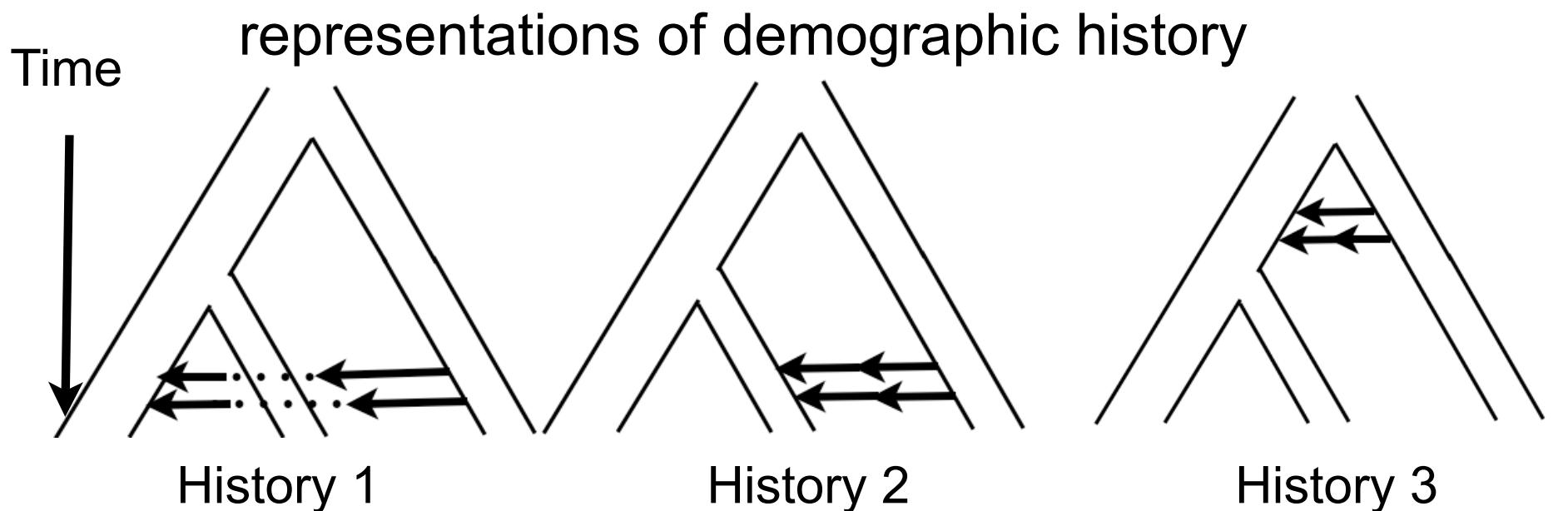


***Lecture 7:***  
Population History and demographic  
Inference 2

# outline

- Summary Statistics, Parameters and Models
- Coalescent Simulations
- Gene Trees vs Species Trees
- Felsenstein Equation
- MCMC, HMM and *maybe* more ABC methods

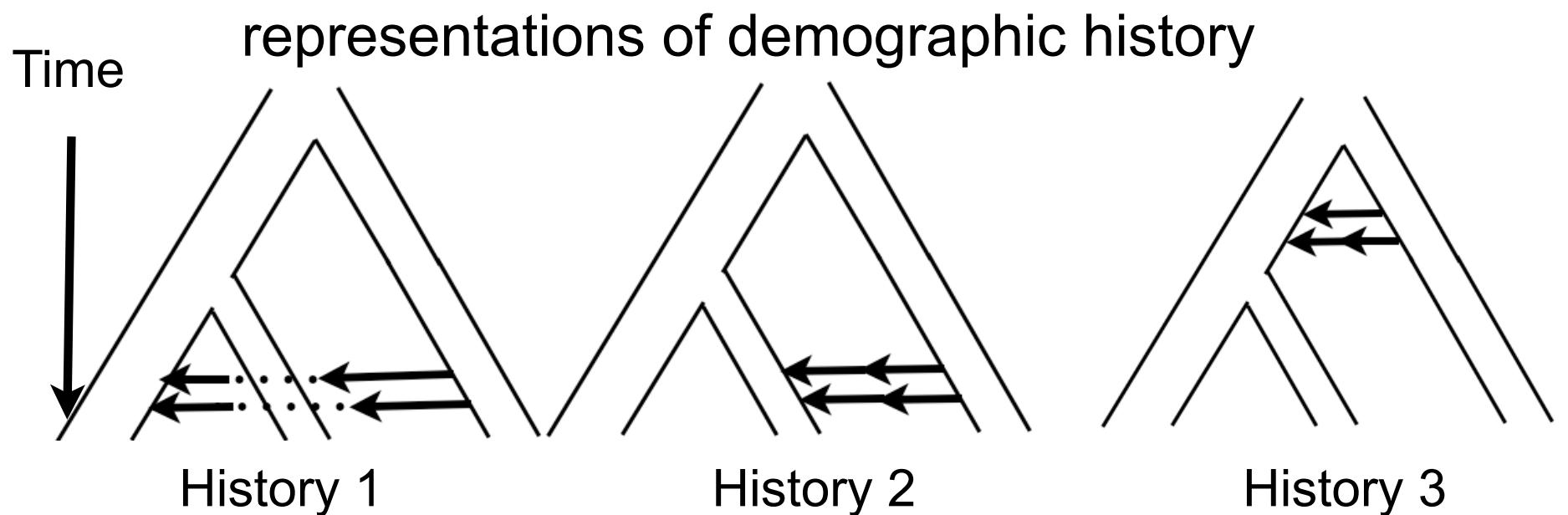
# Population History and demographic Inference



Coalescent Theory allows for statistical inference

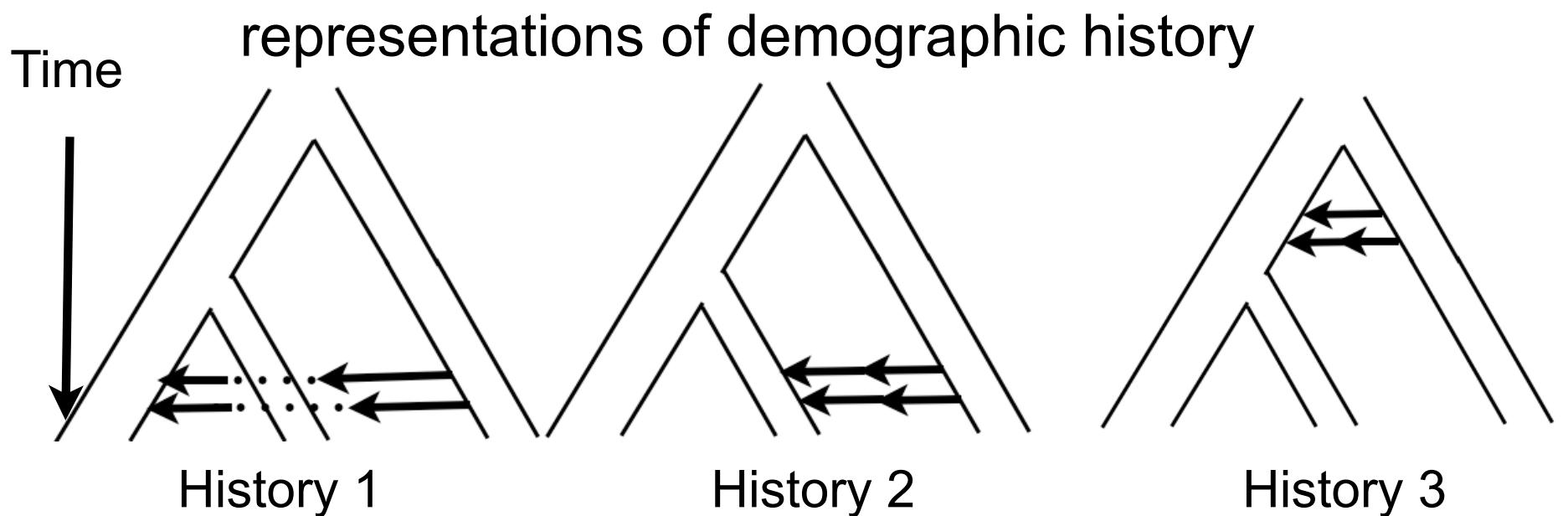
Testing historical demographic models

# Population History and demographic Inference



Population  
Genetic  
data

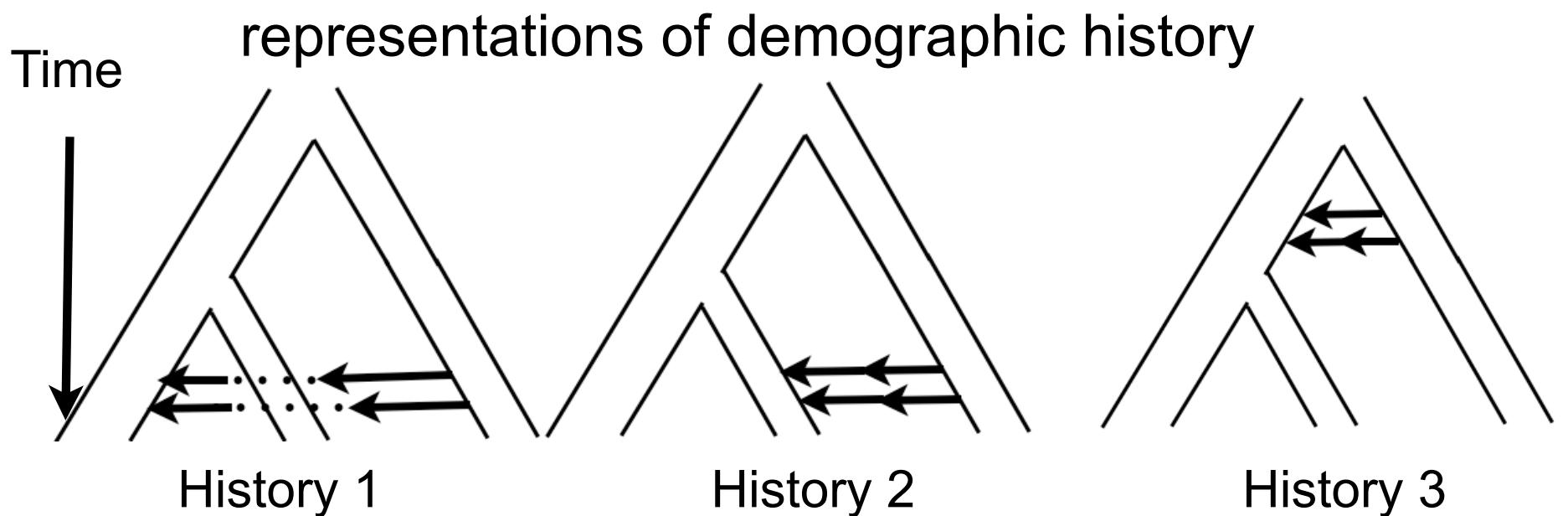
# Population History and demographic Inference



Population  
Genetic  
data

demographic models  
parameters given  
model (s)

# Population History and demographic Inference

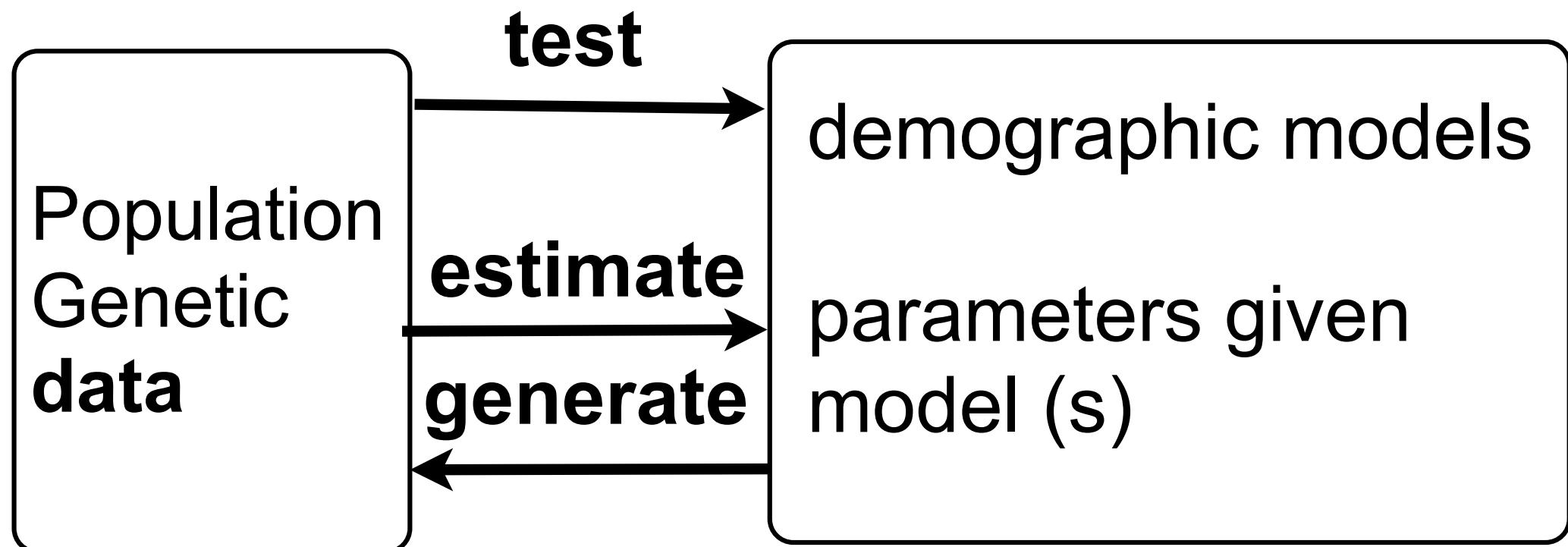
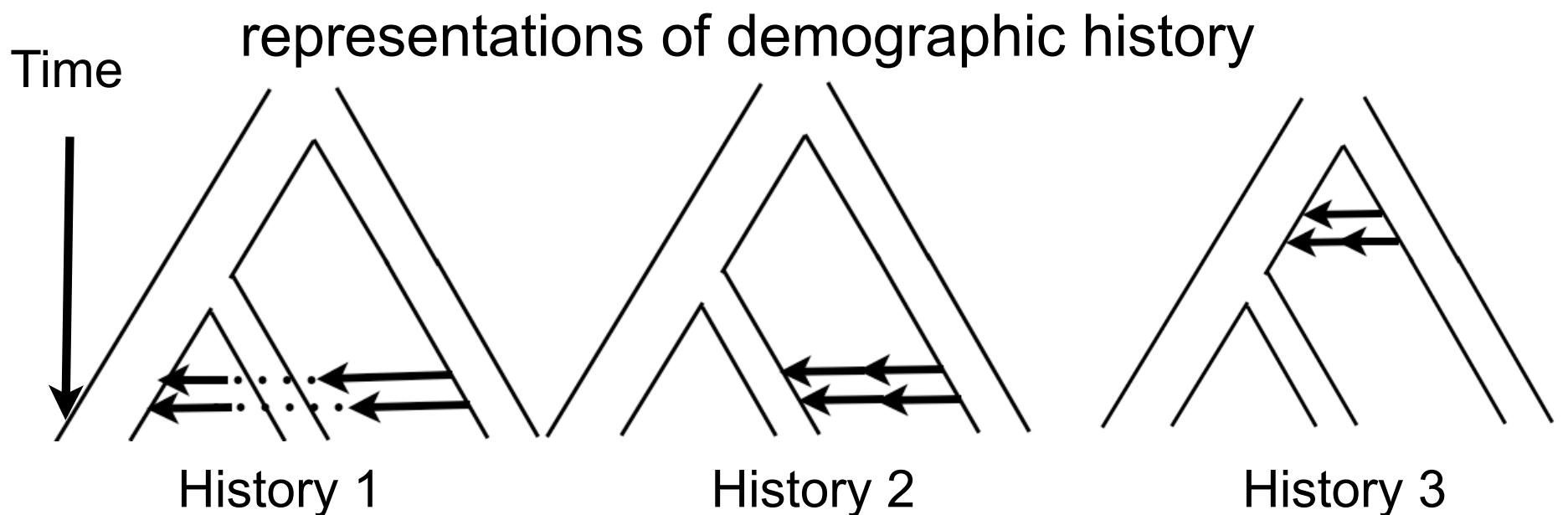


Population  
Genetic  
data

**generate**

demographic models  
parameters given  
model (s)

# Population History and demographic Inference



**summary stats**

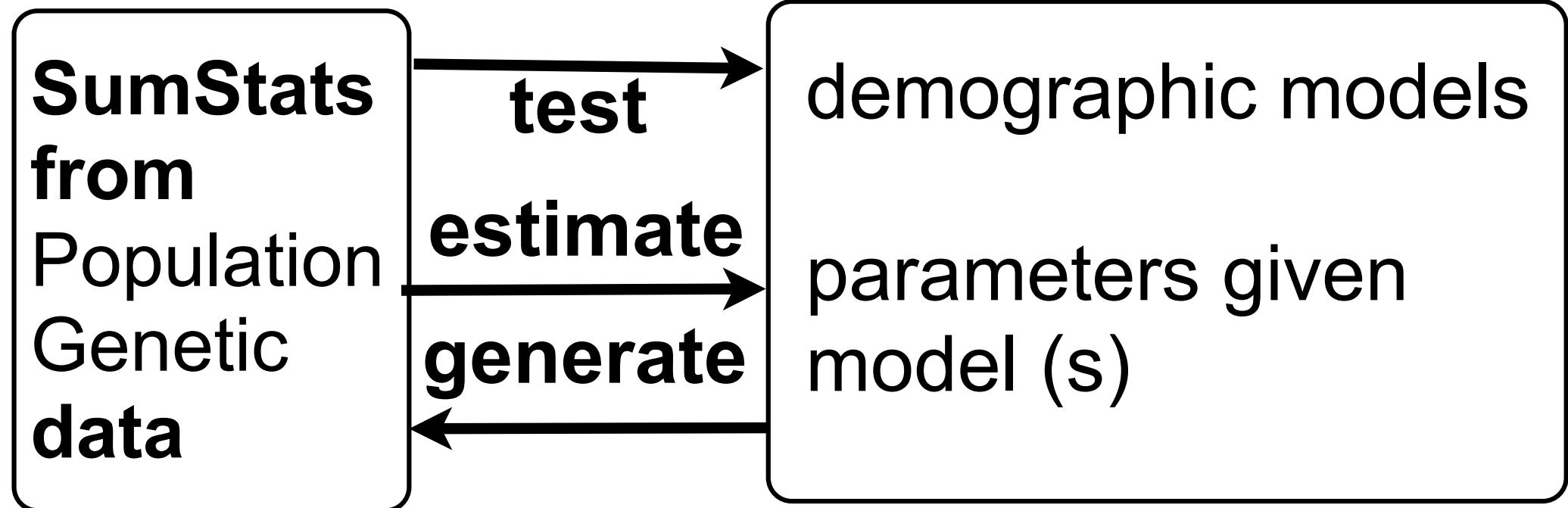
**model parameters**

$$(1 - F_{st})/4F_{st} \longleftrightarrow M = Nm$$

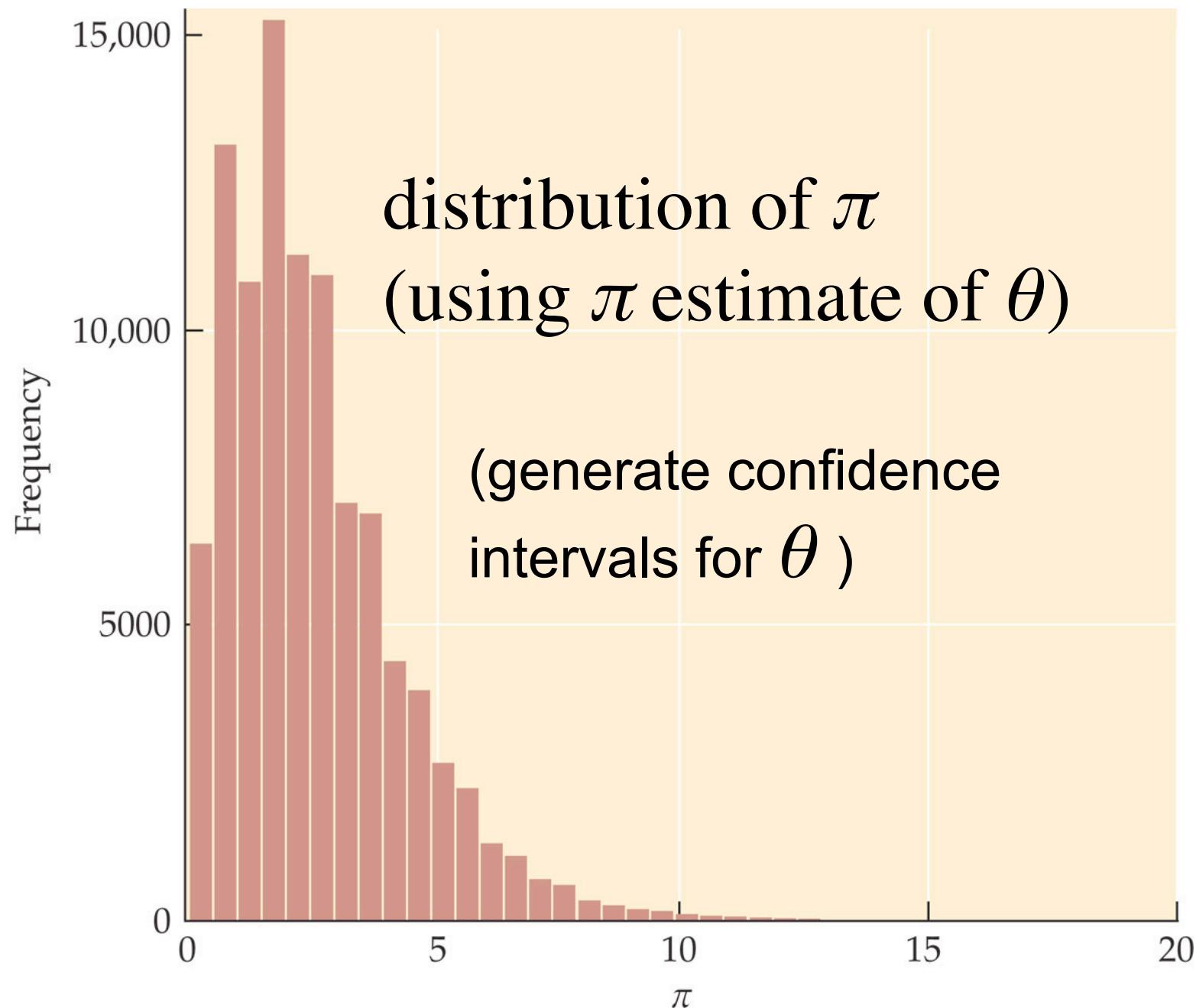
$\pi$  (average  
# pairwise  
differences)

$S$  (# of SNPs)

$$\theta = 4N\mu$$



Simulating  $\pi$  in 100,000 coalescence simulations under the standard coalescence model with infinite sites mutation and  $\theta = 2.6$



## The Tajima's Estimator

To estimate  $\theta$ , we can use the assumption of an infinite sites model and the expected number of mutations separating two individuals. An **estimate** is an educated guess of the true value of a parameter based on information obtained from data. In our case, the parameter is  $\theta$  and the data are the DNA sequences shown above. The data can be summarized in different ways. A popular way of summarizing DNA sequence data is in terms of the **average number of pairwise differences**, or  $\pi$ . The value of  $\pi$  is obtained by calculating the number of sites in which each pair of sequences differ, and then taking the average among all pairs of sequences. We can write this as

$$\pi = \frac{\sum_{i < j} d_{ij}}{n(n-1)/2} \quad (3.5)$$

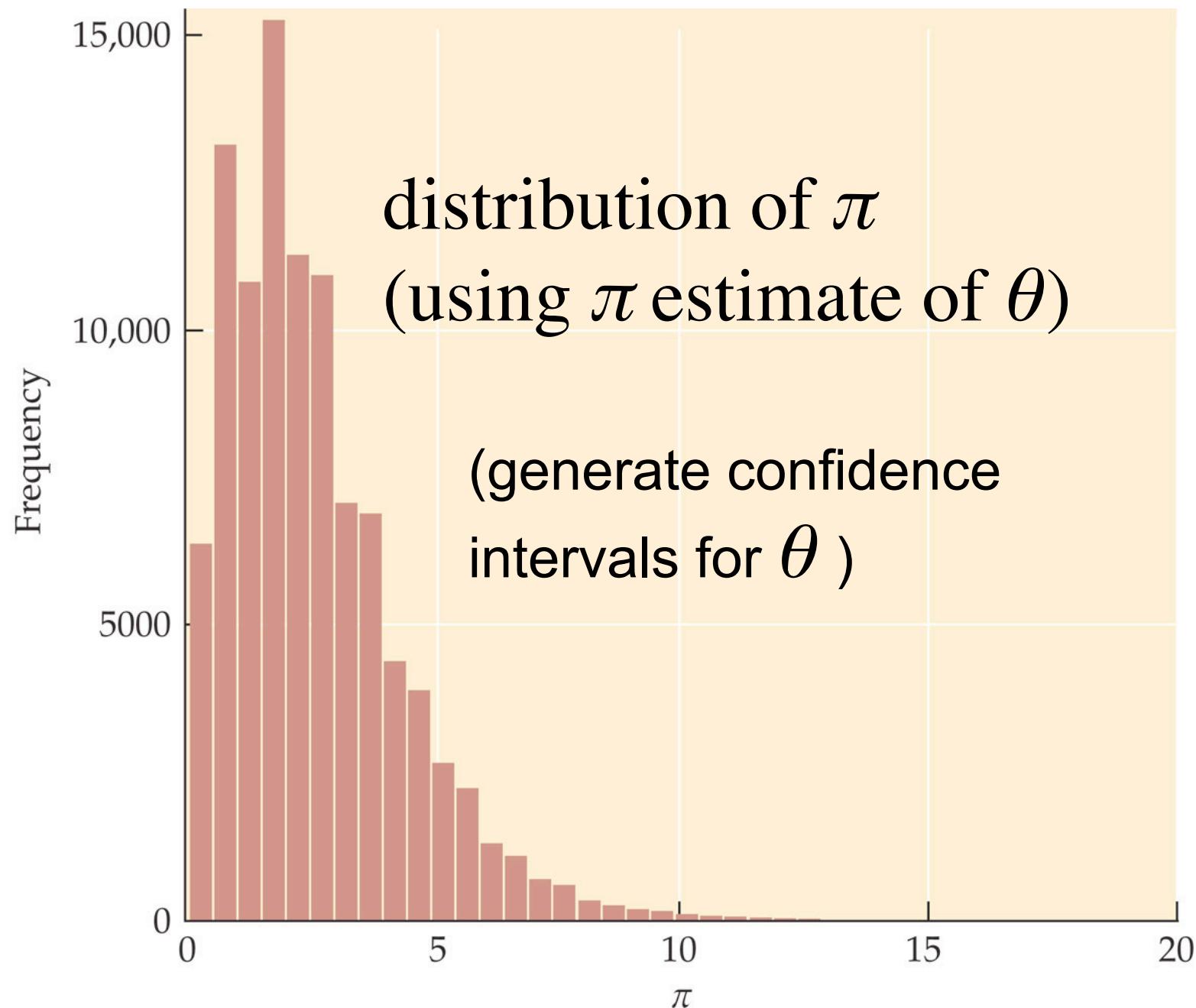
Nielsen and Slatkin 2013

We now have a convenient statistical method for estimating  $\theta$ . A quantity for estimating parameters from the data is called an **estimator**. The estimator of  $\theta$  based on  $\pi$  is sometimes called **Tajima's estimator** after the Japanese population geneticist F. Tajima. We denote an estimator of a quantity such as  $\theta$ , by putting a circumflex, or "hat" on  $\theta$ . We could, for example, write  $\hat{\theta}_T = \pi$ . The subscript  $T$  indicates that this is the particular estimator of  $\hat{\theta}_T$  named after F. Tajima. It is important to distinguish the parameter  $\theta$ , with a true value we may never know, from the estimator  $\hat{\theta}_T$ , which takes on a different value for each data set we consider.

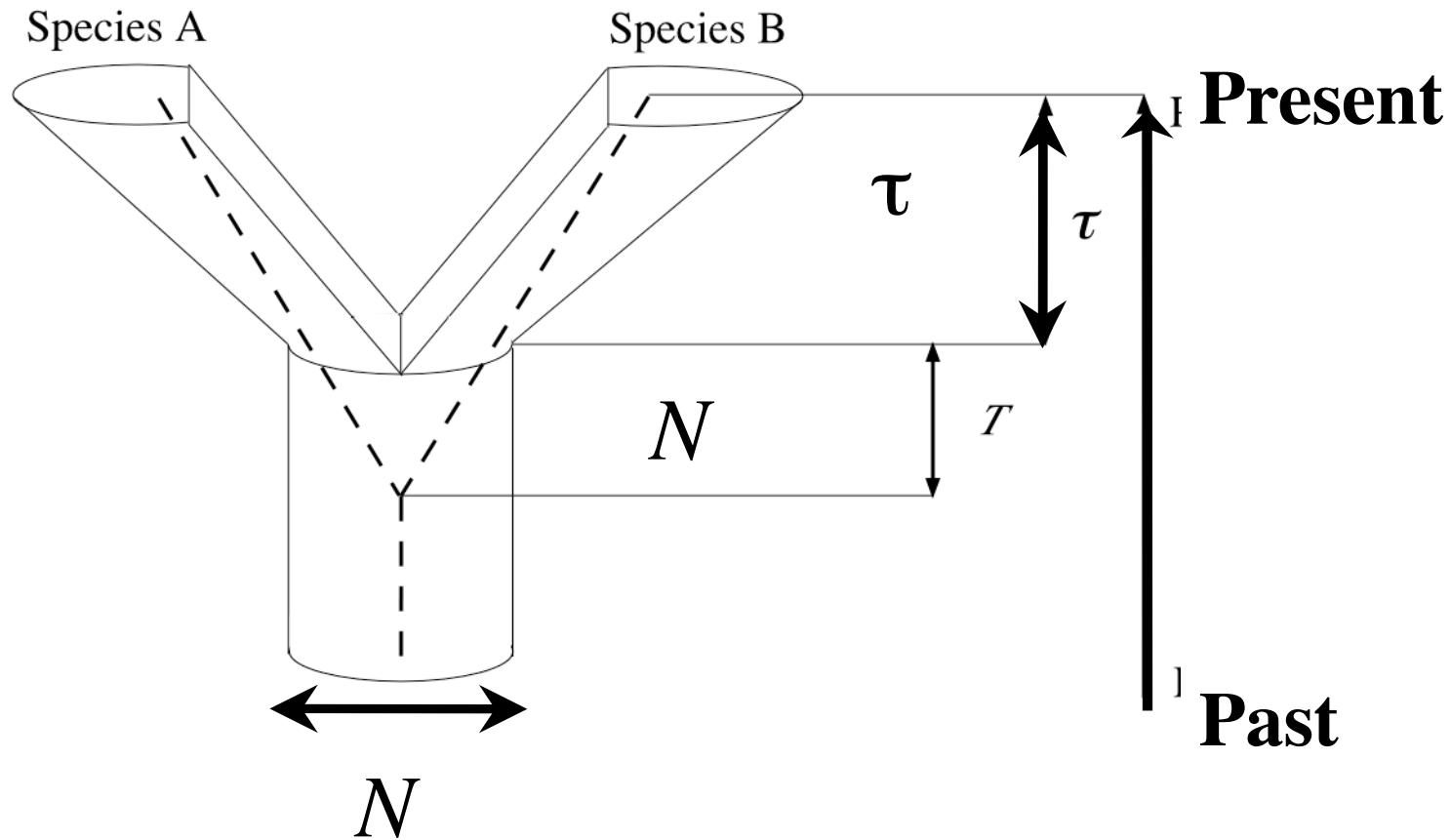


Nielsen and Slatkin 2013

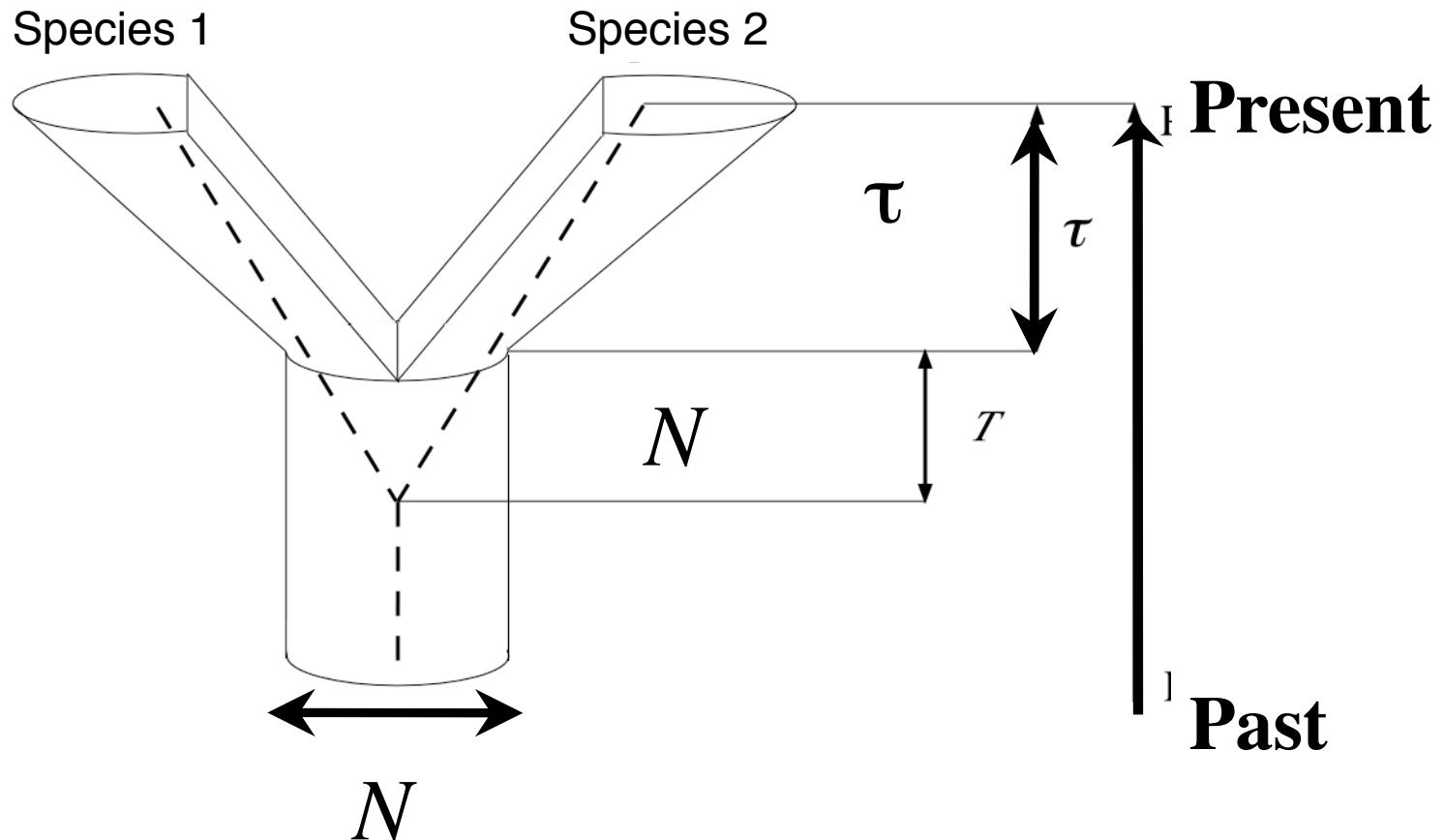
Simulating  $\pi$  in 100,000 coalescence simulations under the standard coalescence model with infinite sites mutation and  $\theta = 2.6$



$$\text{Genetic Divergence} = 2\tau\mu + 2N\mu$$



$$\pi_{net} = \pi_b - (\pi_1 + \pi_2)/2 = 2\tau\mu + 2N\mu$$



# Simulation Assignment: comparing single and double population models

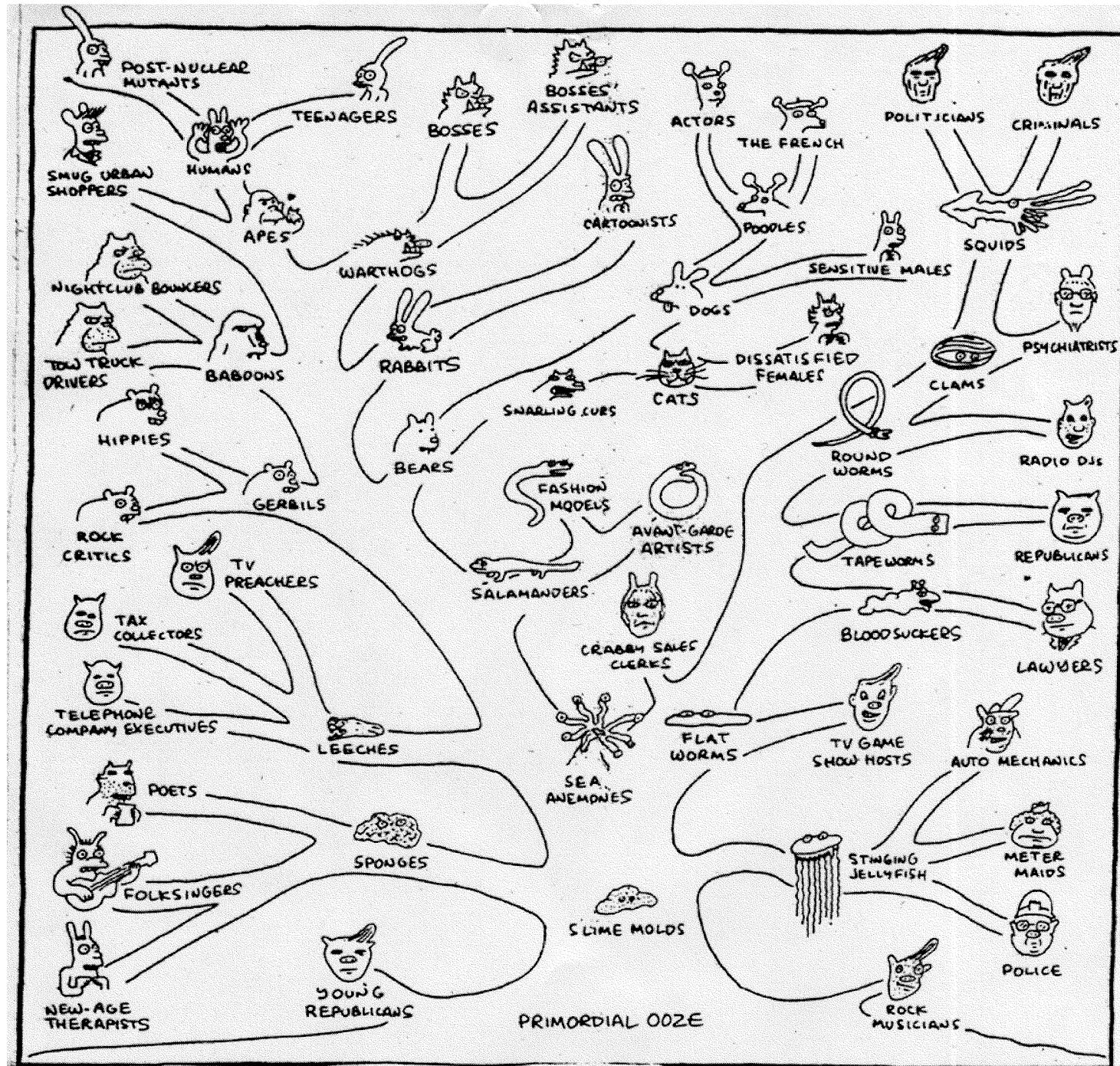
Use msPrime, SLiM or PipeMaster

1 locus, 1000 base pairs  
 $\mu=1e-7$  per base pair

Simulate a 2 population model  
 $\tau = 5,000$  generations  
 $N_1 = 10,000, N_2 = 10,000, N_A = 10,000$   
 $n_1 = 10, n_2 = 10$

1. Simulate once and calculate an observed  $\pi$  (ignoring which of the 2 populations the 20 individuals come from).
2. Simulate a single population model with that observed  $\pi$  value 1,000 times (pretend you know the mutation rate exactly) and plot the distribution of simulated  $\pi$  values
3. indicate the 'true'  $4N\mu$  on the plot. How off is it?
4. Simulate the 2 pop model 1,000 times and choose N values that you think will give you a  $\pi$  distribution that matches the one from the single pop model
5. plot results of 3 and 5

# Using all the data - estimating a gene tree



# Constructing Phylogenetic Trees using Molecular characters

Outgroup

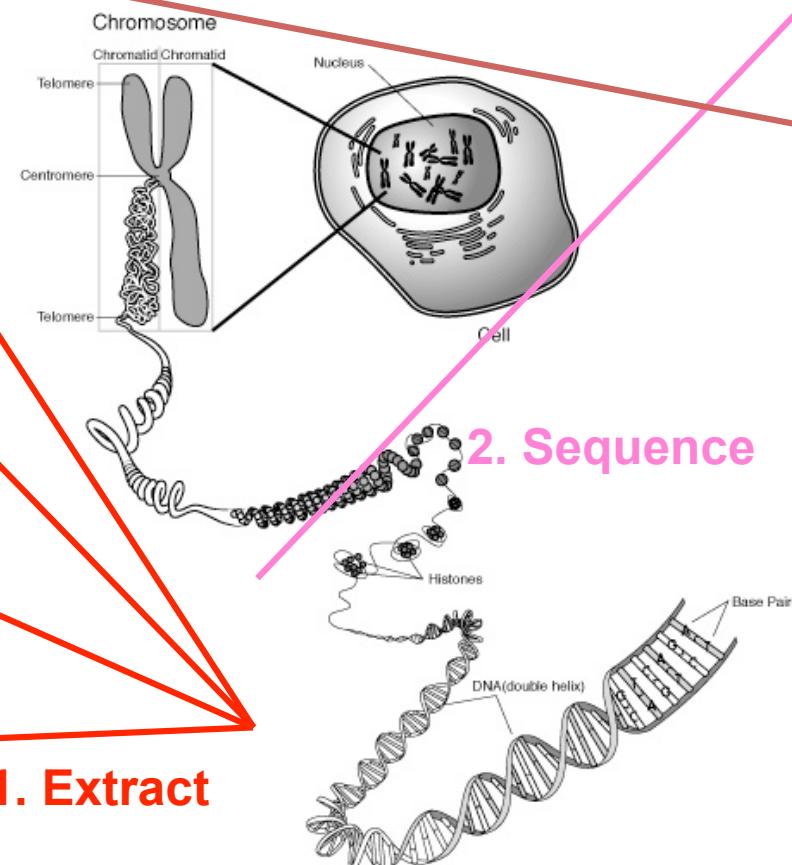
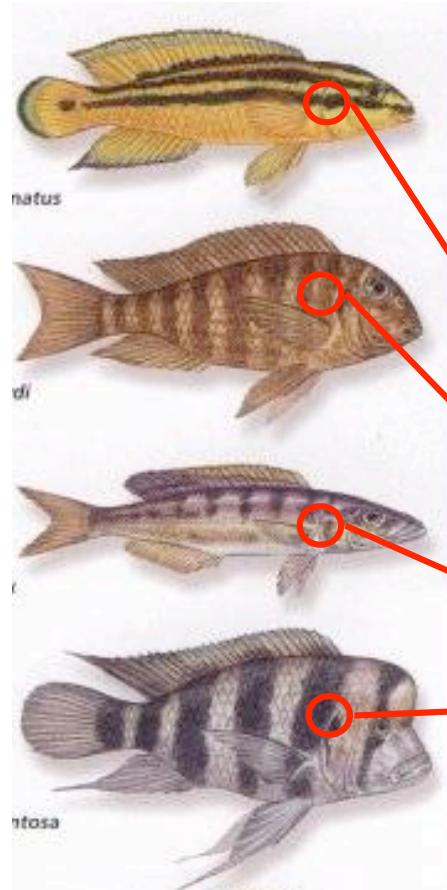
Species A

Species B

Species C

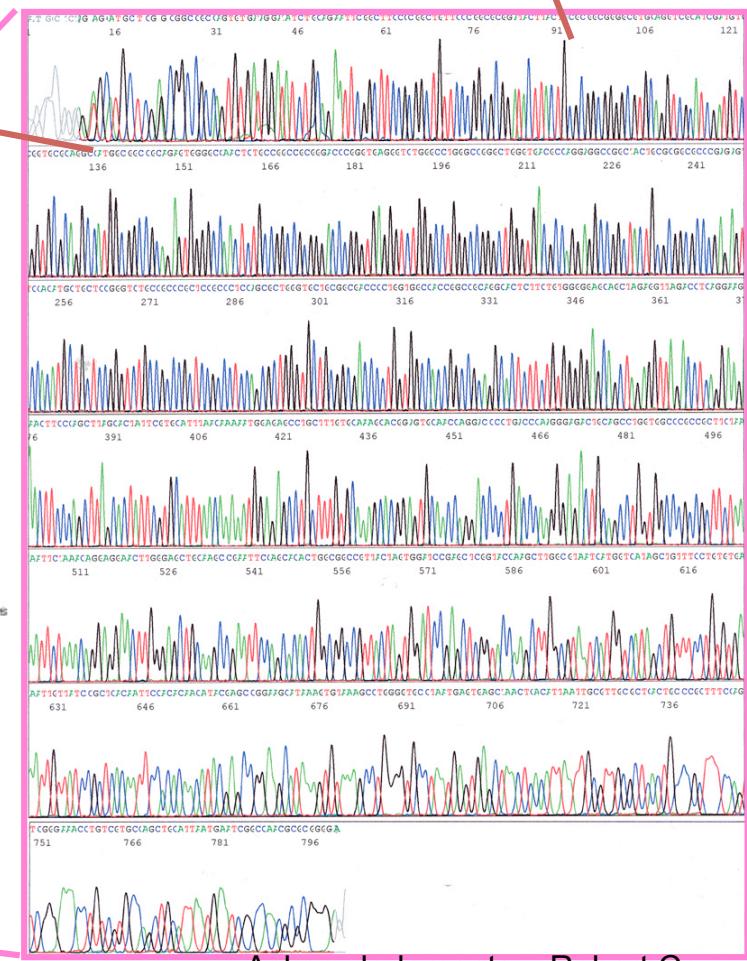
AAGCTTCATAAGGAGCAACCATTCTAATAATAAGCCTCATAAAGCC  
AAGCTTCACCGGCGCAGTTATCCTCATAATATGCCCTCATAAATGCC  
GTGCTTCACCGACGCAGTTGTCCTCATAATGTGCCCTCACTATGCC  
GTGCTTCACCGACGCAGTTGCCCTCATGATGAGCCTCACTATGCA

3. Align



1. Extract

<http://www.accessexcellence.org/AB/GG/chromosome.html>



Acknowledgments – Robert Cox

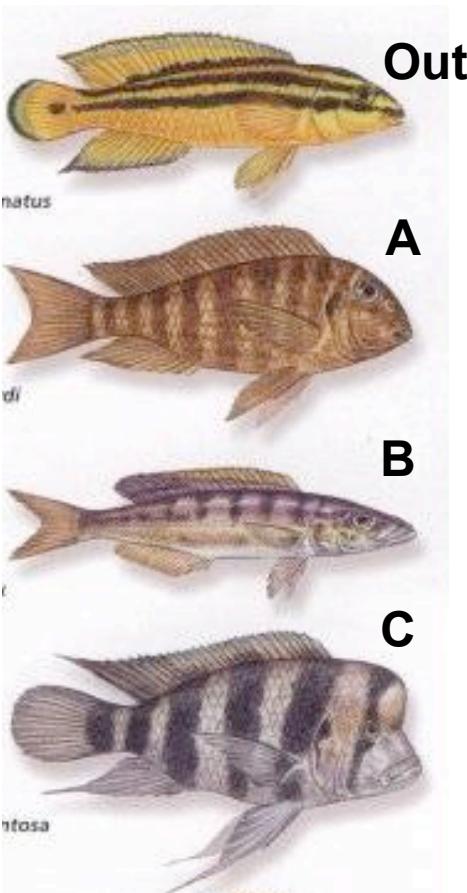
# Molecular characters

Outgroup

Species A

Species B

Species C



AAGCTTCATA

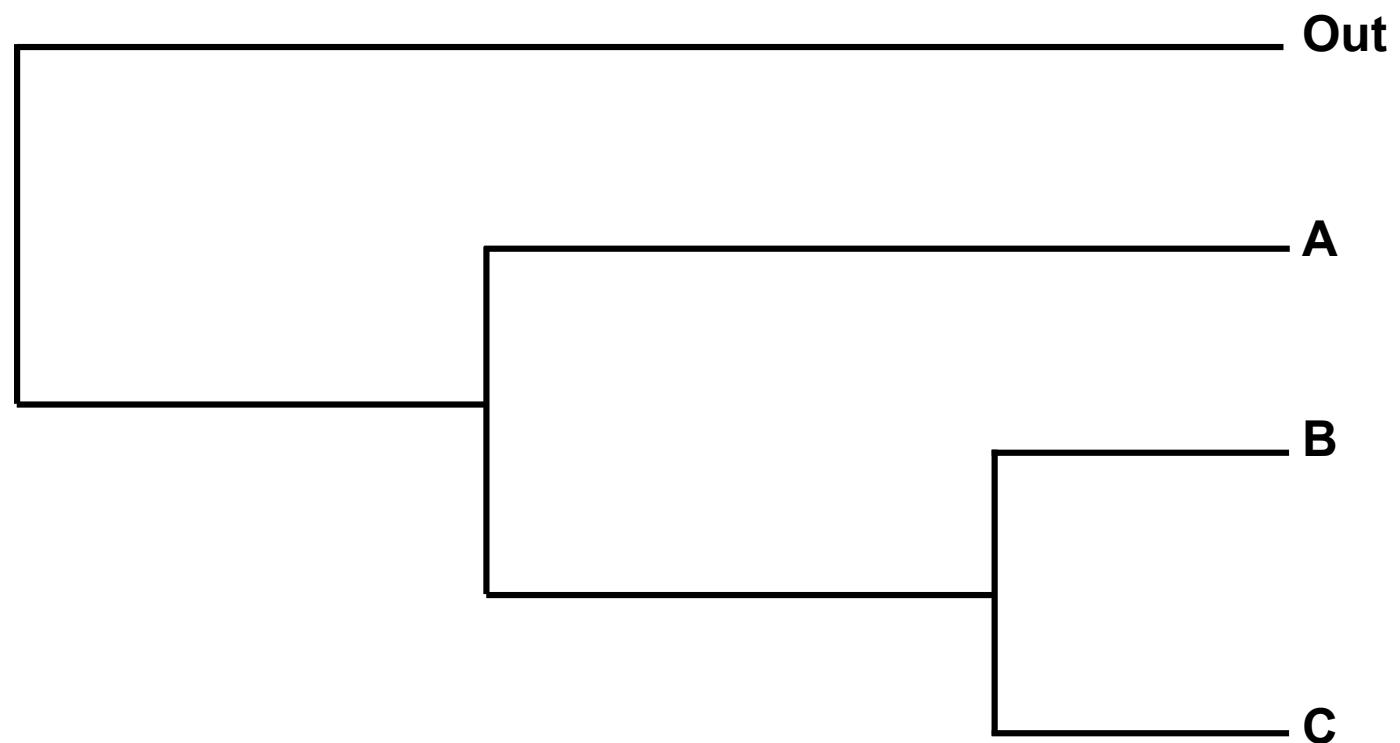
GAGCTTCACA

GTGCTTCACG

GTGCTTCACG

Invariable sites

These are not useful phylogenetic characters



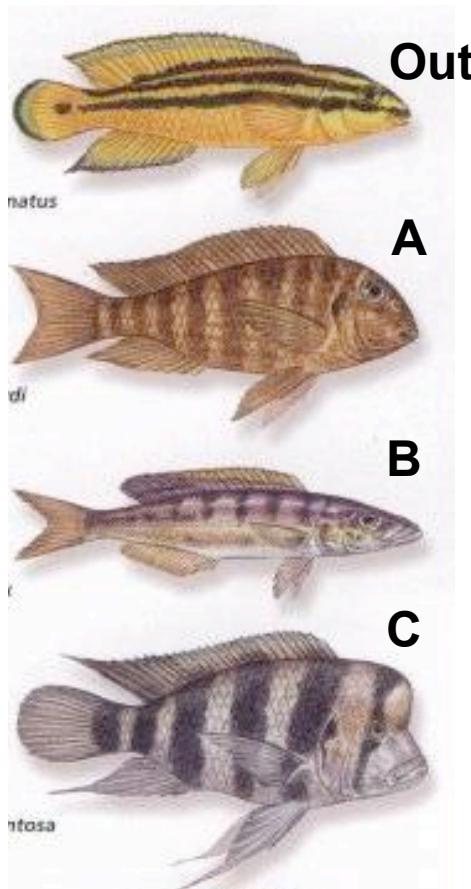
# Molecular characters

Outgroup

Species A

Species B

Species C



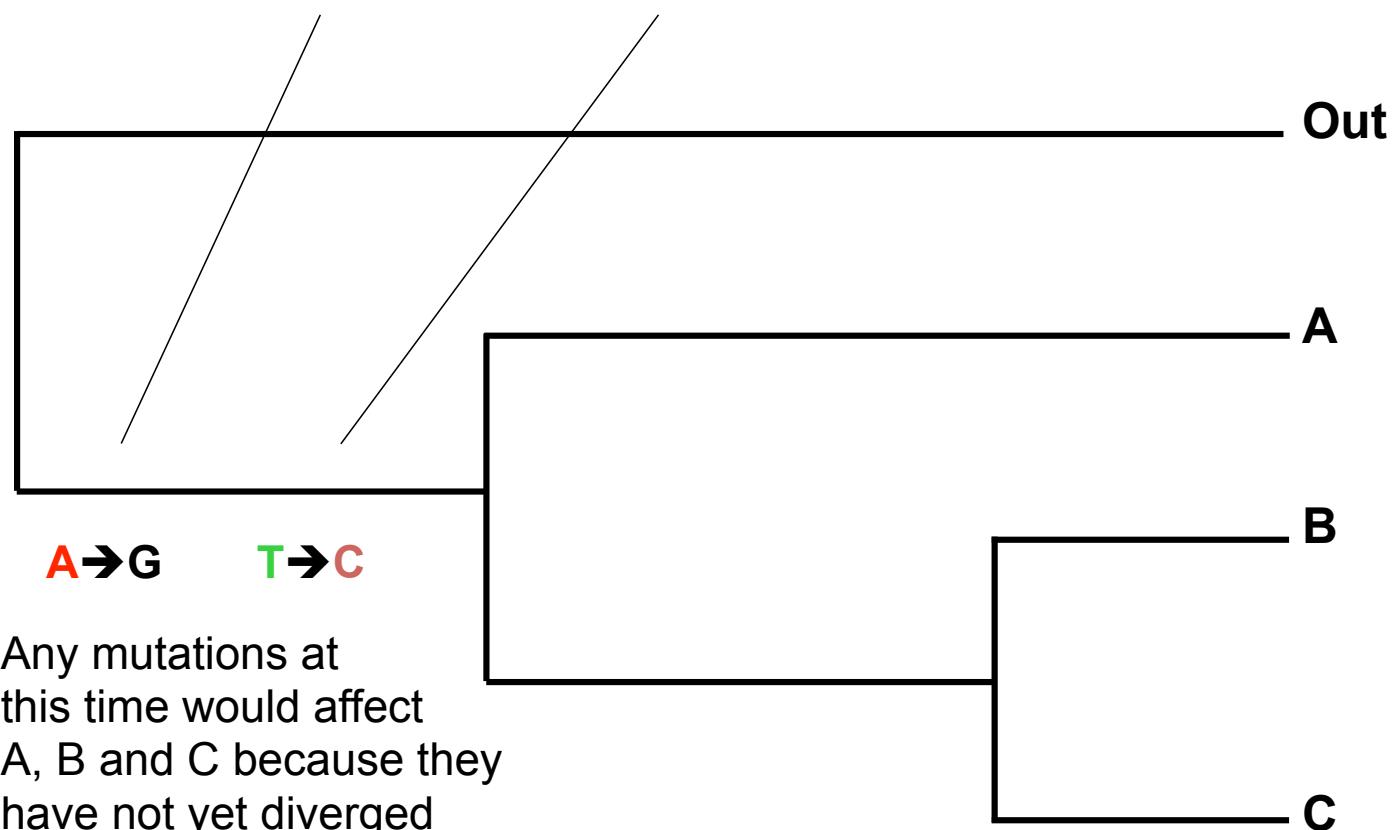
AAGCTTCATA

GAGCTTCACA

GTGCTTCACG

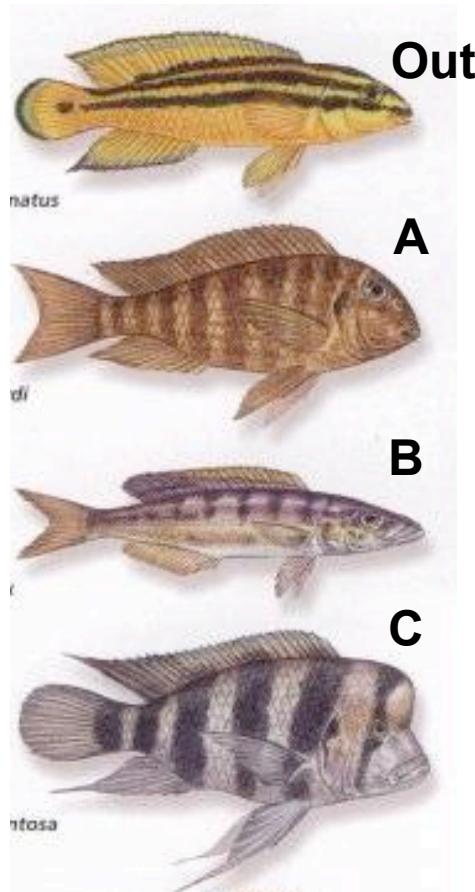
GTGCTTCACG

Synapomorphies  
supporting A+B+C



# Molecular characters

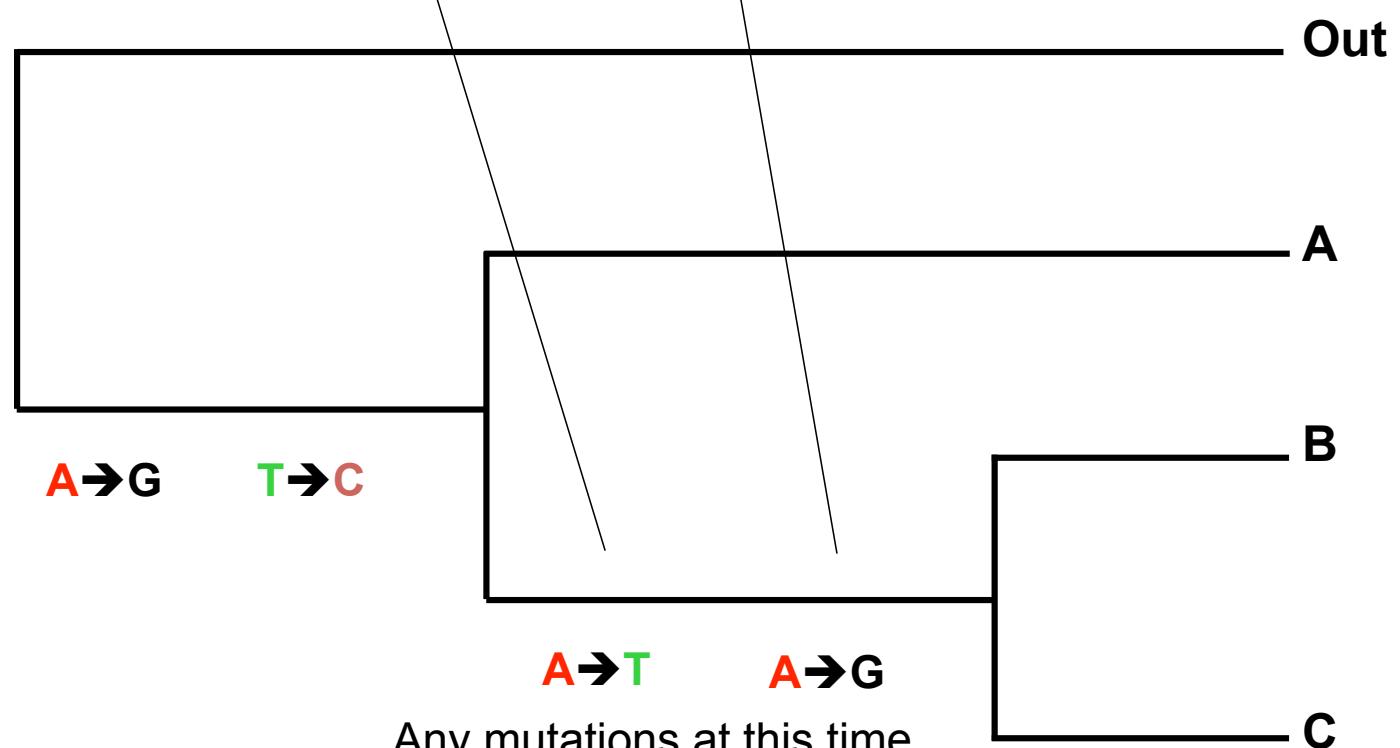
Outgroup  
Species A  
Species B  
Species C



**AAGCTTCATA**  
**GAGCTTCACA**  
**GTGCTTCACG**  
**GTGCCTCACG**

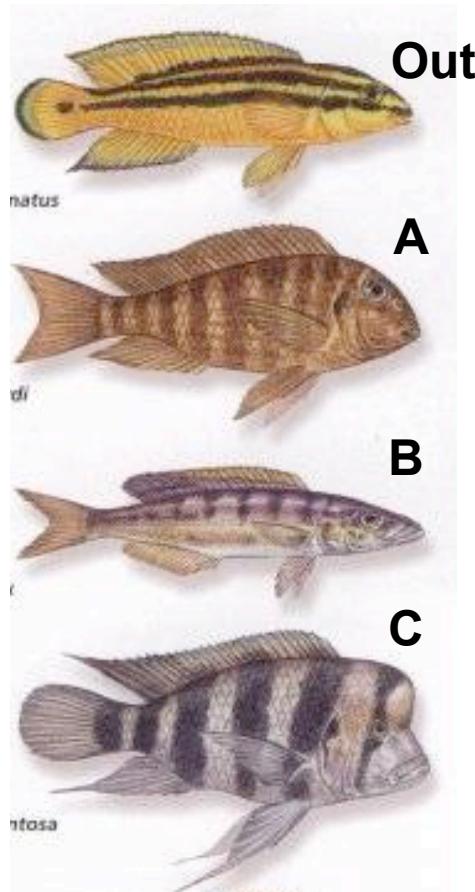
Synapomorphies supporting A+B+C

Synapomorphies supporting B+C



# Molecular characters

Outgroup  
Species A  
Species B  
Species C

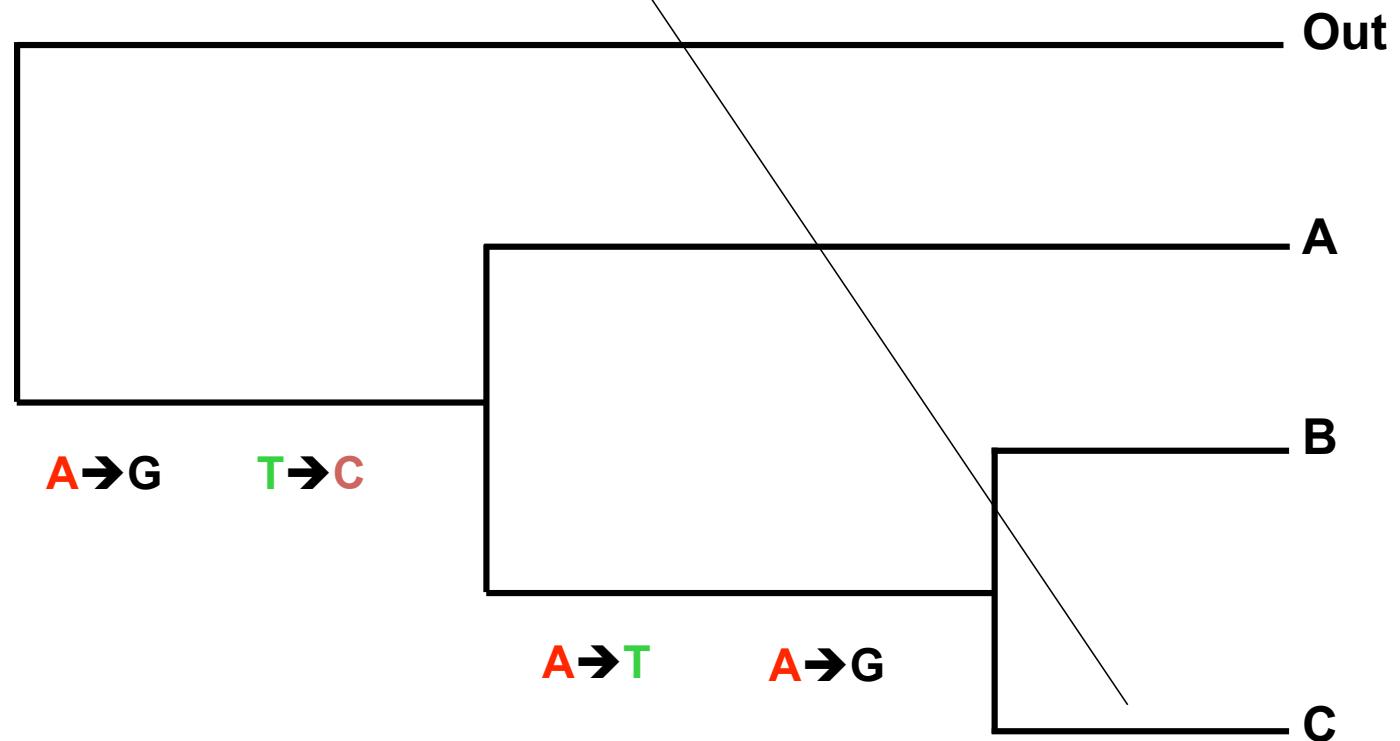


**AAGCTTCATA**  
**GAGCTTCACA**  
**GTGCTTCACG**  
**GTGCCTCACG**

Synapomorphies supporting A+B+C

Synapomorphies supporting B+C

Apomorphy for C



Any mutations at this time would only affect C     $T \rightarrow C$

# Methods in Phylogenetic Reconstruction

Distance

Maximum Parsimony

Maximum Likelihood/Bayesian

\* All algorithms are implemented using available software, eg. PAUP, PHYLIP, McClade, Mr. Bayes etc.

# Comparison of Methods (the old school circa 1997)

<b>Distance</b>	<b>Maximum parsimony</b>	<b>Maximum likelihood</b>
Uses only pairwise distances	Uses only shared derived characters, infinite sites model	Uses all data
Minimizes distance between nearest neighbors	Minimizes total number of changes	Maximizes tree likelihood given specific parameter values
Very fast	Slow	<b>Very slow</b>
Easily trapped in local optima	only uses infinite sites model	can use any evolution model (and use likelihood criteria to choose)
Good for generating tentative tree	Best option when homoplasy rare	optimal method if model is correct (but can be difficult in practice for large data sets)

$$L = \Pr(\text{data} \mid \Theta)$$

$\Theta$  = model, tree & parameters

parameters = e.g. rates, times etc

model = infinite sites

tree = among all possible trees

Likelihood ( $L$ ) is the probability of the data (alignment), given a tree (with topology *and* branch lengths specified) and a probabilistic model of evolution.

$$L = \Pr(\text{data} \mid \Theta)$$

$\Theta$  = model, tree & parameters

parameters = e.g. rates

model = infinite sites

tree = among all possible trees

Likelihood ( $L$ ) is the probability of the data (alignment), given a tree (with topology *and* branch lengths specified) and a probabilistic model of evolution.

Assumptions (the fine print):

- The tree is correct
- The probability that a position has a certain state at time 1 depends only on the state at time 0; knowing that it had some state prior to time 0 is irrelevant
  - (this is called a *Markov process*)
- Data (individual sites) are independent
- A uniform evolutionary process operated across the entire tree (why might this be false? endosymbiosis? loss of function?),
  - i.e., the process of evolution is a *homogeneous Markov process*.

Probability and prediction

Know parameters -> Predict Outcome (predict data)

Likelihood

Observe Data -> Estimate parameters

# Likelihood and coin tosses

$n = 100$  (# of coin tosses)

$h = 56$  (# of heads)

If  $p = 0.5$

$$\Pr(\text{data} | p = 0.50) = \frac{100!}{56!44!} 0.5^{56} 0.5^{44} = 0.0389$$

This starts the count of  
number of ways event can  
occur.

This is the probability  
of success for  $x$  trials.

This ends the count of  
number of ways event can  
occur.

This deletes  
duplications.

This is the probability  
of failure for the  $x$  trials.



$P(\text{heads})$

$$= \frac{1}{2} = 0.5$$

# Likelihood and coin tosses

$n = 100$  (# of coin tosses)

$h = 56$  (# of heads)

If  $p = 0.5$

$$\Pr(\text{data} | p = 0.50) = \frac{100!}{56!44!} 0.5^{56} 0.5^{44} = 0.0389$$

But if  $p = 0.52$

$\Pr(\text{data} | p = 0.52) = \text{ie , bad old coin}$



This starts the count of  
number of ways event can  
occur.

This ends the count of  
number of ways event can  
occur.

This deletes  
duplications.

This is the probability  
of success for x trials.

This is the probability  
of failure for the x trials.

$n = 100$  (# of coin tosses)

$h = 56$  (# of heads)

If  $p = 0.5$

$$\Pr(\text{data}|p = 0.50) = \frac{100!}{56!44!} 0.5^{56} 0.5^{44} = 0.0389$$

But if  $p = 0.52$

$$\Pr(\text{data}|p = 0.52) = \frac{100!}{56!44!} 0.52^{56} 0.52^{44} = 0.0581$$

$$\Pr(\text{data} | \Theta) = \frac{n!}{(n-x)!x!} p^x q^{n-x}$$

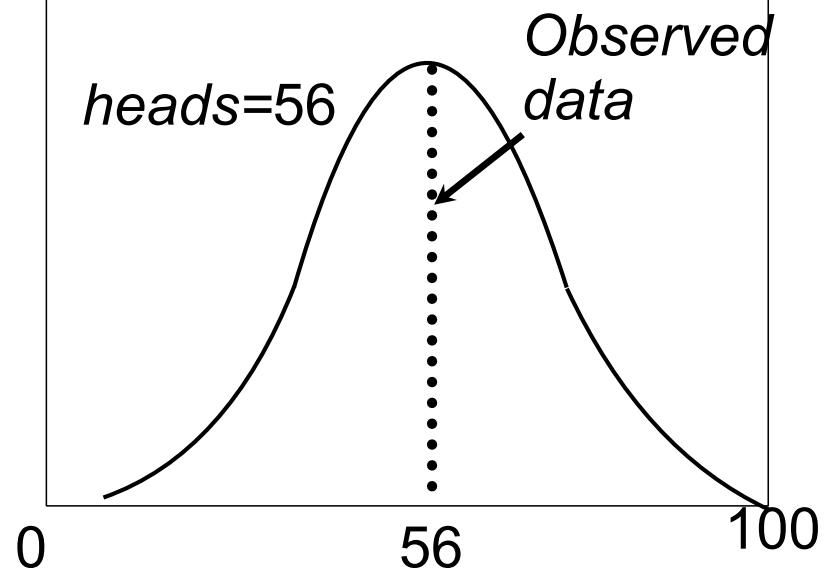
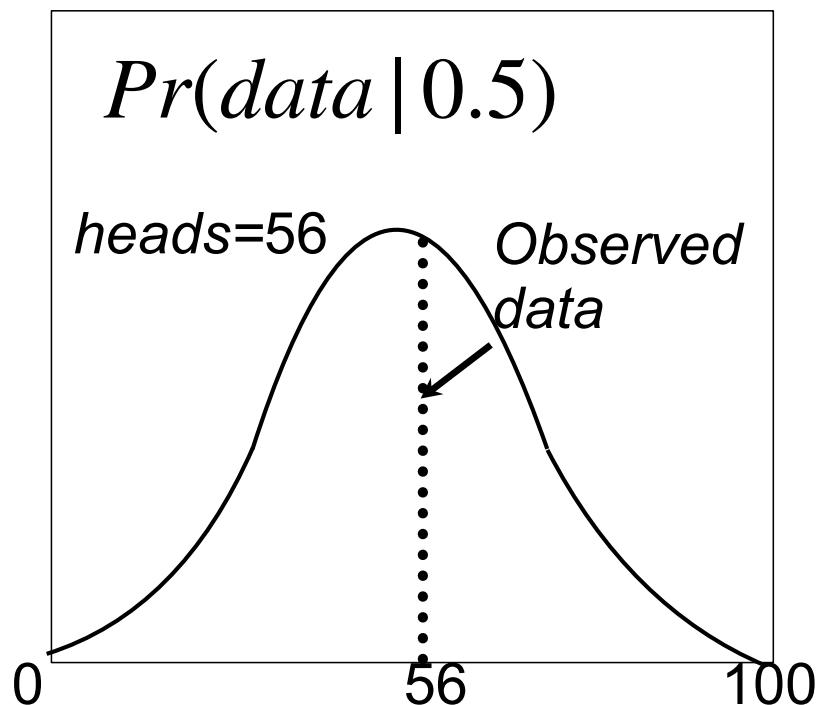
This starts the count of number of ways event can occur.

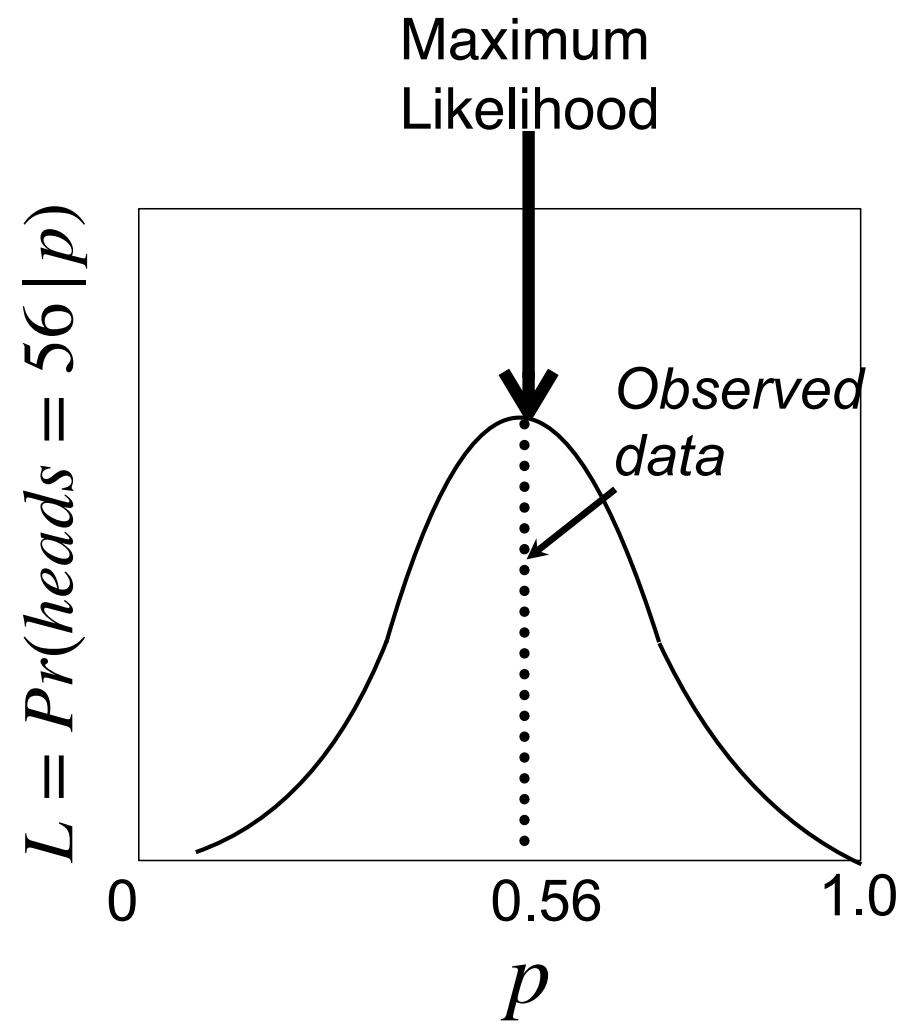
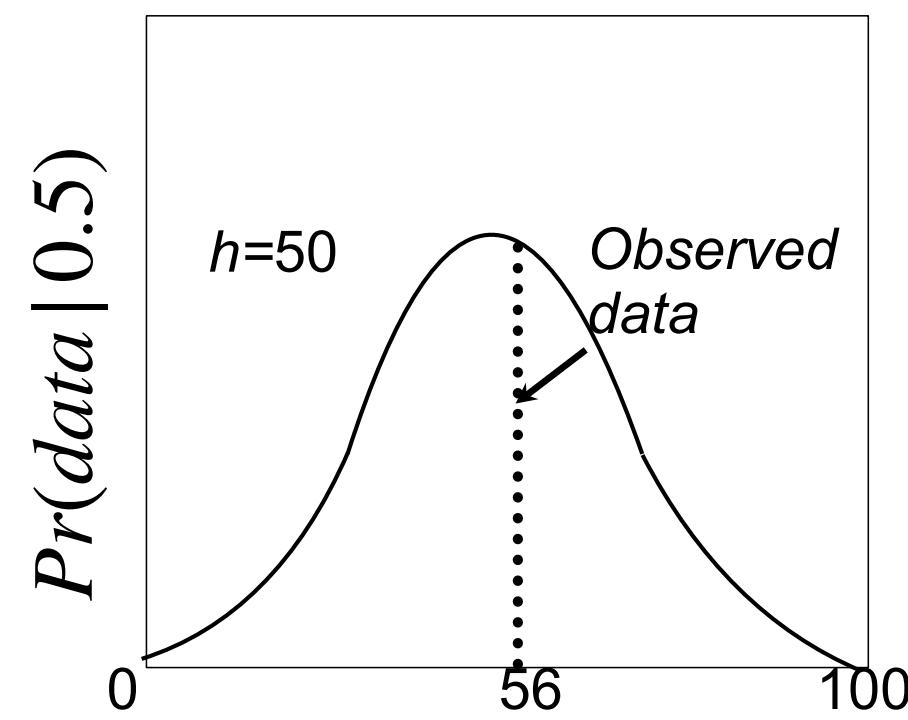
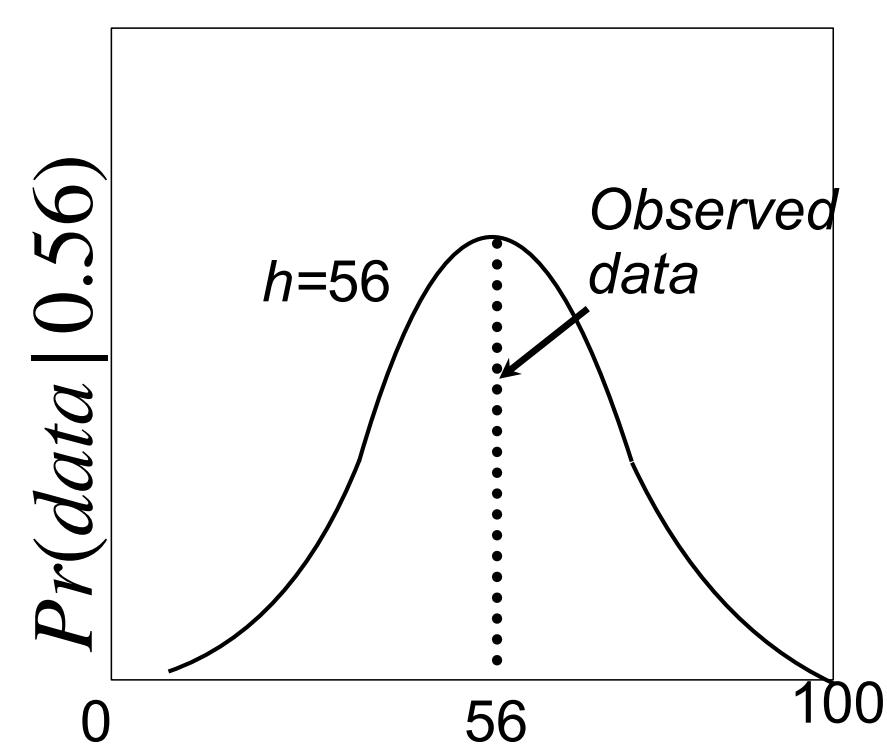
This ends the count of number of ways event can occur.

This deletes duplications.

This is the probability of success for  $x$  trials.

This is the probability of failure for the  $x$  trials.

$$Pr(data | 0.56)$$

$$Pr(data | 0.5)$$


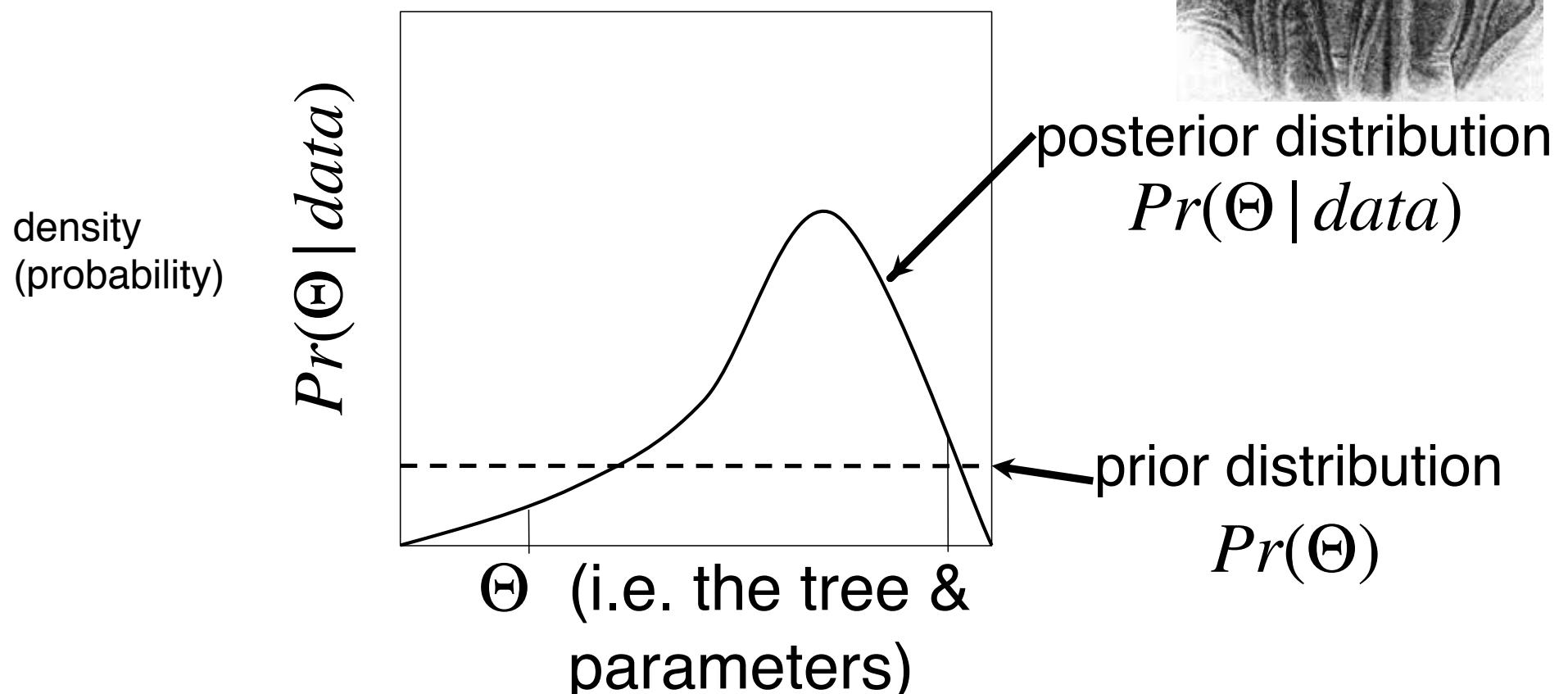


# Bayesian

prior distribution for parameters,  $Pr(\Theta)$

we can obtain posterior distribution  $Pr(\Theta | data)$

$$Pr(\Theta | data) = \frac{Pr(data | \Theta)Pr(\Theta)}{Pr(data)}$$



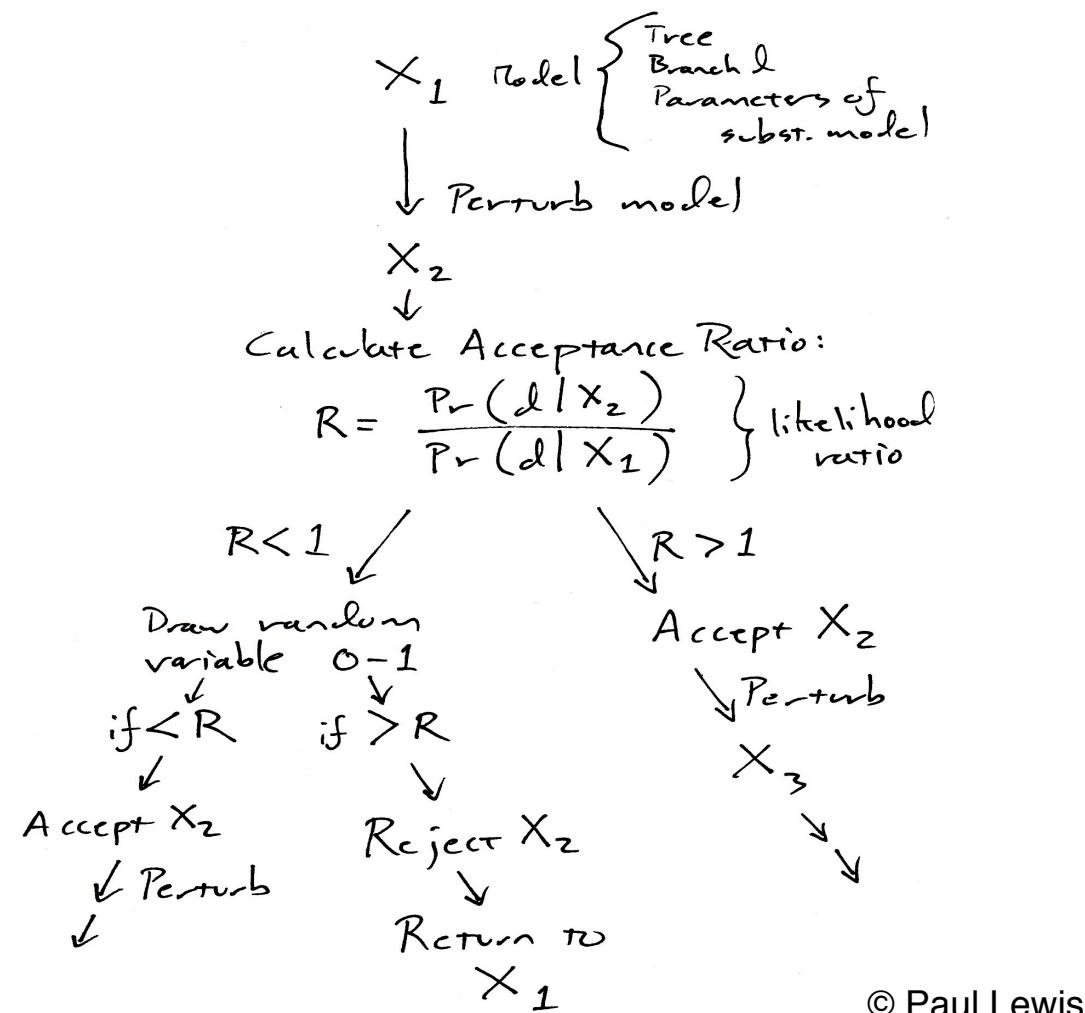
Calculating  $Pr(data | \Theta)$  can be difficult

Can only do it in closed form for the simplest models

Complex models can in principle be solved via integration, but this is often intractable

***must resort to numerical or simulation-based methods  
(MCMC, ABC, supervised machine learning etc)***

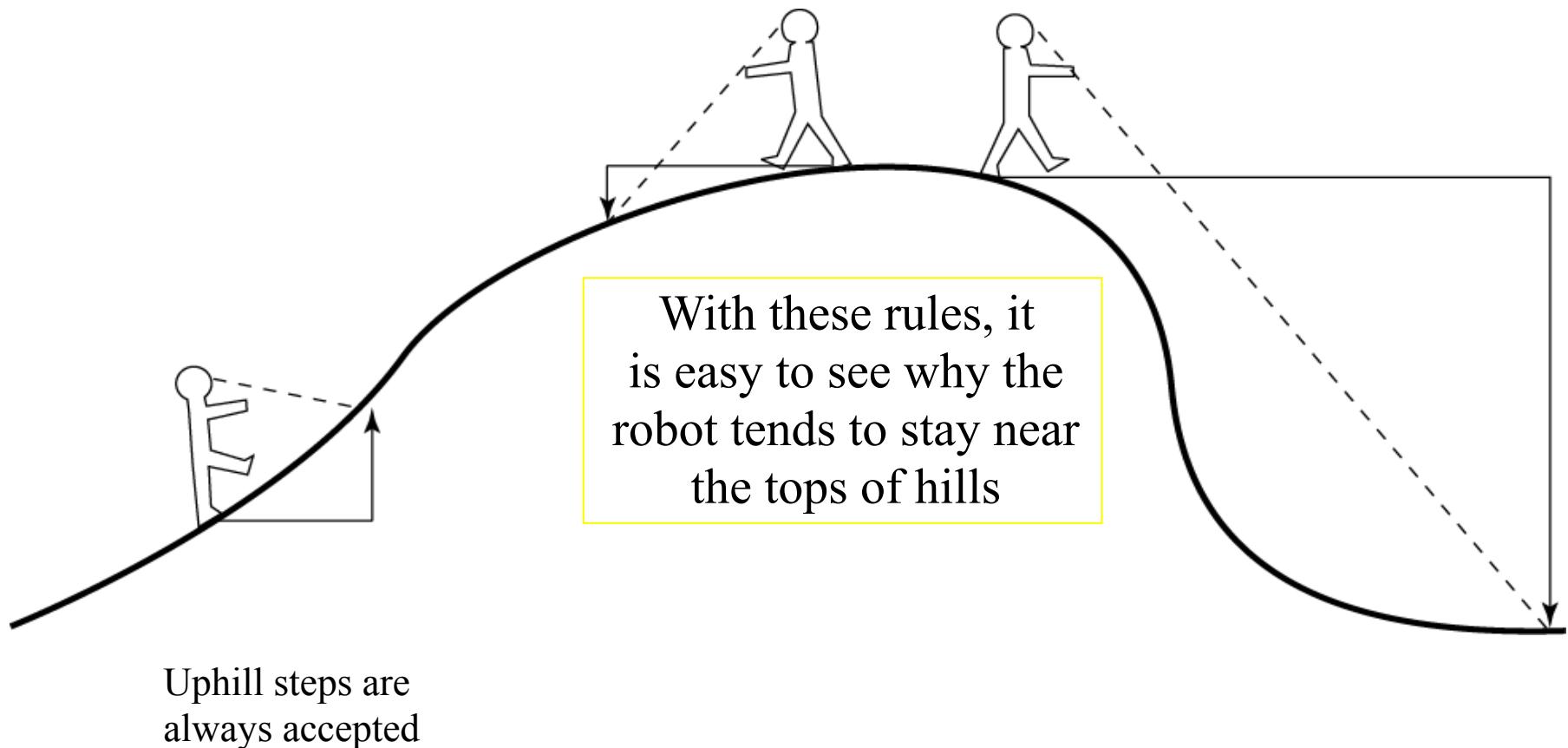
# Markov Chain Monte Carlo (MCMC)

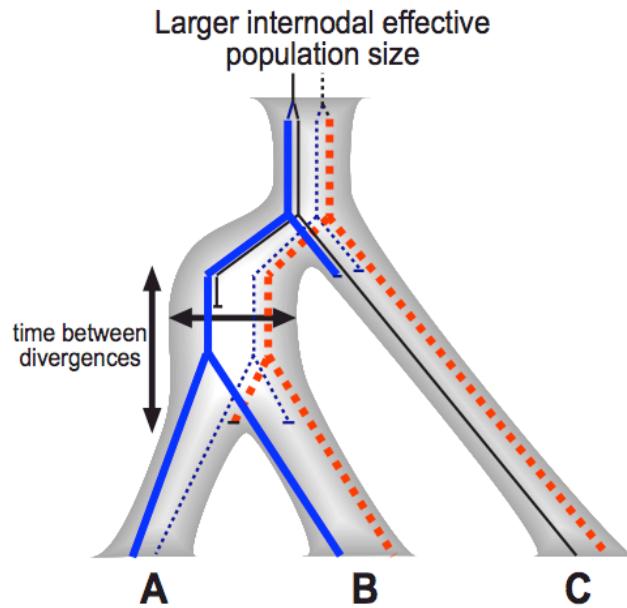


# MCMC robot's rules

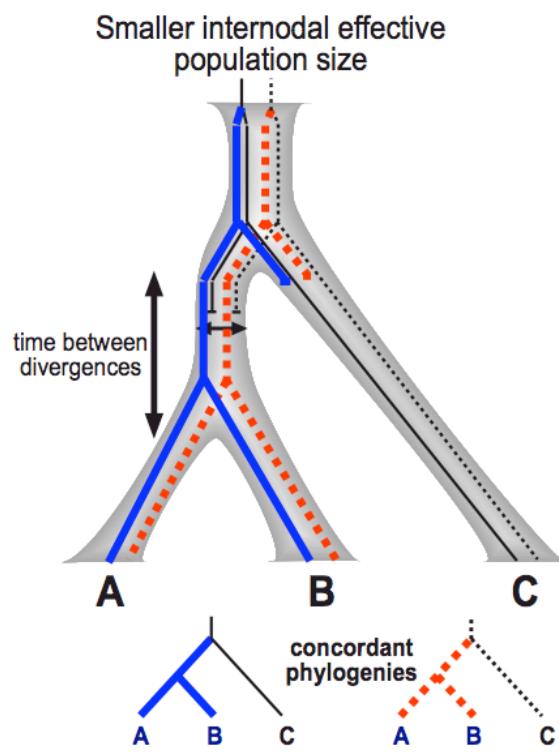
Slightly downhill steps  
are usually accepted

Drastic “off the cliff”  
downhill steps are almost  
never accepted



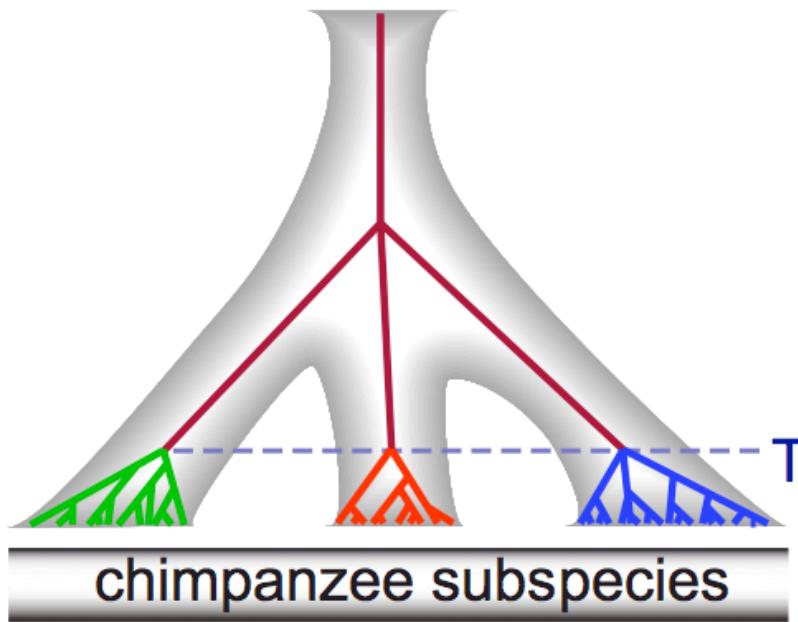


Often the gene tree  
Does not equal the  
“species” tree



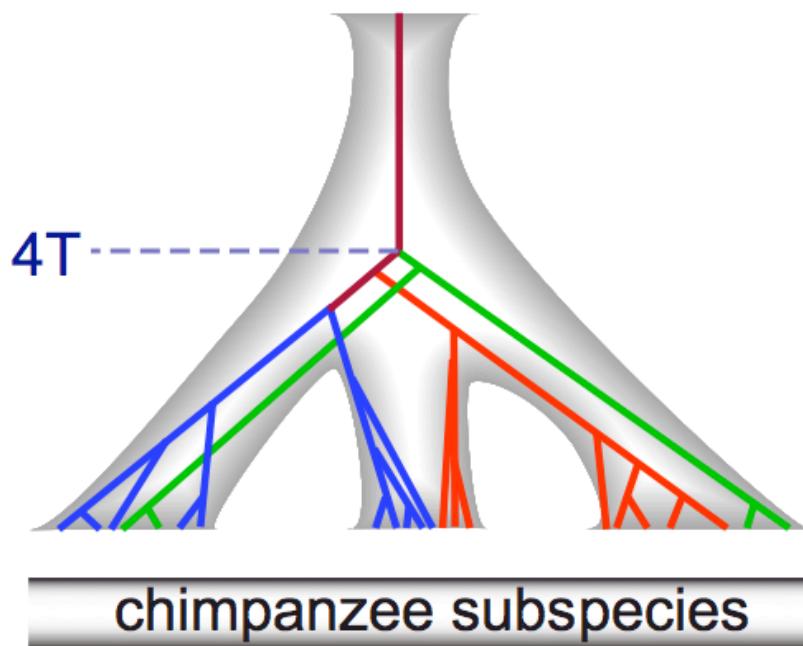
At any one particular  
gene, time to  
Co-ancestry can be  
long

## mitochondrial DNA tree



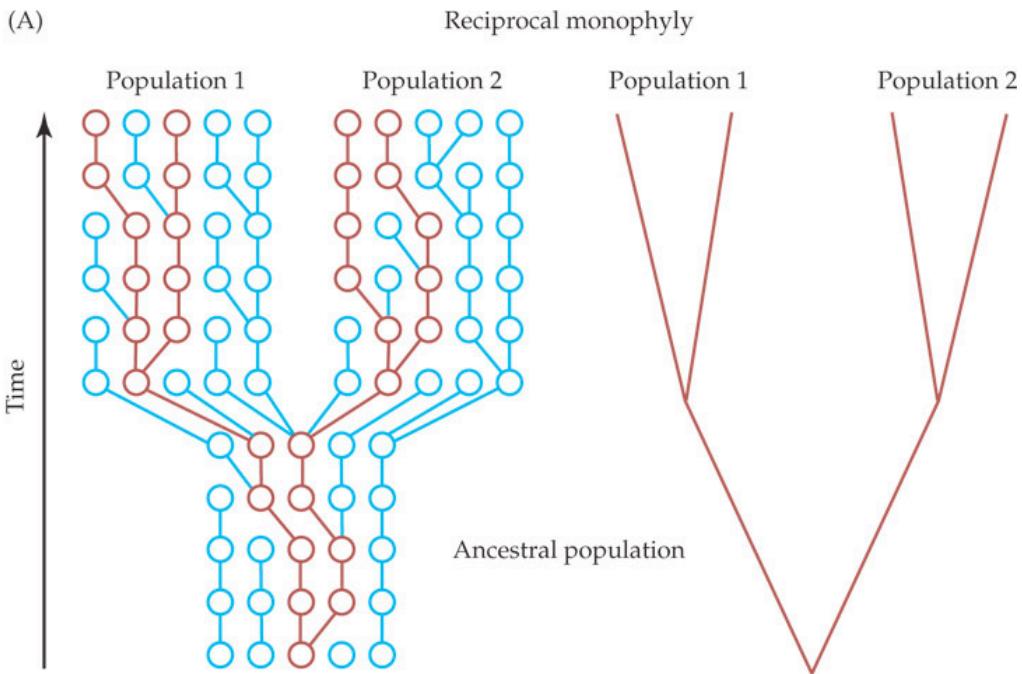
Time to Co-ancestry and  
gene tree/species tree  
Congruence depends on  
Ne (effective population  
sizes)

## Autosomal gene trees

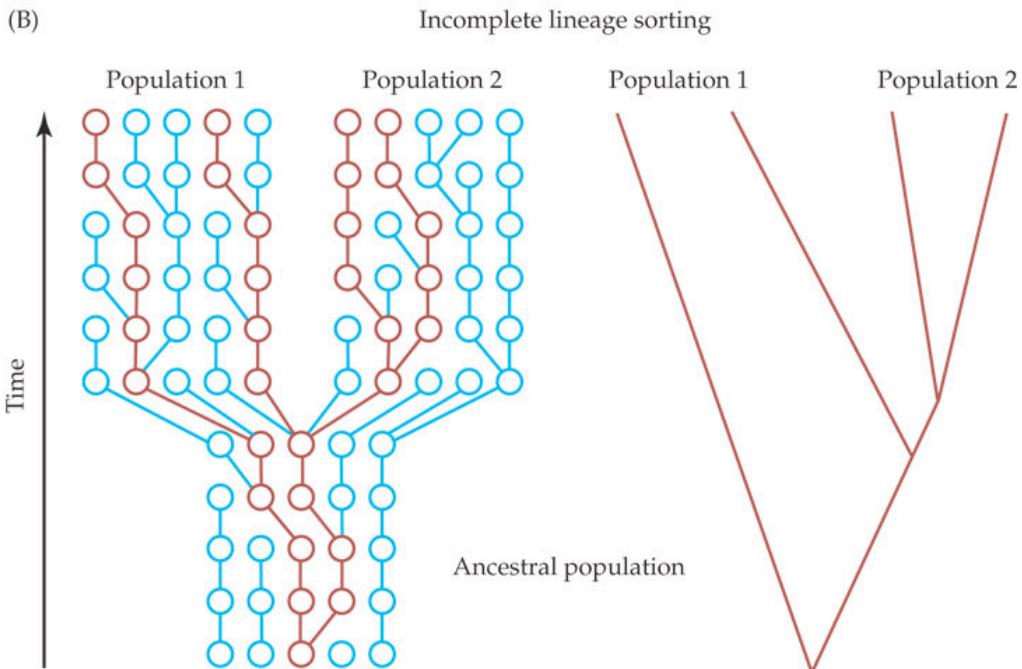


Autosomes have 4x size  
as Y and mtDNA

(A)

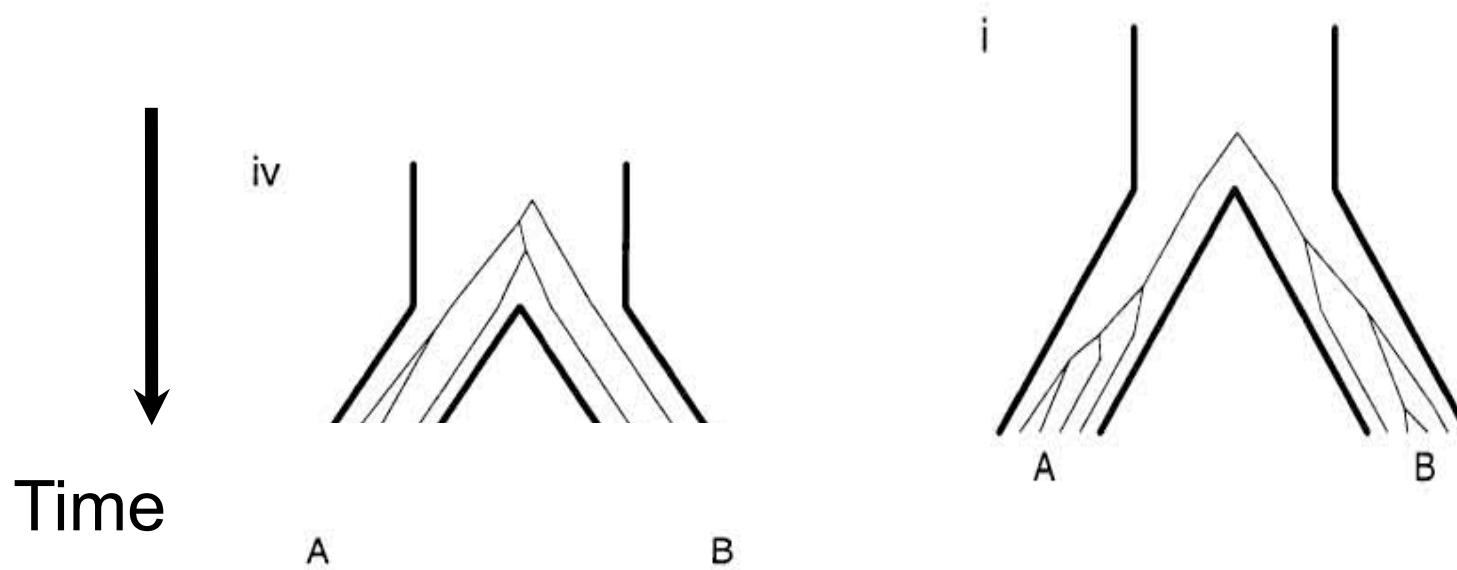


(B)



# Incomplete Lineage Sorting ILS

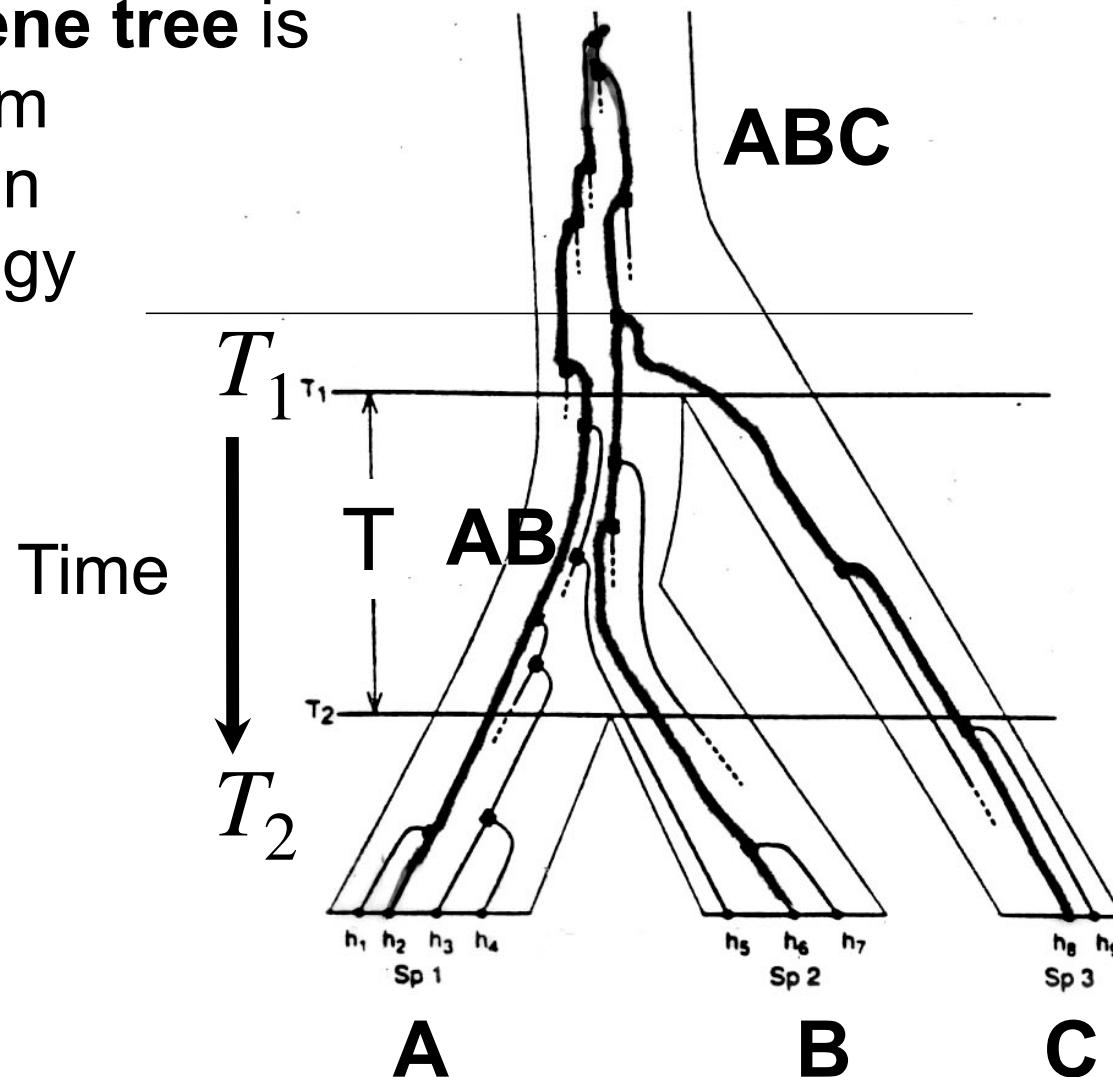
## 2 taxon species tree

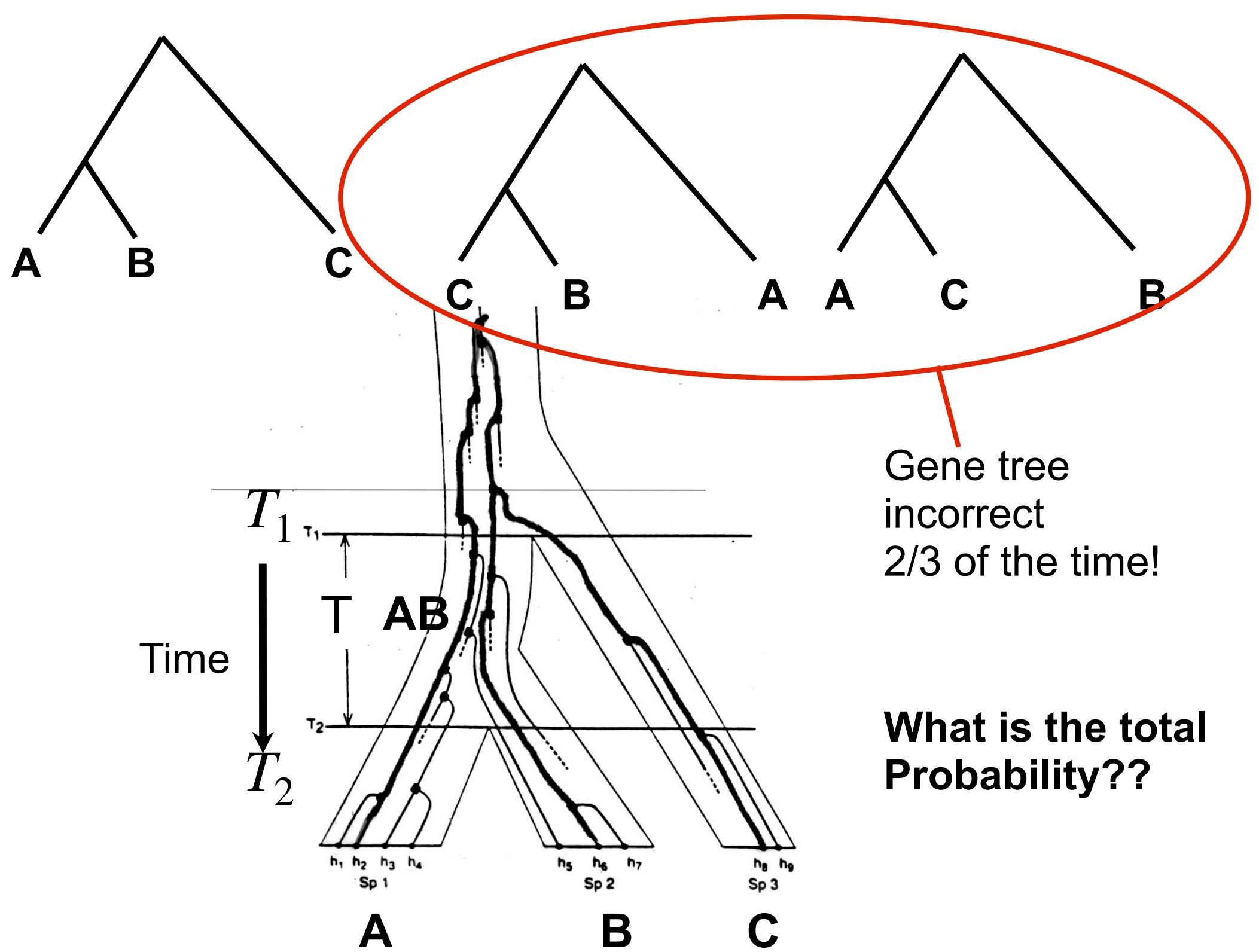


$\Pr(\text{Gene tree} = \text{Species tree})$   
Increases with time

If no coalescent in ancestral species AB,  
(T generations)  
the **gene tree** is  
random  
3 taxon  
topology

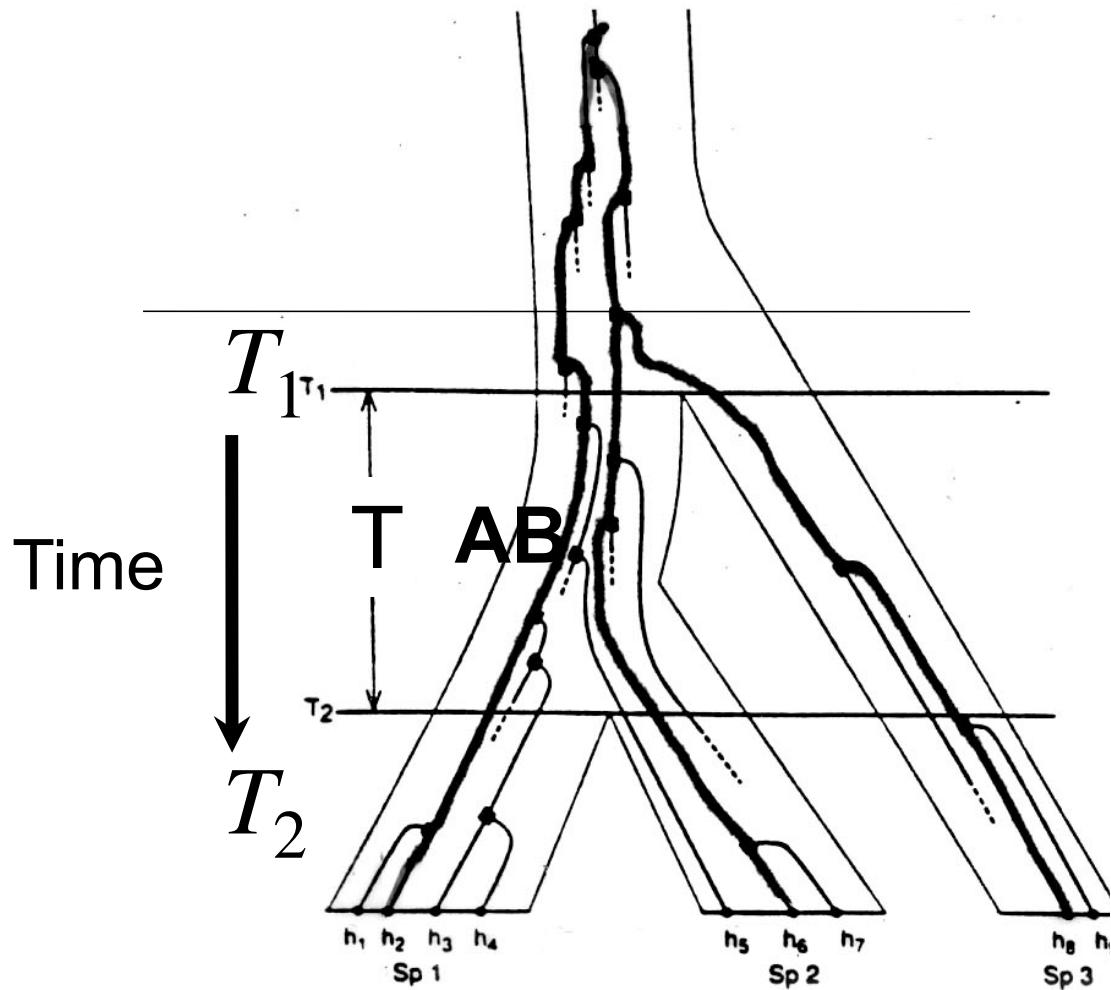
and Gene Tree  
incorrect  
2/3 of the time!





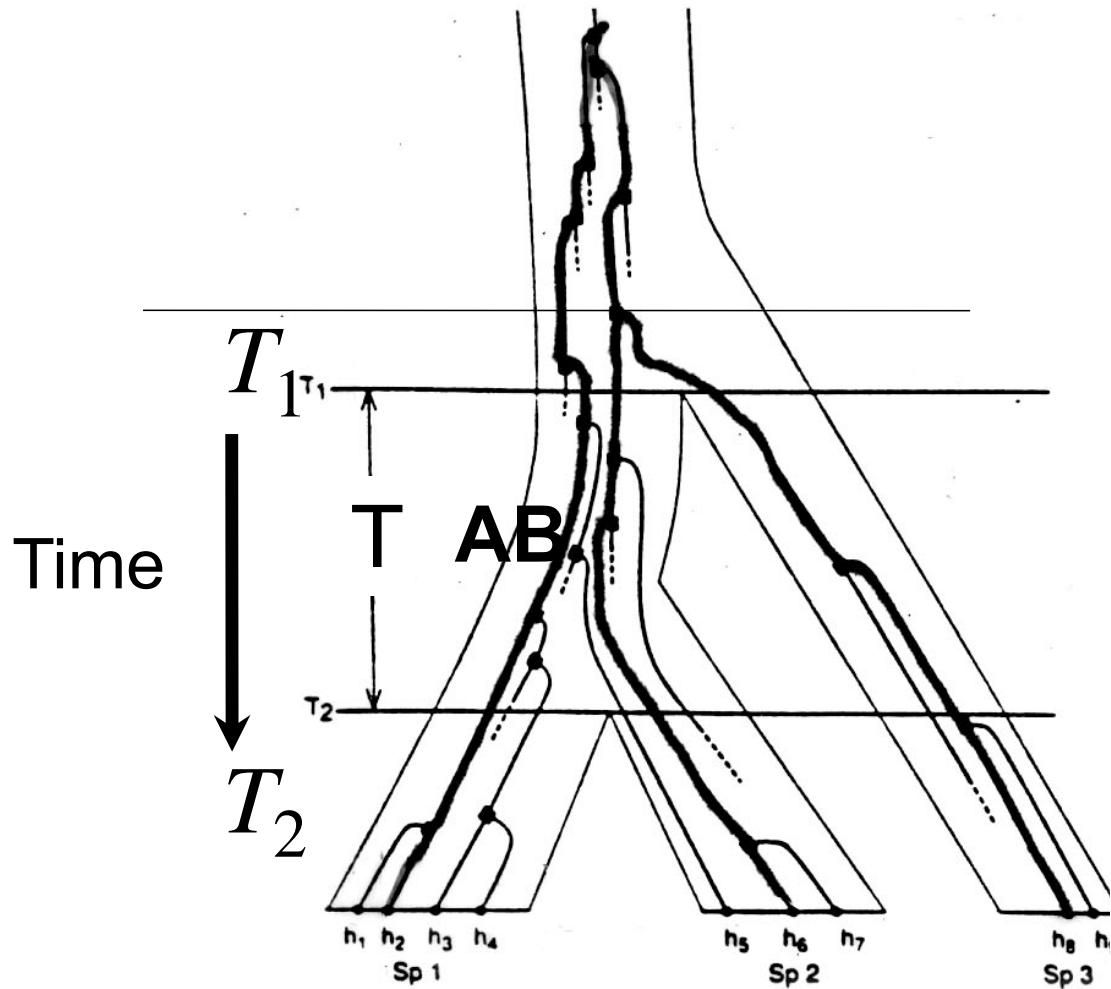
Because probability of no coalescent events  
in  $t$  generations is

$$Pr(k_t = k_0) = \left[1 - \frac{k_0(k_0 - 1)}{4N}\right]^t$$



Because probability of no coalescent events  
in  $t$  generations is

$$Pr(k_t = k_0) = \left[1 - \frac{k_0(k_0 - 1)}{4N}\right]^t$$



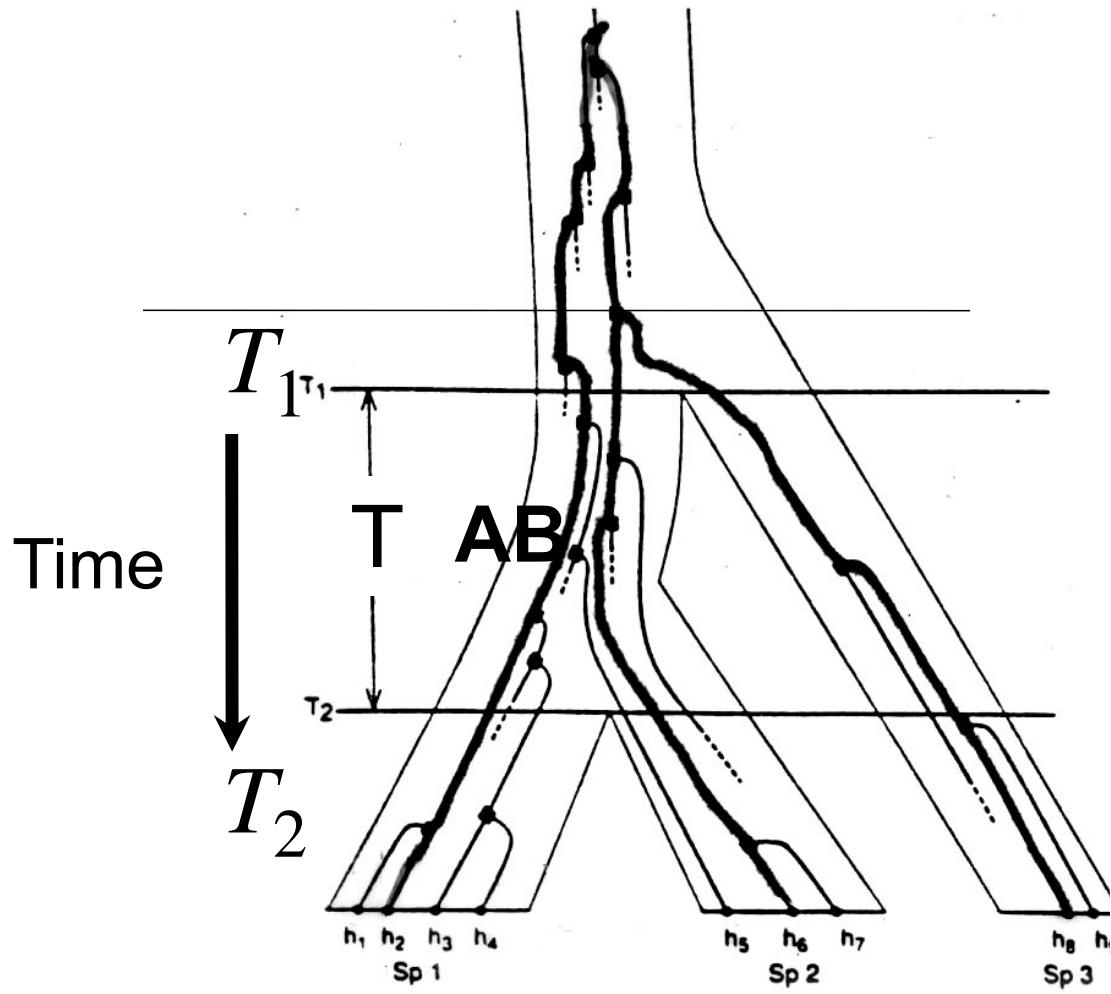
$$\left[1 - \frac{k_0(k_0 - 1)}{4N}\right]^t$$

$$\left[1 - \frac{2(2 - 1)}{4N}\right]^t$$

$$\left[1 - \frac{1}{2N}\right]^t$$

Because probability of no coalescent events  
in  $t$  generations is

$$Pr(k_t = k_0) = \left[1 - \frac{k_0(k_0 - 1)}{4N}\right]^t$$



$$\left[1 - \frac{k_0(k_0 - 1)}{4N}\right]^t$$

$$\left[1 - \frac{2(2 - 1)}{4N}\right]^t$$

$$\left[1 - \frac{1}{2N}\right]^t$$

$$e^{-t/2N}$$

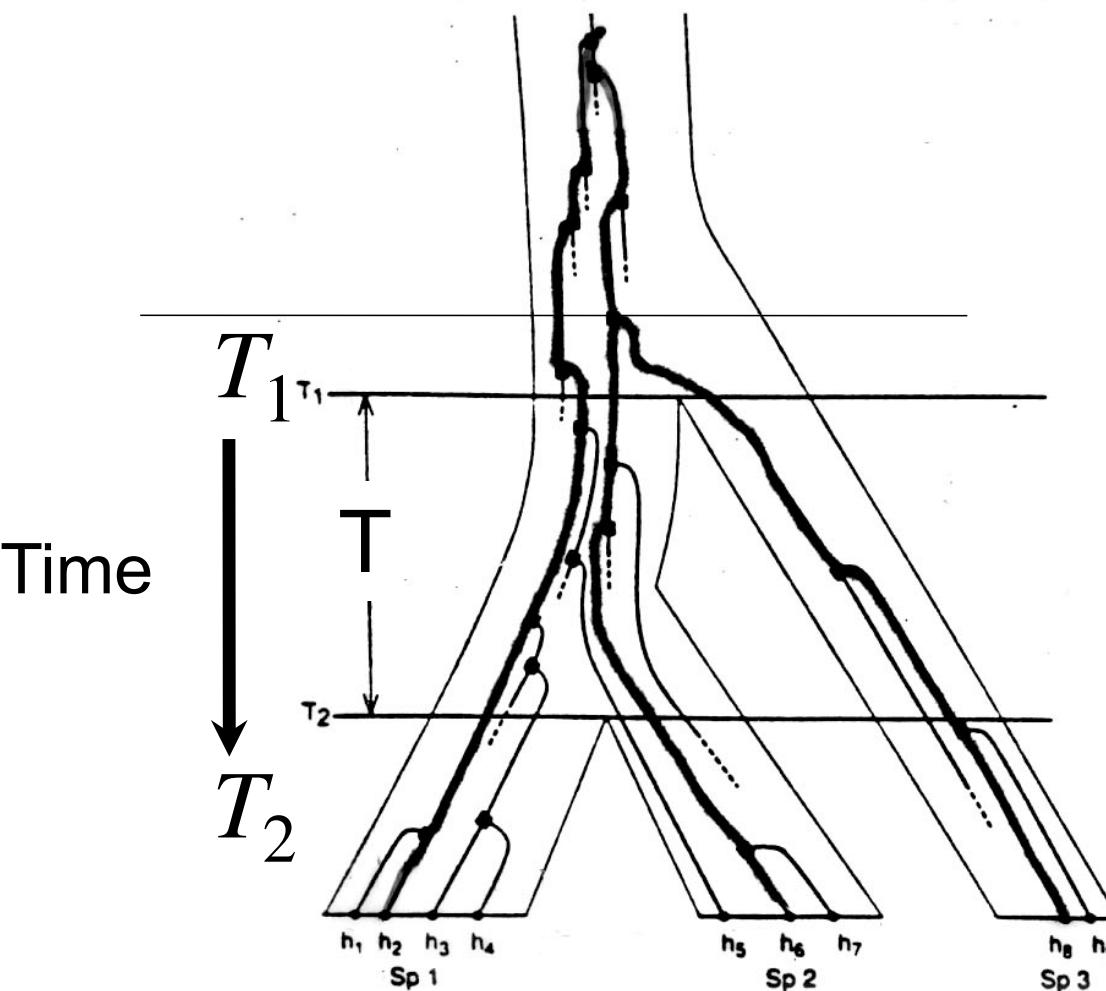
So, .... probability of  
NO coal events between  
 $T_1$  and  $T_2$  ( $T$ ) =  $e^{-t/2N}$

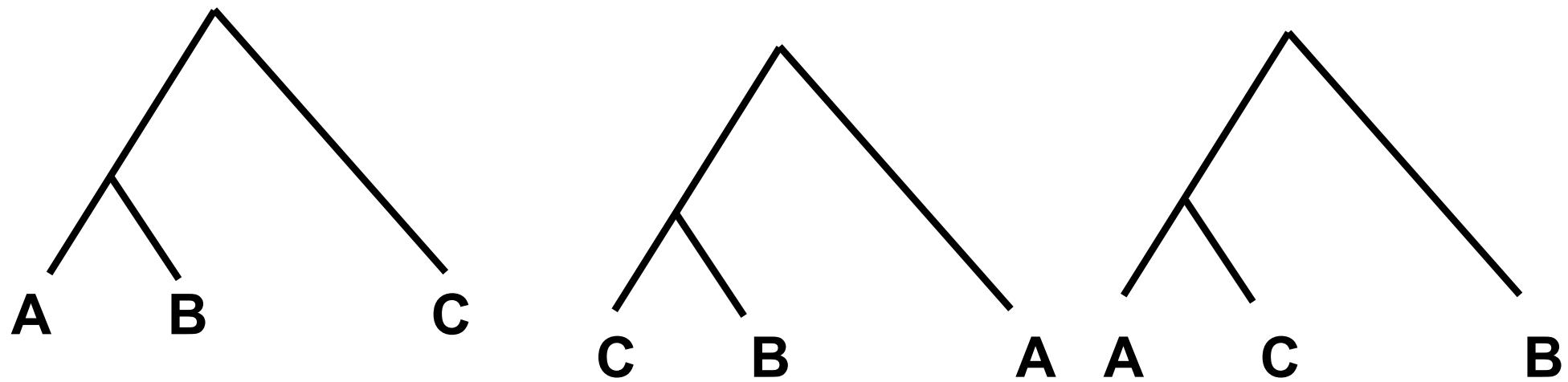
AND

Probability of a random  
3-taxon topology =  
 $2/3$

Therefore probability of  
Gene tree  $\neq$  Species tree =

$$\frac{2}{3}e^{-t/2N}$$





congruent gene tree

incongruent gene tree

incongruent gene tree

Nei 1987. *Molecular Evolutionary Genetics*

Wu 1991. *Genetics* 127:429-435

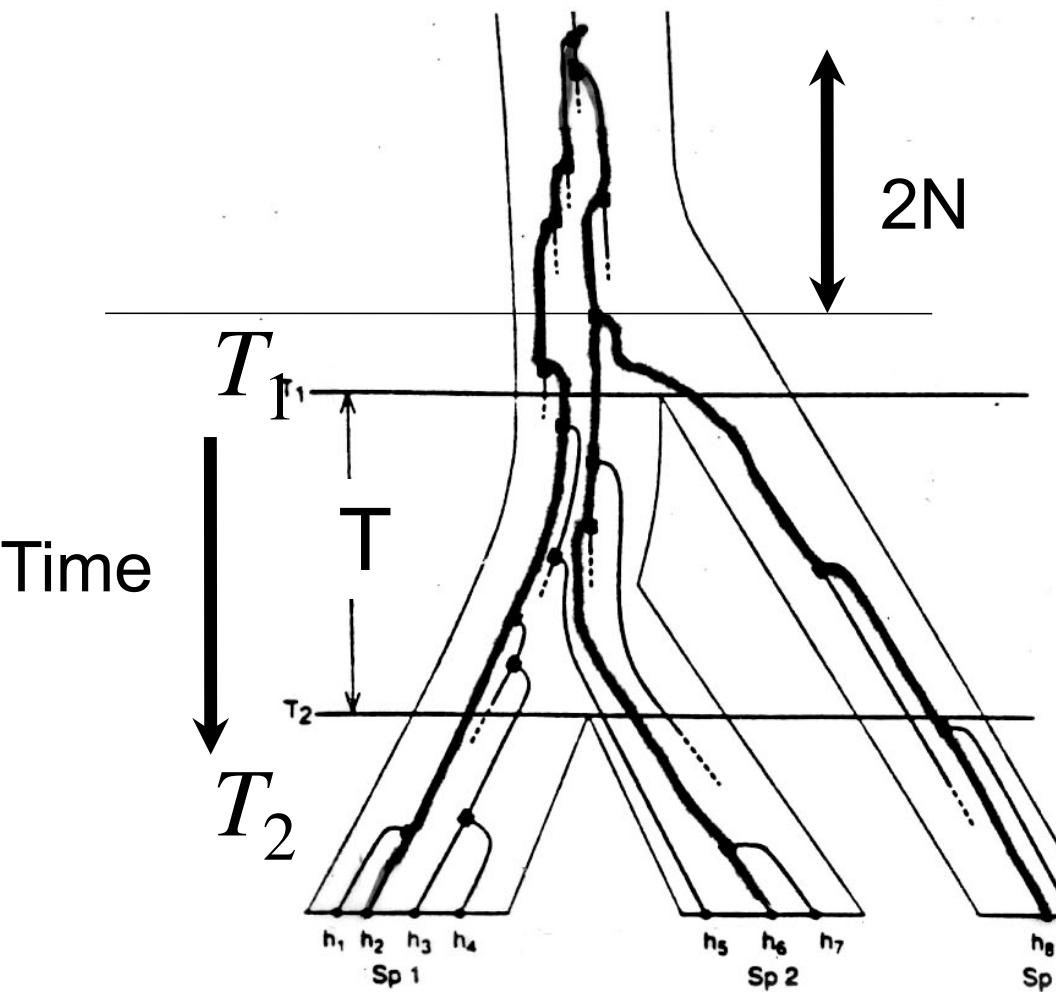
Hudson 1992. *Genetics* 131:509-512

Yang 2002. *Genetics* 162:1811-1823

Rannala B, Yang Z. 2003. *Genetics* 164, 1645-1656.

$$Pr(\text{incongruence}) = \frac{2}{3}e^{-t/2N}$$

$$P_{incongruence} = \frac{2}{3} * e^{-t/2N}$$

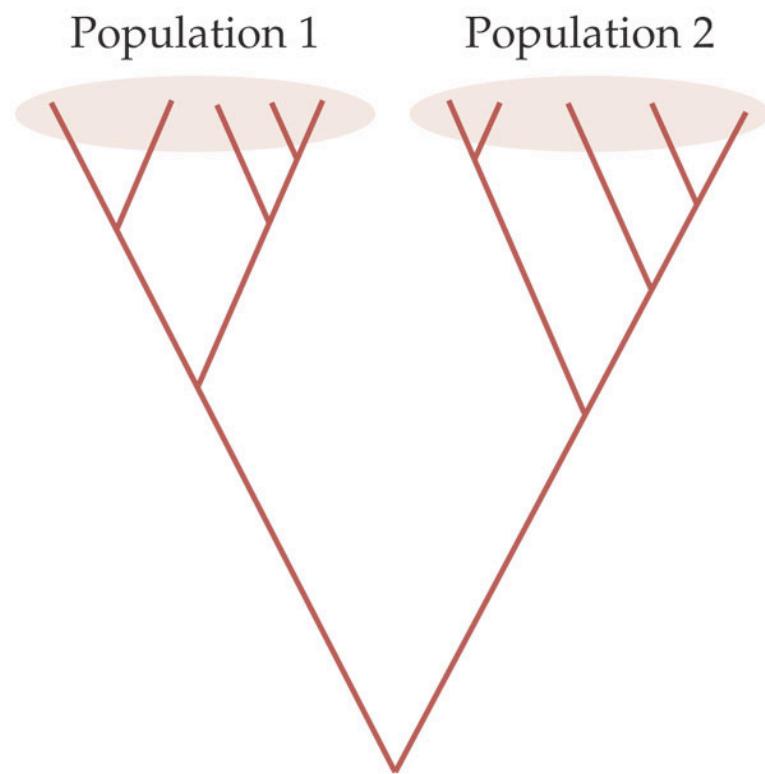


1.  $\Pr(A \& B \text{ do not coalesce in interval } T)$

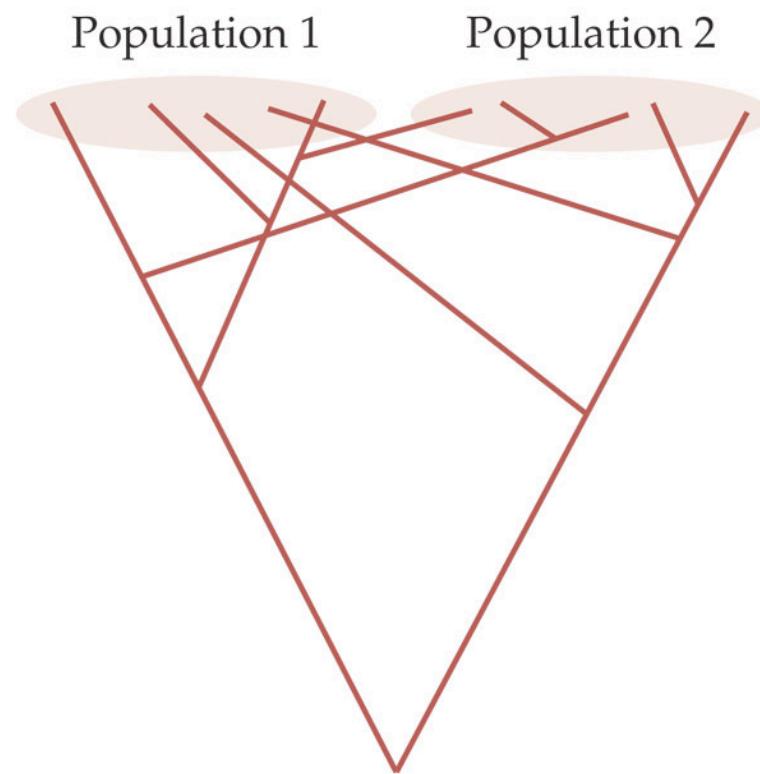
and

2.  $\Pr(A \& B \text{ do not coalesce before } A \& C \text{ or } B \& C)$

(A) Old divergence, little gene flow



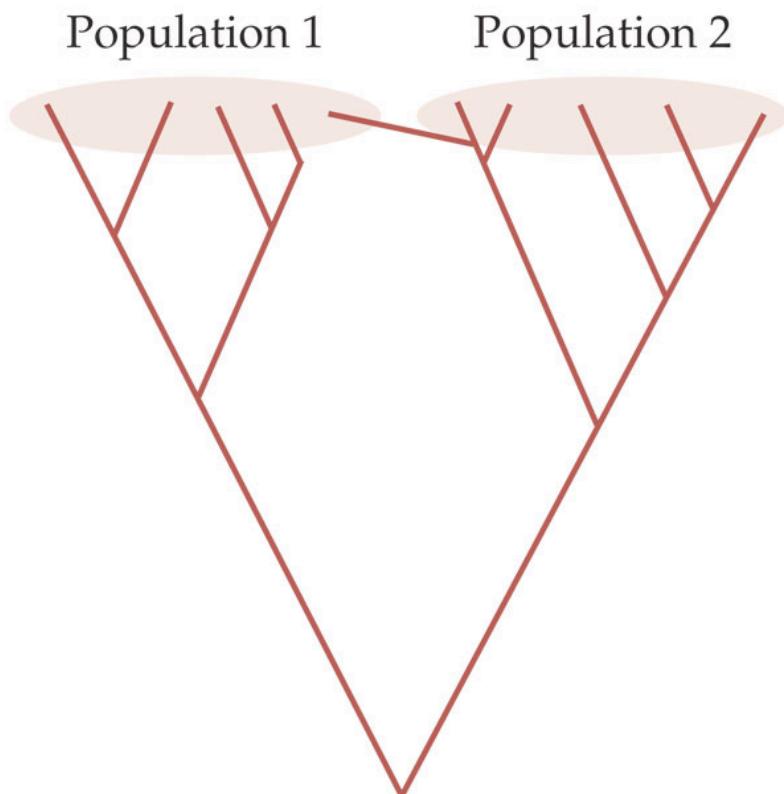
(B) Strong gene-flow, panmixia, very recent divergence



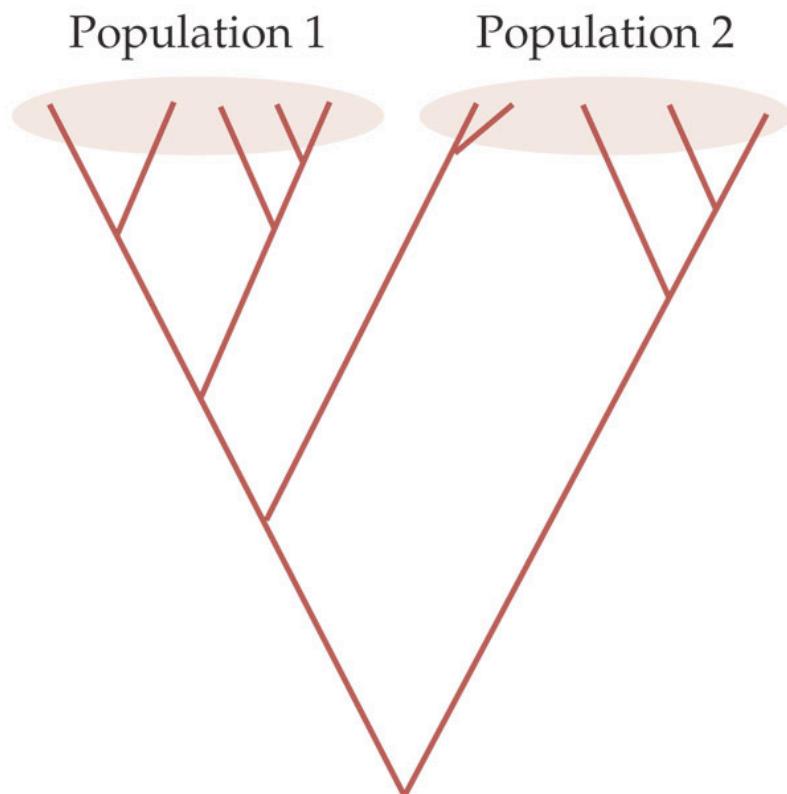
**INTRODUCTION TO POPULATION GENETICS, Figure 5.8 (Part 1)**

© 2013 Sinauer Associates, Inc.

(C) Old divergence, recent gene-flow



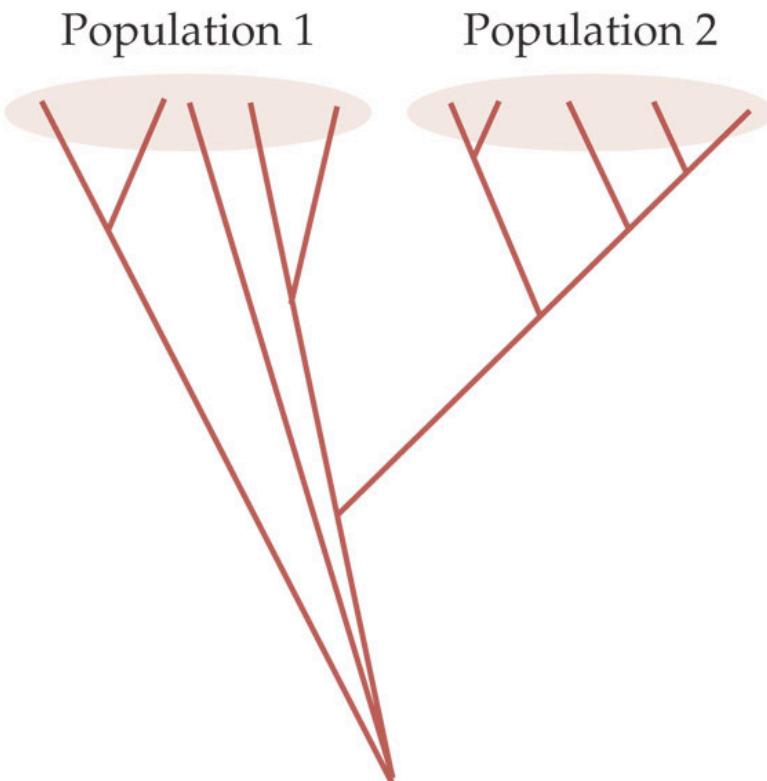
(D) Ongoing gene-flow, old divergence or recent divergence



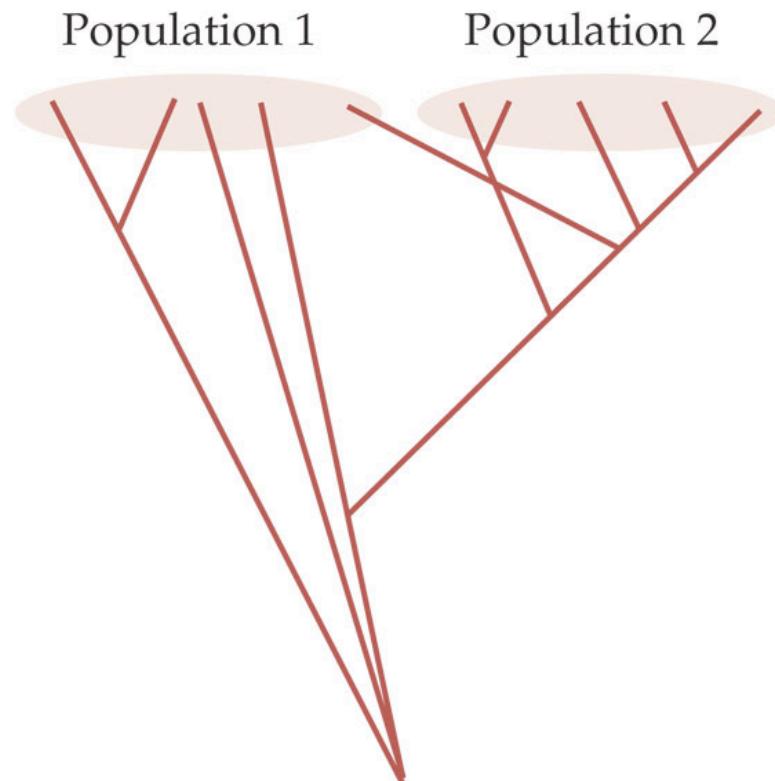
**INTRODUCTION TO POPULATION GENETICS, Figure 5.8 (Part 2)**

© 2013 Sinauer Associates, Inc.

(E) Old divergence or ongoing gene-flow,  
low  $N_e$  in population 2



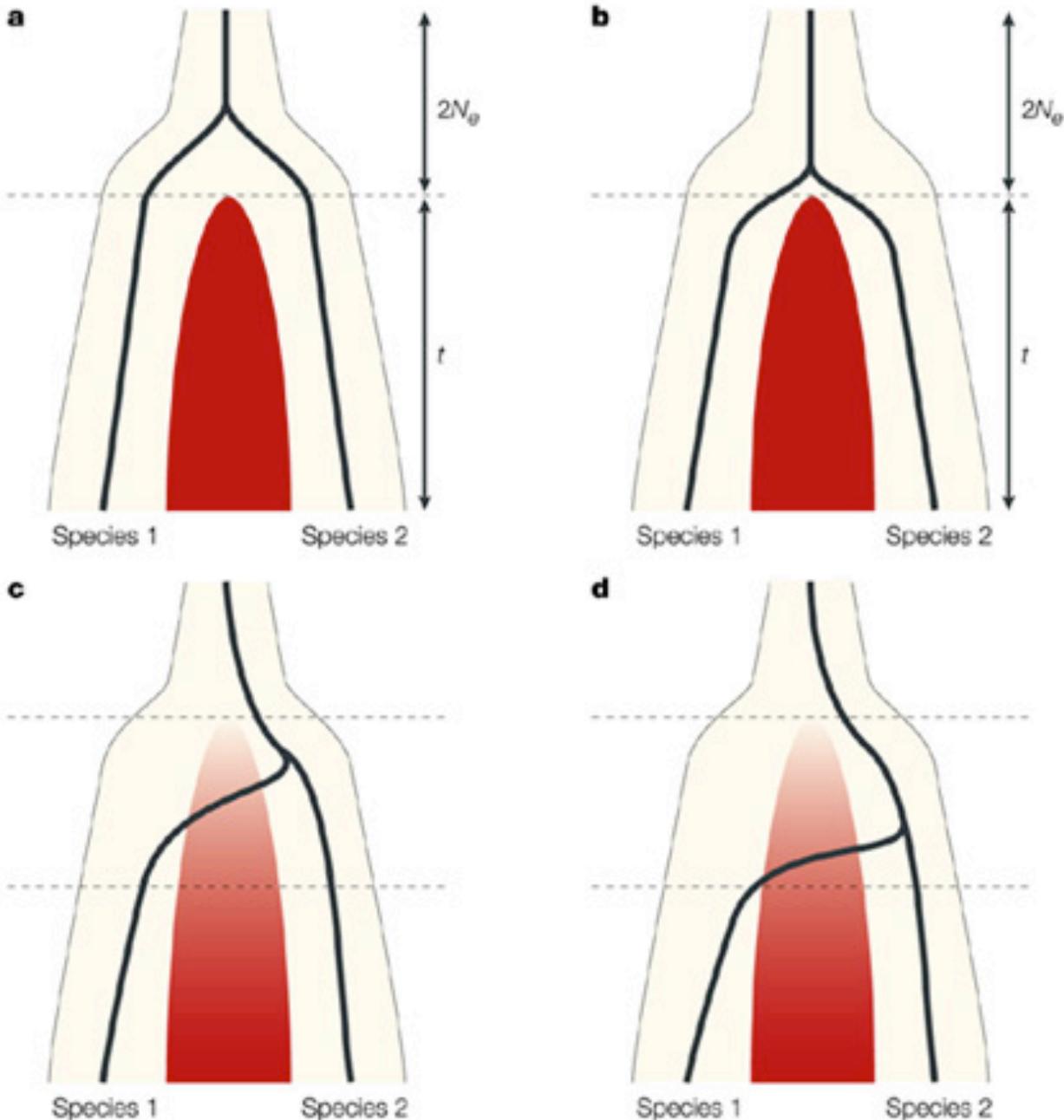
(F) Recent divergence or ongoing gene-flow,  
low  $N_e$  in population 2



**INTRODUCTION TO POPULATION GENETICS, Figure 5.8 (Part 3)**

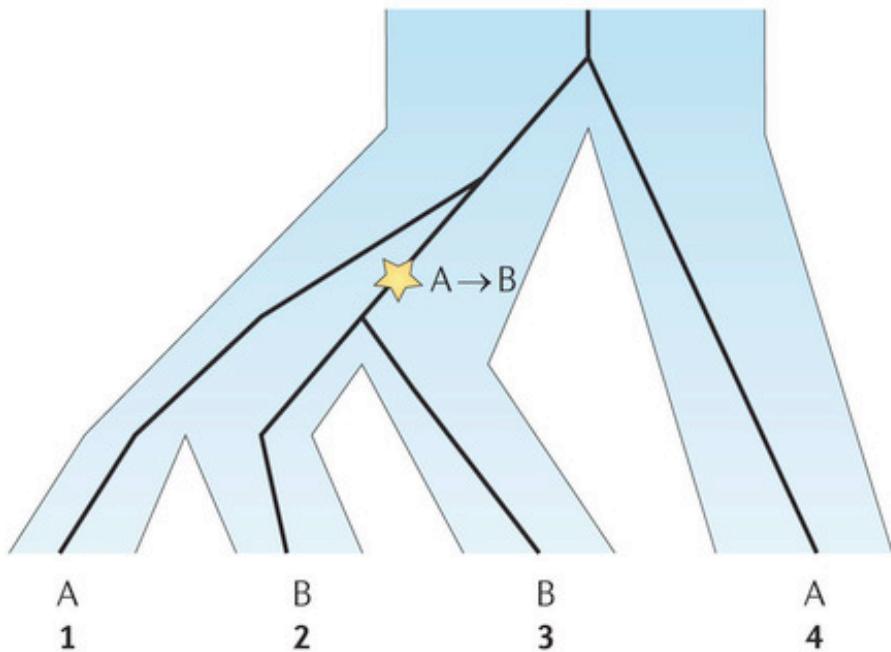
© 2013 Sinauer Associates, Inc.

# Speciation w/ gene flow

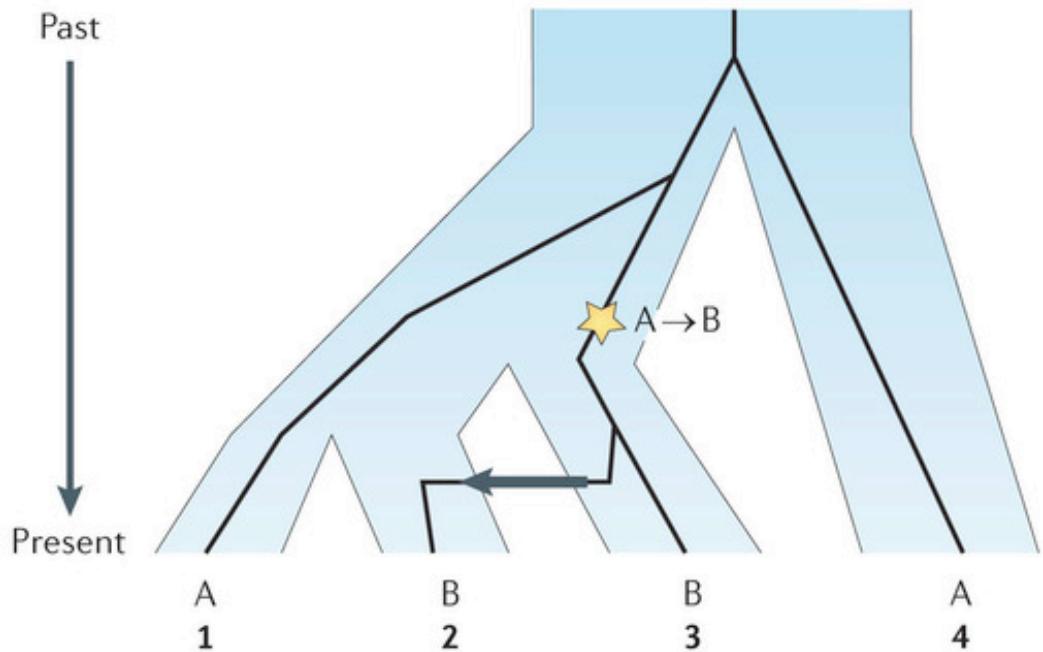


# ABBA/BABA test

a Ancestral polymorphism

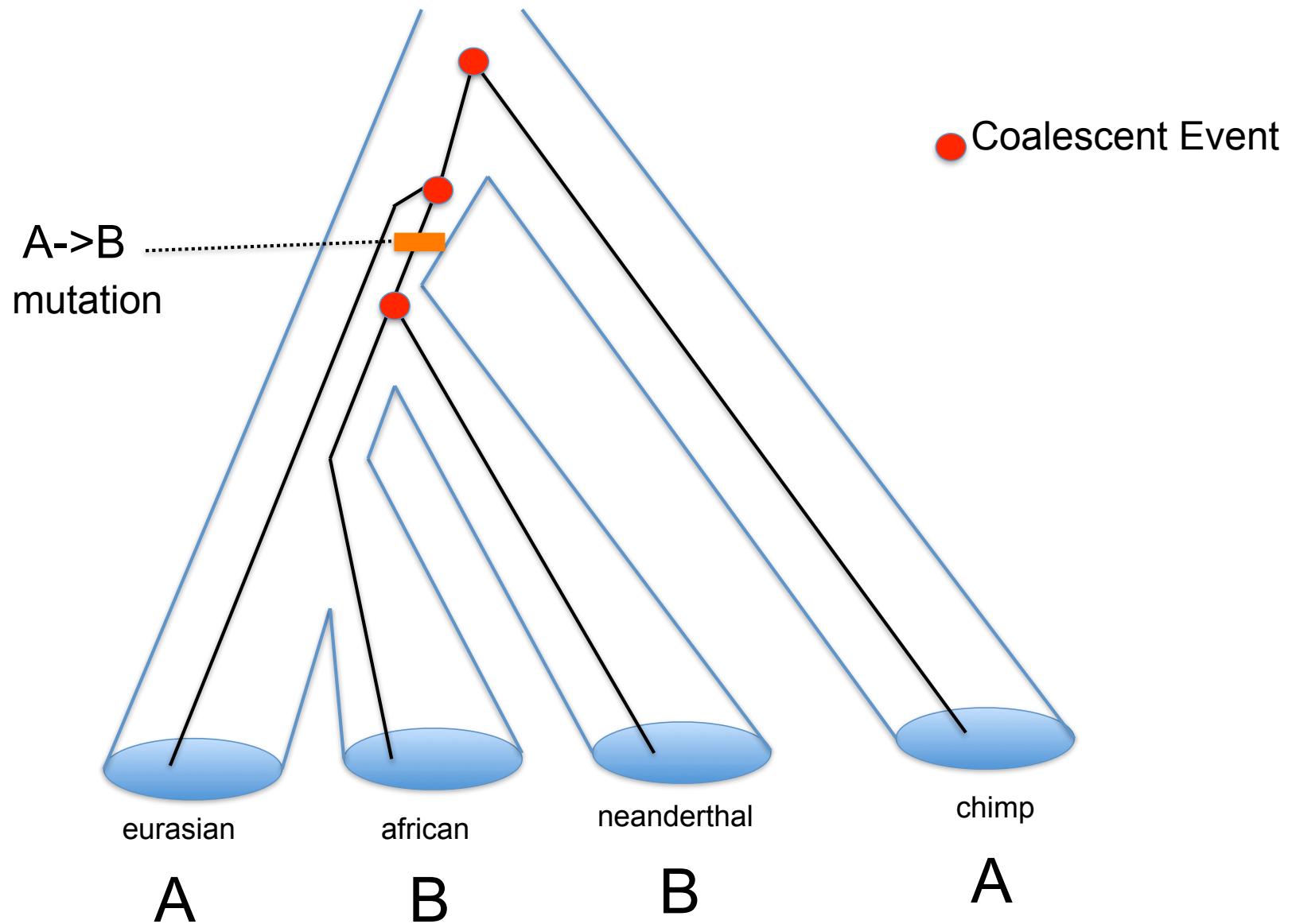


b Introgression (gene flow)

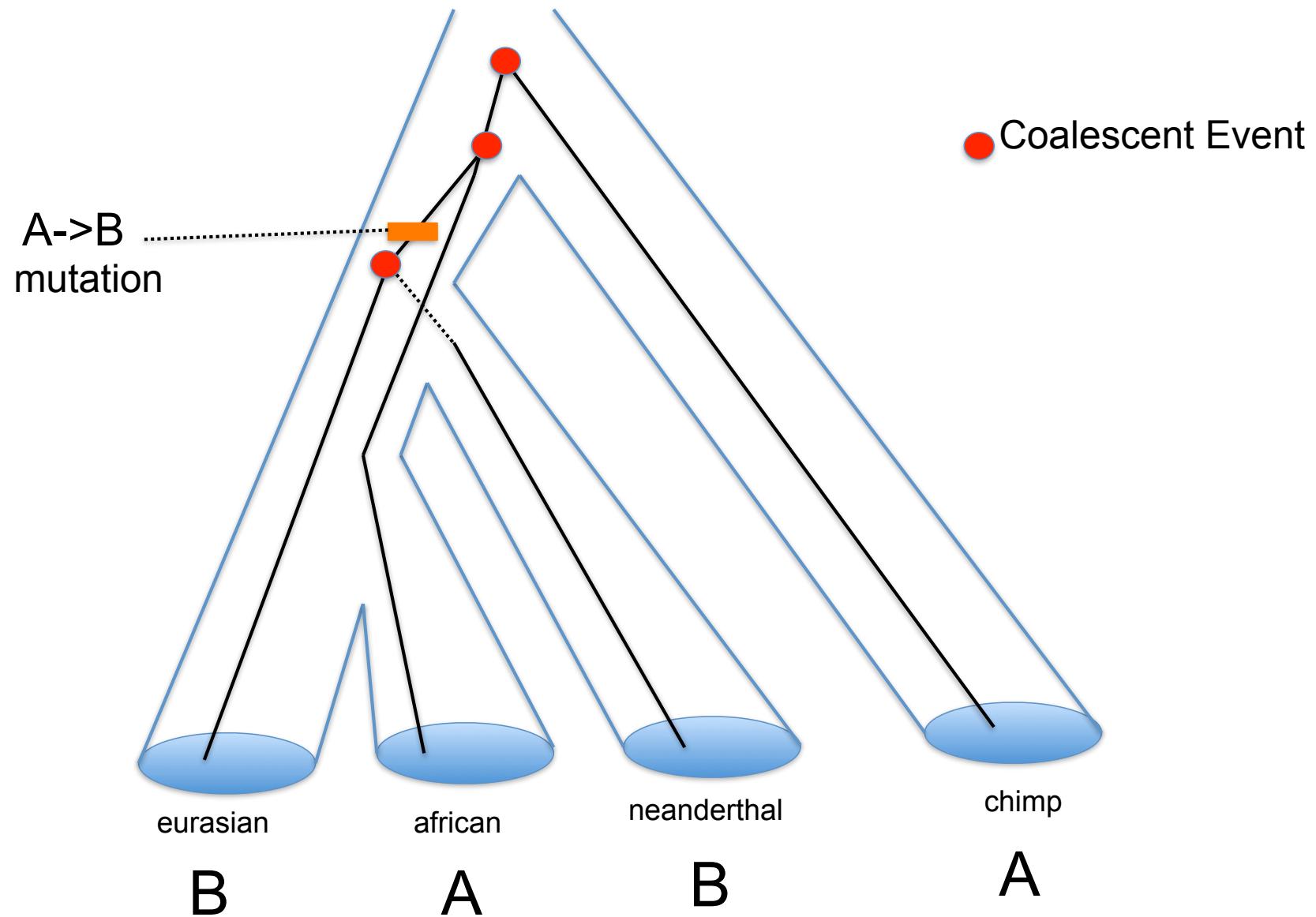


# Did Humans – Neanderthals mix?

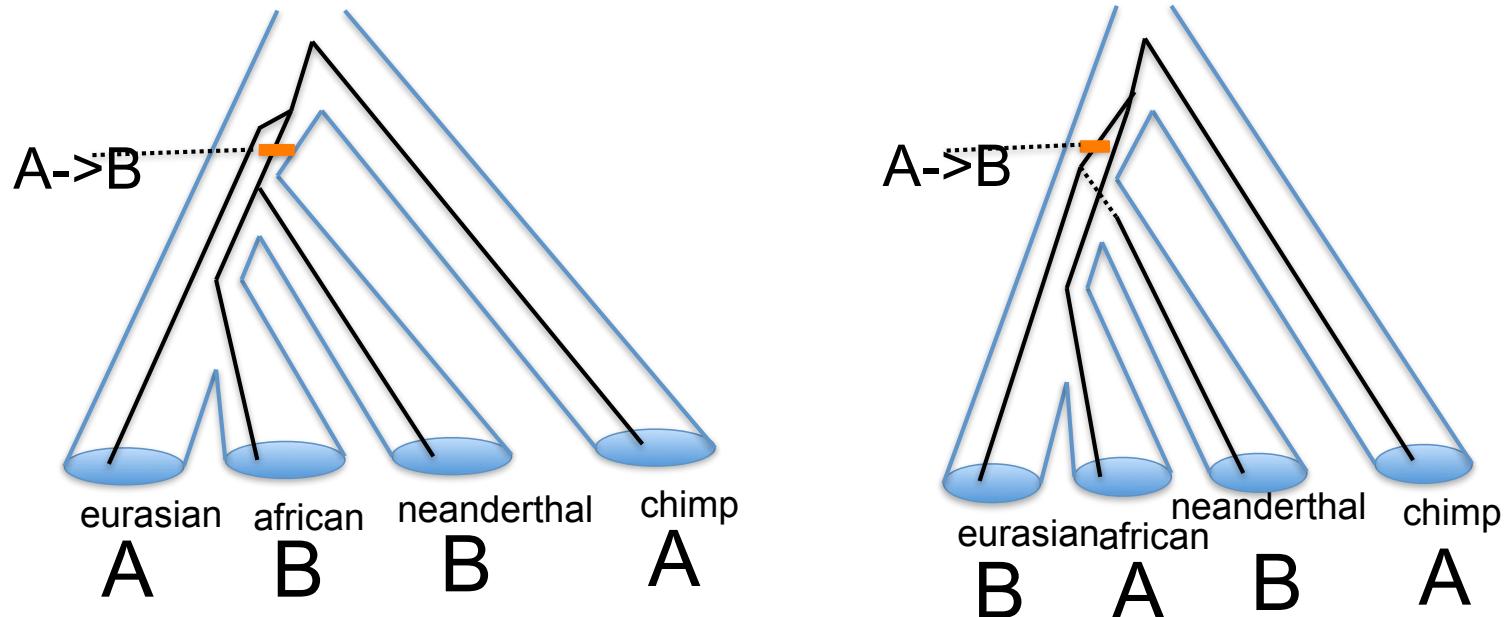
focus on mutations along the internal branch



# Did Humans – Neanderthals mix?



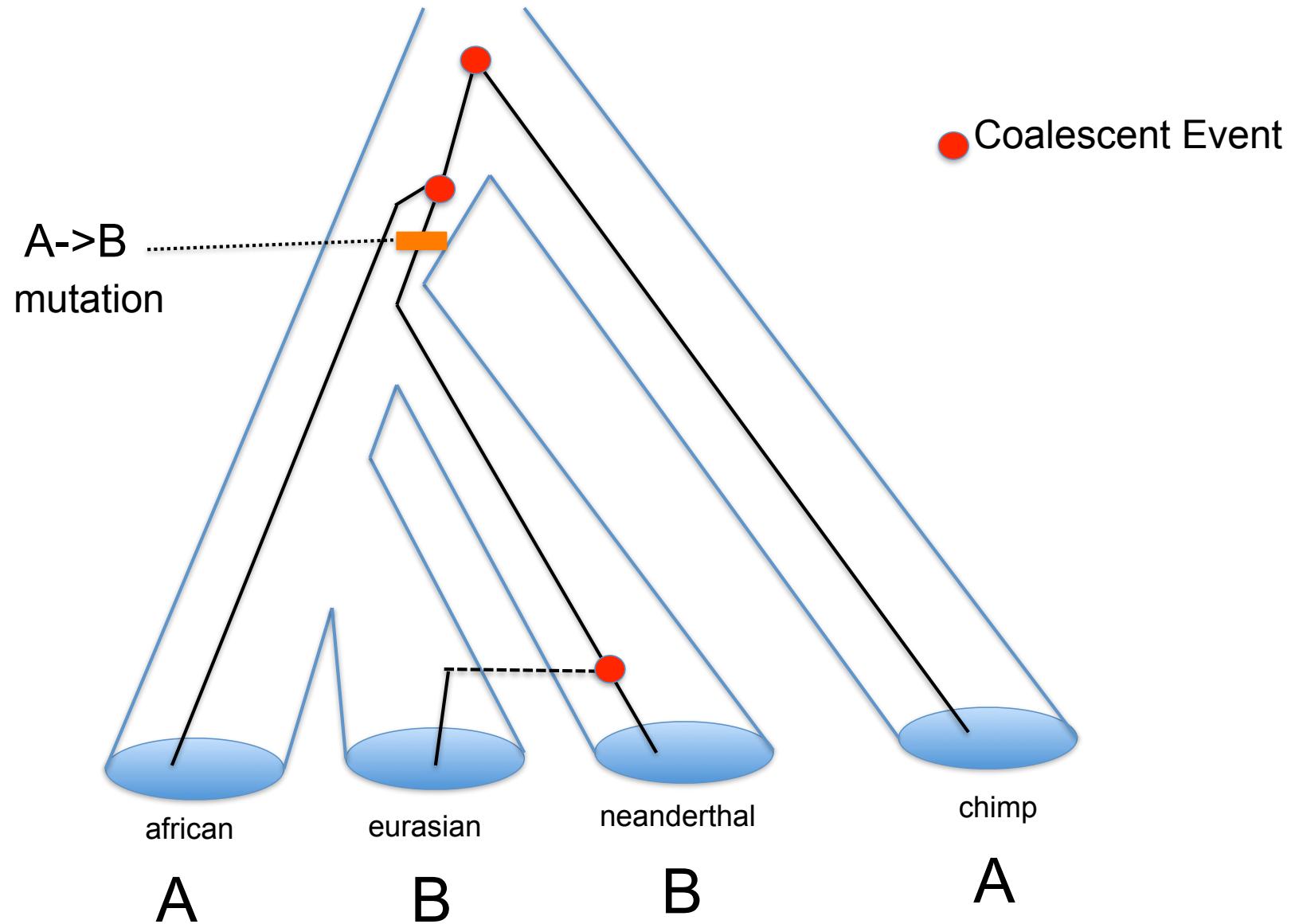
Given no admixture, should be 50% - 50%



ABBA and BABA equally likely if it is Ancestral polymorphism

50%/50% expectation under ancestral polymorphism  
Out of 176,000 SNPs, one compares the observed proportion to the binomial expectation

# Did Humans – Neanderthals mix?



Bottom line - inference of population history from only 1 or a few loci or 1 sumstat is hazardous

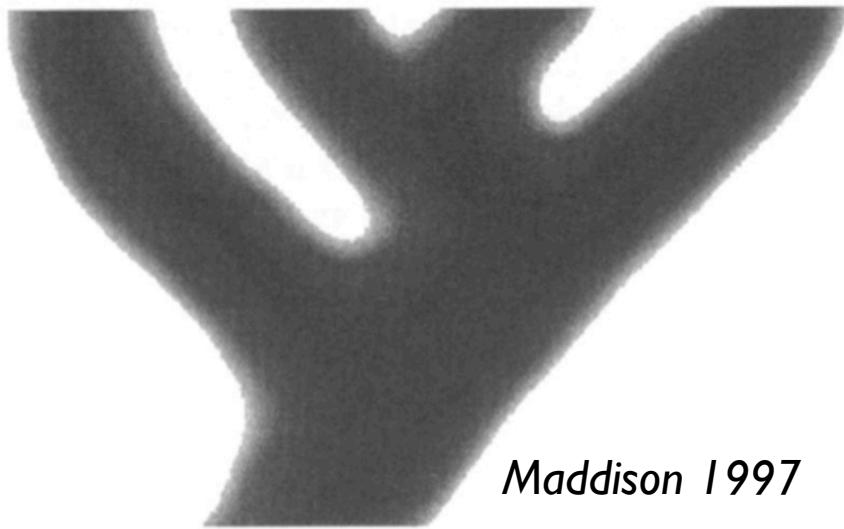
<b>approach</b>	<b>problems</b>
Summary statistic/ simulation	lose a lot of information Can not distinguish histories based on one sumstat
using the estimated gene tree	Can not interpret the history based on one or a few unlinked gene trees  statistical uncertainty in gene tree estimation step (inference only as good as the tree is )

Bottom line - inference of population history from only 1 or a few loci or 1 sumstat is hazardous

make likelihood function,

$$L = \Pr(\text{data} | \Theta)$$

where  $\Theta$  = population history parameters



Maddison 1997

FIGURE 9. Phylogeny as a cloud of gene histories. Phylogeny is more like a statistical distribution than a simple tree of discrete thin branches. It has a central tendency, but it also has a variance because of the diversity of gene trees. Gene trees that disagree with the central tendency are not wrong; rather, they are part of the diffuse pattern that is the genetic history.

Treat the population history like  
you would a phylogenetic tree  
in a likelihood frame work

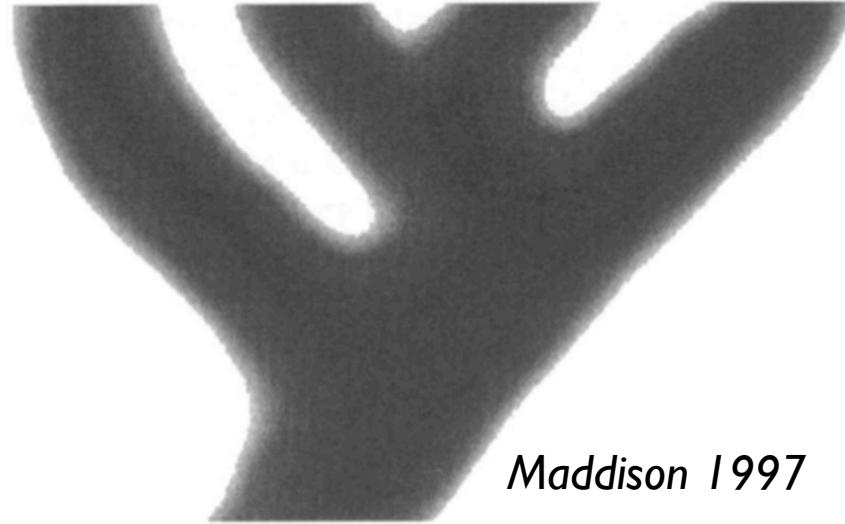


FIGURE 9. Phylogeny as a cloud of gene histories. Phylogeny is more like a statistical distribution than a simple tree of discrete thin branches. It has a central tendency, but it also has a variance because of the diversity of gene trees. Gene trees that disagree with the central tendency are not wrong; rather, they are part of the diffuse pattern that is the genetic history.

account for coalescent stochasticity

$$\prod_{\text{loci}} \sum_{\substack{\text{possible} \\ \text{gene trees}}} [P(\text{sequences} | \text{gene tree}) \\ \cdot P(\text{gene tree} | \text{species tree})]$$

Maddison 1997

$$\Pr(X | \Theta) = \int_G \Pr(X | G) p(G | \Theta) dG$$

G = coalescence gene trees

X = sequence alignments

$\Theta$  = species tree & parameters

$$\Pr(X | \Theta) = \int_G \Pr(X | G) p(G | \Theta) dG$$

$G$  = coalescent gene trees

considers all possible trees  
(weighted by their likelihood)

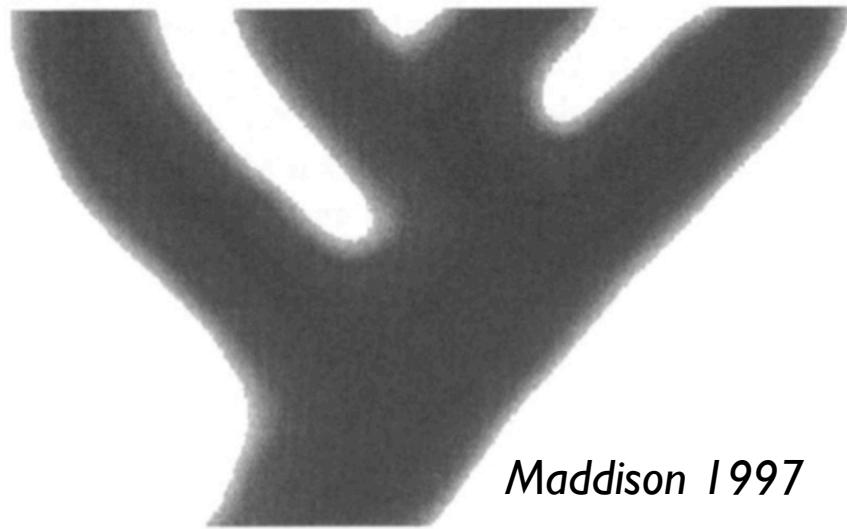


FIGURE 9. Phylogeny as a cloud of gene histories. Phylogeny is more like a statistical distribution than a simple tree of discrete thin branches. It has a central tendency, but it also has a variance because of the diversity of gene trees. Gene trees that disagree with the central tendency are not wrong; rather, they are part of the diffuse pattern that is the genetic history.

can be used to test hypotheses re:  
population history (ie  $\Theta_1$  vs  $\Theta_2$ )

can be very hard to calculate for all  
possible  $G$  (often use MCMC)

sometimes even MCMC is too difficult

## **Too many parameters to calculate likelihood function**

..... **but we can simulate and approximate (ABC; approximate Bayesian computation or now, supervised machine learning)**

**simulate data**

**compress to vector of sumstats**

**parameters that produce data close to observed approximates**       $\Pr(X | \Theta)$

To the limit of  $D_{simulated} - D_{observed} = 0$

and

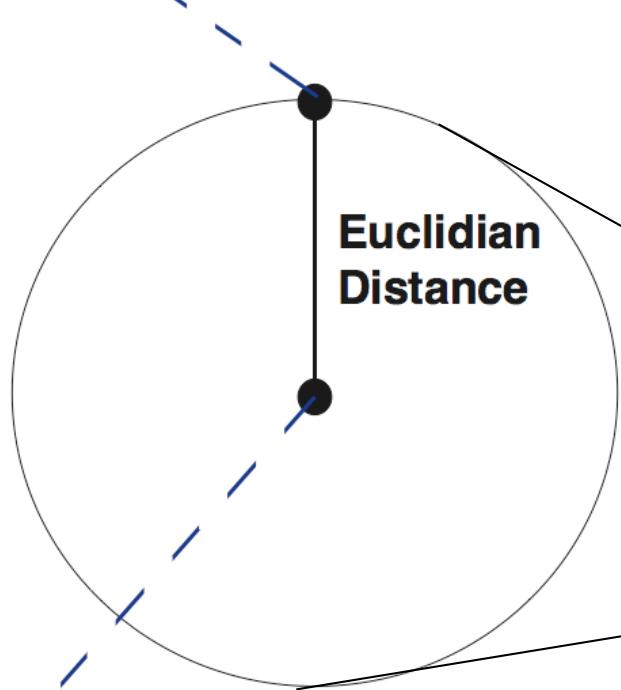
as # of simulations  $\longrightarrow \infty$

the ABC posterior approaches true posterior

$P(\Phi^* | \text{data}) \longrightarrow P(\Phi | \text{data})$

# ABC (simulation based inference)

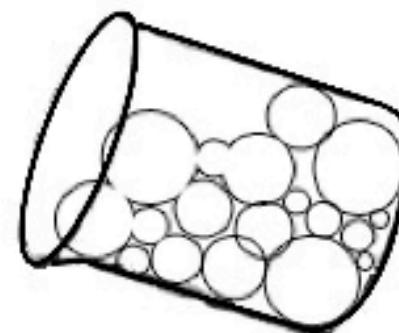
**simulated  
Sum Stat  
Vector**



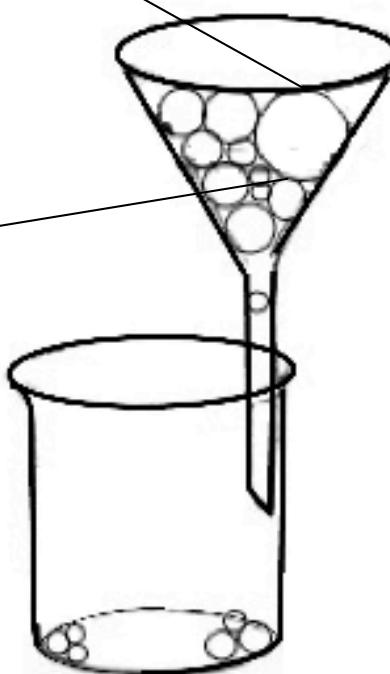
**observed  
Sum Stat  
Vector**



Wen Huang 2008



Simulate data  
w/ randomly  
drawn  
parameter  
values

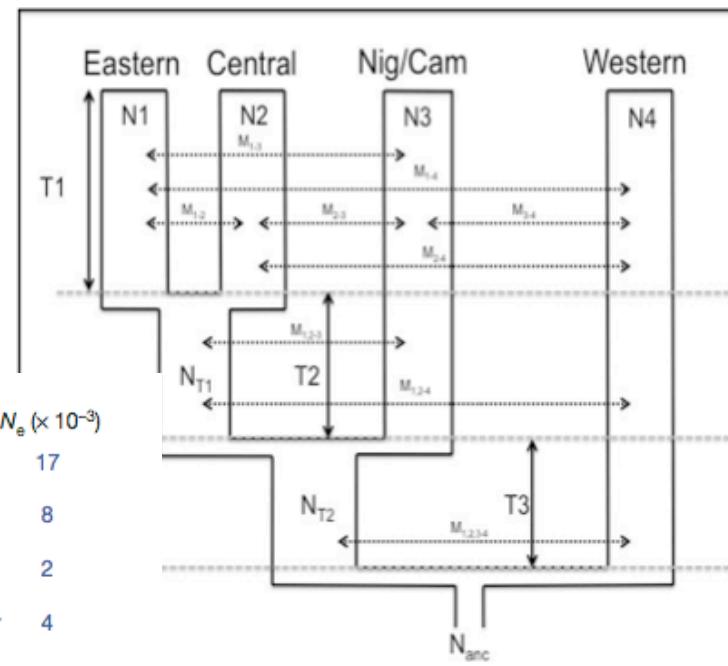
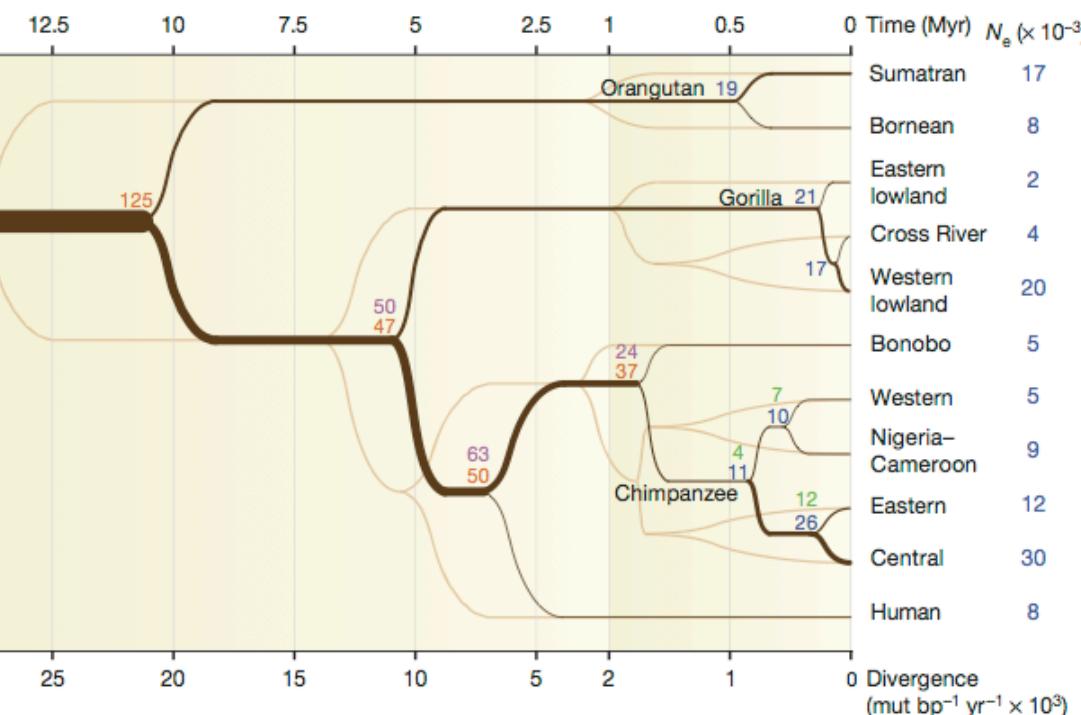


Accept?(if close to  
real data)

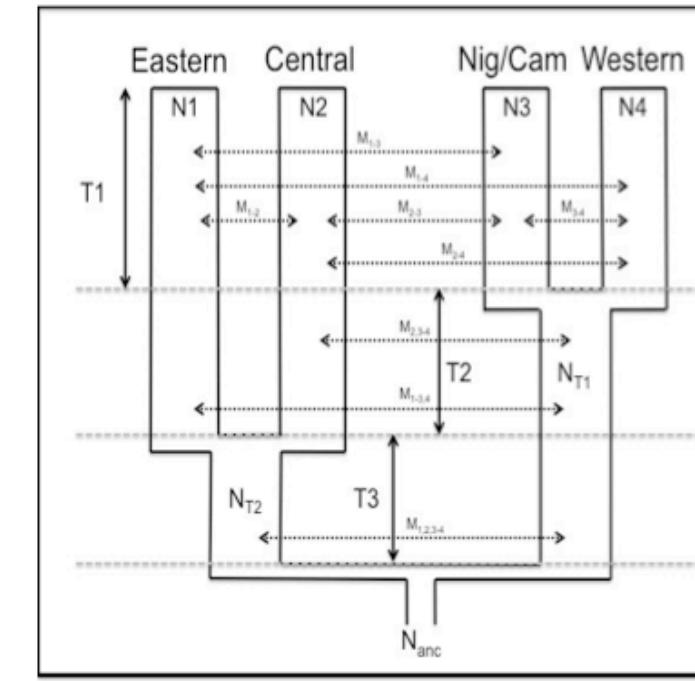
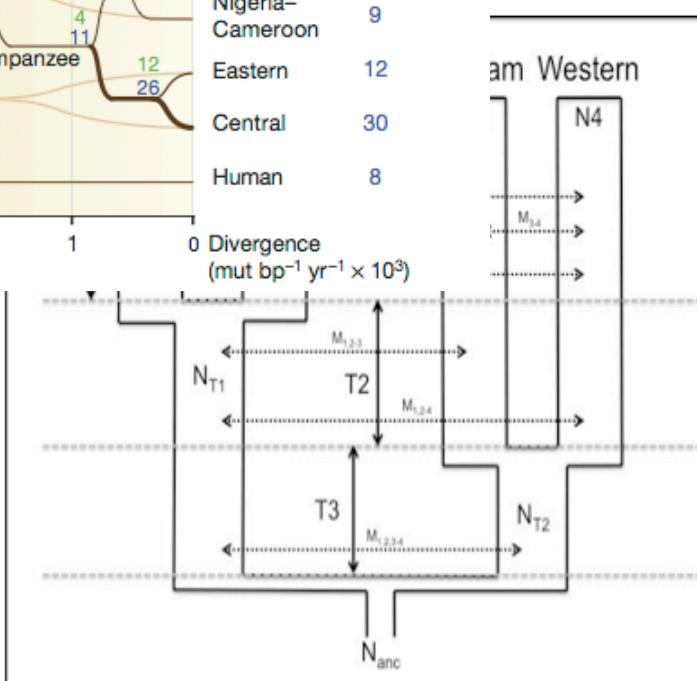
Posterior  
distribution

# ABC

dealing with models that have many many parameters



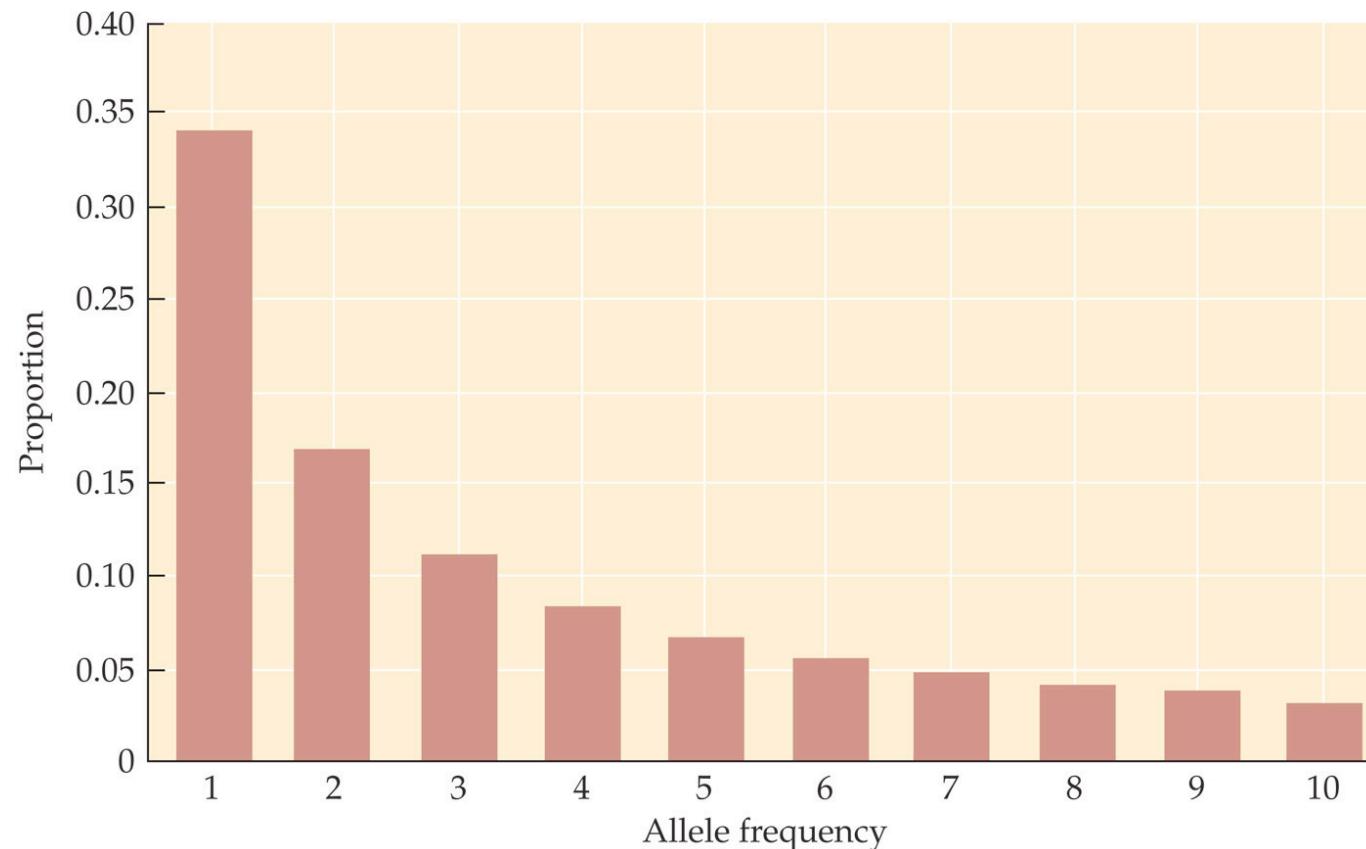
Model 3

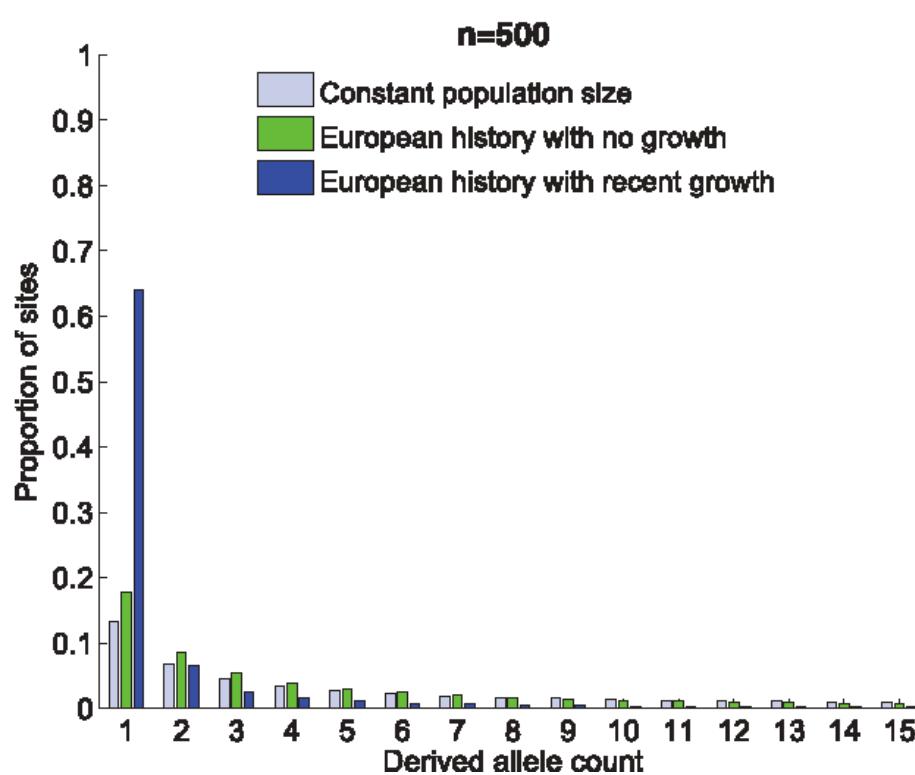
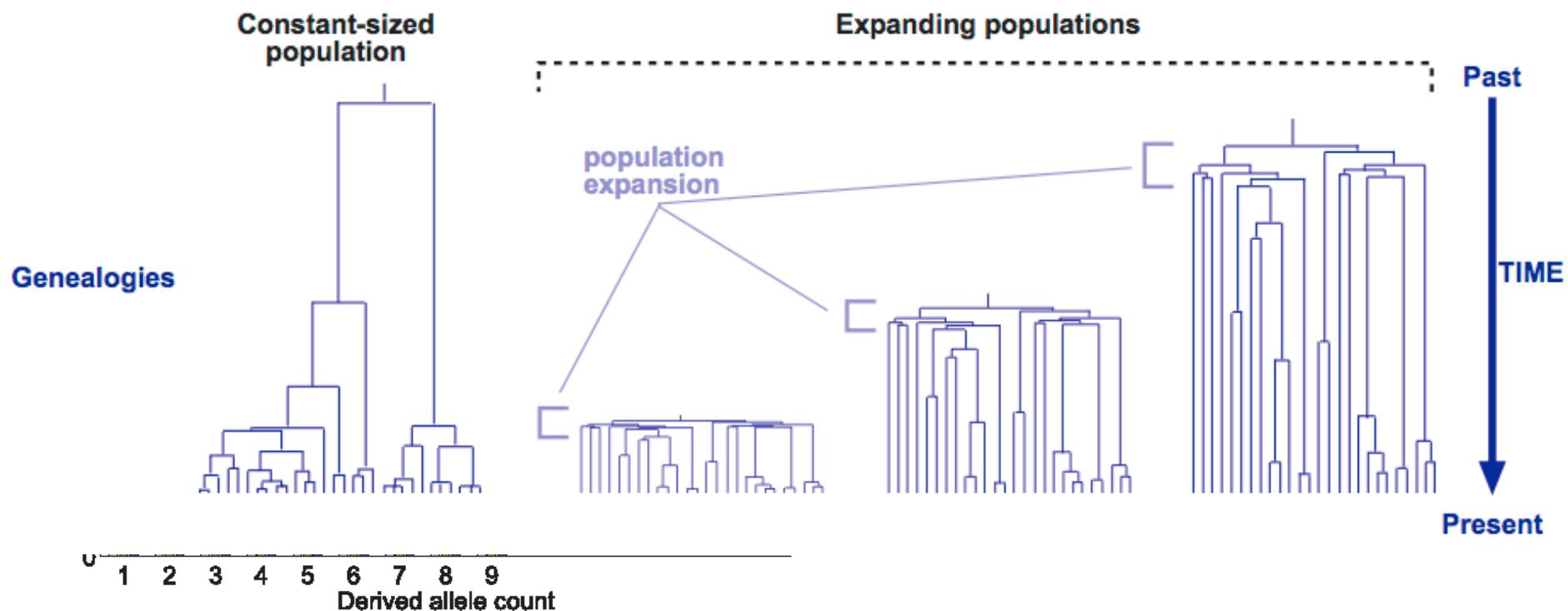


# high-throughput data analysis

illumina SNP data (explicitly incorporating 10,000s of different coalescent trees becomes impractical)

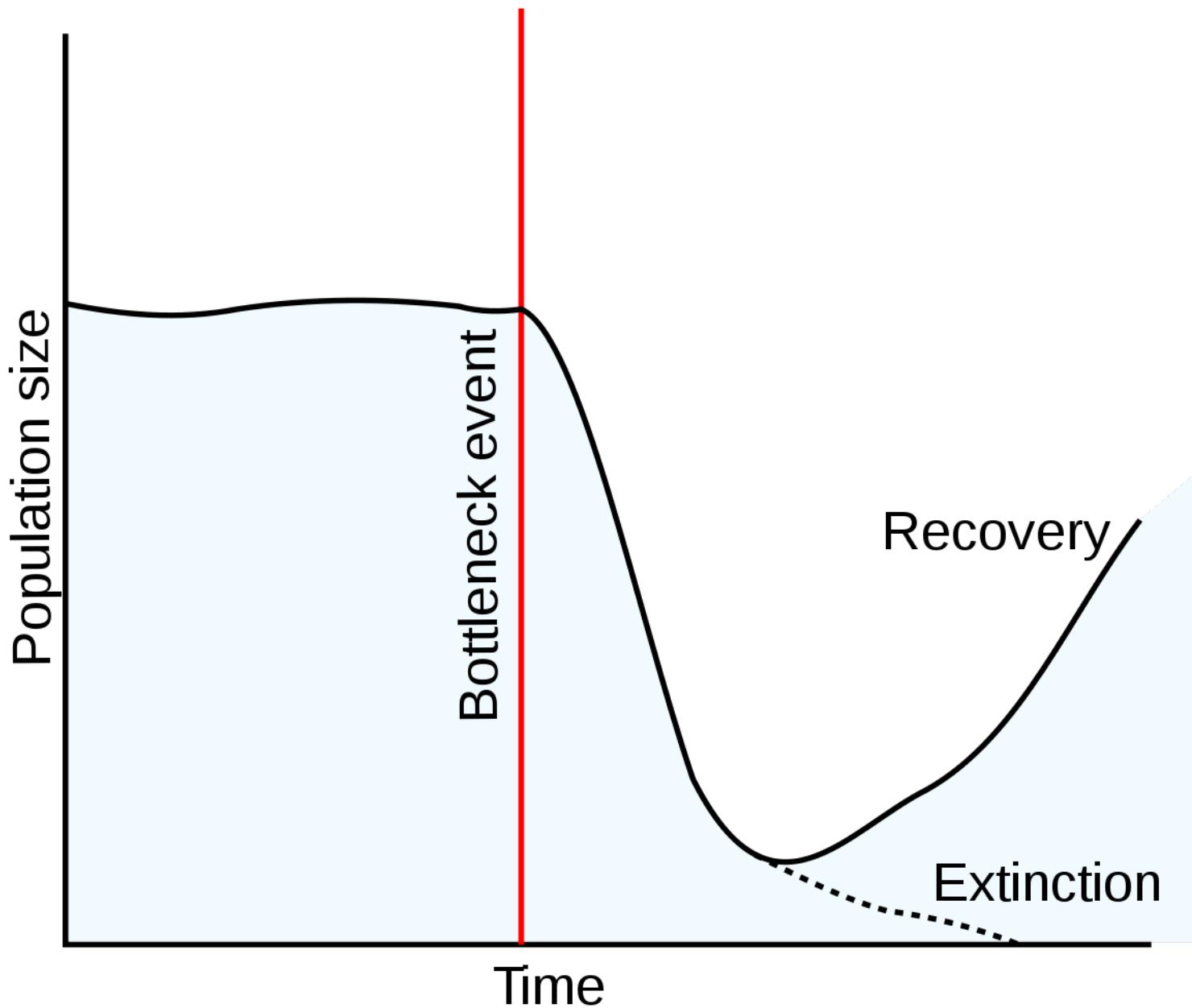
## SFS - Site Frequency Spectrum



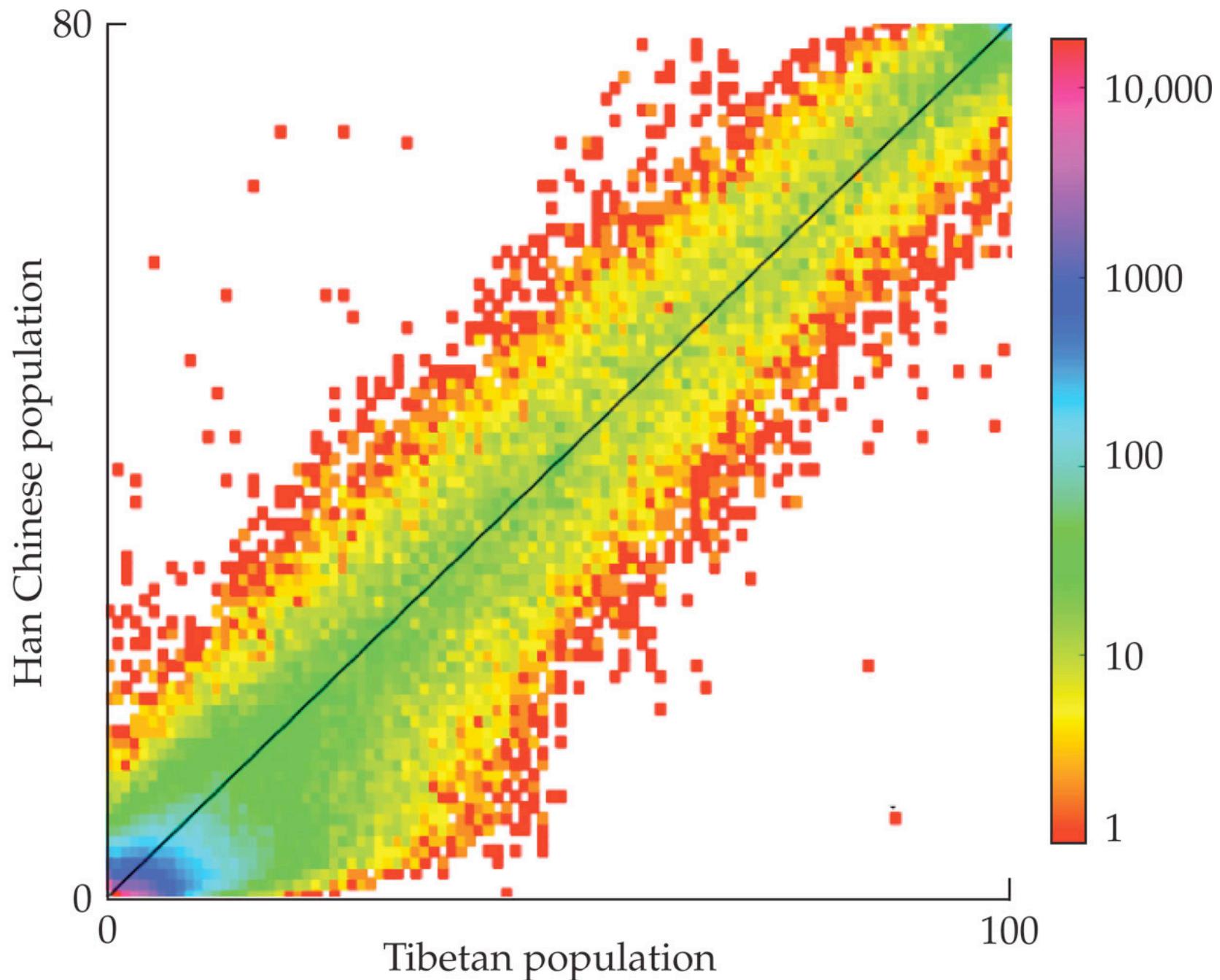


## SFS - Site Frequency Spectrum

Bottleneck?, expansion? both ? (it depends)



# Joint Site Frequency Spectrum



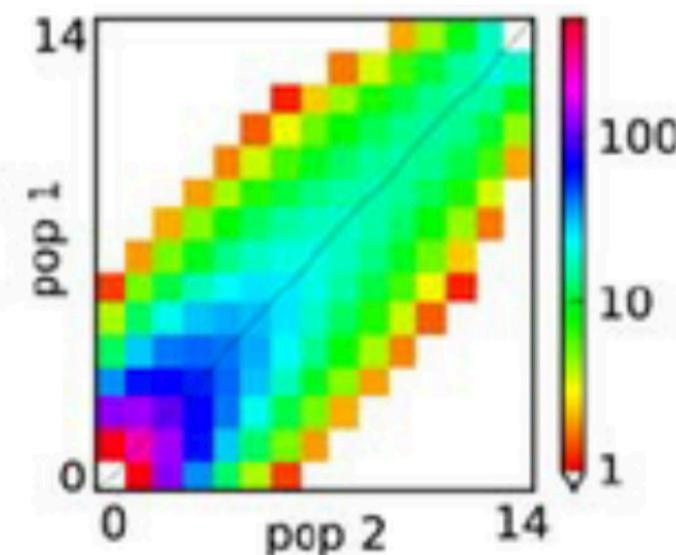
INTRODUCTION TO POPULATION GENETICS, Figure 5.14

© 2013 Sinauer Associates, Inc.

# The joint site frequency spectrum of 2 populations

---

- If an allele has frequency  $f_1$  in population 1 and frequency  $f_2$  in population 2, its joint frequency distribution between the two populations is  $(f_1, f_2)$
- The joint site frequency spectrum (of 2D SFS) is the 2-dimensional histogram of these joint site frequencies

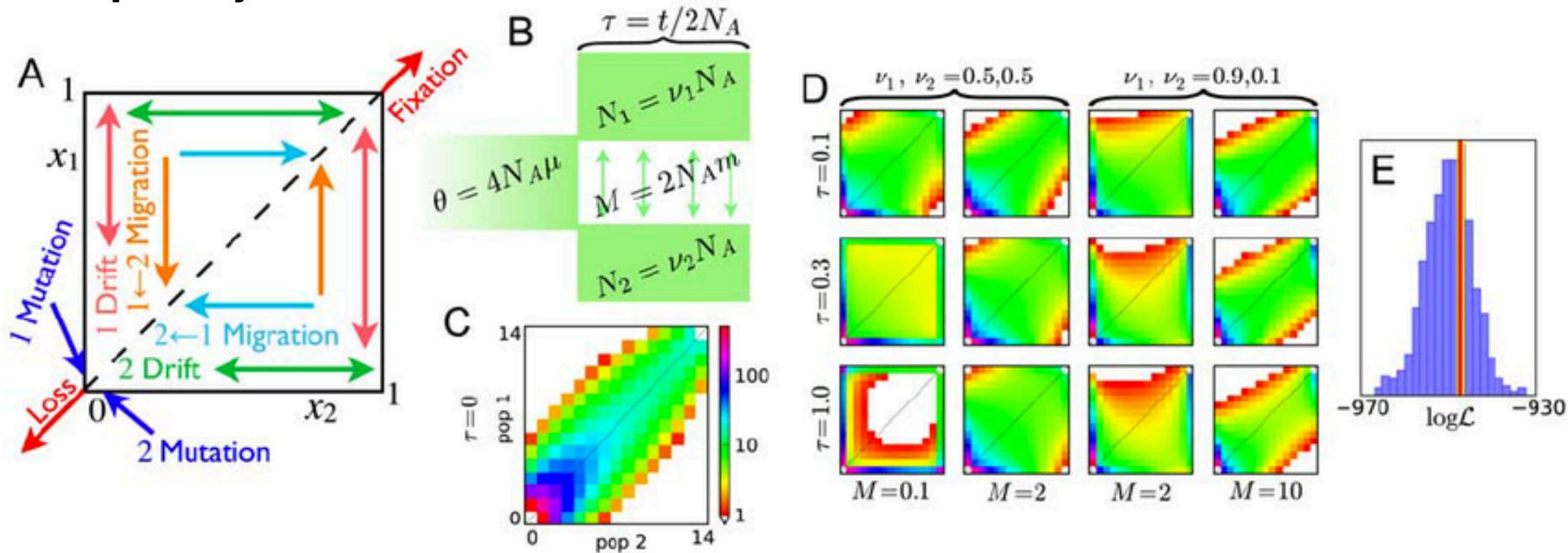


**Does it look like these populations diverged very long ago?**

**Gutenkunst, et al., 2009**

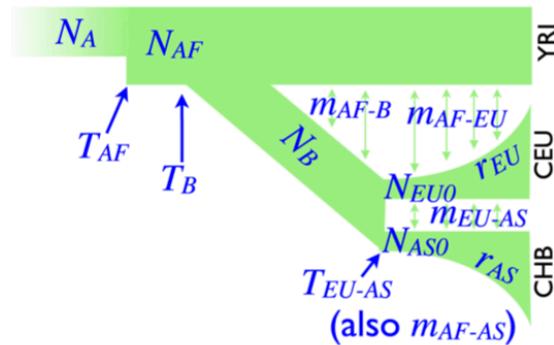
% K. Harris

# relationship btw history and shape of joint SFS

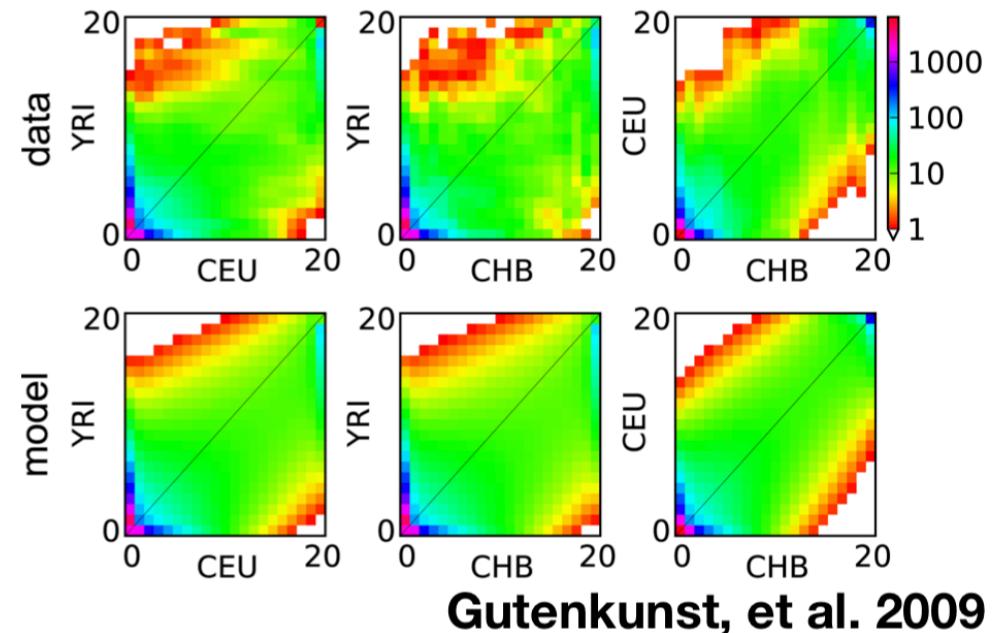


Since small populations diverge faster than large ones, must infer the divergence time jointly with the effective population sizes of the ancestral and present day populations

# An estimate of the joint demographic history of Europeans, Africans, and East Asians



Programs like “dadi” calculate an approximate likelihood of the observed site frequency spectrum given a demographic model and optimize the parameters of the model



**Gutenkunst, et al. 2009**

**YRI = Yorubans from Ibadan, Nigeria**

**CEU = Individuals of European Descent from Utah**

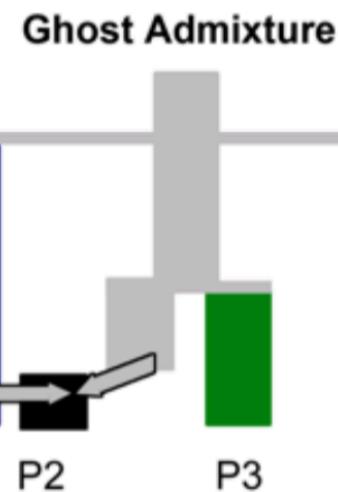
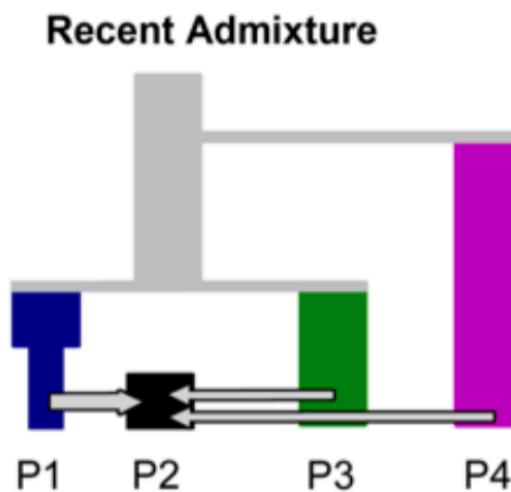
**CHB = Han Chinese from Beijing**

**Why do you think a sample of Yorubans from Ibadan, Nigeria are being used as a proxy for all Africans?**

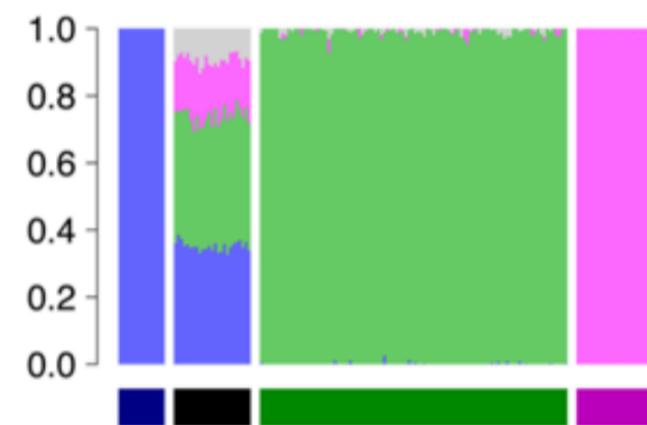
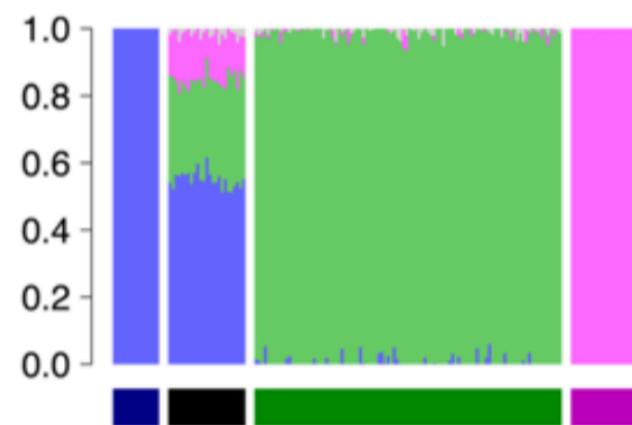
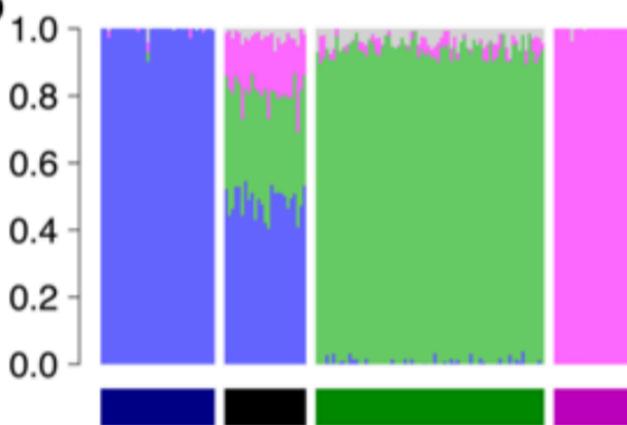
% K. Harris

Remember, STRUCTURE is agnostic to specific demographic histories. (i.e. you should use explicit demographic models for inference after assignment)

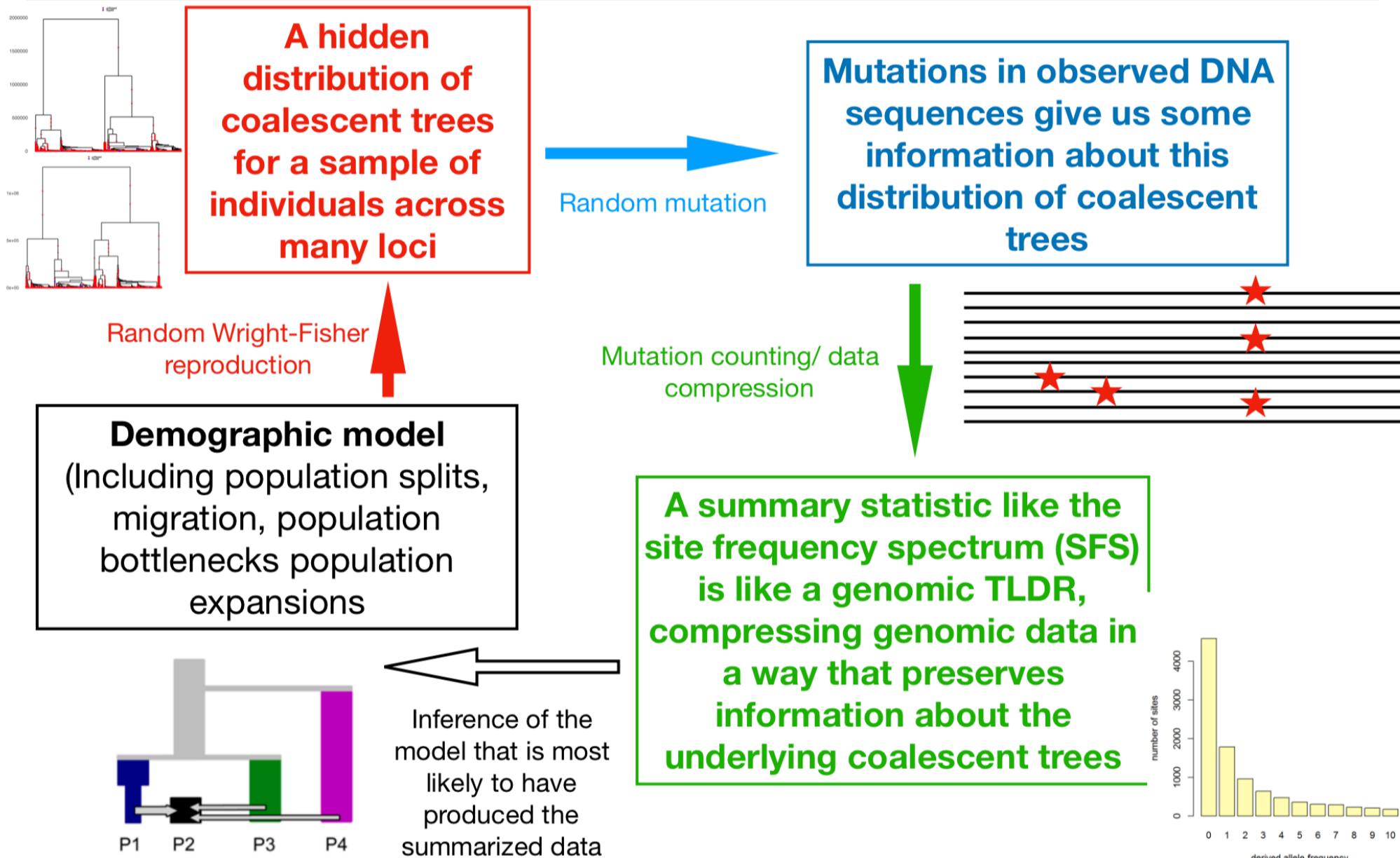
a



b



# weird circularity “work-flow”

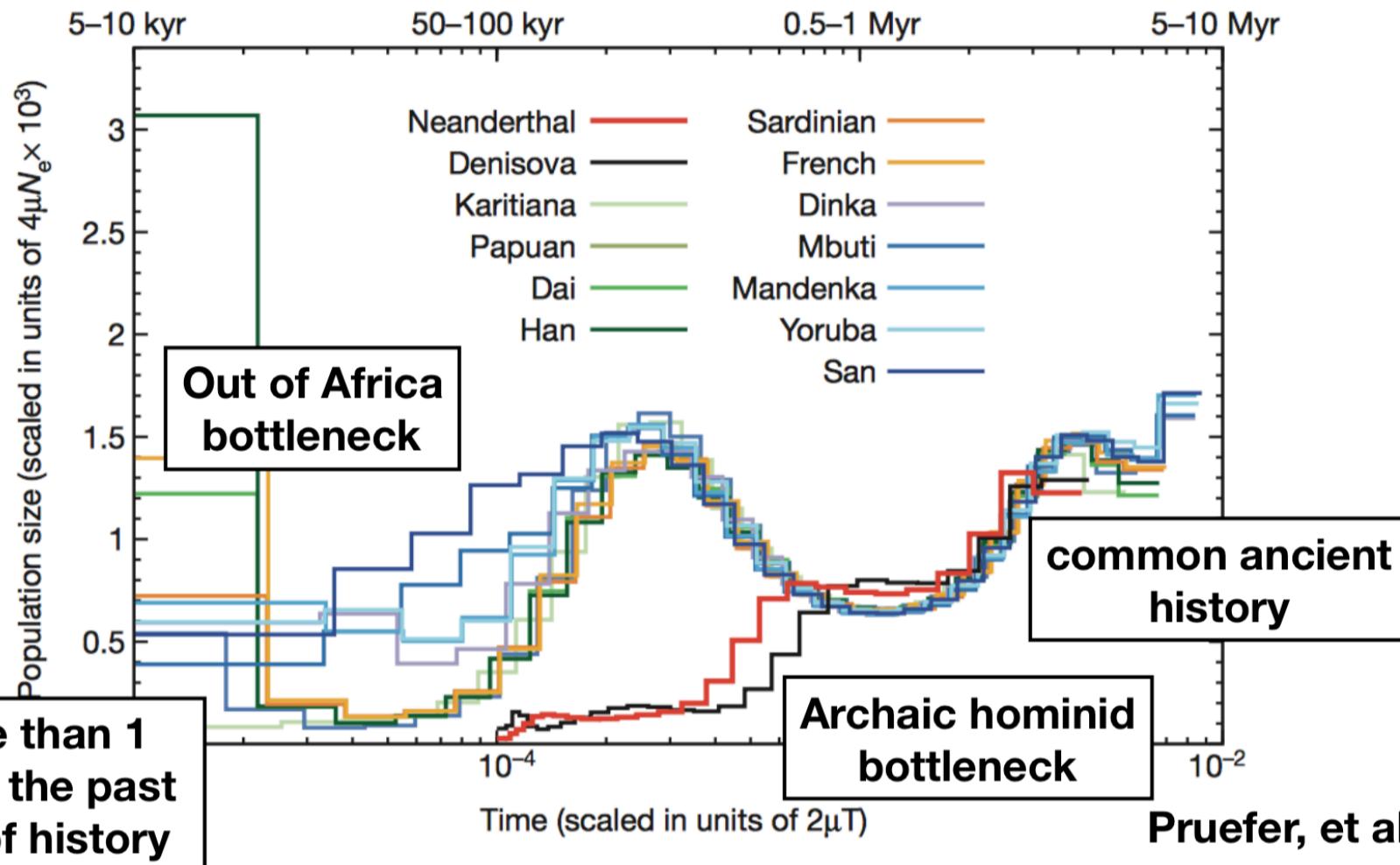


# Inference of human population history from individual whole-genome sequences

Heng Li ✉ & Richard Durbin ✉

% K. Harris

Nature **475**, 493–496 (28 July 2011) | Download Citation ↓



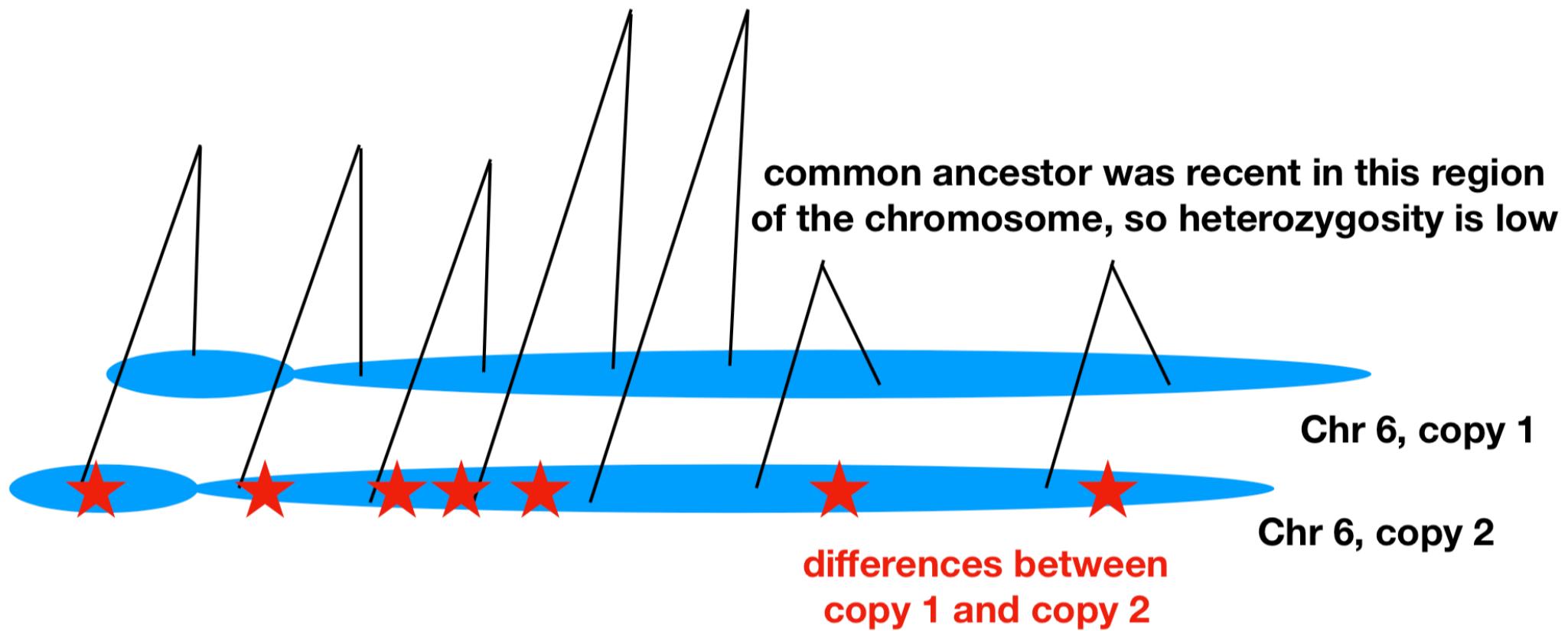
# The PSMC works by leveraging the variation of TMRCA along 1 diploid genome

---

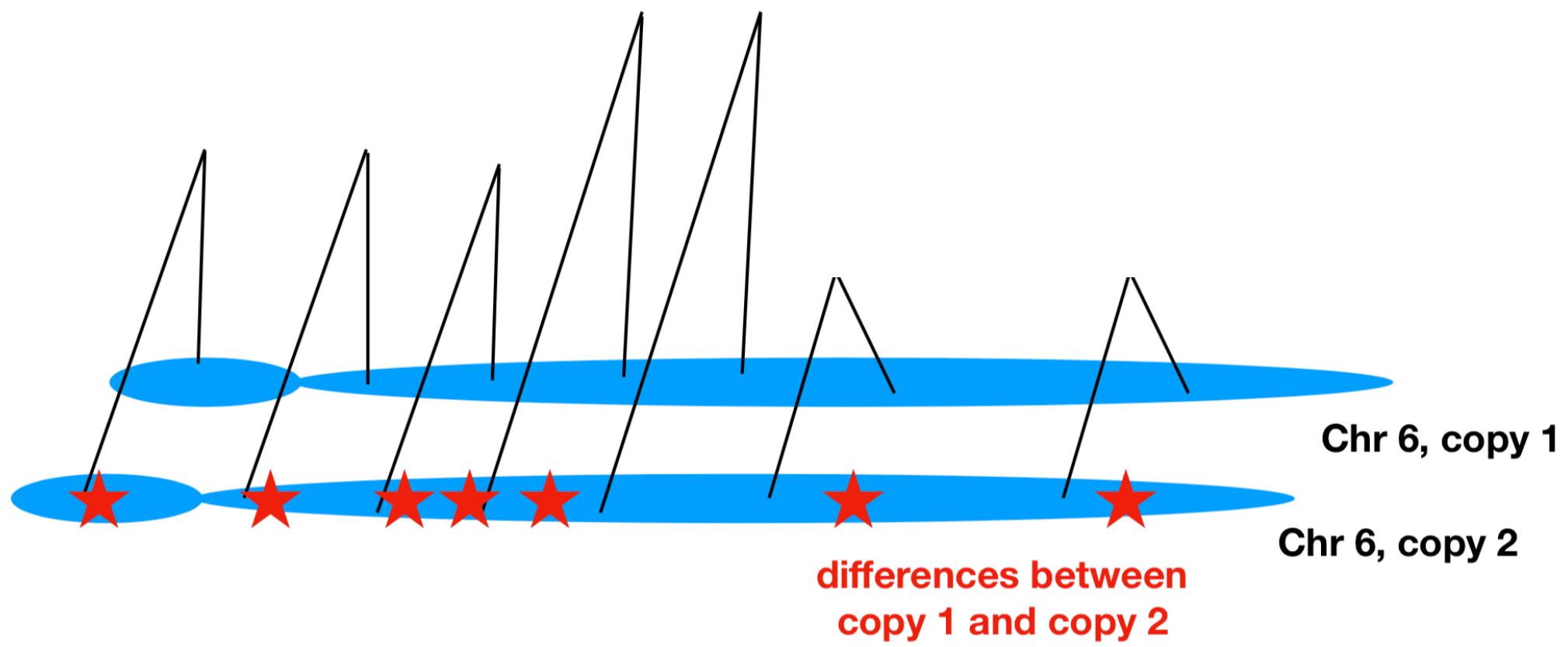
- The genealogical history of 1 genetic locus is just one of many histories that are possible given the population history
- Think about siblings: 1/4 of their genome coalesces 1 generation back, but the rest does not
- To accurately infer demographic history from genetic data, including changes in effective population size over time, information from many independent genetic loci is needed

# Variation of coalescent history along the genome

common ancestors was more ancient in this half of the chromosome, so heterozygosity is higher



HMM, a hidden Markov model run along the chromosome to infer the time periods in which each locus coalesces (the “hidden” states)



# Inference of human population history from one genome

---

- Data from many individuals is required to estimate demographic history from the site frequency spectrum

# Inference of human population history from one genome

---

- Data from many individuals is required to estimate demographic history from the site frequency spectrum
- But heterozygosity in 1 genome is enough to estimate the average effective population size, and it's actually also enough to infer changes in  $N_e$  over time!

# Inference of human population history from one genome

---

- Data from many individuals is required to estimate demographic history from the site frequency spectrum
- But heterozygosity in 1 genome is enough to estimate the average effective population size, and it's actually also enough to infer changes in  $N_e$  over time!
- Recall: if we go back in time to a distant past generation, everyone alive then either has no present-day descendants or is an ancestor of everyone alive today

# Inference of human population history from one genome

---

- Data from many individuals is required to estimate demographic history from the site frequency spectrum
- But heterozygosity in 1 genome is enough to estimate the average effective population size, and it's actually also enough to infer changes in  $N_e$  over time!
- Recall: if we go back in time to a distant past generation, everyone alive then either has no present-day descendants or is an ancestor of everyone alive today
- DNA from a single individual is enough to estimate the effective size of this distant past population

# Demographic change interacts with selection in multiple ways

---

- When the population size grows, genetic drift slows down and selection becomes more efficient

# Demographic change interacts with selection in multiple ways

---

- When the population size grows, genetic drift slows down and selection becomes more efficient
- A population bottleneck speeds up genetic drift and makes natural selection get less efficient

# Demographic change interacts with selection in multiple ways

---

- When the population size grows, genetic drift slows down and selection becomes more efficient
- A population bottleneck speeds up genetic drift and makes natural selection get less efficient
- Population size changes shift the site frequency spectrum away from the shape  $SFS(n) \sim 1/n$  that is expected under neutrality

# Demographic change interacts with selection in multiple ways

---

- When the population size grows, genetic drift slows down and selection becomes more efficient
- A population bottleneck speeds up genetic drift and makes natural selection get less efficient
- Population size changes shift the site frequency spectrum away from the shape  $SFS(n) \sim 1/n$  that is expected under neutrality
- Natural selection also shifts the SFS away from the equilibrium  $1/n$  shape

# Demographic change interacts with selection in multiple ways

---

- When the population size grows, genetic drift slows down and selection becomes more efficient
- A population bottleneck speeds up genetic drift and makes natural selection get less efficient
- Population size changes shift the site frequency spectrum away from the shape  $SFS(n) \sim 1/n$  that is expected under neutrality
- Natural selection also shifts the SFS away from the equilibrium  $1/n$  shape
- The SFS shape can be used to test for selection, but only after correcting for nonequilibrium demography

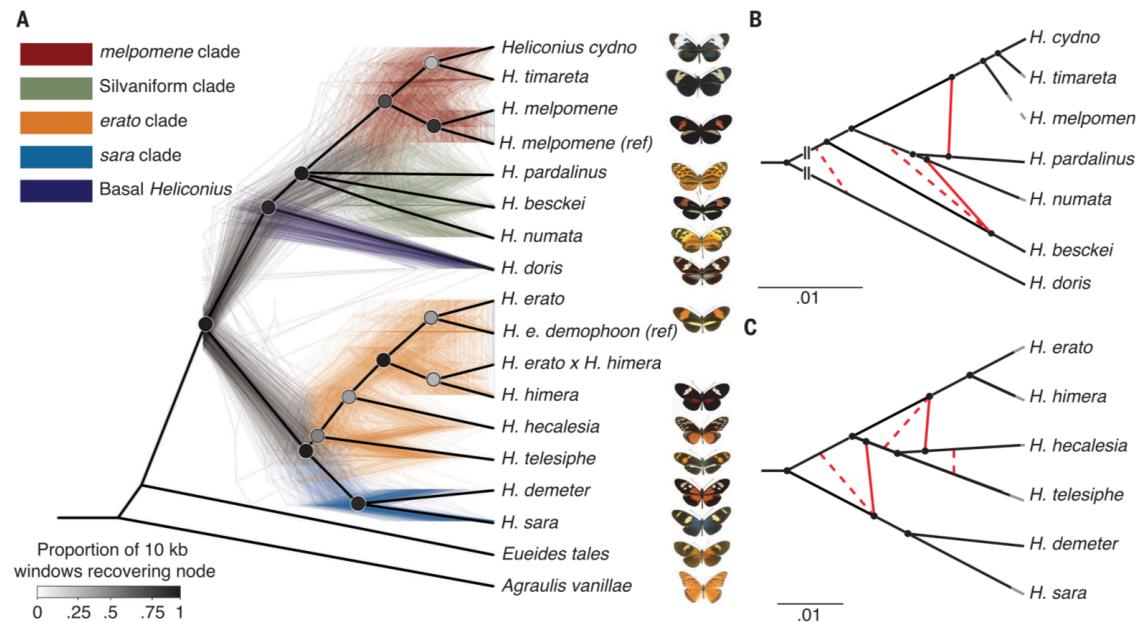
# Next Week:

1. Simulation Assignment
2. Discuss paper

## BUTTERFLY GENOMICS

# Genomic architecture and introgression shape a butterfly radiation

Nathaniel B. Edelman<sup>1\*</sup>, Paul B. Frandsen<sup>2,3</sup>, Michael Miyagi<sup>1</sup>, Bernardo Clavijo<sup>4</sup>, John Davey<sup>5,20</sup>, Rebecca B. Dikow<sup>3</sup>, Gonzalo García-Accinelli<sup>4</sup>, Steven M. Van Belleghem<sup>6</sup>, Nick Patterson<sup>7,8</sup>, Daniel E. Neafsey<sup>8,9</sup>, Richard Challis<sup>10</sup>, Sujai Kumar<sup>11</sup>, Gilson R. P. Moreira<sup>12</sup>, Camilo Salazar<sup>13</sup>, Mathieu Chouteau<sup>14</sup>, Brian A. Counterman<sup>15</sup>, Riccardo Papa<sup>6,16</sup>, Mark Blaxter<sup>10</sup>, Robert D. Reed<sup>17</sup>, Kanchon K. Dasmahapatra<sup>5</sup>, Marcus Kronforst<sup>18</sup>, Mathieu Joron<sup>19</sup>, Chris D. Jiggins<sup>20</sup>, W. Owen McMillan<sup>21</sup>, Federica Di Palma<sup>4</sup>, Andrew J. Blumberg<sup>22</sup>, John Wakeley<sup>1</sup>, David Jaffe<sup>8,23</sup>, James Mallet<sup>1\*</sup>



# Simulation Assignment: comparing single and double population models

Use msPrime, SLiM or PipeMaster

1 locus, 1000 base pairs  
 $\mu=1e-7$  per base pair

Simulate a 2 population model  
 $\tau = 5,000$  generations  
 $N_1 = 10,000, N_2 = 10,000, N_A = 10,000$   
 $n_1 = 10, n_2 = 10$   
no migration

1. Simulate once and calculate an observed  $\pi$  (ignoring which of the 2 populations the 20 individuals come from).
2. Simulate a single population model with that observed  $\pi$  value 1,000 times (pretend you know the mutation rate exactly) and plot the distribution of simulated  $\pi$  values
3. indicate the 'true'  $4N\mu$  ( $\pi$ ) on the plot. How off is it?
4. Simulate the 2 pop model 1,000 times and choose N values that you think will give you a  $\pi$  distribution that matches the one from the single pop model.
5. plot results of 3 and 4