

SSB 2018 Comparative Phylogeography Workshop

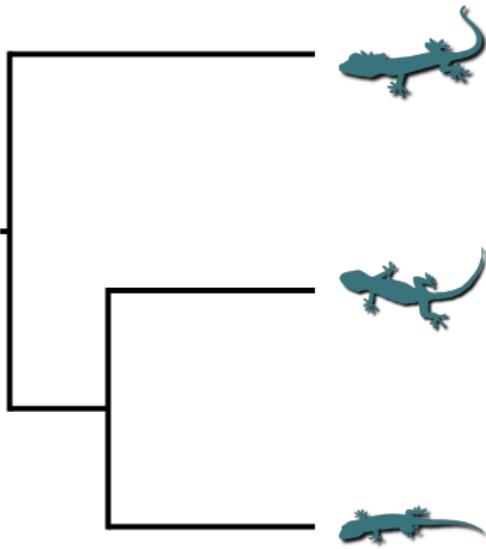
Jamie R. Oaks

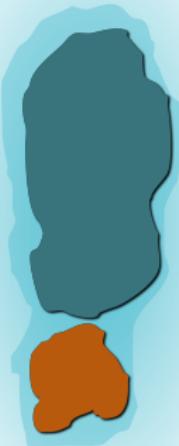
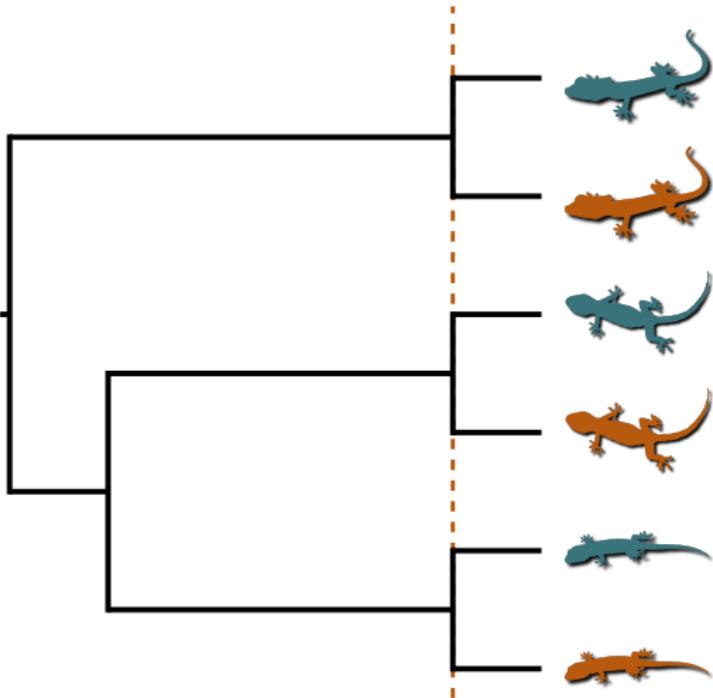
Department of Biological Sciences &
Museum of Natural History, Auburn
University

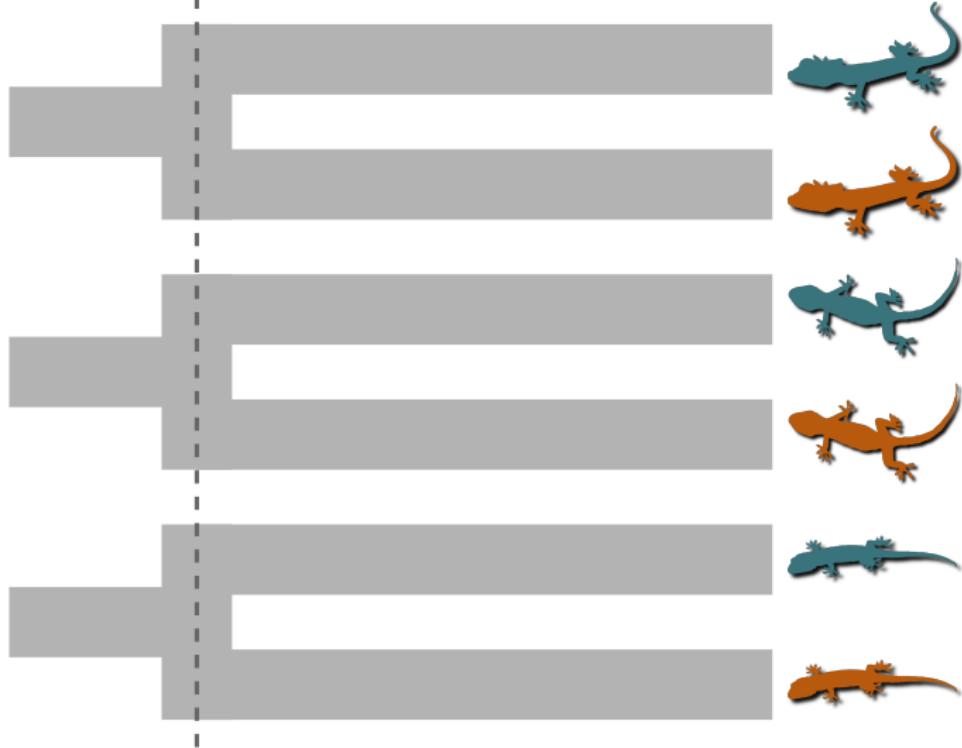
June 4, 2018

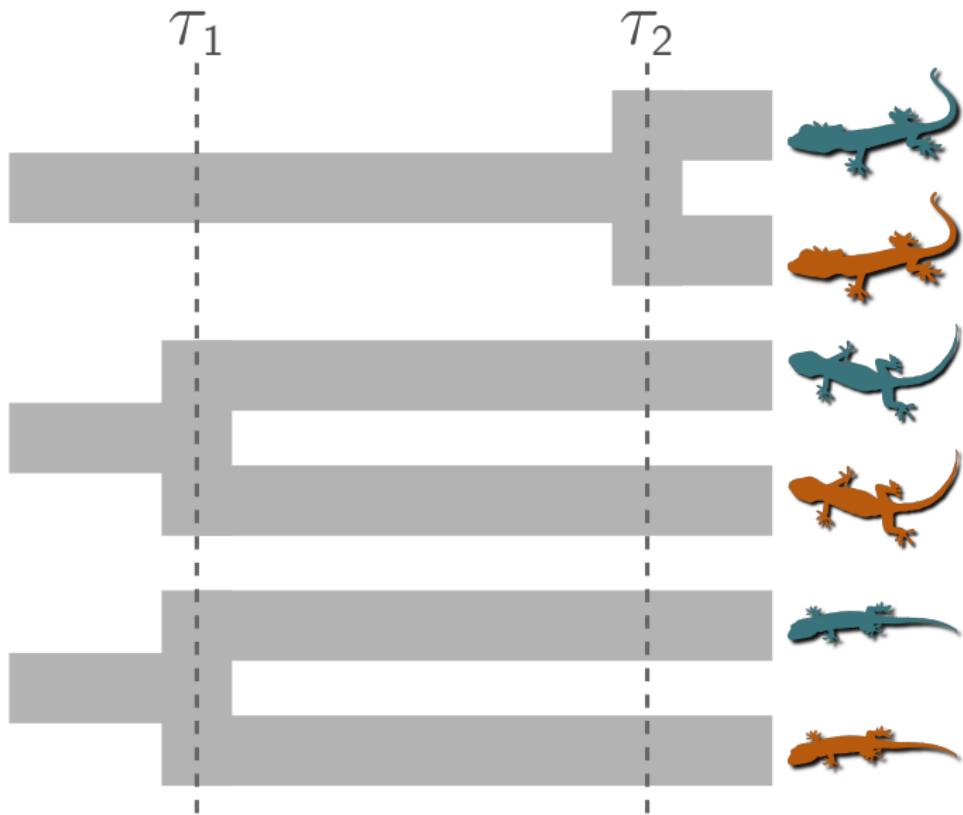


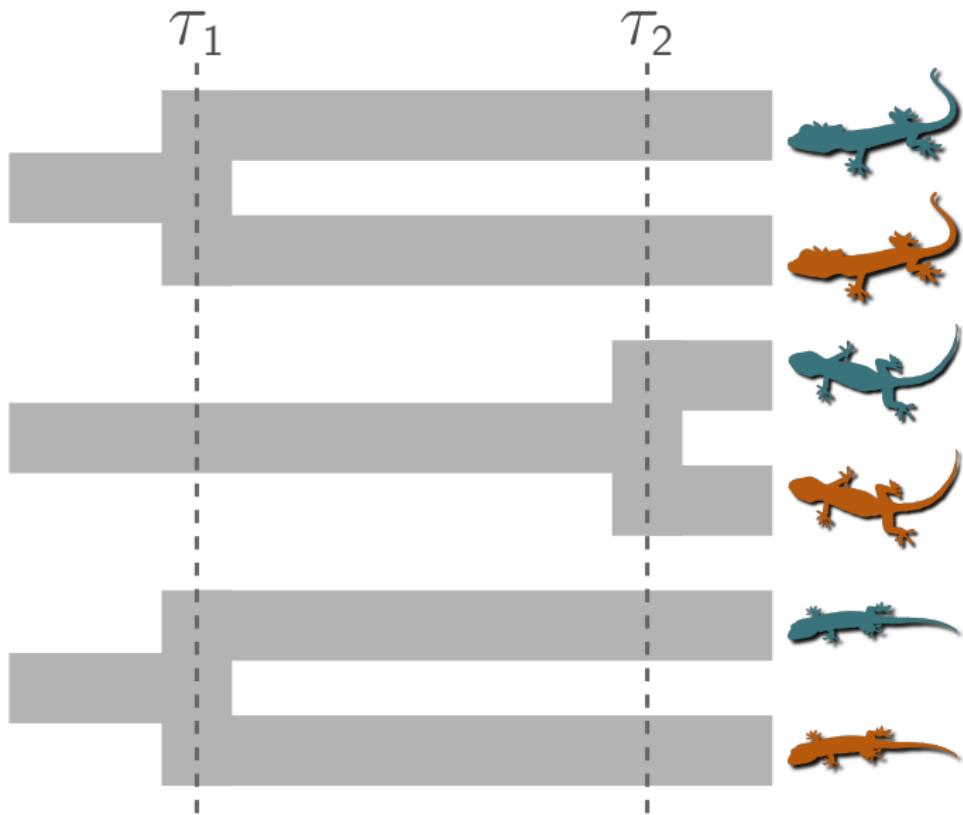
© 2007 Boris Kulikov boris-kulikov.blogspot.com

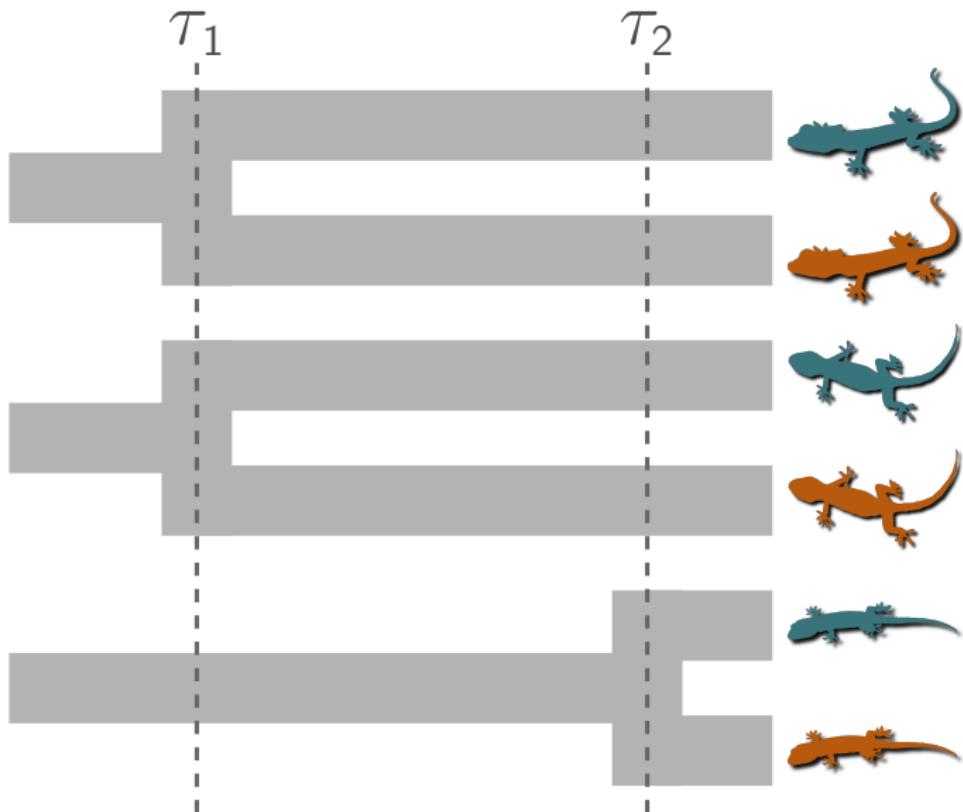


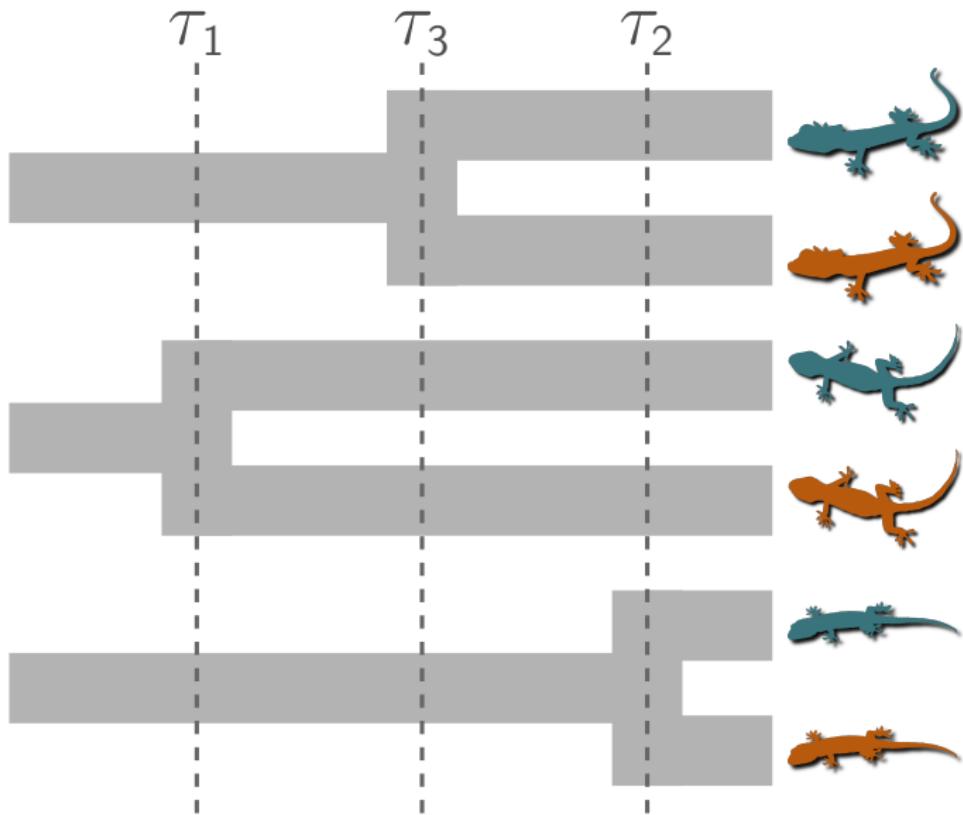


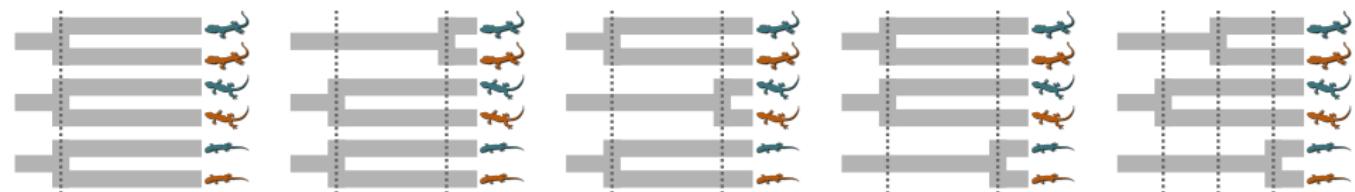
τ_1 

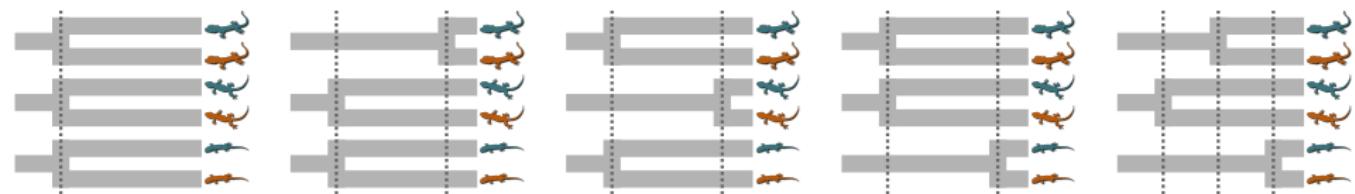








m_1 m_2 m_3 m_4 m_5 

m_1 m_2 m_3 m_4 m_5 

We want to infer the model and divergence times given genetic data

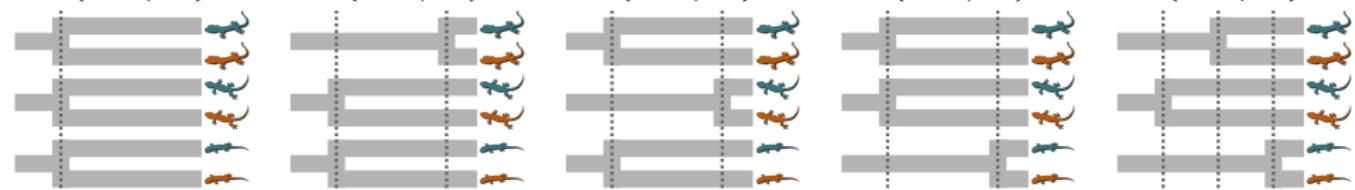
$p(m_1 | \mathbf{X})$

$p(m_2 | \mathbf{X})$

$p(m_3 | \mathbf{X})$

$p(m_4 | \mathbf{X})$

$p(m_5 | \mathbf{X})$



We want to infer the model and divergence times given genetic data

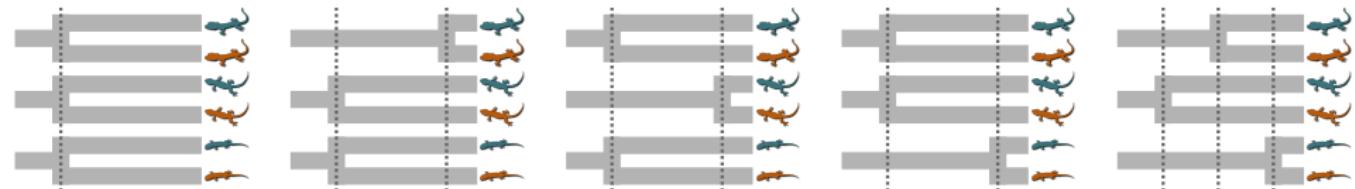
$p(m_1 | \mathbf{X})$

$p(m_2 | \mathbf{X})$

$p(m_3 | \mathbf{X})$

$p(m_4 | \mathbf{X})$

$p(m_5 | \mathbf{X})$



We want to infer the model and divergence times given genetic data

$$p(m_i | \mathbf{X}) \propto p(\mathbf{X} | m_i) p(m_i)$$

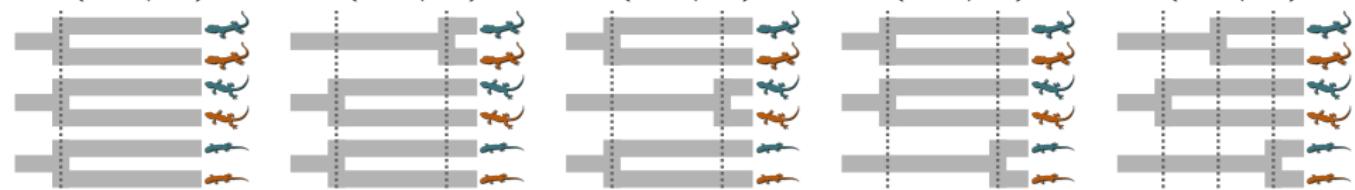
$p(m_1 | \mathbf{X})$

$p(m_2 | \mathbf{X})$

$p(m_3 | \mathbf{X})$

$p(m_4 | \mathbf{X})$

$p(m_5 | \mathbf{X})$



We want to infer the model and divergence times given genetic data

$$p(m_i | \mathbf{X}) \propto p(\mathbf{X} | m_i)p(m_i)$$

$$p(\mathbf{X} | m_i) = \int_{\theta} p(\mathbf{X} | \theta, m_i)p(\theta | m_i)d\theta$$

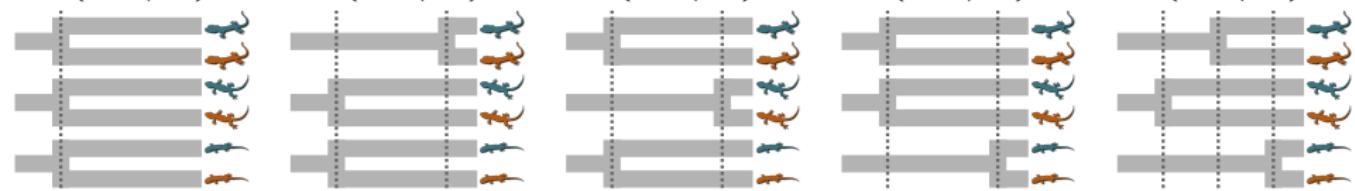
$p(m_1 | \mathbf{X})$

$p(m_2 | \mathbf{X})$

$p(m_3 | \mathbf{X})$

$p(m_4 | \mathbf{X})$

$p(m_5 | \mathbf{X})$



We want to infer the model and divergence times given genetic data

$$p(m_i | \mathbf{X}) \propto p(\mathbf{X} | m_i) p(m_i)$$

$$p(\mathbf{X} | m_i) = \int_{\theta} p(\mathbf{X} | \theta, m_i) p(\theta | m_i) d\theta$$

- ▶ Divergence times
- ▶ Substitution parameters
- ▶ Gene trees
- ▶ Demographic parameters

- ▶ Let's assume we are interested in the probability of a coin we have not seen landing heads-side up when it is flipped (θ)

- ▶ Let's assume we are interested in the probability of a coin we have not seen landing heads-side up when it is flipped (θ)
- ▶ Before flipping, we decide to compare two models that vary in our prior assumptions about the probability of the coin landing heads up

- ▶ Let's assume we are interested in the probability of a coin we have not seen landing heads-side up when it is flipped (θ)
- ▶ Before flipping, we decide to compare two models that vary in our prior assumptions about the probability of the coin landing heads up
- ▶ We assume:
 1. The coin is probably fair
 $M_1: \theta \sim \text{Beta}(5.0, 5.0)$
 2. the coin is weighted to land tails side up most of time
 $M_2: \theta \sim \text{Beta}(1.0, 5.0)$

- We do 100 flips and 50 land heads up

- ▶ We do 100 flips and 50 land heads up
- ▶ Now, we can calculate the posterior distribution for the probability of landing heads up under both our models

- ▶ We do 100 flips and 50 land heads up
- ▶ Now, we can calculate the posterior distribution for the probability of landing heads up under both our models

$$p(\theta | D, M_i) = \frac{p(D | \theta, M_i)p(\theta | M_i)}{p(D | M_i)}$$

- ▶ We do 100 flips and 50 land heads up
- ▶ Now, we can calculate the posterior distribution for the probability of landing heads up under both our models

$$p(\theta | D, M_i) = \frac{p(D | \theta, M_i)p(\theta | M_i)}{p(D | M_i)}$$

- ▶ We see the posterior distribution of θ is very robust to our prior assumptions
- ▶ <https://kerrycobb.github.io/beta-binomial-web-demo/>

- ▶ However, we want to compare the ability of the models to explain the data

- ▶ However, we want to compare the ability of the models to explain the data
- ▶ We need to average (integrate) the likelihood density function over all possible values of θ , weighting by the prior

- ▶ However, we want to compare the ability of the models to explain the data
- ▶ We need to average (integrate) the likelihood density function over all possible values of θ , weighting by the prior

$$p(D | M_1) = \int_{\theta} p(D | \theta, M_1) p(\theta | M_1) d\theta$$

- ▶ However, we want to compare the ability of the models to explain the data
- ▶ We need to average (integrate) the likelihood density function over all possible values of θ , weighting by the prior

$$p(D | M_1) = \int_{\theta} p(D | \theta, M_1)p(\theta | M_1)d\theta$$

- ▶ The marginal likelihoods of our model is sensitive to the prior no matter how much we flip the coin.

- ▶ Why do we care about the marginal likelihood?

- ▶ Why do we care about the marginal likelihood?
- ▶ It's *the evidence* that updates our prior to give us the posterior probability of the model

- ▶ Why do we care about the marginal likelihood?
- ▶ It's *the evidence* that updates our prior to give us the posterior probability of the model

$$p(M_1 | D) = \frac{p(D | M_1)p(M_1)}{p(D | M_1)p(M_1) + p(D | M_2)p(M_2)}$$

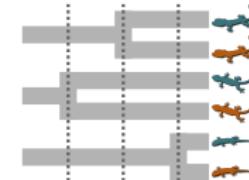
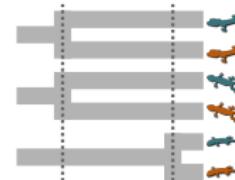
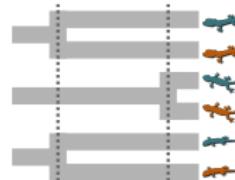
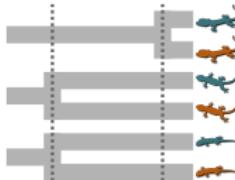
$$p(m_1 | \mathbf{X})$$

$$p(m_2 | \mathbf{X})$$

$$p(m_3 | \mathbf{X})$$

$$p(m_4 | \mathbf{X})$$

$$p(m_5 | \mathbf{X})$$



We want to infer the model and divergence times given genetic data

$$p(m_i | \mathbf{X}) \propto p(\mathbf{X} | m_i) p(m_i)$$

$$p(\mathbf{X} | m_i) = \int_{\theta} p(\mathbf{X} | \theta, m_i) p(\theta | m_i) d\theta$$

- ▶ Divergence times
- ▶ Substitution parameters
- ▶ Gene trees
- ▶ Demographic parameters

Ecoevoluty: Estimating evolutionary coevality

-
- ¹ R. M. Neal (2000). *Journal of Computational and Graphical Statistics* 9: 249–265
 - ² D. Bryant et al. (2012). *Molecular Biology and Evolution* 29: 1917–1932

Ecoevoluty: Estimating evolutionary coevality

- ▶ CTMC model of characters evolving along genealogies
- ▶ Coalescent model of genealogies branching within populations
- ▶ Dirichlet-process prior across divergence models
- ▶ Gibbs sampling¹ to numerically sample models
- ▶ Analytically integrate over genealogies²

¹ R. M. Neal (2000). *Journal of Computational and Graphical Statistics* 9: 249–265

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

Ecoevolity: Estimating evolutionary coevolution

- ▶ CTMC model of characters evolving along genealogies
- ▶ Coalescent model of genealogies branching within populations
- ▶ Dirichlet-process prior across divergence models
- ▶ Gibbs sampling¹ to numerically sample models
- ▶ Analytically integrate over genealogies²
- ▶ *Goal: Fast, full-likelihood Bayesian method to infer patterns of co-diversification from genome-scale data*

¹ R. M. Neal (2000). *Journal of Computational and Graphical Statistics* 9: 249–265

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

Sampling divergence models

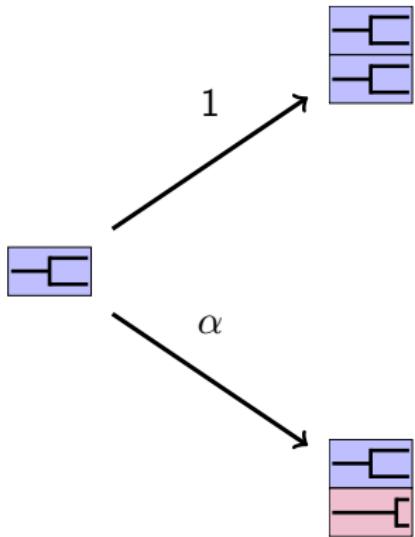
Sampling divergence models

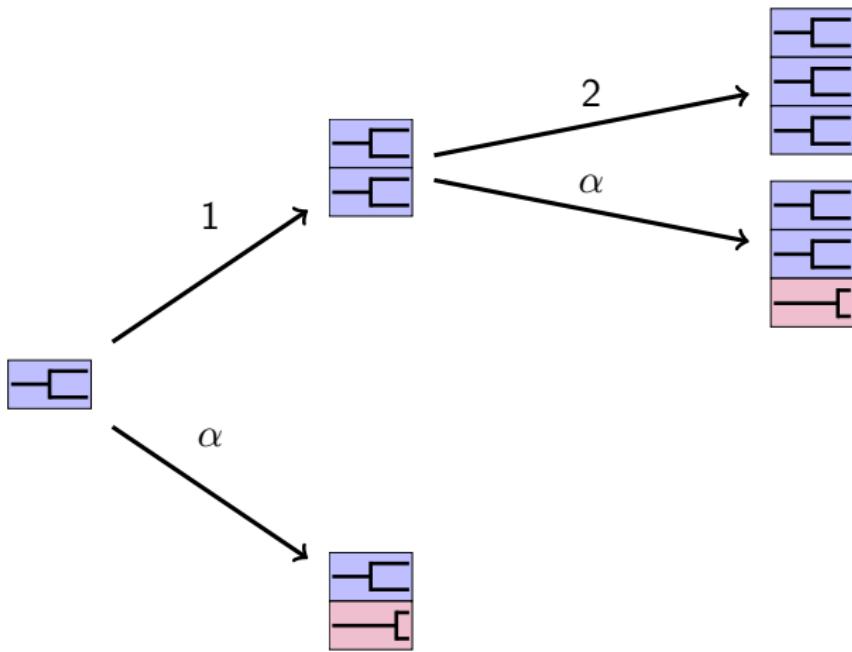
- ▶ The divergence models are ways of assigning our taxa to events
- ▶ A Dirichlet process prior (DPP) model is a convenient and flexible solution
 - ▶ Common Bayesian approach to assigning variables to an unknown number of categories
 - ▶ Controlled by “concentration” parameter: α

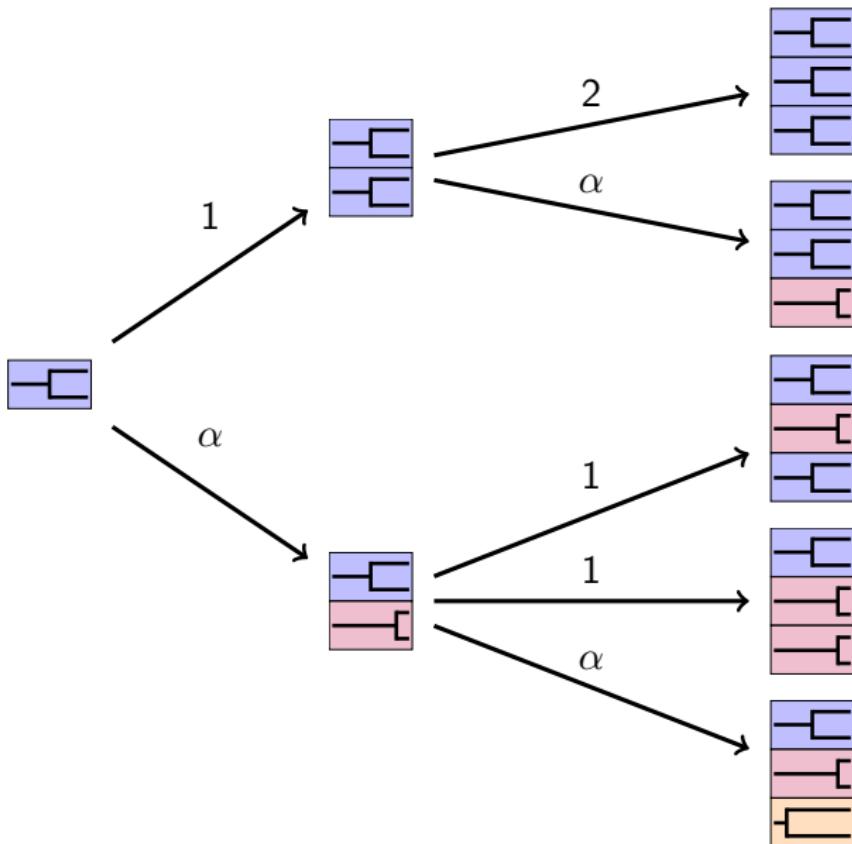


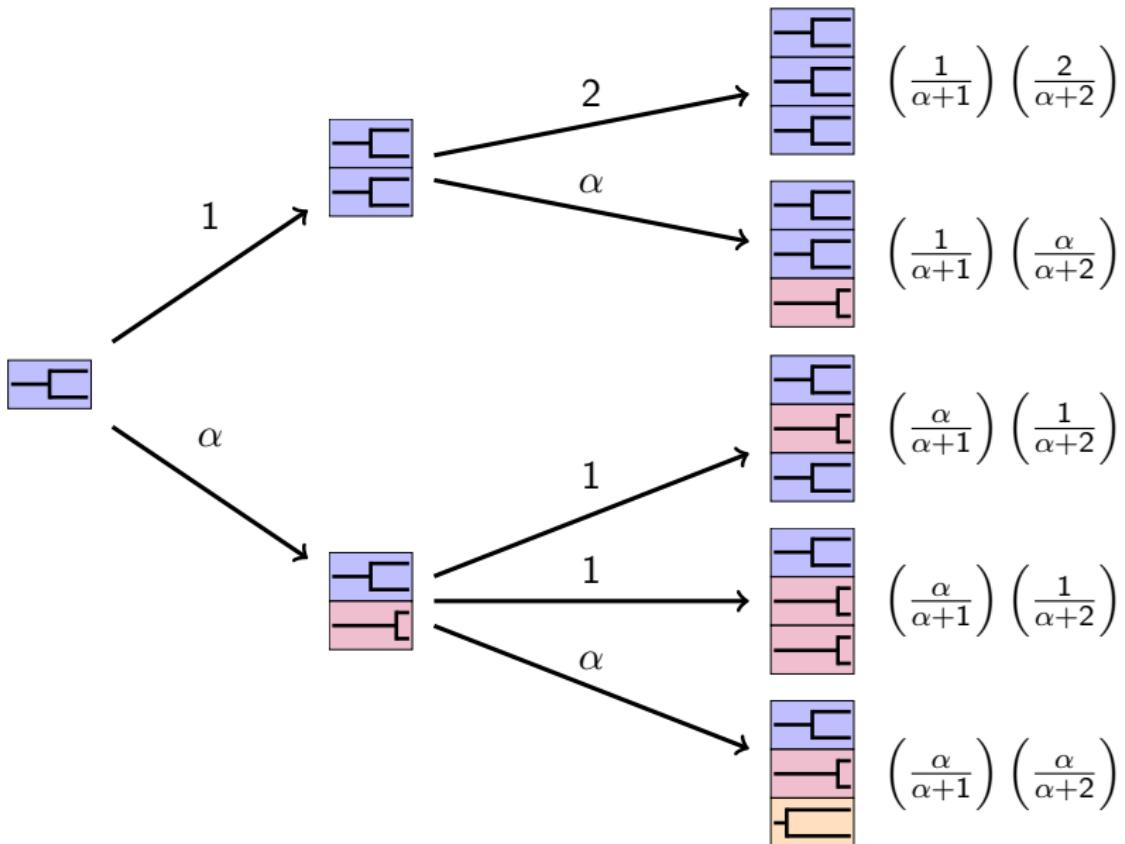
Peter Dirichlet



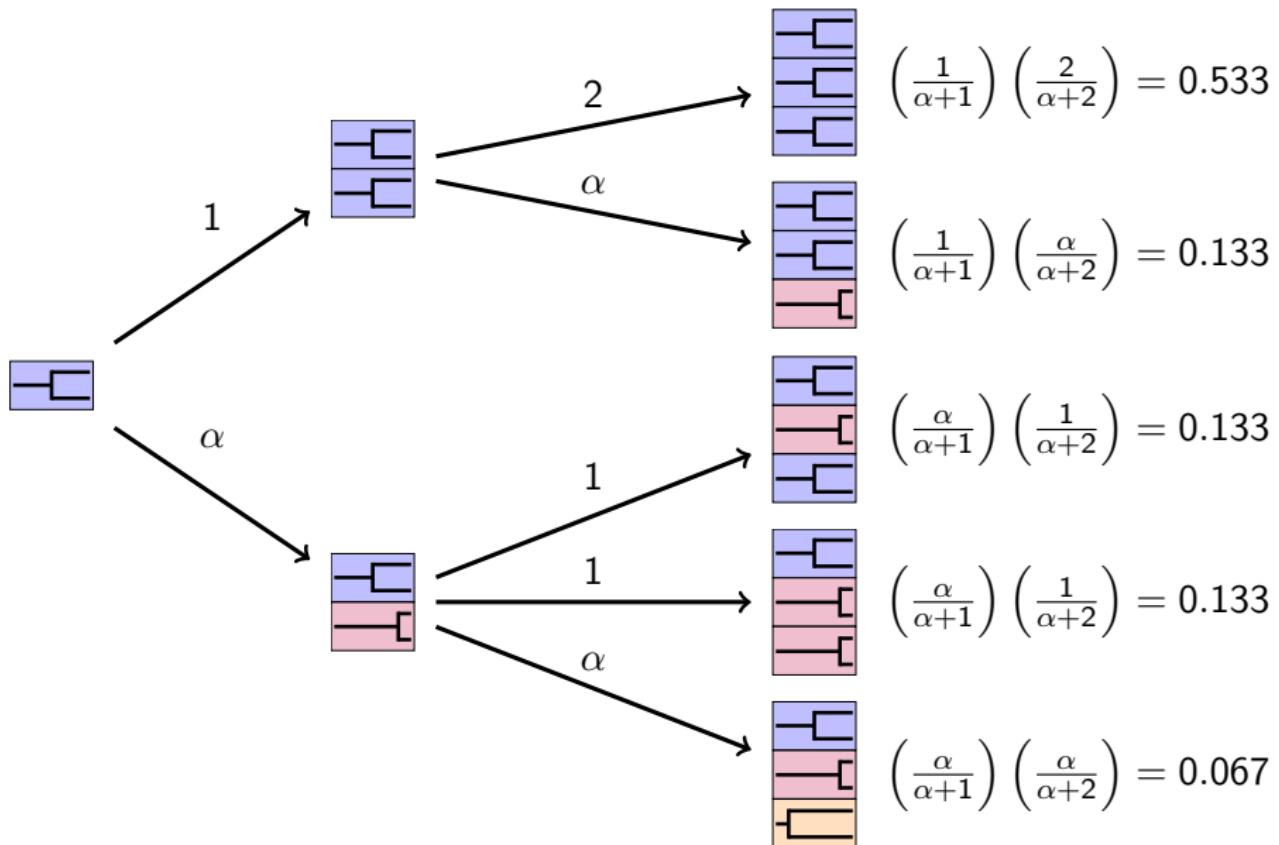




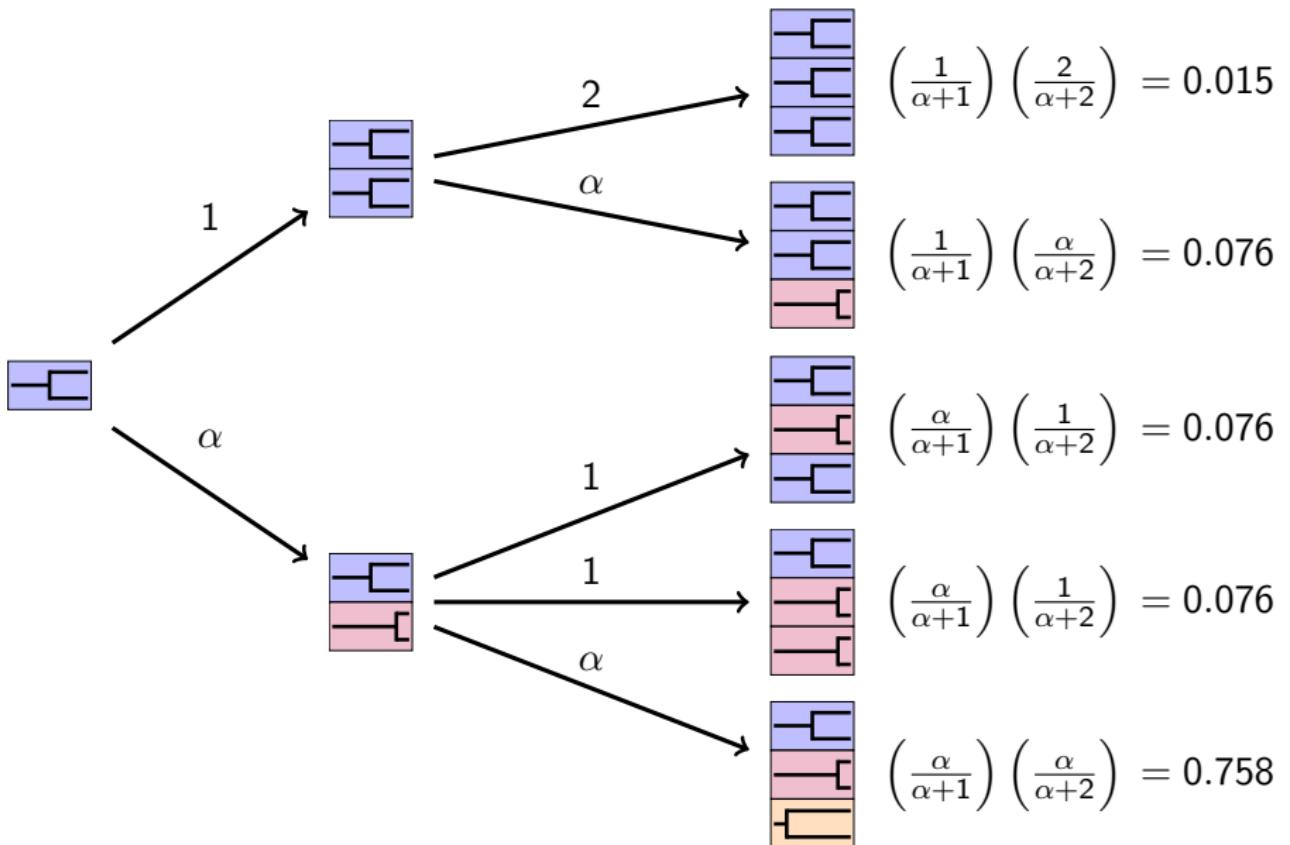




$$\alpha = 0.5$$

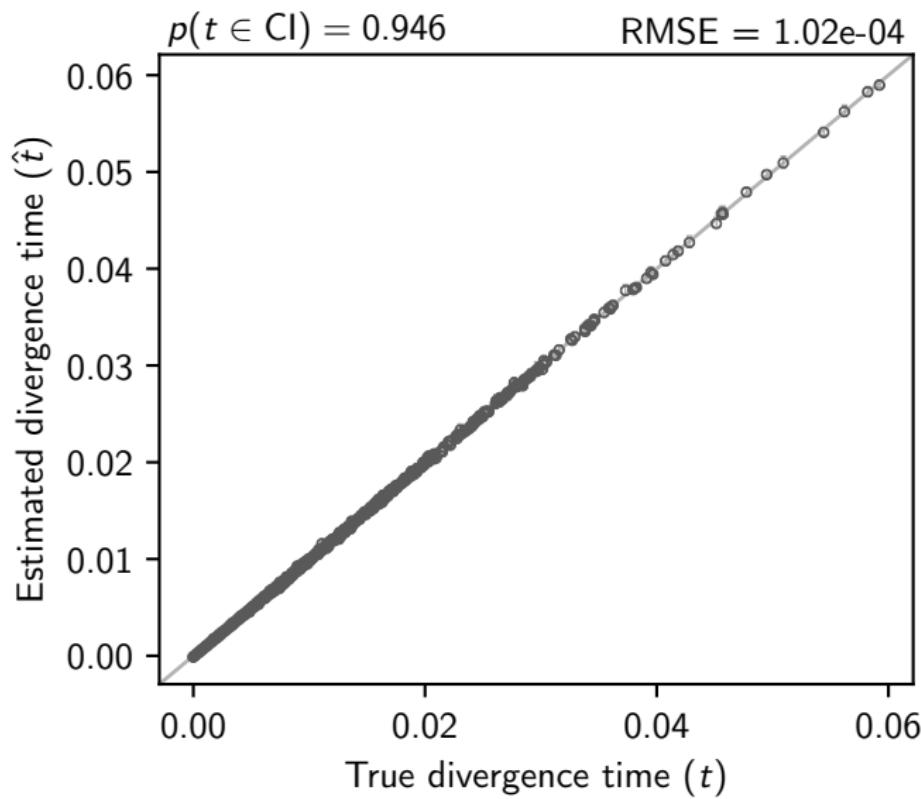


$$\alpha = 10.0$$

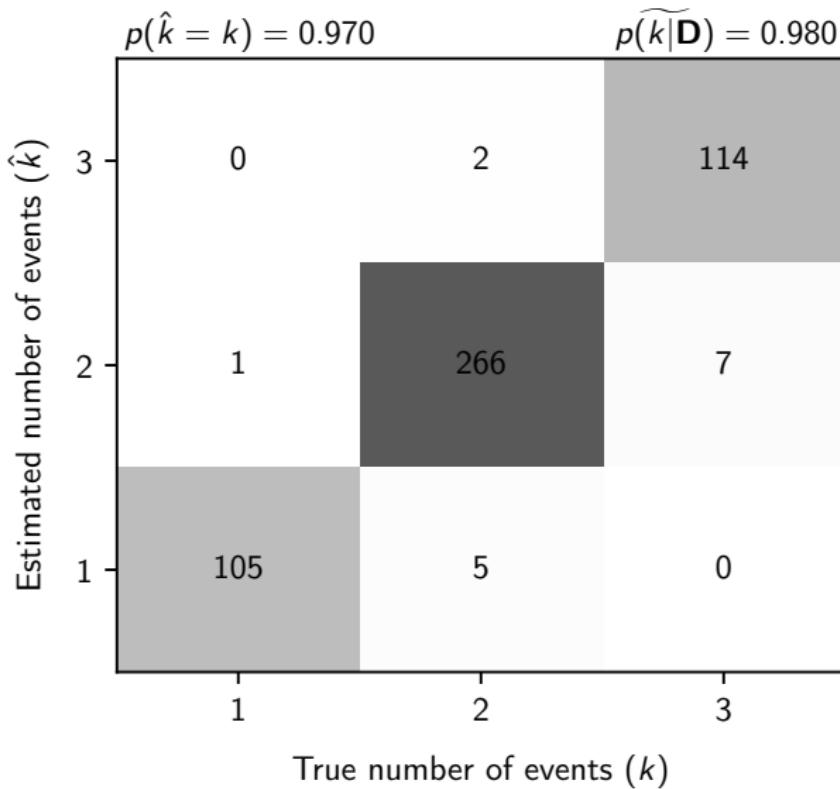


Does it work?

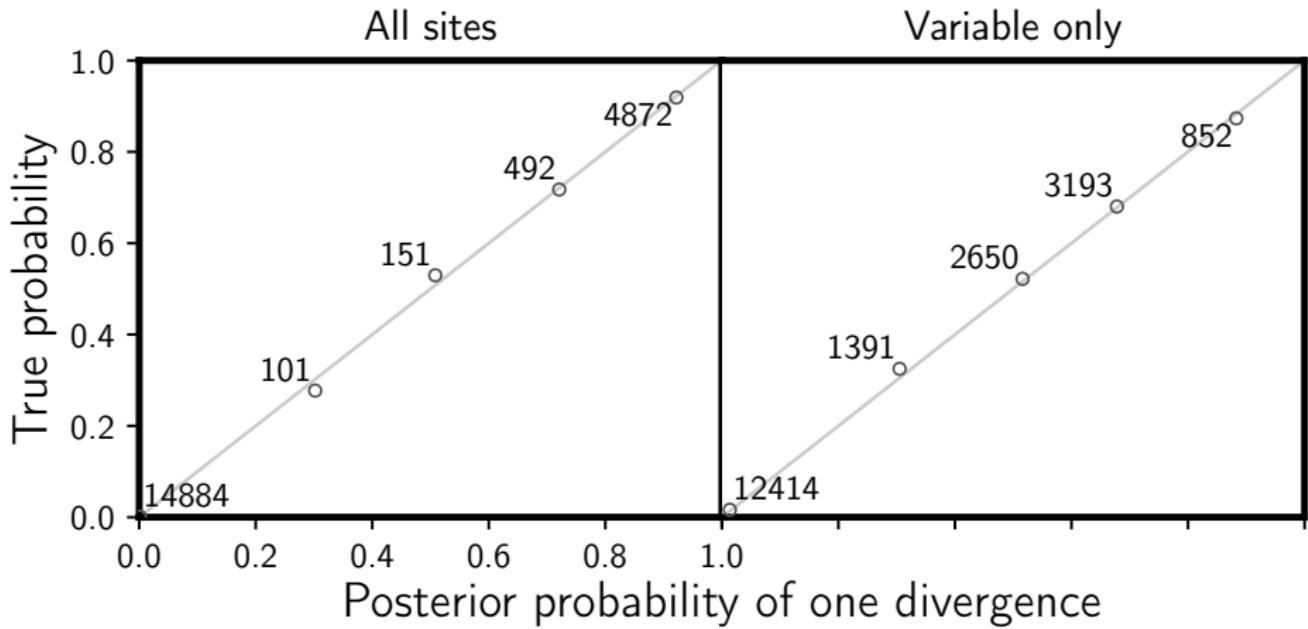
Ecoevolvity: Simulation-based assessment



Ecoevolvity: Simulation-based assessment



Ecoevolvity: Simulation-based assessment

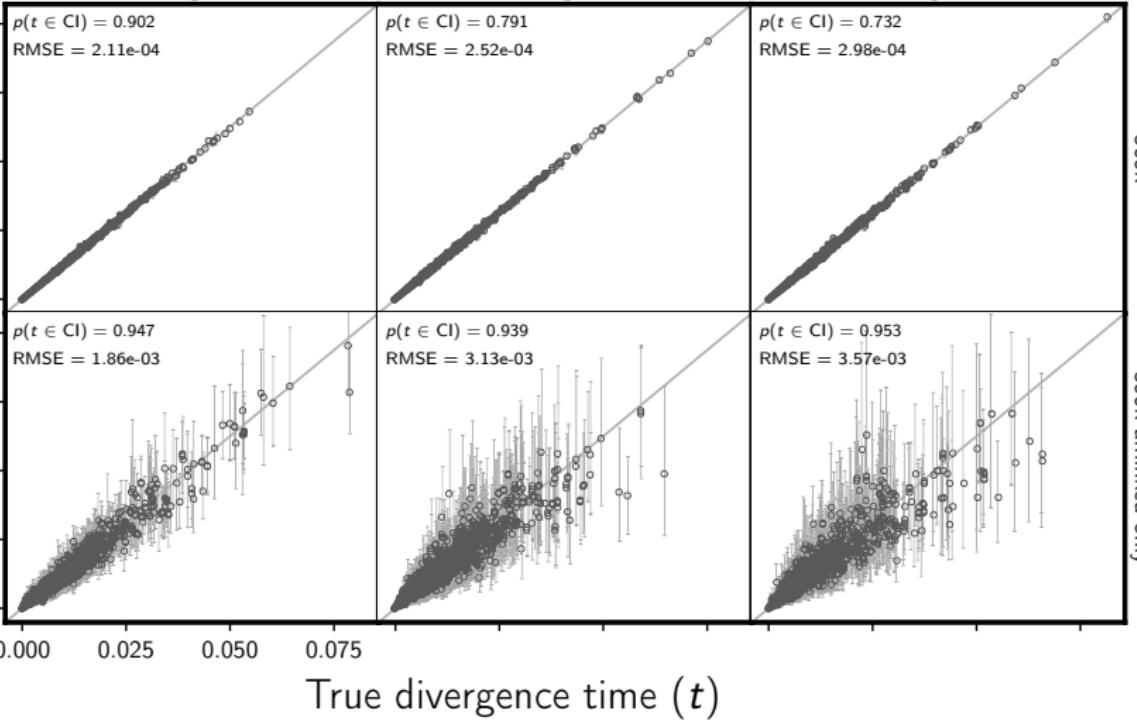


Estimated divergence time (\hat{t})

Locus length = 100

Locus length = 500

Locus length = 1000



Estimated number of events (\hat{k})

Locus length = 100

Locus length = 500

Locus length = 1000

Locus length = 100			Locus length = 500			Locus length = 1000		
$p(\hat{k} = k) = 0.946$			$p(\hat{k} D) = 0.965$			$p(\hat{k} = k) = 0.878$		
0	3	110	0	13	112	0	44	119
3	264	12	24	224	17	34	216	5
99	9	0	103	7	0	78	4	0
$p(k \in CS) = 0.990$			$p(k \in CS) = 0.966$			$p(k \in CS) = 0.946$		
$p(\hat{k} = k) = 0.684$			$p(\hat{k} D) = 0.636$			$p(\hat{k} = k) = 0.668$		
0	12	49	0	8	38	0	16	41
16	200	67	18	209	78	43	194	64
93	54	9	87	53	9	91	41	10
$p(k \in CS) = 0.988$			$p(k \in CS) = 0.998$			$p(k \in CS) = 0.992$		

True number of events (k)

Caveats

- ▶ We have to make strong prior assumptions about the relative rates of mutation among our taxa.
- ▶ The model assumes no migration after divergence