

# Species Tree Estimation

Laura Kubatko  
Departments of Statistics and  
Evolution, Ecology, and Organismal Biology  
The Ohio State University

kubatko.2@osu.edu  
twitter: Laura\_Kubatko

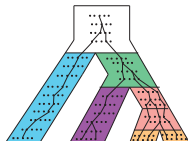
January 31, 2017

## Relationship between population genetics and phylogenetics

- **Population genetics:** Study of genetic variation within a population
- **Phylogenetics:** Use genetic variation between taxa (species, populations) to infer evolutionary relationships
- Previously:
  - ▶ Each taxon is represented by a single sequence – “exemplar sampling”
  - ▶ We have data for a single gene and wish to estimate the evolutionary history for that gene (the **gene tree** or **gene phylogeny**)

## Relationship between population genetics and phylogenetics

- Given current technology, we can do much more:
  - ▶ Sample many individuals within each taxon (species, population, etc.)
  - ▶ Sequence many genes for all individuals
- Need models at two levels:
  - ▶ Model what happens within each population  
[population genetics – coalescent model]
  - ▶ Link each within-population model on a phylogeny  
[phylogenetics]



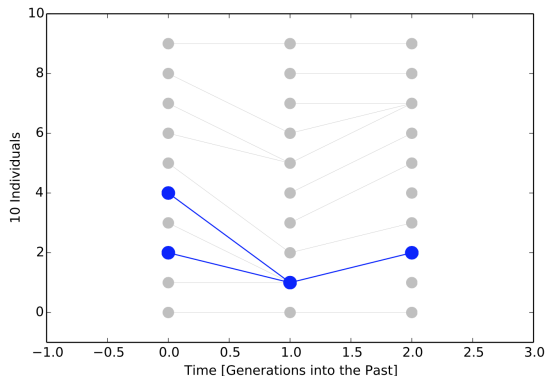
- Assumptions:

- ▶ Population of  $2N$  gene copies
- ▶ Discrete, non-overlapping generations of equal size
- ▶ Parents of next generation of  $2N$  genes are picked randomly with replacement from preceding generation (genetic differences have no fitness consequences)
- ▶ Probability of a specific parent for a gene in the next generation is  $\frac{1}{2N}$



## Wright-Fisher model

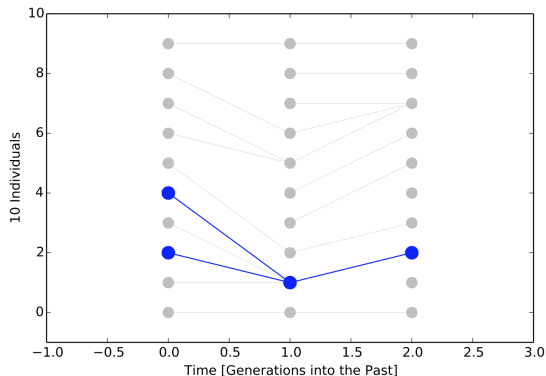
- Consider a population of  $2N$  gene copies at time  $t$
- **Question:** What is the probability that two randomly chosen individuals share a common ancestor in the previous generation?



Figures from PopVizard, by Peter Beerli

## Wright-Fisher model

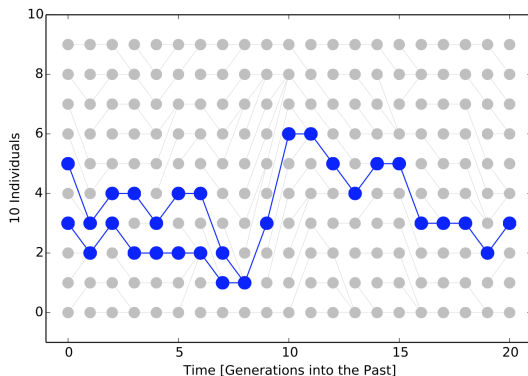
- Consider a population of  $2N$  gene copies at time  $t$
- Question:** What is the probability that two randomly chosen individuals **DO NOT** share a common ancestor in the previous generation?



Figures from PopVizard, by Peter Beerli

## Wright-Fisher model

- Consider a population of  $2N$  gene copies at time  $t$
- Question:** How many generations do we need to wait until two randomly chosen individuals share a common ancestor?



Figures from PopVizard, by Peter Beerli

- Consider a population of  $2N$  gene copies at time  $t$
- **Question:** How many generations do we need to wait until two randomly chosen individuals share a common ancestor?
  - ▶ The number of generations,  $T$ , until two individuals share a common ancestor follows a **geometric distribution** with parameter  $\frac{1}{2N}$ , e.g.,

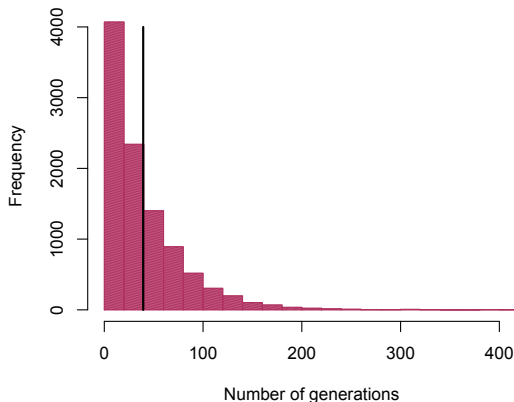
$$P(T = \tau) = \left(1 - \frac{1}{2N}\right)^{\tau-1} \left(\frac{1}{2N}\right)$$

- ▶ The expected number of generations until coalescence is

$$E(T) = 2N$$

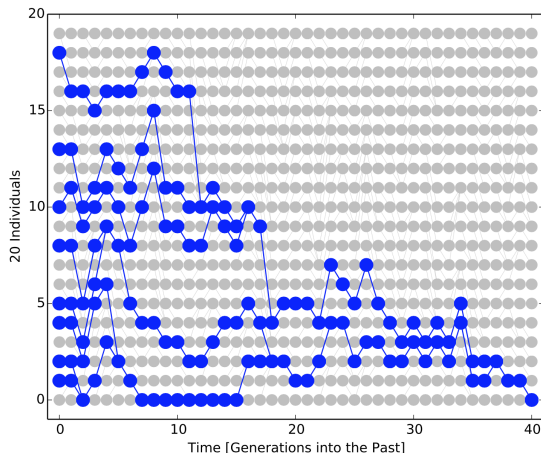
## Wright-Fisher model

- Distribution of time to coalescence for two randomly sampled gene copies in a population of size  $2N = 40$
- Observed mean (black line) = 39.40



## Wright-Fisher model

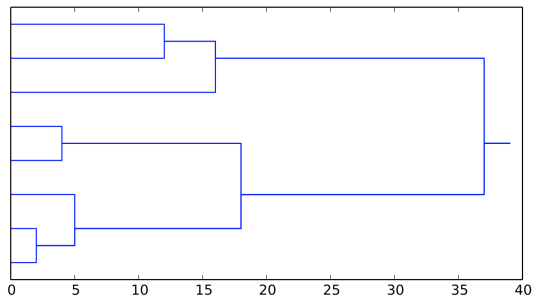
- Consider a population of  $2N$  gene copies at time  $t$
- Question:** What about more than two gene copies?



Figures from PopVizard, by Peter Beerli

## Wright-Fisher model

- Consider a population of  $2N$  gene copies at time  $t$
- **Question:** What about more than two gene copies?



Figures from PopVizard, by Peter Beerli

- Consider a population of  $2N$  gene copies at time  $t$
- **Question:** What about more than two gene copies?
  - ▶ Note from the previous that time to coalescence seems to vary with sample size
  - ▶ Specifically, the probability that two gene copies in a sample of  $k$  coalesce in the previous generation is

$$\binom{k}{2} \frac{1}{2N}$$

- ▶ The distribution of the time to the first coalescent event is **geometric** with parameter  $\binom{k}{2} \frac{1}{2N}$

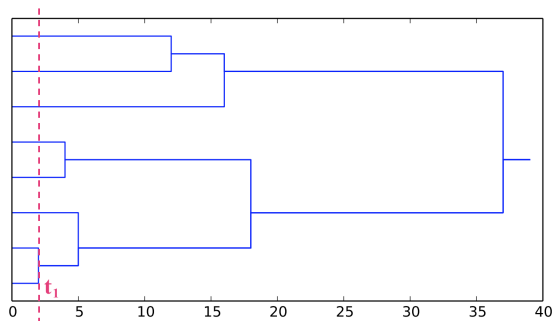


## Kingman's coalescent

- Cumbersome to work in this “discrete” setting where we think of things generation by generation
- **Kingman's approximation**: consider continuous time and a sample of  $k$  lineages. Then, the time back into the past until two lineages coalesce,  $U$ , is exponentially distributed with rate  $\binom{k}{2} \frac{1}{2N}$ 
  - ▶ The probability density function is  $g(u) = \binom{k}{2} \frac{1}{2N} e^{-\binom{k}{2} \frac{u}{2N}}$ , for  $u > 0$
  - ▶ The mean is  $\frac{4N}{k(k-1)}$

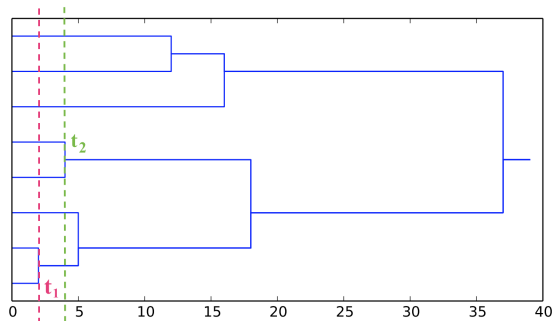


## Computing the probability of a population tree under the coalescent



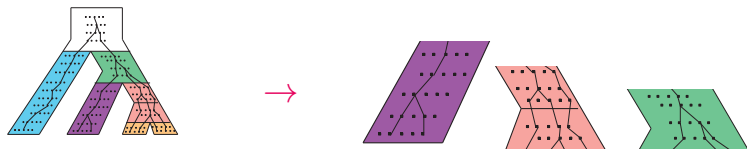
$$P(G) = \left( \frac{1}{2N} e^{-\frac{8(7)}{4N} t_1} \right) \times \dots$$

## Computing the probability of a population tree under the coalescent



$$P(G) = \left( \frac{1}{2N} e^{-\frac{8(7)}{4N} t_1} \right) \left( \frac{1}{2N} e^{-\frac{7(6)}{4N} t_2} \right) \times \dots$$

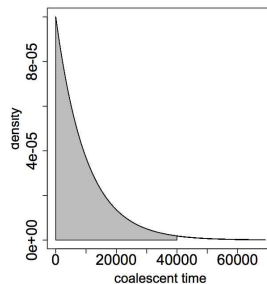
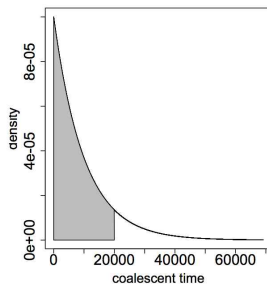
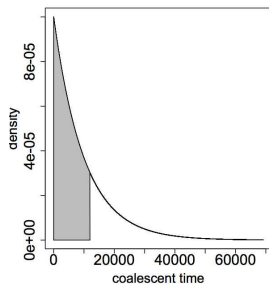
## Fitting population trees into a phylogeny



- Focus on just one **speciation interval** and a sample of  $k = 2$  lineages.
- Then,  $\binom{k}{2} = 1$  and we have an exponential distribution with rate  $\frac{1}{2N}$  and mean  $2N$ .
- Suppose  $N = 5,000$ . Let's find the probability that the two lineages coalesce in an interval of a particular length.

## Fitting population trees into a phylogeny

- $N = 5,000$  and consider the times: 12,000, 20,000 and 40,000 generations

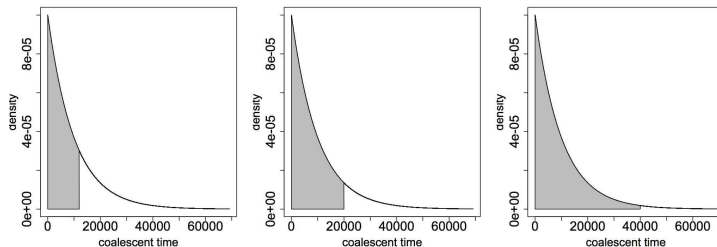


## Fitting population trees into a phylogeny

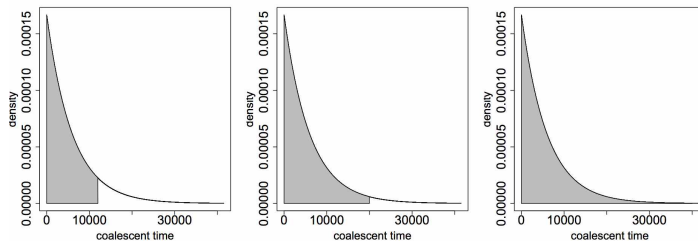
- What happens if we change the population size,  $N$ ?
- Recall that we have an exponential distribution with rate  $\frac{1}{2N}$  and mean  $2N$ .
- Now suppose  $N = 3,000$  and look at the same speciation interval lengths.

## Fitting population trees into a phylogeny

- $N = 5,000$

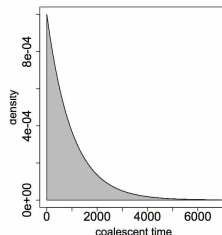
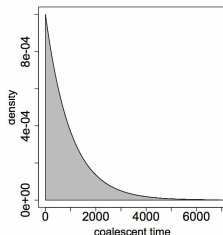
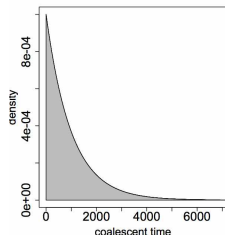


- $N = 3,000$



## Fitting population trees into a phylogeny

- What about the effect of sample size,  $k$ ?
- Consider  $N = 5,000$  again, but now use  $k = 5$ .
  - ▶ Rate is  $\binom{5}{2} \frac{1}{2N} = \frac{10}{2N}$  (was  $\frac{1}{2N}$ )
  - ▶ Mean is  $\frac{4N}{k(k-1)} = \frac{2N}{10}$  (was  $2N$ )



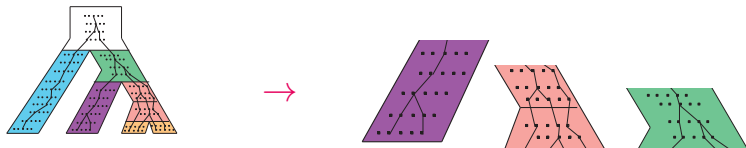


## Fitting population trees into a phylogeny

- Define a common unit of time: **coalescent unit**,  $t = \frac{u}{2N}$
- Examples:
  - ▶  $k = 2$  — exponential distribution with rate 1 and mean 1
  - ▶  $k = 5$  — exponential distribution with rate 10 and mean 0.1
- $t$  “large” is now relative to population size, but the trends are the same:
  - ▶ Longer times lead to a higher probability of coalescence having occurred.
  - ▶ Coalescent events happen more quickly when the population size is smaller.
  - ▶ Coalescent events happen more quickly when the sample size is larger.
- What does this mean for species trees estimation ???

## Fitting population trees into a phylogeny

- Recall our goal to integrate the population process with the phylogeny:



- Can use our previous results to get the following:
  - The probability that  $u$  lineages coalesce into  $v$  lineages in time  $t$  is given by (Tavare, 1984; Watterson, 1984; Takahata and Nei, 1985; Rosenberg, 2002)

$$P_{uv}(t) = \sum_{j=v}^u e^{-j(j-1)t/2} \frac{(2j-1)(-1)^{j-v}}{v!(j-v)!(v+j-1)} \prod_{y=0}^{j-1} \frac{(v+y)(u-y)}{u+y}$$

## Fitting population trees into a phylogeny

- When  $u$  and  $v$  are small, these are easy to compute. For example,

$$\begin{aligned}P_{21}(t) &= \text{probability that 2 lineages coalesce to 1 lineage in time } t \\&= \text{probability of 1 coalescent event in time } t \text{ when } k = 2 \\&= P(T \leq t), \text{ where } T \sim \text{Exp}(\mu = 1) \\&= \int_0^t e^{-x} dx = 1 - e^{-t}\end{aligned}$$

[Note: this is the formula for the gray area in the graphs]

- Similarly,

$$\begin{aligned}P_{22}(t) &= \text{prob. of no coalescence in time } t \text{ for 2 lineages} \\&= P(T > t) \\&= \int_t^\infty e^{-x} dx = e^{-t}\end{aligned}$$

### • Assumptions:

- ▶ Events that occur in one population are independent of what happens in other populations within the phylogeny.
- ▶ More specifically, given the number of lineages entering and leaving a population, coalescent events within populations are independent of other populations.
- ▶ It is also important to recall an assumption we “inherit” from our population genetics model: all pairs of lineages are equally likely to coalesce within a population.
- ▶ No gene flow occurs following speciation.
- ▶ No other evolutionary processes (e.g., horizontal gene flow, duplication, . . .) have led to incongruence between gene trees and the species tree.

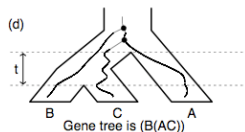
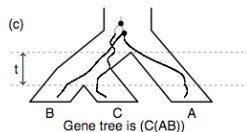
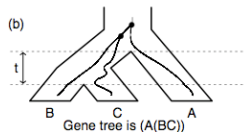
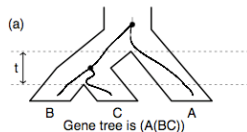
## Putting it together ... the coalescent model along a species tree

- When talking about gene tree distributions, there are two cases of interest:
  - ▶ The gene tree topology distribution
  - ▶ The joint distribution of topologies and branch lengths
- Start with the simple case of 3 species with 1 lineage sampled in each and look at the **gene tree topology distribution**

## Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

Example of gene tree probability computation:

(a)  $\text{Prob} = 1 - e^{-t}$ ; (b), (c), (d)  $\text{Prob} = \frac{1}{3}e^{-t}$



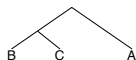
## Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

- Thus, we have the following probabilities:
  - ▶ Gene tree (A,(B,C)):  $\text{prob} = 1 - e^{-t} + \frac{1}{3}e^{-t} = 1 - \frac{2}{3}e^{-t}$
  - ▶ Gene tree (B,(A,C)):  $\text{prob} = \frac{1}{3}e^{-t}$
  - ▶ Gene tree (C,(A,B)):  $\text{prob} = \frac{1}{3}e^{-t}$
- Note: There are two ways to get the first gene tree. We call these **histories**.
- The probability associated with a gene tree topology will be the sum over all histories that have that topology.

## Example: Computation of Gene Tree Topology Probabilities for the 3-taxon Case

- What are these probabilities like as a function of  $t$ , the length of time between speciation events?

(b)



$$\text{prob} = 1 - \exp(-t)$$



$$\text{prob} = (1/3)\exp(-t)$$

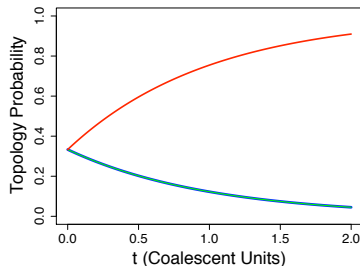


$$\text{prob} = (1/3)\exp(-t)$$



$$\text{prob} = (1/3)\exp(-t)$$

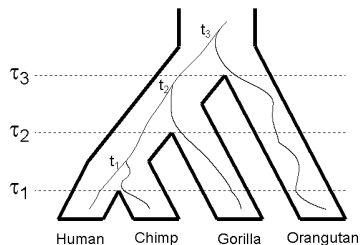
(c)





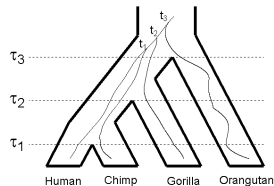
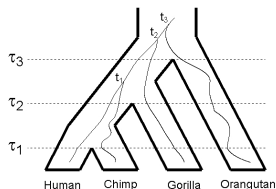
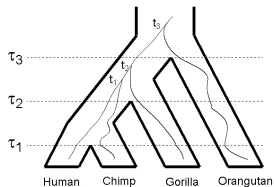
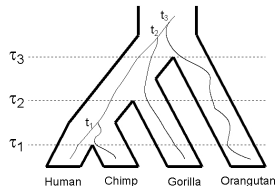
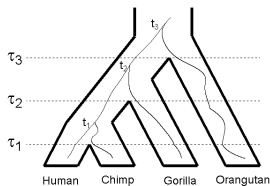
## Example: a slightly larger case

- Consider 4 taxa – the human-chimp-gorilla problem



## Coalescent histories for the 4-taxon example

- There are 5 possible histories for this example:



## Enumerating Histories

TABLE 3. The number of valid coalescent histories when the gene tree and species tree have the same topology. The number of histories is also the number of terms in the outer sum in equation (12).

Taxa	Number of histories		Number of topologies
	Asymmetric trees	Symmetric trees	
4	5	4	15
5	14	10	105
6	42	25	945
7	132	65	10,395
8	429	169	135,135
9	1430	481	2,027,025
10	4862	1369	34,459,425
12	58,786	11,236	13,749,310,575
16	9,694,845	1,020,100	$6.190 \times 10^{15}$
20	1,767,263,190	100,360,324	$8.201 \times 10^{21}$

Degnan and Salter, *Evolution*, 2005

- In the general case, we have the following:

The probability of a gene tree  $g$  gives the species tree  $\mathcal{S}$  is given by

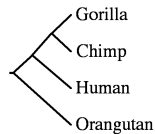
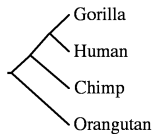
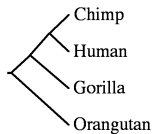
$$P\{G = g|\mathcal{S}\} = \sum_{\text{histories}} P\{G = g, \text{history}|\mathcal{S}\}$$

- Implemented in the software COAL (Degnan and Salter, *Evolution*, 2005)
- A more efficient method has been proposed (Wu, *Evolution*, 2012)

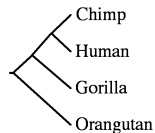
## Applications of the topology distribution - example 1

- Motivation: Paper by Ebersberger et al. 2007. *Mol. Biol. Evol.* 24:2266-2276
- Examined 23,210 distinct alignments for 5 primate taxa: Human, Chimp, Gorilla, Orangutan, Rhesus
- Looked at distribution of gene trees among these taxa - observed strongly supported incongruence only among the Human-Chimp-Gorilla clade.

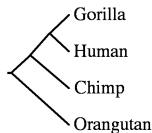
## Applications of the topology distribution - example 1



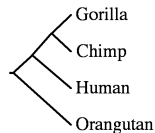
## Applications of the topology distribution - example 1



76.6%



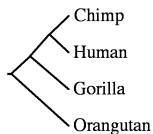
11.4%



11.5%

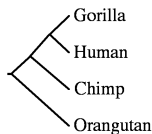
Observed proportions of each  
gene tree among ML phylogenies

## Applications of the topology distribution - example 1



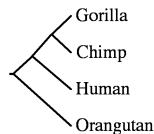
76.6%

79.1%



11.4%

9.9%



11.5%

9.9%

Observed proportions of each gene tree among ML phylogenies

Predicted proportions using parameters from Rannala & Yang, 2003.

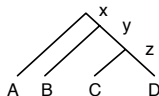


## Applications of the topology distribution - example 2

- In the previous example, one topology is clear preferred
- Must the distribution always look this way?
- Examine the entire distribution when the number of taxa is small

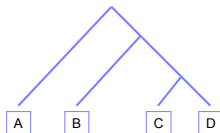
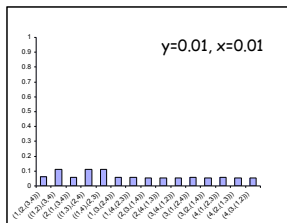
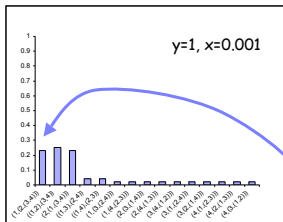
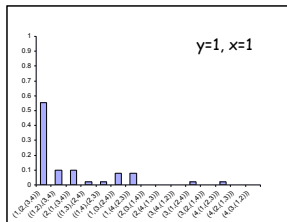
## Applications of the topology distribution - example 2

- Consider 4 taxa: A, B, C, and D
- Species tree:

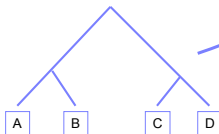
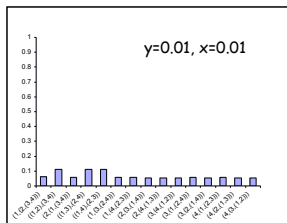
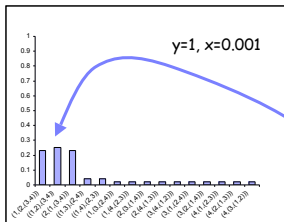
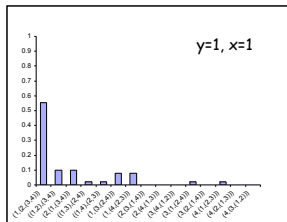


- Look at probabilities of all 15 tree topologies for values of  $x$ ,  $y$ , and  $z$

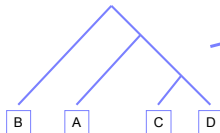
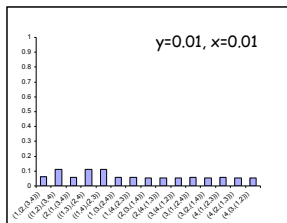
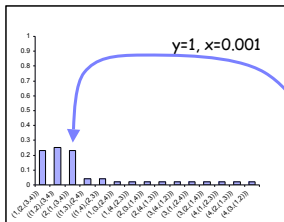
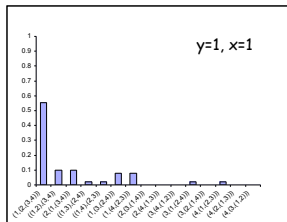
## Applications of the topology distribution - example 2



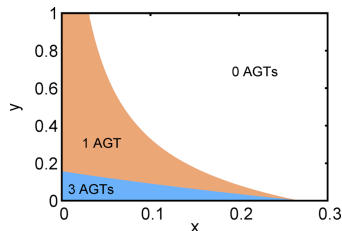
## Applications of the topology distribution - example 2



## Applications of the topology distribution - example 2



## Applications of the topology distribution - example 2



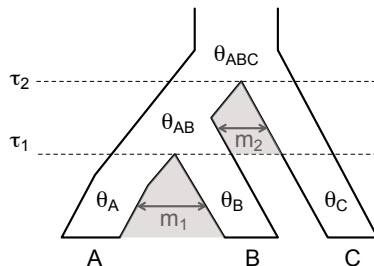
- The existence of **anomalous gene trees** has implications for the inference of species trees

Degnan and Rosenberg, *PLoS Genetics*, 2006

Rosenberg and Tao, *Systematic Biology*, 2008

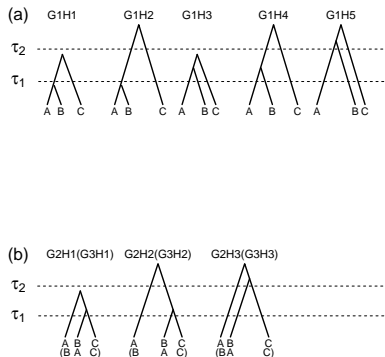
## What about gene flow?

**Question:** What happens to gene tree topology probabilities under a model with gene flow?



## What about gene flow?

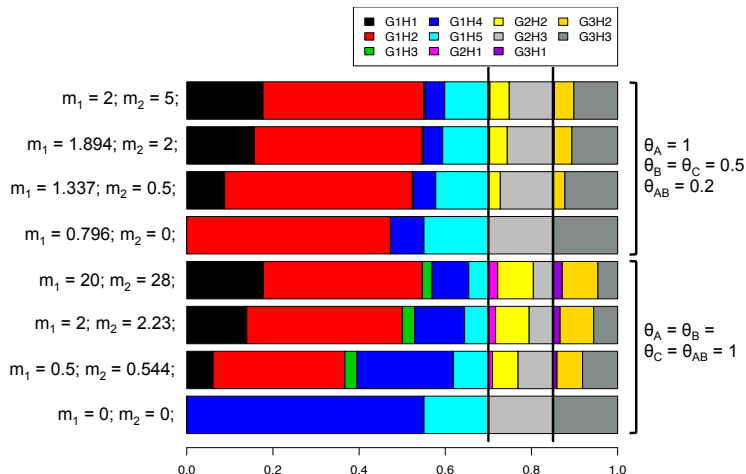
**Complication:** More histories are possible, because coalescent events can happen “before” speciation





## What about gene flow?

With A LOT of work, the probability distribution can be computed!



What about gene flow?

What do we learn from this?

- Many different choices of gene flow and coalescent parameters lead to an identical distribution on the *gene tree topologies*.

The rate of gene flow is not identifiable from the topology distribution!

- The distribution of *gene tree histories* appears to be distinct for different parameter choices.

Conjecture: The gene tree history distribution identifies the rate of gene flow.

- For some choices for the rates of gene flow, all three gene tree topologies are equally likely (and thus the species tree cannot be identified).

Contrast this with the situation in the absence of gene flow

Implications for triplet-based species tree inference methods

Reference: Tian and Kubatko (2016)

## What about mutation?

- What about mutation? How does this affect data analysis?
- The coalescent gives a model for determining gene tree probabilities for **each gene**.
- View DNA sequence data as the results of a two-stage process:
  - ▶ Coalescent process generates a gene tree topology.
  - ▶ Given this gene tree topology, DNA sequences evolve along the tree.

## What about mutation?

### Given this model, how should inference be carried out?

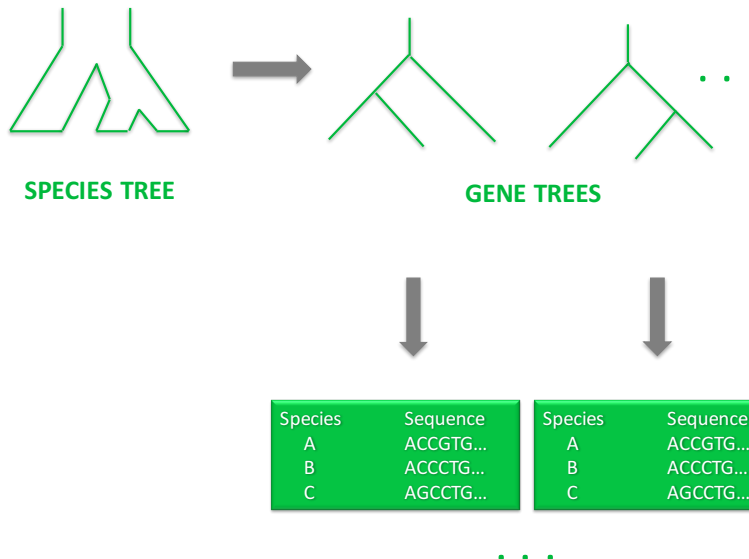
- As more data (genes) are added, the process of estimating species trees from concatenated data can be **statistically inconsistent**
  - May fail to converge to any single tree topology if there are many equally likely trees.
  - May converge to the wrong tree when a gene tree that is topologically incongruent with the species tree has the highest probability.
  - The bootstrap may be **positively misleading** – show strong support for an incorrect clade
- Important note: This is NOT a failing of the bootstrap methodology; the observed “poor” performance is due to the use of an incorrect model (concatenation)

Kubatko and Degnan, 2007

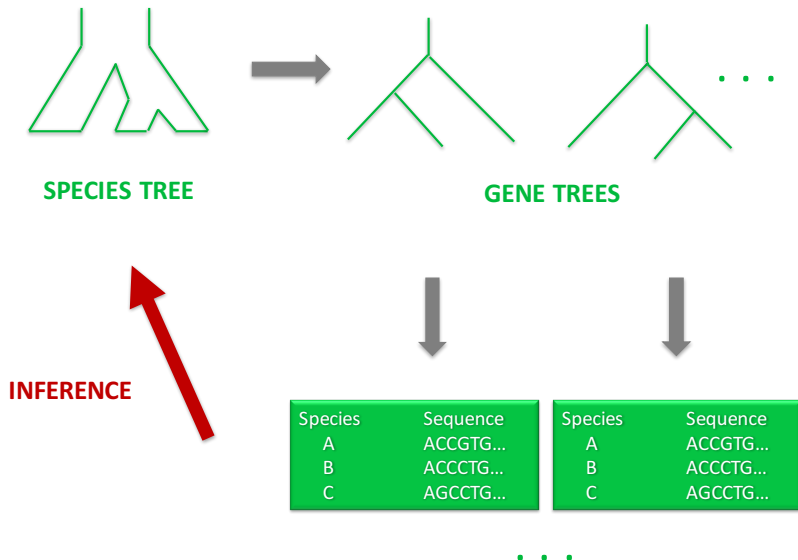
Is there a better way to estimate species phylogenies?

**Explicitly model the coalescent process!**

# Model Underlying Coalescent-based Species Tree Inference



# Model Underlying Coalescent-based Species Tree Inference



## The likelihood function

- Suppose that we have available alignments for  $N$  genes, denoted by  $D_1, D_2, \dots, D_N$
- We would like to find the likelihood of the species phylogeny given these  $N$  alignments, assuming that
  - ▶ individual gene trees are randomly generated according to the coalescent
  - ▶ evolution of sequences along fixed gene trees occurs following a standard nucleotide-based Markov model
  - ▶ the data for the genes are independent given the species tree and associated parameters



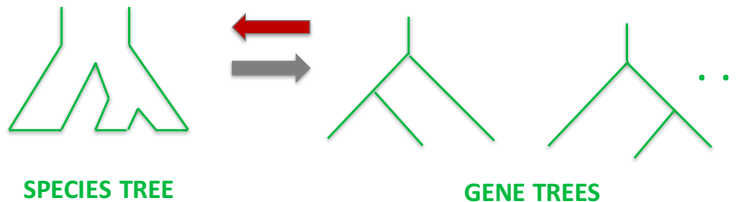
- Computation of the likelihood requires evaluation of the **Felsenstein equation** for a multi-locus data set:

$$\begin{aligned} L(S|D_1, D_2, \dots D_N) &= \prod_{i=1}^N P(D_i|S) \text{ [loci conditionally independent]} \\ &= \prod_{i=1}^N \sum_{j=1}^G P(D_i|g_j) f(g_j|S) \end{aligned}$$

where  $S$  is the species tree (topology and branch lengths) and  $g_j$  represents a gene tree.

- This likelihood is difficult to evaluate directly, because of the dimension of the inner sum (which is really an integral)
- To deal with this, either assume gene trees are known (**summary statistics methods**), use Bayesian techniques (**most full data approaches**), or think about small problems ☹. But new methods based on different ways of summarizing data are being developed!

# Model Underlying Coalescent-based Species Tree Inference



**SUMMARY  
STATISTICS  
METHODS**



Species	Sequence
A	ACCGTG...
B	ACCCTG...
C	AGCCTG...

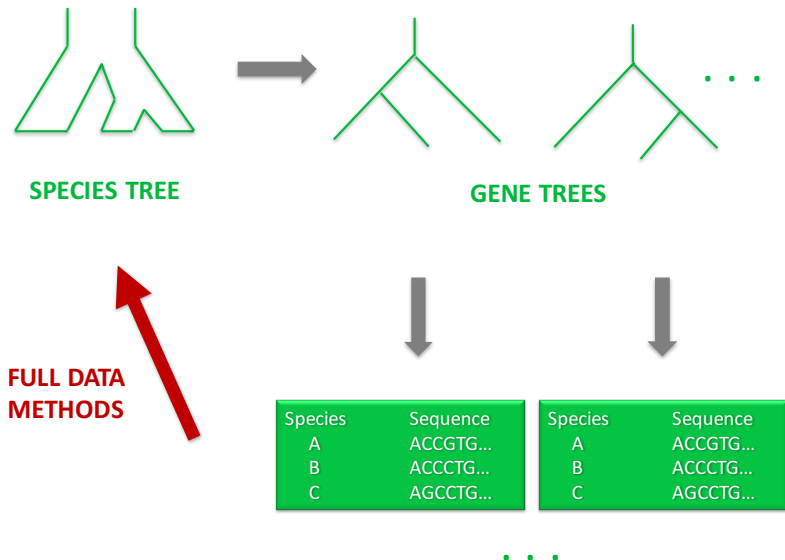
Species	Sequence
A	ACCGTG...
B	ACCCTG...
C	AGCCTG...

...

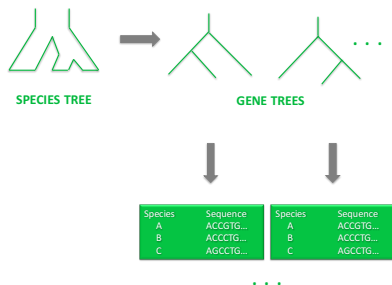
- **Summary statistics methods:** Start with estimated gene trees
  - ▶ Using estimated branch lengths:
    - ★ STEM (Kubatko et al. 2009)
    - ★ STEAC (Liu et al. 2009)
  - ▶ Using topology information only:
    - ★ STAR (Liu et al. 2009)
    - ★ Minimize Deep Coalescences (PhyloNet; Than & Nakhleh 2009)
    - ★ MP-EST (Liu et al. 2010)
    - ★ ST-ABC (Fan and Kubatko 2011)
    - ★ STELLS (Wu 2011)
    - ★ ASTRAL (Mirarab et al. 2014)
    - ★ Statistical binning (Bayzid et al. 2014)

- **Summary statistics methods:** Start with estimated gene trees
  - ▶ Using estimated branch lengths:
    - ★ STEM (Kubatko et al. 2009)
    - ★ STEAC (Liu et al. 2009)
  - ▶ Using topology information only:
    - ★ STAR (Liu et al. 2009)
    - ★ Minimize Deep Coalescences (PhyloNet; Than & Nakhleh 2009)
    - ★ **MP-EST** (Liu et al. 2010)
    - ★ ST-ABC (Fan and Kubatko 2011)
    - ★ STELLS (Wu 2011)
    - ★ **ASTRAL** (Mirarab et al. 2014)
    - ★ Statistical binning (Bayzid et al. 2014)

## Summary of Model Underlying Coalescent-based Species Tree Inference

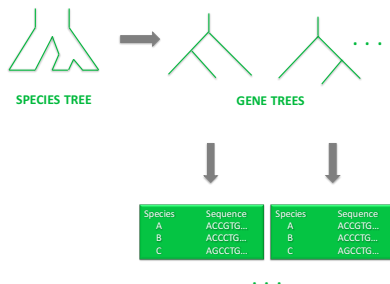


- Model the entire process of data generation
- Goal of these methods is to estimate the posterior distribution of the gene trees and species tree and associated model parameters



- BEST, \*BEAST, and BPP use MCMC by considering both gene trees and the species tree, but their implementations are different
- SNAPP uses a clever two-step peeling algorithm to carry out the integration over gene trees, allowing it to consider a reduced space – but currently limited to biallelic data.

- Model the entire process of data generation
- Avoid computing the likelihood by using algebraic structure in the distribution of site pattern probabilities under the model
- SVDQuartets is implemented in PAUP\*
- SVDQuartets will be discussed in detail in this afternoon's lab



- Comparison of approaches:

- ▶ Summary statistics methods

- ★ Advantage: Quick
- ★ Disadvantage: Ignore information in the data
- ★ Most current implementations do not easily allow assessment of uncertainty (but bootstrap can be used, at the expense of computational efficiency)

- ▶ Full data methods

- ★ Advantage: Fully model-based framework
- ★ Disadvantage: Computationally intensive, sometimes prohibitively so
- ★ BEST, \*BEAST, BPP, and SNAPP utilize a Bayesian framework and involve MCMC



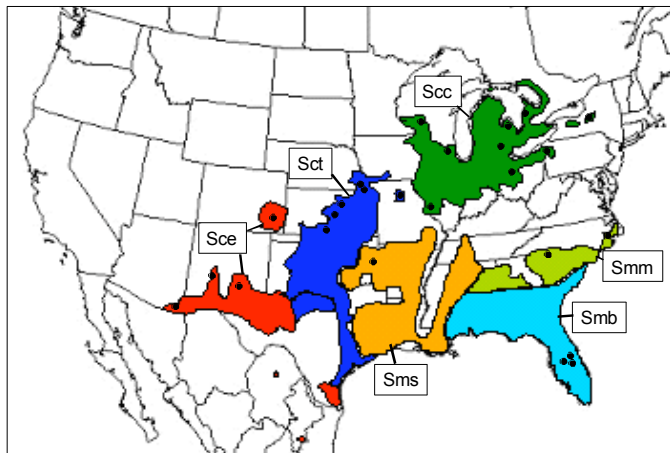
## Multilocus data example: *Sistrurus rattlesnakes*



- North American Rattlesnakes - Joint work with Dr. Lisle Gibbs (EEOB at OSU)
- Of interest evolutionarily because of the diversity of venoms present in the various species and subspecies.
- Of conservation interest because population sizes in the eastern subspecies are very small.

[Pictures by Jimmy Chiucchi and Brian Fedorko]

## Geographic Distribution of Snake Populations



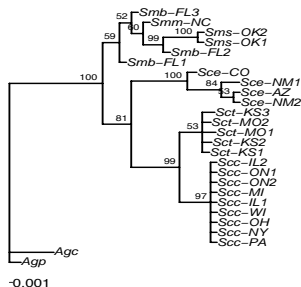
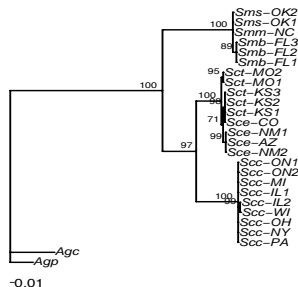


- Data: 7 (sub)species, 26 individuals (52 sequences), 19 genes

Species	Location	No. of individuals per gene
<i>S. catenatus catenatus</i>	Eastern U.S. and Canada	9
<i>S. c. edwardsii</i>	Western U.S.	4
<i>S. c. tergeminus</i>	Western and Central U.S.	5
<i>S. miliarius miliarius</i>	Southeastern U.S.	1
<i>S. m. barbouri</i>	Southeastern U.S.	3
<i>S. m. streckerii</i>	Southeastern U.S.	2
<i>Agkistrodon</i> sp. (outgroup)	U.S.	2

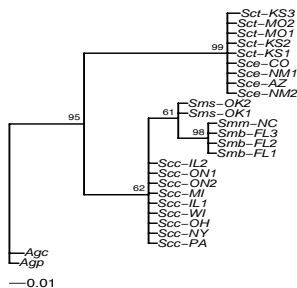
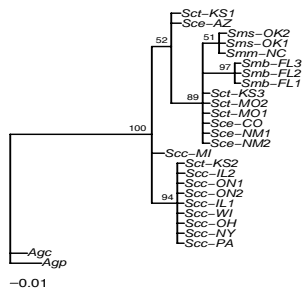
# Individual Gene Tree Estimates

Some are very informative:



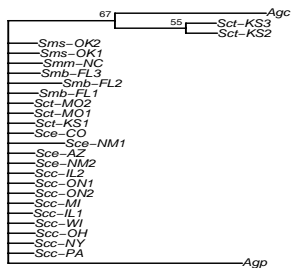
# Individual Gene Tree Estimates

Some are a little informative:

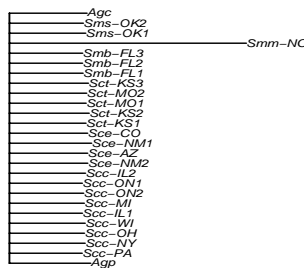


## Individual Gene Tree Estimates

And then there are others .....



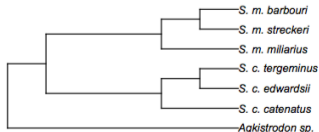
0.001



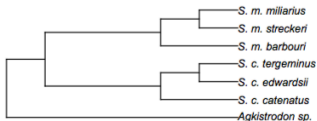
0.001

## Example: *Sistrurus rattlesnakes*

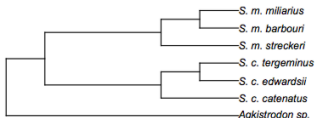
### STEM, STEAC



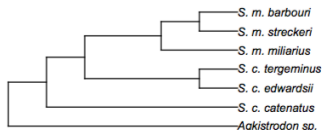
### BEAST (concatenated data), \*BEAST



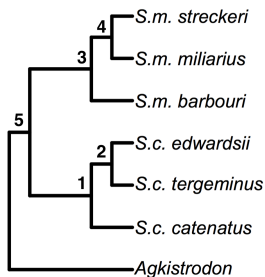
### BEST, Parsimony & MrBayes (concatenated data), Astral



### PhyloNet, STAR



## Example: *Sistrurus rattlesnakes*



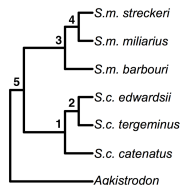
Node	1	2	3	4	5
*BEAST	100	100	100	46*	100
BPP	100	99	100	33*	100
SVDQ	93	100	100	46	100

\* = This clade was not in the maximum clade credibility (*S. m. miliarius* and *S. m. barbouri* received 48.78% posterior probability with \*BEAST and 59% posterior probability with BPP)



## Example: *Sistrurus* rattlesnakes

### Very rough ideas of computational time ...



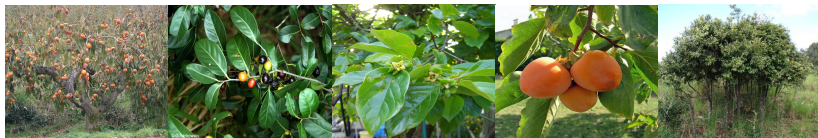
Program	Time	Details
BEST	~3 days	11,770,000 iterations (not converged)
*BEAST	16.8 hours	100,000,000 iterations all ESS > 200 except 1 (>100)
BPP	4 days	500,000 iterations
SVDQ	11 minutes	all quartets sampled 100 bootstrap reps
ASTRAL	2.215 sec	given gene trees! also need bootstrap

## SNP data example: *Diospyros*



From Wikipedia:

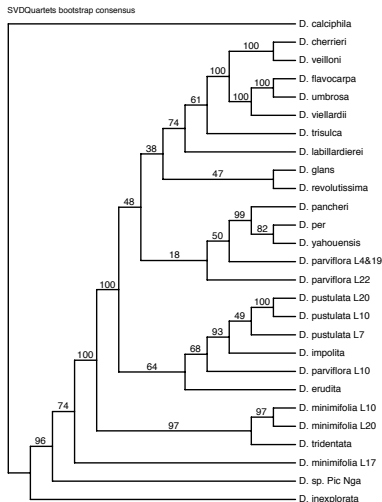
"Diospyros is a genus of over 700 species of deciduous and evergreen trees, shrubs and small bushes. The majority are native to the tropics, with only a few species extending into temperate regions. Depending on their nature, individual species are commonly known as ebony or persimmon trees. Some are valued for their hard, heavy, dark timber, and some for their fruit. Some are useful as ornamentals and many are of local ecological importance."



- Data: samples from New Caledonia archipelago – Ovidiu et al., Syst. Biol. 65(2):212-227, 2016
- 84 individuals, sampled from 39 populations, representing 21 species
- 26 tips on species tree
- Data set 1 (PAUP\*) : 8,488 SNPs
- Data set 2 (SNAPP) : 1,506 SNPs (one per locus)

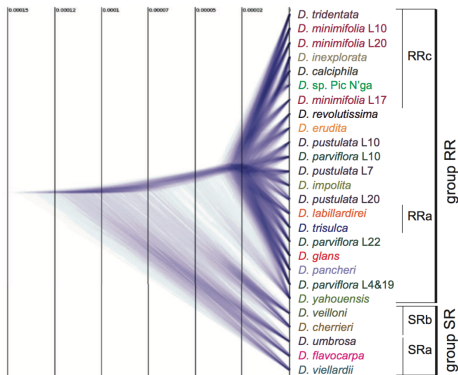
## Diospyros data

SVDQ: 15 hours, 27 minutes for 100 bootstrap reps,  
sampling all quartets



SNAPP: Ovidiu et al. (2016) used 5,000,000 iterations  
(all ESS >100, most > 200)

My analysis: 1500 iterations took ~ 2 days

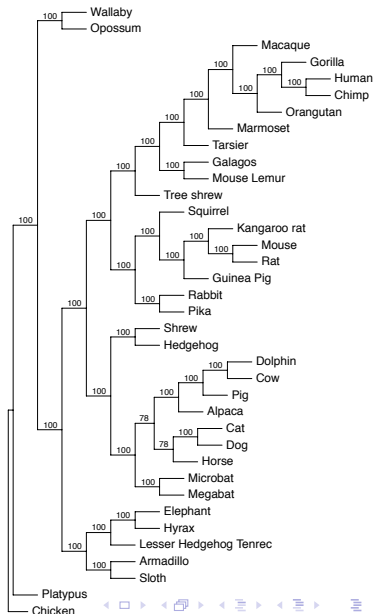


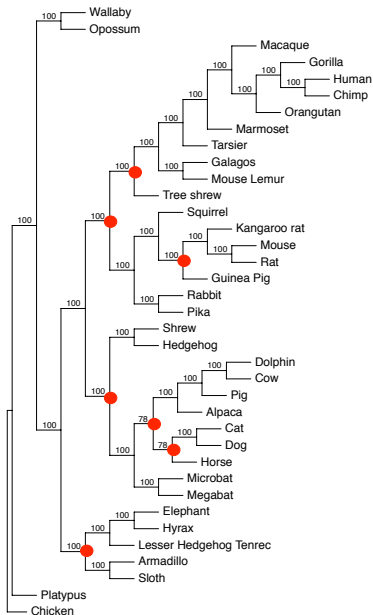
## Mammals example

- **Series of papers** in the literature debating proper phylogenetic relationships among a group of mammals
  - ▶ **Meredith RW, et al. (Science, 2011)** criticized by **Song et al. (PNAS, 2012)**:
    - ★ Amount of data “insufficient” (26 genes, 35,603 bp, 164 mammals)
    - ★ Concatenation not appropriate
  - ▶ Response by **Gatesy and Springer (PNAS, 2013)** criticizing Song et al.:
    - ★ Loci chosen not representative (“concatalescence” – exons ‘pasted’ together)
    - ★ Many nodes still not well supported
    - ★ Subset of 36 species
  - ▶ **Wu et al. (PNAS, 2013)** criticize Gatesy and Springer’s response:
    - ★ Concatenation of all genes is worse than within a few genes
    - ★ The approach of treating exons from a single gene with introns stripped has worked well in other cases
  - ▶ etc. ...
- **Question:** How does SVDQ do?

## Mammals example

- **Dataset:** obtained from Liang Liu, 36 mammal species + outgroup, ~ 1.4 million bp from 447 genes
- **SVDQ** run on 8-year old dual-core linux machine – **27 hours** required to estimate the tree and obtain bootstrap support from 100 replicates





“Historically problematic nodes” identified by McCormack et al. (Genome Research, 2012) are identified with a red circle

SVDQ addresses many criticisms:

- Strong support for many difficult nodes
- Recombination within genes is not a substantial model violation
- SVDQ can accommodate biased base frequencies, which may apply to the placement of bats

- Overall, the SVDQ analysis agrees with the analysis of Song et al. (2012), who used the coalescence-based method MP-EST
- The SVDQ analysis differs from analyses based on concatenation for some of the difficult nodes, but agrees with concatenation for the two nodes with lower bootstrap support
- The total time required for the analysis and the generality of the potential underlying model is a strength



## Species Tree Inference Summary – Comparison of Methods

Software	Data Type	Measure of Uncertainty	Computation Time	Models Included
BEST	multilocus	posterior probability	long; can be run in parallel	coalescent; all reversible substitution models
*BEAST	multilocus	posterior probability	intermediate; can be run in parallel	coalescent; all reversible substitution models; relaxed clock; variable population sizes
BPP	multilocus	posterior probability	long	coalescent; JC69 model only; species delimitation
SVDQ	multilocus; SNP	bootstrap	short	coalescent; all reversible substitution models; parameter estimation ?
SNAPP	biallelic SNP; AFLP	posterior probability	long; can be run in parallel	coalescent; two-state substitution model; Bayes factor delimitation
ASTRAL	unrooted gene trees	bootstrap	short given gene trees	no specific model assumed
MP-EST	rooted gene trees	bootstrap	short given gene trees	coalescent model

## Species Tree Inference Summary

- Failure to incorporate the coalescent model in estimation of the species tree can lead to statistical inconsistency, even when a method that is statistically consistent is applied.
- Many new methods for inferring species trees are being developed – each has its advantages and disadvantages.
- In addition, we should continue to think about other ways of using multi-locus data to its full advantage .... and we should be thinking beyond estimation of the species tree.
- Lots of areas emerging: species delimitation, incorporating horizontal events along the phylogeny, etc. – get involved and have fun!