

Generalizing phylogenetics to  
infer patterns predicted by  
processes of diversification

**Jamie Oaks**

Auburn University

[phyletica.org](http://phyletica.org)

 [@jamoaks](https://twitter.com/jamoaks)

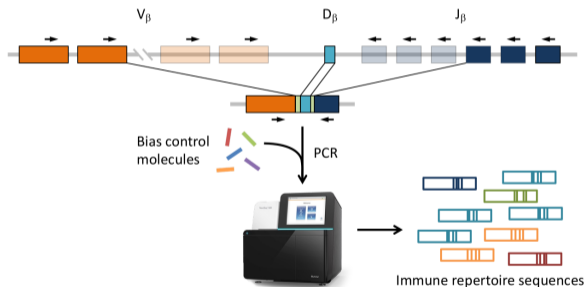
Scan for slides:



[phyletica.org/slides/duke-cbb.pdf](http://phyletica.org/slides/duke-cbb.pdf)



© 2007 Boris Kulikov [boris-kulikov.blogspot.com](http://boris-kulikov.blogspot.com)





AUBURN

UNIVERSITY





Generalizing phylogenetics to  
infer patterns of shared  
evolutionary events

## The Phyleticians

### Postdocs

- ▶ Perry Wood, Jr
- ▶ *Brian Folt*
- ▶ *Jesse Grismer*

### Graduate students

- ▶ Tashitso Anamza
- ▶ Matt Buehler
- ▶ Kerry Cobb
- ▶ Kyle David
- ▶ Saman Jahangiri
- ▶ Randy Klabacka
- ▶ Morgan Muell
- ▶ Tanner Myers
- ▶ Claire Tracy
- ▶ *Breanna Siple*
- ▶ *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*

- Phylogenetics is rapidly becoming the statistical foundation of biology

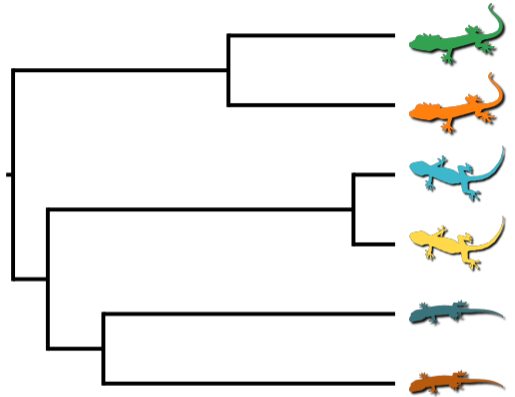


- ▶ Phylogenetics is rapidly becoming the statistical foundation of biology
- ▶ “Big data” present exciting possibilities and challenges

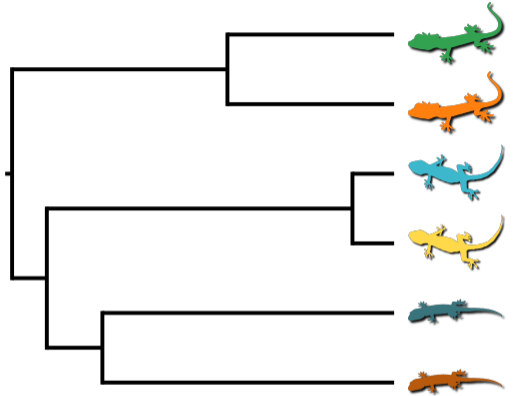


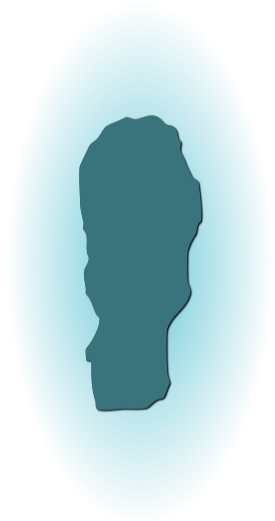
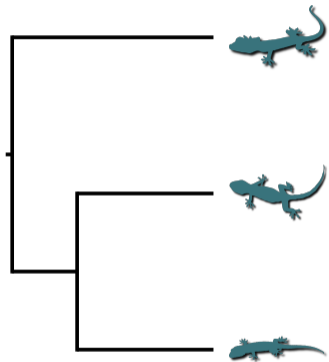
- ▶ Phylogenetics is rapidly becoming the statistical foundation of biology
- ▶ “Big data” present exciting possibilities and challenges
- ▶ Many opportunities to develop new ways to study biology in light of phylogeny

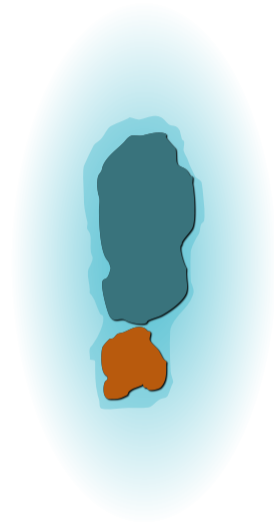
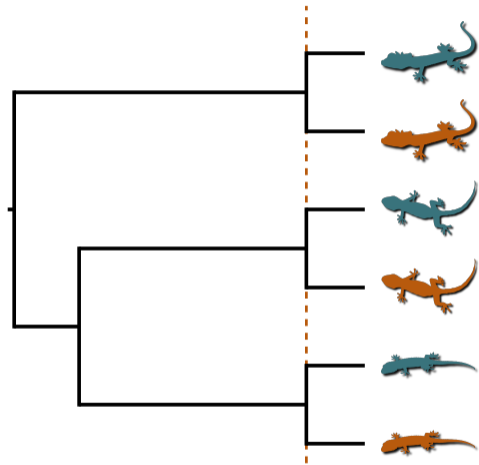




- **Assumption:** All processes of diversification affect each lineage independently

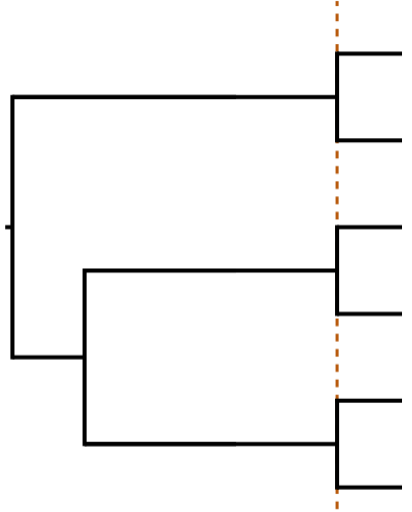






## Biogeography

- ▶ Environmental changes that affect whole communities of species

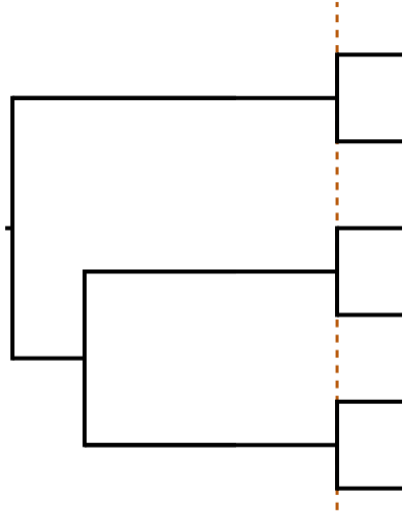


## Biogeography

- ▶ Environmental changes that affect whole communities of species

## Genome evolution

- ▶ Duplication of a chromosome segment harboring gene families



## Biogeography

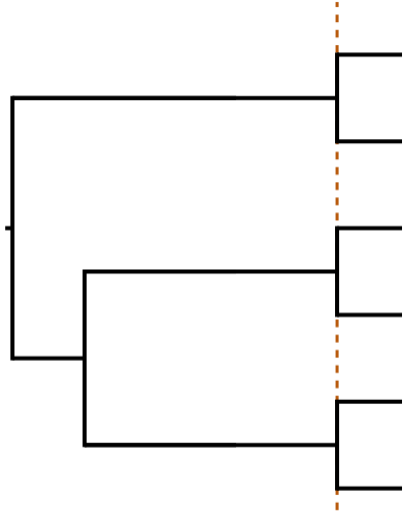
- ▶ Environmental changes that affect whole communities of species

## 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

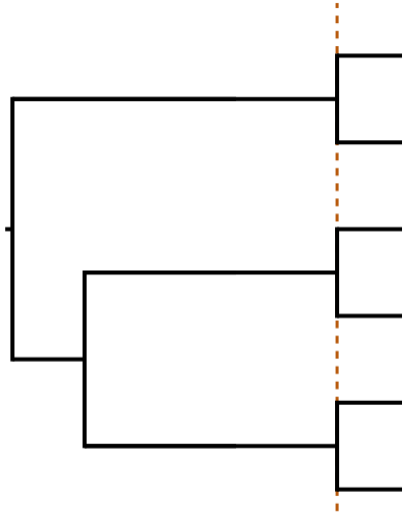
- ▶ Duplication of a chromosome segment harboring gene families

## Epidemiology

- ▶ Transmission at social gatherings

## Endosymbiont evolution (e.g., parasites, microbiome)

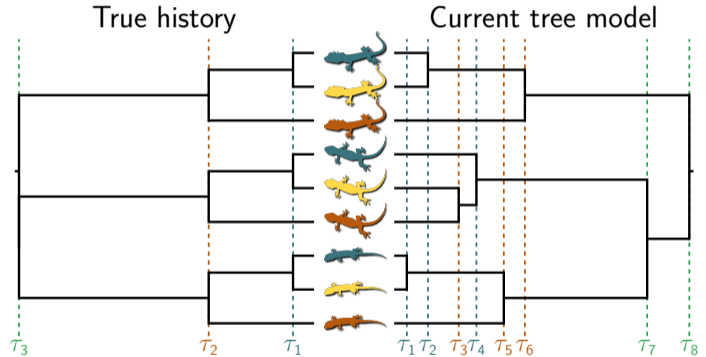
- ▶ Speciation of the host
- ▶ Co-colonization of new host species



# Why account for shared divergences?

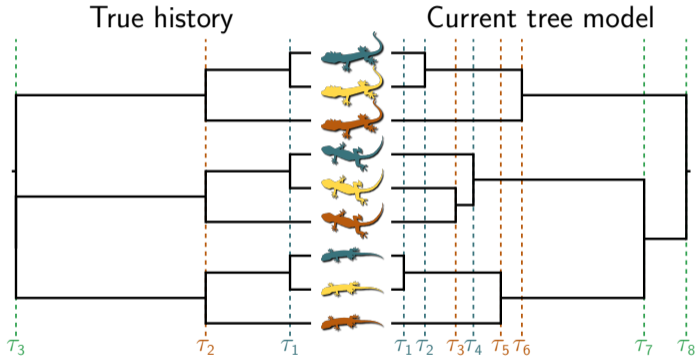
# Why account for shared divergences?

## 1. Improve inference



# 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

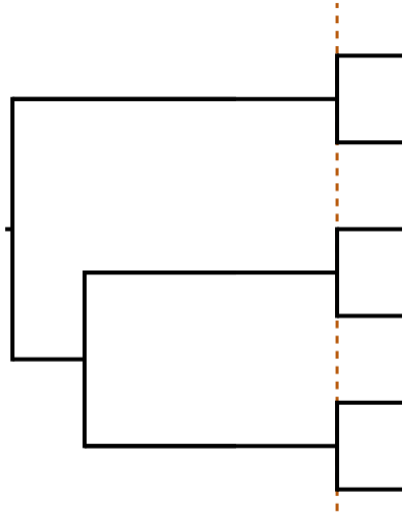
- ▶ Duplication of a chromosome segment harboring gene families

## Epidemiology

- ▶ Transmission at social gatherings

## Endosymbiont evolution (e.g., parasites, microbiome)

- ▶ Speciation of the host
- ▶ Co-colonization of new host species



## Approaches to the problem

A pairwise approach (keep it “simple”)

A fully phylogenetic approach

Tashitso Anamza



Tanner Myers



Randy Klabacka



Perry Wood, Jr.



Claire Tracy



Kerry Cobb



Matt Buehler



Nadia L'bahy



## Approaches to the problem

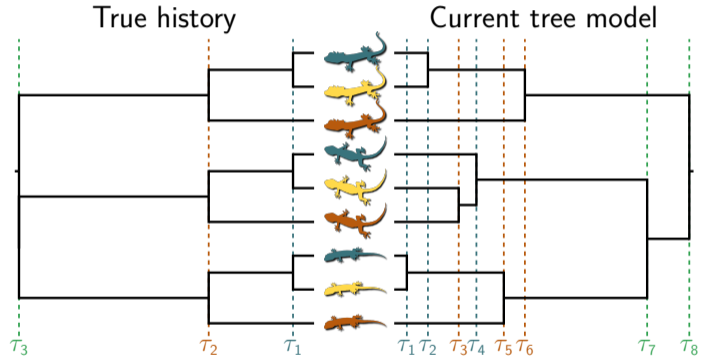
A pairwise approach (keep it “simple”)

A fully phylogenetic approach



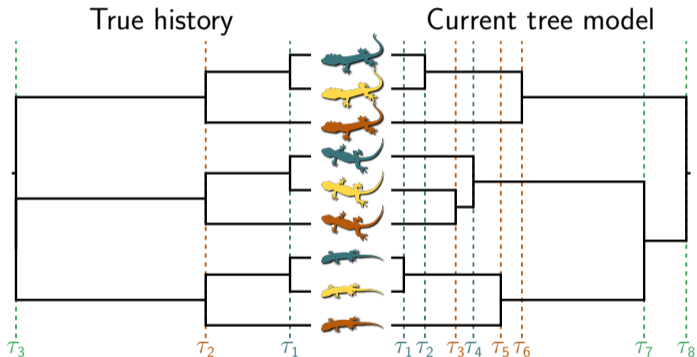
Dr. Perry Wood, Jr.

# Challenges to accounting for shared divergences



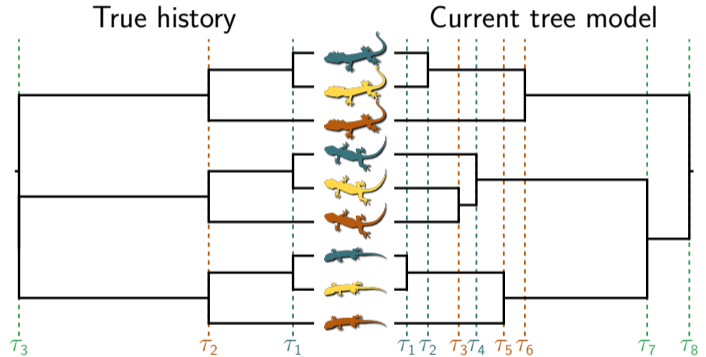
# Challenges to accounting for shared divergences

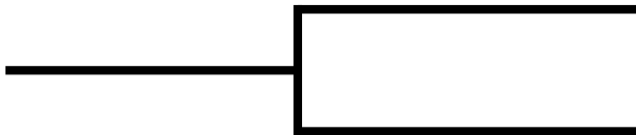
1. Likelihood for genomic data is tricky



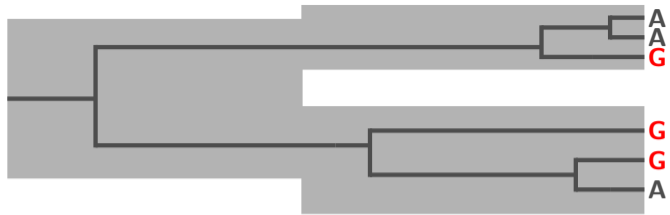
# Challenges to accounting for shared divergences

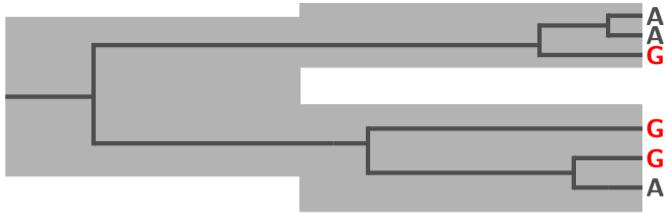
1. Likelihood for genomic data is tricky
2. Lots of possible trees of different dimensions



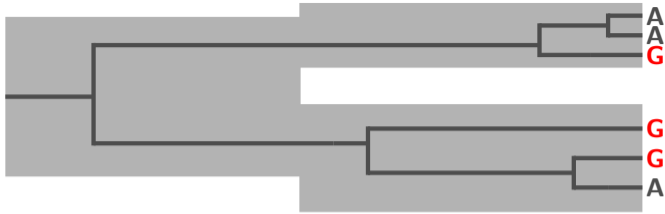




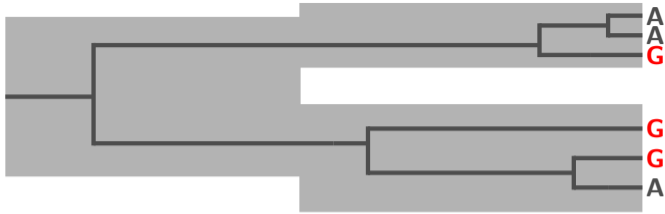




- Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent



- ▶ Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent
  - ▶ Coalescent is a stochastic model of shared inheritance (continuous-time Markov chain = CTMC)



- ▶ Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent
  - ▶ Coalescent is a stochastic model of shared inheritance (continuous-time Markov chain = CTMC)
  - ▶ Gene tree branching patterns are a function of population size



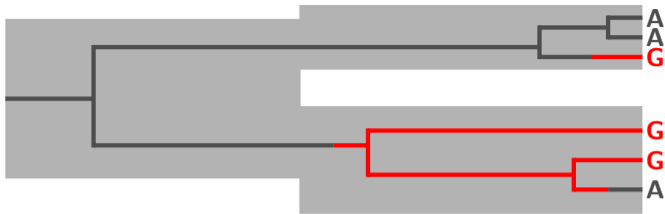
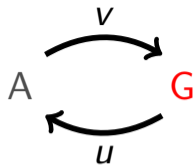
- ▶ Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent
  - ▶ Coalescent is a stochastic model of shared inheritance (continuous-time Markov chain = CTMC)
  - ▶ Gene tree branching patterns are a function of population size
- ▶ Conditional on  $G$ , model mutation as a CTMC



- ▶ Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent
  - ▶ Coalescent is a stochastic model of shared inheritance (continuous-time Markov chain = CTMC)
  - ▶ Gene tree branching patterns are a function of population size
- ▶ Conditional on  $G$ , model mutation as a CTMC
- ▶ Genetic characters provide information about  $G$



- ▶ Conditional on “population tree” ( $T$ ), model “gene trees” ( $G$ ) using coalescent
  - ▶ Coalescent is a stochastic model of shared inheritance (continuous-time Markov chain = CTMC)
  - ▶ Gene tree branching patterns are a function of population size
- ▶ Conditional on  $G$ , model mutation as a CTMC
- ▶ Genetic characters provide information about  $G$
- ▶  $G$  informs  $T$  (population sizes, divergence times, and relationships)





► “Standard” hierarchical approach



- ▶ “Standard” hierarchical approach
  - ▶ Calculate  $p(\text{genetic data} \mid G) \times p(G \mid T)$



- ▶ “Standard” hierarchical approach
  - ▶ Calculate  $p(\text{genetic data} \mid G) \times p(G \mid T)$
  - ▶ Use numerical integration (MCMC) to co-estimate both



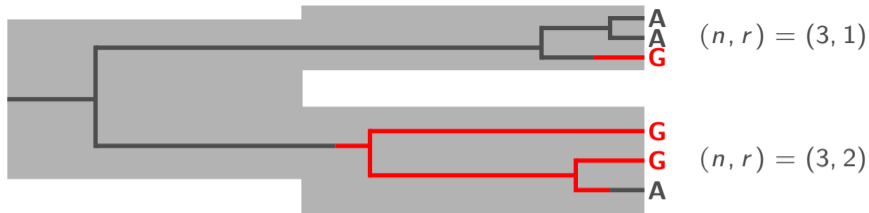
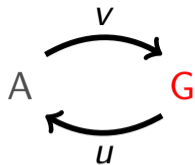
- ▶ “Standard” hierarchical approach
  - ▶ Calculate  $p(\text{genetic data} \mid G) \times p(G \mid T)$
  - ▶ Use numerical integration (MCMC) to co-estimate both
- ▶ But,  $G$  and  $T$  are highly correlated



- ▶ “Standard” hierarchical approach
  - ▶ Calculate  $p(\text{genetic data} \mid G) \times p(G \mid T)$
  - ▶ Use numerical integration (MCMC) to co-estimate both
- ▶ But,  $G$  and  $T$  are highly correlated
- ▶ As the number of loci (gene trees) increases, MCMC falls apart

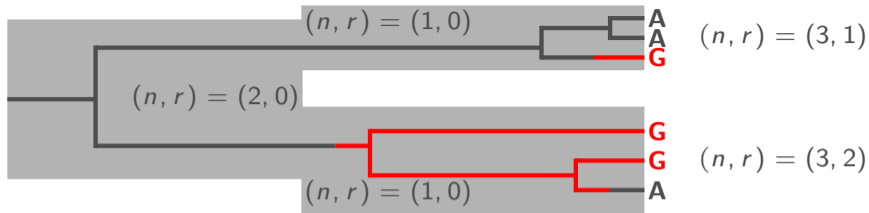
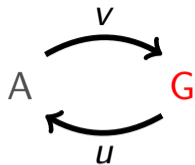


- ▶ “Standard” hierarchical approach
  - ▶ Calculate  $p(\text{genetic data} \mid G) \times p(G \mid T)$
  - ▶ Use numerical integration (MCMC) to co-estimate both
- ▶ But,  $G$  and  $T$  are highly correlated
- ▶ As the number of loci (gene trees) increases, MCMC falls apart
- ▶ Can we integrate  $G$  analytically?



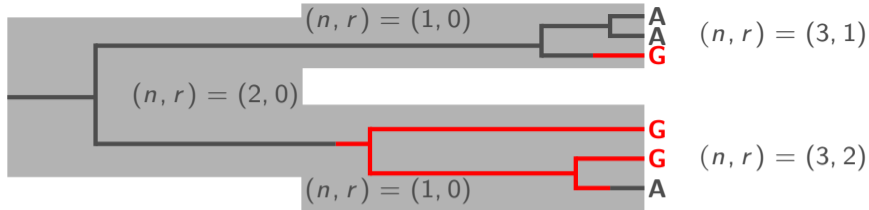
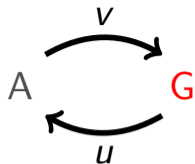
<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932



<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

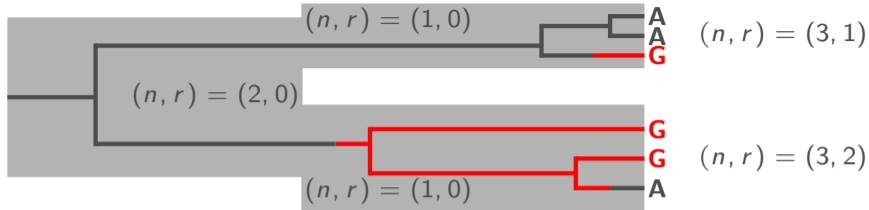
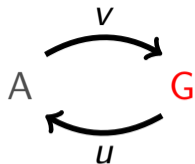
<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932



$$Q = \begin{array}{c|cccccc} & (1, 0) & (1, 1) & (2, 0) & (2, 1) & \cdots & (n, n) \\ \hline (1, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (1, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ \vdots & & & & & & \\ (n, n) & \cdot & \cdot & \cdot & \cdot & & \cdot \end{array}$$

<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

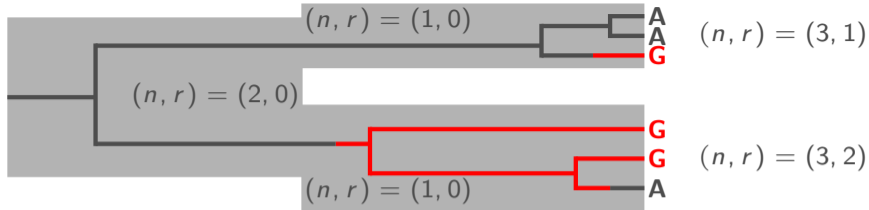
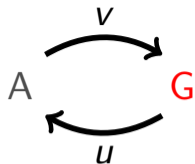


$$Q = \begin{array}{c|cccccc} & (1, 0) & (1, 1) & (2, 0) & (2, 1) & \cdots & (n, n) \\ \hline (1, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (1, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ \vdots & & & & & & \\ (n, n) & \cdot & \cdot & \cdot & \cdot & & \cdot \end{array}$$

$$\begin{aligned} Q_{(n,r);(n,r-1)} &= (n-r+1)v, && \text{mutation,} \\ Q_{(n,r);(n,r+1)} &= (r+1)u, && \text{mutation,} \\ Q_{(n,r);(n-1,r)} &= \frac{(n-1-r)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n-1,r-1)} &= \frac{(r-1)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n,r)} &= -\frac{(n-1)n}{2N_e(u+v)} - (n-r)v - ru. \end{aligned}$$

<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932



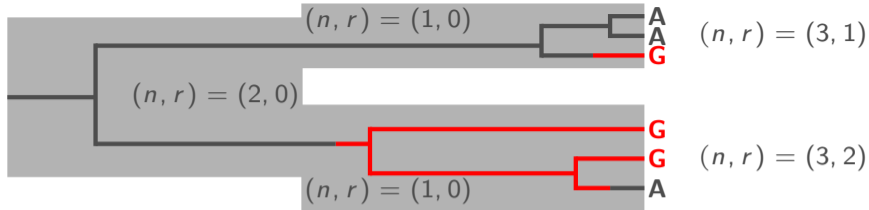
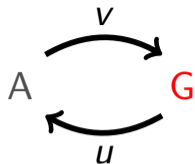
$$Q = \begin{array}{c|cccccc} & (1, 0) & (1, 1) & (2, 0) & (2, 1) & \cdots & (n, n) \\ \hline (1, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (1, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ \vdots & & & & & & \\ (n, n) & \cdot & \cdot & \cdot & \cdot & & \cdot \end{array}$$

$$\begin{aligned} Q_{(n,r);(n,r-1)} &= (n-r+1)v, && \text{mutation,} \\ Q_{(n,r);(n,r+1)} &= (r+1)u, && \text{mutation,} \\ Q_{(n,r);(n-1,r)} &= \frac{(n-1-r)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n-1,r-1)} &= \frac{(r-1)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n,r)} &= -\frac{(n-1)n}{2N_e(u+v)} - (n-r)v - ru. \end{aligned}$$

►  $e^{Qt}$  to keep track of all conditional probabilities along each branch (Carathéodory-Fejér method<sup>1</sup>)

<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932



$$Q = \begin{array}{c|cccccc} & (1, 0) & (1, 1) & (2, 0) & (2, 1) & \cdots & (n, n) \\ \hline (1, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (1, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 0) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ (2, 1) & \cdot & \cdot & \cdot & \cdot & & \cdot \\ \vdots & & & & & & \\ (n, n) & \cdot & \cdot & \cdot & \cdot & & \cdot \end{array}$$

$$\begin{aligned} Q_{(n,r);(n,r-1)} &= (n-r+1)v, && \text{mutation,} \\ Q_{(n,r);(n,r+1)} &= (r+1)u, && \text{mutation,} \\ Q_{(n,r);(n-1,r)} &= \frac{(n-1-r)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n-1,r-1)} &= \frac{(r-1)n}{2N_e(u+v)}, && \text{coalescence,} \\ Q_{(n,r);(n,r)} &= -\frac{(n-1)n}{2N_e(u+v)} - (n-r)v - ru. \end{aligned}$$

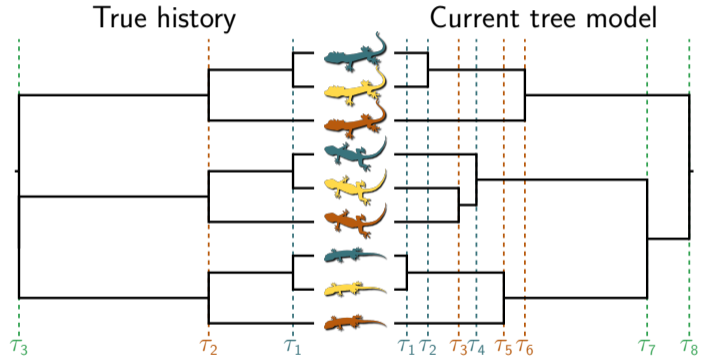
- $e^{Qt}$  to keep track of all conditional probabilities along each branch (Carathéodory-Fejér method<sup>1</sup>)
- At root, get likelihood of population tree integrated over all possible gene trees and mutational histories<sup>2</sup>

<sup>1</sup> T. Schmelzer and L. N. Trefethen (2007). *Electronic Transactions on Numerical Analysis* 29: 1–18

<sup>2</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

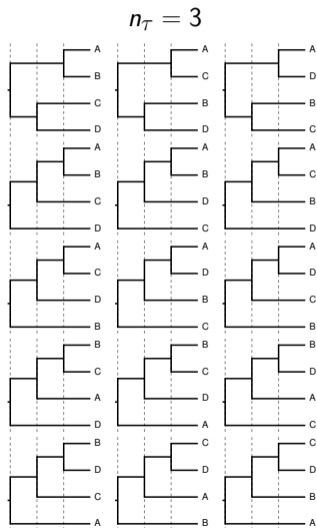
# Challenges to accounting for shared divergences

1. Likelihood for genomic data is tricky
2. Lots of possible trees of different dimensions

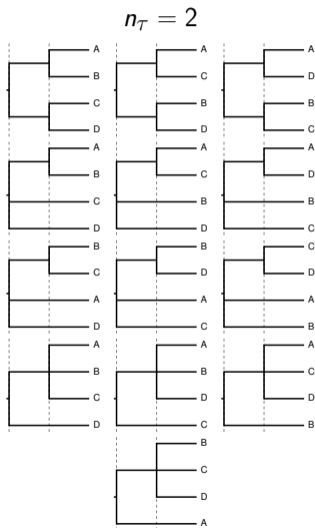
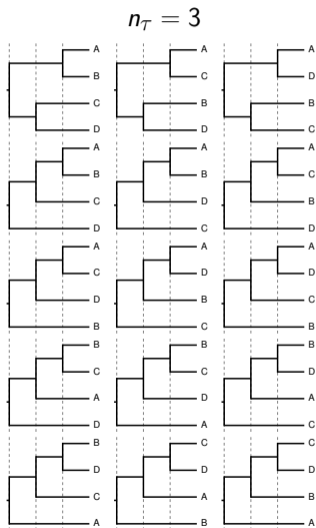


# Generalizing tree space

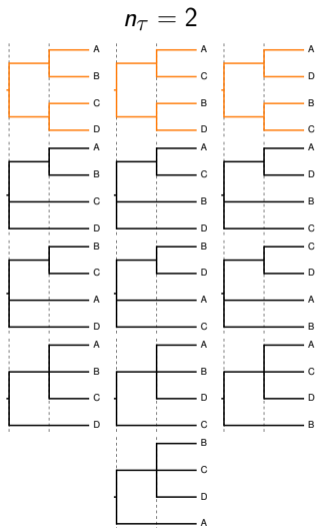
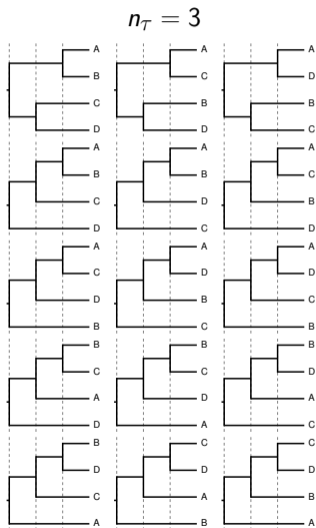
# Generalizing tree space



# Generalizing tree space



# Generalizing tree space

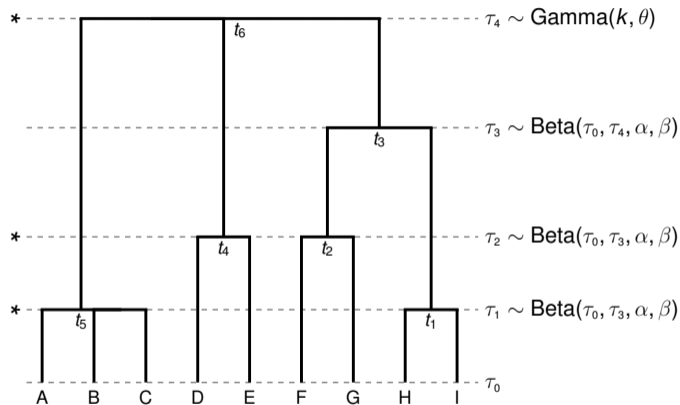


$n_T = 1$

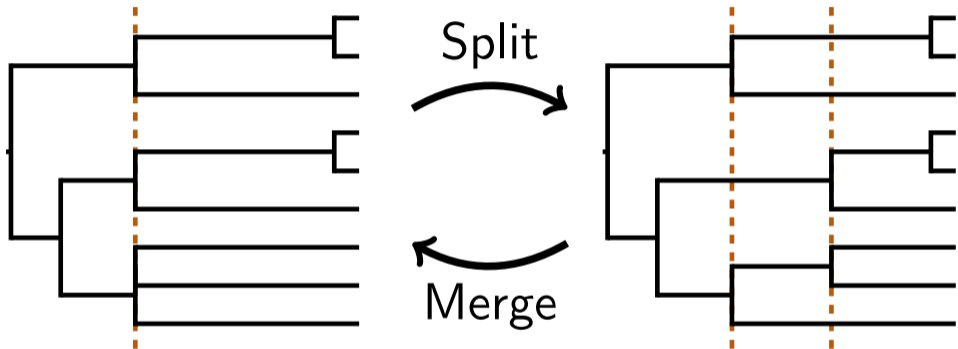


# Generalized tree distribution

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

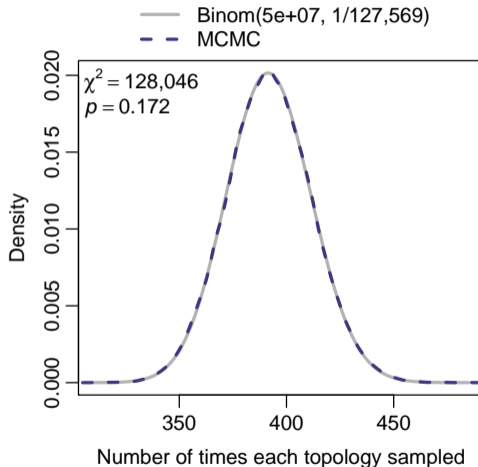


# Inferring trees with shared divergences



Reversible-jump MCMC

# Validating rjMCMC with 7-leaf tree



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

# PhycoEval

Phylogenetic coevality

J. R. Oaks et al. (2022). *PNAS* 119: e2121036119

# Ecoevolity

Estimating evolutionary coevality

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

- ▶ **Tree model**

- ▶ rjMCMC sampling of generalized tree distribution

# PhycoEval

Phylogenetic coevality

J. R. Oaks et al. (2022). *PNAS* 119: e2121036119

# Ecoevolity

Estimating evolutionary 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<sup>1</sup>

<sup>1</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

- ▶ **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<sup>1</sup>

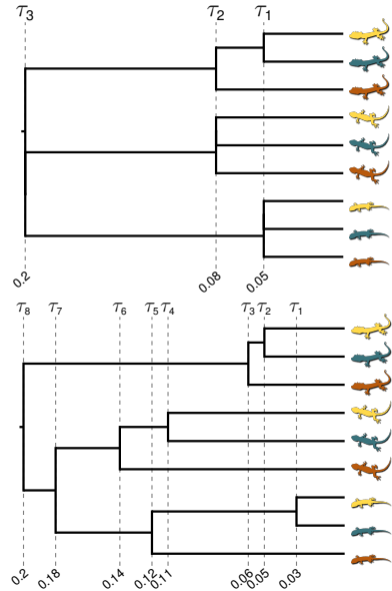
- ▶ *Goal: Co-estimation of phylogeny and shared divergences from genomic data*

<sup>1</sup> D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

Does it work?

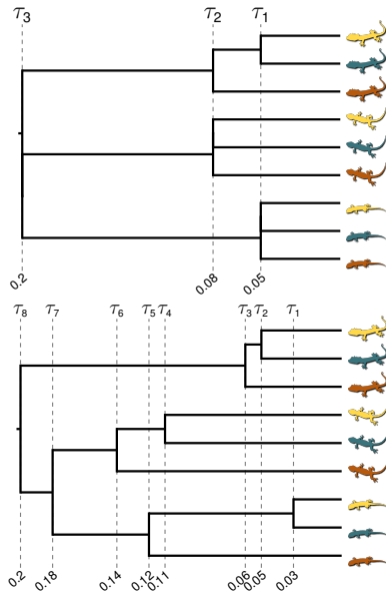
# Methods: Simulations

- ▶ Simulated 100 data sets with 50,000 base pairs



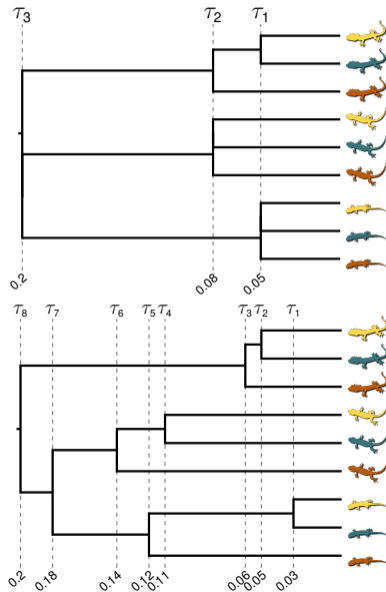
# Methods: Simulations

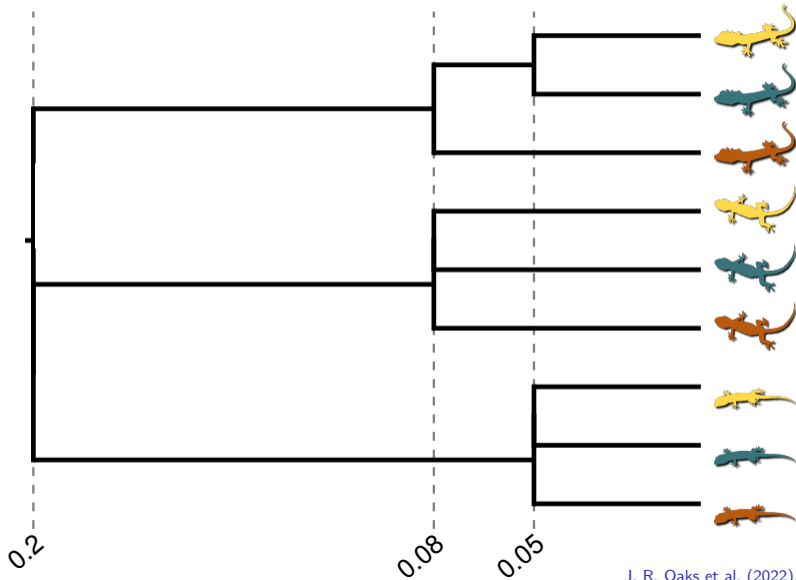
- ▶ Simulated 100 data sets with 50,000 base pairs
- ▶ Analyzed each data set with:
  - ▶  $M_G$  = Generalized tree model
  - ▶  $M_{IB}$  = Independent-bifurcating tree model

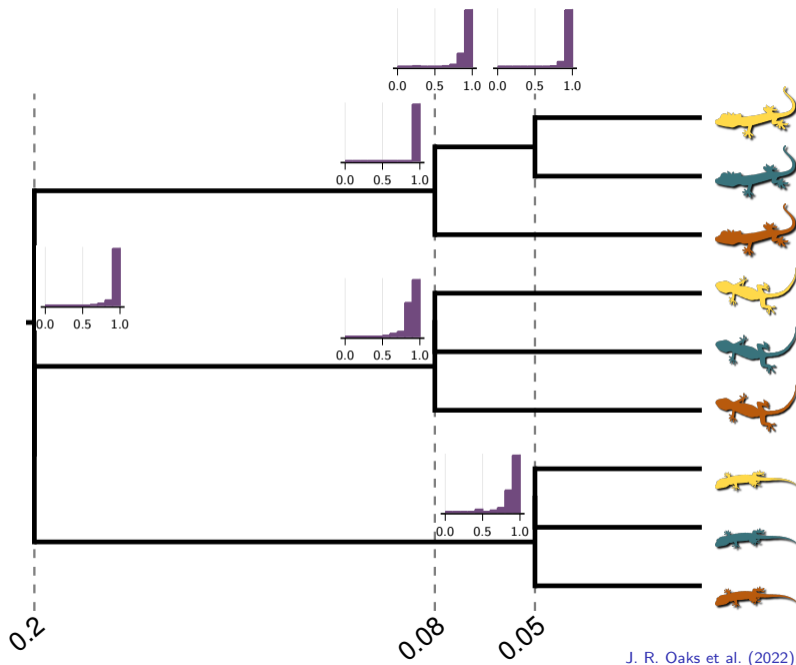


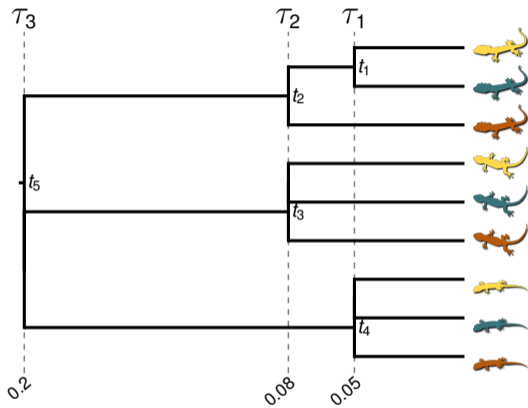
# Methods: Simulations

- ▶ Simulated 100 data sets with 50,000 base pairs
- ▶ 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  $M_{IB}$



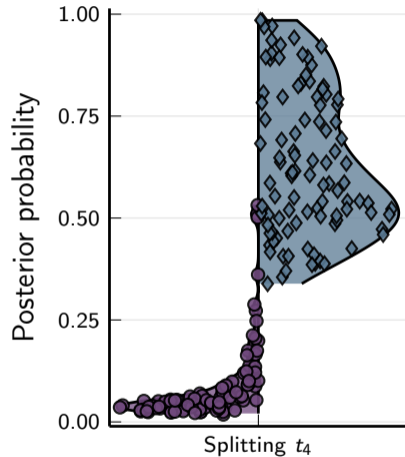
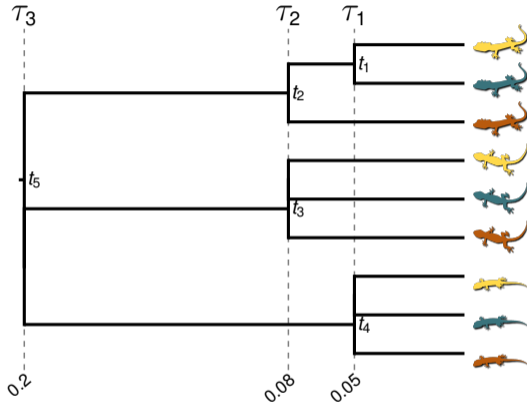




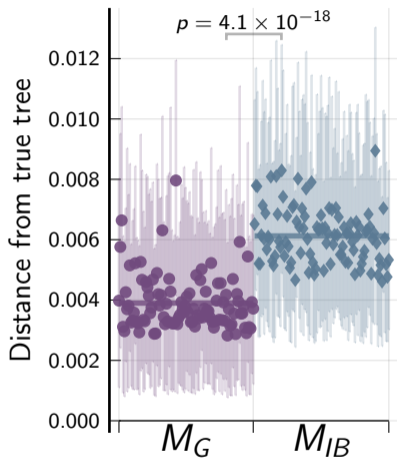
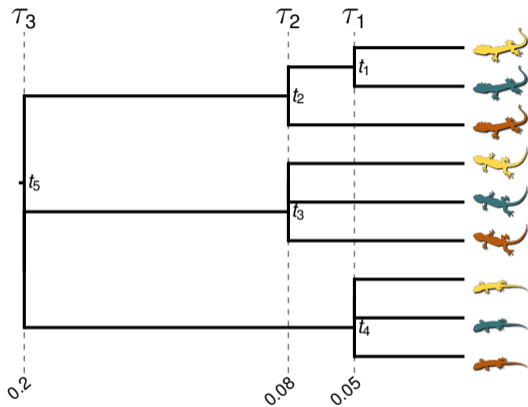


●  $M_G$  = Generalized model

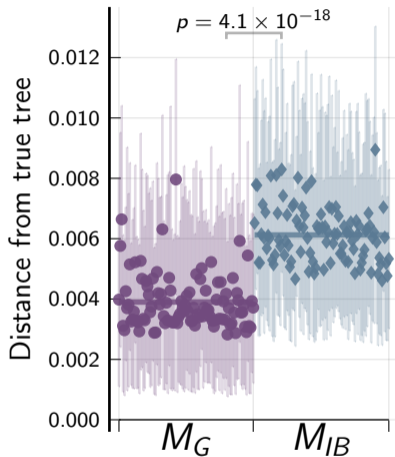
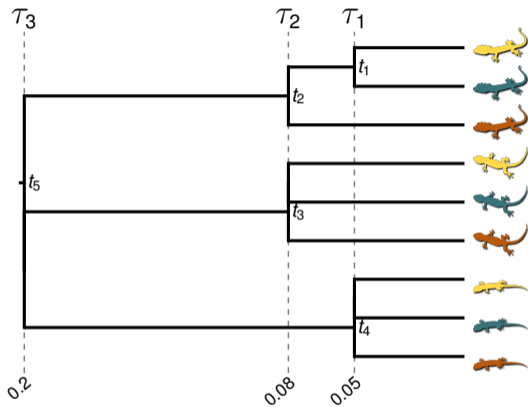
◆  $M_{IB}$  = Independent-bifurcating model



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

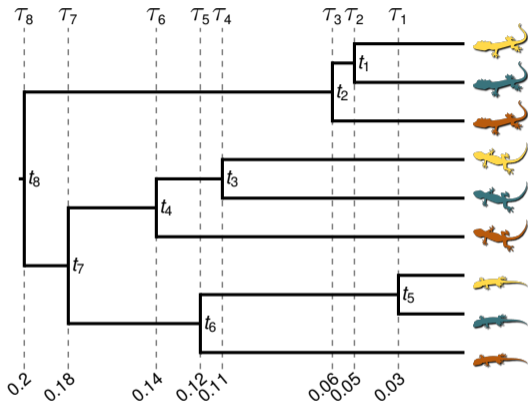


●  $M_G$  = Generalized model      ◆  $M_{IB}$  = Independent-bifurcating model

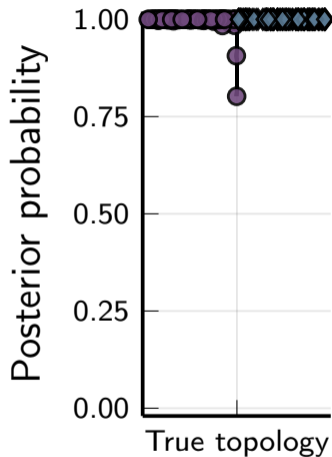
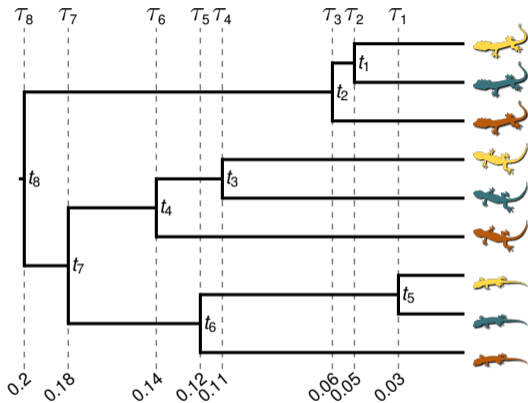


$M_G$  significantly better at inferring trees with shared divergences

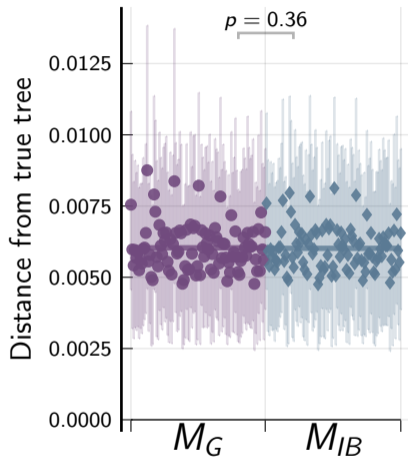
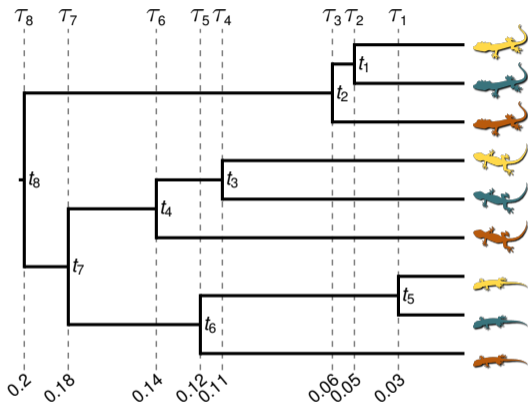
  $M_G$  = Generalized model
   $M_{IB}$  = Independent-bifurcating model



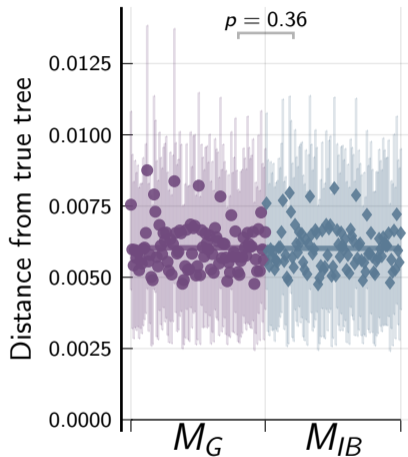
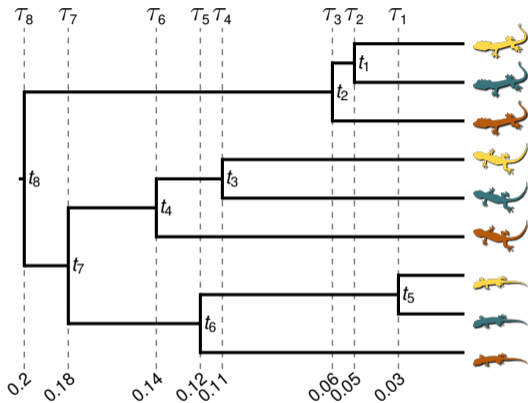
●  $M_G$  = Generalized model
 ◆  $M_{IB}$  = Independent-bifurcating model



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

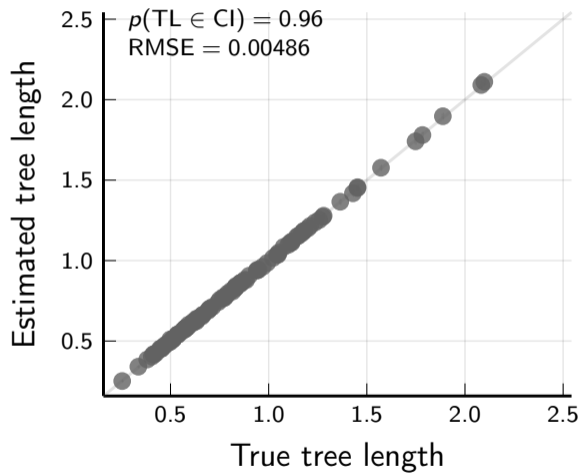


●  $M_G$  = Generalized model      ◆  $M_{IB}$  = Independent-bifurcating model

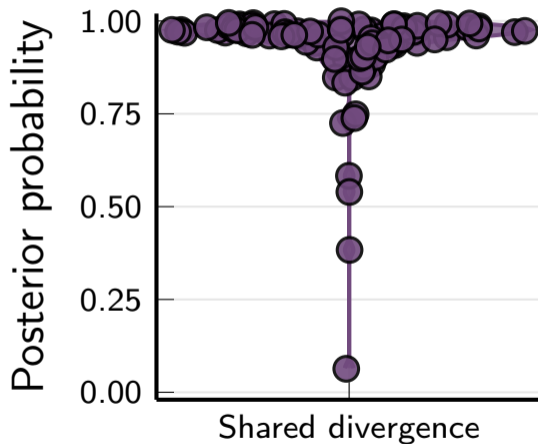
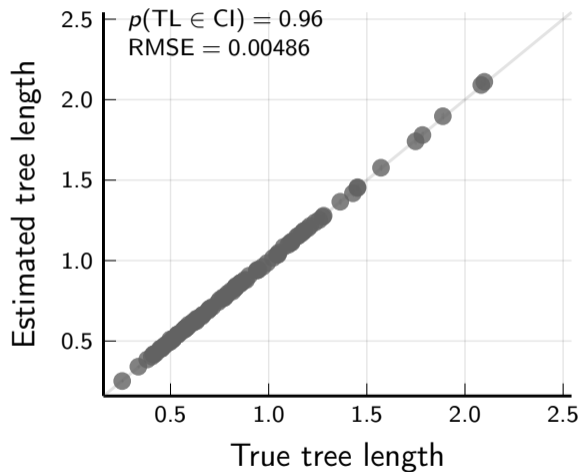


**$M_G$  performs as well as true model when divergences are independent**

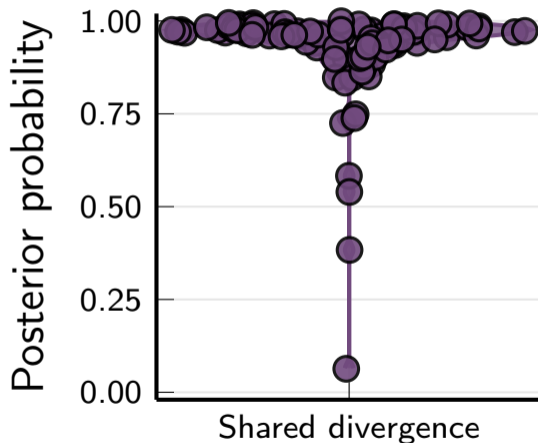
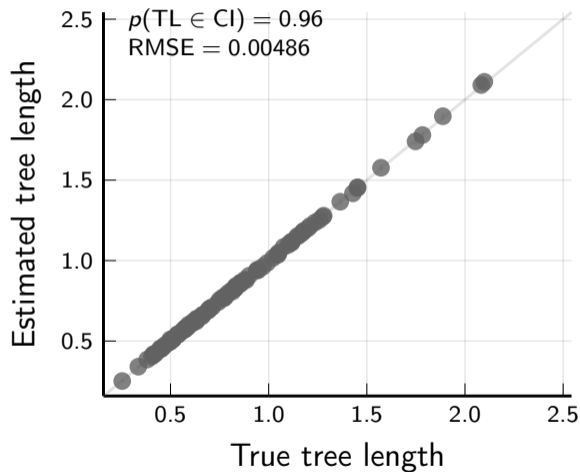
## Results: random $M_G$ trees



## Results: random $M_G$ trees



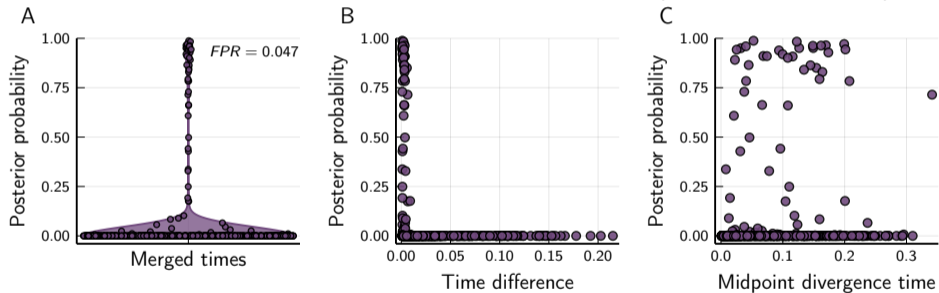
## Results: random $M_G$ trees



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

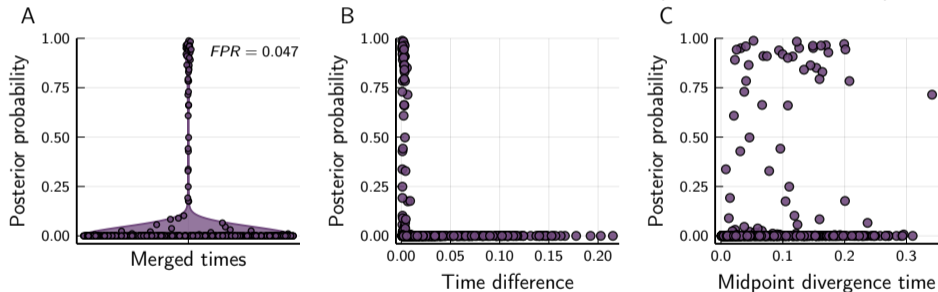
# Results: random $M_{IB}$ trees

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



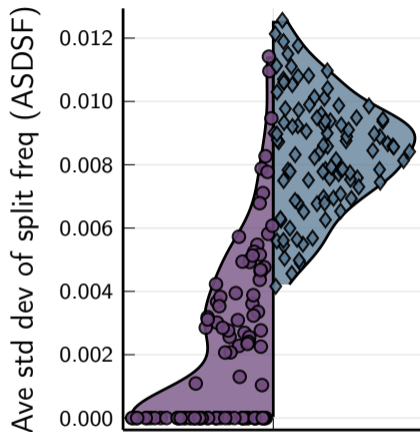
# Results: random $M_{IB}$ trees

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



$M_G$  has low false positive rate

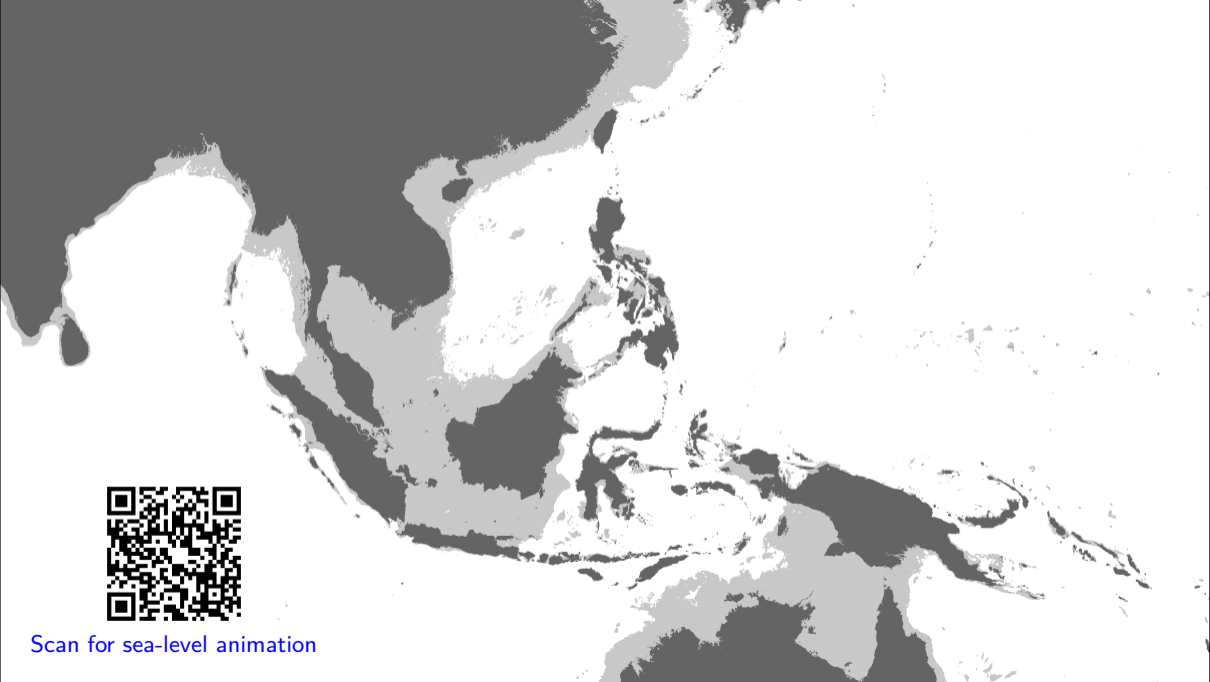
●  $M_G$  = Generalized model      ◆  $M_{IB}$  = Independent-bifurcating model



**Generalizing tree space improves MCMC convergence and mixing**



Scan for sea-level animation



Scan for sea-level animation



**Did fragmentation of islands  
promote diversification?**

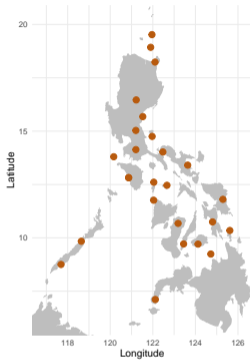


Scan for sea-level animation

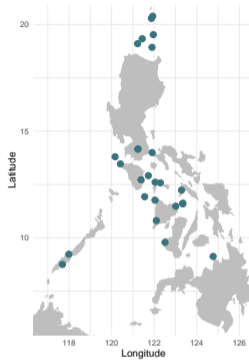
## *Cyrtodactylus*



©Rafe M. Brown



## *Gekko*

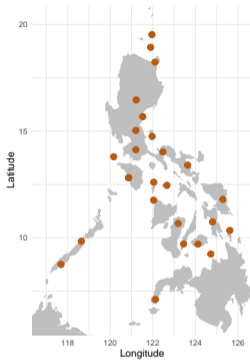


©Rafe M. Brown

## *Cyrtodactylus*

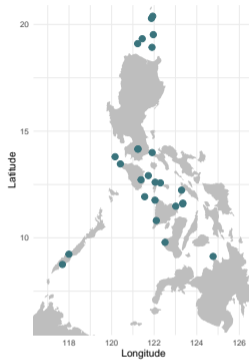


©Rafe M. Brown



1702 loci  
155,887 sites

## *Gekko*

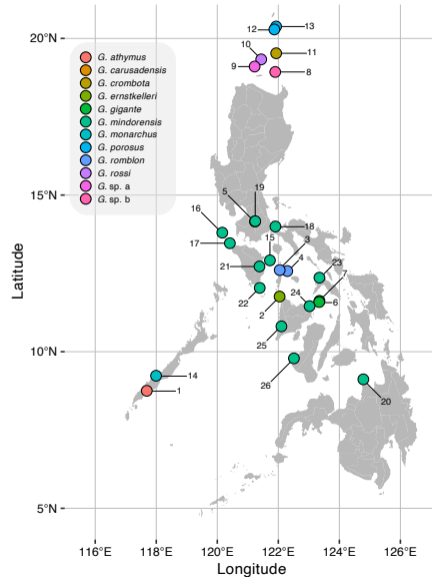
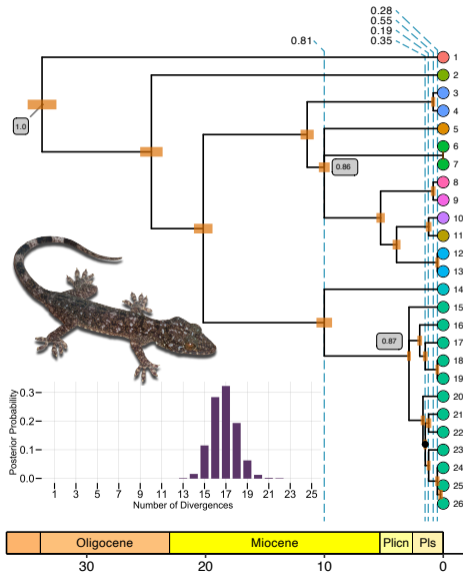


1033 loci  
94,813 sites

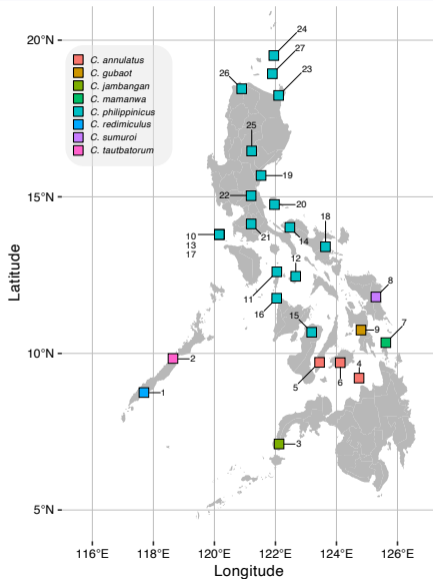
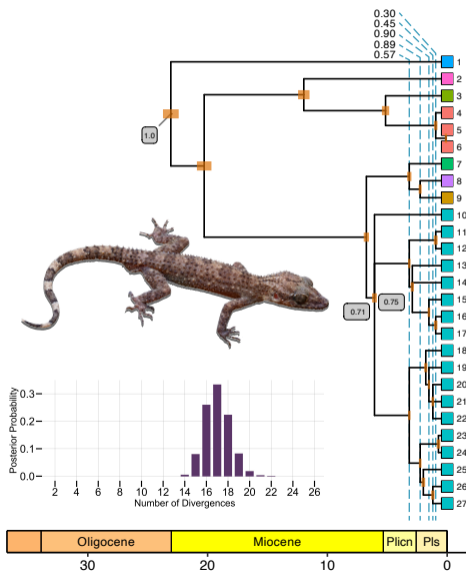


©Rafe M. Brown

# Gekko



# Cyrtodactylus



## Take-home points

- ▶ We can accurately infer phylogenies with shared divergences

## Take-home points

- ▶ We can accurately infer phylogenies with shared divergences
- ▶ Generalizing tree space avoids spurious support incorrect relationships and improves MCMC mixing

## Take-home points

- ▶ We can accurately infer phylogenies with shared divergences
- ▶ Generalizing tree space avoids spurious support incorrect relationships 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: [github.com/phyletica/ecoevolity](https://github.com/phyletica/ecoevolity)  
(release coming soon)

### **Open-Science Notebooks:**

- ▶ Phycoeval analyses: [github.com/phyletica/phycoeval-experiments](https://github.com/phyletica/phycoeval-experiments)
- ▶ Gecko RADseq: [github.com/phyletica/gekgo](https://github.com/phyletica/gekgo)

Moving forward: Theory/methods

## Moving forward: Theory/methods

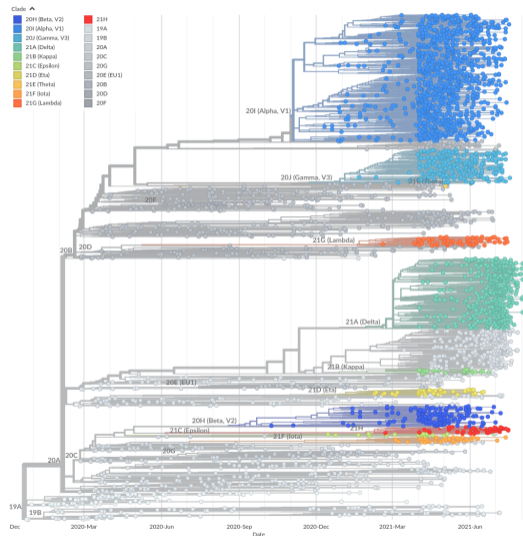
- ▶ Develop process-based and trait-dependent distributions over the space of generalized trees
  - ▶ “Birth-death-burst” model

## Moving forward: Theory/methods

- ▶ Develop process-based and trait-dependent distributions over the space of generalized trees
  - ▶ “Birth-death-burst” model
- ▶ Extend generalized tree distribution to trees that are not ultrametric

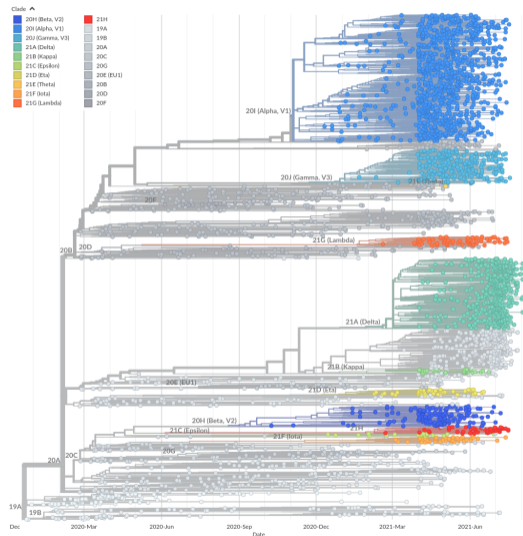
# Moving forward: Theory/methods

- ▶ Develop process-based and trait-dependent distributions over the space of generalized trees
  - ▶ “Birth-death-burst” model
- ▶ Extend generalized tree distribution to trees that are not ultrametric

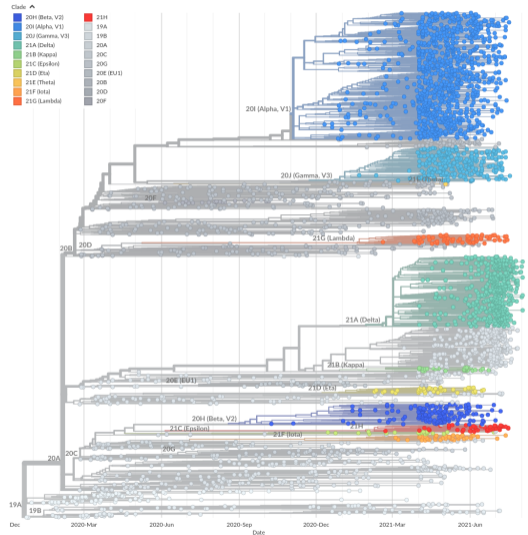


# Moving forward: Theory/methods

- ▶ Develop process-based and trait-dependent distributions over the space of generalized trees
  - ▶ “Birth-death-burst” model
- ▶ Extend generalized tree distribution to trees that are not ultrametric
- ▶ Couple generalized tree distribution with other phylogenetic likelihood models

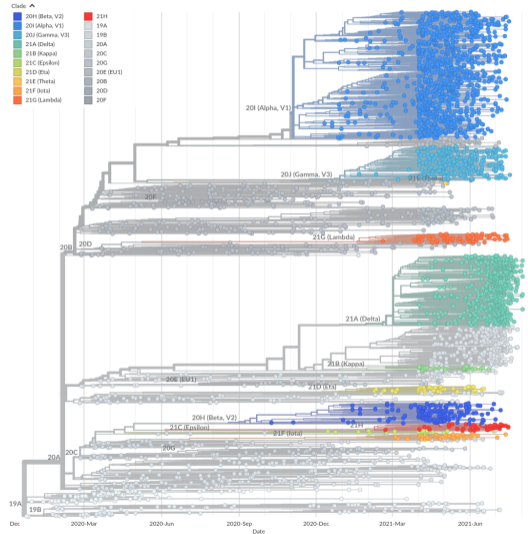


# Moving forward: Applications



# Moving forward: Applications

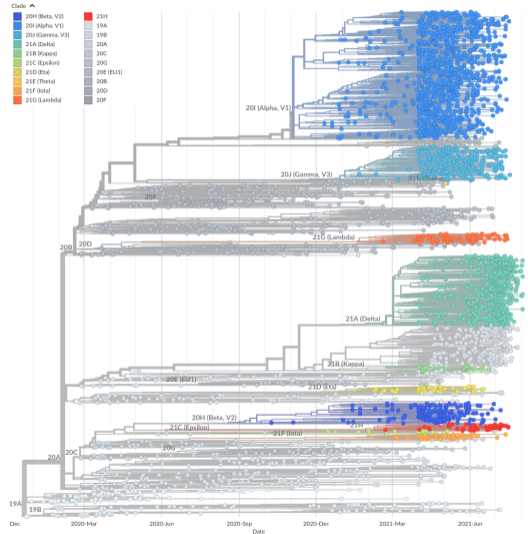
Epidemiological dynamics of “super-spreading” events during the COVID-19 pandemic



# Moving forward: Applications

Epidemiological dynamics of “super-spreading” events during the COVID-19 pandemic

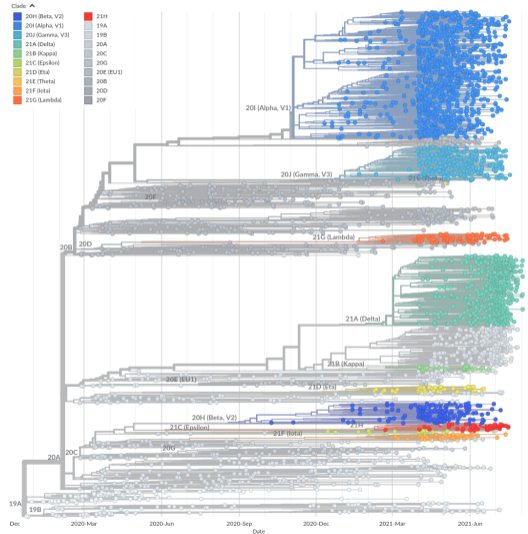
- Spread at social gatherings creates shared and multifurcating divergences in the viral “transmission tree”



# Moving forward: Applications

Epidemiological dynamics of “super-spreading” events during the COVID-19 pandemic

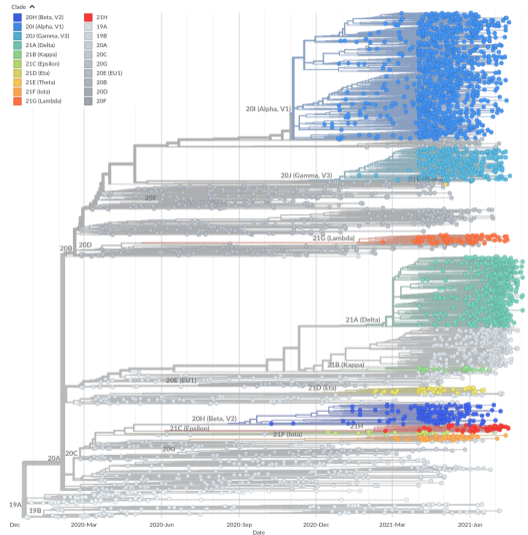
- ▶ Spread at social gatherings creates shared and multifurcating divergences in the viral “transmission tree”
- ▶ Estimate rate of shared divergences as proxy for spread via social gatherings

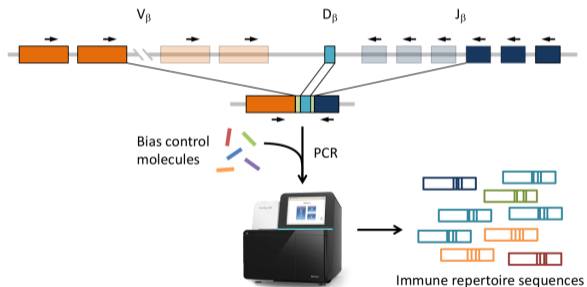


# Moving forward: Applications

Epidemiological dynamics of “super-spreading” events during the COVID-19 pandemic

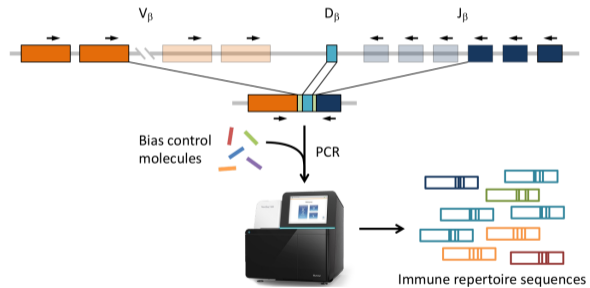
- ▶ Spread at social gatherings creates shared and multifurcating divergences in the viral “transmission tree”
- ▶ Estimate rate of shared divergences as proxy for spread via social gatherings
- ▶ Test if this varies over time, among regions, and among variants of SARS-CoV-2





# Adaptive

biotechnologies™



# Stats & Algorithms

- Our team is very cross-functional
  - We solve statistical and computational challenges for many stakeholders across Adaptive

**Erin Calfee**

Computational Biologist I



**Travers Ching**

Sr Computational Biologist



**Amy Ko**

Computational Biologist II



**Jamie Oaks**

Sr Manager, Computational Bi...



# Come work with us!

## Adaptive Comp Bio Internships!

Scan to internship listing:



[www.adaptivebiotech.com/career-listings](http://www.adaptivebiotech.com/career-listings)



### Internships

**Intern, BCR Discovery**

Seattle

**Intern, Cloud Engineer**

Remote (WFH)

**Intern, Computational Biology**

Remote (WFH)

**Intern, Computational Biology**

Remote (WFH)

**Intern, Computational Biology, Antigen Map**

Remote (WFH)

**Intern, Computational Biology (CRI)**

Remote (WFH)

**Intern, Computational Biology, Stats and Algorithms**

Remote (WFH)

# Thanks everyone!

- ▶ Thanks Devang and Duke GCB & CBB!
- ▶ Bryan Howie and my team at Adaptive
- ▶ Phyletica Lab (the Phyleticians)
- ▶ Mark Holder
- ▶ Rafe Brown
- ▶ Cam Siler
- ▶ Lee Grismer

## **Computation:**

- ▶ Alabama Supercomputer Authority
- ▶ Auburn University Hopper Cluster

## **Funding:**



## **Photo credits:**

- ▶ Rafe Brown
- ▶ Perry Wood, Jr.
- ▶ [PhyloPic](#)

# Questions?

[joaks@auburn.edu](mailto:joaks@auburn.edu)

 [@jamoaks](https://twitter.com/jamoaks)

[phyletica.org](http://phyletica.org)

Scan for slides:



[phyletica.org/slides/duke-cbb.pdf](http://phyletica.org/slides/duke-cbb.pdf)



© 2007 Boris Kulikov [boris-kulikov.blogspot.com](http://boris-kulikov.blogspot.com)