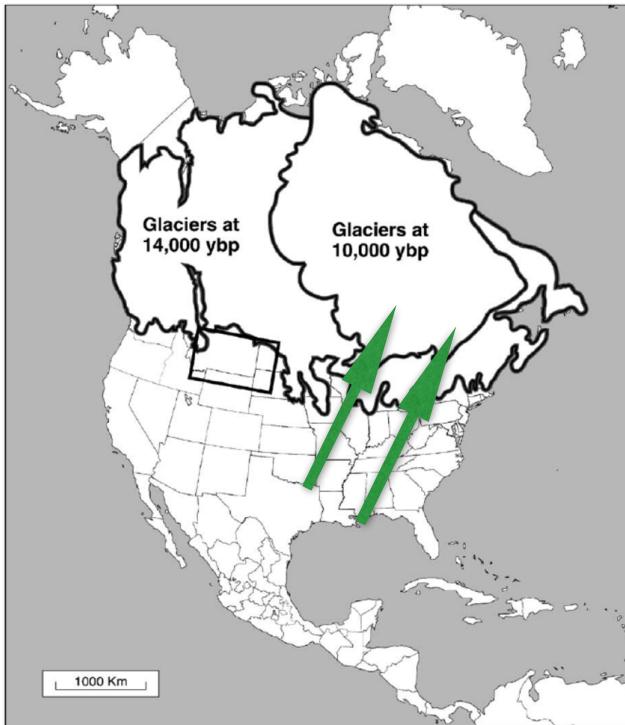


# HABC - step by step part 1



Michael Hickerson  
City University of New York  
City College of New York  
American Museum of Natural History

# HABC - worked example



- How did a temperate assemblage respond to Holocene warming?

*Ecology Letters*, (2016)

doi: 10.1111/ele.12695

LETTER

## Asynchronous demographic responses to Pleistocene climate change in Eastern Nearctic vertebrates

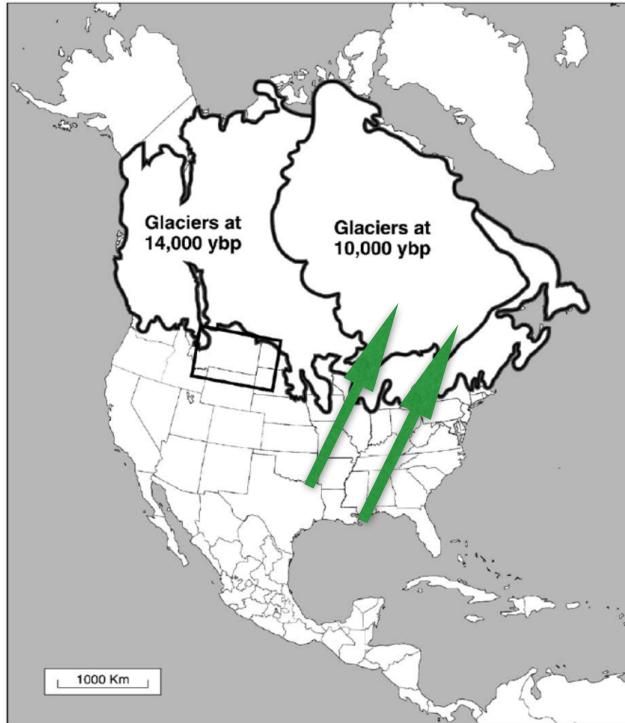
Frank T. Burbrink,<sup>1\*</sup> Yvonne L. Chan,<sup>2</sup> Edward A. Myers,<sup>3,4</sup> Sara Ruane,<sup>5</sup> Brian Tilston Smith<sup>6</sup> and Michael J. Hickerson<sup>4,7,8</sup>



Yvonne Chan - U.Hawaii



Brian Smith - AMNH



Frank Burbrink - AMNH

- How did a temperate assemblage respond to Holocene warming?

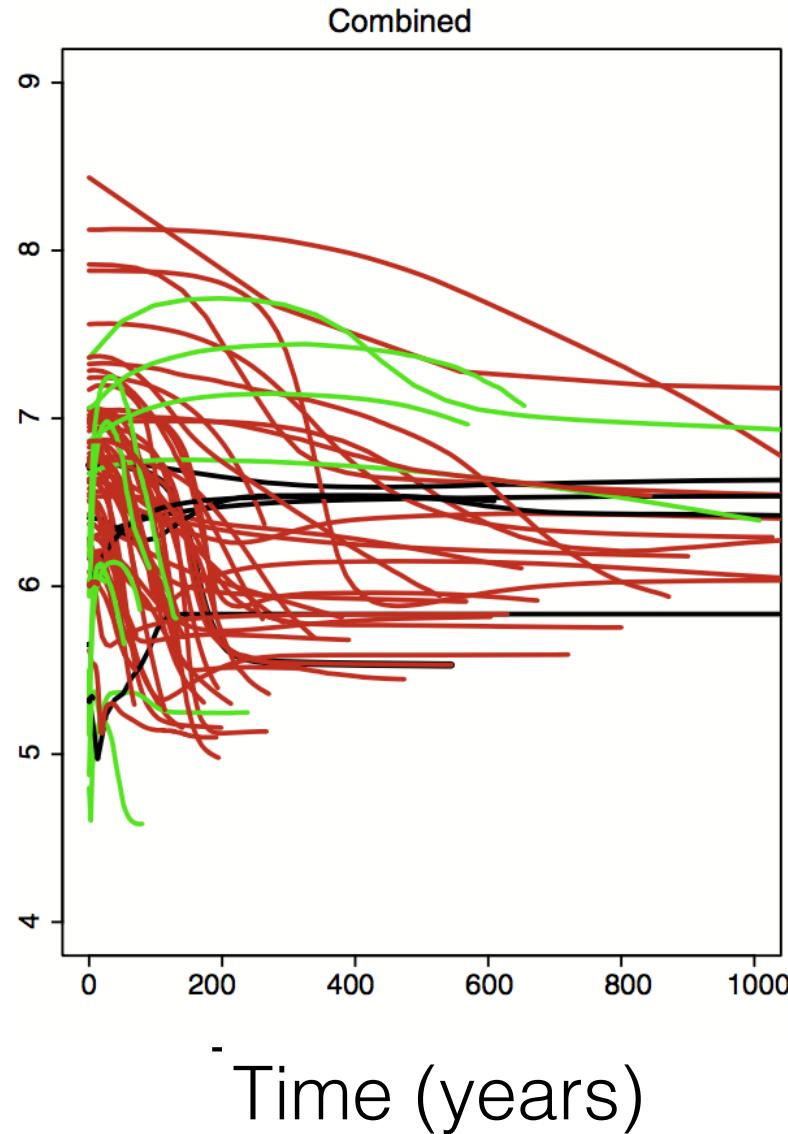
## 90 lineages:

birds  
mammals  
frogs  
salamanders  
lizards  
snakes

mtDNA  
10 - 400 individuals per taxon  
hierarchical Bayesian model

> 90 Skyline plots (eastern Nearctic - w/ mtDNA)

Effective  
Pop Size

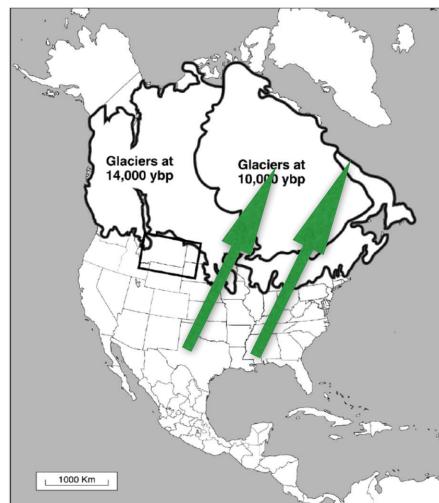


- Impractical for large number of species
- Difficult to statistically compare results

# hierarchical Bayesian model

$$P(\zeta, \tau, \phi \mid \text{Data}) \propto p(\text{Data} \mid \zeta, \tau, \phi) = p(\text{Data} \mid \phi, \tau) p(\phi) p(\tau \mid \zeta) p(\zeta)$$

flexible  
within species uncertainty  
across species estimates



Chan et al. 2014

Xue and Hickerson 2014

Xue and Hickerson 2017

# hierarchical Bayesian model

$$P(\zeta, \tau, \phi \mid \text{Data}) \propto p(\text{Data} \mid \zeta, \tau, \phi) = p(\text{Data} \mid \phi, \tau) p(\phi) p(\tau \mid \zeta) p(\zeta)$$

hyperparameters

$\zeta$  = Co-expansion Proportion

$\tau_s$  = Co-expansion time



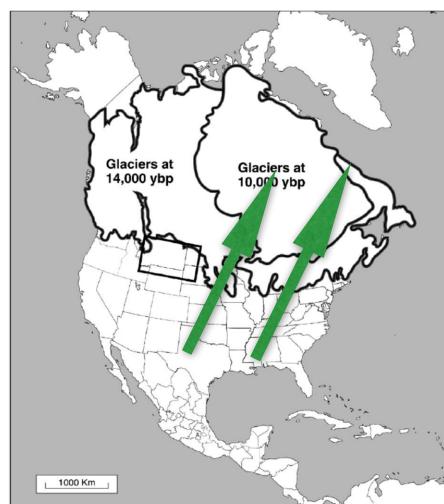
species-specific parameters ( $\phi, \tau$ )

$\tau$  = expansion times

$\{\tau_1, \dots, \tau_n\}$

$\phi = N$ , expansion magnitudes (independent)

$\{\phi_1, \dots, \phi_n\}$      $\{N_1, \dots, N_n\}$



Chan et al. 2014

Xue and Hickerson 2014

# hierarchical Bayesian model

$$P(\zeta, \tau, \phi \mid \text{Data}) \propto p(\text{Data} \mid \zeta, \tau, \phi) = p(\text{Data} \mid \phi, \tau) p(\phi) p(\tau \mid \zeta) p(\zeta)$$

hyperparameters

$\zeta$  = Co-expansion Proportion

$\tau_s$  = Co-expansion time

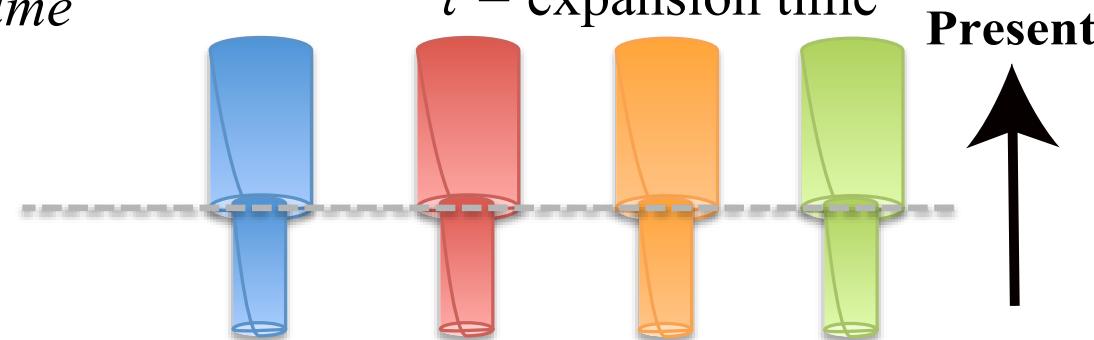
A) Synchronous expansion

$\zeta = 1.0$

species-specific parameters ( $\phi, \tau$ )

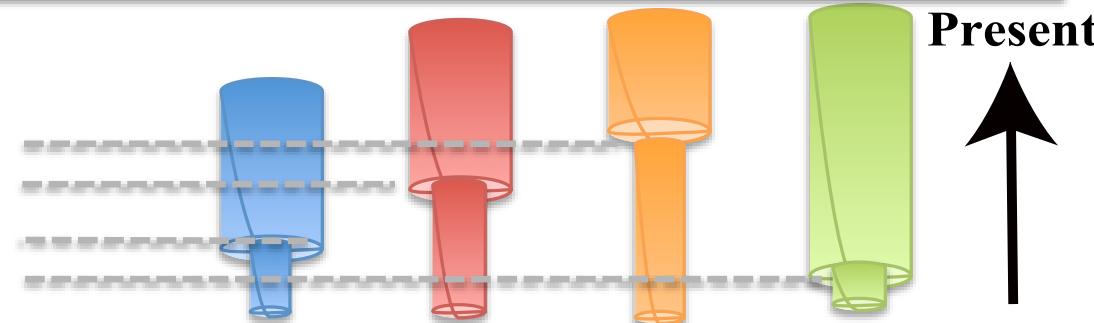
$\phi = N$ , expansion magnitude

$\tau =$  expansion time



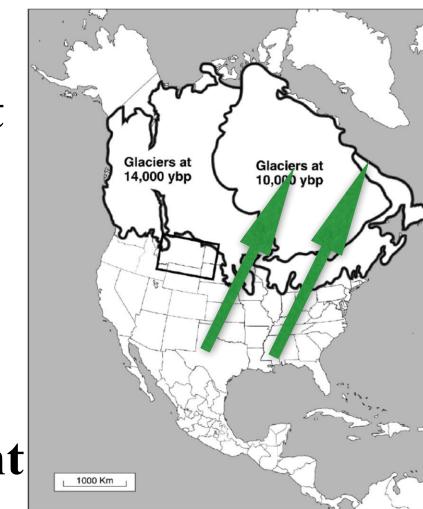
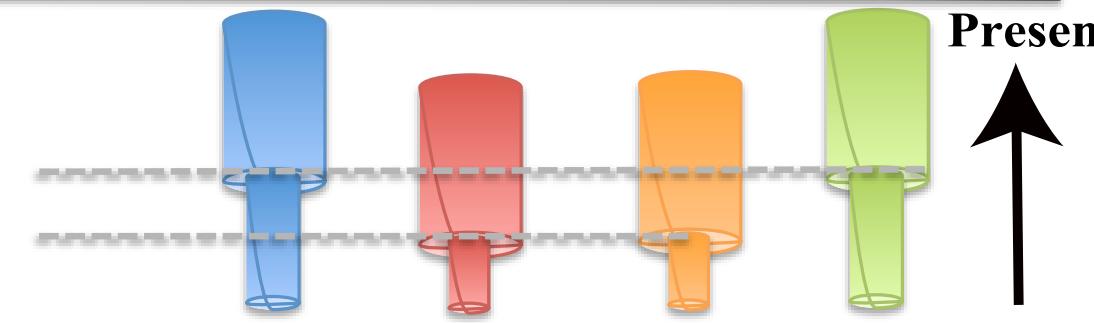
B) Asynchronous expansion

$\zeta = 0.0$



C) Expansion mixture:  
asynchronous  
&  
synchronous

$\zeta = 0.5$



Chan et al. 2014

Xue and Hickerson 2014

$$P(\zeta, \tau, \phi \mid \text{Data}) \propto p(\text{Data} \mid \zeta, \tau, \phi) = p(\text{Data} \mid \phi, \tau) p(\phi) p(\tau \mid \zeta) p(\zeta)$$

hard to solve analytically (and we want flexibility)



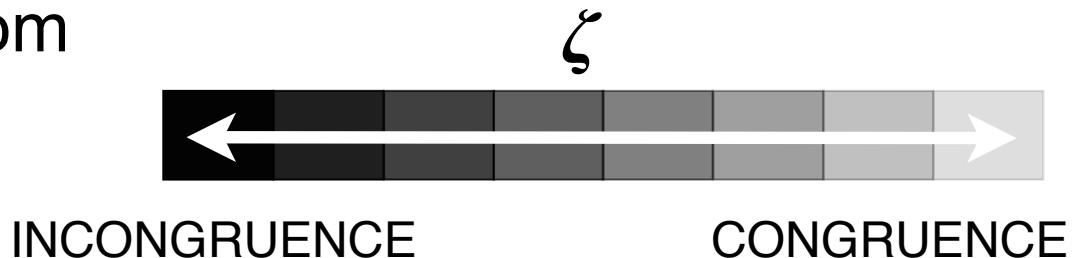
## Approximate Bayesian Computation - ABC

Likelihood of HyperParameter  
value is inversely proportional to  
the difference

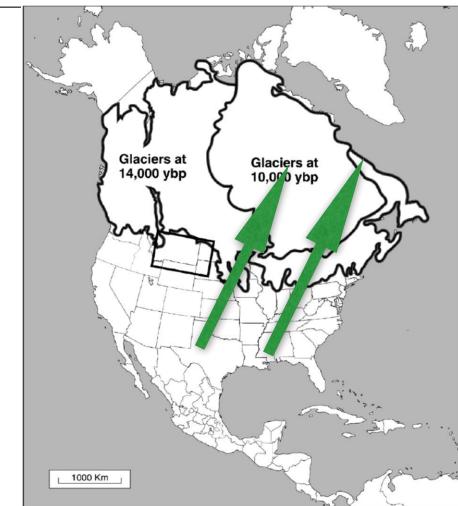
$$\text{Data}_{\text{simulated}} - \text{Data}_{\text{observed}}$$

# Approximate Bayesian Computation - ABC

Simulate Genetic **Data** from  
random values from the  
hyperprior  $p(\phi) p(\tau | \zeta) p(\zeta)$

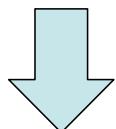


simulate coalescent  
& mutation across  
assemblage

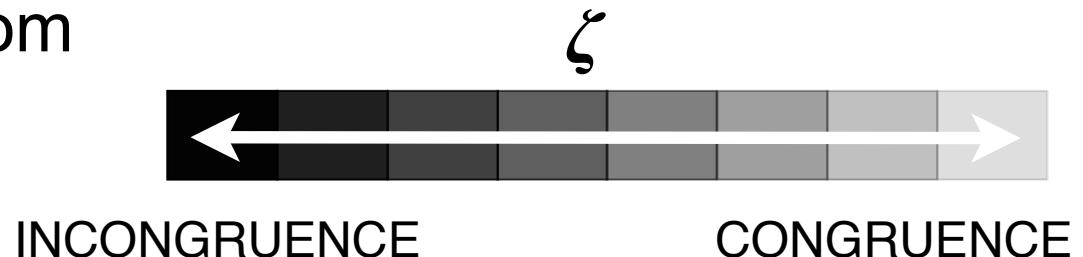


# Approximate Bayesian Computation - ABC

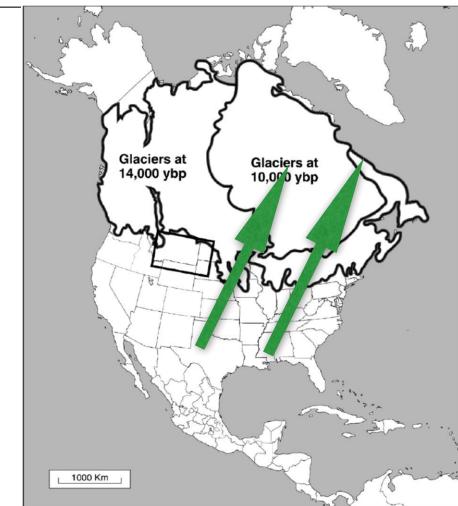
Simulate Genetic **Data** from random values from the hyperprior  $p(\phi) p(\tau | \zeta) p(\zeta)$



Compress **Data** (summary statistics)

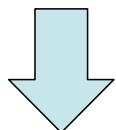


simulate coalescent & mutation across assemblage

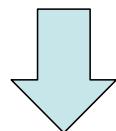


# Approximate Bayesian Computation - ABC

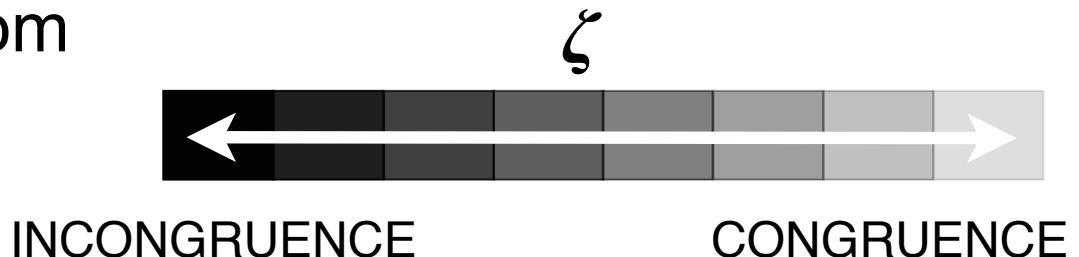
Simulate Genetic **Data** from random values from the hyperprior  $p(\phi) p(\tau | \zeta) p(\zeta)$



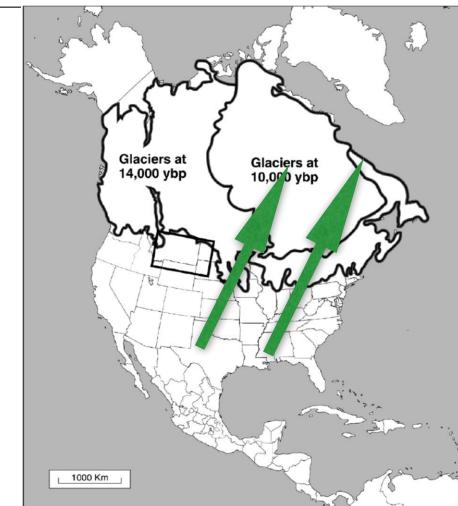
Compress **Data** (summary statistics)



Compare Compressed Simulated **Data** to Compressed Observed **Data**

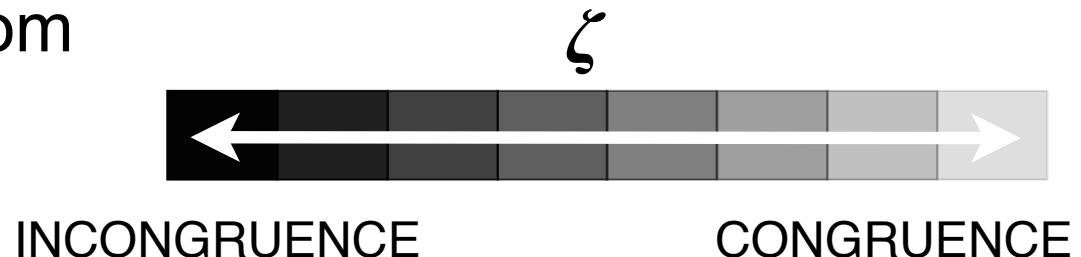
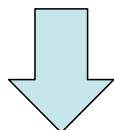


simulate coalescent & mutation across assemblage

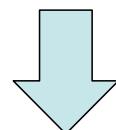


# Approximate Bayesian Computation - ABC

Simulate Genetic **Data** from random values from the hyperprior  $p(\phi) p(\tau | \zeta) p(\zeta)$

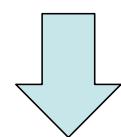


Compress **Data** (summary statistics)



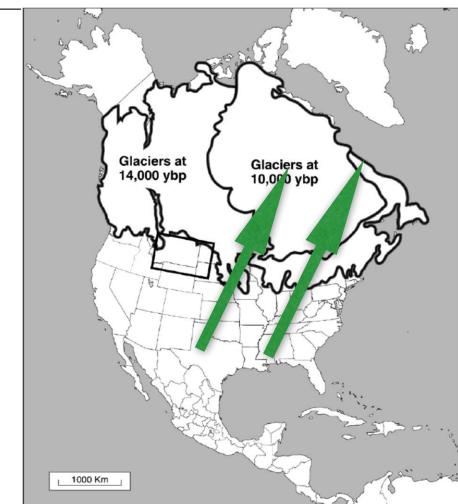
simulate coalescent & mutation across assemblage

Compare Compressed Simulated **Data** to Compressed Observed **Data**



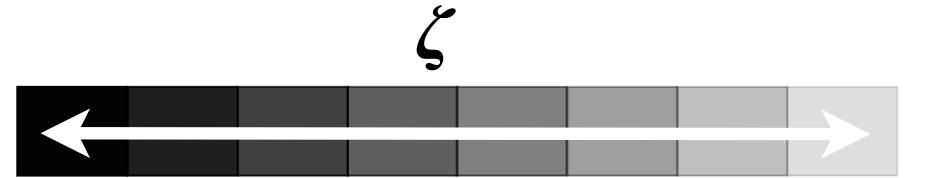
Keep if the Match is Close Enough

$$Data_{simulated} - Data_{observed} \approx 0 ?$$



# Approximate Bayesian Computation - ABC

Simulate Genetic **Data** from random values from the hyperprior  $p(\phi) p(\tau | \zeta) p(\zeta)$



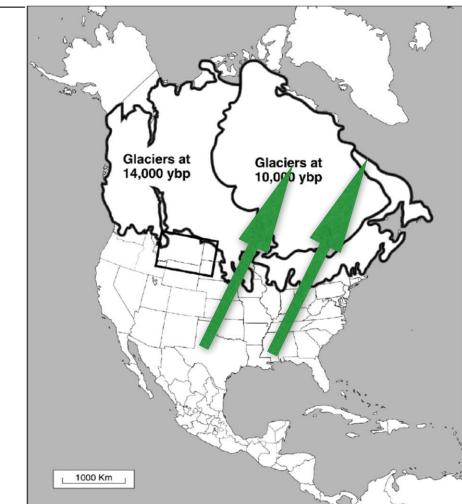
INCONGRUENCE

CONGRUENCE

Compress **Data** (summary statistics)

simulate coalescent & mutation across assemblage

Compare Compressed Simulated **Data** to Compressed Observed **Data**



Keep if the Match is Close Enough

$$Data_{simulated} - Data_{observed} \approx 0 ?$$

# Approximate Bayesian Computation - ABC

Simulate Genetic **Data** from  
random values from the  
hyperprior  $p(\zeta)$

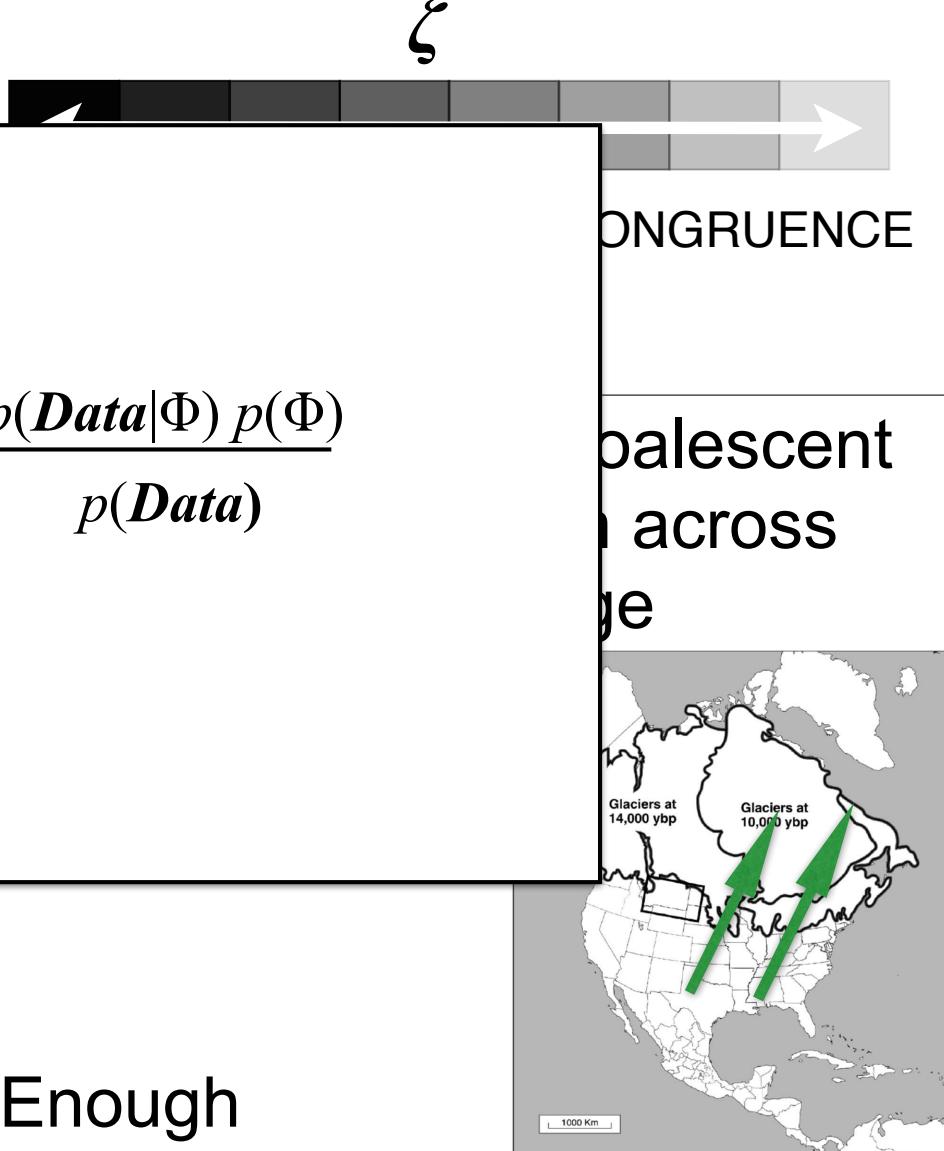
Compress Data

$$P(\Phi | \text{Data}) = \frac{p(\text{Data}|\Phi) p(\Phi)}{p(\text{Data})}$$

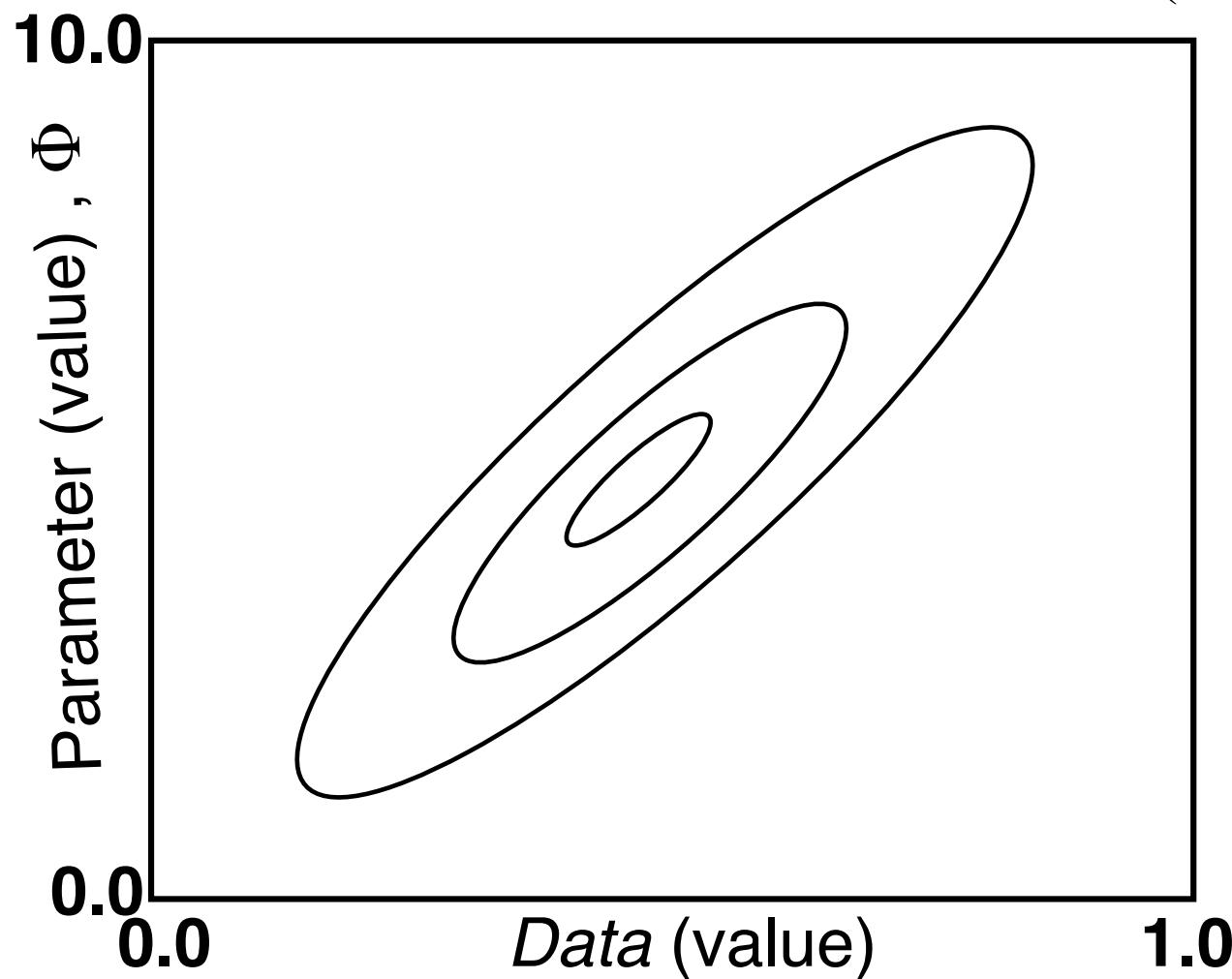
Compare Compressed  
to Compressed Observed

Keep if the Match is Close Enough

$$\text{Data}_{\text{simulated}} - \text{Data}_{\text{observed}} \approx 0 ?$$



$$P(\Phi \mid Data) = \frac{p(Data \mid \Phi) p(\Phi)}{p(Data)}$$

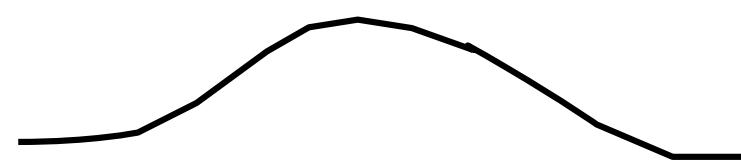
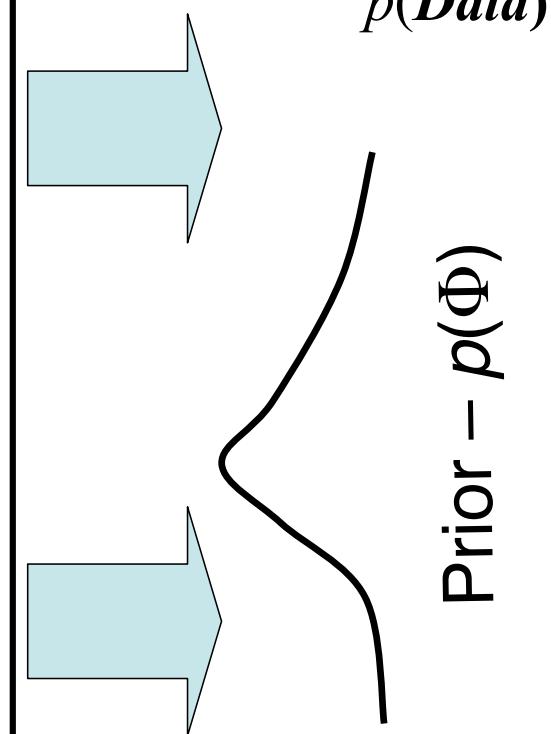
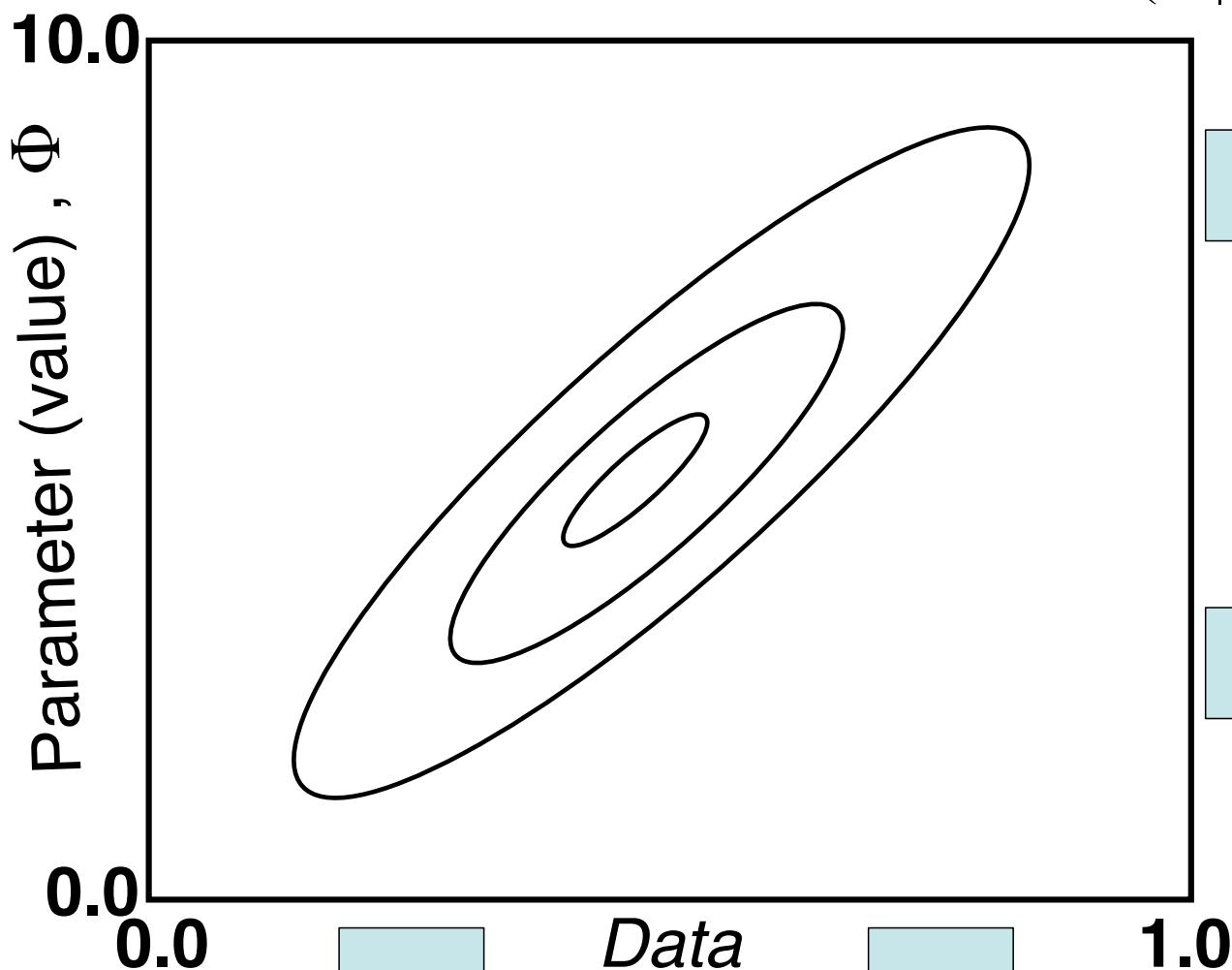


Mark Beaumont's  
ABC schematic

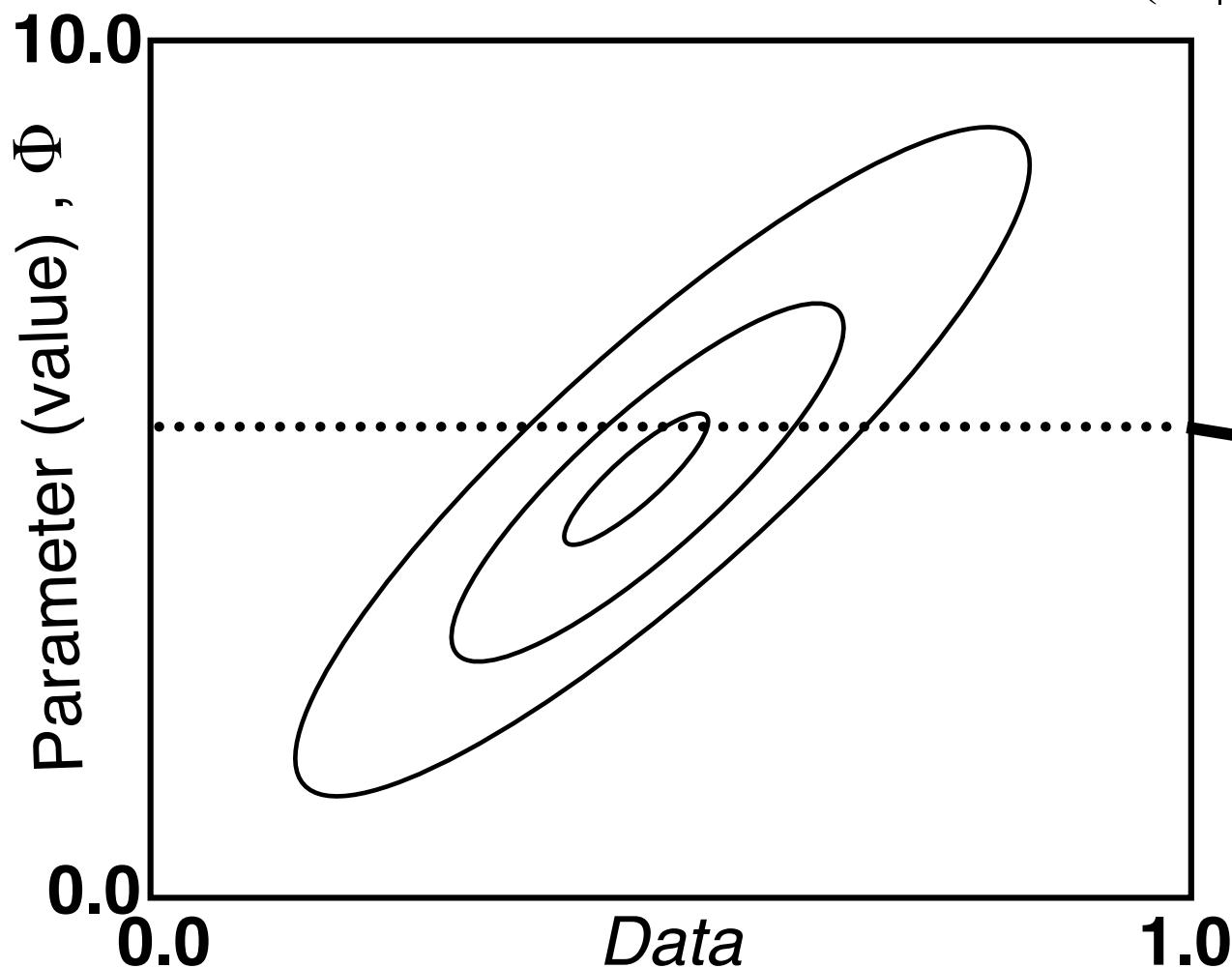
**if**  $Data_{simulated} - Data_{observed} \approx 0 ?$

**then**  $\Phi_{simulated} \approx \Phi_{observed}$

$$P(\Phi \mid Data) = \frac{p(Data \mid \Phi) p(\Phi)}{p(Data)}$$

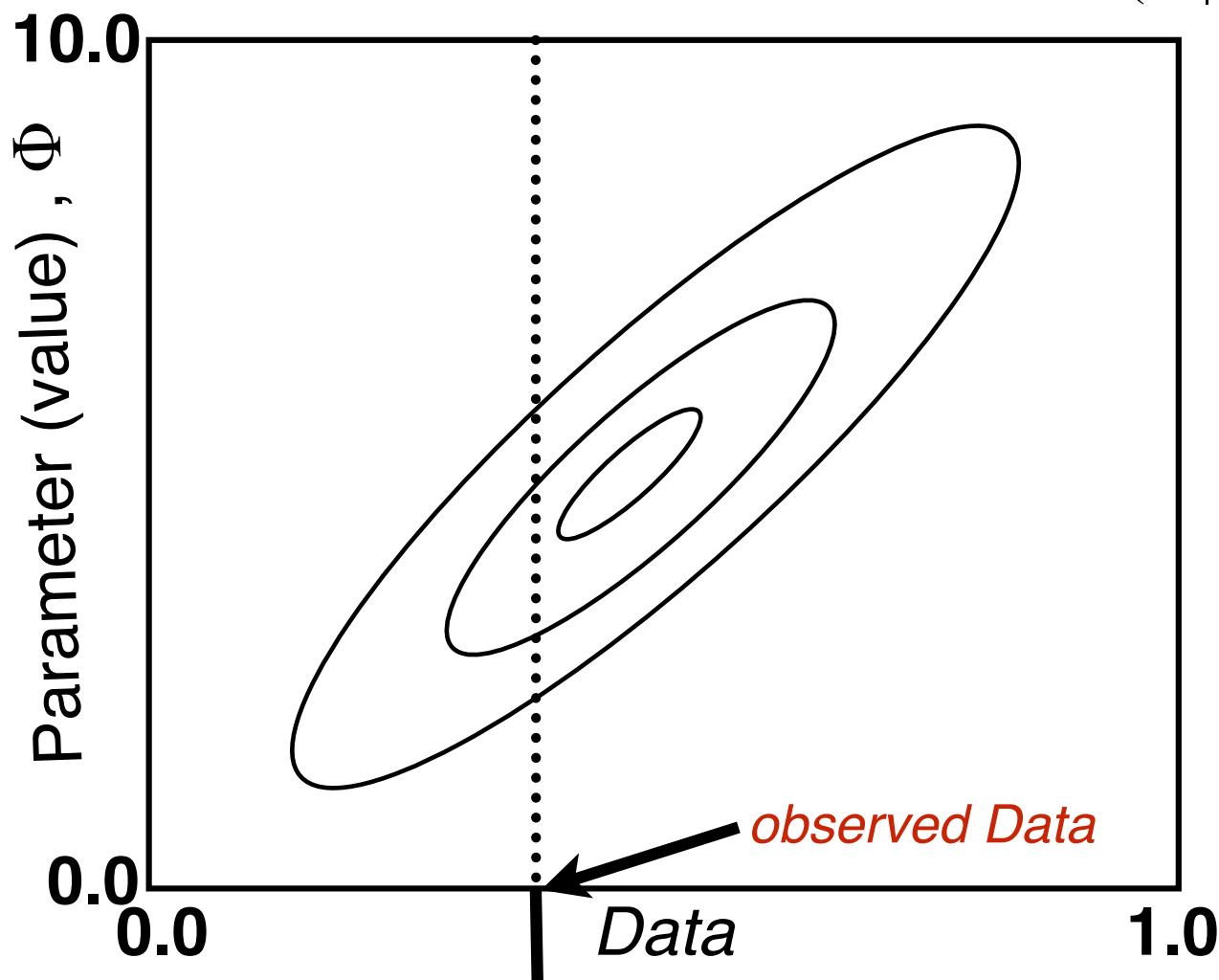


$$P(\Phi \mid \text{Data}) = \frac{p(\text{Data} \mid \Phi) p(\Phi)}{p(\text{Data})}$$



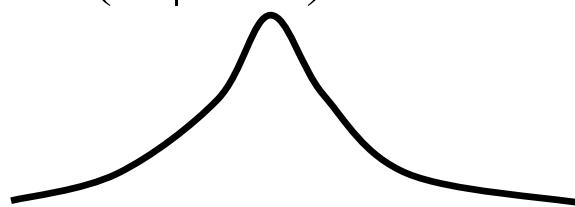
**Likelihood—**  
 $P(\text{Data} \mid \Phi)$

$$P(\Phi \mid Data) = \frac{p(Data \mid \Phi) p(\Phi)}{p(Data)}$$

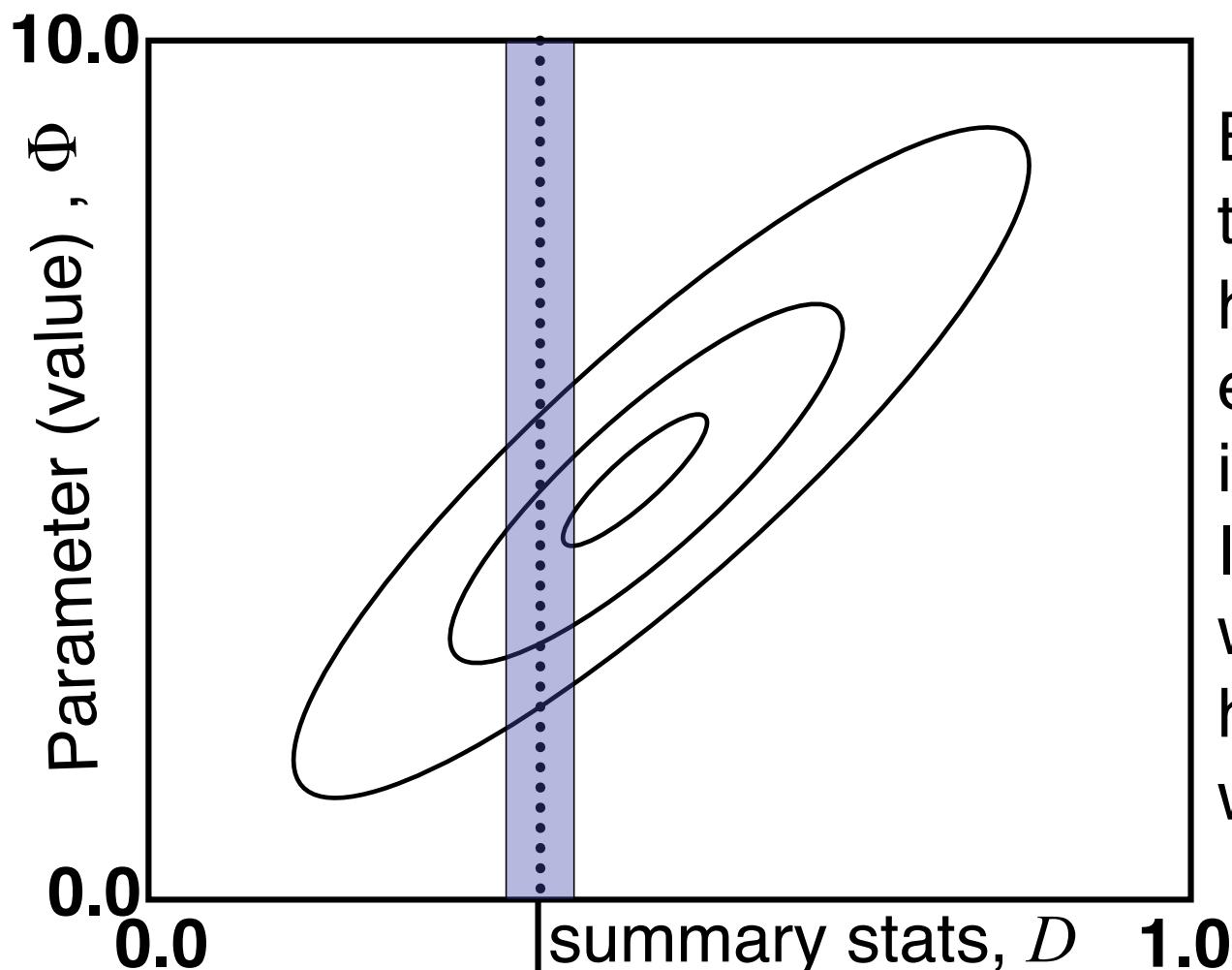


**Posterior distribution –**

$$P(\Phi \mid Data)$$



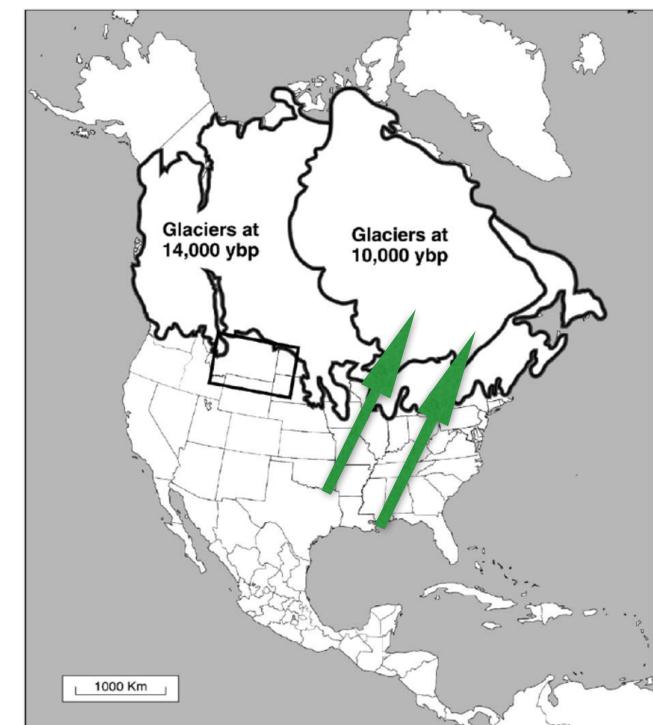
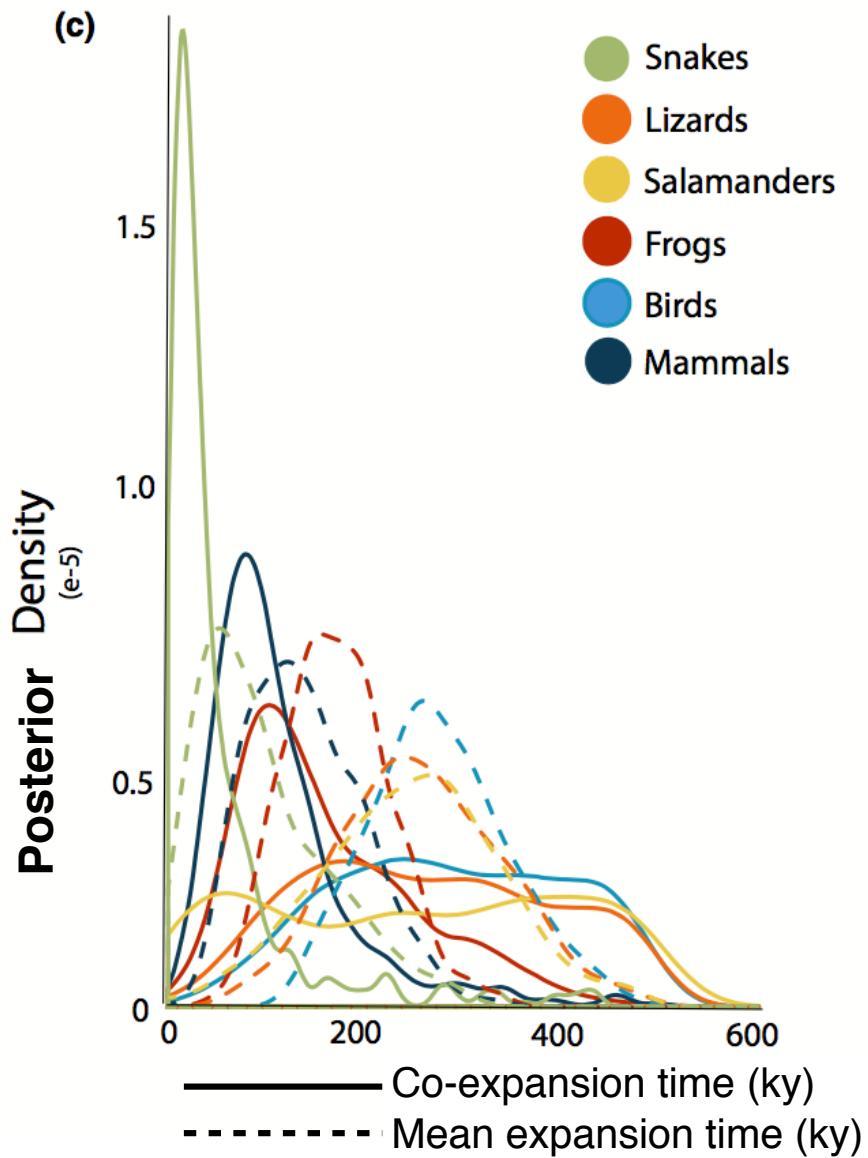
$$P(\Phi \mid \text{Data}) = \frac{p(\text{Data} \mid \Phi) p(\Phi)}{p(\text{Data})}$$



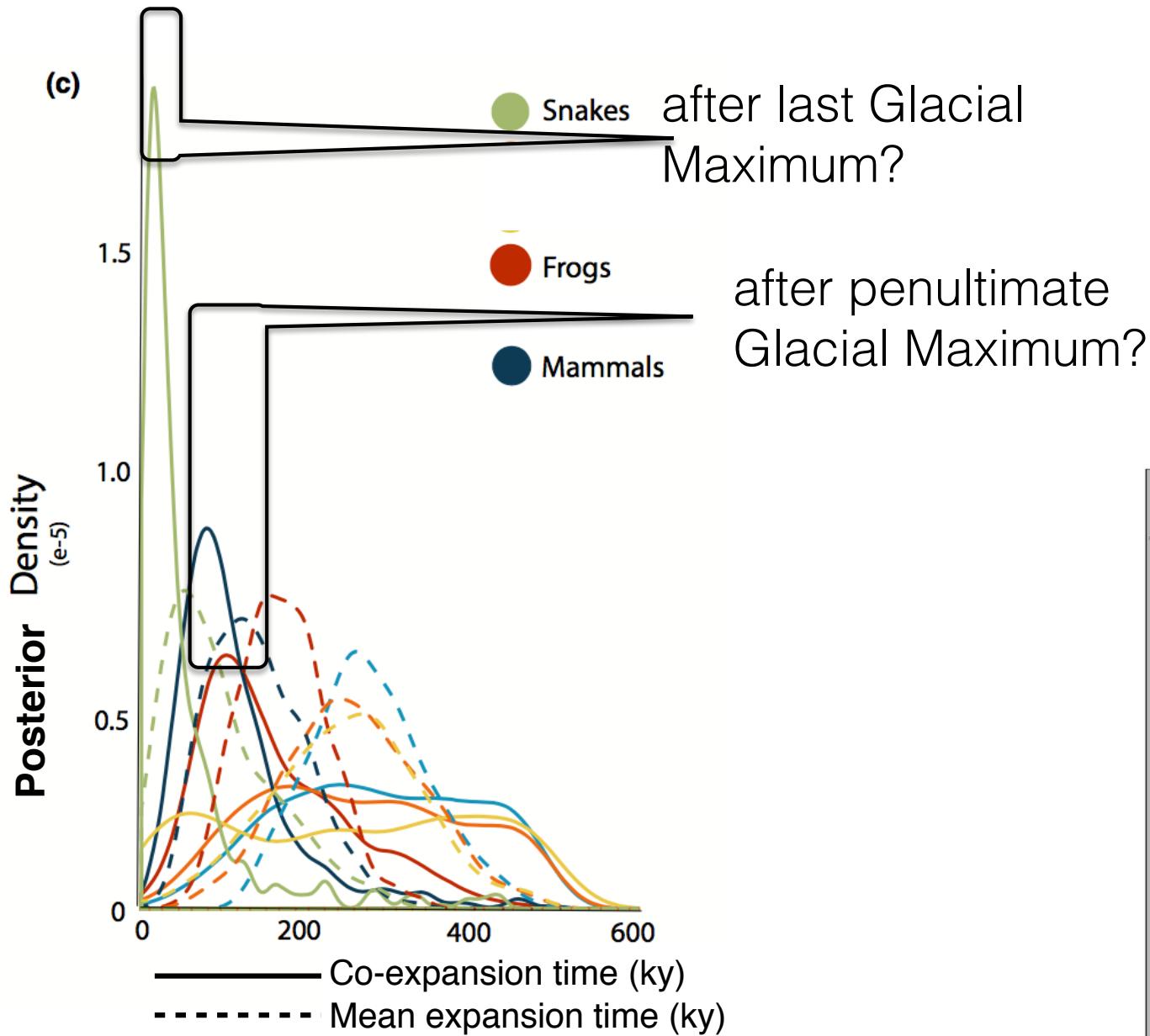
By converting  $\text{Data}$  to summary stats  $D$ , hitting the data exactly is not required. Instead we we need to only hit  $\|D_i - D_{obs}\| < \varepsilon$  within a “threshold”

# Results (hABC)

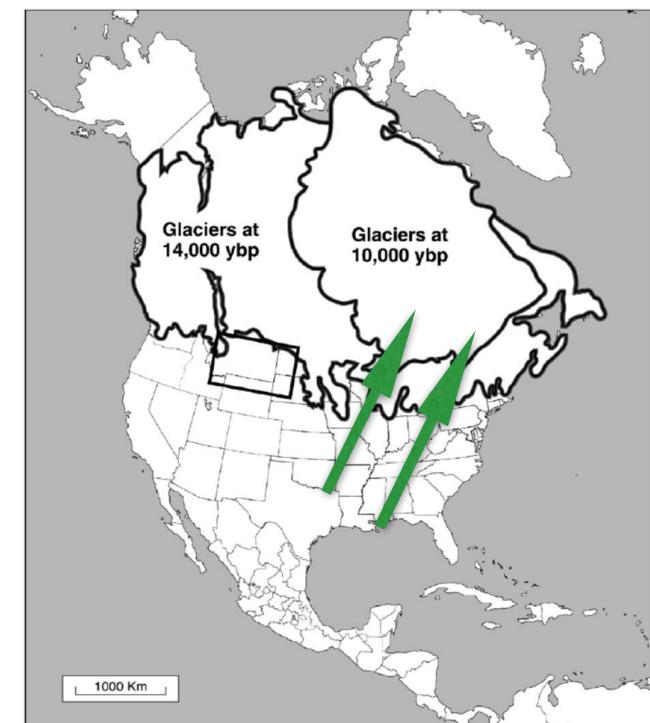
overall, extremely asynchronous, but ....



# Results (hABC)



Strong to moderate recent co-expansion in **Snakes, Frogs and Mammals**



# **Development of ABC inferential strategy**

- 1. Summary statistic Selection**
- 2. Prior Selection & Model Checking**
- 3. Estimator Validation with PODS  
(pseudo-observed data sets)**
- 4. Estimation**

# **Development of ABC inferential strategy**

1. Summary statistic Selection
- 2. Prior Selection & Model Checking**
3. Validation
4. Estimation

## 2. Prior Selection

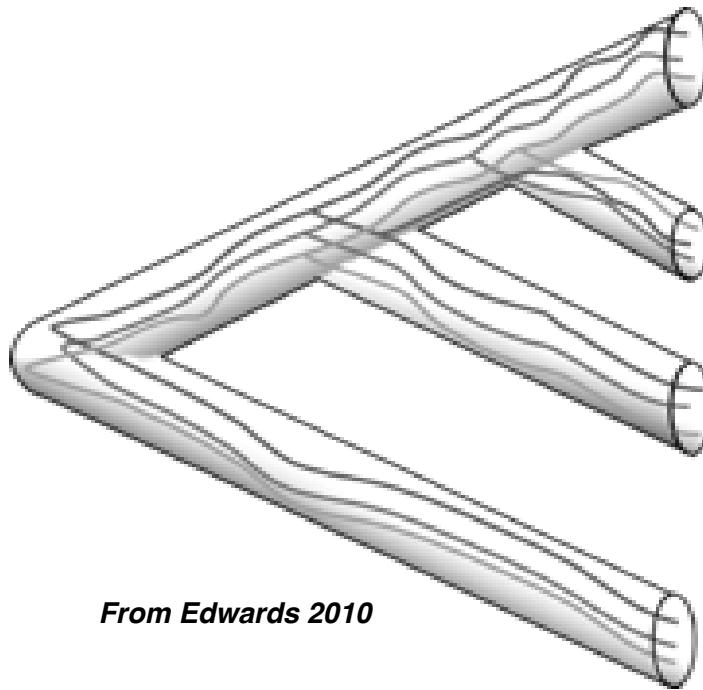
### How do we choose the prior?

1. look at the data
2. graphical checks from trial runs
3. Bayesian model averaging

## 2. Prior Selection

### How do we choose the prior?

1. look at the data

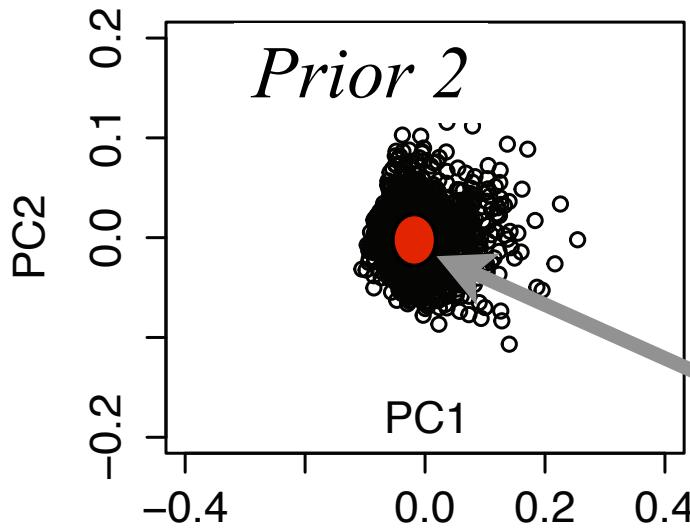


**use gene tree divergences as the maximum  
(population divergence < gene divergence)**

## 2. Prior Selection

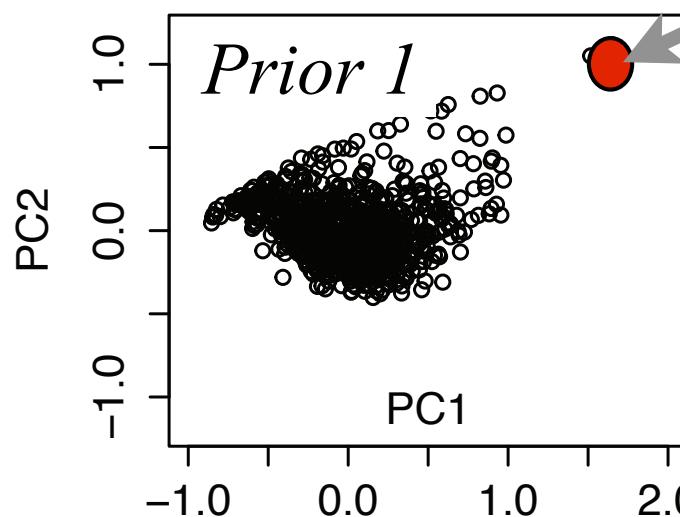
### How do we choose the prior?

2. graphical checks from trial runs (small prior samples)



**well** sampled prior

observed  
data



**poorly** sampled prior

(Hickerson et al. 2014)

## 2. Prior Selection

**How do we choose the prior?**

3. Bayesian model averaging

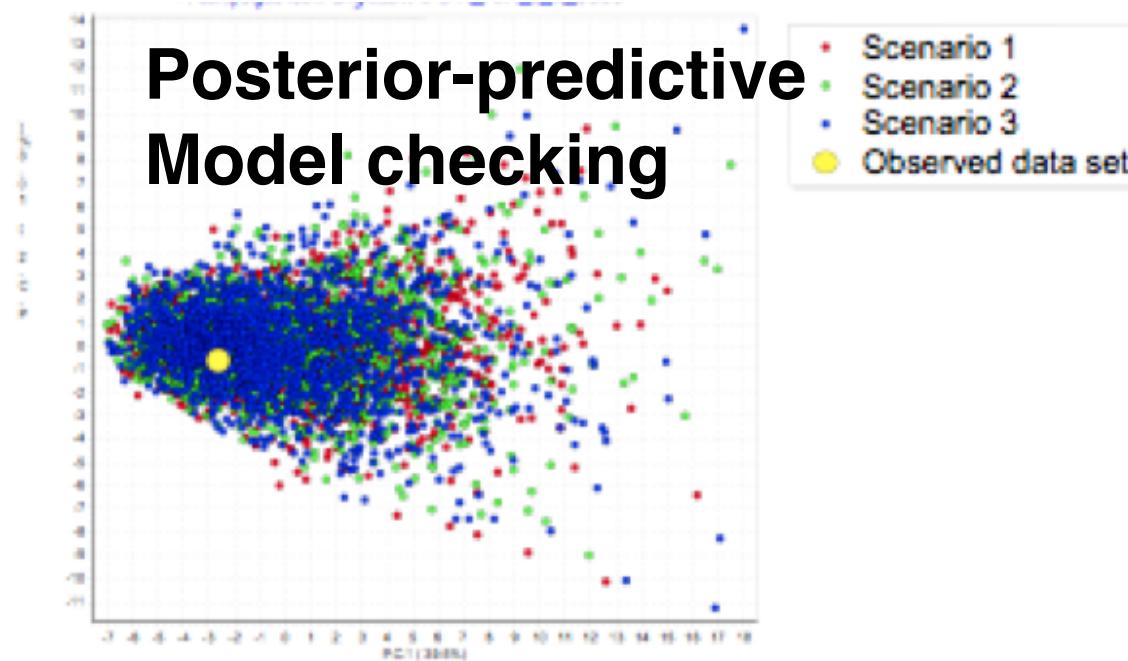
Sample from **multiple** candidate priors

$$\{M_1, \dots, M_8\}$$

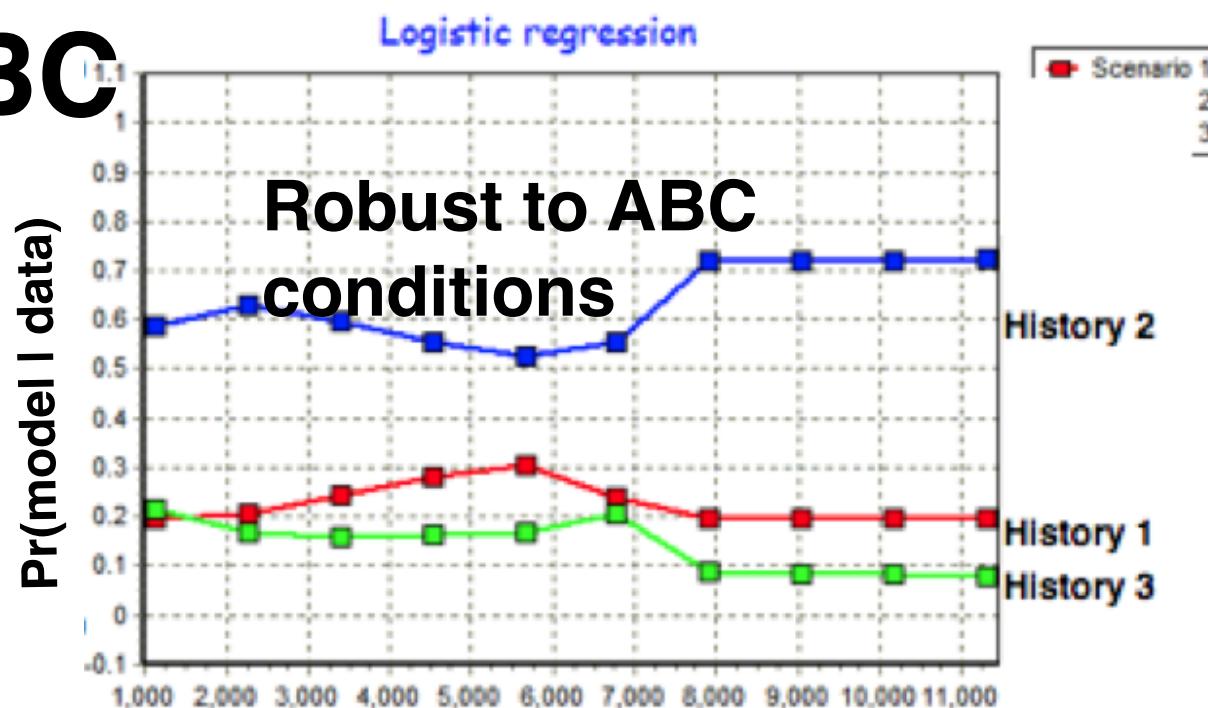
Do regular ABC estimation

ABC posterior estimate is weighted by posterior strength  
of each prior

# Posterior-predictive Model checking



# DIY-ABC



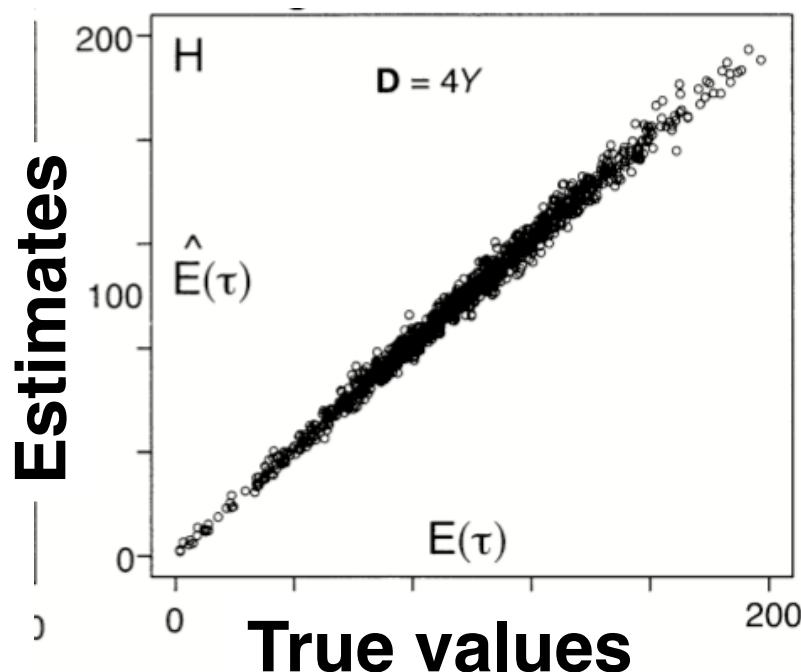
Proportion of accepted simulations

# **Development of ABC inferential strategy**

1. Summary statistic Selection
2. Prior Selection and Model Checking
- 3. Estimator Validation with PODS  
(pseudo-observed data sets)**
4. Estimation

### 3. Estimator Validation with PODS (pseudo-observed data sets)

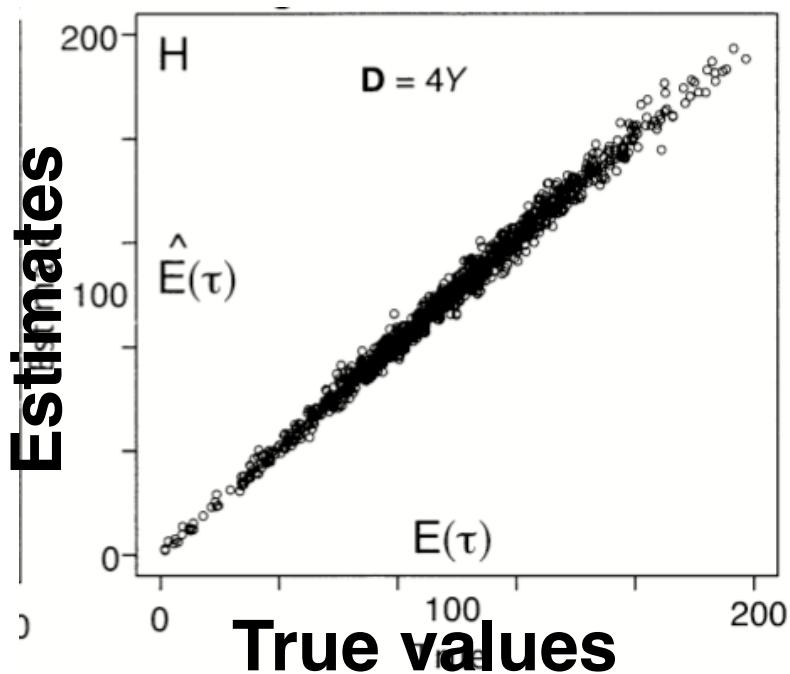
- 1 simulate “reference table” (the prior sample)
- 2 simulate PODS (pseudo-observed data sets with recorded parameter values)
- 3 make ABC estimates from the PODS



requires CPU

CUNY high  
performance computing  
facility

Simulation validation of ABC



> 300,000 Credit card  
numbers later

Beware  
[www.citibank.com](http://www.citibank.com)

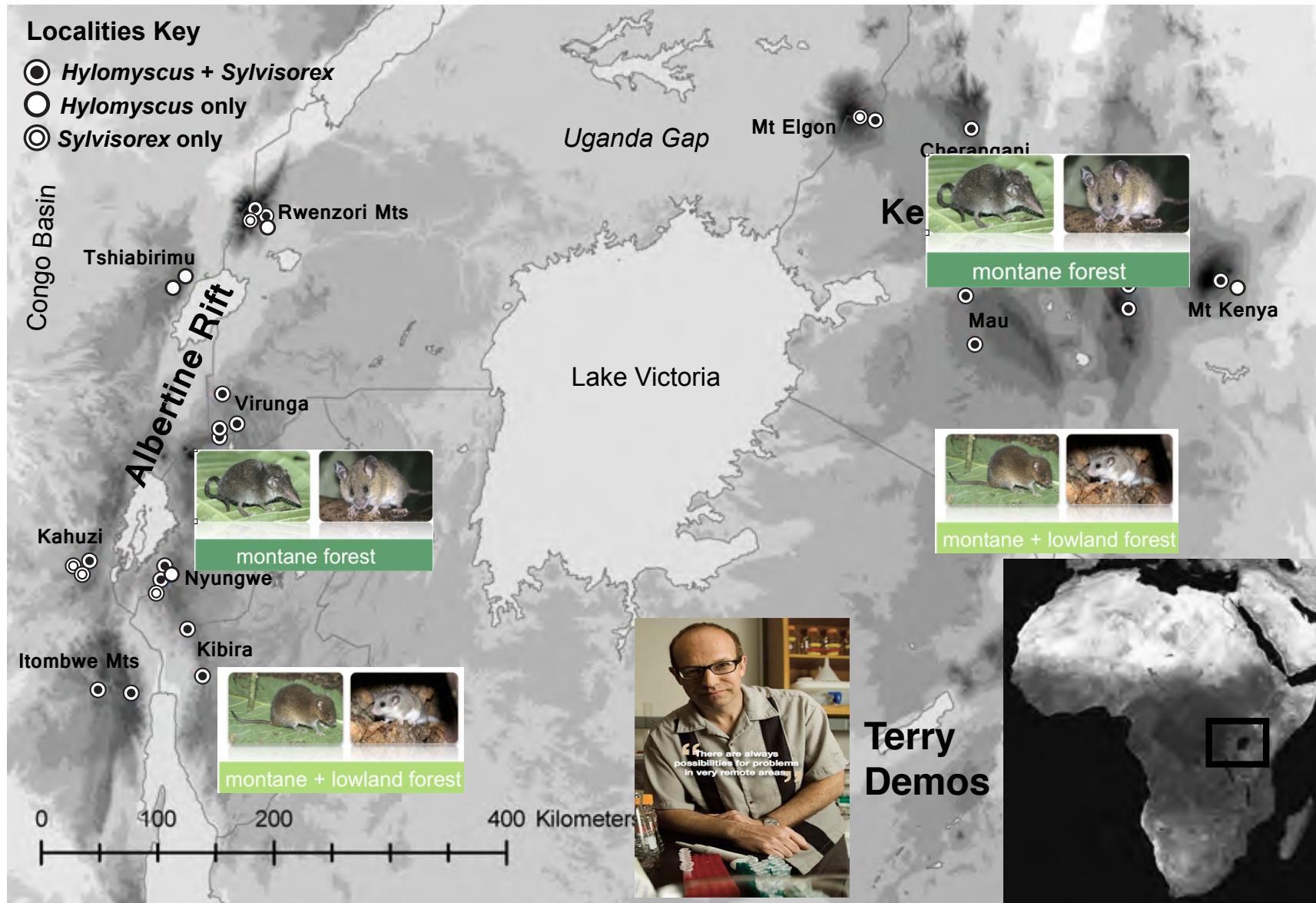


## **Software - MTML - msBayes**

### **6. PODS validation and/or power analysis**

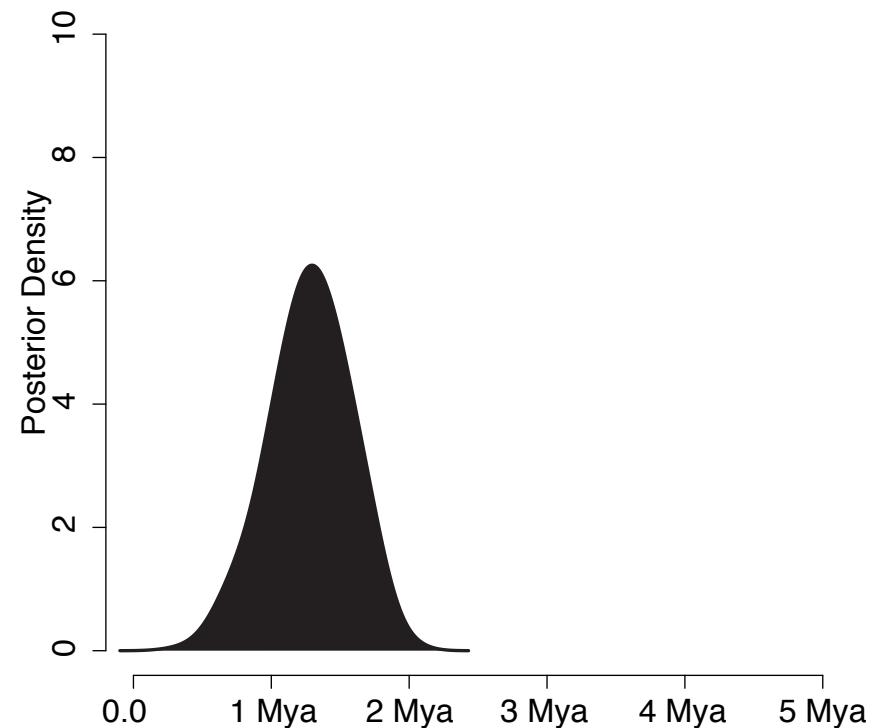
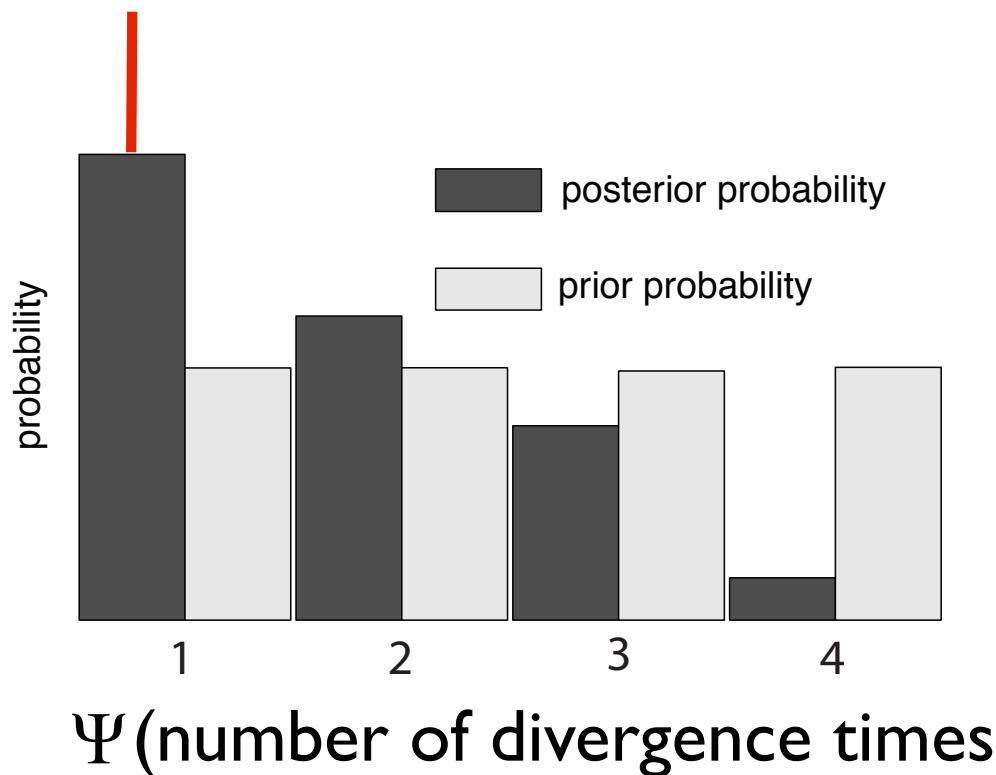
# Software - MTML - msBayes

## example - 4 East/West species pairs of African rodents



# Software - MTML - msBayes

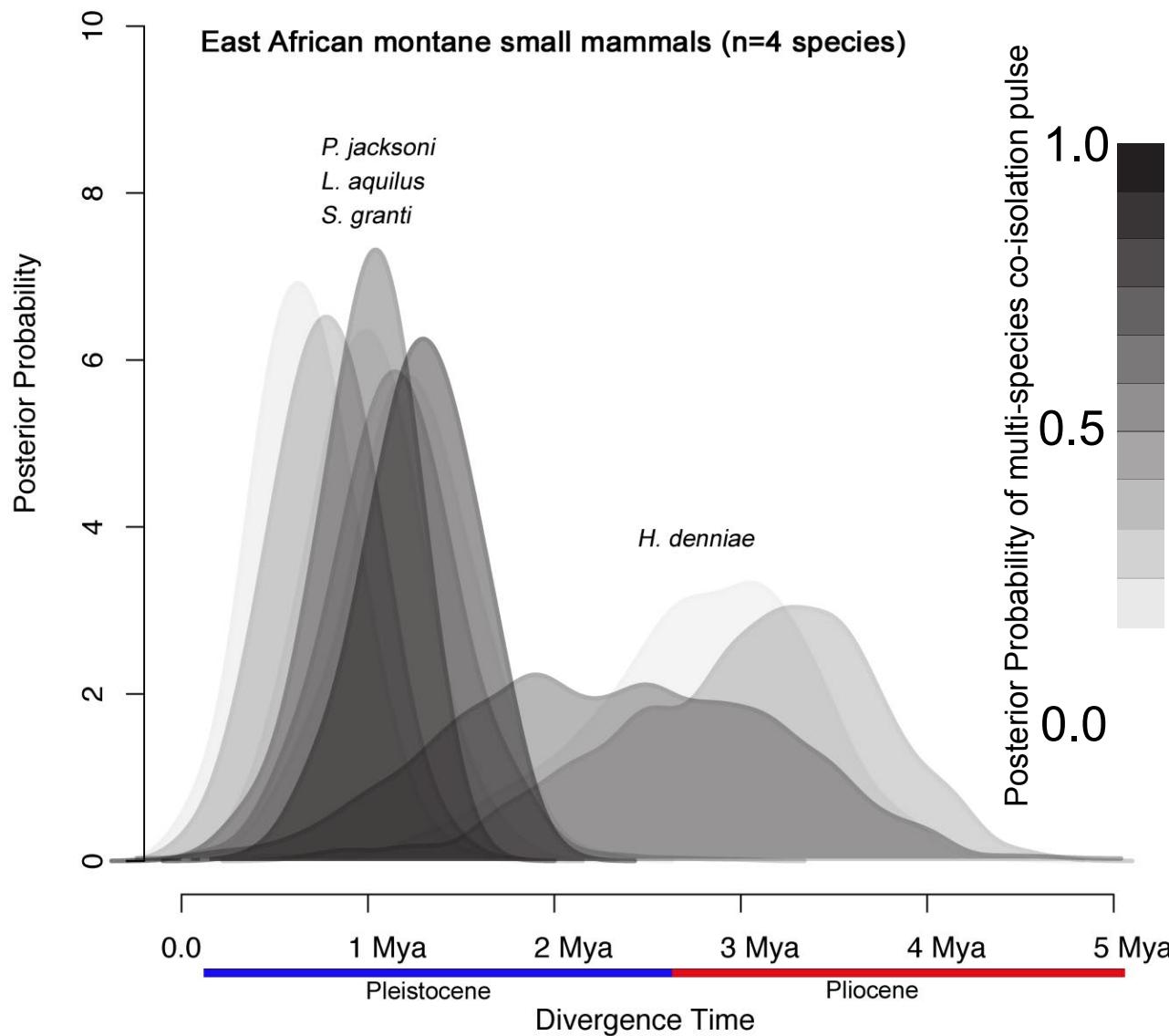
divergence times conditional on  $\Psi$  with  
**highest posterior probability** ( $\Psi=1$  in this case)



“simultaneous” divergence time  
for 4 pairs

# Software - MTML - msBayes

estimate divergence times averaged across posterior for number of times ( $\Psi$ )



# Thank you!

