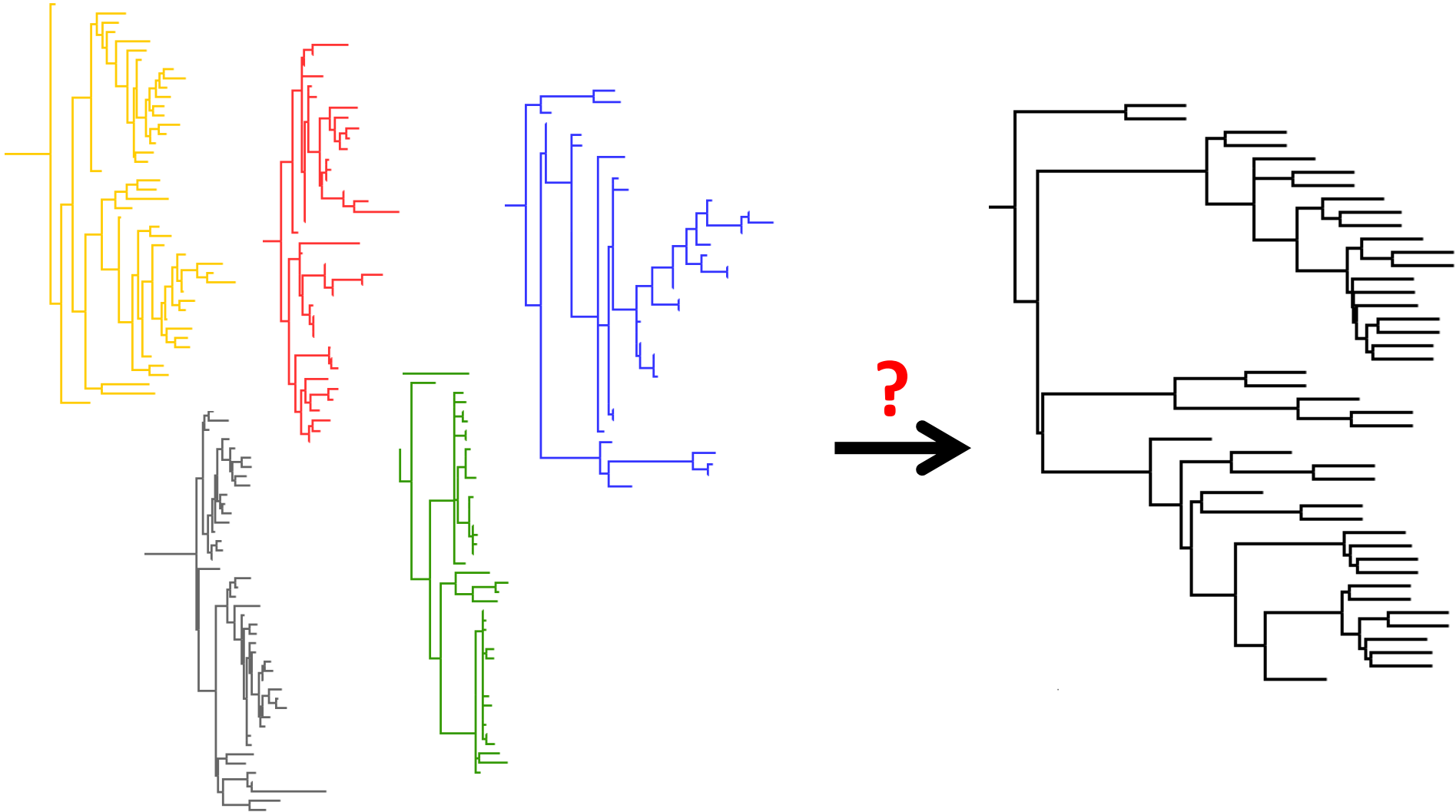# Gene trees vs species tree

## Phylogenetic methods
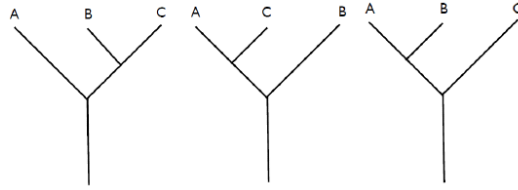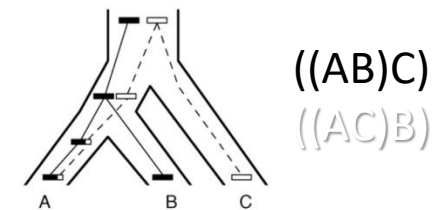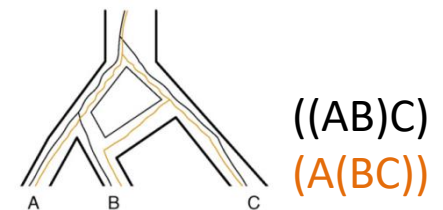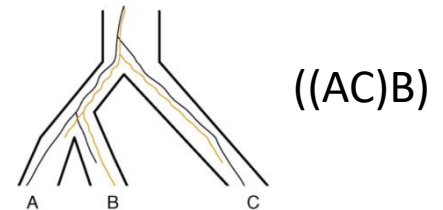
Tomáš Fér

2024

# Species tree from gene trees

# Incongruencies among loci: gene trees vs species tree



- incomplete lineage sorting (ILS)

- horizontal gene transfer (HGT)
  - affects small DNA segments

- gene duplication and loss (GDL)
  - orthology problem

- hybridization
  - affects whole genomes

- recombination
  - different histories for neighboring segments in genes



(A(BC))

((AC)B)

((AB)C)
(A(BC))

((AB)C)
((AC)B)

Degnan & Rosenberg, 2009

# Gene duplications and losses



A   B   C    D

Gene tree

A   B C    D
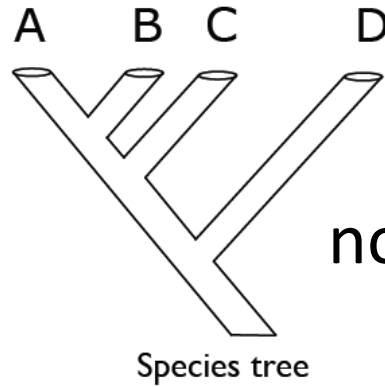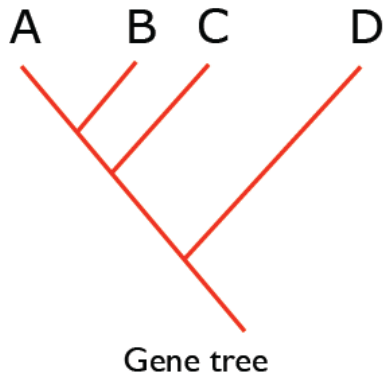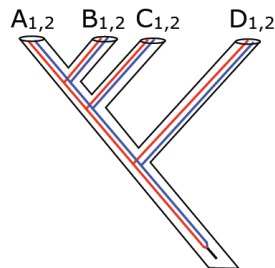
Species tree

no gene duplication

A   B C    D

Gene tree = Species tree

$A_1$   $B_1$ $C_1$   $D_1$ $D_2$   $C_2$   $B_2$   $A_2$

$A_{1,2}$   $B_{1,2}$ $C_{1,2}$    $D_{1,2}$

gene duplications and losses (GDL)

$A_1$   $B_1$ $C_1$   $D_1$ $D_2$   $C_2$ $B_2$   $A_2$

Gene loss

← Gene loss

← gene loss

← gene duplication

true species tree

inferred species tree

A   B   C    D

$A_1$   $B_1$ $C_1$   $D_1$ $D_2$   $C_2$ $B_2$   $A_2$

A   B   D    C

M. Popp, Oslo

# Incomplete lineage sorting
## Coalescence processes



https://frederikleliaert.wordpress.com/green-algae/dna-based-species-delimitation-in-algae/

M. Popp, Oslo

# Species tree estimation



- **concatenation**

- multispecies coalescence
  - *BEAST (**coestimation** of gene trees and species tree)
  - **summary** methods (combining gene trees)

- **supertree** methods
  - MRL (maximum representation using likelihood)

- Bayesian **concordance** analysis (BUCKy)
  - quartet-based Bayesian species tree estimation

- **site-based** methods
  - SNAPP, SVDquartets

# Species tree estimation



- **concatenation**

- multispecies coalescence
  - *BEAST (**coestimation** of gene trees and species tree)
  - **summary** methods (combining gene trees)

- **supertree** methods
  - MRL (maximum representation using likelihood)

- Bayesian **concordance** analysis (BUCKy)
  - quartet-based Bayesian species tree estimation

- **site-based** methods
  - SNAPP, SVDquartets

# Concatenation

- put all the loci after each other (superalignment, supermatrix)
- very good accuracy under low ILS model conditions
- i.e., good approach unless strong ILS

- **single** partition model
  - the whole alignment analyzed with the same parameters
  - statistically inconsistent

- **multiple** partitions model (ML or Bayesian)
  - each alignment (or even codon position) analyzed with separate parameters
  - best partitioning scheme by, e.g., PartitionFinder or ModeltestNG or IQtree
  - fully partitioned analysis
  - maximum likelihood (CA-ML) –  RAxML-ng, ExaML
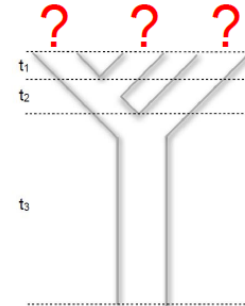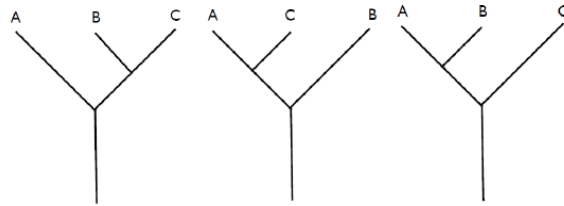  - or Bayesian inference – MrBayes, ExaBayes

# Species tree estimation



- **concatenation**

- multispecies coalescence
  - *BEAST (**coestimation** of gene trees and species tree)
  - **summary** methods (combining gene trees)

- **supertree** methods
  - MRL (maximum representation using likelihood)

- Bayesian **concordance** analysis (BUCKy)
  - quartet-based Bayesian species tree estimation
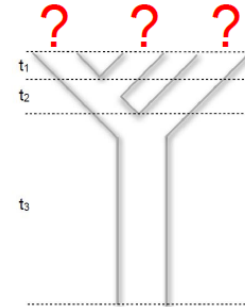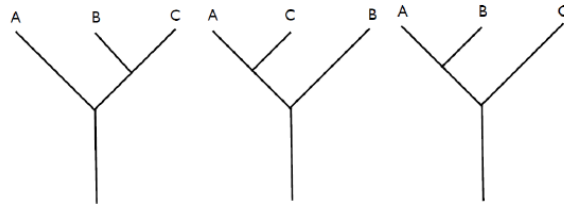
- **site-based** methods
  - SNAPP, SVDquartets

# Multispecies coalescent

- coalescent model applied to gene trees in a species tree
  - combines coalescent and birth-death models

coalescence

birth-death model

# Multispecies coalescent

- used to assemble separate coalescent processes occurring in populations connected by an evolutionary tree
  - coalescent tree distribution (probability of sharing common ancestor t generations back)
  - birth-death model with stochastic rate of birth and death
  - describes probability of gene tree(s) within a species tree

coalescence

birth-death model

# Multispecies coalescent

| | |
|---|---|
| **dot** | individual gene copy |
| **row** | generation |
| **lines** | connection to gene ancestors in previous generations |

Species tree

**the coalescent in several populations**
pink pops – only one lineage (gene copy)

Incomplete lineage sorting

fail to coalesce within populations

A   B   C   D

**populations arranged by evolutionary relationships**

**gene tree ((AB)(CD)) in a species tree ((AB)C)D)**

A   B   C   D

Degnan & Rosenberg, 2009

# Multispecies coalescent

- (incomplete) lineage sorting
  - particular types of genealogical pattern
  - process explaining gene tree discordance
  - failure of lineages in a population to coalesce

Incomplete
lineage sorting

fail to coalesce within populations

A     B     C     D

# *BEAST

## STAR-BEAST = **S**pecies **T**ree **A**ncestral **R**econstruction

- Bayesian framework for species tree reconstruction

- assumptions
  - no recombination within locus
  - free recombination between loci
  - no hybridization
  - each sample mapped to appropriate species

prior distribution

likelihood

posterior distribution

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

marginal likelihood

original Bayesian theorem

probability of the species tree S given the data (D)

gene tree likelihood

prior on species tree

multispecies coalescent likelihood (prior on gene tree given species tree)

$$f(\mathrm{g}, S|D) = \frac{f(S)}{\Pr(D)} \prod_{i=1}^{m} \Pr(D_i|g_i)f(g_i|S),$$

gene tree

marginal likelihood

Drummond & Bouckaert, 2015

# *BEAST

## STAR-BEAST = **S**pecies **T**ree **A**ncestral **R**econstruction

- co-estimates gene trees and species tree

- most accurate species tree method

- computationally intensive

- not suitable for large datasets, i.e.
  - no more than ~50 loci
  - no more than ~20-30 species

- BBCA – divide-and-conquer technique (Zimmerman et al., 2014)

# Summary methods

Estimate each gene tree independently -> summarize them

require rooted gene trees

- MP-EST – **m**aximum **p**seudo-likelihood approach for **e**stimating **s**pecies **t**rees

- STAR – **s**pecies **t**ree estimation using **a**verage **r**anks of coalescences

unrooted gene trees

- STEAC – **s**pecies **t**ree **e**stimation using **a**verage **c**oalescence times

- ASTRAL – **A**ccurate **S**pecies **T**ree **R**econstruction **AL**gorithm

- ASTRID – **A**ccurate **S**pecies **TR**ees from **I**nternode **D**istances (reimplementation of NJ$_{st}$ method)

site-based methods (estimate species trees from the distribution on site pattern within unlinked loci)

- SNAPP – SNP and AFLP Package for Phylogenetic analysis

- SVDquartets

# Tree reconstruction from quartets

- quartet – unrooted tree over 4 taxa
- three possible quartets
- only one quartet *q* is consistent with final tree **T**



Reaz et al. (2015): *Accurate Phylogenetic Tree Reconstruction from Quartets: A Heuristic Approach*. PLoS ONE 9, e104008.

# Tree reconstruction from quartets

- quartet – unrooted tree over 4 taxa
- three possible quartets
- only one quartet $q$ is consistent with final tree **T**



Reaz et al. (2015): *Accurate Phylogenetic Tree Reconstruction from Quartets: A Heuristic Approach*. PLoS ONE 9, e104008.
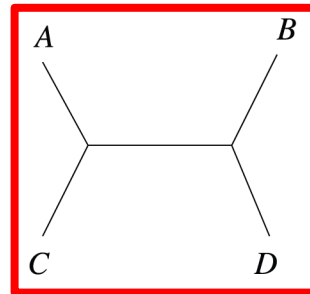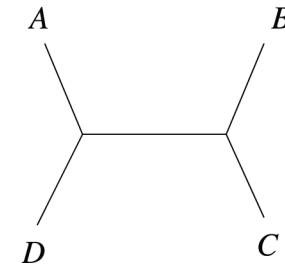
# ASTRAL

## Accurate Species Tree Reconstruction Algorithm
https://github.com/smirarab/ASTRAL

- unrooted gene trees

- species tree that agrees with the largest number of quartet trees induced by the set of gene trees

- weighting all three alternative quartet topologies according to their relative frequencies within gene trees
  - much more frequent topology – trees without this topology are penalized
  - similar frequencies (i.e., close to 0.33) – the quartet has little impact to optimization

- final species tree with
  - local posterior probability that the branch is in the species tree
  - the length of internal branches in coalescent units

Siavash Mirarab

# Unrooted quartets under MSC model

- **for a quartet (4 species) –** the most probable unrooted quartet tree (among the gene trees) is the unrooted species tree topology

- **for 5 or more species –** the unrooted species tree topology can be different from the most probable gene tree (called "anomaly zone")
  - break gene trees into quartets of species
  - find the species tree with the maximum number of induced quartet trees shared with the collection of input gene trees (NP-hard optimization problem)
  - statistically consistent under the multispecies coalescent model with error-free input
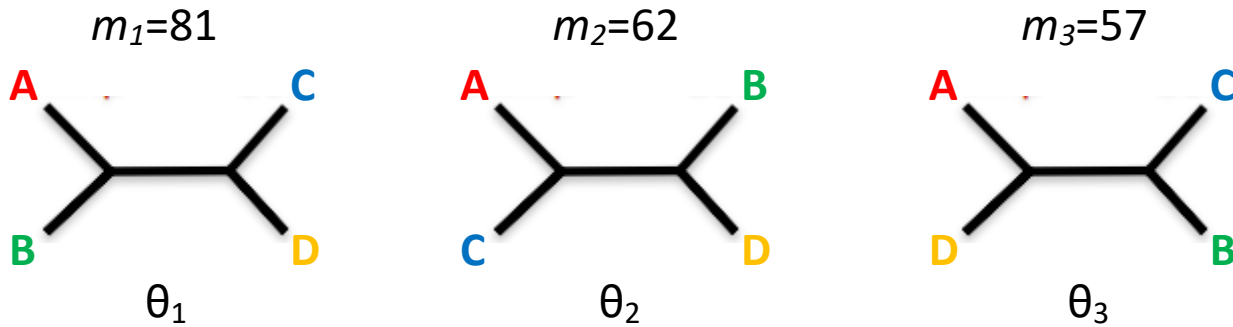  - solved by dynamic programming – ASTRAL

# ASTRAL input/output

- input – unrooted gene trees
  - missing data allowed
  - polytomies allowed
  - multiple alleles per species allowed

- output – estimated unrooted species tree
  - **branch lengths** in coalescent units (on internal branches)
  - measure **of branch support** (LPP, local posterior probability)

# Local posterior probability

- quartet frequencies follow a multinomial distribution

$m_1 = 81$

A — C

B — D

$\theta_1$

$m_2 = 62$

A — B

C — D

$\theta_2$

$m_3 = 57$

A — C

D — B

$\theta_3$

- $P$ (gene tree seen $m_1/m$ times = species tree) = $P(\theta_1 > 1/3)$

  - possible to solve analytically
  - resulting measure is localPP
  - for $n>4$ – averaging quartet scores

- more accurate and faster than multi-locus bootstrap (MLBS; Seo 2008)

increased number of genes  =  increased support
decreased discordance       =  increased support

http://tandy.cs.illinois.edu/astral-apro.pdf

# Branch length of ASTRAL trees

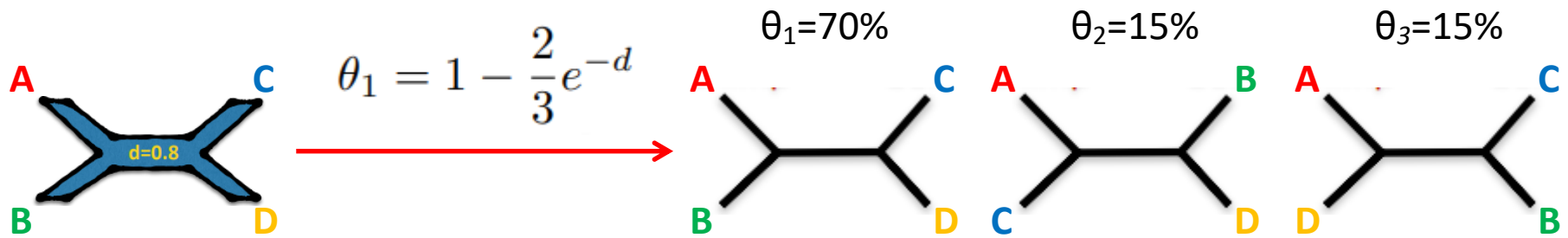- branch length in coalescence units = the level of discordance
- for a single quartet (i.e., *n*=4) – reverse the discordance formula to get multilocus estimate
- for n>4 – average frequencies around the branch



$$\theta_1 = 1 - \frac{2}{3}e^{-d}$$

$\theta_1 = 70\%$  $\theta_2 = 15\%$  $\theta_3 = 15\%$

Sayyari & Mirarab, 2016, MBE
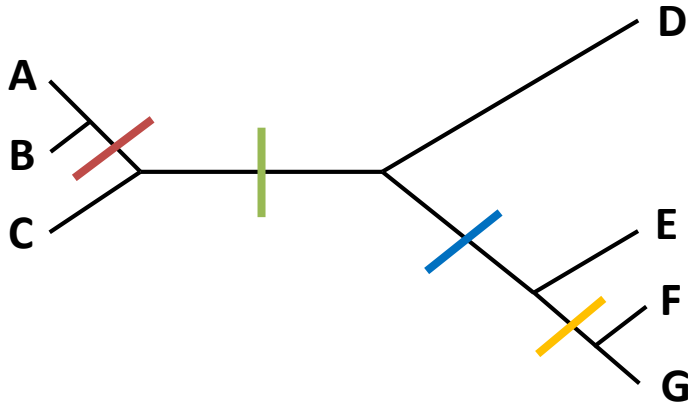
# ASTRAL problems

- assumption for statistical consistency
  - randomly distributed sample of gene trees
    - recombination-free
    - reticulation-free
    - error-free (i.e., topology correctly estimated)
    - orthologous
- in practice: reduced accuracy with low accuracy gene trees
- branch length
  - only for internal branches (unless multiple individuals per species)
  - in coalescent units, i.e., "true value" is a function of population size and generation time
- local posterior probability (LPP)
  - better than MLBS (empirically) but based on many assumptions

http://tandy.cs.illinois.edu/astral-apro.pdf

# MRL
**M**aximum **R**epresentation with **L**ikelihood; Nguyen et al. 2012

- supertree method – estimates species tree on full taxon sets from sets of smaller trees (i.e., with missing species)
- encodes a set of gene trees by a large randomized matrix
  - using mrp.jar; https://github.com/smirarab/mrpmatrix
- each edge (branch) in each gene tree
  - '0' for the taxa that are on one side of the edge
  - '1' for the taxa on the other side
  - '?' for all the remaining taxa (i.e., the ones that do not appear in the tree)
- MRL matrix is analyzed using heuristics for a symmetric 2-state Maximum Likelihood
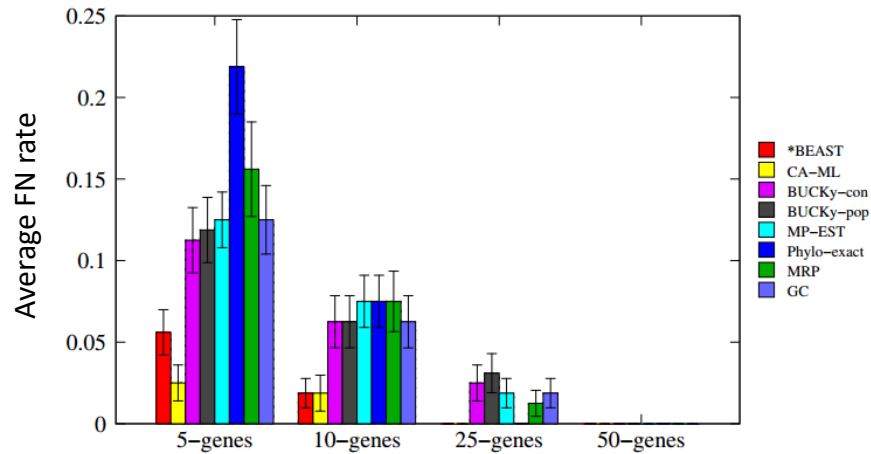  - in RAxML as 'BINGAMMA' model

# MRL binary matrix

# Methods comparison



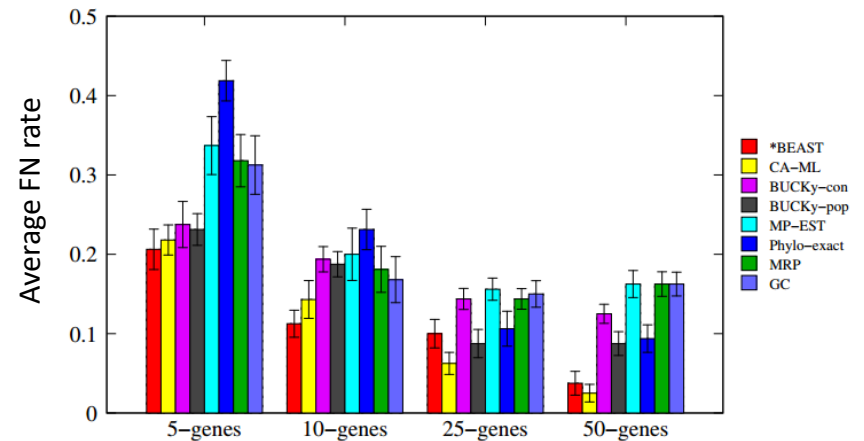## Results on 11-taxon datasets with weak ILS

*BEAST more accurate than summary methods (MP-EST, BUCKy, etc)
CA-ML: concatenated analysis) most accurate

Datasets from Chung and Ané, 2011
Bayzid & Warnow, Bioinformatics 2013

## Results on 11-taxon datasets with strongILS

*BEAST more accurate than summary methods (MP-EST, BUCKy, etc)
CA-ML: (concatenated analysis) also very accurate

Datasets from Chung and Ané, 2011
Bayzid & Warnow, Bioinformatics 2013

T. Warnow, The University of Texas
https://www.cs.utexas.edu/users/tandy/394C-nov20-2013.pdf

# Quartet Sampling (QS)

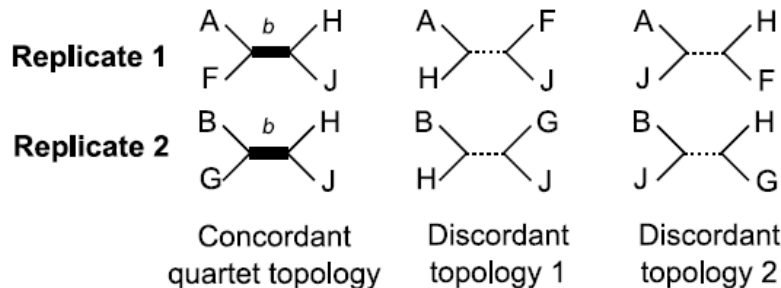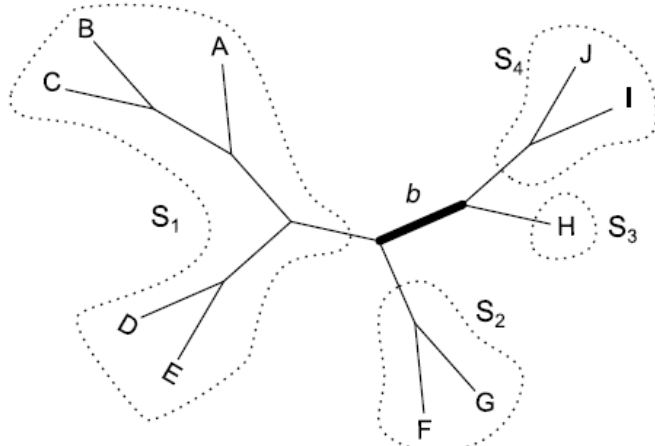### Replacement for bootstrap in phylogenomic studies...

- quartet-based evaluation system

- synthetizes several phylogenetic and genomic analytical approaches

- discordance testing

- distinguishes strong conflict from weak support

- three different scores per branch
  - Quartet Concordance (QC)
  - Quartet Differential (QD)
  - Quartet Informativeness (QI)

- terminal node score
  - Quartet Fidelity (QF)

Pease et al. (2018): Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. *American Journal of Botany* 105(3): 385–403.

# Quartet Sampling

- takes an existing phylogenetic topology and a molecular dataset
- evaluates internal branches – likelihood for all three possible phylogenies for the randomly selected quartets spanning particular branch



Pease et al. (2018): Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. *American Journal of Botany* 105(3): 385–403.

# Quartet Sampling
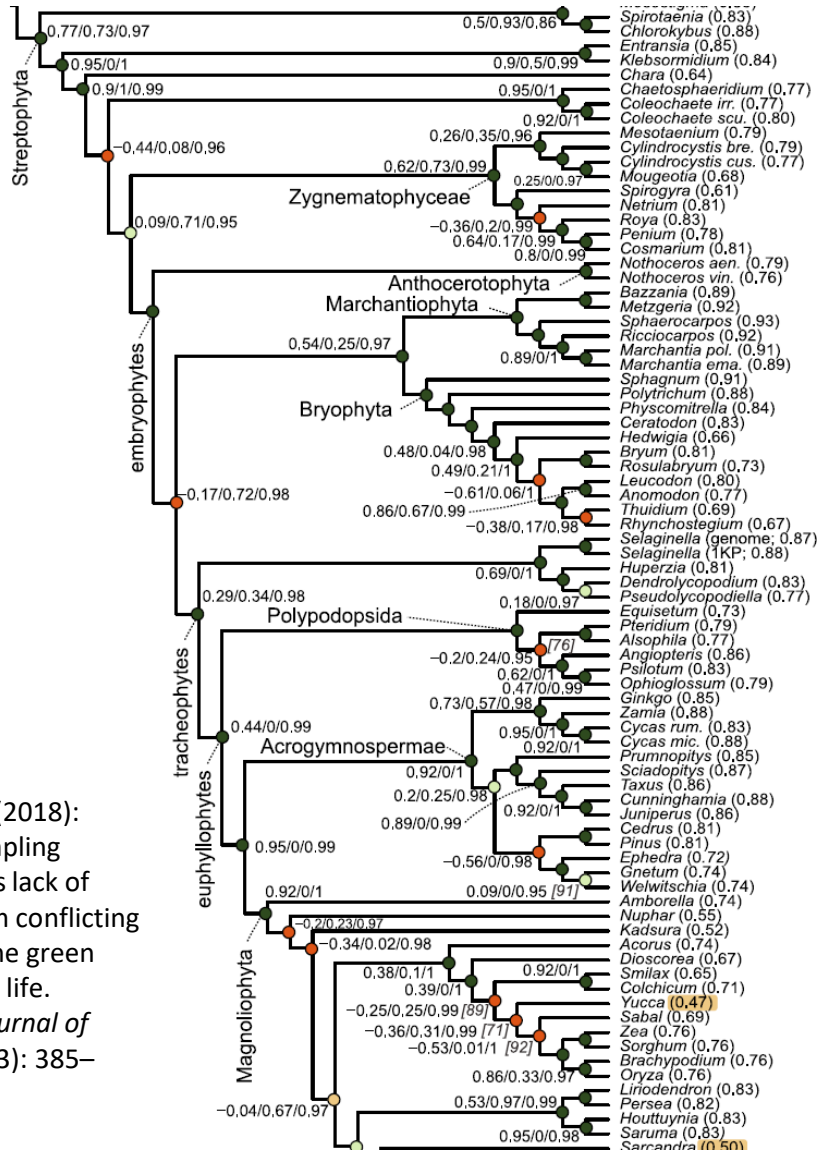## Replacement for bootstrap in phylogenomic studies…

**TABLE 1.** Quartet Sampling (QS) score interpretation.

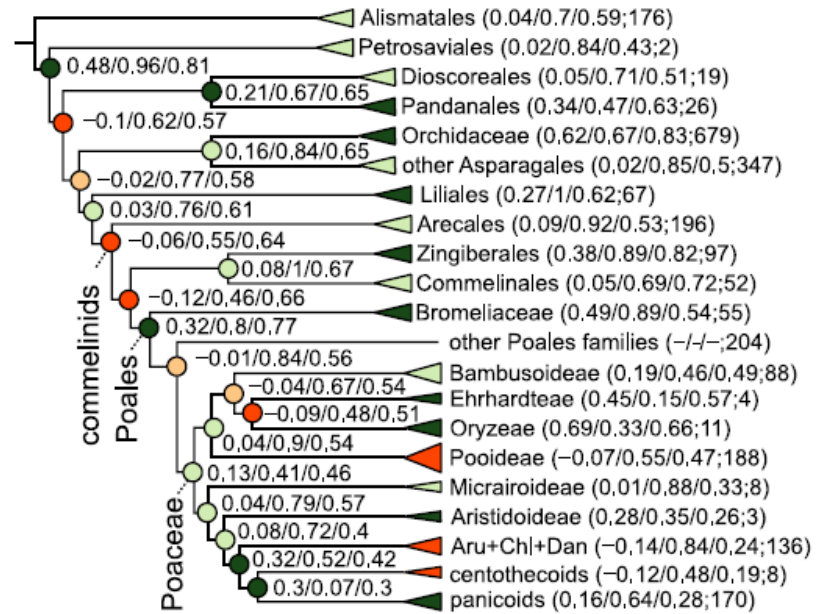| Example QS score (QC/QD/QI) | Interpretation |
|---|---|
| 1.0/–/1.0 | Full support: All sampled quartet replicates support the focal branch (QC = 1) with all trees informative when likelihood cutoffs are used (QI = 1). |
| 0.5/0.98/0.97 | Strong support: A strong majority of quartets support the focal branch (QC = 0.5), and the low skew in discordant frequencies (QD ≈ 1) indicate no alternative history is favored. |
| 0.7/0.1/0.97 | Strong support with discordant skew: A strong majority of quartets support the focal branch (QC = 0.7), but the skew in discordance (QD = 0.1) indicates the possible presence of a supported secondary evolutionary history. |
| 0.05/0.96/0.97 | Weak support: Only a weak majority of quartets support the focal branch (QC = 0.05), and the frequency of all three possible topologies is similar (QD ≈ 1). |
| 0.1/0.1/0.97 | Weak support with discordant skew: Only a weak majority of quartets support the focal branch (QC = 0.1), and the skew in discordance (QD = 0.1) indicates the possible presence of a supported secondary evolutionary history. |
| −0.5/0.1/0.93 | Counter-support: A strong majority of quartets support one of the alternative discordant quartet arrangement history (QC < 0; QD expected to be low). |
| 1/0.97/0.05 | Poorly informed: Despite supportive QC/QD values, only 5% of quartets passed the likelihood cutoff (QI = 0.05), likely indicating few informative sites. |
| 0.0/0.0/1.0 | Perfectly conflicted: The (unlikely) case where the frequencies of all three possible trees are equal and all trees are informative, which indicates a rapid radiation or highly complex conflict. |

*Notes:* QC = Quartet Concordance; QD = Quartet Differential; QI = Quartet Informativeness.

Pease et al. (2018): Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. *American Journal of Botany* 105(3): 385–403.

# Quartet Sampling – land plants



Pease et al. (2018): Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. *American Journal of Botany* 105(3): 385–403.
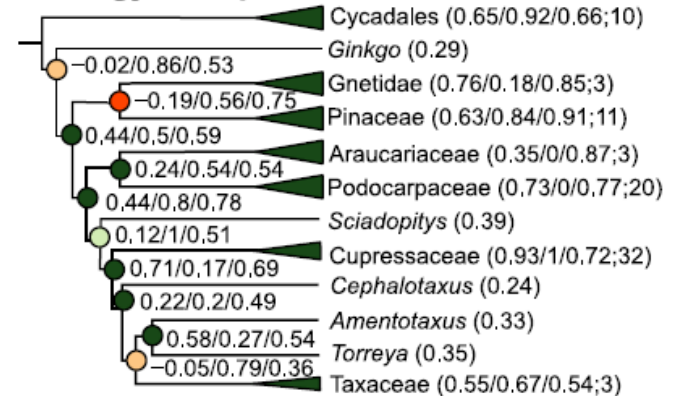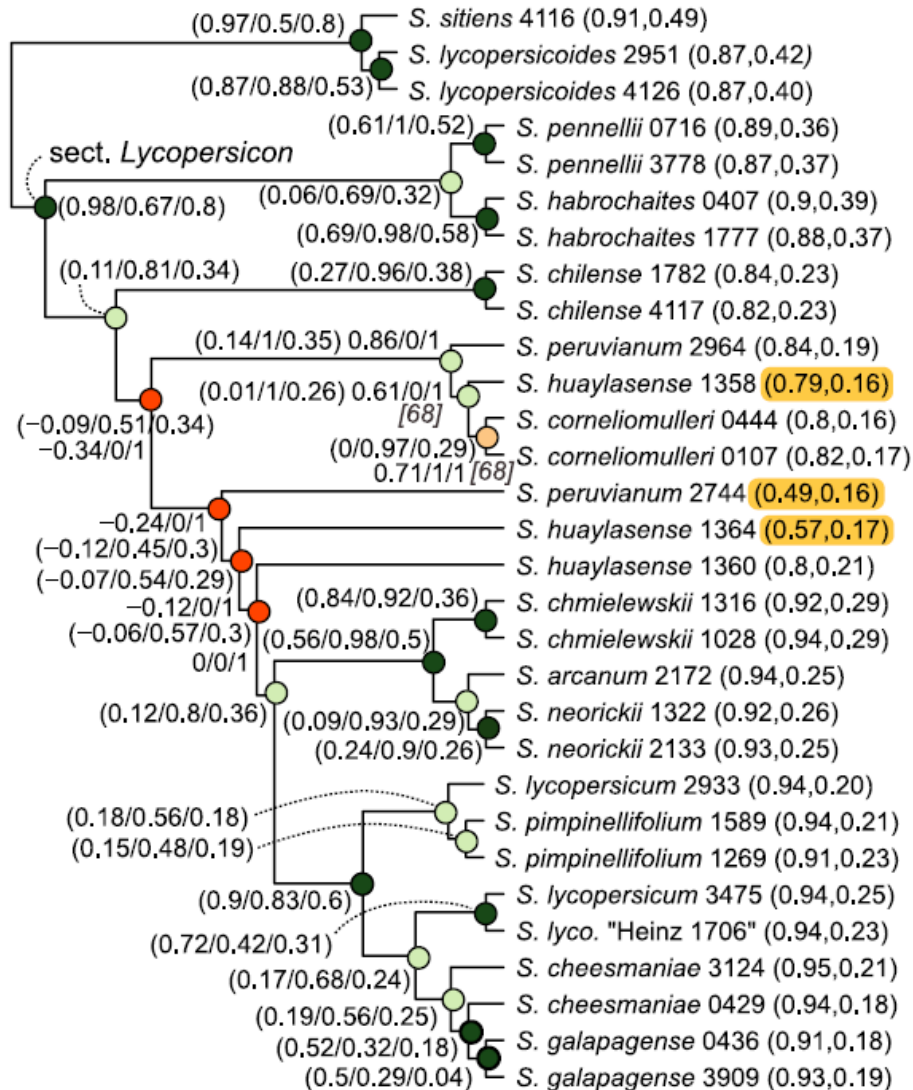
# Quartet Sampling – generic level



A. *Solanum* sect. *Lycopersicon*

Pease et al. (2016)

# References

Degnan, J.H., Rosenberg, N.A., 2009. Gene tree discordance, phylogenetic inference and the multispecies coalescent. Trends Ecol. Evol. 24, 332–340.

Drummond, A.J. & Bouckaert, R.R., 2015. Bayesian Evolutionary Analysis with BEAST. Cambridge University Press.

Mirarab, S., Reaz, R., Bayzid, M.S., Zimmermann, T., S. Swenson, M., Warnow, T., 2014. ASTRAL: Genome-scale coalescent-based species tree estimation. Bioinformatics 30, i541–i548

Sayyari, E., Mirarab, S., 2016. Fast Coalescent-Based Computation of Local Branch Support from Quartet Frequencies. Mol. Biol. Evol. 33, 1654–68.

ASTRAL presentation: http://tandy.cs.illinois.edu/astral-apro.pdf

Nguyen, N., Mirarab, S., Warnow, T., 2012. MRL and SuperFine+MRL: New supertree methods. Algorithms Mol. Biol. 7, 3.

Pease, J.B. et al., 2018. Quartet Sampling distinguishes lack of support from conflicting support in the green plant tree of life. American Journal of Botany 105, 385–403.