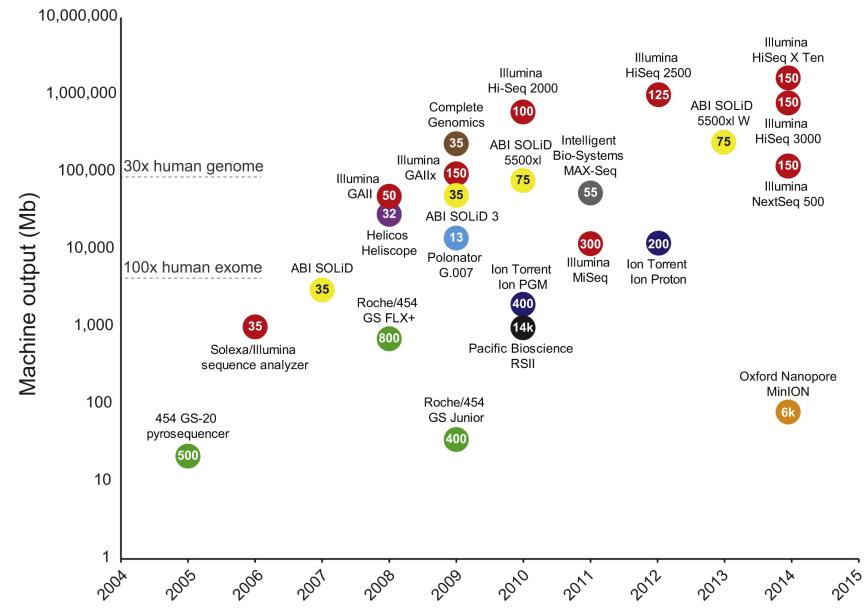


## Lecture 3.1

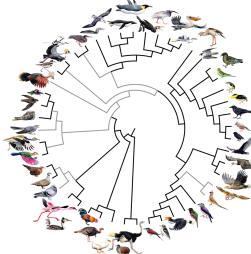
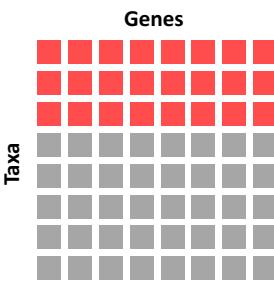
# Phylogenomics

Simon Ho



2

## Large data sets

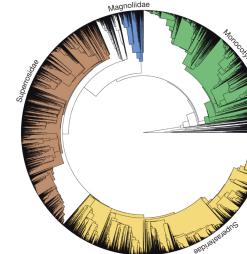
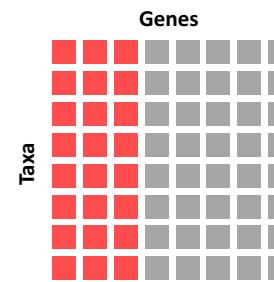


- Calculation of likelihood is expensive
  - Speed up by grouping sites with identical patterns
  - Approximate likelihood calculation
  - Multithreading/parallelisation

48 taxa  
8,295 genes  
Jarvis et al. (2014) *Science*

3

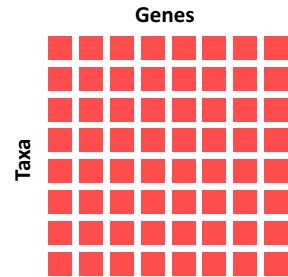
## Large data sets



32,223 taxa  
7 genes  
Zanne et al. (2014) *Nature*

4

## Large data sets



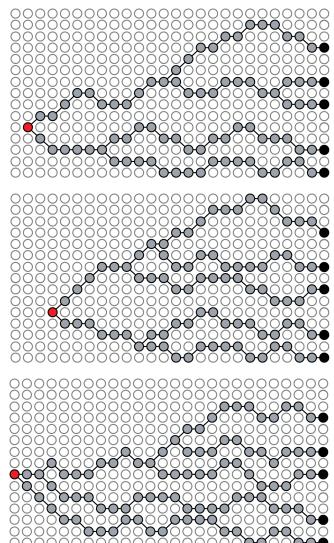
- Analysis is computationally expensive
- Consider filtering the data
  - Phylogenetic signal
  - Mutational saturation
  - Missing data
  - Model fit

5

## Genome-Tree Incongruence

## Gene trees in a species

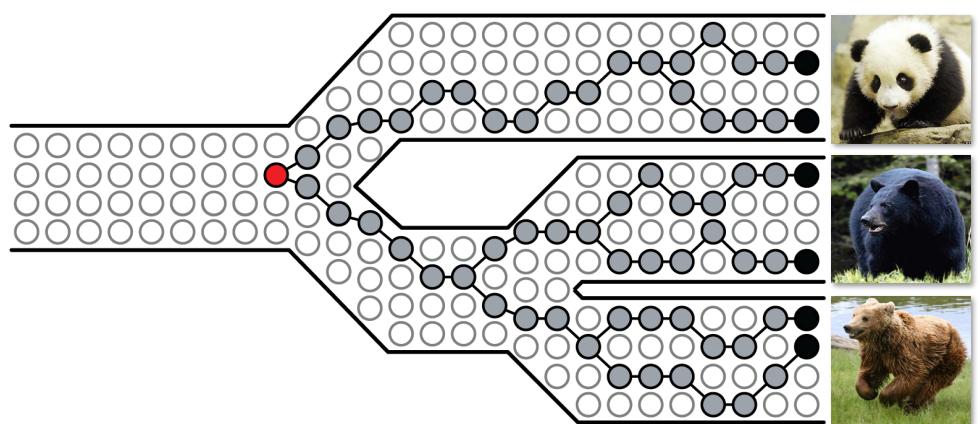
- Genealogies vary stochastically among unlinked loci
- Should not concatenate independent loci in a phylogenetic analysis of intraspecific data
  - Different trees
  - Different coalescence times



7

## Gene trees in multiple species

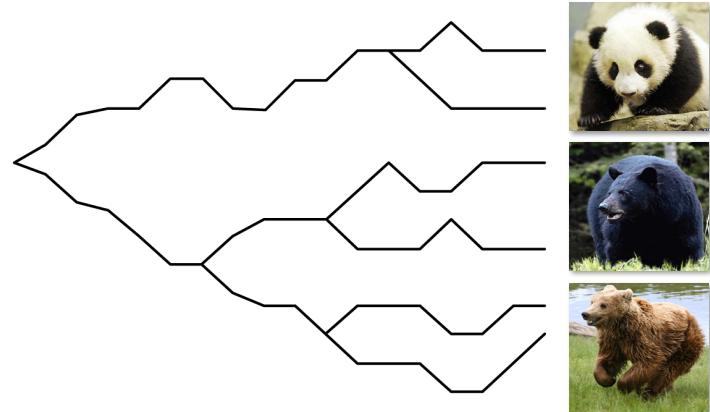
- Gene trees are embedded in the species tree



8

## Gene trees in multiple species

- Gene trees are embedded in the species tree

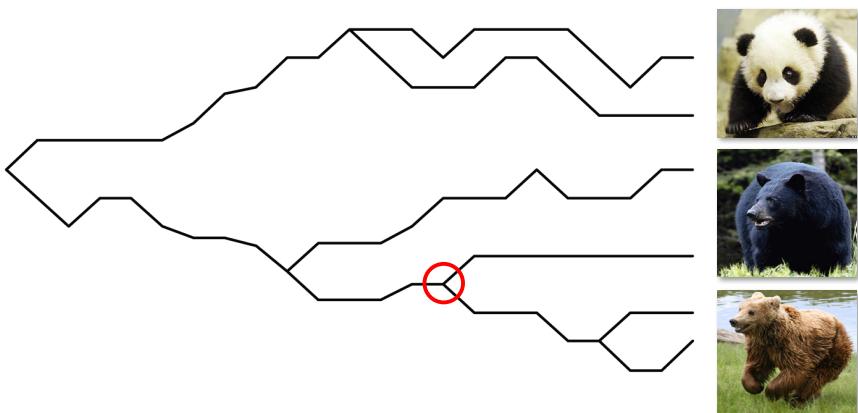


9

10

## Gene trees in multiple species

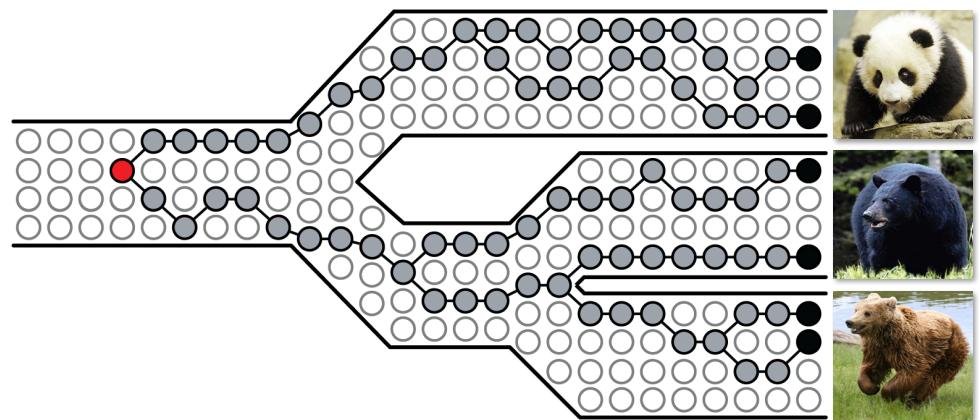
- Incomplete lineage sorting



11

## Gene trees in multiple species

- Incomplete lineage sorting



10

## Species tree

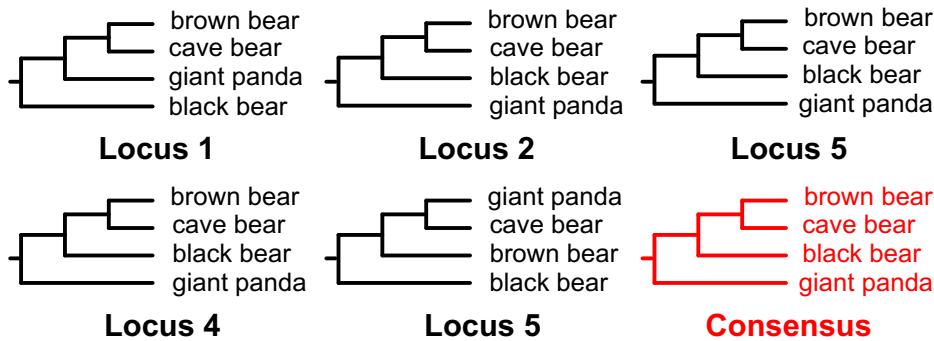
- Incomplete lineage sorting** can lead to gene trees that do not match the species tree
- We can infer the species tree from multiple gene trees
- Three approaches
  1. Consensus
  2. Concatenation
  3. Species-tree methods

12

## Species tree

### 1. Consensus

Estimate genealogy from each locus and find the consensus



But the most frequent gene tree does not always  
match the true species tree (anomaly zone)

13

## Analysing multiple loci

### 2. Concatenation

Assume that all loci share the same evolutionary history



But this ignores the occurrence of different gene trees

14

## Species tree

### 3. Species-tree methods

Estimate the species tree based on gene trees

- Gene trees are independent realisations of a stochastic process (the coalescent) on the same species tree
- Various methods
  - \*BEAST co-estimates gene trees and the species tree

BIOINFORMATICS

Vol. 30 ECCB 2014, pages i541–i548  
doi:10.1093/bioinformatics/btu462

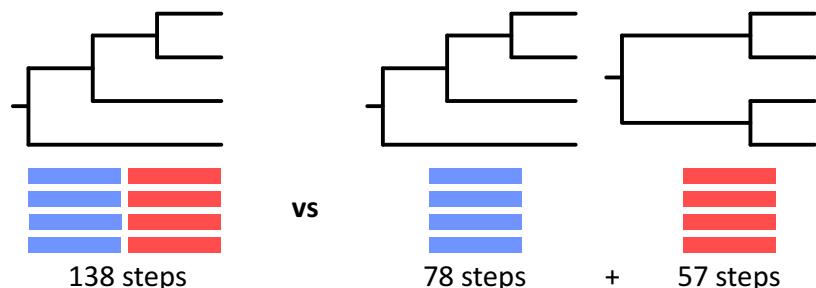
### ASTRAL: genome-scale coalescent-based species tree estimation

S. Mirarab<sup>1</sup>, R. Reaz<sup>1</sup>, Md. S. Bayzid<sup>1</sup>, T. Zimmermann<sup>1,2</sup>, M. S. Swenson<sup>3</sup> and T. Warnow<sup>1,\*</sup>

15

## Partition-homogeneity test

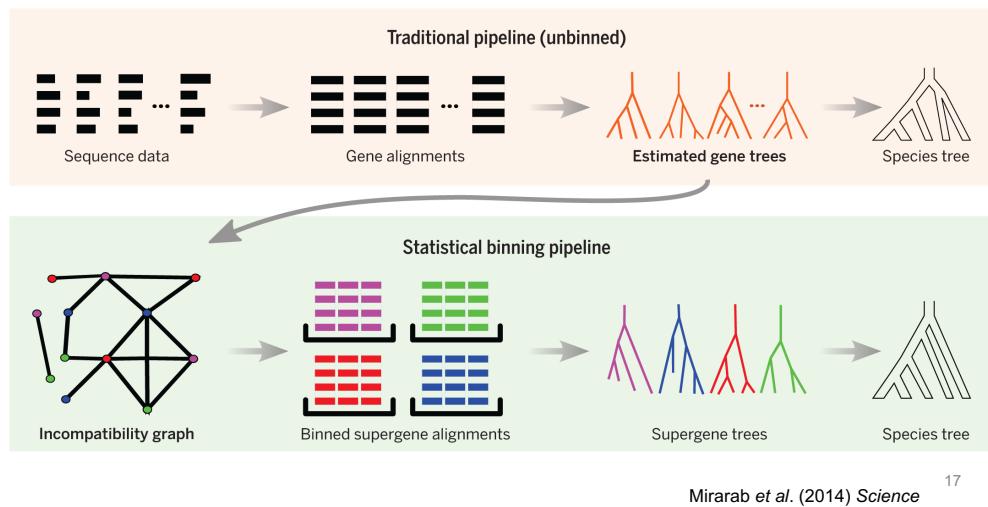
- Test for phylogenetic congruence across markers
- Partition-homogeneity (incongruence length difference) test



16

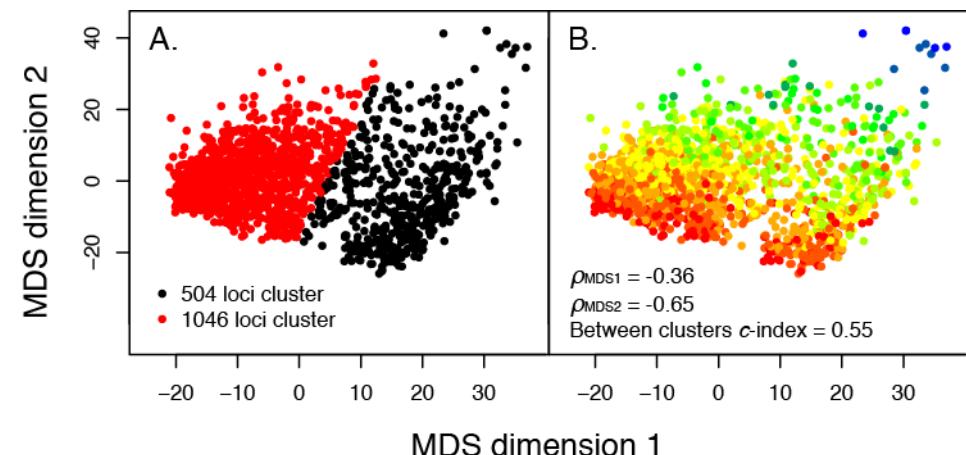
## Species tree

- Statistical binning

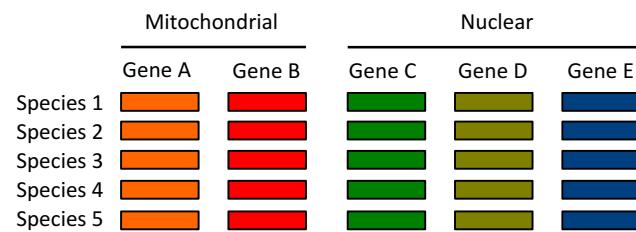


## Species tree

- Topology clustering



## Data partitioning



Lanfear et al. *BMC Evolutionary Biology* 2014, **14**:82  
http://www.biomedcentral.com/1471-2148/14/82



METHODOLOGY ARTICLE

Open Access

## Selecting optimal partitioning schemes for phylogenomic datasets

Robert Lanfear<sup>1,2\*</sup>, Brett Calcott<sup>3†</sup>, David Kainer<sup>1</sup>, Christoph Mayer<sup>4</sup> and Alexandros Stamatakis<sup>5,6</sup>

## Genome-Scale Dating

## Infinite-sites theory

- In Bayesian dating, the calibration priors provide the only information about absolute times
- With increasing data, there is no guarantee that posterior node ages will converge to the true values

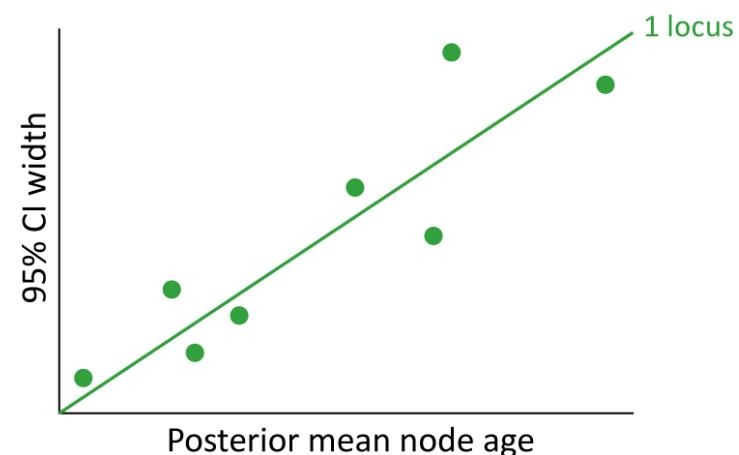
**Infinite-sites plot**

Plot of 95% credibility interval widths vs posterior means of node ages

Yang & Rannala (2006) *Mol Biol Evol* 21

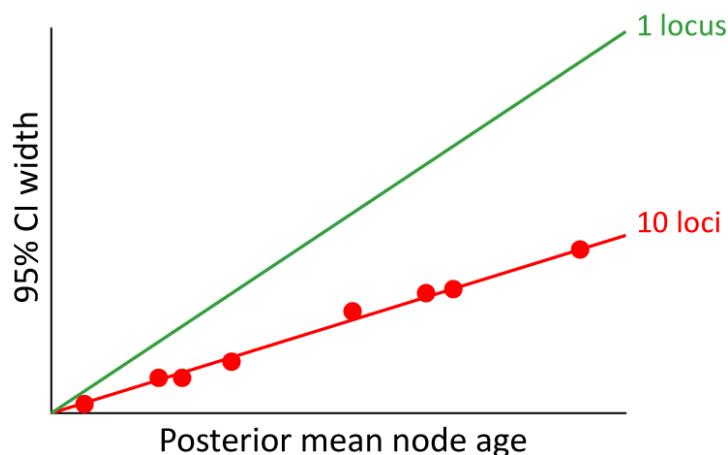
21

## Infinite-sites theory



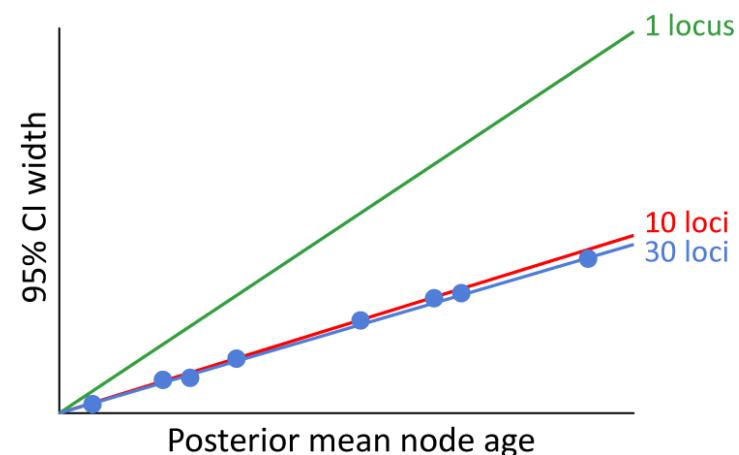
22

## Infinite-sites theory



23

## Infinite-sites theory



24

## Infinite-sites theory

- As the amount of data increases:
  - Approaches a straight line
  - Decline in the slope
- Moderate amounts of data behave similarly to infinite data
- Note that the slope does not reach zero with infinite data

### Even with infinite data

Uncertainty in date estimates depends on uncertainty in calibrations

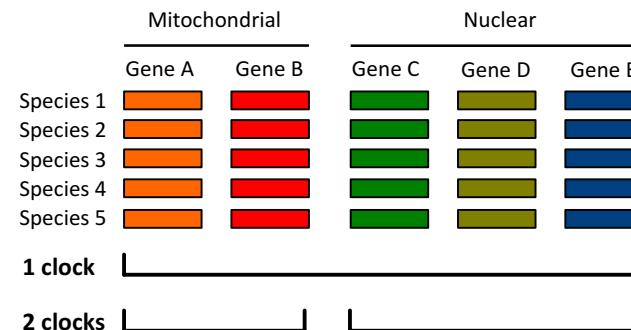
25

## Whole-genome analyses resolve early branches in the tree of life of modern birds

- 8295 genes from 52 taxa
- Identified 1156 clocklike genes
- 20 fossil-based calibrations
- Bayesian dating in *MCMCTree* with approximate likelihood calculation
- Able to test impact of choices of priors and calibrations

Jarvis et al. (2014) *Science* 27

## Data partitioning



### BIOINFORMATICS APPLICATIONS NOTE

Vol. 30 no. 7 2014, pages 1017–1019  
doi:10.1093/bioinformatics/btt665

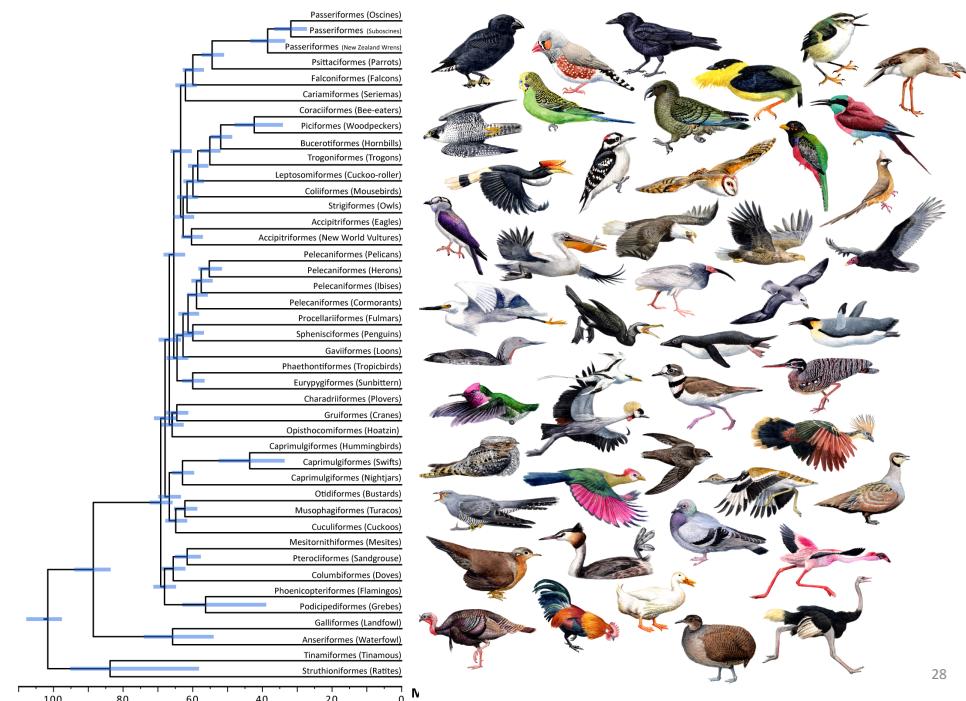
Phylogenetics

Advance Access publication November 14, 2013

### ClockstaR: choosing the number of relaxed-clock models in molecular phylogenetic analysis

Sebastián Duchêne\*,†, Martyna Molak and Simon Y. W. Ho†

26



28

## Data filtering

- Genome-scale projects have produced effectively infinite data
- We can be selective about the data that we use
- **Filter the data** according to some criterion
  - Sequence quality
  - Minimise missing data
  - Clocklike evolution
  - Fewest clock models
  - Model adequacy

29

## Useful references

- **Phylogenomic subsampling: a brief review**  
Edwards (2016) *Zoologica Scripta*, 45: 63–74.
- **Resolving conflict in eutherian mammal phylogeny using phylogenomics and the multispecies coalescent model**  
Song *et al.* (2012) *Proceedings of the National Academy of the USA*, 109: 14942–14947.
- **The changing face of the molecular evolutionary clock**  
Ho (2014) *Trends in Ecology and Evolution*, 29: 496–503.
- **Reconstructing evolutionary timescales using phylogenomics**  
Tong, Lo, & Ho (2016) *Zoological Systematics*, 41: 343–351.

30