

# GOPEN ACCESS

**Citation**: Parsons DJ, Green AE, Carstens BC, Pelletier TA (2024) Predicting genetic biodiversity in salamanders using geographic, climatic, and life history traits. PLoS ONE 19(10): e0310932. <a href="https://doi.org/10.1371/journal.pone.0310932">https://doi.org/10.1371/journal.pone.0310932</a>

**Editor:** Christopher Nice, Texas State University, UNITED STATES OF AMERICA

Received: February 22, 2024 Accepted: September 9, 2024 Published: October 18, 2024

Copyright: © 2024 Parsons et al. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Data are uploaded to Dryad and can be accessed using the following link: https://doi.org/10.5061/dryad.m63xsj474.

Code related to this manuscript, including data cleaning, imputation, predictive modeling, and significance testing has been deposited in GitHub (https://github.com/parsons463/HiddenSalamanders). All remaining data are available in the manuscript and/or Supporting information.

RESEARCH ARTICLE

# Predicting genetic biodiversity in salamanders using geographic, climatic, and life history traits

Danielle J. Parsons<sup>1,2</sup>, Abigail E. Green<sup>3</sup>, Bryan C. Carstens<sup>1,2</sup>, Tara A. Pelletier<sup>3</sup>

- Museum of Biological Diversity, The Ohio State University, Columbus, Ohio, United States of America,
   Department of Evolution, Ecology and Organismal Biology, The Ohio State University, Columbus, Ohio, United States of America,
   Department of Biology, Radford University, Radford, Virginia, United States of America
- \* carstens.12@osu.edu

# **Abstract**

The geographic distribution of genetic variation within a species reveals information about its evolutionary history, including responses to historical climate change and dispersal ability across various habitat types. We combine genetic data from salamander species with geographic, climatic, and life history data collected from open-source online repositories to develop a machine learning model designed to identify the traits that are most predictive of unrecognized genetic lineages. We find evidence of hidden diversity distributed throughout the clade Caudata that is largely the result of variation in climatic variables. We highlight some of the difficulties in using machine-learning models on open-source data that are often messy and potentially taxonomically and geographically biased.

### Introduction

Documenting biodiversity is an important first step in understanding both ecological and evolutionary processes [1], particularly the functional roles that act to connect processes functioning at both shallow and deep time scales [2]. Notably, any such documentation of biodiversity implicitly assumes that the units (e.g., species) are comparable across different geographic regions. Given that a Linnean shortfall (i.e., the ratio of recognized to unrecognized species [3]) exists in most clades and may be substantial across Eukaryota [4], it is not clear that this assumption is reasonable. An alternative approach is to utilize evolutionary significant units [5], or genetic lineages, in place of species in broad analyses of biodiversity (e.g., [6]). This may be particularly useful in clades with relatively high degrees of morphological and ecological conservatism. One such clade is Caudata (i.e., salamanders and newts), which exhibits high frequencies of cryptic species (e.g., [7–9]).

Identifying hidden genetic lineages in Caudata can have important conservation implications. For example, Mead *et al.* [10] discovered a new species of western *Plethodon* salamander that was originally thought to be either *P. elongatus* or *P. stormi* [10]. All three of these species are listed on the IUCN Red List as either near threatened (*P. elongatus*), vulnerable (*P. asupak*),

**Funding:** This work was supported by the Directorate for Biological Sciences (DEB-1910623 to BCC: DEB-1911293 to TAP).

**Competing interests:** The authors have declared that no competing interests exist.

or endangered (*P. stormi*). More recently, Parra Olea *et al.* [11] discovered five cryptic lineages in *Chiropterotriton* from Mexico, several of which are threatened due to their restricted ranges [11]. Species with small ranges and/or limited dispersal capabilities can be harder to protect because their distributions often do not fall within protected areas [12] and small ranges are often used as a factor in assigning conservation priorities [13]. Therefore, it is important to identify these hidden lineages, as they could easily go unnoticed and unprotected. Many other species of salamander that would have otherwise gone unnoticed and have been recognized using molecular data have small ranges and likely need protection [14–18]. The presence of cryptic diversity has been recently highlighted as a key component of undescribed biodiversity that requires greater attention [19,20].

Efforts to conserve undescribed genetic diversity can be facilitated using computational methods that identify genetic lineages representing potentially hidden diversity in need of further investigation. The use of data science techniques has allowed biodiversity studies to expand their geographic and taxonomic focus to explore broader patterns of evolution, which can be difficult to assess using traditional meta-analysis methods [21]. Macrogenetics, a relatively new field that merges biodiversity data with genetic data [22,23], has been used to explore how human impacts influence levels of intraspecific genetic diversity [24,25], to study past and future climate refugia [26,27], and to quantify latitudinal biodiversity gradients [28– 31]. Macrogenetic methods, particularly in combination with predictive modeling, can be used to inform conservation policies by identifying species, taxonomic groups, or geographic areas in need of further investigation [32,33]. Machine learning models excel in processing and analyzing large datasets automatically, making them highly efficient at identifying complex, non-linear patterns within extensive data. While these models are powerful tools for making predictions with high accuracy, the process behind many machine learning models can be difficult to interpret [34]. Random forest models, a type of predictive machine learning model, enhance interpretability by providing clear insights into feature importance, thus representing a powerful tool for identifying factors influencing biodiversity. Recently, such analyses have expanded to taxonomic work.

Parsons *et al.* [35] analyzed mitochondrial DNA sequences from over 4000 species of mammals, representing roughly 66% of currently described species, and found that mammal diversity is largely under-described using molecular species delimitation methods on publicly available barcode data. This is useful for several reasons. A comprehensive list of undescribed genetic lineages that may represent species now exists that can help focus taxonomic efforts. Parsons *et al.* [35] also found that taxa with small bodies, and large geographic distributions with variation in precipitation and isothermality, were more likely to contain cryptic diversity. While some of this might seem obvious (morphological differences are harder to observe in small-bodied animals and these animals may be harder to find), it does allow researchers to document characteristics of species, higher taxonomic groups, or even geographic regions that contribute to diversification and therefore biodiversity patterns. When done in disparate taxonomic groups (e.g., vertebrates, invertebrates, plants, and fungi) and at different levels (e.g., Class, Order, Family) this furthers our understanding of core evolutionary processes.

A similar approach was taken in birds. Using a tree-based molecular species delimitation method, Smith *et al.* [36] found that latitude explained variation in phylogeographic breaks, while other traits pertaining to habitat and life history explained very little. In this case, phylogeographic structure was higher in the tropics. Conversely, in other organisms, isolation-by-distance within species is often higher at higher latitudes (multiple taxonomic groups: [32]; amphibians: [37]). Further, genetic variation within amphibians was best explained by range size and elevation, rather than latitude, in the neotropics [37], while latitude was an important predictor of genetic diversity in the nearctic [30]. This suggests that differences exist in how

genetic variation is distributed within species depending on which taxonomic groups are being examined, and at what spatial scale.

In order to expand these approaches, we conducted an assessment of genetic lineages in roughly 100 described salamander species using the phylogatR database [38]. PhylogatR aggregates DNA sequence data from both GenBank and BOLD into sequence alignments, providing associated GBIF occurrence records (i.e., GPS coordinates) for each sequence. There are over 700 nominal species of salamanders belonging to nine families [39], most located in the northern hemisphere. While salamanders contain a wide variety of life history strategies and habitats, they are likely to have high levels of cryptic diversity due to their moisture requirements and similar body forms. However, their eco-evolutionary processes can vary from species to species and sometimes oppose our expectations [40-45]. We follow methods from Parsons et al. [35] and use molecular species delimitation methods to estimate the number of genetic lineages present in previously collected data that is both openly available and easily tractable. We utilize these delimitation results to identify species that are likely to harbor undescribed diversity. Species for which delimitation reveals multiple genetic lineages are classified as hidden species. The individual genetic lineages that comprise these hidden species are referred to as hidden genetic lineages. We then use a random forest classification to determine whether any variables pertaining to geography, the environment, or life history traits contribute to the presence of hidden genetic lineages within species. We also discuss some of the difficulties in using open-source data that are often messy and potentially taxonomically and geographically biased.

### Materials and methods

# Collection of genetic and geographic data

We downloaded all available data from the *phylogatR* database (<a href="https://phylogatr.org/">https://phylogatr.org/</a>) using the search term 'Caudata' on 2/4/22. The uncleaned data represented four families, 93 different species, and 14 loci with a total of 3768 DNA sequences. To begin cleaning the data, we calculated nucleotide diversity (pi) values for each locus in every species and found outliers by setting lower and upper bounds of 2.5% (0) and 97.5% (0.2193634) respectively. For each of the four outliers and two species with missing pi values, we opened the DNA sequence file in Mesquite v3.7 [46] and removed any extremely short or non-overlapping sequences (S1 Data). Additionally, we discovered a typo for the species *Batrachuperus karlschmidti* causing there to be two different species folders for the same species. Both the sequence and occurrence files were merged for the species and the sequence files were realigned to correct the error. Two species complexes were present in the dataset, and these were kept named as downloaded: *Triturus cristatus x dobrogicus macrosomus* and *Ambystoma laterale jeffersonianum* complex. A review of the available loci indicated that two genes, *COI* and *cytb*, were the most well-represented in both total number of sequences and species coverage. Consequently, we opted to utilize these two genes for downstream analysis.

Species alignments from the download for both the mitochondrial genes Cytochrome oxidase I (COI) and Cytochrome b (cytb) were merged for all salamander species and aligned using MAFFT v7.5 [47] with the default settings and including the–adjustdirection command to account for reverse complement sequences. We visually inspected alignment files for both genes and removed all short sequences, which we classified as those missing 50% or more of the second half of the sequence. Twenty-one sequences were removed from the COI alignment and 99 were removed from the cytb alignment, leaving totals of 768 and 908 sequences for COI and cytb, respectively. The sequences for seven species were completely removed from further analysis due to their short length (missing 50% or more of the second half of the sequences). In

total, eighty-three species remained with an average of approximately 20 sequences per nominal species (see <u>S2 Data</u> for a list of identifiers corresponding to the sequences used in this study).

# Species delimitation

We used three methods of species delimitation to determine the number of genetic lineages present in our samples. The GMYC is a tree-based method that takes a phylogenetic tree as input and finds a point in the tree where branching changes from within to between species [48]. The ABGD [49] and ASAP [50] methods are distance-based delimitation methods that use pairwise genetic distances to establish the threshold between intra- and inter-species divergence. Because each method is based on a specific set of assumptions, it is best to use multiple methods and compare their results in order to achieve a more accurate delimitation [51]. By looking for concordance across methods, we can increase our confidence in the identified lineage boundaries and minimize the potential impact of bias introduced by any single method. While we report delimitation results from the genes *COI* and *cytb* for all methods, we used a consensus of these results—reflecting agreement among the GMYC, ABGD, and ASAP delimitation methods for both *COI* and *cytb*—for assessing the influence of geography, environment, and life history traits on predicting salamander genetic diversity.

To estimate a species tree for input into the GMYC, we used BEAST v2.5.1 [52]. We used the default parameters except for conducting 100,000,000 million generations, sampling every 5,000, and setting the model of sequence evolution to GTR+I+G [53]. The log files were checked by eye using Tracer v1.7.2 [54]. ESS values were all over 1000 for both *cytb* and *COI*. We removed 10% as burnin and retained the maximum clade credibility tree using TreeAnnotator. After checking that the tree was binary and ultrametric, we used the R package *splits* [55] to conduct GMYC analyses. In each case we used the single threshold model and all other default settings. We conducted both ABGD and ASAP delimitation analyses via their web portals (<a href="https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html">https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html</a> and <a href="https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html">https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html</a> and <a href="https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html">https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html</a> and <a href="https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html">https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html</a> and <a href="https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html">https://bioinfo.mnhn.fr/abi/public/abgd/abgdweb.html</a> are settings.

### **Predictor variables**

For each nominal salamander species, we examined numerous geographic, environmental, morphological, and life history variables to identify traits predictive of undescribed salamander diversity using a classification model based on our species delimitation results. A variety of predictor variables were collected, including geographic and environmental values derived from georeferenced locality data (see <a href="S3 Data">S3 Data</a>). In addition, three life history traits were available from AmphiBIO, a global database for amphibian ecological traits <a href="[56]">[56]</a>, for most of the species in our study: reproductive strategy (direct developing, larval phase), habitat (terrestrial, fossorial, aquatic, or some combination of these), and body size (total length). To supplement this dataset and fill in any missing trait values, we used AmphibiaWeb <a href="[57]">[57]</a> and other online sources (<a href="S4 Data">S4 Data</a>).

To extract species specific data related to its environmental distribution, we utilized 42 GIS data layers (see <u>S4 Data</u> for data layer details), meant to capture various aspects of the ecology and habitat of each species. These include all 19 BIOCLIM layers from the CHELSA database [58,59] at 1 km resolution, elevation [60], terrestrial habitat heterogeneity [61], global land cover classification [62], global river classification [63], disaster risk [64], and various indicators of seasonal growth 58–59]. In addition to traits meant to capture various ecological factors we also gathered data for several traits relating to human impact, in order to measure levels of human disturbance and activity to the species environment. While ecological factors directly

affect levels of biodiversity by influencing species biology, anthropogenic factors can influence the way we find and describe this diversity (e.g., increased sampling in wealthier, more populated locations). We extracted species specific data from several GIS layers, including anthropogenic biome [65], human population density [66], and gross domestic product [67] in order to evaluate how anthropogenic factors impact undescribed diversity.

We utilized the R packages 'raster' [68], 'rgdal' [69], 'geosphere' [70], and 'plyr' [71] to extract species specific information from each layer using geographic occurrence records obtained from *phylogatR*. To represent the environmental variation within the occupied range of each species, we extracted the value of each environmental layer for each GPS coordinate associated with each species. We then took the mean and standard deviation for each environmental variable. To obtain species specific data related to geographic distribution we extracted the minimum, maximum, mean, and length of latitude and longitude from the GPS points of each species.

We used the R package 'mice' [72] to impute trait values missing from our dataset (see S1 Fig in S1 File for distribution of missing data and specific trait values imputed). The imputation method 'pmm' was used for all numeric variables and 'polyreg' was used for categorical variables (i.e., reproductive strategy and habitat). We ran the imputation 15 times (S2 Fig in S1 File) and then pooled the iterations to generate the final imputed values. The final database containing all trait values (both imputed and original) is available in S4 Data.

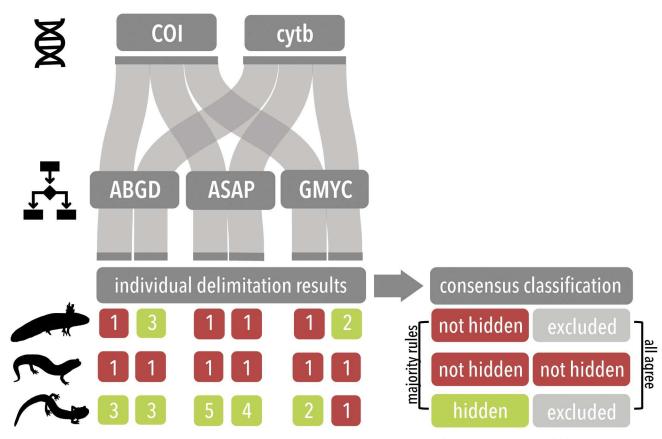
# **Predictive modeling**

We used the R package 'caret' [73] to generate a random forest classification model [34] based on our previously generated database of predictor variables and a consensus of our species delimitation results. Two separate sets of consensus models were generated to assess the role of geography, environment, and life history traits on the presence of hidden diversity (Fig 1A). The first model (all agree) represents a strict consensus of delimitation results from species in which results from all methods of species delimitation agree (Fig 1B). Any species with conflicting delimitation results were excluded from analysis. The second model (majority rules) represents a majority rule consensus in which species are assigned to a response category based on relative support of delimitation results (Fig 1C). For each model, we used 70% of the data to train the model and the remaining 30% was set aside as a test set. Models were generated using 10-fold cross validation with five repeats to tune the parameter 'mtry', the number of variables randomly sampled at each split, and optimize the area under the receiver operating characteristic curve, ROC. After training, we extracted the variable importance measures mean decrease accuracy (MDA) and Gini impurity (Gini) from the final models. We then used the final models on the test set data to evaluate model performance. Model performance was evaluated across a variety of metrics including model accuracy, which reflects how well the predicted classifications agree with the observed classifications, and both positive and negative predictive value, which indicate the how the model performs on observations from each class. Additionally, we calculated the no information rate (NIR), the proportion of observations that fall into the majority class, and the p-value [Accuracy>NIR], to test for model significance. The top important predictor variables from our best model were compared using a Kruskal-Wallis test to determine if these variables are significantly different between species that do or do not contain hidden diversity.

### Results

# Genetic and geographic dataset

Our final dataset consisted of 1676 DNA barcoding sequences (Fig 2). Of these, 768 sequences were from the Cytochrome oxidase I gene (*COI*), and 908 sequences were from the

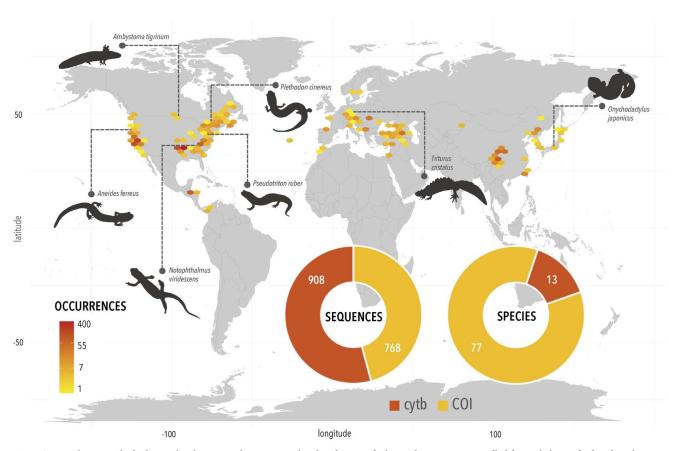


**Fig 1. Consensus classification of species delimitation results.** A, Flowchart describing the process of generating a consensus of delimitation results (among different methods and loci). B, C, Pipeline for classifying nominal species as either containing or not containing hidden genetic diversity in each consensus analysis (*all agree* and *majority rules*, respectively).

Cytochrome b gene (*cytb*). These sequences were derived from 83 nominal species of salamanders, which were distributed among 26 distinct genera occurring across the globe. The dataset contained 13 species with sequences from the gene *cytb*. Comparatively, *COI* exhibited notably broader taxonomic coverage, with 77 nominal species represented. Out of the 83 species analyzed, only seven were shared between *COI* and *cytb*. Of the remaining 76 species, 70 were unique to *COI* and six were unique to *cytb*. To supplement the genetic data collected, a total of 1676 georeferenced occurrence records from *phylogatR* were utilized to collect a combination of geographic, environmental, and life history trait values for each nominal species present in the dataset.

### Species delimitation and consensus assignment

Species delimitation results were generated by analyzing *COI* and *cytb* sequences from each nominal species under three different delimitation methods, ABGD, ASAP, and GMYC. We classified each nominal species as either containing undescribed genetic lineages or not containing undescribed genetic lineages based on the number of genetic groups predicted by each delimitation analysis. While taxonomic overlap between *COI* and *cytb* was narrow, delimitation results for species shared by both loci were mostly congruent with respect to species classification. Of the seven species with sequences from both genes, only two species produced conflicting results regarding the presence of undescribed genetic lineages within a specific



**Fig 2. Geographic spread of salamander data.** Map shows geographic distribution of salamander occurrences pulled from *phylogatR* [38] and used in these analyses. Pie charts show the total number of *cytb* and *COI* sequences used (left) and the number of species represented by those *cytb* and *COI* sequences (right). Basemap created with world map data from the public domain Natural Earth project (<a href="http://www.naturalearthdata.com">http://www.naturalearthdata.com</a>). Salamander figures in black were obtained from Phylopic [74] and are licensed under public domain.

taxon based on loci. Delimitation results across different methods showed slightly less agreement. Classifications resulting from the GMYC and ASAP methods were similar across species. These methods, on average, resulted in slightly fewer predicted species per nominal species than the ABGD method (see Fig 3 for predicted species numbers).

To account for this variation in our final random forest classification models, we generated two consensus classifications to evaluate concordance between delimitation results from different methods and loci. The results of our consensus models indicate that roughly 2/3rds of the nominal salamander species used in this analysis are likely to contain genetic lineages that may be unexplored diversity. The strictest of these classifications produced a consensus model (*all agree*) consisting of 51 total species, 41 of which were classified as containing hidden diversity and 10 of which were classified as not containing hidden diversity. The remaining consensus model (*majority rules*) consisted of 83 total species, of which 51 were classified as containing undescribed genetic lineages and 32 were not (Fig 3).

### **Predictive modeling**

For our *majority rules* and *all agree* consensus classifications, we developed random forest classification models using all available predictor data. To assess potential correlation between variables in our dataset we used the R package 'corrplot' [75] to generate a correlation matrix of

		COI		cytb					COI			cytb		
Species	GMYC	ABGD	ASAP	GMYC	ABGD	ASAP	Species (cont.)	GMYC	ABGD	ASAP	GMYC	ABGD	ASA	
Ambystoma annulatum	1	1	1	-	-	-	Lissotriton boscai	1	1	1		,	-	
Ambystoma californiense	1	2	1	-	-	-	Lissotriton helveticus	1	1	1			-	
Ambystoma laterale	1	2	1	-	-	-	Lissotriton montandoni	1	3	1	-	-	-	
Ambystoma laterale jeffersonianum complex	6	9	4	-	-	-	Lissotriton vulgaris	2	6	2	-	-	-	
Ambystoma opacum	1	2	1	-	-	-	Mertensiella caucasica	4	13	3	-	-	-	
Ambystoma talpoideum	1	3	1_	2	-	-	Neurergus crocatus	1	3	1	-	3	-	
Ambystoma texanum	2	3	2	-	-	-	Notophthalmus viridescens	2	3	1	-	-	-	
Ambystoma tigrinum	1	2	1	-	-	-	Nototriton mime	1	3	1	-	- 4	-	
Aneides ferreus	5	7	4	-	-	-	Ommatotriton nesterovi	3	16	2	-	-	-	
Aneides flavipunctatus	7	7	5	-	-	-	Ommatotriton ophryticus	5	11	5	-	-	-	
Aneides lugubris	4	5	2	-	-	-	Ommatotriton vittatus	4	14	4	-	-	-	
Aneides vagrans	3	3	1	-	-	-	Onychodactylus japonicus	2	3	2	-	-	-	
Batrachoseps attenuatus	1	4	1	3	10	9	Plethodon cinereus	2	9	2		-		
Batrachoseps major	1	3	1	-	-	-	Plethodon fourchensis	-	-	-	3	6	9	
Batrachuperus karlschmidti	4	8	2	4	8	10	Plethodon glutinosus	1	3	1	-	-	-	
Batrachuperus londongensis	1	2	1	2	6	4	Plethodon hubrichti	1	2	1		-		
Batrachuperus pinchonii	4	6	3	5	13	11	Plethodon montanus	3	6	3	-	-	-	
Batrachuperus taibaiensis	5	10	3	5	6	9	Plethodon ouachitae	-	-	-	13	27	42	
Batrachuperus tibetanus	4	9	3	7	21	19	Plethodon richmondi	1	1	1	-		-	
Batrachuperus yenyuanensis	-	-	-	3	4	5	Plethodon serratus	2	3	2	3	10	11	
Bolitoglossa medemi	2	2	2	-		-	Plethodon sherando	1	4	1	-	-	-	
Bolitoglossa porrasorum	3	24	2				Plethodon shermani	-	-		3	3	3	
Bolitoglossa rufescens	3	4	2	-	-	-	Plethodon vehiculum	1	1	1	-	-	-	
Bolitoglossa taylori	2	5	2			-	Plethodon wehrlei	2	4	2			-	
Desmognathus fuscus	5	8	2		_	-	Pleurodeles waltl	2	2	2			-	
Desmognathus monticola	3	7	3	_	-	_	Pseudotriton ruber	2	4	2	_		-	
Desmognathus ochrophaeus	2	3	2	-	-	-	Ranodon sibiricus	1	1	1	-	-	-	
Desmognathus orestes	1	4	1	-		-	Salamandra salamandra	2	8	1			-	
Desmognathus organi	1	1	1		-	-	Salamandrella keyserlingii	2	2	1	_	-	-	
Desmognathus quadramaculatus	2	4	2	-	-	-	Salamandrella schrenckii	3	8	3	-	-	-	
Dicamptodon ensatus	1	1	1	-	-	-	Triturus carnifex	1	3	2			-	
Ensatina eschscholtzii		-		38	107	92	Triturus cristatus	2	5	1	_	-	-	
Eurycea bislineata	1	7	1	30	- 107	-	28 8 10 2 10 00000 10	1	2	1	-	5	-	
	3	9	3	-		-	Triturus cristatus x dobrogicus macrosomus	1	1	1	-	-	-	
Eurycea cirrigera	-		1	-	-	-	Triturus dobrogicus	1	-	1	-		⊢	
Eurycea guttolineata	2	3		- 1	- 1	- 1	Triturus karelinii		5	l l	•	-	-	
Eurycea subfluvicola	-	-	-	1	de	1	В							
Eurycea wilderae	1	3	1	-	-	-	all agree		ma	aiorit	y rule	c		
Gyrinophilus porphyriticus	3	4	1	-	-	-	a ag. cc		1110	ajoint,	y ruic	.3		
Hemidactylium scutatum	3	4	3	-	-	-	■ NA ■ NOT HIDDEN ■ HIDDEN	= NA	N ■ N	OT HID	DEN 📕	HIDD	EN	
Hynobius amjiensis	1	2	1	-	-	-								
Hynobius arisanensis	1	4	1	-	-	-								
Hynobius formosanus	1	3	1	-	-	-								
Hynobius fuca	3	3	2	-	-	-	32				32	2		
Hynobius leechii	2	5	2	-	-	-	41			(				
Hynobius retardatus	1	2	1	-	-	-			51					
Hynobius sonani		3	1	-	-	- 1	10		/		<b>Y</b>			
Hynobius tsuensis	2	2	2		<del>                                     </del>	-				\				

**Fig 3. Species delimitation results.** A, Graphs show the results of ABGD, ASAP, and GMYC species delimitation analyses of the genes *cytb* and *COI* for each nominal species. Numbers represent the predicted genetic lineages from each analysis. Results highlighted in red indicate no hidden genetic lineages were predicted (i.e., number of genetic lineages = 1). Results highlighted in green indicate hidden genetic lineages were predicted (i.e., number of genetic lineages > 1). Grey highlighting indicates that specific analysis was not performed due to a lack of data. B, Pie charts display the number of nominal species classified as either containing or not containing hidden diversity in each consensus analysis (i.e., *all agree* and *majority rules*).

<b>Majority Rules Models</b>	Original	Correlation  > 0.75	Correlation  > 0.85	Correlation  > 0.90
Accuracy	0.75	0.75	0.75	0.8333
Accuracy (95% CI)	(0.5329, 0.9023)	(0.5329, 0.9023)	(0.5329, 0.9023)	(0.6262, 0.9526)
No Information Rate	0.625	0.625	0.625	0.625
Pos Pred Value	0.7368a	0.8	0.7647	0.7895
Neg Pred Value	0.8	0.6667	0.7143	1
P-Value [Acc > NIR]	0.1453	0.1453	0.1453	0.02435

**Table 1. Results of** *majority rules* **consensus random forest models.** Model metrics for each random forest predictive model generated using the *majority rules* consensus classifications are shown.

our predictor variables (S3 Fig in S1 File). Due to the presence of strong correlations between several of the geographic and environmental variables in our dataset we performed multiple random forest models with progressive sets of correlated variables removed at different cutoff values (i.e., |correlation coefficient| > 0.75; 0.85; 0.9). The results of these random forest models are presented below (Table 1).

All random forest models were found to have high predictive accuracy, with the *majority* rules and all agree models achieving accuracies of 75-85% and 87-93%, respectively, in identifying nominal species likely to contain hidden diversity. Although these results may initially seem to suggest that all our models are able to make meaningful predictions, further examination of additional model evaluation metrics reveals potential overfitting and inflation of predictive power. For example, despite the high accuracy of the models, the 95% confidence intervals for these values are broad with an average length of nearly 40% for most of the models (Tables 1 and 2). Additionally, the no information rates (NIRs), a measure of prediction significance based on the underlying dataset that needs to be exceeded in order for model results to be significant, are particularly high for the all agree consensus models, where the class frequencies are more skewed towards species predicted to harbor hidden diversity. The high NIR values combined with wide confidence intervals result in a p-value [Accuracy > NIR] greater than 0.05 in all models, except for the majority rules consensus using a correlation cutoff of 0.90. While all our models show high accuracy, when the additional model evaluation metrics are considered only one has strong predictive power. Therefore, we only used the majority rules consensus using a correlation cutoff of 0.90 for interpreting variable importance of our data.

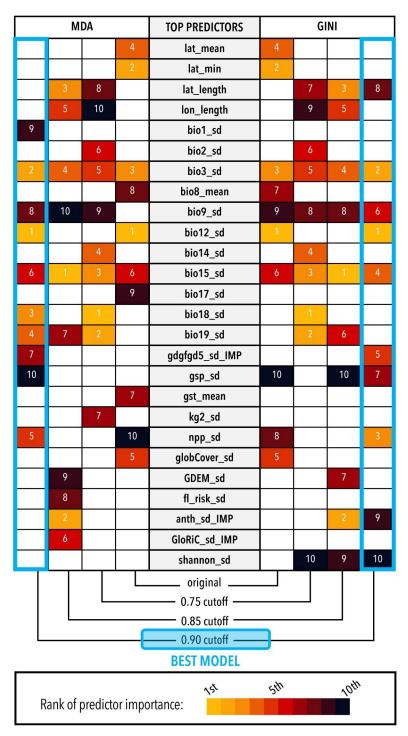
# **Evaluation of variable importance**

We extracted variable importance measurements from each random forest classification using the variable importance metrics MDA and Gini. While there was some overlap of top predictors between different models (Fig 4; S4 Fig in S1 File), no specific predictors were

**Table 2. Results of** *all agree* **consensus random forest models.** Model metrics for each random forest predictive model generated using the *all agree* consensus classifications are shown.

All Agree Models	Original	Correlation  > 0.75	Correlation  > 0.85	Correlation  > 0.90
Accuracy	0.8667	0.9333	0.8667	0.8667
Accuracy (95% CI)	(0.5954, 0.9834)	(0.6805, 0.9983)	(0.5954, 0.9834)	(0.5954, 0.9834)
No Information Rate	0.8	0.8	0.8	0.8
Pos Pred Value	0.8571	0.9231	0.8571	0.8571
Neg Pred Value	1	1	1	1
P-Value [Acc > NIR] 0.398		0.1671	0.398	0.398

https://doi.org/10.1371/journal.pone.0310932.t002



**Fig 4.** Variable importance for random forest classification models generated using the *majority rules* consensus. Variables ranked among the top ten most important variables (based on MDA and Gini) from the classification model generated at different correlation cut-offs are included. Blue highlighting indicates the best consensus model (*majority rules*-correlation cutoff 0.90).

(a)

		Med	dian	Kruskal-Wallis Test		
ES .	Top Predictors	Not Hidden	Hidden	chi-squared	p-value	
MAJORITY RULES (cutoff = 0.90)	annual precipitation sd	22.3	246	22.288	2.35E-06	
<b>3 TY</b>	isothermality sd	0.354	2.1	15.854	6.84E-05	
	warmest quarter precipitation sd (mm)	10.2	53.9	10.943	0.0009397	
M/	coldest quarter precipitation sd (mm)	14.7	63.8	13.49	0.0002398	
	net primary productivity sd	33.4	140	13.562	0.0002308	

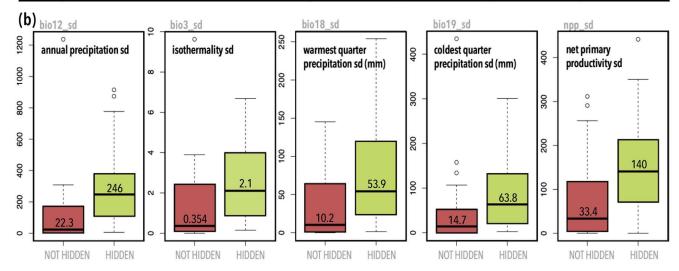


Fig 5. Comparison of hidden vs not hidden trait values for the top five most important predictors of the best consensus model (*majority rules*-correlation cutoff 0.90). A, Columns 1–2 of the table identifys the specific model and predictors (i.e., traits). Columns 3–4 show the median trait values for each group (i.e., hidden vs not hidden). Columns 5–6 show the results of Kruskal-Wallis significance tests, which determine if the difference in median trait values for each group is statistically significant. B, Corresponding boxplots of the median trait values for the top five most important predictors show a significant difference in the range of values between hidden and non-hidden genetic lineages.

https://doi.org/10.1371/journal.pone.0310932.g005

consistently predicted to be of significantly higher importance than other predictors in the model. Instead, importance was split across numerous predictors that were found to be unstable between models. This instability supports previous indications that many of the random forest models are likely prone to overfitting. Despite the lack of a strong set of standout predictors across models, one pattern does emerge that is applicable to the species in our dataset. Of the top ten most important predictors in each model, approximately 85% are measurements of standard deviation (vs. measurements of mean values or life history traits) (\$5 Data). This is supported by further examination of our one model that was able to predict significantly better than random, the *majority rules* consensus with a correlation coefficient cutoff of 0.90, in which the top five most important predictors are measurements of standard deviation. Significance testing indicates that species identified as containing hidden genetic lineages often have ranges characterized by a larger variance in annual and seasonal precipitation, isothermality, and net primary productivity than species not identified as harboring hidden genetic lineages (Fig 5).

### **Discussion**

When identifying genetic lineages or delimiting species, it is important to recognize that species concepts are complex and often differ based on various factors, such as geographic location, reproductive isolating mechanisms, genetic markers, and taxonomic practices. Therefore, it is essential to approach species delimitations with caution and to recognize that they represent a hypothesis or starting point rather than a definitive answer [76]. In addition, while mitochondrial data can be suitable for preliminary assessments of species diversity [77], these assessments should be considered in tandem with other species information and relevant data when describing species boundaries. However, with recent advances in technology rapidly increasing the quantity of publicly accessible genetic and geographic datasets, these data offer a cost effective and efficient way to explore large-scale patterns and predictors of intraspecific genetic variation (e.g., [24,29,78]).

Our results suggest that there are undescribed genetic lineages that may warrant further investigation distributed within Caudata. Adequately documenting biodiversity, both at the species and population level, is a first step in understanding the eco-evolutionary processes generating this diversity. However, in most clades, the Linnean shortfall is likely to influence broad scale patterns detected using macrogenetic approaches [13], making it essential to consider how the taxonomic designations used to inform these approaches influence the patterns detected. This is particularly important when dealing with clades suspected of harboring high levels of cryptic diversity. For example, Miraldo et al. [24] generated the first global map of genetic diversity within species of mammals and amphibians. One of their main conclusions was that amphibians displayed lower levels of genetic variation in areas with higher human impact. Similarly, in amphibians, several recent studies have found within species genetic diversity to be lower in temperate regions in species with smaller ranges and at higher elevations [30,37]. The methods used to detect these patterns are based on current taxonomic knowledge, and as such, rely on the assumption that the species designations used are accurate. However, if species descriptions inaccurately reflect biological diversity, nominal species that contain cryptic species will display higher levels of genetic diversity, while not reflecting true within species variation, potentially skewing our interpretation of any patterns that result.

### **Evaluating support for identified genetic lineages**

While our delimitation of genetic lineages are a starting point, or hypothesis generation step, for evaluating a species in nature where complex processes, such as hybrid zones, and adequate sampling must be considered [e.g., 75–77], we believe these computational approaches are useful for targeting species in further need of examination. We conducted a literature search to explore whether the nominal species in our dataset have been previously explored from a species delimitation approach. We used the online American Museum of Natural History taxonomic and nomenclatural database, Amphibian Species of the World [79]), to evaluate current taxonomic research in each nominal species of salamander predicted to contain hidden diversity in our consensus model. Species in which we were able to identify research-based support for the potential of undescribed diversity were recorded, along with the related articles in which the diversity was described as well as the type of data used (see S7 Data). Nearly 70% of species the majority rules consensus suggests harbor hidden lineages, contain results that also support the potential splitting of species into separate lineages. Out of these about 38% were explored using mt DNA only, 10% with nuclear DNA only, 35% using a combination of both nuclear and mt DNA and 17% using mt DNA, nuclear DNA and morphology. Just under 10% of the species display a complex history of hybridization, making delimitations difficult, a situation not uncommon in salamanders [e.g., 44,80,81]. We were unable to find results for

roughly 25% of our species data. We encountered 5 species in which the results of previous delimitation work were either unclear or considered highly contested (e.g., *Ichthyosaura alpestris*, *Batrachuperus karlschmidti*, *Batrachuperus taibaiensis*, and *Salamandrella schrenckii*). Taxonomy is dynamic field [33] and given our search, it can be difficult to use current opensource data relying solely on species names. However, the current literature largely supports the delimitation results found here and suggests a number of species in further need of investigation (see citations in \$\frac{S7 Data}{2}\$, formal name changes, and an ability to update current opensource databases to reflect these changes). Additionally, even though there are limitations to using current open-source data that might not keep up to date with current taxonomy, we can still determine what factors might predict the presence of species likely to possess undescribed genetic diversity.

# Significant predictors of diversity

Significance testing of the most important predictors from our best model (*majority rules* consensus with a correlation coefficient cutoff of 0.90) indicates that the species which our analysis identified as containing hidden genetic lineages often have ranges characterized by a larger variance in annual and seasonal precipitation, isothermality, and net primary productivity when compared to species that were not identified as containing hidden genetic lineages by our analysis (Fig 5B). And while the order of the most important traits is unstable across different models, across all models most of the traits found to be important were measurements of standard deviation (vs. measurements of mean values or life history traits) (S5 Data). This suggests that the presence of variation in climate, rather than any species-specific trait or characteristic is the most identifiable driving force of within species genetic diversity for salamanders at this scale. Our findings align with similar studies of amphibians using a different measure of genetic variation within species (nucleotide diversity), which concluded that species traits were not a predictor of intraspecific genetic diversity [30,37]. Using similar methods, our results in salamanders differ from that found in mammals, where body size and range size were the most important predictors [35].

These findings are somewhat consistent with other studies of salamander diversification. Reproductive mode (larval stages, direct development) and habitat (combinations of terrestrial, aquatic, arboreal) vary across species and have evolved multiple times but have not been found to directly correlate with speciation, though being a direct developer might increase diversification rates [82]. In vertebrate clades, terrestrial organisms tend to have higher diversification rates than aquatic organisms [83], but we did not have a large number of fully terrestrial species in our dataset, which might have limited our ability to detect this as an important predictor. Alternatively, in one species which has intraspecific variation in habit, *Salamandra salamandra*, terrestrial-breeding individuals exhibited greater geographic genetic differentiation compared to aquatic-breeding individuals [84]. Not surprisingly, this species showed conflicting results in our delimitation analyses. Because various species delimitation methods are not similarly sensitive to differing levels of population structure, we would expect these methods to perform more inconsistently within species with highly variable genetic and geographic distance across different life histories [80,85].

Given that salamanders are relatively constrained in body form and ecological niches, variation in climatic variables seems like a reasonable explanation for species containing cryptic diversity. This follows the suggestion that change in climatic niche variables increases diversification rates in plethodontid salamanders [86]). Diversification rates in frogs and salamanders have been shown to be higher near the tropics [83], so one might expect latitude to be an

•	•				
Mammal Models	ABGD COI	ABGD cytb	GMYC COI	GMYC cytb	consensus
Accuracy	0.737	0.68	0.6429	0.6517	0.781
Accuracy (95% CI)	(0.6802, 0.7885)	(0.6333, 0.7241)	(0.5821, 0.7004)	(0.6014, 0.6996)	(0.7273, 0.8285)
No Information Rate	0.7222	0.6235	0.6128	0.5488	0.6533
Pos Pred Value	0.56667	0.6304	0.17271	0.6624	2.85E-06
Neg Pred Value	0.75833	0.6937	0.5571	0.6345	0.807
P-Value [Acc > NIR]	0.32	0.008792	0.6735	3.00E-05	2.85E-06

Table 3. Summary of results of mammal random forest classification models presented in Parsons et al., 2022 [35]). Model metrics for each random forest classification model generated using data from the class Mammalia are shown.

important predictor. However, latitude was not included in the list of predictor variables that were likely to be important (Fig 4).

# Predictive modeling as a tool to address the Linnean shortfall

Recently, Parsons *et al.* [35] used publicly available genetic barcoding data to develop a predictive framework to identify mammalian clades most likely to contain hidden species and determine specific trait complexes that indicate where hidden mammal diversity is likely to exist. We adopted a similar approach to evaluate undescribed genetic lineages in the clade Caudata, a group which differs from mammals in several key aspects, including species richness and sampling intensity. We focused on a lower taxonomic level so there are fewer recognized species of salamanders (<1000; [57]) compared to the mammal dataset, making the ability to produce robust predictive models more challenging. Additionally, there was a smaller proportion of available data for salamanders than mammals (~10% compared to 60% of described species). However, these smaller datasets might be more realistic in that they are more representative of the type of data most likely to be available for the taxonomic groups that are in greatest need of attention from taxonomists.

While the random forest models generated in this study actually have a higher overall accuracy than those used in Parsons et al. [35] (see Table 3), relying on this metric alone to evaluate the performance of predictive models can be misleading [87–89]. For classification models, model accuracy depends on how well the predicted classifications match the observed classifications. While seemingly straightforward, accuracy does not account for other model characteristics that may be influencing model behavior, such as the class frequencies of the underlying dataset [87]). In cases where one class occurs at a much higher frequency than the other, a predictive model can attain a high accuracy by simply always predicting the higher class. Therefore, an important benchmark to consider when interpreting overall model accuracy is the frequency at which the majority class occurs, the no information rate (NIR) [88]. If a model's accuracy is not significantly higher than the NIR (i.e., p-value [Accuracy > NIR]), it can remain unclear whether the model is making meaningful decisions. In our models, the overall accuracy was found to be high, but the 95% confidence intervals for the accuracy values are very wide for most of the models. In addition, because the dataset is skewed towards species classified as containing hidden diversity, the p-value [Accuracy > NIR] was found to be significant in only one model. This is important to point out because even though there are large datasets available, choosing the right analytical tools can remain challenging depending on the use of the predictive models. Beyond analytical tools, it's also important to consider your dataset, and how the characteristics of your dataset are affecting the results you obtain. Considering the scale of not only the dataset, but also the analytical methods used and the

pattern one is attempting to examine is especially important in meta-analyses, as different patterns emerge at different scales [89].

### **Conclusions**

Here, we chose to utilize biodiversity data from *phylogatR* (i.e., genetic data for which directly associated specimen locality information is available) to avoid potential discrepancies between the distribution of the genetic and geographic data analyzed. By doing so we hoped to gain a more fine-grain understanding of how species genetic diversity is influenced by geographic and environmental factors [23]. However, making this choice significantly decreased the amount of data available and led to a greatly reduced dataset. Our study included 1676 DNA barcoding sequences from the genes *COI* and *cytb* (768 and 908 sequences each, respectively). However, a 3/31/23 search of GenBank for salamander barcoding sequences from the genes *COI* and *cytb* returned a total of 17097 sequences (4468 and 12629 sequences each, respectively; see S6 Data). Similarly, while we were able to obtain 1676 occurrence records tied to the genetic sequences used in this study, a GBIF search for geographic occurrences tied to salamander preserved specimens and material samples returned 675243 records (see S6 Data). This study highlights the lack of genetic data with easily-associated geographic information.

Despite limitations in dataset size and geographic coverage, our framework effectively identified salamander species likely to contain undescribed genetic diversity with agreement across multiple delimitation methods. These species likely represent good candidates for further taxonomic evaluation. While we were unable to pinpoint a specific predictor variable as the most important for predicting undescribed diversity, our findings suggest that hidden diversity in salamanders is likely higher in species with broad geographic ranges characterized by significant climatic variability. This insight serves as a starting point for future integrative taxonomic work and underscores that much diversity remains undiscovered. Although our study was constrained by data availability, the framework we used could further elucidate these relationships with access to more comprehensive genetic and geographic data, highlighting the crucial importance of such data in biodiversity studies.

The numerous benefits of making biological data more broadly available have been repeatedly demonstrated [90]. And recent years have seen a significant increase in the amount of available specimen and biodiversity data. The utility of these data to address large scale patterns of biodiversity, such as those examined in this study, is enhanced by our ability to integrate and synthesize data across different data sources, types, and taxonomic groups [91]. Our study highlights the importance of not just making these data available but making them available in a way that is standardized and will facilitate integration and re-use for future generations to come (e.g., [92,93]).

# Supporting information

S1 File. (Contains Figure S1: Distribution of missing data in the salamander trait database; Figure S2: Distribution of imputed trait data; Figure S3: Correlation matrix of predictor variables; Figure S4: Variable importance for predictive models). (DOCX)

S1 Table. Comparison of salamander data available for each loci on the phylogatR database.

(DOCX)

S2 Table. Comparison of model accuracy confidence intervals between salamander and mammal predictive models.

(DOCX)

S1 Data. Nucleotide diversity of Caudata sequences from phylogatR.

(XLSX)

S2 Data. *PhylogatR* identification numbers (accession\_sourceID; see Pelletier *et al.*, 2022 [38]) for records analyzed.

(XLSX)

S3 Data. Final dataset of response and predictor variables.

(XLSX)

S4 Data. Variable specifics and source information.

(XLSX)

S5 Data. Variable importance extended results.

(XLSX)

S6 Data. Results of search for publicly available genetic and geographic salamander data.

(XLSX)

S7 Data. Results of literature search for genetic lineages in recognized salamander species.

(XLSX)

# **Acknowledgments**

We thank Radford University Office of Undergraduate Research for supporting AEG through the Accelerated Research Opportunities and Research Rookies programs. We thank members of the Carstens lab for their comments on a draft of this manuscript.

### **Author Contributions**

Conceptualization: Danielle J. Parsons, Bryan C. Carstens, Tara A. Pelletier.

Data curation: Danielle J. Parsons, Abigail E. Green, Bryan C. Carstens, Tara A. Pelletier.

Formal analysis: Danielle J. Parsons, Abigail E. Green, Bryan C. Carstens, Tara A. Pelletier.

Funding acquisition: Bryan C. Carstens, Tara A. Pelletier.

Methodology: Tara A. Pelletier.

**Project administration:** Tara A. Pelletier.

Resources: Bryan C. Carstens.

Software: Danielle J. Parsons, Abigail E. Green, Tara A. Pelletier.

**Supervision:** Tara A. Pelletier. **Validation:** Danielle J. Parsons.

Visualization: Danielle J. Parsons, Tara A. Pelletier.

Writing - original draft: Danielle J. Parsons, Bryan C. Carstens, Tara A. Pelletier.

Writing - review & editing: Danielle J. Parsons, Bryan C. Carstens, Tara A. Pelletier.

### References

- Gadelha LMR, Siracusa PC, Dalcin EC, Silva LAE, Augusto DA, Krempser E, et al. 2021. A survey of biodiversity informatics: Concepts, practices, and challenges. WIREs Data Mining and Knowledge Discovery 11.
- Guralnick R, Hill A. 2009. Biodiversity informatics: automated approaches for documenting global biodiversity patterns and processes. Bioinformatics 25: 421–428. <a href="https://doi.org/10.1093/bioinformatics/btn659">https://doi.org/10.1093/bioinformatics/btn659</a> PMID: 19129210
- 3. Whittaker RJ, Araújo MB, Jepson P, Ladle RJ, Watson JEM, Willis KJ. 2005. Conservation Biogeography: assessment and prospect: Conservation Biogeography. Diversity and Distributions 11: 3–23.
- Mora C, Tittensor DP, Adl S, Simpson AGB, Worm B. 2011. How Many Species Are There on Earth and in the Ocean? (Mace GM, Ed.). PLoS Biology 9: e1001127. <a href="https://doi.org/10.1371/journal.pbio.1001127">https://doi.org/10.1371/journal.pbio.1001127</a> PMID: 21886479
- Moritz C. 1994. Defining 'Evolutionarily Significant Units' for conservation. Trends in Ecology & Evolution 9: 373–375. https://doi.org/10.1016/0169-5347(94)90057-4 PMID: 21236896
- Mable BK. 2019. Conservation of adaptive potential and functional diversity: integrating old and new approaches. Conservation Genetics 20: 89–100.
- Jockusch EL, Martínez-Solano I, Hansen RW, Wake DB. 2012. Morphological and molecular diversification of slender salamanders (Caudata: Plethodontidae: Batrachoseps) in the southern Sierra Nevada of California with descriptions of two new species. Zootaxa 3190: 1.
- 8. Camp CD, Wooten JA. 2016. Hidden in Plain Sight: Cryptic Diversity in the Plethodontidae. Copeia 104: 111–117.
- Bernardes M, Le MD, Nguyen TQ, Pham CT, Pham AV, Nguyen TT, et al. 2020. Integrative taxonomy reveals three new taxa within the Tylototriton asperrimus complex (Caudata, Salamandridae) from Vietnam. ZooKeys 935: 121–164. <a href="https://doi.org/10.3897/zookeys.935.37138">https://doi.org/10.3897/zookeys.935.37138</a> PMID: 32508505
- Mead LS, Clayton DR, Nauman RS, Olson DH, Pfrender ME. 2005. Newly discovered populations of salamanders from Siskiyou County California represent a species distinct from Plethodon stormi. Herpetologica 61: 158–177.
- Parra Olea G, Garcia-Castillo MG, Rovito SM, Maisano JA, Hanken J, Wake DB. 2020. Descriptions of five new species of the salamander genus Chiropterotriton (Caudata: Plethodontidae) from eastern Mexico and the status of three currently recognized taxa. PeerJ 8: e8800. <a href="https://doi.org/10.7717/peerj.8800">https://doi.org/10.7717/peerj.8800</a> PMID: 32518712
- Nauman RS, Olson DH. 2008. Distribution and Conservation Of Plethodon Salamanders On Federal Lands In Siskiyou County, California. Northwestern Naturalist 89: 1.
- **13.** Hortal J, de Bello F, Diniz-Filho JAF, Lewinsohn TM, Lobo JM, Ladle RJ. 2015. Seven Shortfalls that Beset Large-Scale Knowledge of Biodiversity. Annual Review of Ecology, Evolution, and Systematics 46: 523–549.
- Steffen MA, Irwin KJ, Blair AL, Bonett RM. 2014. Larval masquerade: a new species of paedomorphic salamander (Caudata: Plethodontidae: Eurycea) from the Ouachita Mountains of North America. Zootaxa 3786: 423. https://doi.org/10.11646/zootaxa.3786.4.2 PMID: 24869544
- Nishikawa K, Matsui M. 2014. Three new species of the salamander genus Hynobius (Amphibia, Urodela, Hynobiidae) from Kyushu, Japan. Zootaxa 3852: 203. <a href="https://doi.org/10.11646/zootaxa.3852.2.3">https://doi.org/10.11646/zootaxa.3852.2.3</a>
   PMID: 25284394
- **16.** Min MS, Baek HJ, Song JY, Chang MH, Poyarkov NAJr. 2016. A new species of salamander of the genus Hynobius (Amphibia, Caudata, Hynobiidae) from South Korea. Zootaxa 4169: 475. PMID: <u>27701288</u>
- 17. Kuchta SR, Brown AD, Highton R. 2018. Disintegrating over space and time: Paraphyly and species delimitation in the Wehrle's Salamander complex. Zoologica Scripta 47: 285–299.
- Okamiya H, Sugawara H, Nagano M, Poyarkov NA. 2018. An integrative taxonomic analysis reveals a new species of lotic Hynobius salamander from Japan. PeerJ 6: e5084. <a href="https://doi.org/10.7717/peerj.5084">https://doi.org/10.7717/peerj.5084</a> PMID: 29942708
- Bickford D, Lohman DJ, Sodhi NS, Ng PKL, Meier R, Winker K, et al. 2007. Cryptic species as a window on diversity and conservation. Trends in Ecology & Evolution 22: 148–155. <a href="https://doi.org/10.1016/j.tree.2006.11.004">https://doi.org/10.1016/j.tree.2006.11.004</a> PMID: 17129636
- Pfenninger M, Schwenk K. 2007. Cryptic animal species are homogeneously distributed among taxa and biogeographical regions. BMC Evolutionary Biology 7: 121. <a href="https://doi.org/10.1186/1471-2148-7-121">https://doi.org/10.1186/1471-2148-7-121</a> PMID: 17640383
- Lyman RA, Edwards CE. 2022. Revisiting the comparative phylogeography of unglaciated eastern North America: 15 years of patterns and progress. Ecology and Evolution 12. <a href="https://doi.org/10.1002/ece3.8827">https://doi.org/10.1002/ece3.8827</a> PMID: 35475178

- Blanchet S, Prunier JG, De Kort H. 2017. Time to Go Bigger: Emerging Patterns in Macrogenetics. Trends in Genetics 33: 579–580. https://doi.org/10.1016/j.tig.2017.06.007 PMID: 28720482
- Leigh DM, van Rees CB, Millette KL, Breed MF, Schmidt C, Bertola LD, et al. 2021. Opportunities and challenges of macrogenetic studies. Nature Reviews Genetics 22: 791–807. <a href="https://doi.org/10.1038/s41576-021-00394-0">https://doi.org/10.1038/s41576-021-00394-0</a> PMID: 34408318
- Miraldo A, Li S, Borregaard MK, Flórez-Rodríguez A, Gopalakrishnan S, Rizvanovic M, et al. 2016. An Anthropocene map of genetic diversity. Science 353: 1532–1535. <a href="https://doi.org/10.1126/science.aaf4381">https://doi.org/10.1126/science.aaf4381</a> PMID: 27708102
- **25.** Millette KL, Fugère V, Debyser C, Greiner A, Chain FJJ, Gonzalez A. 2020. No consistent effects of humans on animal genetic diversity worldwide (A Mooers, Ed.). Ecology Letters 23: 55–67.
- Carstens BC, Morales AE, Field K, Pelletier TA. 2018. A global analysis of bats using automated comparative phylogeography uncovers a surprising impact of Pleistocene glaciation. Journal of Biogeography 45: 1795–1805.
- 27. Baranzelli MC, Cosacov A, Sede SM, Nicola MV, Sérsic AN. 2022. Anthropocene refugia in Patagonia: A macrogenetic approach to safeguarding the biodiversity of flowering plants. Biological Conservation 268: 109492.
- Gratton P, Marta S, Bocksberger G, Winter M, Keil P, Trucchi E, et al. 2017. Which Latitudinal Gradients for Genetic Diversity? Trends in Ecology & Evolution 32: 724–726. <a href="https://doi.org/10.1016/j.tree.2017.07.007">https://doi.org/10.1016/j.tree.2017.07.007</a> PMID: 28807398
- Pelletier TA, Carstens BC. 2018. Geographical range size and latitude predict population genetic structure in a global survey. Biology Letters 14: 20170566. <a href="https://doi.org/10.1098/rsbl.2017.0566">https://doi.org/10.1098/rsbl.2017.0566</a> PMID: 29343561
- **30.** Barrow LN, Fonseca EM, Thompson CEP, Carstens BC. 2021. Predicting amphibian intraspecific genetic diversity with machine learning: Challenges and prospects for integrating traits, geography, and genetic data, Molecular Ecology Resources 21: 2718–2831.
- Fonseca EM, Pelletier TA, Decker SK, Parsons DJ, Carstens BC. 2023. Pleistocene glaciations caused the latitudinal gradient of within-species genetic diversity, Evolution Letters 7: 331–338. <a href="https://doi.org/10.1093/evlett/grad030">https://doi.org/10.1093/evlett/grad030</a> PMID: 37829497
- Pelletier TA, Carstens BC, Tank DC, Sullivan J, Espíndola A. 2018. Predicting plant conservation priorities on a global scale. Proceedings of the National Academy of Sciences 115: 13027–13032. <a href="https://doi.org/10.1073/pnas.1804098115">https://doi.org/10.1073/pnas.1804098115</a> PMID: 30509998
- **33.** Raposo MA, Kirwan GM, Lourenço ACC, Sobral G, Bockmann FA, Stopiglia R. 2021. On the notions of taxonomic 'impediment', 'gap', 'inflation' and 'anarchy', and their effects on the field of conservation. Systematics and Biodiversity 19: 296–311.
- 34. Breiman L. 2001. Random Forests. Machine Learning 45: 5–32.
- **35.** Parsons DJ, Pelletier TA, Wieringa JG, Duckett DJ, Bryan C. Carstens. 2022. Analysis of biodiversity data suggests that mammal species are hidden in predictable places. Proceedings of the National Academy of Sciences 119: e2103400119. https://doi.org/10.1073/pnas.2103400119 PMID: 35344422
- **36.** Smith BT, Seeholzer GF, Harvey MG, Cuervo AM, Brumfield RT. 2017. A latitudinal phylogeographic diversity gradient in birds. PLOS Biology 15: e1002610.
- Amador L, Arroyo-Torres I, Lisa N. Barrow LN. 2023. Machine learning and phylogenetic models identify predictors of genetic variation in Neotropical amphibians. bioRxiv.
- Pelletier TA, Parsons DJ, Decker SK, Crouch S, Franz E, Ohrstrom J, et al. 2022. PhylogatR: Phylogeographic data aggregation and repurposing. Molecular Ecology Resources 22: 2830–2842. <a href="https://doi.org/10.1111/1755-0998.13673">https://doi.org/10.1111/1755-0998.13673</a> PMID: 35748425
- 39. Bánki, O., Roskov, Y., Döring, M., Ower, G., Vandepitte, L., Hobern, D., et al. 2022. Catalogue of Life Checklist (Y. Roskov, Ed.; Version 2022-05-20).
- **40.** Pelletier TA, Duffield DA, DeGrauw EA. 2011. Rangewide Phylogeography of the Western Red-Backed Salamander (Plethodon vehiculum). Northwestern Naturalist 92: 200–210.
- Pelletier TA, Crisafulli C, Wagner S, Zellmer AJ, Carstens BC. 2015. Historical Species Distribution Models Predict Species Limits in Western Plethodon Salamanders. Systematic Biology 64: 909–925. https://doi.org/10.1093/sysbio/syu090 PMID: 25414176
- Pelletier TA, Carstens BC. 2016. Comparing range evolution in two western Plethodon salamanders: glacial refugia, competition, ecological niches, and spatial sorting. Journal of Biogeography 43: 2237– 2249.
- 43. Jones KS, Weisrock DW. 2018. Genomic data reject the hypothesis of sympatric ecological speciation in a clade of Desmognathus salamanders. Evolution 72: 2378–2393. <a href="https://doi.org/10.1111/evo.13606">https://doi.org/10.1111/evo.13606</a> PMID: 30246244

- 44. Pyron RA, O'Connell KA, Lemmon EM, Lemmon AR, Beamer DA. 2020. Phylogenomic data reveal reticulation and incongruence among mitochondrial candidate species in Dusky Salamanders (Desmognathus). Molecular Phylogenetics and Evolution. 146: 1055–7903. <a href="https://doi.org/10.1016/j.ympev.2020.106751">https://doi.org/10.1016/j.ympev.2020.106751</a> PMID: 32028035
- 45. Dufresnes C, Brelsford A, Jeffries DL, Mazepa G, Suchan T, Canestrelli D, et al. 2021. Mass of genes rather than master genes underlie the genomic architecture of amphibian speciation. Proceedings of the National Academy of Sciences 118: e2103963118. <a href="https://doi.org/10.1073/pnas.2103963118">https://doi.org/10.1073/pnas.2103963118</a>
  PMID: 34465621
- 46. Maddison WP, Maddison DR. Mesquite: a modular system for evolutionary analysis. Version 3.81. 2023.
- Katoh K, Standley DM. 2013. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. Molecular Biology and Evolution 30: 772–780. <a href="https://doi.org/10.1093/molbev/mst010">https://doi.org/10.1093/molbev/mst010</a> PMID: 23329690
- **48.** Pons J, Barraclough TG, Gomez-Zurita J, Cardoso A, Duran DP, Hazell S, et al. 2006. Sequence-Based Species Delimitation for the DNA Taxonomy of Undescribed Insects (Hedin M, Ed.). Systematic Biology 55: 595–609. https://doi.org/10.1080/10635150600852011 PMID: 16967577
- Puillandre N, Lambert A, Brouillet S, Achaz G. 2012. ABGD, Automatic Barcode Gap Discovery for primary species delimitation: ABGD, AUTOMATIC BARCODE GAP DISCOVERY. Molecular Ecology 21: 1864–1877.
- Puillandre N, Brouillet S, Achaz G. 2021. ASAP: assemble species by automatic partitioning. Molecular Ecology Resources 21: 609–620. https://doi.org/10.1111/1755-0998.13281 PMID: 33058550
- Carstens BC, Pelletier TA, Reid NM, Satler JD. 2013. How to fail at species delimitation. Molecular Ecology 22: 4369–4383. https://doi.org/10.1111/mec.12413 PMID: 23855767
- 52. Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, et al. 2019. BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis (Pertea M, Ed.). PLOS Computational Biology 15: e1006650. https://doi.org/10.1371/journal.pcbi.1006650 PMID: 30958812
- 53. Abadi S, Azouri D, Pupko T, Mayrose I. 2019. Model selection may not be a mandatory step for phylogeny reconstruction. Nature Communications 10: 934. <a href="https://doi.org/10.1038/s41467-019-08822-w">https://doi.org/10.1038/s41467-019-08822-w</a> PMID: 30804347
- 54. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. 2018. Posterior Summarization in Bayesian Phylogenetics Using Tracer 1.7 (Susko E, Ed.). Systematic Biology 67: 901–904. <a href="https://doi.org/10.1093/sysbio/syy032">https://doi.org/10.1093/sysbio/syy032</a> PMID: 29718447
- **55.** Ezard T. Fujisawa T. Barraclough T. SPLITS: species' limits by threshold statistics. R package version 1.0–20. 2009. https://rdrr.io/rforge/splits/.
- Oliveira BF, São-Pedro VA, Santos-Barrera G, Penone C, Costa GC. 2017. AmphiBIO, a global database for amphibian ecological traits. Scientific Data 4: 170123. <a href="https://doi.org/10.1038/sdata.2017.123">https://doi.org/10.1038/sdata.2017.123</a>
   PMID: 28872632
- 57. University of California, Berkeley. AmphibiaWeb. 2023. https://amphibiaweb.org/.
- Karger DN, Conrad O, Böhner J, Kawohl T, Kreft H, Soria-Auza RW, et al. 2017. Climatologies at high resolution for the earth's land surface areas. Scientific Data 4: 170122. <a href="https://doi.org/10.1038/sdata.2017.122">https://doi.org/10.1038/sdata.2017.122</a> PMID: 28872642
- 59. Karger, Dirk Nikolaus, Conrad, Olaf, Böhner, Jürgen, Kawohl, Tobias, Kreft, Holger, Soria-Auza, Rodrigo Wilber, et al. 2021. Climatologies at high resolution for the earth's land surface areasCHELSA V2.1 (current).: 2.1 KB.
- 60. NASA. ASTER global digital elevation model version 2. 2011. https://asterweb.jpl.nasa.gov/gdem.asp.
- **61.** Tuanmu MN, Jetz W. 2015. A global, remote sensing-based characterization of terrestrial habitat heterogeneity for biodiversity and ecosystem modelling: Global habitat heterogeneity. Global Ecology and Biogeography 24: 1329–1339.
- **62.** European Space Agency. 2009. ESA GlobCover project.
- **63.** Ouellet Dallaire C, Lehner B, Sayre R, Thieme M. 2019. A multidisciplinary framework to derive global river reach classifications at high spatial resolution. Environmental Research Letters 14: 024003.
- 64. Peduzzi P. 2019. The Disaster Risk, Global Change, and Sustainability Nexus. Sustainability 11: 957.
- **65.** Ellis EC, Klein Goldewijk K, Siebert S, Lightman D, Ramankutty N. 2010. Anthropogenic transformation of the biomes, 1700 to 2000: Anthropogenic transformation of the biomes. Global Ecology and Biogeography: no-no.
- Center for International Earth Science Information Network—CIESIN. Socioeconomic Data and Applications Center (SEDAC) Gridded Populations of the World (GPW). 2016.
- United NationsUNISDR. World Bank Development Economics Research Group (DECRG) Gross Domestic Product. 2010. https://datacatalog.worldbank.org/search/dataset/0037850.

- **68.** Hijmans RJ. raster: Geographic data analysis and modeling. R package version 3.6–26. 2016. <a href="https://cran.r-project.org/web/packages/raster/index.html">https://cran.r-project.org/web/packages/raster/index.html</a>.
- Keitt KH, Bivand R. rgdal: Bindings for the Geospatial Data Abstraction Library. R package version 0.2– 5. 2017. https://cran.r-project.org/package=rgdal.
- **70.** Hijmans RJ, Karney C, Williams E, Vennes C. Geosphere: Spherical trigonometry. R package version 1.5–5. 2016. https://cran.r-project.org/package=geosphere.
- Wickham H. 2011. The Split-Apply-Combine Strategy for Data Analysis. Journal of Statistical Software 40.
- van Buuren S, Groothuis-Oudshoorn K. 2011. mice: Multivariate Imputation by Chained Equations in R. Journal of Statistical Software 45.
- Kuhn M. 2008. Building Predictive Models in R Using the caret Package. Journal of Statistical Software 28.
- 74. Keesey M. PhyloPic. https://www.phylopic.org.
- Taiyun Wei, Viliam Simko. 2021. R package 'corrplot': Visualization of a Correlation Matrix. (Version 0.92).
- 76. Hillis DM. 2019. Species Delimitation in Herpetology. Journal of Herpetology 53: 3.
- Gostel MR, Kress WJ. 2022. The Expanding Role of DNA Barcodes: Indispensable Tools for Ecology, Evolution, and Conservation. Diversity 14: 213.
- 78. Yiming L, Siqi W, Chaoyuan C, Jiaqi Z, Supen W, Xianglei H, et al. 2021. Latitudinal gradients in genetic diversity and natural selection at a highly adaptive gene in terrestrial mammals. Ecography 44: 206–218.
- 79. Darrel FR. 2024. Amphibian Species of the World: an Online Reference. Version 6.2 (December 2023). Electronic Database accessible at <a href="https://amphibiansoftheworld.amnh.org/index.php">https://amphibiansoftheworld.amnh.org/index.php</a>. American Museum of Natural History, New York, USA.
- Luo A, Ling C, Ho SYW, Zhu C. 2018. Comparison of Methods for Molecular Species Delimitation Across a Range of Speciation Scenarios. Systematic Biology 67: 830–846. <a href="https://doi.org/10.1093/sysbio/syy011">https://doi.org/10.1093/sysbio/syy011</a> PMID: 29462495
- Denton RD, Morales AE, Gibbs HL. 2018. Genome-specific histories of divergence and introgression between an allopolyploid unisexual salamander lineage and two ancestral sexual species. Evolution. 72: 1689–1700. PMID: 29926914
- Liedtke HC, Wiens JJ, Gomez-Mestre I. 2022. The evolution of reproductive modes and life cycles in amphibians. Nature Communications 13: 7039. <a href="https://doi.org/10.1038/s41467-022-34474-4">https://doi.org/10.1038/s41467-022-34474-4</a> PMID: 36396632
- Wiens JJ. 2015. Explaining large-scale patterns of vertebrate diversity. Biology Letters 11: 20150506. https://doi.org/10.1098/rsbl.2015.0506 PMID: 26202428
- 84. Lourenço A, Gonçalves J, Carvalho F, Wang IJ, Velo-Antón G. 2019. Comparative landscape genetics reveals the evolution of viviparity reduces genetic connectivity in fire salamanders. Molecular Ecology 28: 4573–4591. https://doi.org/10.1111/mec.15249 PMID: 31541595
- **85.** Burbrink FT, Sara Ruane S. 2021. Contemporary Philosophy and Methods for Studying Speciation and Delimiting Species. Ichthyology & Herpetology 109: 874–894.
- **86.** Kozak KH, Wiens JJ. 2010. Accelerated rates of climatic-niche evolution underlie rapid species diversification: Niche evolution and rapid diversification. Ecology Letters 13: 1378–1389.
- 87. Kuhn M. Johnson K. 2013. Applied Predictive Modeling. New York, NY: Springer New York.
- 88. Foster J. Provost, Tom Fawcett, Ron Kohavi. 1998. The Case against Accuracy Estimation for Comparing Induction Algorithms. Machine learning: proceedings of the fifteenth international conference, Madison, Wisconsin, July 24–27, 1998. San Francisco, Calif: Morgan Kaufmann.
- **89.** Gurevitch J, Koricheva J, Nakagawa S, Stewart G. 2018. Meta-analysis and the science of research synthesis. Nature 555: 175–182. https://doi.org/10.1038/nature25753 PMID: 29517004
- **90.** Wüest RO, Zimmermann NE, Zurell D, Alexander JM, Fritz SA, Hof C, et al. 2020. Macroecology in the age of Big Data—Where to go from here? Journal of Biogeography 47: 1–12.
- 91. Heberling JM, Miller JT, Noesgaard D, Weingart SB, Schigel D. 2021. Data integration enables global biodiversity synthesis. Proceedings of the National Academy of Sciences 118: e2018093118. <a href="https://doi.org/10.1073/pnas.2018093118">https://doi.org/10.1073/pnas.2018093118</a> PMID: 33526679
- **92.** Colella JP, Stephens RB, Campbell ML, Kohli BA, Parsons DJ, Mclean BS. 2021. The Open-Specimen Movement. BioScience 71: 405–414.
- 93. Hardisty AR, Ellwood ER, Nelson G, Zimkus B, Buschbom J, Addink W, et al. 2022. Digital Extended Specimens: Enabling an Extensible Network of Biodiversity Data Records as Integrated Digital Objects on the Internet. BioScience 72: 978–987. https://doi.org/10.1093/biosci/biac060 PMID: 36196222