# Generalizing Bayesian phylogenetics to infer shared evolutionary events

**Jamie Oaks**

Auburn University

phyletica.org

🐦 @jamoaks

phyletica.org/slides/lib.pdf

# Phyletica Lab

## The Phyleticians



### Postdocs
- Perry Wood, Jr
- *Brian Folt*
- *Jesse Grismer*

### Graduate students
- Tashitso Anamza
- Matt Buehler
- Kerry Cobb
- Kyle David
- Randy Klabacka
- Morgan Muell
- Tanner Myers
- Claire Tracy
- *Branna Sipley*
- *Aundrea Westfall*

### Undergraduate students
- Laura Lewis
- Mary Wells
- Hailey Whitaker
- Noah Yawn
- *Charlotte Benedict*
- *Eric Carbo*
- *Ryan Cook*
- *Andrew DeSana*
- *Miles Horne*
- *Jacob Landrum*
- *Nadia L'Bahy*
- *Jorge Lopez-Perez*
- *Holden Smith*
- *Virginia White*
- *Kayla Wilson*

*The last 5 years*

 Generalizing Bayesian phylogenetics to infer shared evolutionary events

*The next 5 years*

▶ My vision for a position at the LIB

- Phylogenetics is rapidly becoming the statistical foundation of biology
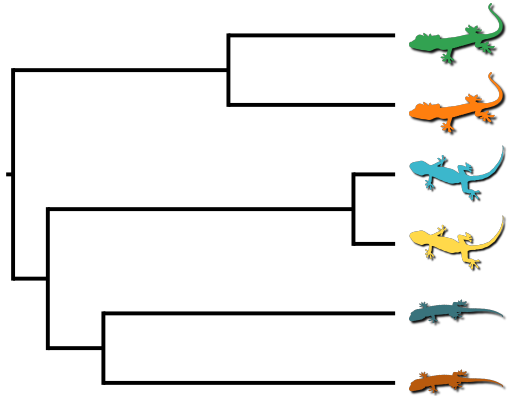


© 2007 Boris Kulikov boris-kulikov.blogspot.com

- Phylogenetics is rapidly becoming the statistical foundation of biology
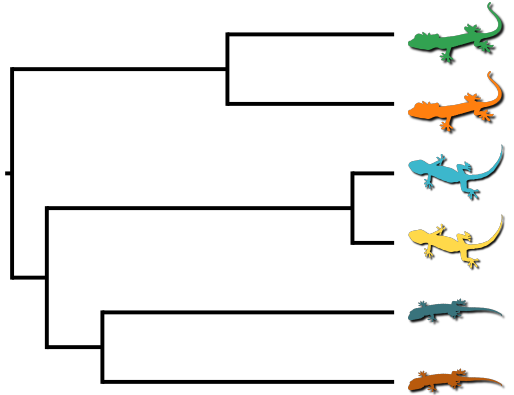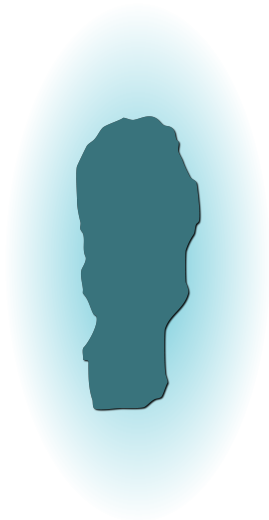- "Big data" present exciting possibilities and challenges
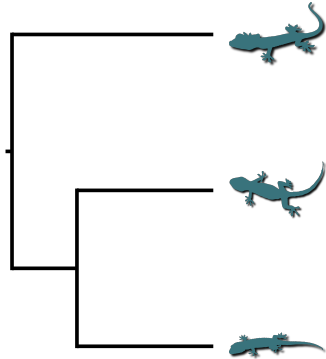


© 2007 Boris Kulikov boris-kulikov.blogspot.com

- Phylogenetics is rapidly becoming the statistical foundation of biology
- "Big data" present exciting possibilities and challenges
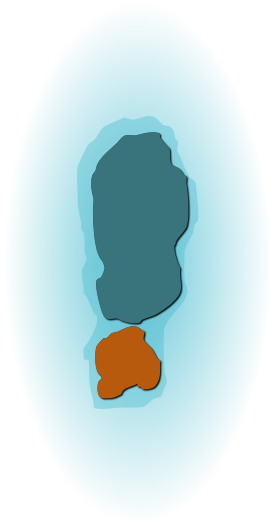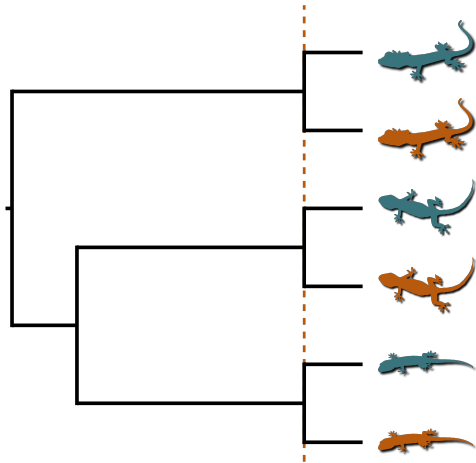- Many opportunities to develop new ways to study biology in light of phylogeny

- ▶ **Assumption:** All processes of diversification affect each lineage independently
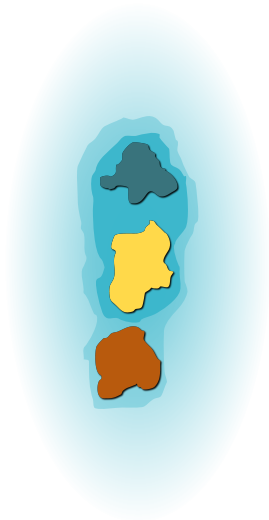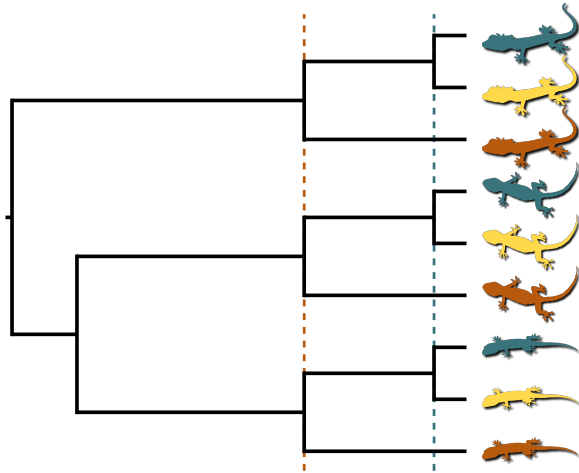
J. R. Oaks (2019). *Systematic Biology* 68: 371–395

J. R. Oaks, C. D. Siler, and R. M. Brown (2019). *Evolution* 73: 1151–1167

J. R. Oaks (2019). *Systematic Biology* 68: 371–395

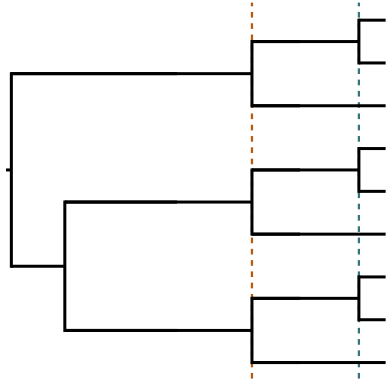J. R. Oaks, C. D. Siler, and R. M. Brown (2019). *Evolution* 73: 1151–1167

J. R. Oaks (2019). *Systematic Biology* 68: 371–395

J. R. Oaks, C. D. Siler, and R. M. Brown (2019). *Evolution* 73: 1151–1167

## Biogeography

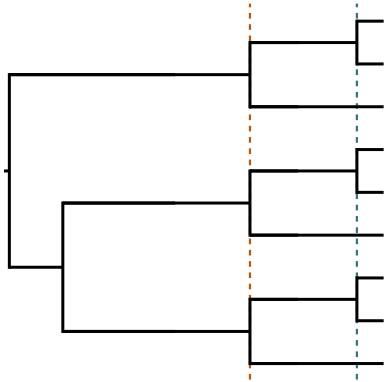▶ Environmental changes that affect whole
communities of species

## Biogeography

▶ Environmental changes that affect whole communities of species

## Genome evolution

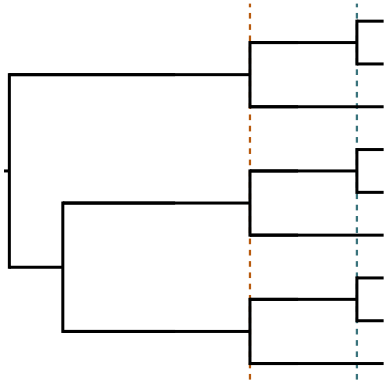▶ Duplication of a chromosome segment harboring gene families

**Biogeography**

► Environmental changes that affect whole communities of species

**Genome evolution**

► Duplication of a chromosome segment harboring gene families

**Epidemiology**

► Transmission at social gatherings

**Biogeography**

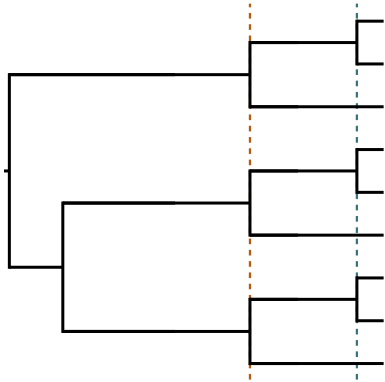▶ Environmental changes that affect whole communities of species

**Genome evolution**

▶ Duplication of a chromosome segment harboring gene families

**Epidemiology**

▶ Transmission at social gatherings

**Endosymbiont evolution** (e.g., parasites, microbiome)
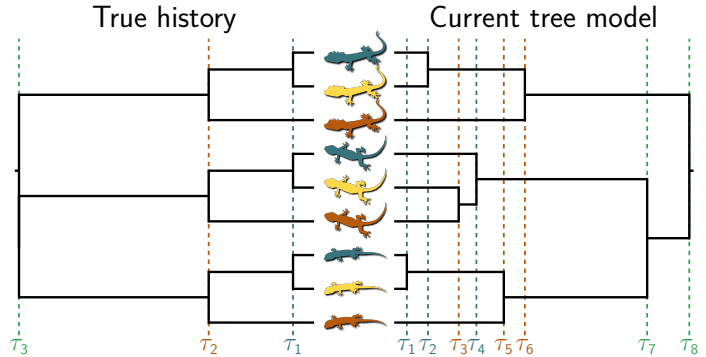
▶ Speciation of the host

▶ Co-colonization of new host species
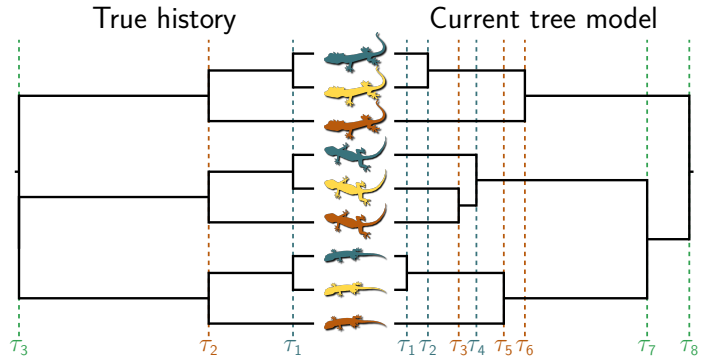
# Why account for shared divergences?

# Why account for shared divergences?

1. Improve inference



True history    Current tree model

$\tau_3$    $\tau_2$    $\tau_1$    $\tau_1 \tau_2$ $\tau_3 \tau_4$ $\tau_5 \tau_6$    $\tau_7$ $\tau_8$

# Why account for shared divergences?

1. Improve inference

2. **Provide a framework for studying processes of co-diversification**

**Biogeography**

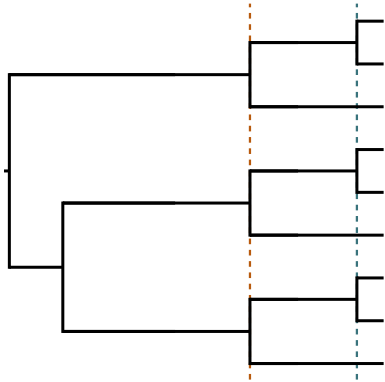▶ Environmental changes that affect whole communities of species

**Genome evolution**

▶ Duplication of a chromosome segment harboring gene families

**Epidemiology**

▶ Transmission at social gatherings
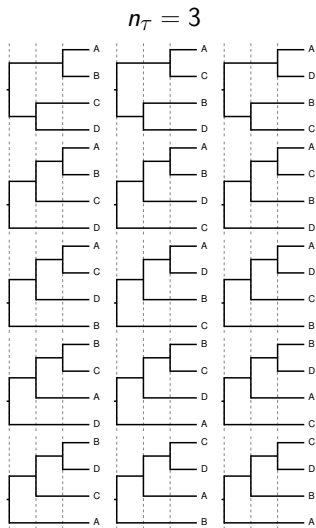
**Endosymbiont evolution** (e.g., parasites, microbiome)

▶ Speciation of the host

▶ Co-colonization of new host species
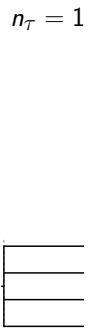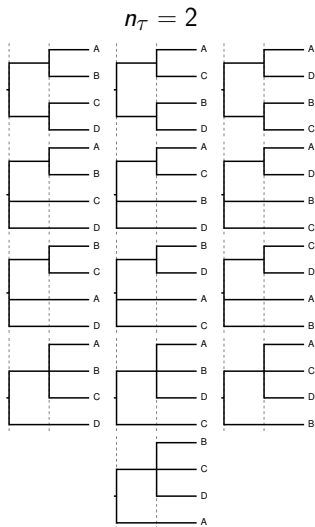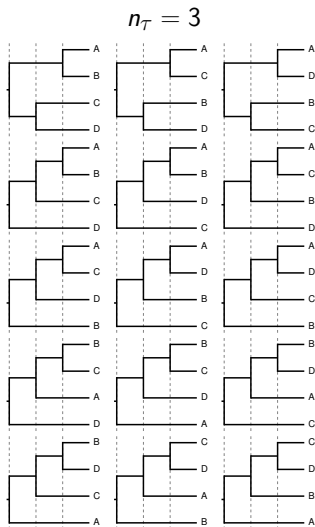
# Generalizing tree space

# Generalizing tree space



$n_\tau = 3$

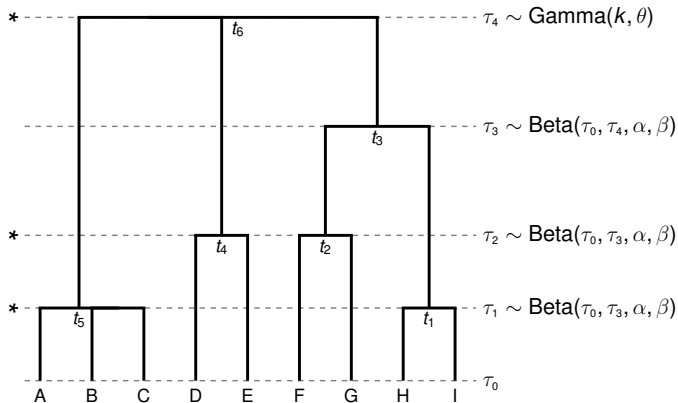# Generalizing tree space

# Generalized tree distribution

- All topologies equally probable
- Parametric distribution on age of root
- Beta distributions on other div times



$\tau_4 \sim \text{Gamma}(k, \theta)$

$\tau_3 \sim \text{Beta}(\tau_0, \tau_4, \alpha, \beta)$

$\tau_2 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$

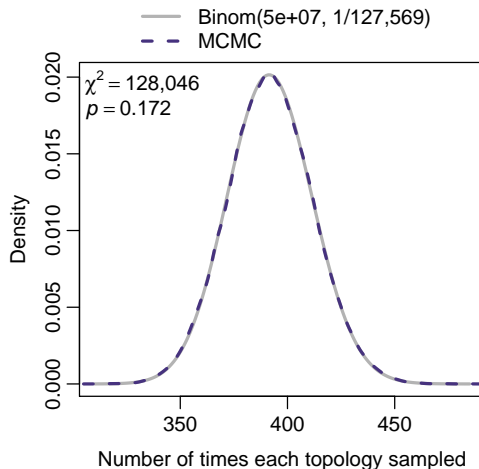$\tau_1 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$

$\tau_0$

# Inferring trees with shared divergences



Reversible-jump MCMC

# Validating rjMCMC with 7-leaf tree



**The rjMCMC algorithms sample the expected generalized tree distribution**

J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

**Phyco Eval**

**Phylogenetic coevality**

**Ecoevolity**

**E**stimating **evo**lutionary **coevality**

▶ **Tree model**

  ▶ rjMCMC sampling of generalized tree distribution

[1] D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

## *Phyco⊑val*
**Phylogenetic coevality**
J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

## *Ecoevolity*
**E**stimating **evo**lutionary **coevality**
J. R. Oaks (2019). *Systematic Biology* 68: 371–395

▶ **Tree model**
  ▶ rjMCMC sampling of generalized tree distribution

▶ **Likelihood model**
  ▶ CTMC model of characters evolving along genealogies
  ▶ Infer species trees by analytically integrate over genealogies[1]

[1] D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

**PhycoEval**

**Phylogenetic coevality**

J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

**Ecoevolity**

**E**stimating **evo**lutionary **coevality**

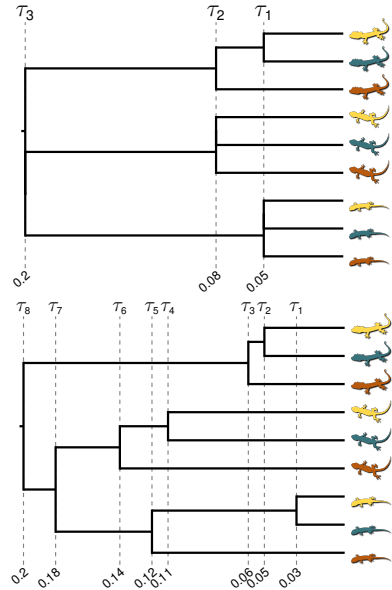J. R. Oaks (2019). *Systematic Biology* 68: 371–395

- ▶ **Tree model**
    - ▶ rjMCMC sampling of generalized tree distribution

- ▶ **Likelihood model**
    - ▶ CTMC model of characters evolving along genealogies
    - ▶ Infer species trees by analytically integrate over genealogies[1]

- ▶ *Goal: Co-estimation of phylogeny and shared divergences from genomic data*

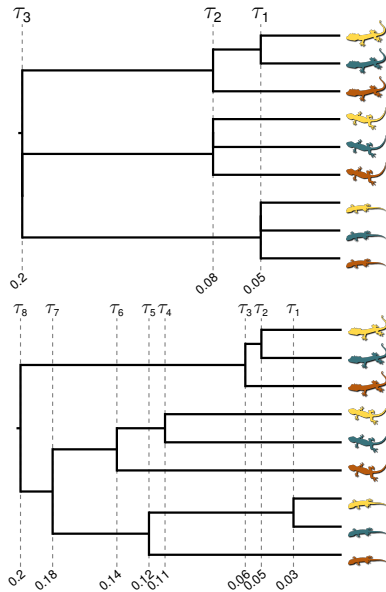[1] D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

# Methods: Simulations
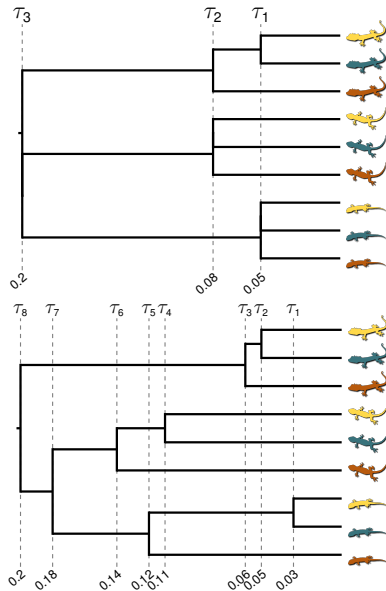
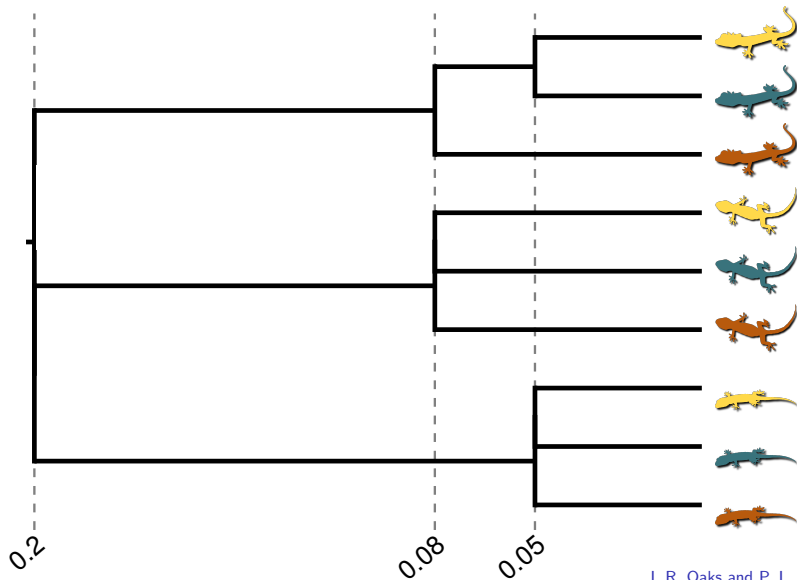▶ Simulated 100 data sets with 50,000 characters

# Methods: Simulations

- Simulated 100 data sets with 50,000 characters
- Analyzed each data set with:
    - $M_G$ = Generalized tree model
    - $M_{IB}$ = Independent-bifurcating tree model

# Methods: Simulations

▶ Simulated 100 data sets with 50,000 characters

▶ Analyzed each data set with:
  ▶ $M_G$ = Generalized tree model
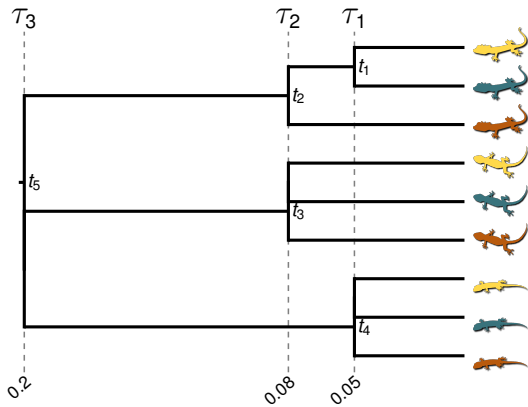  ▶ $M_{IB}$ = Independent-bifurcating tree model

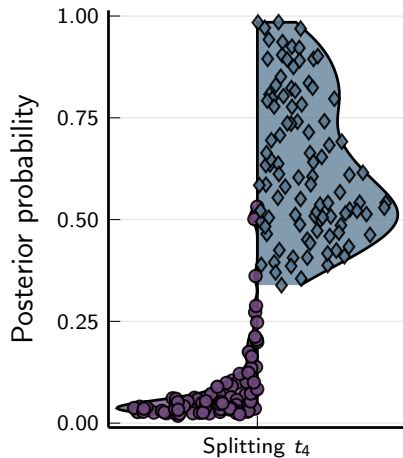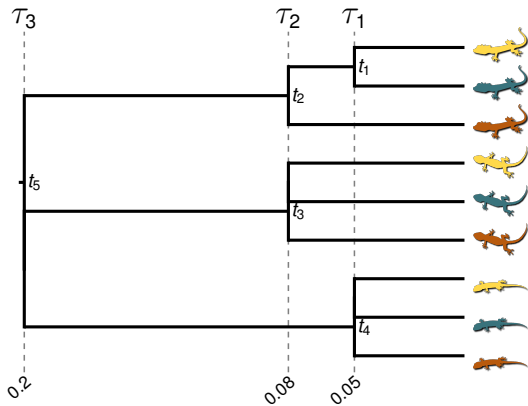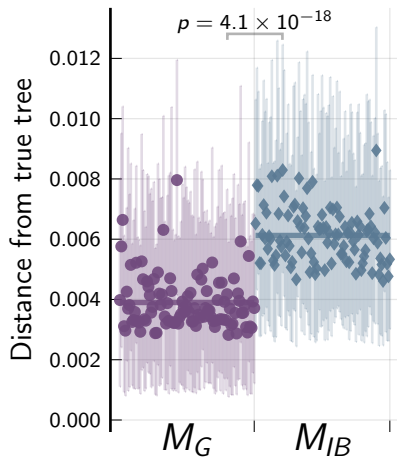▶ Simulated 100 data sets where topology and div times randomly drawn from $M_G$ and $M_{IB}$
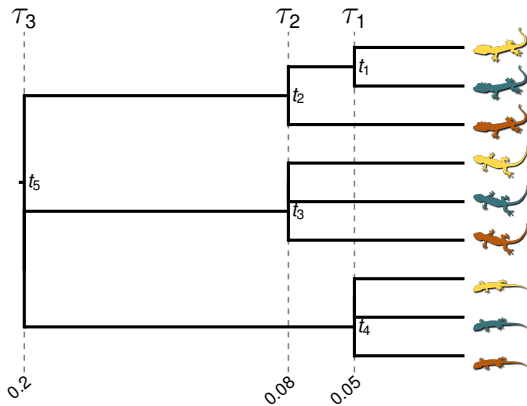
0.2     0.08    0.05

J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

$M_G$ = Generalized model  $M_{IB}$ = Independent-bifurcating model

**$M_G$ significantly better at inferring trees with shared divergences**

$M_G$ = Generalized model   $M_{IB}$ = Independent-bifurcating model
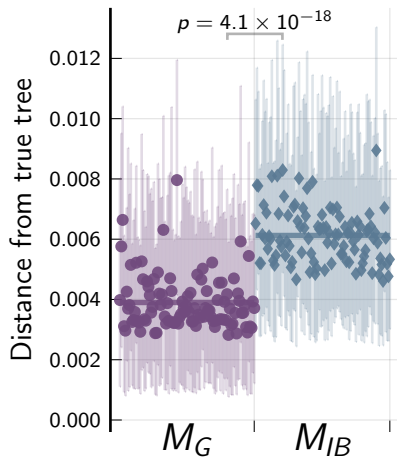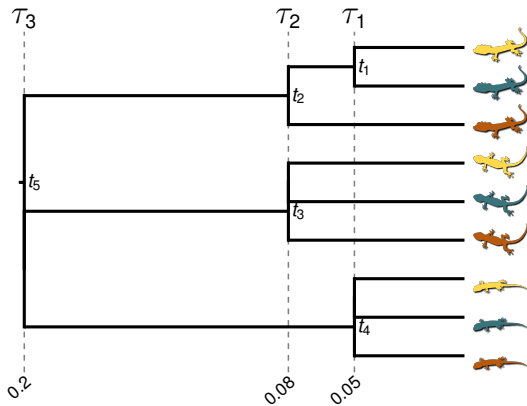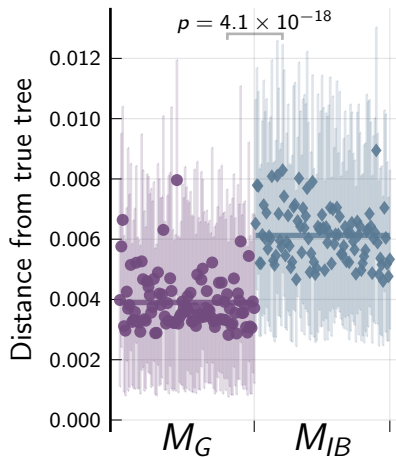
$M_G$ = Generalized model  $M_{IB}$ = Independent-bifurcating model

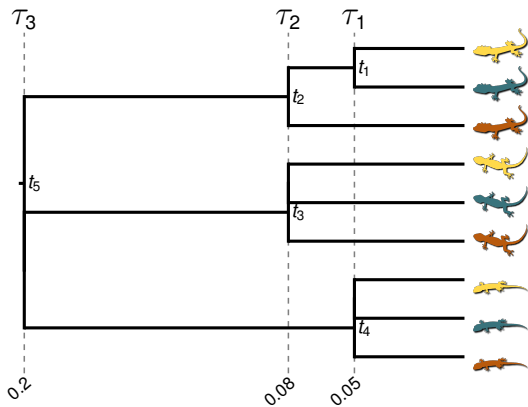$M_G$ = Generalized model   ◆ $M_{IB}$ = Independent-bifurcating model

$M_G$ **performs as well as true model when divergences are independent**

## Results: random $M_G$ trees

# Results: random $M_G$ trees

**$M_G$ performs well with data simulated on random trees with shared divergences**

Probability of incorrectly merged div times (true model $= M_{IB}$)

Probability of incorrectly merged div times (true model $= M_{\text{IB}}$)



**$M_G$ has low false positive rate**

$M_G$ = Generalized model          $M_{IB}$ = Independent-bifurcating model
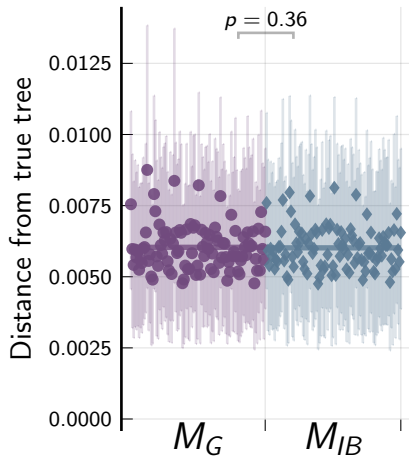
**Generalizing tree space improves MCMC convergence and mixing**

J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

**Did fragmentation of islands promote diversification?**

## Cyrtodactylus



©Rafe M. Brown

## Gekko



©Rafe M. Brown

J. R. Oaks, C. D. Siler, and R. M. Brown (2019). *Evolution* 73: 1151–1167

## Cyrtodactylus


©Rafe M. Brown

1702 loci
155,887 sites

## Gekko


©Rafe M. Brown

1033 loci
94,813 sites

J. R. Oaks, C. D. Siler, and R. M. Brown (2019). *Evolution* 73: 1151–1167

# Gekko

# Cyrtodactylus



J. R. Oaks and P. L. Wood, Jr. (2021). *bioRxiv*

## Take-home points

► We can accurately infer phylogenies with shared divergences with moderately sized data sets

## Take-home points

▶ We can accurately infer phylogenies with shared divergences with moderately sized data sets

▶ Generalizing tree space avoids spurious support and improves MCMC mixing

## Take-home points

▶ We can accurately infer phylogenies with shared divergences with moderately sized data sets

▶ Generalizing tree space avoids spurious support and improves MCMC mixing

▶ Among Philippine gekkonids, we found support for shared divergences predicted by sea-level changes

## Open science: everything is available...

**Software:**
- Phycoeval:
  https://github.com/phyletica/ecoevolity
  (release coming soon)

**Open-Science Notebooks:**
- Phycoeval analyses: https://github.com/phyletica/phycoeval-experiments
- Gecko RADseq:
  https://github.com/phyletica/gekgo



phyletica.org/codiv-sanger-bake-off

# Vision for LIB position

**Phylogenetic theory/methods**

▶ Develop process-based and trait-dependent distributions over the space of generalized trees

**Empirical work**

▶ Did the evolution of habitat preference affect the diversification of bent-toed geckos?

▶ Epidemiological dynamics of "super-spreading" events during the COVID-19 pandemic

**Teaching**

▶ Coding to learn evolution

# Generalized tree distribution

▶ Our current distribution over trees is
motivated by mathematical convenience

# Generalized tree distribution

▶ Our current distribution over trees is motivated by mathematical convenience

▶ A process-based distribution would allow us to learn about parameters that control diversification processes



$\tau_4 \sim \text{Gamma}(k, \theta)$

$\tau_3 \sim \text{Beta}(\tau_0, \tau_4, \alpha, \beta)$

$\tau_2 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$

$\tau_1 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$

$\tau_0$

# Generalized tree distribution

- ▶ Our current distribution over trees is motivated by mathematical convenience
- ▶ A process-based distribution would allow us to learn about parameters that control diversification processes
- ▶ Goal: port $M_G$ algorithms to RevBayes and develop generalized birth-death model



$$\tau_4 \sim \text{Gamma}(k, \theta)$$

$$\tau_3 \sim \text{Beta}(\tau_0, \tau_4, \alpha, \beta)$$

$$\tau_2 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$$

$$\tau_1 \sim \text{Beta}(\tau_0, \tau_3, \alpha, \beta)$$

$$\tau_0$$

# Generalized tree distribution

▶ Our current distribution over trees is motivated by mathematical convenience

▶ A process-based distribution would allow us to learn about parameters that control diversification processes

▶ Goal: port $M_G$ algorithms to RevBayes and develop generalized birth-death model

Sebastian Höhna
LMU Munich

# Generalizing the birth-death process

**Birth-death basics:**
- ▶ Lineages speciate at rate $\lambda$
- ▶ Lineages go extinct at rate $\mu$
- ▶ We sample extant lineages with probability $\rho$

## Generalizing the birth-death process

**Birth-death basics:**

▶ Lineages speciate at rate $\lambda$
▶ Lineages go extinct at rate $\mu$
▶ We sample extant lineages with probability $\rho$

**"Birth-death-burst" (BDB) process:**

▶ Include "burst events" that occur at rate $\lambda_\beta$
▶ Each lineage diverges with probability $\beta$



$t_\beta$

**Birth-death basics:**

- Lineages speciate at rate $\lambda$
- Lineages go extinct at rate $\mu$
- We sample extant lineages with probability $\rho$

**"Birth-death-burst" (BDB) process:**

- Include "burst events" that occur at rate $\lambda_\beta$
- Each lineage diverges with probability $\beta$
- Allow $\lambda$, $\mu$, $\lambda_\beta$, & $\beta$ to vary depending on the traits of lineages across the tree



$t_\beta$

# Generalizing the birth-death process

**Birth-death basics:**

- Lineages speciate at rate $\lambda$
- Lineages go extinct at rate $\mu$
- We sample extant lineages with probability $\rho$

**"Birth-death-burst" (BDB) process:**

- Include "burst events" that occur at rate $\lambda_\beta$
- Each lineage diverges with probability $\beta$
- Allow $\lambda$, $\mu$, $\lambda_\beta$, & $\beta$ to vary depending on the traits of lineages across the tree
- Bayesian model-averaging to infer set of trait-dependent BDB models that best explain data



$t_\beta$

# Birth-death-burst validation

# Birth-death-burst validation



**We have correctly derived the likelihood of trees under the BDB model**

*Cyrtodactylus*
≈ 380 species

India

Indoburma

Indochina

Sundaland

Pacific Ocean

Indian Ocean

Wallacea

Oceania

From L. Grismer et al. (2021). *Diversity* 13:

# Karst endemism in *Cyrtodactylus*

▶ *Cyrtodactylus* are ecologically diverse, ranging from generalists to microhabitat specialists



L. Grismer et al. (2021). *Diversity* 13:

# Karst endemism in *Cyrtodactylus*

- *Cyrtodactylus* are ecologically diverse, ranging from generalists to microhabitat specialists
- Karst-specificity evolved 24 times



L. Grismer et al. (2021). *Diversity* 13:

# Karst endemism in *Cyrtodactylus*

- *Cyrtodactylus* are ecologically diverse, ranging from generalists to microhabitat specialists
- Karst-specificity evolved 24 times
- Comprise 25% of species despite tiny fraction of landscape being karst



L. Grismer et al. (2021). *Diversity* 13:

# Karst endemism in *Cyrtodactylus*

- *Cyrtodactylus* are ecologically diverse, ranging from generalists to microhabitat specialists
- Karst-specificity evolved 24 times
- Comprise 25% of species despite tiny fraction of landscape being karst
- Karst-specific species show remarkable levels of micro-endemism



L. Grismer et al. (2021). *Diversity* 13:

©Perry Wood, Jr.

# Why high levels of diversity and endemism on karst?

► "Rapid" fragmentation of karst habitat caused by the uplift and subsequent erosion of limestone sediment over the last 30my



L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

- "Rapid" fragmentation of karst habitat caused by the uplift and subsequent erosion of limestone sediment over the last 30my

- *E.g.*, Major river systems carved through and isolated limestone karst formations (Ayeyarwady, Chiang Mai, Mekong, Red, and Salween)



L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

- "Rapid" fragmentation of karst habitat caused by the uplift and subsequent erosion of limestone sediment over the last 30my

- *E.g.*, Major river systems carved through and isolated limestone karst formations (Ayeyarwady, Chiang Mai, Mekong, Red, and Salween)

- **Hypothesis**: The fragmentation of limestone karst habitat drove diversification of karst-specific lineages of *Cyrtodactylus*



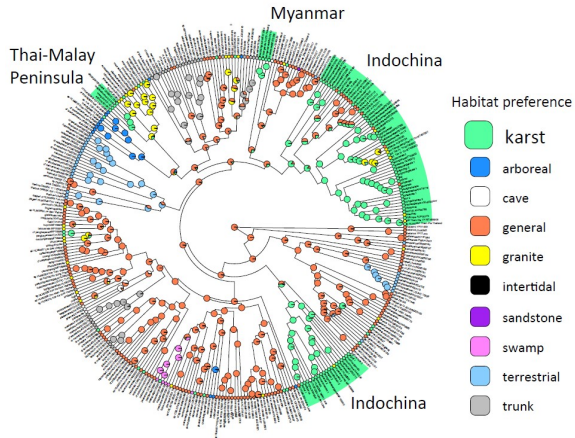L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

- "Rapid" fragmentation of karst habitat caused by the uplift and subsequent erosion of limestone sediment over the last 30my
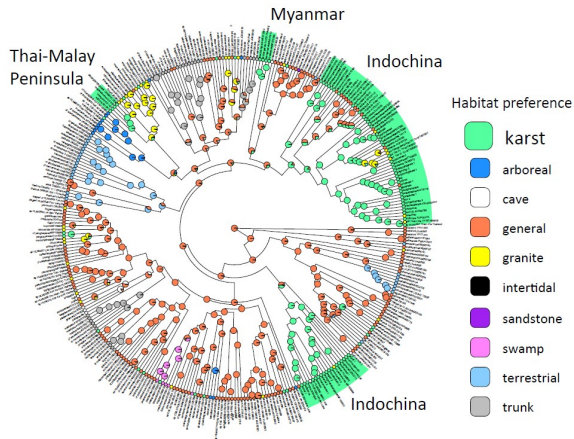
- *E.g.*, Major river systems carved through and isolated limestone karst formations (Ayeyarwady, Chiang Mai, Mekong, Red, and Salween)

- **Hypothesis**: The fragmentation of limestone karst habitat drove diversification of karst-specific lineages of *Cyrtodactylus*

- **Prediction**: Increased rate of ***shared divergences*** in karst-adapted lineages



L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

**Plan**:

▶ Access to tissue samples of 368 of the 380 *Cyrtodactylus* species



L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

**Plan**:

- ▶ Access to tissue samples of 368 of the 380 *Cyrtodactylus* species
- ▶ Sequence 5,060 UCE loci



L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

**Plan**:

- Access to tissue samples of 368 of the 380 *Cyrtodactylus* species
- Sequence 5,060 UCE loci
- Apply habitat-dependent BDB model: Model averaging to infer the posterior set of habitat-dependent models
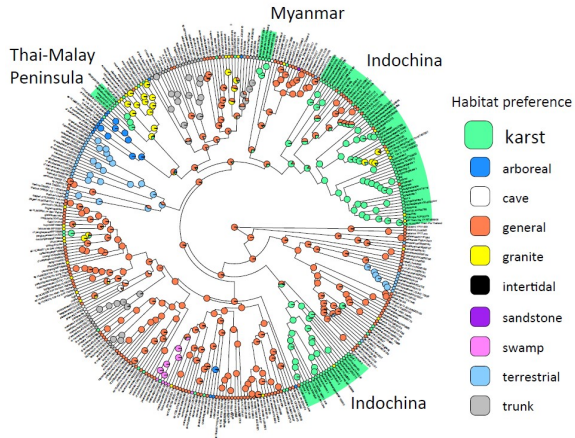


L. Grismer et al. (2021). *Diversity* 13:

# Why high levels of diversity and endemism on karst?

**Plan**:

- ▶ Access to tissue samples of 368 of the 380 *Cyrtodactylus* species
- ▶ Sequence 5,060 UCE loci
- ▶ Apply habitat-dependent BDB model: Model averaging to infer the posterior set of habitat-dependent models
- ▶ Approximate posterior probability that karst-specific lineages have higher rate of shared divergences ($\lambda_\beta$)



L. Grismer et al. (2021). *Diversity* 13:

# Epidemiological dynamics of COVID-19 pandemic

**Questions**:

▶ What is the relative contribution of social
gatherings to the spread of SARS-CoV-2?



nextstrain.org  J. Hadfield et al. (2018). *Bioinformatics* 34: 4121–4123

# Epidemiological dynamics of COVID-19 pandemic

**Questions**:

- ▶ What is the relative contribution of social gatherings to the spread of SARS-CoV-2?
- ▶ Does this vary among variants of the virus?



nextstrain.org  J. Hadfield et al. (2018). *Bioinformatics* 34: 4121–4123

# Epidemiological dynamics of COVID-19 pandemic

**Questions**:

- ▶ What is the relative contribution of social gatherings to the spread of SARS-CoV-2?
- ▶ Does this vary among variants of the virus?
- ▶ Does this increase during holidays?



nextstrain.org  J. Hadfield et al. (2018). *Bioinformatics* 34: 4121–4123

# Divergence patterns predicted by gatherings

# Divergence patterns predicted by gatherings

▶ Multiple infected people spreading
SARS-CoV-2 at a gathering will create shared
divergences across "transmission tree"

# Divergence patterns predicted by gatherings

- Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- **Shared divergences** are a good proxy for spread at gatherings

- Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- **Shared divergences** are a good proxy for spread at gatherings

**Plan**:

# Divergence patterns predicted by gatherings

- Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- **Shared divergences** are a good proxy for spread at gatherings

**Plan**:

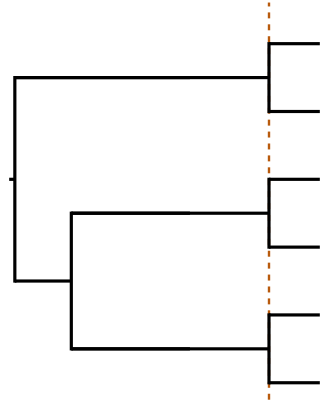- Apply strain-dependent BDB model to regional SARS-CoV-2 sequence datasets

# Divergence patterns predicted by gatherings
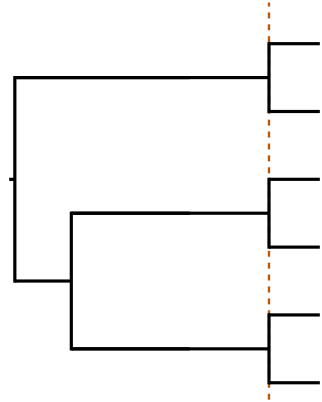
- Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- **Shared divergences** are a good proxy for spread at gatherings

**Plan**:

- Apply strain-dependent BDB model to regional SARS-CoV-2 sequence datasets
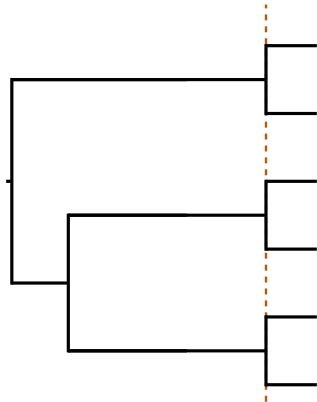- Estimate relative rate of shared divergences $(\lambda_\beta/\lambda)$

# Divergence patterns predicted by gatherings

- ▶ Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- ▶ **Shared divergences** are a good proxy for spread at gatherings

**Plan**:

- ▶ Apply strain-dependent BDB model to regional SARS-CoV-2 sequence datasets
- ▶ Estimate relative rate of shared divergences $(\lambda_\beta/\lambda)$
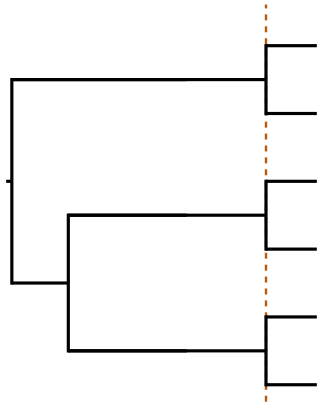- ▶ Approximate posterior probability that $\lambda_\beta$ varies among variants
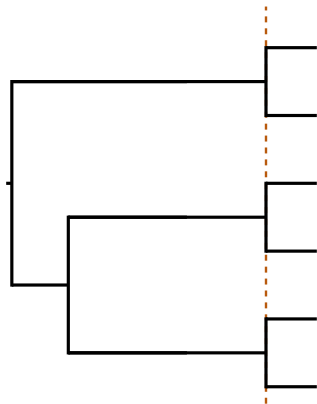
# Divergence patterns predicted by gatherings

- ▶ Multiple infected people spreading SARS-CoV-2 at a gathering will create shared divergences across "transmission tree"
- ▶ **Shared divergences** are a good proxy for spread at gatherings

**Plan**:

- ▶ Apply strain-dependent BDB model to regional SARS-CoV-2 sequence datasets
- ▶ Estimate relative rate of shared divergences $(\lambda_\beta/\lambda)$
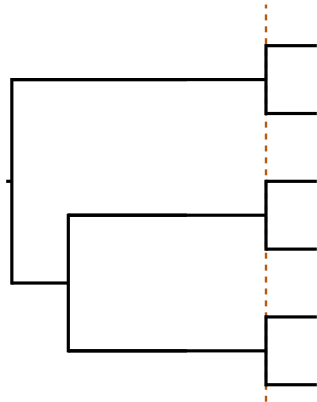- ▶ Approximate posterior probability that $\lambda_\beta$ varies among variants
- ▶ Summarize $\lambda_\beta$ over time to quantify the effect of holidays

# Teaching: Coding to learn evolution

► Develop coding-to-learn evolution course



© 2019 Philipp Messer  messerlab.org/slim

# Teaching: Coding to learn evolution

- Develop coding-to-learn evolution course
- Students use graphical modeling software, like SLiM, to gain intuition for how processes of evolution work and interact



© 2019 Philipp Messer messerlab.org/slim

# Teaching: Coding to learn evolution

- Develop coding-to-learn evolution course
- Students use graphical modeling software, like `SLiM`, to gain intuition for how processes of evolution work and interact
- **Capstone activity**: Design macroevolutionary scenarios and write code to simulate genetic data under them



© 2019 Philipp Messer messerlab.org/slim
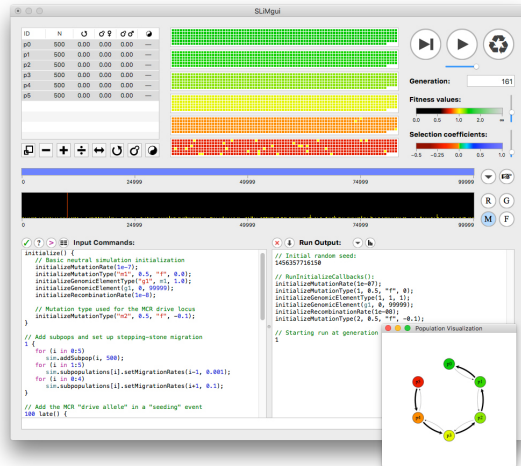
# Teaching: Coding to learn evolution

- Develop coding-to-learn evolution course
- Students use graphical modeling software, like SLiM, to gain intuition for how processes of evolution work and interact
- **Capstone activity**: Design macroevolutionary scenarios and write code to simulate genetic data under them
- Use these data for phylogenetic "bake-off"



© 2019 Philipp Messer messerlab.org/slim

# Teaching: Coding to learn evolution

- Develop coding-to-learn evolution course
- Students use graphical modeling software, like SLiM, to gain intuition for how processes of evolution work and interact
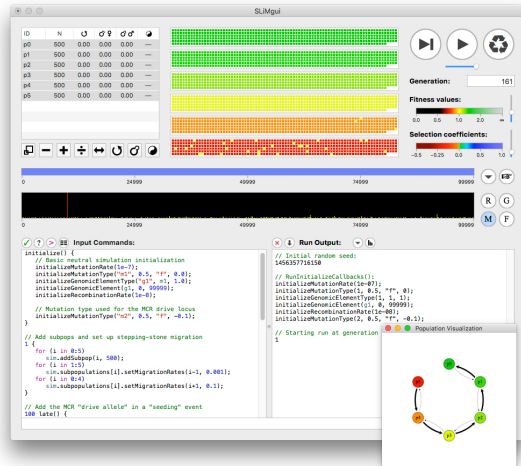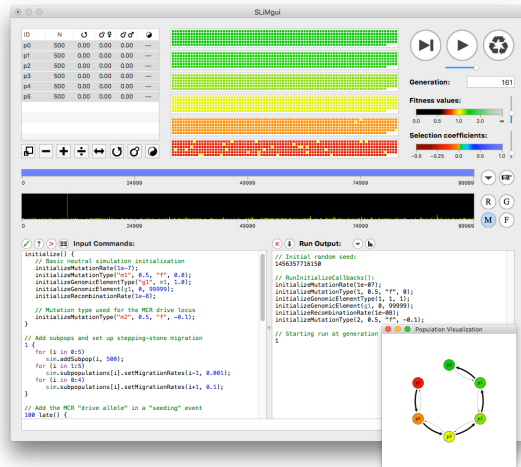- **Capstone activity**: Design macroevolutionary scenarios and write code to simulate genetic data under them
- Use these data for phylogenetic "bake-off"
- How robust are estimates under the BDB model when applied to data generated under *very* different models?
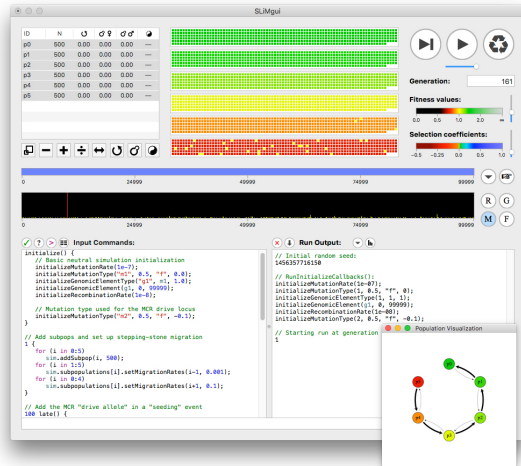


© 2019 Philipp Messer messerlab.org/slim

# Teaching: Coding to learn evolution

- Develop coding-to-learn evolution course

- Students use graphical modeling software, like SLiM, to gain intuition for how processes of evolution work and interact

- **Capstone activity**: Design macroevolutionary scenarios and write code to simulate genetic data under them

- Use these data for phylogenetic "bake-off"

- How robust are estimates under the BDB model when applied to data generated under *very* different models?

- Students co-author paper



© 2019 Philipp Messer messerlab.org/slim

## Acknowledgments

- Phyletica Lab (the Phyleticians)
- Mark Holder
- Rafe Brown
- Cam Siler
- Lee Grismer

**Computation:**

- Alabama Supercomputer Authority
- Auburn University Hopper Cluster

**Photo credits:**

- Rafe Brown
- Perry Wood, Jr.
- PhyloPic

# Questions?

joaks@auburn.edu

phyletica.org