



DBI - 1661029

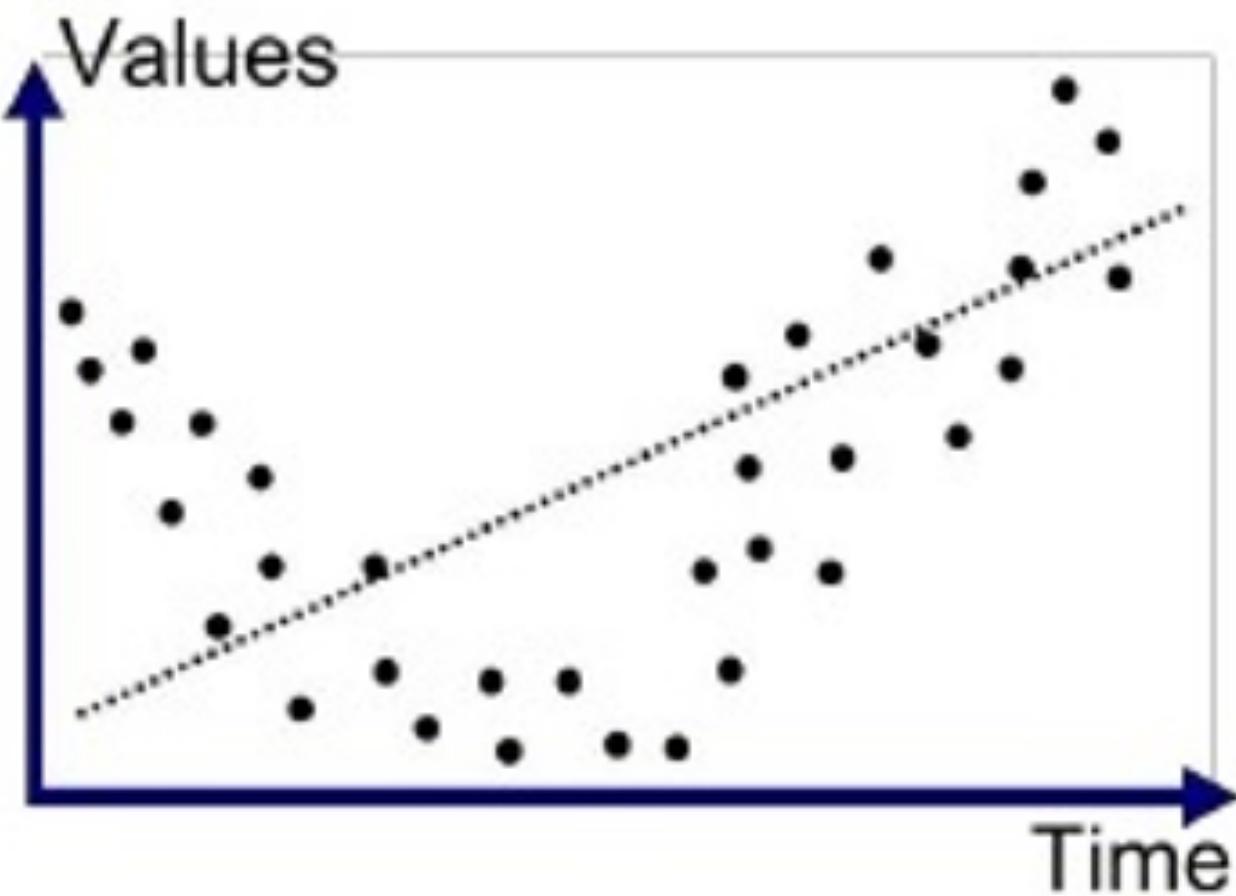
Posterior Predictive Checks in Bayesian Phylogenetics

R package: P2C2M.SNAPP

Drew Duckett, Tara A Pelletier, Bryan C Carstens



Models matter!





Bayesian Inference of Phylogeny and Its Impact on Evolutionary Biology

John P. Huelsenbeck,^{1*} Fredrik Ronquist,² Rasmus Nielsen,³ Jonathan P. Bollback¹

As a discipline, phylogenetics is becoming transformed by a flood of molecular data. These data allow broad questions to be asked about the history of life, but also present difficult statistical and computational problems. Bayesian inference of phylogeny brings a new perspective to a number of outstanding issues in evolutionary biology, including the analysis of large phylogenetic trees and complex evolutionary models and the detection of the footprint of natural selection in DNA sequences.

ity of a tree (Fig. 1). Bayes's theorem

$$\Pr[\text{Tree} \mid \text{Data}] = \frac{\Pr[\text{Data} \mid \text{Tree}] \times \Pr[\text{Tree}]}{\Pr[\text{Data}]}$$

(where the vertical bar should be read as “given”) is used to combine the prior probability of a phylogeny ($\Pr[\text{Tree}]$) with the likelihood ($\Pr[\text{Data} \mid \text{Tree}]$) to produce a posterior probability distribution on trees ($\Pr[\text{Tree} \mid \text{Data}]$). The posterior probability of a tree



Bayesian Inference of Phylogeny and Its Impact on Evolutionary Biology

John P. Huelsenbeck,^{1*} Fredrik Ronquist,² Rasmus Nielsen,³ Jonathan P. Bollback¹

As a discipline, phylogenetics is becoming transformed by a flood of molecular data. These data allow broad questions to be asked about the history of life, but also present difficult statistical and computational problems. Bayesian inference of phylogeny brings a new perspective to a number of outstanding issues in evolutionary biology, including the analysis of large phylogenetic trees and complex evolutionary models and the detection of the footprint of natural selection in DNA sequences.

ity of a tree (Fig. 1). Bayes's theorem

$$\Pr[\text{Tree} \mid \text{Data}] = \frac{\Pr[\text{Data} \mid \text{Tree}] \times \Pr[\text{Tree}]}{\Pr[\text{Data}]}$$

(where the vertical bar should be read as “given”) is used to combine the prior probability of a phylogeny ($\Pr[\text{Tree}]$) with the likelihood ($\Pr[\text{Data} \mid \text{Tree}]$) to produce a posterior probability distribution on trees ($\Pr[\text{Tree} \mid \text{Data}]$). The posterior probability of a tree

“The posterior probability of a tree can be interpreted as the probability that the tree is correct”, given the model.

Multi-species coalescent model (MSCM)

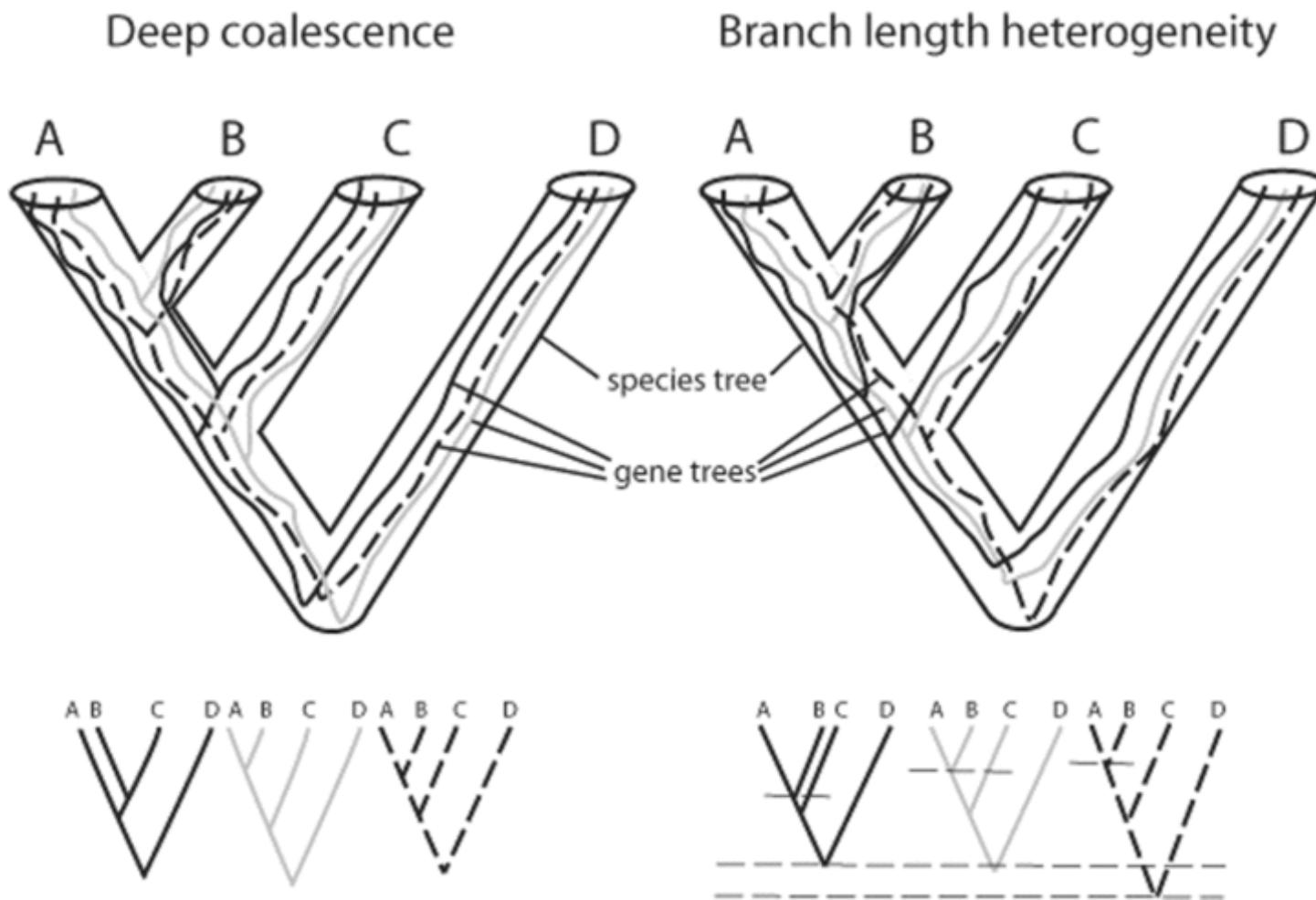
- developed to estimate species trees, while accounting for the coalescent process that can lead to incongruence among gene trees
(due to **incomplete lineage sorting**)

Multi-species coalescent model (MSCM)

Deep coalescence

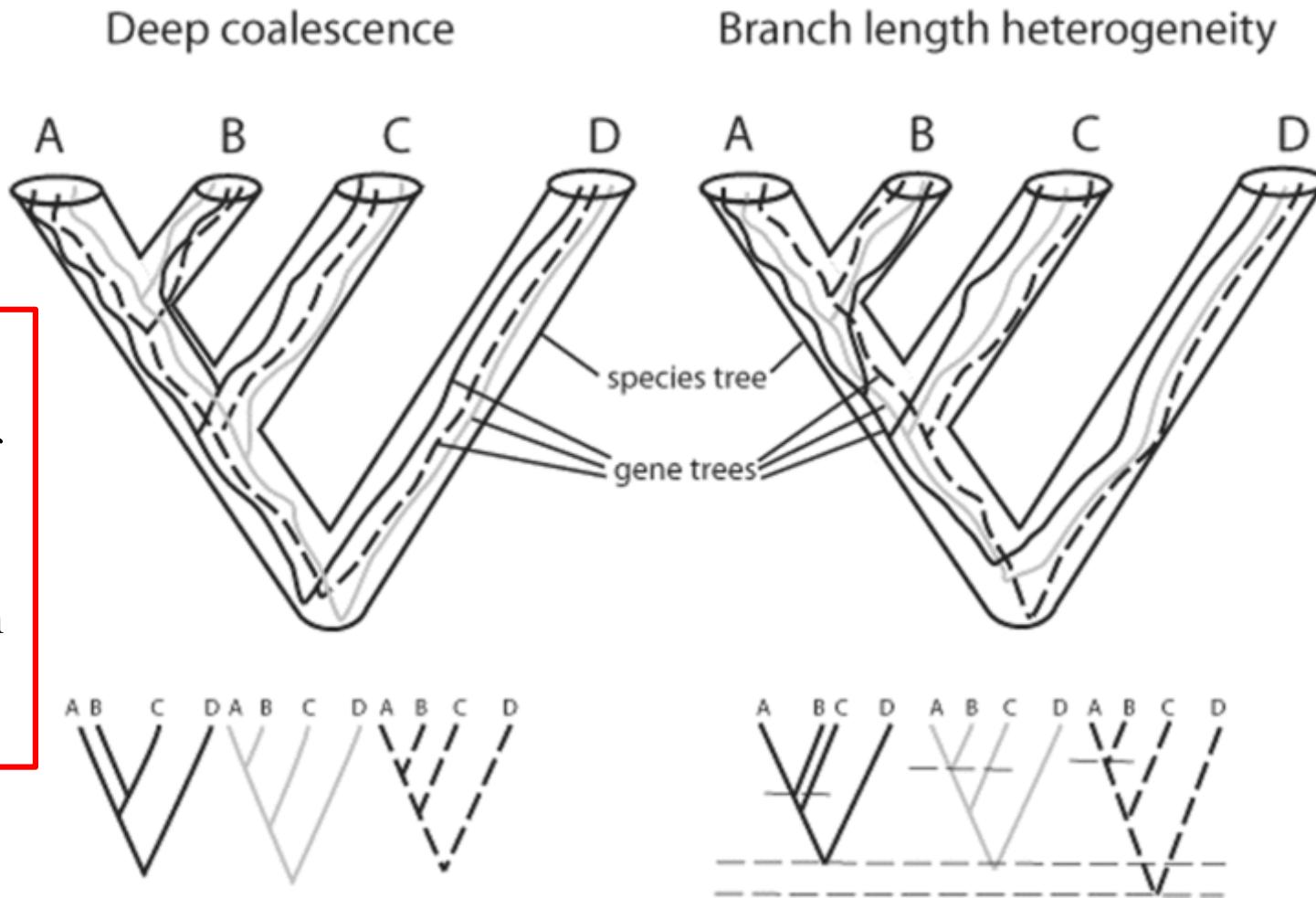
Branch length heterogeneity

Multi-species coalescent model (MSCM)

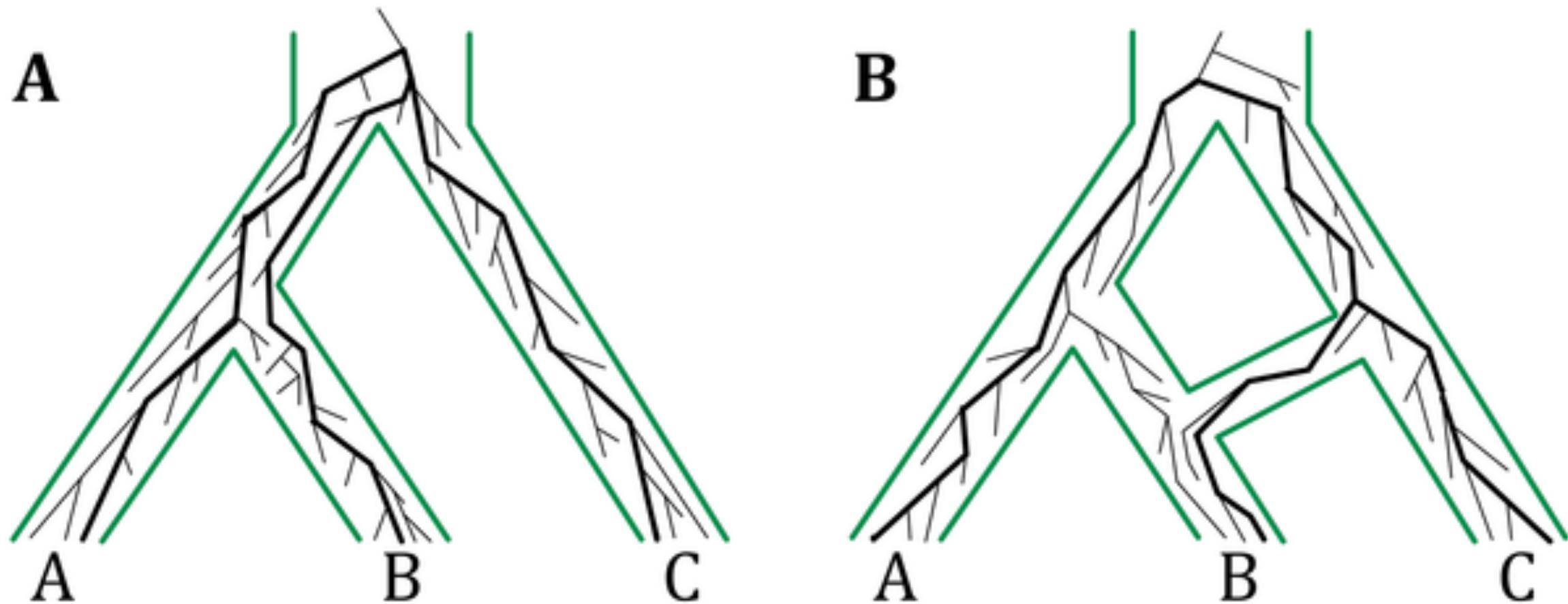


Multi-species coalescent model (MSCM)

Given a species tree, we can calculate the probability distribution of gene trees, then use this information to estimate the best species tree given our data under the coalescent model.



Multi-species coalescent model (MSCM)



Syst. Biol. 63(1):17–30, 2014

© The Author(s) 2013. 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/syt049

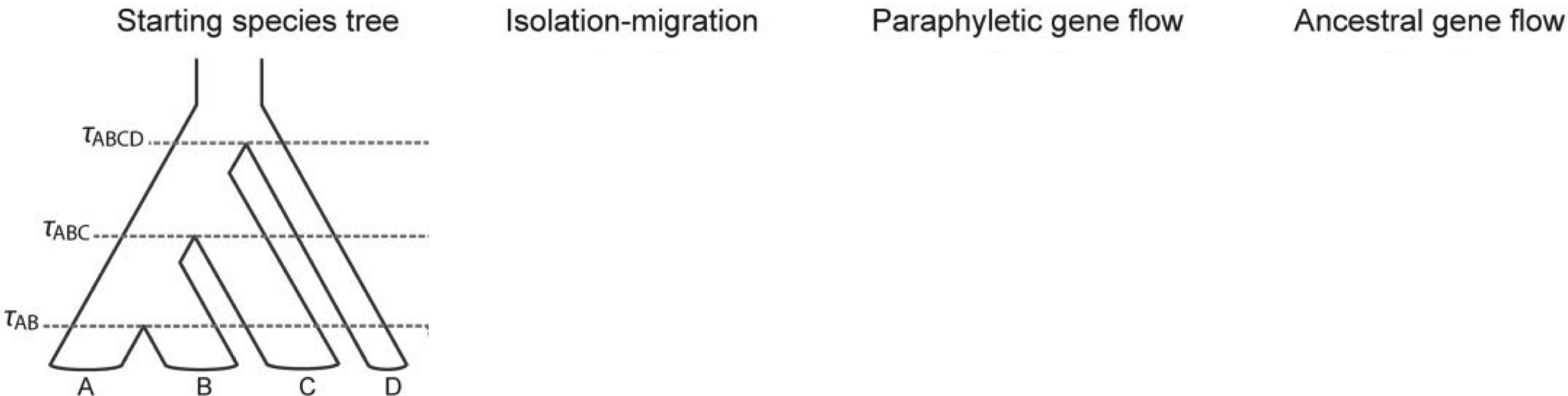
Advance Access publication August 13, 2013

The Influence of Gene Flow on Species Tree Estimation: A Simulation Study

ADAM D. LEACHÉ^{1,*}, REBECCA B. HARRIS¹, BRUCE RANNALA^{2,3}, AND ZIHENG YANG^{3,4}

The Influence of Gene Flow on Species Tree Estimation: A Simulation Study

ADAM D. LEACHÉ^{1,*}, REBECCA B. HARRIS¹, BRUCE RANNALA^{2,3}, AND ZIHENG YANG^{3,4}



The Influence of Gene Flow on Species Tree Estimation: A Simulation Study

ADAM D. LEACHE^{1,*}, REBECCA B. HARRIS¹, BRUCE RANNALA^{2,3}, AND ZIHENG YANG^{3,4}

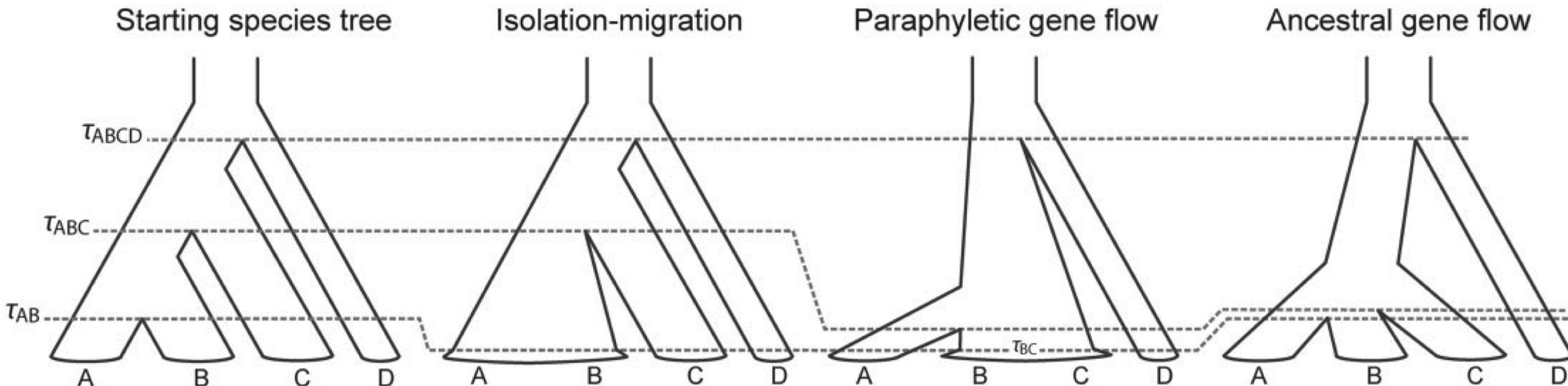


Figure 8 Species tree distortions caused by gene flow that can result from coalescent methods that only model ILS. Dashed lines illustrate species tree compression, and the widening of branches illustrates species tree dilation in relation to the starting species tree.

Syst. Biol. 63(3):322–333, 2014

© The Author(s) 2013. 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/syt057

Advance Access publication August 28, 2013

Poor Fit to the Multispecies Coalescent is Widely Detectable in Empirical Data

NOAH M. REID^{1,*}, SARAH M. HIRD¹, JEREMY M. BROWN¹, TARA A. PELLETIER², JOHN D. MCVAY¹, JORDAN D. SATLER²,
AND BRYAN C. CARSTENS²

Syst. Biol. 67(2):269–284, 2018

© The Author(s) 2017. 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/syx073

Advance Access publication September 4, 2017

Impact of Model Violations on the Inference of Species Boundaries Under the Multispecies Coalescent

ANTHONY J. BARLEY^{1,*}, JEREMY M. BROWN², AND ROBERT C. THOMSON¹

Posterior predictive simulation (PPS)

- We can approximate the **posterior predictive distribution** of a model by **simulating new observations** from parameter values sampled from the posterior distribution of your Bayesian phylogenetic analysis.

Posterior predictive simulation (PPS)

- We can approximate the **posterior predictive distribution** of a model by **simulating new observations** from parameter values sampled from the posterior distribution of your Bayesian phylogenetic analysis.
- If an evolutionary model is a good fit to your data (i.e., it does a good job of explaining patterns in the DNA), then data simulated under that model (**PPD**) should be similar to the **empirical data**.

Posterior predictive simulation (PPS)

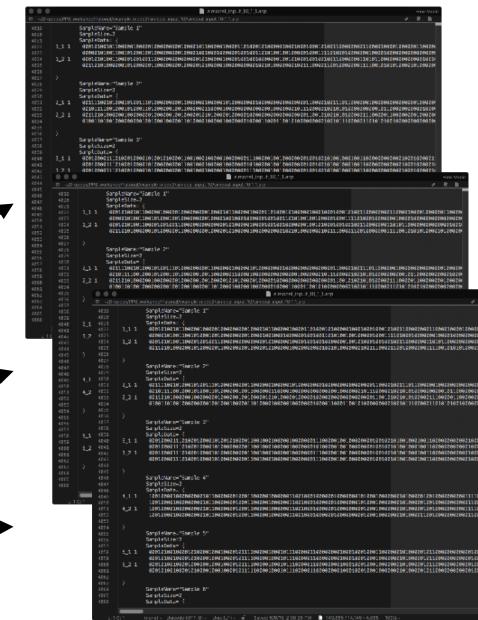
- We can approximate the **posterior predictive distribution** of a model by **simulating new observations** from parameter values sampled from the posterior distribution of your Bayesian phylogenetic analysis.
- If an evolutionary model is a good fit to your data (i.e., it does a good job of explaining patterns in the DNA), then data simulated under that model (**PPD**) should be similar to the **empirical data**.
- We can ask: does a particular model adequately describe an individual empirical data set?

Posterior predictive simulation (PPS)

1) Sample tree from the posterior distribution

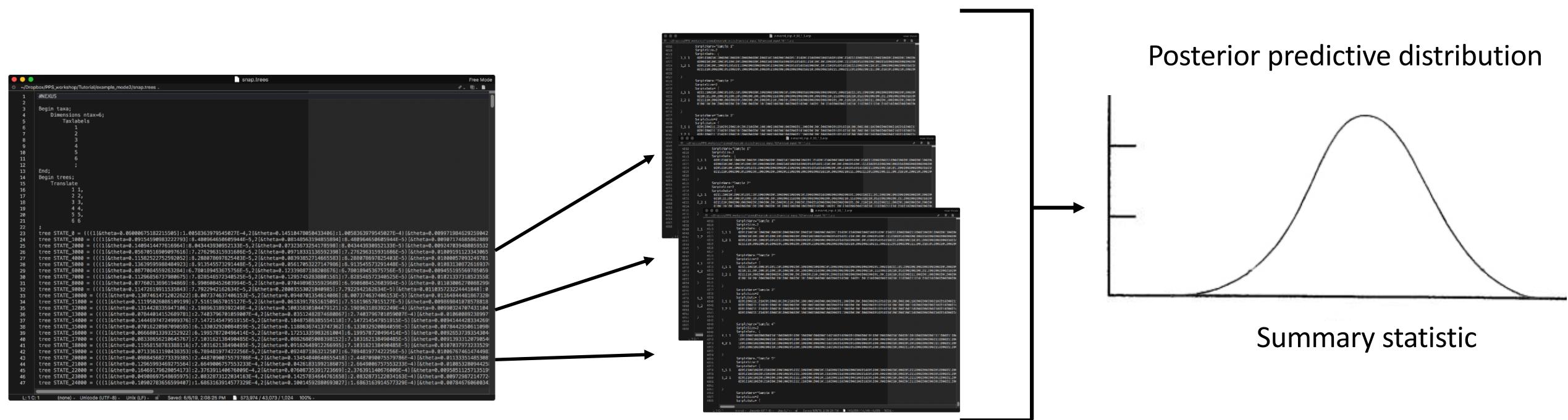
Posterior predictive simulation (PPS)

- 1) Sample tree from the posterior distribution
 - 2) Simulate data using this tree/parameters under the MSCM



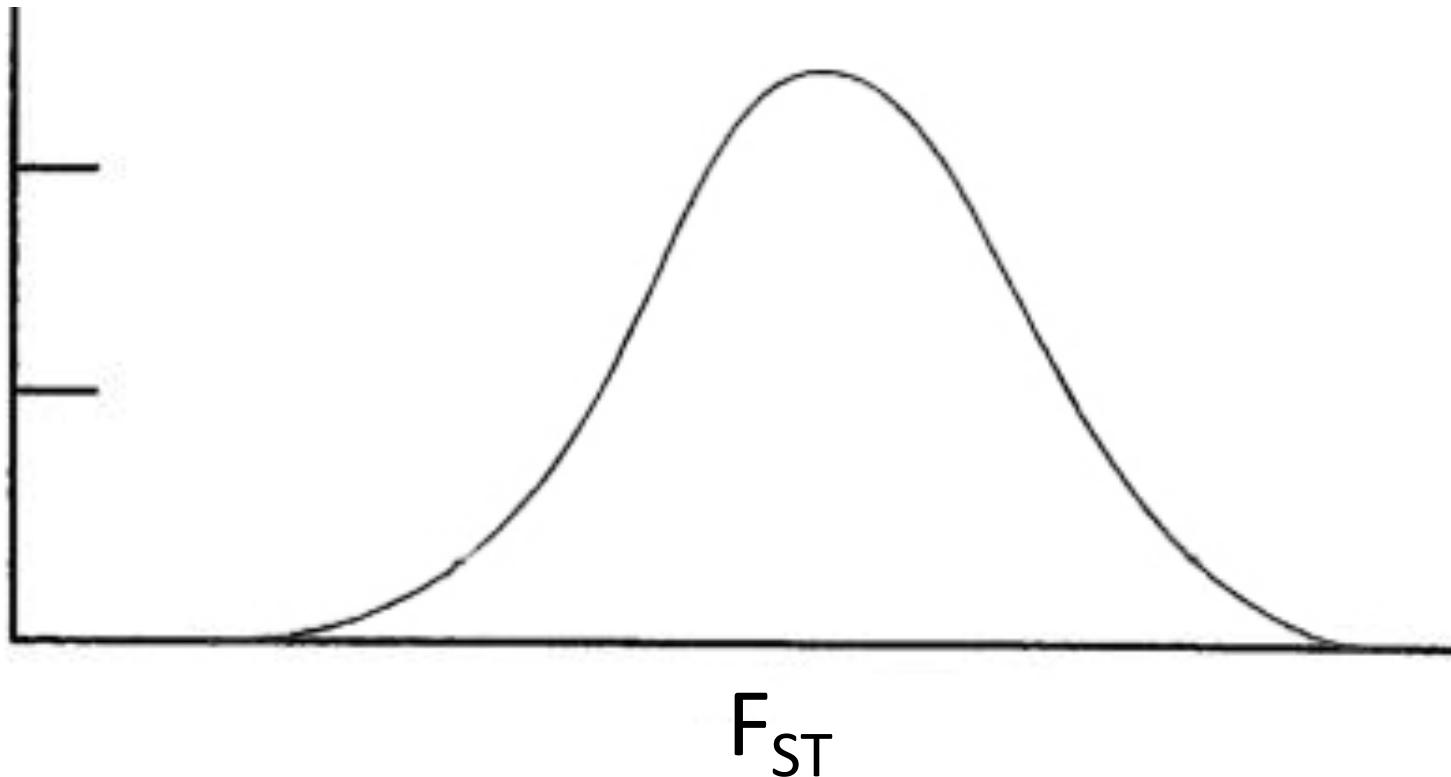
Posterior predictive simulation (PPS)

- 1) Sample tree from the posterior distribution
 - 2) Simulate data using this tree/parameters under the MSCM
 - 3) Summarize simulated datasets with test statistic



Posterior predictive distribution (PPD)

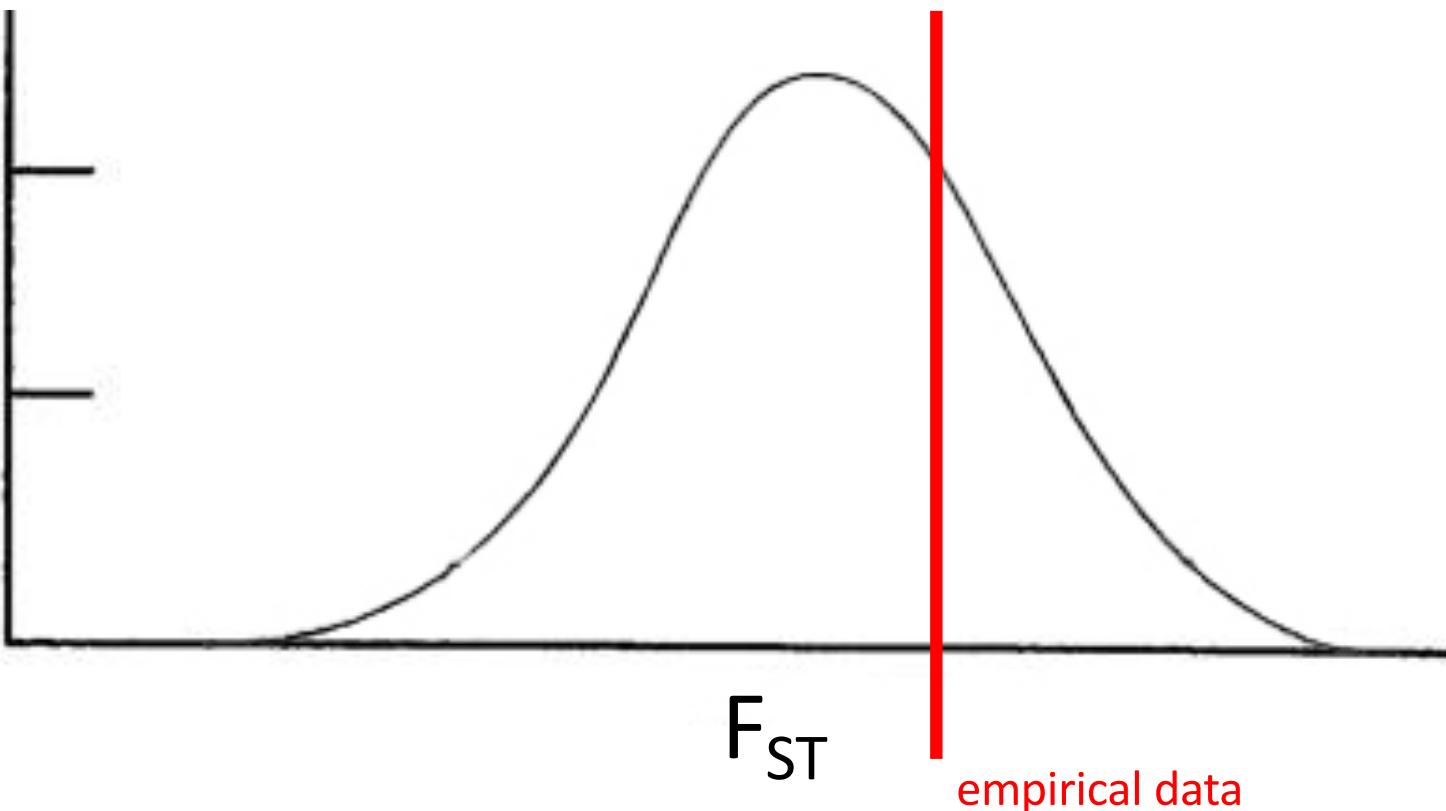
- How well does the empirical data fit this distribution?



This distribution is a representation of your data if the model were true – if the MSC were in fact generating your empirical data

Posterior predictive distribution (PPD)

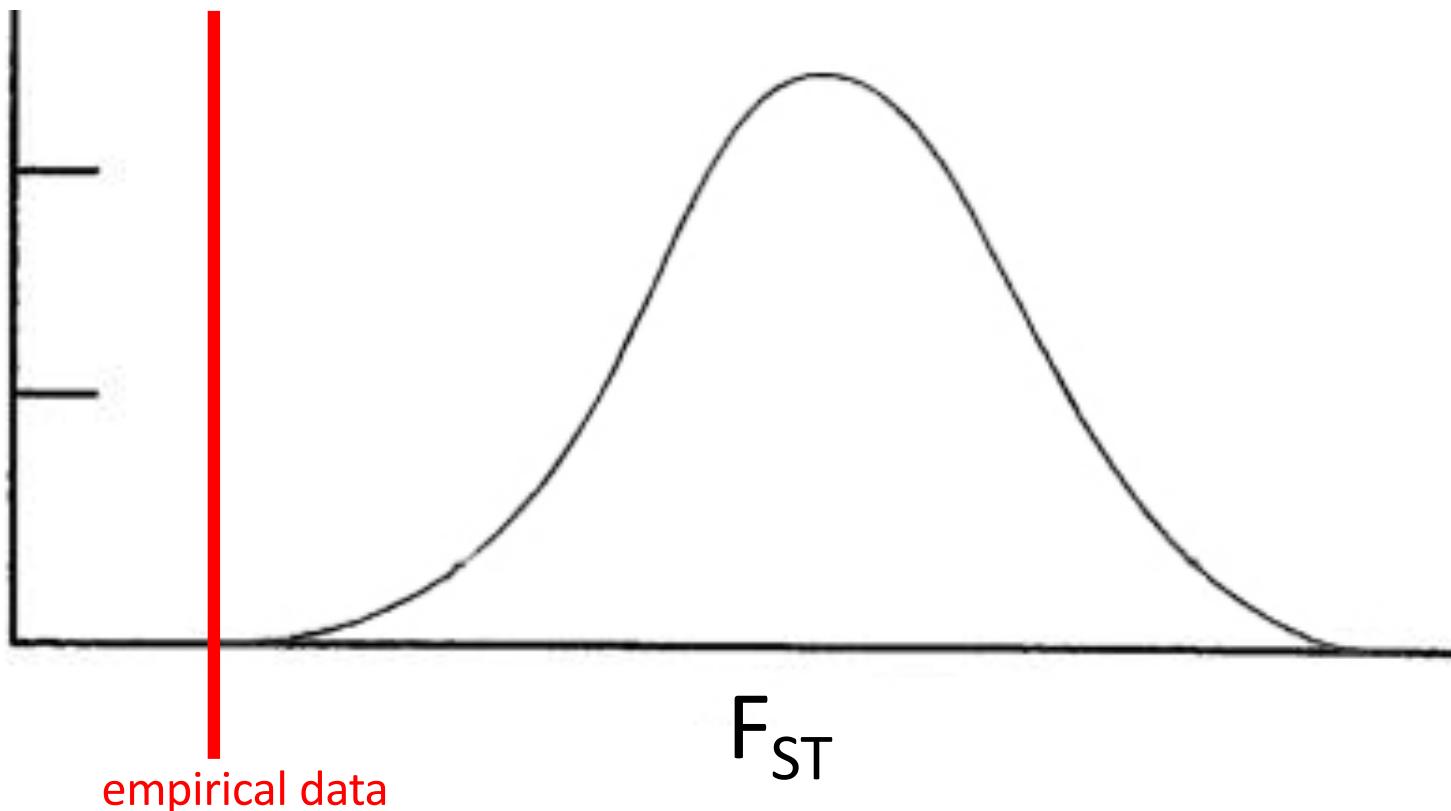
- How well does the empirical data fit this distribution?



This distribution is a representation of your data if the model were true – if the MSC were in fact generating your empirical data

Posterior predictive distribution (PPD)

- How well does the empirical data fit this distribution?



This distribution is a representation of your data if the model were true – if the MSC were in fact generating your empirical data

Posterior predictive checks of coalescent models: P2C2M, an R package

MICHAEL GRUENSTAEUDL,*‡ NOAH M. REID,† GREGORY L. WHEELER* and BRYAN C. CARSTENS*

**Department of Evolution, Ecology & Organismal Biology, Ohio State University, Columbus, OH 43210, USA* †*Department of Environmental Toxicology, University of California, Davis, CA 95616, USA*

*BEAST – multi-locus sequence data

SNAPP (SNP and AFLP Phylogenies)

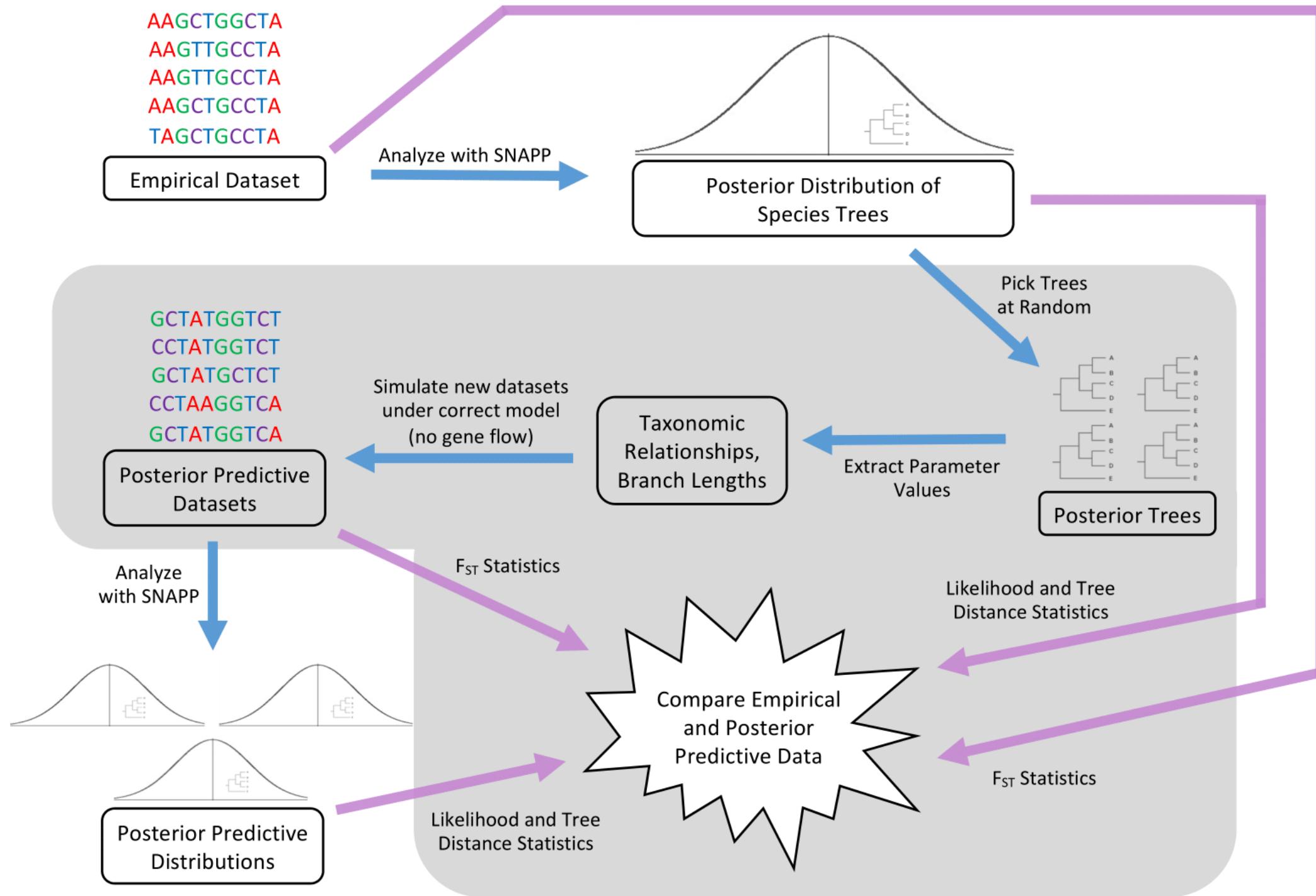
- Interfaces with the BEAST package
- Algorithm bypasses the gene trees and computes species tree likelihoods directly from the markers
- Returns a sample of species trees with (relative) divergence times and population sizes (posterior distribution)

SNAPP assumptions

- Those of the coalescent process (shared polymorphism is due to ILS)
- Each marker is a single biallelic character
- The genealogies for separate markers are conditionally independent (satisfied for SNPs that are well spaced along the genome)

R package: P2C2M.SNAPP

- Conducts posterior predictive checks on your analysis from the program SNAPP
- We are about to submit both this paper and the package to CRAN



P2C2M.SNAPP simulation testing

- Which summary statistics work best?
- We simulated two trees 100X: **MSCM** and **MSCM+*m***
 - 6 species (symmetrical tree)
 - 2 individuals per species
 - 2000 SNPs
 - $N_e = 100,000$
 - Speciation times at 5N, 10N, and 20N
 - When gene flow was incorporated into the model it happened at 2.5N generations in the past and m was drawn from a uniform prior between 0.5 and 5 migrants per generation

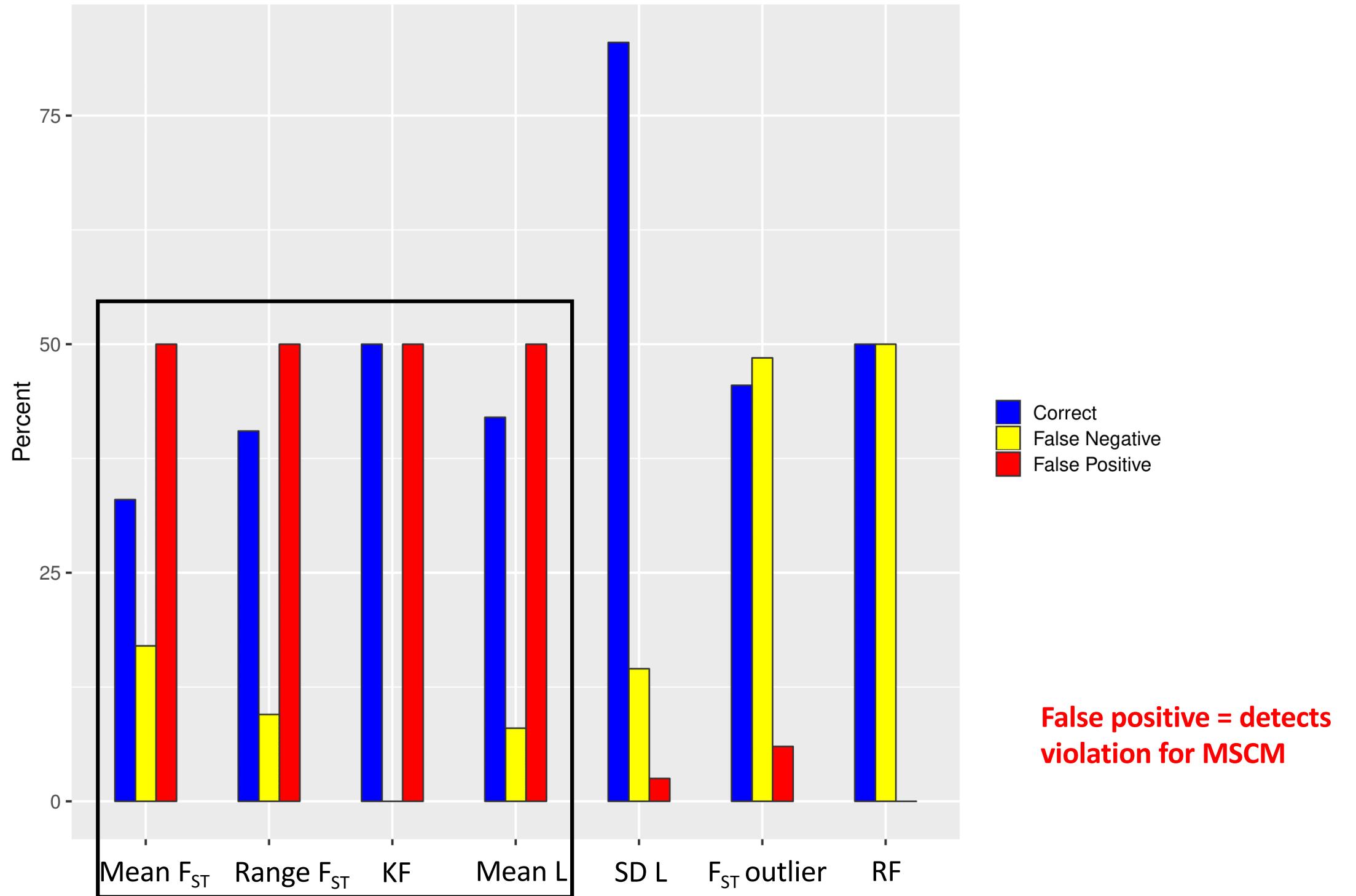
P2C2M.SNAPP simulation testing

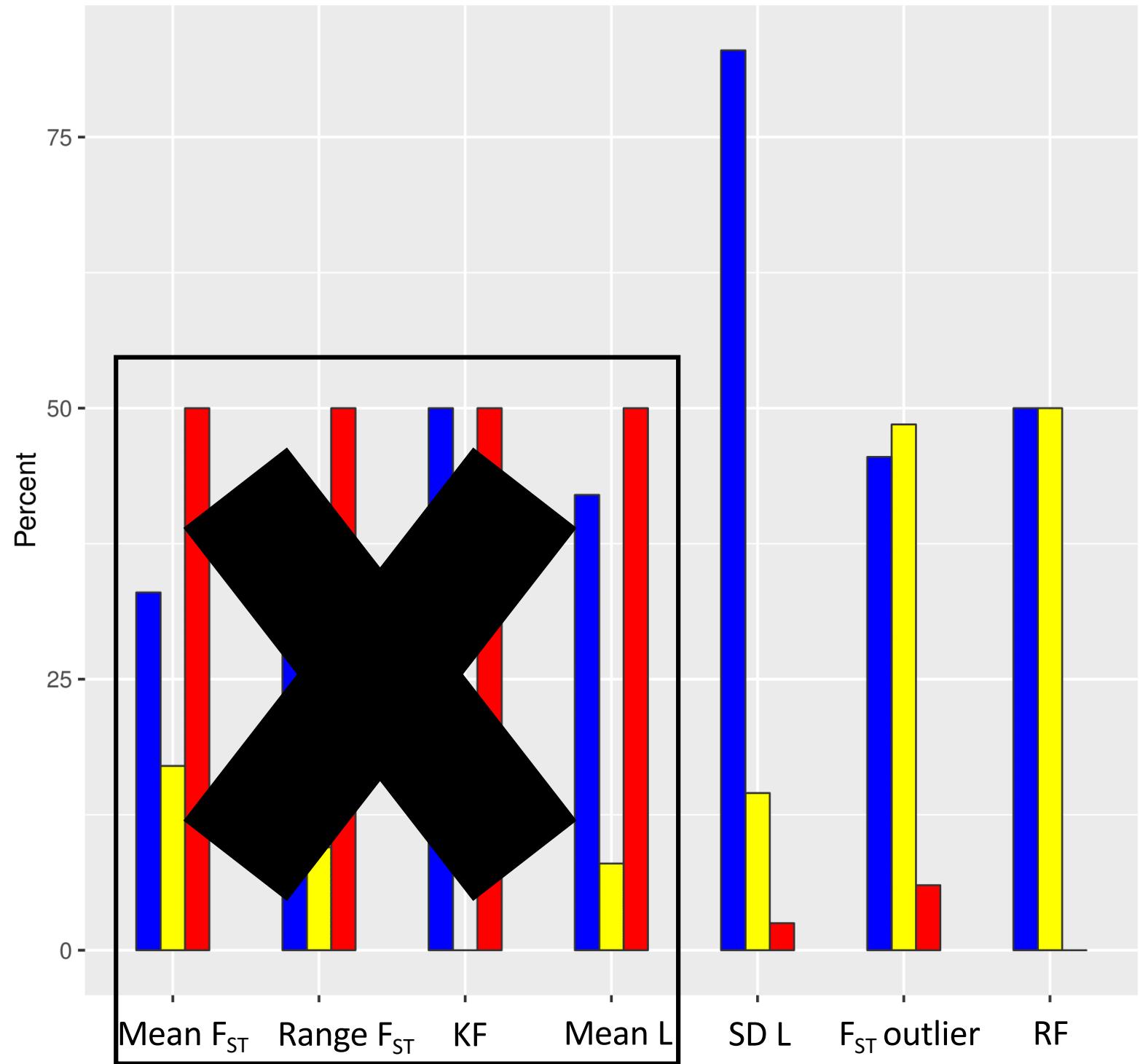
- **Summary statistics tested:**

- Mean pairwise F_{ST}
- Range of pairwise F_{ST}
- F_{ST} outlier test
- Robinson-Foulds distance (topological distance only)
- Kuhner-Felsenstein distance (includes branch lengths)
- Mean of tree likelihood
- Standard deviation of tree likelihood

P2C2M.SNAPP simulation testing

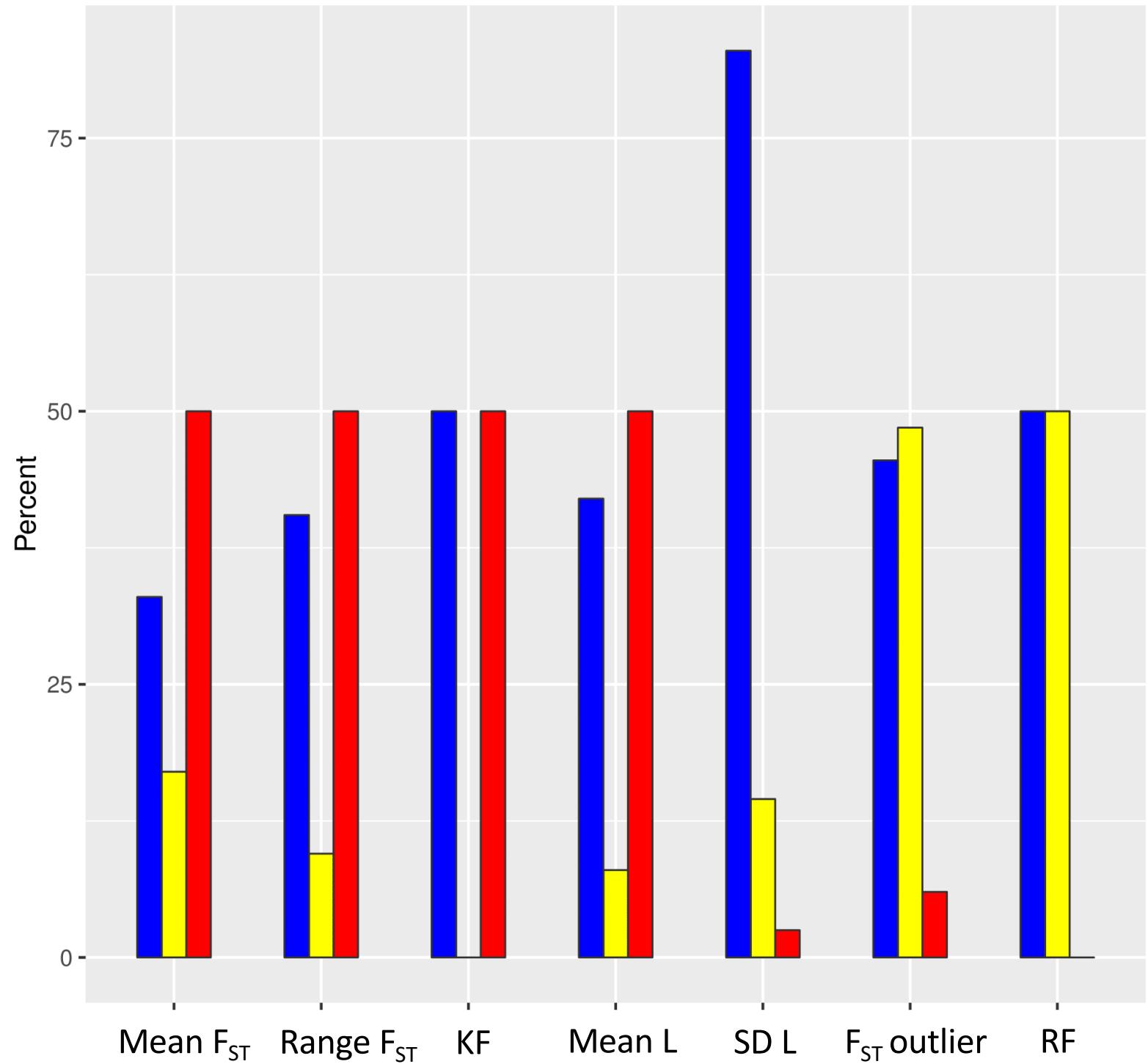
- **Summary statistics tested:**
 - Mean pairwise F_{ST}
 - Range of pairwise F_{ST}
 - F_{ST} outlier test
 - Robinson-Foulds distance (**topological distance only**)
 - Kuhner-Felsenstein distance (includes branch lengths)
 - Mean of tree likelihood
 - Standard deviation of tree likelihood





Correct
False Negative
False Positive

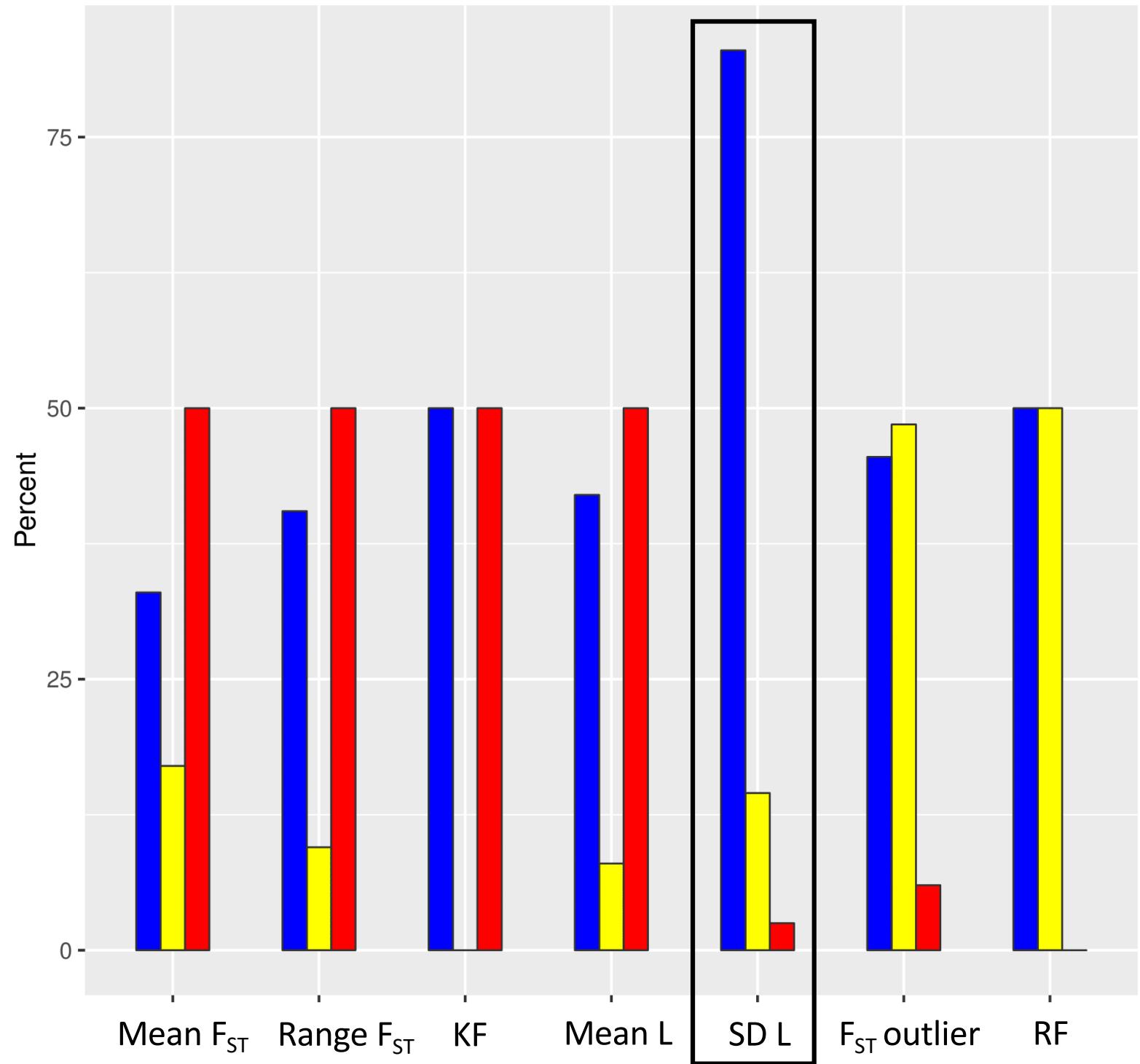
False positive = detects violation for MSCM



Correct
False Negative
False Positive

False negative = does not identify violation for MSCM+ m

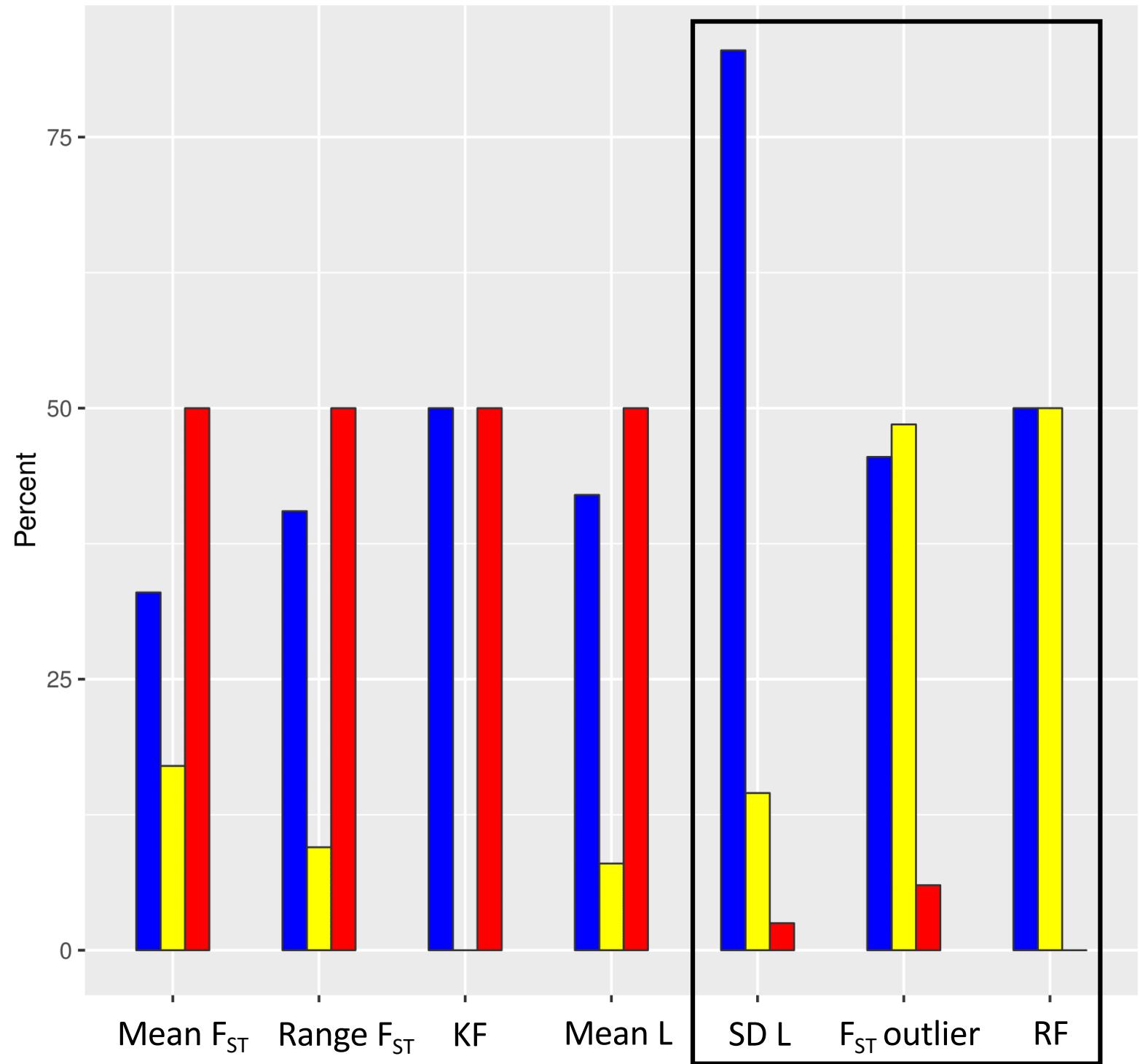
False positive = detects violation for MSCM



Correct
False Negative
False Positive

False negative = does not identify violation for MSCM+ m

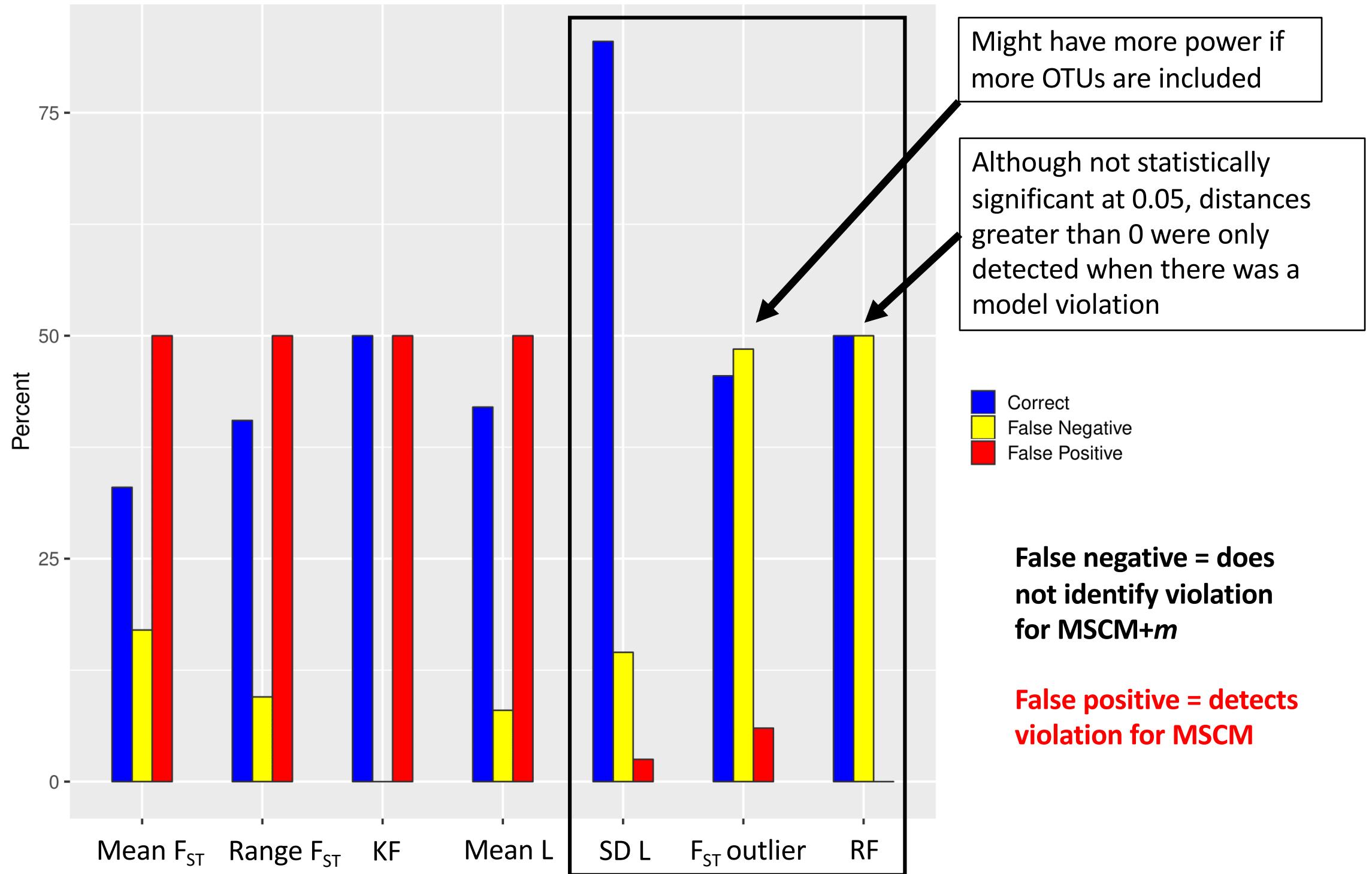
False positive = detects violation for MSCM



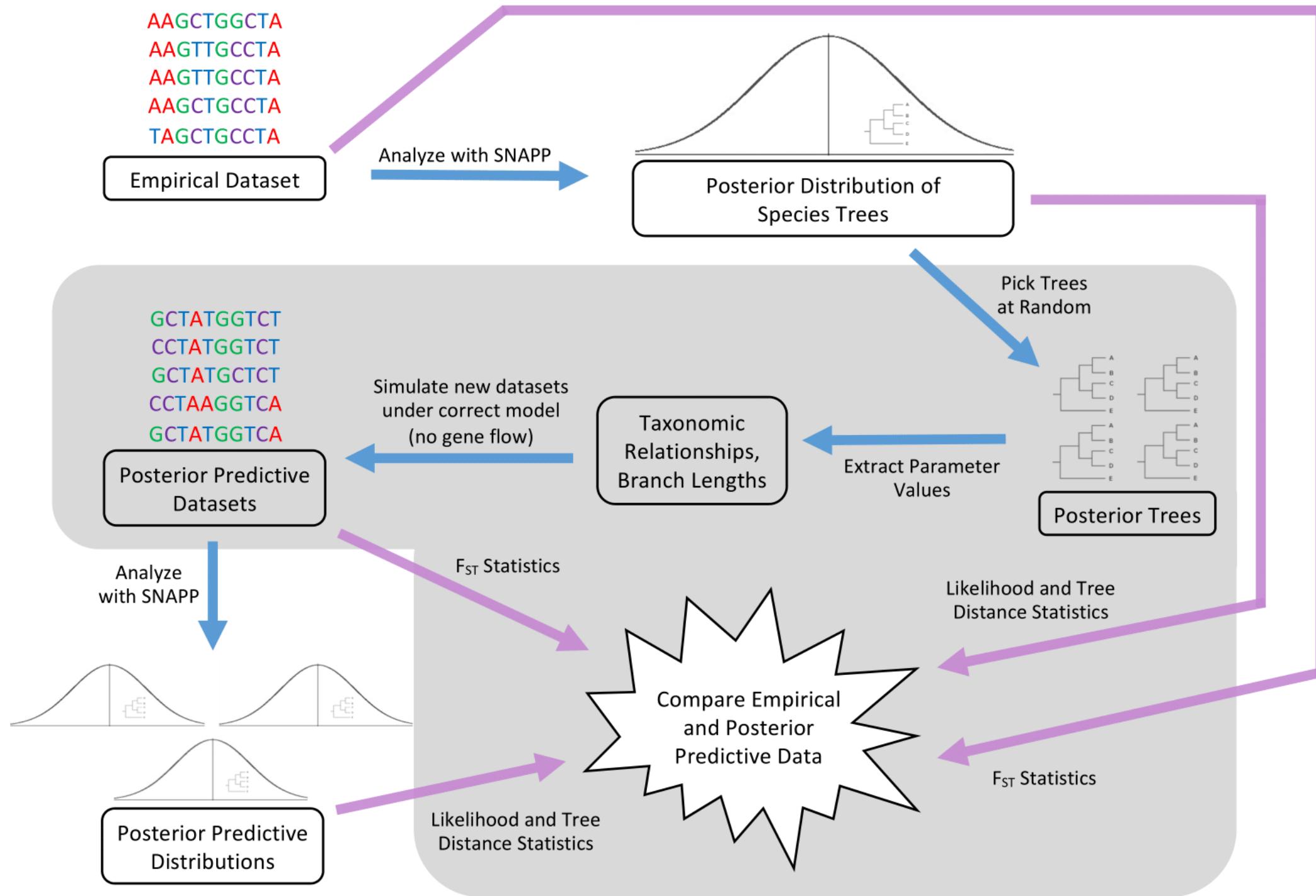
Correct
False Negative
False Positive

False negative = does not identify violation for MSCM+ m

False positive = detects violation for MSCM



- **Data-based test statistics:** assess whether the observed and posterior predictive data sets exhibit similar characteristics (estimated directly from the data)
 - F_{ST} outliers
- **Inference-based test statistics:** where inferences drawn from the observed posterior distribution and posterior predictive distributions are compared
 - SD likelihood, RF



What if you detect violations?

- acknowledge the model violation and the effects it could have on your phylogeny estimate (we used a cut-off of 0.05 for simulation testing but you should consider the p-value)
- conduct additional analyses to examine the cause of the model violation, as such violations indicate interesting evolutionary processes not accounted for by the MSCM model
 - PhyloNet (Wen et al. 2018): MSNC

Running P2C2M.SNAPP

- First make sure your SNAPP analysis has converged (next presentation)! Otherwise you get false positives.
- Getting started (see tutorial)
 - Install P2C2M.SNAPP
 - Install Fastsimcoal
 - Put SNAPP output and fastsimcoal executable in working directory