
What Do You See in Common?

Learning Hierarchical Prototypes over Tree-of-Life to Discover Evolutionary Traits

Anonymous Author(s)

Affiliation

Address

email

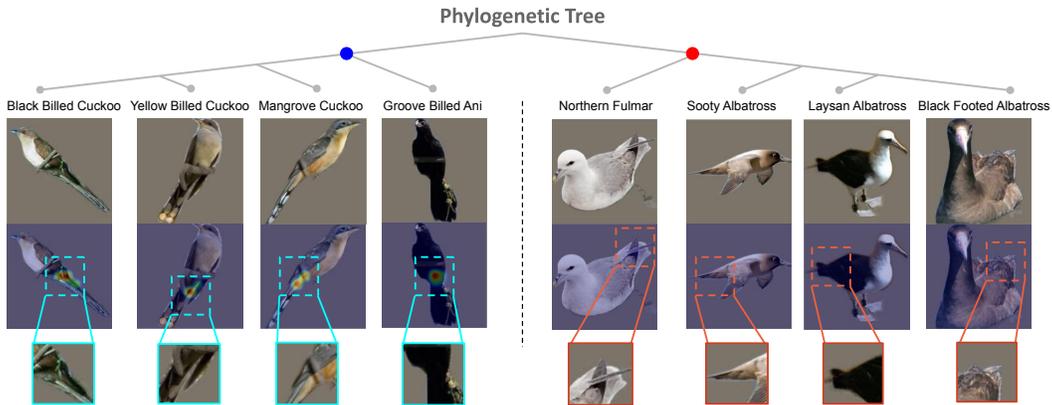


Figure 1: Sample images of bird species with zoomed-in views of learned prototypes along with their associated score maps. We consider the problem of finding evolutionary traits common to a group of species derived from the same ancestor (blue) that are absent in other species from a different ancestor (red). We can infer that descendants of the blue node share a common trait: *long tail*, absent from descendants of the red node.

Abstract

1 A grand challenge in biology is to discover evolutionary traits—features of organ-
2 isms common to a group of species with a shared ancestor in the tree of life (also
3 referred to as phylogenetic tree). With the growing availability of image repositories
4 in biology, there is a tremendous opportunity to discover evolutionary traits directly
5 from images in the form of a hierarchy of prototypes. However, current prototype-
6 based methods are mostly designed to operate over a flat structure of classes and
7 face several challenges in discovering hierarchical prototypes, including the issue of
8 learning over-specific features at internal nodes. To overcome these challenges, we
9 introduce the framework of **H**ierarchy aligned **C**ommonality through **P**rototypical
10 **N**etworks (**HComP-Net**). We empirically show that HComP-Net learns prototypes
11 that are accurate, semantically consistent, and generalizable to unseen species in
12 comparison to baselines on birds, butterflies, and fishes datasets.

13 1 Introduction

14 A central goal in biology is to discover the observable characteristics of organisms, or *traits* (e.g., beak
15 color, stripe pattern, and fin curvature), that help in discriminating between species and understanding

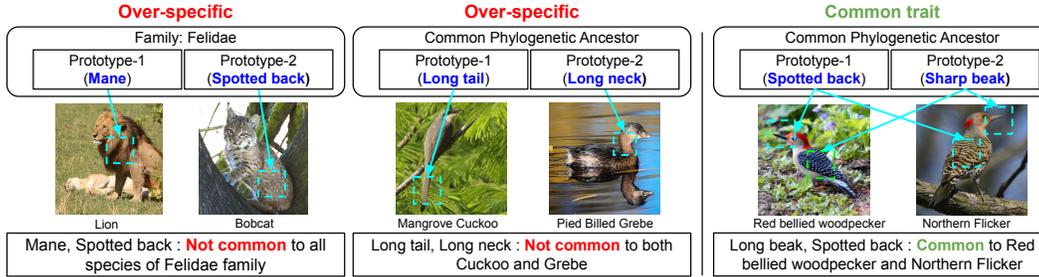


Figure 2: Examples to illustrate the problem of learning “over-specific” prototypes at internal nodes, which only cover one descendant species of the node instead of learning prototypes *common* to all descendants.

16 how organisms evolve and adapt to their environment [1]. For example, discovering traits inherited by
 17 a group of species that share a common ancestor on the tree of life (also referred to as the *phylogenetic*
 18 *tree*, see Figure 1) is of great interest to biologists to understand how organisms diversify and evolve
 19 [2]. The measurement of such traits with evolutionary signals, termed *evolutionary traits*, is not
 20 straightforward and often relies on subjective and labor-intensive human expertise and definitions
 21 [3, 4], hindering rapid scientific advancement [5].

22 With the growing availability of large-scale image repositories in biology containing millions of
 23 images of organisms [6, 7, 8], there is an opportunity for machine learning (ML) methods to discover
 24 evolutionary traits automatically from images [5, 9]. This is especially true in light of recent advances
 25 in the field of explainable ML, such as the seminal work of ProtoPNet [10] and its variants [11, 12, 13]
 26 which find representative patches in training images (termed *prototypes*) capturing discriminatory
 27 features for every class. We can thus cast the problem of discovering evolutionary traits into asking
 28 the following question: *what image features or prototypes are common across a group of species*
 29 *with a shared ancestor in the tree of life that are absent in species with a different shared ancestor?*

30 For example, in Figure 1, we can see that the four species of birds on the left descending from the
 31 blue node show the common feature of having “long tails,” unlike any of the descendant species of
 32 the red node. Learning such common features at every internal node as a hierarchy of prototypes can
 33 help biologists generate novel hypotheses of species diversification (e.g., the splitting of blue and red
 34 nodes) and accumulation of evolutionary trait changes.

35 Despite the success of ProtoPNet [10] and its variants in learning prototypes over a flat structure of
 36 classes, applying them to discover a hierarchy of prototypes is challenging for three main reasons.
 37 *First*, existing methods that learn multiple prototypes for every class are prone to learning “over-
 38 specific” prototypes at internal nodes of a tree, which cover only one (or a few) of its descendant
 39 species. Figure 2 shows a few examples to illustrate the concept of over-specific prototypes. Consider
 40 the problem of learning prototypes common to descendant species of the Felidae family: Lion and
 41 Bobcat. If we learn one prototype focusing on the feature of the mane (specific only to Lion) and
 42 another prototype focusing on the feature of spotted back (specific only to Bobcat), then these two
 43 prototypes taken together can classify all images from the Felidae family. However, they do not
 44 represent *common* features shared between Lion and Bobcat and hence are not useful for discovering
 45 evolutionary traits. Such over-specific prototypes should be instead pushed down to be learned at
 46 lower levels of the tree (e.g., the species leaf nodes of Lion and Bobcat).

47 *Second*, while existing methods such as ProtoPShare [11], ProtoPool [12], and ProtoTree [13] allow
 48 prototypes to be shared across classes for re-usability and sparsity, in the problem of discovering
 49 evolutionary traits, we want to learn prototypes at an internal node n that are not just shared across
 50 all its descendant species but are also absent in the *contrasting set* of species (i.e., species descending
 51 from sibling nodes of n representing alternate paths of diversification). *Third*, at higher levels of the
 52 tree, finding features that are common across a large number of diverse species is challenging [14, 15].
 53 In such cases, we should be able to abstain from finding common prototypes without hampering
 54 accuracy at the leaf nodes—a feature missing in existing methods.

55 To address these challenges, we present **Hierarchy aligned Commonality through Prototypical**
 56 **Networks (HComp-Net)**, a framework to learn hierarchical prototypes over the tree of life for
 57 discovering evolutionary traits. Here are the main contributions of our work:

- 58 1. HComP-Net learns common traits shared by all descendant species of an internal node and
59 avoids the learning of over-specific prototypes in contrast to baseline methods using a novel
60 *overspecificity loss*.
- 61 2. HComP-Net uses a novel *discriminative loss* to ensure that the prototypes learned at an
62 internal node are absent in the contrasting set of species with different ancestry.
- 63 3. HComP-Net includes a novel *masking module* to allow for the exclusion of over-specific
64 prototypes at higher levels of the tree without hampering classification performance.
- 65 4. We empirically show that HComP-Net learns prototypes that are accurate, semantically
66 consistent, and generalizable to unseen species compared to baselines on data from 190
67 species of birds (CUB-200-2011 dataset) [8], 38 species of fishes [9], and 30 species of
68 butterflies [16]. We show the ability of HComP-Net to generate novel hypotheses about
69 evolutionary traits at different levels of the phylogenetic tree of organisms.

70 2 Related Works

71 One of the seminal lines of work in the field of prototype-based interpretability methods is the
72 framework of ProtoPNet [10] that learns a set of “prototypical patches” from training images of every
73 class to enable case-based reasoning. Following this work, several variants have been developed
74 such as ProtoPShare [11], ProtoPool [12], ProtoTree [13], and HPnet [17] suiting to different
75 interpretability requirements. Among all these approaches, our work is closely related to HPnet [17],
76 the hierarchical extension of ProtoPNet that learns a prototype layer for every parent node in the
77 tree. Despite sharing a similar motivation as our work, HPnet is not designed to avoid the learning of
78 over-specific prototypes or to abstain from learning common prototypes at higher levels of the tree.

79 Another related line of work is the framework of PIPNet [18], which uses self-supervised learning
80 methods to reduce the “semantic gap” [19, 20] between the latent space of prototypes and the space
81 of images, such that the prototypes in latent space correspond to the same visual concept in the image
82 space. In HComP-Net, we build upon the idea of self-supervised learning introduced in PIPNet to
83 learn semantically consistent hierarchy of prototypes. Our work is also related to ProtoTree [13],
84 which structures the prototypes as nodes in a decision tree to offer more granular interpretability.
85 However, ProtoTree differs from our work in that it learns the tree-based structure of prototypes
86 automatically from data and cannot handle a known hierarchy. Moreover, the prototypes learned in
87 ProtoTree are purely discriminative and allow for negative reasoning, which is not aligned with our
88 objective of finding common traits of descendant species.

89 Other related works that focus on finding shared features are ProtoPShare [11] and ProtoPool [12].
90 Both approaches aim to find common features among classes, but their primary goal is to reduce
91 the prototype count by exploiting similarities among classes, leading to a sparser network. This is
92 different from our goal of finding a hierarchy of prototypes to find evolutionary traits common to a
93 group of species (that are absent from other species).

94 Outside the realm of prototype-based methods, the framework of Phylogeny-guided Neural Networks
95 (PhyloNN) [9] shares a similar motivation as our work to discover evolutionary traits by representing
96 biological images in feature spaces structured by tree-based knowledge (i.e., phylogeny). However,
97 PhyloNN primarily focuses on the tasks of image generation and translation rather than interpretability.
98 Additionally, PhyloNN can only work with discretized trees with fixed number of ancestor levels per
99 leaf node, unlike our work that does not require any discretization of the tree.

100 3 Proposed Methodology

101 3.1 HComP-Net Model Architecture

102 Given a phylogenetic tree with N internal nodes, the goal of HComP-Net is to jointly learn a set of
103 prototype vectors \mathbf{P}_n for every internal node $n \in \{1, \dots, N\}$. Our architecture as shown in Figure 3
104 begins with a CNN that acts as a common feature extractor $f(x; \theta)$ for all nodes, where θ represents
105 the learnable parameters of f . f converts an image x into a latent representation $Z \in \mathbb{R}^{H \times W \times C}$,
106 where each “patch” at location (h, w) is, $\mathbf{z}_{h,w} \in \mathbb{R}^C$. Following the feature extractor, for every node
107 n , we initialize a set of K_n prototype vectors $\mathbf{P}_n = \{\mathbf{p}_i\}_{i=1}^{K_n}$, where $\mathbf{p}_i \in \mathbb{R}^C$. Here, the number of

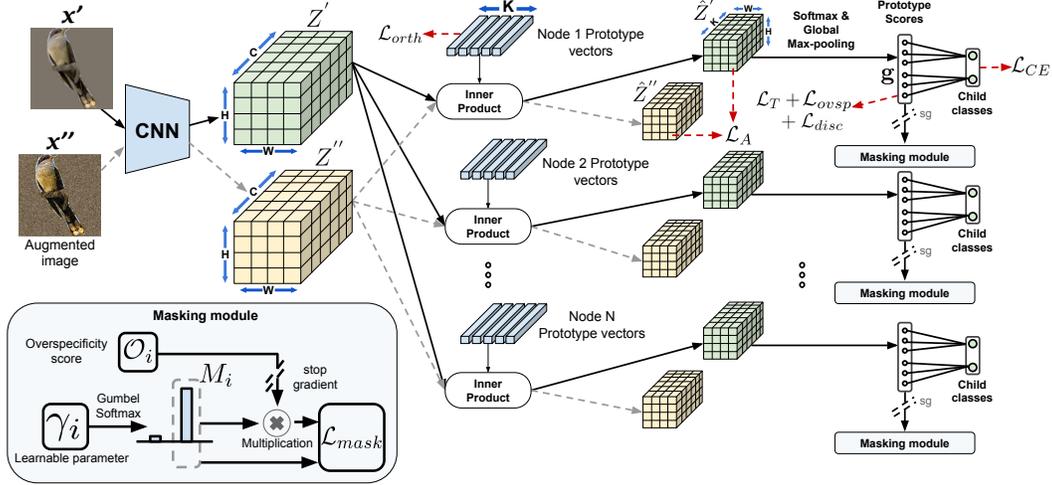


Figure 3: Schematic illustration of HCompP-Net model architecture.

108 prototypes K_n learned at node n varies in proportion to the number of children of node n , with β
 109 as the proportionality constant, i.e., at each node n we assign β prototypes for every child node. To
 110 simplify notations, we drop the subscript n in \mathbf{P}_n and K_n while discussing the operations occurring
 111 in node n .

112 We consider the following sequence of operations at every node n . We first compute the similarity
 113 score between every prototype in \mathbf{P} and every patch in Z . This results in a matrix $\hat{Z} \in \mathbb{R}^{H \times W \times K}$,
 114 where every element represents a similarity score between image patches and prototype vectors. We
 115 apply a softmax operation across the K channels of \hat{Z} such that the vector $\hat{z}_{h,w} \in \mathbb{R}^K$ at spatial
 116 location (h, w) in \hat{Z} represents the probability that the corresponding patch $\mathbf{z}_{h,w}$ is similar to the K
 117 prototypes. Furthermore, the i^{th} channel of \hat{Z} serves as a prototype score map for the prototype
 118 vector \mathbf{p}_i , indicating the presence of \mathbf{p}_i in the image. We perform global max-pooling across the
 119 spatial dimensions $H \times W$ of \hat{Z} to obtain a vector $\mathbf{g} \in \mathbb{R}^K$, where the i^{th} element represents the
 120 highest similarity score of the prototype vector \mathbf{p}_i across the entire image. \mathbf{g} is then fed to a linear
 121 classification layer with weights ϕ to produce the final classification scores for every child node of
 122 node n . We restrict the connections in the classification layer so that every child node n_c is connected
 123 to a distinct set of β prototypes, to ensure that every prototype uniquely maps to a child node. ϕ is
 124 restricted to be non-negative to ensure that the classification is done solely through positive reasoning,
 125 similar to the approach used in PIP-Net [18]. We borrow the regularization scheme of PIP-Net to
 126 induce sparsity in ϕ by computing the logit of child node n_c as $\log((\mathbf{g}\phi)^2 + 1)$. \mathbf{g} and ϕ here are
 127 again unique to each node.

128 3.2 Loss Functions Used to Train HCompP-Net

129 **Contrastive Losses for Learning Hierarchical Prototypes:** PIP-Net [18] introduced the idea of
 130 using self-supervised contrastive learning to learn semantically meaningful prototypes. We build
 131 upon this idea in our work to learn semantically meaningful hierarchical prototypes at every node
 132 in the tree as follows. For every input image \mathbf{x} , we pass in two augmentations of the image, \mathbf{x}' and
 133 \mathbf{x}'' to our framework. The prototype score maps for the two augmentations, \hat{Z}' and \hat{Z}'' , are then
 134 considered as positive pairs. Since $\hat{z}_{h,w} \in \mathbb{R}^K$ represents the probabilities of patch $\mathbf{z}_{h,w}$ being similar
 135 to the prototypes from \mathbf{P} , we align the probabilities from the two augmentations $\hat{z}'_{h,w}$ and $\hat{z}''_{h,w}$ to be
 136 similar using the following alignment loss:

$$137 \mathcal{L}_A = -\frac{1}{HW} \sum_{(h,w) \in H \times W} \log(\hat{z}'_{h,w} \cdot \hat{z}''_{h,w}) \quad (1)$$

137 Since $\sum_{i=1}^K \hat{z}_{h,w,i} = 1$ due to softmax operation, \mathcal{L}_A is minimum (i.e., $\mathcal{L}_A = 0$) when both $\hat{z}'_{h,w}$
 138 and $\hat{z}''_{h,w}$ are identical one-hot encoded vectors. A trivial solution that minimizes \mathcal{L}_A is when all

139 patches across all images are similar to the same prototype. To avoid such representation collapse, we
 140 use the following tanh-loss \mathcal{L}_T of PIP-Net [18], which serves the same purpose as uniformity losses
 141 in [21] and [22]:

$$\mathcal{L}_T = -\frac{1}{K} \sum_{i=1}^K \log(\tanh(\sum_{b=1}^B \mathbf{g}_{b,i})), \quad (2)$$

142 where $\mathbf{g}_{b,i}$ is the prototype score for prototype i with respect to image b of mini-batch. \mathcal{L}_T encourages
 143 each prototype \mathbf{p}_i to be activated at least once in a given mini-batch of B images, thereby helping to
 144 avoid the possibility of representation collapse. The use of tanh ensures that only the presence of a
 145 prototype is taken into account and not its frequency.

146 **Over-specificity Loss:** To achieve the goal of learning prototypes common to all descendant species
 147 of an internal node, we introduce a novel loss, termed *over-specificity loss* \mathcal{L}_{ovsp} that avoids learning
 148 over-specific prototypes at any node n . \mathcal{L}_{ovsp} is formulated as a modification of the tanh-loss such
 149 that prototype \mathbf{p}_i is encouraged to be activated at least once in every one of the descendant species
 150 $d \in \{1, \dots, D_i\}$ of its corresponding child node in the mini-batch of images fed to the model, as
 151 follows:

$$\mathcal{L}_{ovsp} = -\frac{1}{K} \sum_{i=1}^K \sum_{d=1}^{D_i} \log(\tanh(\sum_{b \in B_d} \mathbf{g}_{b,i})), \quad (3)$$

152 where B_d is the subset of images in the mini-batch that belong to species d .

153 **Discriminative loss:** In order to ensure that a learned prototype for a child node n_c is not activated
 154 by any of its *contrasting set* of species (i.e., species that are descendants of child nodes of n other
 155 than n_c), we introduce another novel loss function, \mathcal{L}_{disc} , defined as follows:

$$\mathcal{L}_{disc} = \frac{1}{K} \sum_{i=1}^K \sum_{d \in \widetilde{D}_i} \max_{b \in B_d}(\mathbf{g}_{b,i}), \quad (4)$$

156 where \widetilde{D}_i is the contrasting set of all descendant species of child nodes of n other than n_c . This is
 157 similar to the separation loss used in other prototype-based methods such as [10], [13], and [23].

158 **Orthogonality loss:** We also apply kernel orthogonality as introduced in [24] to the prototype vectors
 159 at every node n , so that the learned prototypes are orthogonal and capture diverse features:

$$\mathcal{L}_{orth} = \|\hat{\mathbf{P}}\hat{\mathbf{P}}^\top - I\|_F^2 \quad (5)$$

160 where $\hat{\mathbf{P}}$ is the matrix of normalized prototype vectors of size $C \times K$, I is an identity matrix, and
 161 $\|\cdot\|_F^2$ is the Frobenius norm. Each prototype $\hat{\mathbf{p}}_i$ in $\hat{\mathbf{P}}$ is normalized as, $\hat{\mathbf{p}}_i = \frac{\mathbf{p}_i}{\|\mathbf{p}_i\|}$.

162 **Classification loss:** Finally, we apply cross entropy loss for classification at each internal node as
 163 follows:

$$\mathcal{L}_{CE} = -\sum_b^B y_b \log(\hat{y}_b) \quad (6)$$

164 where y is ground truth label and \hat{y} is the prediction at every node of the tree.

165 3.3 Masking Module to Identify Over-specific Prototypes

166 We employ an additional masking module at every node n to identify over-specific prototypes without
 167 hampering their training. The learned mask for prototype \mathbf{p}_i simply serves as an indicator of whether
 168 \mathbf{p}_i is over-specific or not, enabling our approach to abstain from finding common prototypes if there
 169 are none, especially at higher levels of the tree. To obtain the mask values, we first calculate the
 170 over-specificity score for prototype \mathbf{p}_i as the product of the maximum prototype scores obtained
 171 across all images in the mini-batch belonging to every descendant species d as:

$$\mathcal{O}_i = -\prod_{d=1}^{D_i} \max_{(b \in B_d)}(\mathbf{g}_{b,i}) \quad (7)$$

172 where $\mathbf{g}_{b,i}$ is the prototype score for prototype \mathbf{p}_i with respect to image b of mini-batch and B_d
 173 is the subset of images in the mini-batch that belong to descendant species d . Since $\mathbf{g}_{b,i}$ takes a
 174 value between 0 to 1 due to the softmax operation, \mathcal{O}_i ranges from -1 to 0, where -1 denotes least

175 over-specificity and 0 denotes the most over-specificity. The multiplication of the prototype scores
 176 ensures that even when the score is less with respect to only one descendant species, the prototype
 177 will be assigned a high over-specificity score (close to 0).

178 As shown in Figure 3, \mathcal{O}_i is then fed into the masking module, which includes a learned mask value
 179 M_i for every prototype \mathbf{p}_i . We generate M_i from a Gumbel-softmax distribution [25] so that the
 180 values are skewed to be very close to either 0 or 1, i.e., $M_i = \text{Gumbel-Softmax}(\gamma_i, \tau)$, where γ_i are
 181 the learnable parameters of the distribution and τ is temperature. We then compute the masking loss,
 182 \mathcal{L}_{mask} , as:

$$\mathcal{L}_{mask} = \sum_{i=1}^K (\lambda_{mask} M_i \circ \text{stopgrad}(\mathcal{O}_i) + \lambda_{L_1} \|M_i\|_1) \quad (8)$$

183 where λ_{mask} and λ_{L_1} are trade-off coefficients, $\|\cdot\|_1$ is the L_1 norm added to induce sparsity in
 184 the masks, and stopgrad represents the stop gradient operation applied over \mathcal{O}_i to ensure that the
 185 gradient of \mathcal{L}_{mask} does not flow back to the learning of prototype vectors and impact their training.
 186 Note that *the learned masks are not used for pruning the prototypes during training*, they are only
 187 used during inference to determine which of the learned prototypes are over-specific and likely to not
 188 represent evolutionary traits. Therefore, even if all the prototypes are identified as over-specific by
 189 the masking module at an internal node, it will not affect the classification performance at that node.

190 3.4 Training HComP-Net

191 We first pre-train the prototypes at every internal node in a self-supervised learning manner using
 192 alignment and tanh-losses as $\mathcal{L}_{SS} = \lambda_A \mathcal{L}_A + \lambda_T \mathcal{L}_T$. We then fine-tune the model using the following
 193 combined loss: $(\lambda_{CE} \mathcal{L}_{CE} + \mathcal{L}_{SS} + \lambda_{ovsp} \mathcal{L}_{ovsp} + \lambda_{disc} \mathcal{L}_{disc} + \lambda_{orth} \mathcal{L}_{orth} + \mathcal{L}_{mask})$, where λ 's are
 194 trade-off parameters. Note that the loss is applied over every node in the tree. We show an ablation of
 195 key loss terms in our framework in Table 6 in the Supplementary Section.

196 4 Experimental Setup

197 **Dataset:** In our experiments, we primarily focus on the 190 species of birds (**Bird**) from the CUB-200-
 198 2011 [8] dataset for which the phylogenetic relationship [26] is known. The tree is quite large with a
 199 total of 184 internal nodes. We removed the background from the images to avoid the possibility of
 200 learning prototypes corresponding to background information such as the bird’s habitat as we are
 201 only interested in the traits corresponding to the body of the organism. We also apply our method on
 202 a fish dataset with 38 species (**Fish**) [9] along with its associated phylogeny [9] and 30 subspecies
 203 of Heliconius butterflies (**Butterfly**) from the Jiggins Heliconius Collection dataset [16] collected
 204 from various sources ¹ along with its phylogeny [52, 53]. The qualitative results of Butterfly and
 205 Fish datasets are provided in the supplementary materials. The complete details of hyper-parameter
 206 settings and training strategy are also provided in the Supplementary Section E.

207 **Baselines:** We compare HComP-Net to ResNet-50 [54], INTR (Interpretable Transformer) [55] and
 208 HPnet [17]. For HPnet, we used the same hyperparameter settings and training strategy as used by
 209 ProtoPNet for CUB-200-2011 dataset. For a fair comparison, we also set the number of prototypes
 210 for each child in HPnet to be equal to 10 similar to our implementation. We follow the same training
 211 strategy as provided by ProtoPNet for CUB-200-2011 dataset.

212 5 Results

213 5.1 Fine-grained Accuracy

214 Similar to HPnet [17], we calculate the fine-grained accuracy for each leaf node by calculating the
 215 path probability over every image. During inference, the final probability for leaf class Y given
 216 an image X is calculated as, $P(Y|X) = P(Y^{(1)}, Y^{(2)}, \dots, Y^{(L)}|X) = \prod_{l=1}^L P(Y^{(l)}|X)$, where
 217 $P(Y^{(l)}|X)$ is the probability of assigning image X to a node at level l , and L is the depth of the
 218 leaf node. Every image is assigned to the leaf class with maximum path probability, which is used
 219 to compute the fine-grained accuracy. The comparison of the fine-grained accuracy calculated for

¹Sources: [27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51]

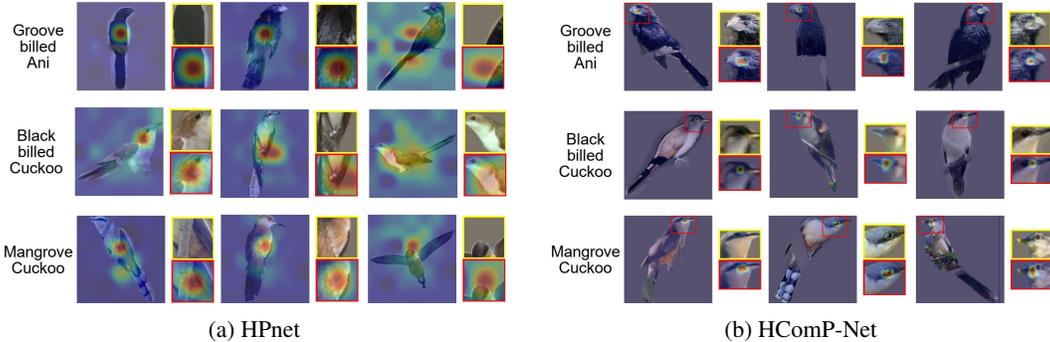


Figure 4: Comparing the part consistency of HPnet and HComP-Net for their prototype learned at an internal node in the bird dataset that corresponds to 3 descendant species (names shown on the rows). For every species, we are visualizing the top-3 images with highest prototype score for both HPnet and HComP-Net, shown as the four columns with zoomed in views of their discovered prototypes. We can see that *HPnet highlights varying parts of the bird* across the 3 species and across multiple images of the same species, making it difficult to associate a consistent semantic meaning to its learned prototype. In contrast, *HComP-Net consistently highlights the head region* of the bird across all four species and their images.

220 HComP-Net and the baselines are given in Table 1. We can see that HComP-Net performs better
 221 than the other interpretable methods, such as INTR and HPNet, and is also able to nearly match the
 222 performance of non-interpretable models, such as ResNet-50, even outperforming it for the Fish
 223 and Butterfly dataset. This shows the ability of our proposed framework to achieve competitive
 224 classification accuracy along with serving the goal of discovering evolutionary traits.

| Model | Hierarchy | Bird | Butterfly | Fish |
|-----------|-----------|--------------|--------------|--------------|
| ResNet-50 | No | 74.18 | 95.76 | 86.63 |
| INTR | | 69.22 | 95.53 | 86.73 |
| HPnet | Yes | 36.18 | 94.69 | 77.51 |
| HComP-Net | | 70.01 | 97.35 | 90.80 |

| Species Name | HComP-Net | HPnet |
|------------------|-----------|-------|
| Fish Crow | 53.33 | 10.55 |
| Rock Wren | 53.33 | 10.22 |
| Indigo Bunting | 96.67 | 49.2 |
| Bohemian Waxwing | 70.00 | 44.9 |

225 5.2 Generalizing to Unseen Species in the Phylogeny

226 We analyze the performance of HComP-Net in generalizing to unseen species that the model hasn't
 227 seen during training. The biological motivation for this experiment is to evaluate if HComP-Net
 228 can situate newly discovered species at its appropriate position in the phylogeny by identifying its
 229 common ancestors shared with the known species. An added advantage of our work is that along with
 230 identifying the ancestor of an unseen species, we can also identify the common traits shared by the
 231 novel species with known species in the phylogeny. Since unseen species cannot be classified to the
 232 finest levels (i.e., up to the leaf node corresponding to the unseen species), we analyze the ability of
 233 HComP-Net to classify unseen species accurately up to one level above the leaf level in the hierarchy.
 234 With this consideration, the final probability of an unseen species for a given image is calculated
 235 as, $P(Y|X_{unseen}) = P(Y^{(1)}, Y^{(2)}, \dots, Y^{(L-1)}|X) = \prod_{l=1}^{L-1} P(Y^{(l)}|X)$. Note that we leave out the
 236 class probability at the L^{th} level, since we do not take into account the class probability of the leaf
 237 level. We leave four species from the Bird training set and calculate their accuracy during inference
 238 in Table 2. We can see that HComP-Net is able to generalize better than HPnet for all four species.

239 5.3 Analyzing the Semantic Quality of Prototypes

240 Following the method introduced in PIPNet [18], we assess the semantic quality of our learned
 241 prototypes by evaluating their part purity. A prototype with high part purity (close to 1) is one that
 242 consistently highlights the same image region in the score maps (corresponding to consistent local
 243 features such as the eye or wing of a bird) across images belonging to the same class. The part

244 purity is calculated using the part locations of
 245 15 parts that are provided in the CUB dataset.
 246 For each prototype, we take the top-10 im-
 247 ages from each leaf descendant. We con-
 248 sider the 32×32 image patch that is centered
 249 around the max activation location of the pro-
 250 totype from the top-10 images. With these
 251 top-10 image patches, we calculate for each
 252 part how frequently the part is present inside
 253 the image patch. For example, a part that is found inside the image patch 8 out of 10 times is given a
 254 score of 0.8. In PIP-Net, the highest value among the values calculated for each part is given as the
 255 part purity of the prototype. In our approach, since we are dealing with a hierarchy and taking the
 256 top-10 from each leaf descendant, a particular part, let’s say the eye, might have a score of 0.5 for
 257 one leaf descendant and 0.7 for a different leaf descendant. Since we want the prototype to represent
 258 the same part for all the leaf descendants, we take the lowest score (the weakest link) among all the
 259 leaf descendants as the score of the part. By following this method, for a given prototype we can
 260 arrive at a value for each part and finally take the maximum among the values as the purity of the
 261 prototype. We take the mean of the part purity across all the prototypes and report the results in Table
 262 3 for different ablations of HComP-Net and also HPnet, which is the only baseline method that can
 263 learn hierarchical prototypes.

Table 3: Part purity of prototypes on **Bird** dataset.

| Model | \mathcal{L}_{ovsp} | Masking | Part purity | % masked |
|-----------|----------------------|---------|-----------------------------------|----------|
| HPnet | - | - | 0.14 ± 0.09 | - |
| HComP-Net | - | - | 0.68 ± 0.22 | - |
| HComP-Net | - | ✓ | 0.75 ± 0.17 | 21.42% |
| HComP-Net | ✓ | - | 0.72 ± 0.19 | - |
| HComP-Net | ✓ | ✓ | 0.77 ± 0.16 | 16.53% |

264 We can see that HComP-Net, even without the use of over-specificity loss performs much better than
 265 HPnet due to the contrastive learning approach we have adopted from PIPNet [18]. The addition
 266 of over-specificity loss improves the part purity because over-specific prototypes tend to have poor
 267 part purity for some of the leaf descendants which will affect their overall part purity score. Further,
 268 for both ablations with and without over-specificity loss, we apply the masking module and remove
 269 masked (over-specific) prototypes during the calculation of part purity. We see that the part purity goes
 270 higher by applying the masking module, demonstrating its effectiveness in identifying over-specific
 271 prototypes. We further compute the purity of masked-out prototypes and notice that the masked-out
 272 prototypes have drastically lower part purity (0.29 ± 0.17) compared to non-masked prototypes
 273 (0.77 ± 0.16). An alternative approach to learning the masking module is to identify over-specific
 274 prototypes using a fixed global threshold over \mathcal{O}_i . We show in Table 9 of Supplementary Section F,
 275 that given the right choice of such a threshold, we can identify over-specific prototypes. However,
 276 selecting the ideal threshold can be non-trivial. On the other hand, our masking module learns the
 277 appropriate threshold dynamically as part of the training process.

278 Figure 4 visualizes the part consistency of prototypes discovered by HComP-Net in comparison to
 279 HPnet for the bird dataset. We can see that HComP-Net is finding a consistent region in the image
 280 (corresponding to the head region) across all three descendant species and all images of a species, in
 281 contrast to HPnet. Furthermore, thanks to the alignment loss, every patch $\hat{\mathbf{z}}_{h,w}$ is encoded as nearly
 282 a one-hot encoding with respect to the K prototypes which causes the prototype score maps to be
 283 highly localized. The concise and focused nature of the prototype score maps makes the interpretation
 284 much more effective compared to baselines.

285 5.4 Analyzing Evolutionary Traits Discovered by HComP-Net

286 We now qualitatively analyze some of the hypothesized evolutionary traits discovered in the hierarchy
 287 of prototypes learned by HComP-Net. Figure 5 shows the hierarchy of prototypes discovered over
 288 a small subtree of the phylogeny from Bird (four species) and Fish (three species) dataset. In the
 289 visualization of bird prototypes, we can see that the two Pelican species share a consistent region in the
 290 learned Prototype labeled 2, which corresponds to the head region of the birds. We can hypothesize
 291 this prototype to be capturing the white colored crown common to the two species. On the other hand,
 292 Prototype 1 finds the shared trait of similar beak morphology (e.g., sharpness of beaks) across the
 293 two Cormorant species. We can see that HComP-Net avoids the learning of over-specific prototypes
 294 at internal nodes, which are pushed down to individual leaf nodes, as shown in visualizations of
 295 Prototype 3, 4, 5, and 6. Similarly, in the visualization of the fish prototypes, we can see that Prototype
 296 1 is highlighting a specific fin (dorsal fin) of the *Carassius auratus* and *Notropis hudsonius* species,
 297 possibly representing their pigmentation and structure, which is noticeably different compared to
 298 the contrasting species of *Alosa chrysochloris*. Note that while HComP-Net identifies the common

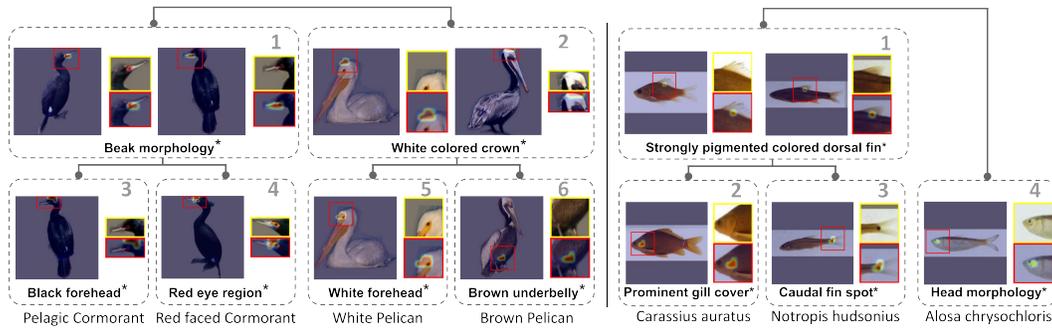


Figure 5: Visualizing the hierarchy of prototypes discovered by HComP-Net for birds and fishes. *Note that the textual descriptions of the hypothesized traits shown for every prototype are based on human interpretation.

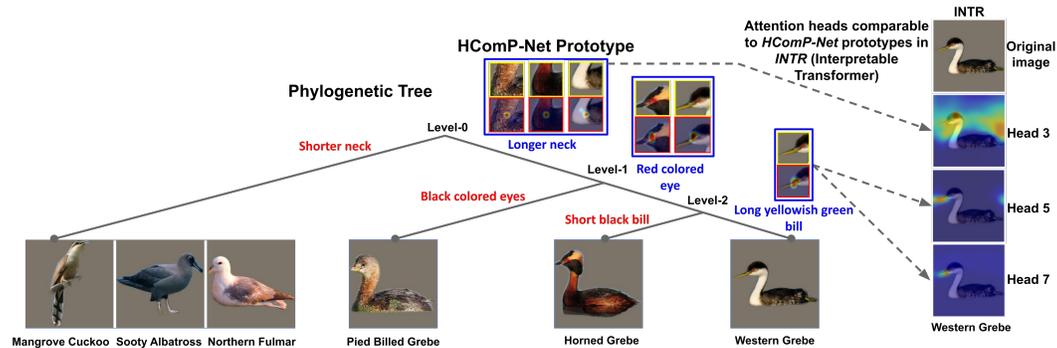


Figure 6: We trace the prototypes learned for Western Grebe at three different levels in the phylogenetic tree (corresponding to different periods of time in evolution). Text in blue is the interpretation of common traits of descendants found by HComP-Net at every ancestor node of Western Grebe.

299 regions corresponding to each prototype (shown as heatmaps), the textual descriptions of the traits
 300 provided in Figure 5 are based on human interpretation.

301 Figure 6 shows another visualization of the sequence of prototypes learned by HComP-Net for the
 302 Western Grebe species at different levels of the phylogeny. We can see that at level 0, we are capturing
 303 features closer to the neck region, indicating the likely difference between the length of necks between
 304 Grebe species and other species (Cuckoo, Albatross, and Fulmar) that diversify at an earlier time in
 305 the process of evolution. At level 1, the prototype is focusing on the eye region, potentially indicating
 306 to difference in the color of red and black patterns around the eyes. At level 2, we are differentiating
 307 Western Grebe from Horned Grebe based on the feature of bills. We also validate our prototypes by
 308 comparing them with the multi-head cross-attention maps learned by INTR [55]. We can see that
 309 some of the prototypes discovered by HComP-Net can be mapped to equivalent attention heads of
 310 INTR. However, while INTR is designed to produce a flat structure of attention maps, we are able
 311 to place these maps on the tree of life. This shows the power of HComP-Net in generating novel
 312 hypotheses about how trait changes may have evolved and accumulated across different branches of
 313 the phylogeny. Additional visualizations of discovered evolutionary traits for butterfly species and
 314 fish species are provided in the supplementary section in Figures 7 to 16.

315 6 Conclusion

316 We introduce a novel approach for learning hierarchy-aligned prototypes while avoiding the learning
 317 of over-specific features at internal nodes of the phylogenetic tree, enabling the discovery of novel
 318 evolutionary traits. Our empirical analysis on birds, fishes, and butterflies, demonstrates the efficacy
 319 of HComP-Net over baseline methods. Furthermore, HComP-Net demonstrates a unique ability
 320 to generate novel hypotheses about evolutionary traits, showcasing its potential in advancing our
 321 understanding of evolution. We discuss the limitations of our work in Supplementary Section I. While
 322 we focus on the biological problem of discovering evolutionary traits, our work can be applied in
 323 general to domains involving a hierarchy of classes, which can be explored in future research.

324 **References**

- 325 [1] David Houle and Daniela M Rossoni. Complexity, evolvability, and the process of adaptation.
326 *Annual Review of Ecology, Evolution, and Systematics*, 53, 2022.
- 327 [2] Maureen A O’Leary and Seth Kaufman. Morphobank: phylophenomics in the “cloud”. *Cladis-*
328 *tics*, 27(5):529–537, 2011.
- 329 [3] Tiago R Simões, Michael W Caldwell, Alessandro Palci, and Randall L Nydam. Giant taxon-
330 character matrices: quality of character constructions remains critical regardless of size. *Cladis-*
331 *tics*, 33(2):198–219, 2017.
- 332 [4] Paul C Sereno. Logical basis for morphological characters in phylogenetics. *Cladistics*,
333 23(6):565–587, 2007.
- 334 [5] Moritz D Lürig, Seth Donoughe, Erik I Svensson, Arthur Porto, and Masahito Tsuboi. Computer
335 vision, machine learning, and the promise of phenomics in ecology and evolutionary biology.
336 *Frontiers in Ecology and Evolution*, 9:642774, 2021.
- 337 [6] Grant Van Horn, Oisín Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig
338 Adam, Pietro Perona, and Serge Belongie. The inaturalist species classification and detection
339 dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
340 pages 8769–8778, 2018.
- 341 [7] Randal A Singer, Kevin J Love, and Lawrence M Page. A survey of digitized data from us fish
342 collections in the idigbio data aggregator. *PloS one*, 13(12):e0207636, 2018.
- 343 [8] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The
344 caltech-ucsd birds-200-2011 dataset. 2011.
- 345 [9] Mohannad Elhamod, Mridul Khurana, Harish Babu Manogaran, Josef C Uyeda, Meghan A
346 Balk, Wasila Dahdul, Yasin Bakis, Henry L Bart Jr, Paula M Mabee, Hilmar Lapp, et al.
347 Discovering novel biological traits from images using phylogeny-guided neural networks. In
348 *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*,
349 pages 3966–3978, 2023.
- 350 [10] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su.
351 This looks like that: deep learning for interpretable image recognition. *Advances in neural*
352 *information processing systems*, 32, 2019.
- 353 [11] Dawid Rymarczyk, Łukasz Struski, Jacek Tabor, and Bartosz Zieliński. Protopshare: Prototypi-
354 cal parts sharing for similarity discovery in interpretable image classification. In *Proceedings of*
355 *the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 1420–1430,
356 2021.
- 357 [12] Dawid Rymarczyk, Łukasz Struski, Michał Górszczak, Koryna Lewandowska, Jacek Tabor, and
358 Bartosz Zieliński. Interpretable image classification with differentiable prototypes assignment.
359 In *European Conference on Computer Vision*, pages 351–368. Springer, 2022.
- 360 [13] Meike Nauta, Ron Van Bree, and Christin Seifert. Neural prototype trees for interpretable
361 fine-grained image recognition. In *Proceedings of the IEEE/CVF Conference on Computer*
362 *Vision and Pattern Recognition*, pages 14933–14943, 2021.
- 363 [14] Luke J Harmon, Jonathan B Losos, T Jonathan Davies, Rosemary G Gillespie, John L Gittleman,
364 W Bryan Jennings, Kenneth H Kozak, Mark A McPeck, Franck Moreno-Roark, Thomas J Near,
365 et al. Early bursts of body size and shape evolution are rare in comparative data. *Evolution*,
366 64(8):2385–2396, 2010.
- 367 [15] Matthew W Pennell, Richard G FitzJohn, William K Cornwell, and Luke J Harmon. Model
368 adequacy and the macroevolution of angiosperm functional traits. *The American Naturalist*,
369 186(2):E33–E50, 2015.
- 370 [16] Christopher Lawrence and Elizabeth G. Campolongo. Heliconius collection (cambridge butter-
371 fly), 2024.

- 372 [17] Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin. Interpretable image recognition with
373 hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and*
374 *Crowdsourcing*, volume 7, pages 32–40, 2019.
- 375 [18] Meike Nauta, Jörg Schlötterer, Maurice van Keulen, and Christin Seifert. Pip-net: Patch-based
376 intuitive prototypes for interpretable image classification. In *Proceedings of the IEEE/CVF*
377 *Conference on Computer Vision and Pattern Recognition*, pages 2744–2753, 2023.
- 378 [19] Adrian Hoffmann, Claudio Fanconi, Rahul Rade, and Jonas Kohler. This looks like that... does
379 it? shortcomings of latent space prototype interpretability in deep networks. *arXiv preprint*
380 *arXiv:2105.02968*, 2021.
- 381 [20] Sunnie SY Kim, Nicole Meister, Vikram V Ramaswamy, Ruth Fong, and Olga Russakovsky.
382 Hive: Evaluating the human interpretability of visual explanations. In *European Conference on*
383 *Computer Vision*, pages 280–298. Springer, 2022.
- 384 [21] Tongzhou Wang and Phillip Isola. Understanding contrastive representation learning through
385 alignment and uniformity on the hypersphere. In *International Conference on Machine Learning*,
386 pages 9929–9939. PMLR, 2020.
- 387 [22] Thalles Silva and Adín Ramírez Rivera. Representation learning via consistent assignment of
388 views to clusters. In *Proceedings of the 37th ACM/SIGAPP Symposium on Applied Computing*,
389 pages 987–994, 2022.
- 390 [23] Jiaqi Wang, Huafeng Liu, Xinyue Wang, and Liping Jing. Interpretable image recognition
391 by constructing transparent embedding space. In *Proceedings of the IEEE/CVF International*
392 *Conference on Computer Vision*, pages 895–904, 2021.
- 393 [24] Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X Yu. Orthogonal convolutional
394 neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern*
395 *recognition*, pages 11505–11515, 2020.
- 396 [25] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax.
397 *arXiv preprint arXiv:1611.01144*, 2016.
- 398 [26] W. Jetz, G. H. Thomas, J. B. Joy, K. Hartmann, and A. O. Mooers. The global diversity of birds
399 in space and time. *Nature*, 491:444–448, 2012.
- 400 [27] Gabriela Montejo-Kovacevich, Eva van der Heijden, Nicola Nadeau, and Chris Jiggins. Cam-
401 bridge butterfly wing collection batch 10, November 2020.
- 402 [28] Patricio A. Salazar, Nicola Nadeau, Gabriela Montejo-Kovacevich, and Chris Jiggins. Sheffield
403 butterfly wing collection - Patricio Salazar, Nicola Nadeau, Ikiam broods batch 1 and 2, Novem-
404 ber 2020.
- 405 [29] Gabriela Montejo-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing
406 collection batch 2, May 2019.
- 407 [30] Chris Jiggins, Gabriela Montejo-Kovacevich, Ian Warren, and Eva Wiltshire. Cambridge
408 butterfly wing collection batch 3, May 2019.
- 409 [31] Gabriela Montejo-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing
410 collection batch 4, May 2019.
- 411 [32] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge
412 butterfly wing collection batch 5, May 2019.
- 413 [33] Ian Warren and Chris Jiggins. Miscellaneous Heliconius wing photographs (2001-2019) Part 1,
414 February 2019.
- 415 [34] Ian Warren and Chris Jiggins. Miscellaneous Heliconius wing photographs (2001-2019) Part 3,
416 February 2019.
- 417 [35] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge
418 butterfly wing collection batch 6, May 2019.

- 419 [36] Chris Jiggins and Ian Warren. Cambridge butterfly wing collection - Chris Jiggins 2001/2
420 broods batch 1, January 2019.
- 421 [37] Chris Jiggins and Ian Warren. Cambridge butterfly wing collection - Chris Jiggins 2001/2
422 broods batch 2, January 2019.
- 423 [38] Joana I. Meier, Patricio Salazar, Gabriela Montejo-Kovacevich, Ian Warren, and Chris Jiggins.
424 Cambridge butterfly wing collection - Patricio Salazar PhD wild specimens batch 3, October
425 2020.
- 426 [39] Gabriela Montejo-Kovacevich, Chris Jiggins, and Ian Warren. Cambridge butterfly wing
427 collection batch 1- version 2, May 2019.
- 428 [40] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, Camilo Salazar, Marianne Elias,
429 Imogen Gavins, Eva Wiltshire, Stephen Montgomery, and Owen McMillan. Cambridge and
430 collaborators butterfly wing collection batch 10, May 2019.
- 431 [41] Patricio Salazar, Gabriela Montejo-Kovacevich, Ian Warren, and Chris Jiggins. Cambridge
432 butterfly wing collection - Patricio Salazar PhD wild and bred specimens batch 1, December
433 2018.
- 434 [42] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge
435 butterfly wing collection batch 7, May 2019.
- 436 [43] Patricio Salazar, Gabriela Montejo-Kovacevich, Ian Warren, and Chris Jiggins. Cambridge
437 butterfly wing collection - Patricio Salazar PhD wild and bred specimens batch 2, January 2019.
- 438 [44] Erika Pinheiro de Castro, Christopher Jiggins, Karina Lucas da Silva-Brandão, Andre Victor
439 Lucci Freitas, Marcio Zikan Cardoso, Eva Van Der Heijden, Joana Meier, and Ian Warren.
440 Brazilian Butterflies Collected December 2020 to January 2021, February 2022.
- 441 [45] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, and Eva Wiltshire. Cambridge
442 butterfly wing collection batch 8, May 2019.
- 443 [46] Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, Eva Wiltshire, and Imogen Gavins.
444 Cambridge butterfly wing collection batch 9, May 2019.
- 445 [47] Gabriela Montejo-Kovacevich, Eva van der Heijden, and Chris Jiggins. Cambridge butterfly
446 collection - GMK Broods Ikiam 2018, November 2020.
- 447 [48] Gabriela Montejo-Kovacevich, Quentin Paynter, and Amin Ghane. *Heliconius erato cyrba*,
448 Cook Islands (New Zealand) 2016, 2019, 2021, September 2021.
- 449 [49] Ian Warren and Chris Jiggins. Miscellaneous *Heliconius* wing photographs (2001-2019) Part 2,
450 February 2019.
- 451 [50] Camilo Salazar, Gabriela Montejo-Kovacevich, Chris Jiggins, Ian Warren, and Imogen Gavins.
452 Camilo Salazar and Cambridge butterfly wing collection batch 1, May 2019.
- 453 [51] Anniina Mattila, Chris Jiggins, and Ian Warren. University of Helsinki butterfly collection -
454 Anniina Mattila bred specimens, February 2019.
- 455 [52] OpenTreeOfLife, Benjamin Redelings, Luna Luisa Sanchez Reyes, Karen A. Cranston, Jim
456 Allman, Mark T. Holder, and Emily Jane McTavish. Open tree of life synthetic tree, 2019.
- 457 [53] Francois Michonneau, Joseph W. Brown, and David J. Winter. *rotl*: an R package to interact
458 with the open tree of life data. *Methods in Ecology and Evolution*, 7(12):1476–1481, 2016.
- 459 [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
460 recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
461 pages 770–778, 2016.
- 462 [55] Dipanjyoti Paul, Arpita Chowdhury, Xinqi Xiong, Feng-Ju Chang, David Carlyn, Samuel
463 Stevens, Kaiya Provost, Anuj Karpatne, Bryan Carstens, Daniel Rubenstein, et al. A simple
464 interpretable transformer for fine-grained image classification and analysis. *arXiv preprint*
465 *arXiv:2311.04157*, 2023.

- 466 [56] Abien Fred Agarap. Deep learning using rectified linear units (relu). *arXiv preprint*
467 *arXiv:1803.08375*, 2018.
- 468 [57] Samuel G Müller and Frank Hutter. Trivialaugment: Tuning-free yet state-of-the-art data
469 augmentation. In *Proceedings of the IEEE/CVF international conference on computer vision*,
470 pages 774–782, 2021.
- 471 [58] R. Farrell. Cub-200-2011 segmentations (1.0) [data set], 2024.

472 **NeurIPS Paper Checklist**

473 **1. Claims**

474 Question: Do the main claims made in the abstract and introduction accurately reflect the
475 paper's contributions and scope?

476 Answer: [\[Yes\]](#)

477 Justification: We claim that HComP-Net can generate novel hypotheses for potential evo-
478 lutionary traits (shared traits among species due to common ancestry in the phylogeny)
479 from image by learning prototypes at each internal node in the phylogenetic tree. We show
480 through various visualizations of the prototypes in Figures 5, 6, and 7 to 16, that the learned
481 prototypes at the internal nodes can identify possible evolutionary traits from images. We
482 also evaluate the improved interpretability of our approach quantitatively in Table 3 by
483 computing part purity metric on Bird dataset.

484 Guidelines:

- 485 • The answer NA means that the abstract and introduction do not include the claims
486 made in the paper.
- 487 • The abstract and/or introduction should clearly state the claims made, including the
488 contributions made in the paper and important assumptions and limitations. A No or
489 NA answer to this question will not be perceived well by the reviewers.
- 490 • The claims made should match theoretical and experimental results, and reflect how
491 much the results can be expected to generalize to other settings.
- 492 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
493 are not attained by the paper.

494 **2. Limitations**

495 Question: Does the paper discuss the limitations of the work performed by the authors?

496 Answer: [\[Yes\]](#)

497 Justification: We discuss the limitations of our work in Supplementary Section I.

498 Guidelines:

- 499 • The answer NA means that the paper has no limitation while the answer No means that
500 the paper has limitations, but those are not discussed in the paper.
- 501 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 502 • The paper should point out any strong assumptions and how robust the results are to
503 violations of these assumptions (e.g., independence assumptions, noiseless settings,
504 model well-specification, asymptotic approximations only holding locally). The authors
505 should reflect on how these assumptions might be violated in practice and what the
506 implications would be.
- 507 • The authors should reflect on the scope of the claims made, e.g., if the approach was
508 only tested on a few datasets or with a few runs. In general, empirical results often
509 depend on implicit assumptions, which should be articulated.
- 510 • The authors should reflect on the factors that influence the performance of the approach.
511 For example, a facial recognition algorithm may perform poorly when image resolution
512 is low or images are taken in low lighting. Or a speech-to-text system might not be
513 used reliably to provide closed captions for online lectures because it fails to handle
514 technical jargon.
- 515 • The authors should discuss the computational efficiency of the proposed algorithms
516 and how they scale with dataset size.
- 517 • If applicable, the authors should discuss possible limitations of their approach to
518 address problems of privacy and fairness.
- 519 • While the authors might fear that complete honesty about limitations might be used by
520 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
521 limitations that aren't acknowledged in the paper. The authors should use their best
522 judgment and recognize that individual actions in favor of transparency play an impor-
523 tant role in developing norms that preserve the integrity of the community. Reviewers
524 will be specifically instructed to not penalize honesty concerning limitations.

525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [NA]

Justification: The assumptions made in our work do not require explicit theoretical proofs. Instead, for the key loss terms that we introduce in this work, we provide ablation results in Supplementary Table 6 to show empirically the importance of each component.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We provide details of hyperparameters in Supplementary Section E. Furthermore, we also provide the full code, data, and necessary data preprocessing pipelines to reproduce all the experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).

578 (d) We recognize that reproducibility may be tricky in some cases, in which case
579 authors are welcome to describe the particular way they provide for reproducibility.
580 In the case of closed-source models, it may be that access to the model is limited in
581 some way (e.g., to registered users), but it should be possible for other researchers
582 to have some path to reproducing or verifying the results.

583 5. Open access to data and code

584 Question: Does the paper provide open access to the data and code, with sufficient instruc-
585 tions to faithfully reproduce the main experimental results, as described in supplemental
586 material?

587 Answer: [Yes]

588 Justification: We provide the full code, data, and necessary data preprocessing pipelines to
589 reproduce the experiments.

590 Guidelines:

- 591 • The answer NA means that paper does not include experiments requiring code.
- 592 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/
593 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 594 • While we encourage the release of code and data, we understand that this might not be
595 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not
596 including code, unless this is central to the contribution (e.g., for a new open-source
597 benchmark).
- 598 • The instructions should contain the exact command and environment needed to run to
599 reproduce the results. See the NeurIPS code and data submission guidelines ([https://
600 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 601 • The authors should provide instructions on data access and preparation, including how
602 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 603 • The authors should provide scripts to reproduce all experimental results for the new
604 proposed method and baselines. If only a subset of experiments are reproducible, they
605 should state which ones are omitted from the script and why.
- 606 • At submission time, to preserve anonymity, the authors should release anonymized
607 versions (if applicable).
- 608 • Providing as much information as possible in supplemental material (appended to the
609 paper) is recommended, but including URLs to data and code is permitted.

610 6. Experimental Setting/Details

611 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-
612 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the
613 results?

614 Answer: [Yes]

615 Justification: The details of the data splits for each dataset used are provided in Table 8 and
616 overview of phylogeny is given in Table 7. We also give details of key hyperparameters
617 and the way they were chosen in Supplementary Section E. Full code also provided for
618 reproducibility.

619 Guidelines:

- 620 • The answer NA means that the paper does not include experiments.
- 621 • The experimental setting should be presented in the core of the paper to a level of detail
622 that is necessary to appreciate the results and make sense of them.
- 623 • The full details can be provided either with the code, in appendix, or as supplemental
624 material.

625 7. Experiment Statistical Significance

626 Question: Does the paper report error bars suitably and correctly defined or other appropriate
627 information about the statistical significance of the experiments?

628 Answer: [Yes]

629 Justification: We have done multiple runs of our model on Bird dataset with different
630 random weight initialization, and report the mean and standard deviation of accuracy in
631 Supplementary Section D

632 Guidelines:

- 633 • The answer NA means that the paper does not include experiments.
- 634 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
635 dence intervals, or statistical significance tests, at least for the experiments that support
636 the main claims of the paper.
- 637 • The factors of variability that the error bars are capturing should be clearly stated (for
638 example, train/test split, initialization, random drawing of some parameter, or overall
639 run with given experimental conditions).
- 640 • The method for calculating the error bars should be explained (closed form formula,
641 call to a library function, bootstrap, etc.)
- 642 • The assumptions made should be given (e.g., Normally distributed errors).
- 643 • It should be clear whether the error bar is the standard deviation or the standard error
644 of the mean.
- 645 • It is OK to report 1-sigma error bars, but one should state it. The authors should
646 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
647 of Normality of errors is not verified.
- 648 • For asymmetric distributions, the authors should be careful not to show in tables or
649 figures symmetric error bars that would yield results that are out of range (e.g. negative
650 error rates).
- 651 • If error bars are reported in tables or plots, The authors should explain in the text how
652 they were calculated and reference the corresponding figures or tables in the text.

653 8. Experiments Compute Resources

654 Question: For each experiment, does the paper provide sufficient information on the com-
655 puter resources (type of compute workers, memory, time of execution) needed to reproduce
656 the experiments?

657 Answer: [Yes]

658 Justification: Details of computer resources used are provided in Supplementary Section E

659 Guidelines:

- 660 • The answer NA means that the paper does not include experiments.
- 661 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
662 or cloud provider, including relevant memory and storage.
- 663 • The paper should provide the amount of compute required for each of the individual
664 experimental runs as well as estimate the total compute.
- 665 • The paper should disclose whether the full research project required more compute
666 than the experiments reported in the paper (e.g., preliminary or failed experiments that
667 didn't make it into the paper).

668 9. Code Of Ethics

669 Question: Does the research conducted in the paper conform, in every respect, with the
670 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

671 Answer: [Yes]

672 Justification: The work abides by NeurIPS Code of Ethics in every aspect.

673 Guidelines:

- 674 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 675 • If the authors answer No, they should explain the special circumstances that require a
676 deviation from the Code of Ethics.
- 677 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
678 eration due to laws or regulations in their jurisdiction).

679 10. Broader Impacts

680 Question: Does the paper discuss both potential positive societal impacts and negative
681 societal impacts of the work performed?

682 Answer: [NA]

683 Justification: There is no negative societal impact of the work performed to the best of our
684 knowledge.

685 Guidelines:

- 686 • The answer NA means that there is no societal impact of the work performed.
- 687 • If the authors answer NA or No, they should explain why their work has no societal
688 impact or why the paper does not address societal impact.
- 689 • Examples of negative societal impacts include potential malicious or unintended uses
690 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
691 (e.g., deployment of technologies that could make decisions that unfairly impact specific
692 groups), privacy considerations, and security considerations.
- 693 • The conference expects that many papers will be foundational research and not tied
694 to particular applications, let alone deployments. However, if there is a direct path to
695 any negative applications, the authors should point it out. For example, it is legitimate
696 to point out that an improvement in the quality of generative models could be used to
697 generate deepfakes for disinformation. On the other hand, it is not needed to point out
698 that a generic algorithm for optimizing neural networks could enable people to train
699 models that generate Deepfakes faster.
- 700 • The authors should consider possible harms that could arise when the technology is
701 being used as intended and functioning correctly, harms that could arise when the
702 technology is being used as intended but gives incorrect results, and harms following
703 from (intentional or unintentional) misuse of the technology.
- 704 • If there are negative societal impacts, the authors could also discuss possible mitigation
705 strategies (e.g., gated release of models, providing defenses in addition to attacks,
706 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
707 feedback over time, improving the efficiency and accessibility of ML).

708 11. Safeguards

709 Question: Does the paper describe safeguards that have been put in place for responsible
710 release of data or models that have a high risk for misuse (e.g., pretrained language models,
711 image generators, or scraped datasets)?

712 Answer: [NA]

713 Justification: The work poses no risk of misuse.

714 Guidelines:

- 715 • The answer NA means that the paper poses no such risks.
- 716 • Released models that have a high risk for misuse or dual-use should be released with
717 necessary safeguards to allow for controlled use of the model, for example by requiring
718 that users adhere to usage guidelines or restrictions to access the model or implementing
719 safety filters.
- 720 • Datasets that have been scraped from the Internet could pose safety risks. The authors
721 should describe how they avoided releasing unsafe images.
- 722 • We recognize that providing effective safeguards is challenging, and many papers do
723 not require this, but we encourage authors to take this into account and make a best
724 faith effort.

725 12. Licenses for existing assets

726 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
727 the paper, properly credited and are the license and terms of use explicitly mentioned and
728 properly respected?

729 Answer: [Yes]

730 Justification: All datasets used have been cited along with their source wherever necessary.
731 We also mention the license details for each dataset used in Supplementary Section E

732 Guidelines:

- 733 • The answer NA means that the paper does not use existing assets.
- 734 • The authors should cite the original paper that produced the code package or dataset.
- 735 • The authors should state which version of the asset is used and, if possible, include a
- 736 URL.
- 737 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 738 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 739 service of that source should be provided.
- 740 • If assets are released, the license, copyright information, and terms of use in the
- 741 package should be provided. For popular datasets, paperswithcode.com/datasets
- 742 has curated licenses for some datasets. Their licensing guide can help determine the
- 743 license of a dataset.
- 744 • For existing datasets that are re-packaged, both the original license and the license of
- 745 the derived asset (if it has changed) should be provided.
- 746 • If this information is not available online, the authors are encouraged to reach out to
- 747 the asset's creators.

748 13. New Assets

749 Question: Are new assets introduced in the paper well documented and is the documentation
750 provided alongside the assets?

751 Answer: [Yes]

752 Justification: We provide the full code and the necessary documentation to reproduce the
753 results.

754 Guidelines:

- 755 • The answer NA means that the paper does not release new assets.
- 756 • Researchers should communicate the details of the dataset/code/model as part of their
- 757 submissions via structured templates. This includes details about training, license,
- 758 limitations, etc.
- 759 • The paper should discuss whether and how consent was obtained from people whose
- 760 asset is used.
- 761 • At submission time, remember to anonymize your assets (if applicable). You can either
- 762 create an anonymized URL or include an anonymized zip file.

763 14. Crowdsourcing and Research with Human Subjects

764 Question: For crowdsourcing experiments and research with human subjects, does the paper
765 include the full text of instructions given to participants and screenshots, if applicable, as
766 well as details about compensation (if any)?

767 Answer: [NA]

768 Justification: No crowdsourcing or human subject involved in this research.

769 Guidelines:

- 770 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 771 human subjects.
- 772 • Including this information in the supplemental material is fine, but if the main contribu-
- 773 tion of the paper involves human subjects, then as much detail as possible should be
- 774 included in the main paper.
- 775 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 776 or other labor should be paid at least the minimum wage in the country of the data
- 777 collector.

778 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human 779 Subjects

780 Question: Does the paper describe potential risks incurred by study participants, whether
781 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
782 approvals (or an equivalent approval/review based on the requirements of your country or
783 institution) were obtained?

784 Answer: [NA]

785
786
787
788
789
790
791
792
793
794
795
796

Justification: No crowdsourcing or human subject involved in this research.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

797 **A Ablation of Over-specificity Loss Trade-off Hyperparameter**

798 We have provided an ablation for the over-specificity loss trade-off hyperparameter (λ_{ovsp}) in Table 4.
 799 We can observe that increasing the weight of over-specificity loss reduces the model’s classification
 800 performance, as the model struggles to find any commonality especially at internal nodes where the
 801 number of leaf descendant species are large in number and quite diverse. It is natural that species that
 802 are diverse and distantly related may share fewer characteristics with each other, in comparison to a
 803 set of species that diverged more recently from a common ancestor [14, 15]. Therefore, forcing the
 804 model to learn common traits with a strong \mathcal{L}_{ovsp} constraint can cause the model to perform bad in
 805 terms of accuracy.

Table 4: Ablation of over-specificity loss trade-off hyperparameter (λ_{ovsp}). Done on Bird dataset.

| λ_{ovsp} | Part purity | Part purity with mask applied | % masked | % Accuracy |
|--------------------------|--------------------|-------------------------------|----------|--------------|
| w/o \mathcal{L}_{ovsp} | 0.68 ± 0.22 | 0.75 ± 0.17 | 21.42% | 58.32 |
| 0.05 | 0.72 ± 0.19 | 0.77 ± 0.16 | 16.53% | 70.01 |
| 0.1 | 0.71 ± 0.18 | 0.74 ± 0.16 | 11.31% | 70.97 |
| 0.5 | 0.71 ± 0.19 | 0.72 ± 0.18 | 4.2% | 68.23 |
| 1.0 | 0.70 ± 0.19 | 0.70 ± 0.2 | 2.13% | 62.68 |
| 2.0 | 0.69 ± 0.19 | 0.69 ± 0.19 | 0.55% | 53.16 |

806 **B Ablation of Number of Prototypes**

807 In Table 5 we vary the number of prototypes per child β for a node to see the impact on model’s
 808 performance. We note that while the accuracy increases marginally with increasing the number of
 809 prototypes per child (β) from 10 to 15, it also considerably increases the overall number of prototypes
 810 initialized. Therefore we continue to work with $\beta = 10$ for all of our experiments.

Table 5: Ablation of number of prototypes per child for a node (β). Done on Bird dataset.

| Number of Prototypes (β) | % Accuracy |
|----------------------------------|--------------|
| 10 | 70.01 |
| 15 | 70.92 |
| 20 | 67.93 |

811 **C Ablation of Individual Losses**

812 In Table 6, we perform an ablation of the various loss terms used in our methodology. As it can be
 813 observed, removal of \mathcal{L}_{ovsp} and \mathcal{L}_{disc} degrades performance in terms of both semantic consistency
 814 (part purity) and accuracy. On the other hand, removal of self supervised contrastive loss \mathcal{L}_{SS}
 815 improves accuracy but at the cost of drastically decreasing the semantic consistency.

Table 6: Ablation of individual losses. Done on Bird dataset.

| Model | Part purity | Part purity with mask applied | % masked | % Accuracy |
|------------------------------------|-------------|-------------------------------|----------|------------|
| HComP-Net | 0.72 ± 0.19 | 0.77 ± 0.16 | 16.53% | 70.01 |
| HComP-Net w/o \mathcal{L}_{ovsp} | 0.68 ± 0.22 | 0.75 ± 0.17 | 21.42% | 58.32 |
| HComP-Net w/o \mathcal{L}_{disc} | 0.69 ± 0.19 | 0.72 ± 0.17 | 10.95% | 65.99 |
| HComP-Net w/o \mathcal{L}_{SS} | 0.53 ± 0.18 | 0.57 ± 0.15 | 8.36% | 81.62 |

816 **D Consistency of Classification Performance Over Multiple Runs**

817 We trained the model using five distinct random weight initializations. The results showed that the
 818 model’s fine-grained accuracy averaged 70.63% with a standard deviation of 0.18%.

819 **E Implementation Details**

820 We have included all the source code and dataset along with the comprehensive instructions to
 821 reproduce the results, in the supplementary material (.zip file).

822 **Model hyper-parameters:** We build HComp-Net on top of a ConvNeXt-tiny architecture as the
 823 backbone feature extractor. We have modified the stride of the max pooling layers of later stages
 824 of the backbone from 2 to 1 similar to PIP-Net such that the backbone produces feature maps of
 825 increased height and width, in order to get more fine-grained prototype score maps. We implement
 826 and experiment our method on ConvNeXt-tiny backbones with 26×26 feature maps. The length
 827 of prototype vectors C is 768. The weights ϕ at every node n of HComp-Net are constrained to be
 828 non-negative by the use of ReLU activation function [56]. Further, the prototype activation nodes are
 829 connected with non-negative weights only to their respective child classes in W while their weights
 830 to other classes are made zero and non-trainable.

831 **Training details:** All models were trained with images resized and appropriately padded to 224×224
 832 pixel resolution and augmented using TrivialAugment [57] for contrastive learning. The prototypes
 833 are pretrained with self-supervised learning similar to PIP-Net for 10 epochs, following which the
 834 model is trained with the entire set of loss functions for 60 epochs. We use a batch size of 256 for
 835 Bird dataset and 64 for Butterfly and Fish dataset. The masking module is trained in parallel and its
 836 training is continued for 15 additional epochs after the training of rest of the model is completed. The
 837 trade-off hyper-parameters for the loss functions are set to be $\lambda_{CE} = 2$; $\lambda_A = 5$; $\lambda_T = 2$; $\lambda_{ovsp} =$
 838 0.05 ; $\lambda_{disc} = 0.1$; $\lambda_{orth} = 0.1$; $\lambda_{mask} = 2.0$; $\lambda_{L1} = 0.5$. λ_{CE} , λ_T and λ_A were borrowed from
 839 PIP-Net [18]. Ablations to arrive at suitable λ_{ovsp} is provided in Table 4. λ_{disc} and λ_{orth} were
 840 chosen empirically and found to work well on all three datasets. Experiment on unseen species was
 841 done by leaving out certain classes from the datasets, so that they are not considered during training.

842 **Dataset and Phylogeny Details:** Dataset statistics and phylogeny statistics are provided in Table
 843 8 and Table 7 respectively. Bird dataset is created by choosing 190 species from CUB-200-2011
 844 ² [8] dataset, which were part of the phylogeny. Background from all images were filtered using
 845 the associated segmentation metadata [58]. For Butterfly dataset we considered each subspecies
 846 as an individual class and considered only the subspecies of genus *Heliconius* from the *Heliconius*
 847 *Collection* (Cambridge Butterfly)³ [16]. There is substantial variation among subspecies of *Heliconius*
 848 species. Furthermore, we balanced the dataset by filtering out the subspecies which did not have
 849 20 or more images. We also sampled a subset of 100 images from each subspecies that had more
 850 than 100 images. For Fish⁴ dataset, we followed the exact same preprocessing steps as outlined in
 851 PhyloNN [9].

852 **Compute Resources:** The models for Bird dataset were trained on two NVIDIA A100 GPUs with
 853 80GB of RAM each. Butterfly and Fish models were trained on single A100 GPU. As a rough
 854 estimate the execution time for training model on Bird dataset is around 2.5 hours. For Butterfly and
 855 Fish datasets, the training completes under 1 hour. We used a single A100 GPU during inference
 856 stage for all other analysis.

Table 7: High level statistics of the phylogenies used for different datasets.

| Phylogeny | # Internal nodes | Max-depth | Min-depth |
|-----------|------------------|-----------|-----------|
| Bird | 184 | 25 | 3 |
| Butterfly | 13 | 5 | 2 |
| Fish | 20 | 11 | 2 |

²License: CC BY

³Note that this dataset is a compilation of images from 25 Zenodo records by the Butterfly Genetics Group at Cambridge University, licensed under Creative Commons Attribution 4.0 International ([27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51]).

⁴License: CC BY-NC

Table 8: Dataset statistics (# train and validation images).

| Dataset | # Classes | Train set | Validation set |
|------------------|-----------|-----------|----------------|
| Bird | 190 | 5695 | 5512 |
| Butterfly | 30 | 1418 | 358 |
| Fish | 38 | 4140 | 1294 |

Table 9: Part purity with post-hoc thresholding approach. Done on **Bird** dataset.

| Threshold | Part purity with mask applied | % masked |
|-----------|-------------------------------|----------|
| 0.2 | 0.74 ± 0.28 | 12.28% |
| 0.3 | 0.75 ± 0.27 | 13.47% |
| 0.4 | 0.76 ± 0.26 | 14.97% |
| 0.5 | 0.77 ± 0.15 | 16.66% |
| 0.6 | 0.77 ± 0.26 | 17.43% |

857 F Post-hoc Thresholding to Identify Over-specific Prototypes

858 An alternative approach to learning masking module is to calculate the over-specificity score for each
 859 prototype on the test set after training the model. We calculate the over-specificity scores for the
 860 prototypes of a trained model as follows,

$$\mathcal{O}_i = - \prod_{d=1}^{D_i} \frac{1}{\text{top}_k} \sum_{i=1}^{\text{top}_k} (g_i) \quad (9)$$

861 For a given prototype, we choose the top_k images with the highest prototype scores from each
 862 leaf descendant. After taking mean of the top_k prototype score, we multiply the values from each
 863 descendant to arrive at the over-specificity score for the particular prototype. Subsequently we choose
 864 a threshold to determine which prototypes are over-specific. We provide the results of post-hoc
 865 thresholding approach that can also be used to identify overspecific prototypes in Table 9. While we
 866 can note that this approach can also be effective, validating the threshold particularly in scenarios
 867 where there is no part annotations available (such as part location annotation of CUB-200-2011) can
 868 be an arduous task. In such cases directly identifying over-specific prototypes as part of the training
 869 through masking module can be the more feasible option.

870 G Additional Visualizations of the Hierarchical Prototypes Discovered by 871 HComp-Net

872 We provide more visualizations of the hierarchical prototypes discovered by HComp-Net for Butterfly
 873 (Figures 7 and 8) and Fish (Figure 9) datasets in this section. For ease of visualization, in each figure
 874 we visualize the prototypes learned over a small sub-tree from the phylogeny. The prototypes at the
 875 lowest level capture traits that are species-specific whereas the prototypes at internal nodes capture the
 876 commonality between its descendant species. For Fish dataset, we have provided textual descriptions
 877 purely based on human interpretation for the traits that are captured by prototypes at different levels.
 878 For Butterfly dataset, since the prototypes are capturing different wing patterns, assigning textual
 879 description for them is not straightforward. Therefore, we refrain from providing any text description
 880 for the highlighted regions of the learned prototypes and leave it to the reader’s interpretation.

881 H Additional Top-K Visualizations of HComp-Net Prototypes

882 We provide additional top-K visualizations of the prototypes from Butterfly (Figures 10 to 13) and
 883 Fish (Figures 14 to 16) datasets, where every row corresponds to a descendant species and the
 884 columns corresponds to the top-K images from the species with the largest prototype activation scores.
 885 A requirement of a semantically meaningful prototype is that it should consistently highlight the
 886 same part of the organisms in various images, provided that the part is visible. We can see in the

887 figures that the prototypes learned by HComP-Net consistently highlight the same part across all
888 top-K images of a species, and across all descendant species. We additionally show that HComP-Net
889 can find common traits at internal nodes with varying number of descendant species, including 4
890 species (Figure 10), 5 species (Figures 11 and 12), and 10 species (Figure 13) of butterflies, and
891 5 species (Figure 14), 8 species (Figure 15) and 18 species (Figure 16) for fish. We also provide
892 several top-k visualizations of prototypes learned for bird species in Figures 17 to 25. This shows the
893 ability of HComP-Net to discover common prototypes at internal nodes of the phylogenetic tree that
894 consistently highlight the same regions in the descendant species images even when the number of
895 descendants is large.

896 **I Limitations of Our Work**

897 A fundamental challenge of every prototype-based interpretability method (including ours) is the
898 difficulty in associating a semantic interpretation to the underlying visual concept of a prototype.
899 While some prototypes can be interpreted easily based on visual inspection of prototype activation
900 maps, other prototypes are harder to interpret and require additional domain expertise of biologists.
901 Also, while we have considered large phylogenies as that of the 190 species from CUB dataset, it may
902 still not be representative of all bird species. This limited scope may cause our method to identify
903 apparent homologous evolutionary traits that could differ with the inclusion of more species into the
904 phylogeny. Therefore, our method can be seen as a system that generates potential hypotheses about
905 evolutionary traits discovered in the form of hierarchical prototypes.

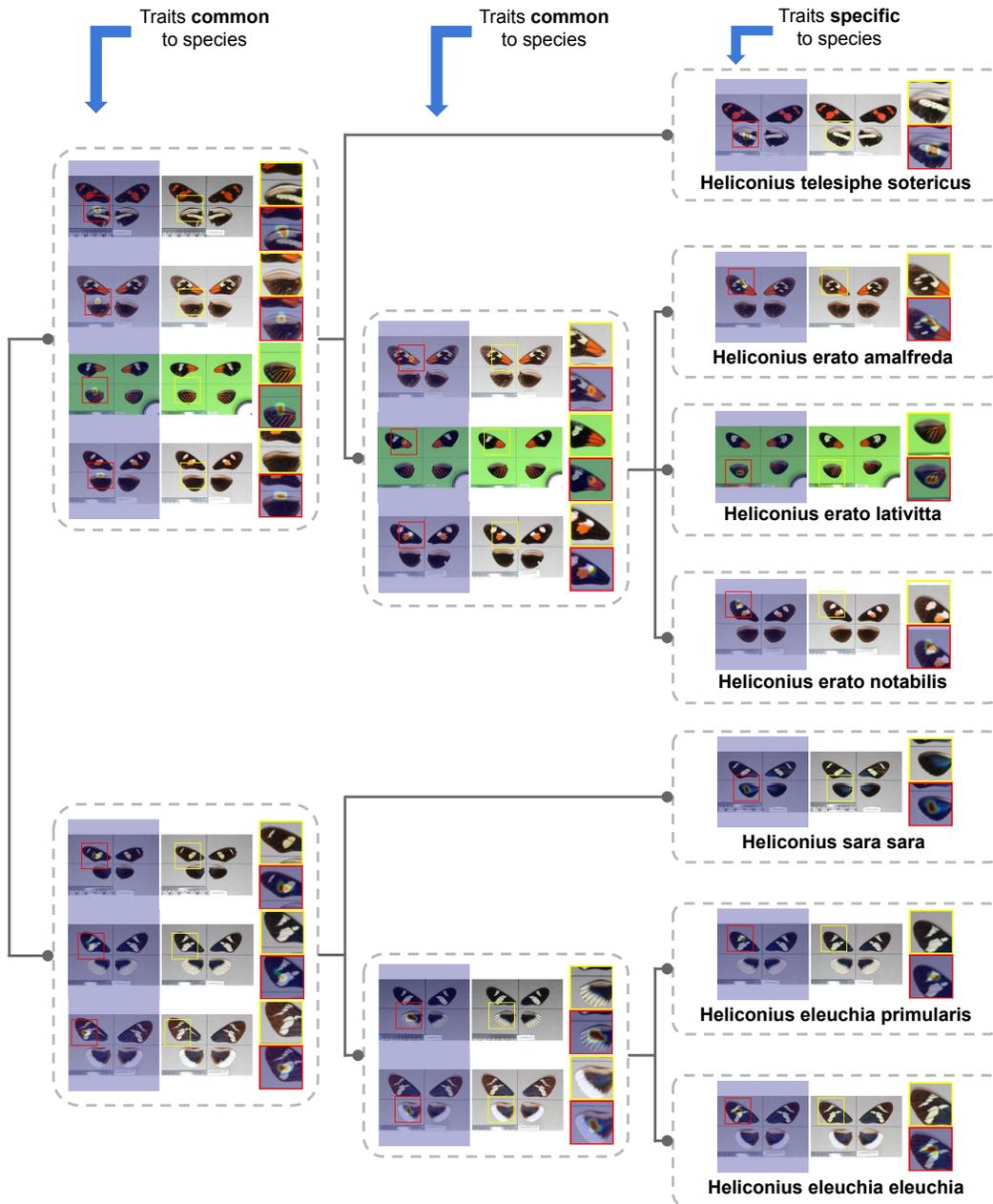


Figure 7: Visualizing the hierarchy of prototypes discovered by HComP-Net over three levels in the phylogeny of seven species from **Butterfly** dataset. For each prototype we visualize one image from each of its leaf descendant. Therefore, for prototypes at species level (rightmost column) we show only one image whereas for prototypes at internal nodes we show multiple images (equal to the number of leaf descendants). For each image, we show the zoomed in view of the original image as well as the heatmap overlaid image in the region of the learned prototype. The prototypes appear to be capturing different wing patterns of the butterflies.

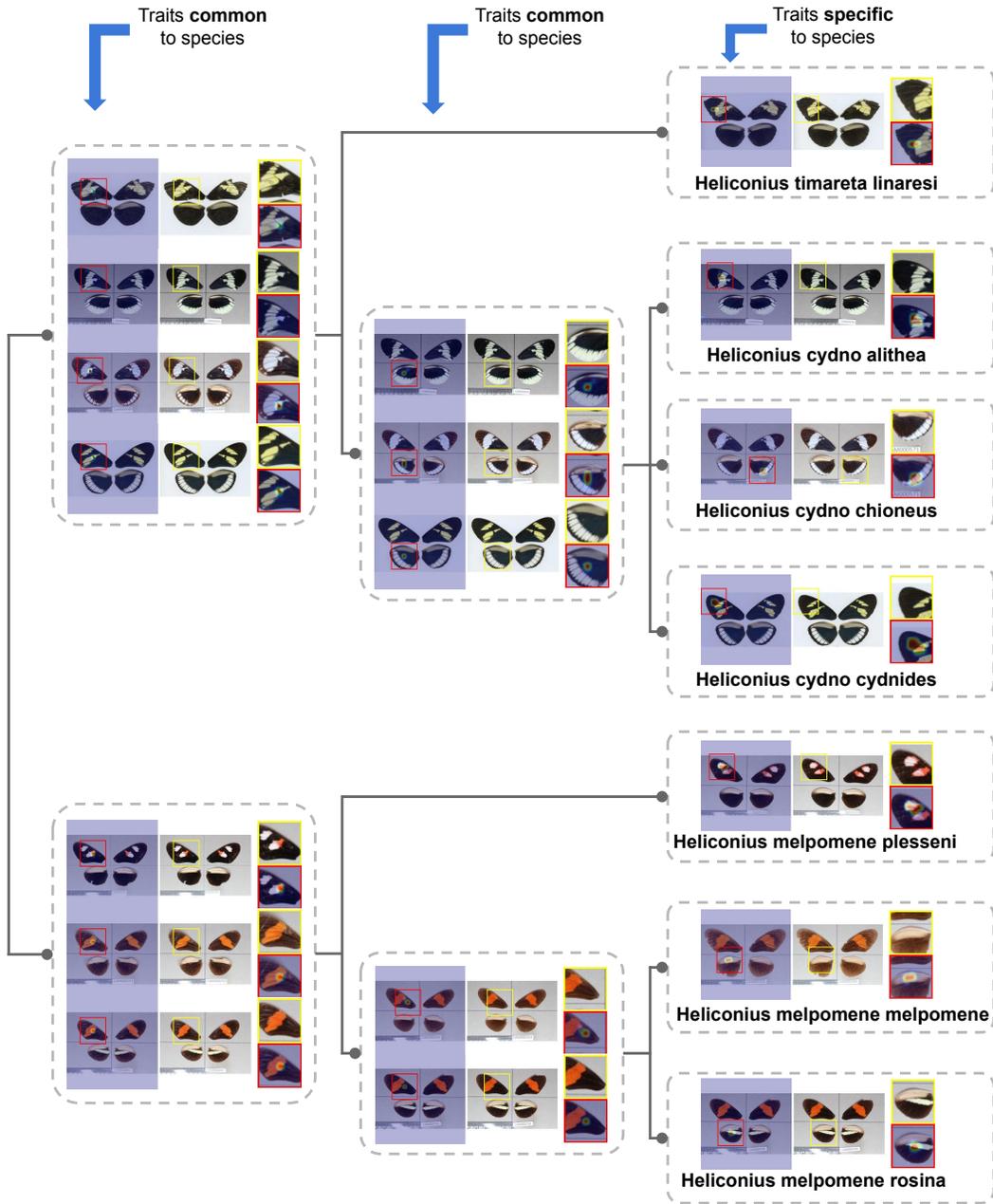


Figure 8: Visualizing the hierarchy of prototypes discovered by HComp-Net over three levels in the phylogeny of seven species from **Butterfly** dataset.

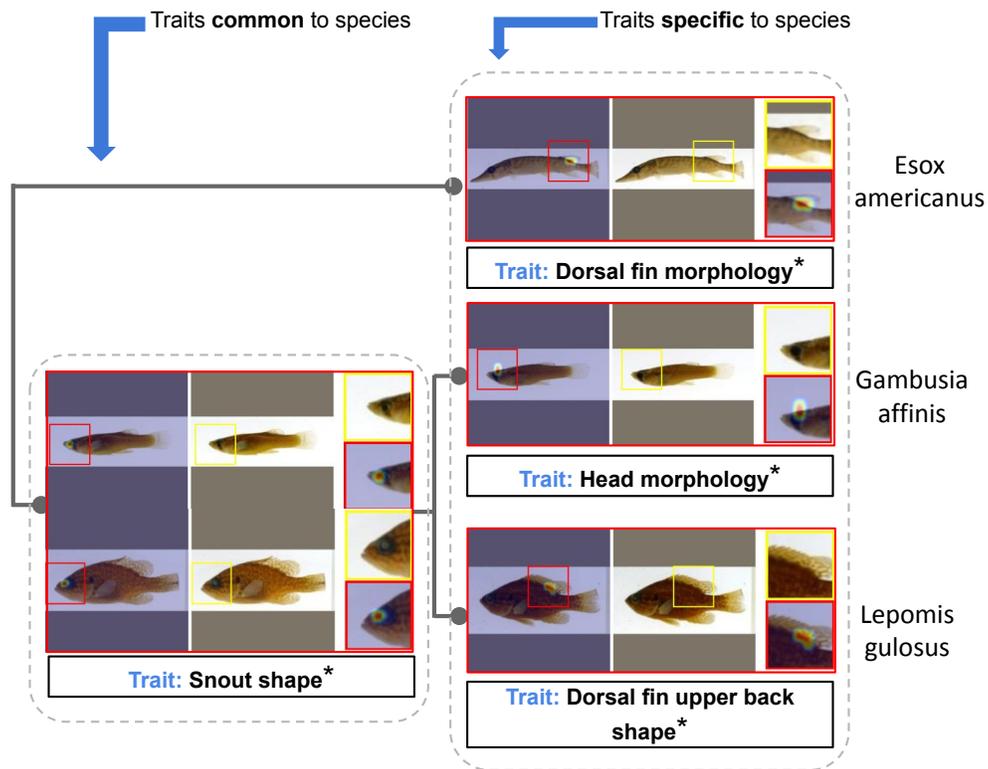


Figure 9: Visualizing the hierarchy of prototypes discovered by HComP-Net for a sub-trees with three species from **Fish** dataset. *Note that the textual descriptions of the hypothesized traits shown for every prototype are based on human interpretation.

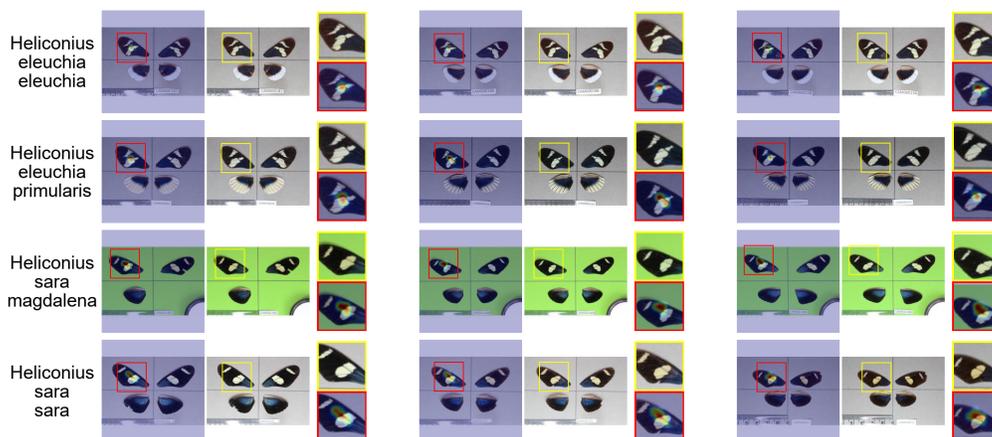


Figure 10: Top-K visualization of a prototype finding commonality between four species of butterfly sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.



Figure 11: Top-K visualization of a prototype finding commonality between nine species of butterfly sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

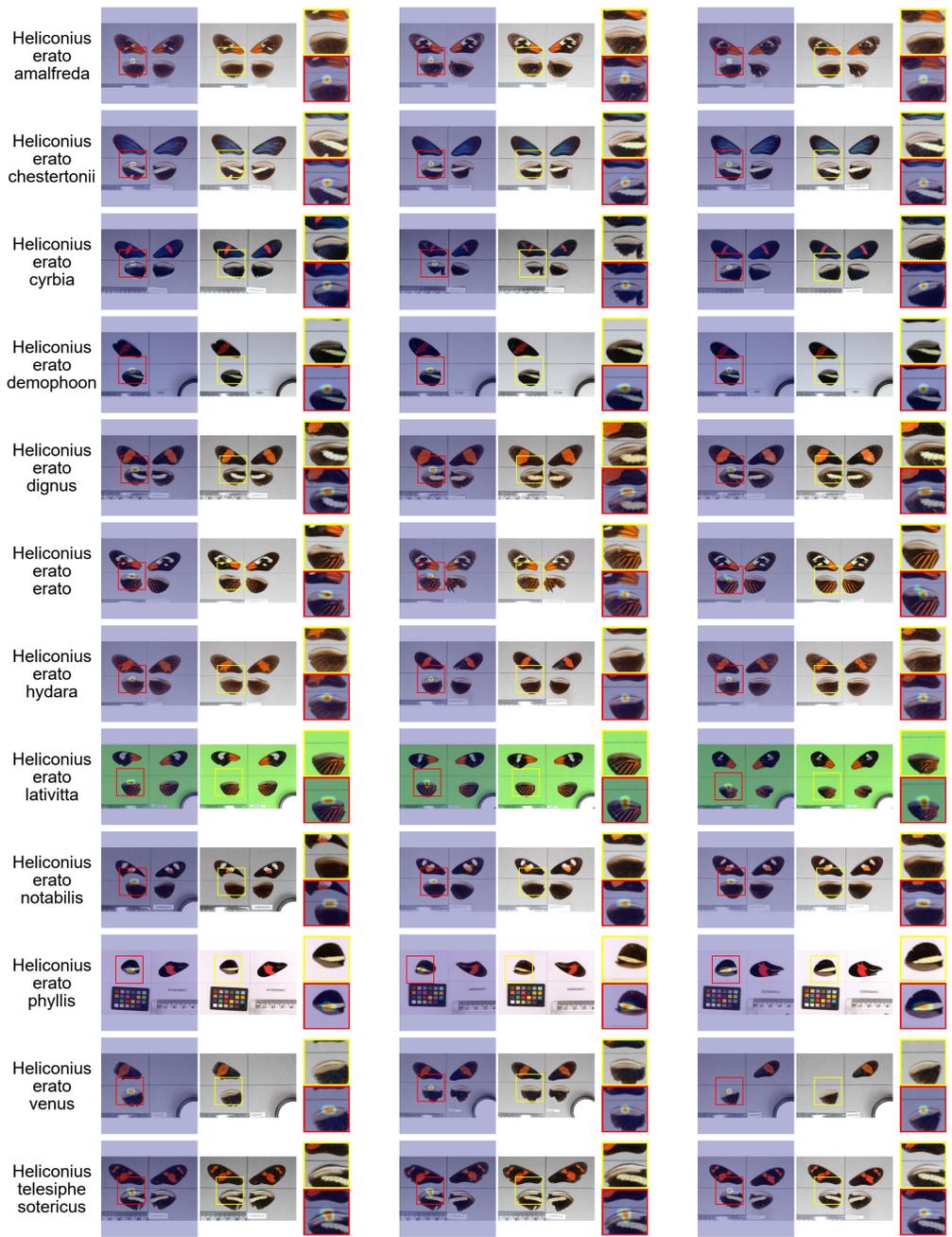


Figure 12: Top-K visualization of a prototype finding commonality between twelve species of butterfly sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.



Figure 13: Top-K visualization of a prototype finding commonality between four species of butterfly sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

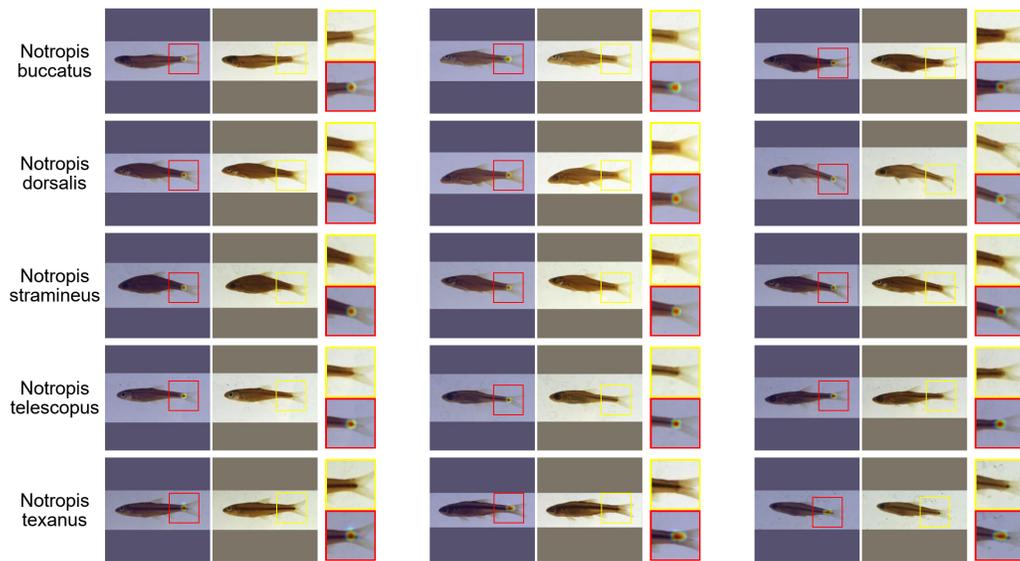


Figure 14: Top-K visualization of a prototype finding commonality between five species of fish sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

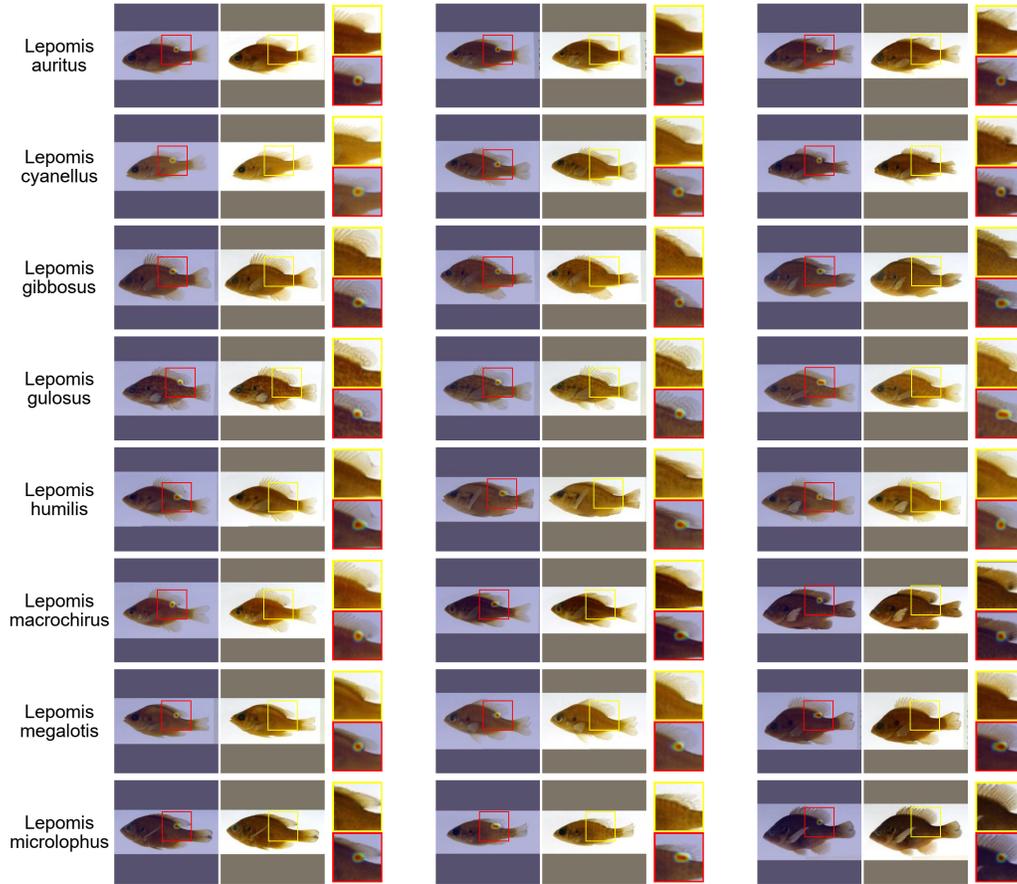


Figure 15: Top-K visualization of a prototype finding commonality between eight species of fish sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

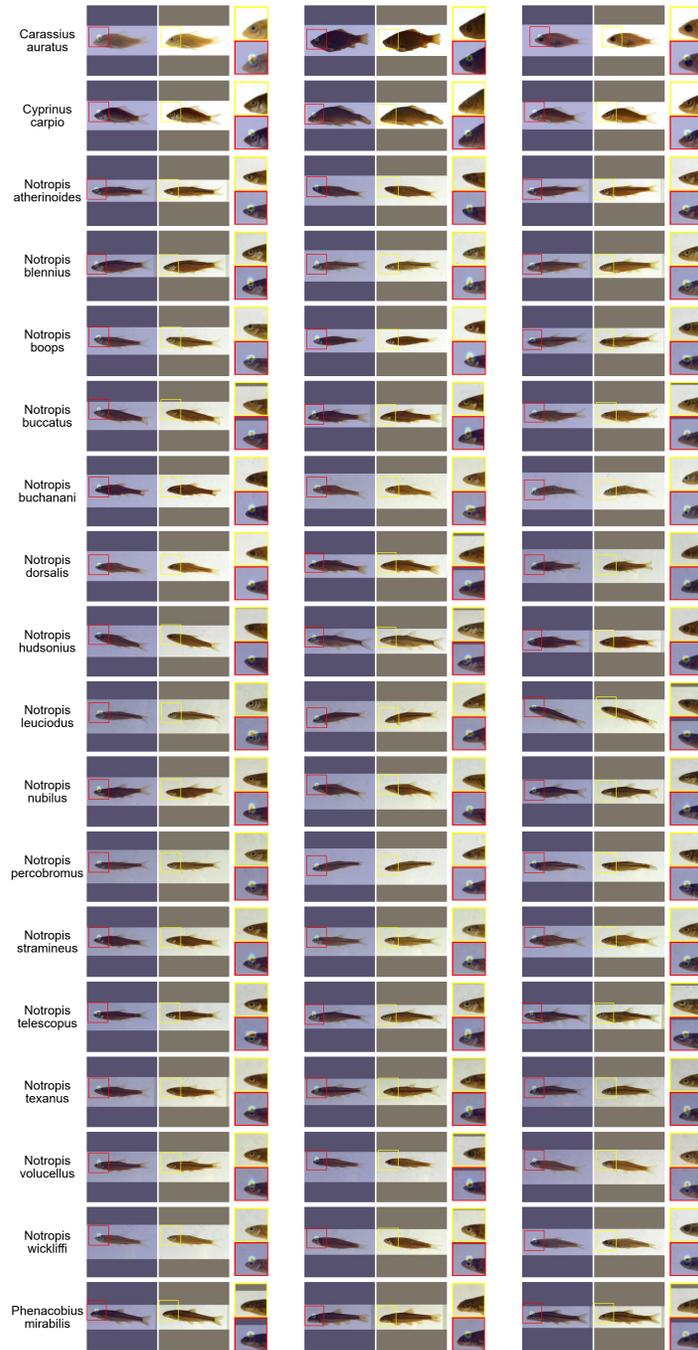


Figure 16: Top-K visualization of a prototype finding commonality between eighteen species of fish sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

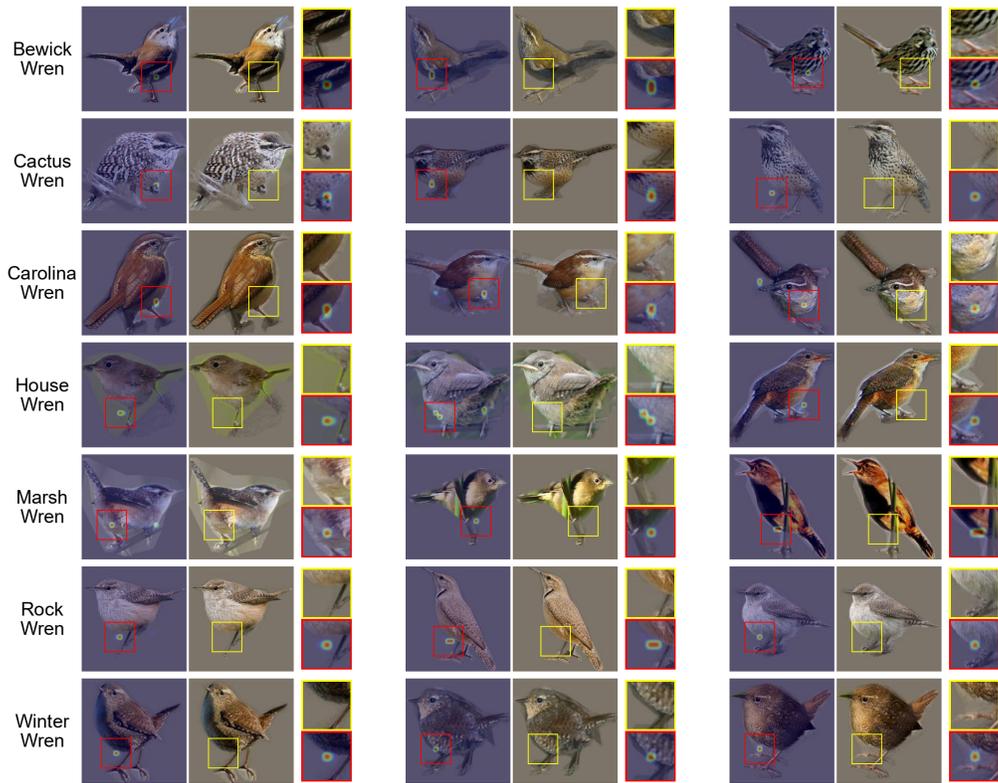


Figure 17: Top-K visualization of a prototype finding commonality between seven species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.



Figure 18: Top-K visualization of a prototype finding commonality between eight species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

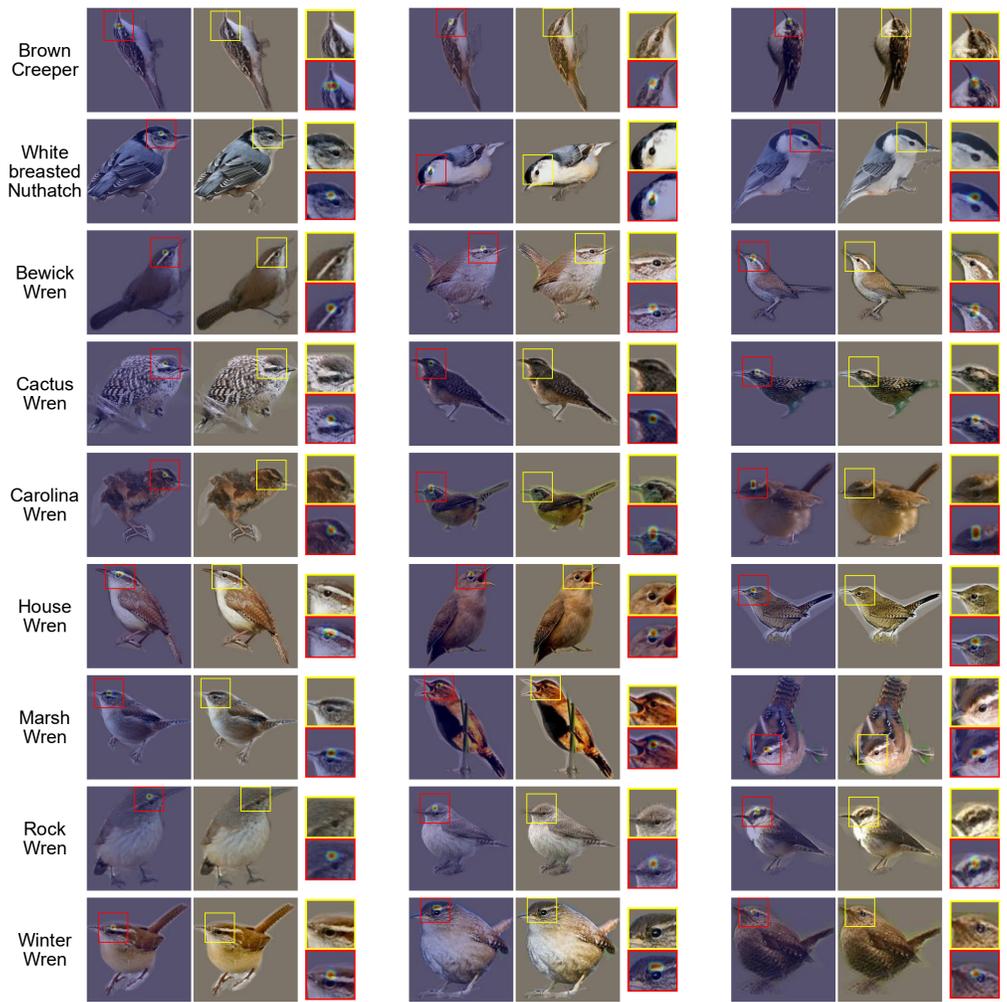


Figure 19: Top-K visualization of a prototype finding commonality between nine species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

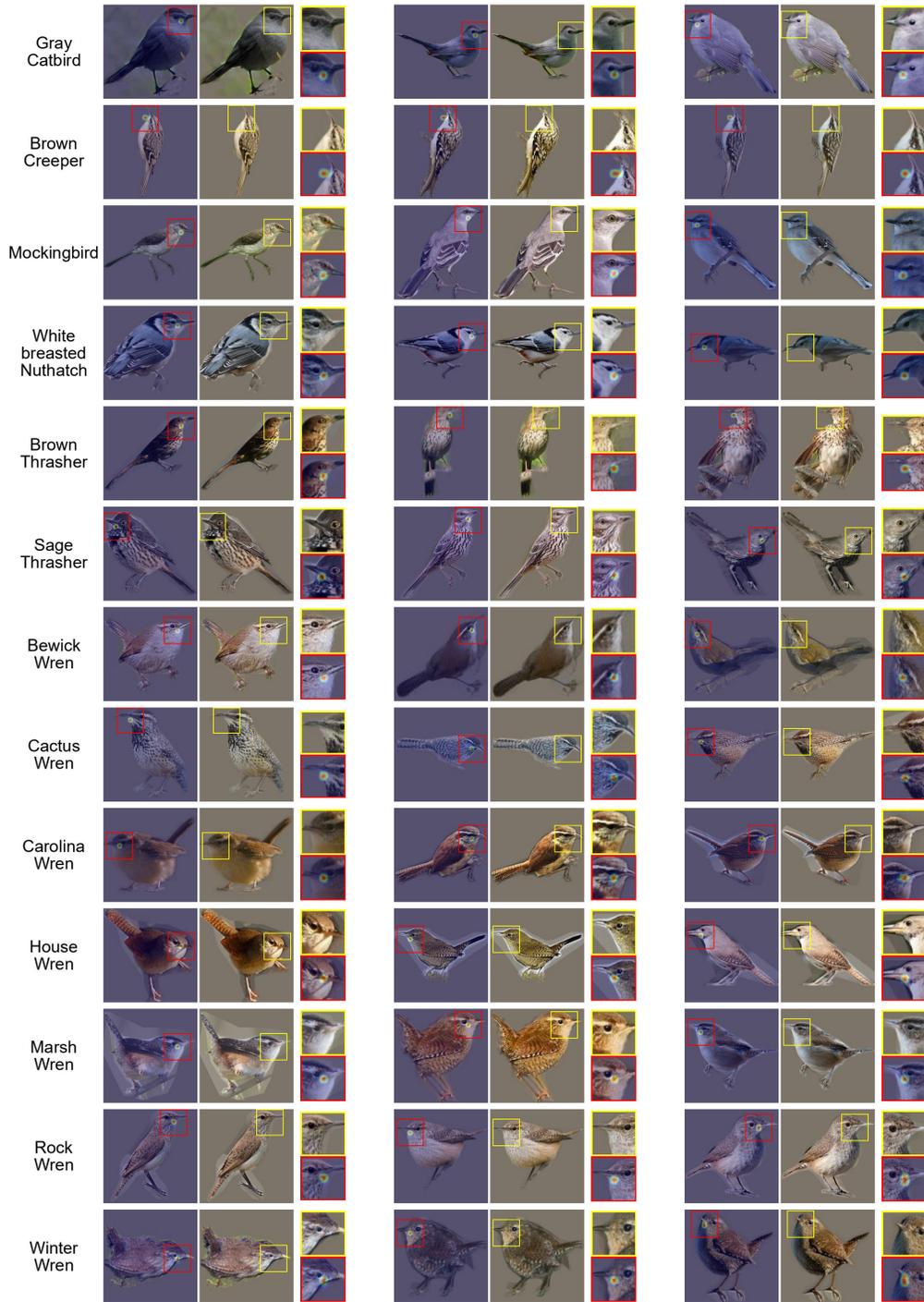


Figure 20: Top-K visualization of a prototype finding commonality between thirteen species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.



Figure 21: Top-K visualization of a prototype finding commonality between five species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.



Figure 22: Top-K visualization of a prototype finding commonality between five species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

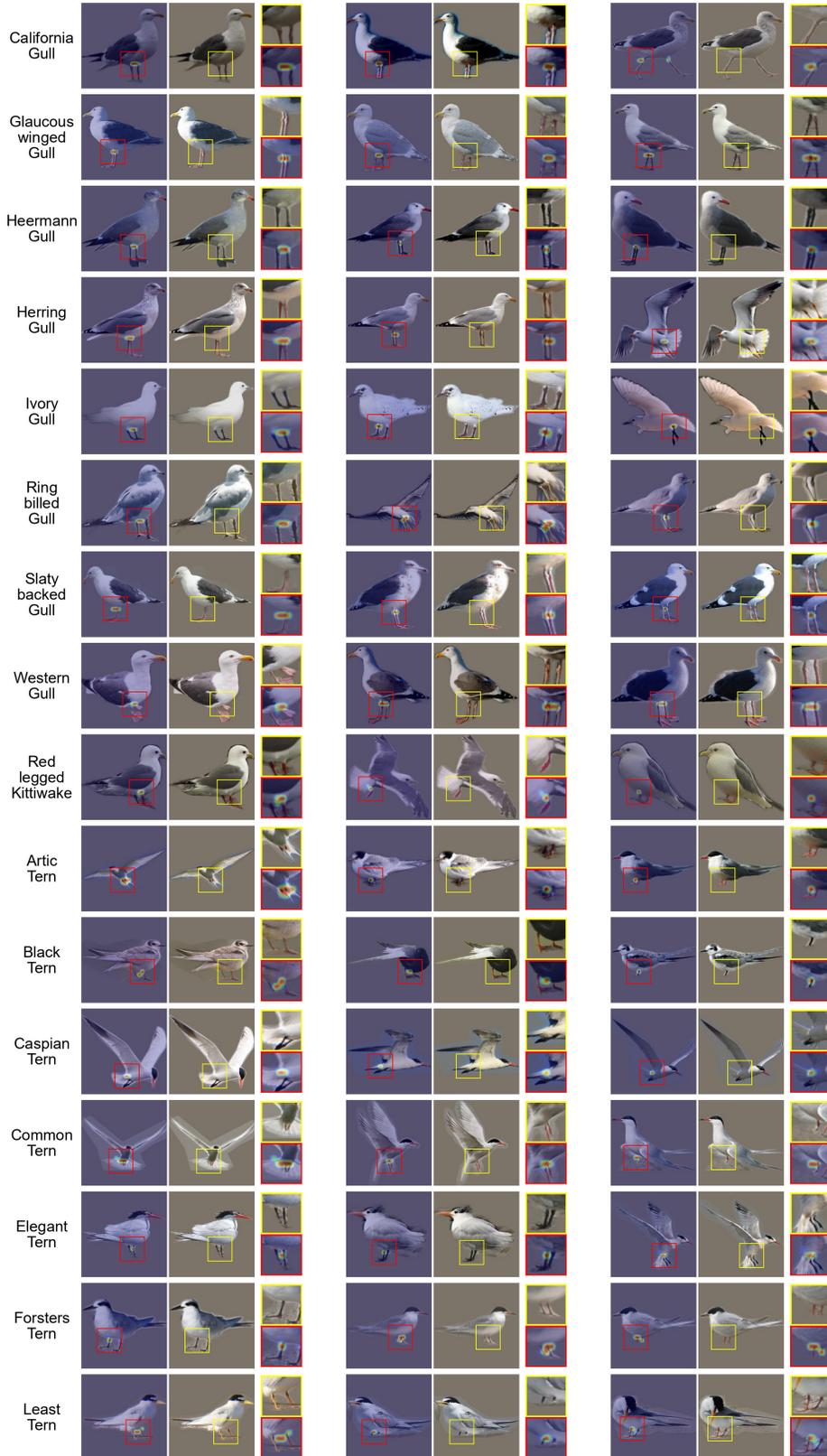


Figure 23: Top-K visualization of a prototype finding commonality between sixteen species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

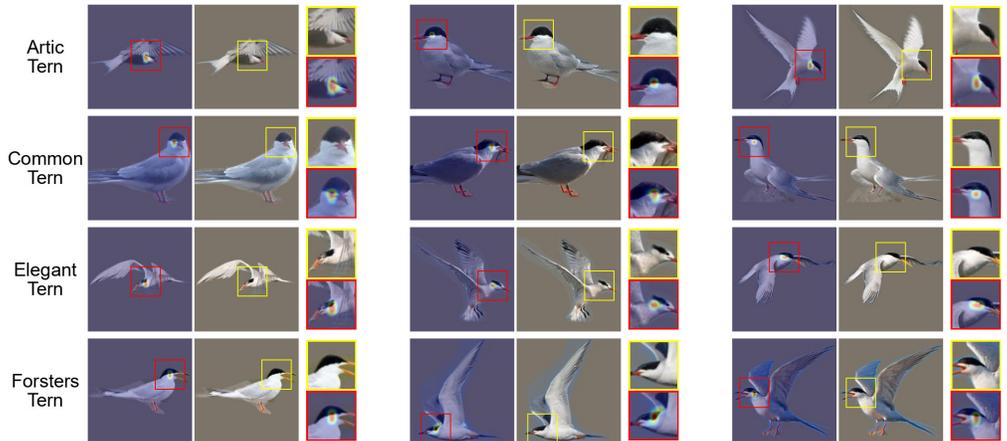


Figure 24: Top-K visualization of a prototype finding commonality between four species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.

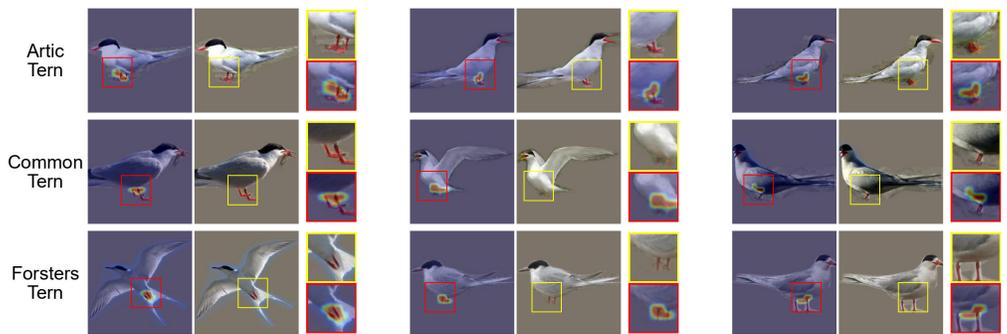


Figure 25: Top-K visualization of a prototype finding commonality between three species of birds sharing a common ancestor. Each row represents the top 3 images from the respective species. For each image we show the zoomed in view of the original image as well as the heatmap overlaid image.