



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

ReproPhylo

Citation for published version:

Szitenberg, A, John, M, Blaxter, ML & Lunt, DH 2015, 'ReproPhylo: An environment for reproducible Phylogenomics', PLoS Computational Biology, vol. 11, no. 9, 1004447.
<https://doi.org/10.1371/journal.pcbi.1004447>

Digital Object Identifier (DOI):

[10.1371/journal.pcbi.1004447](https://doi.org/10.1371/journal.pcbi.1004447)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Publisher's PDF, also known as Version of record

Published In:

PLoS Computational Biology

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



RESEARCH ARTICLE

ReproPhylo: An Environment for Reproducible Phylogenomics

Amir Szitenberg^{1*}, Max John¹, Mark L. Blaxter², David H. Lunt¹

1 Evolutionary Biology Group, School of Biological, Biomedical & Environmental Sciences, The University of Hull, Hull, United Kingdom, **2** Institute of Evolutionary Biology, The University of Edinburgh, Edinburgh, United Kingdom

* A.Szitenberg@hull.ac.uk



OPEN ACCESS

Citation: Szitenberg A, John M, Blaxter ML, Lunt DH (2015) ReproPhylo: An Environment for Reproducible Phylogenomics. *PLoS Comput Biol* 11(9): e1004447. doi:10.1371/journal.pcbi.1004447

Editor: Paul P Gardner, University of Canterbury, NEW ZEALAND

Received: May 27, 2015

Accepted: July 13, 2015

Published: September 3, 2015

Copyright: © 2015 Szitenberg et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: ReproPhylo is distributed under the CC0 license and uses open access dependencies. It is under active development within a publicly accessible GitHub repository (<http://goo.gl/s6EdVM>). Documentation is provided as a version tracked publicly-editable Google Docs manual (<http://goo.gl/yW6J1J>). A frozen version of the programme (Version 1.0), utilizing Jupyter Notebook as interface, is available as a self contained environment in a Docker image (<http://goo.gl/JcHMGn>). Use cases discussed in this manuscript are also available as Git repositories on GitHub (use case 1: <https://goo.gl/BsOxfl>, nbviewer: <http://goo.gl/KzFAvj>, use case 2: <https://goo.gl/26laiF>, nbviewer:

Abstract

The reproducibility of experiments is key to the scientific process, and particularly necessary for accurate reporting of analyses in data-rich fields such as phylogenomics. We present ReproPhylo, a phylogenomic analysis environment developed to ensure experimental reproducibility, to facilitate the handling of large-scale data, and to assist methodological experimentation. Reproducibility, and instantaneous repeatability, is built in to the ReproPhylo system and does not require user intervention or configuration because it stores the experimental workflow as a single, serialized Python object containing explicit provenance and environment information. This ‘single file’ approach ensures the persistence of provenance across iterations of the analysis, with changes automatically managed by the version control program Git. This file, along with a Git repository, are the primary reproducibility outputs of the program. In addition, ReproPhylo produces an extensive human-readable report and generates a comprehensive experimental archive file, both of which are suitable for submission with publications. The system facilitates thorough experimental exploration of both parameters and data. ReproPhylo is a platform independent CC0 Python module and is easily installed as a Docker image or a WinPython self-sufficient package, with a Jupyter Notebook GUI, or as a slimmer version in a Galaxy distribution.

This is a *PLOS Computational Biology* Software paper.

Introduction

Experimental reproducibility has become a widely discussed issue in many areas of science [1,2]. Strict experimental reproducibility is not common in any area of the biological sciences and while the reasons for this may be varied they include the technical challenges in routine and robust implementation. Phylogenetic analyses are very widely used across the biological sciences [3], and, even in studies that are not primarily phylogenetic, the understanding of phylogenetic relationships is almost always required for a meaningful statistical inference [4–6]. Despite this importance, the reproducibility of phylogenetic experiments is low, and Magee et al. [7] estimated that 60% of published phylogenetic analyses are “lost to science” due to the

<http://goo.gl/g3XP5B>), and in FigShare (<http://dx.doi.org/10.6084/m9.figshare.1409426>).

Funding: The Science of the Environment Council grant (<http://www.nerc.ac.uk/>) NE/J011355/1 was awarded to DHL and MLB. The Science of the Environment Council grant (<http://www.nerc.ac.uk/>) R8/H10/56 was awarded to GenPool, University of Edinburgh. The Medical Research Council grant (<http://www.mrc.ac.uk/>) G0900740 was awarded to GenPool, University of Edinburgh. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

unavailability of the underlying data, an outcome also predicted in other areas of biology [8]. However, even the public archiving of all data does not ensure reproducibility, since complete knowledge of the analytical software, software versions, software parameters, dependencies and operating system versions can be very challenging to both discover and recreate from published manuscripts. The increasing quantity of DNA sequence data available, and the proliferation of analytic toolkits, makes phylogenetics carried out on a genomic scale (“phylogenomics”) both especially powerful, and especially problematic to reproduce. Reproducibility in phylogenomics requires tracking of data provenance of multiple loci from many taxa, and, frequently, deeply nested analyses that explore, sift and partition data to achieve the end goals of biological understanding.

Here we introduce ReproPhylo, a Python package designed to deliver reproducible phylogenomic analyses. ReproPhylo promotes reproducibility on two levels. First, it eases the complex phylogenomic pipeline design process by providing a simple and concise scripting syntax for the execution of complex and forked phylogenetic workflows. Second, it automates reproducibility by employing well trusted containerization, versioning and provenance programs. In ReproPhylo, management of the experiment’s reproducibility and version control is carried out in a ‘frictionless’ manner in the background, without a need for user attention (although users have the option to access and tailor these aspects). Third, it ensures persistence and availability of metadata throughout the workflow, and in all the final products. With these three components of the analysis process considerably simplified, major important practices are addressed [9], and time and effort can be directed towards the core goals of understanding phylogenetic relationships by experimental parameter selection and data exploration, as the examples described here show (See [Results](#) section).

ReproPhylo is not the first package to provide phylogenetic workflow or pipeline tools [10–13]. A pipeline approach is a step forward from the point of view of reproducibility, as pipelines can serve as machine-readable records of analyses. Existing solutions [10–13] typically focus on the analysis itself, and do not attempt to provide complete reproducibility solutions. Several phylogenomic pipelines exist as web services [14–16], however, server-based analysis introduces additional complexities and reproducibility challenges, the main one of which is the dependency on a remote software environment. Osiris [17] achieves reproducibility through use of the Galaxy [18–20] reproducible bioinformatics environment, which can easily be used locally. Within the Galaxy framework, Osiris offers tools and format converters for widely used phylogenetic analysis programs, with user friendly and flexible GUI.

ReproPhylo explores an alternative, more generalised, approach to reproducibility, as it avoids dependency on any single high level software environment. It unifies the different components of a flexible, convenient, platform-independent, user friendly and reproducible workflow, drawing on the many advantages of standard data formats and community standard Biopython [21] code classes. ReproPhylo is simply accessed within a Jupyter Notebook (formerly IPython Notebook) [22]. We have also designed several basic ReproPhylo Galaxy tools, which produce self-contained and fully reproducible outputs, even outside the Galaxy system, as a proof of concept.

Design and Implementation

ReproPhylo interfaces with existing phylogenetic analysis tools *via* standard data structures, such as SeqRecord or MultipleSeqAlignment Biopython objects. In addition, it imports and exports data as text files in all standard formats supported by Biopython [21], and does not itself implement any novel data formats.

ReproPhylo can be run using Jupyter Notebook [22], where it is interacted with using a simple and self-explanatory Python syntax (examples in [S1 Methods](#)). We provide a range of notebooks for different types of analysis with the ReproPhylo distribution, including one for the Lepidoptera case analysis presented below. These notebooks are examples of ‘literate programming’ [23] in that they combine instructions, documentation, and code. The user may modify these Notebook pipelines either trivially (e.g. just changing the input data and executing), or more substantially (by altering the nature or sequence of analyses *via* Python code). Our testing with undergraduates, postgraduates, and academics without coding experience indicates that Jupyter Notebook is an effective GUI for scientists lacking a background in programming.

The ReproPhylo pipeline

ReproPhylo aids processes through the complete arc of a phylogenomics study: dataset collation, data analysis and visualisation/exploration. [Table 1](#) lists the data classes in ReproPhylo and their associated methods and functions. [Fig 1](#) illustrates a typical ReproPhylo workflow, and code snippets associated with each of the workflow steps are demonstrated in [S1 Methods](#). The ReproPhylo module uses a set of Python packages to control the pipeline and report results and quality statistics. The workflow is carried out by Biopython [21] and ETE2 [24], the latter of which also powers tree annotation. The primary output data file format is PhyloXML, although other formats can be produced. Graphics other than phylogenetic trees, such as alignment statistics and sequence statistics box-plots, are produced using Matplotlib [25].

Dataset collation in ReproPhylo has three components: harvesting, selection and filtering. An example of *data harvest* would be importing all GenBank records for a specific taxonomic group from a Genbank format text file, and adding unpublished sequences from a fasta or ab1 format sequence file. Exonerate [26] can be deployed within ReproPhylo to harvest loci of interest from genome or transcript data *via* specialized functions. *Data selection* exploits ReproPhylo’s loci report to automatically include or exclude specific genes and coding sequences present in an input Genbank file. *Data filtering* automatically excludes or includes sequences, or loci, based on user specifications—length, GC content, sequence number or taxonomic coverage—informed by ReproPhylo’s sequence and alignment summary statistics reports.

The analysis workflow in ReproPhylo includes sequence alignment, alignment trimming, and tree reconstruction. These steps can be forked to explore alternative analytic approaches while tracking data provenance in each branch and step. We have included commonly used analysis tools for each step, and additional algorithms can be suggested, or included by modifying the ReproPhylo module code (described in the manual, <http://goo.gl/yW6J11>). The first release of ReproPhylo can utilise the sequence aligners MAFFT [27], MUSCLE [28,29] and Pal2Nal [30]. Trimming of alignments to remove poorly aligned ‘gappy’ regions can improve analyses [31], and is carried out based on explicit trimming criteria using TrimAl [32]. Tree reconstruction programs accessible through ReproPhylo include RAxML [33] and PhyloBayes [34].

ReproPhylo facilitates phylogenetic output visualisation and exploration. Tree annotation, and creation of publication quality figures, is powered by ETE2 [24] and informed by metadata from the data harvest step provided to it by ReproPhylo. BayesTraits [35,36] is included for comparative phylogenetic analyses, and is invoked by a function which accepts a ReproPhylo Project object as the source of both the tree and trait information. Pairwise tree distances between trees in the Project can be computed and visualized (see [Results](#) section).

Table 1. Summary of the Python module structure.

Module feature	Description
Class Locus	Descriptor of the name, aliases, feature type and sequence type of an analysed locus
Class Project	Container for the input, intermediate and output datasets, and their metadata. Structured using Locus and Concatenation objects
method categories	
Read	Read data and metadata in any Biopython compatible format or tabular format for metadata
Filter	Filter sequences based on length, GC content or ID
edit_metadata	Programmatically manipulate sequence metadata
Align	Conduct sequence alignment(s) configured by a Conf object
Trim	Conduct alignment trimming configured by a Conf object
Tree	Conduct tree reconstruction(s) configured by a Conf object
Annotate	Annotate and root trees based on metadata stored in the Project
Write	Write files containing sequences, alignments, trees or metadata in any Biopython format
View	View alignments, statistics plots, occupancy tables etc. in the browser
Fetch	Copy a Project attribute (e.g. a tree or alignment object) into an independent variable
Conf Classes	A set of classes for configuring the different analytic steps
Class LociStats	Contains alignment and sequence parameters of the data in the Project
Methods	
Sort	sort the loci based on one of the available parameters
Plot	plot parameter boxplots
Slice	produce a supermatrix with certain parameter limits
Slide	create supermatrices by a sliding window approach along a gradient of a given parameter
Class Concatenation	Descriptor of the locus and OTU composition of a supermatrix
method categories	
Add	Add the concatenation to the analysis
Make	Prepare a supermatrix based on the instructions
Function categories	
list_loci	List loci found in a gb file, synonymize and choose from
Report	Write human readable report containing detailed methods and results
Pickle	Serialize/ Unserialize a Project object
Exonerate	Functions to run exonerate yielding metadata rich gb files
Bayestraits	Invokes BayesTraits using a Project object as the input source for both trees and traits

doi:10.1371/journal.pcbi.1004447.t001

Data provenance and reproducibility

Data provenance, the recording of the input and transformation of information used to generate a result, is a key issue in reproducibility. To maintain phylogenomic data provenance, ReproPhylo keeps the full workflow in a single instance of the Project ReproPhylo class (Fig 2A). This object contains all the analytical steps and their outputs, together with machine and human readable unique process IDs that describe the provenance of each data object for both the programme and the user. In addition, the Project instance contains the metadata associated with each sequence of each locus, with a unique ID, which allows it to associate the metadata

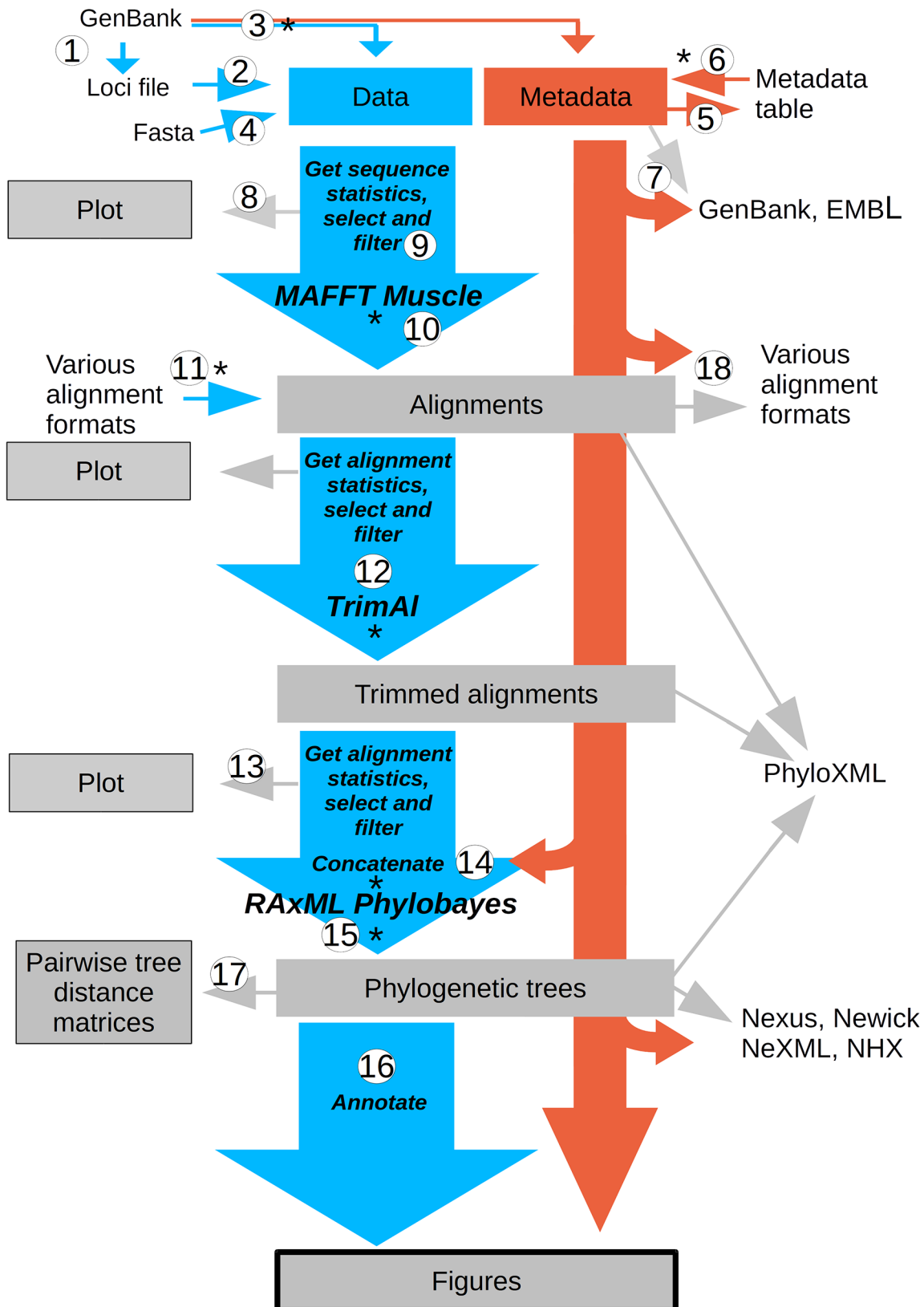
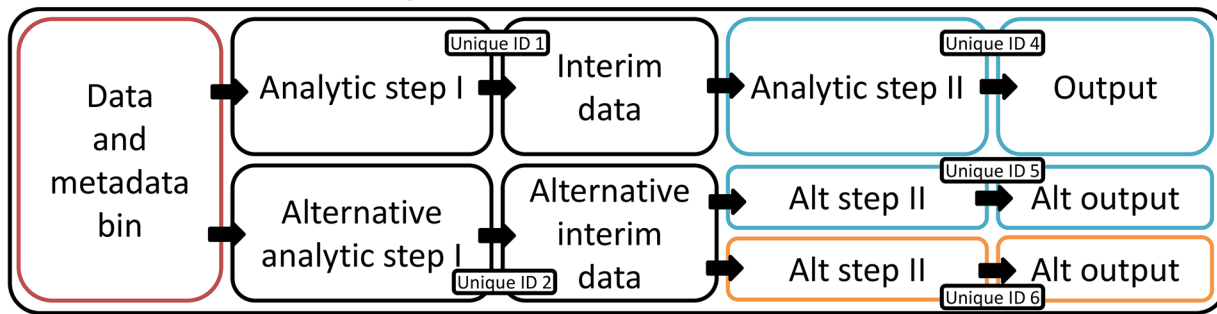


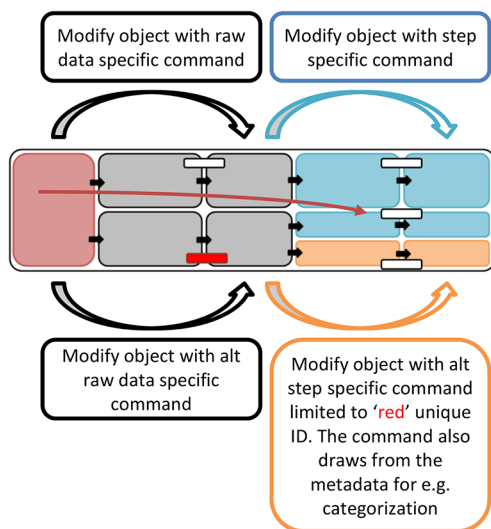
Fig 1. A typical ReproPhylo workflow. This illustration demonstrates the flow of data (blue arrows) and metadata (red arrows) through the phylogenetic analysis. Numbers on arrows correspond with code snippets in [S1 Methods](#). Asterisks indicate an automatic pickle and Git checkpoint. The user can toggle between these checkpoints indefinitely using a built in ReproPhylo function.

doi:10.1371/journal.pcbi.1004447.g001

A The workflow as an object



B Analytic steps as object modifiers



C Single file approach to provenance

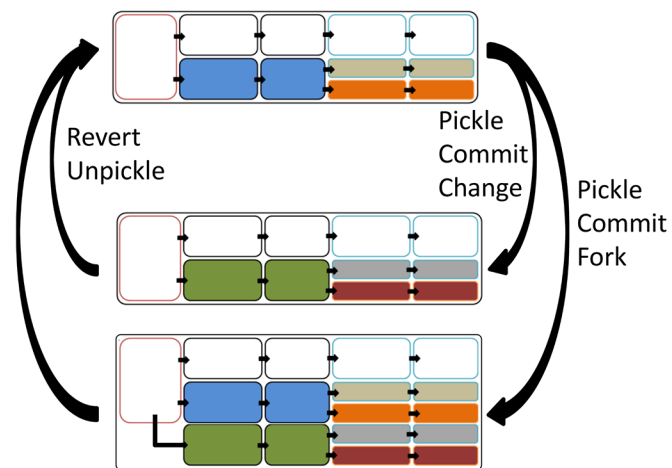


Fig 2. The phylogenetic workflow as a single Python object. (A) The workflow is contained as a single object with bins (attributes) for the raw data and metadata, as well as for the various workflow analyses and forks. These are made provenance-explicit with unique IDs and names. (B) Analyses are invoked via commands that modify the workflow object. A command can invoke batch analysis for all the relevant data in the object. For example, the command 'align' will apply for all the unaligned datasets. Commands can be limited to certain datasets using IDs. Commands can be customized using options. (C) Provenance survives version changes. The workflow object can be serialized (pickled) and then committed to a version control repository as a single file. Reverting to previous output version will also revert to the intermediate steps leading to it. Forks can be done post-hoc using the all-inclusive and provenance explicit workflow (pickled) object.

doi:10.1371/journal.pcbi.1004447.g002

with its sequence or tree leaf in any of the existing data objects (the SeqRecord, MultipleSeqAlignment and Tree objects). Analysis is invoked by Project class methods, which modify the data (e.g. align the sequences), place the resulting data object (e.g. MultipleSeqAlignment) in the appropriate Project attribute (e.g. Project.alignments) under a unique ID (Fig 2B), update the binary file storing the Project, and commit it to the Git repository. In each analytical step metadata can be retrieved using unique sequence identifiers, and alternative analytic approaches (forks) can be stored within a single Project through their unique process IDs.

Since the complete workflow is represented as a single Python object, provenance can be maintained across different versions of the analysis (Fig 2C). ReproPhylo serializes (“pickles”) the Project object and maintains it as a binary file that allows the user to pause and resume the analysis seamlessly. ReproPhylo uses the version control program Git (git-scm.com) to record a version of the binary Project file each time it is modified, and thus allows forwards and

backwards toggling of file versions. When an older version is restored, the full chain of intermediate results and the records detailing their production are restored throughout the workflow and across forks. ReproPhylo's version control and reproducibility are implemented passively in the background and are frictionless for the user, requiring neither specialist knowledge nor action to produce a reproducible phylogenomics experiment. The integration of Git in ReproPhylo is demonstrated in [S1 Example](#) (also in <http://dx.doi.org/10.6084/m9.figshare.1419590> and in nbviewer, <http://goo.gl/g3XP5B>).

To facilitate publication of the reproducible experiment, ReproPhylo produces a compressed experiment directory (.zip format) suitable for upload to a data repository such as FigShare (<http://figshare.com/>) or Dryad (<http://datadryad.org/>). This file contains trees and sequence alignments (in standard phyloXML format [37]), all analysis scripts, tree figure files, and a complete, human-readable report. The report includes a methods section ready for inclusion in a manuscript, which contains program versions, accession numbers, references etc., to which the digital object identifier of the full experimental record can be added. The compressed experiment directory also contains the binary file in which the serialized Project object is stored. This object contains all the data, metadata, method descriptions and results, and includes explicit provenance information. It can be used to revive the entire analysis, either in the ReproPhylo Docker container, in a local ReproPhylo installation or independently of ReproPhylo, and instantly repeat it or extend it. Another product of ReproPhylo is a Git repository, which can be published on websites such as Github (<http://github.com/>) and Figshare (<http://figshare.com/>). Both the compressed experiment directory and the Git repository satisfy all the Minimum Information about a Phylogenetic Analysis (MIAPA) goal [38], but the requirement for a description of the research objectives, by providing data files, data objects and human readable reports. They supersede the MIAPA requirements by also providing full software environment details and the machine readable scripts which have produced the intermediate and final files.

Version 1 of ReproPhylo is distributed as a Docker image (See [Availability and Future Directions](#) section). Using Docker as a work environment also facilitates reproducibility and reusability, as all relevant files can be committed to the image, generating a single Docker image file containing the computer environment, specific program copies, and data components of the finished analysis. Such containerisation approaches, which deliver both reproducible and easily reusable experiments, are powerful development and delivery tools [39].

Example use case

Several examples of use of the ReproPhylo phylogenomic analytical pipeline are provided as Jupyter notebooks in the distribution files. We focus here on parameter space exploration using ReproPhylo to demonstrate the advantages of phylogenomic analysis delivered by a fully scripted, reproducible environment. In this use case we demonstrate exploration of the effect of the median residue conservation (gene variability level) in each locus on a resulting species topology, using an existing multigene dataset of lepidopteran species [40]. Loci with different levels of conservation may hold phylogenetic signal of events that occurred in different times in the past, or may be too conserved, or too rapidly evolving and saturated with homoplasies, to provide any signal at all [41]. We utilise Shannon Entropy (SE) [42] as a conservation scoring method [43]. The script generating this analysis is available as [S2 Methods](#). The original Jupyter Notebook, together with the input and output files and figures, has been archived on FigShare (doi:[10.6084/m9.figshare.1409423](https://doi.org/10.6084/m9.figshare.1409423), goo.gl/KzFAvj), and has also been included as one of the tutorials in the current distribution of ReproPhylo (see ReproPhylo documentation at <http://goo.gl/aZeRXf>). A report with supplementary results generated by ReproPhylo is

provided as [S1 Results](#). Instructions on accessing the Project file in order to reproduce this demonstration are provided in the manual.

We obtained a nucleotide sequence alignment of 465 loci from 26 Lepidoptera species [40]. Using a built-in function ([S2 Methods](#), section 2.6.1), SE values [42], ignoring gap characters, were calculated for each residue in each locus. An entropy distribution plot ([Fig 3A](#), centre) illustrates the differences in SE among the loci. This plot is typical of alignment statistics and representations produced by the *ReproPhylo* LociStats class (see Section 2.6.3 of [S2 Methods](#) for code generating this plot). Six supermatrices were extracted, each from a sliding window of 200 loci, starting with the highest entropy loci and ending with the lowest entropy loci, and shifting the window by 50 loci between subsets ([Fig 3A](#)). Lastly, following the original analysis, all 26 species were included in all of the supermatrices, which contained no missing data ([S1 Results](#), [S1 Methods](#) section 2.7). Trees ([Fig 2](#)) were reconstructed as described in [S2 Methods](#), sections 2.5–2.10. Note that data partition information is utilised by *ReproPhylo* automatically. The trees were formally compared using the Symmetric Distance of Robinson-Foulds [44] ([Fig 3B](#)), the Branch Distance [45,46] ([Fig 3C](#)), and a modified Branch Distance [45] ([Fig 3D](#)), with standardized evolutionary rate ([S1 Methods](#), section 2.11).

Reproducibility statement

The entire project workflow for our analysis was saved as a pickle file ([S1 Results](#)), a Git repository generated by *ReproPhylo* (doi:10.6084/m9.figshare.1409423), and a publishable archive file ([S1 Results](#)). The pickled workflow can most productively be used within the *ReproPhylo* environment, where it is possible to add data and repeat the analysis or extend the analysis without the need to repeat any previous step. Importantly, the data within the pickled workflow is accessible using Biopython, even in the absence of *ReproPhylo*. The archive file represents a more traditional approach to reproducibility, as it includes alignment and tree text files, the tree figures ([Fig 3A](#)), and a human readable report containing complete methods and results information.

Results

We explored the partitioned Lepidoptera data for support for the clade Rhopalocera (butterflies) in loci with different SE values. Butterfly taxa are indicated in [Fig 3A](#) with dark blue highlight. The resulting topologies depend on the median entropy values in the dataset, with loci possessing low entropy values providing most support for Rhopalocera monophyly ([Fig 3A](#) trees 5–6). The result is similar for three other clades identified by Kawahara and Breinholt [40] (their clades I, III and IV; [Fig 3A](#) insets, light blue, yellow and gray highlights respectively). The entropy calculations were shown to be unbiased by the GC content or missing data ([S1 Fig](#); generated by section 2.4.6, [S2 Methods](#)). Formal tree comparisons ([Fig 2B–2D](#)), showing the topological differences ([Fig 3B](#)), the branch length differences ([Fig 3C](#)), and a combination of both ([Fig 3D](#)), also illustrate the effect of entropy on the topology and branch-lengths. This reaffirms the importance of analytic control over confounding effects.

The key novelty in the *ReproPhylo* environment is the ease and flexibility with which a complex phylogenetic investigation such as this can be set up, and be instantaneously repeatable and reproducible without compromising the user's control over parameter choice and configuration. *ReproPhylo* facilitates informed parameter choices and data filtering based on clearly documented and reproducible experimentation. Additional use cases are included with the package and they demonstrate the usage of additional components of the module and their interaction with Git and Docker.

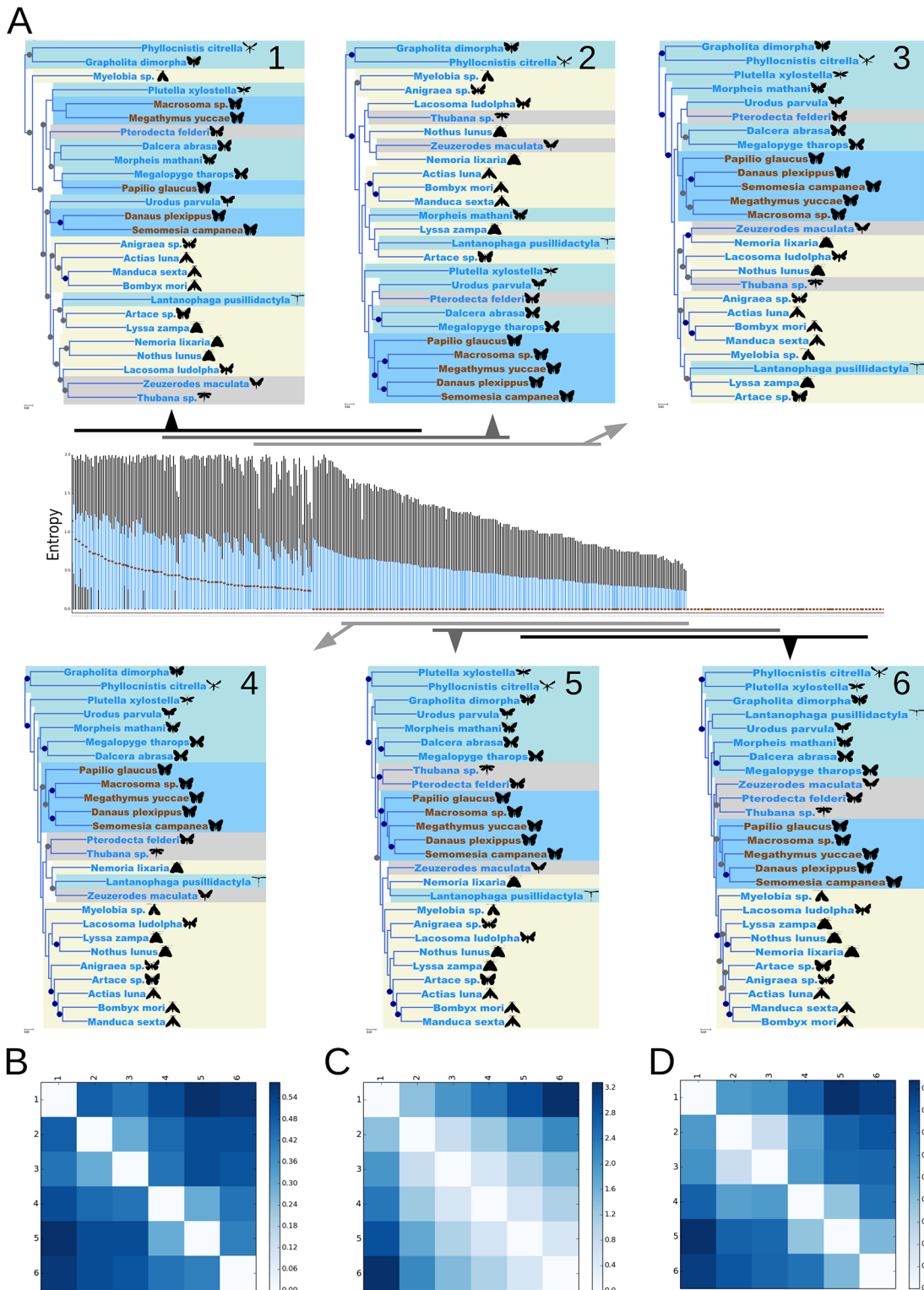


Fig 3. Exploratory phylogenomic analysis of a Lepidoptera dataset. (A) A nucleotide dataset from 26 species from Kawahara and Breinholt [40] was reanalyzed. Loci were sorted by their median, 75 percentile and 25 percentile entropy values (centre panel). For each locus, a box plot was generated. The medians are denoted by brown dots. The boxes (blue) represent the 25–75 percentiles. Whiskers (black) represent values that are found within a range outside the box, 1.5 times as long as the box (which is null, when the box itself has a null range) Trees (insets A 1–6) were reconstructed from 200-locus windows with 50 locus overlap between neighbouring windows. The windows are represented by black and gray horizontal bars, each with an arrow pointing to the tree generated from it. In trees 1–6, dark blue highlights denote Rhopalocera (butterfly) taxa, and light blue, gray and yellow highlights denote clades I,

III and IV respectively (*sensu* Kawahara and Breinholt [40]). Bullets on nodes represent Bootstrap percentages (BP). Blue bullets represent maximal support. Other support values above 80% are denoted by gray bullets. (B-D) Three pairwise tree divergence metrics were calculated and presented as heatmaps, with the most divergent tree pairs denoted by dark blue and identical tree pairs by a white box. While the scales are not comparable among the metrics, the relative differences are. The metrics are (B) the Symmetric Distance of Robinson-Foulds [44], (C) the Branch Distance [45] and (D) evolutionary rate corrected Branch Distance [45].

doi:10.1371/journal.pcbi.1004447.g003

ReproPhylo is an integrated environment for performing fully reproducible, platform independent, phylogenomics analyses that is highly accessible for scientists even without a strong computational background. ReproPhylo, by dealing with input and output formatting of data and results, can improve the accessibility and integration of existing computational tools. Phylogenetic analyses focussing on a single locus are becoming rarer as the power of modern genomics makes the *de novo* generation of large-scale data for multiple species feasible, especially with targeted sequencing approaches [47]. The rapid growth of public databases provides a resource that can be mined for new sets of loci across wide taxonomic spans, offering a second source of very large phylogenomic datasets. To exploit these new data, and at the same time deliver fully reproducible science that can lead to a truly incremental synthesis of evolution of life on earth, toolkits such as ReproPhylo that are large-data-ready, and natively reproducible will be essential.

Availability and Future Directions

ReproPhylo is open source, using strictly open source dependencies, and is under active development within a publicly accessible Github repository (<https://github.com/HullUni-bioinformatics/ReproPhylo>). Documentation is provided as a version tracked publicly-editable Google Docs manual at <http://goo.gl/yW6J1J>, allowing corrections and expansions by the user community. A frozen version of the module (Version 1), utilizing Jupyter Notebook as interface, is available as a self-contained environment in a Docker image (<http://goo.gl/JcHMGN>). Bioinformatics pipelines may often be challenging to install but the use of a Docker image for distribution eliminates such difficulties, and facilitates installation on any system. The Docker image is accompanied by a shell script that will install and deploy the ReproPhylo image as a Docker container, with a local web browser based GUI. We also provide ReproPhylo as a Win-Python version (see manual), and currently develop a Vagrant box solution (<https://www.vagrantup.com/>) for OSX. These will address any issues with the X11 server within Docker on Windows and Mac OSs. A repository containing the data and script for the analysis presented here is available on FigShare (<http://dx.doi.org/10.6084/m9.figshare.1409423>), as well as a repository containing the script and data for a demonstration of version control in ReproPhylo (<http://dx.doi.org/10.6084/m9.figshare.1419590>). The notebook containing the version control demonstration (<http://goo.gl/g3XP5B>) is also provided here as [S1 Example](#). As a proof of concept, ReproPhylo is also provided as a Galaxy distribution (<http://goo.gl/udsS3Q>) containing ReproPhylo Galaxy tools. This version utilises the Galaxy framework, while retaining completely reproducible results even outside the Galaxy GUI.

Future development is intended to include an extended suite of quality control indices, allowing better control over large datasets. Specifically, ReproPhylo can benefit from analyses that allow one to detect misleading signal in phylogenies [48]. In addition, we would like to include Resource Description Framework (RDF) outputs and parsers that will allow interactions with online repositories utilizing formal ontology descriptions [49] of phylogenetic experiments (e.g. CDAO-store [50]). Finally, ReproPhylo is intended to be a community tool, and we hope its future development will be guided by input from users, either by pull requests or issue reporting and suggestions in the Github repository.

Supporting Information

S1 Fig. Loci statistics boxplots for data derived from [40]. For each locus, the plots illustrate the distributions of (from top to bottom) per-position entropy, per-position gap score [32], per position conservation score [32], sequence length and GC content. <http://dx.doi.org/10.6084/m9.figshare.1409424>

(TIFF)

S1 Methods. An example code. The code snippets in this supplementary file are those associated with the numbered steps in the workflow illustrated in Fig 1. <http://dx.doi.org/10.6084/m9.figshare.1502477>.

(PDF)

S2 Methods. Scripts used in this research. A static HTML representation of the code that was used to create all the analyses in this study. <http://dx.doi.org/10.6084/m9.figshare.1409427> (HTML). Also in nbviewer: <http://goo.gl/KzFAvj>.

S1 Results. ReproPhylo report. A results archive produced by ReproPhylo, containing the serialized Project, input and output files, scripts and an HTML report. <http://dx.doi.org/10.6084/m9.figshare.1409488>

(ZIP)

S1 Example. A Jupyter notebook demonstrating version control in ReproPhylo (also available in FigShare (<http://dx.doi.org/10.6084/m9.figshare.1419590>) and nbviewer (<http://goo.gl/g3XP5B>)).

(HTML)

Acknowledgments

We thank Dr. Africa Gómez, Dr. Christoph Hahn, Dr. Stephen Moss, Daniel Jeffries, and Claudia Scavariello for useful comments on the program and the manuscript. The silhouettes in Fig 2 are distributed here under the [Creative Commons Attribution 3.0 Unported](https://creativecommons.org/licenses/by/3.0/) and are credited as follows: Geometroidea, Bombycoidea: Gareth Monger, Cossioidea: Didier Descouens (vectorized by T. Michael Keesey), Gelechioidea: Caroline Harding, MAF (vectorized by T. Michael Keesey).

Author Contributions

Conceived and designed the experiments: DHL AS MLB. Performed the experiments: AS DHL MJ. Analyzed the data: AS MJ. Contributed reagents/materials/analysis tools: DHL MLB. Wrote the paper: AS DHL MLB MJ.

References

1. McNutt M. Journals unite for reproducibility. *Science*. 2014; 346: 679. PMID: [25383411](https://pubmed.ncbi.nlm.nih.gov/25383411/)
2. Begley CG, Ioannidis JPA. Reproducibility in science improving the standard for basic and preclinical research. *Circ Res*. 2015; 116: 116–126. doi: [10.1161/CIRCRESAHA.114.303819](https://doi.org/10.1161/CIRCRESAHA.114.303819) PMID: [25552691](https://pubmed.ncbi.nlm.nih.gov/25552691/)
3. Eales JM, Pinney JW, Stevens RD, Robertson DL. Methodology capture: discriminating between the “best” and the rest of community practice. *BMC Bioinformatics*. 2008; 9: 359. doi: [10.1186/1471-2105-9-359](https://doi.org/10.1186/1471-2105-9-359) PMID: [18761740](https://pubmed.ncbi.nlm.nih.gov/18761740/)
4. Penny D. The comparative method in evolutionary biology. *J Classification*. 1992; 9: 169–172.
5. Whitney KD, Baack EJ, Hamrick JL, Godt MJW, Barringer BC, Bennett MD, et al. A role for nonadaptive processes in plant genome size evolution? *Evolution*. 2010; 64: 2097–2109. doi: [10.1111/j.1558-5646.2010.00967.x](https://doi.org/10.1111/j.1558-5646.2010.00967.x) PMID: [20148953](https://pubmed.ncbi.nlm.nih.gov/20148953/)

6. Ågren JA, Wang W, Koenig D, Neuffer B, Weigel D, Wright SI. Mating system shifts and transposable element evolution in the plant genus *Capsella*. *BMC Genomics*. 2014; 15: 602. doi: [10.1186/1471-2164-15-602](https://doi.org/10.1186/1471-2164-15-602) PMID: [25030755](https://pubmed.ncbi.nlm.nih.gov/25030755/)
7. Magee AF, May MR, Moore BR. The dawn of open access to phylogenetic data. *PLoS ONE*. 2014; 9: e110268. doi: [10.1371/journal.pone.0110268](https://doi.org/10.1371/journal.pone.0110268) PMID: [25343725](https://pubmed.ncbi.nlm.nih.gov/25343725/)
8. Vines TH, Albert AYK, Andrew RL, Débarre F, Bock DG, Franklin MT, et al. The availability of research data declines rapidly with article age. *Curr Biol*. 2014; 24: 94–97. doi: [10.1016/j.cub.2013.11.014](https://doi.org/10.1016/j.cub.2013.11.014) PMID: [24361065](https://pubmed.ncbi.nlm.nih.gov/24361065/)
9. Cranston K, Harmon LJ, O’Leary MA, Lisle C. Best practices for data sharing in phylogenetic research. *PLoS Curr*. 2014; 6.
10. Huerta-Cepas J, Bork P, Gabaldon T. ETE-NPR: A portable application for Nested Phylogenetic Reconstruction and workflow design http://etetoolkit.org/ete_npr/.
11. Pearse WD, Purvis A. phyloGenerator: an automated phylogeny generation tool for ecologists. *Methods Ecol Evol*. 2013; 4: 692–698.
12. Grant JR, Katz LA. Building a phylogenomic pipeline for the eukaryotic tree of life—addressing deep phylogenies with genome-scale data. *PLoS Curr*. 2014; 6.
13. Dunn CW, Howison M, Zapata F. Agalma: an automated phylogenomics workflow. *BMC Bioinformatics*. 2013; 14: 330. doi: [10.1186/1471-2105-14-330](https://doi.org/10.1186/1471-2105-14-330) PMID: [24252138](https://pubmed.ncbi.nlm.nih.gov/24252138/)
14. Sánchez R, Serra F, Tárraga J, Medina I, Carbonell J, Pulido L, et al. Phylemon 2.0: a suite of web-tools for molecular evolution, phylogenetics, phylogenomics and hypotheses testing. *Nucleic Acids Res*. 2011; 39: W470–4. doi: [10.1093/nar/gkr408](https://doi.org/10.1093/nar/gkr408) PMID: [21646336](https://pubmed.ncbi.nlm.nih.gov/21646336/)
15. Dereeper A, Guignon V, Blanc G, Audic S, Buffet S, Chevenet F, et al. Phylogeny.fr: robust phylogenetic analysis for the non-specialist. *Nucleic Acids Res*. 2008; 36: W465–9. doi: [10.1093/nar/gkn180](https://doi.org/10.1093/nar/gkn180) PMID: [18424797](https://pubmed.ncbi.nlm.nih.gov/18424797/)
16. Miller MA, Wayne P, Terri S. Creating the CIPRES Science Gateway for inference of large phylogenetic trees. 2010 Gateway Computing Environments Workshop (GCE). 2010.
17. Oakley TH, Alexandrou MA, Ngo R, Pankey MS, Churchill CKC, Chen W, et al. Osiris: accessible and reproducible phylogenetic and phylogenomic analyses within the Galaxy workflow management system. *BMC Bioinformatics*. 2014; 15: 230. doi: [10.1186/1471-2105-15-230](https://doi.org/10.1186/1471-2105-15-230) PMID: [24990571](https://pubmed.ncbi.nlm.nih.gov/24990571/)
18. Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, et al. Galaxy: A platform for interactive large-scale genome analysis. *Genome Res*. 2005; 15: 1451–1455. PMID: [16169926](https://pubmed.ncbi.nlm.nih.gov/16169926/)
19. Blankenberg D, Kuster GV, Coraor N, Ananda G, Lazarus R, Mangan M, et al. Galaxy: A web-based genome analysis tool for experimentalists. *Current Protocols in Molecular Biology*. John Wiley & Sons, Inc.; 2001.
20. Goecks J, Nekrutenko A, Taylor J. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol*. 2010; 11: R86. doi: [10.1186/gb-2010-11-8-r86](https://doi.org/10.1186/gb-2010-11-8-r86) PMID: [20738864](https://pubmed.ncbi.nlm.nih.gov/20738864/)
21. Cock PJA, Antao T, Chang JT, Chapman BA, Cox CJ, Dalke A, et al. Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics*. 2009; 25: 1422–1423. doi: [10.1093/bioinformatics/btp163](https://doi.org/10.1093/bioinformatics/btp163) PMID: [19304878](https://pubmed.ncbi.nlm.nih.gov/19304878/)
22. Pérez F, Granger BE. IPython: a system for interactive scientific computing. *Comput Sci Eng*. 2007; 9: 21–29.
23. Knuth DE. Literate programming. *Comput J*. 1984; 27: 97–111.
24. Huerta-Cepas J, Dopazo J, Gabaldón T. ETE: a python environment for tree exploration. *BMC Bioinformatics*. 2010; 11: 24. doi: [10.1186/1471-2105-11-24](https://doi.org/10.1186/1471-2105-11-24) PMID: [20070885](https://pubmed.ncbi.nlm.nih.gov/20070885/)
25. Hunter JD. Matplotlib: A 2D graphics environment. *Comput Sci Eng*. 2007; 9: 90–95.
26. Slater GSC, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics*. 2005; 6: 31. PMID: [15713233](https://pubmed.ncbi.nlm.nih.gov/15713233/)
27. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol*. 2013; 30: 772–780. doi: [10.1093/molbev/mst010](https://doi.org/10.1093/molbev/mst010) PMID: [23329690](https://pubmed.ncbi.nlm.nih.gov/23329690/)
28. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*. 2004; 32: 1792–1797. PMID: [15034147](https://pubmed.ncbi.nlm.nih.gov/15034147/)
29. Edgar RC. MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics*. 2004; 5: 1–19.
30. Suyama M, Torrents D, Bork P. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. *Nucleic Acids Res*. 2006; 34: W609–W612. PMID: [16845082](https://pubmed.ncbi.nlm.nih.gov/16845082/)

31. Talavera G, Castresana J. Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Syst Biol*. 2007; 56: 564–577. PMID: [17654362](#)
32. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*. 2009; 25: 1972–1973. doi: [10.1093/bioinformatics/btp348](#) PMID: [19505945](#)
33. Stamatakis A. RAxML Version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics*. 2014; btu033.
34. Lartillot N, Lepage T, Blanquart S. PhyloBayes 3: a Bayesian software package for phylogenetic reconstruction and molecular dating. *Bioinformatics*. 2009; 25: 2286–2288. doi: [10.1093/bioinformatics/btp368](#) PMID: [19535536](#)
35. Pagel M. Detecting correlated evolution on phylogenies: a general method for the comparative analysis of discrete characters. *Proc R Soc B*. 1994; 255: 37–45.
36. Pagel M, Meade A, Barker D. Bayesian estimation of ancestral character states on phylogenies. *Syst Biol*. 2004; 53: 673–684. PMID: [15545248](#)
37. Han MV, Zmasek CM. phyloXML: XML for evolutionary biology and comparative genomics. *BMC Bioinformatics*. 2009; 10: 356. doi: [10.1186/1471-2105-10-356](#) PMID: [19860910](#)
38. Leebens-Mack J, Vision T, Brenner E, Bowers JE, Cannon S, Clement MJ, et al. Taking the first steps towards a standard for reporting on phylogenies: Minimum Information About a Phylogenetic Analysis (MIAPA). *OMICS*. 2006; 10: 231–237. PMID: [16901231](#)
39. Boettiger C. An introduction to Docker for reproducible research. *Oper Syst Rev. ACM*; 2015; 49: 71–79.
40. Kawahara AY, Breinholt JW. Phylogenomics provides strong evidence for relationships of butterflies and moths. *Proc R Soc B*. 2014; 281: 20140970. doi: [10.1098/rspb.2014.0970](#) PMID: [24966318](#)
41. Higgs PG. RNA secondary structure: physical and computational aspects. *Q Rev Biophys*. 2000; 33: 199–253. PMID: [11191843](#)
42. Shannon CE. A Mathematical Theory of Communication. *SIGMOBILE Mob Comput Commun Rev*. 2001; 5: 3–55.
43. Valdar WSJ. Scoring residue conservation. *Proteins*. 2002; 48: 227–241. PMID: [12112692](#)
44. Robinson DF, Foulds LR. Comparison of phylogenetic trees. *Math Biosci*. 1981; 53: 131–147.
45. Kuhner MK, Felsenstein J. A simulation comparison of phylogeny algorithms under equal and unequal evolutionary rates. *Mol Biol Evol*. 1994; 11: 459–468. PMID: [8015439](#)
46. Sukumaran J, Holder MT. DendroPy: a Python library for phylogenetic computing. *Bioinformatics*. 2010; 26: 1569–1571. doi: [10.1093/bioinformatics/btq228](#) PMID: [20421198](#)
47. Lemmon AR, Emme SA, Lemmon EM. Anchored hybrid enrichment for massively high-throughput phylogenomics. *Syst Biol*. 2012; 61: 727–744. doi: [10.1093/sysbio/sys049](#) PMID: [22605266](#)
48. Struck TH. TreSpEx-Detection of misleading signal in phylogenetic reconstructions based on tree information. *Evol Bioinform Online*. 2014; 10: 51–67. doi: [10.4137/EBO.S14239](#) PMID: [24701118](#)
49. Schulze-Kremer S. Ontologies for molecular biology and bioinformatics. *In Silico Biol*. 2002; 2: 179–193. PMID: [12542404](#)
50. Chisham B, Wright B, Le T, Son TC, Pontelli E. CDAO-store: ontology-driven data integration for phylogenetic analysis. *BMC Bioinformatics*. 2011; 12: 98. doi: [10.1186/1471-2105-12-98](#) PMID: [21496247](#)