

# Inference of the tree of life with applications to plants

Siavash  
Mirarab

Electrical and  
Computer Engineering

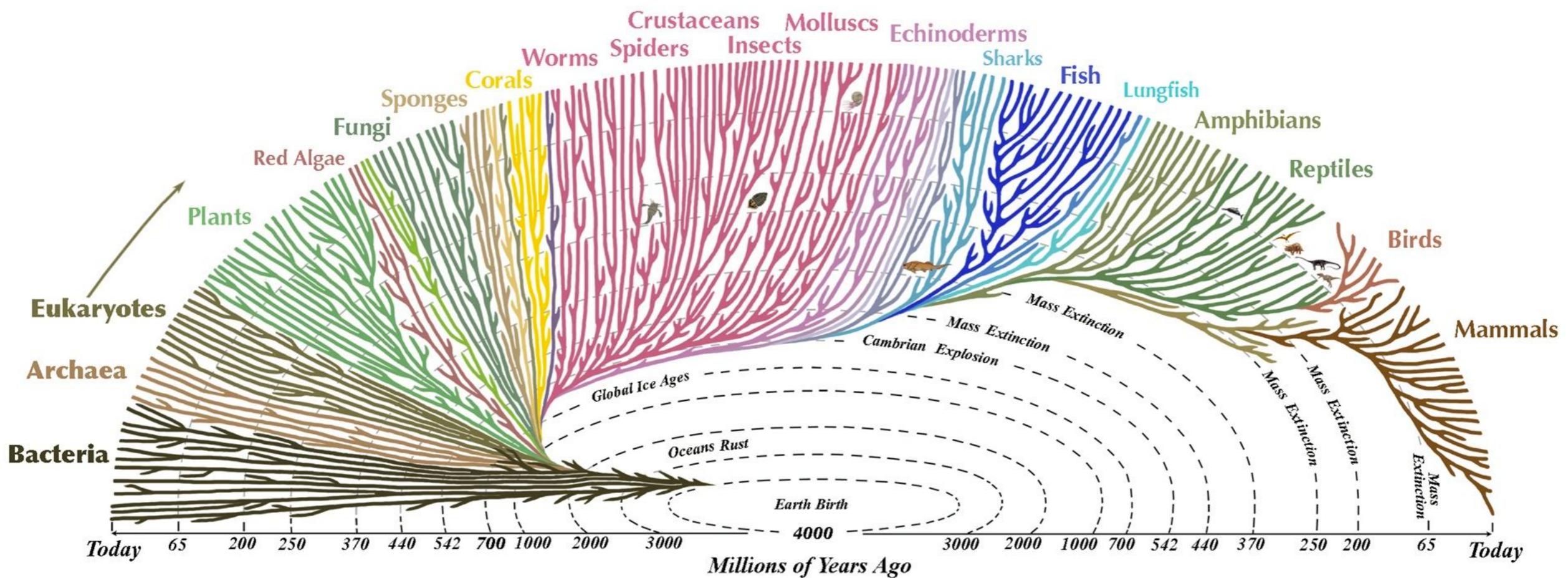
University of  
California, San Diego



# Tree of life

***“Nothing in biology makes sense except in the light of evolution.”***

Dobzhansky, 1973

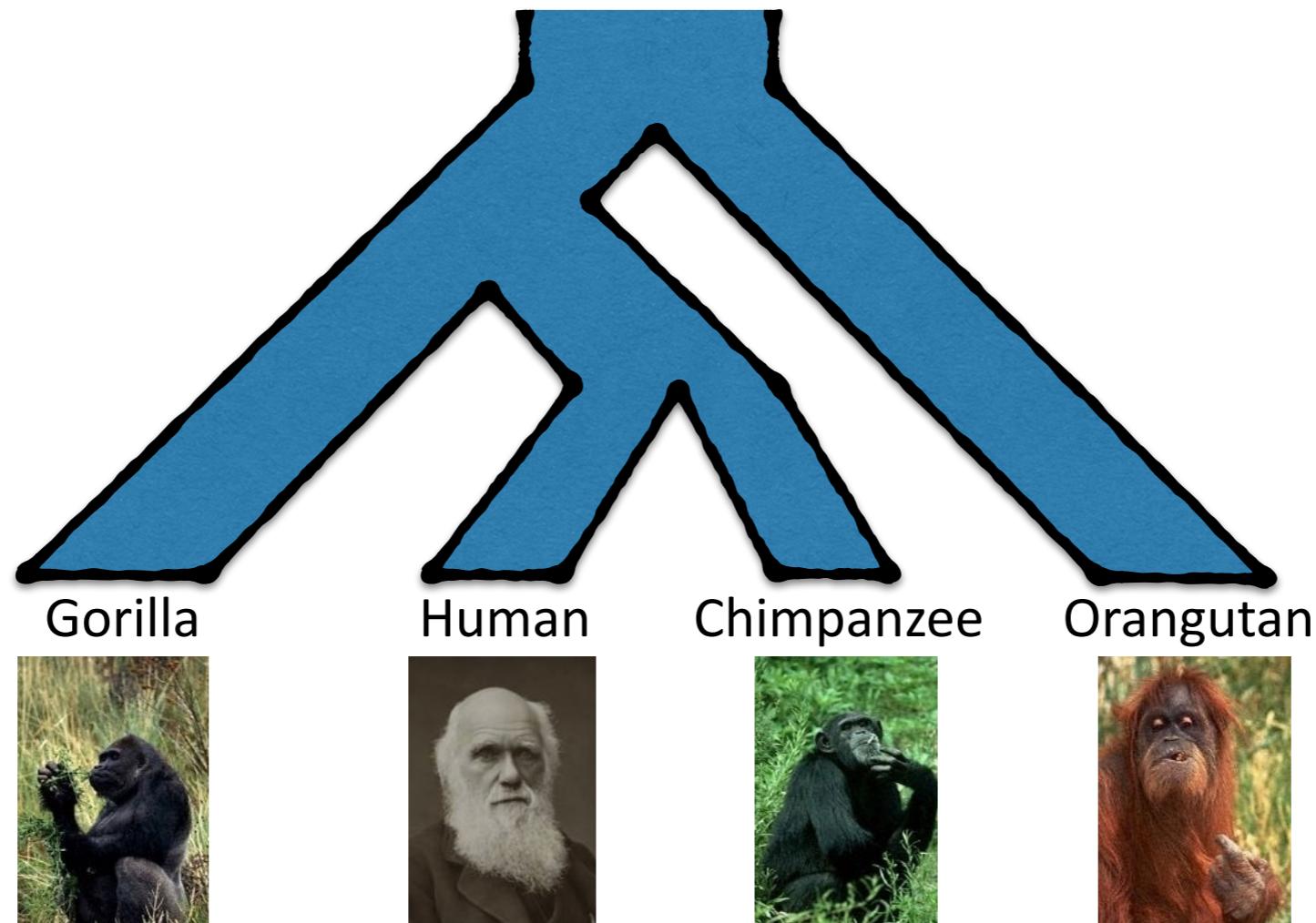
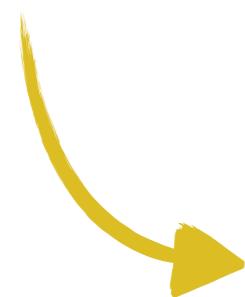


source: <http://www.evogeneao.com/>

© Leonard Eisenberg 2008  
evogeneao.com

# Phylogeny

Leaves are  
Species  
(a.k.a. **taxa**)



The branching structure shows evolutionary relationships

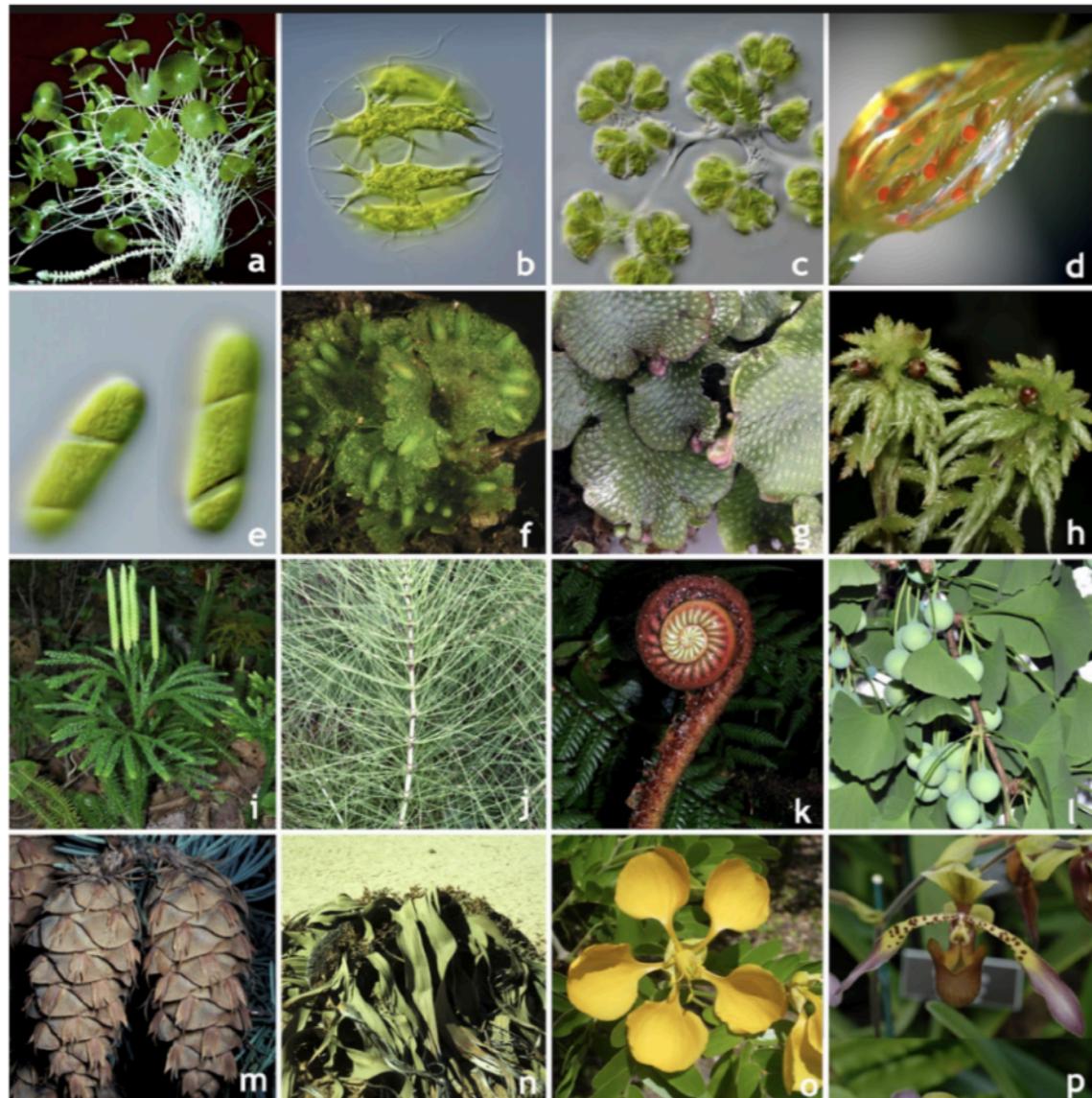
# Why useful?

It ENABLEs downstream  
analyses

# Order of biological innovation

Zym.: algae-like

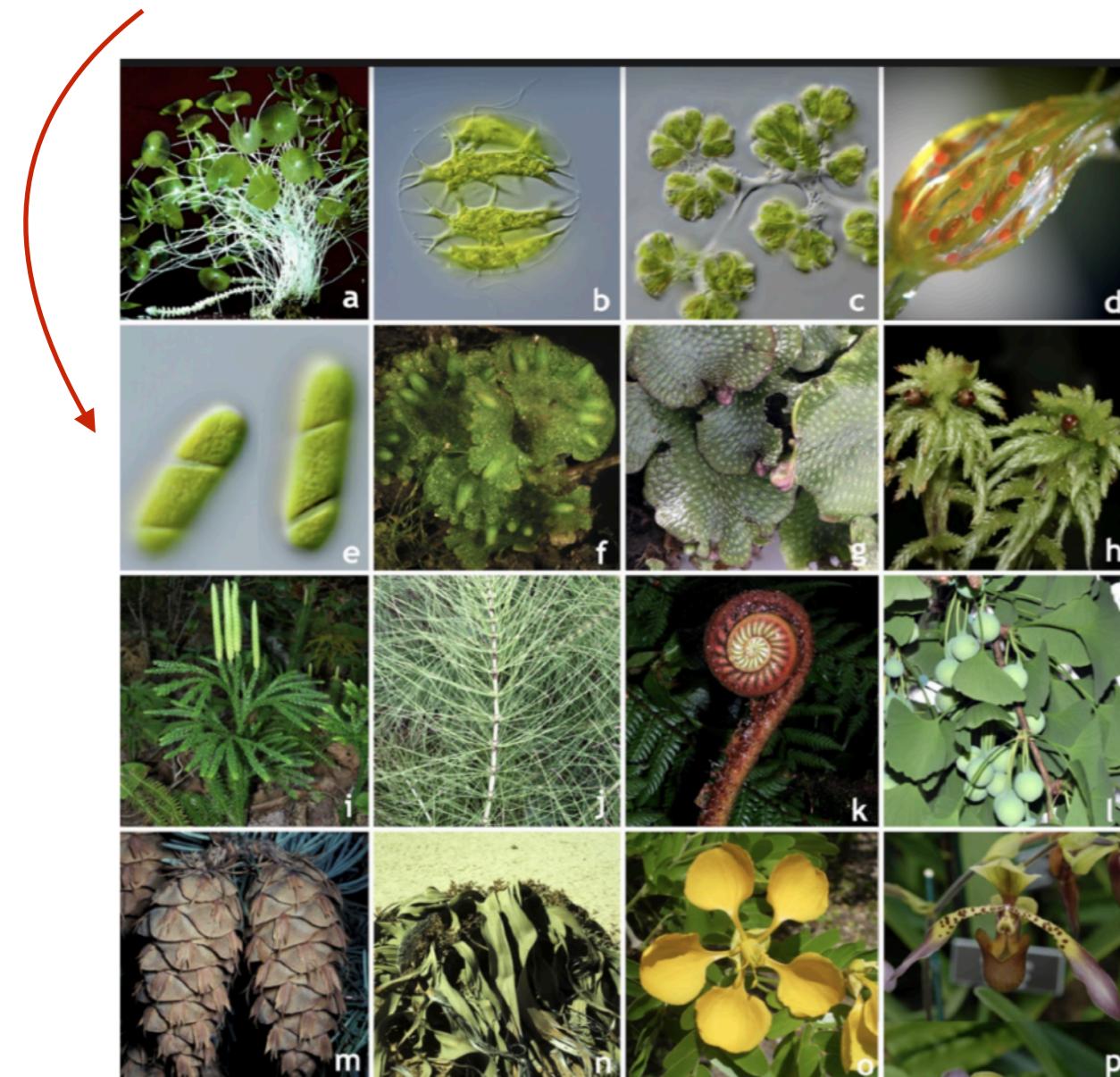
Chara: land plant-like



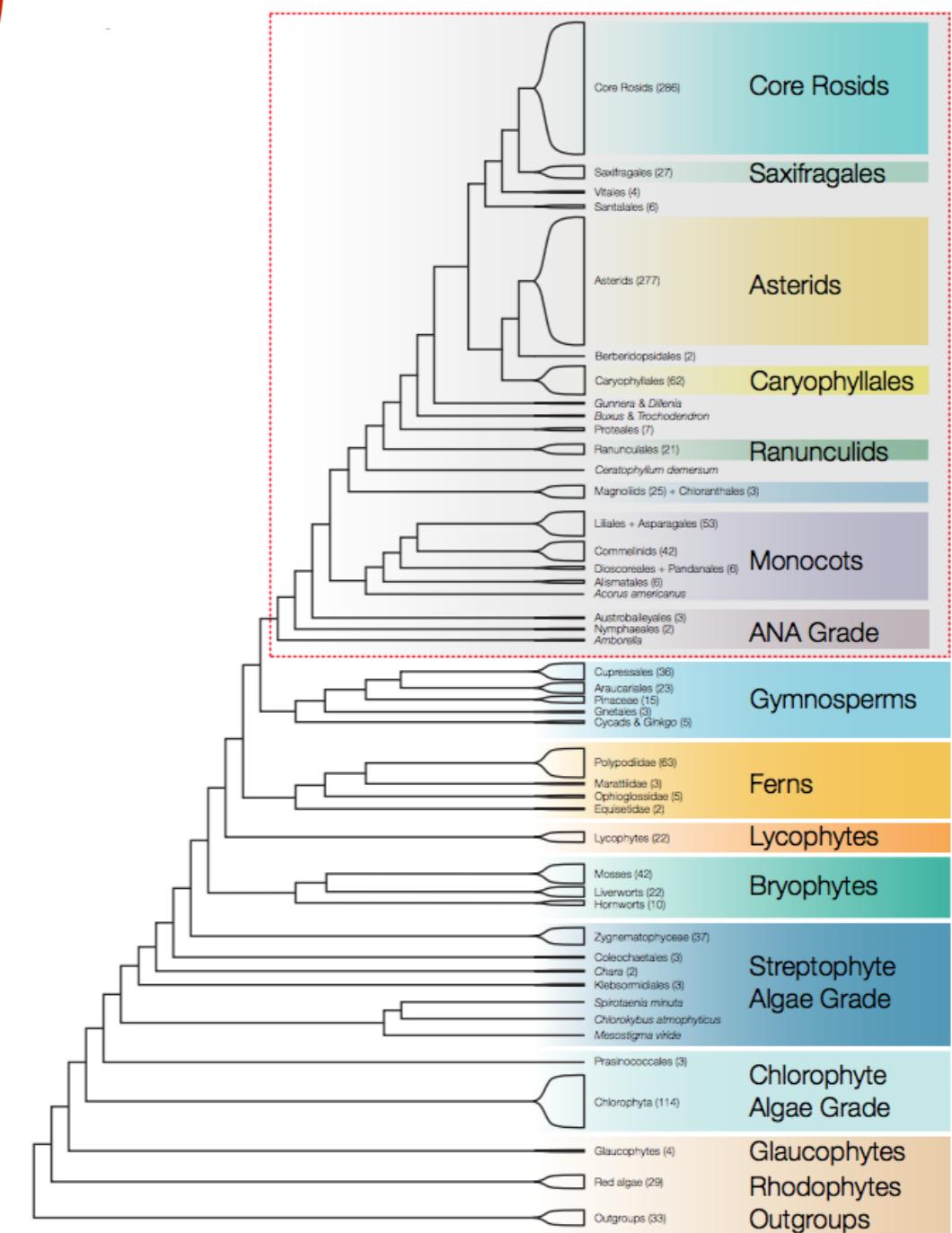
1KP project [unpublished results]

# Order of biological innovation

Zym.: algae-like



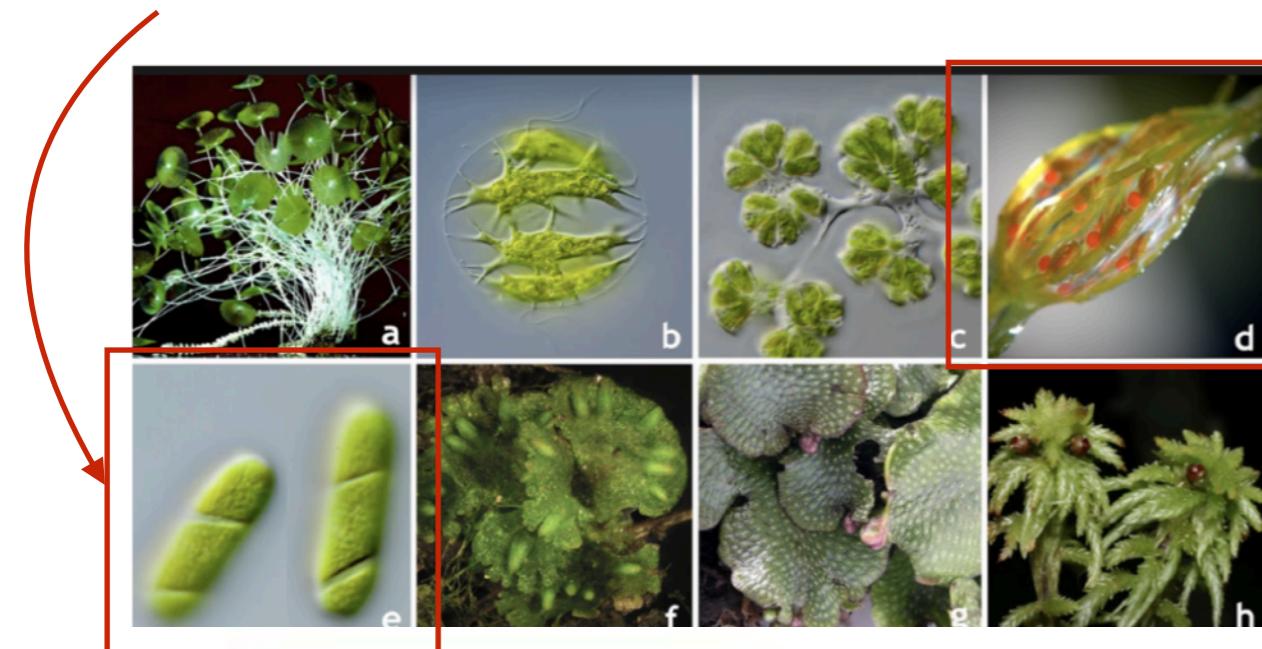
Chara: land plant-like



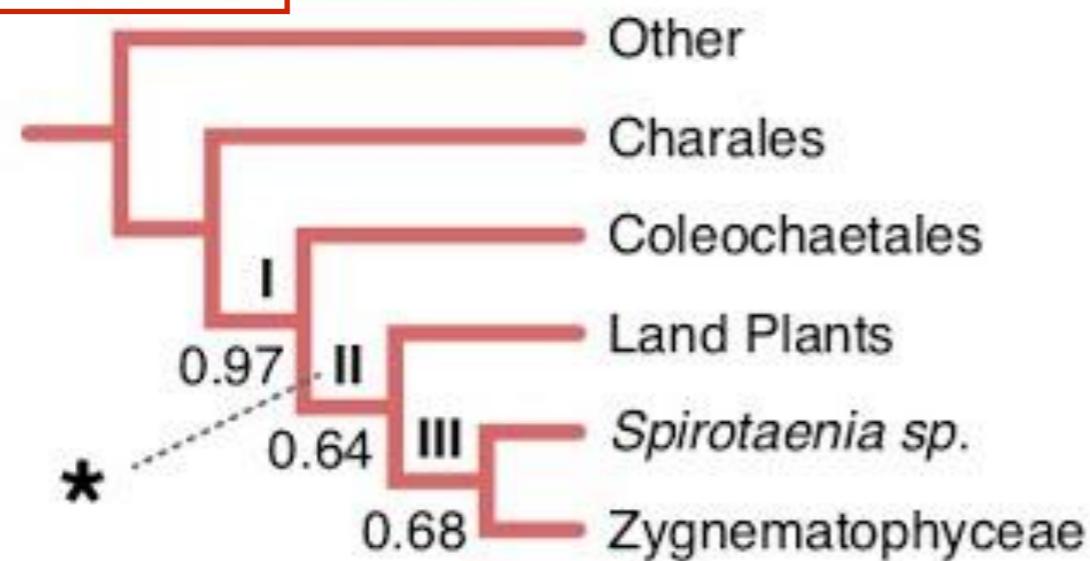
1KP project [unpublished results]

# Order of biological innovation

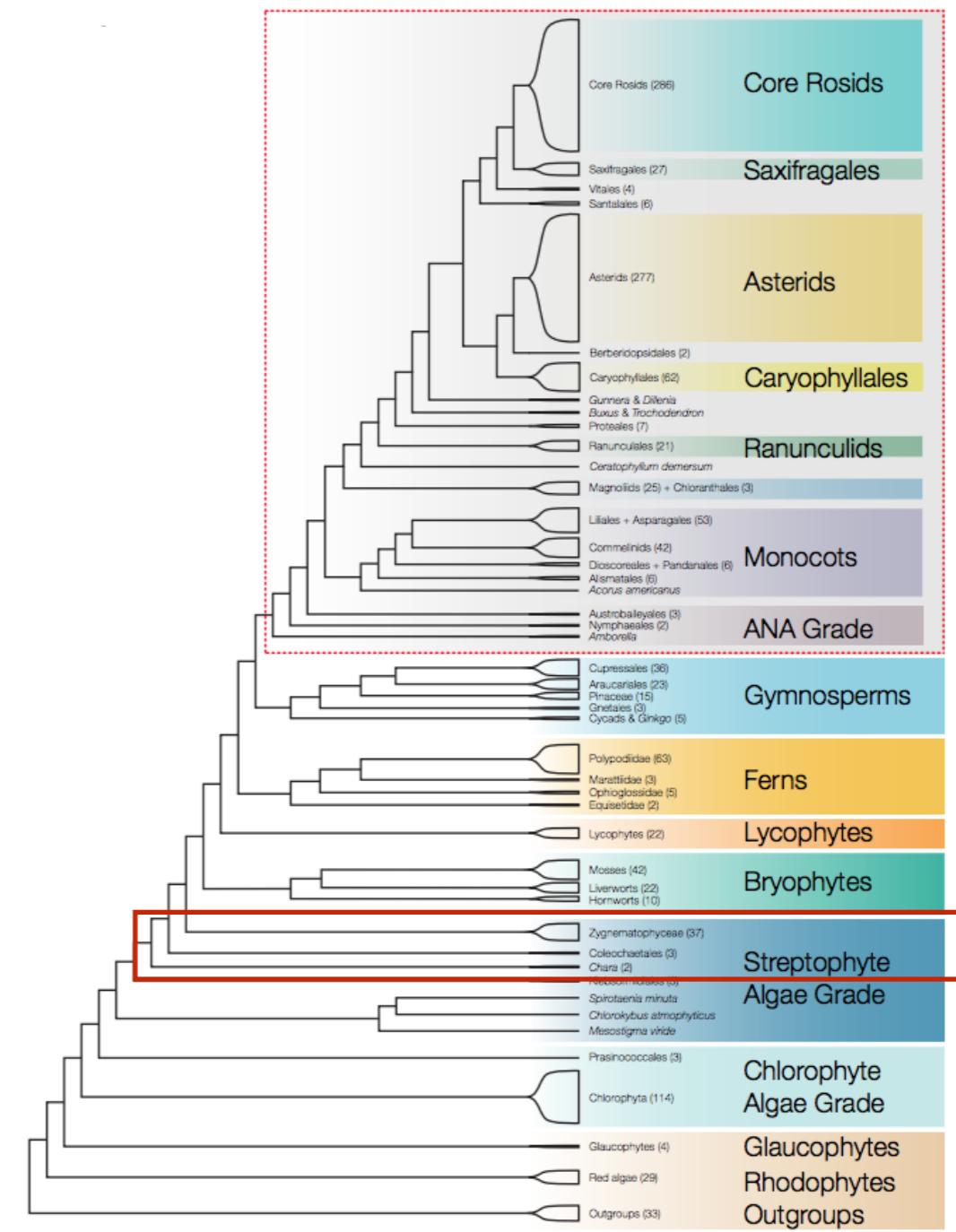
Zym.: algae-like



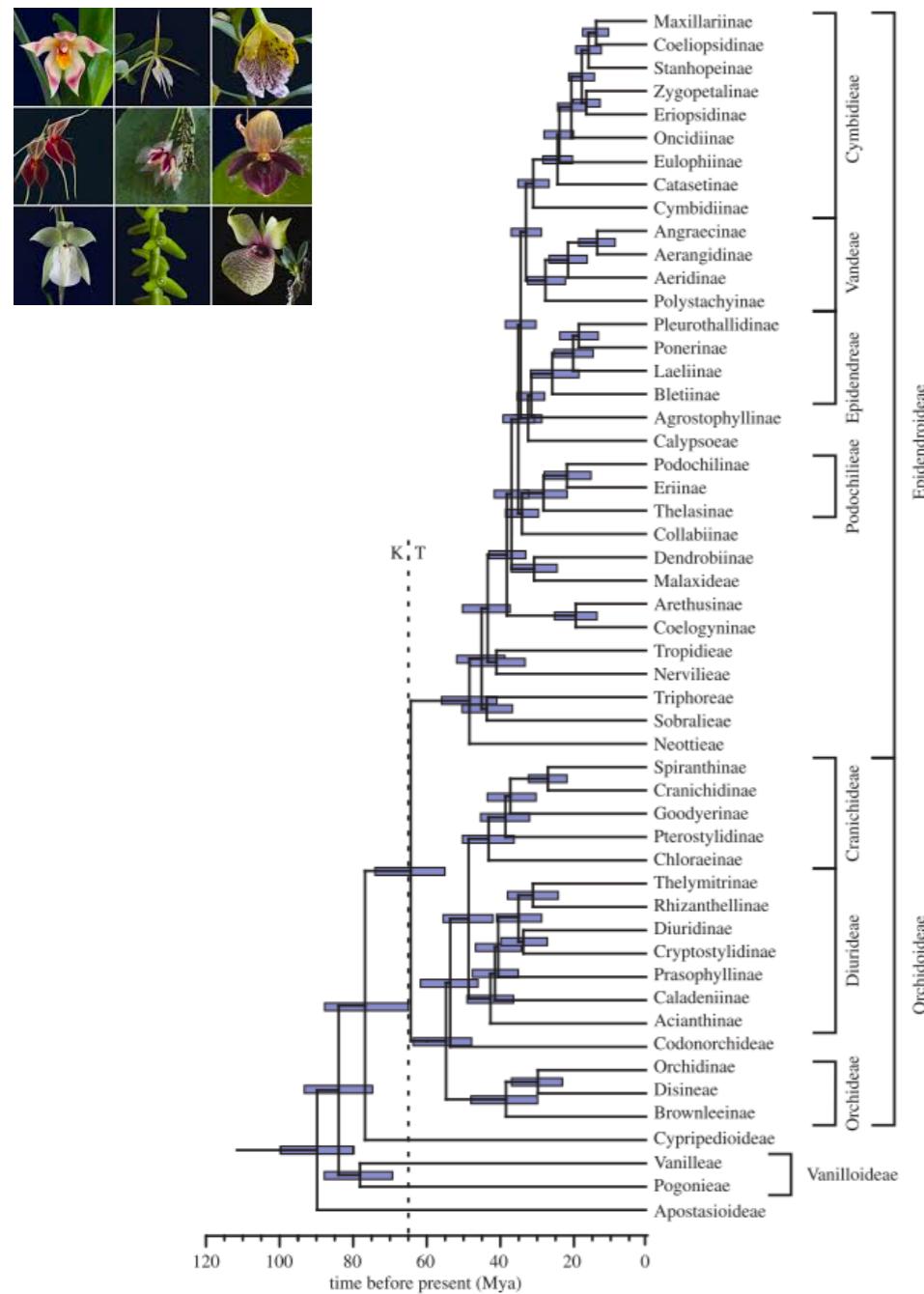
Chara: land plant-like



1KP project [unpublished results]

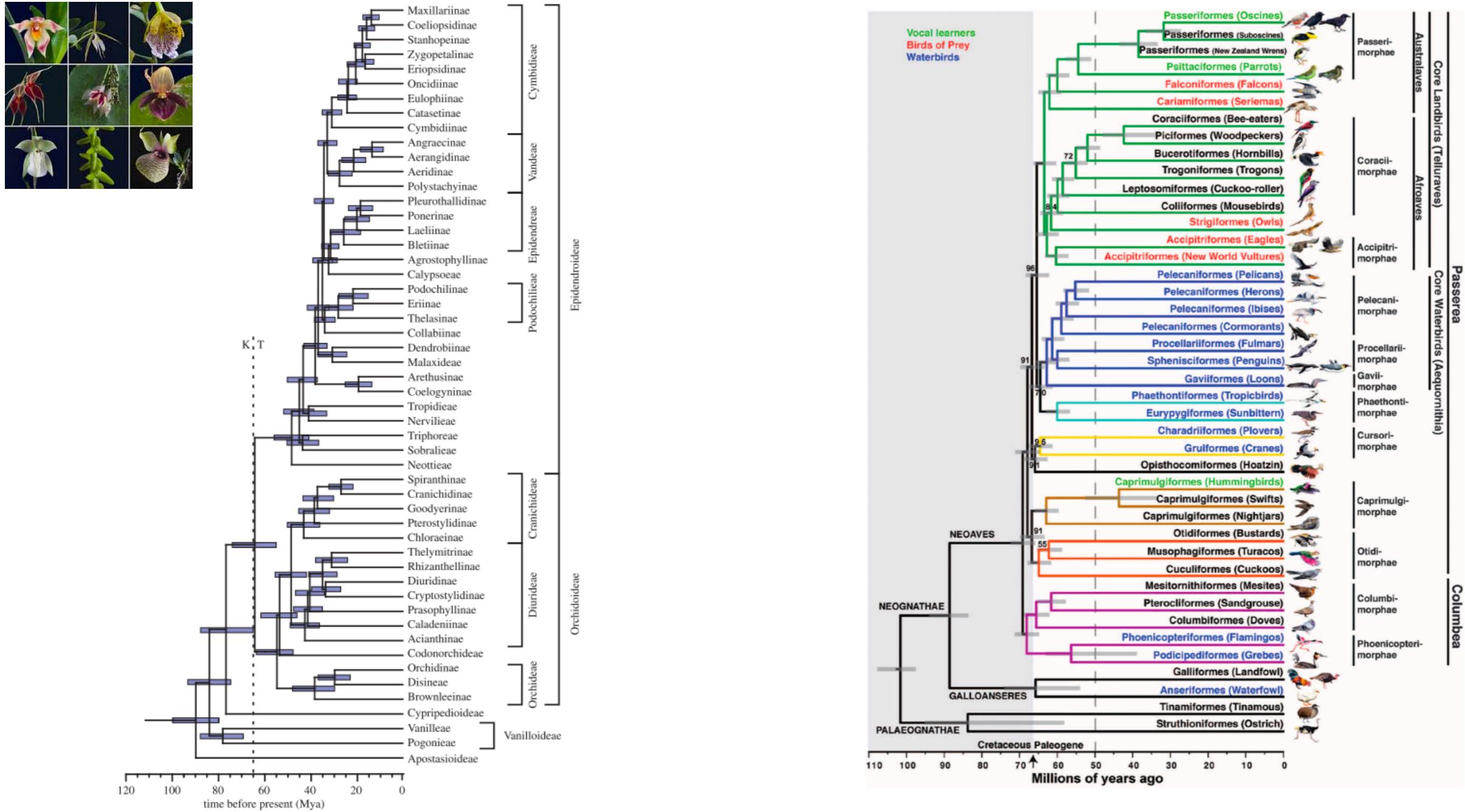


# Understand timing of events



Givnish *et al*, Royal Society B, 2015

# Understand timing of events



Givnish *et al*, Royal Society B, 2015

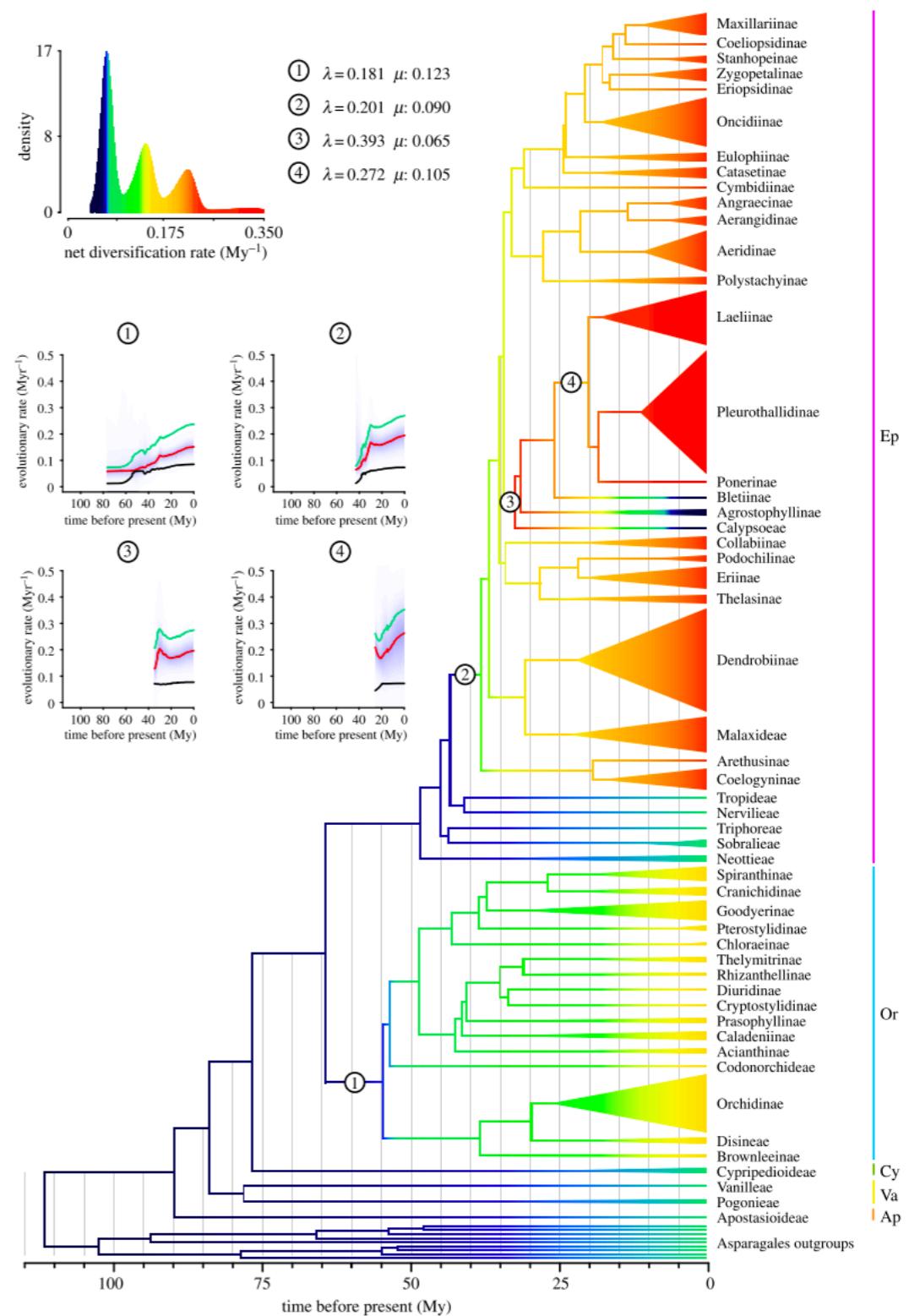
Jarvis, Mirarab *et al*, Science, 2014

# Understand drivers of diversity

- Orchids have more species than mammals, birds and reptiles combined

# Understand drivers of diversity

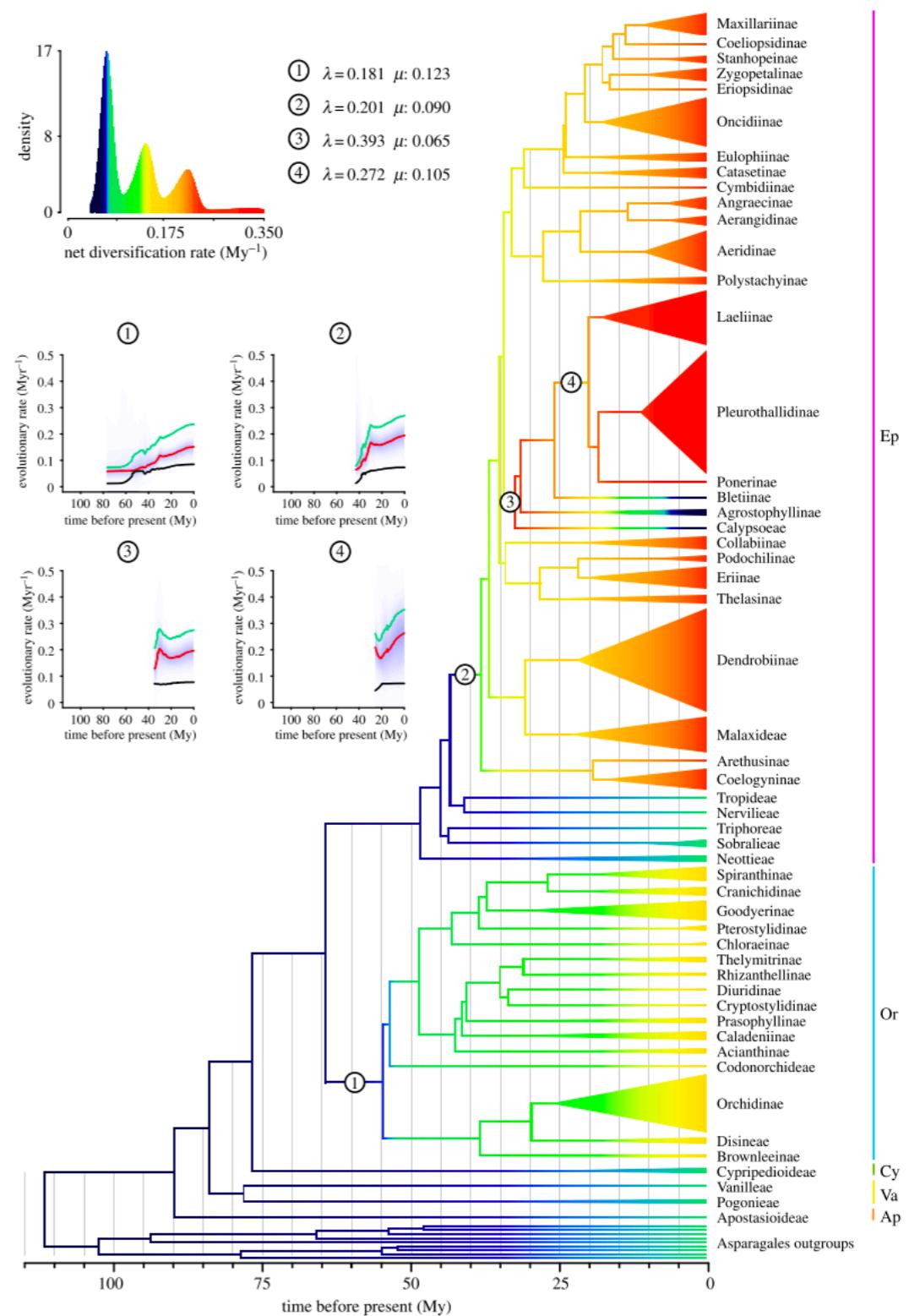
- Orchids have more species than mammals, birds and reptiles combined
- Three significant shifts of diversification rate



Givnish et al, Royal Society B, 2015

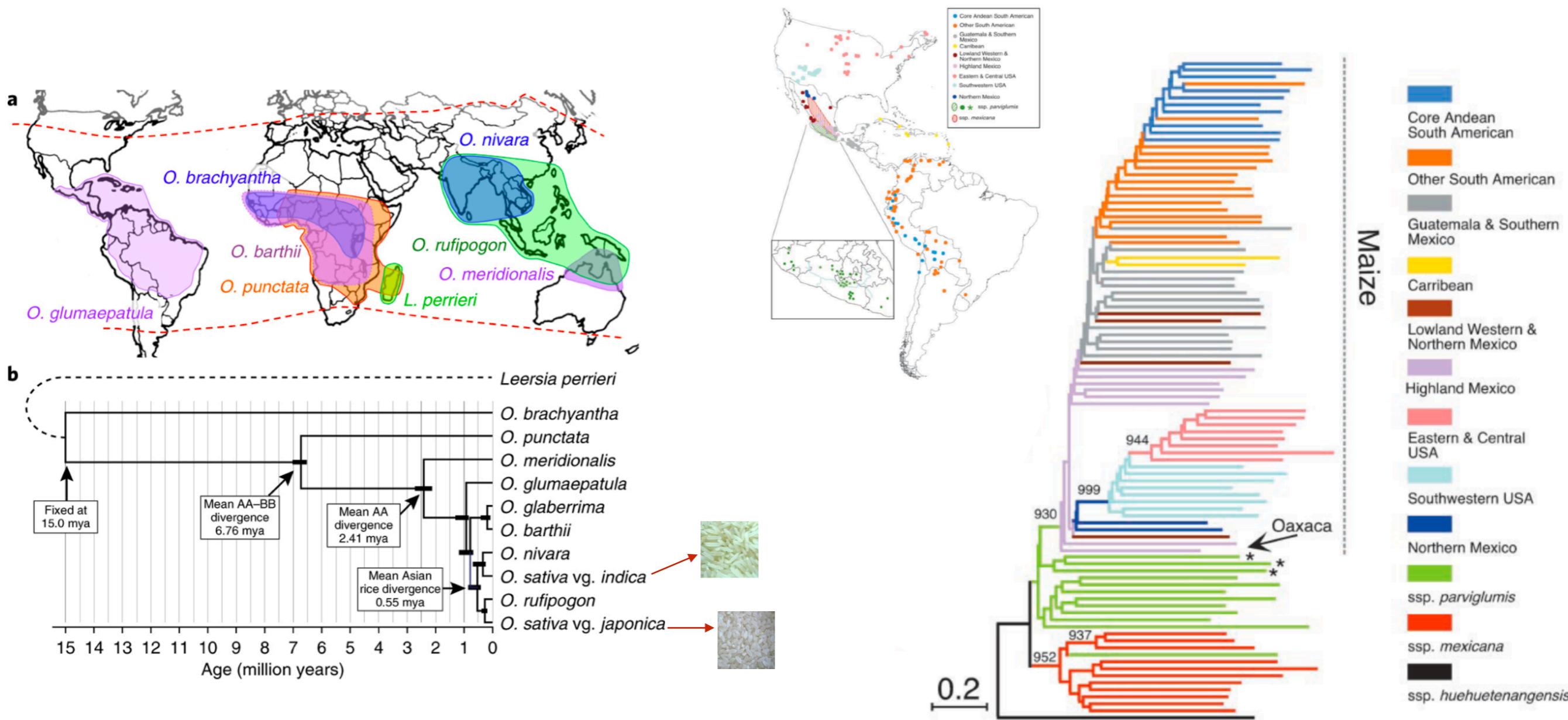
# Understand drivers of diversity

- Orchids have more species than mammals, birds and reptiles combined
- Three significant shifts of diversification rate
- These coincide with evolution of pollinia, the epiphytic habit, CAM photosynthesis, tropical distribution, and pollination via butterflies or bees



Givnish et al, Royal Society B, 2015

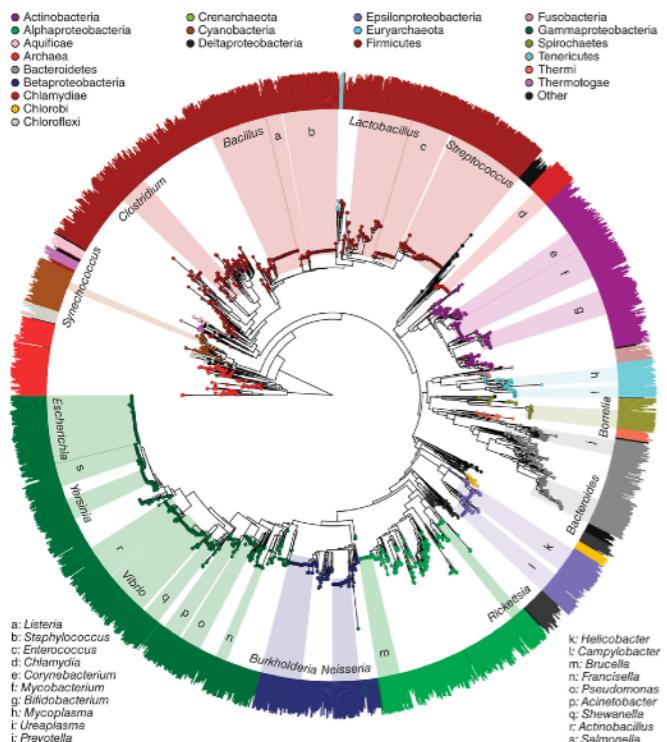
# Domestication history



Stein et al, Nature genetics, 2018

Matsuoka et al, PNAS, 2002

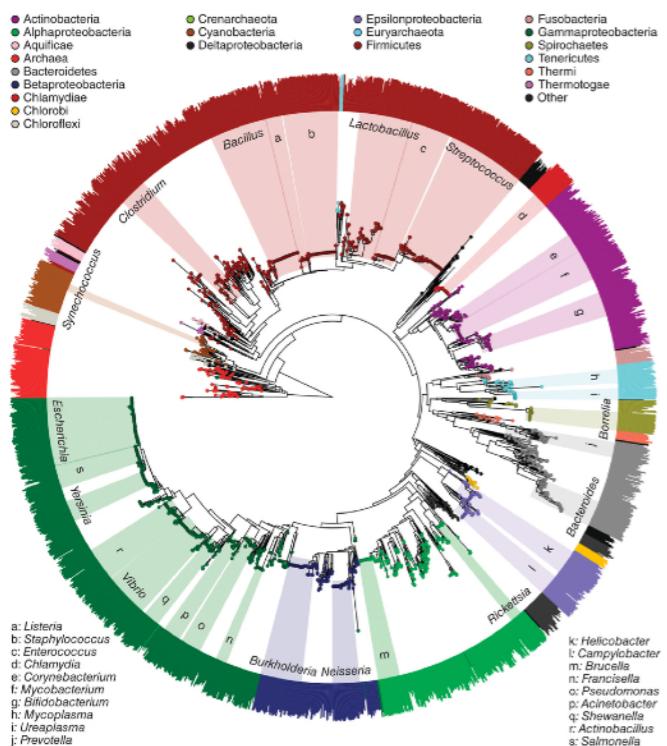
# Biomedical Applications



source: Langille et al., Nature, 2013

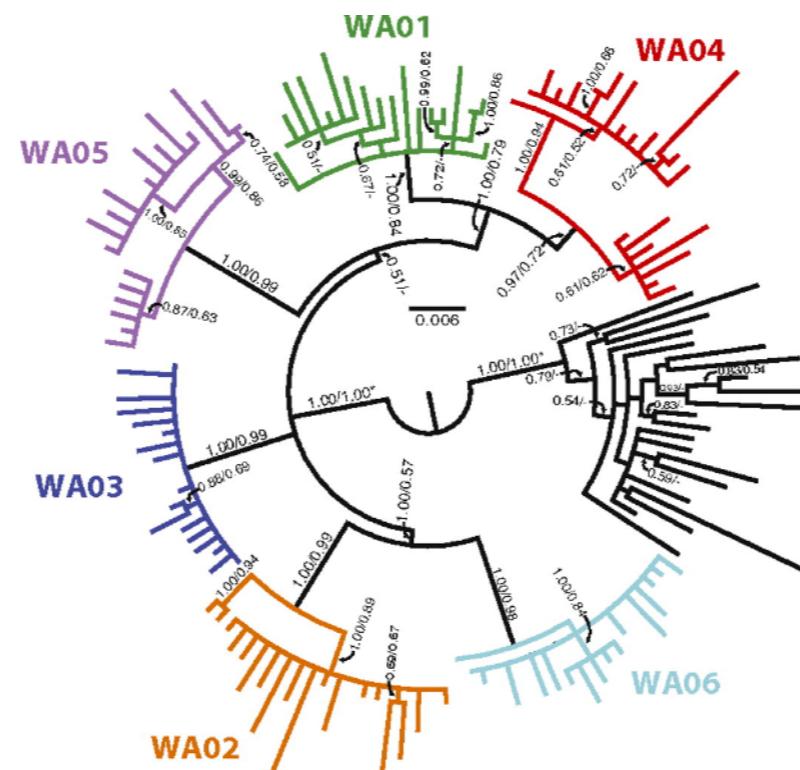
**Microbiome**

# Biomedical Applications



source: Langille et al., Nature, 2013

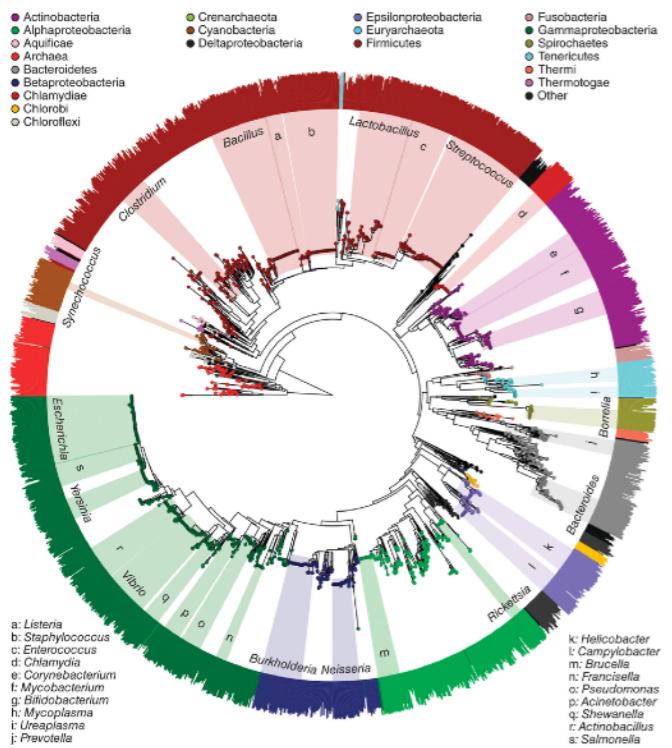
**Microbiome**



source: Scaduto et al., PNAS, 2010

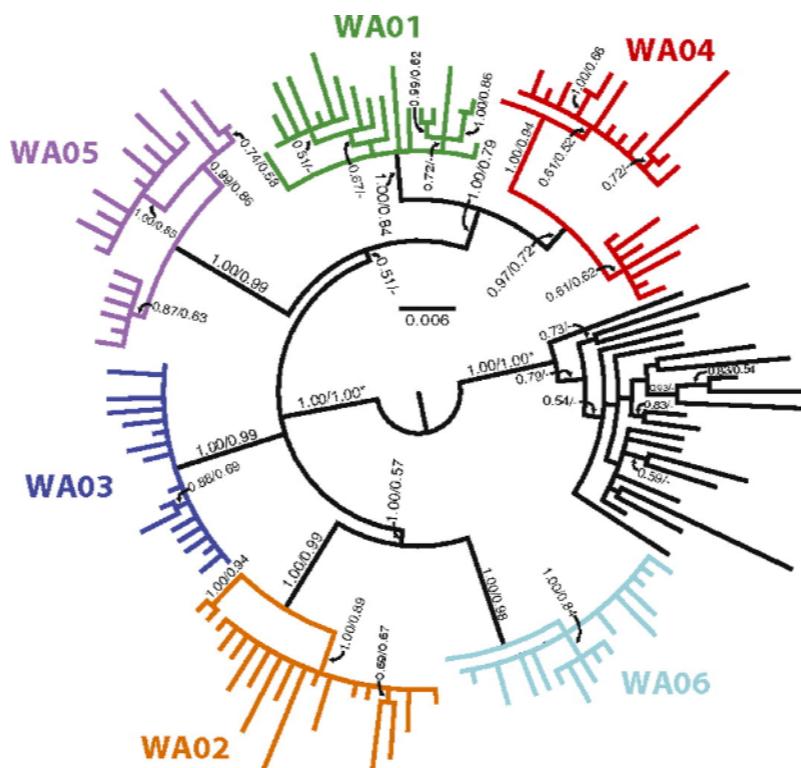
**Forensic**

# Biomedical Applications



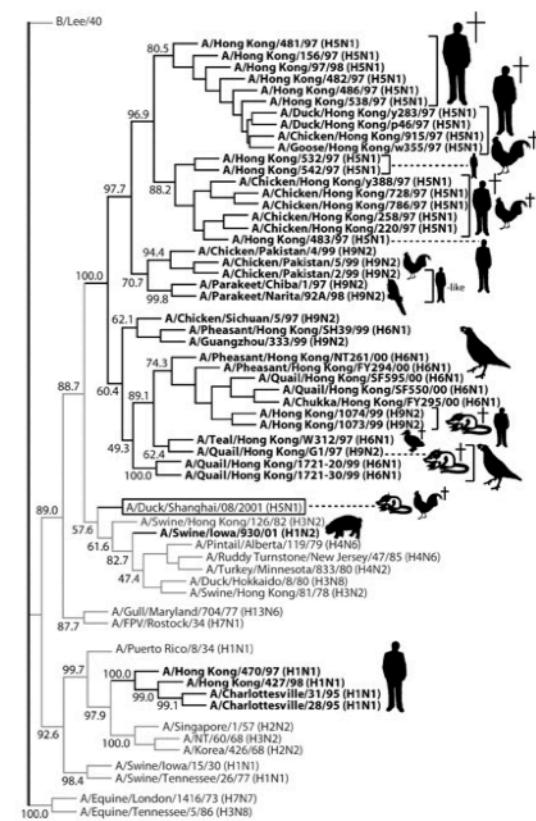
source: Langille et al., Nature, 2013

**Microbiome**



source: Scaduto et al., PNAS, 2010

**Forensic**

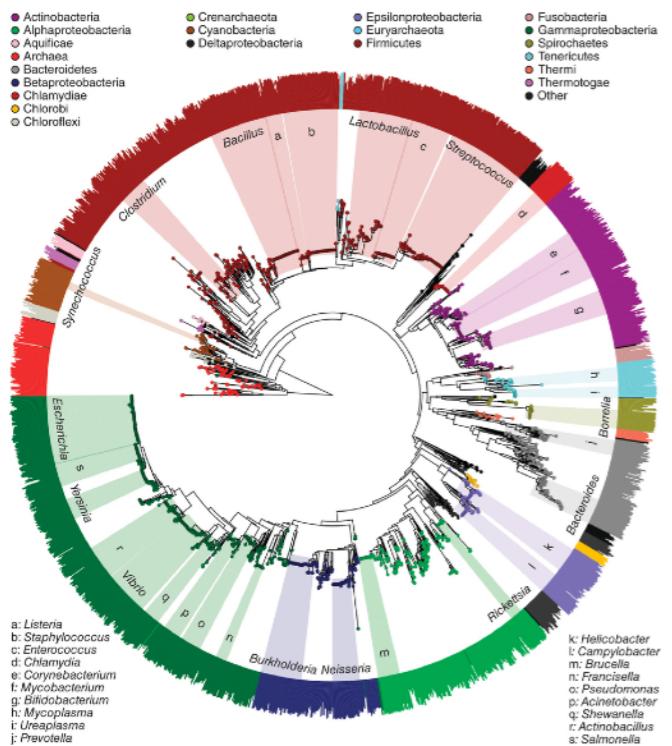


**Avian influenza**

# Biomedical Applications

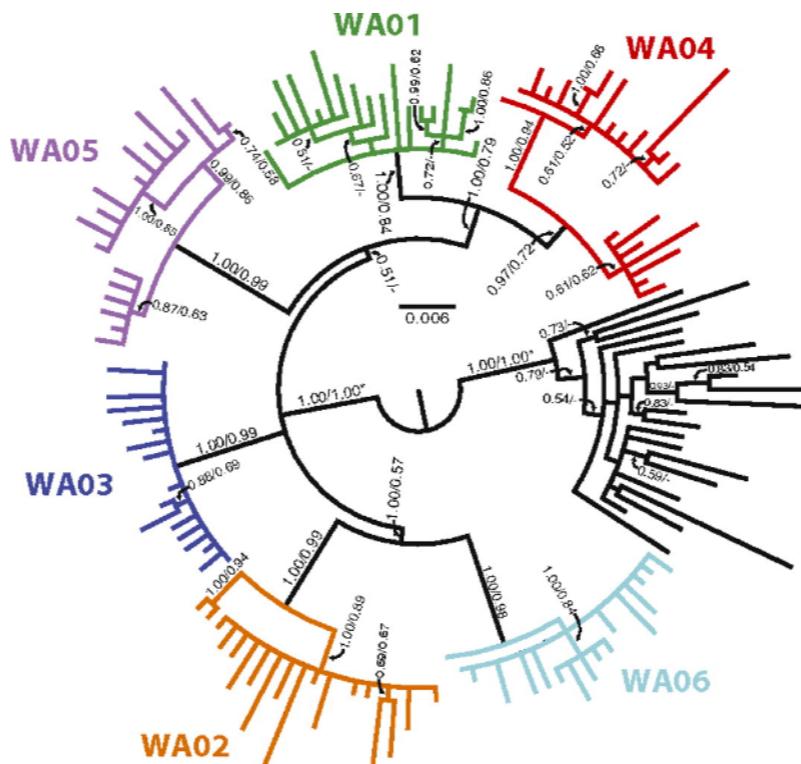
**“Nothing in evolution makes sense except in the light of phylogeny.”**

multiple coinage



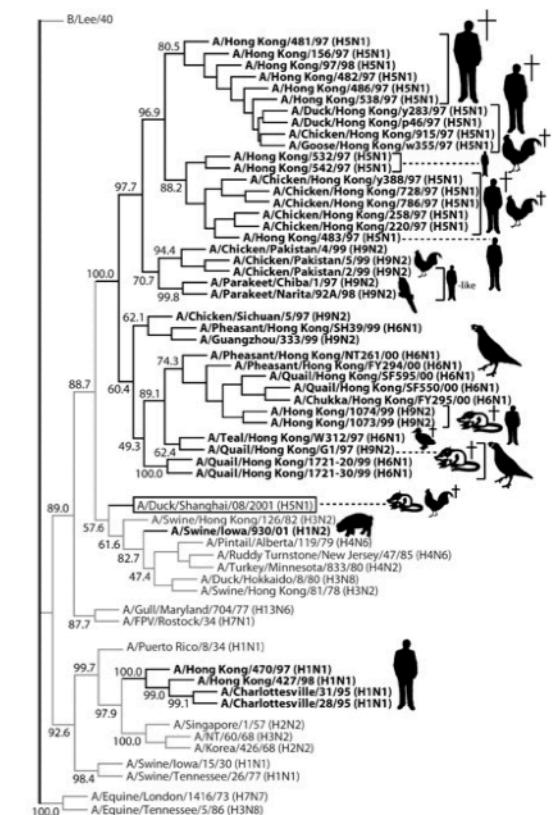
source: Langille et al., Nature, 2013

Microbiome



source: Scaduto et al., PNAS, 2010

Forensic



Avian influenza

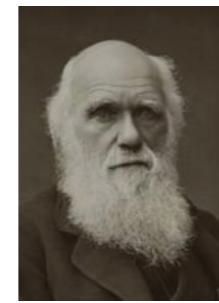
Why related to engineering?

Computational challenges?

# Phylogenetic reconstruction from data



Gorilla  
ACTGCACACCG



Human  
ACTGCCCG

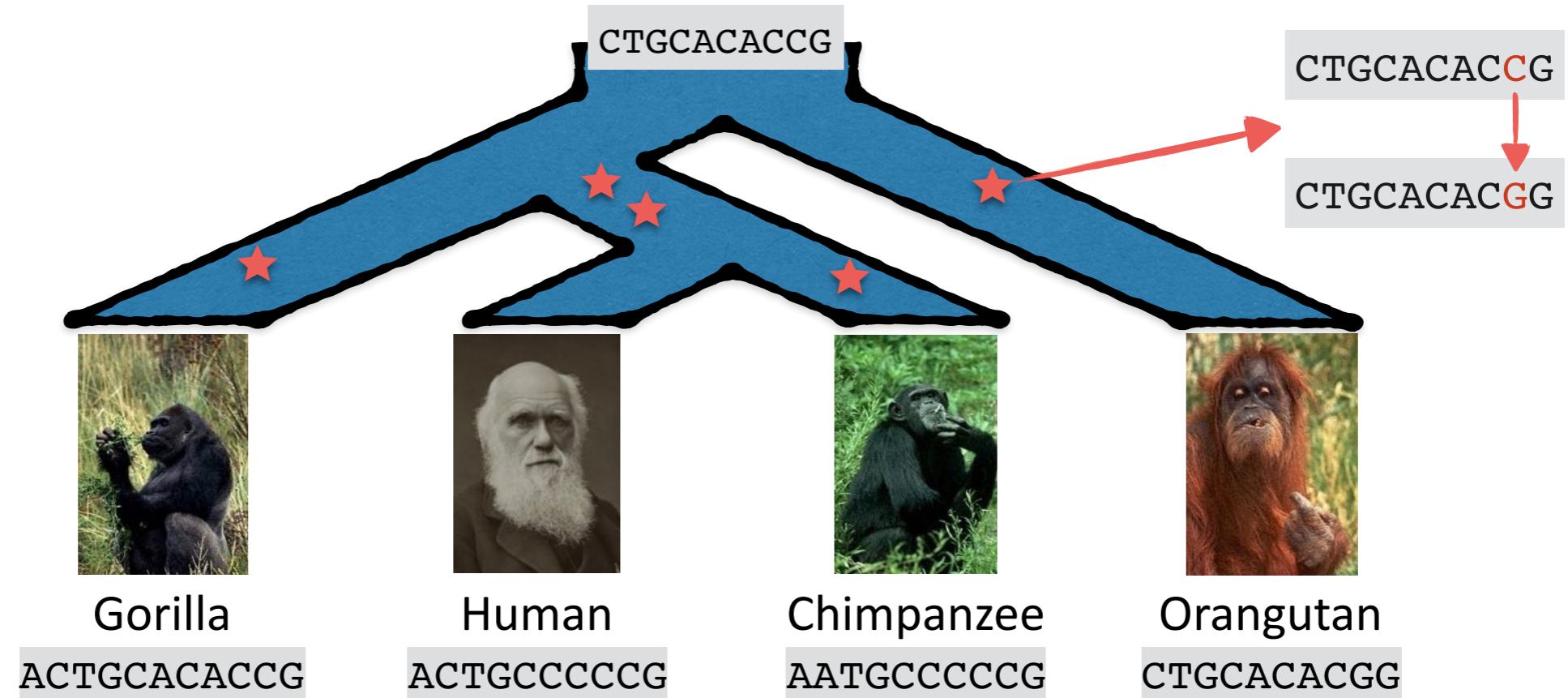


Chimpanzee  
AATGCCCG



Orangutan  
CTGCACACGG

# Phylogenetic reconstruction from data



# Phylogenetic reconstruction from data



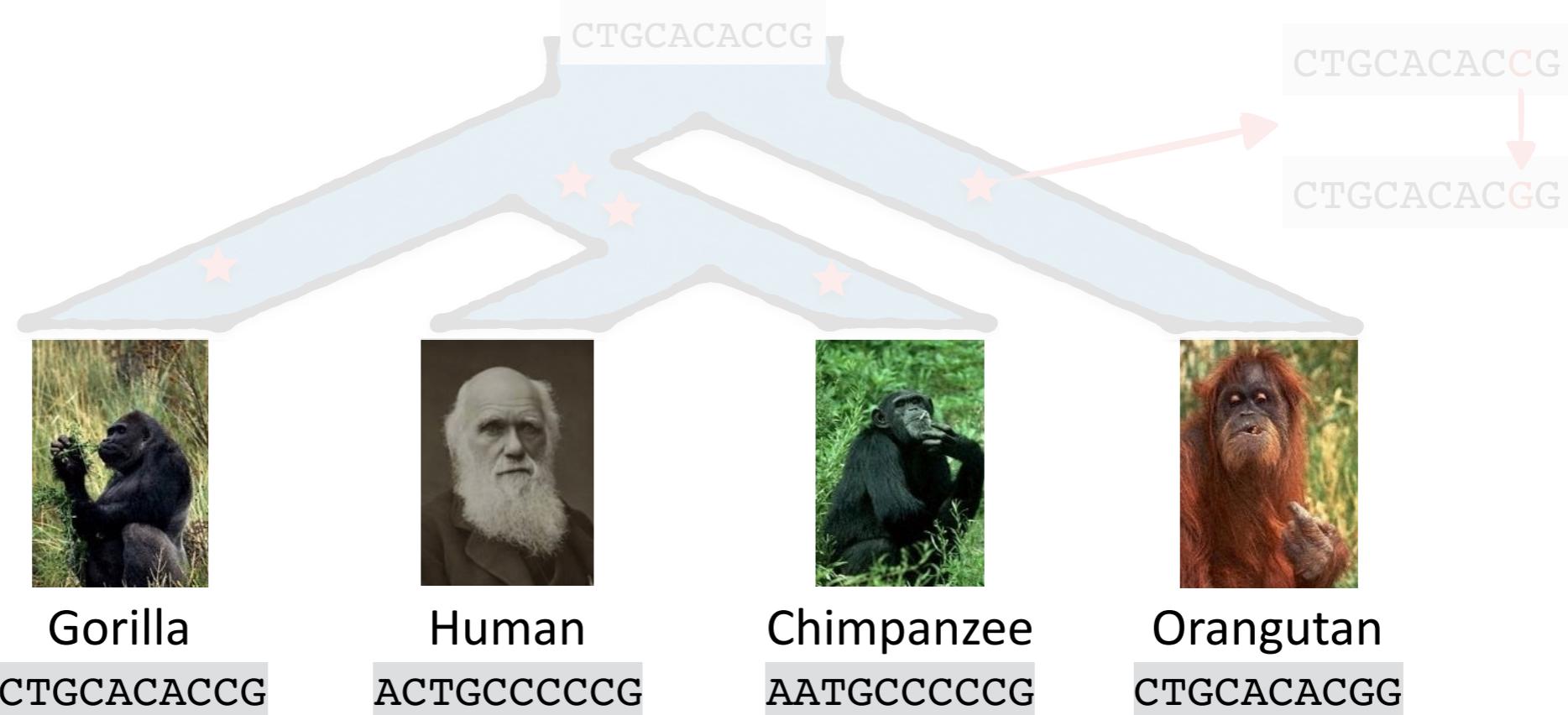
# Phylogenetic reconstruction from data



Gorilla	ACTGCACACCG
Human	ACTGC-CCCG
Chimpanzee	AATGC-CCCG
Orangutan	-CTGCACACGG

*D*

# Phylogenetic reconstruction from data

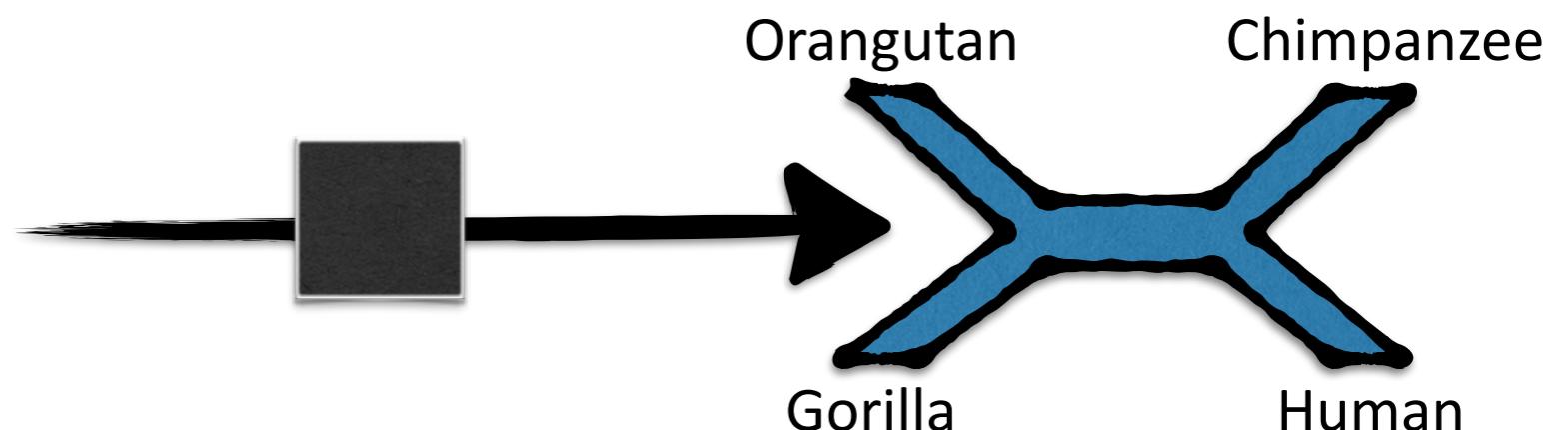


Gorilla	ACTGCACACCCG
Human	ACTGC-CCCG
Chimpanzee	AATGC-CCCG
Orangutan	-CTGCACACGG

$D$

$P(D|T)$

$T$



# Computational inference

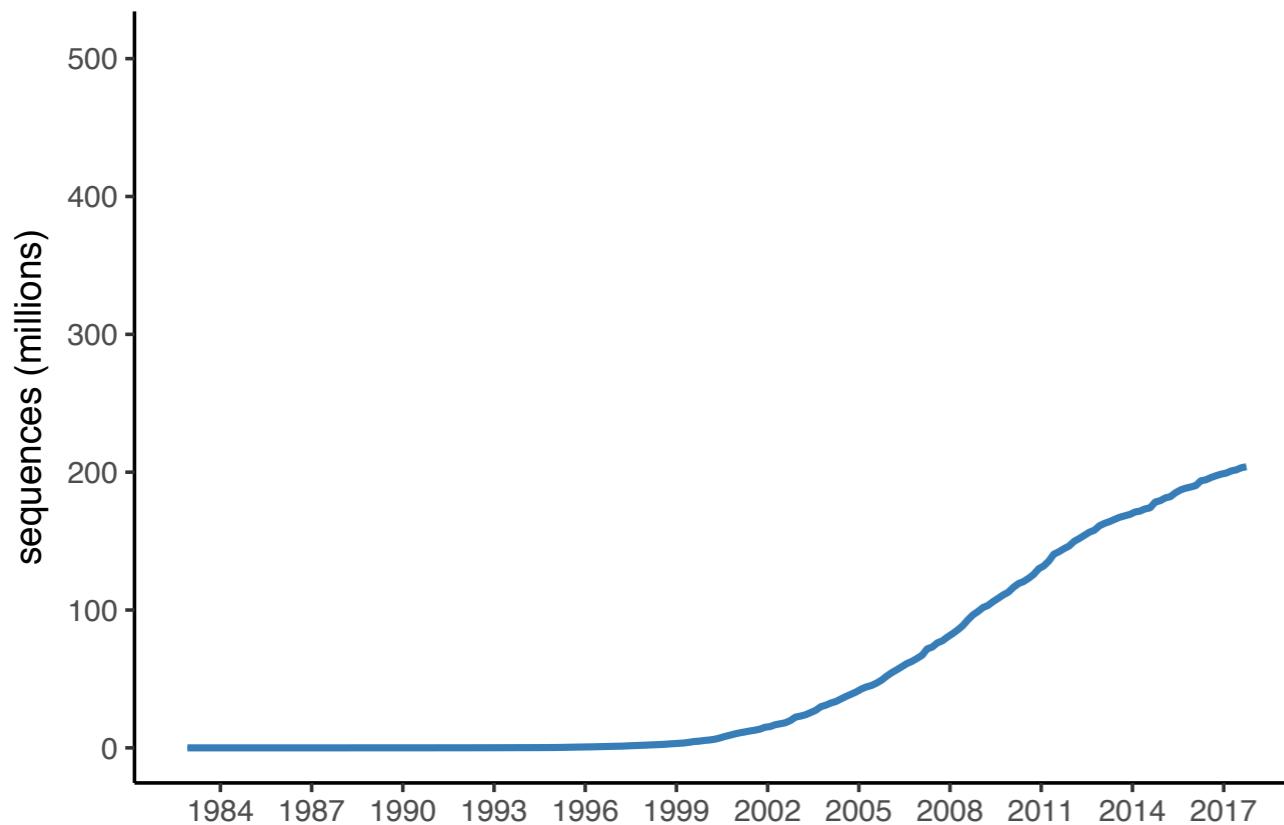
- Really no “experimental” way to validate results.
  - Phylogenetics is ENABLEd only due to computation
  - We use computation-heavy procedures to calculate statistical support

# Computational inference

- Really no “experimental” way to validate results.
  - Phylogenetics is ENABLEd only due to computation
  - We use computation-heavy procedures to calculate statistical support
- Genome evolution is complex. We need complex statistical models for accuracy
  - And need large data to infer under complex models

# Sequence data growth

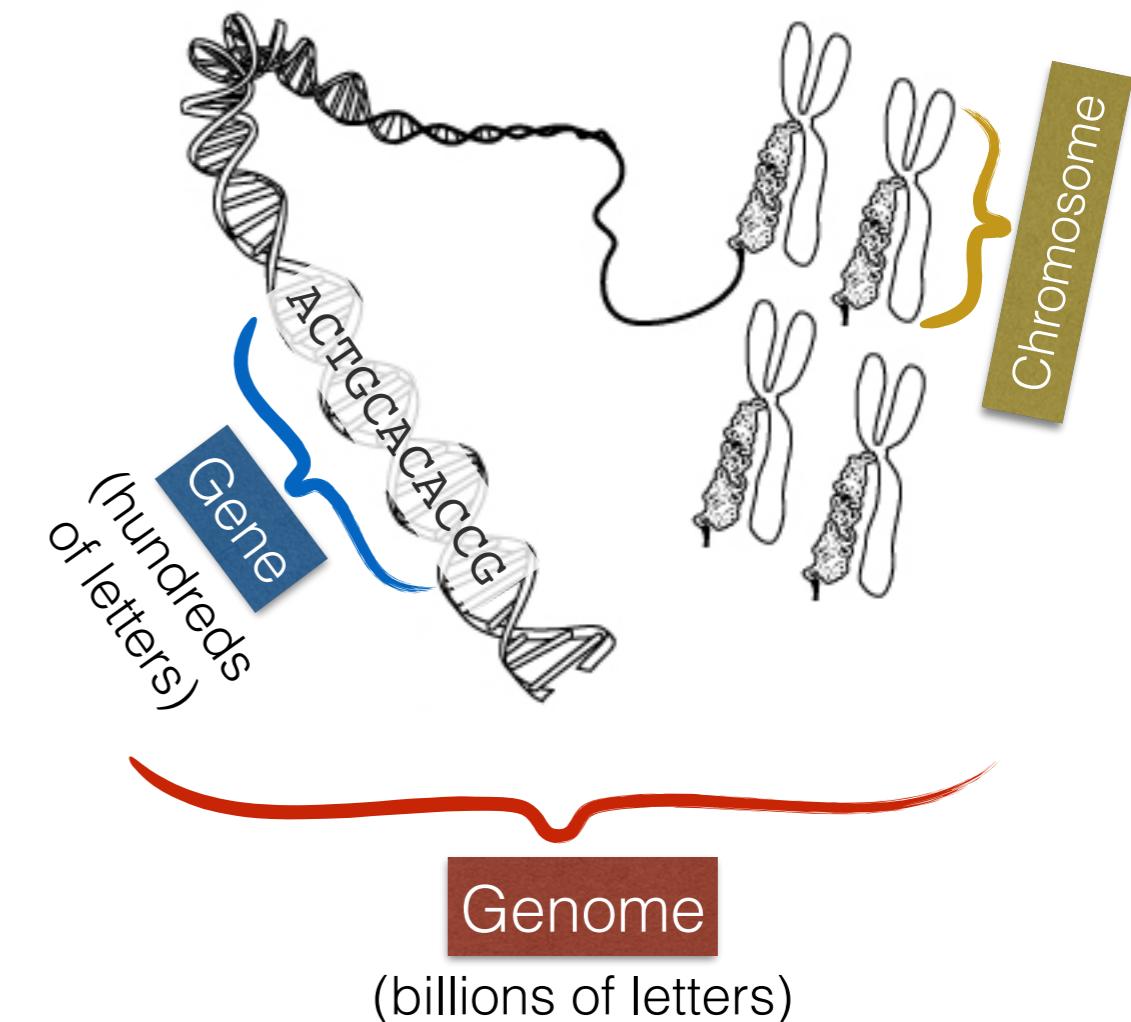
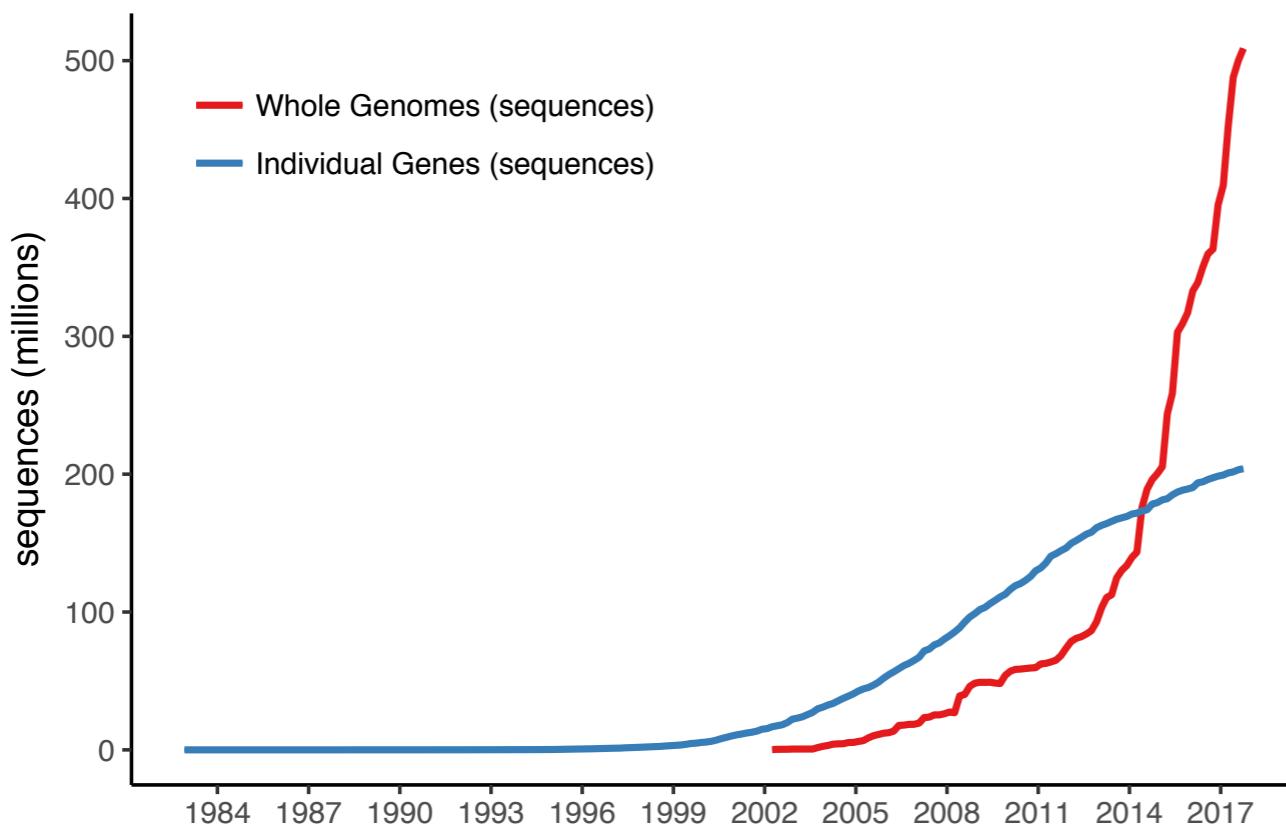
data source: NCBI (<http://www.ncbi.nlm.nih.gov/genbank/statistics>)



- Rapid growth in the available sequences

# Sequence data growth

data source: NCBI (<http://www.ncbi.nlm.nih.gov/genbank/statistics>)



- Rapid growth in the available sequences
- Our focus has shifted to “whole genomes”

**Billions of columns**

More genomic regions

SYKVKL	I	T	P	D	G	P	I	E	F	D	C	P	D	D	V	Y	I	L	D	Q	A	E	E	A	G	H	D	L	P	Y	
SYKVKL	I	T	P	D	G	P	I	E	F	D	C	P	D	N	V	Y	I	L	D	Q	A	E	E	A	G	H	D	L	P	Y	
SYKVKL	I	T	P	E	G	P	I	E	F	E	C	P	D	D	V	Y	I	L	D	Q	A	E	E	E	G	H	D	L	P	Y	
SYKVKL	I	T	P	D	G	P	I	E	F	E	C	P	D	D	V	Y	I	L	D	Q	A	E	E	E	G	H	D	L	P	Y	
SYKVKL	L	V	T	P	D	G	T	Q	E	F	E	C	P	S	D	V	Y	I	L	D	H	A	E	E	V	G	I	D	L	P	Y
TYKVKL	I	T	P	E	G	P	Q	E	F	D	C	P	D	D	V	Y	I	L	D	H	A	E	E	V	G	I	E	L	P	Y	
AYKVT	L	V	T	P	E	G	K	Q	E	L	E	C	P	D	D	V	Y	I	L	D	A	A	E	E	A	G	I	D	L	P	Y
AYKVT	L	V	T	P	T	G	N	V	E	F	Q	C	P	D	D	V	Y	I	L	D	A	A	E	E	E	G	I	D	L	P	Y
TYKVKF	I	T	P	E	G	E	Q	E	V	E	C	D	D	V	V	V	L	D	A	A	E	E	A	G	I	D	L	P	Y		
TYKVKF	I	T	P	E	G	E	L	E	V	E	C	D	D	V	V	V	L	D	A	A	E	E	A	G	I	D	L	P	Y		
TYKVKF	I	T	P	E	G	E	Q	E	V	E	C	D	D	V	V	V	L	D	A	A	E	E	A	G	I	D	L	P	Y		
TYKVKF	I	T	P	E	G	E	Q	E	V	E	C	E	E	D	V	V	V	L	D	A	A	E	E	A	G	L	D	L	P	Y	
TYKVKF	I	T	P	E	G	E	Q	E	V	E	C	E	E	D	V	V	V	L	D	A	A	E	E	A	G	L	D	L	P	Y	
TYNVKL	I	T	P	E	G	E	V	E	L	Q	V	P	D	D	V	Y	I	L	D	Q	A	E	E	D	G	I	D	L	P	Y	

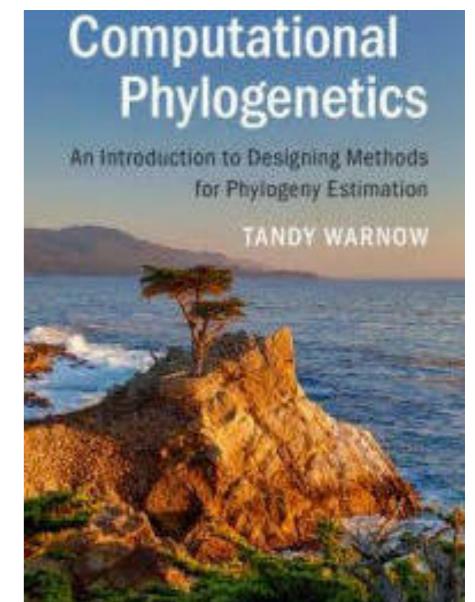
More sequences



**A million rows**

# Phylogenetic reconstruction : a computational nightmare

- Almost all problems are NP-Hard!
- The search space is huge.
  - Focusing on topology alone,  
there are  $(2n-5)!!$  trees with n leaves
  - We also care about branch lengths:  $\mathbb{R}^{2n-3}$
  - We are interested in  $n$  in  $10^1 \dots 10^6$  range



Tandy Warnow

# To scale to large datasets . . .

- Approximate and heuristic solutions

# To scale to large datasets . . .

- Approximate and heuristic solutions
- Make the problem easier
  - Divide-and-conquer
  - Constrained search

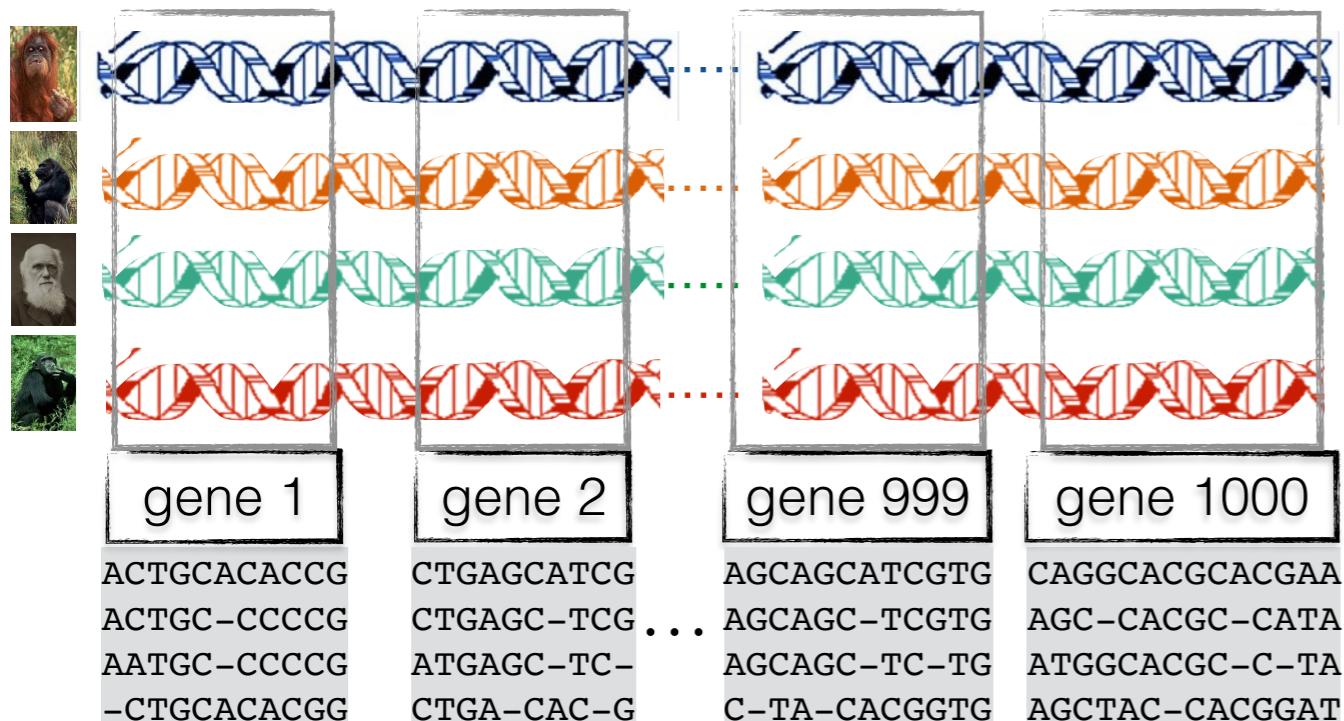
# To scale to large datasets . . .

- Approximate and heuristic solutions
- Make the problem easier
  - Divide-and-conquer
  - Constrained search
- Develop optimized code.
  - Take advantage of High Performance Computing

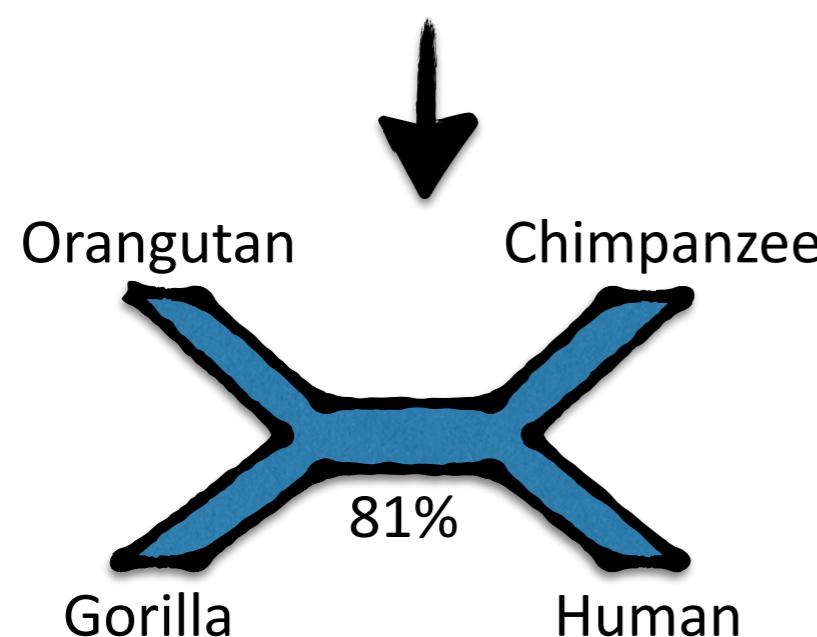
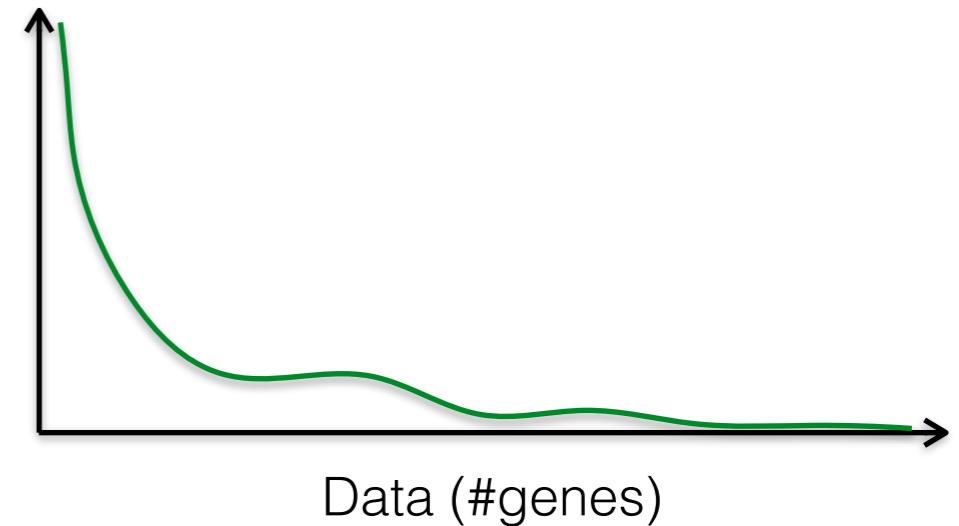
# To scale to large datasets . . .

- Approximate and heuristic solutions
- Make the problem easier
- But how about accuracy?
- Do large data help?
- Develop optimized code.
- Take advantage of High Performance Computing

# Phylogenomics: more data → better inference



Species tree error



# Phylogenomics: more data → better inference



Species tree error

**nature**

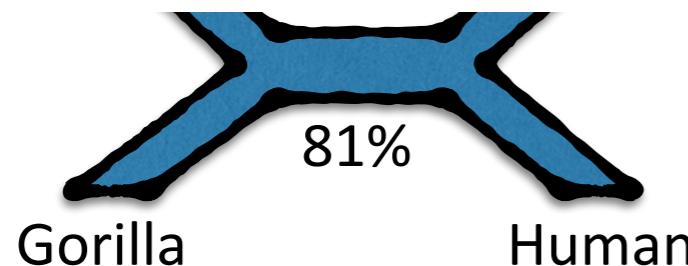
International journal of science

Evolution

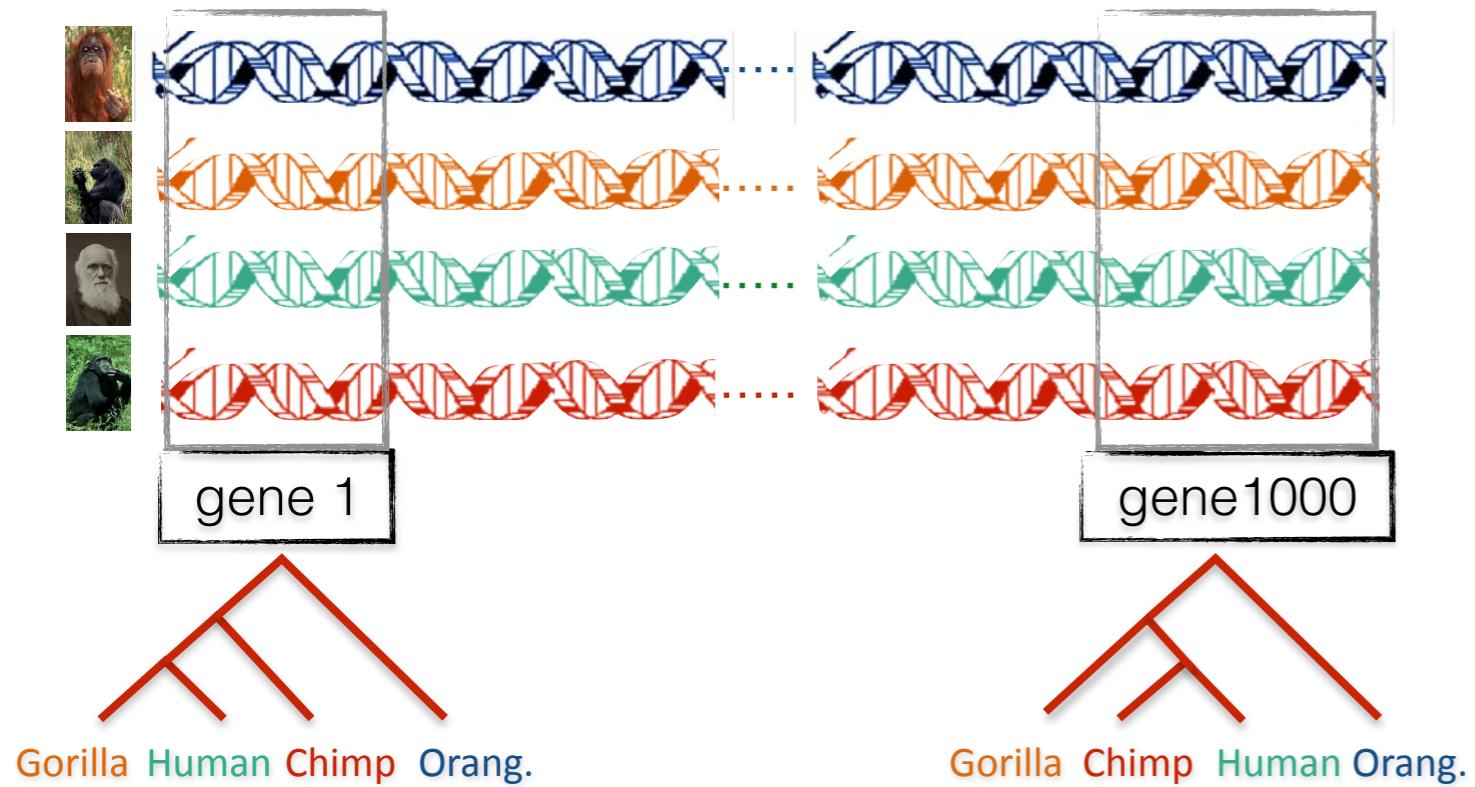
## Ending incongruence

Henry Gee

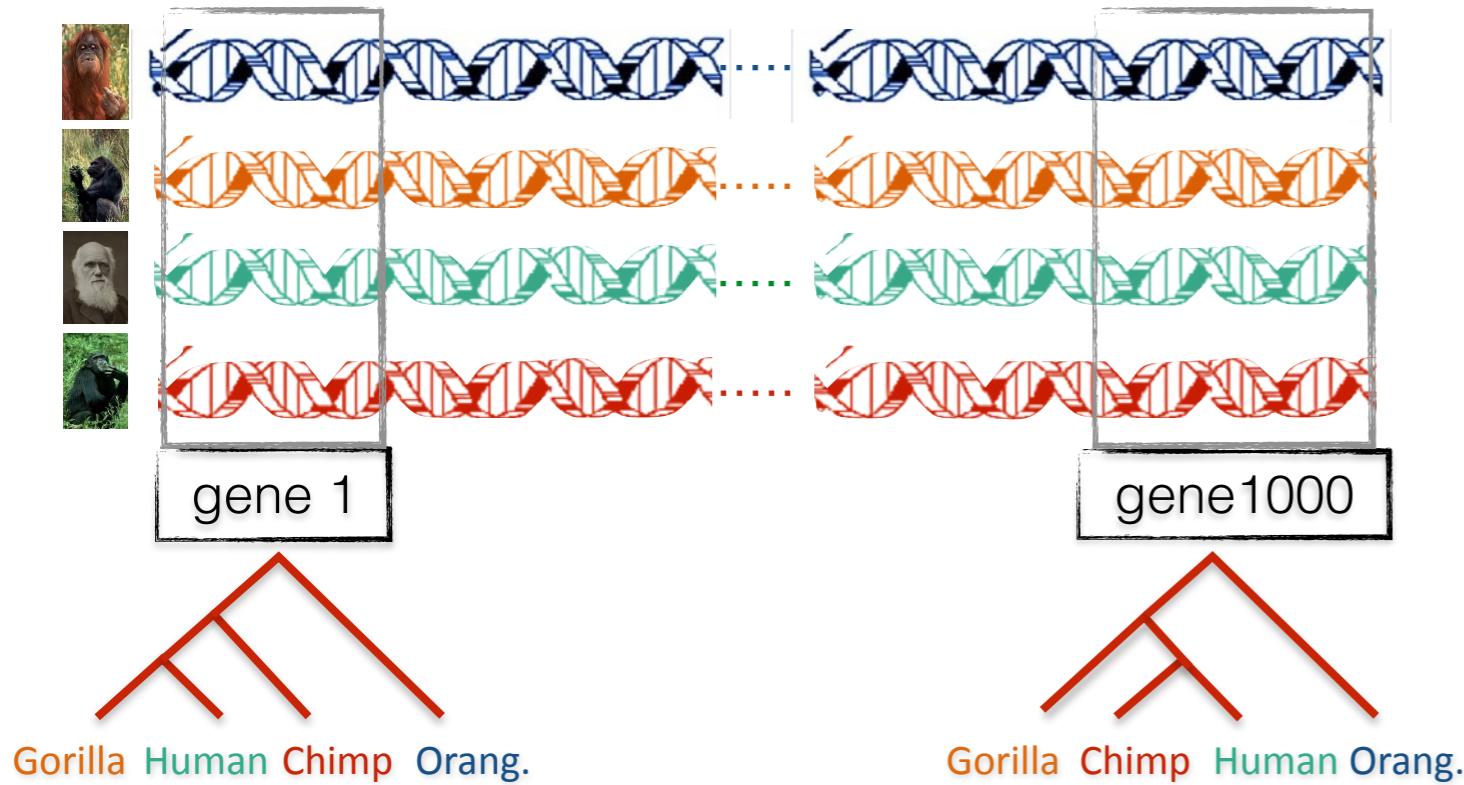
Recovering the true evolutionary history of any group of organisms has seemed impossible. The availability of large amounts of genomic data promises an era in which the uncertainties are better constrained.



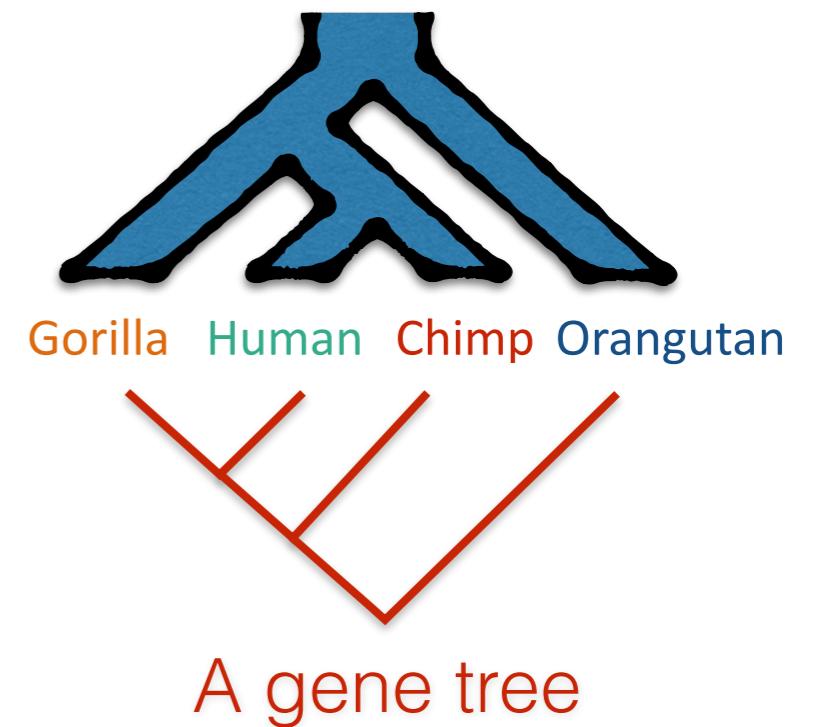
# Gene tree discordance



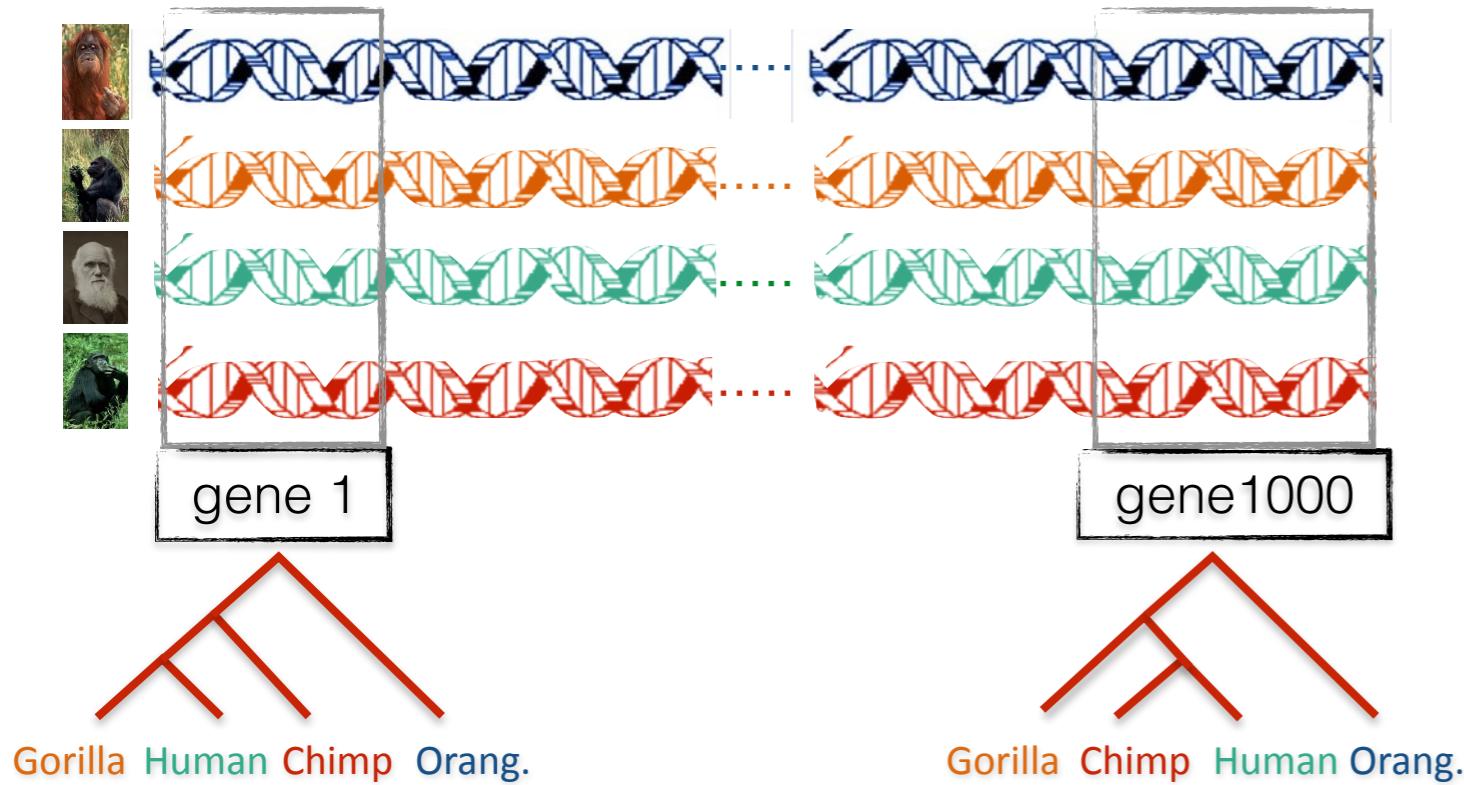
# Gene tree discordance



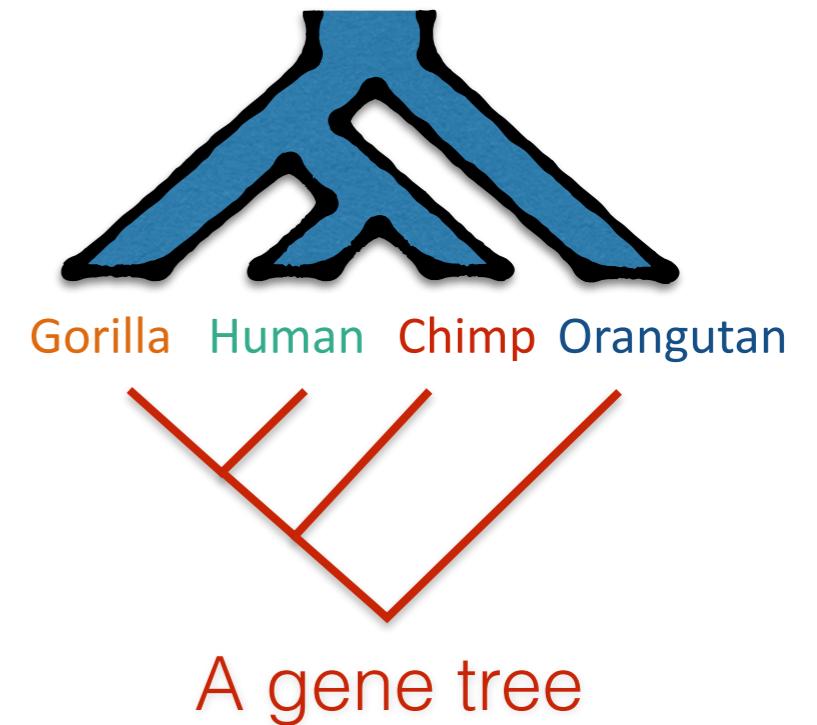
**The species tree**



# Gene tree discordance



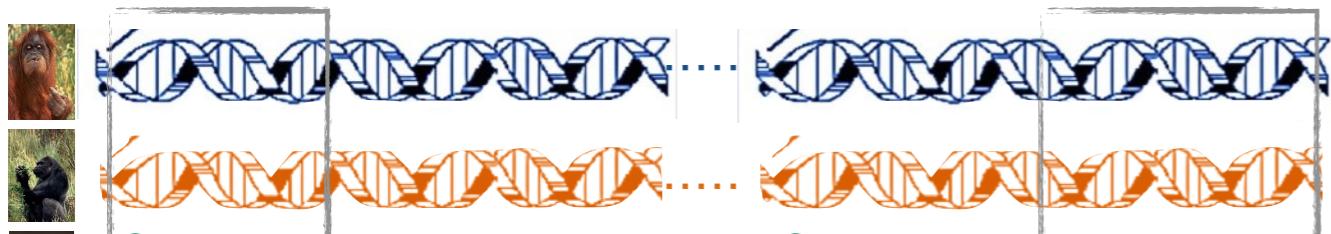
The species tree



## Causes of gene tree discordance include:

- Incomplete Lineage Sorting (ILS)
- Duplication and loss
- Horizontal Gene Transfer (HGT)

# Gene tree discordance



The species tree



## Trends in Genetics



Volume 22, Issue 4, April 2006, Pages 225–231

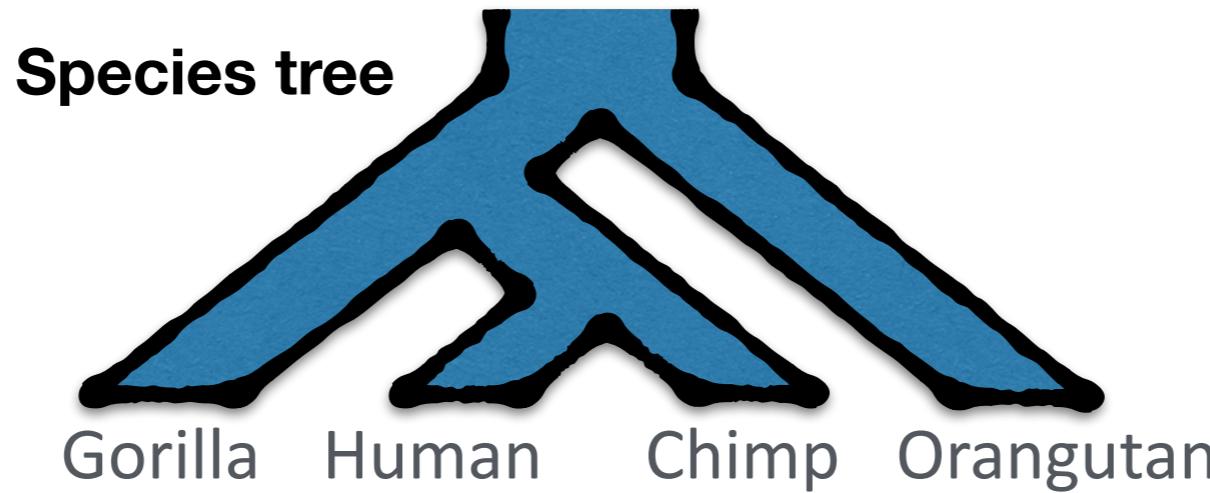
### Phylogenomics: the beginning of incongruence?

Olivier Jeffroy, Henner Brinkmann, Frédéric Delsuc, Hervé Philippe

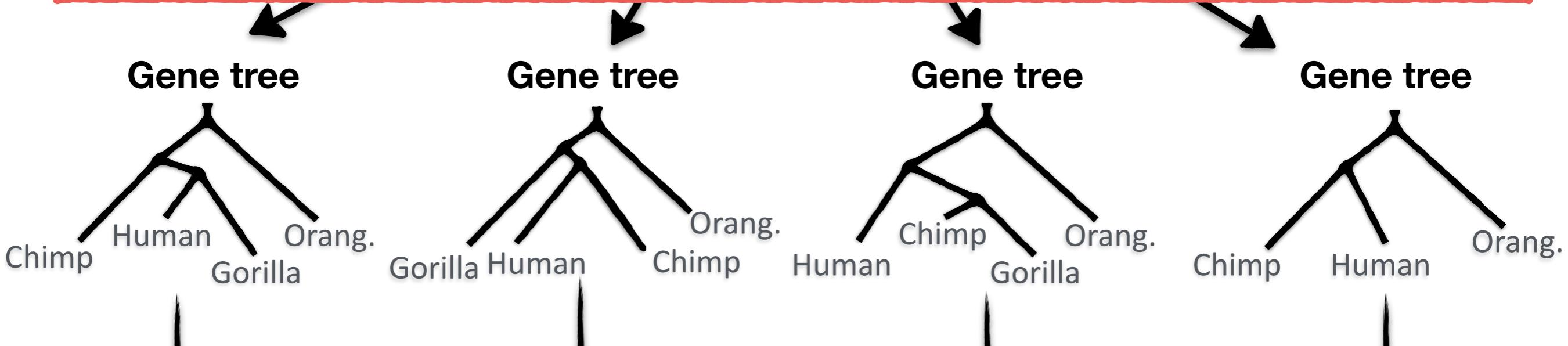
Show more

<https://doi.org/10.1016/j.tig.2006.02.003>

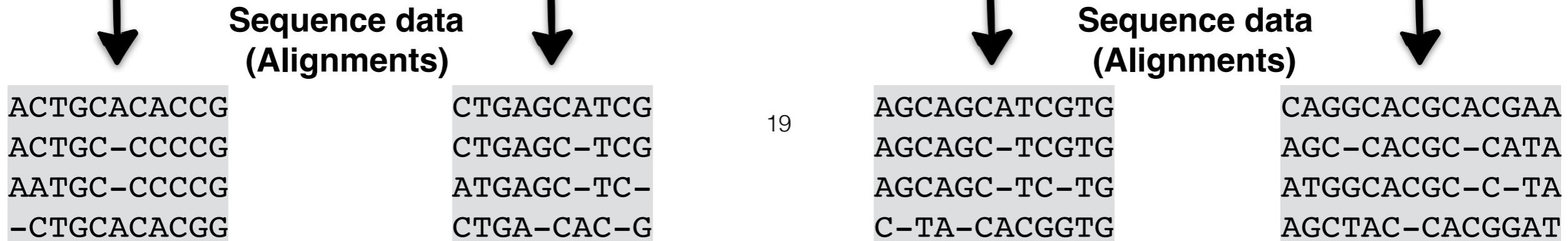
[Get rights and content](#)

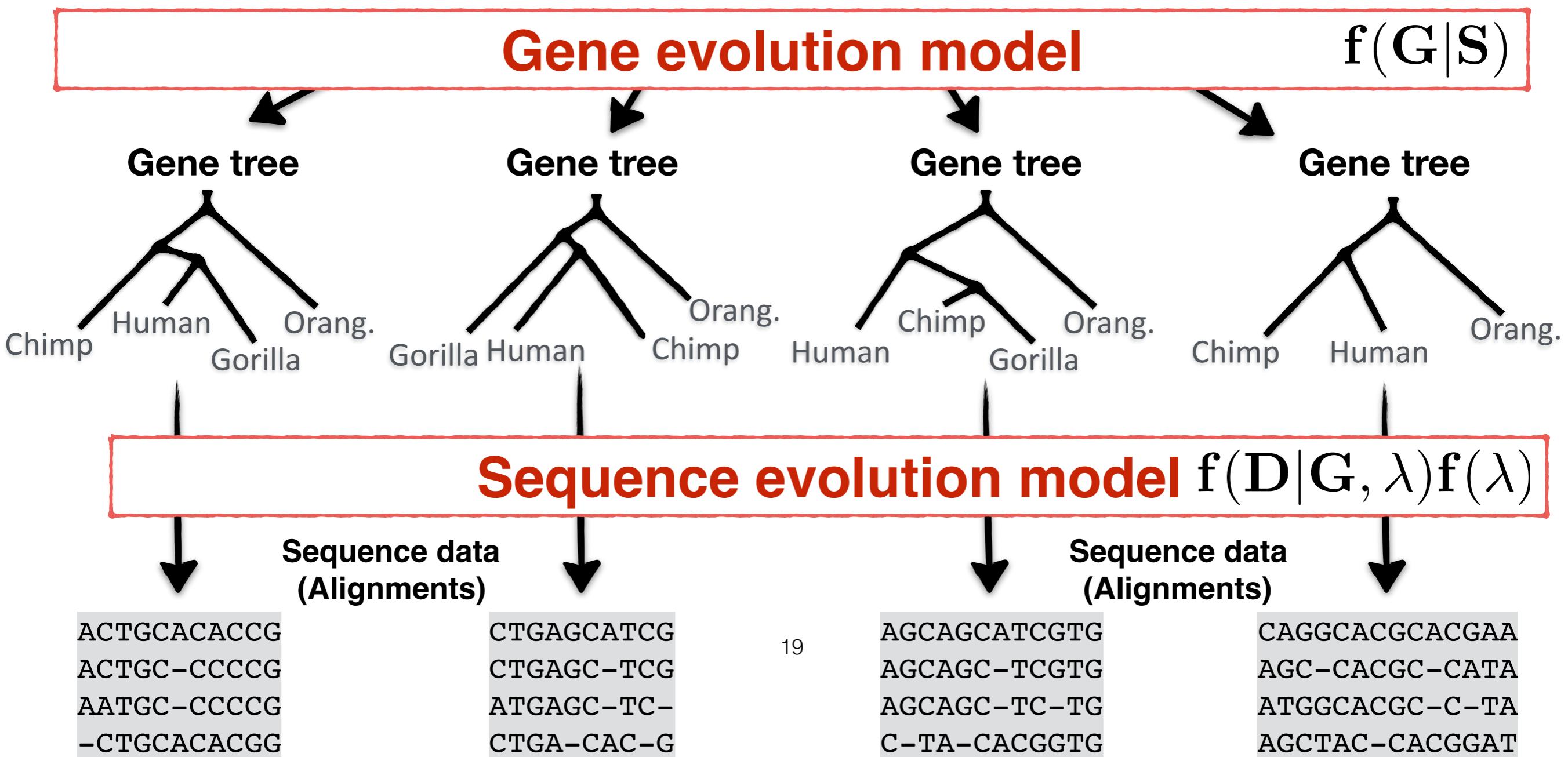
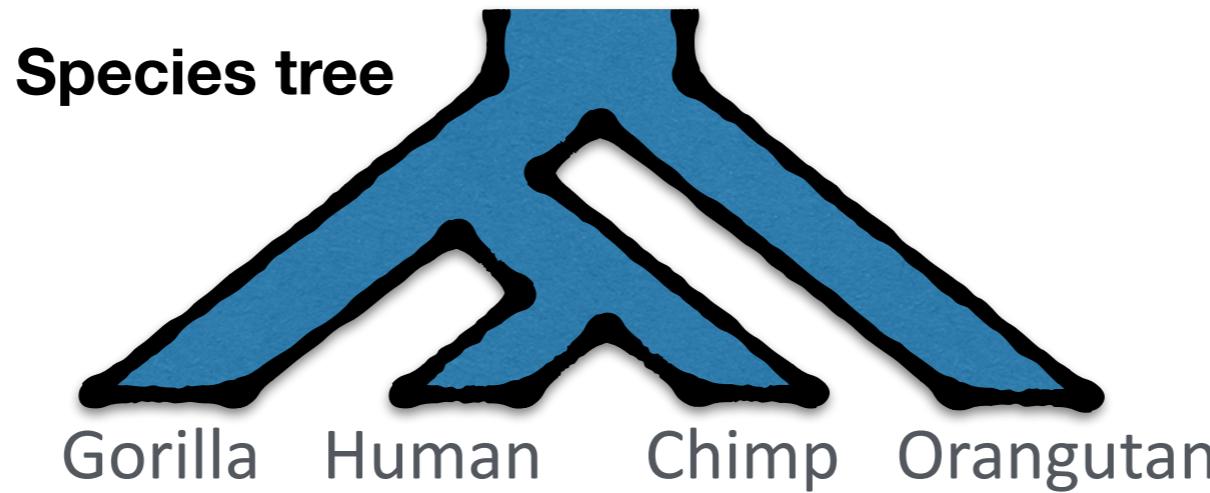


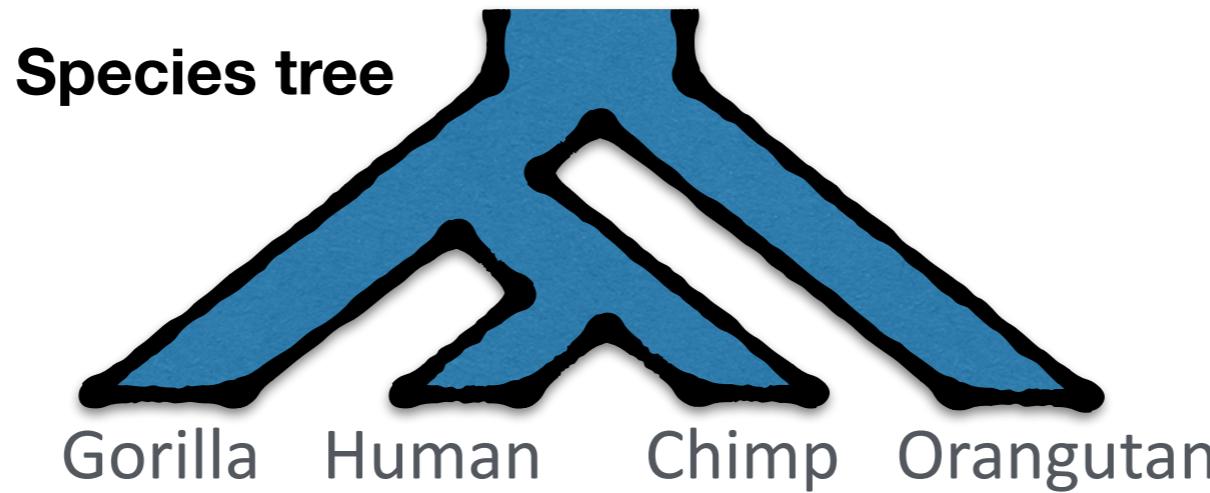
## Gene evolution model



## Sequence evolution model



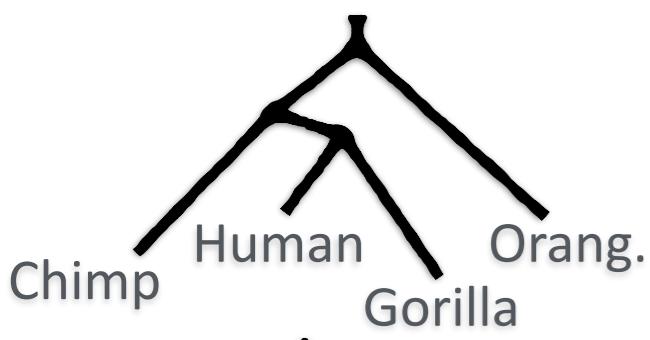




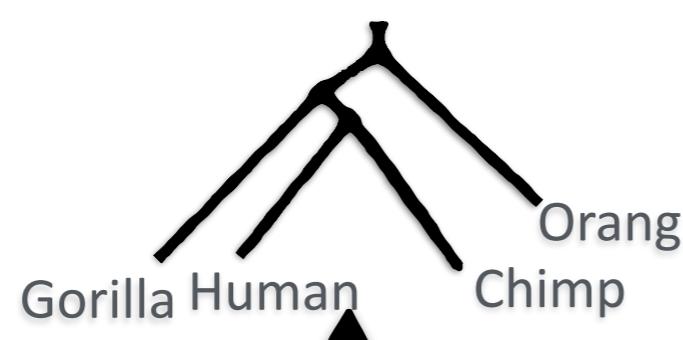
## Gene evolution model

$$f(G|S)$$

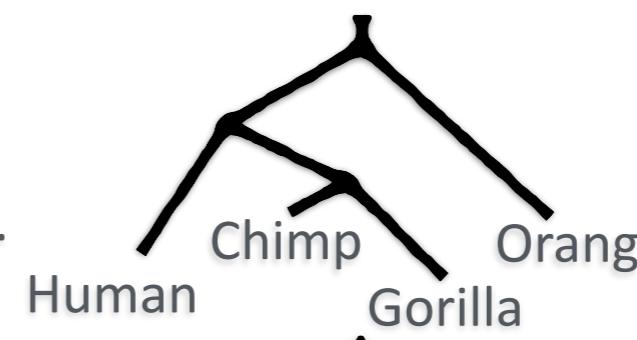
Gene tree



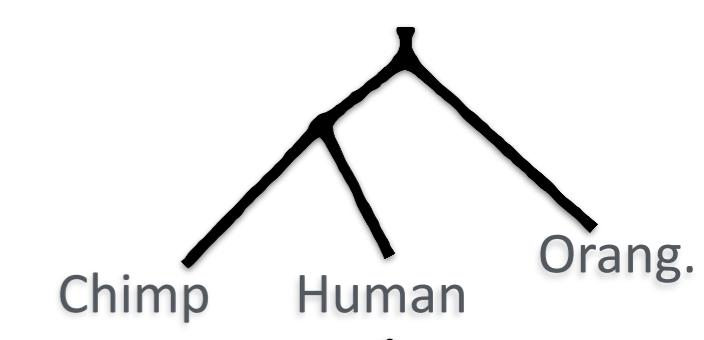
Gene tree



Gene tree



Gene tree



## Sequence evolution model

$$f(D|G, \lambda) f(\lambda)$$

Sequence data  
(Alignments)

```
ACTGCACACCCG
ACTGC-CCCCG
AATGC-CCCCG
-CTGCACACGG
```

```
CTGAGCATCG
CTGAGC-TCG
ATGAGC-TC-
CTGA-CAC-G
```

20

```
AGCAGGCATCGTG
AGCAGC-TCGTG
AGCAGC-TC-TG
C-TA-CACGGTG
```

```
CAGGCACGCACGAA
AGC-CACGC-CATA
ATGGCACGC-C-TA
AGCTAC-CACGGAT
```

# 1KP: Plant whole transcriptomes

[Wickett, Mirarab, *et al.*, PNAS, 2014]



gane wong

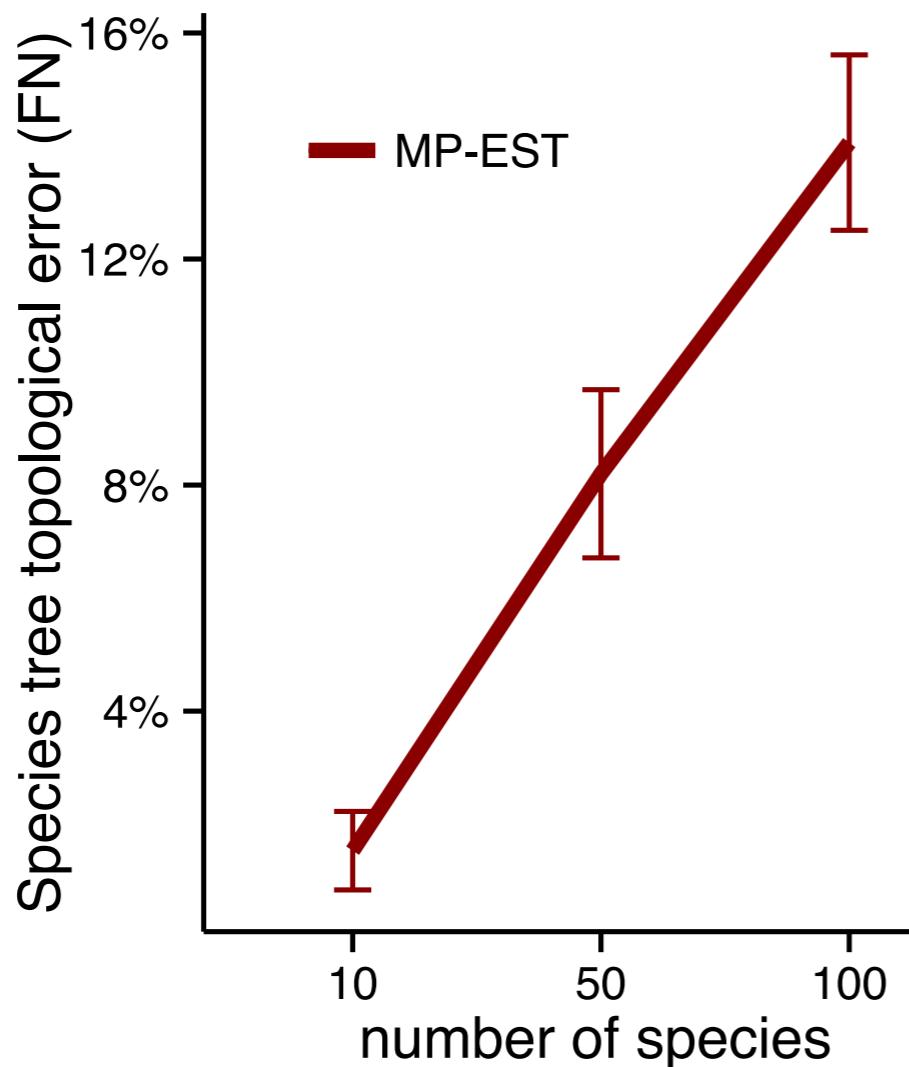
华大基因  
**BGI**

- Two phases
  - Pilot: Whole transcriptomes for **103 plant species** and ~800 single copy genes
  - Capstone: Whole transcriptomes for **1178 plant species** and ~410 single copy genes
- Spans all plants and algae: ~1 billion years of evolution
- Many unanswered questions about plant evolution



Jim  
leebens-mack

# Number of species impacts error

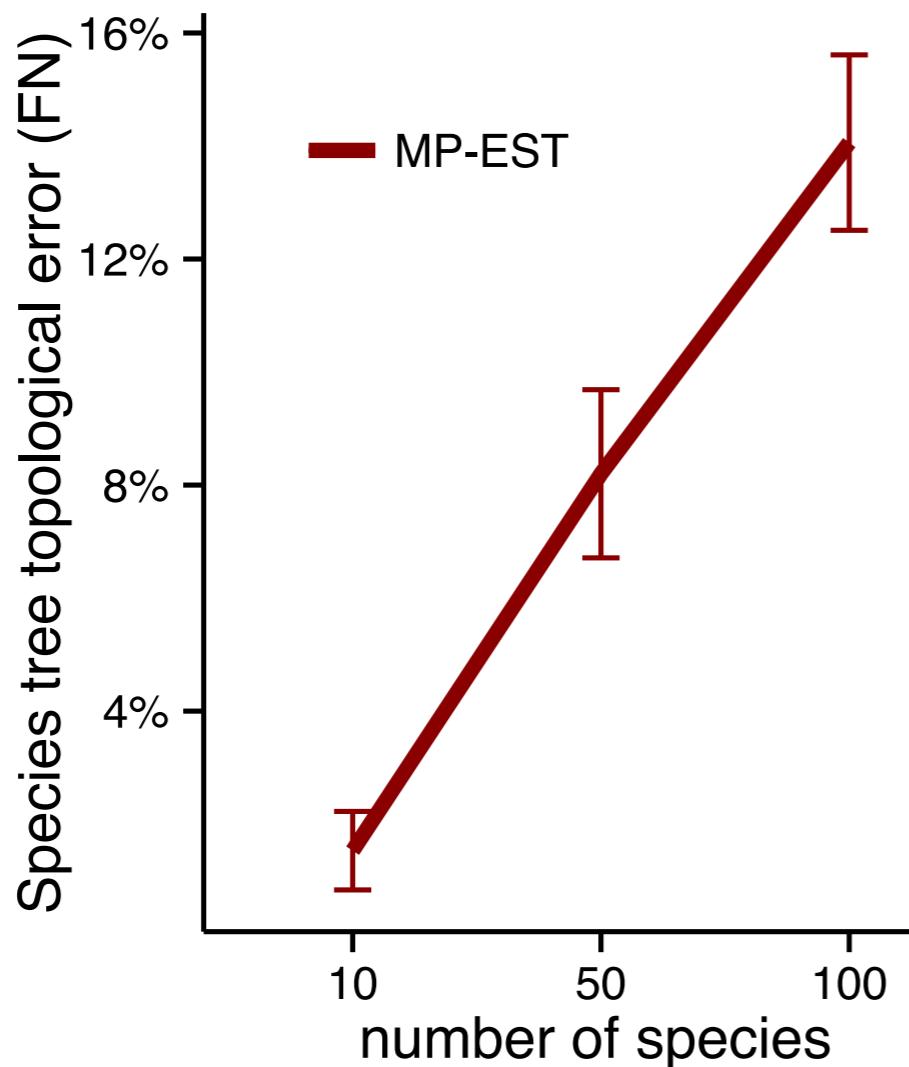


1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# Our response: a new method called ASTRAL

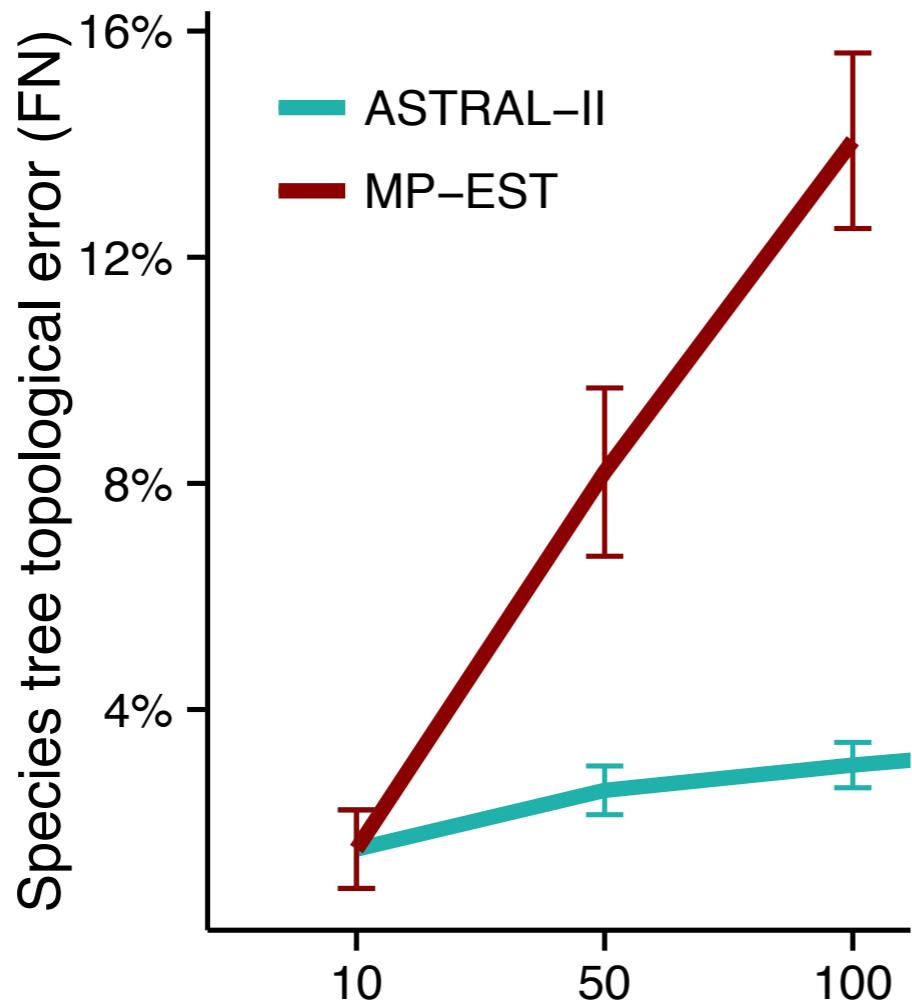
- Solves a simpler problem:
  - Instead of computing likelihood (difficult), it introduces a **distance-based** consistent estimator
- Uses divide-and-conquer
  - Each tree on  $n$  species can be broken into  $\binom{n}{4}$  “quartet trees”

# Number of species impacts error



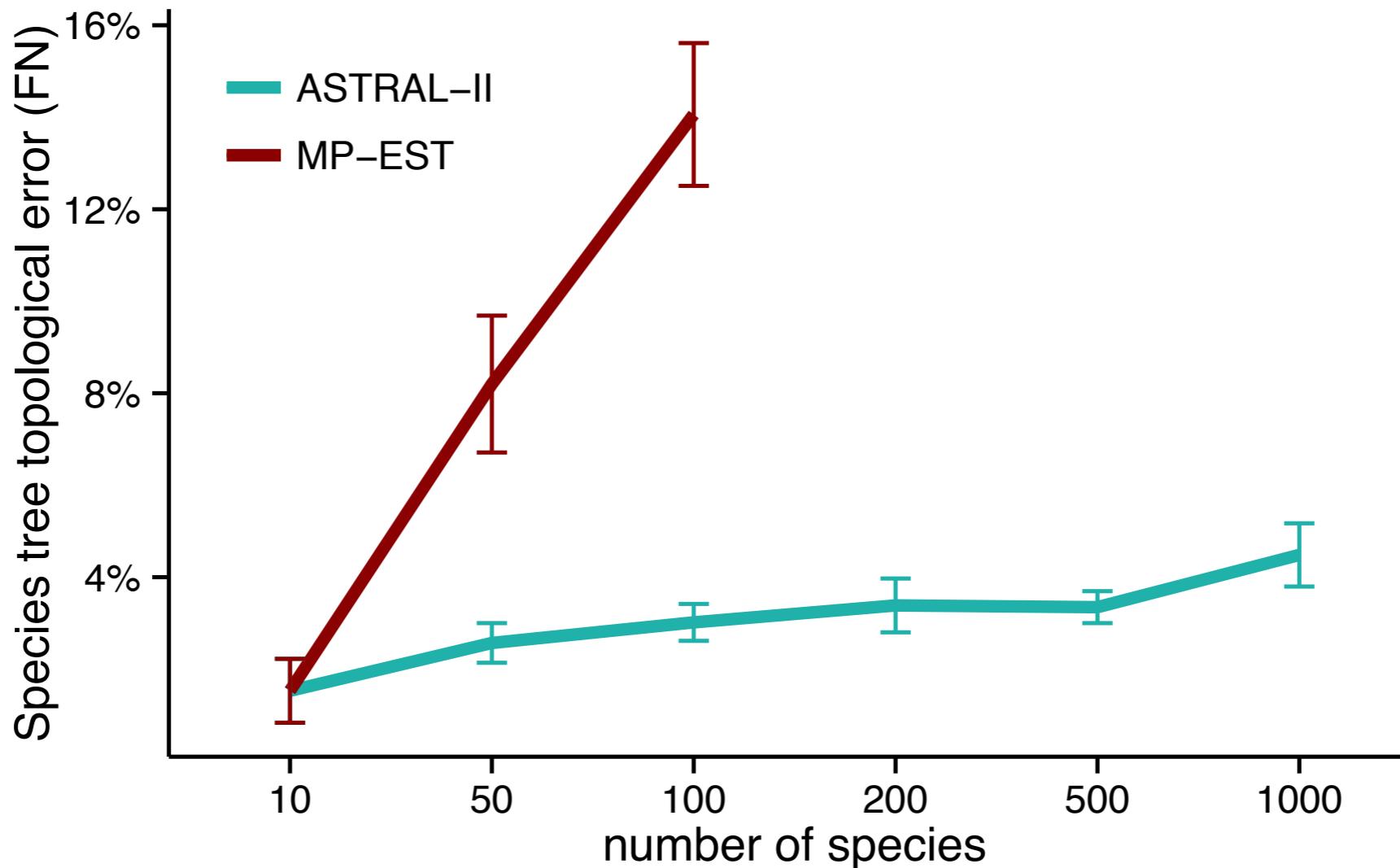
1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# ASTRAL: accurate and scalable



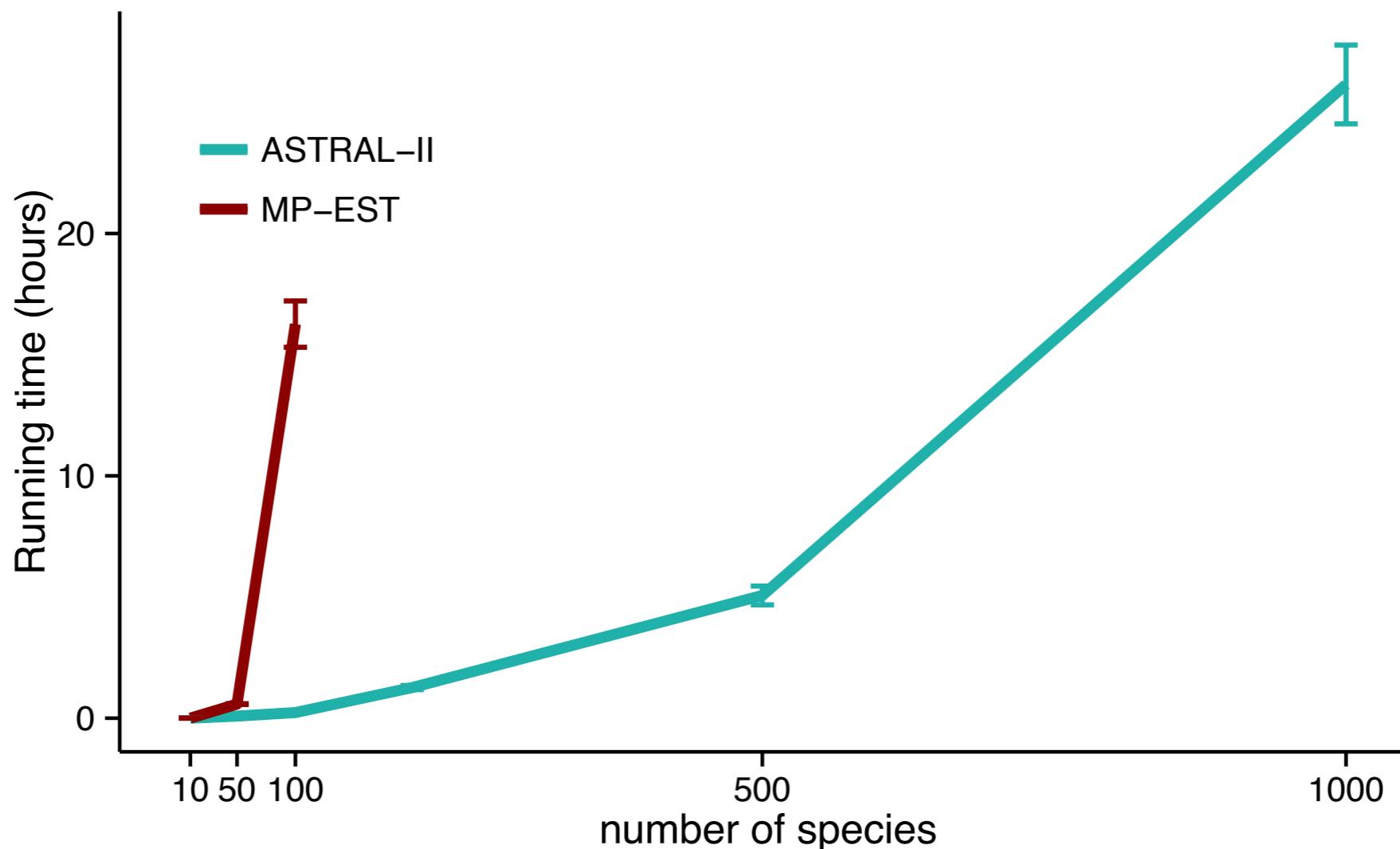
1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# ASTRAL: accurate and scalable



1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# Running time as function of # species



1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

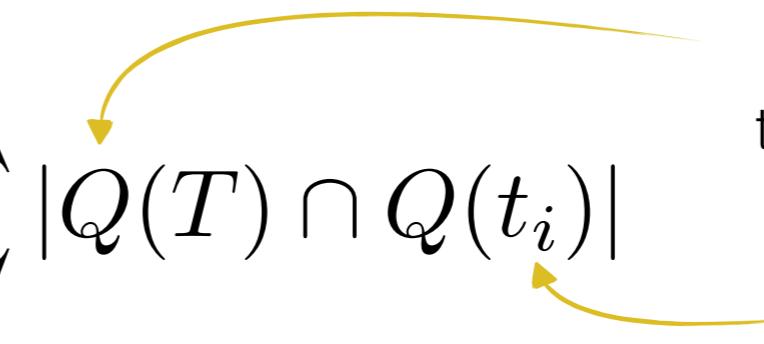
# Maximum Quartet Support Species Tree

- Median tree problem:

Find the species tree with the maximum number of induced quartet trees shared with the collection of input gene trees

$$Score(T) = \sum_1^m |Q(T) \cap Q(t_i)|$$

the set of quartet trees induced by T  
a gene tree



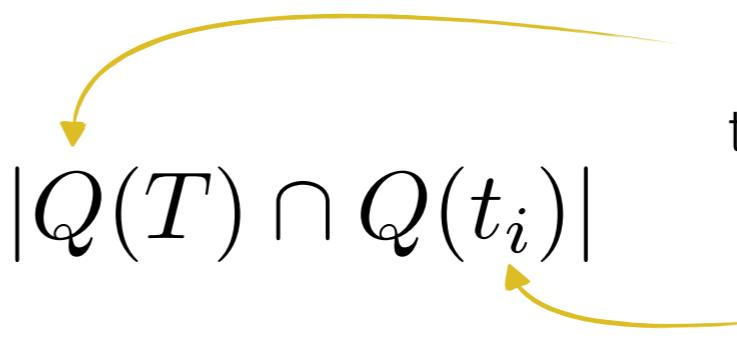
# Maximum Quartet Support Species Tree

- Median tree problem:

Find the species tree with the maximum number of induced quartet trees shared with the collection of input gene trees

$$Score(T) = \sum_1^m |Q(T) \cap Q(t_i)|$$

the set of quartet trees induced by T  
a gene tree



- Theorem:** Statistically consistent under the multi-species coalescent model when solved exactly

# ASTRAL

[Mirarab, et al., 2014] [Mirarab and Warnow, 2015] [Zhang et. al., 2018]

- Solve the Maximum Quartet Support problem exactly using [dynamic programming](#)

# ASTRAL

[Mirarab, et al., 2014] [Mirarab and Warnow, 2015] [Zhang et. al., 2018]

- Solve the Maximum Quartet Support problem exactly using [dynamic programming](#)
- [Constrains](#) the search space using gene trees
  - The constrained version remains [statistically consistent](#)
  - Running time of the constrained version increases polynomially with the input size:  $O(\mathcal{D}(nk)^{1.73})$   
 $\mathcal{D}=O(nk)$  is the sum of degrees of all unique gene tree nodes

# ASTRAL

[Mirarab, et al., 2014] [Mirarab and Warnow, 2015] [Zhang et. al., 2018]

- Solve the Maximum Quartet Support problem exactly using [dynamic programming](#)
- [Constrains](#) the search space using gene trees
  - The constrained version remains [statistically consistent](#)
  - Running time of the constrained version increases polynomially with the input size:  $O(\mathcal{D}(nk)^{1.73})$   
 $\mathcal{D}=O(nk)$  is the sum of degrees of all unique gene tree nodes
- [Sample complexity](#): Can find the correct species tree with arbitrarily high probability given  $O(\log(n)f^{-2})$  gene trees where  $f$  is the length of the shortest branch in the tree [Shekhar, Roch, Mirarab, TCBB, 2017]

# ASTRAL on the “pilot” 1KP

[Wickett, Mirarab, *et al.*, PNAS, 2014]

- The ASTRAL tree had high support and was one of the accepted trees in the resulting publication
- ASTRAL took only 10 minutes on 103 taxa and 400 genes
- Solved many but not all questions of interest



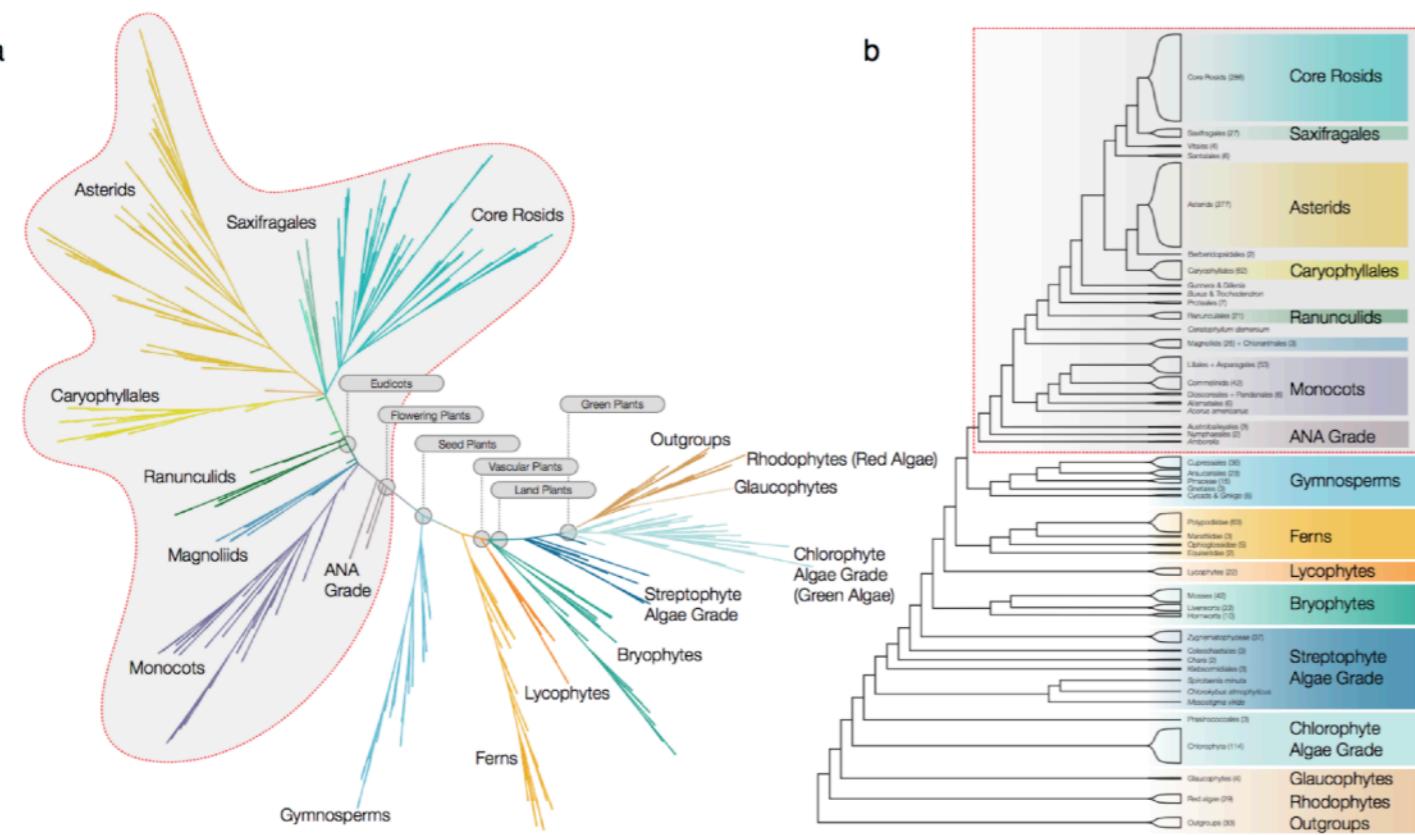
## Phylotranscriptomic analysis of the origin and early diversification of land plants

Norman J. Wickett<sup>a,b,1,2</sup>, Siavash Mirarab<sup>c,1</sup>, Nam Nguyen<sup>c</sup>, Tandy Warnow<sup>c</sup>, Eric Carpenter<sup>d</sup>, Naim Matasci<sup>e,f</sup>, Saravanaraj Ayyampalayam<sup>g</sup>, Michael S. Barker<sup>f</sup>, J. Gordon Burleigh<sup>h</sup>, Matthew A. Gitzendanner<sup>h,i</sup>, Brad R. Ruhfel<sup>h,j,k</sup>, Eric Wafula<sup>l</sup>, Joshua P. Der<sup>l</sup>, Sean W. Graham<sup>m</sup>, Sarah Mathews<sup>n</sup>, Michael Melkonian<sup>o</sup>, Douglas E. Soltis<sup>h,i,k</sup>, Pamela S. Soltis<sup>h,i,k</sup>, Nicholas W. Miles<sup>k</sup>, Carl J. Rothfels<sup>p,q</sup>, Lisa Pokorny<sup>p,r</sup>, A. Jonathan Shaw<sup>p</sup>, Lisa DeGironimo<sup>s</sup>, Dennis W. Stevenson<sup>s</sup>, Barbara Surek<sup>o</sup>, Juan Carlos Villarreal<sup>t</sup>, Béatrice Roure<sup>u</sup>, Hervé Philippe<sup>u,v</sup>, Claude W. dePamphilis<sup>l</sup>, Tao Chen<sup>w</sup>, Michael K. Deyholos<sup>d</sup>, Regina S. Baucom<sup>x</sup>, Toni M. Kutchan<sup>y</sup>, Megan M. Augustin<sup>y</sup>, Jun Wang<sup>z</sup>, Yong Zhang<sup>y</sup>, Zhijian Tian<sup>z</sup>, Zhixiang Yan<sup>z</sup>, Xiaolei Wu<sup>z</sup>, Xiao Sun<sup>z</sup>, Gane Ka-Shu Wong<sup>d,z,aa,2</sup>, and James Leebens-Mack<sup>g,2</sup>

# ASTRAL-III on the capstone 1KP

[unpublished, 2018, 2019]

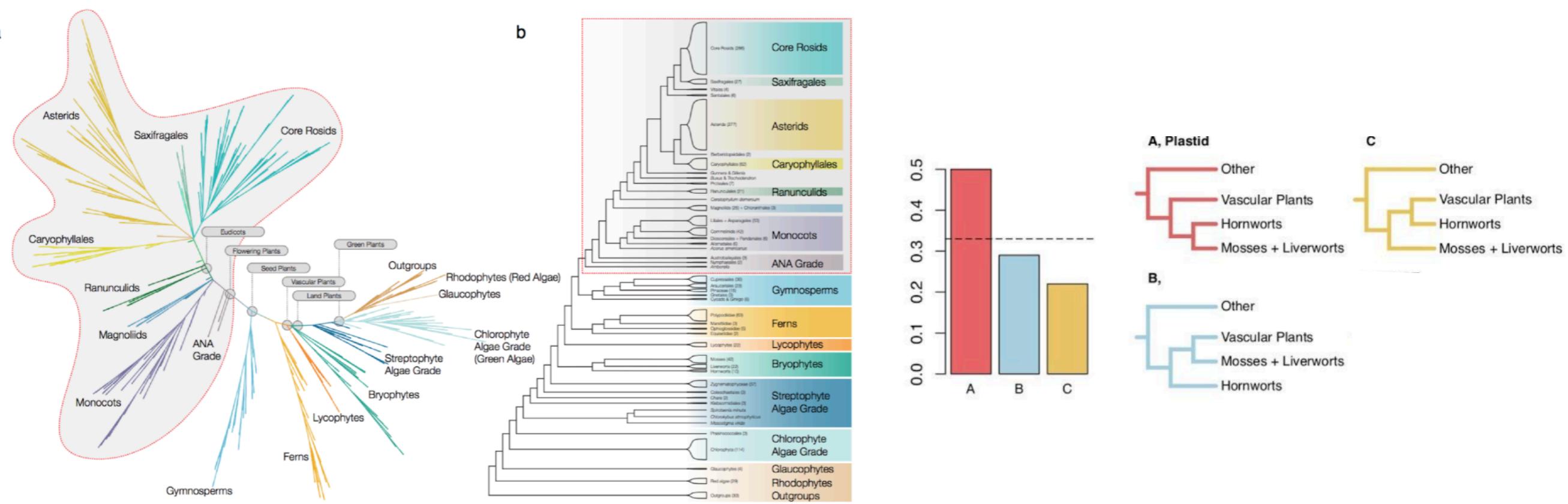
- ASTRAL is the main tree used
  - Took eight hours for 1178 species and 410 genes



# ASTRAL-III on the capstone 1KP

[unpublished, 2018, 2019]

- ASTRAL is the main tree used
  - Took eight hours for 1178 species and 410 genes
  - Sheds light on trait evolution and differences among genes

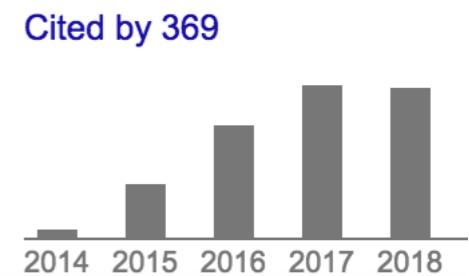


# ASTRAL used widely

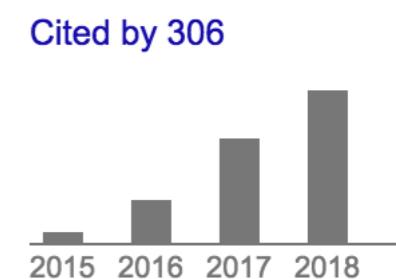
Early use:

- Plants: Wickett, et al., 2014, PNAS
- Birds: Prum, et al., 2015, Nature
- Xenoturbella, Cannon et al., 2016, Nature
- Xenoturbella, Rouse et al., 2016, Nature
- Flatworms: Laumer, et al., 2015, eLife
- Shrews: Giarla, et al., 2015, Syst. Bio.
- Frogs: Yuan et al., 2016, Syst. Bio.
- Tomatoes: Pease, et al., 2016, PLoS Bio.
- Angiosperms: Huang et al., 2016, MBE
- Worms: Andrade, et al., 2015, MBE

ASTRAL



ASTRAL-II

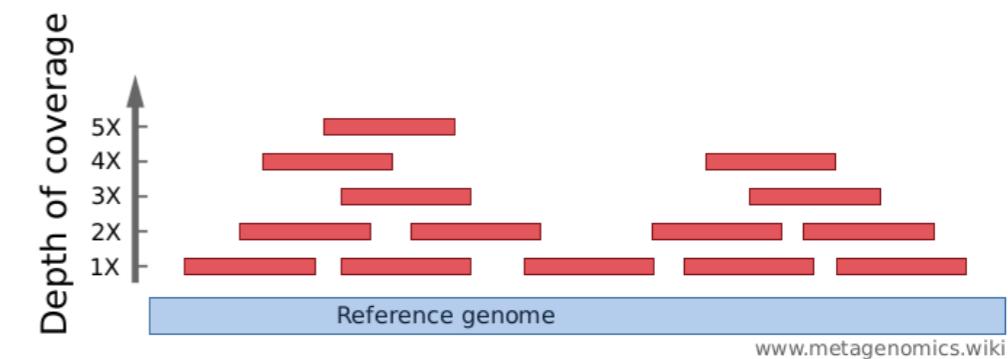


ASTRAL-III



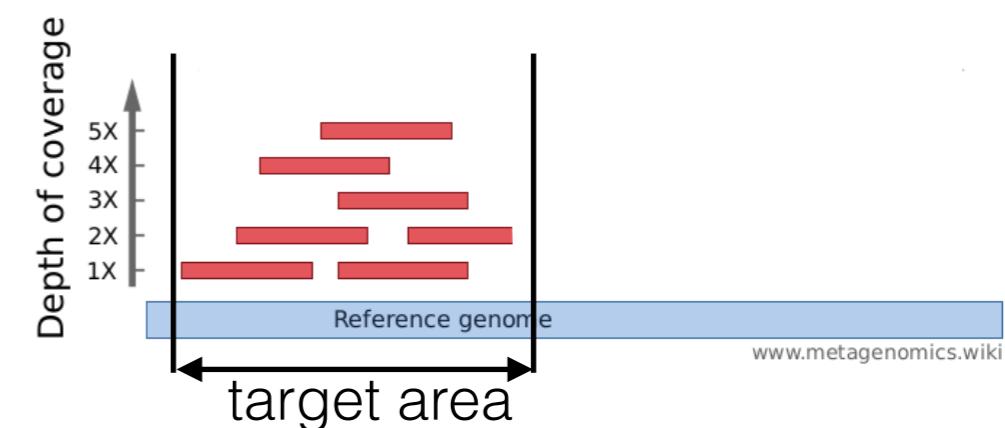
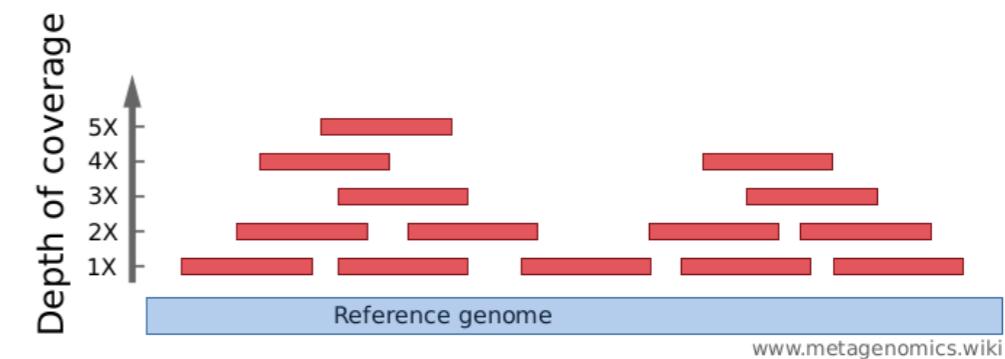
# How do we get genome-wide data?

- **Whole-genome** sequencing:  
very expensive but complete



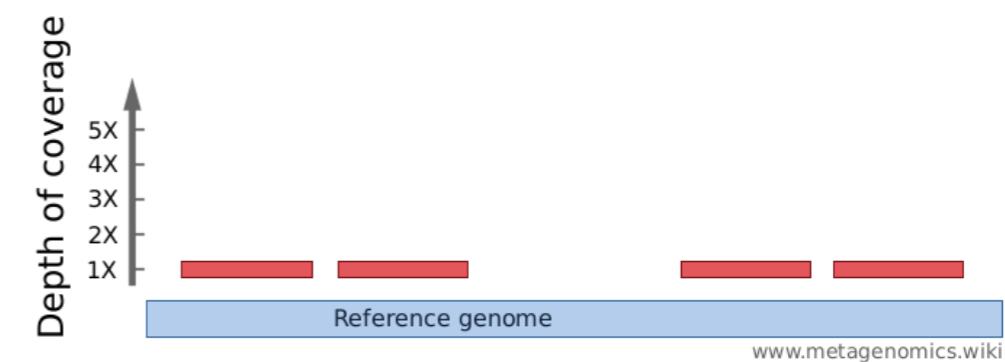
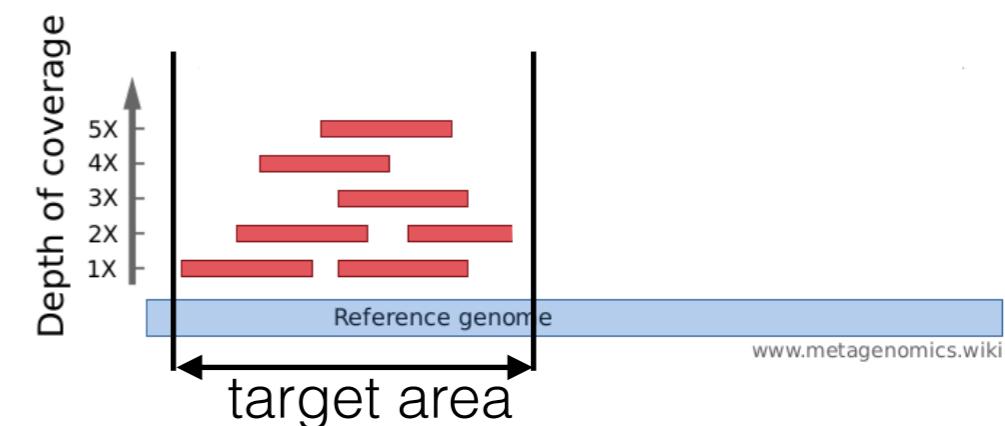
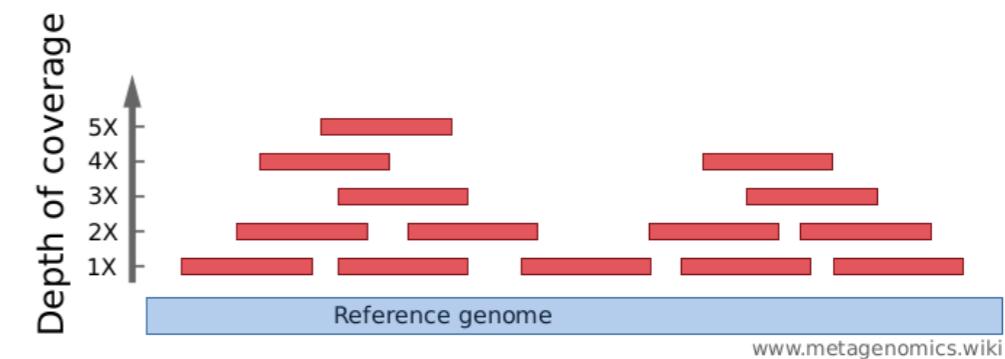
# How do we get genome-wide data?

- **Whole-genome** sequencing: very expensive but complete
- **Transcriptome** sequencing: still expensive but much less data
- **Targeted capture**: a bit less expensive but less useful data



# How do we get genome-wide data?

- **Whole-genome** sequencing: very expensive but complete
- **Transcriptome** sequencing: still expensive but much less data
- **Targeted capture**: a bit less expensive but less useful data
- **Genome skimming**: cheap (\$50) but little data and hard to use



# Genome skimming

- Sequence genomes at very **low coverage**
  - Not enough for assemblies
  - May be enough for **sample identification**: find the closest species to a query skim
  - Build phylogenetic trees with reduced cost



Tom Gilbert,  
Copenhagen



guojie zhang

华大基因  
**BGI**

# Genome skimming

- Sequence genomes at very **low coverage**
  - Not enough for assemblies
  - May be enough for **sample identification**: find the closest species to a query skim
  - Build phylogenetic trees with reduced cost
  - Both applications boil down to **computing distances** between genomes



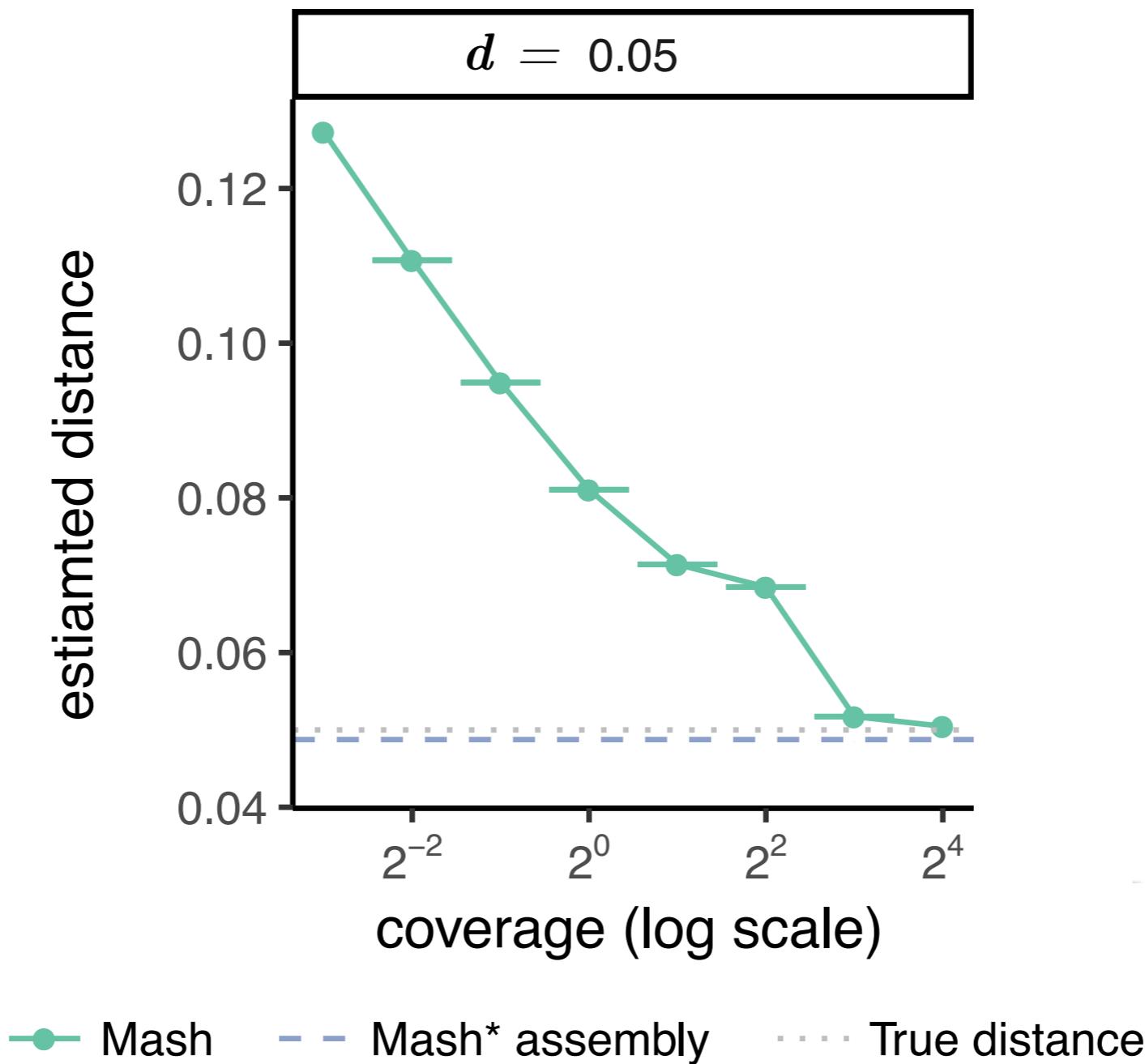
Tom Gilbert,  
Copenhagen



guojie zhang

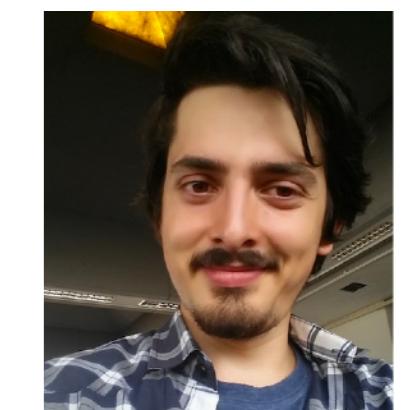
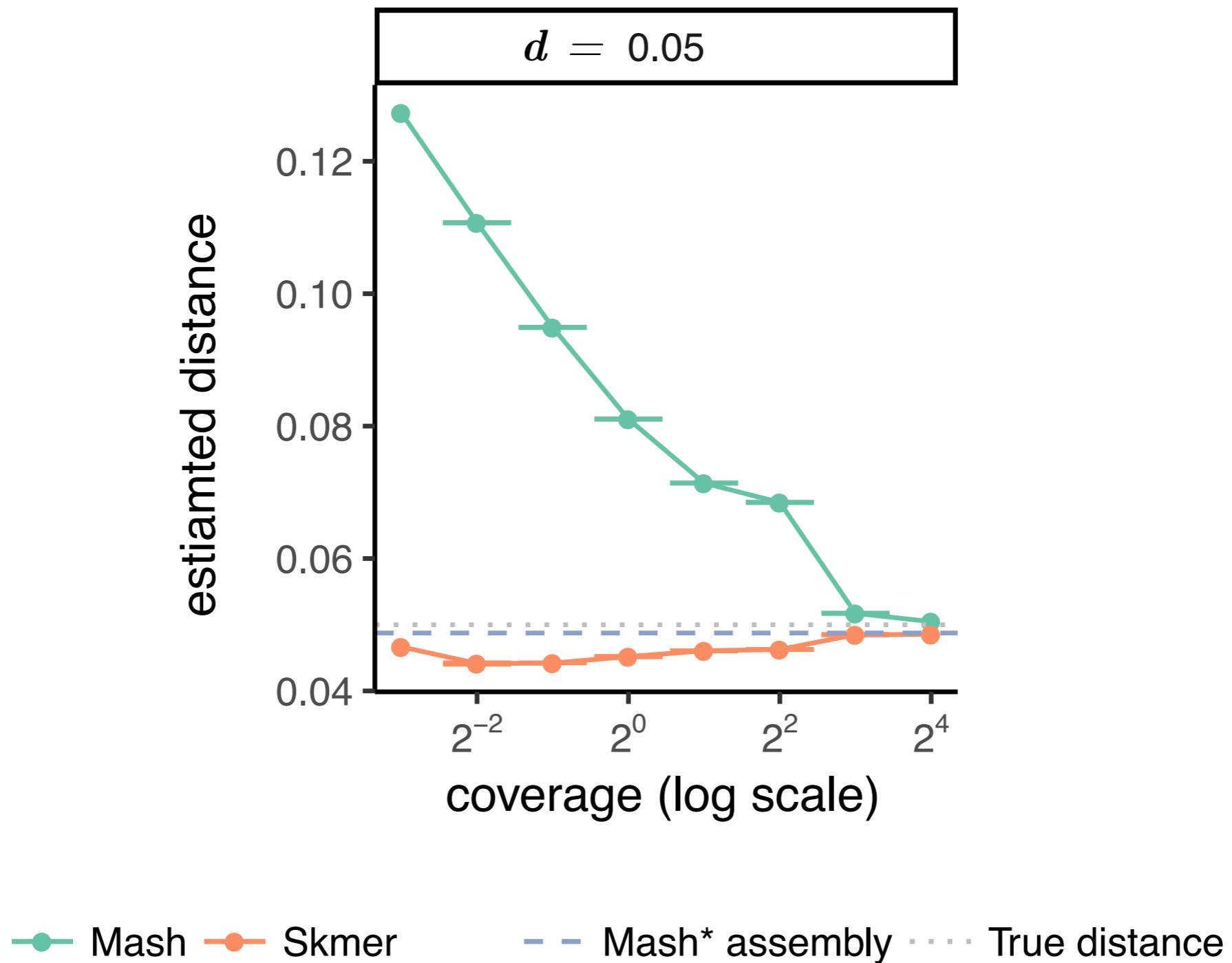


# What coverage is needed?

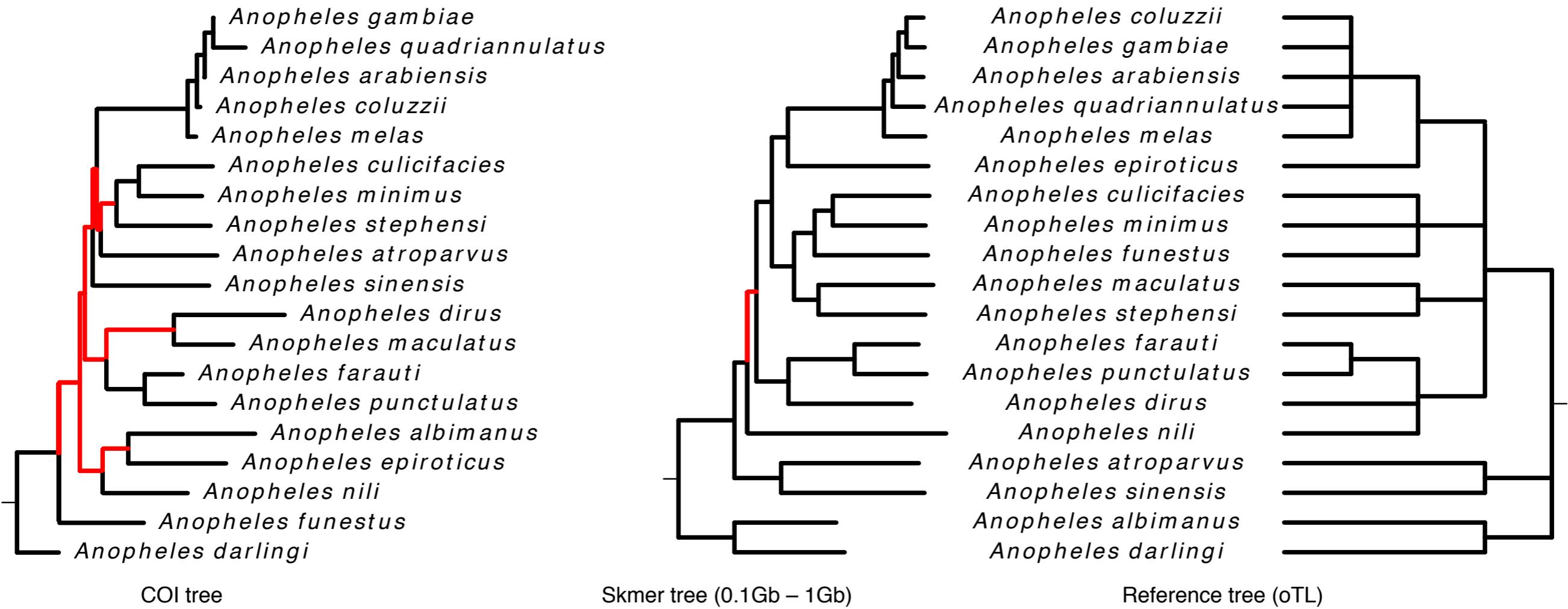


# Skmer: a better translation from $J$ to $D$

[Sarmashghi, et al, Genome biolog, 2018]



# Building Phylogenies using Skimmed data



# Conclusions

- Evolution can be reconstructed using genomic data
  - It ENABLEs answering basic science questions in biology and has many downstream applications
- We constantly get bigger and **bigger data**, but
  - With more data, we often want to solve **harder problems**
  - Large data **increase the running time** but *can reduce accuracy of methods* not designed for large size
  - There is a need for developing **new algorithms** that can handle large dataset

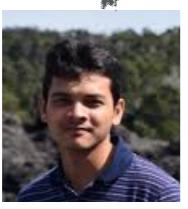
# Acknowledgments



Tandy Warnow



Siavash  
Mirarab



S.M. Bayzid



Nam Nguyen  
(now at UIUC)



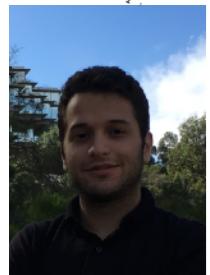
Jim Leebens-mack  
(UGA)



Norman Wickett  
(U Chicago)



Gane Wong  
(U of Alberta)



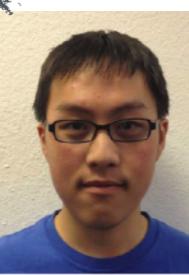
Erfan  
Sayyari



Chao Zhang



Maryam  
Hashemi



John Yin



Uyen Mai



Guojie Zhang  
(BGI, China)



Tom Gilbert  
(U Copenhagen)



Erich Jarvis  
(Duke, HMMI)



Bastien Boussau  
(Université Lyon)



Ed Braun  
(U Florida)