 1.1	<u>78.0</u> \pm 0.6	33.9 \pm 0.6	54.3 \pm 2.2	85.7 \pm 1.7	61.6 \pm 1.9	81.7 \pm 1.1	96.7 \pm 0.6	<u>65.7</u> \pm 0.4	68.5	
BIOCLIP 2	92.4 \pm 0.2	<u>50.5</u> \pm 1.1	89.3 \pm 0.4	<u>44.3</u> \pm 1.1	<u>67.7</u> \pm 1.9	<u>94.4</u> \pm 0.8	<u>85.0</u> \pm 1.1	83.9 \pm 0.9	<u>98.4</u> \pm 0.4	77.2 \pm 0.4	78.3	

example is 55,085 taxa of *Fungi* in TREEOF LIFE-200M, close to 4 \times of that in TREEOF LIFE-10M (14,793). The improved diversity facilitates accurate species classification of these clades, as evidenced in Table 1 (42.9% absolute improvement over BIOCLIP on zero-shot Fungi benchmark).

4 BIOCLIP 2 and Species Classification

BIOCLIP adopts a hierarchical multi-modal contrastive training framework, where images are associated with their corresponding hierarchical labels including taxonomic labels, scientific names, and common names [13]. Different from one-hot labels, taxonomic labels inherently encode hierarchical biological information from different levels [15]. In combination with an auto-regressive text encoder, BIOCLIP yielded superior species classification performance on both zero- and few-shot settings. In this work, we stick with the hierarchical contrastive training recipe and focus on the impact of scale.

Modifications. In addition to the significantly larger and more diverse dataset, we also scale model capacity by adopting a larger vision transformer (ViT-L/14 pre-trained on LAION-2B [12, 53, 54]). An auxiliary replay mechanism is introduced [55, 56] to maintain general-domain understanding for broader applications [33, 57]; a portion of CLIP training data (LAION-2B) is interleaved simultaneously with species contrastive learning. We ablate this decision and find that the experience replay improves biological understanding and performance across diverse tasks in §6.

We train BIOCLIP 2 on 32 NVIDIA H100 GPUs for 10 days on 214M organismal biology images with hierarchical labels and 26M randomly-sampled image-text pairs from LAION-2B for 30 epochs. We provide the training details in §D.

4.1 Species Classification Performance

We evaluate BIOCLIP 2 on species classification tasks in Table 1. We use the same benchmarks as BioCLIP [13], including seven tasks from Meta-Album [58] and Rare Species [59]. We substitute

Table 2: Biological visual tasks beyond species classification. **Bold** and underlined entries indicate the **best** and second best accuracies. See §H for task and evaluation methodology details.

Model	Animals			Plants		Mean
	FishNet	NeWT	AwA2	Herb. 19	PlantDoc	
CLIP (ViT-L/14)	27.9 \pm 0.2	83.4 \pm 0.1	61.6 \pm 0.6	18.2 \pm 0.1	22.3 \pm 3.3	42.7
SigLIP	31.9 \pm 0.1	83.2 \pm 0.1	<u>67.3</u> \pm 0.6	18.6 \pm 0.2	28.2 \pm 5.3	45.8
Supervised-IN21K	29.4 \pm 0.1	75.8 \pm 0.2	52.7 \pm 1.6	14.9 \pm 0.1	25.1 \pm 1.1	39.6
DINOv3	37.9 \pm 0.1	<u>85.7</u> \pm 0.0	48.0 \pm 2.8	<u>31.2</u> \pm 0.2	40.3 \pm 1.2	48.6
BioTrove-CLIP	22.1 \pm 0.0	82.5 \pm 0.1	45.7 \pm 0.7	20.4 \pm 0.2	37.7 \pm 1.2	41.7
BIOCLIP	30.1 \pm 0.2	82.7 \pm 0.1	65.9 \pm 0.3	26.8 \pm 0.4	39.5 \pm 2.3	<u>49.0</u>
BIOCLIP 2	39.8 \pm 0.4	89.1 \pm 0.1	69.5 \pm 1.1	48.6 \pm 0.6	40.4 \pm 3.7	57.5

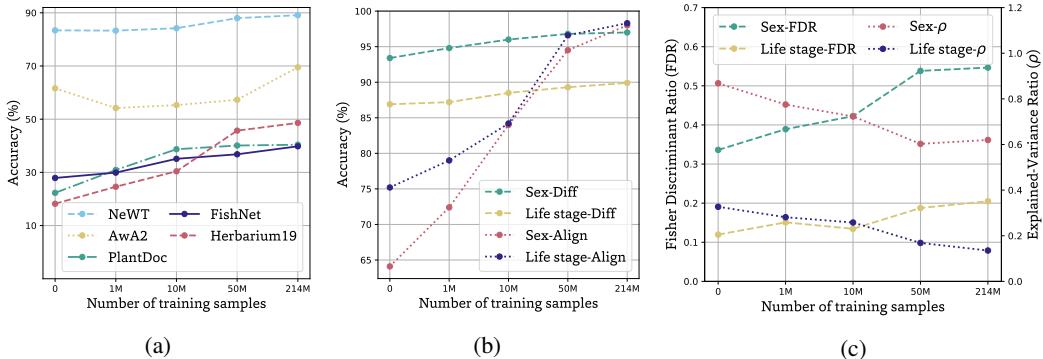


Figure 3: (a) The model performance on five downstream tasks under different scales of training data. (b) The model performance on differentiating and aligning different life stages and sexes. (c) The separation and orthogonality evaluation of models trained with different amounts of data.

Birds-525 [60] with NABirds [61] due to data inaccessibility. Additionally, we collect a test set of IDLE-OO Camera Traps from the Labeled Information Library of Alexandria: Biology and Conservation (LILA-BC) [62, 63, 64, 65, 66, 67] to illustrate more realistic species classification applications in the wild. The results suggest substantial improvements of BIOCLIP 2 over BIOCLIP. Particularly, attributed to more comprehensive species and image type coverage of TREEOFLIFE-200M, we observe over 20% zero-shot improvement on Camera Trap, Fungi, and Rare Species. On average, BIOCLIP 2 surpasses the second-best model by 16.1% and provides a 30.1% improvement over the original CLIP model that serves as weight initialization. Information on baselines is in §E.

5 Emergent Properties from Scaling Hierarchical Contrastive Learning

5.1 Beyond Species Classification

Biology’s organization extends beyond species taxonomies; if scaling truly induces emergent behavior, model representations learned through species-level supervision should transfer to problems far removed from species classification. We collect and benchmark models on five visual benchmarks that push past species ID: habitat classification (ecological context) [16], trait prediction (evolutionary studies) [18, 17], new-species identification (biodiversity monitoring) [19], and agricultural disease detection [20]. For each task, we keep the evaluated models frozen and extract the corresponding sample embeddings. The embeddings are subsequently processed using machine-learning techniques (*e.g.*, support vector classifiers). Detailed evaluation procedures are listed in §H.

Table 2 presents the performance comparison among BIOCLIP 2, vision-language baselines, and vision-only models. Although no information on these tasks is explicitly described during training, BIOCLIP 2 yields an average performance improvement of 14.8% over the original CLIP baseline. DINOv3 is commonly believed to capture fine-grained visual features and is adopted for diverse visual tasks [22, 68, 69]. Nevertheless, BIOCLIP 2 yields an 8.9% performance gap over DINOv3.

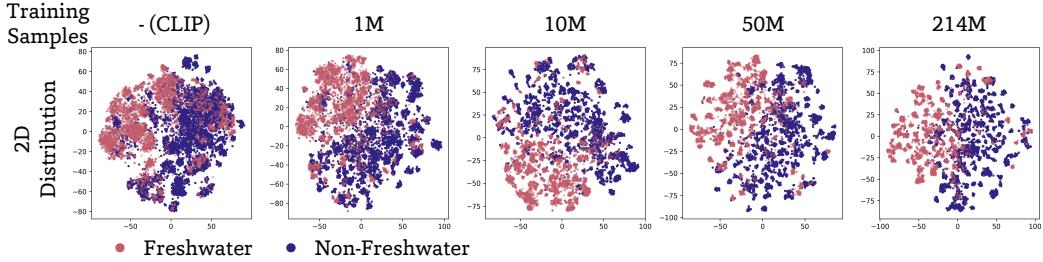


Figure 4: t-SNE embedding visualization of FishNet test set for models trained with different amounts of data. The leftmost plot is the original LAION-2B CLIP ViT-L/14. As the training data scales, freshwater fish become more distinct from saltwater fish and brackish fish, despite no explicit supervision, demonstrating that data scale contributes to emergent properties in model representations.

5.2 Scaling Trends

Better species classification with more species-labeled data is expected, but its effect on other, non-species classification tasks is unexpected. To better understand the relationship between training data scale and non-species classification performance, we apply the same hierarchical training with varying sizes of data: 1M, 10M, 50M, and 214M samples. The smaller datasets are randomly sampled from the complete set without losing taxa representativeness. We compare the performance of the baseline CLIP ViT-L/14 and the four obtained models on five non-species classification tasks in [Figure 3a](#). A consistent improvement is observed as the volume of training data increases from 1M to 214M. In AwA2, for example, the 1M model is worse than the baseline. However, the model gradually learns more generalized representations for different attributes across species and obtains improved performance as data scales up.

Next, we investigate how scale affects representations within species. We collect two groups of images with intra-species appearance variations: life stage variations from NeWT [18] and sex variations from NABirds [61]. We ask whether scaling hierarchical contrastive training collapses all images of one species onto a single prototype or still distinguishes juveniles from adults and males from females. We accordingly design two complementary tasks for each type of variation: (i) *alignment*, where a species classifier trained on one variant (*e.g.*, juvenile images) is expected to recognize the species on the other variant (*e.g.*, adult images), and (ii) *differentiation*, where the task is to tell the variants apart (*e.g.*, juvenile vs. adult). As illustrated in [Figure 3b](#), data scale steadily improves cross-variant species recognition. But at the same time, the model also becomes *better* at distinguishing the variants themselves.

5.3 Emergent Properties and Qualitative Analysis

Why does scaling data boost tasks that are never supervised during hierarchical contrastive training? We look deeper into BIOCLIP 2’s embedding space and identify two emergent properties that generalize beyond species classification.

First, the embedding distribution of different species *aligns with their ecological and functional relationships*. [Figure 4](#) shows t-SNE plots [70] of FishNet test set embeddings at four training scales, colored by whether the fish can live in freshwater or not. In the baseline CLIP plot (left), freshwater and non-freshwater fish have a large portion of overlap. Larger training sets progressively separate the two groups in the embedding space. We note that there is no explicit constraint to arrange meaningful distribution across species in contrastive loss, highlighting the emergence at scale.

Second, the intra-species variations are *preserved and separated*. [Figure 5](#) shows t-SNE plots of BIOCLIP 2 embeddings of three species from NeWT exhibiting life-stage variation. In the 2D plots (top), different species (shown in different colors) tighten progressively from left to right, which is a direct consequence of scaled contrastive training. At the same time, the intra-species variations are preserved and better separated than the baseline CLIP model (leftmost sub-figure). We further project the embeddings onto 3D spaces created by singular value decomposition (SVD, bottom), which reveals that the intra-species variations lie in subspaces roughly orthogonal to the species span. Therefore, the existence of intra-species variation does not interfere with inter-species distinctions. §F contains more empirical observations.

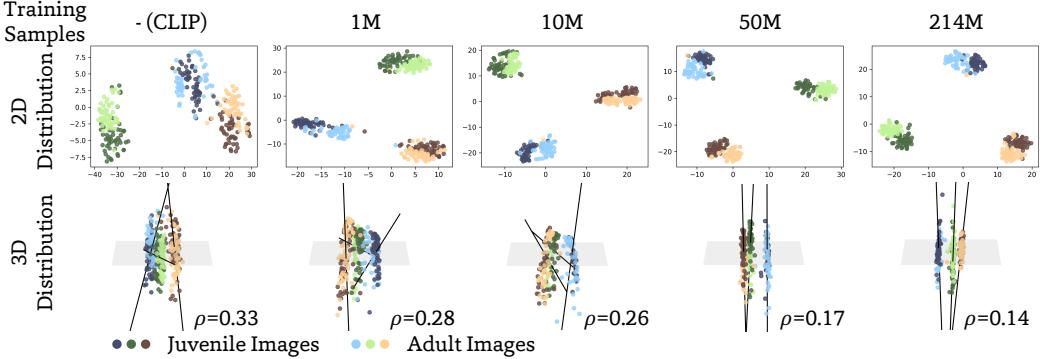


Figure 5: The embedding distribution of life stage variations under different scales of training data. The 2D distributions are obtained using t-SNE. For the 3D distributions, we first run SVD with the mean embedding of each species. The first two singular vectors are used to construct the gray plane that captures most inter-species differences. The embeddings are then projected into the 3D space with an additional orthogonal dimension. The straight lines point from the mean embedding of juvenile images to that of the adult images. As the training scales up, the intra-species variations are preserved in the subspace orthogonal to the inter-species differences. Orthogonality improves with data scale, as evidenced by the decreasing explained-variance ratio ρ .

5.4 Formal Analysis

We next ask *why* these two properties of inter-species ecological alignment and intra-species variation separation emerge with scale.

Inter-species ecological alignment. As training data increases, species that share proximal taxonomic labels are pulled toward common textual prototypes at multiple levels. As related taxa typically share morphology, behavior, and ecological characteristics, this multi-level supervision aligns visual similarity with functional similarity [15]. With more samples per species providing supervision at scale, embeddings of species in the same family or genus form coherent macro-clusters. In effect, the embeddings extracted by BioCLIP 2 are more separable for different ecological groups.

Intra-species variation separation. While the inter-species ecological alignment can be explained by hierarchical supervision, the intra-species variations are *not* encoded in taxonomic labels. The preservation of intra-species variations also contradicts the common intuition of contrastive training effects. Therefore, we investigate the optimization of contrastive loss. We propose the following theorem to suggest that subspace orthogonal to inter-species differences is allowed after extensive training to accommodate intra-species variations.

Theorem 5.1. *Let μ be the prototypes of species, with μ_s as the prototype of species s . Let τ be temperature. If different μ_k are nearly orthogonal (i.e., species are well separated), the intra-species variation δ for species s is constrained by*

$$\delta^\top \left[\frac{1}{2\tau^2} \left(\sum_k w_k \mu_k \mu_k^\top - \mu_s \mu_s^\top \right) \right] \delta, \quad \text{where } w_k = \frac{\exp(\mu_s^\top \mu_k)/\tau}{\sum_k (\exp(\mu_s^\top \mu_k)/\tau)}.$$

Proof. See §C. □

Thus, as long as the variation δ is distributed in a subspace orthogonal to the inter-species distinctions, the scale of δ won't interfere with the overall contrastive optimization. The orthogonality is qualitatively supported by Figure 5. To further quantify it, we calculate the explained-variance ratio, *i.e.*, the ratio in the intra-species variation that is captured by the species span [71]. We first obtain an orthonormal basis for the species prototypes \mathbf{U} using QR decomposition. Let $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_n] \in \mathbb{R}^{d \times n}$ be the matrix stacking n intra-species variation difference vectors with dimension d . The explained-variance ratio calculates the energy fraction inside species span by $\rho = \|\mathbf{U}^\top \mathbf{D}\|_F^2 / \|\mathbf{D}\|_F^2$, where $\|\cdot\|_F$ denotes the Frobenius norm [72]. We show the ratio change as the data scales up in Figure 3c. The results suggest that the intra-species variations are increasingly orthogonal to the species differences. Due to the orthogonality, the existence of intra-species variations will not interfere with the inter-species distinctions. The observation of smaller projection areas in Figure 5 at the species span also indicates better species classification accuracy after extensive contrastive training.

Table 3: Ablation study with different training settings. BIOCLIP 2 adopts hierarchical contrastive training with taxonomic labels. We ablate using scientific names solely without taxonomic labels, and one-hot labels with cross-entropy loss instead of the contrastive objective.

Dataset	Hierarchical Contrastive	Contrastive w/ Scientific Name	Cross-entropy Loss
FishNet	35.1 \pm 0.1	33.8 \pm 0.1	33.0 \pm 0.1
PlantDoc	38.7 \pm 3.7	37.3 \pm 3.3	30.9 \pm 1.9
Life stage-Diff	88.0 \pm 0.1	88.5 \pm 0.2	85.5 \pm 0.2
Life stage-Align	84.1 \pm 0.1	84.5 \pm 0.1	78.6 \pm 0.1
Sex-Diff	97.0 \pm 0.1	96.6 \pm 0.1	95.5 \pm 0.1
Sex-Align	84.1 \pm 0.2	82.7 \pm 0.3	74.9 \pm 0.2

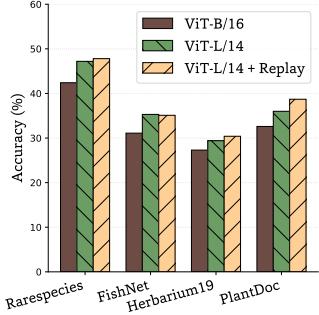


Figure 6: Ablation study on model size and experience replay.

Furthermore, we look at the separation between these variations using the Fisher Discriminant Ratio (FDR) metric [73]. Given two variant classes A and B , FDR is defined by their embedding mean μ and standard deviation σ : $FDR = \|\mu_A - \mu_B\|^2 / (\sigma_A^2 + \sigma_B^2)$. The increasing trends in Figure 3c support that, in addition to orthogonality, the intra-species variations are more separable as the training scales up. This trend is also qualitatively evidenced by the top row of Figure 5. These theoretical and empirical insights validate that BIOCLIP 2 learns to preserve and separate intra-species variations without explicit training constraints. We provide more qualitative analyses in §F.2.

6 Ablation Study and Analysis

The necessity of contrastive loss and taxonomic labels. The previous analyses highlight the effectiveness of hierarchical contrastive training. Table 3 ablates two key modeling decisions using TREEOFLIFE-10M as a feasible testbed: scientific names vs. hierarchical labels and contrastive learning vs. cross-entropy loss on one-hot labels. The experiments are conducted on TREEOFLIFE-10M data with ViT-L/14 as the visual encoder. When trained solely with scientific names, the model loses some of the hierarchical supervision embedded in taxonomic labels. As a result, the performance on FishNet drops by 1.3%. However, contrastive training still leads to separation for both species and intra-species variations. Training a 952K-class softmax classifier with cross-entropy loss is optimization-heavy and leads to inferior performance on all benchmarks. The adopted hierarchical contrastive supervision leverages the advantages of both aspects and yields the best overall performance across benchmarks.

Architecture and replay. BIOCLIP 2 scales up the visual encoder of BIOCLIP from a ViT-B/16 to ViT-L/14 and introduces experience replay of CLIP training data (LAION-2B). We ablate the effects of these two changes, again using TREEOFLIFE-10M as a testbed. Figure 6 shows that increasing model capacity improves performance across all benchmarks. Comparatively, experience replay leads to better species classification accuracy and improved performance on some of the other visual tasks. We provide a more detailed empirical study of experience replay in §F.1.

7 Conclusion

In this work, we curate TREEOFLIFE-200M, the largest and most diverse biological organism dataset to date, and train BIOCLIP 2 with hierarchical taxonomic labels. BIOCLIP 2 achieves state-of-the-art accuracy on species classification. More importantly, large-scale training gives rise to two emergent properties not described during training. At the inter-species level, the embedding distribution of different species aligns with their ecological relationships. At the intra-species level, the appearance variations within species are preserved and well separated in the embedding space. We demonstrate that combining the effort of domain-specific scaling and structured supervision leads to effective generalization beyond the initial training objectives. BIOCLIP 2 serves as a strong foundation model for biological research and simultaneously evidences the effectiveness of scale-driven scientific discovery.

Acknowledgments and Disclosure of Funding

We would like to thank Zhiyuan Tao, Shuheng Wang, Ziheng Zhang, Zhongwei Wang, and Leanna House for their help with the TREEOFLIFE-200M dataset, Charles (Chuck) Stewart, Sara Beery, and other **Imageomics Team** members for their constructive feedback and Sergiu Sanielevici, Tom Maiden, and TJ Olesky for their dedicated assistance with arranging the necessary computational resources.

We are grateful to Kakani Katija and Dirk Steinke for helpful conversations regarding use and integration of FathomNet and BIOSCAN-5M, respectively, as well as Stephen Formel and Markus Döring for GBIF. We thank Marie Grosjean for comparative methods for filtering citizen science images and Dylan Verheul for assistance with acquiring images from observation.org from GBIF. We thank Suren Byna for a helpful conversation on early dataset design decisions. We thank Doug Johnson for his collaboration in hosting this large dataset on the Ohio Supercomputer Center research storage file system.

Our research is supported by NSF OAC 2118240 and resources from the Ohio Supercomputer Center [74]. This work used the Bridges-2 system, which is supported by NSF award number OAC-1928147 at the Pittsburgh Supercomputing Center (PSC) [75], under the auspices of the NAIRR Pilot program.

References

- [1] Hanchen Wang, Tianfan Fu, Yuanqi Du, Wenhao Gao, Kexin Huang, Ziming Liu, Payal Chandak, Shengchao Liu, Peter Van Katwyk, Andreea Deac, et al. Scientific discovery in the age of artificial intelligence. *Nature*, 620(7972):47–60, 2023. [2](#)
- [2] Yongjun Xu, Xin Liu, Xin Cao, Changping Huang, Enke Liu, Sen Qian, Xingchen Liu, Yanjun Wu, Fengliang Dong, Cheng-Wei Qiu, et al. Artificial intelligence: A powerful paradigm for scientific research. *The Innovation*, 2(4), 2021. [2](#)
- [3] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021. [2](#)
- [4] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023. [2](#)
- [5] Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gowoon Cheon, and Ekin Dogus Cubuk. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, 2023. [2](#)
- [6] Huan Tran, Rishi Gurnani, Chihio Kim, Ghanshyam Pilania, Ha-Kyung Kwon, Ryan P Lively, and Rampi Ramprasad. Design of functional and sustainable polymers assisted by artificial intelligence. *Nature Reviews Materials*, pages 1–21, 2024. [2](#)
- [7] Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Willson, Peter Wirnsberger, Meire Fortunato, Ferran Alet, Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen, Weihua Hu, et al. Learning skillful medium-range global weather forecasting. *Science*, 382(6677):1416–1421, 2023. [2, 25](#)
- [8] Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen, Xiaotao Gu, and Qi Tian. Accurate medium-range global weather forecasting with 3d neural networks. *Nature*, 619(7970):533–538, 2023. [2](#)
- [9] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *TMLR*, pages 1–30, 2022. [2, 3](#)
- [10] Daniil A Boiko, Robert MacKnight, Ben Kline, and Gabe Gomes. Autonomous chemical research with large language models. *Nature*, 624(7992):570–578, 2023. [2](#)
- [11] Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle Ott, C Lawrence Zitnick, Jerry Ma, et al. Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of Sciences*, 118(15):e2016239118, 2021. [2](#)

- [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763. PMLR, 2021. 2, 4, 5, 17, 18
- [13] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song, David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, et al. Bioclip: A vision foundation model for the tree of life. In *CVPR*, pages 19412–19424, 2024. 2, 3, 5, 20, 26
- [14] 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. In *NeurIPS*, volume 37, pages 102101–102120, 2024. 2, 3, 18, 24
- [15] José M Padial, Aurélien Miralles, Ignacio De la Riva, and Miguel Vences. The integrative future of taxonomy. *Frontiers in zoology*, 7:1–14, 2010. 2, 5, 8
- [16] Faizan Farooq Khan, Xiang Li, Andrew J Temple, and Mohamed Elhoseiny. Fishnet: A large-scale dataset and benchmark for fish recognition, detection, and functional trait prediction. In *ICCV*, pages 20496–20506, 2023. 2, 3, 6, 22
- [17] Yongqin Xian, Christoph H Lampert, Bernt Schiele, and Zeynep Akata. Zero-shot learning—a comprehensive evaluation of the good, the bad and the ugly. *IEEE TPAMI*, 41(9):2251–2265, 2018. 2, 6, 19, 23
- [18] Grant Van Horn, Elijah Cole, Sara Beery, Kimberly Wilber, Serge Belongie, and Oisin Mac Aodha. Benchmarking representation learning for natural world image collections. In *CVPR*, pages 12884–12893, 2021. 2, 3, 6, 7, 23, 26
- [19] Kiat Chuan Tan and Yulong Liu. Herbarium challenge 2019 - fgvc6. <https://kaggle.com/competitions/herbarium-2019-fgvc6>, 2019. Kaggle. 2, 6, 18, 23
- [20] Davinder Singh, Naman Jain, Pranjali Jain, Pratik Kayal, Sudhakar Kumawat, and Nipun Batra. Plantdoc: A dataset for visual plant disease detection. In *Proceedings of the 7th ACM IKDD CoDS and 25th COMAD*, pages 249–253, 2020. 2, 3, 6, 18, 23
- [21] Xiaohua Zhai, Basil Mustafa, Alexander Kolesnikov, and Lucas Beyer. Sigmoid loss for language image pre-training. In *ICCV*, pages 11975–11986, 2023. 2, 18
- [22] Oriane Siméoni, Huy V Vo, Maximilian Seitzer, Federico Baldassarre, Maxime Oquab, Cijo Jose, Vasil Khalidov, Marc Szafraniec, Seungeun Yi, Michaël Ramamonjisoa, et al. Dinov3. *arXiv preprint arXiv:2508.10104*, 2025. 2, 6, 18
- [23] Sangeet Lamichhaney, Jonas Berglund, Markus Sällman Almén, Khurram Maqbool, Manfred Grabherr, Alvaro Martinez-Barrio, Marta Promerová, Carl-Johan Rubin, Chao Wang, Neda Zamani, et al. Evolution of darwin’s finches and their beaks revealed by genome sequencing. *Nature*, 518(7539):371–375, 2015. 2, 19
- [24] Jinxin Zhou, Chong You, Xiao Li, Kangning Liu, Sheng Liu, Qing Qu, and Zhihui Zhu. Are all losses created equal: A neural collapse perspective. In *NeurIPS*, volume 35, pages 31697–31710, 2022. 2, 21
- [25] Jianfeng Lu and Stefan Steinerberger. Neural collapse under cross-entropy loss. *Applied and Computational Harmonic Analysis*, 59:224–241, 2022. 2, 21
- [26] Leonardo Berti, Flavio Giorgi, and Gjergji Kasneci. Emergent abilities in large language models: A survey. *arXiv preprint arXiv:2503.05788*, 2025. 3
- [27] Zheda Mai, Arpita Chowdhury, Zihe Wang, Sooyoung Jeon, Lemeng Wang, Jiacheng Hou, and Wei-Lun Chao. Ava-bench: Atomic visual ability benchmark for vision foundation models. *arXiv preprint arXiv:2506.09082*, 2025. 3
- [28] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 3
- [29] Jiarui Xu, Shalini De Mello, Sifei Liu, Wonmin Byeon, Thomas Breuel, Jan Kautz, and Xiaolong Wang. Groupvit: Semantic segmentation emerges from text supervision. In *CVPR*, pages 18134–18144, 2022. 3
- [30] Morris Alper and Hadar Averbuch-Elor. Emergent visual-semantic hierarchies in image-text representations. In *ECCV*, pages 220–238. Springer, 2024. 3

- [31] Reza Abbasi, Mohammad Hossein Rohban, and Mahdieh Soleymani Baghshah. Deciphering the role of representation disentanglement: Investigating compositional generalization in clip models. In *ECCV*, pages 35–50. Springer, 2024. 3
- [32] Valentin Gabeff, Marc Rußwurm, Devis Tuia, and Alexander Mathis. Wildclip: Scene and animal attribute retrieval from camera trap data with domain-adapted vision-language models. *IJCV*, 132(9):3770–3786, 2024. 3
- [33] Julian D Santamaria, Claudia Isaza, and Jhony H Giraldo. Catalog: A camera trap language-guided contrastive learning model. In *WACV*, pages 1197–1206. IEEE, 2025. 3, 5
- [34] Srikumar Sastry, Subash Khanal, Aayush Dhakal, Adeel Ahmad, and Nathan Jacobs. Taxabind: A unified embedding space for ecological applications. In *WACV*, pages 1765–1774. IEEE, 2025. 3
- [35] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song, David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, Wei-Lun Chao, and Yu Su. TreeOfLife-10M. 2023. doi: 10.57967/hf/1972. URL <https://huggingface.co/datasets/imageomics/TreeOfLife-10M>. 3
- [36] Grant Van Horn and Oisin Mac Aodha. inat challenge 2021 - fgvc8, 2021. URL <https://kaggle.com/competitions/inaturalist-2021>. 3, 26
- [37] Z. Gharaee, Z. Gong, N. Pellegrino, I. Zarubiieva, J. B. Haurum, S. C. Lowe, J. T. A. McKeown, C. Y. Ho, J. McLeod, Y. C. Wei, J. Agda, S. Ratnasingham, D. Steinke, A. X. Chang, G. W. Taylor, and P. Fieguth. A step towards worldwide biodiversity assessment: The BIOSCAN-1M insect dataset. In *NeurIPS*, volume 36, pages 43593–43619. Curran Associates, Inc., 2023. 3
- [38] GBIF.org. GBIF occurrence download. May 2024. doi: 10.15468/DL.BFV433. URL <https://doi.org/10.15468/dl.bfv433>. 3, 24
- [39] Encyclopedia of Life (EOL). URL <https://eol.org>. Accessed August 2024. 3
- [40] Zahra Gharaee, Scott C. Lowe, ZeMing Gong, Pablo Millan Arias, Nicholas Pellegrino, Austin T. Wang, Joakim Bruslund Haurum, Iuliia Zarubiieva, Lila Kari, Dirk Steinke, Graham W. Taylor, Paul Fieguth, and Angel X. Chang. BIOSCAN-5M: A multimodal dataset for insect biodiversity. In *NeurIPS*, volume 37, pages 36285–36313, 2024. 3
- [41] Kakani Katija, Eric Orenstein, Brian Schlining, Lonny Lundsten, Kevin Barnard, Giovanna Sainz, Oceane Boulais, Megan Cromwell, Erin Butler, Benjamin Woodward, and Katherine L. C. Bell. FathomNet: A global image database for enabling artificial intelligence in the ocean. *Scientific Reports*, 12(1):15914, September 2022. ISSN 2045-2322. doi: 10.1038/s41598-022-19939-2. URL <https://www.nature.com/articles/s41598-022-19939-2>. 3
- [42] iNaturalist contributors and iNaturalist. iNaturalist Research-grade Observations. Accessed via GBIF.org (Occurrence Snapshot: 10.15468/DL.BFV433) 2024-05-01. 4, 24
- [43] Smithsonian Institution. URL <https://www.si.edu/>. Accessed May 15, 2025. 4
- [44] Flickr. URL <https://www.flickr.com/>. Accessed May 15, 2025. 4
- [45] Andrei Kopanov, Matthew J Thompson, and Elizabeth G Campolongo. Distributed-downloader. 2025. doi: 10.5281/zenodo.17418004. URL <https://github.com/Imageomics/distributed-downloader>. 4
- [46] Dmitry Mozzherin. GNVerifier – a reconciler and resolver of scientific names against more than 100 data sources. November 2024. doi: 10.5281/zenodo.10070488. URL <https://github.com/gnames/gnVerifier>. 4
- [47] Sara Beery, Dan Morris, and Siyu Yang. Efficient pipeline for camera trap image review, 2019. 4, 25
- [48] Andres Hernandez, Zhongqi Miao, Luisa Vargas, Sara Beery, Rahul Dodhia, and Juan Lavista. Pytorch-wildlife: A collaborative deep learning framework for conservation, 2024. 4, 25
- [49] Kaipeng Zhang, Zhanpeng Zhang, Zhifeng Li, and Yu Qiao. Joint face detection and alignment using multitask cascaded convolutional networks. *IEEE signal processing letters*, 23(10):1499–1503, 2016. 4, 26
- [50] Facebook, Inc. Threatexchange. <https://github.com/facebook/ThreatExchange>, 2019. 4, 27
- [51] IUCN. IUCN Red List Summary Table 1a, March 2025. URL https://nc.iucnredlist.org/redlist/content/attachment_files/2025-1_RL_Table_1a.pdf. 4, 15, 28

- [52] IUCN. The IUCN Red List of Threatened Species. version 2025-1. <https://www.iucnredlist.org/search?dl=true&permalink=fd473452-551c-4561-a2b2-fb51fe51d7b7>, 2025. Accessed on 14 May 2025. 4, 27
- [53] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, pages 1–21, 2021. 5, 18
- [54] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade W Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, Patrick Schramowski, Srivatsa R Kundurthy, Katherine Crowson, Ludwig Schmidt, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5b: An open large-scale dataset for training next generation image-text models. In *NeurIPS*, pages 1–17, 2022. 5, 17
- [55] Sylvestre-Alvise Rebuffi, Alexander Kolesnikov, Georg Sperl, and Christoph H Lampert. icarl: Incremental classifier and representation learning. In *CVPR*, pages 2001–2010, 2017. 5
- [56] David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. In *NeurIPS*, volume 32, 2019. 5
- [57] Edward Vendrow, Omiros Pantazis, Alexander Shepard, Gabriel Brostow, Kate Jones, Oisin Mac Aodha, Sara Beery, and Grant Van Horn. Inquire: A natural world text-to-image retrieval benchmark. In *NeurIPS*, volume 37, pages 126500–126514, 2024. 5, 19
- [58] Ihsan Ullah, Dustin Carrion, Sergio Escalera, Isabelle M Guyon, Mike Huisman, Felix Mohr, Jan N van Rijn, Haozhe Sun, Joaquin Vanschoren, and Phan Anh Vu. Meta-album: Multi-domain meta-dataset for few-shot image classification. In *NeurIPS*, 2022. URL <https://meta-album.github.io/>. 5
- [59] Samuel Stevens, Jiaman Wu, Matthew J Thompson, Elizabeth G Campolongo, Chan Hee Song, David Edward Carlyn, Li Dong, Wasila M Dahdul, Charles Stewart, Tanya Berger-Wolf, Wei-Lun Chao, and Yu Su. Rare species. 2023. doi: 10.57967/hf/1981. URL <https://huggingface.co/datasets/imageomics/rare-species>. 5, 20, 26
- [60] Gerald Piosenka. Birds 525 species - image classification, 05 2023. URL <https://www.kaggle.com/datasets/gpiosenka/100-bird-species>. 6
- [61] 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 *CVPR*, pages 595–604, 2015. 6, 7, 23
- [62] Labeled information library of alexandria: Biology and conservation (LILA-BC). <https://lila.science/datasets>. 6, 20, 22
- [63] Island Conservation. Island conservation camera traps. <https://lila.science/datasets/island-conservation-camera-traps/>. 6, 20, 22
- [64] Desert Lion Conservation. Desert lion conservation camera traps. <https://lila.science/datasets/desert-lion-conservation-camera-traps/>, July 2024. 6, 20, 22
- [65] Juliana Vélez, Paula J Castiblanco-Camacho, Michael A Tabak, Carl Chalmers, Paul Fergus, and John Fieberg. Choosing an appropriate platform and workflow for processing camera trap data using artificial intelligence. *arXiv preprint arXiv:2202.02283*, 2022. 6, 20, 22
- [66] S. Balasubramaniam. Optimized classification in camera trap images: An approach with smart camera traps, machine learning, and human inference. Master’s thesis, The Ohio State University, 2024. URL http://rave.ohiolink.edu/etdc/view?acc_num=osu1721417695430687. 6, 20, 22
- [67] Hayder Yousif, Roland Kays, and Zhihai He. Dynamic programming selection of object proposals for sequence-level animal species classification in the wild. *IEEE Transactions on Circuits and Systems for Video Technology*, 2019. 6, 20, 22
- [68] Guangxing Han and Ser-Nam Lim. Few-shot object detection with foundation models. In *CVPR*, pages 28608–28618, 2024. 6
- [69] Simon Damm, Mike Laszkiewicz, Johannes Lederer, and Asja Fischer. Anomalydino: Boosting patch-based few-shot anomaly detection with dinov2. In *WACV*, pages 1319–1329. IEEE, 2025. 6

- [70] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008. 7
- [71] Hervé Abdi and Lynne J Williams. Principal component analysis. *Wiley interdisciplinary reviews: computational statistics*, 2(4):433–459, 2010. 8
- [72] Gene H. Golub and Charles F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, MD, 4 edition, 2013. ISBN 978-1-4214-0794-4. 8
- [73] Ronald A Fisher. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936. 9
- [74] Ohio Supercomputer Center. Ohio supercomputer center, 1987. URL <http://osc.edu/ark:/19495/f5s1ph73>. 10
- [75] Shawn T. Brown, Paola Buitrago, Edward Hanna, Sergiu Sanielevici, Robin Scibek, and Nicholas A. Nystrom. Bridges-2: A platform for rapidly-evolving and data intensive research. In *Practice and Experience in Advanced Research Computing 2021: Evolution Across All Dimensions*, PEARC ’21, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450382922. doi: 10.1145/3437359.3465593. URL <https://doi.org/10.1145/3437359.3465593>. 10
- [76] Devis Tuia, Benjamin Kellenberger, Sara Beery, Blair R Costelloe, Silvia Zuffi, Benjamin Risse, Alexander Mathis, Mackenzie W Mathis, Frank van Langevelde, Tilo Burghardt, et al. Perspectives in machine learning for wildlife conservation. *Nature communications*, 13(1):792, 2022. 15
- [77] David Lindenmayer and Ben Scheele. Do not publish. *Science*, 356(6340):800–801, 2017. doi: 10.1126/science.aan1362. 16
- [78] Xi Chen, Xiao Wang, Soravit Changpinyo, AJ Piergiovanni, Piotr Padlewski, Daniel Salz, Sebastian Goodman, Adam Grycner, Basil Mustafa, Lucas Beyer, et al. Pali: A jointly-scaled multilingual language-image model. In *ICLR*, pages 1–33, 2023. 18
- [79] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, pages 248–255. IEEE, 2009. 18
- [80] Vignesh Kothapalli. Neural collapse: A review on modelling principles and generalization. *Transactions on Machine Learning Research*, pages 1–43, 2022. 21
- [81] Vardan Papyan, XY Han, and David L Donoho. Prevalence of neural collapse during the terminal phase of deep learning training. *Proceedings of the National Academy of Sciences*, 117(40):24652–24663, 2020. 21
- [82] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. In *NeurIPS*, volume 33, pages 18661–18673, 2020. 21
- [83] Jaidev Gill, Vala Vakilian, and Christos Thrampoulidis. Engineering the neural collapse geometry of supervised-contrastive loss. In *ICASSP*, pages 7115–7119. IEEE, 2024. 21
- [84] Zhihui Zhu, Tianyu Ding, Jinxin Zhou, Xiao Li, Chong You, Jeremias Sulam, and Qing Qu. A geometric analysis of neural collapse with unconstrained features. In *NeurIPS*, volume 34, pages 29820–29834, 2021. 21
- [85] Sagar Vaze, Kai Han, Andrea Vedaldi, and Andrew Zisserman. Generalized category discovery. In *CVPR*, pages 7492–7501, 2022. 23
- [86] Yan Wang, Wei-Lun Chao, Kilian Q Weinberger, and Laurens Van Der Maaten. Simpleshot: Revisiting nearest-neighbor classification for few-shot learning. *arXiv preprint arXiv:1911.04623*, 2019. 23
- [87] GBIF Secretariat. GBIF Backbone Taxonomy. 2023. doi: 10.15468/39omei. URL <https://doi.org/10.15468/39omei>. 24
- [88] J. Doubt, R. Blades, and A. Tyrell. Canadian Museum of Nature Herbarium. 2025. doi: 10.15468/6e8nje. Version 1.283. Canadian Museum of Nature. Occurrence dataset <https://doi.org/10.15468/kowta4> accessed via GBIF.org on 2025-05-22. 26, 27
- [89] Natural History Museum and Benjamin Drew. Natural History Museum (London) Collection Specimens. 2025. Occurrence dataset <https://doi.org/10.5519/qd.x8gq8f5t> accessed via GBIF.org on 2025-05-14. <https://www.gbif.org/occurrence/1056329684>. Collected in Suriname by The Trustees of the Natural History Museum, London (licensed under <http://creativecommons.org/licenses/by/4.0/>). 27

Appendix

The Appendix is organized as follows:

- In §A we discuss limitations of our work.
- In §B we discuss the broader impacts of our work.
- In §C we present the proof of [Theorem 5.1](#).
- In §D we describe the training implementation for BIOCLIP 2 in detail.
- In §E we introduce the baselines employed for the numerical experiments.
- In §F we demonstrate more empirical observations of BIOCLIP 2.
- In §G we discuss the relationship between the adopted training scheme and neural collapse.
- In §H we describe both the biological visual benchmarks and the implementation details.
- In §I we provide the details of data processing for TREEOFLIFE-200M.
- In §J we list the contribution of each author.

A Limitations

Theoretical limitation. In this work, we have proved that the intra-species variations are preserved in subspaces orthogonal to the inter-species difference. We also have empirical observations that the variations are more separable as the training data scales up (See [Figure 3c](#)). However, we haven't formally proved this. It will be our future work to have deeper theoretical analyses of the separation of intra-species variations to better understand BIOCLIP 2's emergent properties.

Data limitation. TREEOFLIFE-200M is an imbalanced dataset in both taxonomic coverage and image type. Specifically, the dataset exhibits a long-tailed distribution across taxa. This is to be expected when working with biological data—not all taxonomic ranks are represented evenly across the tree of life. For instance, though TREEOFLIFE-200M has a balanced representation (at the kingdom level) between plants and animals, animals represent a larger proportion of described species [51].

Some of this is due to the nature of the image type distribution, which we provide for GBIF (Camera-trap, Citizen Science, and eleven Museum Specimen types: Fungi, Insect, Invertebrate Zoology, Microbiology, Plant, Vertebrate Zoology - [Amphibians, Birds, Fishes, Mammals, Others], as well as Unidentified). EOL contains the same categories, but we do not have precise numbers; BIOSCAN-5M are essentially all insect museum specimens, though the images are taken by researchers, so will skew toward their area of study; FathomNet contains a mix. Citizen Science images are the vast majority (151M from GBIF alone); these will skew toward more charismatic species and plants. Our next largest category is museum specimen images (51.8M from GBIF), which are limited more to representatives of the species, so may not contain as many images per taxa, though more taxa are represented. Finally, camera trap images make up the smallest portion of the dataset (617.8K in GBIF), and when filtered, these are only images of animals and generally those large enough to trigger a motion sensor or be detected in the primary provider's post-processing.

Further emphasis on the impact of citizen science images in amassing larger representations of species: the most prevalent taxonomic classes, flowering plants, insects, birds, mushrooms, and mammals, have millions of representative images, while the least-represented, microscopic organisms (e.g., bacteria, viruses), have a dozen or fewer.

B Broader Impacts

BIOCLIP 2 and TREEOFLIFE-200M provide great potential to improve and enhance existing conservation efforts, in particular by facilitating recognition of threatened species. As noted in [§3.3](#), TREEOFLIFE-200M has expansive coverage of threatened species, as classified by IUCN. It additionally builds on the coverage of species considered to be Data Deficient. Based on the emergent properties displayed by BIOCLIP 2, there is potential to add to the effort to understand the risks facing these species that cannot currently be classified by IUCN due to lack of available information, as suggested in [76]. These designations are crucial to the international effort to protect biodiversity across the planet.

Unfortunately, as with many open-source efforts to further conservation goals, there is potential for bad actors to make use of these tools for malicious purposes. Though the improvement on threatened species *could* make it easier for poachers to identify protected species, these types of tools are a force-multiplier to monitor illicit trade and sales of these same species. The primary risk to endangered species comes from disclosure of precise location information rather than improved classification capability [77]. Our data does not provide geo-tagged information on the organisms included, minimizing the vulnerabilities that could be used in poaching.

C Proof of Theorem 5.1

Proof. The contrastive loss for one visual embedding \mathbf{z} belonging to the class s and the corresponding text embedding \mathbf{c}_s is:

$$l(\mathbf{z}, \mathbf{c}_s) = -\log \frac{\exp(\mathbf{z}^\top \mathbf{c}_s / \tau)}{\sum_k \exp(\mathbf{z}^\top \mathbf{c}_k / \tau)}. \quad (1)$$

Let $h_\phi(L(\cdot))$ be the text encoder. Assume the representation is already close to the species prototype $\boldsymbol{\mu}_s = h_\phi(L(s)) = \mathbf{c}_s$, and there is a residual $\boldsymbol{\delta}$ representing the intra-species variance:

$$\boldsymbol{\delta} := \mathbf{z} - \boldsymbol{\mu}_s, \quad \|\boldsymbol{\delta}\| \ll \|\boldsymbol{\mu}_s - \boldsymbol{\mu}_{k \neq s}\|.$$

Define:

$$a_k := \frac{\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_k}{\tau}, \quad Z = \sum_k \exp(a_k), \quad w_k := \frac{\exp(a_k)}{Z}$$

Substituting $\mathbf{z} = \boldsymbol{\mu}_s + \boldsymbol{\delta}$ into Equation 1:

$$l(\boldsymbol{\mu}_s + \boldsymbol{\delta}, \mathbf{c}_s) = -\frac{(\boldsymbol{\mu}_s + \boldsymbol{\delta})^\top \boldsymbol{\mu}_s}{\tau} + \log \left[\sum_k \exp \left(\frac{(\boldsymbol{\mu}_s + \boldsymbol{\delta})^\top \boldsymbol{\mu}_k}{\tau} \right) \right] \quad (2)$$

$$= -\frac{\boldsymbol{\mu}_s^\top (\boldsymbol{\mu}_s + \boldsymbol{\delta})}{\tau} + \log \left[\sum_k \exp \left(\frac{\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_k}{\tau} \right) \cdot \exp \left(\frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right) \right] \quad (3)$$

$$= -\frac{\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_s}{\tau} - \frac{\boldsymbol{\mu}_s^\top \boldsymbol{\delta}}{\tau} + \log \left[\sum_k \exp \left(\frac{\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_k}{\tau} \right) \sum_k \frac{\exp(\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_k / \tau)}{\sum_k \exp(\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_k / \tau)} \cdot \exp \left(\frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right) \right] \quad (4)$$

$$= -\frac{\boldsymbol{\mu}_s^\top \boldsymbol{\mu}_s}{\tau} - \frac{\boldsymbol{\mu}_s^\top \boldsymbol{\delta}}{\tau} + \log Z + \log \left[\sum_k w_k \exp \left(\frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right) \right]. \quad (5)$$

Use the Taylor expansion to the second order for the argument of the logarithm of the last term $\Psi(\boldsymbol{\delta}) = \log [\sum_k w_k \exp(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k / \tau)]$:

$$\sum_k w_k \exp \left(\frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right) = 1 + \sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} + \frac{1}{2} \sum_k w_k \frac{(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k)^2}{\tau^2} + O(\|\boldsymbol{\delta}\|^3), \quad (6)$$

$$\text{obtaining} \quad \Psi(\boldsymbol{\delta}) = \sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} + \frac{1}{2} \left[\sum_k w_k \frac{(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k)^2}{\tau^2} - \left(\sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right)^2 \right] + O(\|\boldsymbol{\delta}\|^3). \quad (7)$$

Insert $\Psi(\boldsymbol{\delta})$ back into Eq.2:

$$\begin{aligned} l &= l(\boldsymbol{\mu}_s, \boldsymbol{\mu}_s) - \frac{\boldsymbol{\mu}_s^\top \boldsymbol{\delta}}{\tau} + \sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \\ &\quad + \frac{1}{2} \left[\sum_k w_k \frac{(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k)^2}{\tau^2} - \left(\sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right)^2 \right] + O(\|\boldsymbol{\delta}\|^3). \end{aligned}$$

Define $\mathbf{m} := \sum_k w_k \boldsymbol{\mu}_k$. Then collecting first-order terms around $\boldsymbol{\delta}$ results in the following:

$$\frac{(\sum_k w_k \boldsymbol{\mu}_k - \boldsymbol{\mu}_s)^\top \boldsymbol{\delta}}{\tau} = \frac{(\mathbf{m} - \boldsymbol{\mu}_s)^\top \boldsymbol{\delta}}{\tau}.$$

Suppose the training resulted in $w_s \approx 1$, leading to $\mathbf{m} \approx \boldsymbol{\mu}_s$. Then the first-order terms will vanish as the embeddings of different species are better separated.

We further rewrite the second-order term of Eq.6 as:

$$\begin{aligned} \frac{1}{2} \left[\sum_k w_k \frac{(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k)^2}{\tau^2} - \left(\sum_k w_k \frac{\boldsymbol{\delta}^\top \boldsymbol{\mu}_k}{\tau} \right)^2 \right] &= \frac{1}{2} \left[\sum_k w_k \frac{(\boldsymbol{\delta}^\top \boldsymbol{\mu}_k)^2}{\tau^2} - \frac{(\boldsymbol{\delta}^\top \mathbf{m})^2}{\tau^2} \right] \\ &= \boldsymbol{\delta}^\top \left[\frac{1}{2\tau^2} \left(\sum_k w_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top - \mathbf{m} \mathbf{m}^\top \right) \right] \boldsymbol{\delta}. \end{aligned}$$

Substituting $\mathbf{m} \approx \boldsymbol{\mu}_s$, we have an approximation as

$$\boldsymbol{\delta}^\top \left[\frac{1}{2\tau^2} \left(\sum_k w_k \boldsymbol{\mu}_k \boldsymbol{\mu}_k^\top - \boldsymbol{\mu}_s \boldsymbol{\mu}_s^\top \right) \right] \boldsymbol{\delta}.$$

The Hessian matrix lies in the span of species, as each term is an outer product of a prototype. Therefore, as long as the residual $\boldsymbol{\delta}$ is distributed in a subspace orthogonal to the species differences, the scale of $\boldsymbol{\delta}$ will not interfere with the overall contrastive optimization. \square

D Training Implementation Details

Table 4: The adopted hyper-parameter setting in training BIOCLIP 2.

Hyper-parameter	Value
Architecture	ViT-L/14
Optimizer	Adam
Batch size/GPU (organism)	2,816
Batch size/GPU (replay)	320
GPUs	32 H100s
Epochs	30
Max learning rate	1×10^{-4}
Warm-up steps	1,875
Weight decay	0.2
Input resolution	224

Table 5: The adopted hyper-parameter setting in the ablation study.

Hyper-parameter	Value
Architecture	ViT-L/14
Optimizer	Adam
Batch size/GPU (organism)	2,816
Batch size/GPU (replay)	320
GPUs	8 H100s
Epochs	100
Max learning rate	1×10^{-4}
Warm-up steps	1,500
Weight decay	0.2
Input resolution	224

We list the adopted hyperparameters in Table 4 and Table 5 for the BIOCLIP 2 training and ablation study, respectively. TREEOF LIFE-10M is used for the ablation study. The batch size presented in both tables is the size per GPU. In addition to the larger GPU number and smaller batch size compared with the training in BIOCLIP due to the larger model size, another important modification is the introduction of experience replay. An additional visual projector is introduced as described in §F.1. Beyond the visual projector, the model architecture is kept the same as that of CLIP [12]. During the evaluation, all the embeddings are extracted with the visual projector for hierarchical label matching to avoid extra influence.

E Baseline Details

In the quantitative experiments, we compare BIOCLIP 2 with the following baseline vision-language models and vision-only models. Without specification, the input image size is 224.

- **CLIP (ViT-L/14).** We compare BIOCLIP 2 with the CLIP model pre-trained on LAION-2B [54], which has the same architecture and patch size. The CLIP model is also used as the weight initialization of the BIOCLIP 2 training. It uses a ViT-large visual encoder with a patch size of 14. We load the weight from the OpenCLIP repository.

Table 6: Performance comparison between different replay designs of CLIP training data. All the models in the bottom three rows are initialized with CLIP (the first row) and trained with TREEOFLIFE-10M data. The CLIP model is pre-trained with LAION-2B data, from which we randomly select 2M samples for this experiment. Δ represents the performance gap over the CLIP baseline. **Bold** entries indicate the best accuracy.

Model	Species		Non-Species						INQUIRE	Mean (Δ)
	NABirds	Rarespecies	FishNet	NeWT	AWA2	Herb. 19	PlantDoc			
CLIP (ViT-L/14)	66.5	35.2	27.9 ± 0.2	83.4 ± 0.1	61.6 ± 0.6	18.2 ± 0.1	22.3 ± 3.3	35.0	43.8	–
No Replay	68.9	46.1	35.3 ± 0.1	83.8 ± 0.1	58.0 ± 2.8	29.4 ± 0.4	36.0 ± 2.8	34.4	49.0	$\uparrow 5.2$
Single-proj	68.8	44.8	34.4 ± 0.2	84.5 ± 0.2	58.1 ± 1.4	31.0 ± 0.2	38.1 ± 2.4	34.7	49.3	$\uparrow 5.5$
Separate-proj (Ours)	71.2	47.2	35.1 ± 0.1	84.2 ± 0.1	57.4 ± 0.9	30.4 ± 0.3	38.7 ± 3.7	37.1	50.2	$\uparrow 6.4$

- **SigLIP.** In addition to the standard CLIP model, we also evaluate the performance of SigLIP [21]. We adopt the SigLIP model pre-trained on WebLI data [78]. The adopted visual encoder is ViT-large [53], the patch size is 16, and the input image resolution is 256. The model weight is also loaded from the [OpenCLIP](#) repository.
- **Supervised-IN21K.** For vision-only models, we first select a ViT-large model [53] trained in a supervised way on ImageNet-21K dataset [79]. As it is a vision-only model, we only run a few-shot and non-species classification tasks with it. The patch size is 32. The model is publicly downloadable in [Hugging Face](#).
- **DINOv3.** Besides supervised pre-training, we also evaluate the performance of DINOv3, which is pre-trained in an unsupervised way [22]. The backbone architecture is ViT-large, and the patch size is 16. The model can be downloaded from [Hugging Face](#). Similarly, we only run few-shot and non-species classification evaluations with DINOv3.
- **BioTrove-CLIP.** The above four models are trained on general knowledge covering a variety of topics. We also compare BioCLIP 2 with domain-specific models. BioTrove-CLIP is trained with BioTrove [14]. The model weights can be downloaded from [Hugging Face](#). Among the provided three models, we use the model initialized with OpenAI CLIP [12] that yields the best average accuracy [14]. The visual backbone is ViT-base with a patch size of 16.
- **BioCLIP.** BioCLIP is trained on TREEOFLIFE-10M. It adopts ViT-base as the visual encoder, with a patch size of 16. The model weight is publicly available in [Hugging Face](#).

F More Empirical Observations

In this section, we present more empirical observations of the training design, detailed numerical results, and qualitative analyses.

F.1 CLIP Training Data Replay

Together with the contrastive supervision of hierarchical labels, we also introduce the replay of CLIP training data to retain the understanding of general knowledge. For each experiment of different training scales, we randomly select a subset from LAION-2B with 10%-20% of the corresponding biological image total. Specifically, for the largest run on 214M biological images, we select 26M samples from LAION-2B. Each training batch consists of 69,312 biological images and 8,192 replay samples. However, the text labels in TREEOFLIFE-200M are primarily different forms of taxonomic names, which have a distributional gap from the CLIP training data. Therefore, we apply a separate visual projector specifically for the replay data to avoid the optimization conflict. Other than this difference, the biological data and replay data share the same visual backbone and text encoder.

We evaluate different replay settings quantitatively in [Table 6](#). The “No Replay” row shows the baseline performance applying biological contrastive training on top of the pre-trained CLIP model, where a 5.2% performance improvement is achieved. When the replay data is added, which shares the visual projector with biological images, we observe a conflicting performance change. On Herb. 19 [19] and PlantDoc [20], more than 2% improvement is acquired. However, the species classification accuracy on Rare Species drops by 1.3%. We attribute the inconsistent performance change to the distribution gap between biological and replay text labels. After applying a separate

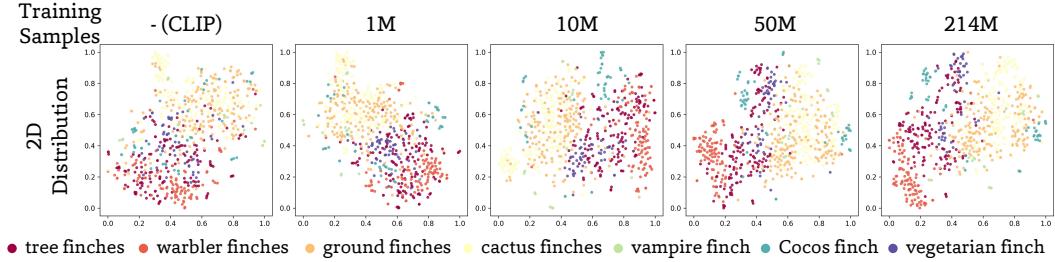


Figure 7: The t-SNE distribution of Darwin’s finches under different scales of training data. Different colors represent different groups of finch species. As the training data scales up, the embeddings form biologically meaningful clusters that align with the phylogenetic tree and their functional traits.

visual projector for the replay data, we observe overall improved performance across multiple benchmarks. At the same time, we also notice that replay has limited influence on benchmarks like AwA2 [17]. Additionally, we evaluate the preservation of general knowledge understanding on INQUIRE-Rerank [57]. Applying contrastive training with taxonomic labels slightly hurts the performance, while the single-projector replay fails to retain the understanding. Comparatively, the adopted separate-projector replay improves the performance by 2.1%. Given that there is still a distribution gap between the taxonomic labels and natural language, we do not treat INQUIRE as one of our main focuses in this work.

F.2 Qualitative Analyses

Embedding distribution of Darwin’s finches. In Figure 1, we show that the embedding distribution of Darwin’s finches aligns with the beak size. Here we further visualize the distribution under different training scales in Figure 7. Based on the genome-based phylogeny, warbler finches were the most ancient branches, while tree finches and ground finches form the recent branches [23]. Among these species, warbler finches have the smallest beak, convenient for extracting tiny arthropods from leaves. Comparatively, ground finches have larger beaks, which are more suitable for cracking seeds and nuts. In the original CLIP embedding space, the warbler finches and tree finches are mixed. As the training data scales up, these two groups are separated, and the relative geometric relationship of all the finches aligns with their phylogeny tree. While the species separation is induced by taxonomic supervision, BIOCLIP 2’s embedding space again illustrates emergent higher-level biological meaningfulness after extensive training.

Embedding distribution of Sex data. Similar to Figure 5, we visualize the 2D and 3D distributions of embeddings from 3 species of the Sex data in Figure 8. We can draw conclusions similar to those obtained from Figure 5. When no training data is incorporated, the embeddings extracted by the original CLIP visual encoder (leftmost sub-figure) demonstrate large portions of overlap between male and female images. After extensive vision-language contrastive training, the embeddings present clear decision boundaries between the two variants. Furthermore, as evidenced in the 3D distribution, the embeddings within each species form more compact clusters when projected onto the species span (the gray plane). Comparatively, instead of being eliminated after contrastive training, the intra-species variations are embedded in the subspace orthogonal to the inter-species differences. The extensive training facilitates BIOCLIP 2 to acquire a biologically meaningful embedding space, highlighting its value in serving as a biology foundation model.

Embedding distribution of PlantDoc data. We further visualize the distribution of 6 classes from the PlantDoc dataset in Figure 9. When no training is processed (leftmost sub-figure), embeddings of different species, as well as diseases, are mixed. As the training scale increases, we observe two trends. First, the margin between different species is enlarged, and the embeddings belonging to the same species are clustered together. Second, the diseased leaves are easier to separate within each species, although not explicitly constrained during training. More interestingly, the embeddings of healthy apple leaves are distributed close to the healthy blueberry leaves. These observations again highlight the biologically meaningful embedding space of BIOCLIP 2 after extensive training.

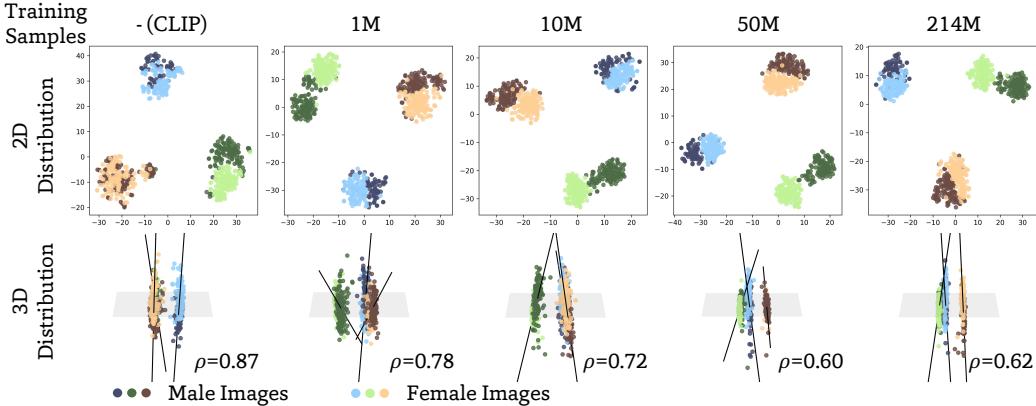


Figure 8: The embedding distribution of sex variations under different scales of training data. The 2D distributions are obtained using t-SNE. For the 3D distributions, we first run SVD with the mean embedding of each species. The first two singular vectors are used to construct the gray plane that captures most inter-species differences. The embeddings are then projected into the 3D space with an additional orthogonal dimension. The straight lines point from the mean embedding of male images to that of the female images. As the training data scales up, the intra-species variations are preserved in the subspace orthogonal to the inter-species differences. Orthogonality improves with data scale, as evidenced by the decreasing explained-variance ratio ρ .

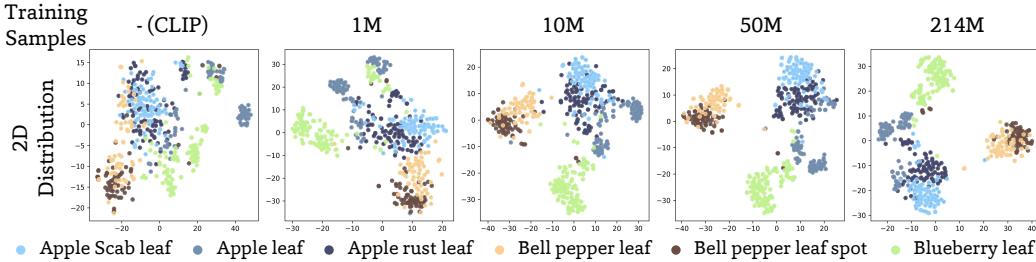


Figure 9: The t-SNE distribution of 6 classes from the PlantDoc dataset, including three species and three different diseases. As the training data scales up, not only are the species better separated, but the intra-species variations also form clusters, making them easier to separate.

E.3 The decay of FDR numerator/denominator along the data scale.

We observe that FDR between two intra-species variation classes is increasing as the training data scales up. More specifically, we look into the numerator term (the difference between two class centers) and the denominator term (the variation of features). We visualize the curve of the numerator and denominator terms in Figure 10a and Figure 10b for life stage data and sex data, respectively. The y-axes are scaled to the same maximum ratio to the minimum value. As the number of training samples scales up from 1M to 214M, the denominator term goes through a larger decay than the numerator term. The numerical results further support the increasing FDR trend.

E.4 Camera Trap Results

In addition to the species classification benchmarks adopted in [13], we further introduce a balanced camera trap image benchmark for species classification, **IDLE-OO Camera Traps**, derived from LILA-BC datasets [62], to construct a more realistic application scenario. Specifically, we select five datasets from LILA-BC that are labeled to the image level to avoid testing on noisy images—those labeled as containing an animal when it is simply the animal’s habitat. The Island Conservation Camera Traps [63] were of particular interest for their stated purpose of assisting in the prevention of endangered island species’ extinction and the varied ecosystems represented. This provides a fine-grained complement to the Rare Species test set [13, 59]. The Desert Lion Conservation Camera Traps dataset [64] is similarly intended to advance conservation efforts. With the Orinoquia Camera Traps [65], Ohio Small Animals [66], and ENA24 [67] camera trap datasets, we can test on camera

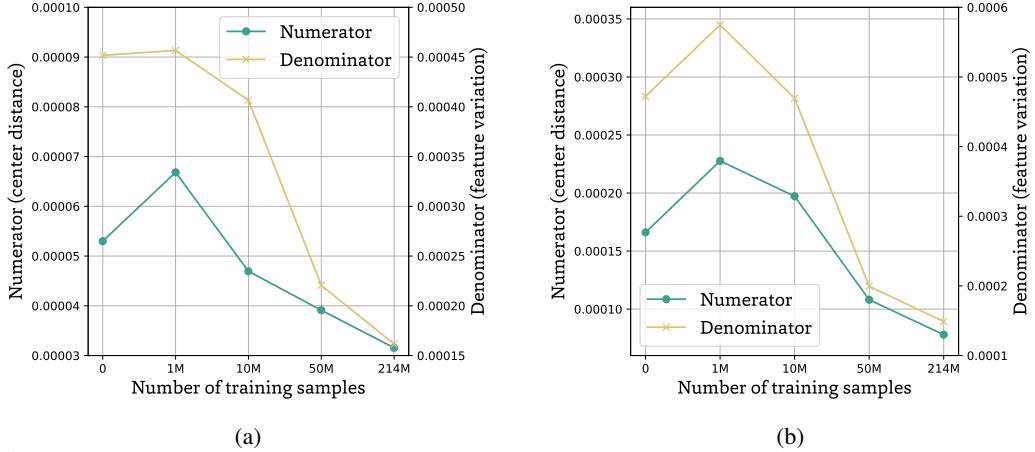


Figure 10: The decay curves for the numerator (difference between two class centers) and the denominator (feature variation) terms of the FDR metric along the increasing data scale. (a) The curves of life stage data. (b) The curves of sex data.

trap images across varied settings. The Island Conservation set uses common names, while the remaining datasets are reduced to just those that are labeled to the species level, and then sampled to create a balanced test set across the remaining classes (these are evaluated on scientific names). The sampling and image manifests are provided in [IDLE-OO Camera Traps](#). We report the detailed accuracy on each camera trap dataset in [Table 7](#).

G Discussion on the Relationship with Neural Collapse

Neural collapse is a status where the intra-class variations collapse to zero, and the embeddings collapse to their corresponding class prototypes [80, 81]. The class prototypes collapse to the vertices of a simplex Equiangular Tight Frame (ETF). It has been empirically observed and theoretically proved that the commonly used loss functions—including cross-entropy loss and supervised contrastive loss [82]—lead to neural collapse at the terminal phase of training [24, 25, 83, 84]. If neural collapse happens, the intra-species appearance variations will be hardly separable. However, we reveal that after extensive training of BIOCLIP 2, the intra-species variations are preserved in the subspace orthogonal to inter-species differences, indicating that BIOCLIP 2 does not suffer neural collapse.

Why BIOCLIP 2 does not lead to Neural Collapse. We summarize two key reasons that allow the existence of the intra-species variation subspace in BIOCLIP 2.

First, the class prototypes in standard classification tasks and supervised contrastive learning (SCL) can both be treated as fully trainable parameters [82]. In standard classification, the prototypes are the weights of the linear classifier. In SCL, they are simply the mean embedding of the corresponding classes. The class prototypes will form a simplex ETF after extensive training. Comparatively, the adopted contrastive training scheme in BIOCLIP 2 employs text embeddings of the taxonomic labels as class prototypes. During the creation of taxonomic labels, hierarchical structures have been naturally embedded into them based on ecological and functional evidence. Different species can share higher taxonomic levels, which makes them hard negatives. Therefore, even if the text encoder is also being trained, the generated prototypes will not become ETF.

Second, TREEOFLIFE-200M poses an enormous label space, where 952K classes are involved in training. Typical SCL is performed upon CIFAR or ImageNet-level datasets, consisting of 100–1,000 classes. Previous works analyzing neural collapse usually assume the dimension of embeddings is larger than the class number and the sample number is balanced across different classes [25, 81]. In contrast, the class number involved in BIOCLIP 2 is much larger than the embedding dimension (768 channels). These conditions restrict the possibility of prototypes forming ETF.

Table 7: Zero-, one-, and five-shot species classification accuracy of **IDLE-OO Camera Traps** from the Labeled Information Library of Alexandria: Biology and Conservation (LILA-BC) [62, 63, 64, 65, 66, 67] for different models. **Bold** and underlined entries indicate the **best** and second best accuracies, respectively.

LILA-BC						
Model	Desert Lion	ENA24	Island	Orinoquia	OH Small Animals	Mean
Random Guessing	3.1	5.0	3.0	3.6	2.6	3.5
<i>Zero-Shot Classification</i>						
CLIP (ViT-L/14)	35.2	38.2	27.1	25.6	21.6	29.5
SigLIP	46.9	41.0	<u>28.1</u>	<u>31.9</u>	20.5	<u>33.7</u>
BioTrove-CLIP	9.7	10.4	13.2	8.0	12.4	10.7
BIOCLIP	<u>47.2</u>	<u>42.5</u>	18.4	27.1	<u>23.1</u>	31.7
BIOCLIP 2	58.8	68.5	42.6	47.9	51.5	53.9
<i>One-Shot Classification</i>						
CLIP (ViT-L/14)	37.6 ± 2.5	35.9 ± 3.8	44.7 ± 3.0	31.3 ± 2.7	30.7 ± 2.0	36.0
SigLIP	41.1 ± 1.8	39.1 ± 4.0	48.5 ± 2.9	32.7 ± 1.6	27.6 ± 2.5	37.8
Supervised-IN21K	32.4 ± 2.0	28.4 ± 3.2	40.1 ± 2.2	26.3 ± 2.4	26.0 ± 2.4	30.6
DINOv3	48.1 ± 1.3	38.5 ± 3.2	50.5 ± 2.8	40.0 ± 2.6	<u>36.1</u> ± 2.3	<u>42.6</u>
BioTrove-CLIP	30.5 ± 1.7	30.5 ± 2.2	37.9 ± 2.1	29.0 ± 2.6	28.3 ± 3.1	31.2
BIOCLIP	39.9 ± 2.2	34.4 ± 2.7	45.7 ± 3.0	27.5 ± 3.7	27.5 ± 2.2	35.0
BIOCLIP 2	54.3 ± 1.1	48.6 ± 2.3	49.5 ± 2.5	43.1 ± 2.1	45.0 ± 3.0	48.1
<i>Five-Shot Classification</i>						
CLIP (ViT-L/14)	58.3 ± 1.7	57.8 ± 2.6	66.1 ± 2.4	43.9 ± 2.9	43.5 ± 1.6	53.9
SigLIP	64.3 ± 2.2	60.0 ± 2.7	<u>71.6</u> ± 1.4	46.2 ± 1.8	42.6 ± 1.7	56.9
Supervised-IN21K	51.3 ± 2.3	48.5 ± 2.6	59.0 ± 1.3	39.3 ± 3.4	41.9 ± 1.6	48.0
DINOv3	66.3 ± 3.3	63.2 ± 3.1	73.4 ± 1.3	59.9 ± 3.7	53.8 ± 2.4	<u>63.3</u>
BioTrove-CLIP	47.6 ± 3.0	46.7 ± 2.1	62.2 ± 1.1	41.1 ± 1.6	41.8 ± 2.4	47.9
BIOCLIP	62.5 ± 1.7	57.0 ± 3.2	63.2 ± 3.8	44.6 ± 1.7	44.0 ± 0.5	54.3
BIOCLIP 2	73.4 ± 2.9	73.4 ± 0.9	70.7 ± 1.6	<u>59.4</u> ± 2.8	61.8 ± 1.3	67.7

H Biological Visual Evaluation Details

Instead of training specialist models, we treat the evaluated models as frozen visual embedding extractors. Standard machine learning algorithms are applied on top of the acquired visual embeddings to predict the corresponding labels. Such a design is adopted to evaluate the quality of the embeddings while avoiding the influence of complicated optimization loops. In the following, we introduce the details of the adopted benchmarks and the evaluation algorithms. All the experiments are conducted with 1 NVIDIA A100 GPU, and the running time for each task is within 30 minutes.

FishNet. FishNet focuses on recognizing, locating, and predicting species and their functional traits [16]. Specifically, 94,532 images are collected with annotations of habitat, ecological role, and nutritional value. In this work, we mainly focus on the prediction of habitats and ecological roles, involving 9 groups of binary labels (*e.g.*, whether the fish can live in freshwater). Following the practice in the original paper, we train a two-layer linear classifier with binary cross-entropy loss to predict the 9 labels. We count a correct prediction only if all the 9 labels are predicted correctly for the sample. The original train-test split is adopted, where 75,631 images are used in training, and the remaining 18,901 images are used for testing. This task evaluates whether the embedding distribution of different species is aligned with their ecological relationships.

NeWT. NeWT comprises 164 binary classification tasks in the natural world [18]. The tasks include appearance, gestalt, context, counting, and behavior concepts. In each task, 50-100 images are assigned per class per train-test split. After extracting visual embeddings, we apply a support vector classifier for each of the tasks. The average accuracy is reported across all the evaluated tasks.

AwA2. AwA2 consists of 37,322 images of 50 animal classes, with annotations of 85 numeric attribute values for each class [17]. The dataset can be used for testing the capability of attribute-based classification and zero-shot transfer learning of trait prediction. In this work, we mainly focus on the transfer learning scheme. The training set of 45 classes is used to train an attribute classifier, and the remaining 5 unseen classes are used for testing. Similar to FishNet, we incorporate a linear classifier on top of the extracted features to predict the binary labels for the 85 attributes. The average F1 score over all the attributes is reported for this benchmark.

Herb. 19. Herbarium 19 is a task for discovering new species [19]. Specifically, given the images of known and unknown species, the model is required to predict the labels for both of them. As there are no fixed labels for unknown species, the task is implemented with a form of semi-supervised clustering [85]. Given the total number of species, a semi-supervised K-means algorithm is conducted on top of the extracted embeddings to cluster images. Clustering accuracy is calculated for the predictions following the original practice [85].

PlantDoc. PlantDoc is a dataset targeting the incorporation of computer vision for scalable and early plant disease detection [20]. 2,598 images of 13 plant species and up to 17 classes of diseases are collected in uncontrolled natural settings. We evaluate the model on PlantDoc in a few-shot learning style. One image per class is randomly selected from the original training split as the support set, and SimpleShot [86] is employed to predict the class labels for the test set. Accuracy over all the testing samples is reported as the performance.

Life stage-Diff/Align. We use the data of Age tasks from NeWT [18] to construct the Life stage-Diff/Align benchmark, where images are labeled with juvenile and adult classes. Specifically, the differentiation tasks aim to separate the binary appearance variations within each species. Conversely, in the alignment task, we train a species classifier with juvenile images while testing it using adult images. It requires the embeddings of two variations within one species to be closer than the embeddings of different species. For both of the tasks, we incorporate a support vector classifier to predict the corresponding labels. Ideally, after extensive contrastive training, the embeddings of different variation classes are expected to collapse to the species prototype. However, we demonstrate that the intra-species variations are well preserved in the embedding space of BIOCLIP 2.

Sex-Diff/Align. NABirds consists of 48,000 images from 400 bird species [61]. Sex and life stage labels are provided for those species with large appearance variations. We manually examine the images and select 81 species with male-female differences to construct the Sex-Diff/Align benchmark, where 13,624 images are used in total. For the differentiation task, we filter out 20 male images and 20 female images per species as the test set. Among the remaining 10,384 images, we filter out at most 20 male and 20 female images for training. For the alignment setting, we use the images of female birds to train a species classifier and use the male images for testing. Similar to the life stage benchmark, we use a support vector classifier for both of the sub-tasks.

I Data Processing Details

In this section, we provide more details on the data curation process.

I.1 Taxonomic Standardization

When GNVerifier returns a result from a query, our taxonomic alignment package combines it with the input taxonomic hierarchy, along with the query parameter, to form a resolution attempt. The query response, based on the most specific taxonomic term available in the input data, along with the remaining input hierarchy (entry’s resolution attempt), is then matched against pre-defined profiles. The algorithm uses these three components to fit against the series of profiles to determine whether a confident resolution is found or if an alternative query strategy is needed (such as using a different query term or data source), iterating until a match is made or alternative approaches are exhausted.

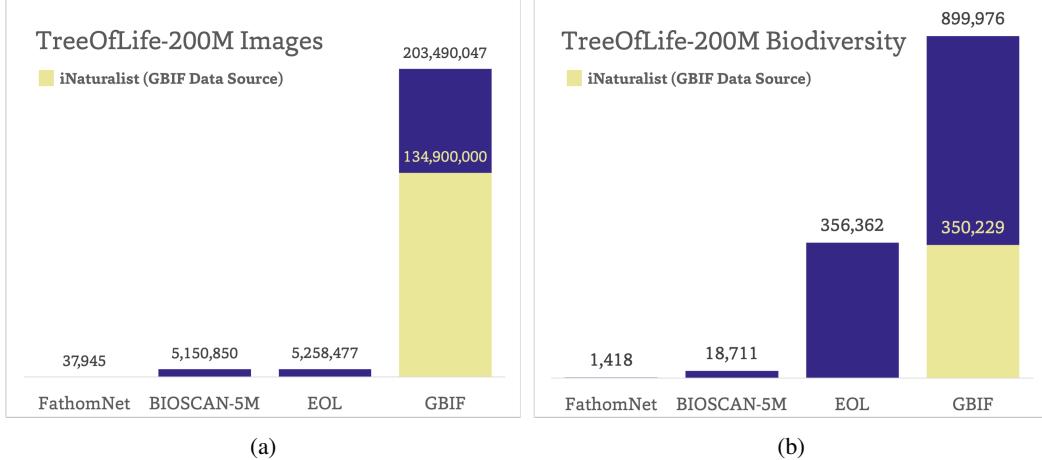


Figure 11: The number of images (a) and unique 7-rank taxa (b) by core data provider in TREEOFLIFE-200M. iNaturalist [42] is the largest GBIF data source (by image count) and the data provider used by BioTrove [14], so it is included here for reference.

Table 8: Top taxa contributors in GBIF [38]. Here, “Distinct Taxa” is used to describe 7-rank taxonomic labels only provided to GBIF by that publisher, while “Total Taxa” is the total number of unique taxa labels within that publisher. Observe that, following iNaturalist [42], the next three most diverse sources are museums, which, together, account for more distinct taxa than iNaturalist alone, despite having significantly fewer images overall.

Publisher (GBIF Data Source)	Distinct Taxa	Total Taxa	Images	% GBIF Taxa
iNaturalist.org	81,671	350,229	134,877,019	9.07%
Natural History Museum	38,535	148,637	3,603,256	4.28%
Museum national d’Histoire nat.	37,201	258,766	6,185,477	4.13%
Naturalis Biodiversity Center	25,421	244,697	5,282,428	2.82%

Table 9: The distribution of unique taxa and the covered images at different taxonomic ranks in TREEOFLIFE-200M. 92.26% images in TREEOFLIFE-200M have taxonomic labels specific to the species level.

Rank	Total Images	Unique Taxa	Percent of Dataset
Species	197,382,628	867,455	92.26%
Genus	206,160,396	135,380	96.36%
Family	207,489,189	13,790	96.99%
Order	210,063,485	1,683	98.19%
Class	211,236,966	382	98.74%
Phylum	212,416,362	127	99.29%
Kingdom	213,932,022	11	100.00%

If all attempts ultimately fail to match a profile, the original input data is used for the entry’s label in the final output. After hierarchical resolution, common names are annotated. For each entry, the common name corresponding to its most specific resolved taxonomic term is selected, when available, using the preferred English vernacular names from the GBIF Backbone Taxonomy [87]. In the output, the proportion of records changed for each data source is EOL: 98.7%, FathomNet: 16.0%, BIOSCAN: 11.8%, and GBIF: 0.3%. The notably low modification rate for GBIF reflects our taxonomic alignment package’s preference for the GBIF Backbone Taxonomy as its primary reference. After standardization, we summarize the image number and taxa number from each source in Figure 11. We also report the distinct taxa and image number provided by top publishers from GBIF in Table 8. We provide the source code for this part in [TaxonoPy](#).

As claimed in §3.3, not all taxa are specific to the species level. We summarize the detailed taxa distribution in Table 9. Although some ambiguous labels were mapped up to higher taxonomic ranks, our overall dataset still predominantly comprises images confidently labeled at the species level. Specifically, out of the total 214M images, approximately 92.26% retain species-level labels.

I.2 Image-quality Screens

At download, all images were checked to be sufficiently large (224 pixels or more on the shortest side), and resized, where needed, so that the largest side does not exceed 720 pixels. The code, support set class embeddings, and further details of the image-quality screens and the subsequent duplicate control are provided in our dataset repository, [TreeOfLife-toolbox](#).

I.2.1 Museum Specimen Processing

Museums often consider multiple specimens of the same species to be connected instead of differentiating them (as they are often collected from the same or similar location). Thus, these occurrences often include images of their metadata, duplicates of the same or different specimens under one occurrence ID, or less informative images. Some common complications to consider when working with museum specimen images that influenced our processing are that

1. Plant and fungi specimens may be stored in envelopes or folders. These are often photographed and digitized under the same occurrence ID as the image of the specimen itself, thus creating extraneous images we do not want to include in training. They also sometimes pair this with living specimen images (which would be worth keeping as well). See [Figure 12](#), which demonstrates variety within a single occurrence for fungi specimens. Plant specimens have a secondary confounding factor in that they are often pressed for preservation with their metadata on the page (see [Figure 13](#)).
2. Fish, worms, and similar will be stored in jars. A single jar is considered an occurrence, so only one specimen may be photographed—perhaps at multiple angles—the jar may be photographed, multiple specimens photographed, etc. In any of these cases, the images will all be labeled with the same metadata that does not include this context.
3. For animal specimens, both of the above cases may occur: there may be simply a close-up of the tag (similar issue to envelopes with plants and fungi). There may also be multiple specimens photographed within the same occurrence (as noted with 2 about fish and worms). Museums often consider multiple specimens of the same species to be connected instead of differentiating them. This is two-fold, in that they are generally, collected from the same or similar location at or about the same time, and a specimen is kept as a representative of its species, so there is not a clear need for distinguishing between them. We also see multiple views of the same specimen within an occurrence. Examples in [Figure 14](#).

In order to appropriately separate these images, we treated them as museums do, specifically by first dividing them into 11 collection areas (Fungi, Insect, Invertebrate Zoology, Microbiology, Plant, Uncategorized, and five classes of Vertebrate Zoology: Amphibians, Birds, Fishes, Mammals, Others) inspired by the [Smithsonian Institution’s categorical subdivisions](#) for their biological museum collections. From here, we further divided each category based on its image type (e.g., fossil or preserved specimen, as specified in GBIF metadata).

For each museum specimen category, we manually curated a small “support set” of representative specimen and non-specimen images. We embedded these examples using CLIP (ViT-L/14@336px) [7]; we chose not to use BIOCLIP since museum specimen labels were not filtered from its training data (e.g., EOL contains many museum specimen images). To capture intra-class diversity, we ran K-Means on each support set and retained the resulting cluster centers. During classification, we L2-normalized both the input image’s embedding and each center, then assigned the image to the nearest center in Euclidean space. This processing was applied to all museum specimen images identified within GBIF. The support set embeddings are included in [TreeOfLife-toolbox](#).

I.2.2 Camera Trap Images

Some occurrences include large-volume camera-trap sequences, with up to 10,000 images. These occurrences have a single taxonomic label applied across all images, though there are different taxa (e.g., the first few frames may have a duck, while later images have a swan, another a goose, but the label for all is a duck). To reduce the risk of introducing such noise while still capturing relevant biodiversity, we filter the dataset to include only occurrences with 15 images or fewer. We then use MegaDetector [47, 48] to filter “empty” frames. In creating a dataset to train a foundation model on the entire tree of life, it is prudent to avoid introducing too many plant images labeled as animals.

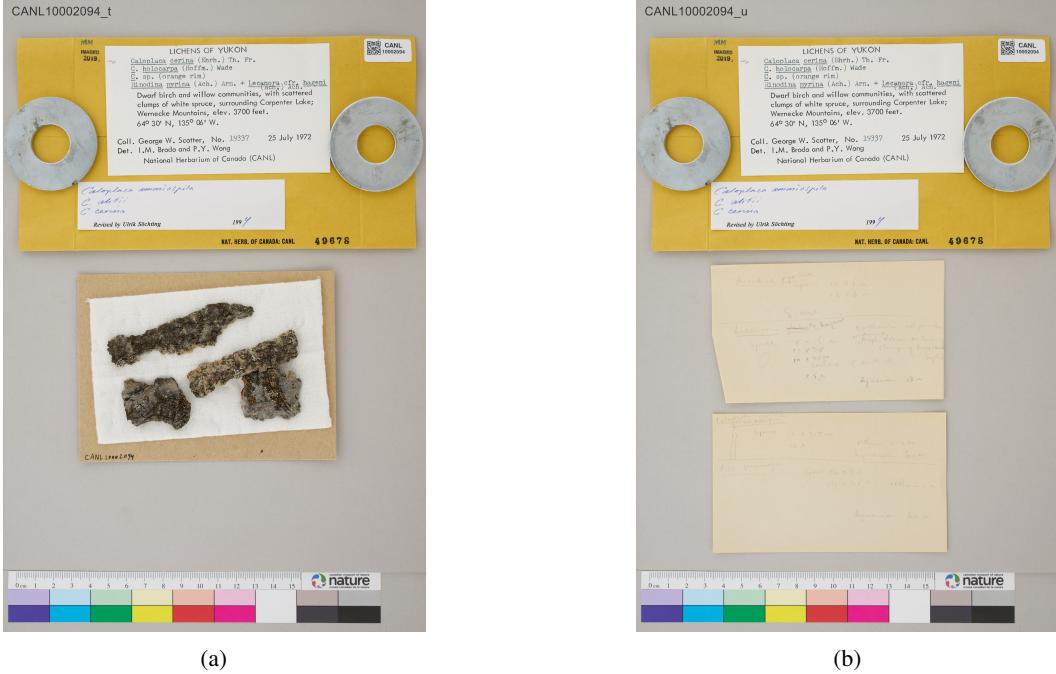


Figure 12: These images all belong to GBIF Occurrence 2301948912 [88], and thus share metadata. (a) Contains the specimen (*Caloplaca cerina*), (b) contains the file metadata from the specimen folder; only (a) should be retained. Collected in Canada by ©Canadian Museum of Nature (licensed under CC BY-NC 4.0).

I.2.3 Citizen Science Images

Similarly, citizen science occurrences may contain many images assigned a single label. For each occurrence, we thus embedded the images using BIOCLIP [13], standardized them, calculated the pair-wise cosine similarity, and used the mean to indicate the occurrence’s “distinctness”. In these instances, we identified three broad categories into which to further subdivide them for reduction:

1. Mixed occurrences: multiple species all labeled as just one of them (low similarity). These were treated as noise and discarded.
2. Multiple image occurrences of the same species: observations suggest these are creature close-ups and environment images under a single occurrence. They may also be larger groups with close-ups. For these, we randomly sample down to 5 images.
3. Camera trap images: likely images collected in backyards (high similarity). Occasionally, people upload images from camera traps to citizen science platforms like iNaturalist (also ex: MammalWeb, a citizen science platform for camera trap images). These are processed under camera trap protocols described above.

The final step in our quality control processing pipeline was to filter out identifiable images of people from the training data. This was done by running MTCNN [49] on all images downloaded from GBIF and EOL.

I.3 Duplicate and Leakage Control

GBIF and EOL are large biodiversity data aggregators, sourcing images that are also used in biological benchmarks such as iNat21 [36] and NeWT [18]. Due to metadata provenance complications, the images in these test sets cannot easily be matched to their sources or copies downloaded from other sources. To prevent the introduction of test data sourced from iNaturalist or “other sources online”, we perform two content deduplication steps. The first is to run MD5 hashes on all downloaded images. During our download, these hash sums are used to distinguish images that have already been downloaded. We record the hash of both the original image and our resized image. These were used to ensure that Rare Species [59] images (sourced from EOL) were not included in the training data.



(a)

(b)

Figure 13: Both images contain specimens included in the dataset, but their dominance in-frame is quite different. The pressed flower sheets are the same size, but (a) *Osmorrhiza longistylis* is dominant in the frame, while (b) *Chara globularis* covers less pixel area than the text. These images belong to [GBIF Occurrence 5132787320](#) and [GBIF Occurrence 5135831095](#) [88], respectively, and were both collected in the United States of America by The New York Botanical Garden (licensed under CC BY 4.0).



(a)

(b)

(c)

(d)

Figure 14: These images all belong to [GBIF Occurrence 1056329684](#) [89], and thus share metadata. (a) through (c) all contain the specimen (*Pristimantis zeuctotylus*), though (c) is primarily focused on its label or tag. Meanwhile, (d) is the label for another specimen of the same species from the collection; there are multiple views of this specimen as well. (c) provides an alternate view and can be retained, while (d) should be removed.

Traditional hash sums are highly sensitive to small changes—a one-pixel difference between two images will produce a different hash. Hence, we applied perceptual hashing [PDQ 50], with distance less than 10, to identify training images that may be in our desired test sets that could not otherwise be filtered out (*i.e.*, by MD5 hash sum or through metadata). Note that PDQ hash evaluation was only run on GBIF citizen science images and EOL images sourced from Flickr since they are not reliable for museum specimens.

IUCN Red List coverage. According to the most recent IUCN Red List assessment[52], TREEOFLIFE-200M contains images of 69.5% (55,512) of all IUCN-assessed species in categories characterized as rare species or data deficient. This coverage was determined by applying our taxonomic alignment package to the IUCN taxonomic data to enable direct comparison with our standardized dataset. TREEOFLIFE-200M demonstrates particularly strong representation of threatened species, with 77.1% coverage across threatened categories (36,370 species), including 79.7% of Vulnerable species (14,038), 79.4% of Endangered species (15,190), and 68.4% of Critically Endangered species (7,142). Coverage extends to 81.7% of Near Threatened species (8,073) and 82.7% of species classified as Extinct in the Wild (67). Data-deficient species have a lower representation at 48.4% (11,002), likely reflecting the challenges in imaging and identifying the

species in this group. Notably, these species are designated in this category because there is not sufficient information about them for IUCN to evaluate their status; only 8% of the 2.14M described species have been evaluated [51]. Thus, including 48% of these species in TREEOFLIFE-200M, along with the threatened species coverage, establishes the approach to integrating diverse data sources used in TREEOFLIFE-200M as a valuable resource for conservation research, providing visual representation for a substantial majority of species prioritized for global conservation action.

J Author Contribution Statement

Jianyang Gu led the research project, ran the experiments, conducted analyses, and wrote the major part of the manuscript. For the dataset processing, he built the face detection pipeline to remove images with recognizable/identifiable people. He also retrieved the common names from GBIF backbone taxonomy alignment to TaxonoPy output.

Samuel Stevens constructed the initial benchmark for non-species classification tasks and significantly contributed to the paper writing. In addition, he also developed the content-based de-duplication pipeline based on PDQ-hash to identify test images that were contained in training data.

Elizabeth G Campolongo analyzed, evaluated, planned, organized, and documented the dataset effort. She developed the data processing approach by categories and supervised the tool development to retrieve and organize the dataset onto the HPC file system and of dataset curation. She curated camera trap test sets for benchmarking and Darwin's Finches for embedding analysis. She also contributed significantly to the dataset part of the paper writing.

Matthew J Thompson planned, organized, and supervised the development of tools to retrieve and organize the dataset onto the HPC file system. He developed and executed the taxonomic standardization (TaxonoPy) and webdataset ML-ready format conversion. He also contributed to the paper writing of the dataset part.

Net Zhang played an essential role in cleaning the noisy data. He dealt with the museum specimen images, citizen science images, and data de-duplication. Besides, he also processed and analyzed the downloaded data, primarily focusing on GBIF.

Jiaman Wu provided detailed advice on designing the training experiments for BIOCLIP 2. She also contributed to the dataset processing by initiating specimen label filtering and informing the taxonomic standardization for hemihomonym cases.

Andrei Kopanov developed the distributed-downloader tool to retrieve and organize the original images in the HPC file system and converted the dataset into ML-ready format (webdataset).

Zheda Mai provided constructive advice on designing the experience replay of CLIP training data.

Alexander E. White provided expertise in consultation on the museum specimen image processing approach.

James Balhoff and **Wasila Dahdul** provided expertise in consultation on the algorithm for taxonomic hierarchy resolution relevant to TaxonoPy development.

Daniel Rubenstein and **Hilmar Lapp** were involved in regular meetings and gave valuable feedback on results and experiments.

Tanya Berger-Wolf and **Wei-Lun Chao** provided insightful discussions and directions throughout the development of the project. They offered constructive advice for both the training design and the model evaluation. They also provided detailed comments on the manuscript.

Yu Su is the senior lead that oversaw the project and was involved in every aspect. He conceived the idea of observing emergent properties from scaling hierarchical contrastive training. He gave insightful guidance in formalizing the research question and directions to verify the new capabilities of BIOCLIP 2. He provided critical feedback and helped shape the manuscript.