

Discovering Novel Biological Traits From Images Using Phylogeny-Guided Neural Networks

Mohannad Elhamod
Virginia Tech
elhamod@vt.edu

Mridul Khurana
Virginia Tech
mridul@vt.edu

Harish Babu Manogaran
Virginia Tech
harishbabu@vt.edu

Josef C. Uyeda
Virginia Tech
juyeda@vt.edu

Meghan A. Balk
Battelle
meghan.balk@gmail.com

Wasila Dahdul
University of California, Irvine
wdahdul@uci.edu

Yasin Bakış
Tulane University
ybakis@tulane.edu

Henry L. Bart Jr.
Tulane University
hbartjr@tulane.edu

Paula M. Mabee
Battelle
mabee@battelleecology.org

Hilmar Lapp
Duke University
Hilmar.Lapp@duke.edu

James P. Balhoff
University of North Carolina at
Chapel Hill
balhoff@renci.org

Caleb Charpentier
Virginia Tech
calebc22@vt.edu

David Carlyn
The Ohio State University
carlyn.1@osu.edu

Wei-Lun Chao
The Ohio State University
chao.209@osu.edu

Charles V. Stewart
Rensselaer Polytechnic Institute
stewart@rpi.edu

Daniel I. Rubenstein
Princeton University
dir@princeton.edu

Tanya Berger-Wolf
The Ohio State University
berger-wolf.1@osu.edu

Anuj Karpatne
Virginia Tech
karpatne@vt.edu

ABSTRACT

Discovering evolutionary traits that are heritable across species on the tree of life (also referred to as a phylogenetic tree) is of great interest to biologists to understand how organisms diversify and evolve. However, the measurement of traits is often a subjective and labor-intensive process, making *trait discovery* a highly label-scarce problem. We present a novel approach for discovering evolutionary traits directly from images without relying on trait labels. Our proposed approach, *Phylo-NN*, encodes the image of an organism into a sequence of quantized feature vectors—or codes—where different segments of the sequence capture evolutionary signals at varying ancestry levels in the phylogeny. We demonstrate the effectiveness of our approach in producing biologically meaningful results in a number of downstream tasks including species image generation and species-to-species image translation, using fish species as a target example.¹

CCS CONCEPTS

- **Computing methodologies** → *Neural networks*; **Computer vision**; **Image representations**; **Neural networks**; *Regularization*;
- **Applied computing** → **Imaging**; **Imaging**.

KEYWORDS

computer vision, neural networks, phylogeny, morphology, knowledge-guided machine learning

1 INTRODUCTION

One of the grand challenges in biology is to find features of organisms—or *traits*—that define groups of organisms, their genetic and developmental underpinnings, and their interactions with environmental selection pressures [18]. Traits can be physiological, morphological, and/or behavioral (e.g., beak color, stripe pattern, and fin curvature) and are integrated products of genes and the environment. The analysis of traits is critical for predicting the effects of environmental change or genetic manipulation, and to understand the process of evolution. For example, discovering traits that are heritable across species on the tree of life (also referred to as the *phylogenetic tree*), can serve as a starting point for linking traits to underlying genetic factors. Traits with such genetic or phylogenetic signal, termed *evolutionary traits*, are of great interest to biologists, as the history of genetic ancestry captured by such traits can guide our understanding of how organisms diversify and evolve. This understanding enables tasks such as estimating the morphological features of ancestors, understanding how species have responded to environmental changes, and even predicting the potential future course of trait changes [7, 30]. However, the measurement of traits is not straightforward and often relies on subjective and labor-intensive human expertise and definitions [43]. Hence, *trait*

¹The code and datasets for running all the analyses reported in this paper can be found at <https://github.com/elhamod/phylo.nn>.

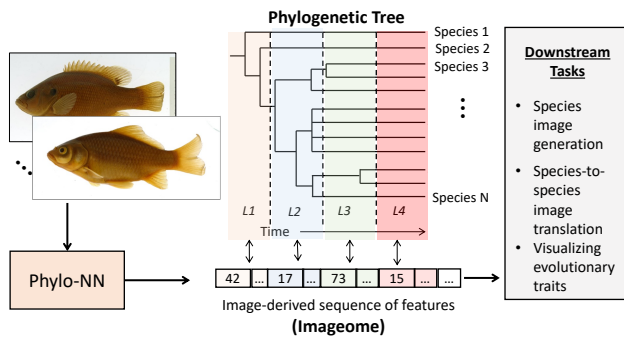


Figure 1: *Phylo-NN* converts images to discrete sequences of features (called *Imageomes*) where different sequence segments (shown in distinct colors) capture evolutionary information at varying ancestry levels of phylogeny (L1 to L4).

discovery has remained a highly label-scarce problem, hindering rapid scientific advancement [29].

With the recent availability of large-scale image repositories containing millions of images of biological specimens [45, 48, 49], there is a great opportunity for machine learning (ML) to contribute to the problem of trait discovery [29]. In particular, advances in deep learning have enabled us to extract useful information from images and to map them to structured feature spaces where they can be manipulated in a number of ways. We ask the question: *how can we develop deep learning models to discover novel evolutionary traits automatically from images without using trait labels?*

Despite the biological relevance of this question, answering it is challenging for two main reasons. First, not all image features extracted by a deep learning model for ML tasks such as image reconstruction or species classification will exhibit evolutionary signals. Hence, it is important to disentangle the image features of an organism that preserve evolutionary information, from remaining features influenced by unrelated factors [7]. Second, information about evolutionary signals is not available as a set of known attributes (or trait labels) but rather in the form of structured knowledge of how species are related to each other in the phylogenetic tree (see Figure 1). Without access to trait labels, current methods for feature disentanglement in deep learning [5, 28] are unfit for discovering evolutionary traits. Furthermore, current standards in deep learning for generative modeling [14, 41] or interpretable ML [4, 32] are unable to leverage structured forms of biological knowledge (e.g., phylogenetic trees) in the learning of image features, and hence are unable to analyze and manipulate learned features in biologically meaningful ways.

We propose a novel approach for discovering evolutionary traits automatically from images, termed *phylogeny-guided neural networks (Phylo-NN)*, which encodes the image of an organism into a sequence of quantized feature vectors or “codes” (see Figure 1). A unique feature of the image-derived sequences learned by *Phylo-NN* is that different segments of the sequence capture evolutionary information at varying ancestry levels in the phylogeny, where every level corresponds to a certain point of time in evolutionary history. Analogous to how the genome of an organism encodes all

its genetic information and structures it as a set of genes, our image-derived sequences encodes all of the visual information contained in the organism’s image and structures it as a set of evolutionary traits shared with ancestor nodes within its lineage (at different levels of phylogeny). We thus refer to the image-derived sequences of *Phylo-NN* as *Imageomes*, a brand-new concept in evolutionary biology. By analyzing and manipulating the codes in the *Imageomes* of organisms, we can perform a number of biologically meaningful downstream tasks such as species image generation, species-to-species image translation, and visualization of evolutionary traits. We demonstrate the effectiveness of *Phylo-NN* in solving these tasks using fish species as a target example.

Our work, for the first time, provides a bridge between the “language of evolution” represented as phylogeny and the “language of images” extracted by *Phylo-NN* as *Imageomes*. This work is part of a larger-scale effort to establish a new field of research in “*Imageomics*” [33], where images are used as the source of information to accelerate biological understanding of traits, ranging from their selective consequences to their causation. Our work also provides a novel methodological advance in the emerging field of knowledge-guided machine learning (KGML) [20–22] by using tree-based knowledge to structure the embedding space of neural networks and produce scientifically meaningful image generation and translation results.

2 BACKGROUND AND RELATED WORK

What is a Phylogenetic Tree? A phylogenetic tree visually characterizes the evolutionary distances among a set of species and their common ancestors represented as nodes of the tree. In this tree, the length of every edge is a value that represents the evolutionary distance between two nodes (measured in time intervals representing thousands or millions of years), which is estimated from living species and time-calibrated ages using dated fossil ancestors or molecular methods. While rates of change along different edges may vary substantially, on average we expect that longer edges will accumulate higher levels of evolutionary trait change than shorter edges. In our work, we consider discretized versions of the phylogenetic tree with $n_1 = 4$ ancestry levels, such that every species class (leaf node in the tree) has exactly $n_1 - 1$ ancestors. Every ancestry level corresponds to a certain point of time in evolutionary history. See Appendix B for details on phylogeny preprocessing.

Generative Modeling for Images: There exists a large body of work in deep learning for image generation, including methods based on Variational Autoencoders (VAEs) [26], Generative Adversarial Networks (GANs) [16, 23–25, 41], Transformer networks [11, 17], and Diffusion models [8]. While some recent advances in this field (e.g., DALL-E 2) have been shown to produce images with very high visual quality, they involve large and complex embedding spaces that are difficult to structure and analyze using tree-based knowledge (e.g., phylogeny). Instead, we build upon a recent line of work in generative modeling using vector-quantized feature representations of images [14, 34] that are easier to manipulate than continuous features. In particular, a recent variant of the VAE, termed Vector-Quantized VAE (VQVAE) [34], uses discrete feature spaces quantized using a learned codebook of feature vectors and employs a PixelCNN [46] model for sampling in the discrete feature

space. This work was extended in [14] to produce VQGAN, which is different from VQVAE in two aspects. First, it adds a discriminator to its framework to improve the quality of the generated images. Second, it uses a Transformer model, namely the GPT architecture [39], to generate images from the quantized feature space instead of a PixelCNN. VQGAN is a state-of-the-art method that generates images of better quality efficiently at higher resolutions than other counterparts such as StyleGAN [23] and Vision Transformers [11, 17]. Our work draws inspiration from VQGAN to embed images in discrete feature spaces (analogous to the discrete nature of symbols used in genome sequences) but with the grounding of biological knowledge available as phylogenetic trees.

Interpretable ML: There is a growing trend in the ML community to focus on the interpretability of deep learning features [12]. Some of the earliest works in this direction include the use of saliency scores [44] and Class Activation Maps (CAMs) [42] that reveal sensitive regions of an image influencing classification decisions. However, these methods are known to be noisy and often imprecise [1]. Recent work includes the ProtoPNet model [4], which first learns a set of template image patches (or prototypes) for each class during training, and then uses those templates to both predict and explain the class label of a test image. These methods suffer from two drawbacks. First, they do not allow for structured knowledge to guide the learning of interpretable features and hence are not designed to produce results that are *biologically meaningful*. Second, they are mostly developed for classification problems and cannot be directly applied to image generation or translation problems.

Disentangling ML Features: Another related line of research involves disentangling the feature space of deep learning models to align the disentangled features with target “concepts.” This includes the approach of “Concept whitening” (CW) [5], where the latent space of a classification model is whitened (i.e., normalized and decorrelated) such that the features along every axis of the latent space corresponds to a separate class. Another approach in this area is that of Latent Space Factorization (LSF) [28], where the latent space of an autoencoder is linearly transformed using matrix subspace projections to partition it into features aligned with target concepts (or attributes) and features that do not capture attribute information. Note that our proposed *Phylo-NN* model can also be viewed as a latent space disentanglement technique, where the disentangled segments of the learned Imageome correspond to different ancestry levels of the phylogeny. We thus use CW and LSF as baselines in our experiments to test if it is possible to discover evolutionary traits just by disentangling the latent space using species classes as orthogonal concepts, without using the structured knowledge of how species are related to one another in the phylogeny.

Knowledge-Guided ML: KGML is an emerging area of research that aims to integrate scientific knowledge in the design and learning of ML models to produce generalizable and scientifically valid solutions [22]. Some examples of previous research in KGML include modifying the architecture of deep learning models to capture known forms of symmetries and invariances [2, 51], and adding loss functions that constrain the model outputs to be scientifically consistent even on unlabeled data [9, 40]. In biology, KGML methods have been developed for species classification that leverage the

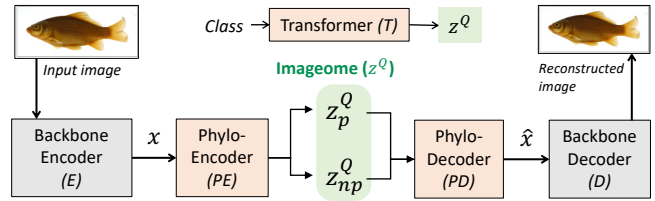


Figure 2: Overview of proposed *Phylo-NN* model architecture.

knowledge of taxonomic grouping of species [10, 13]. KGML methods have also been developed for generative modeling of images using domain knowledge available as knowledge graphs or ontologies [15, 37]. In contrast to these prior works, we focus on structuring the embedding space of neural networks using tree-based knowledge (i.e., phylogeny) to enable the discovery and analysis of novel evolutionary traits automatically from images.

3 PROPOSED APPROACH: *PHYLO-NN*

We consider the problem of discovering novel (or “unknown”) evolutionary traits automatically from images without using any trait labels or knowledge of how the unknown traits correspond to known concepts in a knowledge graph or ontology. We only use the “distant” supervision of how these unknown traits have evolved over time and are shared across species, available in the form of the phylogenetic tree. Figure 2 provides an overview of our proposed *Phylo-NN* model. Our method can operate on the latent space of any backbone encoder model E that takes in images as input and produces continuous feature maps x as output. There are three computing blocks in *Phylo-NN* as shown in Figure 2. The first block, *Phylo-Encoder (PE)*, takes continuous feature maps x as input and generates quantized feature sequences (or Imageomes) as output. Imageome sequences z^Q comprise of two *disentangled* segments: z_p^Q , which captures phylogenetic information (p) at varying ancestry levels, and z_{np}^Q , which captures non-phylogenetic information (np) that is still important for image reconstruction but is unrelated to the phylogeny. The second block, *Phylo-Decoder (PD)*, maps the Imageome sequences back to the space of feature maps \hat{x} , such that \hat{x} is a good reconstruction of x . We then feed \hat{x} into a backbone decoder model D that reconstructs the original image. Note that in the training of *PE* and *PD* models, both the backbone models E and D are kept frozen, thus requiring low training time. *Phylo-NN* can thus be plugged into the latent space of any powerful encoder-decoder framework. The third block of *Phylo-NN* is a transformer model T that takes in the species class variable as input, and generates a distribution of plausible Imageome sequences corresponding to the class as output. These sequences can be fed to the *PD* model to generate a distribution of synthetic images. In the following, we provide details on each of the three blocks of *Phylo-NN*.

3.1 *Phylo-Encoder (PE)* Block

Figure 3 shows the sequence of operations that we perform inside the *PE* block. We first apply a convolutional layer on x to produce feature maps of size $(H \times W \times C)$, where C is the number of channels. We split these C feature maps into two sets. The first C_p maps are fed into an MLP layer to learn a global set of feature vectors z_p capturing phylogenetic information. The size of z_p is kept equal to

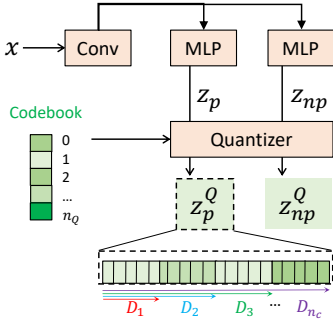


Figure 3: Detailed view of the Phylo-Encoder block.

$(n_1 n_p \times d)$, where n_1 is the number of phylogeny levels, n_p is the number of feature vectors we intend to learn at every phylogeny level, and d is the dimensionality of feature vectors. Similarly, the remaining $C - C_p$ maps are fed into an MLP layer to produce a set of feature vectors z_{np} capturing non-phylogenetic information of size $(n_{np} \times d)$.

Vector Quantization: Both z_p and z_{np} are converted to *quantized* sequences of feature vectors, z_p^Q and z_{np}^Q , respectively, using the approach developed in VQVAE [34]. The basic idea of this quantization approach is to learn a set (or codebook) of n_q distinct feature vectors (or codes), such that every feature vector in z_p and z_{np} is replaced by its nearest counterpart in the codebook. This is achieved by minimizing the *quantization loss*, $L_q = |z - z^Q|$. The advantage of working with quantized vectors is that every feature vector in z_p^Q and z_{np}^Q can be referenced just by its location (or index) in the codebook. This allows for faster feature manipulations in the space of discrete code positions than continuous feature vectors.

Using phylogenetic knowledge in z_p^Q : Here, we describe our approach to ensure that the quantized feature sequence z_p^Q contains phylogenetic information. Note that z_p^Q contains n_1 sub-sequences of length n_p , where every sub-sequence corresponds to a different ancestry level in the phylogeny. While the first sub-sequence S_1 should ideally capture information contained in x that is necessary for identifying ancestor nodes at level 1 of the phylogeny, S_2 should contain additional information that when combined with S_1 is sufficient to identify the correct ancestor node of x at level 2. In general, we define the concept of a *Phylo-descriptor* $D_i = \{S_1, S_2, \dots, S_i\}$ of x that contains the necessary information for identifying nodes at level i (see Figure 3). We feed D_i to an MLP layer that predicts the class probabilities of nodes at level i , which are then matched with the correct node class of x at level i , $c_i(x)$, by minimizing the following *phylogeny-guided loss*, L_p :

$$L_p = \sum_{i=0}^{n_1} \beta_i \text{CE}(\text{MLP}_i(D_i(x)), c_i(x)), \quad (1)$$

where CE is the cross-entropy loss and β_i is the weighting hyperparameter for level i .

Disentangling z_p^Q and z_{np}^Q : While minimizing L_p guides the learning of z_p^Q to contain phylogenetic information, we still need a way to ensure that z_{np}^Q focuses on complementary features and does not contain phylogenetic information. To achieve this, we first apply an orthogonal convolution loss L_o (originally proposed in [50]) to

the convolutional layer of Phylo-Encoder, to constrain the C convolutional kernels to be orthogonal to each other. To further ensure that z_{np}^Q has no phylogenetic information, we also employ an adversarial training procedure to incrementally remove phylogenetic information from z_{np}^Q . In particular, we apply an MLP layer MLP_{adv} on z_{np}^Q , and then train the parameters of MLP_{adv} to minimize the following *adversarial loss*:

$$L_{adv} = \sum_{i=0}^{n_1} \beta_i \text{CE}(\text{MLP}_i(\text{MLP}_{adv}(z_{np}^Q(x))), c_i(x)), \quad (2)$$

This is aimed at training MLP_{adv} to detect any phylogenetic information contained in z_{np}^Q . Simultaneously, we train the rest of *Phylo-NN*'s parameters to maximize L_{adv} , such that z_{np}^Q becomes irrelevant for the task of identifying nodes in the phylogeny and only contains non-phylogenetic information.

4 PHYLO-DECODER (PD) BLOCK

The goal of the PD block is to convert the space of Imageome sequences, $z^Q = \{z_p^Q, z_{np}^Q\}$, back to the space of original feature maps, x . The sequence of operations in *PD* is almost a mirror image of those used in *PE*. We first pass z_p^Q and z_{np}^Q through two MLPs, and then concatenate their outputs to create feature maps of size $(H \times W \times C)$. These feature maps are then fed into a convolutional layer to produce \hat{x} . Minimizing the reconstruction loss, $L_{rec} = |\hat{x} - x|$, ensures that \hat{x} is a good approximation of x . Finally, *PE* and *PD* are jointly trained using a weighted summation of all the losses mentioned above.

4.1 Transformer (T) Block

Once *PE* and *PD* are trained, we can extract Imageome sequences z^Q for every image in the training set. The goal of the Transformer block is to learn the patterns of codes in the extracted Imageome sequences of different classes (e.g., species class or ancestor node class), and use these patterns to generate synthetic Imageome sequences for every class. To achieve this task, we follow the approach used by VQGAN [14] and train a GPT transformer model [39] T_i to generate plausible sequences of z^Q for every node class at level i . The generated Imageome sequences can then be converted into synthesized specimen images using *PD* and D .

5 EVALUATION SETUP

5.1 Data

We used a curated dataset of teleost fish images from five ichthyological research collections that participated in the Great Lakes Invasives Network Project (GLIN). After obtaining the raw images from these collections, we handpicked a subset of about 11,000 images and pre-processed them by resizing and appropriately padding each image to be of a 256×256 pixel resolution. Finally, we partitioned the images into a training set and a validation set using an 80 – 20 split. See Appendix A for details on data pre-processing.

Our dataset includes images from 38 species of teleost fishes with an average number of 200 images per species. We discretized the phylogenetic tree to have $n_1 = 4$ ancestry levels, where the last

level is the species class. See Appendix B for details on phylogeny selection and discretization.

5.2 Backbone Encoder and Decoder

Since *Phylo-NN* can operate on the feature space \mathbf{x} of any backbone encoder E and produce reconstructed feature maps $\hat{\mathbf{x}}$ that can be decoded back to images by a corresponding backbone decoder D , we tried different encoder-decoder choices including pix2pix [19], ALAE [35], and StyleGAN [25]. However, we found VQGAN [14] feature maps to produce images of better visual quality than other encoder-decoder models. Hence, we used the embeddings of a base VQGAN encoder E as inputs in *Phylo-NN* for all our experiments. The reconstructed feature maps of *Phylo-NN* were then fed into a base VQGAN quantizer serving as the backbone decoder D . Note that while training *Phylo-NN*, we kept the parameters of the backbone models fixed, thus saving training time and resources.

5.3 Baseline Methods

Since no direct baselines exist for structuring the embedding space of neural networks using tree-based knowledge or discovering novel evolutionary traits from images, we considered the following baselines that are closest in motivation to *Phylo-NN*:

Vanilla VQGAN [14]: The first baseline that we consider is a vanilla VQGAN model trained to generate and reconstruct images on the fish dataset. By comparing the learned embeddings and generated images of *Phylo-NN* with vanilla VQGAN, we aim to demonstrate the importance of using biological knowledge to structure the embedding space of neural networks for trait discovery, rather than solely relying on information contained in data.

Concept whitening (CW) [5]: For this second baseline, we replaced the last normalization layer in the encoder block of vanilla VQGAN with the concept whitening (CW) module, where we used species class labels as concept definitions. This is intended to evaluate if CW is capable of disentangling the evolutionary traits of species automatically from images without using the phylogeny. The whitened embeddings \mathbf{z}_{CW} produced by the CW module are fed into the quantizer module of vanilla VQGAN for converting the embeddings to images. While training the CW module, we optimized the whitening and rotation matrices for all concepts every 30 batches. We used the VQGAN’s transformer to generate plausible feature sequences \mathbf{z}_{CW} conditioned on the species label, which are then decoded into specimen images using the VQGAN’s decoder.

Latent Space Factorization (LSF) [28]: The third baseline that we considered is the LSF method, which is another approach for feature disentanglement given concept attribute labels. Specifically, we introduced a variational autoencoder (VAE) model between the encoder and the quantization layer of the base VQGAN model. Similar to CW, we used the species class of each image as the concept attribute for factorizing the latent space in LSF. The LSF module was trained to optimize VAE’s KL-divergence loss and recreation loss along with the attribute and non-attribute losses, as originally defined in the LSF method [28].

6 RESULTS

In the following, we analyze the results of *Phylo-NN* from multiple angles to assess the quality of its learned embeddings and generated images in comparison with baseline methods.

6.1 Validating Species Distances in the Embedding Space

In order to evaluate the ability of *Phylo-NN* to extract novel (or unknown) evolutionary traits from images without using trait labels, we show that distances between species pairs in the embedding space of *Phylo-NN* are biologically meaningful and are correlated with ground-truth values better than baseline methods. In the following, we describe the two types of ground-truths used, the approach used for computing distances in the embedding space of comparative methods, and the comparison of correlations with ground-truth values.

Phylogenetic Ground-truth (GT): The first ground-truth distance between pairs of species is the *evolutionary distance* between their corresponding nodes in the phylogenetic tree. In particular, for any two species, we can calculate the total sum of edge lengths in the path between their nodes in the phylogenetic tree. The longer the path, the more distant the species are on the evolutionary scale. Hence, if *Phylo-NN* indeed captures evolutionary traits in its embedding space, we would expect it to show higher correlations with evolutionary distances computed from the phylogeny as compared to baselines. We applied min-max scaling of evolutionary distances so that they range from 0 to 1.

Morphological Ground-truth (GT): Another type of ground-truth distance between species was computed based on measurements of known morphological traits obtained from the FishShapes v1.0 dataset [36], which contains expert-measured traits known to carry evolutionary signals, defined and collected using traditional methods that are subjective and labor-intensive. We specifically used 8 functionally relevant traits from this dataset for every fish species. Some species were not available in this dataset, so when possible, either the closest relative was substituted or the species was dropped. The species were then matched to a time-calibrated phylogeny of fishes [3, 38] and the log-transformed measurements were rotated with phylogenetically-aligned components analysis (PACA) [7], which rotates the traits to the axis with the highest level of phylogenetic signal. After correcting for overall size and allometry, the principal components of PACA were used to compute the Mahalanobis distance between every species-pair, using a covariance matrix proportional to the evolutionary rate matrix. See Appendix D for details on PACA calculations.

Computing Embedding Distances: To compute pair-wise distances in the embedding space of *Phylo-NN*, we first compute the probability distributions (or histograms) of quantized codes at every position of the Imageome sequence (i.e., \mathbf{z}_{p}^Q and \mathbf{z}_{np}^Q) in the test images for every species. We then compute the *Jensen-Shannon (JS) divergence* [31] between the probability distributions of codes at a pair of species to measure the dissimilarity of their learned embeddings. We adopt a similar approach for computing the JS-divergence of species-pairs in the quantized feature space of vanilla VQGAN. For baseline methods that operate in continuous feature spaces (CW and LSF), we first calculate the mean feature vector for every

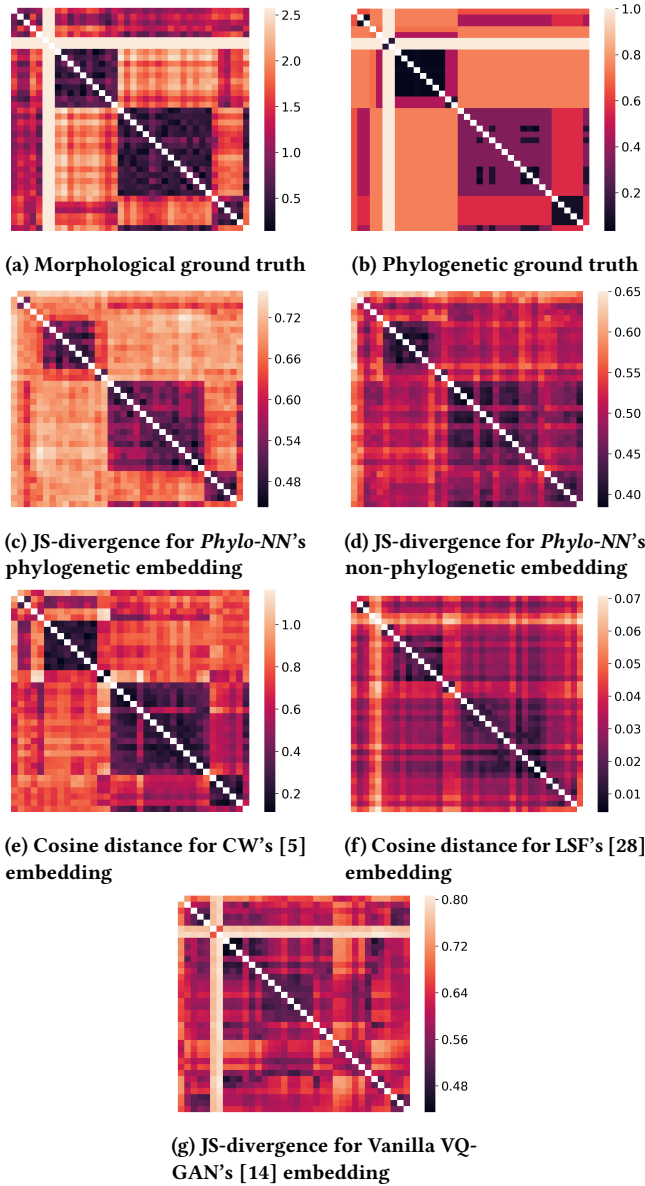


Figure 4: Comparing embedding distance matrices of methods with morphological and phylogenetic ground-truths.

species and then compute the *cosine distance* ($1 - \text{cosine similarity}$) of vectors for a pair of species. For both metrics, JS-divergence and cosine distance, a value closer to 0 represents higher similarity.

Comparing Correlations with Ground Truth: Figure 4(a) and Figure 4(b) show the pair-wise species distance matrices for morphological and phylogenetic GTs, respectively. Note that the rows and columns of all matrices in Figure 4 are species ordered according to their position in the phylogeny (see Appendix B for details) and the diagonal values (correlation with self) are removed so that they do not affect the colormap scale. We can see that both ground-truths show a similar clustering structure of species. However, there are differences too; while phylogenetic GT is solely based on phylogeny, the morphological GT uses both the phylogeny and

Table 1: Correlations between GT and embedding distances

	Morphological	Phylogenetic	
PhyloNN	level0	0.86	0.83
	level1	0.87	0.85
	level2	0.78	0.83
	species	0.70	0.78
LSF	0.38	0.55	
CW	0.70	0.67	
vanilla VQGAN	0.31	0.24	

information about “known” traits. Figure 4(c) and Figure 4(d) show the JS-divergences among species computed separately for the two disentangled parts of *Phylo-NN*’s embeddings (\mathbf{z}_p^Q and \mathbf{z}_{np}^Q). We can see that the embeddings containing phylogenetic information show a similar clustering structure of distances as the GT matrices, in contrast to the non-phylogenetic embeddings. This shows the ability of *Phylo-NN* to disentangle features related to phylogeny from other unrelated features. Figure 4 also shows the embedding distance matrices of the baseline methods, which are not as visually clean as *Phylo-NN* in terms of matching with the GT matrices.

To quantitatively evaluate the ability of *Phylo-NN* to match with GT distances, we compute the Spearman correlation between the GT distance matrices and embedding distance matrices for different methods as shown in Table 1. We can see that *Phylo-NN* shows higher correlations at the species level with both GTs. Furthermore, since *Phylo-NN* learns a different descriptor for every ancestry level in contrast to baseline methods that learn a flat representation, we can also compute *Phylo-NN*’s distance matrix at any ancestry level and compare it with GT matrices at the same level. Table 1 shows that *Phylo-NN* shows significantly higher correlations with GT matrices at higher ancestry levels than the species level.

6.2 Evaluating Species-to-species Image Translations

To further assess how well *Phylo-NN*’s embeddings capture evolutionary traits, we investigate how altering the learned Imageome sequence of an image specimen incrementally in a phylogenetically meaningful ordering affects the observed traits when the altered embeddings are decoded back as an image. To do that, we set up the following experiment. We pick two specimen images from a pair of species. By encoding the two images using *Phylo-NN*, we obtain their corresponding Imageome encodings, \mathbf{z}_1^Q and \mathbf{z}_2^Q . We then start to replace the codes in the Imageome sequence \mathbf{z}_1^Q with the corresponding codes in \mathbf{z}_2^Q iteratively, until \mathbf{z}_1^Q transforms completely into \mathbf{z}_2^Q . The order of this iterative replacement is by first replacing the codes representing the non-phylogenetic part of the embedding $\mathbf{z}_{np,1}^Q$, then the part capturing evolutionary information at the earliest ancestry level (level 0), to the next ancestry level (level 1), till we eventually reach the last level of the phylogeny, which is the species level. At the final point, the entire Imageome sequence \mathbf{z}_1^Q has been replaced with \mathbf{z}_2^Q . This phylogeny-driven ordering of code replacements helps us capture key “snapshots” of the species-to-species translation process that are biologically

meaningful. In particular, by observing the traits that appear or disappear at every ancestry level of code replacement, we can infer and generate novel hypotheses about the biological timing of trait changes as they may have happened in evolutionary history.

Figure 5 shows an example of such a translation process between a specimen of the species *Carassius auratus* to a specimen of the species *Lepomis cyanellus*. We can see that although the two specimens look similar on the surface, there are several subtle traits that are different in the two species that are biologically interesting. For example, the source species has a V-shaped tail fin (termed caudal fin), while the target species has a rounded caudal fin. By looking at their place of occurrence in the translation process of *Phylo-NN*, we can generate novel biological hypotheses of whether they are driven by phylogeny or not, and whether they appeared earlier or later in the target species in the course of evolution. For example, we can see that the rounded tail feature of the target species appears right after replacing the non-phylogenetic part of the embedding (see blue circle), indicating that this feature may not be capturing evolutionary signals and instead maybe affected by unrelated factors (e.g., environment). On the other hand, if we observe another fin (termed pectoral fin) that appears on the side of the body just behind the gill cover, compared to the fin’s lower position closer to the underside of the specimen in the source image, we can see that it seems to get sharper and compact only in the later levels (it is faintly visible in level 1 but shows up prominently as a white region in level 2, see green circle). This suggests that the change in the position and shape of the pectoral fin occurred later in fish evolution, which is supported by the phylogeny and is in fact the case. Our work opens novel opportunities for generating such biological hypotheses, which can be further investigated by biologists to potentially accelerate scientific discoveries. Figure 5 also shows the translations obtained by baseline methods for the same pair of species specimens. We can see that the baselines are mostly performing a smooth interpolation between the source and target images. This is in contrast to the discrete and non-smooth nature of changes observed in the translation of *Phylo-NN*, which is indeed the desired behavior since the appearance or disappearance of traits at every ancestry level are expected to be orthogonal to those at other levels. Furthermore, the transition points in the translation process of baseline methods do not correspond to biologically meaningful events as opposed to *Phylo-NN*.

6.3 Generalization to Unseen Species

As *Phylo-NN* aims to encode specimen images into their corresponding phylogenetic and non-phylogenetic sequences of codes, we expect specimens that belong to the same species to largely share the same phylogenetic code in terms of the species descriptor D_{n_i} , while varying in terms of the non-phylogenetic codes. More generally, specimens belonging to species sharing a common ancestor at phylogenetic level i should largely share the codes with the descriptor at level i , D_i , while varying at the rest of codes. This should also apply for specimens of *unseen* (or newly discovered) species that we have not yet observed in the training set. We posit that by looking at the similarity of codes generated for an unseen species, we should be able to infer its ancestral lineage in terms of the species sampled during training. In other words, by analyzing

the distribution of codes generated from the image of an unseen species, we should be able to locate it on the phylogenetic tree and position it to next to the subset of known species (seen during training) that share a common ancestor.

To quantify this phenomenon, for a given species or ancestor node, we construct two sets of histograms, H_p and H_{np} , of sizes $[n_1 \times n_p]$ and $[n_{np}]$, respectively. Each value in the histograms, $H_p^{i,j}$ and H_{np}^k , describes the distribution of codes at a certain location in the Imageomes across all specimens belonging to the species or ancestor node. See Appendix F for an example. We then compute the entropy of each sequence location in H_p and H_{np} to measure the “purity” of codes used at every location. If the entropy is low for a certain location, it means only a few possible codes occur at that location, suggesting that those specific codes are key at characterizing the species or ancestor node in question. On the other hand, higher entropy means a variety of codes occur at that location, implying that such a location is not discriminative to the species or ancestor node. Finally, to compare the code distributions for two species, we use the JS-divergence metric for calculating the difference between two histograms of a sequence location. Similar to Section 6.1, such a metric can be aggregated to quantify the coding differences between species-pairs.

To assess *Phylo-NN*’s ability to generalize to unseen species, we train it on a subset of the species and then evaluate the quality of the embedding space when the model is introduced to species it has never seen before during training. In our experiment, we chose to train on the same dataset as before while only excluding three species. Once the model is trained, we look at the average JS-divergence distance between these missing species and three other species in the tree. These three other species were selected such that each missing species has one seen species that is close to it phylogenetically (i.e., both species share the same ancestor at the immediate ancestry level) while others are relatively far from it.

Table 2 shows the average distance of the phylogenetic codes among the six aforementioned species. We can see that the distance is smallest for each unseen species and its counterpart that shares the same immediate ancestor (shown as the diagonal in the table). This confirms that even though the model has not seen the former species, it was able to characterize it using an Imageome sequence that is significantly closer to that of its seen counterpart than the other species’ Imageomes.

While Table 2 highlights the phylogenetic matching in the embedding space at the species descriptor level, D_{n_i} , Table 3 does the same but for the descriptor at a distant ancestry level (level 0), i.e., D_0 . Based on the phylogenetic tree we have used in this example, both the *Notropis* and *Noturus* species share the same distant ancestor at that descriptor level. On the other hand, *Lepomis* species does not share that ancestor. Hence, we find that the JS-divergences increase for the *Lepomis* unseen species with seen species that are not *Lepomis* as compared to Table 2. On the other hand, the JS-divergences decrease for the other two unseen species w.r.t. seen species that are on the off-diagonals of the table. This confirms that D_0 specifically captures the phylogenetic information of that distant ancestor that is common across *Notropis* and *Noturus* seen and unseen species. Finally, to confirm that this phylogenetic correlation is mainly constrained only to the phylo-descriptors, we

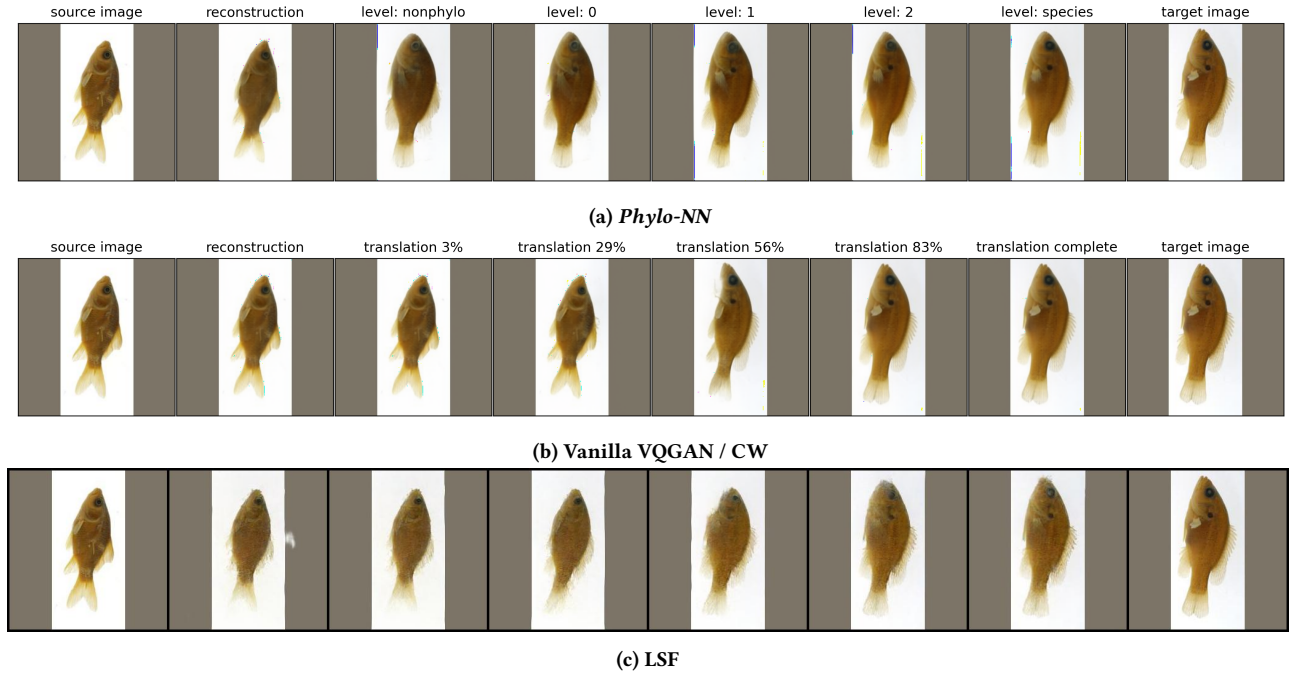


Figure 5: Comparing species-to-species image translations from a *Carassius auratus* specimen to a *Lepomis cyanellus* specimen.

Table 2: JS-divergence of the phylogenetic codes at the species level between unseen and seen species

		Seen species		
		<i>Notropis nubilus</i>	<i>Lepomis macrochirus</i>	<i>Noturus flavus</i>
Unseen species	<i>Notropis percobromus</i>	0.47	0.71	0.62
	<i>Lepomis megalotis</i>	0.73	0.43	0.72
	<i>Noturus miurus</i>	0.62	0.71	0.48

calculate the same distances but using the non-phylogenetic part of the sequences. The result is shown in Table 4. We can see that the distances are much closer to each other, implying that the non-phylogenetic embedding is not specialized at differentiating among different species, and hence cannot be used to phylogenetically categorize the unseen species.

6.4 Assessing the Clustering Quality of the Embedding Space Using t-SNE Plots

In this section, we qualitatively assess the quality of generated images by visualizing their embedding space. Visualization tools such as loss landscape visualizations [27] and t-SNE plots [47], have been frequently used as investigative tools in deep learning in recent years as they help gauge a model’s generalization power. To that end, we are interested in understanding how *Phylo-NN* clusters the

Table 3: JS-divergence of the phylogenetic codes at the earliest ancestral level between unseen and seen species

		Seen species		
		<i>Notropis nubilus</i>	<i>Lepomis macrochirus</i>	<i>Noturus flavus</i>
Unseen species	<i>Notropis percobromus</i>	0.26	0.81	0.50
	<i>Lepomis megalotis</i>	0.81	0.27	0.81
	<i>Noturus miurus</i>	0.52	0.80	0.31

Table 4: JS-divergence of the non-phylogenetic codes between unseen and seen species

		Seen species		
		<i>Notropis nubilus</i>	<i>Lepomis macrochirus</i>	<i>Noturus flavus</i>
Unseen species	<i>Notropis percobromus</i>	0.39	0.45	0.39
	<i>Lepomis megalotis</i>	0.46	0.36	0.48
	<i>Noturus miurus</i>	0.40	0.42	0.36

embedding space compared to other baselines by analyzing these models’ t-SNE plots. To construct the t-SNE plot for each model,

we iterate through its generated images, encode them, obtain the quantized embedding vector for each image (z_p^Q and z^Q for *Phylo-NN* and vanilla VQGAN, respectively), and finally create the t-SNE plots. For CW, we use the whitened embeddings z_{cw} instead.

Figure 6 shows these constructed t-SNE plots with two different color-coding schemes. The first one (left column) color-codes the data-points based on the grouping of species at the second phylogenetic level (i.e., the direct ancestor of the specimen’s species). Using this color-coding scheme allows us to inspect how different species cluster in the embedding space. The second color-coding (right column) is the average phylogenetic distance between the data-point and its k -nearest neighbors (KNN), where $k = 5$ in this setup. The higher the average distance (i.e., the darker the data-point’s color), the more distant the specimen is from those k specimen’s that are closest to it in the quantized embedding space. This color-coding helps us spot how well the different species are separated from each other in the embedding space, which generally characterizes the quality of the encoding and its propensity for downstream tasks, such as classification.

From Figure 6, we can see that *Phylo-NN* (top row) clusters the generated images better than vanilla VQGAN and CW as evident from its hierarchical clustering where the specimens belonging to the same species clump into small clusters and these clusters in turn clump into larger clusters (representing ancestor nodes) that have a singular color. This demonstrates that *Phylo-NN* is able to learn a phylogenetically-meaningful encoding, whereas the other base models’ clustering is quite fuzzy and poorly characterizes any biological knowledge. Also, by looking at the right column, we can see that *Phylo-NN* commits very little clustering error in terms of its phylogenetic constraints because the average phylogenetic distance is low (almost zero) for the majority of points. This is in contrast to the other baselines where there is quite a high clustering error as seen from the “heat” of its scatter plot.

7 CONCLUSIONS AND FUTURE WORK

In this work, we presented a novel approach of *Phylo-NN* for discovering biological traits related to evolution automatically from images in an unsupervised manner without requiring any trait labels. The key novelty of our approach is to leverage the biological knowledge of phylogeny to structure the quantized embedding space of *Phylo-NN*, where different parts of the embedding capture phylogenetic information at different ancestry levels of the phylogeny. This enables our method to perform a variety of tasks in a biologically meaningful way such as species-to-species image translation and identifying the ancestral lineage of newly discovered unseen species.

In the future, our work can be extended to include a larger number of embedding dimensions to improve the visual quality of generated images and can be applied to other image datasets beyond the fish dataset. Future work can explore extensions of *Phylo-NN* to generate images of ancestor species or to predict images of species that are yet to be evolved. Future work can also focus on making the discovered Imageome sequences more explainable by understanding the correspondence of each quantized code with a region in the image space. Our work opens a novel area of research in grounding image representations using tree-based knowledge, which can lead

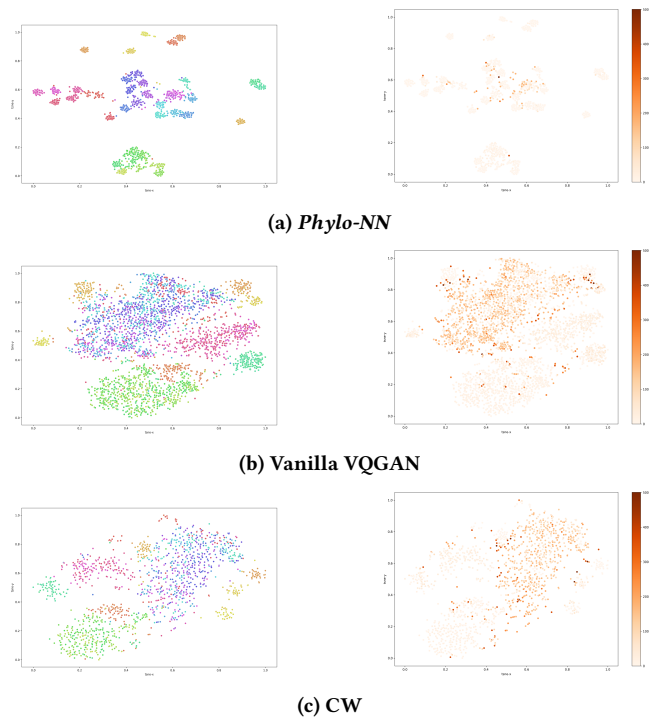


Figure 6: t-SNE plots of the images generated by *Phylo-NN* and other baselines.

to new research paradigms in other fields of science where images are abundant but labels are scarce.

ACKNOWLEDGMENTS

This work was supported, in part, by NSF awards for the HDR Imageomics Institute (Award # 2118240) and the Biology-guided Neural Network (BGNN) projects (Award # 1940247, # 1940322, # 1940233, # 2022042, # 1939505). Access to computing facilities was provided by the Advanced Research Computing (ARC) Center at Virginia Tech.

REFERENCES

- [1] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian Goodfellow, Moritz Hardt, and Been Kim. 2018. Sanity checks for saliency maps. *Advances in neural information processing systems* 31 (2018).
- [2] Brandon Anderson, Truong Son Hy, and Risi Kondor. 2019. Cormorant: Covariant Molecular Neural Networks. *Advances in Neural Information Processing Systems* 32 (2019), 14537–14546.
- [3] Jonathan Chang, Daniel L Rabosky, Stephen A Smith, and Michael E Alfaro. 2019. An R package and online resource for macroevolutionary studies using the ray-finned fish tree of life. *Methods in Ecology and Evolution* 10, 7 (2019), 1118–1124.
- [4] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. 2019. This Looks Like That: Deep Learning for Interpretable Image Recognition. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/adf7ee2dcf142b0e11888e72b43fcb75-Paper.pdf>
- [5] Zhi Chen, Yijie Bei, and Cynthia Rudin. 2020. Concept whitening for interpretable image recognition. *Nature Machine Intelligence* 2, 12 (2020), 772–782.
- [6] Julien Clavel, Gilles Escarguel, and Gildas Merceron. 2015. mvMORPH: an R package for fitting multivariate evolutionary models to morphometric data. *Methods in Ecology and Evolution* 6, 11 (2015), 1311–1319.
- [7] Michael L. Collyer and Dean C. Adams. 2021. Phylogenetically aligned component analysis. *Methods in Ecology and Evolution*

- 12, 2 (2021), 359–372. <https://doi.org/10.1111/2041-210X.13515> arXiv:<https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13515>
- [8] Florinel-Alin Croitoru, Vlad Hondru, Radu Tudor Ionescu, and Mubarak Shah. 2023. Diffusion models in vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2023).
- [9] Arka Daw, Anuj Karpatne, William D Watkins, Jordan S Read, and Vipin Kumar. 2017. Physics-guided neural networks (pgnn): An application in lake temperature modeling. In *Knowledge-Guided Machine Learning*. Chapman and Hall/CRC, 353–372.
- [10] Anderson Aparecido dos Santos and Wesley Nunes Gonçalves. 2019. Improving Pantanal fish species recognition through taxonomic ranks in convolutional neural networks. *Ecological Informatics* 53 (2019), 100977.
- [11] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xi-aohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).
- [12] Mengnan Du, Ninghao Liu, and Xia Hu. 2019. Techniques for interpretable machine learning. *Commun. ACM* 63, 1 (2019), 68–77.
- [13] Mohannad Elhamod, Kelly M. Diamond, A. Murat Maga, Yasin Bakis, Henry L. Bart Jr., Paula Mabee, Wasila Dahdul, Jeremy Leipzig, Jane Greenberg, Brian Avants, and Anuj Karpatne. 2022. Hierarchy-guided neural network for species classification. *Methods in Ecology and Evolution* 13, 3 (2022), 642–652. <https://doi.org/10.1111/2041-210X.13768> arXiv:<https://besjournals.onlinelibrary.wiley.com/doi/pdf/10.1111/2041-210X.13768>
- [14] Patrick Esser, Robin Rombach, and Bjorn Ommer. 2021. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 12873–12883.
- [15] Raissa Garozzo, Cettina Santagati, Concetto Spampinato, and Giuseppe Vecchio. 2021. Knowledge-based generative adversarial networks for scene understanding in Cultural Heritage. *Journal of Archaeological Science: Reports* 35 (2021), 102736.
- [16] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.
- [17] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. 2022. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 16000–16009.
- [18] David Houle and Daniela M Rossoni. 2022. Complexity, Evolvability, and the Process of Adaptation. *Annual Review of Ecology, Evolution, and Systematics* 53 (2022).
- [19] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. 2017. Image-to-Image Translation with Conditional Adversarial Networks. *CVPR* (2017).
- [20] George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. 2021. Physics-informed machine learning. *Nature Reviews Physics* 3, 6 (2021), 422–440.
- [21] Anuj Karpatne, Gowtham Atluri, James H Faghmous, Michael Steinbach, Arindam Banerjee, Auroop Ganguly, Shashi Shekhar, Nagiza Samatova, and Vipin Kumar. 2017. Theory-guided data science: A new paradigm for scientific discovery from data. *IEEE Transactions on knowledge and data engineering* 29, 10 (2017), 2318–2331.
- [22] Anuj Karpatne, Ramakrishnan Kannan, and Vipin Kumar. 2022. *Knowledge Guided Machine Learning: Accelerating Discovery using Scientific Knowledge and Data*. CRC Press.
- [23] Tero Karras, Miika Aittala, Samuli Laine, Erik Härkönen, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2021. Alias-free generative adversarial networks. *Advances in Neural Information Processing Systems* 34 (2021), 852–863.
- [24] Tero Karras, Samuli Laine, and Timo Aila. 2019. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 4401–4410.
- [25] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. 2020. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 8110–8119.
- [26] Diederik P Kingma and Max Welling. 2013. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013).
- [27] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the Loss Landscape of Neural Nets. In *Advances in Neural Information Processing Systems*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Eds.), Vol. 31. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2018/file/a41b3bb3e6b050b6c9067c67f663b915-Paper.pdf>
- [28] Xiao Li, Chenghua Lin, Ruizhe Li, Chaozheng Wang, and Frank Guerin. 2020. Latent space factorisation and manipulation via matrix subspace projection. In *International Conference on Machine Learning*. PMLR, 5916–5926.
- [29] Moritz D Lürig, Seth Donoughe, Erik I Svensson, Arthur Porto, and Masahito Tsuboi. 2021. Computer vision, machine learning, and the promise of phenomics in ecology and evolutionary biology. *Frontiers in Ecology and Evolution* 9 (2021), 642774.
- [30] Michael Lynch. 1991. Methods for the analysis of comparative data in evolutionary biology. *Evolution* 45, 5 (1991), 1065–1080.
- [31] ML Menéndez, JA Pardo, L Pardo, and MC Pardo. 1997. The jensen-shannon divergence. *Journal of the Franklin Institute* 334, 2 (1997), 307–318.
- [32] Meike Nauta, Ron van Bree, and Christin Seifert. 2021. Neural prototype trees for interpretable fine-grained image recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14933–14943.
- [33] NSF HDR Imageomics Institute. 2021. Imageomics: A new frontier of biological information powered by knowledge-guided machine learning. <https://imageomics.osu.edu/>.
- [34] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. 2017. Neural Discrete Representation Learning. <https://doi.org/10.48550/ARXIV.1711.00937>
- [35] Stanislav Pidhorskyi, Donald A Adjeroh, and Gianfranco Doretto. 2020. Adversarial latent autoencoders. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14104–14113.
- [36] Samantha A Price, Sarah T Friedman, Katherine A Corn, Olivier Larouche, Kasey Brockelsby, Anna J Lee, Maya Nagaraj, Nick G Bertrand, Maile Danao, Megan C Coyne, et al. 2022. FishShapes v1: Functionally relevant measurements of teleost shape and size on three dimensions.
- [37] Mengshi Qi, Yunhong Wang, Jie Qin, and Annan Li. 2019. Ke-gan: Knowledge embedded generative adversarial networks for semi-supervised scene parsing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5237–5246.
- [38] Daniel L Rabosky, Jonathan Chang, Peter F Cowman, Lauren Sallan, Matt Friedman, Kristin Kaschner, Cristina Garilao, Thomas J Near, Marta Coll, Michael E Alfaro, et al. 2018. An inverse latitudinal gradient in speciation rate for marine fishes. *Nature* 559, 7714 (2018), 392–395.
- [39] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. 2019. Language models are unsupervised multitask learners. *OpenAI blog* 1, 8 (2019), 9.
- [40] Maziar Raissi, Paris Perdikaris, and George E Karniadakis. 2019. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational physics* 378 (2019), 686–707.
- [41] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. 2021. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. 2287–2296.
- [42] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.
- [43] Tiago R Simões, Michael W Caldwell, Alessandro Palci, and Randall L Nydam. 2017. Giant taxon-character matrices: quality of character constructions remains critical regardless of size. *Cladistics* 33, 2 (2017), 198–219.
- [44] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2013. Deep inside convolutional networks: Visualising image classification models and saliency maps. *arXiv preprint arXiv:1312.6034* (2013).
- [45] Randal A Singer, Kevin J Love, and Lawrence M Page. 2018. A survey of digitized data from US fish collections in the iDigBio data aggregator. *PLoS one* 13, 12 (2018), e0207636.
- [46] Aaron Van den Oord, Nal Kalchbrenner, Lasse Espeholt, Oriol Vinyals, Alex Graves, et al. 2016. Conditional image generation with pixelcnn decoders. *Advances in neural information processing systems* 29 (2016).
- [47] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9, 86 (2008), 2579–2605. <http://jmlr.org/papers/v9/vandermaaten08a.html>
- [48] Grant Van Horn, Oisin Mac Aodha, Yang Song, Yin Cui, Chen Sun, Alex Shepard, Hartwig Adam, Pietro Perona, and Serge Belongie. 2018. The inaturalist species classification and detection dataset. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 8769–8778.
- [49] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. 2011. The caltech-ucsd birds-200-2011 dataset. (2011).
- [50] Jiayun Wang, Yubei Chen, Rudrasis Chakraborty, and Stella X. Yu. 2020. Orthogonal Convolutional Neural Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [51] Rui Wang, Robin Walters, and Rose Yu. 2020. Incorporating symmetry into deep dynamics models for improved generalization. *arXiv preprint arXiv:2002.03061* (2020).

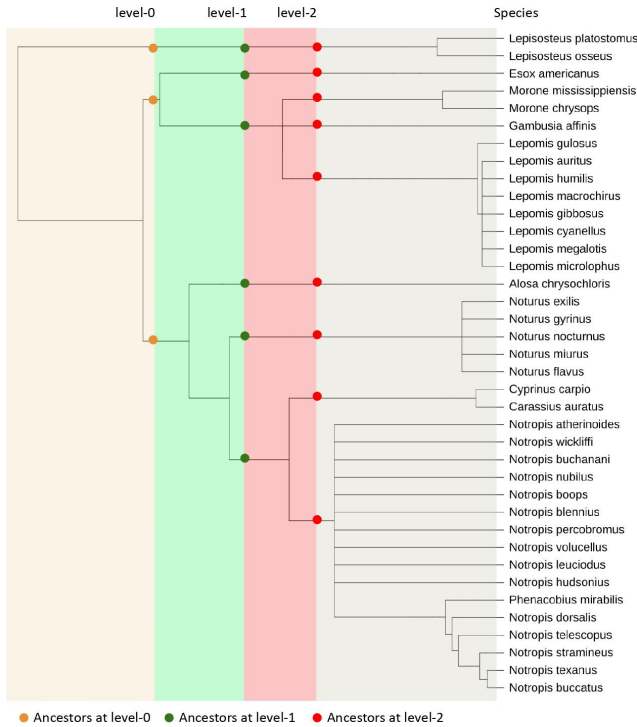


Figure 7: Ancestors at different levels of phylogeny

A DATASET

As mentioned in Section 5.1, the images we use in this work come from a 38-species semi-balanced subset of a larger collection that participated in the Great Lakes Invasives Network Project (GLIN). After further splitting this subset into training and test sets, we apply some pre-processing that is necessary for the neural network to yield best results. This pre-processing includes padding the images with the ImageNet mean RGB color. We also use data augmentation when training the base VQGAN model. This includes random horizontal flips, spatial shifts and rotations, and brightness and contrast changes.

B PHYLOGENY PREPROCESSING

As mentioned in Section 2, we use a phylogenetic tree to characterize the evolutionary distances between the species in our dataset. The phylogeny corresponding to the dataset described in Section 5.1 was obtained using *opentree* (<https://opentree.readthedocs.io/en/latest/>) python package. Phylogeny processing and manipulation were done using *ete3* (<http://etetoolkit.org/>) python package.

In our application, we quantize our 38-species tree into $m_1 = 4$ distinct phylogenetic levels. Each level groups the 38 species based on their common ancestry within that level. Figure 7 outlines each level with its corresponding species groupings.

C HYPER-PARAMETER SELECTION

In terms of hyper-parameter tuning, we used the following settings for each of the trained models:

Vanilla VQGAN: We trained a VQGAN with a codebook of 1024 possible codes and an embedding sequence of 256 codes. We trained

the model for 836 epochs with a learning rate of 4.5×10^{-6} . We used this VQGAN as the base model for the rest of the models. A batch size of 32 was used.

Phylo-NN Taking the base VQGAN model described above, we trained a *Phylo-NN* that has z_p^Q of dimensions ($n_l = 4, n_p = 8$). z_{np}^Q also has the same dimensionality. The dimensionality of the embedding itself is $d = 16$, and the size of the codebook $n_q = 64$. A batch size of 32 was used.

Concept Whitening (CW) Taking the base VQGAN model described above, we trained CW for 20 epochs using the same hyper-parameters as Vanilla VQGAN. We used a batch size of 20 for the concepts.

Latent Space Factorization (LSF): With the base VQGAN model described above, a variational autoencoder was introduced between the base encoder and the quantization layer. The model was trained for 200 epochs with a learning rate of 1×10^{-4} . We used an embedding dimension of 1024 for the variational autoencoder.

D DETAILS OF MORPHOLOGICAL DISTANCE PROCESSING

The 8 functionally relevant traits that we used from the FishShapes dataset include: standard length, maximum body depth, maximum fish width, lower jaw length, mouth width, head depth, minimum caudal peduncle depth, and minimum caudal peduncle. Some species were not available in the FishShapes dataset, so when possible, the closest relative was substituted. (*Notropis percobromus* was replaced with *Notropis rubellus*, and *Carassius auratus* was replaced with *Carassius carassius*). Also, two species of *Lepisosteus* had no close relatives and were thus removed from the Spearman correlation analysis. To correct for overall size and allometry, each measurement was log transformed and regressed against Standard Length (SL) using a phylogenetic regression in the R package *phylolm*, with the residuals from the regression being the inputs into PACA. Distances in the principal components of PACA were measured as the Mahalanobis distance between the multivariate means using a covariance matrix proportional to the evolutionary rate matrix from the multivariate Brownian Motion fit in the R package *mvMORPH* [6].

E T-SNE PLOTS FOR TEST IMAGES

As we have shown in Section 6.4, the embedding of the images generated by our *Phylo-NN* model are more meaningful than those generated by a vanilla VQGAN. In this section, we run the same analysis on the test images. Clearly, a similar case can be made here, namely that the encoding of the images using *Phylo-NN* is superior to other models' in terms of its clustering. Figure 8 illustrates this point clearly when comparing *Phylo-NN*'s (top row) t-SNE plots with those of the other models'. Both vanilla VQGAN and CW perform worse at clustering the dataset compared to *Phylo-NN*.

F EXAMPLES OF PHYLO HISTOGRAMS

In Section 6.3, by means of calculating the average JS-divergence of sequence histograms, we investigated how well the Imageome sequences match for species that share a common ancestor, as opposed to those that don't. In this section, we show some of these



Figure 9: *Notropis nubilus*

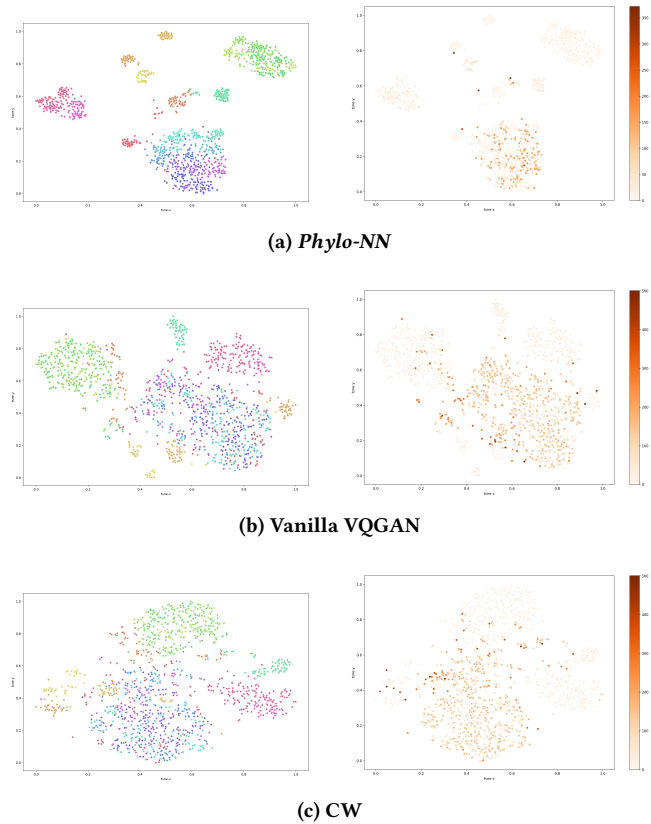


Figure 8: t-SNE plots of the test set images using different models

histogram plots to illustrate their value and the insight they could provide.

In figures 9 to 11, each histogram represents a code location in the phylogenetic sequence. Each column represents one of the $n_l = 4$ phylogenetic levels into which the phylogeny was quantized, starting with the species level from right and climbing the phylogeny all the way to the earliest ancestral level on the left. Each column has $n_p = 8$ codes. Each histogram shows the relative frequency of each code of the learned $n_q = 64$ codes for its corresponding sequence location. The lower a histogram's entropy (i.e., when there is only one or a couple of codes that dominate the histogram's frequency spectrum), the more important that code location hypothetically is for characterizing the species at its corresponding phylogenetic level.

As can be seen, both *Notropis* species share many codes at many sequence locations up to, but not including, the species level. This is because these species share an immediate ancestor. In contrast, we can see that the *Lepomis* species has a distinct histogram signature and does not share almost no codes with the *Notropis* species, except for the earlier ancestral level (i.e., the left column).

Discovering Novel Biological Traits From Images Using Phylogeny-Guided Neural Networks

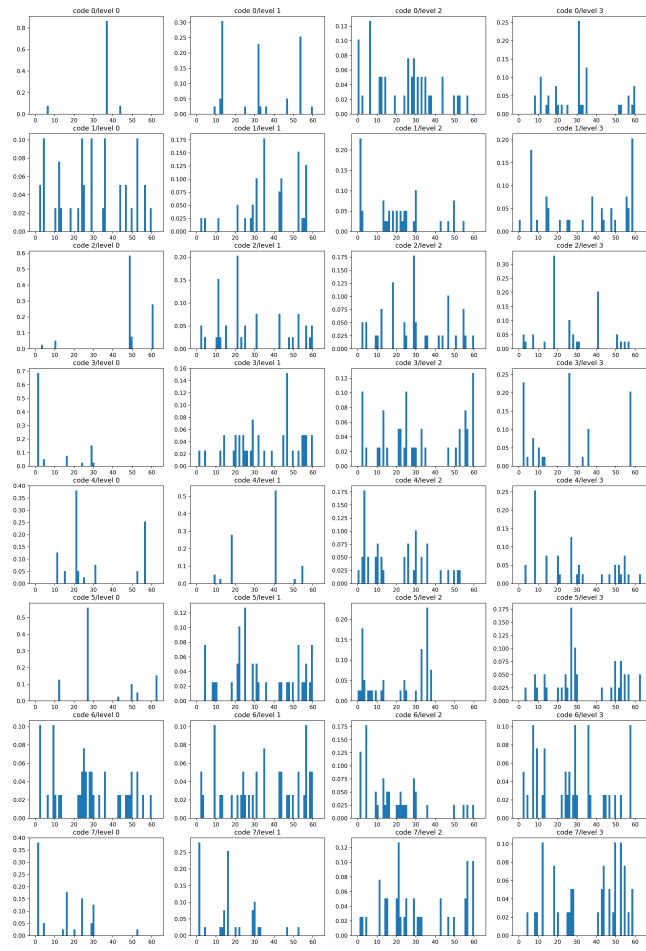


Figure 10: *Notropis percobromus*



Figure 11: *Lepomis macrochirus*