Recent Advances in General Phylogenetic Inference and Energyefficient Computing

Alexandros Stamatakis^{1,2,3}

1. Institute of Computer Science, Foundation for Research and Technology - Hellas

2. Heidelberg Institute for Theoretical Studies

3. Dept. of Informatics, Karlsruhe Institute of Technology

www.biocomp.gr (Crete lab)

www.exelixis-lab.org (Heidelberg lab)

Outline

- Recent Advances in Phylogenetic Inference
- Energy-efficient Computing

The number of trees

 $3 \text{ taxa} \rightarrow 1$ tree

The number of trees



4 taxa \rightarrow 3 trees

The number of trees



5 taxa \rightarrow 15 trees





possible trees with 2000 taxa

stamatak@exelixis:~/Desktop/GIT/TreeCounter\$./treeCounter -n 2000

GNU GPL tree number calculator released June 2011 by Alexandros Stamatakis

Number of unrooted binary trees for 2000 taxa: 30049638174211656151632910065681814981377232074237013089504954043012636525258308210827685996688247000464352

Problem Complexity



Problem Complexity



Finding the best tree under Maximum Likelihood is **NP-hard**!

Problem Complexity



Starting Trees



Starting Trees



A Tree with Support Values





Easy Dataset



Difficult Dataset



Phylogenetic Inference





Which data is more difficult to analyze?

Thousands of sequences, short sequence length

Which data is more difficult to analyze?







Intuitively it is this dataset here, as it contains much less information for telling apart more sequences

JOURNAL ARTICLE

Phylogenetic Analysis of SARS-CoV-2 Data Is Difficult d

Benoit Morel, Pierre Barbera, Lucas Czech, Ben Bettisworth, Lukas Hübner, Sarah Lutteropp, Dora Serdari, Evangelia-Georgia Kostaki, Ioannis Mamais, Alexey M Kozlov, Pavlos Pavlidis, Dimitrios Paraskevis, Alexandros Stamatakis ☎ Author Notes

Molecular Biology and Evolution, Volume 38, Issue 5, May 2021, Pages 1777–1791, https://doi.org/10.1093/molbev/msaa314

Published: 15 December 2020

Easy & Difficult Likelihood Surfaces



Easy & Difficult Likelihood Surfaces



Inferred 20 ML trees

125 taxa, 34 genes Inferred 20 ML trees

Now we can quantify this

- In past years these slides about easy and hard datasets were very hand-wavy
- Since 2022 we can quantify & predict difficulty

JOURNAL ARTICLE

From Easy to Hopeless—Predicting the Difficulty of Phylogenetic Analyses 👌

Julia Haag ™, Dimitri Höhler, Ben Bettisworth, Alexandros Stamatakis

Molecular Biology and Evolution, Volume 39, Issue 12, December 2022, msac254, https://doi.org/10.1093/molbev/msac254 Published: 17 November 2022

Easy & Difficult Likelihood Surfaces



Pythia Usage

- Pythia predicts difficulty of phylogenetic analysis via boosted tree regressor
- Input: molecular multiple sequence alignment or binary dataset
- **Output:** difficulty between **0.0** (easy) and **1.0** (hopeless)
- Invocations for our example datasets:

pythia --msa 125.phy --raxmlng ~/bin/raxml-ng

pythia --msa 7764.phy --raxmlng ~/bin/raxml-ng

Anecdotal Observations

- Good correlation between the difficulty score and the average bootstrap support values
- "Apparent convergence" speed of MCMC analyses can be predicted

Small SARS-CoV-2 dataset

"Phylogenetic Analysis of SARS-CoV-2 Data Is Difficult" (https://doi.org/10.1093/molbev/msaa314)

The predicted difficulty for MSA examples/covid.fasta is: 0.84.
FEATURES:
num_taxa: 4869
num_sites: 28361
[]
num_sites/num_taxa: 5.82
[]
avg_rfdist_parsimony: 0.79
proportion_unique_topos_parsimony: 1.0
Feature computation runtime: 1830.182 seconds
Г Т

PYTHIA Features

Table 1. Importance of the Subset of Features we use to Train Pythia.

	Feature	Impurity Importance
Parsimony = 76%	% Unique topologies parsimony trees	42.9%
	RF-distance parsimony trees	33.2%
	Entropy	17.0%
	Patterns-over-taxa	13.6%
	% Gaps	2.5%
	Bollback	2.3%
	Sites-over-taxa	1.5%
	% Invariant	0.6%

Empirical Difficulty Distributions

#MSAs



Difficulty: Binary Biological & Language Data



Difficulty score

Difficulty: Binary Biological & Language Data



Difficulty score

Difficult Datasets

- On difficult datasets
 - Infer a plausible tree set
 - And summarize it \rightarrow summary statistics

plausible tree set



SARS-CoV-2 consensus tree colored by country

Difficult Datasets

- On difficult datasets
 - Infer a plausible tree set
 - And summarize it \rightarrow summary statistics

plausible tree set If there is insufficient signal in the data to just infer a tree, one should not attempt to analyze the data under more complex (parameter-rich) models kelihood surface

Difficult Datasets

- On difficult datasets
 - Infer a plausible tree set
 - And summarize it \rightarrow summary statistics

Use/Produce a set of plausible Glottolog reference trees instead of just one tree? plausible tree set **Random resolutions of polytomies?**
Conflicting Signal

Quantify strength of conflicting signal?

JOURNAL ARTICLE

Novel Information Theory-Based Measures for Quantifying Incongruence among Phylogenetic Trees @

Leonidas Salichos, Alexandros Stamatakis, Antonis Rokas 🐱 👘 Author Notes

Molecular Biology and Evolution, Volume 31, Issue 5, May 2014, Pages 1261–1271, https://doi.org/10.1093/molbev/msu061

Published: 07 February 2014

Educated Bootstrap Guesser (EBG)





THE PREPRINT SERVER FOR BIOLOGY

New Results

Predicting Phylogenetic Bootstrap Values via Machine Learning

Julius Wiegert, Dimitri Höhler, Dulia Haag, Alexandros Stamatakis doi: https://doi.org/10.1101/2024.03.04.583288

This article is a preprint and has not been certified by peer review [what does this mean?].

A Tree with computationally expensive Standard Bootstrap Values



Educated Bootstrap Guesser (EBG)

- One order of magnitude faster than existing fast methods (UFBoot2: UltraFast Bootstrap version 2)
- Median error of 5 when predicting bootstrap values between 0-100
- 1654 SARS-CoV2 sequences
 - Bootstrap prediction in 3 hours on mid-class laptop

The Renaissance of Parsimony for Machine Learning

Light Gradient-Boosting Model from tree-based boosting ensemble framework.

	Feature	Importance in %
Parsimonv: 85%	PBS	82.2
	PS	3.1
	Normalized branch length	2.0
	# child inner branches	1.7
	Skewness PBS	1.5

PBS = **P**arsimony **B**ootstrap **S**upport from *200* parsimony bootstraps PS = **P**arsimony **S**upport from *1000* parsimony starting trees

Accuracy – Simulated Data



42

But ...



Accuracy on simulated data from UFBoot2 paper



Inferred bootstrap support



Inferred bootstrap support





Simulated Data Suck!

JOURNAL ARTICLE

Simulations of Sequence Evolution: How (Un)realistic They Are and Why 👌

Johanna Trost, Julia Haag ⊠, Dimitri Höhler, Laurent Jacob, Alexandros Stamatakis, Bastien Boussau Author Notes

Molecular Biology and Evolution, Volume 41, Issue 1, January 2024, msad277, https://doi.org/10.1093/molbev/msad277 Published: 20 December 2023 Article history •

We can distinguish between empirical and simulated molecular data with high accuracy using two distinct and independently developed machine learning based classification approaches!

Outline

- Recent Advances in Phylogenetic Inference
- Energy-efficient Computing

Motivation Rural Community Southern Crete

Crete

- Water temperature around Crete was 3 degrees above average in April
- Local fire chief was extremely nervous beginning of April
- Medium term summer temperature forecast is worrisome (European Centre for Medium-Range Weather Forecasts)
- Fire danger is **extreme** in many regions of Greece



Τρίτη - Τετάρτη, 11 – 12.06.2024

Σε περίπτωση εκδήλωσης πυρκαγιών, οι πυρομετεωρολογικές συνθήκες θα είναι ευνοϊκές για την γρήγορη εξάπλωση τους, καθιστώντας δύσκολο τον έλεγχο τους.

> Αποφυγή οποιασδήποτε χρήσης φωτιάς σε εξωτερικό χώρο. Παραμείνετε σε επαγρύπνηση. Μείνετε ενημερωμένοι. #ΜιαΦωτιάΛιγότερη





Efficiency is insufficient

- Efficiency gains != savings (Jevons 1865)
 - "But we optimize jet turbines!" \rightarrow see above & below



FIGURE 16

Renewable energy purchasing compared with total electricity use



Google Environment report (2023)

Decarbonizing the energy system

- Data centers have a role to play
 - Renewable Power Purchase Agreements (PPAs)
 - Heat reuse (partially)
 - Flexible demand

Clean-Power Shopping Spree

Top corporate buyers of clean energy globally through PPAs



Existing capacity Newly contract capacity

Source: BloombergNEF

Note: Chart only includes offsite, publicly announced power purchase agreements. Data is displayed in gigawatts of direct current capacity, with the exception of Microsoft's latest announcement. PPAs refers to power purchase agreements.

34

BloombergNEF

Timing the electricity market



energy-charts.info

Timing the electricity market



energy-charts.info

Carbon-aware scheduling





Dynamic carbon-aware power scaling



Carbon signal(s)

- Carbon intensity: CI [gCO₂/kWh]
- Share of renewable

+ ELECTRICITY MAPS





Alternatively: Carbon "traffic light"

- Discrete signal, e.g. green / yellow / red
- Ideally, reflects regional CI (Carbon Intensity)



https://energy-charts.info



https://carbonintensity.org.uk/





https://www.stromgedacht.de/

Test Carbon Intensities 2023 Historical Data

	Carbon intensity	Price	
Region	T y Signal ¢	N a r Type k e t	
Germany	a v e r continuous a g e	V F C I E national S a I E	

Carbon and price profiles: Germany (2023)



Carbon and price profiles: London (2023)



Carbon and price profiles: N. California (2023)



Dynamic carbon-aware power scaling



Experimental Setup

- Energy data
 - Historical data from 2023
- Hardware
 - 2x Intel Xeon CascadeLake
 - 2x Intel Xeon IceLake + 4x NVIDIA A40
- Software
 - SpecHPC 2021
 - + real-world tools: Arepo, Gromacs, RAxML-NG
 - 14 workloads (total)

Power profiles: RAPL powercap - CPU



System: Xeon Platinum 8260, Cascade Lake, 48C, 768 GB RAM

Scaling policy

Germany & North California

$$P_{lim}(CI) = \begin{cases} 100\% \text{ TDP} & \text{if } 0\% < CI <= 33\% \\ 80\% \text{ TDP} & \text{if } 33\% < CI <= 66\% \\ 60\% \text{ TDP} & \text{if } 66\% < CI <= 100\% \end{cases}$$

London



Evaluation



2023 electricity market data; 3-step scaling policy (100%/80%/60%); 14 HPC workloads, Xeon Cascade Lake

Features & Future Work

- Features & Setup
 - Specify Power Policy
 - Specify Electricity Data Provider
 - Put hardware to sleep on idle
- Future Work
 - Adapt to PV systems that are using net-metering



Contact



Oleksiy Kozlov Staff scientist, HITS gGmbH alexey.kozlov@h-its.org https://cme.h-its.org/exelixis



- Code: https://github.com/amkozlov/eco-freq
- Paper:

https://ieeexplore.ieee.org/document/10528928

Thank you for your attention



Listaros village, Crete



Difficult


What does Difficulty mean?

Difficulty = ruggedness of the tree space



- Few highly similar tree topologies
- Single likelihood peak

 Highly distinct topologies, statistically indistinguishable

Difficult

• Multiple likelihood peaks

Predicting Difficulty with Pythia

- Pythia = Boosted Tree Regressor
- Supervised Regression Task
 - Predict difficulty between **0** (easy) and **1** (difficult)
 - Ground truth difficulty as training target based on 100 distinct Maximum Likelihood tree inferences
- Initially trained on 4K empirical MSAs
 - Mean absolute error: 2.5%

Pythia developments

- New release (May 19, 2023)
 - Trained on 12K datasets
 - 11,108 DNA MSAs
 - 979 Protein MSAs
 - 460 Morphological MSAs
 - Two new features
 - Improved accuracy
 - Mean absolute error: 0.07 (previously 0.09)
 - Mean absolute percentage error: 1.7% (previously 2.5%)
- Using Pyhtia
 - See next slides

Definition of Difficulty



Prediction Features

• Eight Features

- 4 MSA attributes
 - Sites-over-taxa
 - patterns-over-taxa
 - % gaps
 - % invariant sites
- 2 MSA information metrics
 - Shannon entropy
 - Bollback multinomial test statistic
- 2 Parsimony-tree-based features
 - Infer 100 parsimony trees
 - \rightarrow average RF-Distance
 - \rightarrow % unique topologies

Using Pythia

- **Prior** to tree inference
 - \rightarrow determine analysis & post-analysis setup
 - → adjust/modify MSA
 - \rightarrow explore data filtering & assembly strategies
 - \rightarrow adjust user expectations about data

