# Challenges and advances in genome-wide species tree reconstruction
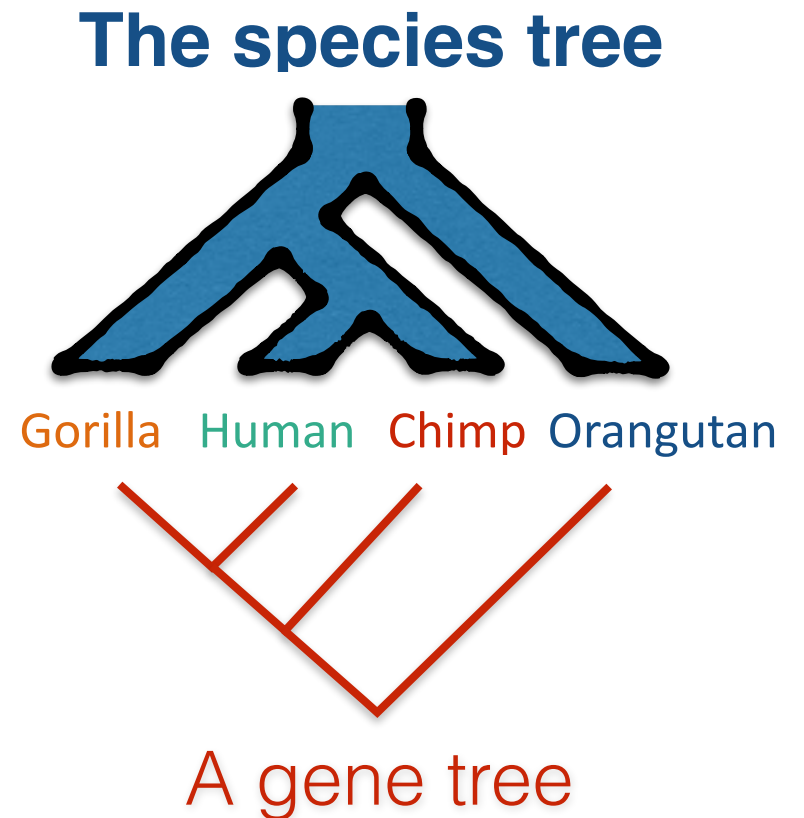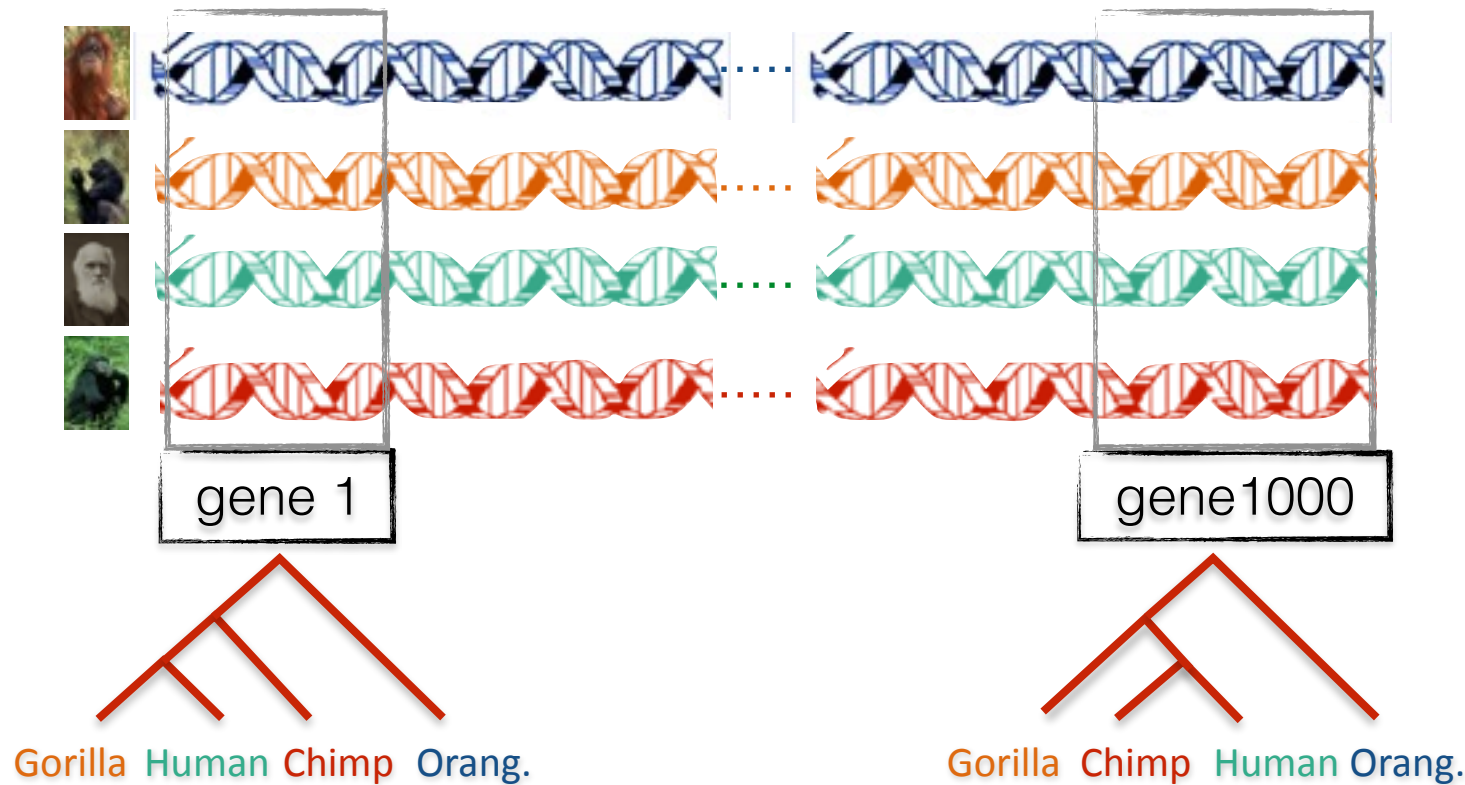
Siavash Mirarab
University of California, San Diego
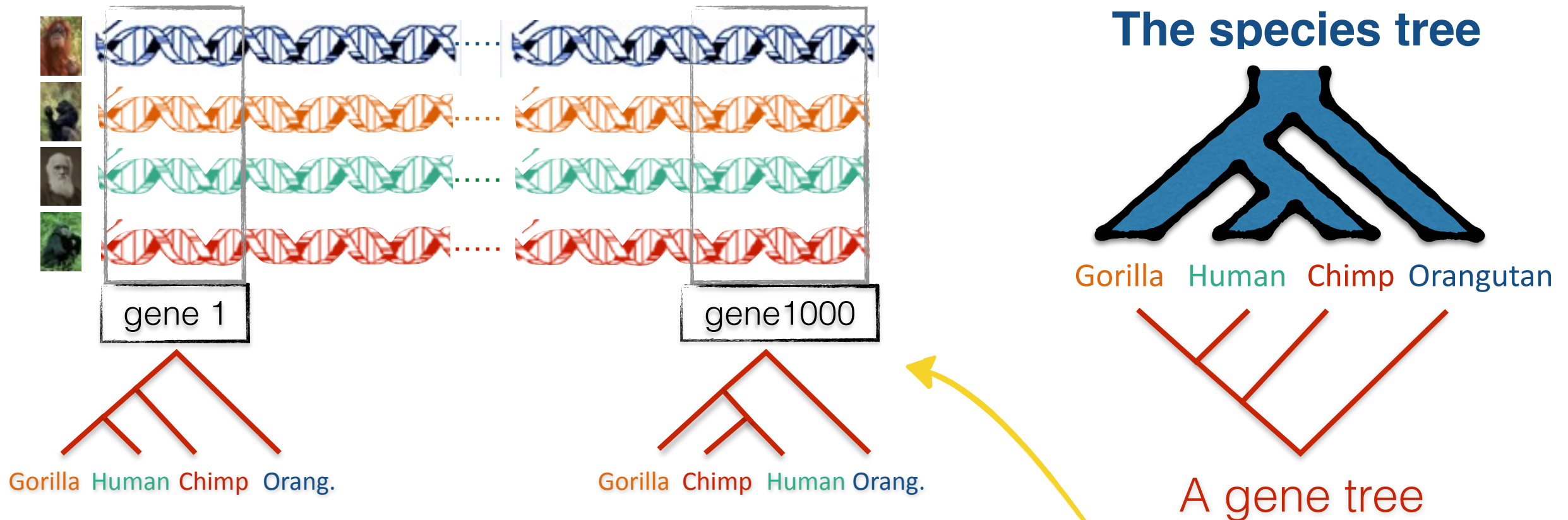
# Gene tree discordance



The species tree

Gorilla   Human   Chimp   Orangutan

A gene tree

gene 1

Gorilla  Human  Chimp  Orang.

gene1000

Gorilla  Chimp  Human  Orang.

**Causes of gene tree discordance include:**

- Incomplete Lineage Sorting (ILS)

- Duplication and loss

- Horizontal Gene Transfer (HGT)

- Hybridization

# Gene tree discordance



The species tree

Gorilla  Human  Chimp  Orangutan

A gene tree

gene 1

Gorilla  Human  Chimp  Orang.

gene1000

Gorilla  Chimp  Human  Orang.

**Causes of gene tree discordance include:**

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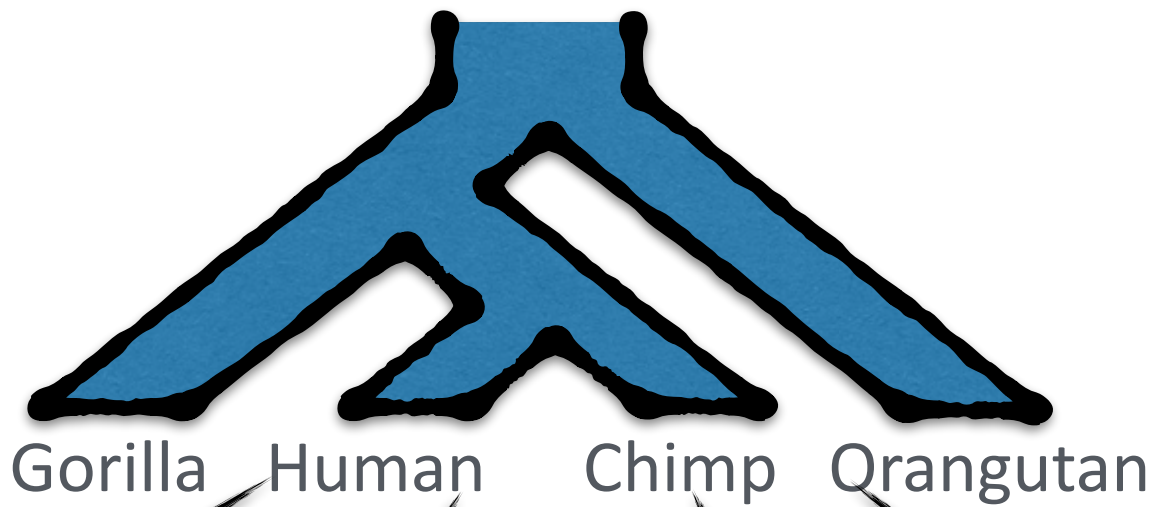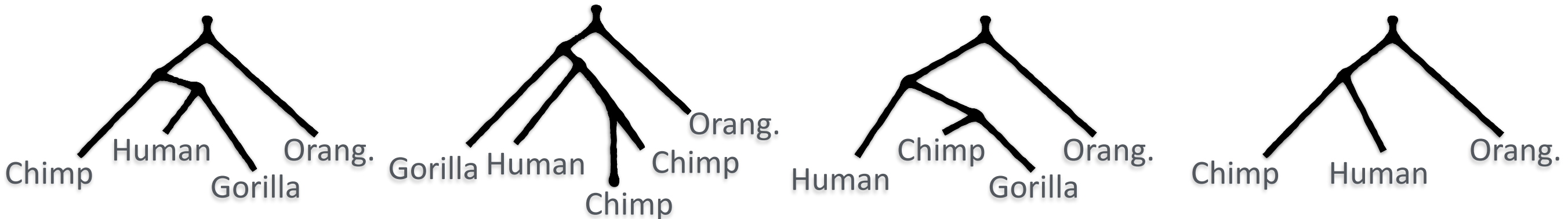- Horizontal Gene Transfer (HGT)
- Hybridization

"c-gene":
recombination-free orthologous
stretches of the genome

Gorilla Human Chimp Orangutan

**Gene evolution model**

Chimp Human Gorilla Orang.

Gorilla Human Chimp Orang. Chimp

Human Chimp Gorilla Orang.

Chimp Human Orang.

**Sequence evolution model**

```
ACTGCACACCG
ACTGC-CCCCG
AATGC-CCCCG
-CTGCACACGG
```
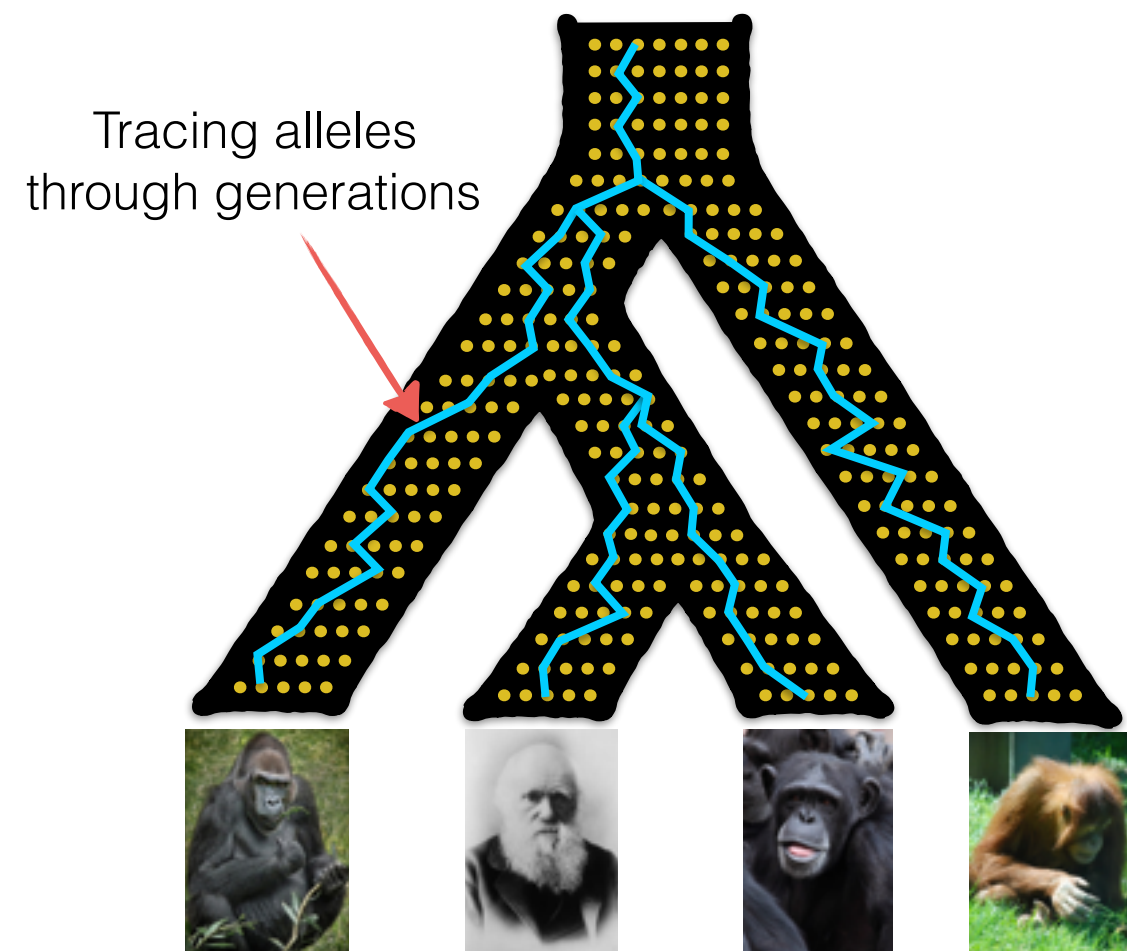
```
CTGAGCATCG
CTGAGC-TCG
ATGAGC-TC-
CTGA-CAC-G
```

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

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

# Incomplete Lineage Sorting (ILS)

- The coalescent process extended to multiple species

  - Omnipresent; most likely for rapid radiations, like birds



Tracing alleles through generations

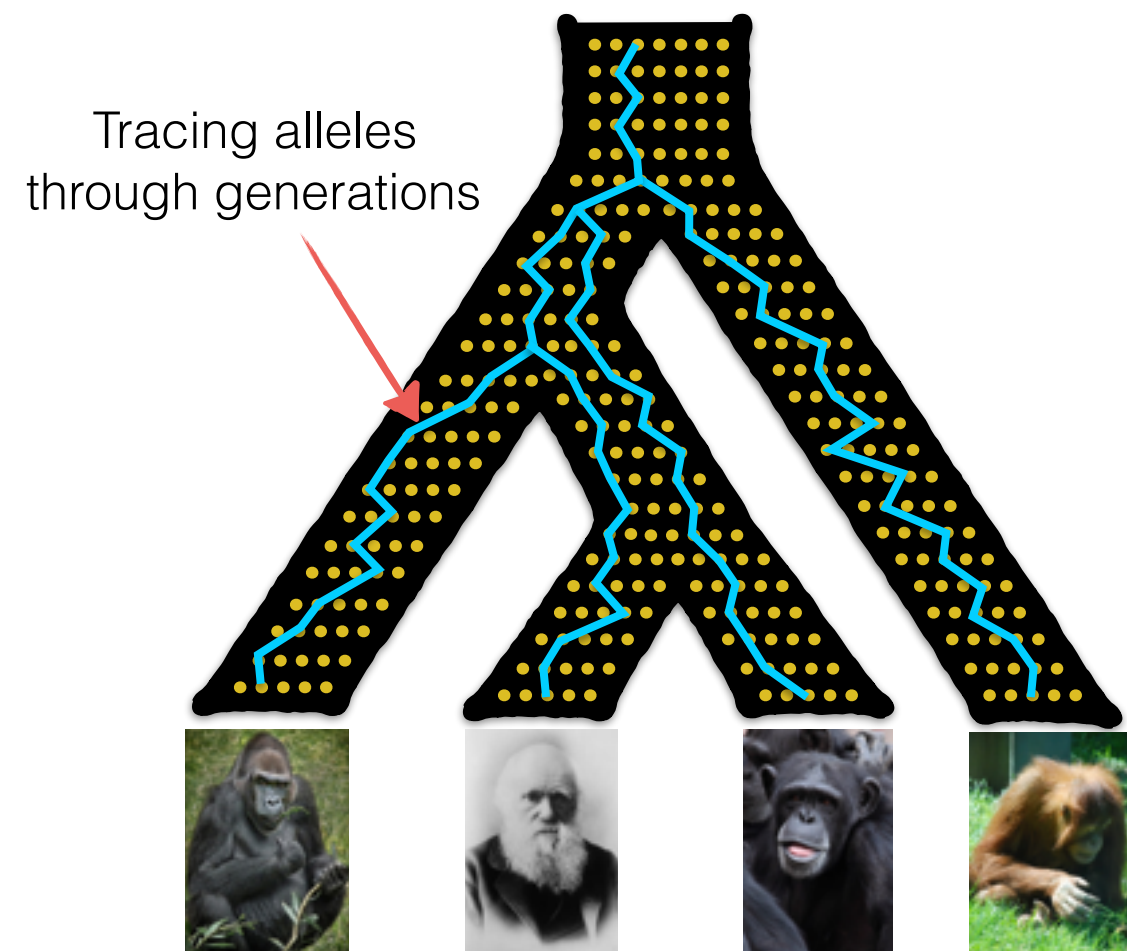# Incomplete Lineage Sorting (ILS)

- The coalescent process extended to multiple species

  - Omnipresent; most likely for rapid radiations, like birds

Tracing alleles through generations

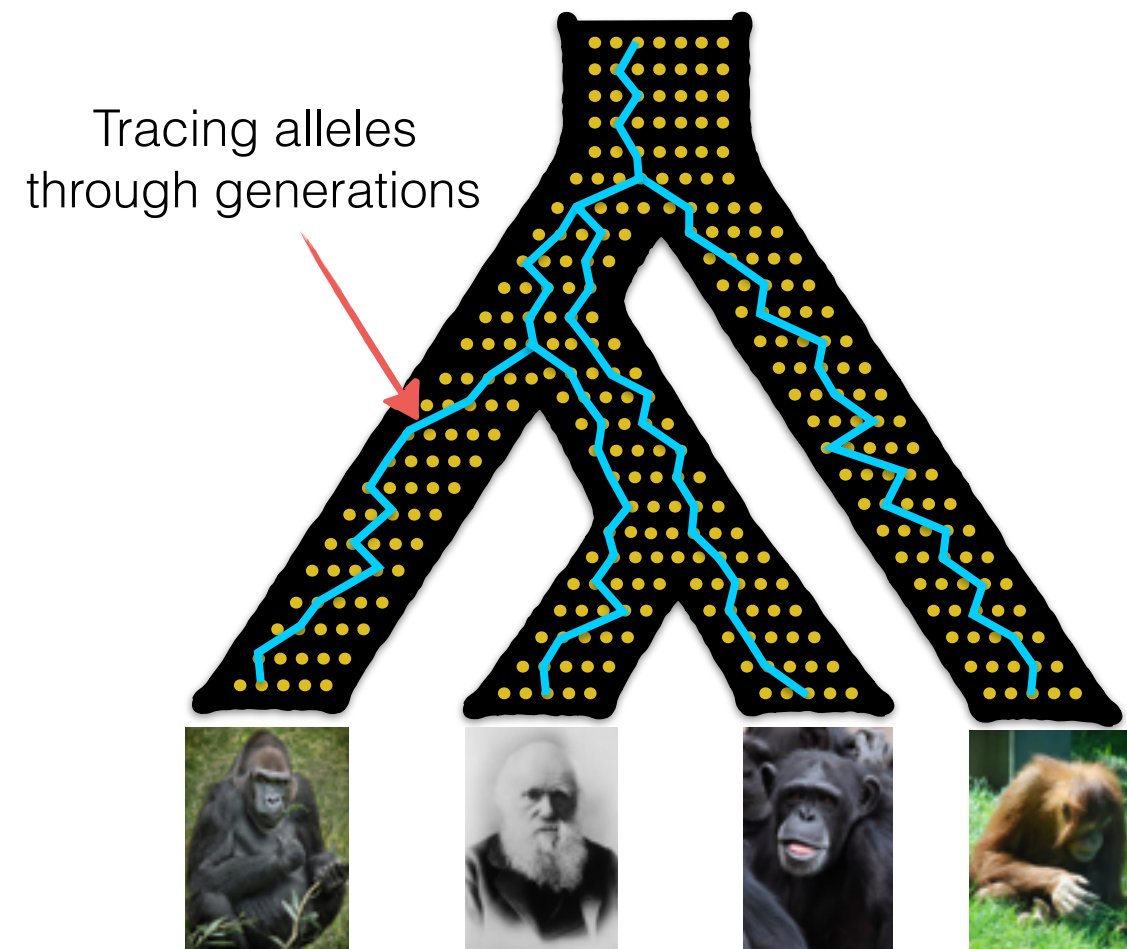# Incomplete Lineage Sorting (ILS)

- The coalescent process extended to multiple species

  - Omnipresent; most likely for rapid radiations, like birds

- Multi-species coalescent. The species tree defines the probability distribution on gene trees, and is identifiable from the distribution on gene tree topologies
[Degnan and Salter 2005]



Tracing alleles through generations

# Challenges

- What is a gene or a species and how do we <u>find</u> them?

# Challenges

- What is a gene or a species and how do we <u>find</u> them?

- Modeling: multiple evolutionary processes operate together, sometimes creating patterns that are hard to distinguish. How do we <u>untangle</u> them?

# Challenges

- What is a gene or a species and how do we <u>find</u> them?

- Modeling: multiple evolutionary processes operate together, sometimes creating patterns that are hard to distinguish. How do we <u>untangle</u> them?

- Inference: phylogenomics is hard. Dealing with multi-locus datasets and complex evolutionary processes often requires approximations.
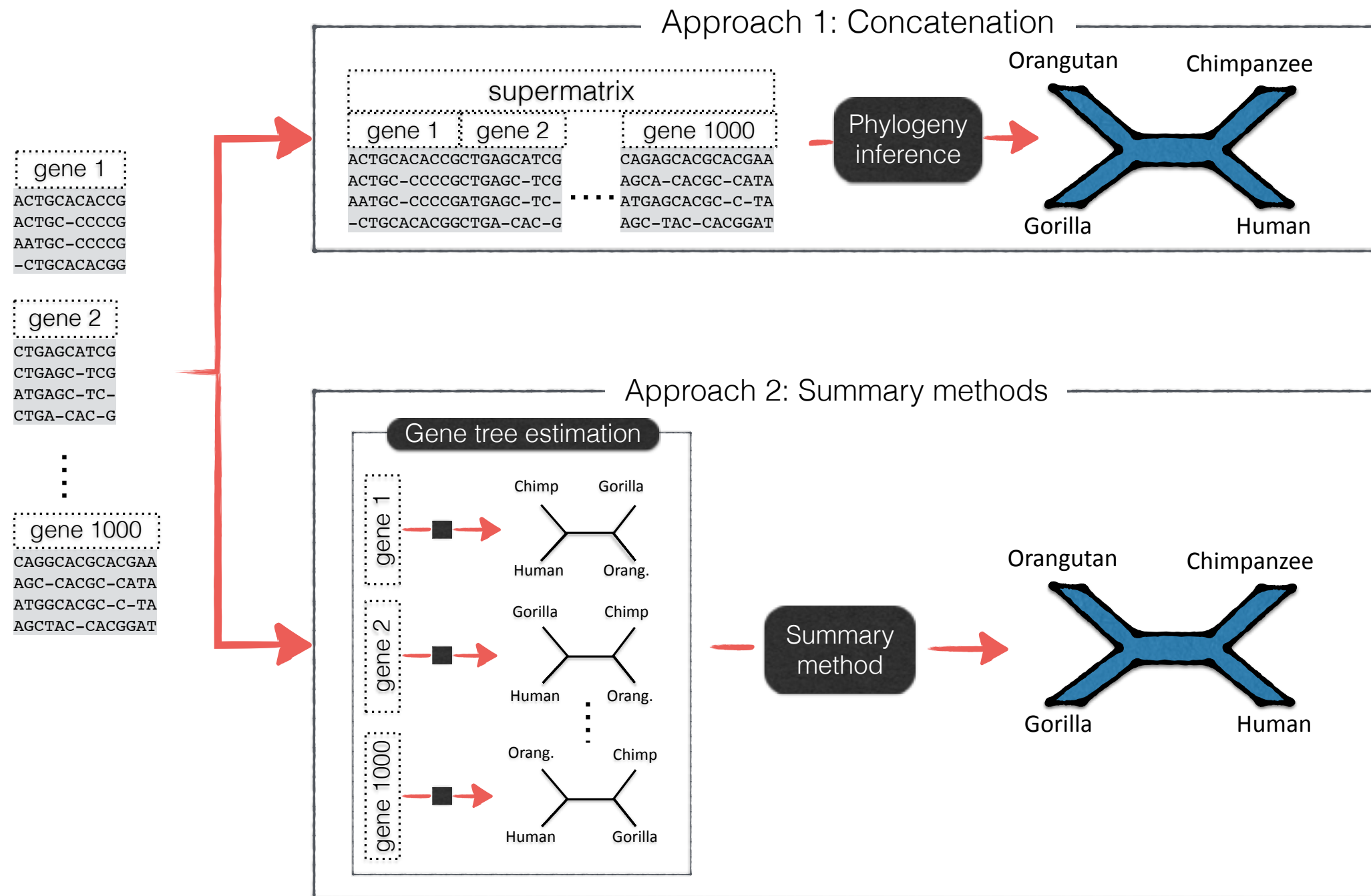
# Challenges

- **What** is a gene or a species and how do we <u>find</u> them?

- **Modeling**: multiple evolutionary processes operate together, sometimes creating patterns that are hard to distinguish. How do we <u>untangle</u> them?

- **Inference**: phylogenomics is hard. Dealing with multi-locus datasets and complex evolutionary processes often requires approximations.
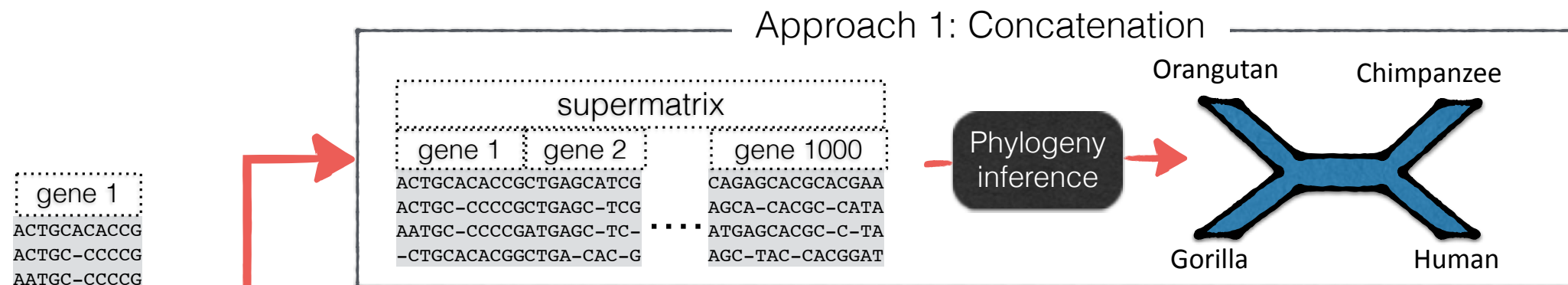
- **Reliability** and interpretation

# Challenges

- What is a gene or a species and how do we <u>find</u> them?

- Modeling: multiple evolutionary processes operate together, sometimes creating patterns that are hard to distinguish. How do we <u>untangle</u> them?

- Inference: phylogenomics is hard. Dealing with multi-locus datasets and complex evolutionary processes often requires approximations.

- Reliability and interpretation

- Catching up with new types of data

# Challenges

- What is a gene or a species and how do we <u>find</u> them?

- Modeling: multiple evolutionary processes operate together, sometimes creating patterns that are hard to distinguish. How do we <u>untangle</u> them?

- **Inference: phylogenetics is hard. Dealing with multi-locus datasets and complex evolutionary processes is often requires approximations.**

- Reliability and interpretation

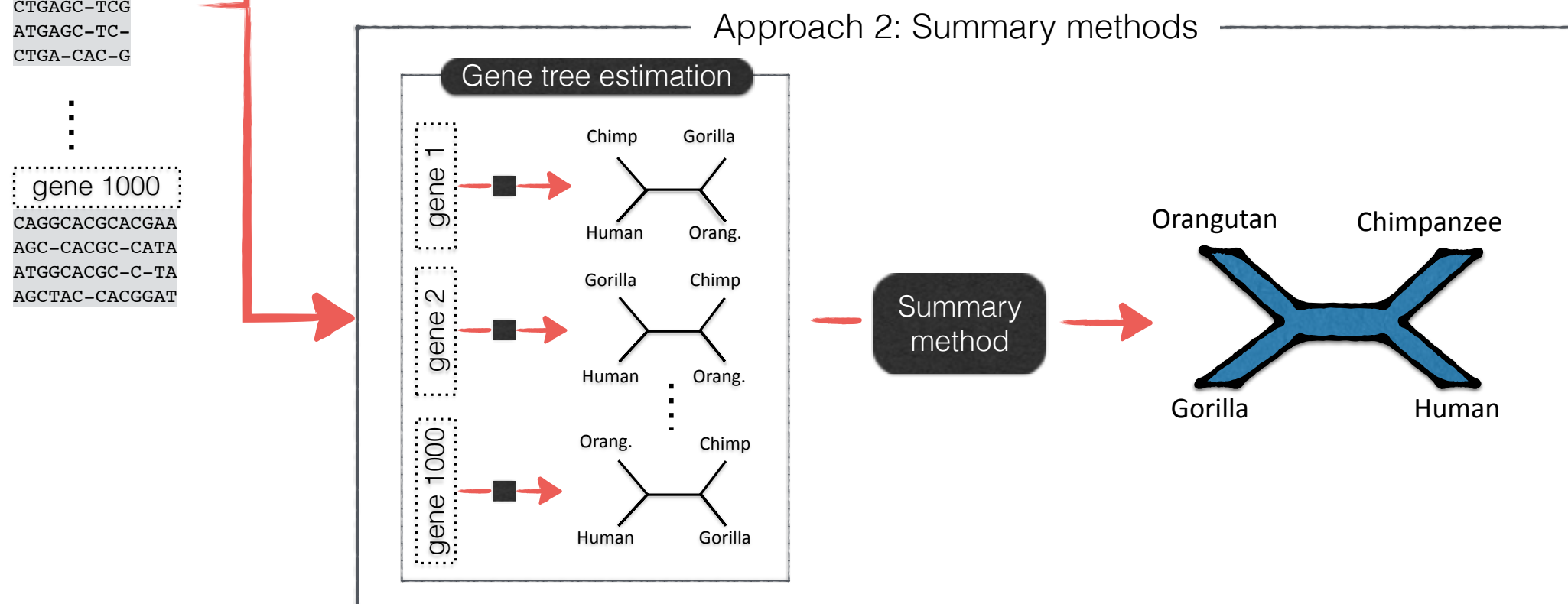- Catching up with new types of data

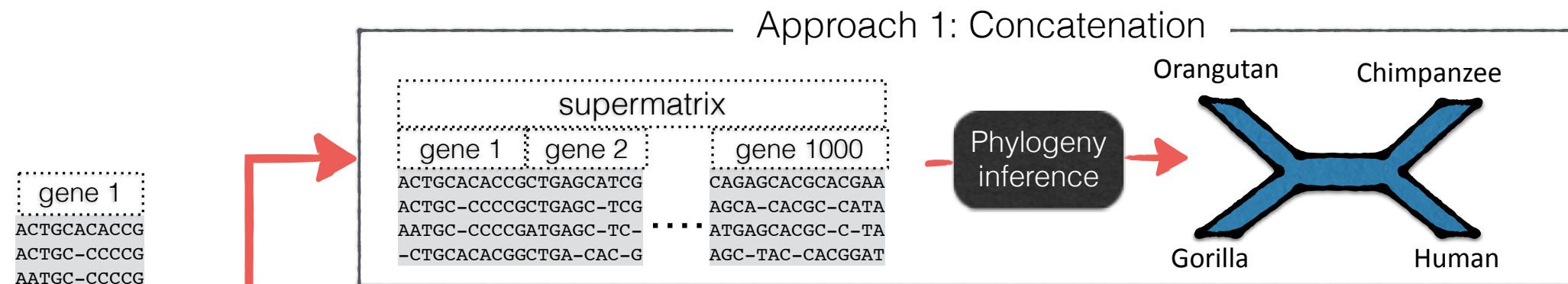# Multi-gene species tree estimation

# Multi-gene species tree estimation
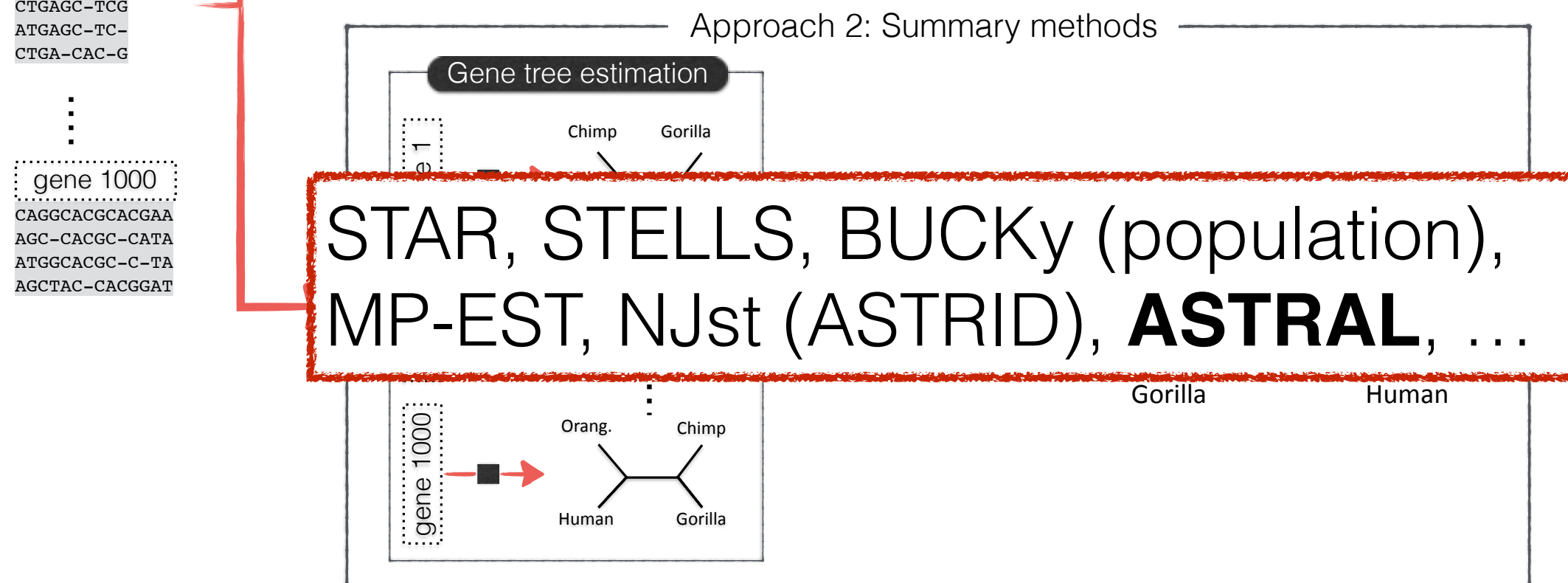


Statistically <u>inconsistent</u> [Roch and Steel,2014]

Can be statistically consistent given <u>true gene trees</u>
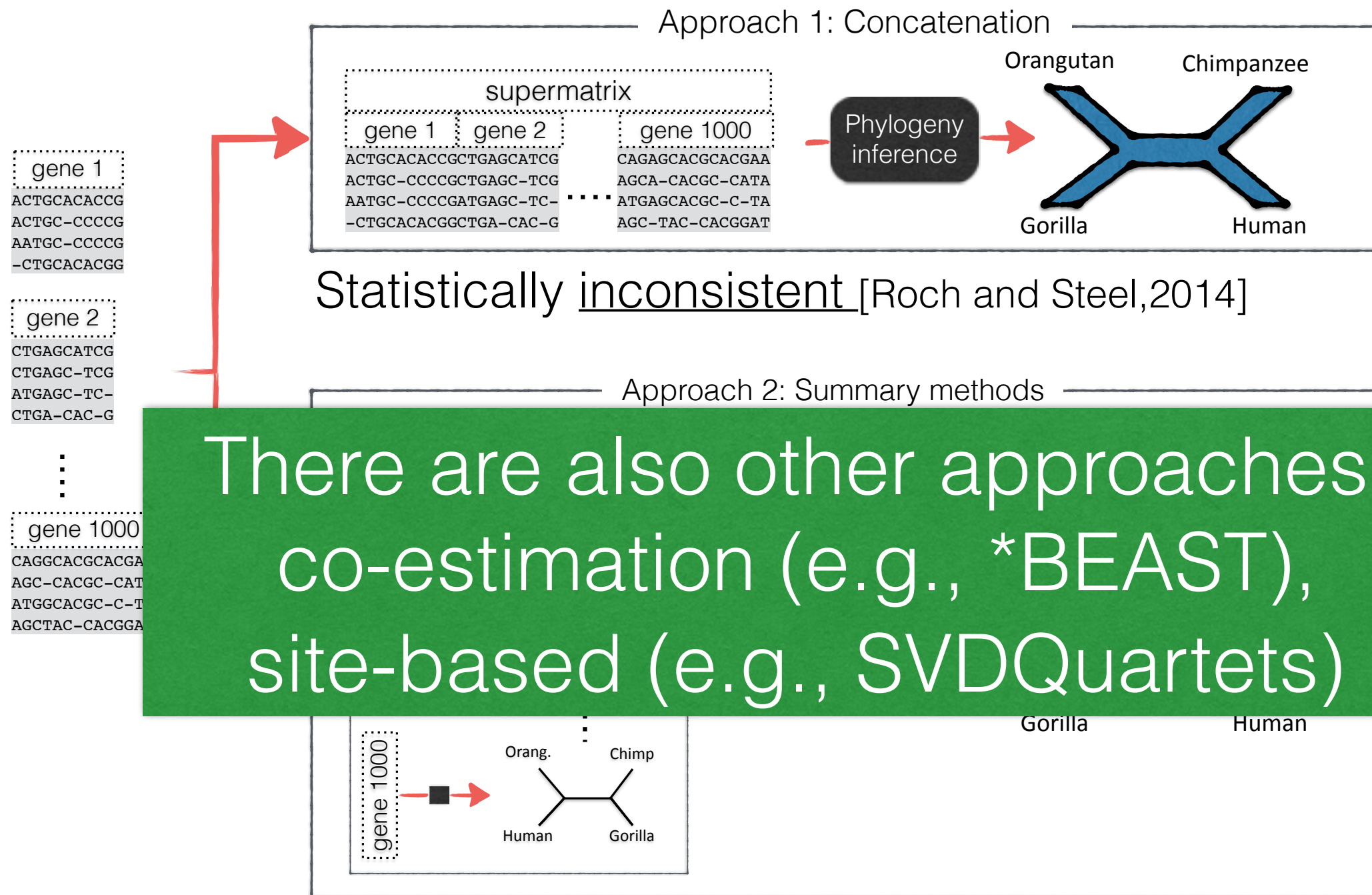
# Multi-gene species tree estimation

# Multi-gene species tree estimation



Statistically <u>inconsistent</u> [Roch and Steel, 2014]

There are also other approaches: co-estimation (e.g., *BEAST), site-based (e.g., SVDQuartets)

Can be statistically consistent given <u>true gene trees</u>

# ASTRAL

- **Input:** A set of inferred unrooted gene trees

- **Output:** A species tree with branch lengths in coalescent units and branch support values

# ASTRAL

- **Input:** A set of inferred unrooted gene trees

- **Output:** A species tree with branch lengths in coalescent units and branch support values

- **Approach**: try to find the species tree that shares the maximum number of quartet trees with input gene trees
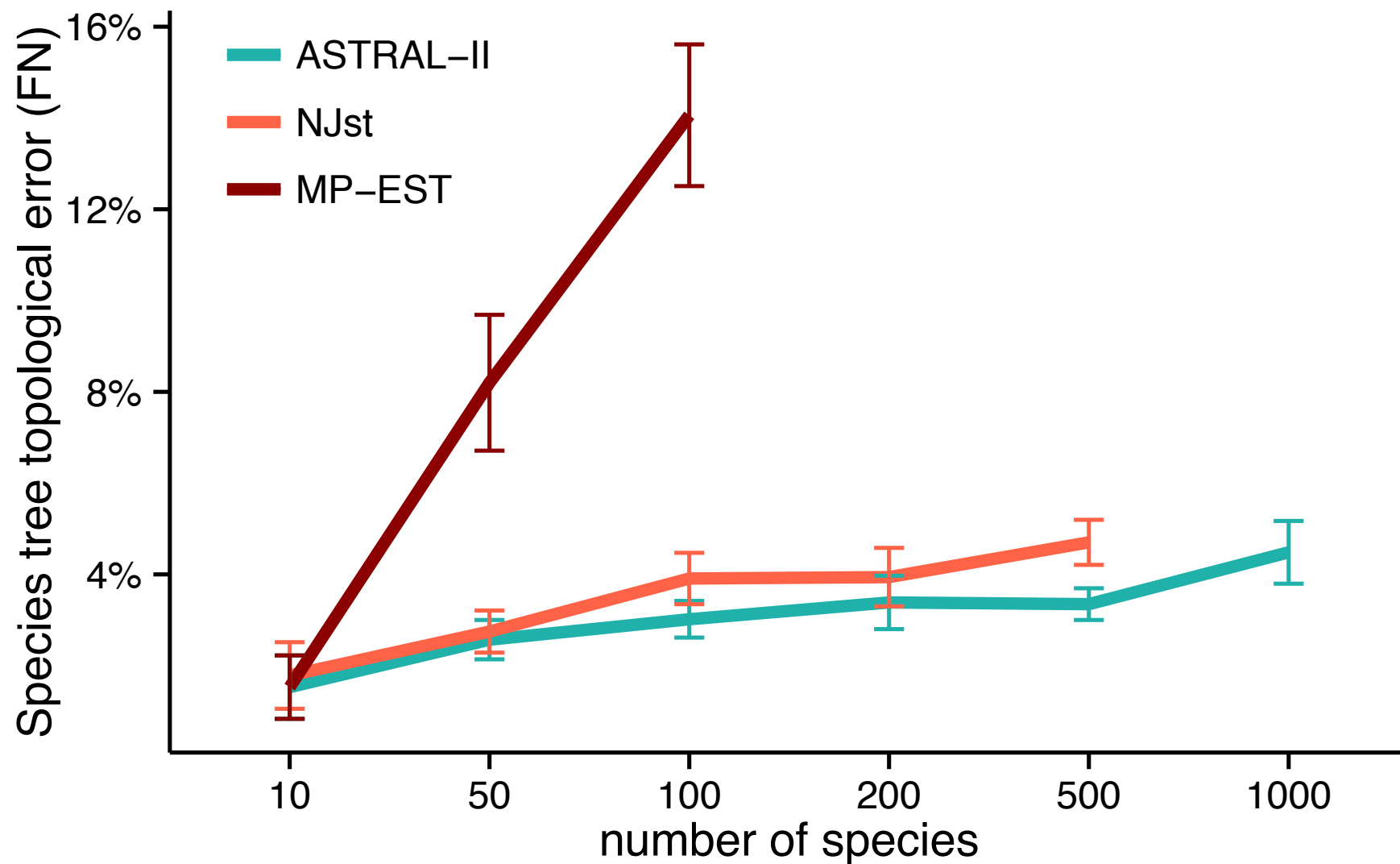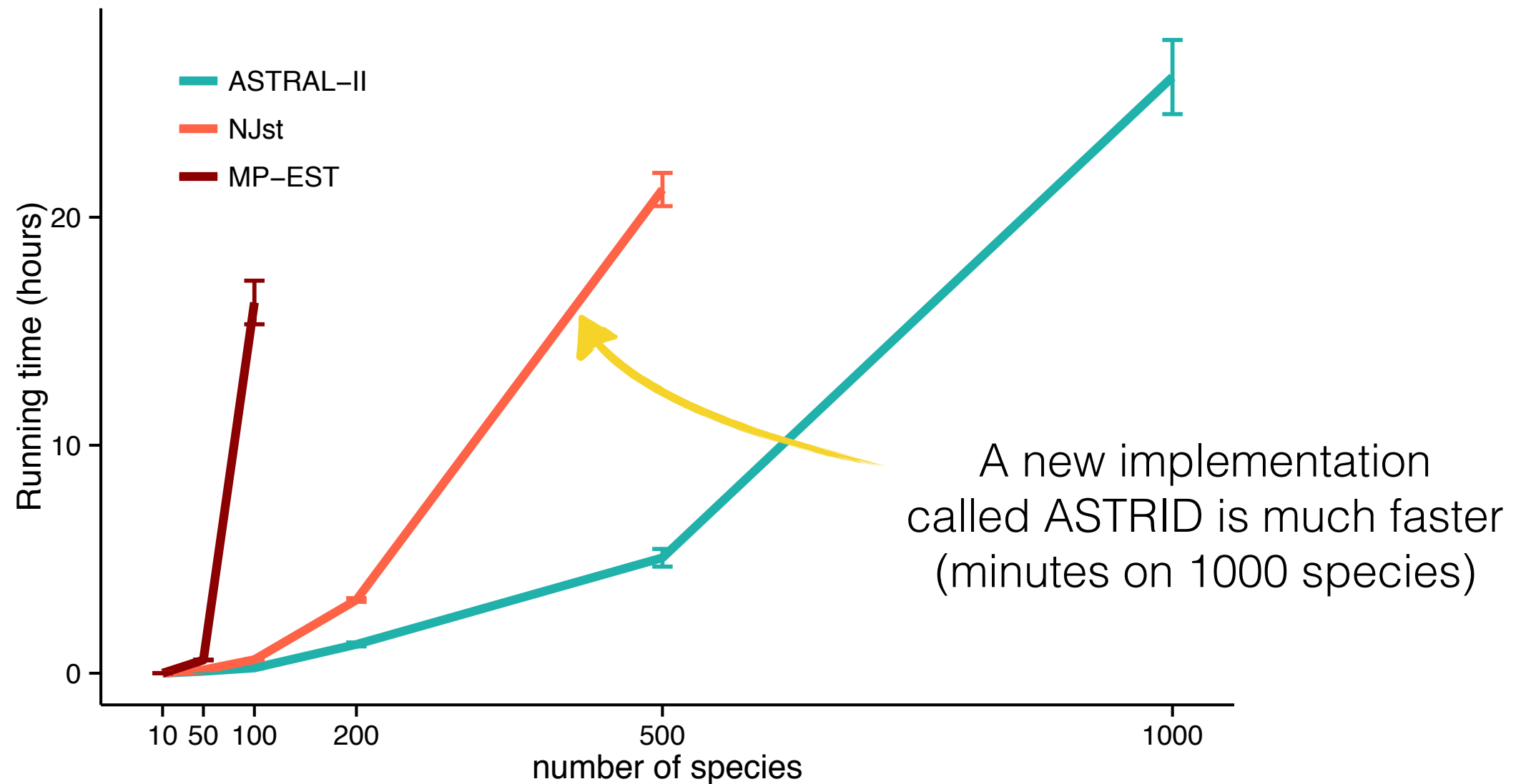
# ASTRAL

- **Input:** A set of inferred unrooted gene trees

- **Output:** A species tree with branch lengths in coalescent units and branch support values

- **Approach**: try to find the species tree that shares the maximum number of quartet trees with input gene trees

- **Designed for:**

  - Accuracy (established in simulation studies)

  - Scalability: the default version runs on a thousand genes from a thousand species in a day
    —> Important for next phases of B10K
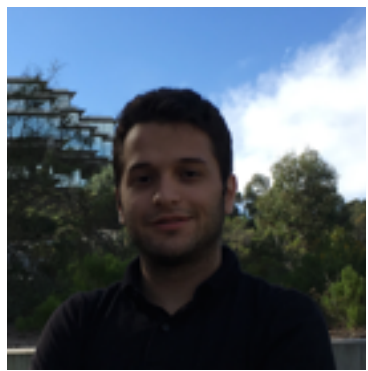
# ASTRAL: accurate and scalable



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

9

# Running time as function of # species



A new implementation
called ASTRID is much faster
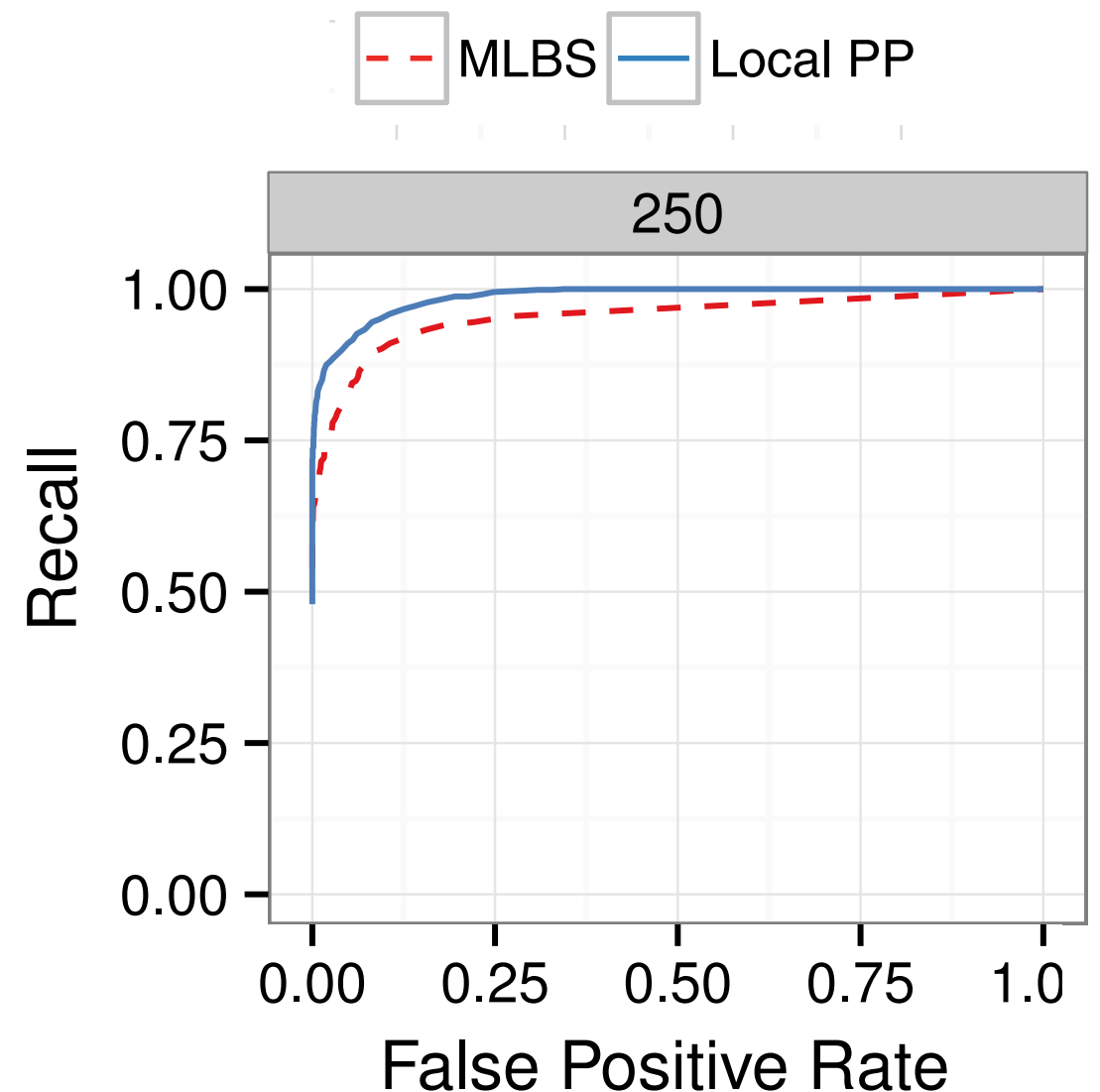(minutes on 1000 species)

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

# Local branch support

Erfan Sayyari

- Use frequency of quartets defined *around* each species tree branch and some strong assumptions to compute statistical support for each species tree branch

- Extremely scalable.

  - Doesn't need bootstrapped gene trees

  - Minutes on 1K species
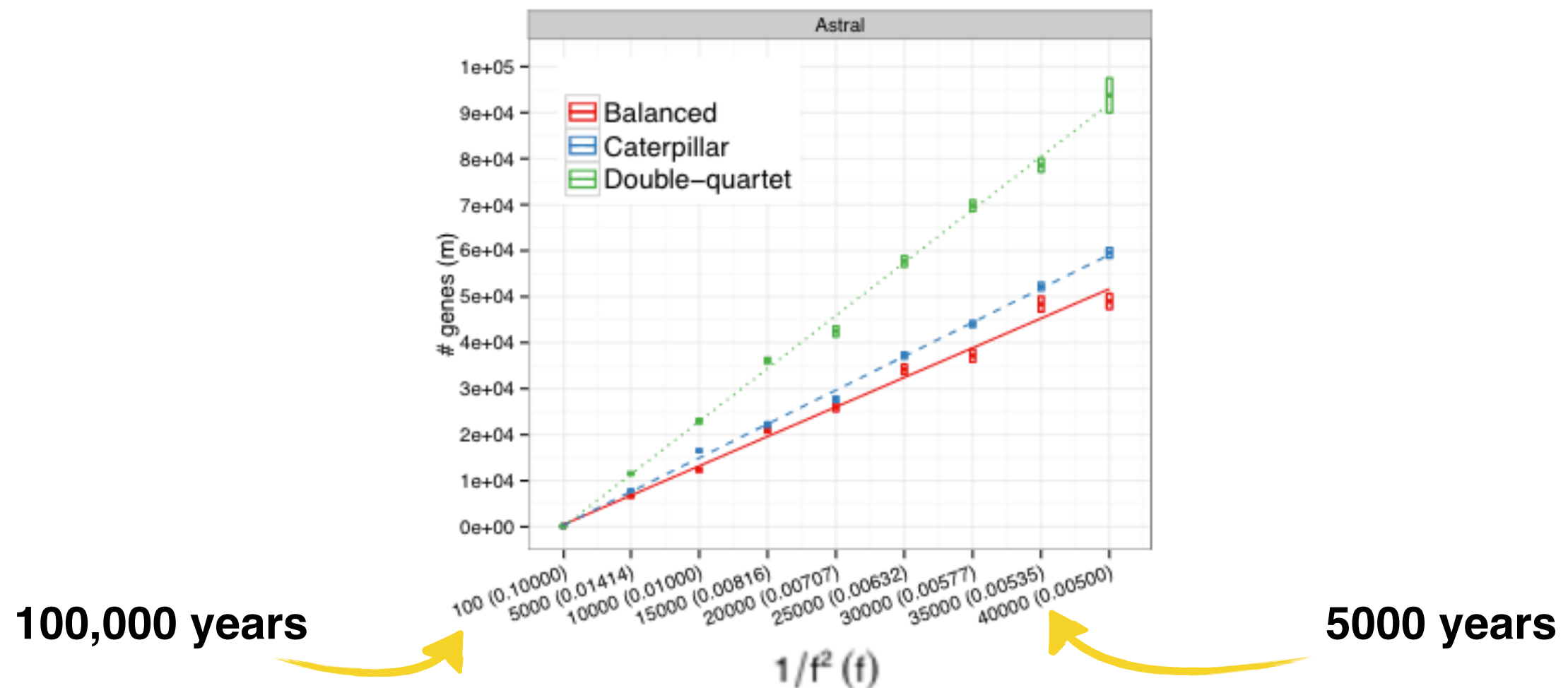
- More accurat than traditional bootstrapping



Avian simulated dataset
(48 taxa, 1000 genes)

[Sayyari, Mirarab, MBE 2016]

# How many genes does ASTRAL need?

Depends on the branch length ($f$) and the number of species ($n$):
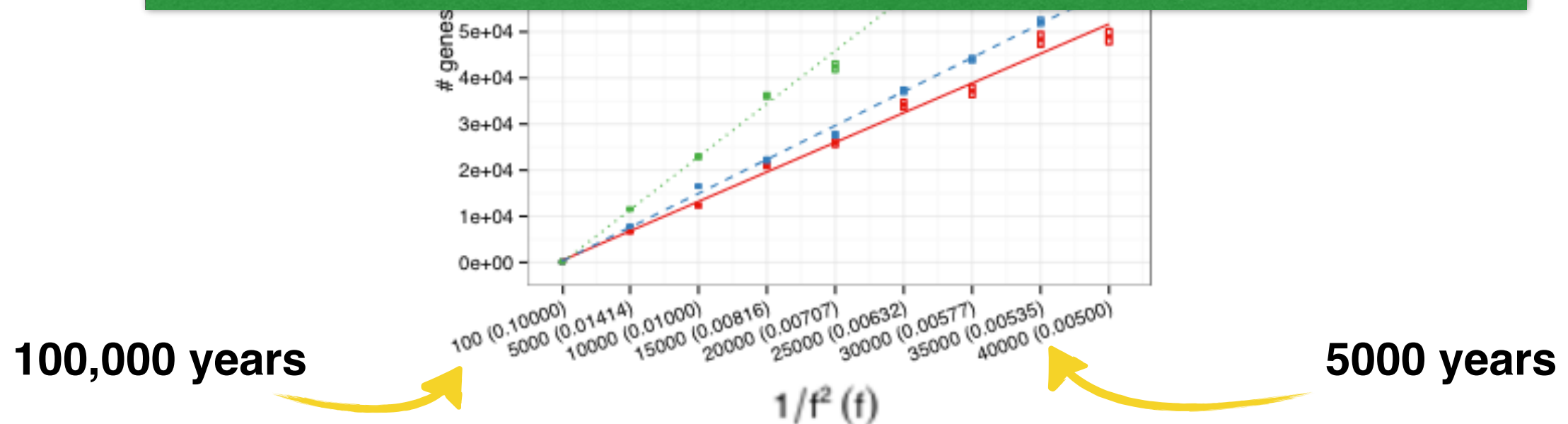
$$log(n)\, f^{\,-2}$$



**100,000 years**                    **5000 years**

# How many genes does ASTRAL need?

Depends on the branch length ($f$) and the number of species ($n$):

$$log(n) \, f^{-2}$$



**100,000 years**          **5000 years**

# Ongoing improvements to ASTRAL

- A GPU implementation is 10-20X faster (less than an hour on 1000 species).

- Improved measures of support (fewer and weaker assumptions)

- Better ways of dealing with multiple individuals from the same species

- Divide-and-conquer to enable analyzing many tens of thousands of species

# Unsolved challenges for the next phase of B10K

- **Hybridization:** Inferring species networks is doable now, but not on the scales targeted by B10K

- **Gene tree error/uncertainty :** gene tree uncertainty is notoriously high for avian genomes.

  - We previously proposed statistical binning for this problem. It is not clear that binning will scale to hundreds of species

- **Standard pipelines for quality control:** for example, for dealing with fragmentary data, codon bias, etc.

- And many more …