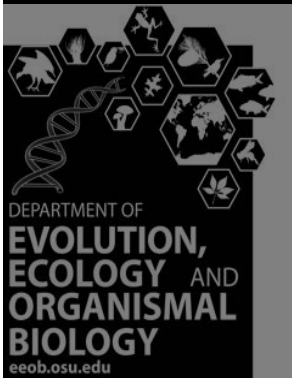# Big data, Bioinformatics, & Biodiversity.

@bryanccarstens
carstens.12@osu.edu
https://carstenslab.osu.edu

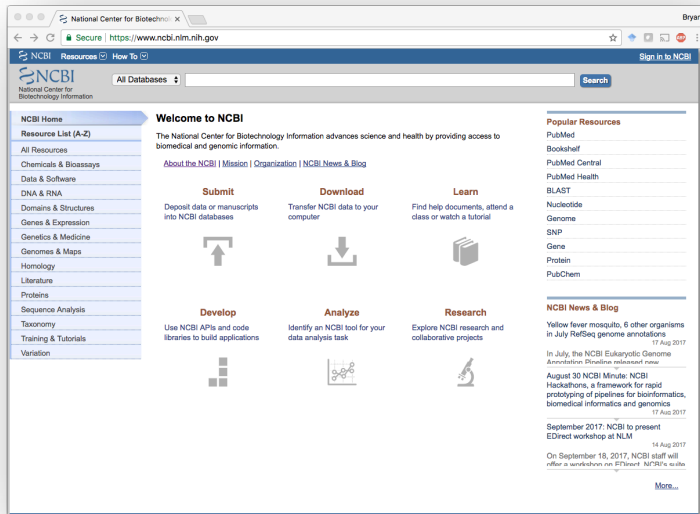DEPARTMENT OF
**EVOLUTION,
ECOLOGY** AND
**ORGANISMAL
BIOLOGY**
eeob.osu.edu

T · H · E
OHIO
STATE
UNIVERSITY

# The biological sciences are data rich…



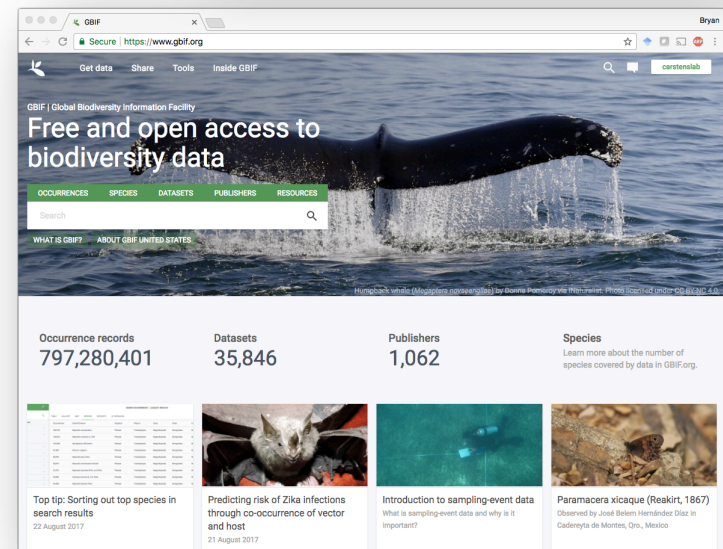- Gene sequences
- Climate data

Bryan Carstens - OSU EEOB
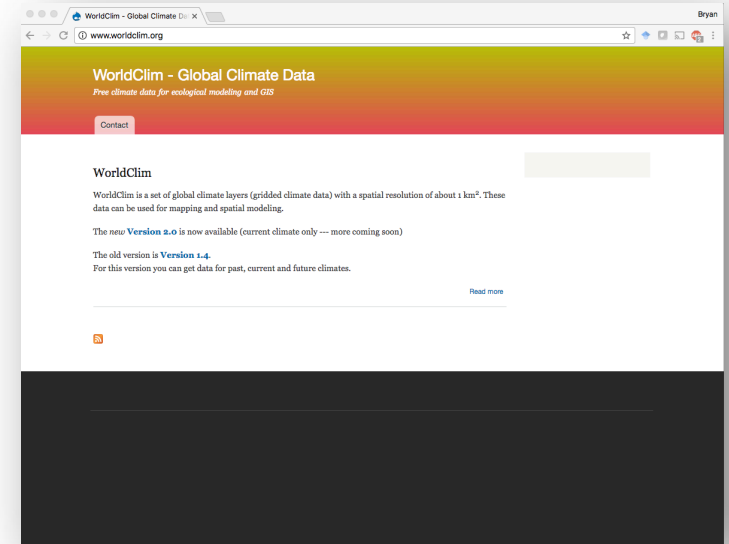
# The biological sciences are data rich…



- Gene sequences
- Climate data

- Collection localities
- Genomes

Bryan Carstens - OSU EEOB

**search = 'phylogeograph*'   -   Web of Science 1987-2015**

~40,000 peer-reviewed papers
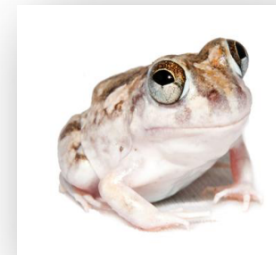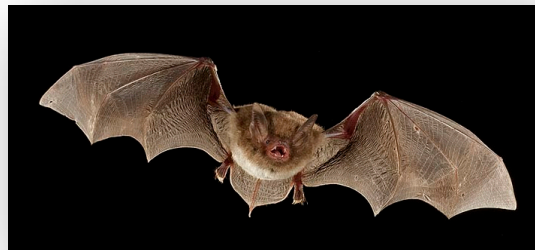


**Publications by Year**

## Speciation with Gene Flow in North American *Myotis* Bats

ARIADNA E. MORALES[1], NATHAN D. JACKSON[2], TANYA A. DEWEY[3], BRIAN C. O'MEARA[2],
AND BRYAN C. CARSTENS[1,*]

[1]Department of Evolution, Ecology and Organismal Biology, Ohio State University, 318 W. 12th Avenue, Columbus, OH 43210, USA;
[2]Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville, 442 Hesler Biology Building, Knoxville, TN 37996, USA;
[3]Department of Biology, Colorado State University, 1878 Campus Delivery, Fort Collins, CO 80523, USA
*Correspondence to be sent to: Department of Evolution, Ecology and Organismal Biology, Ohio State University, 318 W. 12th Avenue,
Columbus, OH 43210, USA; E-mail: carstens.12@osu.edu.

---

www.nature.com/hdy

ORIGINAL ARTICLE

## Model-based analysis supports interglacial refugia over long-dispersal events in the diversification of two South American cactus species

MF Perez[1,3], IAS Bonatelli[1,3], EM Moraes[1] and BC Carstens[2]

---

## Phylogeographic model selection leads to insight into the evolutionary history of four-eyed frogs

Maria Tereza C. Thomé[a] and Bryan C. Carstens[b,1]

[a]Departamento de Zoologia, Instituto de Biociências, Universidade Estadual Paulista, Campus Rio Claro, 13506900 Rio Claro, SP, Brazil; and [b]Department of Evolution, Ecology and Organismal Biology, The Ohio State University, Columbus, OH 43210

Phylogeographic research investigates biodiversity at the interface between populations and species, in a temporal and geographic data, particularly models that incorporate coalescent theory (7) to estimate parameters of interest under a formal framework.

---

## Biogeographic barriers drive co-diversification within associated eukaryotes of the *Sarracenia alata* pitcher plant system

Jordan D. Satler[1], Amanda J. Zellmer[2] and Bryan C. Carstens[1]

[1] Department of Evolution, Ecology and Organismal Biology, The Ohio State University, Columbus, OH, United States
[2] Department of Biology, Occidental College, Los Angeles, CA, United States

---

## Historical Species Distribution Models Predict Species Limits in Western *Plethodon* Salamanders

TARA A. PELLETIER[1], CHARLIE CRISAFULLI[2], STEVE WAGNER[3], AMANDA J. ZELLMER[4], AND BRYAN C. CARSTENS[1,*]

[1]Department of Evolution, Ecology and Organismal Biology, Ohio State University, Columbus, OH 43201; [2]U.S. Forest Service, Pacific Northwest Research Station, Olympia, WA 98512; [3]Department of Biological Sciences, Central Washington University, Ellensburg, WA 98926; and [4]Department of Biology, Occidental College, Los Angeles, CA 90041, USA
*Correspondence to be sent to: Department of Evolution, Ecology and Organismal Biology, Ohio State University, Columbus, OH 43201, USA;
E-mail: carstens.12@osu.edu.

# Phylogeographic methods development…

- **PHRAPL** (DEB-1257784)

### PHRAPL: Phylogeographic Inference Using Approximate Likelihoods

NATHAN D. JACKSON[1], ARIADNA E. MORALES[2], BRYAN C. CARSTENS[2] AND BRIAN C. O'MEARA[1,*]

[1]*Department of Ecology and Evolutionary Biology, University of Tennessee, 442 Hesler Biology Building, Knoxville, TN 37996, USA and*
[2]*Department of Evolution, Ecology and Organismal Biology, Ohio State University, 318 W. 12th Avenue, Columbus, OH 43210, USA*
*Correspondence to be sent to: Department of Ecology and Evolutionary Biology, University of Tennessee, Knoxville, TN 37996, USA; E-mail: bomeara@utk.edu*

- **P2C2M** (DBI-1661029)

## Posterior predictive checks of coalescent models: P2C2M, an R package

MICHAEL GRUENSTAEUDL,*‡ NOAH M. REID,† GREGORY L. WHEELER* and BRYAN C. CARSTENS*

*Department of Evolution, Ecology & Organismal Biology, Ohio State University, Columbus, OH 43210, USA†Department of Environmental Toxicology, University of California, Davis, CA 95616, USA*

# Phylogeographic methods development…

- **Predictive phylogeography** (DEB-1457519)

Identifying cryptic diversity with predictive phylogeography

Anahí Espíndola[1,2], Megan Ruffley[1,2], Megan L. Smith[3], Bryan C. Carstens[3], David C. Tank[1,2] and Jack Sullivan[1,2]

[1]Department of Biological Sciences, University of Idaho, 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[2]Biological Sciences, Institute for Bioinformatics and Evolutionary Studies (IBEST), 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[3]Department of Evolution, Ecology, and Organismal Biology, The Ohio State University, 318 W. 12th Avenue, 300 Aronoff Labs, Columbus, OH 43210-1293, USA

AE, 0000-0001-9128-8836

# Phylogeographic methods development…

- **Predictive phylogeography** (DEB-1457519)

A **framework** that seeks to answer key questions relevant to organismal biology.

- **Multiple species** (regional to global, particular clades to broader taxonomic groups)

- **Integrative** – b/c incorporates all sorts of data (environmental, organismal, genetic)

- relies on **machine learning** to identify key variables

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

Identifying cryptic diversity with predictive phylogeography

Anahí Espíndola[1,2], Megan Ruffley[1,2], Megan L. Smith[3], Bryan C. Carstens[3], David C. Tank[1,2] and Jack Sullivan[1,2]

[1]Department of Biological Sciences, University of Idaho, 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[2]Biological Sciences, Institute for Bioinformatics and Evolutionary Studies (IBEST), 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[3]Department of Evolution, Ecology, and Organismal Biology, The Ohio State University, 318 W. 12th Avenue, 300 Aronoff Labs, Columbus, OH 43210-1293, USA
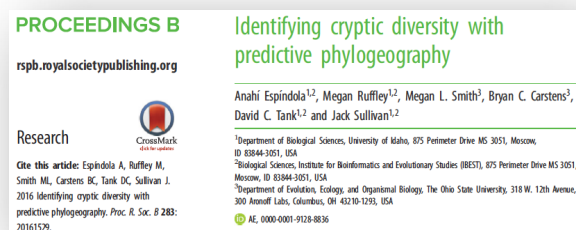
AE, 0000-0001-9128-8836

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

**1. Train the model using existing data.**

Anahí Espíndola[1,2], Megan Ruffley[1,2], Megan L. Smith[3], Bryan C. Carstens[3], David C. Tank[1,2] and Jack Sullivan[1,2]

[1]Department of Biological Sciences, University of Idaho, 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[2]Biological Sciences, Institute for Bioinformatics and Evolutionary Studies (IBEST), 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[3]Department of Evolution, Ecology, and Organismal Biology, The Ohio State University, 318 W. 12th Avenue, 300 Aronoff Labs, Columbus, OH 43210-1293, USA
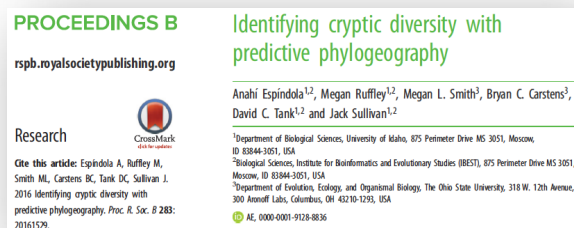
AE, 0000-0001-9128-8836

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.
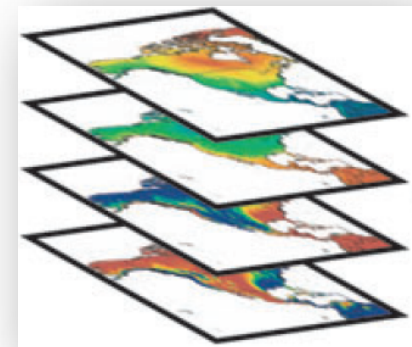
**1. Train the model using existing data.**

- climate data from WorldClim

- species distribution models for all taxa

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.
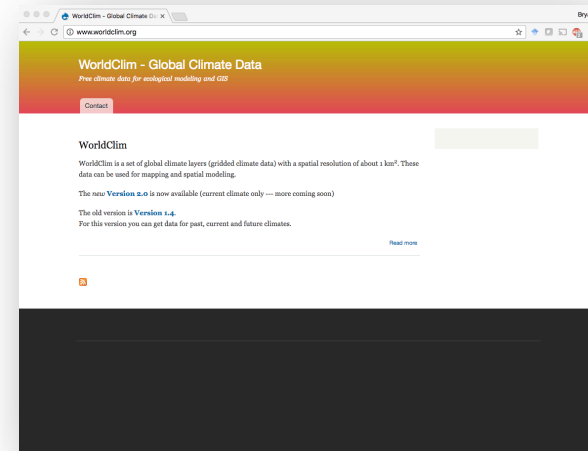
**1. Train the model using existing data.**

- genetic data (sequence, SNPs, MSATs)

- ABC used to calculate posterior probability of historical demographic models for all taxa

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

**2. Build data table**

## For each species in focal ecosystem:

1. species distribution model (climate layers)

2. evolutionary model (posterior probability)

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

**2. Build data table**

## For each species in focal ecosystem:

1. species distribution model (climate layers)

2. evolutionary model (posterior probability)

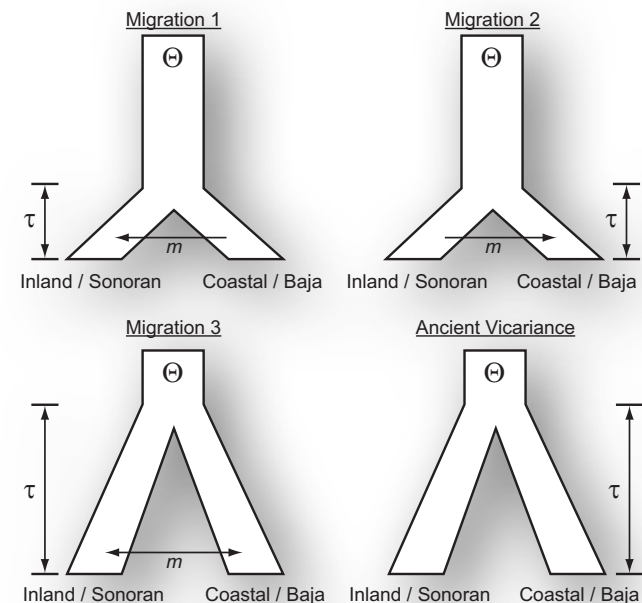3. life history traits (categorical)

4. taxonomic traits (categorical)

Bryan Carstens - OSU EEOB

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

**2. Build data table**

How do we analyze
these disparate data?

**For each species in focal ecosystem:**

1.  species distribution model (climate layers)

2.  evolutionary model (posterior probability)

3.  life history traits (categorical)

4.  taxonomic traits (categorical)

# Random Forest Analysis

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

**3. Random Forest Analysis**

# Random Forest Analysis



Build decision tree by choosing predictor variables at random from data table.

One tree is likely to be a bad predictor of the response variables (cryptic, noncryptic)

# Random Forest Analysis



Repeating this process once may produce a better decision tree…

Bryan Carstens - OSU EEOB

# Random Forest Analysis

# Random Forest Analysis

Bryan Carstens - OSU EEOB

# Random Forest Analysis



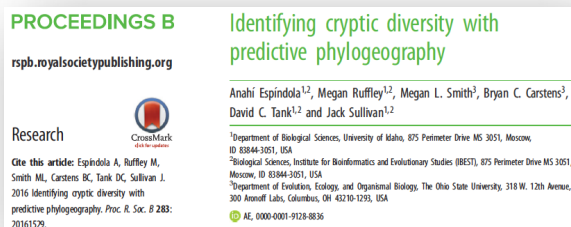Repeating this process many times and using a consensus tree produces the best decision tree.
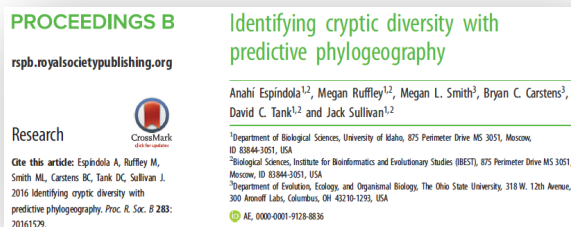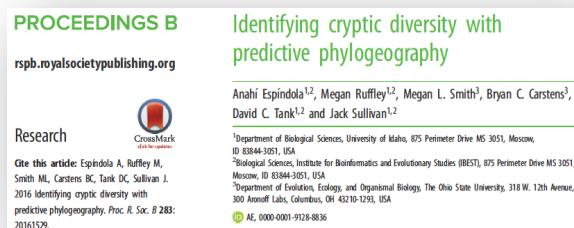
# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

## 4. Evaluation



Table 1. Prediction accuracies (in %), based on the full, the downsampled and the resampled datasets. Values indicate overall and category-based (i.e. cryptic versus non-cryptic) accuracies.

| dataset | overall | cryptic | non-cryptic |
|---|---|---|---|
| **PNW** | | | |
| full | 98.78 | 98.52 | 100.00 |
| downsampling | 98.78 | 98.52 | 100.00 |
| resampling 141 | 77.52 | 83.14 | 51.78 |
| resampling 1500 | 98.78 | 98.52 | 100.00 |
| resampling 4500 | 98.78 | 98.52 | 100.00 |
| resampling 9000 | 98.78 | 98.52 | 100.00 |

Identifying cryptic diversity with predictive phylogeography

Anahí Espíndola[1,2], Megan Ruffley[1,2], Megan L. Smith[3], Bryan C. Carstens[3], David C. Tank[1,2] and Jack Sullivan[1,2]

[1]Department of Biological Sciences, University of Idaho, 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[2]Biological Sciences, Institute for Bioinformatics and Evolutionary Studies (IBEST), 875 Perimeter Drive MS 3051, Moscow, ID 83844-3051, USA
[3]Department of Evolution, Ecology, and Organismal Biology, The Ohio State University, 318 W. 12th Avenue, 300 Aronoff Labs, Columbus, OH 43210-1293, USA
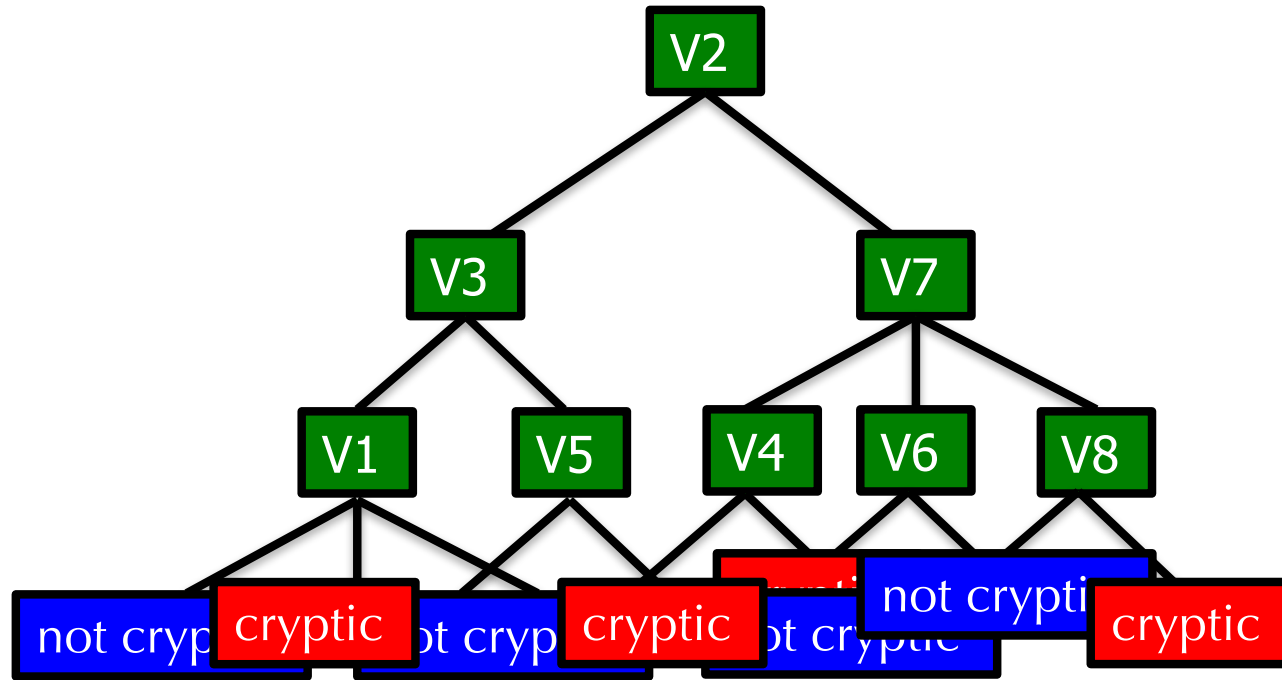
AE, 0000-0001-9128-8836

# Building a predictive framework

**Goal**: to develop a predictive framework to identify species that are likely to contain cryptic diversity.

## 5. Prediction



### (e) Predicting diversity in unknown taxa
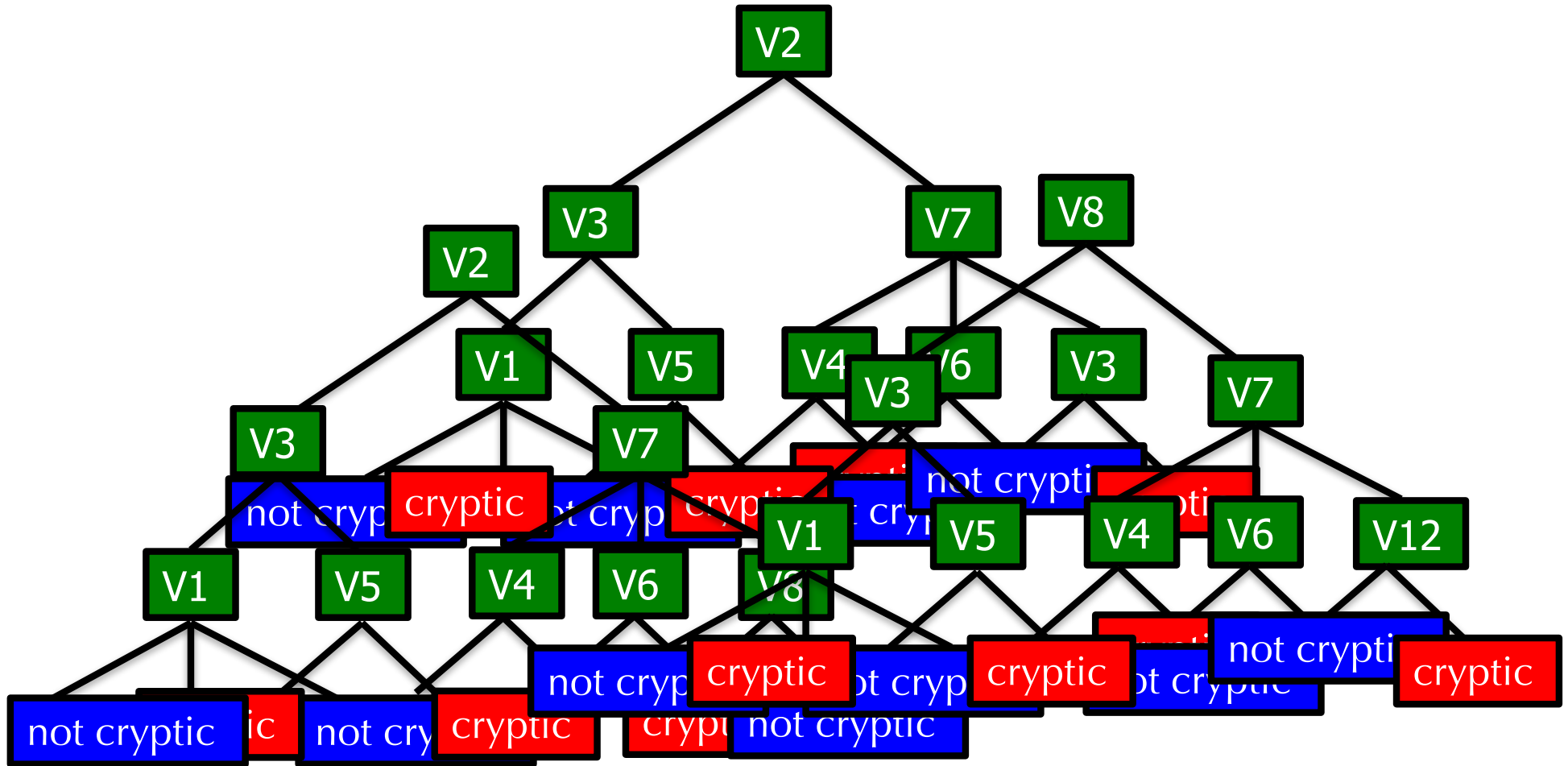
To demonstrate the application of our method, we used the RF approach to predict the presence or absence of cryptic diversity in a set of taxa for which the presence of cryptic diversity has not been assessed with genetic data, so that we could prioritize future work. We assessed three taxa per biome; the three taxa from the PNW (i.e. red alder *Alnus rubra*, Western red cedar *Thuja plicata* and robust lancetooth *Haplotrema vancouverense*) were predicted to lack cryptic diversity with relatively high probabilities (98.06%, 97.91% and 98.24%, respectively). Two of the three taxa selected from the SAL (Costa's hummingbird, *C. costae* and the desert woodrat, *N. lepida*) were predicted to contain cryptic diversity, whereas the Gila woodpecker *M. uropygiales* was predicted to lack cryptic diversity (55.28%, 68.48% and 51.23%, respectively). Interestingly, *N. lepida* has been recently shown to possess cryptic diversity based on published revisionary data [72].

# Building a predictive framework

*Haplotrema vancouverense* was predicted to be non-cryptic absent genetic data…

**Megan Smith** used SNPs from RADseq, developed a novel approach to analyzing these data, and confirmed that the snail was non-cryptic.

## Demographic model selection using random forests and the site frequency spectrum

Megan L. Smith[1] | Megan Ruffley[2,3] | Anahí Espíndola[2,3] | David C. Tank[2,3] |
Jack Sullivan[2,3] | Bryan C. Carstens[1]

# Predictive frameworks
# **recycle** and **repurpose** existing data.



Publications by Year

Web of Science search term = 'phylogeograph*'

Program notes & review papers…

## Web of Science search term = 'phylogeograph*'

…and LOTS of data papers.

Citations

Regression of Citations on Year

$r^2$ ~0.99

mode = 3

median = 8

Bryan Carstens - OSU EEOB

# Ecologists and Evolutionary Biologists…

- …want to learn about interesting ecosystems and species

- … hope to understand how biodiversity evolves

# ...should use big data!

1. **aggregate** available data on a global scale



2. **analyze** these data using predictive modeling

WorldClim - Global Climate Data
Free climate data for ecological modeling and GIS

19 bioclimatic variables at ~1km global resolution

NCBI    Resources    How To

GenBank            Nucleotide

>200,000,000 sequences

Global Biodiversity Information Facility
Free and Open Access to Biodiversity Data

| 651,297,925 | 1,634,951 | 18,828 |
|---|---|---|
| OCCURRENCES | SPECIES | DATASETS |

**Greg Wheeler** helped develop the initial versions of the scripts to aggregate available data.

# Phylogatr

**Tara Pelletier** developed the complete set of Python & R scripts to aggregate available data.

GBIF

BIOCLIM

GenBank

GPS

Format data

Sequence alignments

**Data analysis**

Bryan Carstens - OSU EEOB

# Big Data!

**Phylogatr** (Tara Pelletier)

- 561,534 – georeferenced sequences
- 42,206 – species w/ georeferenced sequence data
- 12,266 – sequence alignments
- 10,991 – species with alignments

Bryan Carstens - OSU EEOB

Bryan Carstens - OSU EEOB

# Big Data!

**Phylogatr** (w/ Tara Pelletier)

- **Global processes** (structure of genetic diversity)
- **Classic questions on global scales** (response to climate change)
- **Quantifying biodiversity** (species limits in major clades)

# What factors promote intraspecific genetic structure?

genetic distance

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

Wright 1943

- **Isolation by distance (IBD)**

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

Wright 1943

- **Isolation by distance (IBD)**

- **Isolation by environment (IBE)**

# What factors promote intraspecific genetic structure?

Wright 1943

- **Isolation by distance (IBD)**

- **Isolation by environment (IBE)**

- correlation within species between environment and geography on a global scale ($r = 0.77$)

- **IBD/E**: multiple matrix regression with randomization (Wang 2013)

*environmental or genetic distance* (vertical axis)

*geographic distance* (horizontal axis)

3/27/18    Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | P-value Geo. | prop.sig Env. | P-value Env. |
|---|---|---|---|---|---|
| Fungi | 23 | 0.04 | | | |
| Mosses | | | | | |
| Ferns | | | | | |
| Gymnosperms | | | | | |
| Angiosperms | | | | | |
| Arthropods | | | | | |
| Vertebrates | | | | | |
| Annelida | | | | | |
| Cnidaria | | | | | |
| Echinodermata | | | | | |
| Mollusca | | | | | |
| Nematoda | | | | | |
| Platyhelminthes | | | | | |
| Total | | | | | |

## each species tested at $P = 0.05$

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | **P-value Geo.** | prop.sig Env. | P-value Env. |
|---|---|---|---|---|---|
| Fungi | 23 | 0.04 | **0.69** | | |
| Mosses | | | | | |
| Ferns | | | | | |
| Gymnosperms | | | | | |
| Angiosperms | | | | | |
| Arthropods | | | | | |
| Vertebrates | | | | | |
| Annelida | | | | | |
| Cnidaria | | | | | |
| Echinodermata | | | | | |
| Mollusca | | | | | |
| Nematoda | | | | | |
| Platyhelminthes | | | | | |
| **Total** | | | | | |

**Exact Binomial test:** Is the proportion of species that are isolated by distance higher than expected by chance?

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | P-value Geo. | **prop.sig Env.** | **P-value Env.** |
|---|---|---|---|---|---|
| Fungi | 23 | 0.04 | 0.69 | **0.04** | **0.69** |
| Mosses | | | | | |
| Ferns | | | | | |
| Gymnosperms | | | | | |
| Angiosperms | | | | | |
| Arthropods | | | | | |
| Vertebrates | | | | | |
| Annelida | | | | | |
| Cnidaria | | | | | |
| Echinodermata | | | | | |
| Mollusca | | | | | |
| Nematoda | | | | | |
| Platyhelminthes | | | | | |
| Total | | | | | |

Because we're using a **multiple matrix regression with randomization** both IBD and IBE are considered.

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | P-value Geo. | prop.sig Env. | P-value Env. |
|---|---|---|---|---|---|
| Fungi | 23 | | | | |
| Mosses | 10 | | | | |
| Ferns | 7 | | | | |
| Gymnosperms | 111 | | | | |
| Angiosperms | 870 | | | | |
| Arthropods | 6015 | | | | |
| Vertebrates | 2723 | | | | |
| Annelida | 33 | | | | |
| Cnidaria | 6 | | | | |
| Echinodermata | 14 | | | | |
| Mollusca | 44 | | | | |
| Nematoda | 6 | | | | |
| Platyhelminthes | 15 | | | | |
| **Total** | **9877** | | | | |

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | P-value Geo. | prop.sig Env. | P-value Env. |
|---|---|---|---|---|---|
| Fungi | 23 | 0.04 | | 0.04 | |
| Mosses | 10 | 0 | | 0 | |
| Ferns | 7 | 0 | | 0 | |
| Gymnosperms | 111 | 0.07 | | 0.06 | |
| Angiosperms | 870 | 0.1 | | 0.1 | |
| Arthropods | 6015 | 0.15 | | 0.13 | |
| Vertebrates | 2723 | 0.29 | | 0.21 | |
| Annelida | 33 | 0.21 | | 0.15 | |
| Cnidaria | 6 | 0.5 | | 0 | |
| Echinodermata | 14 | 0.21 | | 0.21 | |
| Mollusca | 44 | 0.16 | | 0.16 | |
| Nematoda | 6 | 0.33 | | 0.33 | |
| Platyhelminthes | 15 | 0 | | 0.2 | |
| **Total** | **9877** | **0.19** | | **0.15** | |

# What factors promote intraspecific genetic structure?

| Group | # datasets | prop.sig Geo. | P-value Geo. | prop.sig Env. | P-value Env. |
|---|---|---|---|---|---|
| Fungi | 23 | 0.04 | 0.69 | 0.04 | 0.69 |
| Mosses | 10 | 0 | 1 | 0 | 1 |
| Ferns | 7 | 0 | 1 | 0 | 1 |
| Gymnosperms | 111 | 0.07 | 0.19 | 0.06 | 0.32 |
| Angiosperms | 870 | 0.1 | <0.01 | 0.1 | <0.01 |
| Arthropods | 6015 | 0.15 | <0.01 | 0.13 | <0.01 |
| Vertebrates | 2723 | 0.29 | <0.01 | 0.21 | <0.01 |
| Annelida | 33 | 0.21 | <0.01 | 0.15 | 0.02 |
| Cnidaria | 6 | 0.5 | <0.01 | 0 | 1 |
| Echinodermata | 14 | 0.21 | 0.03 | 0.21 | 0.03 |
| Mollusca | 44 | 0.16 | 0.01 | 0.16 | 0.01 |
| Nematoda | 6 | 0.33 | 0.03 | 0.33 | 0.03 |
| Platyhelminthes | 15 | 0 | 1 | 0.2 | 0.0362 |
| **Total** | **9877** | **0.19** | **<0.01** | **0.15** | **<0.01** |

Bryan Carstens - OSU EEOB

What explains this variation in IBD/E across biological groups?

# Building predictive frameworks for big data analysis

**Isolation by distance / environment analyses grouped categorically by result:**

- not significant / significant (at the species level)

- data table containing general information about organisms and environment

# Building predictive frameworks for big data analysis

**Data table: 33 variables used in machine learning analysis**

- environmental characteristics (canopy cover, wetlands, habitat type)

- organismal traits (metabolism, taxonomy, type of gene)

- geographic (max. latitude, range area, mid-point latitude)

Random Forest Analysis

Proportion significant for IBDE

0%　　　　　　　　　　　　　　　　　100%

Fungi　Mosses　Ferns　Gymnosperms　Angiosperms　Insects　Other Arthropods　Fish　Amphibia　Reptile　Birds　Bats　Terrestrial Mammals*　Marine Mammals*　Annelida　Cnidaria　Echinodermata　Mollusca　Nematoda　Platyhelminthes

# What factors promote intraspecific genetic structure?

Variable performance quantified by measuring the mean decrease in accuracy (MDA) of the predictive function that occurs when that variable is omitted from analysis.

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

Variable performance quantified by measuring the mean decrease in accuracy (MDA) of the predictive function that occurs when that variable is omitted from analysis.



**Mean Decrease in Accuracy**

Legend:
- Geographic
- Intrinsic
- Environmental
- Sampling

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

| Variable | mean with IBD | mean without IBD | t-test p-value |
|----------|---------------|------------------|----------------|
| Area (km$^2$) | $6.57 \times 10^6$ | $2.98 \times 10^6$ | $3.05 \times 10^{-11}$ |
| Minimum distance from equator | 30.04982 | 31.88308 | 2.92E-06 |
| Mid-point latitude of range | 31.86497 | 33.42718 | 0.0009313 |
| Length of latitude° | 14.66759 | 9.74694 | 2.20E-16 |

# Geography!

Bryan Carstens - OSU EEOB

# What factors promote intraspecific genetic structure?

- If IBD/E is a precursor to local adaptation, organismal traits represent evolved responses to aspects of the environment

- more precise organismal traits are needed…

# What factors promote intraspecific genetic structure?

- If IBD/E is a precursor to local adaptation, organismal traits represent evolved responses to aspects of the environment
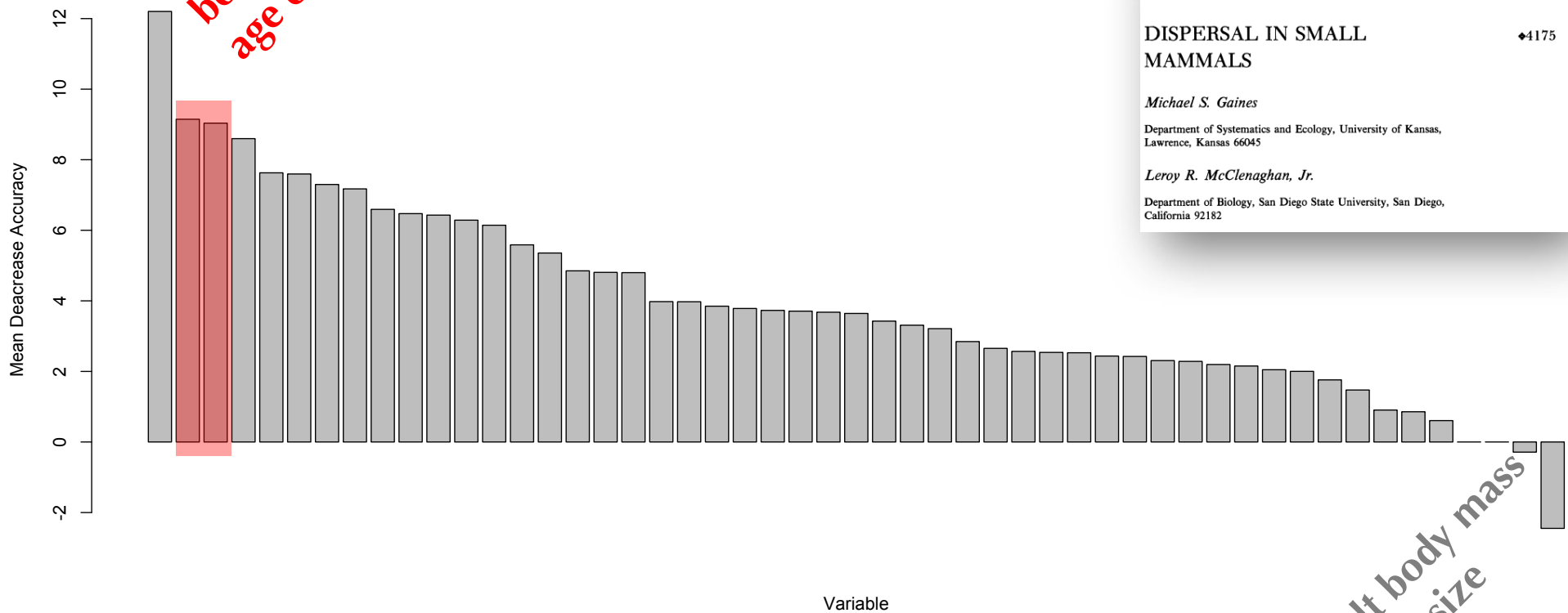
- more precise organismal traits are needed...

**panTHERIA database of Class Mammalia**
**(Jones et al. 2009)**

- 55 organismal traits by 4630 species.

- repeat the RF analysis using only mammals (954 species)

Bryan Carstens - OSU EEOB

**INVITED REVIEW**

### Advances in our understanding of mammalian sex-biased dispersal

L. J. LAWSON HANDLEY*† and N. PERRIN†

*Theoretical and Molecular Population Genetics Group, Department of Genetics, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK, †Department of Ecology and Evolution, University of Lausanne, CH-1015 Lausanne, Switzerland

### DISPERSAL IN SMALL MAMMALS ◆4175

*Michael S. Gaines*

Department of Systematics and Ecology, University of Kansas, Lawrence, Kansas 66045

*Leroy R. McClenaghan, Jr.*

Department of Biology, San Diego State University, San Diego, California 92182

n (number of samples)
body mass at weaning
age of sexual maturity

adult body mass
group size

# What factors promote intraspecific genetic structure?

1. on a global scale, environmental and geographic distance are broadly correlated within species

2. geographic attributes such as maximum latitude and range size are the best predictors of which species are likely to exhibit IBD/E

3. organismal traits may be difficult to compare across the Tree of Life

Traditional comparative phylogeography to particular regions:

- SE US (*Avise 2000*)
- Europe (*Hewitt 2000*)
- Pacific Northwest of NA (*Carstens et al. 2005*)

One traditional goal of phylogeographic work has been to understand how particular species respond to large scale climatic shifts (e.g., such as that of the end Pleistocene).

Traditional comparative phylogeography to particular regions:

- SE US (*Avise 2000*)
- Europe (*Hewitt 2000*)
- Pacific Northwest of NA (*Carstens et al. 2005*)

One traditional goal of phylogeographic work has been to understand how particular species respond to large scale climatic shifts (e.g., such as that of the end Pleistocene).
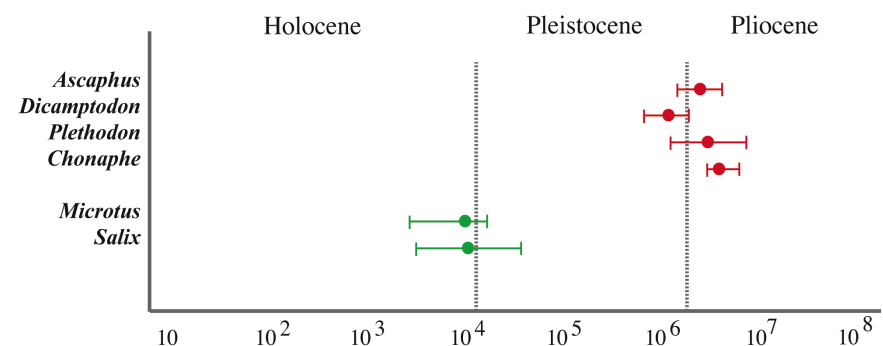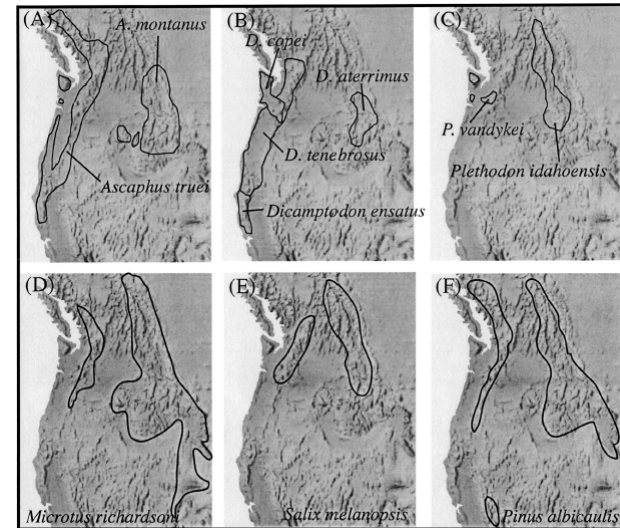
How did bats respond to climate change….
    …on a global scale

- downloaded >30,000 sequences from 123 species with greater than 15 georeferenced samples

- **A**pproximate **B**ayesian **C**omputation used to calculate the probability of two models (expansion, bottleneck) in all species

- species distribution modeling to compare predicted size of current range to predicted range size at end Pleistocene (thanks to Ariadna Morales!)
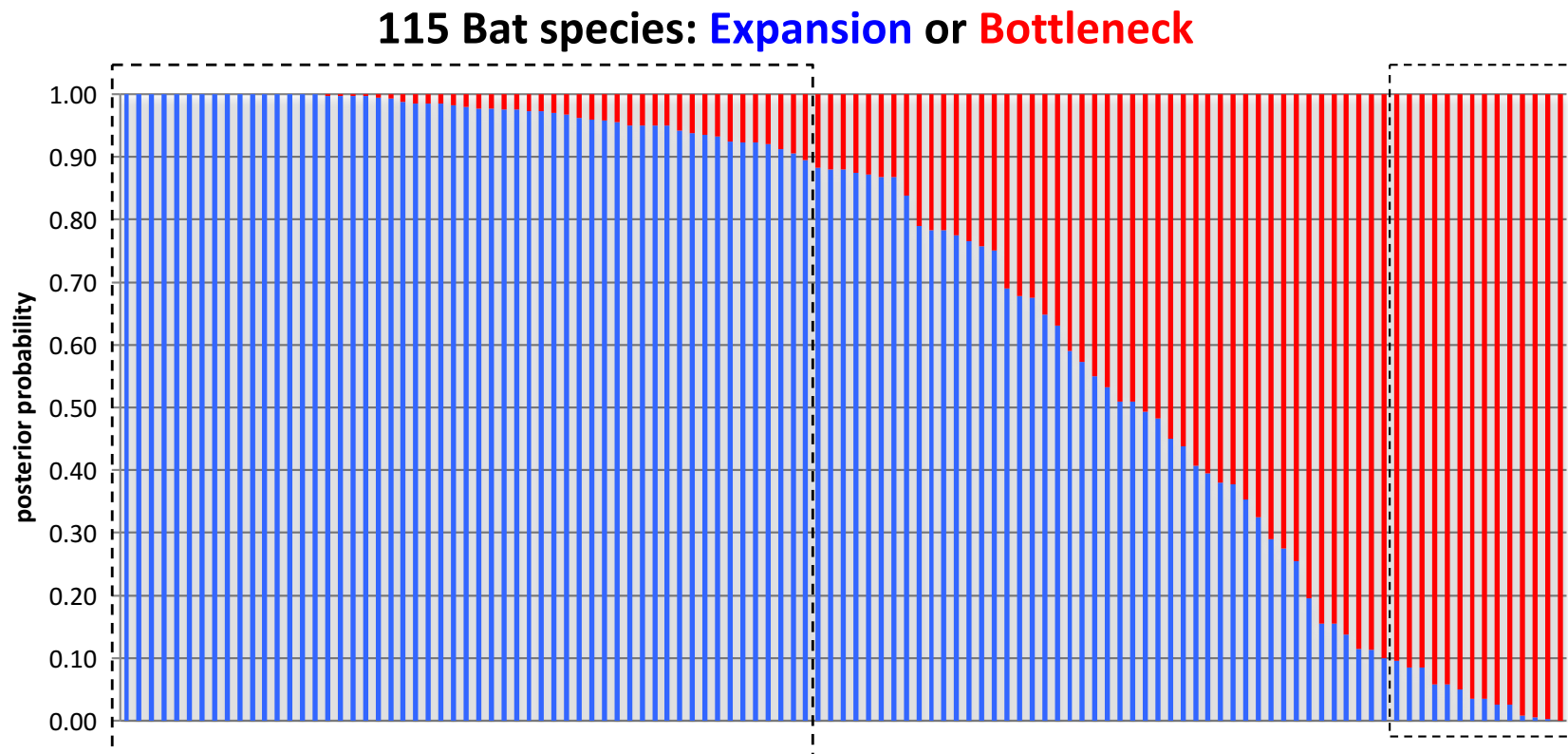
# How did bats respond to climate change on a global scale?

**115 Bat species: <span style="color:blue">Expansion</span> or <span style="color:red">Bottleneck</span>**

## How did bats respond to climate change on a global scale?

- Adult body mass is significantly correlated with $PP_{expansion}$.

- Expansion species are nearly twice the size as bottleneck species (19.9 : 10.8)

- Related to dietary niche, also significantly correlated.

Bryan Carstens - OSU EEOB

**How did bats respond to climate change on a global scale?**

**Organismal Traits**
- body size
- wing shape
- breeding strategy
- roosting location
- dietary niche

**Environmental**
- predicted size of current range
- predicted size of range at LGM
- maximum latitude
- mid point latitude
- average temperature in observed range

Identify factors (intrinsic, environmental) that predict the observed response:

# How did bats respond to climate change on a global scale?

| RF P=0.9 | OOB error | Expansion error | Bottleneck error | n | n Expansion | n Bottleneck |
|---|---|---|---|---|---|---|
| all spatial variables | 0.2115 | 0.089 | 1.0 | 52 | 45 | 7 |
| change spatial variables | 0.1923 | 0.067 | 1.0 | 52 | 45 | 7 |
| all variables | 0.1522 | 0.000 | 1.0 | 46 | 39 | 7 |
| other variables | 0.1379 | 0.000 | 1.0 | 58 | 50 | 8 |

| RF P=0.7 | OOB error | Expansion error | Bottleneck error | n | n Expansion | n Bottleneck |
|---|---|---|---|---|---|---|
| all spatial variables | 0.292 | 0.121 | 1.0 | 72 | 58 | 14 |
| change spatial variables | 0.278 | 0.121 | 0.9 | 72 | 58 | 14 |
| all variables | 0.194 | 0.020 | 1.0 | 62 | 51 | 11 |
| other variables | 0.190 | 0.015 | 1.0 | 79 | 65 | 14 |

Random forest prediction error rates unacceptably high due to disparity in response variables…

Most available DNA sequence data lack georeferencing.

(previous analysis based on 13% of total mtDNA data)

1.  downloaded all mtDNA from bats

2.  aligned by gene
    - ~ 20,000 barcoding loci (10,421 *cyt b* seqs, 9552 *COI* seqs)
    - 842 nominal species (75% of described)
    - 1116 total species in Chiroptera (Wilson & Reeder, 3rd Ed)

3.  estimated distributions of gene trees (by family)

4.  used GMYC model (Pons et al. 2006) to estimate number of cryptic bat species

General **M**ixed **Y**ule – **C**oalescent model (Pons et al. 2006)

# General Mixed Yule – Coalescent model (Pons et al. 2006)

- similar to *lineage through time plots* in that it considers rates of branching in rooted ultrametric trees
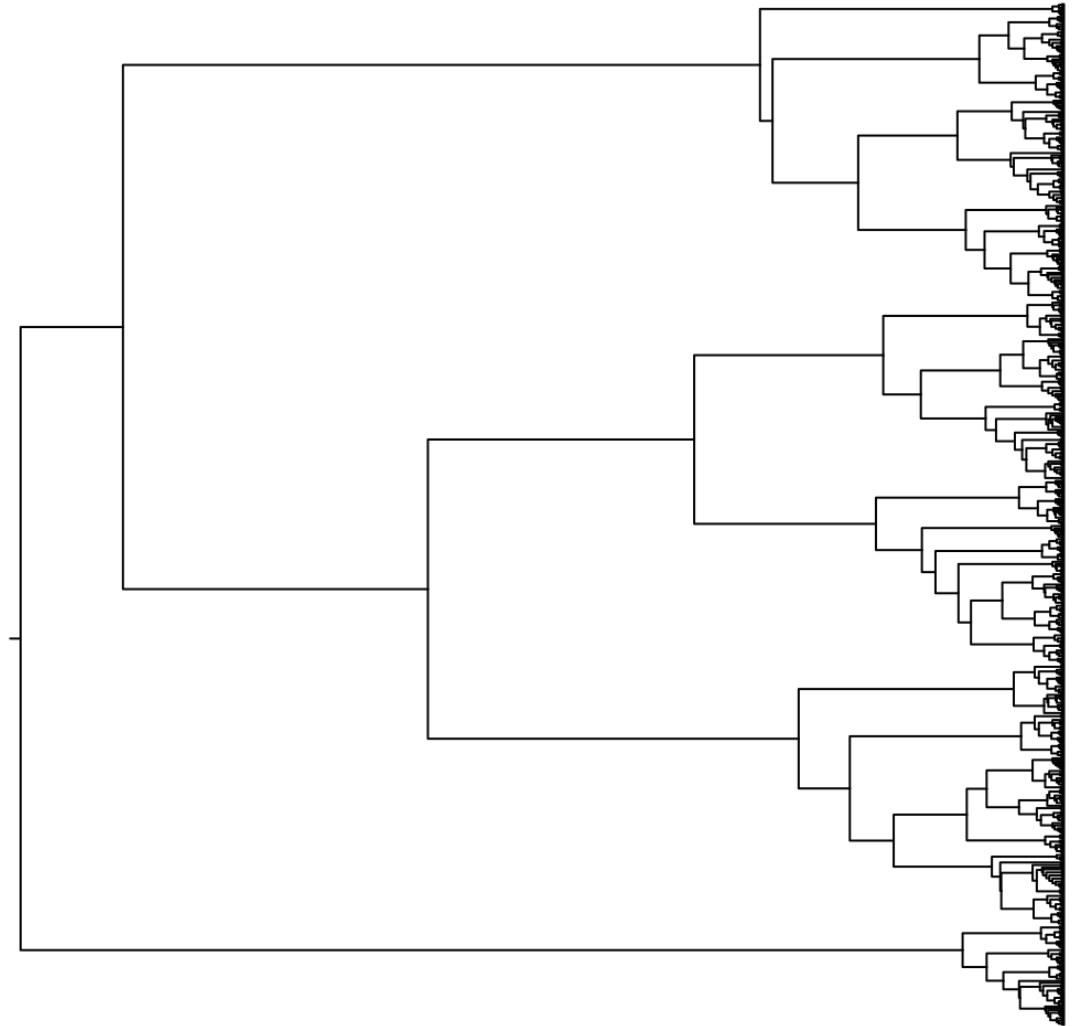
# General Mixed Yule – Coalescent model (Pons et al. 2006)

- similar to *lineage through time plots* in that it considers rates of branching in rooted ultrametric trees

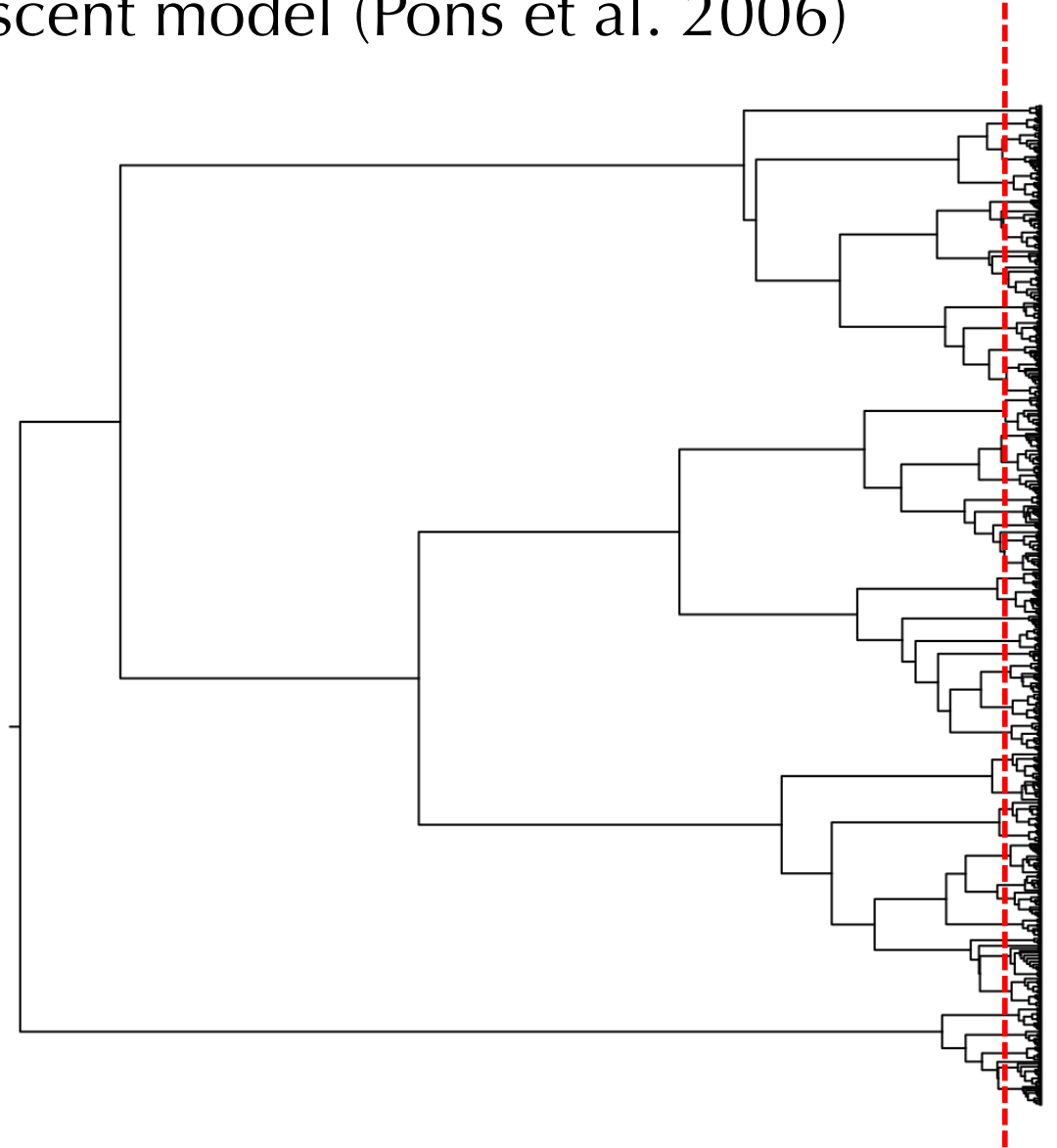- attempts to find point where rate of branching transitions from slow to fast

## General Mixed Yule – Coalescent model (Pons et al. 2006)

- similar to *lineage through time plots* in that it considers rates of branching in rooted ultrametric trees

- attempts to find point where rate of branching transitions from slow to fast

- assumes rate of speciation is slow compared to rate of allele coalescence

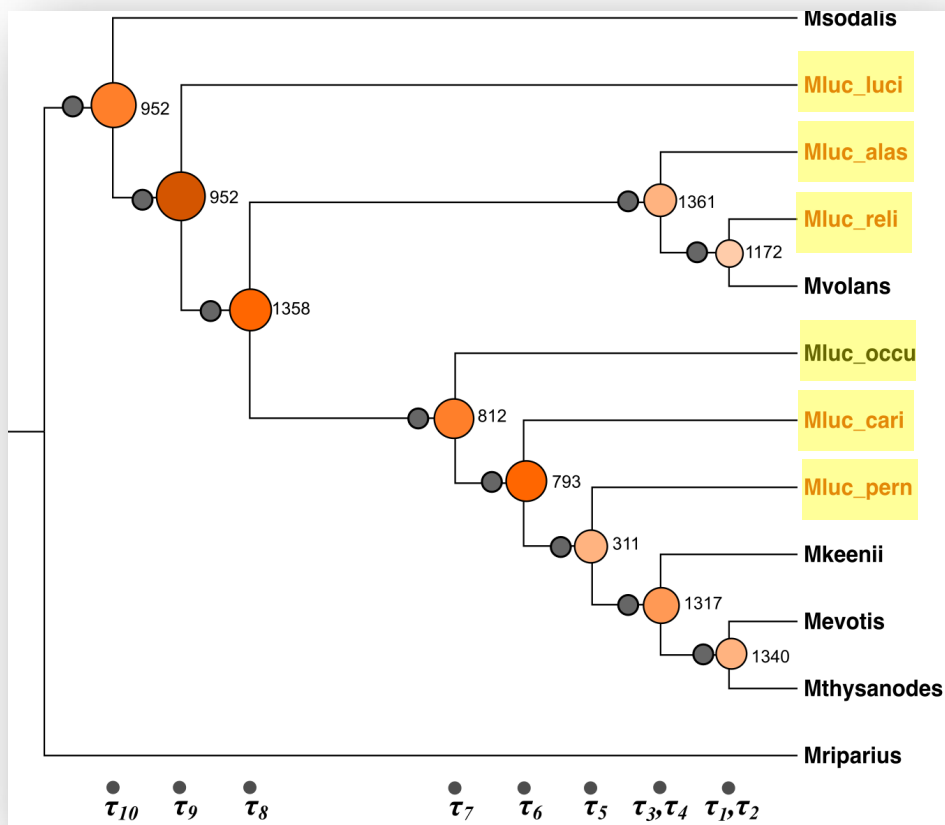Species Diversity - Chiroptera

2161 discrete GMYC entities were detected.

GMYC has been criticized as being biased towards overestimation (e.g., Esselstyn et al. 2012), but…
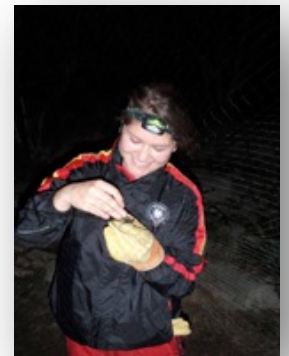
…follow up investigations often confirm GMYC results.

## Species Diversity - Chiroptera



- *Myotis lucifugus* contains multiple described subspecies.

- Our results delimit 4 GMYC entities in *M. lucifugus*

- Morales et al. (in rev) collected ~800 UCE loci from *Myotis* bats.
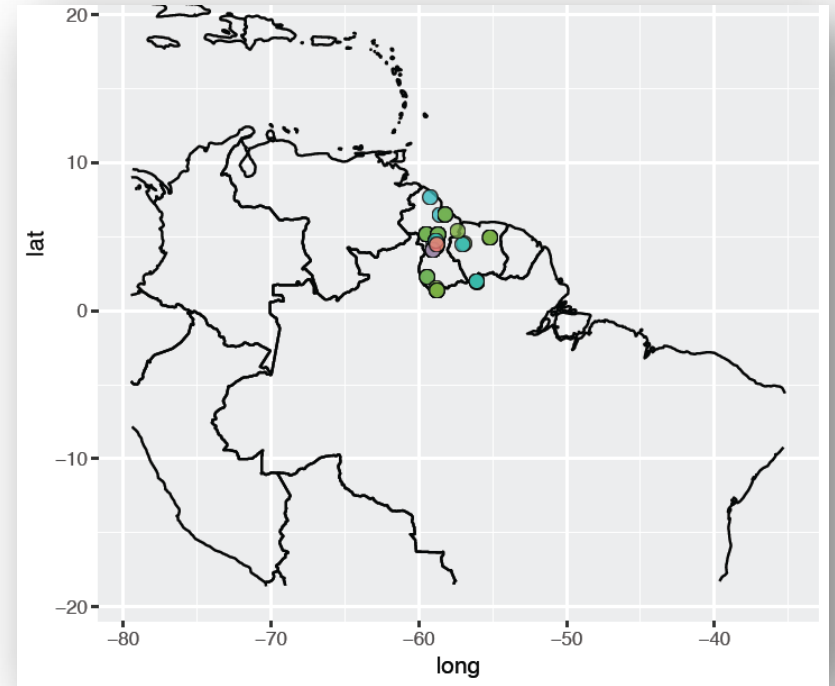
- *M. lucifugus* subspecies are not monophyletic…

data from Ariadna's DDIG….

# Species Diversity - Chiroptera

Geographic data are available from 332 nominal species, and 134 contained data for >1 of the GMYC entities.

- GMYC entities contradicted by geographic distribution of samples (0.13).
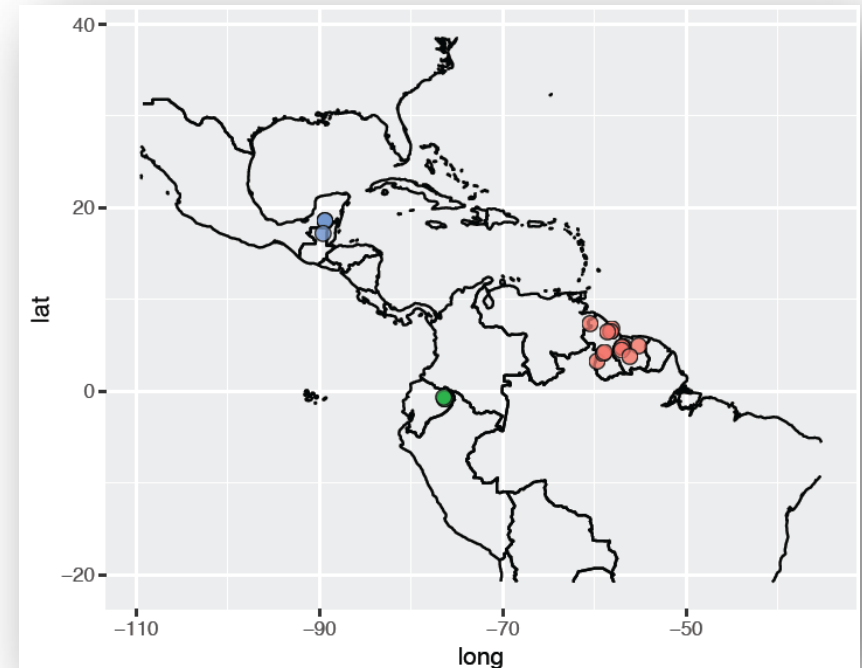


contradicts GMYC: *Artibeus concolor*

Bryan Carstens - OSU EEOB

# Species Diversity - Chiroptera

Geographic data are available from 332 nominal species, and 134 contained data for >1 of the GMYC entities.

- GMYC entities contradicted by geographic distribution of samples (0.13).

- GMYC entities corresponded to discrete geographic clusters for all (0.29) or some (0.31) of the delimited groups
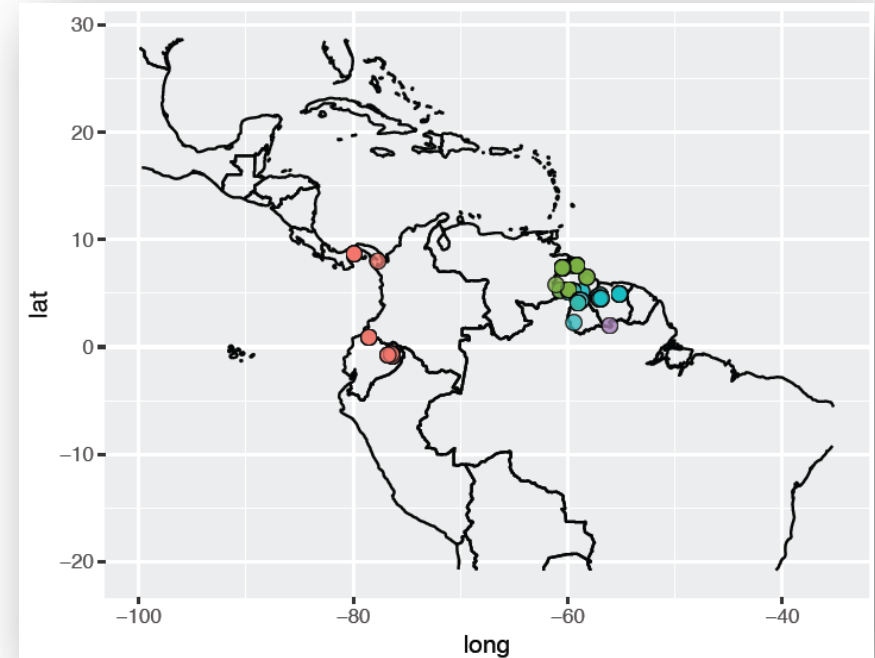


all discrete: *Eptesicus furinalis*

## Species Diversity - Chiroptera

Geographic data are available from 332 nominal species, and 134 contained data for >1 of the GMYC entities.

- GMYC entities contradicted by geographic distribution of samples.

- GMYC entities corresponded to discrete geographic clusters for all (0.29) or some (0.31) of the delimited groups
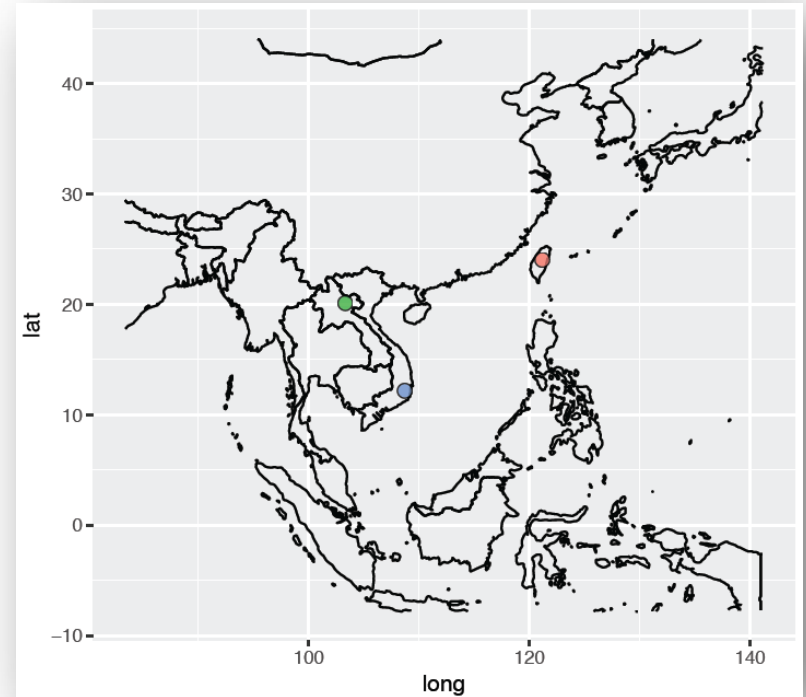


some discrete: *Carollia brevicada*

# Species Diversity - Chiroptera

Geographic data are available from 332 nominal species, and 134 contained data for >1 of the GMYC entities.

- GMYC entities contradicted by geographic distribution of samples.

- GMYC entities corresponded to discrete geographic clusters for all (0.29) or some (0.31) of the delimited groups

- Sampling inadequate to draw conclusions (0.27)



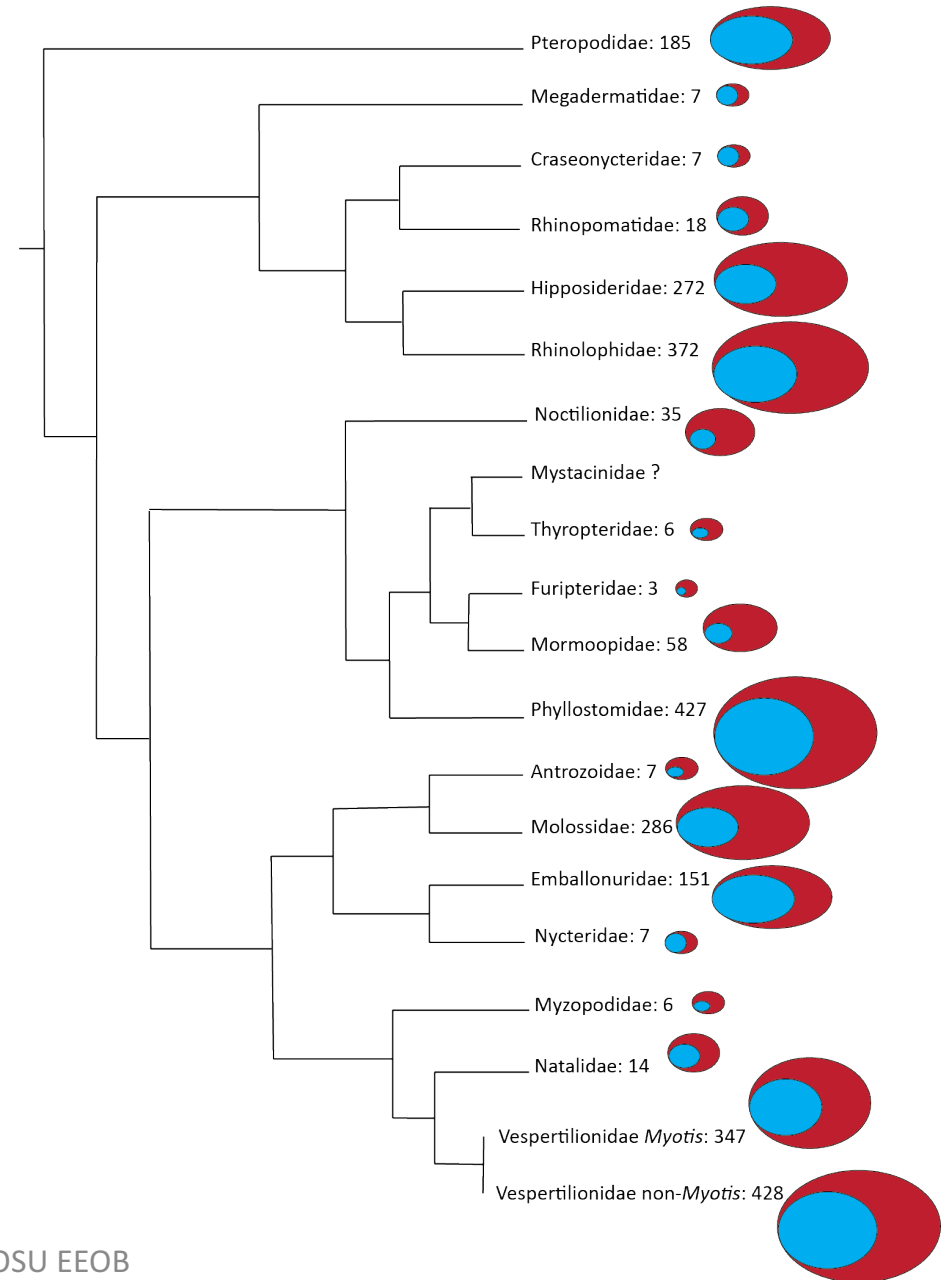Inadequate sampling: *Coelops frithii*

## Species Diversity - Chiroptera
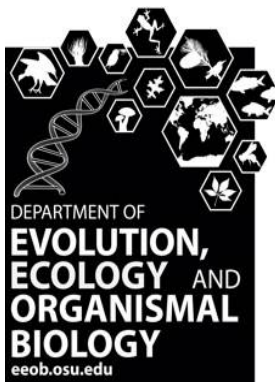
**Nominal** / **Predicted**

- 842 nominal species in analysis
- 1116 described in Chiroptera
- 75% of species represented

- **2073 GMYC entities**

- **~2700 bat species!?!**

# Big data!

*Conclusions*

> ~300 years of taxonomic work
>
> ~ 6 million genetic data points
>
> ~600 million occurrence records on GBIF

- NSF Proposal to fund **phylogatr**: database aggregator w/ R pipelines to facilitate meta-analysis of phylogeographic data.

DEPARTMENT OF
**EVOLUTION,
ECOLOGY** AND
**ORGANISMAL
BIOLOGY**
eeob.osu.edu

**Ohio Supercomputer Center**
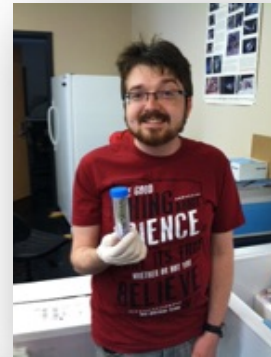An **OH·TECH** Consortium Member

T · H · E
**OHIO
STATE
UNIVERSITY**

**Postdocs**
Margaret Koopman
Yi-Hsin Erica Tsai
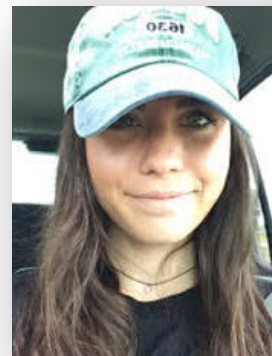Amanda Zellmer
Tereza Thomé
Michael Gruenstaeudl
Tara Pelletier

**Grad Students**
Sarah Hird
Noah Reid
John McVay
Tara Pelletier
Jordan Satler
Ariadna Morales
Greg Wheeler
Megan Smith
Cole Thompson

**Undergrads**
Danielle Fuselier
Holly Stoute
Dan Ence
Matt Demarest
Maxim Kim
Edwin Rice
Ivy Larkin
Katy Field

**NSF Funding**
DBI - 1661029
DEB - 1701810
DBI - 1560116
DEB - 1501474
DEB - 1457519
DEB - 1403034
DEB - 1500774
DEB - 1257784
OISE -1118408
DEB - 0918212
DEB - 0956069

3/27/18

Bryan Carstens - OSU EEOB

NSF