## PhyloBayes - what, why, and how

Xiaofan Zhou Integrative Microbiology Research Centre South China Agricultural University 2019-01-25

# **PhyloBayes**

#### a **Bayes**ian **phylogenetic** software based on *mixture* models

#### **Bayesian phylogenetics**

P(M,D)

**Posterior** probability

Likelihood

**Prior** 

# Among site variation in sequence evolution

- Evolutionary rates
- Substitution model
  - Exchange rate matrix
  - Equilibrium frequency profile (**PhyloBayes-CAT!!!**)

#### **Components of substitution model**



#### **Dirichlet process for site assignment**

#### **Long-Branch Attraction**



#### **PhyloBayes-CAT can suppress LBA**

![](_page_7_Figure_1.jpeg)

### Models available in PhyloBayes

Exchangeability matrix		Equilibrium frequency		<b>Evolutionary Rate</b>			Data partition		
•	Poisson Empirical model	•	Site homogeneous - similar to "+F" in ML	•	Gamma distribution - discrete - continuous	•	Partially linked branch-length		
•	- JTT - LG 	•	Site heterogeneous - CAT - empirical mixture models (e.g. C10-	•	CAT model of rate	•	Partition-specific substitution matrix not supported (?)		
•	GTR		C60)						
•	Mixture model								

#### **Partitioned PhyloBayes-MPI**

#### phylobayes mpi

记 141 commits	<b>}∕</b> 4 br	anches 🔊 <b>0</b> releases			4 2 contributors			<b>র্কু</b> GPL-2.0			
Branch: partition - New pull request				Crea	ate new file	Upload files	Find file	Clone or download -			
Switch branches/tags	×	behind master.					🎁 Pull	request 🖹 Compare			
Filter branches/tags	le) Latest commit 5d621ee on 23 Jul 201										
Branches Tags		pilation error (cons	t double)					6 months ago			
marginallogl		14781a2						3 years ago			
master		mit						4 years ago			
mutsel		titionfinder						8 months ago			
✓ partition	pcomp2 files m			3 years ago							
E README.md											
<b>pbmpi</b> partitioned phylobayes mpi Options unique to the partiti	oned v	ersion:									

-p <partition-file> partitioning scheme file in PartitionFinder format

### **Limitations of PhyloBayes**

- CAT+GTR model is **very** slow.....
  - Even with PhyloBayes-MPI
- Hard to converge
  - Particularly CAT+GTR
- CAT-Poisson is often used instead of CAT-GTR
- Results of unconverged analyses were reported

Syst. Biol. 66(2):232–255, 2017 © The Author(s) 2016. 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/syw084 Advance Access publication September 14, 2016

#### Who Let the CAT Out of the Bag? Accurately Dealing with Substitutional Heterogeneity in Phylogenomic Analyses

NATHAN V. WHELAN $^*$  AND KENNETH M. HALANYCH

![](_page_11_Figure_3.jpeg)

"..... CAT-F81 consistently performed worse than other models (*partitioning and CAT-GTR*) in inferring the correct branching patterns ....."

#### **Alternative solution in ML**

- Empirical mixture models in ML
  - C10 to C60

CAT-C10 profile mixture model of Le, Gascuel & Lartillot (2008)

frequency (10pi1 = 0.4082573125 0.0081783015 0.0096285438 0.0069870889 0.0349388179 0.0075279735 0.0097846653 0.1221613215 0.0039151830 0.0125784287 0.0158338663 0.0059670150 0.00813132 frequency (10pi2 = 0.1027763487 0.0418664491 0.0213272051 0.0155943616 0.0149663448 0.0440685478 0.0419667447 0.0138805792 0.0158864807 0.1066076641 0.1131944125 0.0436343681 0.04378003 frequency (10pi3 = 0.0351766018 0.0019678632 0.0016591476 0.0006768741 0.0078706538 0.0016559557 0.0019686768 0.0022420602 0.0012878339 0.3515819591 0.1278183107 0.0018856550 0.02426317 frequency (10pi4 = 0.0408513927 0.0269887074 0.2185648186 0.2333814790 0.0037602852 0.0380451418 0.0901238869 0.1158332065 0.0373197176 0.0025523644 0.0052164616 0.0485017266 0.000227517 frequency (10pi5 = 0.0185492661 0.0062362395 0.0024895723 0.0009775062 0.0070416514 0.0083539447 0.0024891617 0.0028952913 0.0040103982 0.1632422345 0.4443079499 0.0043570878 0.12028156 frequency (10pi5 = 0.1106750119 0.0352190043 0.0405186210 0.1636437899 0.0014834855 0.0877962201 0.2638456592 0.0325228293 0.0163803600 0.0068334902 0.0140679579 0.0677158208 0.002489881 frequency (10pi7 = 0.0522657662 0.0668294648 0.0714836849 0.0297745257 0.0143324928 0.0736540298 0.0388386669 0.0228101108 0.1551638111 0.0187406149 0.0653779932 0.0043969345 0.020718924 frequency (10pi7 = 0.0522657662 0.0668294648 0.0714836849 0.0207745257 0.0143324928 0.0736540298 0.0388386669 0.0228101108 0.1551638111 0.0187406149 0.0653779932 0.0043969345 0.020718934 frequency (10pi7 = 0.0522657662 0.0668294648 0.0714836849 0.0207142457 0.016334743 0.004203734 0.0024251122 0.003470413 0.0366787049 0.0187185330 0.067648974 0.02270839355 0.01243548 frequency (10pi9 = 0.0627195947 0.2038782162 0.002609174 0.0052662886 0.1098111767 0.068628494 0.0256174957 0.0332612124 0.0128968249 0.035627740 0.02270839355 0.01240359 frequency (10pi19 = 0.1145518598 0.0324008908 0.0750614981 0.0416192189 0.0098549497 0.0339624663 0.0364907910 0.05503817581 0.0165233329 0.0092949460 0.0139153707 0.0423026886 0.0082240 model C10 = P

#### **Alternative solution in ML**

- Empirical mixture models in ML
  - C10 to C60
- Fast implementation in IQ-TREE, further accelerated using PMSF

Tree,  
Alignment,  
Model (e.g.,  
C10)
$$P(j|x) = \frac{w_j \times P(x|j)}{\sum_j w_j \times P(x|j)} \longrightarrow f_a(x) = \sum_j f_{aj} \times P(j|x)$$

Wang et al. (2017) Syst. Biol.

#### Model comparison in PhyloBayes

- Bayes factor
- Cross-Validation
  - How well can the models predict future data?
- Posterior predictive analysis
  - How well does simulated data approximate observed data?

#### **Cross Validation**

![](_page_15_Figure_1.jpeg)

•  $P(D_{test}|D_{learn}, Ma)$  vs.  $P(Dtes_t|Dl_{earn}, Mb)$ 

#### **Posterior predictive analysis**

- Compare observed data vs. simulate data (using trained model)
- Resemble parametric bootstrap in ML
  - Parameters sampled from posterior distributions instead of being fixed
- Available tests:
  - Biochemical specificity
  - Compositional homogeneity
  - Saturation

![](_page_17_Picture_0.jpeg)

#### Tutorial

- 1. Run PhyloBayes analyses on a toy dataset with four taxa:
  - Check for convergence (e.g., TRACER)
  - Compare inferred tree with true tree (e.g., phylo.io)
  - Try different models (e.g., F81, GTR, CAT+F81, CAT+GTR...)
- 2. Perform model comparison
  - Cross-Validation
  - Posterior predictive analysis
- 3. Compare with empirical mixture models in ML (IQ-TREE)

#### The dataset was simulated under LBA

![](_page_19_Figure_1.jpeg)

- 4 taxa, ~21,000 amino acid sites
- Simulated under LG+C60+F
- *n* ranges from 1 to 10 to model different levels of LBA artefact

True tree

Wang et al. (2018) Syst. Biol.

#### Markov chain Monte Carlo

![](_page_20_Picture_1.jpeg)

(a) Symmetrical moves

![](_page_20_Picture_3.jpeg)

(b) Asymmetrical moves

#### When does it finish?

- Check for convergence between runs
- Visually:
  - TRACER (http://beast.community/tracer)
  - AWTY (http://king2.scs.fsu.edu/CEBProjects/awty/awty\_start.php)
  - Graphylo (https://github.com/wrf/graphphylo)
- bpcomp (in PhyloBayes)

#### When does it finish?

![](_page_22_Figure_1.jpeg)

![](_page_23_Figure_0.jpeg)