

# Fundamentals of Evolution

Session 14 - 10/23/2017

Large-scale phylogenetics and Genomics

# Recap of Phylogenetics

- Tree-thinking: Reading history from trees.
- Tree inference:
  - Algorithmic: Minimize number of derived character changes.
  - Model based: Account for multiple changes (homoplasy)
  - Why use one method versus another?
- Probability and Likelihood
- Bayes theorem and Bayesian phylogenetics

# Where are we in the textbook?

You should have now read:

- Ch. 2 -- Tree of Life
- Ch. 16 -- Phylogeny
- Articles: Tree-thinking & Bird phylogeny

Read for next session:

- [Article: Heliconius introgression 2012](#)

# Recap of last session

**Bayesian statistics** asks “what is the probability of my hypothesis given the data?” by incorporating our prior belief as probability.

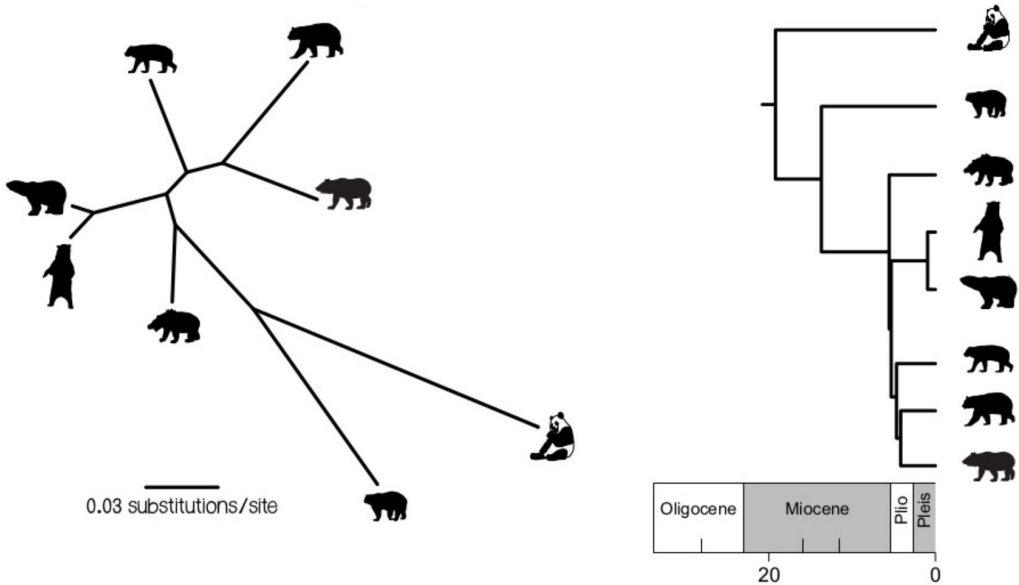
The diagram illustrates Bayes' theorem with the following components:

- Likelihood of hypothesis  $\theta$** : Points to  $\Pr(D|\theta)$  in the numerator.
- Prior probability of hypothesis  $\theta$** : Points to  $\Pr(\theta)$  in the numerator.
- Posterior probability of hypothesis  $\theta$** : Points to  $\Pr(\theta|D)$  on the left side of the equation.
- Marginal probability of the data (marginalizing over hypotheses)**: Points to the denominator  $\sum_{\theta} \Pr(D|\theta) \Pr(\theta)$ .

$$\Pr(\theta|D) = \frac{\Pr(D|\theta) \Pr(\theta)}{\sum_{\theta} \Pr(D|\theta) \Pr(\theta)}$$

# Why is Bayesian analysis useful for phylogenetics?

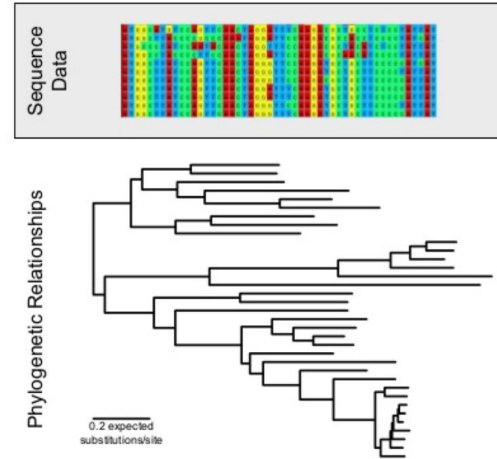
Phylogenies with branch lengths in units of time provide more information than unrooted trees with branch lengths in units of rate\*time (substitutions/site).



Sequence data provide information about **branch lengths**

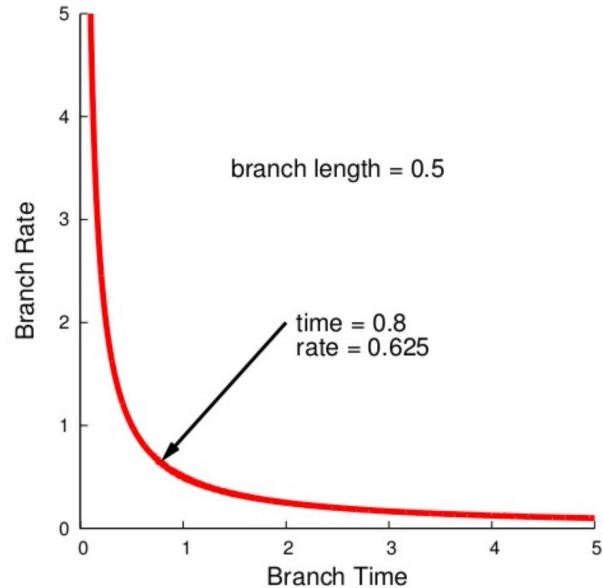
In units of **the expected # of substitutions per site**

branch length = rate  $\times$  time



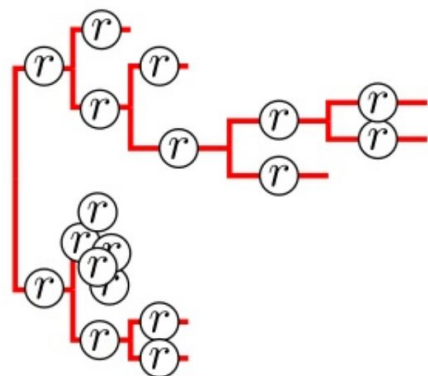
The sequence data provide information about branch length

for any possible rate, there's a time that fits the branch length perfectly

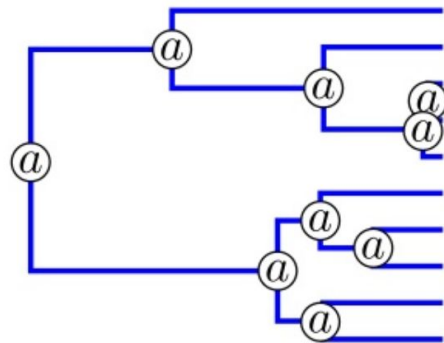


(figure based on Thorne & Kishino, 2005)

Methods for dating species divergences estimate the substitution rate and time separately



length = rate



length = time

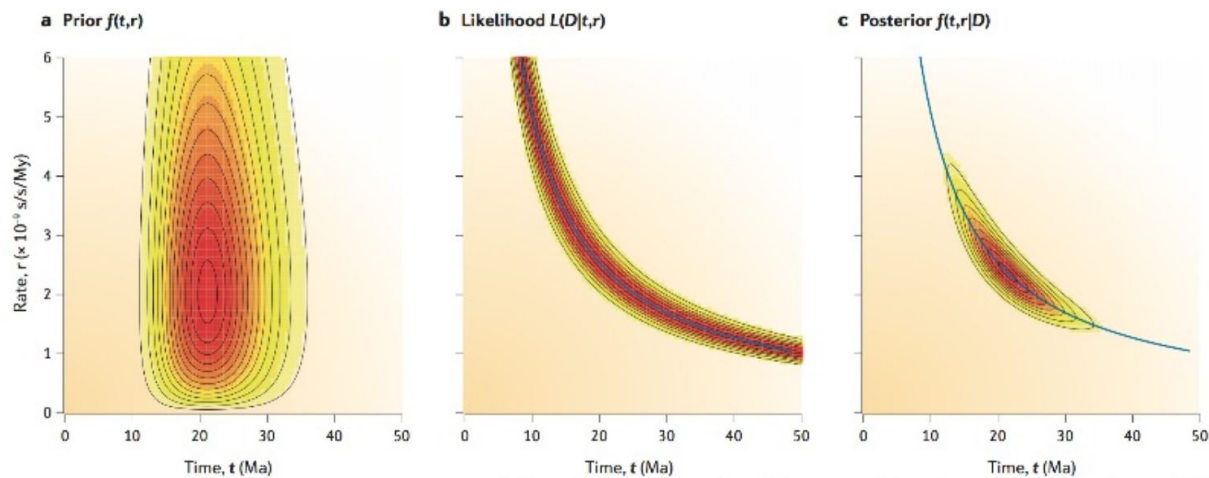
$$\mathcal{R} = (r_1, r_2, r_3, \dots, r_{2N-2})$$

$$\mathcal{A} = (a_1, a_2, a_3, \dots, a_{N-1})$$

$$N = \text{number of tips}$$



Methods for dating species divergences estimate the **substitution rate** and **time** separately



(dos Reis et al. *Nature Reviews Genetics*, 2016)

Tree-time priors for molecular phylogenies are only informative on a **relative** time scale

# Recap of last session

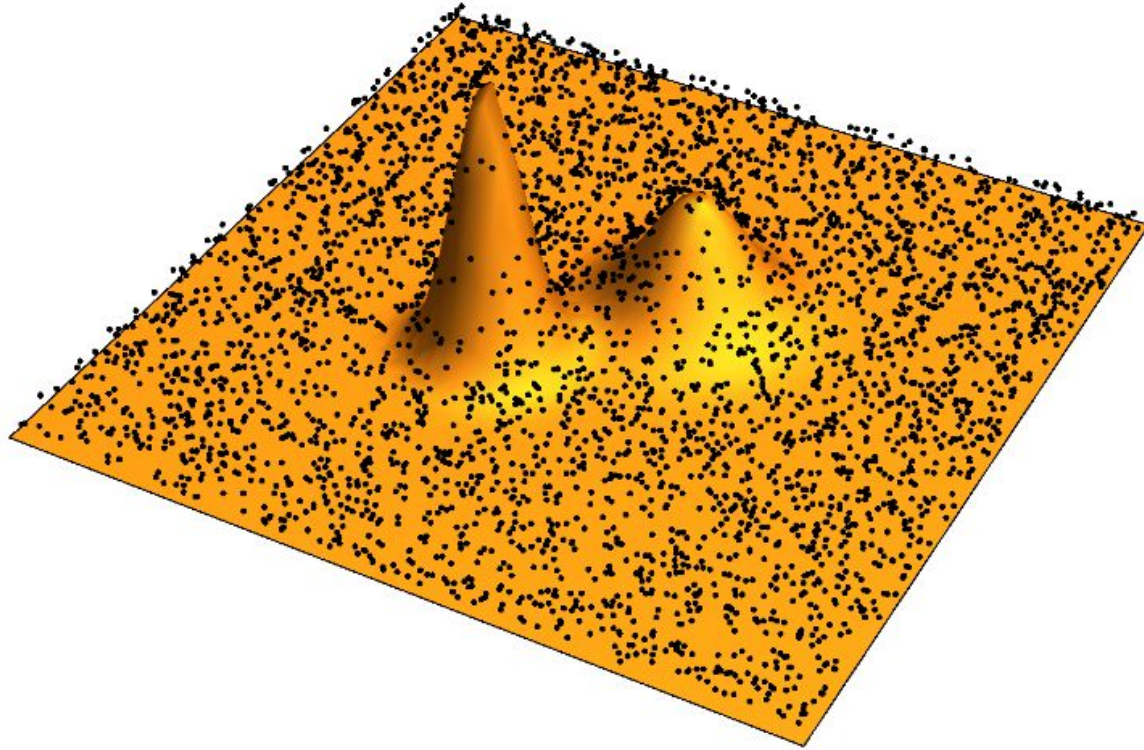
**Bayesian statistics** asks “what is the probability of my hypothesis given the data?” by incorporating our prior belief as probability.

The diagram illustrates Bayes' theorem with the following components and labels:

- Likelihood of hypothesis  $\theta$** : Points to  $\Pr(D|\theta)$  in the numerator.
- Prior probability of hypothesis  $\theta$** : Points to  $\Pr(\theta)$  in the numerator.
- Posterior probability of hypothesis  $\theta$** : Points to  $\Pr(\theta|D)$  on the left side of the equation.
- Marginal probability of the data (marginalizing over hypotheses)**: Points to the denominator  $\sum_{\theta} \Pr(D|\theta) \Pr(\theta)$ .

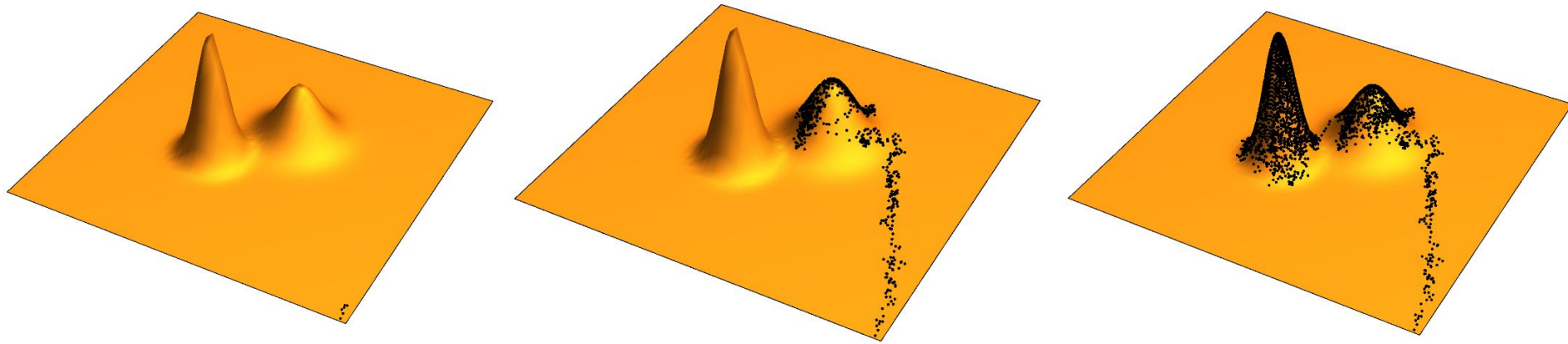
$$\Pr(\theta|D) = \frac{\Pr(D|\theta) \Pr(\theta)}{\sum_{\theta} \Pr(D|\theta) \Pr(\theta)}$$

## Naive integration approach

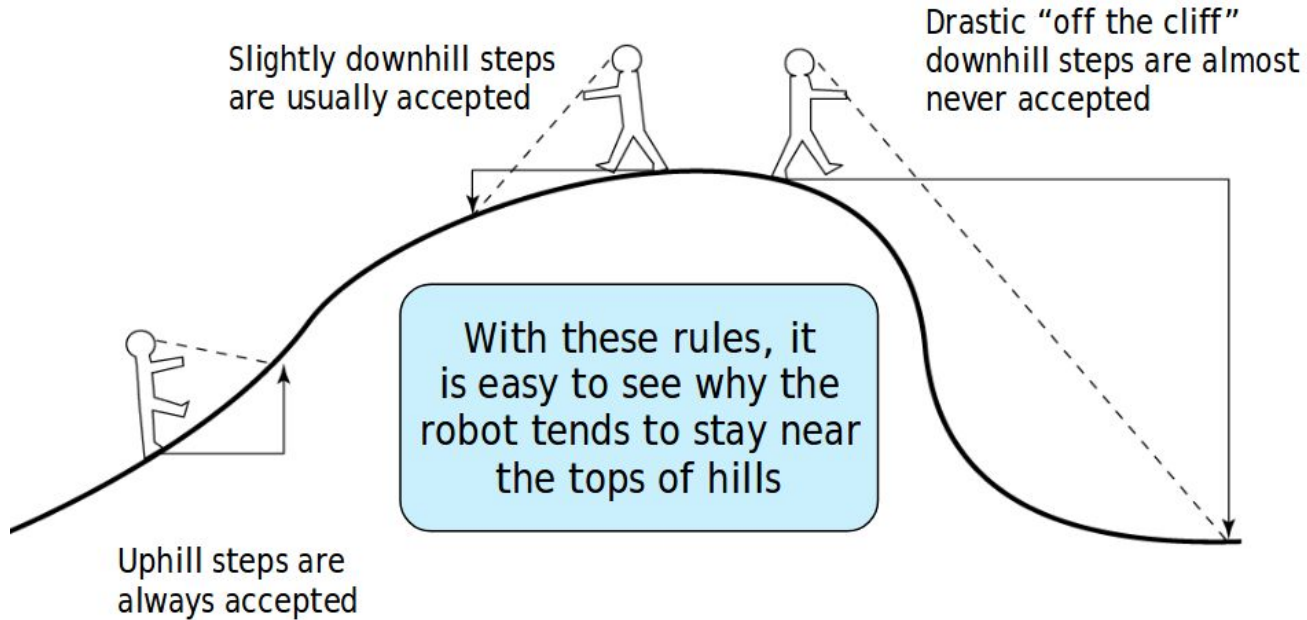


## Markov chain Monte-Carlo (MCMC)

Heuristic method of integrating across marginal probabilities. Mechanistic algorithm to search parameter space where the proportion of steps spent in any part of search space reflects the posterior probability support for that parameter. The result is a **posterior probability distribution**.



# MCMC robot's rules

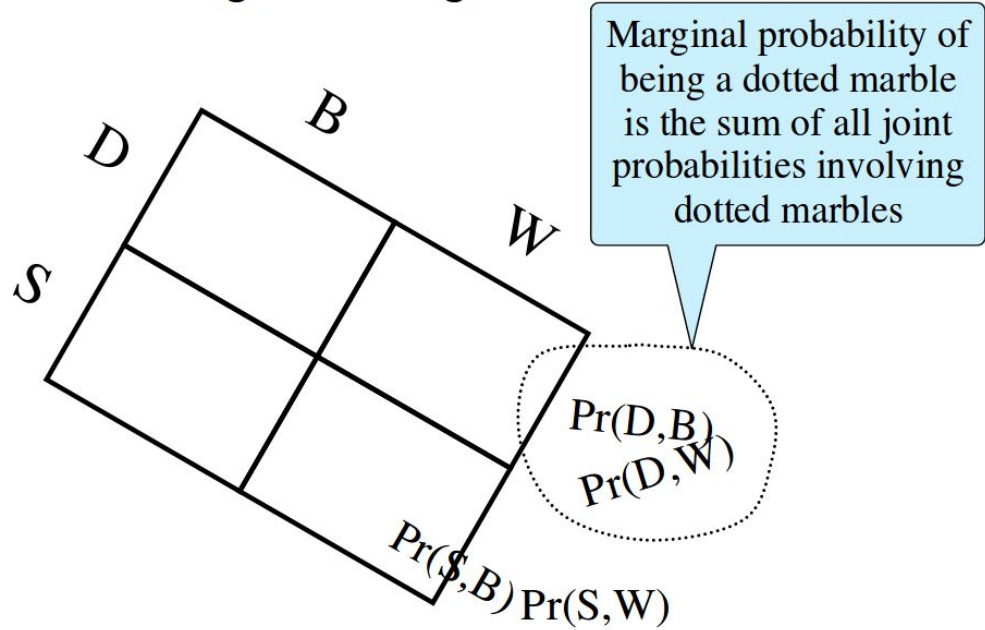


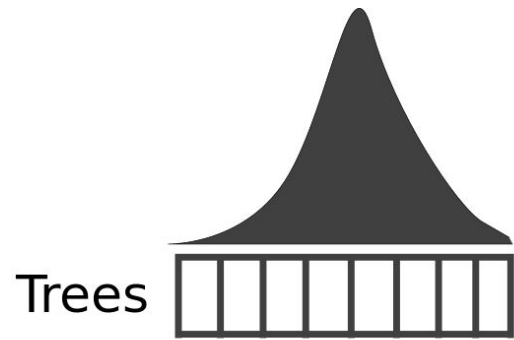
$$f(\mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s | D) =$$

$$\frac{f(D | \mathcal{R}, \mathcal{A}, \theta_s) f(\mathcal{R} | \theta_{\mathcal{R}}) f(\mathcal{A} | \theta_{\mathcal{A}}) f(\theta_s)}{f(D)}$$

$f(D   \mathcal{R}, \mathcal{A}, \theta_{\mathcal{R}}, \theta_{\mathcal{A}}, \theta_s)$	Likelihood
$f(\mathcal{R}   \theta_{\mathcal{R}})$	Prior on rates
$f(\mathcal{A}   \theta_{\mathcal{A}})$	Prior on node ages
$f(\theta_s)$	Prior on substitution parameters
$f(D)$	Marginal probability of the data

# Marginalizing over colors





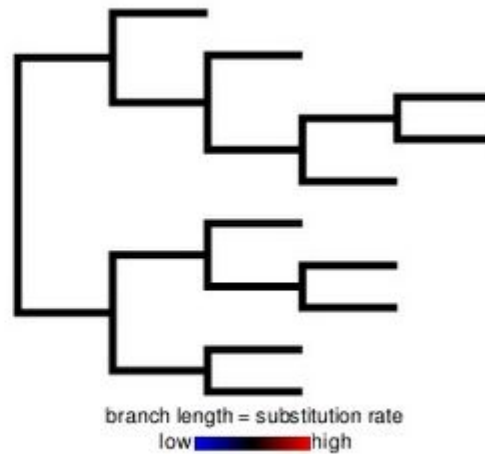


- **Global clock** (Zuckerkandl & Pauling, 1962)
- **Local clocks** (Hasegawa, Kishino & Yano 1989; Kishino & Hasegawa 1990; Yoder & Yang 2000; Yang & Yoder 2003, Drummond and Suchard 2010)
- **Punctuated rate change model** (Huelsenbeck, Larget and Swofford 2000)
- **Log-normally distributed autocorrelated rates** (Thorne, Kishino & Painter 1998; Kishino, Thorne & Bruno 2001; Thorne & Kishino 2002)
- **Uncorrelated/independent rates models** (Drummond et al. 2006; Rannala & Yang 2007; Lepage et al. 2007)
- **Mixture models on branch rates** (Heath, Holder, Huelsenbeck 2012)

The substitution rate is  
constant over time

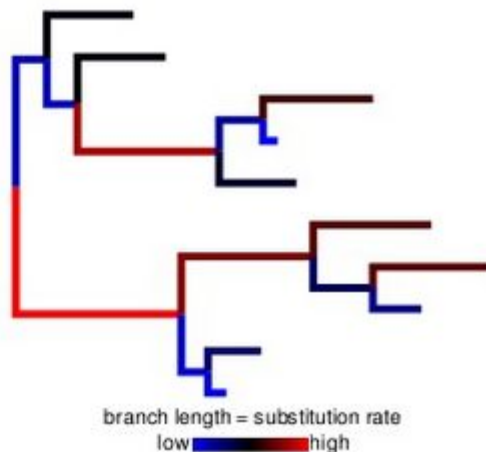
All lineages share the same  
rate

(Zuckerkandl & Pauling, 1962)



Lineage-specific rates are uncorrelated when the rate assigned to each branch is independently drawn from an underlying distribution

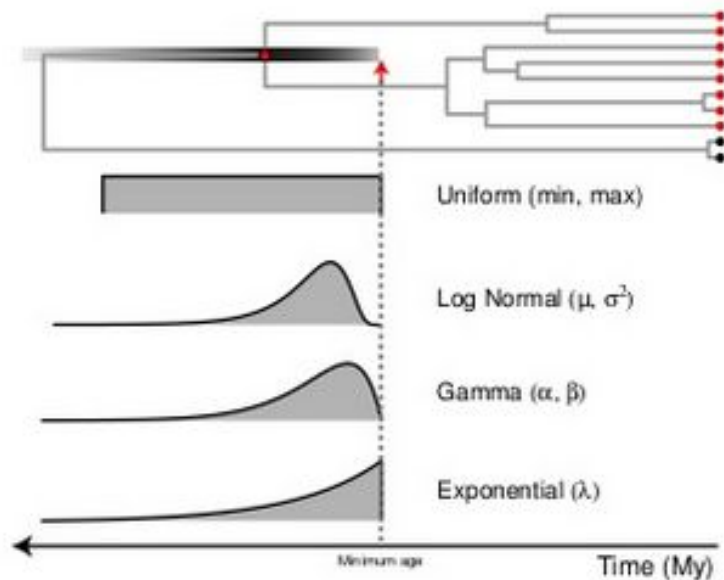
(Drummond et al. 2006; Rannala & Yang 2007; Lepage et al. 2007)



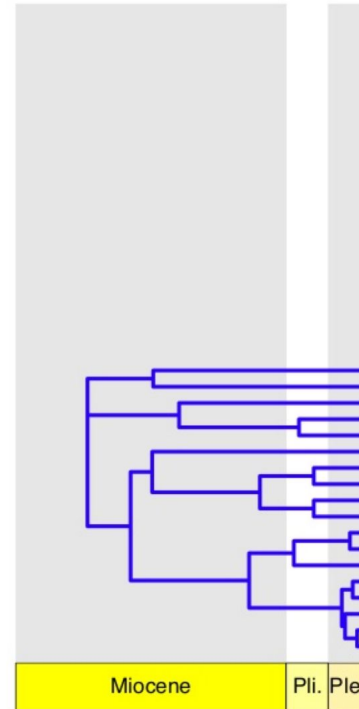
## Common practice in Bayesian divergence-time estimation:

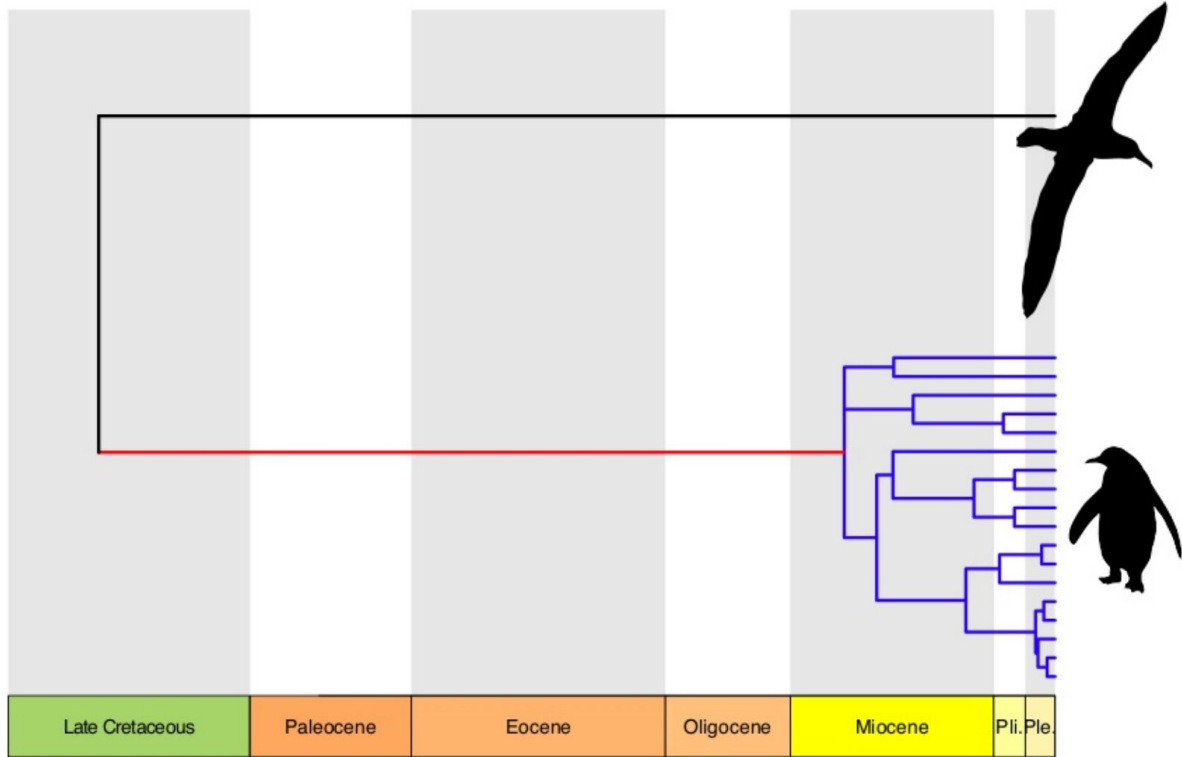
Parametric distributions are typically off-set by the age of the oldest fossil assigned to a clade

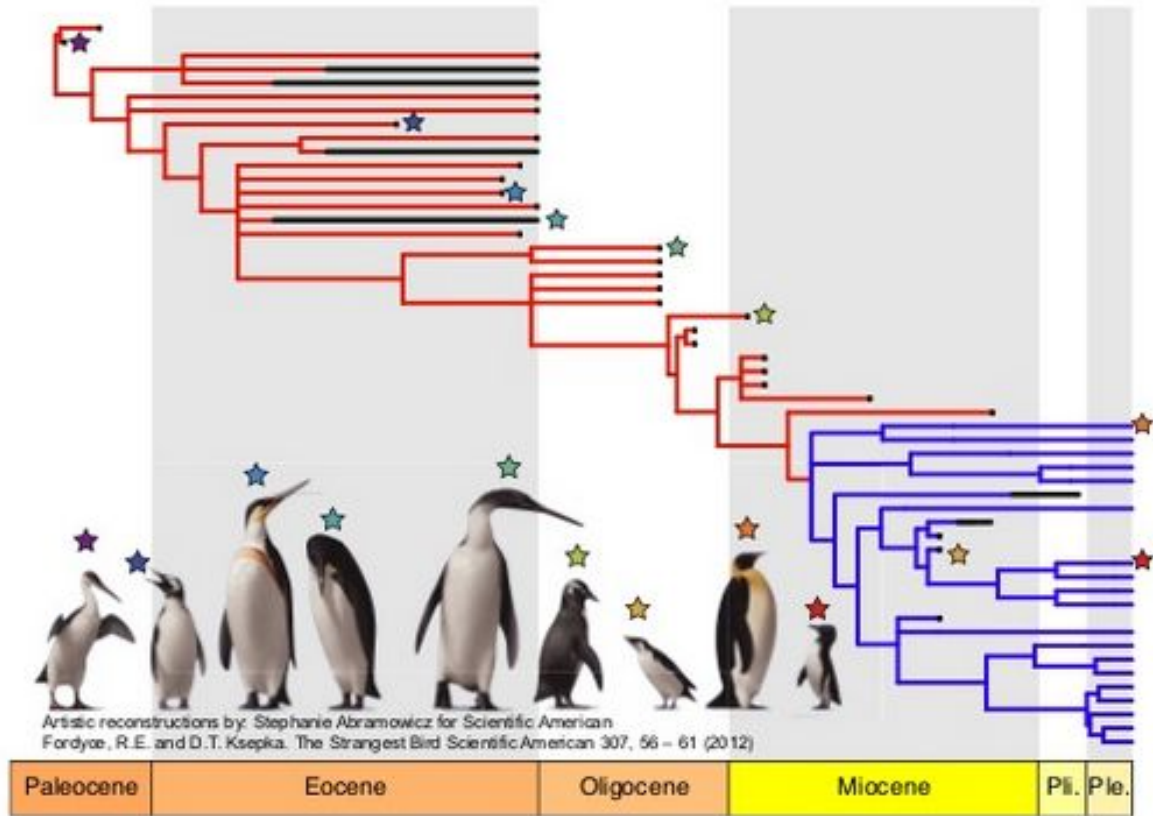
These prior densities do not (necessarily) require specification of maximum bounds



# PENGUIN DIVERSITY

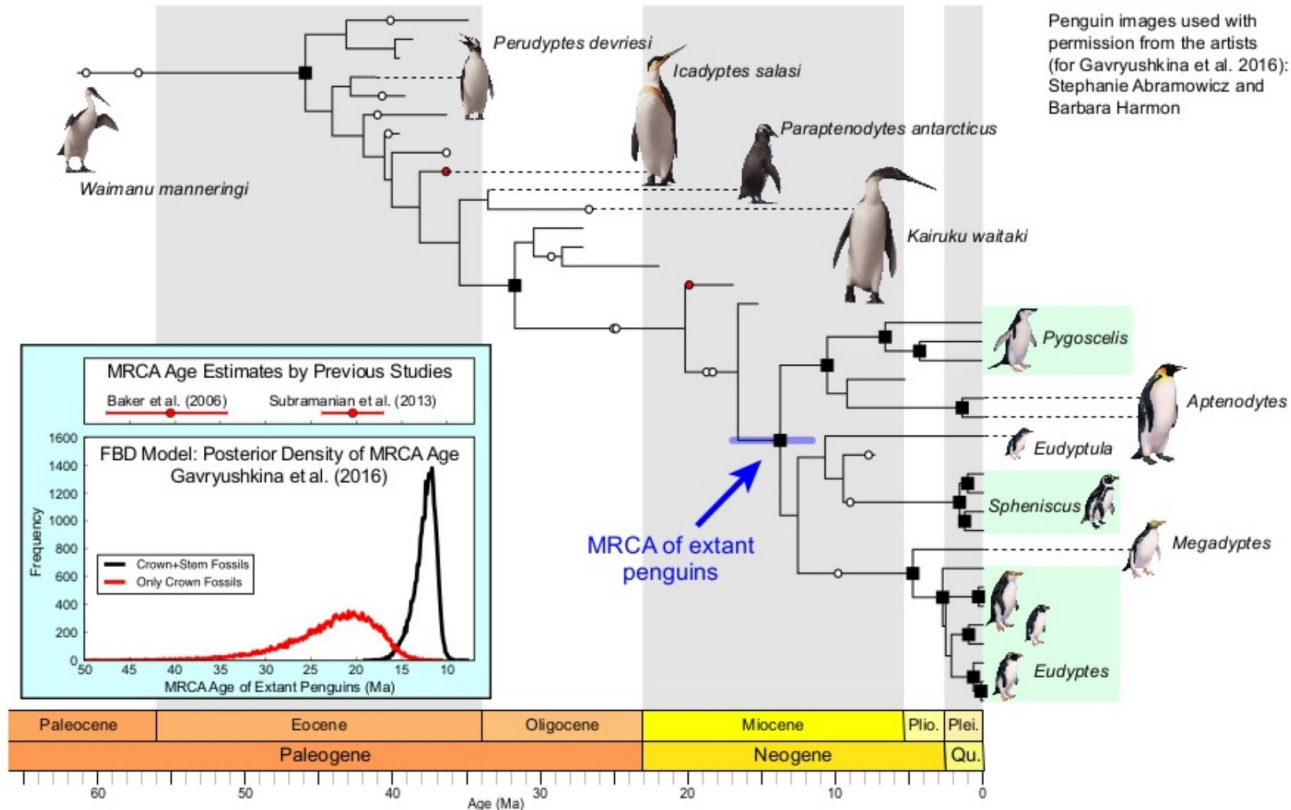






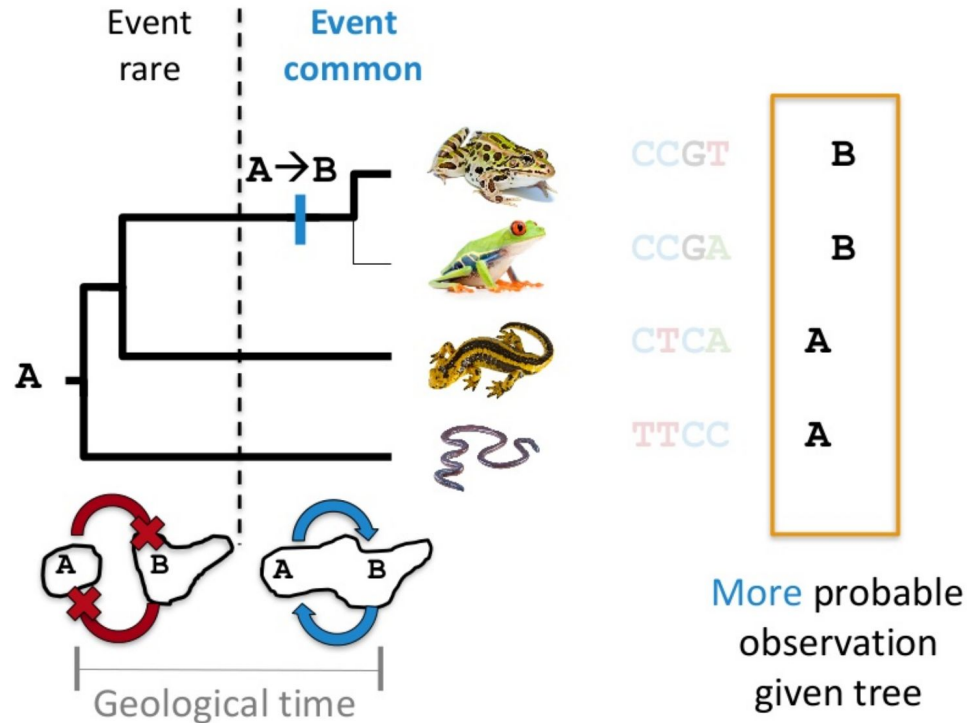
Artistic reconstructions by: Stephanie Abramowicz for Scientific American  
 Fordyce, R.E. and D.T. Ksepka. The Strangest Bird Scientific American 307, 56 – 61 (2012)

Incorporating both fossils and DNA sequences, and informed priors on the fossil placements, Gavryushkina et al. (2016) found the crown age of extant penguins is much younger than previously thought.





# Even without fossils, time-informed priors

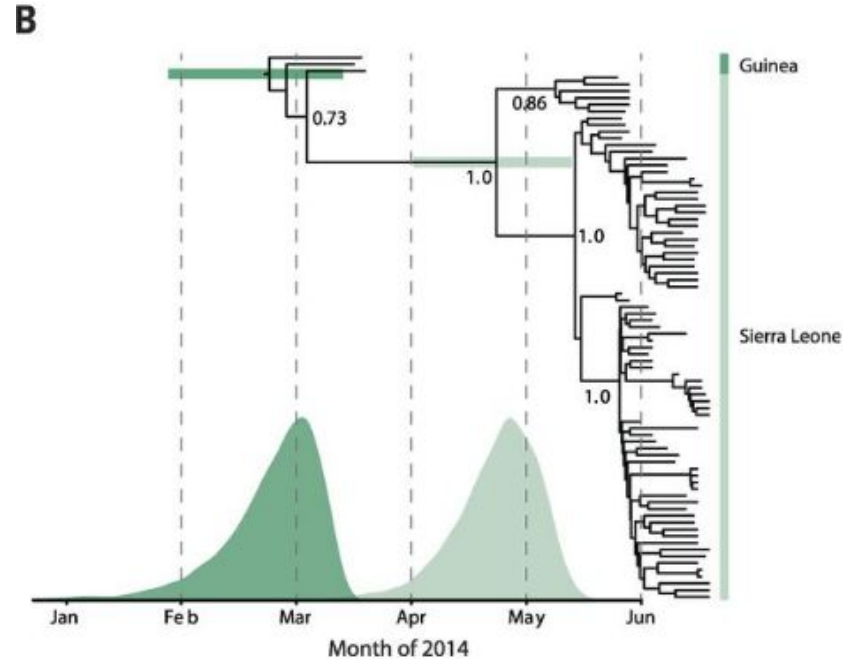


# Phylodynamics

- The study of how epidemiological, immunological, and evolutionary processes act and potentially interact to shape viral phylogenies.
- Bayesian phylogenetics is highly important because *rate* varies dramatically during viral outbreaks

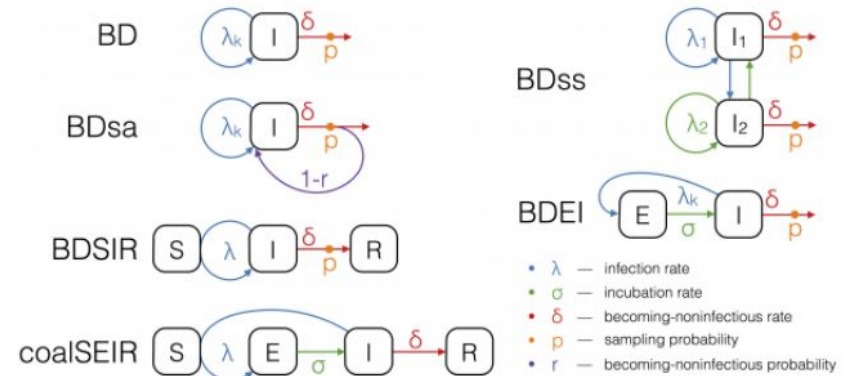
## Estimating the rate of infection of Ebola

- The 2013 West African Ebola virus epidemic spread primarily through Guinea, Sierra Leone and Liberia and killed over 11,000 people
- Estimated that strain began at a funeral in Guinea December 2013
- Phylogenetic analysis shows MRCA was February 2014 with 2 strains introduced to Sierra Leone.



## Estimating the rate of infection of Ebola

- Multiple birth-death model approaches were used to estimate epidemiological parameters across a Bayesian phylogeny.
- **Birth** is the rate of transmission, **death** is recovery or death of host.
- Incubation time: 4.92 days
- Infectious period: 2.58 days
- RO: 2.18 people



Stadler T *et al.* PLOS Currents Outbreaks. 2014

# Summary of Bayesian phylogenetics

- Broadly applicable statistical framework that allows one to combine data from many different sources through defining priors.
- In practice, often used for dated phylogenies because with priors on ages or rates you can better differentiate age from rate (which cannot be done in ML)
- However, it can be rather slow (MCMC search)
- And if you define too strict of priors and your data are not very informative then your results may just return what you put in. Requires careful testing/refining.

# Large-scale phylogenetics

- Increasingly, phylogenetic and phylogenomics is a field of **informatics**, or **data science**, and *computer science*.
- Data archiving and mining. Researchers focus on specific groups and over time accumulate enough data to span deeper and deeper in time.
- Methods for combining knowledge and minimizing the need to optimization + tree search.

# How many species are there?

nature  
microbiology

Article | [OPEN](#)

Recovery of nearly 8,000 metagenome-assembled genomes substantially expands the tree of life

# How many species are there?

- Globally, our best approximation to the total number of species, based on taxonomic expertise, is 3-100 million species (May 2010).
- Many methods are employed to estimate the number of undiscovered/described species: e.g., body-size distribution, species-area relationship, ratios between taxa, time-series relationships (Mora et al. 2011)



Species	Earth			Ocean		
	Catalogued	Predicted	±SE	Catalogued	Predicted	±SE
<b>Eukaryotes</b>						
Animalia	953,434	7,770,000	958,000	171,082	2,150,000	145,000
Chromista	13,033	27,500	30,500	4,859	7,400	9,640
Fungi	43,271	611,000	297,000	1,097	5,320	11,100
Plantae	215,644	298,000	8,200	8,600	16,600	9,130
Protozoa	8,118	36,400	6,690	8,118	36,400	6,690
<i>Total</i>	1,233,500	8,740,000	1,300,000	193,756	2,210,000	182,000
<b>Prokaryotes</b>						
Archaea	502	455	160	1	1	0
Bacteria	10,358	9,680	3,470	652	1,320	436
<i>Total</i>	10,860	10,100	3,630	653	1,320	436
<b>Grand Total</b>	<b>1,244,360</b>	<b>8,750,000</b>	<b>1,300,000</b>	<b>194,409</b>	<b>2,210,000</b>	<b>182,000</b>

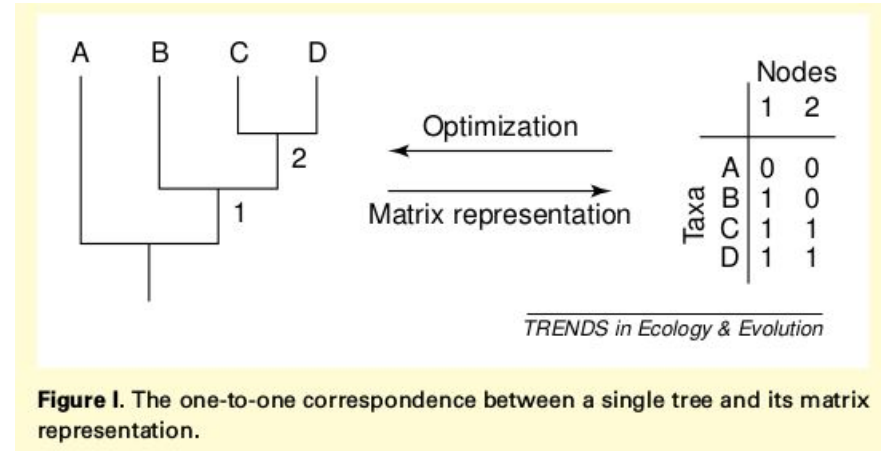
Predictions for prokaryotes represent a lower bound because they do not consider undescribed higher taxa. For protozoa, the ocean database was substantially more complete than the database for the entire Earth so we only used the former to estimate the total number of species in this taxon. All predictions were rounded to three significant digits.

doi:10.1371/journal.pbio.1001127.t002

(Mora et al. 2011)

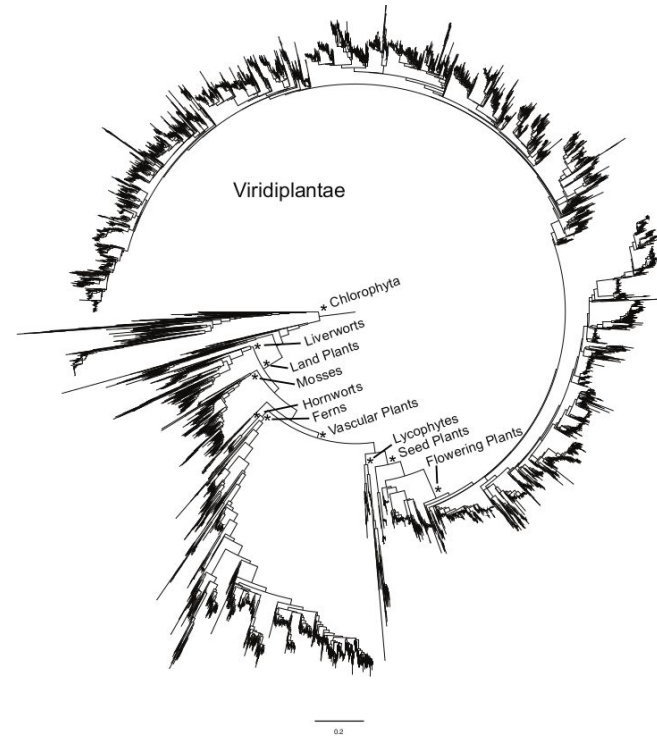
# Large-scale phylogenetics

- Super trees:
- Inferring large trees is difficult and time consuming, it is easier to join together smaller trees. Several techniques.
- The largest phylogenies that we have are all supertrees.



# Large-scale phylogenetics

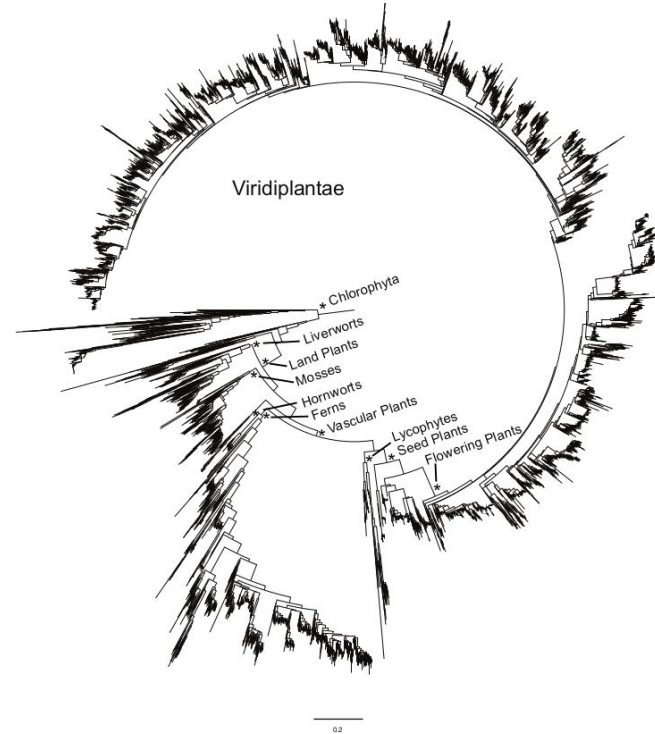
- Supermatrices:
- Around the early 2000s common markers were discovered that could be sequenced reliably across many organisms, which made it possible to combine their data into larger analyses. Faster inference methods developed.
- Hundreds of taxa sequenced at one or more of the same genes.



**Figure 3**  
Maximum-likelihood phylogeny for 13,533 species of green plants based on *rbcL* DNA sequences. The data matrix was constructed using the mega-phylogeny method; major clades are labeled and denoted with a star.

# Large-scale phylogenetics

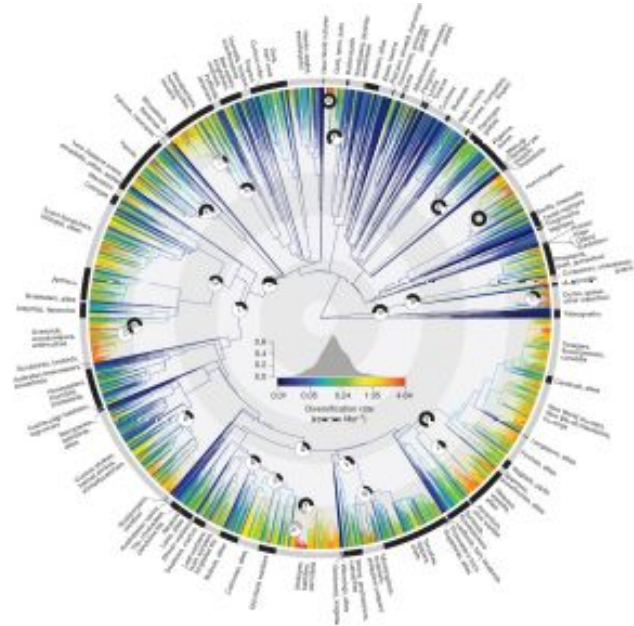
- Megaphylogeny pipelines: Automated procedures to build supermatrices by finding sequences in databases and aligning them at multiple hierarchical levels.
- Example: >13K species of plants analyzed for one gene.



**Figure 3**  
Maximum-likelihood phylogeny for 13,533 species of green plants based on *rbcL* DNA sequences. The data matrix was constructed using the mega-phylogeny method; major clades are labeled and denoted with a star.

# Large-scale phylogenetics

- Time-scaled megaphylogenies:
- Bayesian relaxed clock analysis on a reduced set of taxa to infer the backbone.
- Many smaller Bayesian relaxed clock analyses of subclades are added to the estimated backbone.



# Large-scale phylogenetics

- National Science Foundation initiatives to support Assembling the Tree of Life programs starting in early 2000s.

## **Assembling the fungal tree of life: progress, classification, and evolution of subcellular traits <sup>1</sup>**

Plant scientists plan massive effort to sequence 10,000 genomes

Genome 10K is a project to sequence the genome of at least one individual from each vertebrate genus, approximately 10,000 genomes. It is a key milestone on

# Large-scale phylogenetics

- Open Tree of Life.
- Compilation of all published phylogenetic knowledge.
- Uses a taxonomy (lists of groups within groups) to stitch trees together where information is missing.
- Stores information about conflict among different published studies as a network.

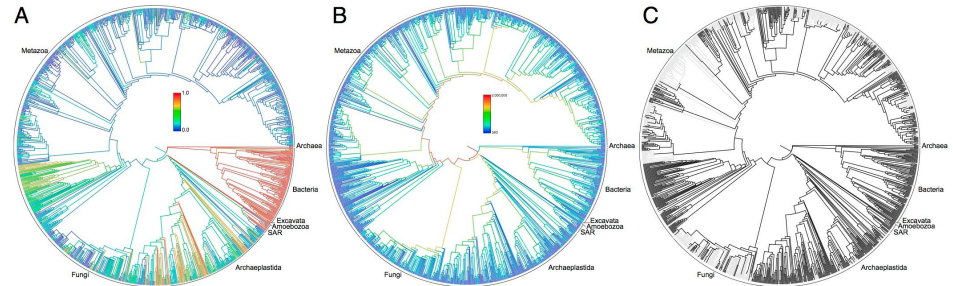


Fig. 1. Phylogenies representing the synthetic tree. The depicted tree is limited to lineages containing at least 500 descendants. (A) Colors represent proportion of lineages represented in NCBI databases. (B) Colors represent the amount of diversity measured by number of descendant tips. (C) Dark lineages have at least one representative in an input source tree.

# Large-scale phylogenetics

- However, some groups are difficult to characterize as ‘species’, and therefore to confirm sampling.
- Most data does not end up in databases.
- Manual curation and ranking remains necessary.

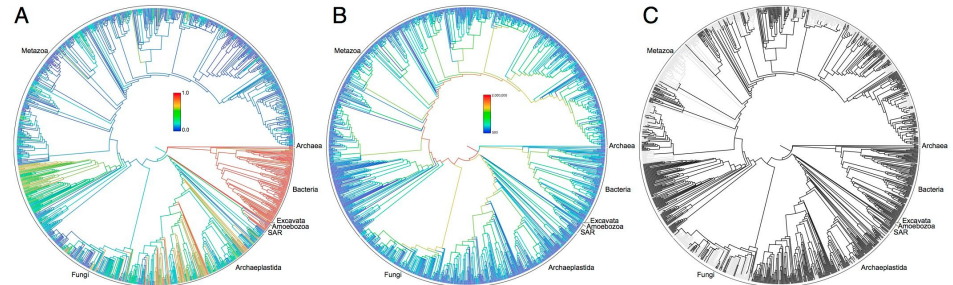


Fig. 1. Phylogenies representing the synthetic tree. The depicted tree is limited to lineages containing at least 500 descendants. (A) Colors represent proportion of lineages represented in NCBI databases. (B) Colors represent the amount of diversity measured by number of descendant tips. (C) Dark lineages have at least one representative in an input source tree.



# Summary of large-scale phylogenetics

- Supermatrix approaches combine huge numbers of taxa for few genes. Often sparse matrices (missing data). Made possible by algorithmic and computational improvements to likelihood calculations.
- Supertree methods aim to combine information from multiple trees without the need to infer the actual sequence data for all samples at once.
- At the largest scale, both approaches are typically combined to *stitch together* the tree of life with both known (inferred) relationships, and estimated (taxonomy) relationships. *A lot of work remains to be done!*

<https://pollev.com/dereineaton004>

## A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing

Richard O. Prum<sup>1,2\*</sup>, Jacob S. Berv<sup>3\*</sup>, Alex Dornburg<sup>1,2,4</sup>, Danil J. Field<sup>2,5</sup>, Jeffrey P. Townsend<sup>1,6</sup>, Emily Moriarty Lemmon<sup>7</sup> & Alan R. Lemmon<sup>8</sup>

Although reconstruction of the phylogeny of living birds has progressed tremendously in the last decade, the evolutionary history of Neoaves—a clade that encompasses nearly all living bird species—remains the greatest unresolved challenge in dinosaur systematics. Here we investigate avian phylogeny with an unprecedented scale of data: >390,000 bases of genomic sequence data from each of 198 species of living birds, representing all major avian lineages, and two crocodilian outgroups. Sequence data were collected using anchored hybrid enrichment, yielding 259 nuclear loci with an average length of 1,523 bases for a total data set of over  $7.8 \times 10^7$  bases. Bayesian and maximum likelihood analyses yielded highly supported and nearly identical phylogenetic trees for all major avian lineages. Five major clades form successive sister groups to the rest of Neoaves: (1) a clade including nightjars, other caprimulgiforms, swifts, and hummingbirds; (2) a clade uniting cuckoos, bustards, and turacos with pigeons, mesites, and sandgrouse; (3) cranes and their relatives; (4) a comprehensive waterbird clade, including all diving, wading, and shorebirds; and (5) a comprehensive landbird clade with the enigmatic hoatzin (*Opisthocomus hoazin*) as the sister group to the rest. Neither of the two main, recently proposed Neoavian clades—Columbea and Passerea<sup>1</sup>—were supported as monophyletic. The results of our divergence time analyses are congruent with the palaeontological record, supporting a major radiation of crown birds in the wake of the

It has long been recognized that phylogenetic confidence depends not only on the number of characters analysed and their rate of evolution, but also on the number and relationships of the taxa sampled relative to the nodes of interest<sup>9–11</sup>. Theory predicts that sampling a single taxon that diverges close to a node of interest will have a far greater effect on phylogenetic resolution than will adding more characters<sup>11</sup>. Despite using an alignment of >40 million base pairs, sparse sampling of 48 species in the recent avian genomic analysis may not have been sufficient to confidently resolve the deep divergences among major lineages of Neoaves. Thus, expanded taxon sampling is required to test the monophyly of neoavian clades, and to further resolve the phylogenetic relationships within Neoaves.

Here, we present a phylogenetic analysis of 198 bird species and 2 crocodilians (Supplementary Table 1) based on loci captured using anchored enrichment<sup>12</sup>. Our sample includes species of 122 avian families in all 40 extant avian orders<sup>2</sup>, with denser representation of non-oscine birds (108 families) than of oscine songbirds (14 families). Effort was made to include taxa that would break up long phylogenetic branches, and provide the highest likelihood of resolving short internodes at the base of Neoaves<sup>11</sup>. We also sampled multiple species within groups whose monophyly or phylogenetic interrelationships have been controversial—that is, tinamous, nightjars, hummingbirds, turacos, cuckoos, pigeons, sandgrouse, mesites, rails, storm petrels, petrels, storks, herons, hawks, hornbills, mousebirds, trogons, king-

## The goal of Prum et al. (2015) was to:

Sequence DNA from every species of bird for phylogenetic inference

Sequence huge amounts of DNA from one species in every major lineage of birds for phylogenetic inference

Use morphology to study the phylogeny of birds

Prum et al. (2015) used a method called "Sequence capture" to generate DNA sequences. Which statement best describes these data?

The entire genomes of each species were sequenced which included billions of DNA sites

A few hundred loci were sequenced which included millions of DNA sites

A few hundred loci were sequenced which included hundreds of thousands of DNA sites

# Bird phylogeny (Prum et al.)

- “Birds (Aves) are the most diverse lineage of extant tetrapod vertebrates. They comprise over 10,000 living species”
- “Here, we present a phylogenetic analysis of 198 bird species and 2 crocodylians based on loci captured using anchored enrichment. Our sample includes species of 122 avian families in all 40 extant avian orders”
- “We targeted 394 loci centred on conserved anchor regions of the genome that are flanked by more variable regions. We performed all phylogenetic analyses on a data set of 259 genes with the highest quality assemblies. The average locus was 1,524 bases in length”

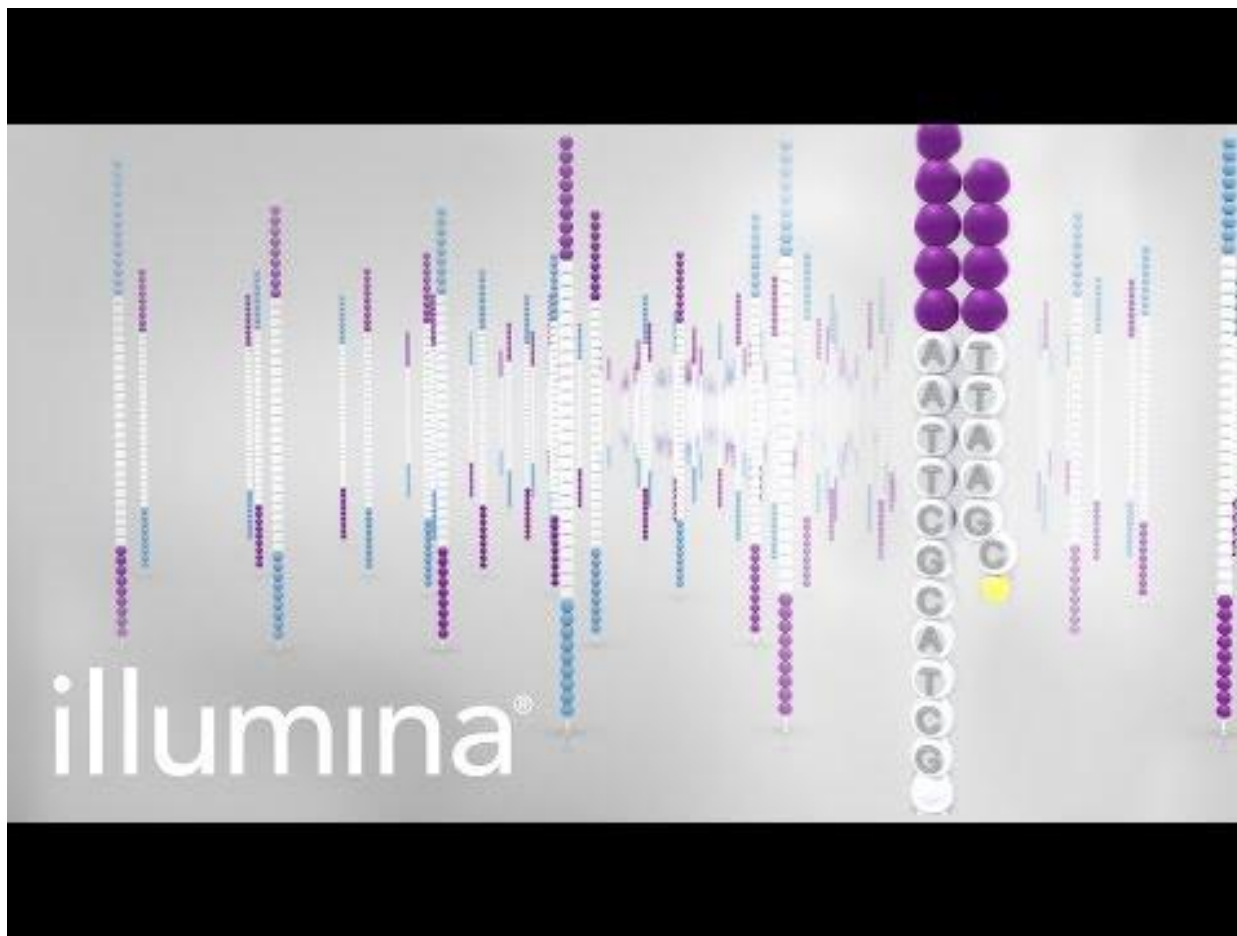
# Bird phylogeny (Prum et al.)

- “Our results indicate that the recent genome phylogeny (Jarvis et al.) may contain some erroneous relationships induced by long branch attraction from sparse taxon sampling”.
- “Differences in tree topology when taxa are excluded are to be expected if early internodes in Neoaves are very short. Adding taxa that have diverged near nodes of interest has been theoretically demonstrated to constrain the possible historical substitution patterns, and increase the accuracy of phylogenetic inference”

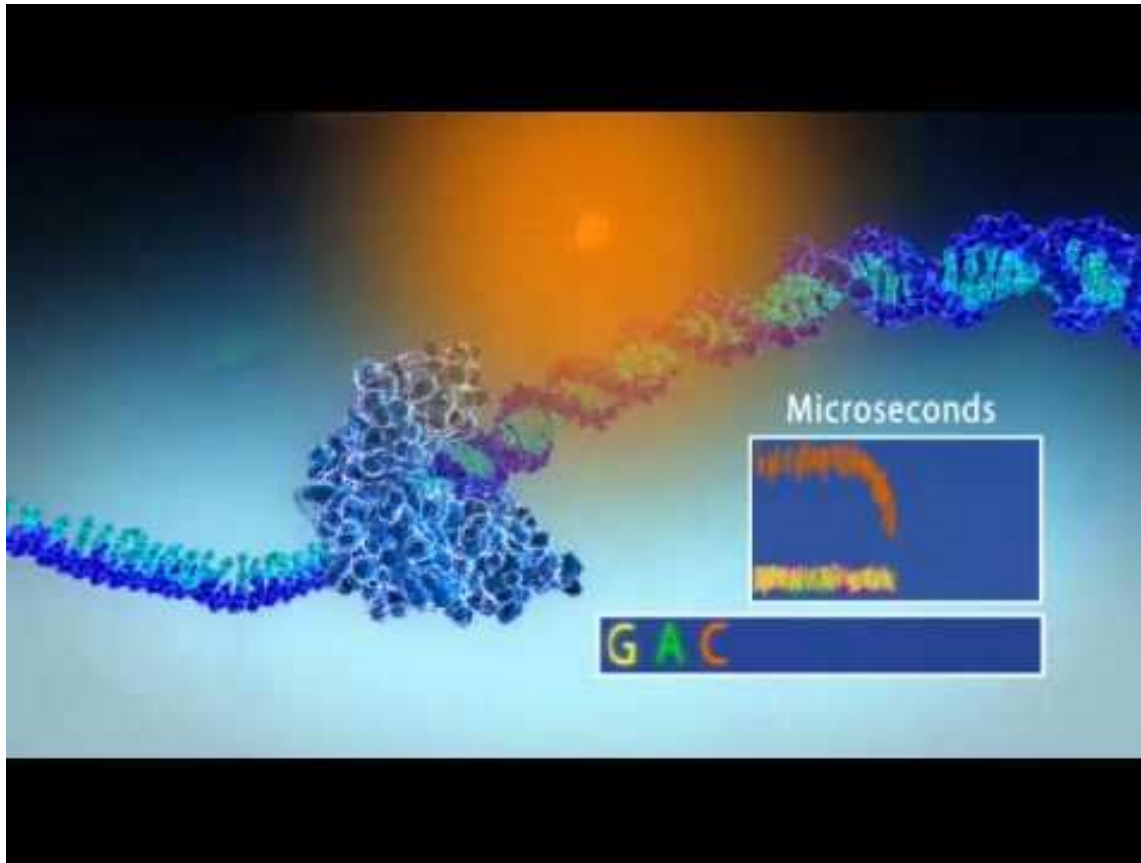
# Phylogenomics

- Of the many ways to sequence genomic data for phylogenetic analyses, how to choose? What methods are available and how do they differ?
  - Whole genome sequencing (WGS)
  - Transcriptome sequencing (RNA-seq)
  - Sequence capture (UCEs; RNA-baits)
  - Amplicon sequencing
  - Restriction-site associated DNA sequencing (RAD-seq)
- It depends on the goal of your study

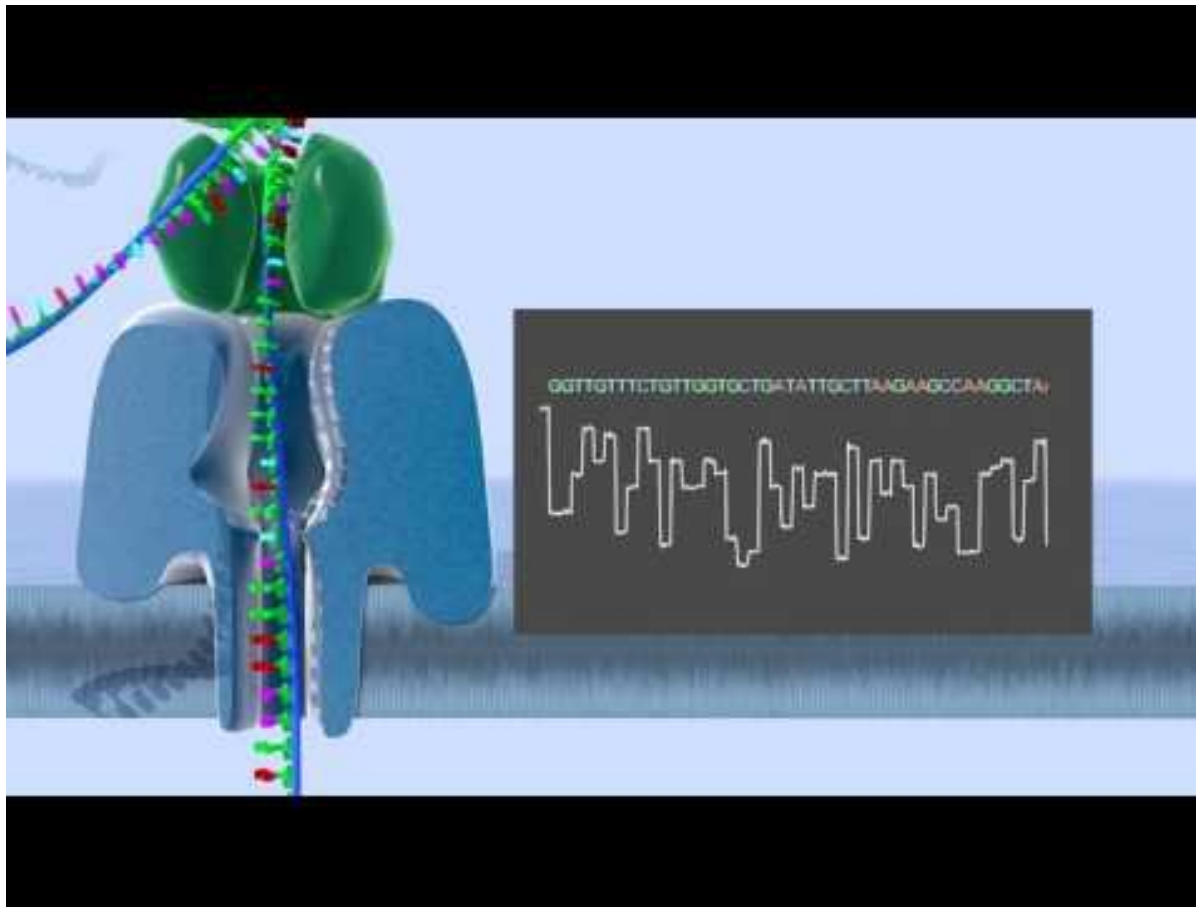




- Illumina short-read technologies
  - <https://www.youtube.com/watch?v=fCd6B5HRaZ8>



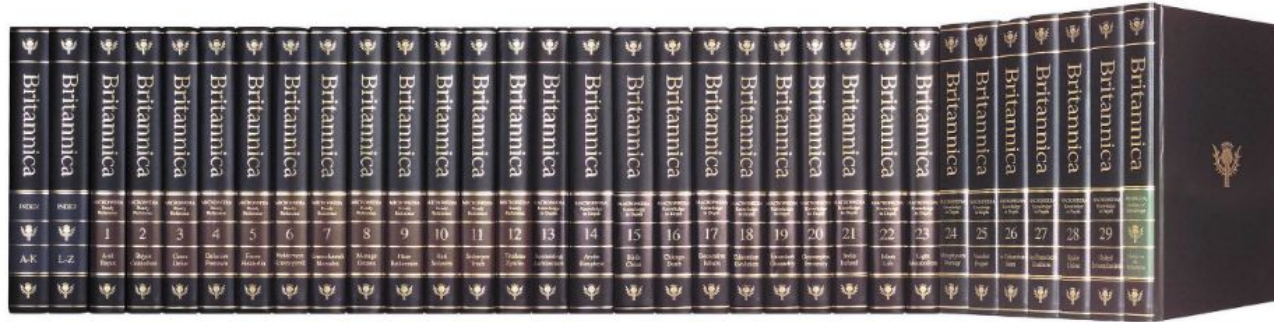
- Pac-Bio SMRT (single molecule real time) sequencing
  - <https://www.youtube.com/watch?v=v8p4ph2MAvI>



- Nano-pore technology
  - <https://www.youtube.com/watch?v=GUb1TZvMWsw>

# Whole genome sequencing vs. ...

- It is easy to sequence small genomes (e.g., *E. coli*; 4.6Mb), but very difficult and expensive to sequence large ones (e.g., *Sequoia sempervirens*; 31,000Mb).
- Studies of organisms with small genomes tend to study the whole genome, while large eukaryotic studies tend to subsample the genome, or sequence it to very low depth (~1x), which can introduce many errors.
- Subsampling methods target fewer regions of the genome and typically analyze loci/genes separately (as gene trees).



Encyclopedia  $\approx$  *Apis* Genome  $\approx$  320 million characters

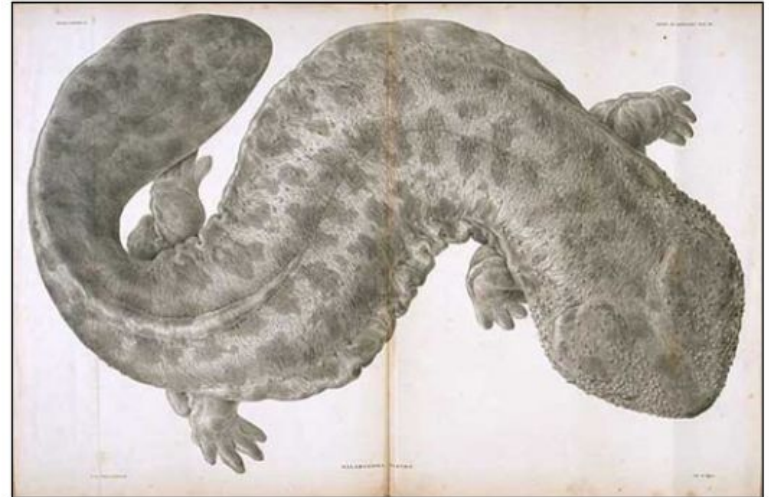
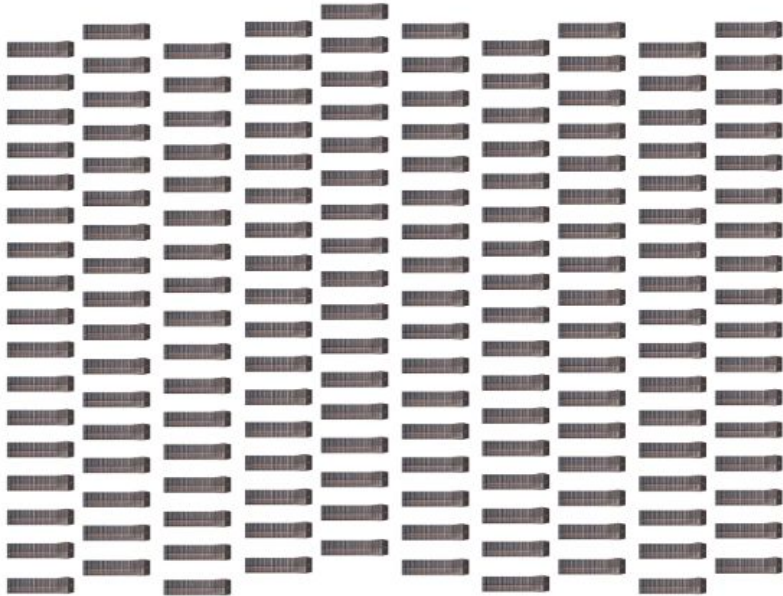




Hominid Genome  $\approx$  3.2 billion bp



Cryptobranchid Genome  $\approx$  55 billion bp

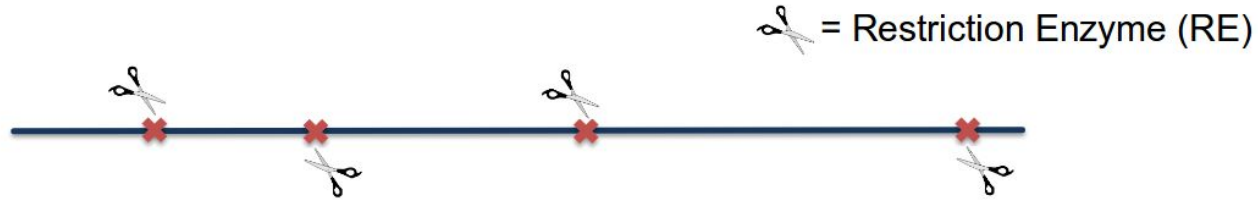


# Restriction-site associated DNA

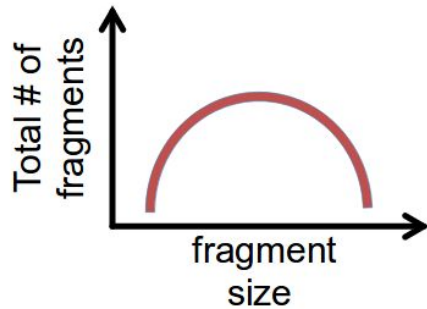
- RAD-sequencing (RAD-seq) and variants (GBS, ddRAD)
- Aims to narrow down the number of sampled regions by targeting a subsample of the genome. In this case, based on the presence of **restriction-enzyme recognition sites**.
- Subsampling targets fewer regions of the genome and typically treats loci/genes as distinct gene trees.



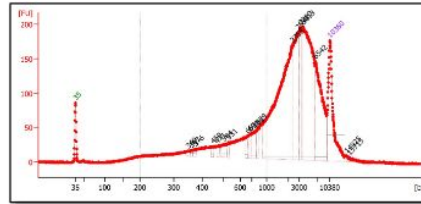
# Restriction enzyme digestion.



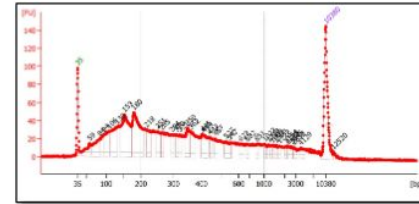
Choice of RE influences the number of fragments produced.



EcoRI



MluCI



In addition to reducing the genome for sequencing, you also need to sequence each RAD locus enough times to accurately identify genetic variation.



Because of the shotgun nature of NGS, there is no expectation of even coverage across all loci.

Insufficient sequencing can lead to missing one allele at a polymorphic locus, or to completely missing a locus.

# Pros and Cons

- RAD-sequencing can target thousands to hundreds of thousands of loci, depending on the enzyme you choose...
- The sequenced loci are short 50-300bp, which may provide little information for inferring gene trees.
- Many new analysis methods make use of only SNPs (single nucleotide polymorphisms), for example, by integrating over all possible gene trees that could produce the SNP (these tend to be slow methods).

```

Individual_1 ATCGA.....TCTTAACGATCCATGC
Individual_1 ATCGA.....TCTTAACGTTCCATGC
Individual_2 ATCGA.....TCTTAACGTTCCATGC
Individual_2 ATCGA.....TCTTAACGTTCCATGC
Individual_3 ATGGA.....TCTTAACGATCCATGC
Individual_3 ATGGA.....TCTTAACGATCCATGC
    
```



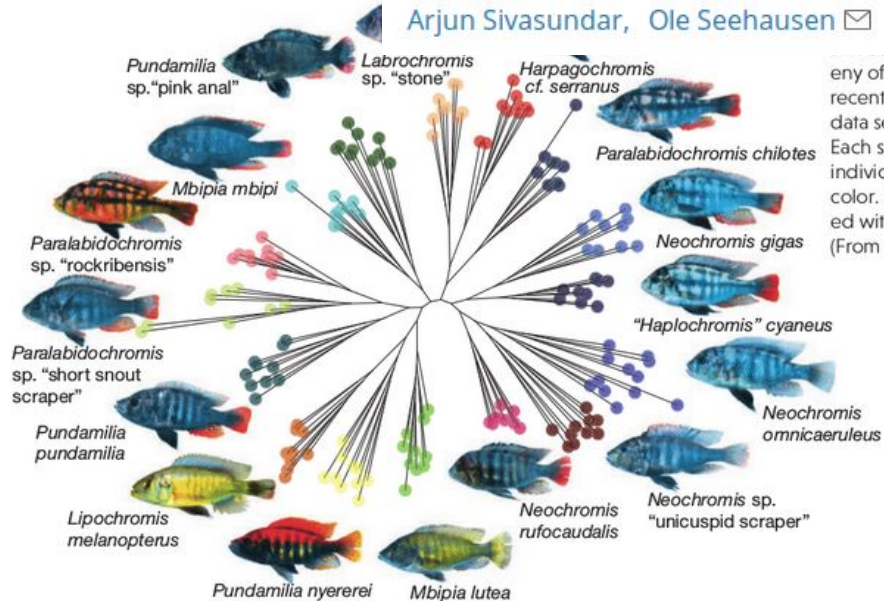
Individual 1	A	T	.	.	.	.	.
Individual 1	T	T	.	.	.	.	.
Individual 2	T	C	.	.	.	.	.
Individual 2	T	C	.	.	.	.	.
Individual 3	A	T	.	.	.	.	.
Individual 3	A	C	.	.	.	.	.

Highly flexible and cheap data type for working across evolutionary scales from very shallow questions to relatively deep-scale questions.

Cichlid radiation includes hundreds of species in just a few thousand years

## Genome-wide RAD sequence data provide unprecedented resolution of species boundaries and relationships in the Lake Victoria cichlid adaptive radiation

Catherine E. Wagner, Irene Keller, Samuel Wittwer, Oliver M. Selz, Salome Mwaiko, Lucie Greuter, Arjun Sivasundar, Ole Seehausen ✉



Each species is represented by several individuals, shown by dots of the same color. Because this tree was constructed without an outgroup, it has no root. (From [36].)

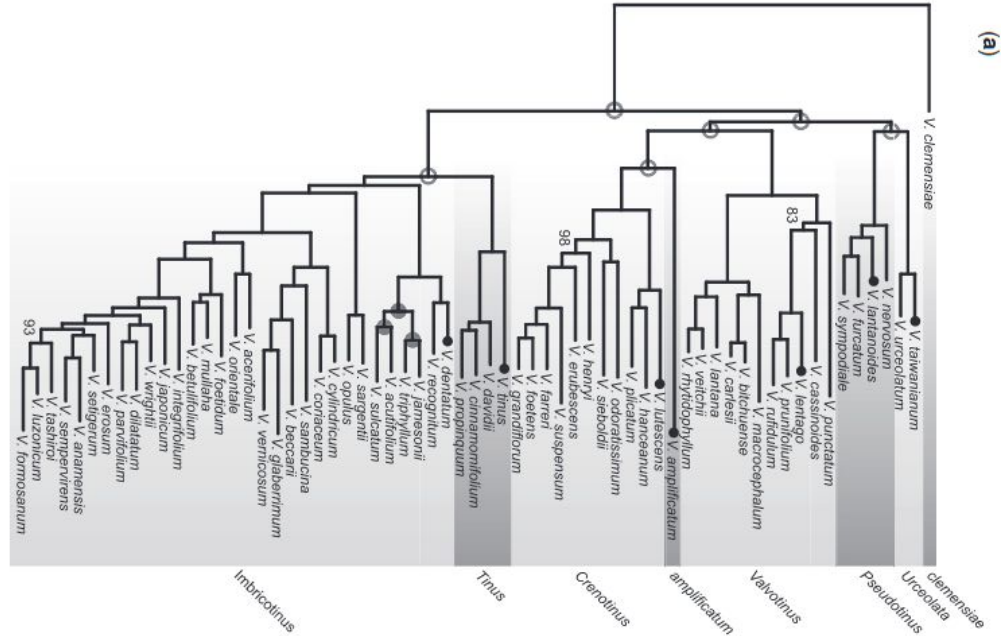
## Misconceptions on Missing Data in RAD-seq Phylogenetics with a Deep-scale Example from Flowering Plants

DEREN A. R. EATON\*, ELIZABETH L. SPRIGGS, BRIAN PARK, AND MICHAEL J. DONOGHUE

Highly flexible and cheap data type for working across evolutionary scales from very shallow questions to relatively deep-scale questions.

Cichlid radiation includes hundreds of species in just a few thousand years

Viburnum radiation includes a few hundred species over >60 million years.



# Common uses of RAD-seq data

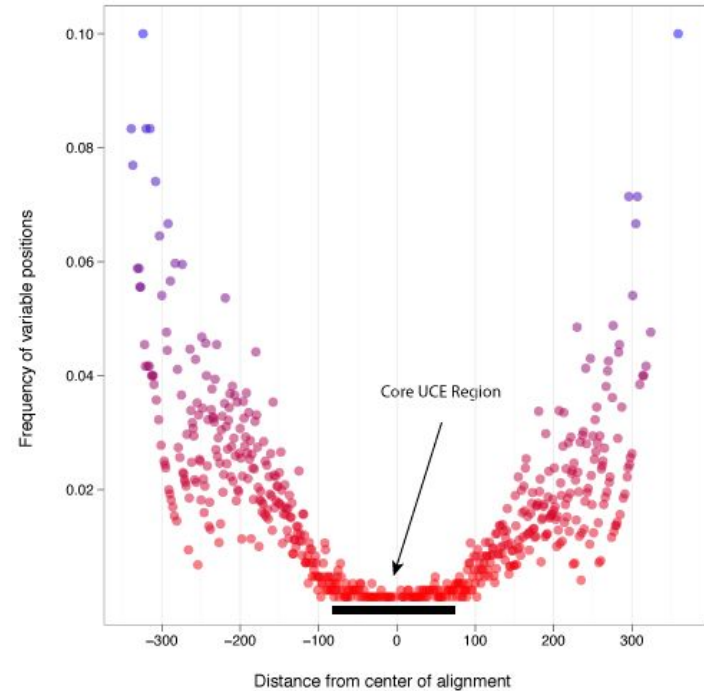
- Phylogenetic inference from hundreds of years to ~100 Ma years.
- QTL mapping, identifying sex-determination loci, loci of large effect.
- Constructing linkage maps (estimating distances between SNPs on chromosomes) based on frequency of recombination throughout the genome.
- Demographic inference: Inferring changes in population sizes and gene flow through time.

# Target capture data

- Similarly to RAD-seq this method aims to sample a reduced portion of the genome. Develops “baits” (RNA probes) which capture certain sequences and let non-matching sequences wash away prior to sequencing.
- Requires prior knowledge to develop the baits, typically from a closely related published genome. Or, uses “universal baits” like [ultra-conserved elements \(UCEs\)](#) , which have been found to change very little across huge amounts of time.
- Targets regions downstream from invariant target site. Can target several overlapping regions to build larger *contigs*, spanning up to several thousand base pairs.
- Very repeatable, and creates reusable data across distant taxa.

# Pros and Cons

- UCEs target up to a few thousand loci.
- Sequenced loci are longer than RAD loci (200-1000bp), although variation is highly heterogenous. Fewer total SNPs but more informative gene trees on average.
- Hugely useful for deeper phylogenetic analyses because many gene trees can be reliably sampled across very distant taxa (e.g., all birds, all vertebrates).





# Ultraconserved elements are novel phylogenomic markers that resolve placental mammal phylogeny when combined with species-tree analysis

John E. McCormack,<sup>1,8</sup> Brant C. Faircloth,<sup>2</sup> Nicholas G. Crawford,<sup>3</sup>  
Patricia Adair Gowaty,<sup>4,5</sup> Robb T. Brumfield,<sup>1,6</sup> and Travis C. Glenn<sup>7</sup>

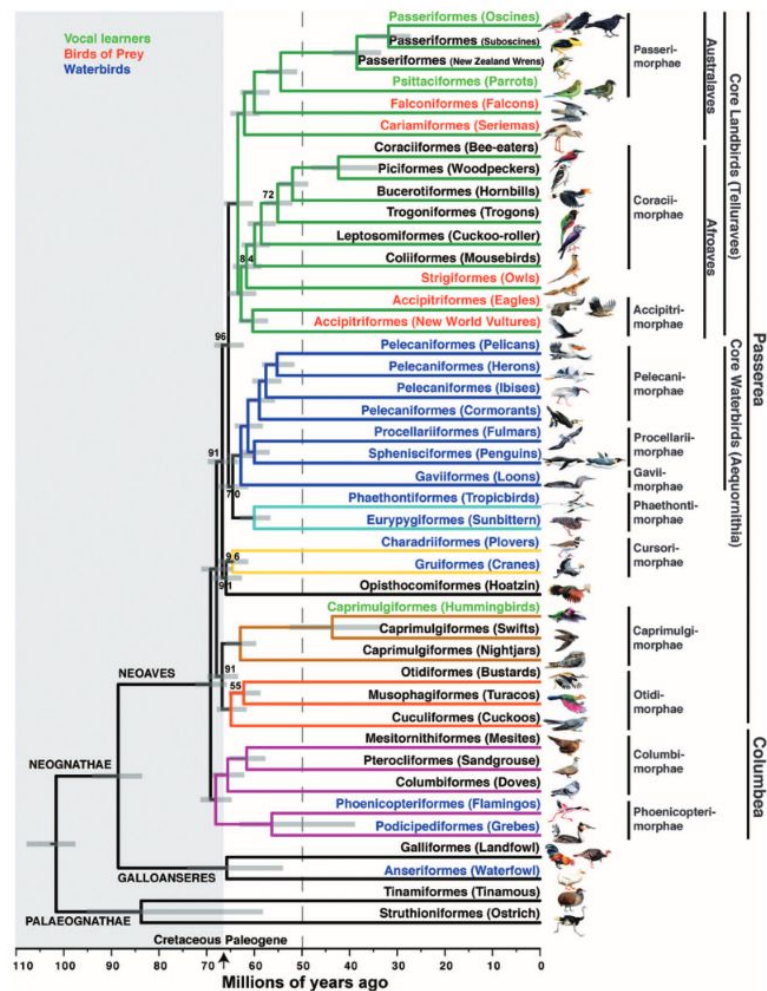
## A Phylogeny of Birds Based on Over 1,500 Loci Collected by Target Enrichment and High-Throughput Sequencing

John E. McCormack , Michael G. Harvey, Brant C. Faircloth, Nicholas G. Crawford, Travis C. Glenn, Robb T. Brumfield

Published: January 29, 2013 • <https://doi.org/10.1371/journal.pone.0054848>

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

Erich D. Jarvis,<sup>1\*</sup> Siavash Mirarab,<sup>2\*</sup> Andre J. Aberer,<sup>3</sup> Bo Li,<sup>4,5,6</sup> Peter Houde,<sup>7</sup> Cai Li,<sup>4,6</sup> Simon Y. W. Ho,<sup>8</sup> Brant C. Faircloth,<sup>9,10</sup> Benoit Nabholz,<sup>11</sup> Jason T. Howard,<sup>1</sup> Alexander Suh,<sup>12</sup> Claudia C. Weber,<sup>12</sup> Rute R. da Fonseca,<sup>6</sup> Jianwen Li,<sup>4</sup> Fang Zhang,<sup>4</sup> Hui Li,<sup>4</sup> Long Zhou,<sup>4</sup> Nitish Narula,<sup>7,13</sup> Liang Liu,<sup>14</sup> Ganesh Ganapathy,<sup>1</sup> Bastien Boussau,<sup>15</sup> Md. Shamsuzzoha Bayzid,<sup>2</sup> Volodymyr Zavidovych,<sup>1</sup> Sankar Subramanian,<sup>16</sup> Toni Gabaldón,<sup>17,18,19</sup> Salvador Capella-Gutiérrez,<sup>17,18</sup> Jaime Huerta-Cepas,<sup>17,18</sup> Bhanu Rekepalli,<sup>20</sup> Kasper Munch,<sup>21</sup> Mikkel Schierup,<sup>21</sup> Bent Lindow,<sup>6</sup> Wesley C. Warren,<sup>22</sup> David Ray,<sup>23,24,25</sup> Richard E. Green,<sup>26</sup> Michael W. Bruford,<sup>27</sup> Xiangjiang Zhan,<sup>27,28</sup> Andrew Dixon,<sup>29</sup> Shengbin Li,<sup>30</sup> Ning Li,<sup>31</sup> Yinhua Huang,<sup>31</sup> Elizabeth P. Derryberry,<sup>32,33</sup> Mads Frost Bertelsen,<sup>34</sup> Frederick H. Sheldon,<sup>33</sup> Robb T. Brumfield,<sup>33</sup> Claudio V. Mello,<sup>35,36</sup> Peter V. Lovell,<sup>35</sup> Morgan Wirthlin,<sup>35</sup> Maria Paula Cruz Schneider,<sup>36,37</sup> Francisco Prosdocimi,<sup>36,38</sup> José Alfredo Samaniego,<sup>6</sup> Amhed Missael Vargas Velazquez,<sup>6</sup> Alonzo Alfaro-Núñez,<sup>6</sup> Paula F. Campos,<sup>6</sup> Bent Petersen,<sup>39</sup> Thomas Sicheritz-Ponten,<sup>39</sup> An Pas,<sup>40</sup> Tom Bailey,<sup>41</sup> Paul Scofield,<sup>42</sup> Michael Bunce,<sup>43</sup> David M. Lambert,<sup>16</sup> Qi Zhou,<sup>44</sup> Polina Perelman,<sup>45,46</sup> Amy C. Driskell,<sup>47</sup> Beth Shapiro,<sup>26</sup> Zijun Xiong,<sup>4</sup> Yongli Zeng,<sup>4</sup> Shiping Liu,<sup>4</sup> Zhenyu Li,<sup>4</sup> Binghang Liu,<sup>4</sup> Kui Wu,<sup>4</sup> Jin Xiao,<sup>4</sup> Xiong Yinqi,<sup>4</sup> Qiumei Zheng,<sup>4</sup> Yong Zhang,<sup>4</sup> Huanming Yang,<sup>48</sup> Jian Wang,<sup>48</sup> Linnea Smeds,<sup>12</sup> Frank E. Rheindt,<sup>49</sup> Michael Braun,<sup>50</sup> Jon Fjeldsa,<sup>51</sup> Ludovic Orlando,<sup>6</sup> F. Keith Barker,<sup>52</sup> Knud Andreas Jønsson,<sup>51,53,54</sup> Warren Johnson,<sup>55</sup> Klaus-Peter Koepfli,<sup>56</sup> Stephen O'Brien,<sup>57,58</sup> David Haussler,<sup>59</sup> Oliver A. Ryder,<sup>60</sup> Carsten Rahbek,<sup>51,54</sup> Eske Willerslev,<sup>6</sup> Gary R. Graves,<sup>51,61</sup> Travis C. Glenn,<sup>62</sup> John McCormack,<sup>63</sup> Dave Burt,<sup>64</sup> Hans Ellegren,<sup>12</sup> Per Alström,<sup>65,66</sup> Scott V. Edwards,<sup>67</sup> Alexandros Stamatakis,<sup>3,68</sup> David P. Mindell,<sup>69</sup> Joel Cracraft,<sup>70</sup> Edward L. Braun,<sup>71</sup> Tandy Warnow,<sup>2,72</sup> Wang Jun,<sup>48,73,74,75,76</sup> M. Thomas P. Gilbert,<sup>6,43</sup> Guojie Zhang<sup>4,77</sup>



## A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing

Richard O. Prum<sup>1,2\*</sup>, Jacob S. Berv<sup>3\*</sup>, Alex Dornburg<sup>1,2,4</sup>, Daniel J. Field<sup>2,5</sup>, Jeffrey P. Townsend<sup>1,6</sup>, Emily Moriarty Lemmon<sup>7</sup> & Alan R. Lemmon<sup>8</sup>

Although reconstruction of the phylogeny of living birds has progressed tremendously in the last decade, the evolutionary history of Neoaves—a clade that encompasses nearly all living bird species—remains the greatest unresolved challenge in dinosaur systematics. Here we investigate avian phylogeny with an unprecedented scale of data: >390,000 bases of genomic sequence data from each of 198 species of living birds, representing all major avian lineages, and two crocodilian outgroups. Sequence data were collected using anchored hybrid enrichment, yielding 259 nuclear loci with an average length of 1,523 bases for a total data set of over  $7.8 \times 10^7$  bases. Bayesian and maximum likelihood analyses yielded highly supported and nearly identical phylogenetic trees for all major avian lineages. Five major clades form successive sister groups to the rest of Neoaves: (1) a clade including nightjars, other caprimulgiforms, swifts, and hummingbirds; (2) a clade uniting cuckoos, bustards, and turacos with pigeons, mesites, and sandgrouse; (3) cranes and their relatives; (4) a comprehensive waterbird clade, including all diving, wading, and shorebirds; and (5) a comprehensive landbird clade with the enigmatic hoatzin (*Opisthocomus hoazin*) as the sister group to the rest. Neither of the two main, recently proposed Neoavian clades—Columbea and Passerea<sup>1</sup>—were supported as monophyletic. The results of our divergence time analyses are congruent with the palaeontological record, supporting a major radiation of crown birds in the wake of the

It has long been recognized that phylogenetic confidence depends not only on the number of characters analysed and their rate of evolution, but also on the number and relationships of the taxa sampled relative to the nodes of interest<sup>9–11</sup>. Theory predicts that sampling a single taxon that diverges close to a node of interest will have a far greater effect on phylogenetic resolution than will adding more characters<sup>11</sup>. Despite using an alignment of >40 million base pairs, sparse sampling of 48 species in the recent avian genomic analysis may not have been sufficient to confidently resolve the deep divergences among major lineages of Neoaves. Thus, expanded taxon sampling is required to test the monophyly of neoavian clades, and to further resolve the phylogenetic relationships within Neoaves.

Here, we present a phylogenetic analysis of 198 bird species and 2 crocodilians (Supplementary Table 1) based on loci captured using anchored enrichment<sup>12</sup>. Our sample includes species of 122 avian families in all 40 extant avian orders<sup>2</sup>, with denser representation of non-oscine birds (108 families) than of oscine songbirds (14 families). Effort was made to include taxa that would break up long phylogenetic branches, and provide the highest likelihood of resolving short internodes at the base of Neoaves<sup>11</sup>. We also sampled multiple species within groups whose monophyly or phylogenetic interrelationships have been controversial—that is, tinamous, nightjars, hummingbirds, turacos, cuckoos, pigeons, sandgrouse, mesites, rails, storm petrels, petrels, storks, herons, hawks, hornbills, mousebirds, trogons, king-

# Transcriptome sequencing

- Sequence only the transcribed portion of the genome by extracting RNA, and from that making cDNA.
- Con: Extracting and storing RNA can be very difficult.
- Con: Assembling transcripts can be difficult due to differential splicing.
- Con: Markers are too conserved, too little variation.
- Pro: Phylogenetic markers may also provide insight into the evolution of functional differences, or convergence.
- Pro: Markers are conserved across very deep time-scales.
- Pro: Paralogs (duplicated genes) are typically easier to detect.

# Summary of phylogenomic methods

- Whole genome sequencing is difficult for organisms with large genomes where it is preferable to sequence fewer regions to high/reliable depth. Whole genomes must be split into non-recombining loci, or analyzed by sliding window.
- RAD-seq targets regions on the basis of restriction-sites, and can generate many thousands of short markers. Used for younger clade analyses (0-80Ma).
- UCEs target regions on the basis of designed baits chosen to target known conserved regions. Generates a few hundred or thousand markers of longer length and informativeness than RAD. Very useful for deep scale phylogenetic analyses (20-300 Ma).

# Reading assignments

- No textbook reading for Thursday
- Article on courseworks:
  - “Butterfly genome reveals promiscuous exchange of mimicry adaptations among species”. The Heliconius Genome Consortium (2012).