

## **Module 3: Phylogenomics**

- Phylogenetic inference using maximum-likelihood and Bayesian analyses
- Phylogenomics and challenges of whole genome inferences

Patrícia H. Brito



## Who am I?



Patrícia H. Brito, Ph.D. **Computational Biologist** Yeast Genomics LAB UCIBIO, NOVA Science and Technology, UNL https://yeastgenomicslab.wordpress.com phbrito@fct.unl.pt

Microbial genome evolution **Evolutionary genomics** Phylogenomics **Population genomics Molecular Evolution Evolutionary Medicine** Silva et al 2022 Molecular Ecology 3 subtilis str. KCTC 1342 subtilis str. gtP20t allismortis str. DV/1-E-3 B. subtilis str. TU-B-10 B. subtilis str. DV1-100 KCTC 13622 100 100 Plant btilis str. Miyad Food htilis str BSP Gut Marine Other Lab ubtilis str. GXA-2 ubtilis str OH-B subtilis str AUSI98 **Bacterial phylogenomics** subtilis str MB73 **Protein family evolution** subtilis str. JH642 subtilis str. JH642 substr. AG1 subtilis str. SMY ubtilis str. QB928 0.02 subtilis str. PY79 Brito et al 2018 GBE



Torcato et al 2019 JBC

WILEY T delhaveekä DE Torulaspora sp. I

T. pretoriensis - PF

T. nypae - NYP T quercuum - QUE

Torulaspora sp. V

3 CC 8929 IF 3987 CC 8930 CC 5321 CC 5321 CC 5323

Yeaflow Alph YCC 7193 YCC 8416 ATF 4303 YCC 8819

AI-2 receptors in Clostridic Firmicutes

Yeast phylogenomics

## Phylogenetic inference – what it is and why do we do it?

"The affinities of all the beings of the same class have sometimes been represented by a great tree. I believe this simile largely speaks the truth."

Charles Darwin, Origin of Species, 1859

#### **Evolution is descent with modification from a common ancestor**



I think





## **Phylogenetic inference – terminology**





root

## **Phylogenetic inference – typical pipeline**

#### Question:

we want to know how a particular group of species/organisms/cells relate to each other

#### Typical pipeline in Molecular Phylogenetics

- 1. Choose the molecular marker (genomic region or specific data type)
- 2. Get the sequences of that molecular marker for all terminals in the tree
- 3. Choose an **optimality criterion** and an **algorithm** to estimate the gene tree

1 gene  $\rightarrow$  some genes  $\rightarrow$  many genes  $\rightarrow$  genomes

More genes  $\rightarrow$  more resolution at different levels of the tree  $\rightarrow$  higher support

In most analyses the implicit assumption is that all genes do in fact have the same gene tree, that these gene trees are congruent with and converge on the species tree

#### Sequence alignment

CACCTGTCGT			TCTGGTGCAG
CAGCTGTCGT	GCTCTTTCTG	TTGAGCCTGG	TCTGGTGCAG
CAGCTGCCGT	GTTTTCTCTG	TTGAGCCTGG	TCTGGTACAG
CAGCTGCCGC	GTTCTCTCCG		TCTGGTGCAA
CTCCTGCCGG	GTGCTCTCAG		
CTCCTGCCGG		CTGAGCCGGG	TCTGGTGCAG
CTCTTGCCGG		CTGAGCCTTG	





not always true!

## **Phylogenetic inference – methods overview**





## Maximum Likelihood



 $\rightarrow$  The same observed result leads to very different likelihoods depending on the assumed model (hypothesis)

 $\rightarrow$  The model that assumes "non-fair dice" is the one that gives higher probability of observing that result, and as such it is the most likely model.



## Maximum Likelihood of a Tree (topology & branch lengths)





Two assumptions:

- 1. Evolution in different sites is independent
- 2. Evolution in different lineages is independent



Designation	Rate params	Base frequencies	of free params
JC	a=b=c=d=e=f	$\pi_A = \pi_C = \pi_G = \pi_T$	1
K80, K2P	a=c=d=f, b=e	$\pi_A = \pi_C = \pi_G = \pi_T$	2
TrNef	a=c=d=f, b, e	$\pi_A = \pi_C = \pi_G = \pi_T$	3
K81, K3ST	a=f, b=e, c=d	$\pi_A = \pi_C = \pi_G = \pi_T$	3
TvMef	a, c, d, f, b=e	$\pi_A = \pi_C = \pi_G = \pi_T$	5
TiMef	a=f, c=d, b, e	$\pi_A = \pi_C = \pi_G = \pi_T$	4
SYM	a, b, c, d, e, f	$\pi_A = \pi_C = \pi_G = \pi_T$	6
F81	a=b=c=d=e	$\pi_A, \pi_C, \pi_G, \pi_T$	4
НКҮ	a=c=d=f, b=e	$\pi_A, \pi_C, \pi_G, \pi_T$	5
TrN	a=c=d=f, b, e	$\pi_A, \pi_C, \pi_G, \pi_T$	6
K81uf	a=f, b=e, c=d	$\pi_A, \pi_C, \pi_G, \pi_T$	6
TVM	a, c, d, f, b=e	$\pi_A, \pi_C, \pi_G, \pi_T$	8
TiM	a=f, c=d, b, e	$\pi_A, \pi_C, \pi_G, \pi_T$	7
GTR, REV	a, b, c, d, e, f	$\pi_A, \pi_C, \pi_G, \pi_T$	9

Number



## Maximum Likelihood of a Tree (topology & branch lengths)





## Phylogenetic inference is a NP-complete problem where exhaustive searches for datasets of 10+ terminals are practically impossible $\rightarrow$ Heuristic methods





Phylogenetic tree space

#### Subtree pruning and regrafting (SPR)





Tree bisection and reconnection (TBR)

- Branch swapping NNI < SPR < TBR</li>
- Multiple replicates with random starting points





**Bootstrap** resampling method used to estimate branch support on a phylogenetic tree. Provides an indication of the robustness (confidence) in each bipartition.



VOVA SCHOOL OF SCIENCE & TECHNOLOGY DEPARTMENT OF LIFE SCIENCES

### **Bayes' theorem**





## **Bayesian phylogenetic inference**



Rev. Thomas Bayes (1701-1761)

#### Given:

- $\tau$  = phylogenetic tree (topology + branch lengths)
- X = data (aligned molecular data)



#### **Prior probability**

- Uniform dist topology
- Exponential dist branch lengths
- Gamma dist rate variation
- Dirichlet dist allele frequency

Typically, impossible to estimate. But by using MCMC chains to sample the posterior distribution we do not need to estimate this quantity



## Markov Chain Monte Carlo sampling (MCMC) <---- Not a "hill-climbing" method!

MCMC approximates  $f(\tau_i | X)$  by sampling a high number of trees  $\tau_i$  from the posterior distribution. The trees with higher probabilities are the ones most likely to be sampled during the MCMC sampling process. Therefore, MCMC focuses most of the sampling effort on sampling the distribution of interest - **the proportion of time that MCMC method samples a give region of the parameter space is proportional to the posterior distribution of that region.** 

#### Strategy for running an MCMC:

- 1. Start at a random point
- 2. Make a small-scale change
- 3. Estimate the ratio (r) of the probabilities of the new and the original state:
  - If r > 1 -> accept change
  - If r < 1 -> accept change with probability r
- 4. Back to step 2

#### Metropolis-Hasting decision criteria:





#### MCMC ...





adapted from J. Felsenstein



## **The Trace Plot**





Assessing Convergence:

- 1. Check for the plateau in the trace plot
- 2. Look at sampling behavior within the run (autocorrelation times, effective sample size etc)
- 3. Compare independent runs with different, randomly chosen starting points



### **Summarizing sampled topologies**

 $f(\boldsymbol{X}|\tau_i)$ 

0.000

0.026

0.000

0.000

0.000

0.000

0.037

0.000

0.001

0.001

0.001

0.004

0.919

0.009

Results are summarized with
 credibility intervals and majority consensus trees







The posterior probability of a clade is simply the sum of the posterior probabilities of all trees that contain that clade.

#### Example:

- A credible 95% interval for these topologies includes trees 14 and 8 trees ->  $f(X|\tau_i)=0.956$
- The probability of the Human-Chimp clade is T13 + T14 + T15 = 0.932



List of sampled topologies

 $\tau_i$ (Gi,Hu,((Ch,Go),Or))

(Gi,(Hu,(Ch,Go)),Or)

(Gi,(Hu,Or),(Ch,Go))

(Gi,((Hu,Or),Go),Ch)

(Gi,((Hu,Or),Ch),Go)

(Gi,Hu,((Ch,Or),Go))

(Gi,(Hu,Go),(Ch,Or))

(Gi,((Hu,Go),Ch),Or)

(Gi,((Hu,Go),Or),Ch)

(Gi,(Hu,(Ch,Or)),Go)

(Gi,Hu,(Ch,(Go,Or)))

(Gi,(Hu,(Go,Or)),Ch)

(Gi,(Hu,Ch),(Go,Or))

(Gi,((Hu,Ch),Go),Or)

(Gi,((Hu,Ch),Or),Go)

5

9

10

11 12

13

14

15

A majority rule consensus tree is formed by combining all the clades with the highest posterior probability that are compatible



Mixing refers to how often proposed changes to parameters are accepted during the MCMC run. High acceptance rate means chain is making too small moves. Low acceptance rate means proposed changes are too large. Optimal acceptance rate: 20-60 percent.

The time it takes for a MCMC to obtain an adequate sample of the posterior depends on its mixing behavior





### Metropolis-coupled Markov chain Monte Carlo aka (MC)<sup>3</sup>







## **Inference of large phylogenies**







#### concatenation





## Inference of large phylogenies

- The implicit assumption is that all genes do in fact have the same gene tree, that these gene trees are congruent with and converge on the species tree
- The use of many genes eliminates **stochastic error** (e.g. insufficient sequence length) and **systematic error** (some gene trees my depart from model assumptions)
- If we add extra requirements such as single copy orthologs and core genes then we might also reduce/eliminate biological causes of incongruence between gene tree and species tree

# Things to be aware of... analysis of large concatenated datasets may lead to **misleading bootstrap support**







# Things to be aware of... sometimes all gene trees differ from each other and from the concatenation phylogeny!!

doi:10.1038/nature12130

#### Inferring ancient divergences requires genes with strong phylogenetic signals

Leonidas Salichos<sup>1</sup> & Antonis Rokas<sup>1</sup>

To tackle incongruence, the topological conflict between different gene trees, phylogenomic studies couple concatenation with practices such as rogue taxon removal or the use of slowly evolving genes. Phylogenomic analysis of 1,070 orthologues from 23 yeast genomes identified 1,070 distinct gene trees, which were all incongruent with the phylogeny inferred from concatenation. Incongruence severity increased for shorter internodes located deeper in the phylogeny. Notably, whereas most practices had little or negative impact on the yeast phylogeny, the use of genes or internodes with high average internode support significantly improved the robustness of inference. We obtained similar results in analyses of vertebrate and metazoan phylogenomic data sets. These results question the exclusive reliance on concatenation and associated practices, and argue that selecting genes with strong phylogenetic signals and demonstrating the absence of significant incongruence are essential for accurately reconstructing ancient divergences.



Salichos and Rokas Nature 2013

# Some recipes for handling incongruence in concatenation analysis:

- Remove all sites containing gaps
- Remove fast-evolving or unstable species
- Selecting genes that recover specific clades
- Selecting the most slow-evolving genes
- Selecting genes whose bootstrap consensus trees have high average support
- Multiple searches using distinct starting trees

## Strategy:

- Apply different phylogenetic methods (different optimality criteria/approaches)
- Assess conflict across gene trees
- Investigate alternative hypotheses for branches showing conflict/assess sensitivity of results



#### References

The Phylogenetic Handbook, A practical approach to phylogenetic analysis and hypothesis testing. Ed, Philippe Lemey, Marco Salemi, Anne-Mieke Vandamme. 2009 Inferring Phylogenies, Joe Felsenstein 2004

#### Software

Iqtree <u>http://www.iqtree.org/</u> MrBayes <u>http://mrbayes.scs.fsu.edu</u> RaxML <u>https://cme.h-its.org/exelixis/web/software/raxml/</u>



