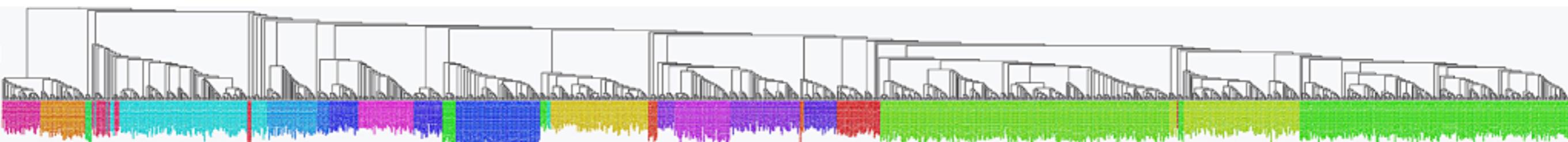
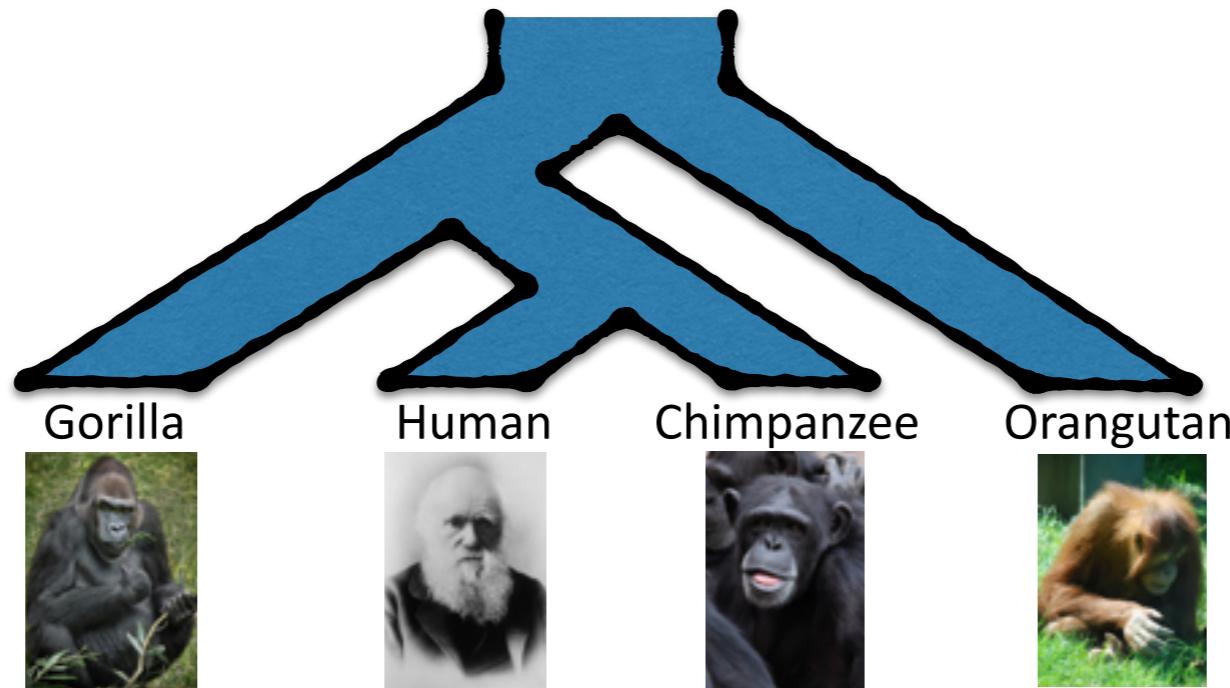


# Reconstruction of species histories using genomic data

Siavash Mirarab  
University of California San Diego



# Statistical inference of phylogenies



Gorilla  
ACTGCACACCG



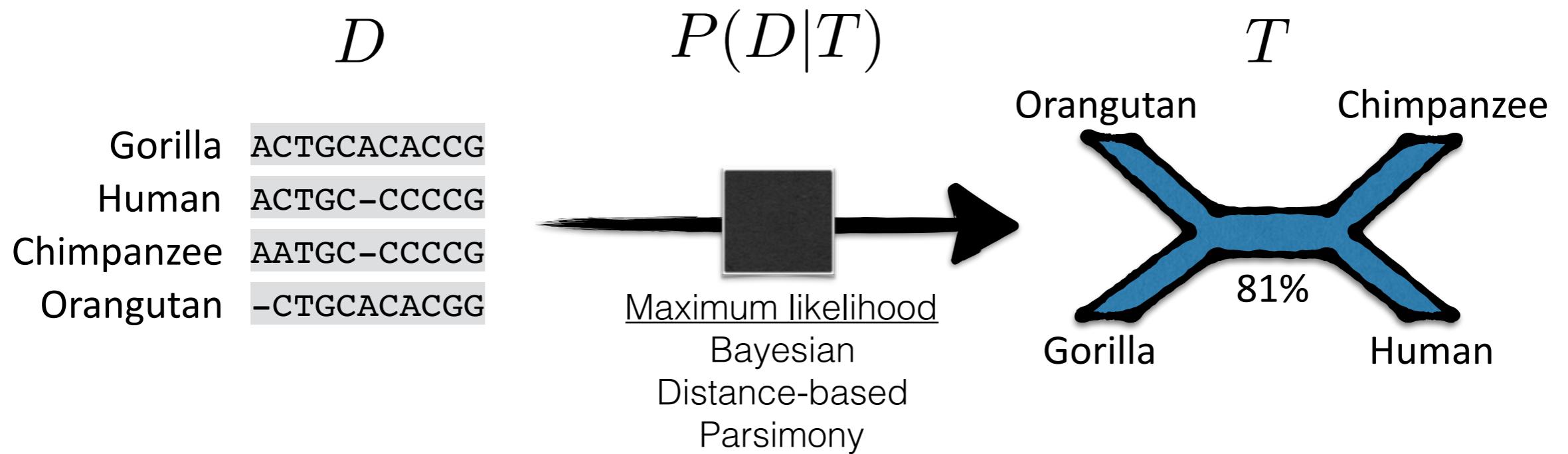
Human  
ACTGCCCG



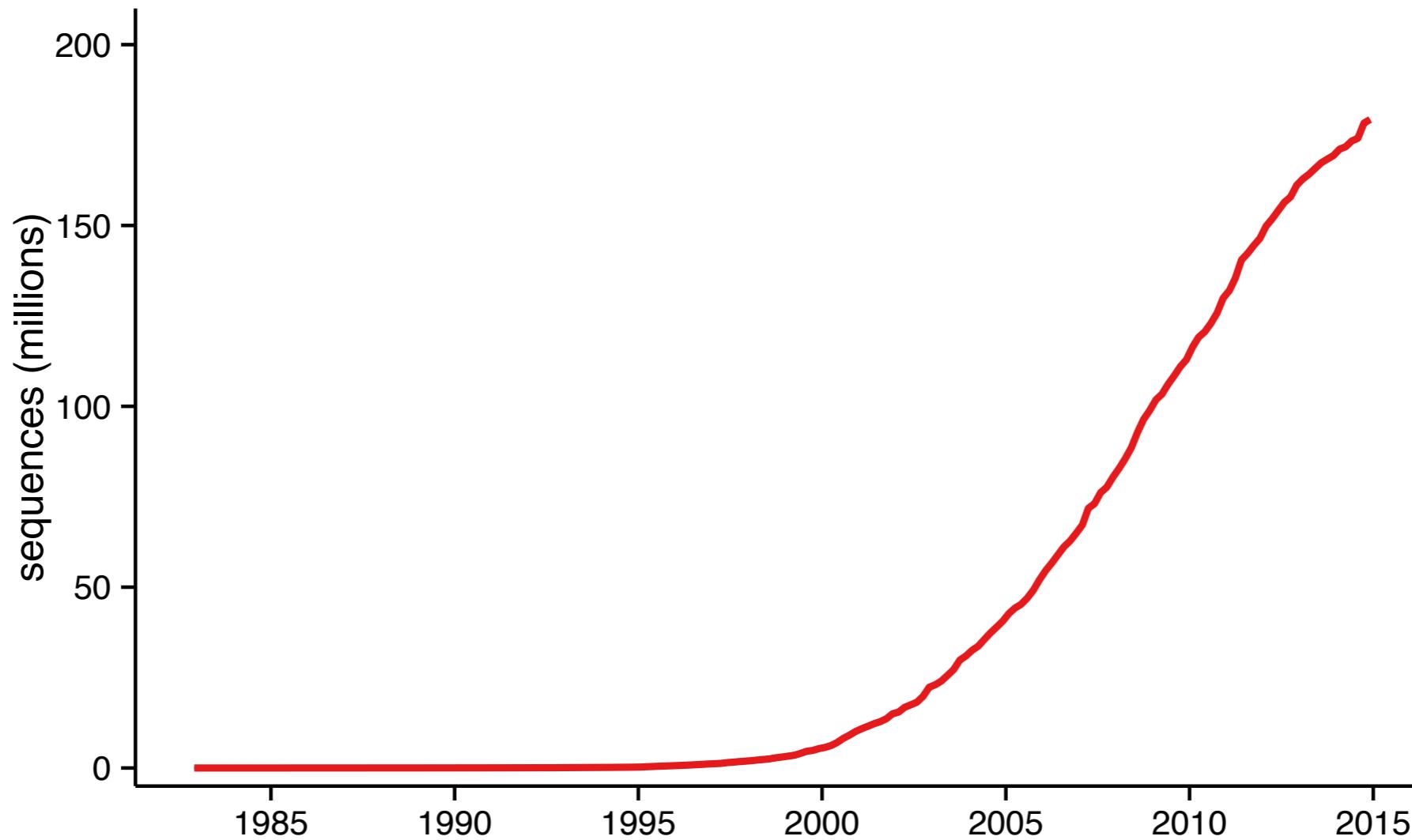
Chimpanzee  
AATGCCCG



Orangutan  
CTGCACACGG

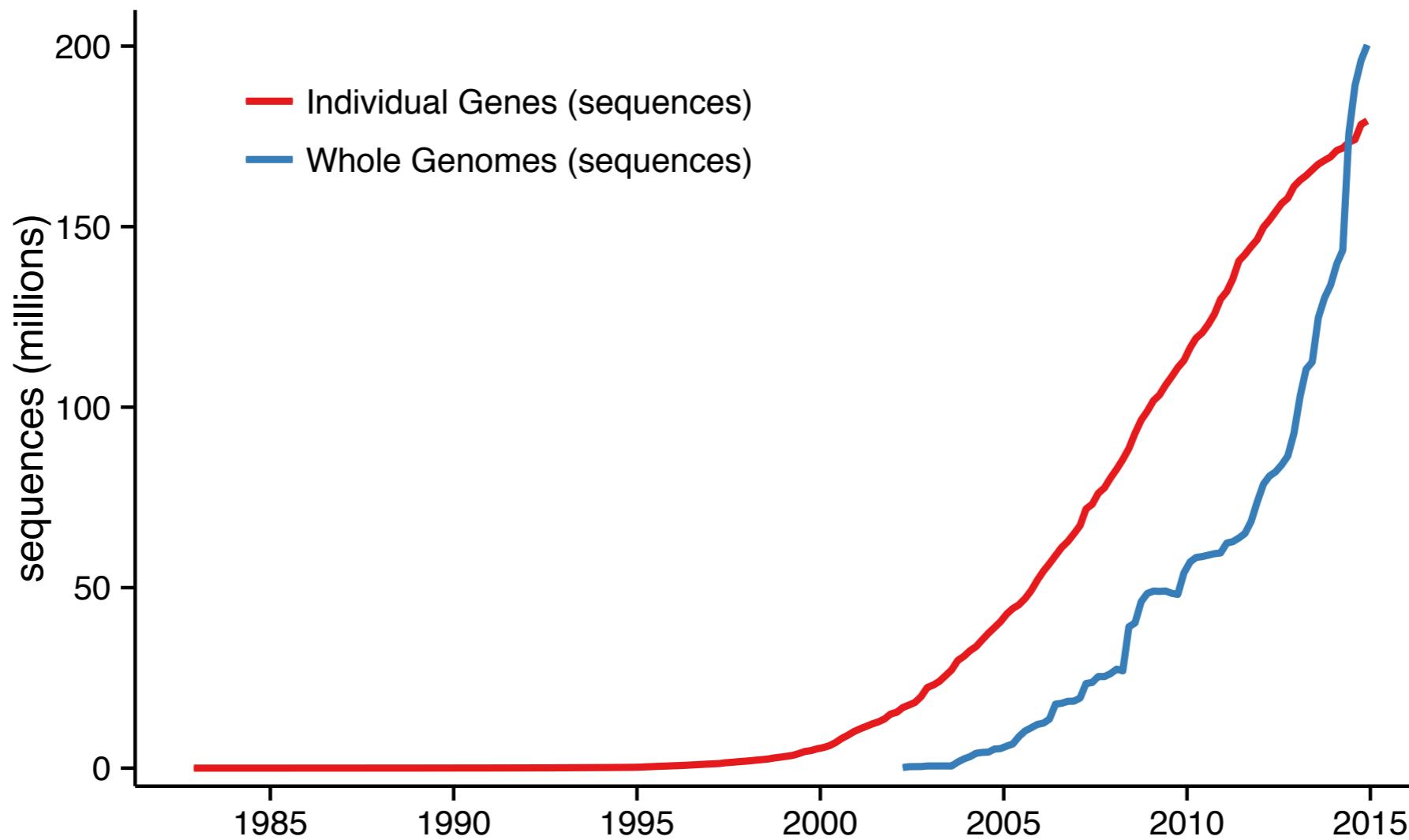


# Sequence data growth



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

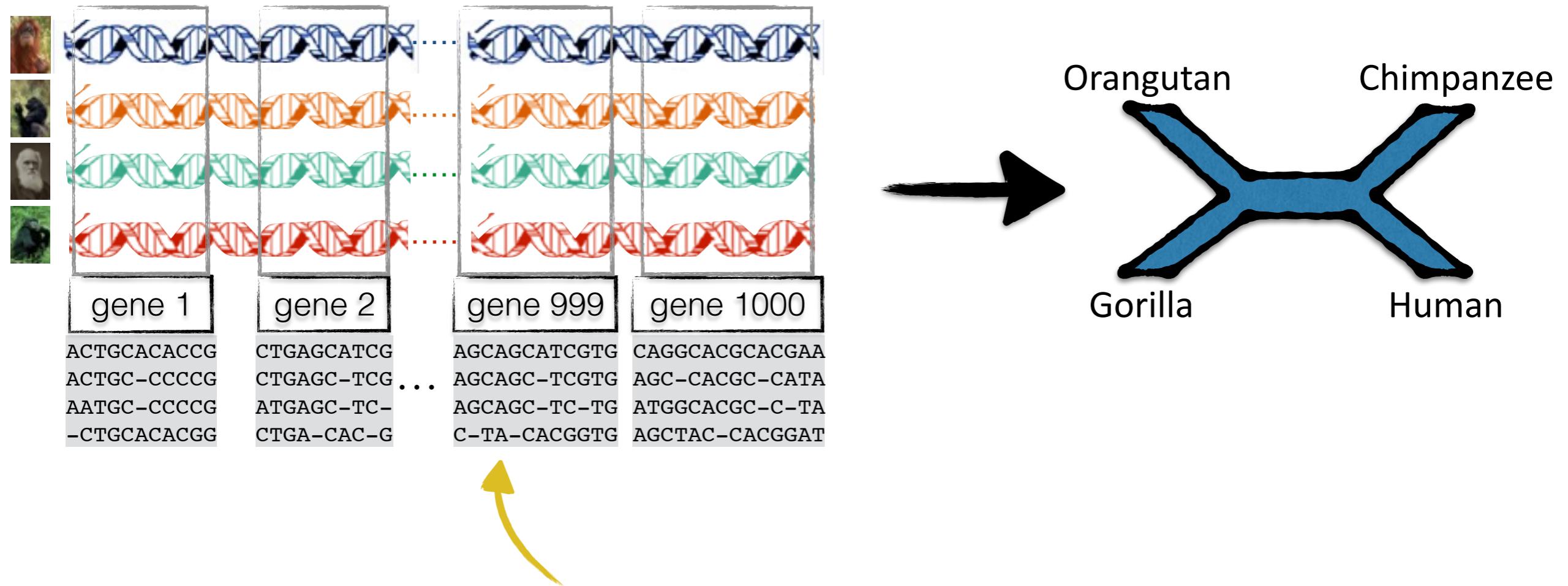
# Sequence data growth



Growth in Genebank sequence data

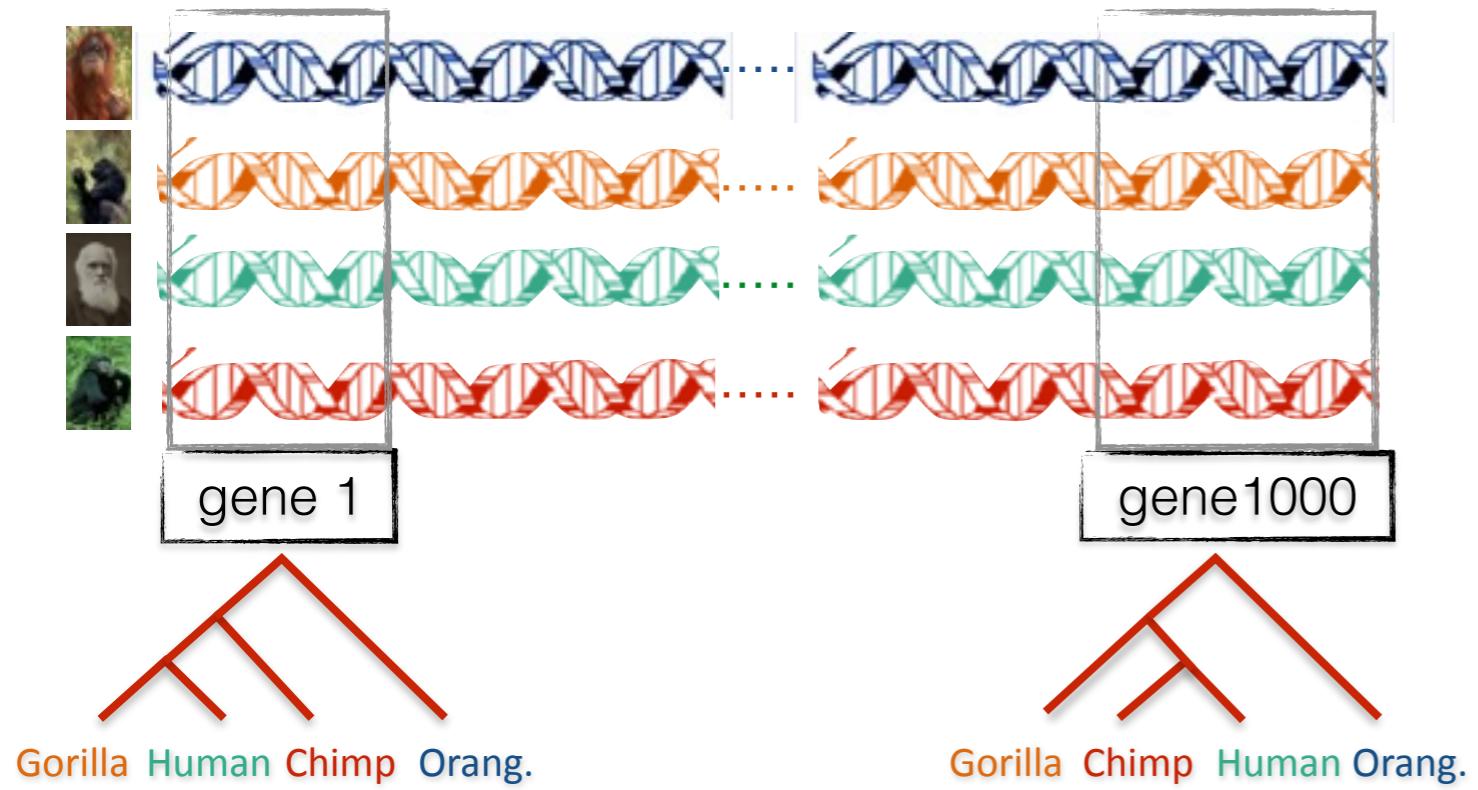
data source: NCBI (<http://www.ncbi.nlm.nih.gov/genbank/statistics>)

# Phylogenomics

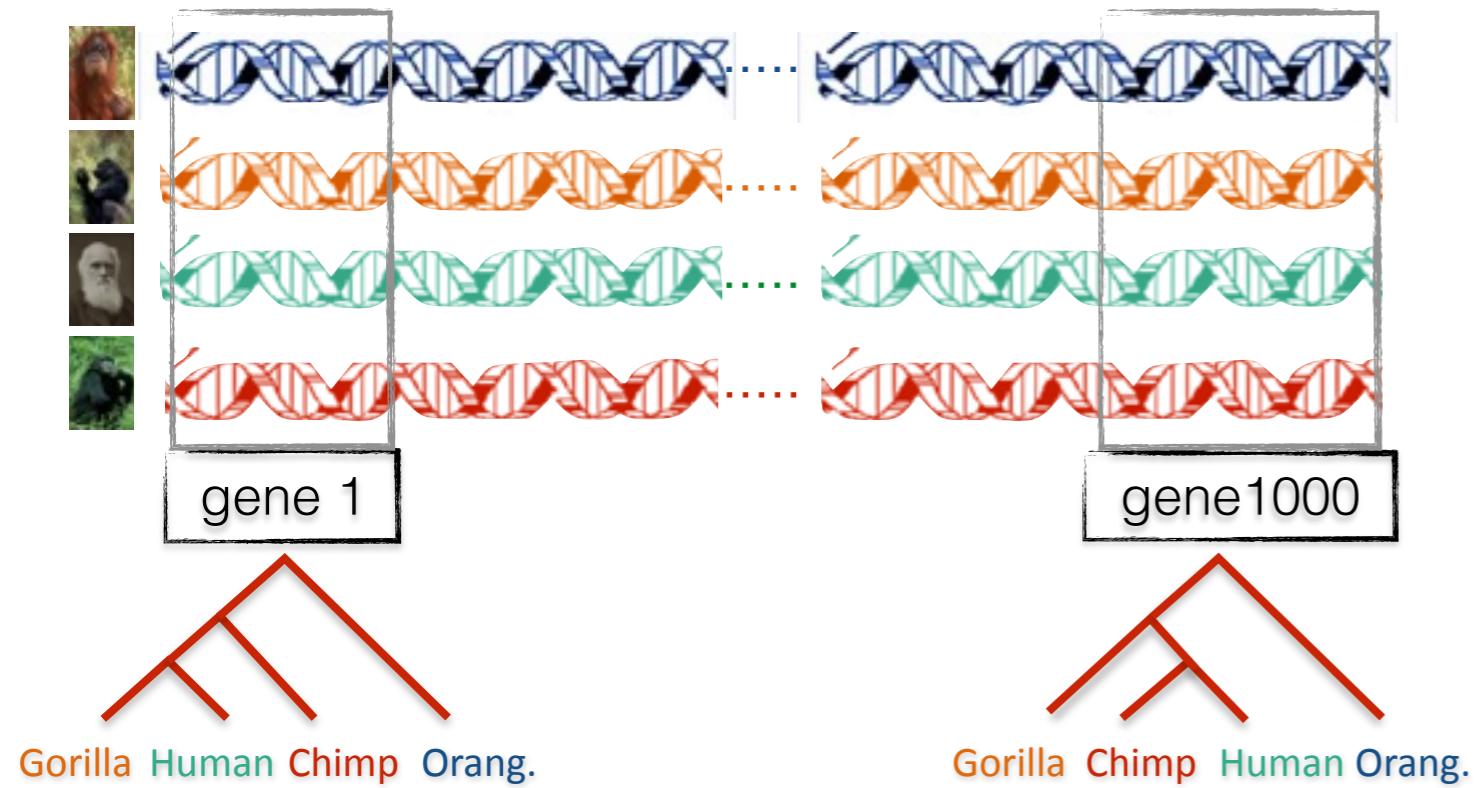


I'll use the term "gene" to refer to "c-genes":  
recombination-free orthologous stretches of the genome

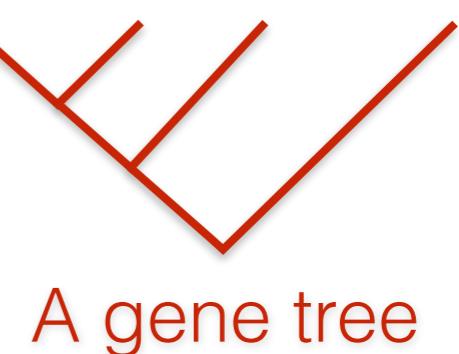
# Gene tree discordance



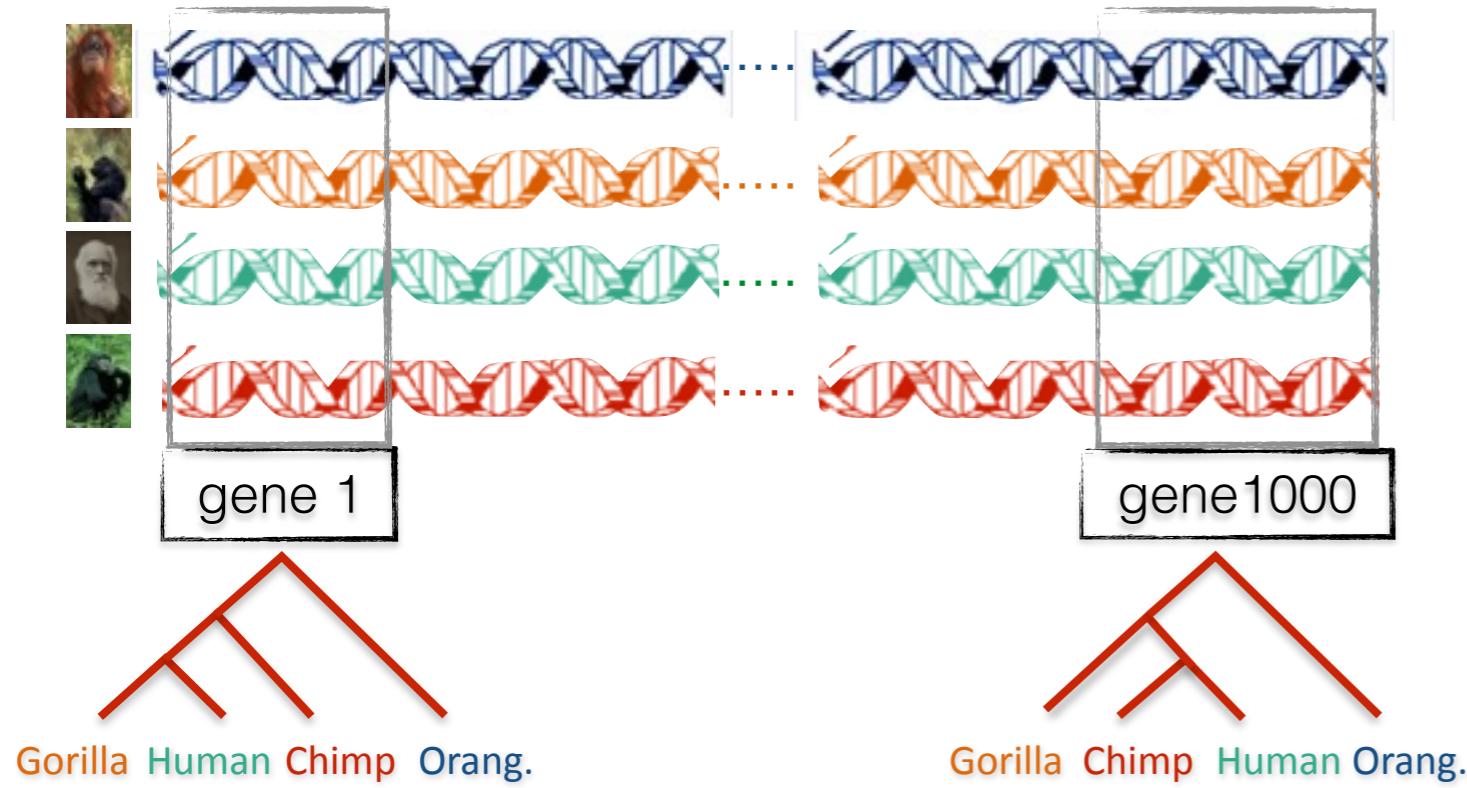
# Gene tree discordance



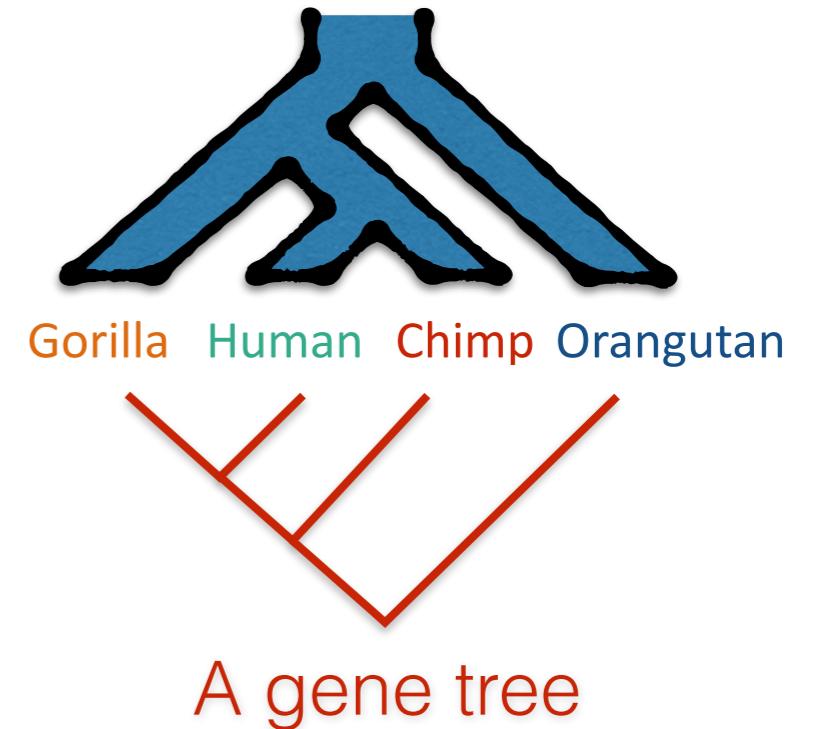
**The species tree**



# Gene tree discordance



The species tree

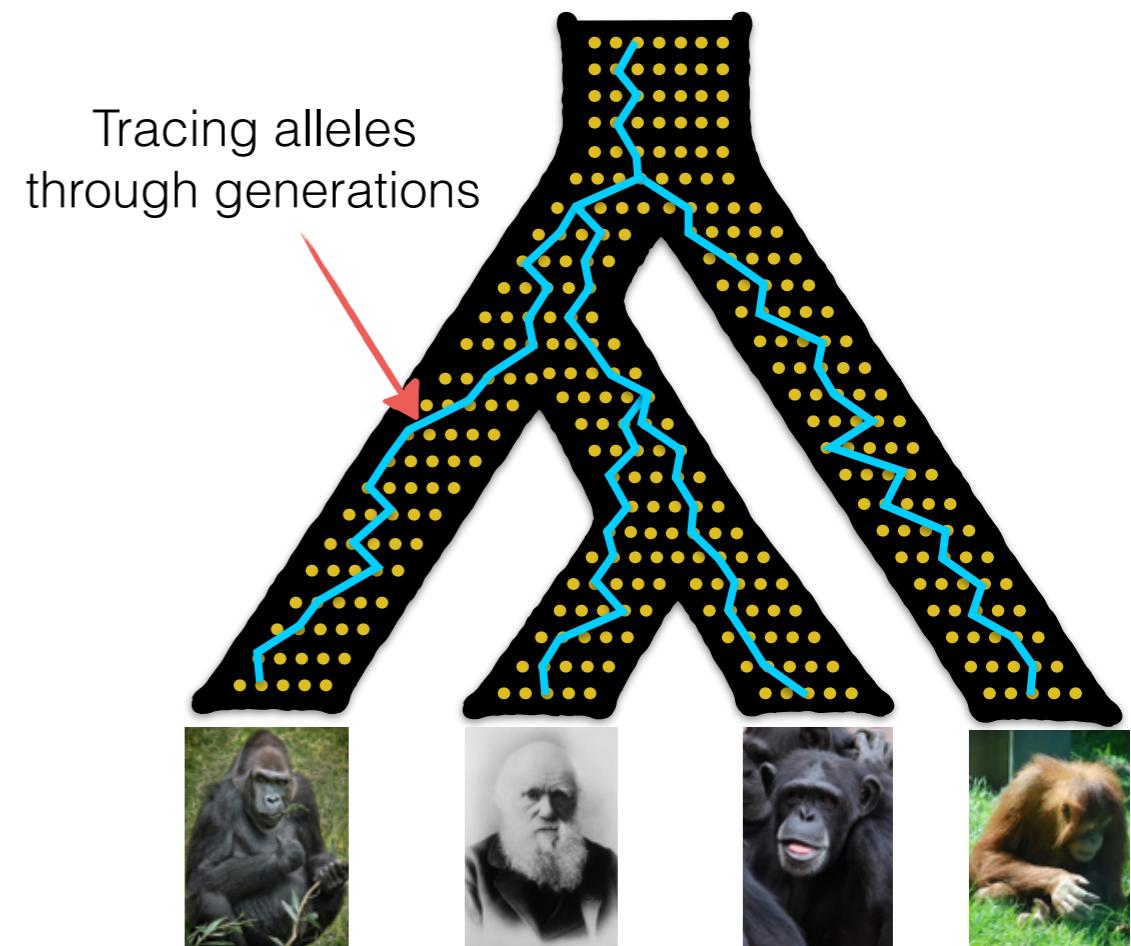


**Causes of gene tree discordance include:**

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

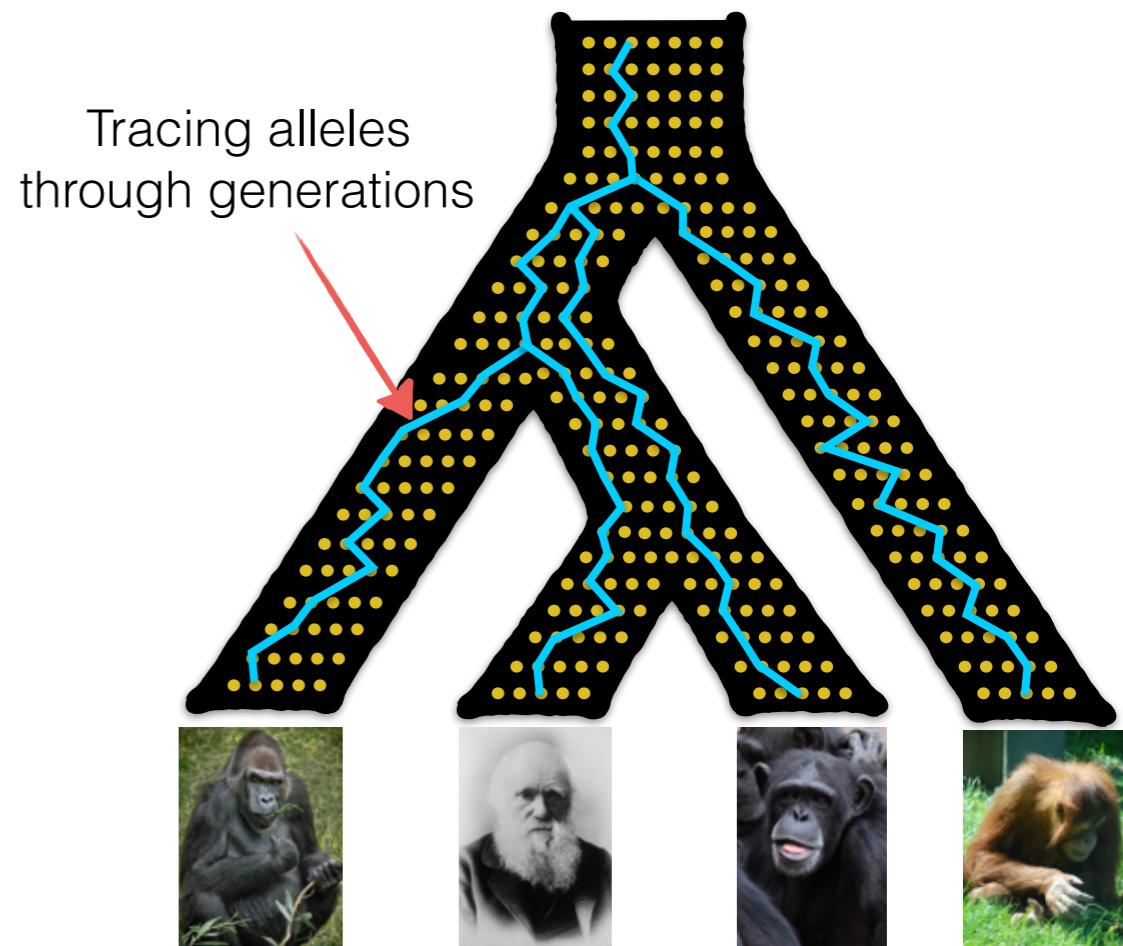
# Incomplete Lineage Sorting (ILS)

- A random process related to having multiple versions of each gene in a population



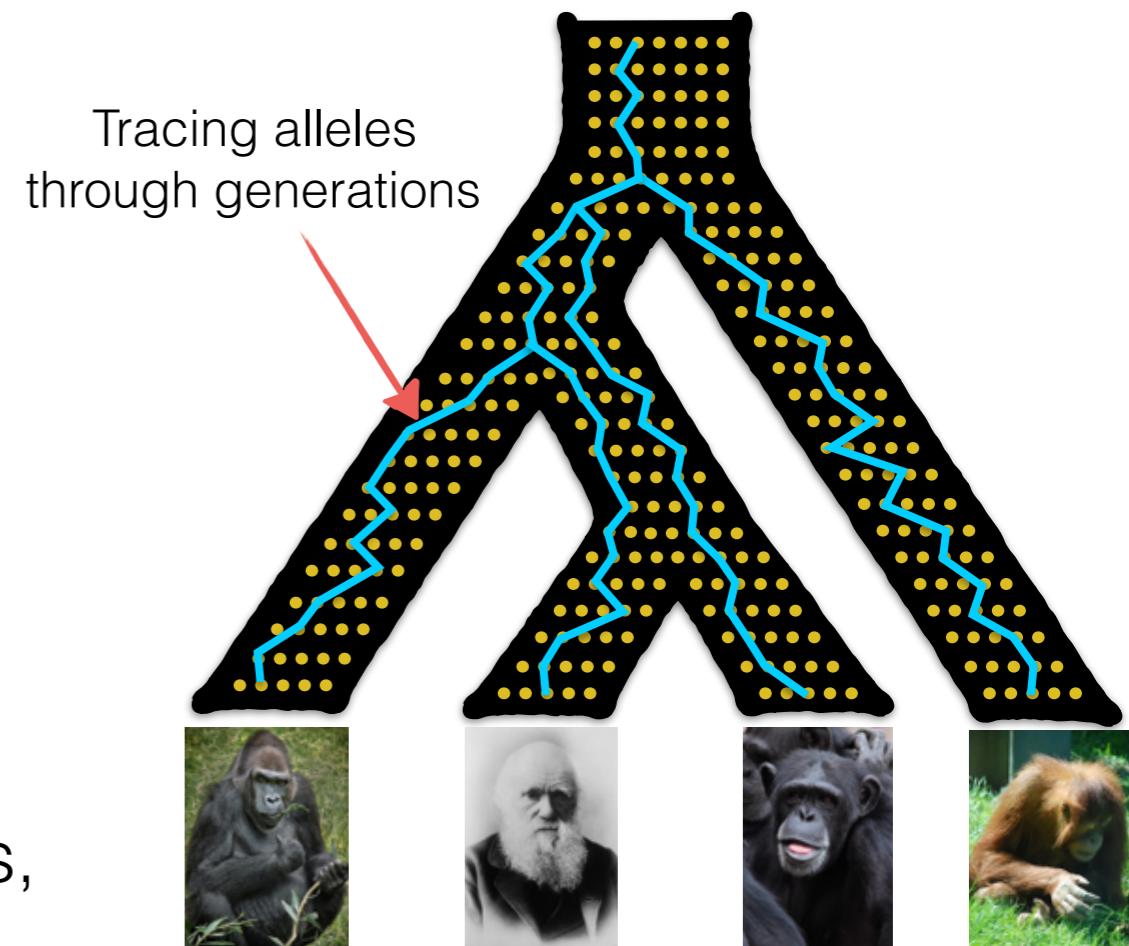
# Incomplete Lineage Sorting (ILS)

- A random process related to having multiple versions of each gene in a population
- Omnipresent; most likely for short branches and large population size

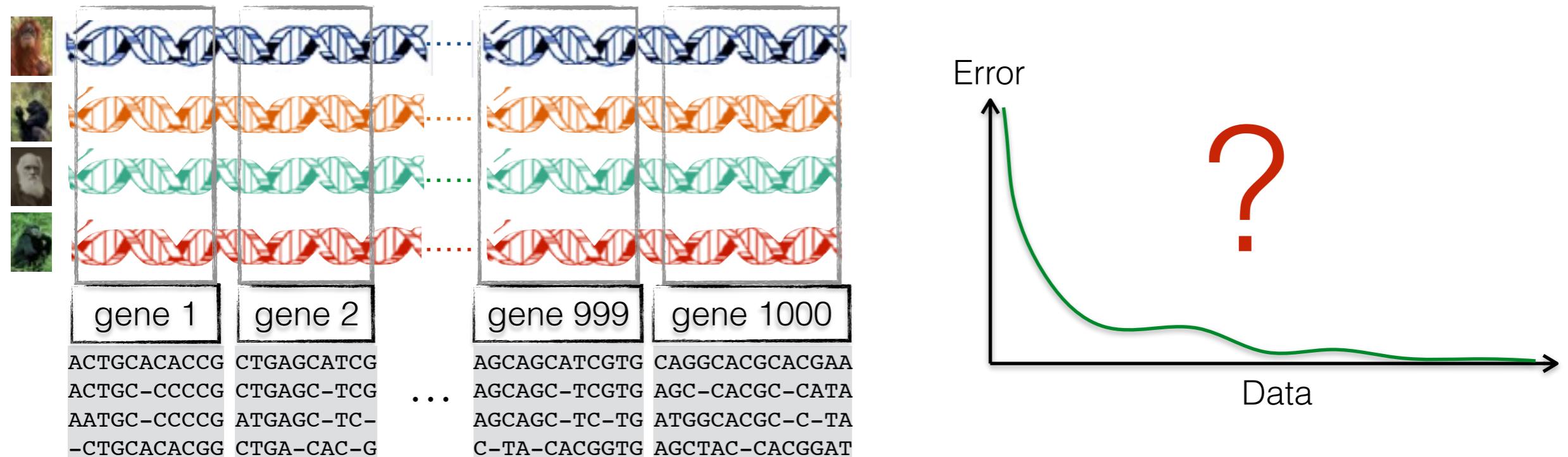


# Incomplete Lineage Sorting (ILS)

- A random process related to having multiple versions of each gene in a population
- Omnipresent; most likely for short branches and large population size
- We have statistical models of ILS (multi-species coalescent)
  - The species tree **defines the probability distribution** on gene trees, and is **identifiable** from the distribution on gene trees  
[Degnan and Salter, Int. J. Org. Evolution, 2005]

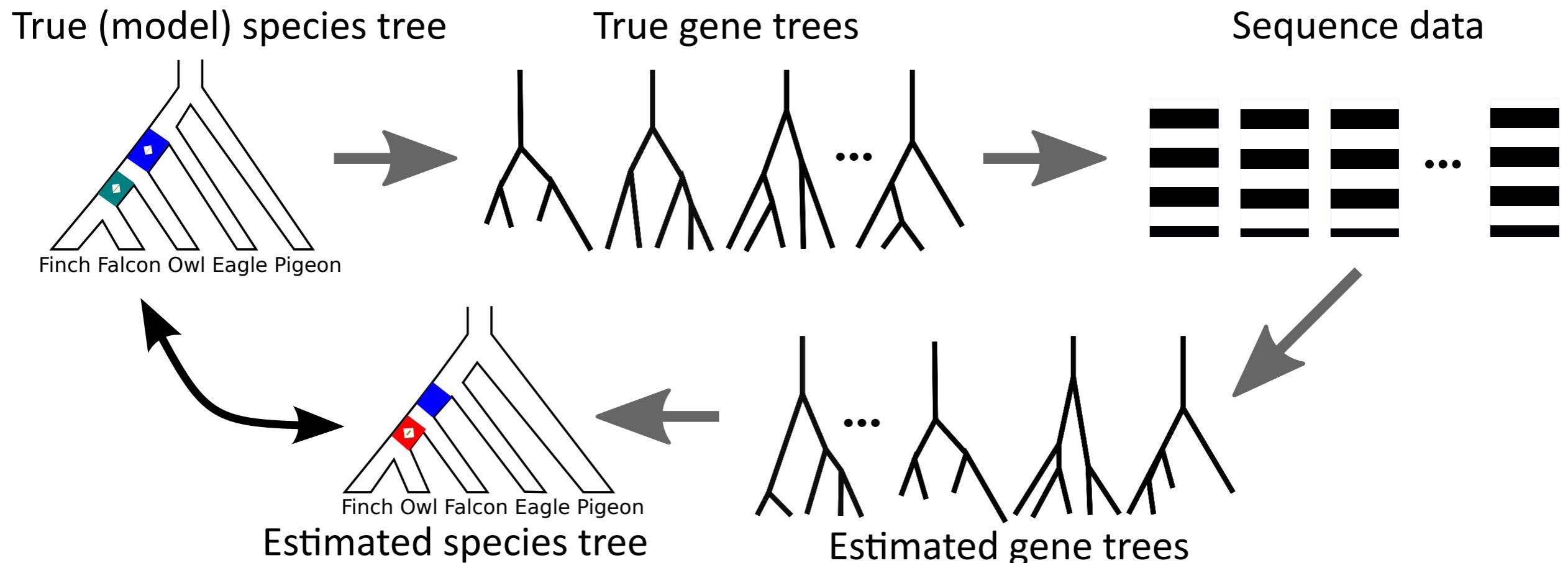


# Increased data → accuracy?



- A. **Theoretical statistical guarantees:** assuming data are generated under the multi-species coalescent model, we aspire to statistical consistency
- B. **Simulation studies**, generating synthetic data according to the multi-species coalescent model and models of sequence evolution

# Simulation studies

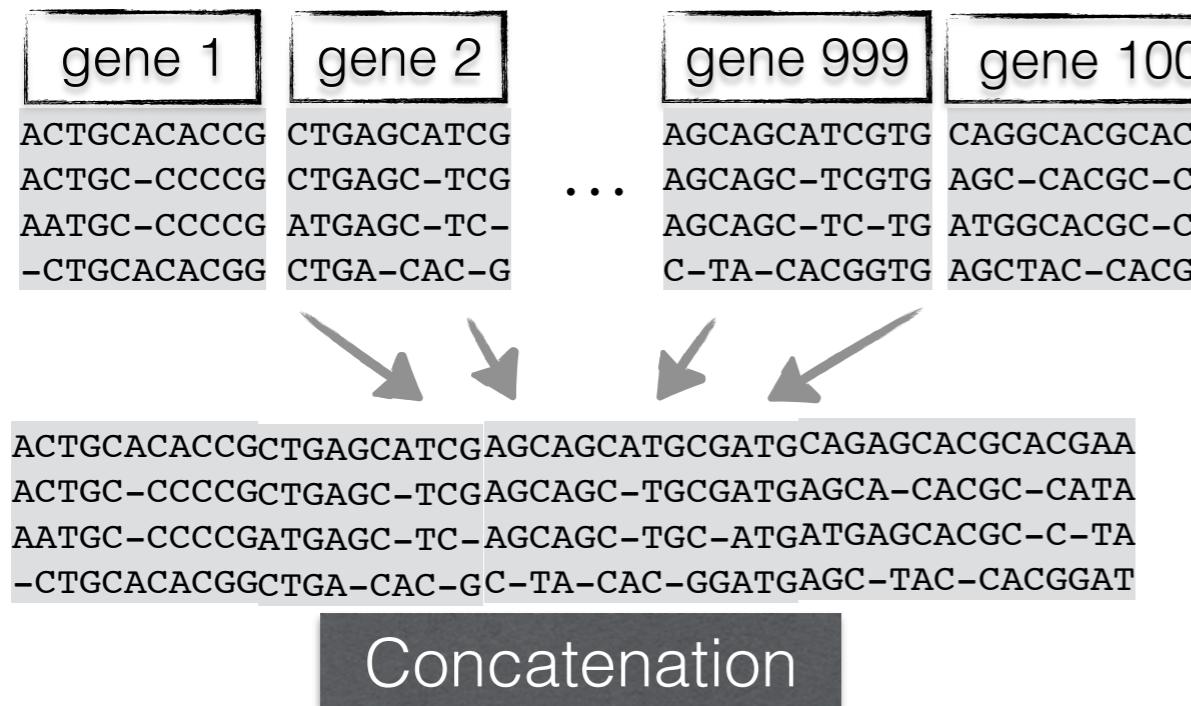


Error metric: percentage of branches in true tree that are missing from the estimated tree

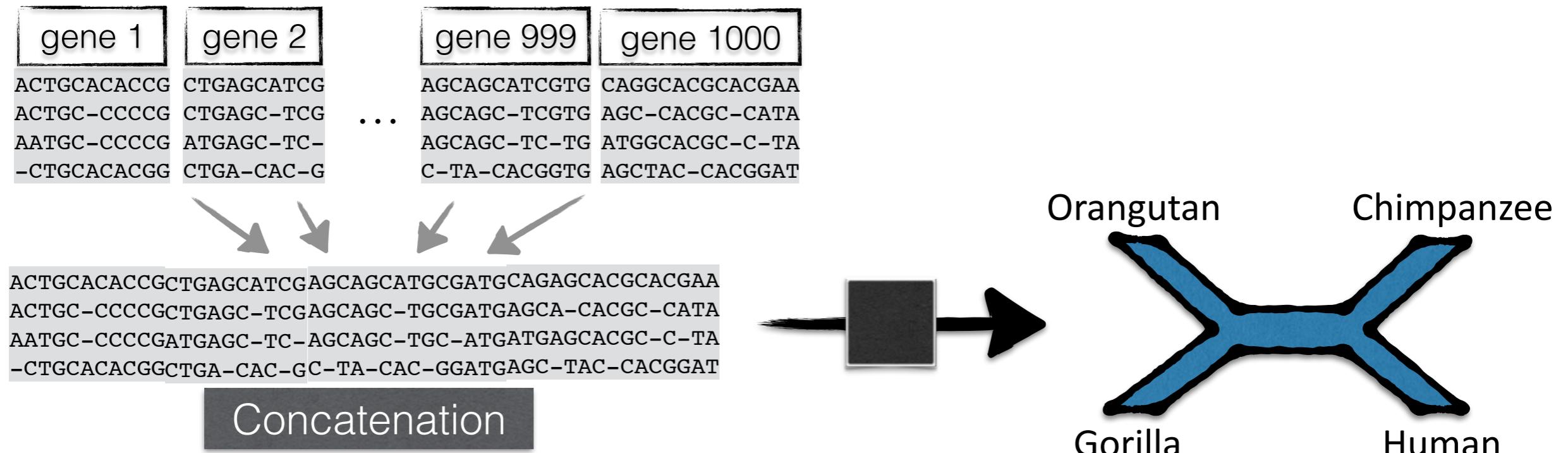
# Species tree estimation from phylogenomic data (Approach 1: concatenation)

gene 1	gene 2	gene 999	gene 1000
ACTGCACACCG	CTGAGCATCG	AGCAGCATTGTG	CAGGCACGCACGAA
ACTGC-CCCCG	CTGAGC-TCG	AGCAGC-TCGTG	AGC-CACGC-CATA
AATGC-CCCCG	ATGAGC-TC-	AGCAGC-TC-TG	ATGGCACGC-C-TA
-CTGCACACGG	CTGA-CAC-G	C-TA-CACGGTG	AGCTAC-CACGGAT

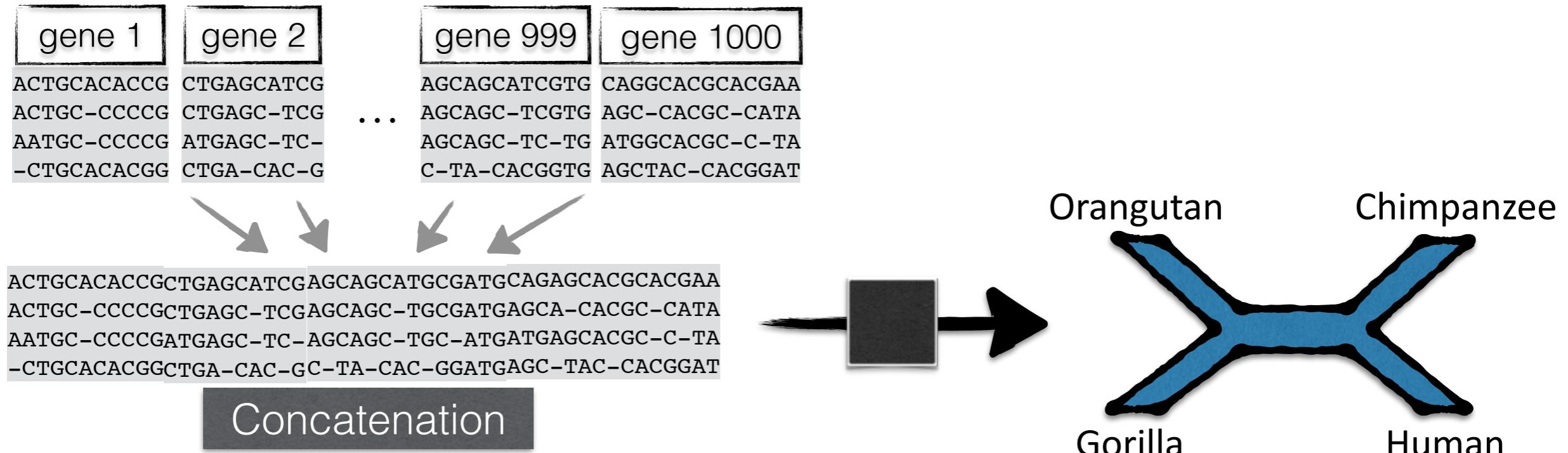
# Species tree estimation from phylogenomic data (Approach 1: concatenation)



# Species tree estimation from phylogenomic data (Approach 1: concatenation)



# Species tree estimation from phylogenomic data (Approach 1: concatenation)



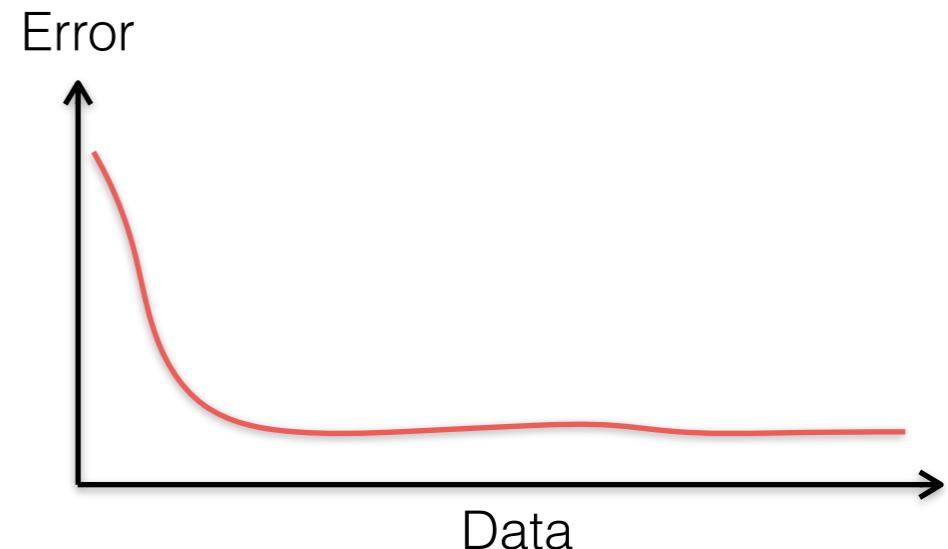
Statistically inconsistent & positively misleading  
(proof for unpartitioned maximum likelihood)

[Roch and Steel, Theo. Pop. Gen., 2014]

Mixed accuracy in simulations

[Kubatko and Degnan, Systematic Biology, 2007]

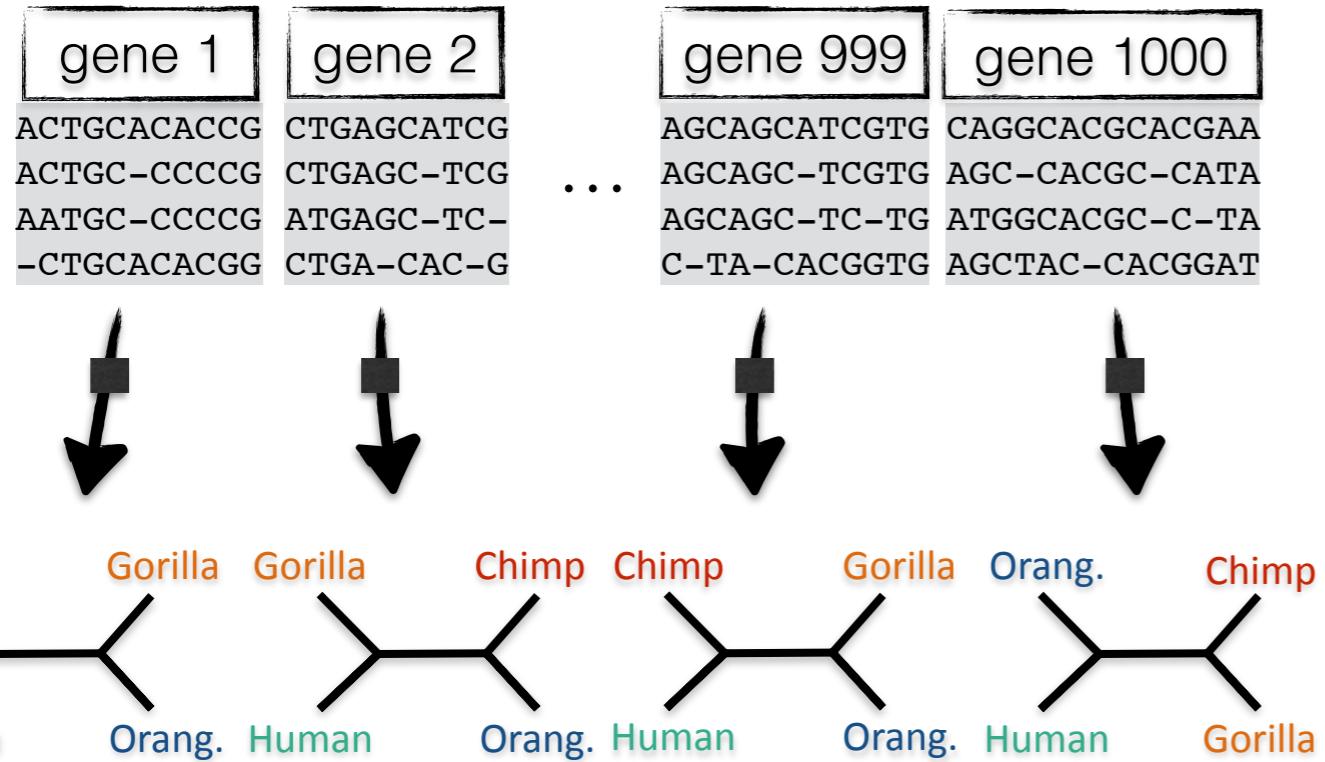
[Mirarab, et al., Systematic Biology, 2014]



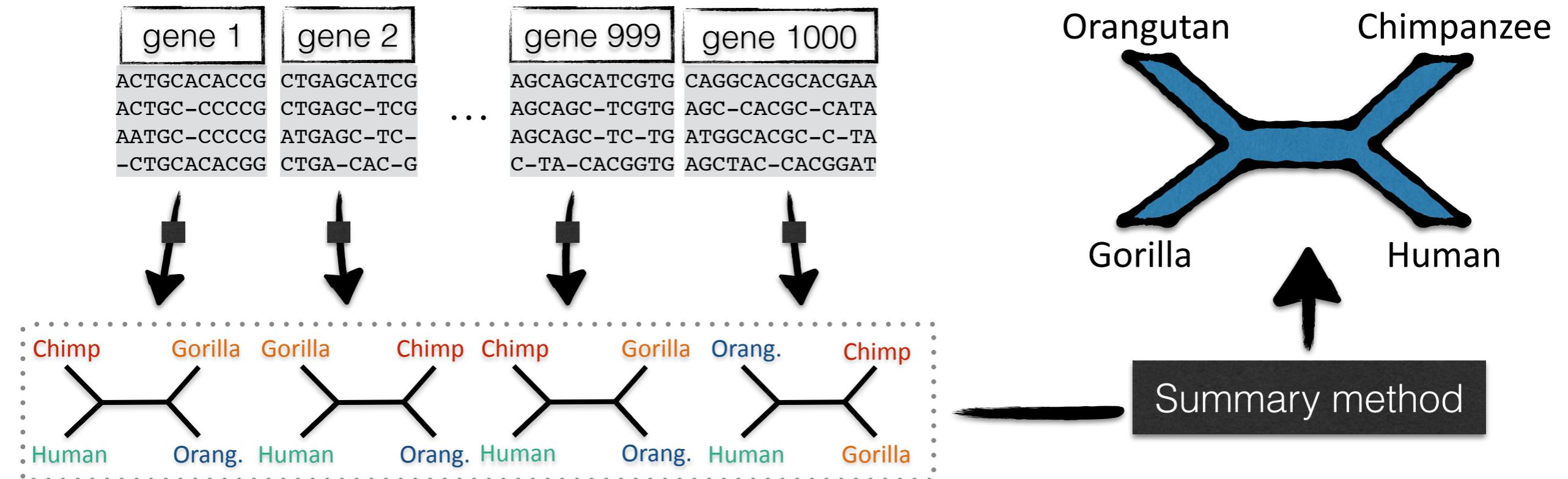
# Species tree estimation from phylogenomic data (Approach 2: summary methods)

gene 1	gene 2	gene 999	gene 1000
ACTGCACACCG	CTGAGCATCG	AGCAGC ATCGTG	CAGGCACGCACGAA
ACTGC-CCCCG	CTGAGC-TCG	AGCAGC-TCGTG	AGC-CACGC-CATA
AATGC-CCCCG	ATGAGC-TC-	AGCAGC-TC-TG	ATGGCACGC-C-TA
-CTGCACACGG	CTGA-CAC-G	C-TA-CACGGTG	AGCTAC-CACGGAT

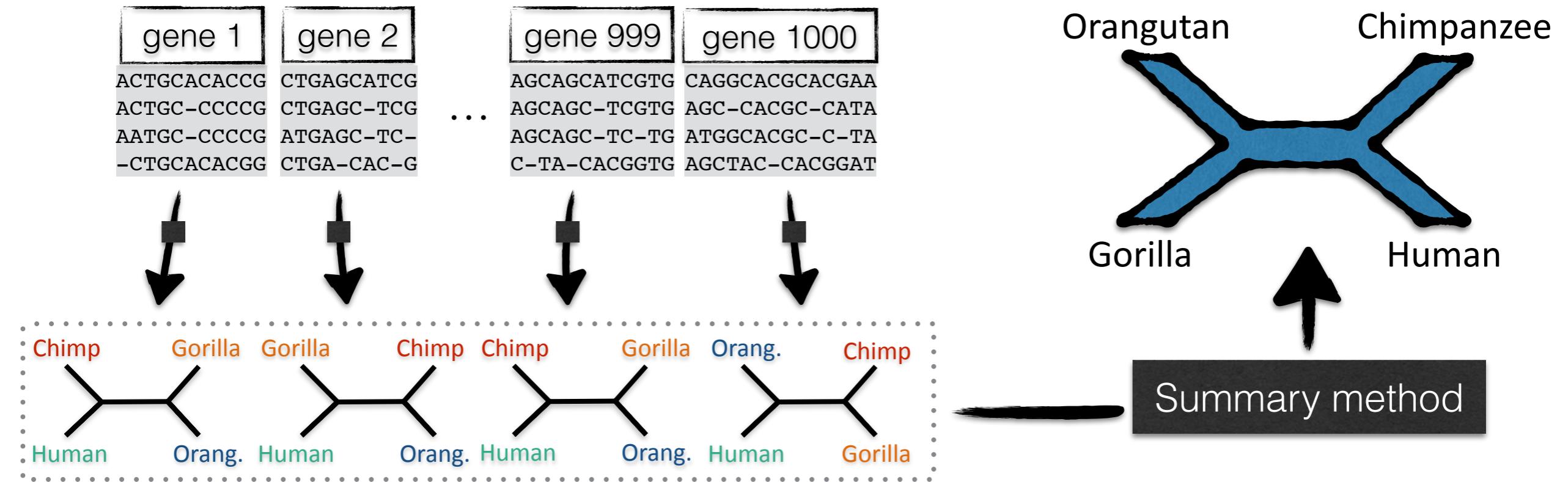
# Species tree estimation from phylogenomic data (Approach 2: summary methods)



# Species tree estimation from phylogenomic data (Approach 2: summary methods)

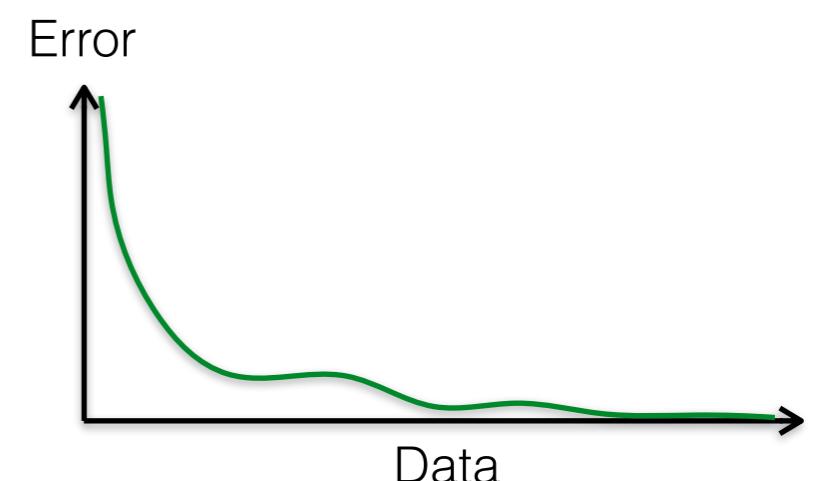


# Species tree estimation from phylogenomic data (Approach 2: summary methods)

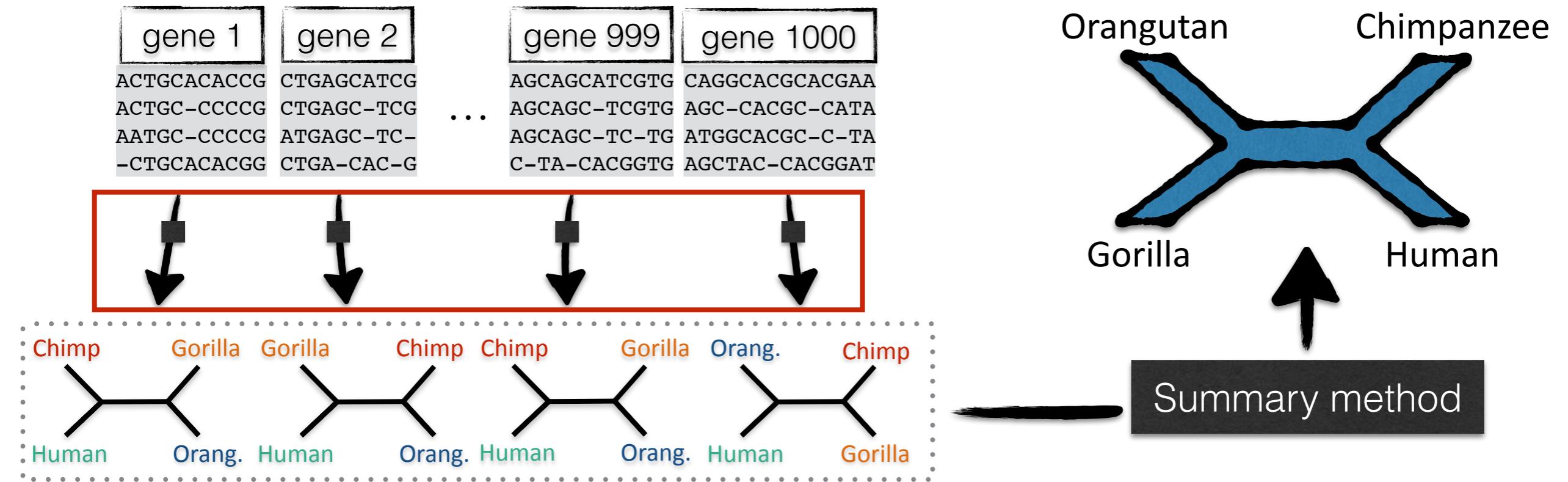


Can be statistically consistent

- **MP-EST** (maximum pseudo-likelihood)  
[Liu, Yu, Edwards, BMC Evol. Bio., 2010]
- **NJst**, STAR, STELLS, ... BUCKy-population

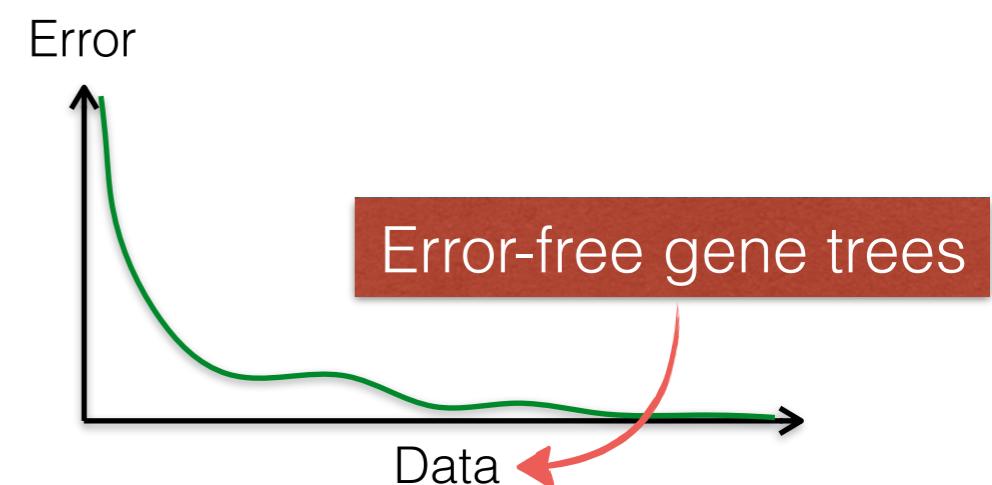


# Species tree estimation from phylogenomic data (Approach 2: summary methods)



Can be statistically consistent

- **MP-EST** (maximum pseudo-likelihood)  
[Liu, Yu, Edwards, BMC Evol. Bio., 2010]
- **NJst**, STAR, STELLS, ... BUCKy-population



# New methods

- Statistical Binning — improves estimation of gene tree distribution (input to summary methods)  
[Mirarab et al., Science, 2014] [Bayzid et al., PLoS ONE, 2015]
- ASTRAL — a summary method with improved accuracy and scalability compared to other summary methods  
[Mirarab et al., Bioinformatics, 2014 and 2015]

# Recent large-scale phylogenomic projects

## ■ Avian phylogenomics [Jarvis, Mirarab, et al., Science, 2014]

### 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 Rekpalli,<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> Qiuemei 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 Jönsson,<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>



## ■ 1K Plants (1KP) [Wickett, Mirarab, et al., PNAS, 2014]

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

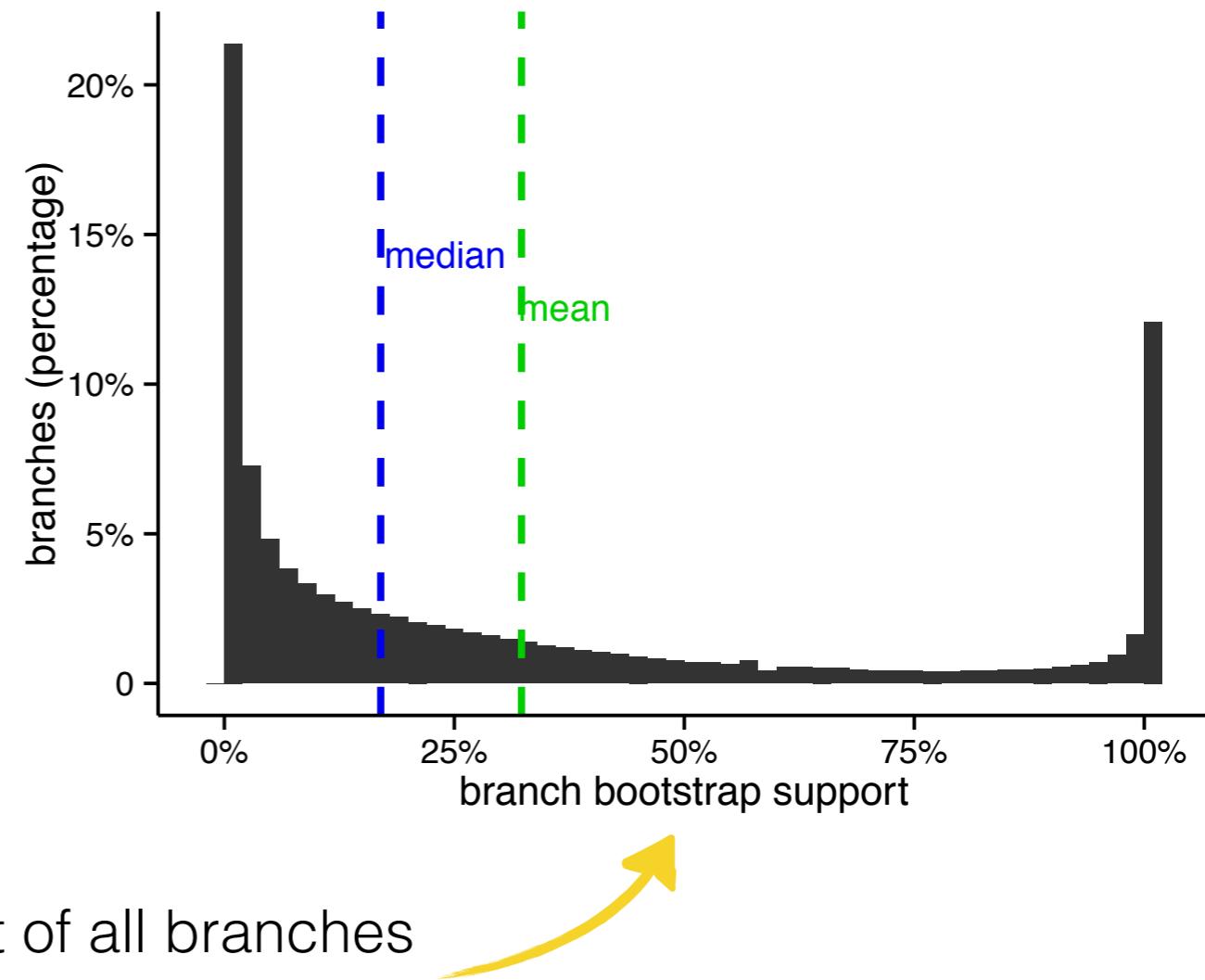
Norman J. Wickett<sup>a,b,1,2</sup>, Siavash Mirarab<sup>c,1</sup>, Nam Nguyen<sup>c</sup>, Tandy Warnow<sup>c</sup>, Eric Carpenter<sup>d</sup>, Naim Matasci<sup>e,f</sup>, Saravaraj Ayyampalayam<sup>g</sup>, Michael S. Barker<sup>f</sup>, J. Gordon Burleigh<sup>h</sup>, Matthew A. Gitzendanner<sup>h,i</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>t</sup>, Barbara Surek<sup>o</sup>, Juan Carlos Villarreal<sup>t</sup>, Béatrice Roure<sup>u,v</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. Kutchari<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</sup>



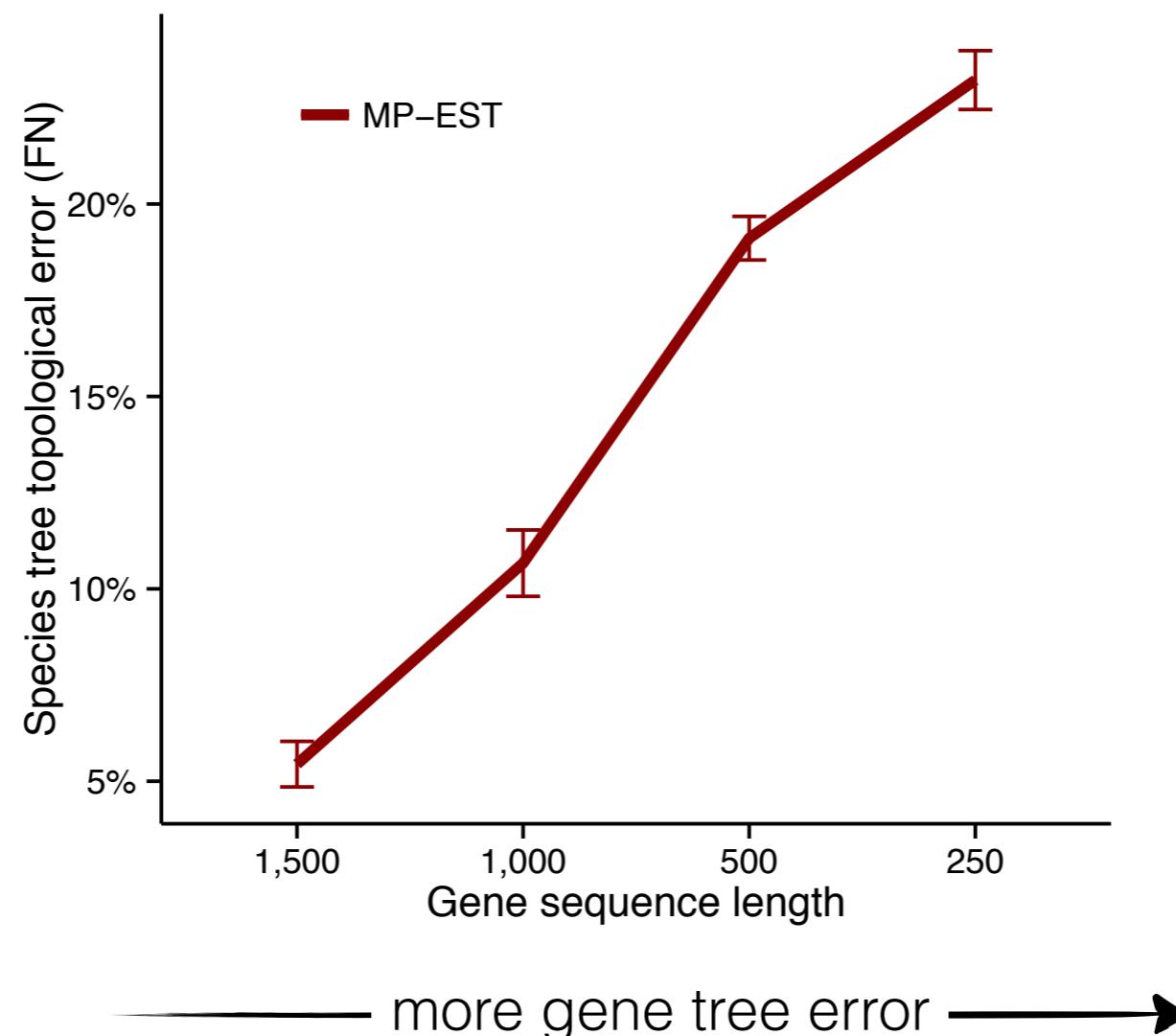
# Challenges of avian phylogenomics

[Jarvis, Mirarab, et al., Science, 2014]

- Whole genomes for 48 bird species (~100m years of evolution)
- Goal: a phylogeny of bird orders
- Extremely challenging due to rampant gene tree incongruence
- 14,000 “noisy” genes (i.e., high gene tree error)

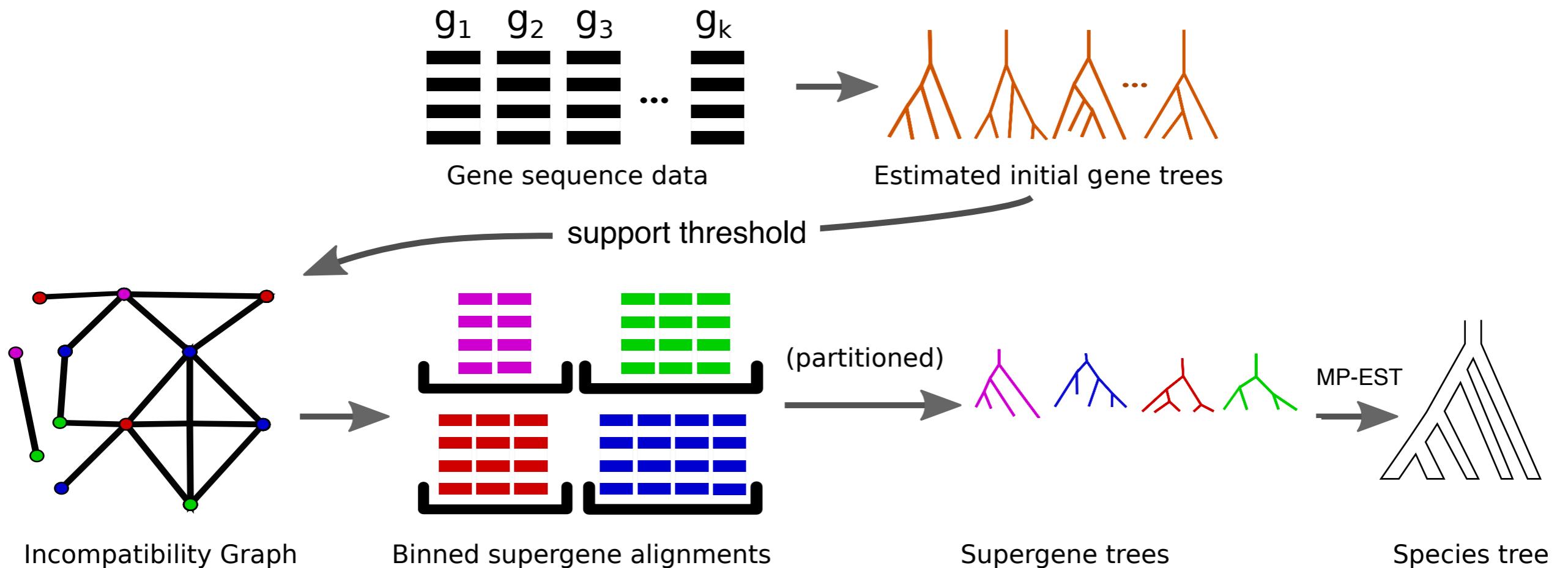


# Gene tree estimation error impacts species tree estimation



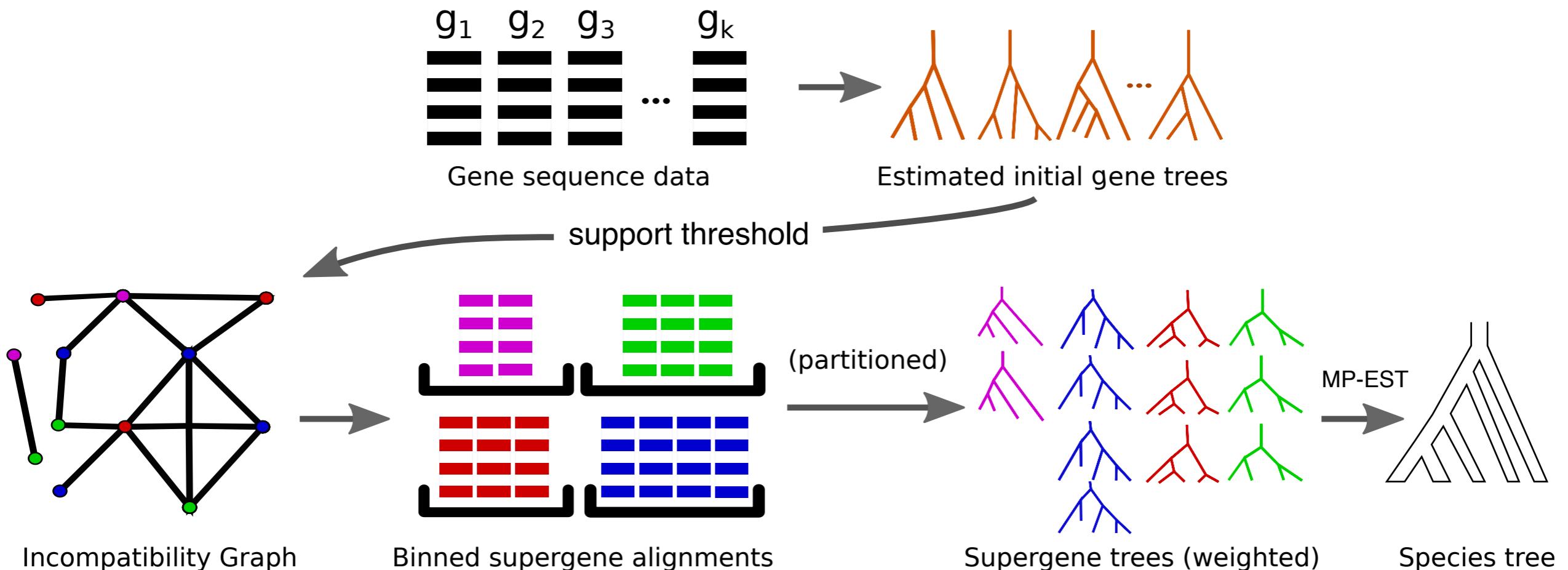
45 avian-like species, 1000 genes  
[Mirarab, et al., Science, 2014]

# Statistical binning: overview



**Original version:** unweighted [Mirarab, et al., Science, 2014]

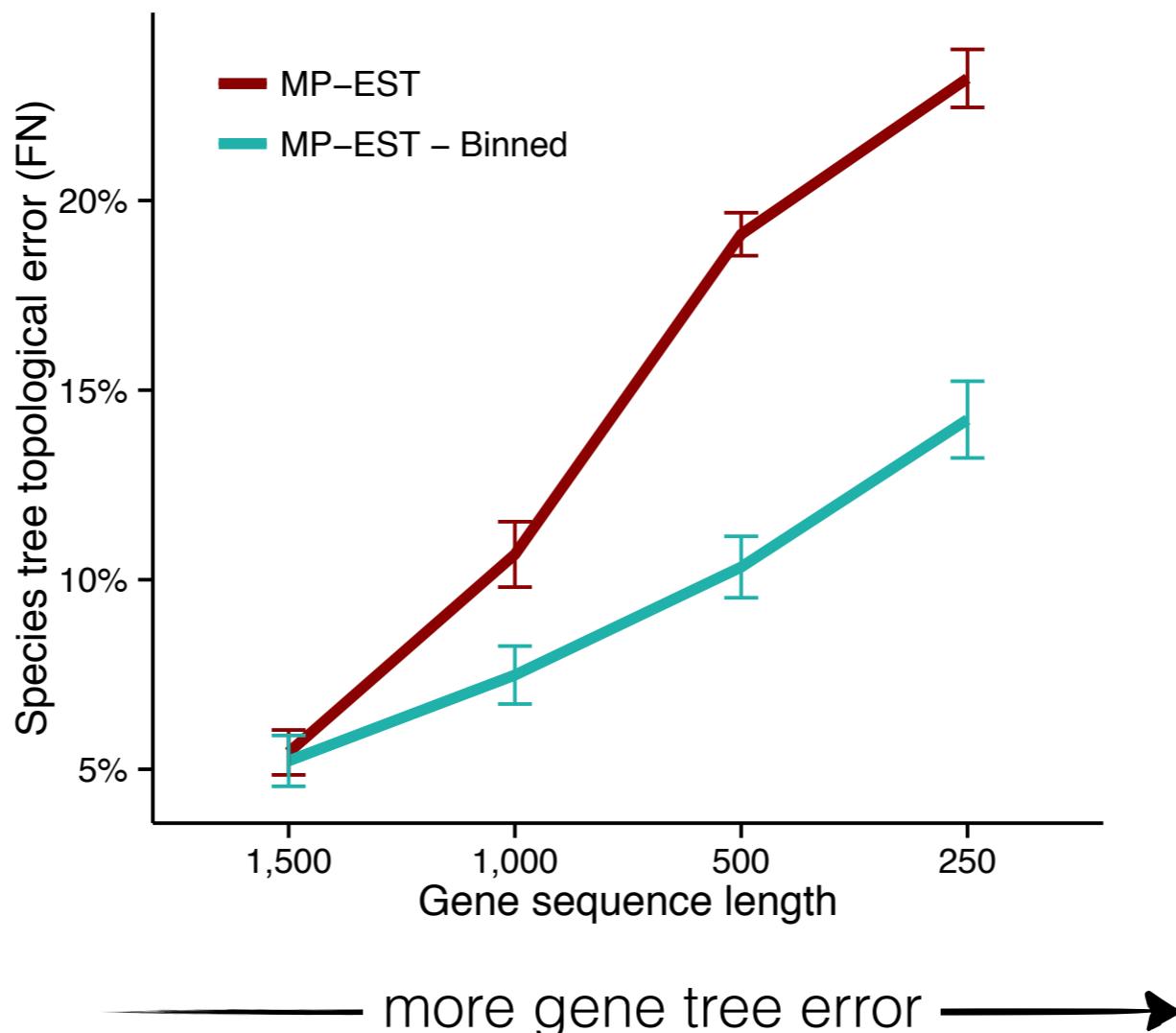
# Statistical binning: overview



**Original version:** unweighted [Mirarab, et al., Science, 2014]

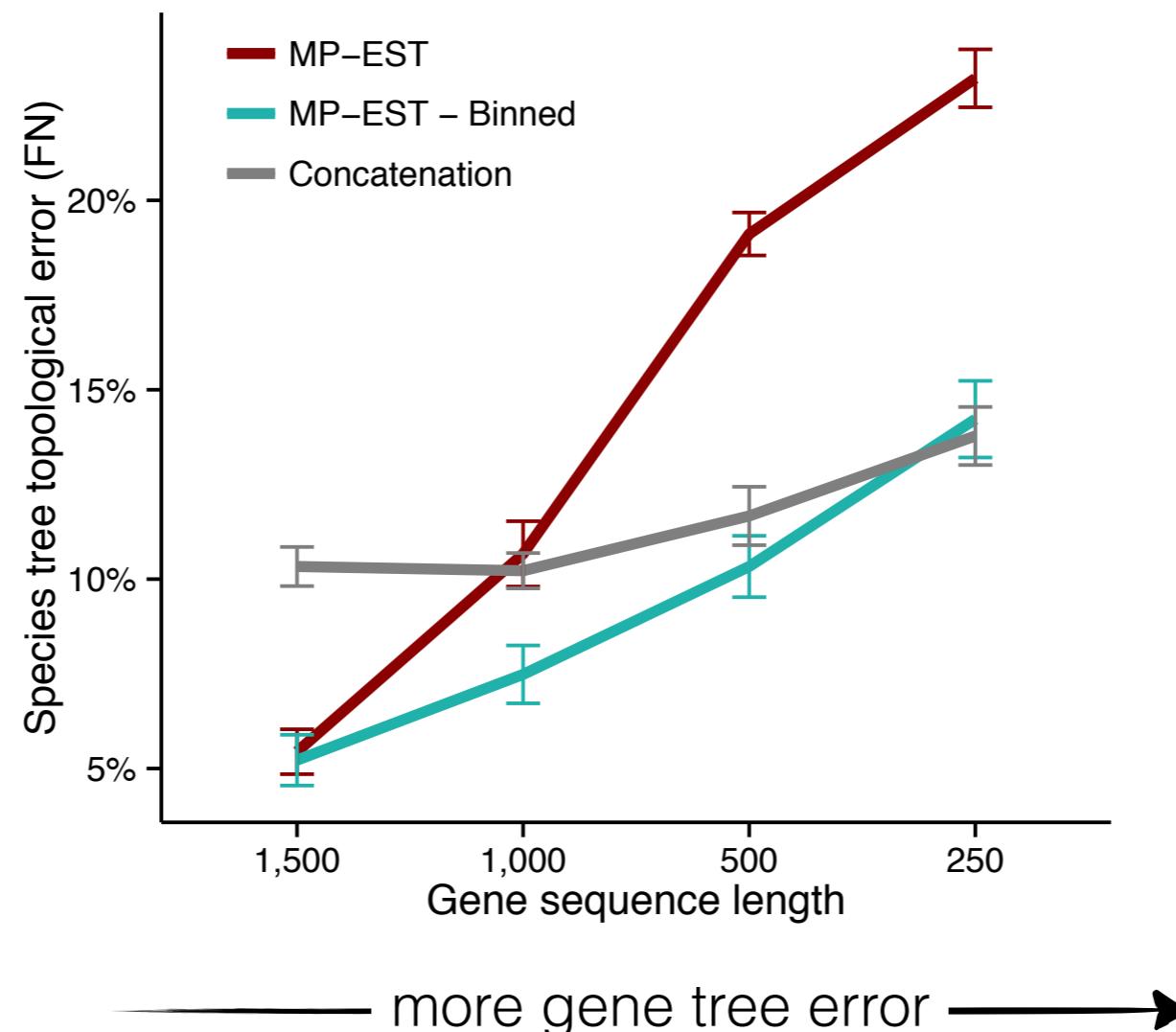
**New version:** weighted; statistically consistent [Bayzid, Mirarab, Warnow, PloS ONE, 2015]

# Statistical binning improves species tree estimation



45 avian-like species, 1000 genes  
[Mirarab, et al., Science, 2014]

# Statistical binning improves species tree estimation



45 avian-like species, 1000 genes  
[Mirarab, et al., Science, 2014]

# Binning also improves other measures of accuracy

- More accurate gene tree distributions

# Binning also improves other measures of accuracy

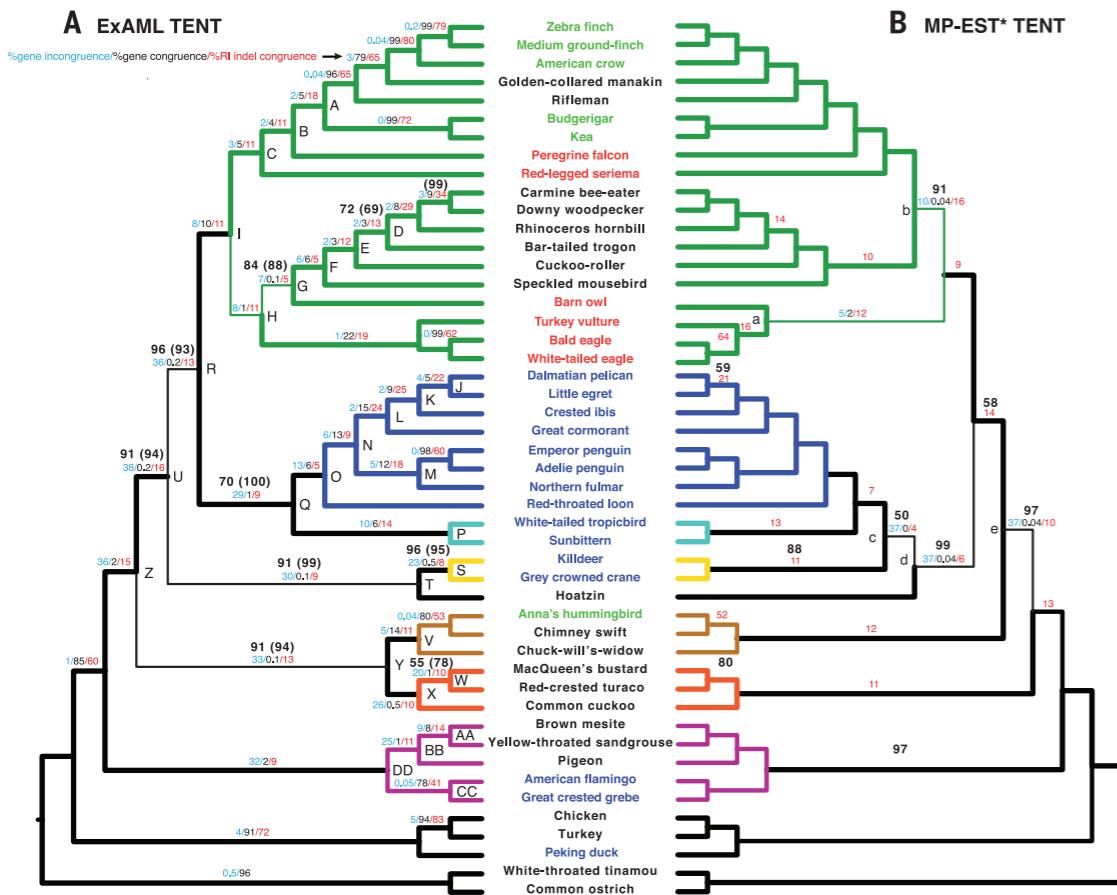
- More accurate gene tree distributions
- Better Species tree bootstrap support
  - Higher support for true positives
  - Fewer false positives with high support

# Binning also improves other measures of accuracy

- More accurate gene tree distributions
- Better Species tree bootstrap support
  - Higher support for true positives
  - Fewer false positives with high support
- More accurate species tree branch lengths
  - Hence, estimates of ILS levels

# Binning on the Avian dataset

Binned analyses results in a tree that was largely congruent with the concatenation



**B MP-EST\* TENT**

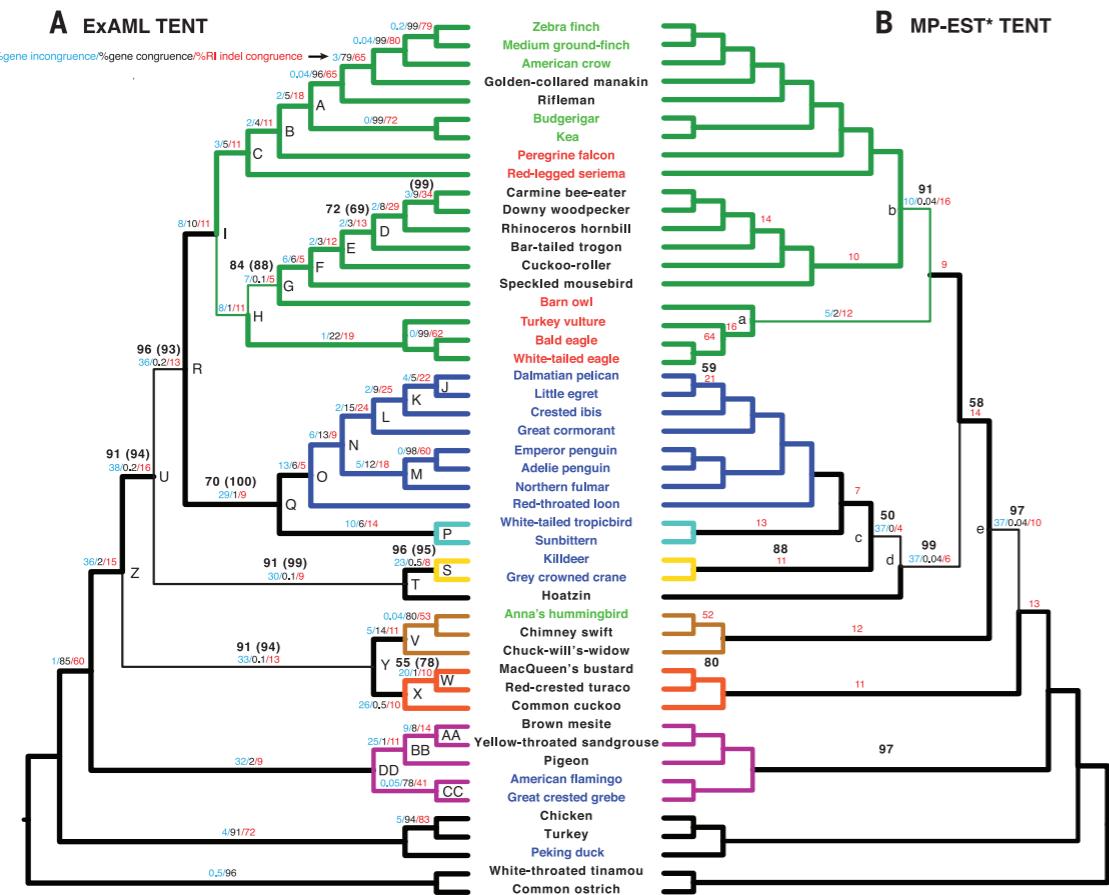


[Jarvis, Mirarab, et al., Science, 2014]

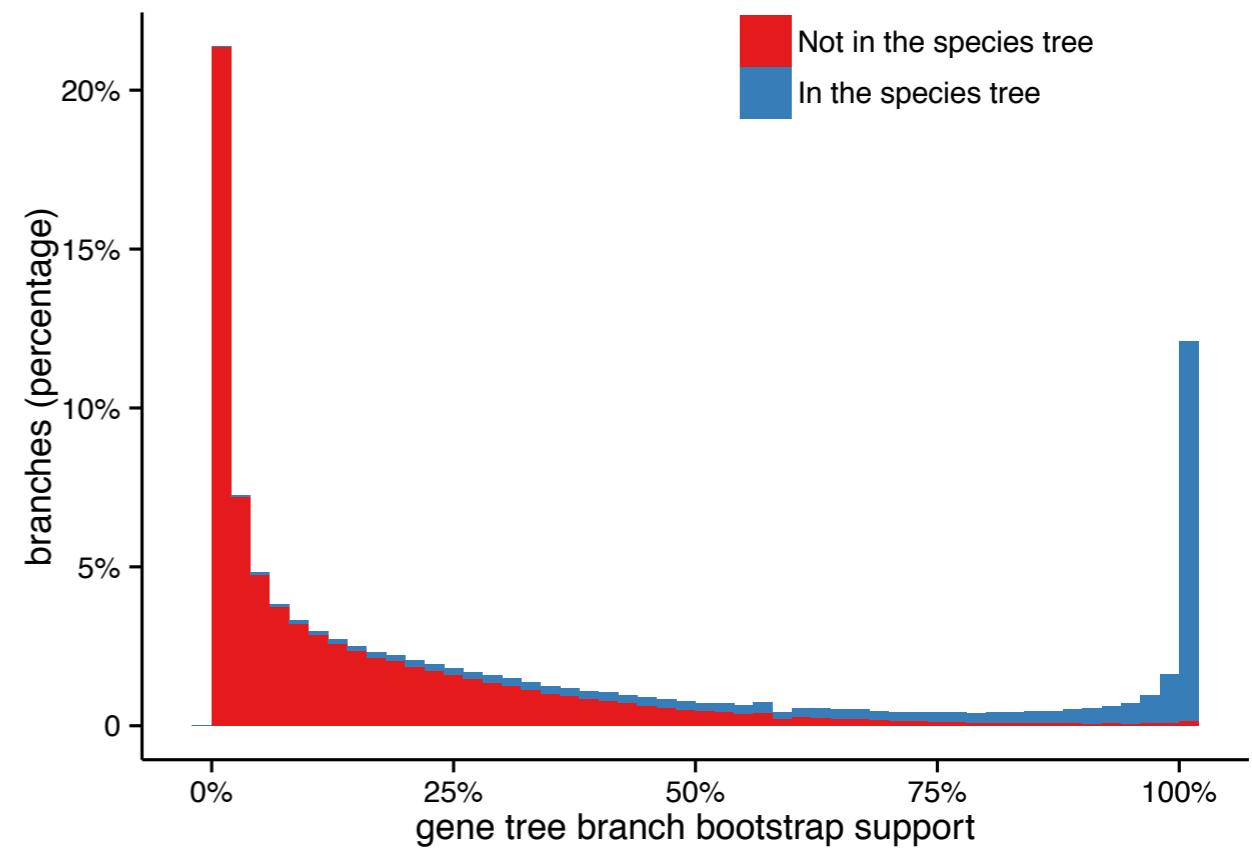
# Binning on the Avian dataset

Binned analyses results in a tree that was largely congruent with the concatenation

The binned analyses recovered almost all highly supported gene tree branches



[Jarvis, Mirarab, et al., Science, 2014]



# New methods

- Statistical Binning — improves estimation of gene tree distribution (input to summary methods)  
[Mirarab et al., Science, 2014] [Bayzid et al., PLoS ONE, 2015]

# 1KP: Plant whole transcriptomes



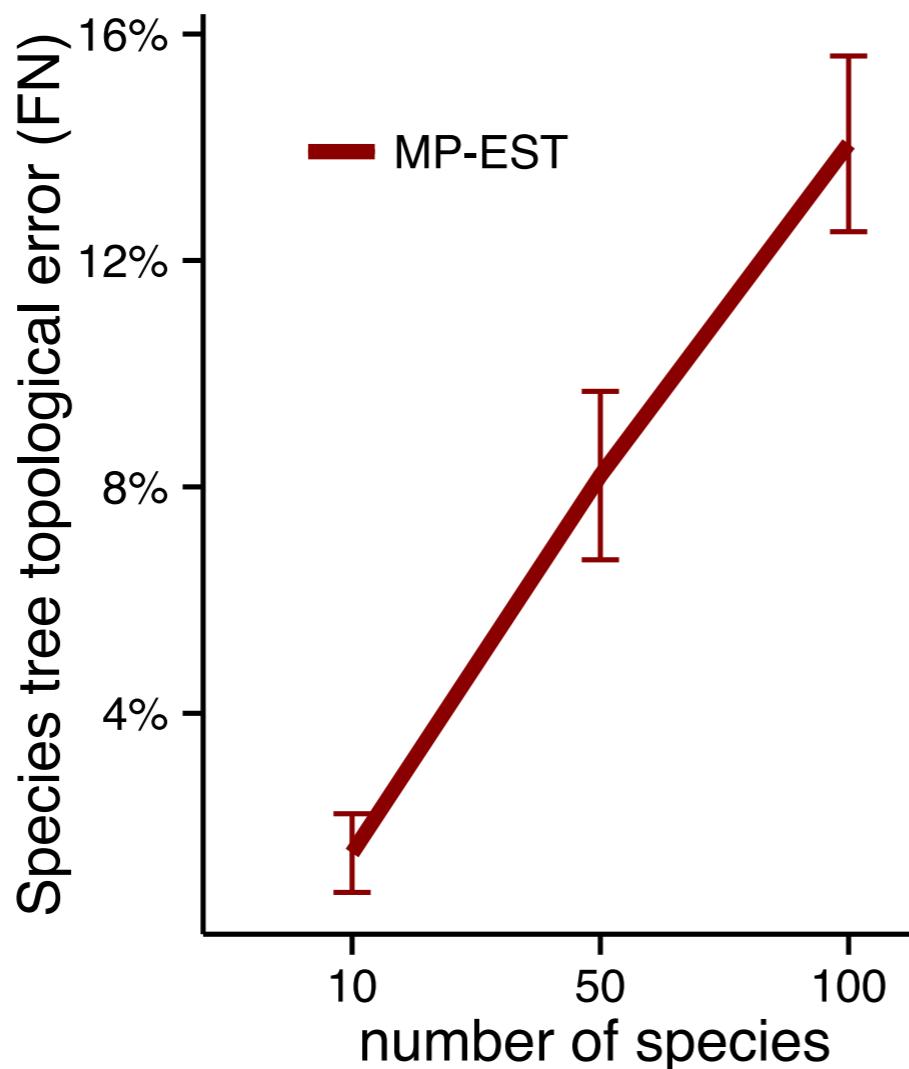
## Phylogenomic analysis of the origin and early diversification of land plants

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>t</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>v</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>



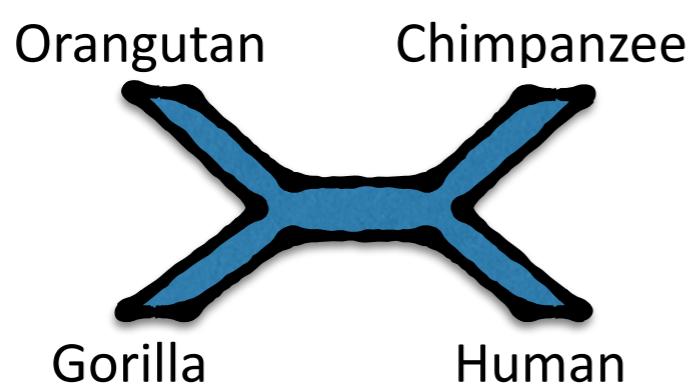
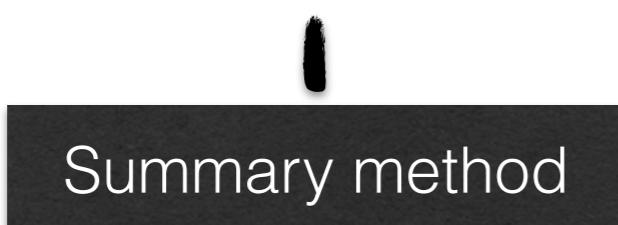
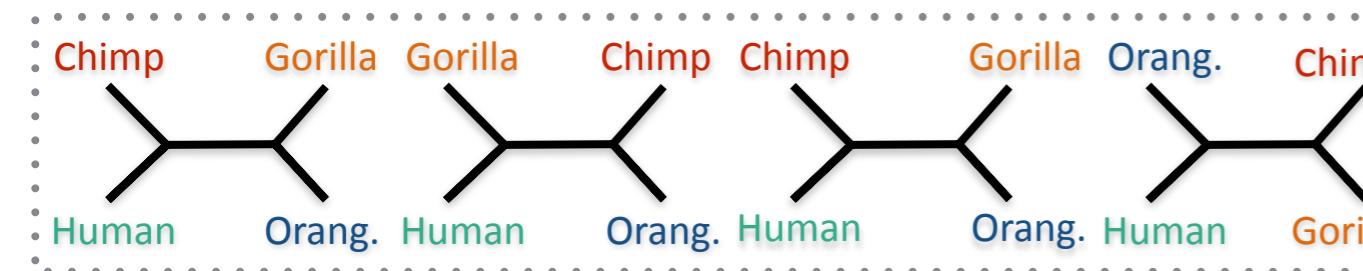
- Whole transcriptomes for 103 plant species
  - 1,200 in the next phase
  - 400-800 single copy “genes”
  - Spans ~1 billion years of evolution
  - Many unanswered questions about plant evolution

# Number of species impacts estimation error in the species tree

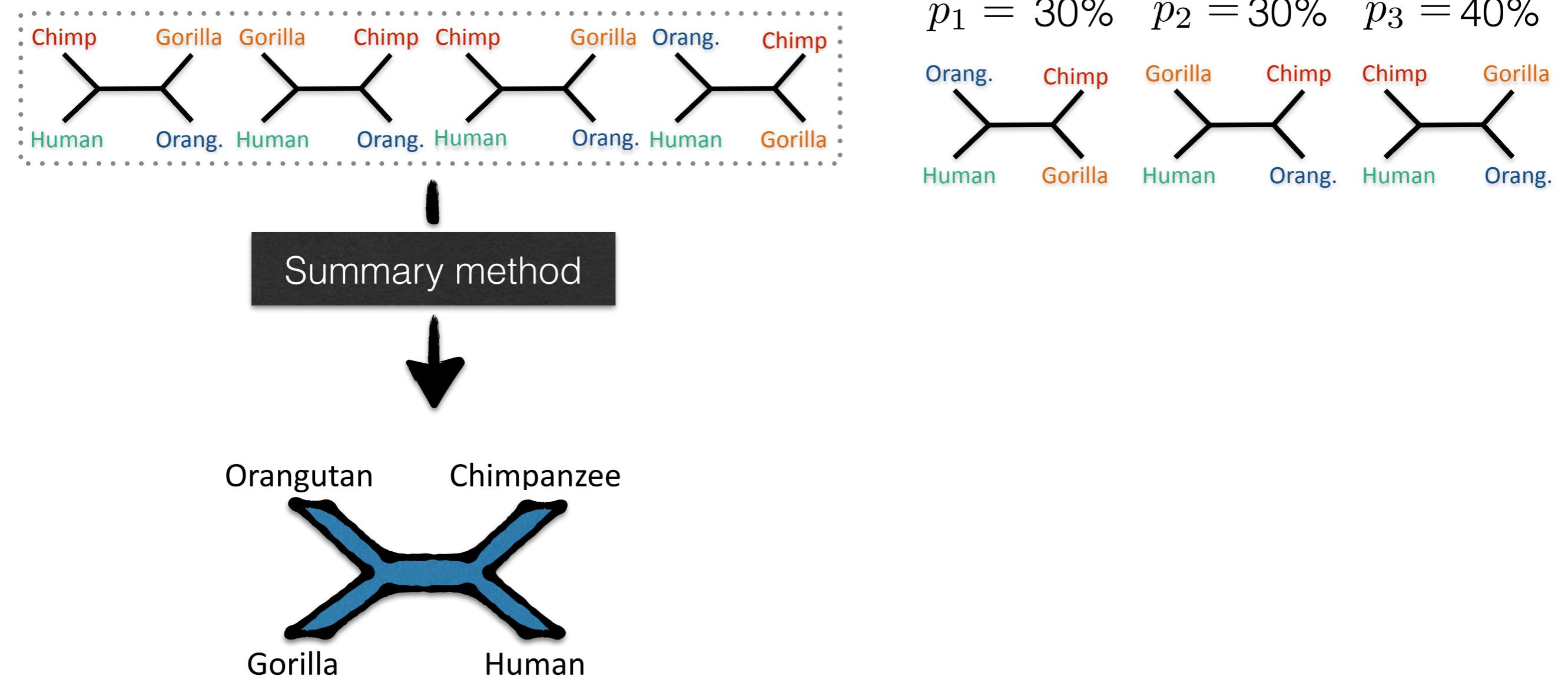


1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

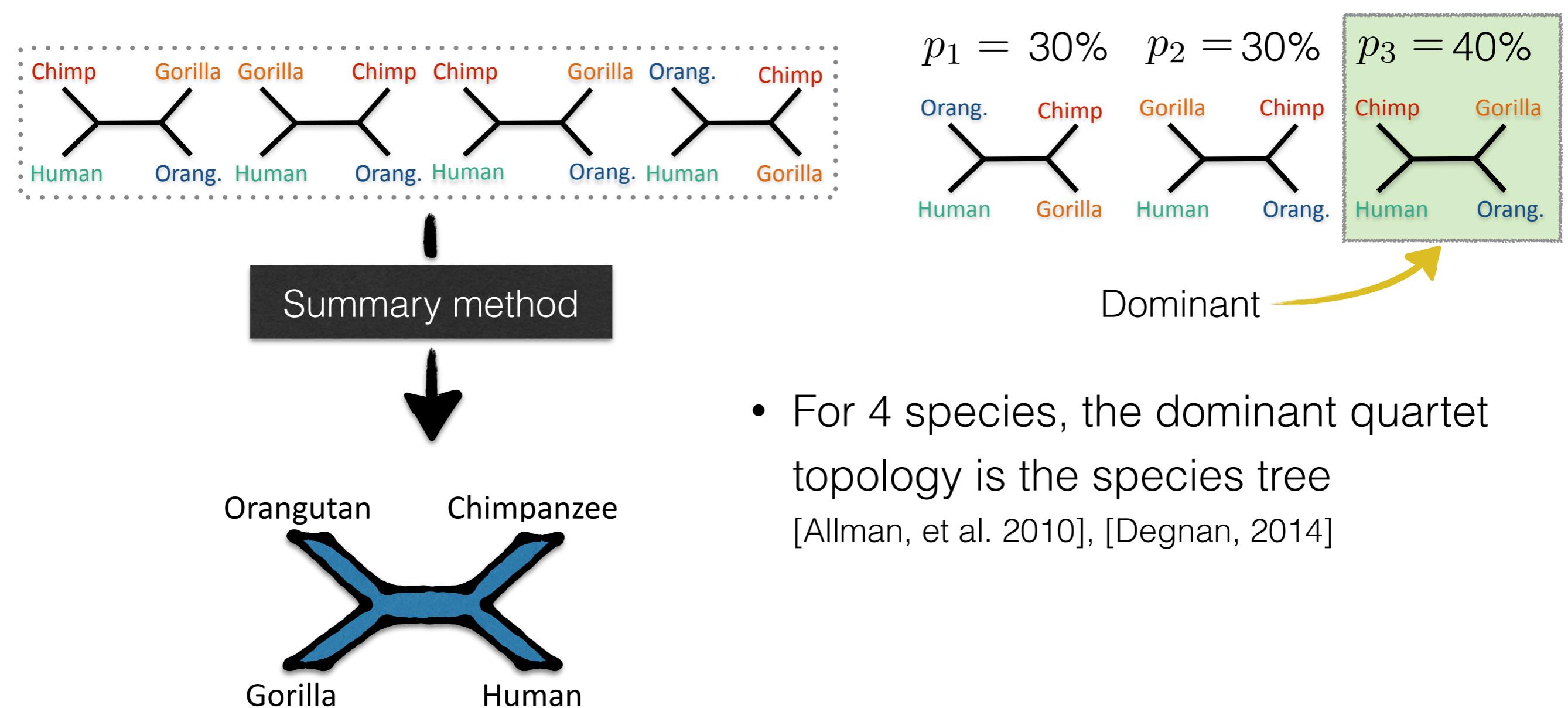
# Designing summary methods for quartets of taxa



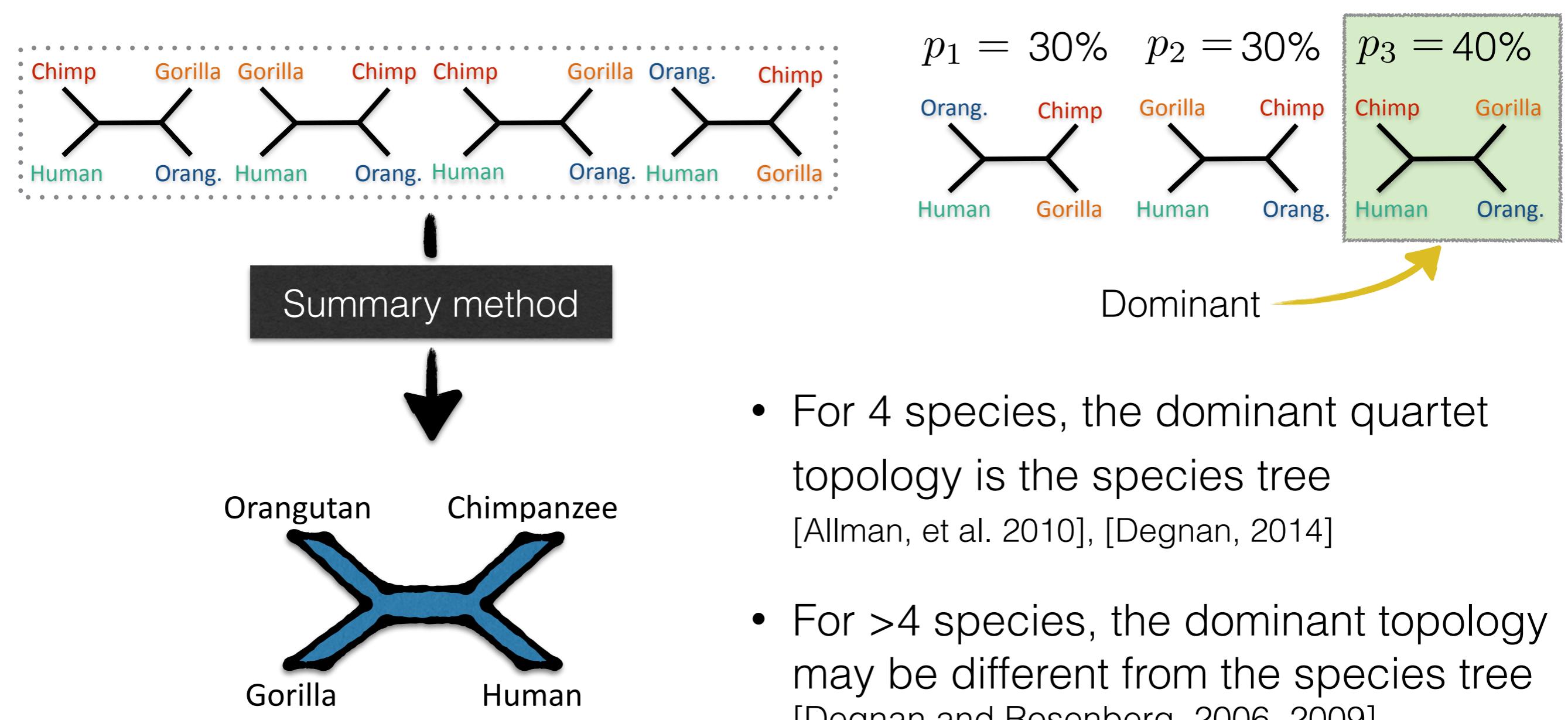
# Designing summary methods for quartets of taxa



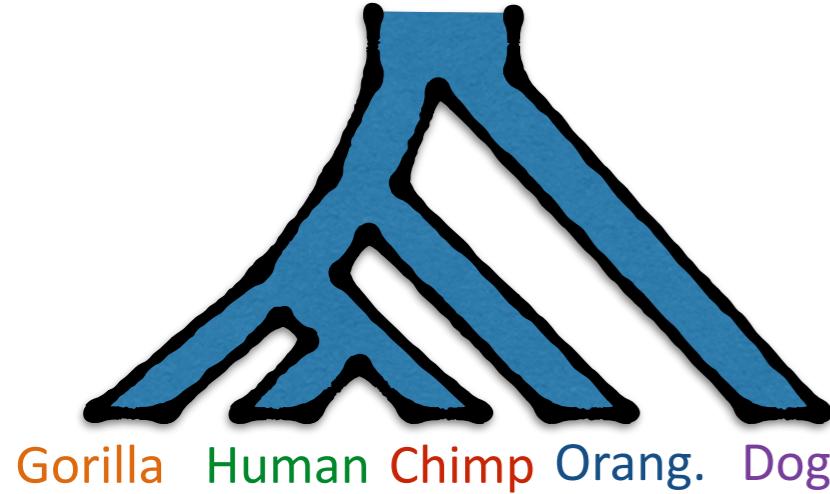
# Designing summary methods for quartets of taxa



# Designing summary methods for quartets of taxa

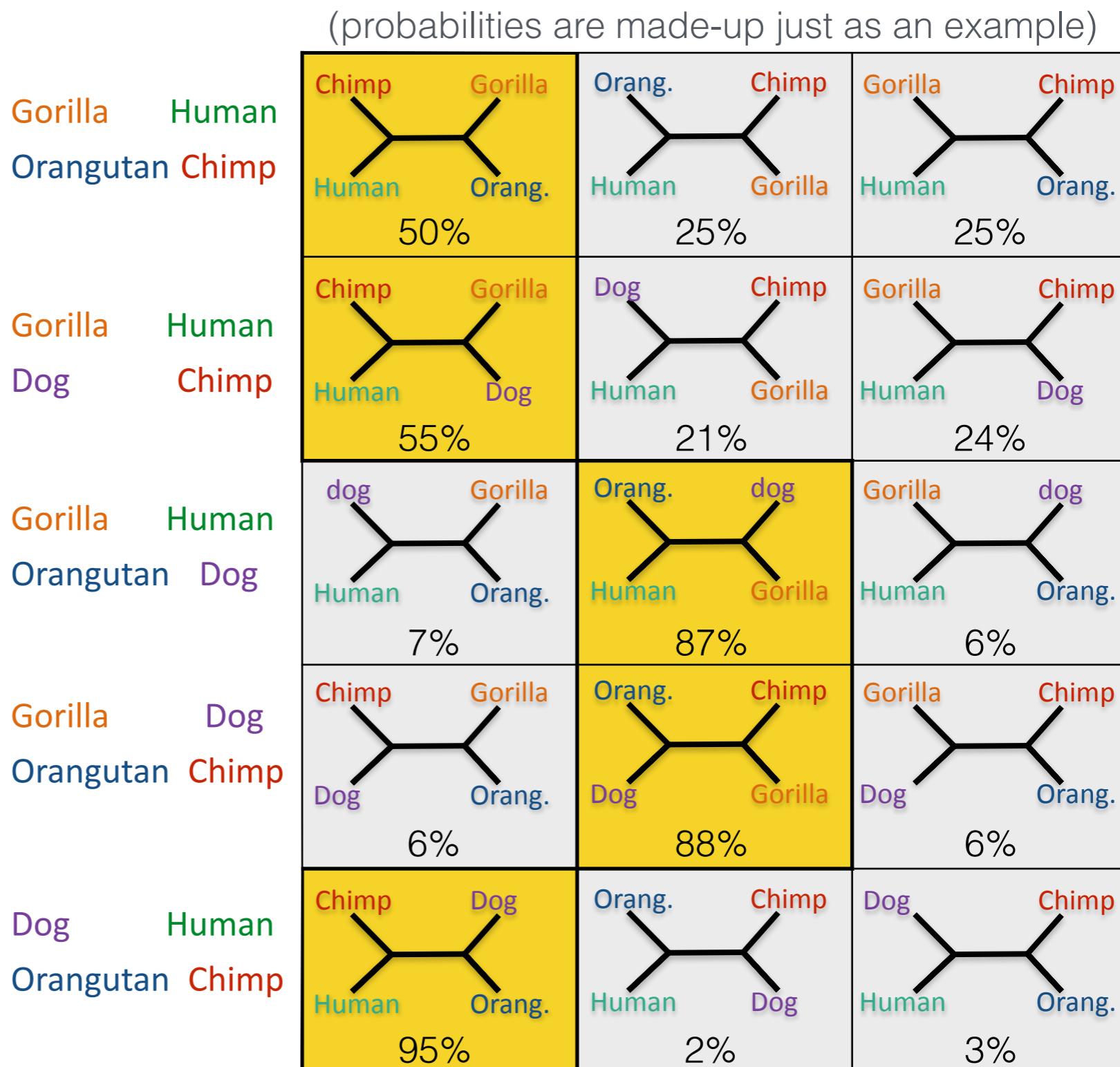


# More than 4 species

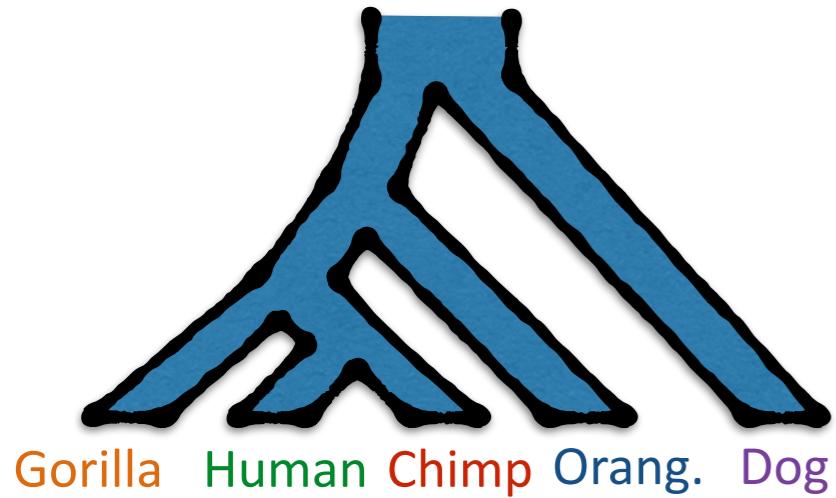


1. Break gene trees into  $\binom{n}{4}$  quartets of species
2. Find the dominant tree for all quartets of taxa
3. Combine quartet trees

**Example:** BUCKy-pop.  
[Larget, et al., Bioinformatics, 2010]



# ASTRAL: weighting by frequency



1. Break gene trees into  $\binom{n}{4}$  quartets of species
2. Find the tree that “satisfies” the maximum number of weighted quartets from gene trees

(probabilities are made-up just as an example)

Gorilla Orangutan	Human Chimp	 50%	 25%	 25%
Gorilla Dog	Human Chimp	 55%	 19%	 26%
Gorilla Orangutan	Human Dog	 7%	 87%	 6%
Gorilla Orangutan	Dog Chimp	 6%	 88%	 6%
Dog Orangutan	Human Chimp	 95%	 2%	 3%

# Maximum Quartet Support Species Tree

[Mirarab, et al., ECCB, 2014]

- Optimization Problem (suspected NP-Hard):

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

$$Score(T) = \sum_{t \in \mathcal{T}} |Q(T) \cap Q(t)|$$

Set of quartet trees  
induced by T

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

# ASTRAL

- We developed a dynamic programming algorithm to solve the problem exactly
  - Exponential running time (still feasible for <18 species)
- Developed a constrained version of the problem that can be solved exactly in polynomial time
  - Constraints search to branches that appear in gene trees ( $\mathcal{X}$ )
  - Runs for 1000 species and 1000 genes in about a day
  - Remains statistically consistent

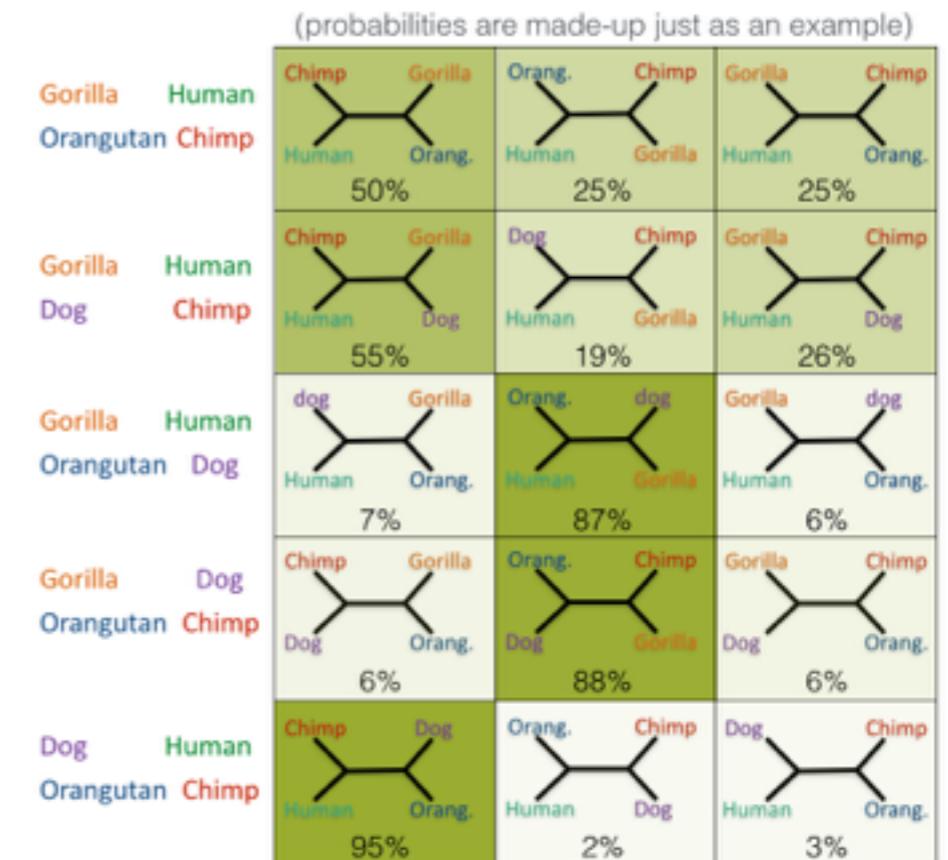
# Asymptotic running time

- $O(nk|\mathcal{X}|^2)$  for  $k$  genes of  $n$  species

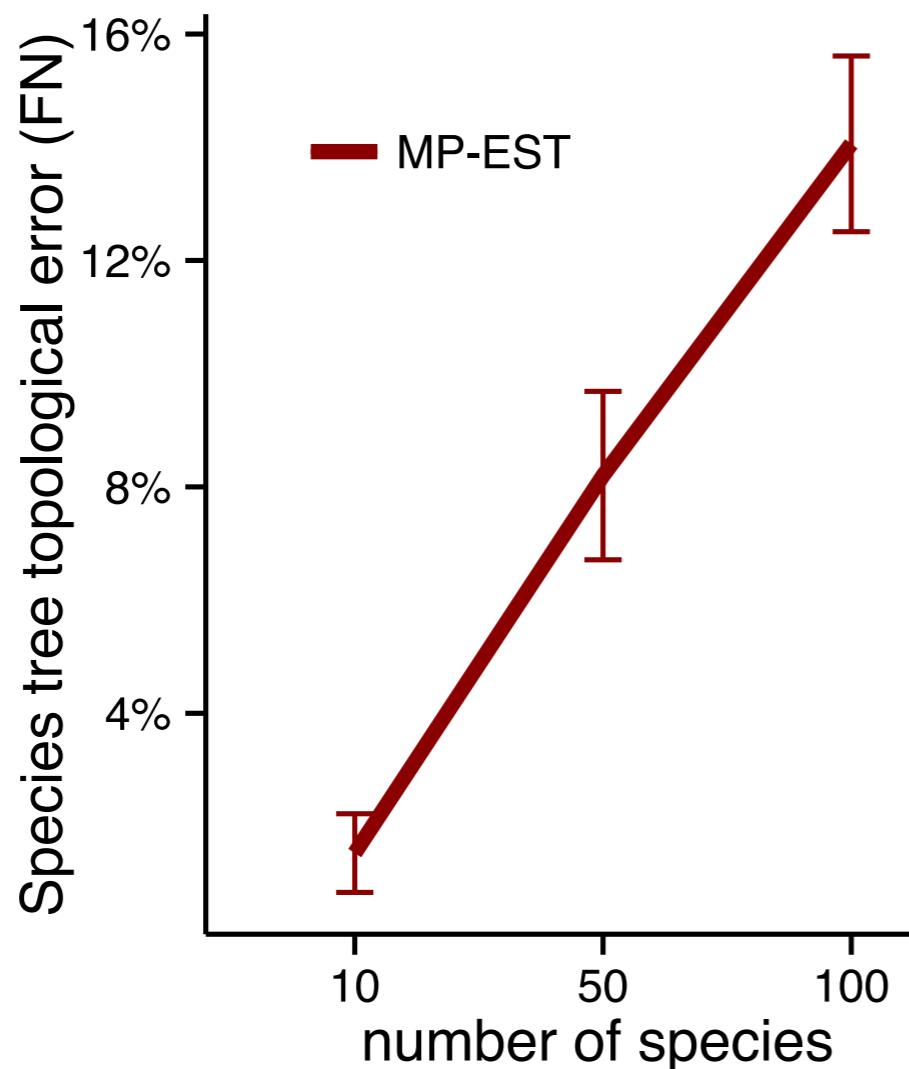
# Asymptotic running time

- $O(nk|\mathcal{X}|^2)$  for  $k$  genes of  $n$  species

- Surprise: running time is better than  $\Theta(n^4)$
- Don't we have to at least list all  $\binom{n}{4}$  quartets?
- No! we calculate scores without listing  $\binom{n}{4}$  quartets

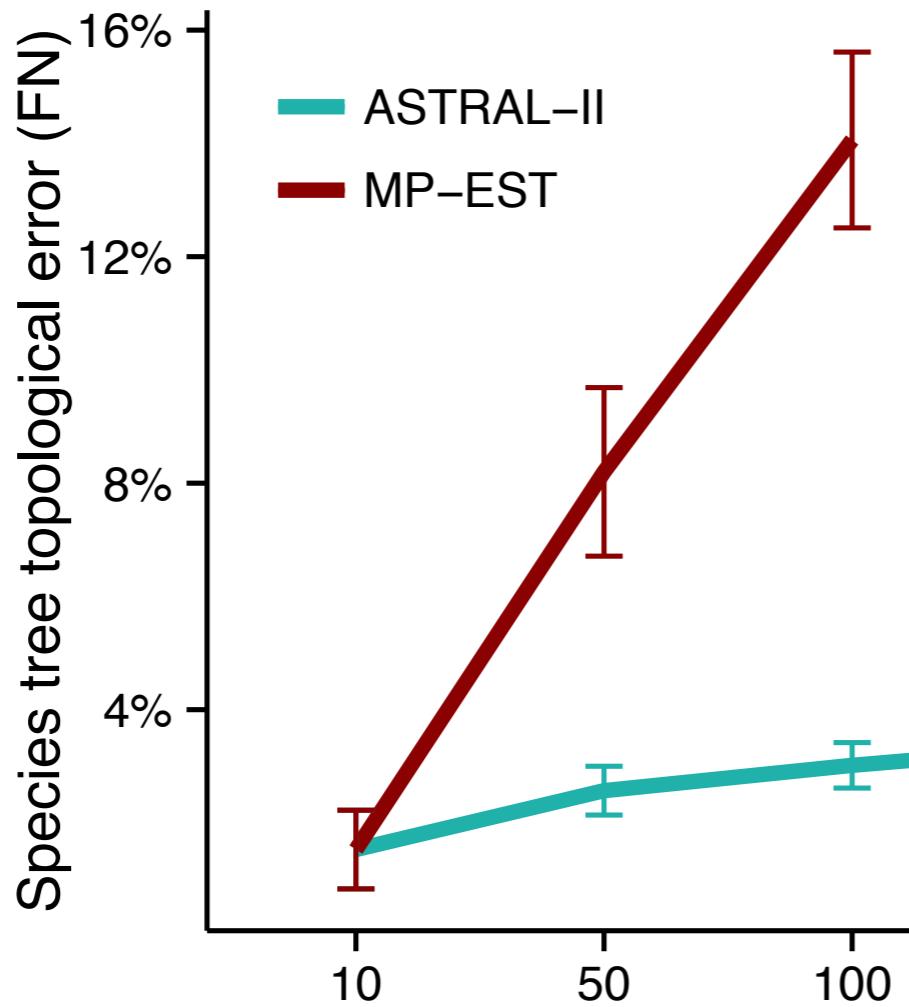


# Number of species impacts estimation error in the species tree



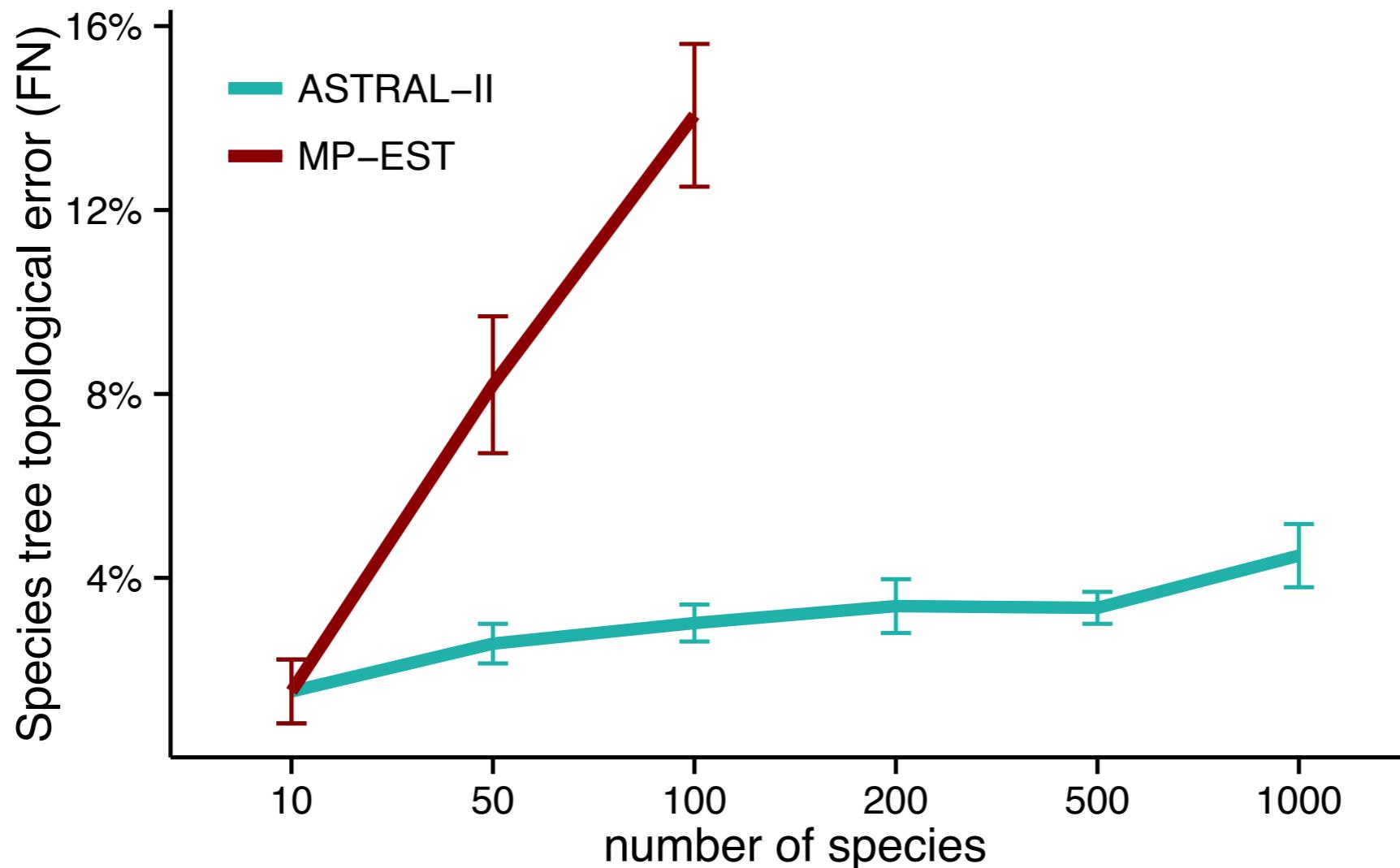
1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# ASTRAL: accurate and scalable



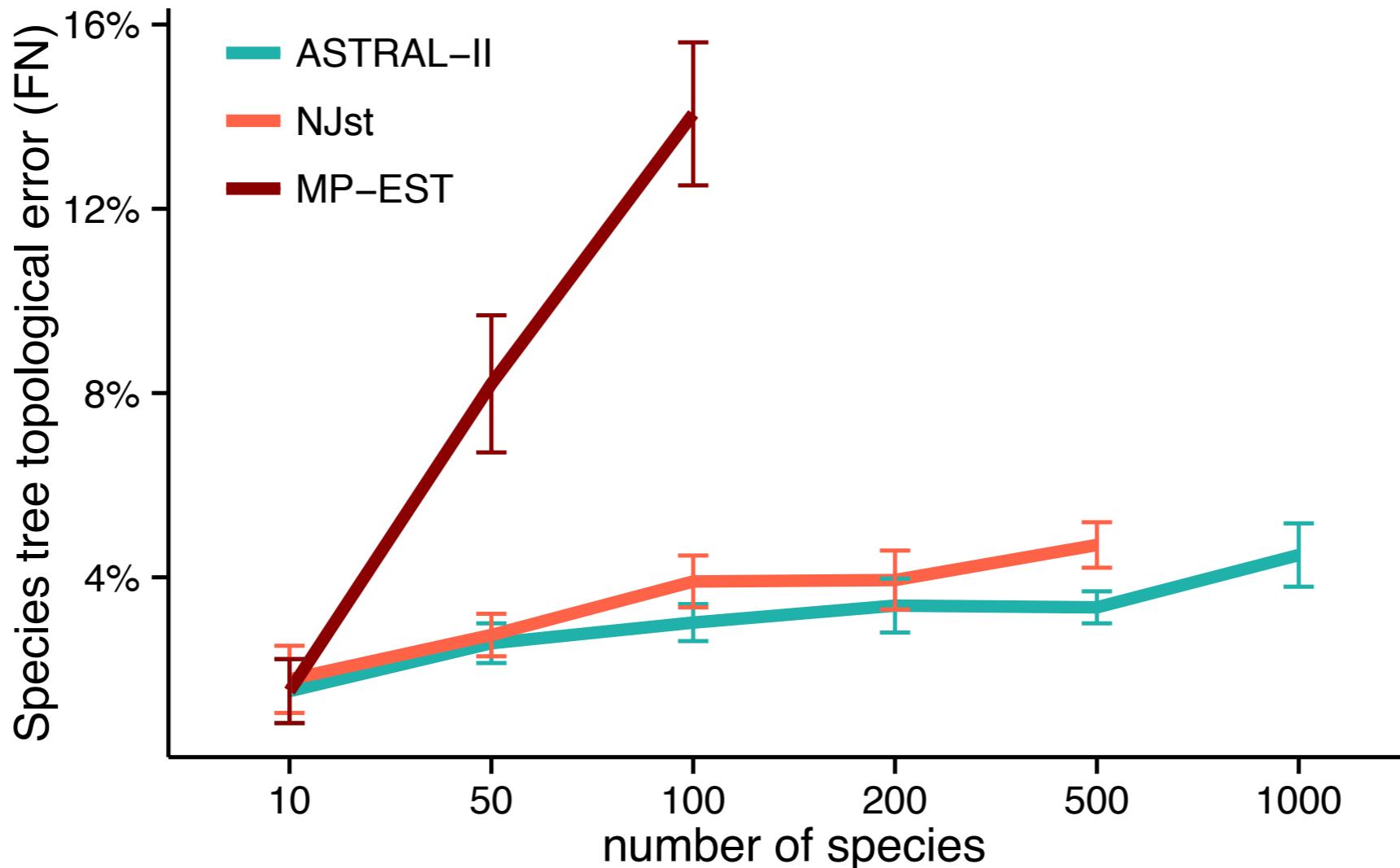
1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# ASTRAL: accurate and scalable



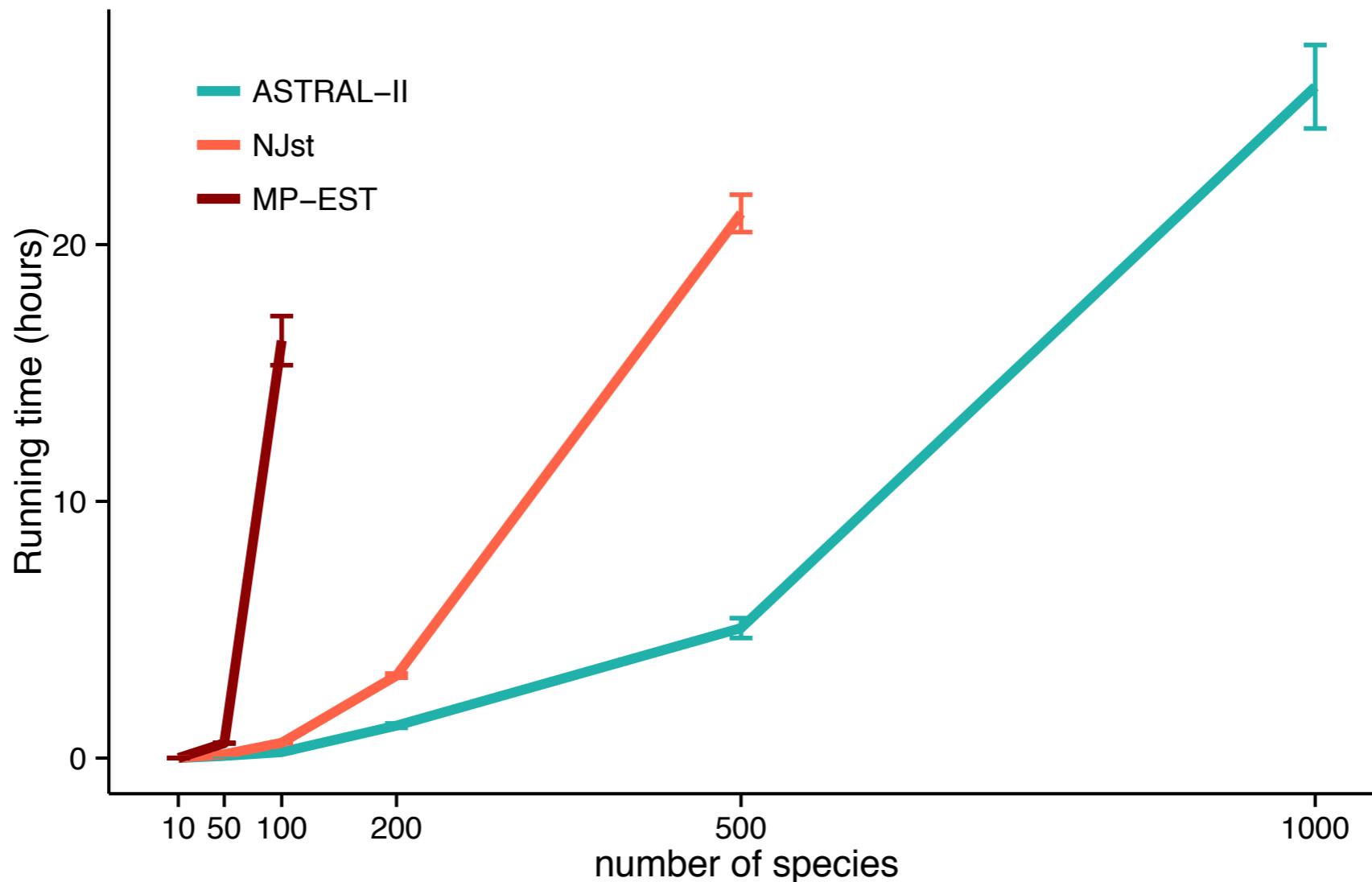
1000 genes, “medium” levels of ILS, simulated species trees  
[Mirarab and Warnow, ISMB, 2015]

# ASTRAL: accurate and scalable



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

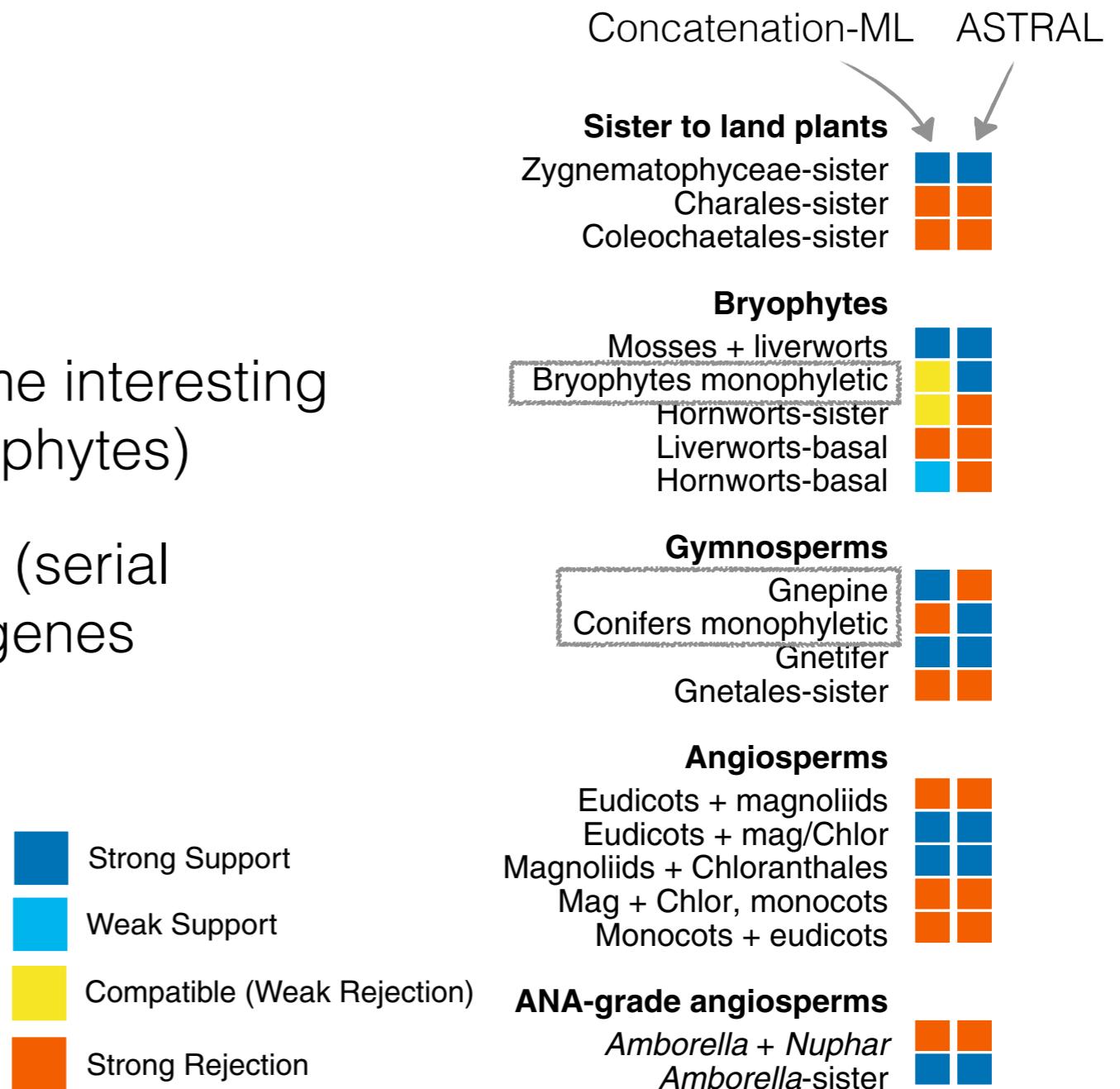
# Running time as function of # species



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

# ASTRAL on plants dataset

- The ASTRAL tree:
  - High support
  - Similar to concatenation with some interesting differences (e.g., recovered bryophytes)
- ASTRAL took only about 10 minutes (serial running time) on 103 taxa and 400 genes



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

# ASTRAL on biological datasets

- 1KP: **103** plant species, 400-800 genes
- Yang, et al. **96** Caryophyllales species, 1122 genes
- Dentinger, et al. **39** mushroom species, 208 genes
- Giarla and Esselstyn. **19** Philippine shrew species, 1112 genes
- Laumer, et al. **40** flatworm species, 516 genes
- Grover, et al. **8** cotton species, 52 genes
- Hosner, Braun, and Kimball. **28** quail species, 11 genes
- Simmons and Gatesy. **47** angiosperm species, 310 genes
- Prum et al, **198** avian species, 259 genes

## Dissecting Molecular Evolution in the Highly Diverse Plant Clade Caryophyllales Using Transcriptome Sequencing

Syst. Biol. 0(0):1–14, 2015  
© The Author(s) 2015. Published by Oxford University Press, on behalf of the Society of Systematic Biologists. All rights reserved.  
For Permissions, please email: journals.permissions@oup.com  
DOI:10.1093/sysbio/syv029



eLIFE  
elifesciences.org

The Challenges of Resolving a Rapid, Recent Radiation: Empirical and Simulated Phylogenomics of Philippine Shrews

Nuclear genomic signals of the 'microturbellarian' roots of platyhelminth evolutionary innovation

Christopher E Laumer<sup>1\*</sup>, Andreas Hejnol<sup>2</sup>, Gonzalo Giribet<sup>1</sup>



## Re-evaluating the phylogeny of allopolyploid *Gossypium* L. <sup>☆</sup>

Corrinne E. Grover<sup>1,2\*</sup>, Joseph P. Gallagher<sup>3</sup>, Josef J. Jareczek<sup>4</sup>, Justin T. Page<sup>5</sup>, Joshua A. Udall<sup>6</sup>, Michael A. Gore<sup>4</sup>, Jonathan F. Wendt<sup>7</sup> *Journal of Biogeography* U. Biogeogr. (2015)



Land connectivity changes and global cooling shaped the colonization history and diversification of New World quail (Aves: Galliformes: Odontophoridae)

Peter A. Hosner<sup>1\*</sup>, Edward L. Braun<sup>1,2,3</sup> and Rebecca T. Kimball<sup>1,2,3</sup>

doi:10.1038/nature15697

## LETTER

## 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 Domburg<sup>1,3,4</sup>, Daniel J. Field<sup>1,5</sup>, Jeffrey P. Townsend<sup>1,6</sup>, Emily Moriarty Lemmon<sup>7</sup> & Alan R. Lemmon<sup>8</sup>

# ASTRAL-II on biological datasets (ongoing collaborations)

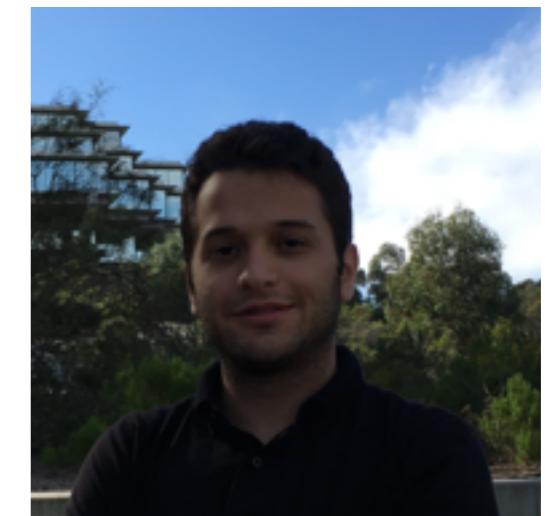
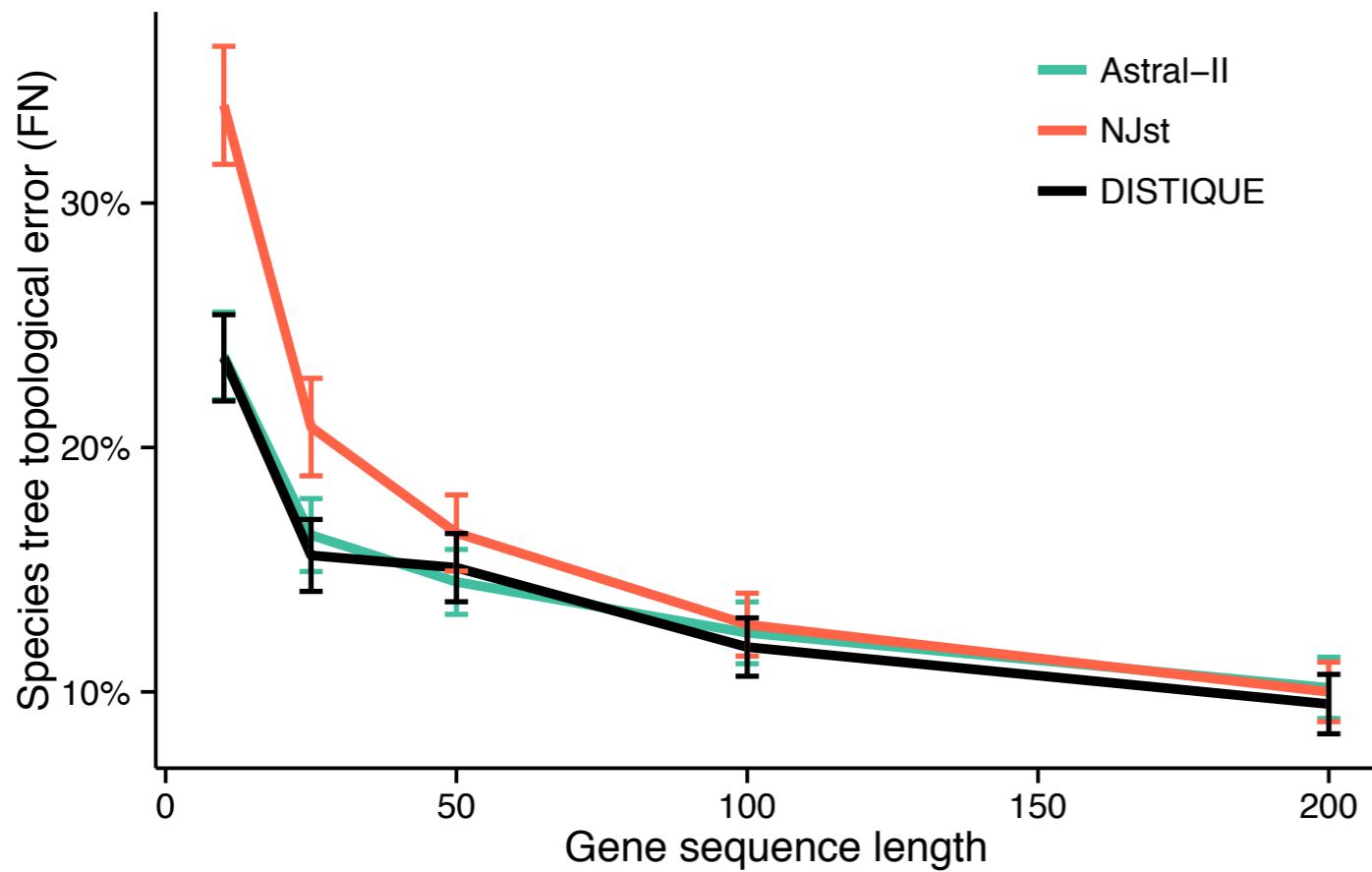
- 1200 plants with ~ 400 genes (1KP consortium)
- 250 avian species with 2000 genes (with LSU, UF, and Smithsonian)
- 200 avian species with whole genomes (with Genome 10K, international)
- 250 suboscine species (birds) with ~2000 genes (with LSU and Tulane)
- 140 Insects with 1400 genes (with U. Illinois at Urbana-Champaign)
- 50 Hummingbird species with 2000 genes (with U. Copenhagen and Smithsonian)
- 40 raptor species (birds) with 10,000 genes (with U. Copenhagen and Berkeley)
- 38 mammalian species with 10,000 genes (with U. of Bristol, Cambridge, and Nat. Univ. of Ireland)

# New methods

- Statistical Binning — improves estimation of gene tree distribution (input to summary methods)  
[Mirarab et al., Science, 2014] [Bayzid et al., PLoS ONE, 2015]
- ASTRAL — a summary method with improved accuracy and scalability compared to other summary methods  
[Mirarab et al., Bioinformatics, 2014 and 2015]

# Ongoing work

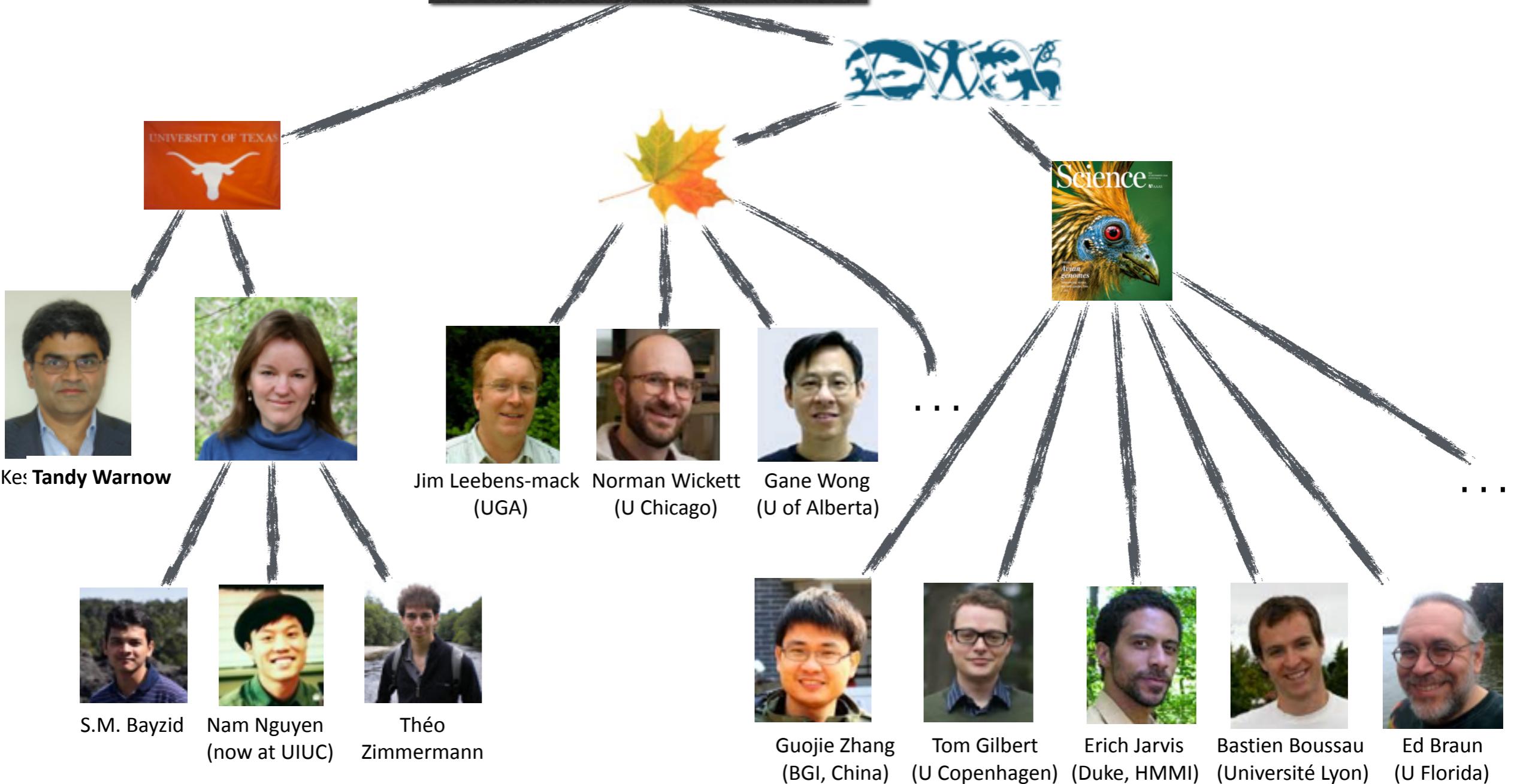
DISTIQUE — a family of summary method, some with quadratic running time and interesting theory  
[unpublished and unfinished work; with Erfan Sayyari]



# Summary

- Genome-scale sequence data provides a wealth of information
- Yet, reconstruction of species phylogenies remains challenging
  - Scalability to many species: ASTRAL
  - Limited data per gene: statistical binning
  - Impact of model violations
  - Recombination
  - Missing data
  - Multiple sources of gene tree discordance
- Many interesting statistical and computational questions and need for method development

# Acknowledgments



HMMI international student fellowship

