

Supplementary Information

palaeoverse: a community-driven R package to support palaeobiological analysis

Lewis A. Jones¹, William Gearty², Bethany J. Allen^{3,4}, Kilian Eichenseer⁵, Christopher D. Dean⁶, Sofia Galván¹, Miranta Kouvari^{6,7}, Pedro L. Godoy^{8,9}, Cecily Nicholl⁶, Lucas Buffan¹⁰, Erin M. Dillon^{11,12}, Joseph T. Flannery-Sutherland¹³, and Alfio Alessandro Chiarenza¹

¹*Centro de Investigación Mariña, Grupo de Ecoloxía Animal, Universidade de Vigo, 36310 Vigo, Spain.*

²*Division of Paleontology, American Museum of Natural History, New York, NY, 10024 USA.*

³*Department of Biosystems Science and Engineering, ETH Zürich, 4058 Basel, Switzerland.*

⁴*Computational Evolution Group, Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland.*

⁵*Department of Earth Sciences, Durham University, South Road, DH1 3LE, Durham, United Kingdom.*

⁶*Department of Earth Sciences, University College London, Gower Street, WC1E 6BT, London, United Kingdom.*

⁷*Life Sciences Department, Natural History Museum, Cromwell Road, SW7 5BD, London, United Kingdom.*

⁸*Laboratório de Paleontologia, Faculdade de Filosofia, Ciências e Letras de Ribeirão Preto, Universidade de São Paulo, Ribeirão Preto, SP, 14040-901 Brazil.*

⁹*Department of Anatomical Sciences, Stony Brook University, Stony Brook, NY, 11794 USA.*

¹⁰*Département de Biologie, École Normale Supérieure de Lyon, Université Claude Bernard Lyon 1, 69342 Lyon Cedex 07, France.*

¹¹*Smithsonian Tropical Research Institute, Balboa, Republic of Panama.*

¹²*Department of Ecology, Evolution, and Marine Biology, University of California, Santa Barbara, CA 93106, USA.*

¹³*School of Earth Sciences, University of Bristol, BS8 1RL, Bristol, UK*

Corresponding author: LewisAlan.Jones@uvigo.es

Survey

We conducted an online public survey to collect the opinions of the palaeobiological community, and determine which tools were most needed. The survey was advertised via social media (i.e. Twitter) and a mailing list (not included here for the sake of privacy). The survey was opened on the 26th of May 2022 via the following link (<http://www.tinyurl.com/palaeoverse>). For the purpose of this manuscript, answers were collected until the 9th of August, 2022. However, the survey remains open for the community to respond to. In total, 35 participants from 14 different countries completed the survey, most of which are affiliated with European and North American institutions. This observation likely reflects a geographical bias in the reach of the survey, and demonstrates the additional effort that will need to be made to reach further afield in the future.

Below, we include the conducted survey for the sake of completeness:

palaeoverse: towards a community-driven R package

palaeoverse is an R package being developed by palaeobiologists, for palaeobiologists.

The aim of palaeoverse is to generate a community-driven software package of generic functions for the palaeobiology community. The package does not aim to provide implementations of statistical approaches, rather it provides auxiliary functions to help streamline analyses, and improve code readability and reproducibility.

As part of the initial development of palaeoverse, we would like to hear from you, the palaeobiology community! What generic functions do you wish existed to streamline your work? What takes up more of your precious time than it should? For example, time and spatial binning of data, palaeorotating fossil occurrences, checking species names for errors... this kind of thing! Let us know what functions you think would be useful for you, and for the wider community.

1. Name (short answer text)
2. Email address (short answer text)
3. Affiliation (long answer text)
4. Do you wish to join our mailing list, and receive updates on palaeoverse? (If so, please include your email above) (select one of the following)
 - Yes
 - No

5. What types of palaeontological data do you typically use or are interested in? (Multiple choice)

- Taxonomic identifications
- Age or biozone
- Geographic information (modern or at time of deposition)
- Geological/palaeoenvironmental context
- Taxon abundance
- Taphonomy
- Trait values, classifications or descriptions
- Phylogenies
- Other

6. What R packages (or other tools) have you used previously to clean or explore palaeontological data sets? (long answer text)

7. What kinds of tasks do you typically carry out when cleaning and exploring palaeontological data? (long answer text)

8. Please specify and detail generic functions that you feel would be useful to include in the palaeoverse package. If you have suggestions for more than three functions, please submit a second form.

- Function 1 (long answer text)
- Function 2 (long answer text)
- Function 3 (long answer text)

9. Do you have pre-existing code (for the above functions or other functions) you would like to contribute to palaeoverse? (select one of the following)

- Yes
- No

10. If yes, what function is your code performing? (long answer text)

Survey responses

Below we provide a summary of the responses to survey questions not documented in the main text (6–8). Some questions (1–4) are omitted for the sake of privacy. Where appropriate, we grouped survey responses into the following distinct categories to aid summary:

- Data access

- 91 • Checking and transforming data
- 92 • Plotting data
- 93 • Time binning
- 94 • Spatial analyses
- 95 • Tree modification and plotting
- 96 • Phylogenetic analyses
- 97 • Other analyses

98 *Responses to questions related to the tasks participants usually carry out*

99 Checking and transforming data

- 100 • Transforming data into portable structures
- 101 • Organising based on multiple data types
- 102 • Grouping things by character state
- 103 • Restructuring and filtering
- 104 • Cross-referencing PBDB with other datasets for taxonomy errors (synonyms/misspellings)
- 105 • Checking for occurrences outside expected taxon duration, or freshwater species in marine
- 106 sediments
- 107 • Checking that column names are identical between imported datasets from the literature
- 108 • Reshaping into tidy/long format
- 109 • Transforming species abbreviations into full species names
- 110 • Formatting and matching data with trees
- 111 • Processing age/taxon information for phylogenetic analysis
- 112 • Merging files using specific specimen ids, cleaning repeat specimens, fixing age ranges across time
- 113 intervals, updating species names and connecting to past species names
- 114 • Checking taxonomic assignments/spellings, verifying geochronology and environmental setting,
- 115 checking for missing data, reviewing spatial distribution, reviewing data summaries by taxonomic
- 116 group
- 117 • Refining chronostratigraphy for PBDB collections
- 118 • Standardising taxonomy
- 119 • Surveying and cleaning biostratigraphic data on PBDB
- 120 • Taxonomic harmonisation, filtering non-pollen and converting to percentages

- 121 • Aligning multiple ecological abundance, paleoenvironmental and depositional datasets (sometimes
- 122 from different csv files, collected by different people or from different core intervals) in time/space
- 123 and with age-depth model outputs
- 124 • Calculating derived variables from raw data
- 125 • Reformatting dataframes for various analyses/packages which might require data to be formatted
- 126 in a different way, especially when some variables are measured per sample (relative abundance)
- 127 and others are measured per individual specimen (taphonomy scores)
- 128 • Tidying taxonomic identifications
- 129 • Trimming datasets above a threshold

130 Plotting data

- 131 • Exploratory figures
- 132 • Plotting census data as a diagram to see faunal trends, nMDS, etc
- 133 • Skimming through the dataset and plotting temporal/morphological data
- 134 • Data visualisation
- 135 • Visualising variables to look for distributions, outliers, and preliminary trends
- 136 • Overlaying multiple paleo time series plots with aligned axes

137 Time binning

- 138 • Temporal subsetting/assignment
- 139 • Binning data by time period
- 140 • Time binning with variable binning rules (midpoint, range-through, random age assignment,
- 141 exclusion of multi-bin spanning data)
- 142 • Binning
- 143 • Time binning

144 Tree modification and plotting

- 145 • Stitching two phylogenetic trees together
- 146 • Acquiring FADs and LADs, matching taxon names to tree tip labels, reading character data with
- 147 missing data or polymorphisms, rooting trees
- 148 • Tree calibration
- 149 • Time-calibrating phylogenies
- 150 • Converting TNT tree files into NEXUS format

- 151 • Tree plotting, pruning, assigning max ages or nodes, exploration of evolutionary rates and shape
152 changes, phylogeny building, etc

153 Phylogenetic methods

- 154 • Estimating phylogenetic signal, phylogenetic regressions and multimodel inference, predictive
155 modelling for paleobiological inference, ancestral state reconstructions, analysis and mapping of
156 evolutionary rates
- 157 • Phylogenetic comparative methods and multivariate statistics, origination analyses, model fitting
158 and ancestral state reconstructions
- 159 • Tree plotting, pruning, assigning max ages or nodes, exploration of evolutionary rates and shape
160 changes, phylogeny building, etc

161 Spatial analyses

- 162 • Geographical analyses
- 163 • Spatial subsetting/assignment
- 164 • Converting extant latitude and longitude to extinct latitude and longitude
- 165 • Aggregate collections within a given radius into palaeocommunities/sites

166 Data Access

- 167 • GBIF guide

168 Other analyses

- 169 • Diversity analyses using different metrics
- 170 • Resampling and appropriate rarefaction
- 171 • Creating tip priors/contrast matrices
- 172 • Applying an age model
- 173 • Comparing morphological/isotope measurements across space and time
- 174 • Paleoenvironmental reconstructions
- 175 • Ecological coupling between aquatic and terrestrial systems
- 176 • Human-environment links
- 177 • Trend analyses
- 178 • Multi-proxy comparison/analysis
- 179 • Age modelling

- 180 • Sensitivity analyses
- 181 • Relative abundance

182 *Responses to questions related to the functions that participants consider useful to be included in*
183 *palaeoverse*

184 Checking and transforming data

- 185 • Removing taxonomic equivalents, species with too little data and uninformative characters
- 186 • Checking for typos in taxa names, especially between two vectors; checking if taxa names are
187 formatted as binomial entities, i.e. checking if they follow a “Genus_species” format and return an
188 error message if not.
- 189 • Cross-matching PBDB and phylogenetic data (bonus: some functionality for visualising this data)
- 190 • Cleaning specimen description data into components when single specimen number has multiple
191 entries
- 192 • Detection of unusual spatiotemporal occurrences of PBDB records
- 193 • Taxonomic harmonization
- 194 • Resolving to taxonomic authorities, such as Paleobiology Database or World Register of Marine
195 Species
- 196 • Prepping datasets for RevBayes analyses
- 197 • Merging duplicate species or lumping specified species

198 Plotting data

- 199 • Plotting data on stratigraphic-geological timeline with appropriate colours
- 200 • Range charts
- 201 • Visualising geographic data
- 202 • Adding geologic timescales to plots in base R graphics rather than ggplot
- 203 • Plotting spindle diagrams given a tree and diversity (or abundance) data through time
- 204 • Perspective and stacked time slice disparity plots
- 205 • coord_geo (similar to deeptime R package) to add geologic timescale to a plot
- 206 • core_photograph to automatically look up core images remotely, download them, shrink them, add
207 them to the plot
- 208 • Faunal diagram (eg % abundance of different species vs age or depth) from census data
- 209 • Heatmaps to overview faunal composition

- 210 • Stratigraphic charts
- 211 • Summary plots (e.g. richness by geologic unit, time interval, and/or taxonomic groups)
- 212 • Adding geological scales to ggplot2
- 213 • Plotting
- 214 • Overlaying multiple time series plots over an aligned time axis
- 215 Time binning
 - 216 • Time binning
 - 217 • Age binning
 - 218 • FADs to LADs ranging with gap info
 - 219 • Age binning of PBDB data
 - 220 • Biostrat or geological epoch to quantitative age
 - 221 • Time binning function with variable binning rules (midpoint age, assign randomised age in FAD-
 - 222 LAD range then bin, range-through binning)
 - 223 • Translating periods/stages into up-to-date absolute ages
 - 224 • Converting and organising input data according to age (across different age types - calBP, 14C,
 - 225 etc)
 - 226 • Updating timescale ages with new chronostratigraphic works
 - 227 • Tools to facilitate time binning and alignment of coeval data points across time intervals, especially
 - 228 when dealing with age-depth model uncertainty, different temporal resolutions, etc
- 229 Tree modification and plotting
 - 230 • Tree slicing
 - 231 • Treestich: attach two trees together
 - 232 • Plotting unambiguous synapomorphies on a tree like MacClade used to do
 - 233 • Plotting up tree spaces
 - 234 • Finding centroid trees
- 235 Phylogenetic methods
 - 236 • Using ranges in ancestral state estimation, rather than a single mean value
 - 237 • Grouping taxa by character trait and work with character matrices easily
 - 238 • Total-evidence phylogeny plot

- 239 • Performing cross-validation for predictive models in a phylogenetic context (e.g. for PGLS)
- 240 • A function that simply removes all fossil taxa based on non-ultrametricity and performs leave-n-
- 241 out cross-validation
- 242 • Tree stats, Templeton Test, double decay indices, positional congruence, leaf stability
- 243 • Comparisons or tests for correlations between the diversification rates of two different clades

244 Spatial analyses

- 245 • Latitudinal bins assignment from palaeolatitudes
- 246 • Palaeo-ocean basin assignment from lat/long and date
- 247 • Converting long/lat with fossil age to past long/lat
- 248 • Spatial binning function using hexagonal grids
- 249 • Palaeocoordinating

250 Data Access

- 251 • Gathering data from the Paleobiology Database
- 252 • Processing private data into a format to be easily uploaded to Neotoma
- 253 • Read TNT files or a way to quickly and easily convert files
- 254 • Bringing TNT into R
- 255 • Data downloading functions from various databases (e.g. PBDB)
- 256 • Getting data
- 257 • Compiling a list of relevant packages that perform useful functions to increase exposure

258 Other analyses

- 259 • Rarefaction
- 260 • Tools for data simulation to generate null models
- 261 • Generalised differencing
- 262 • Fossil frequency
- 263 • Oldest record
- 264 • Parsing stratigraphic range data for use in phylogenetic packages and PyRate
- 265 • Age models
- 266 • Converting element descriptions to a standardize form which the user can define, using a few
- 267 specific shared terms in those descriptions

- Integrating paleobiological and stratigraphic data
- Breakpoint and smooth community data analyses
- Age-depth modelling
- Assessing distributions of sampling rates among taxa
- Easy extraction of some results (e.g. clock rates) from Bayesian analyses
- Age model
- Beta diversity time-series analyses. For instance, a function that computes cumulative (from a reference point) and successive (between time points) assemblage turnover
- With pollen data, deciphering relative abundance of terrestrial vs all species

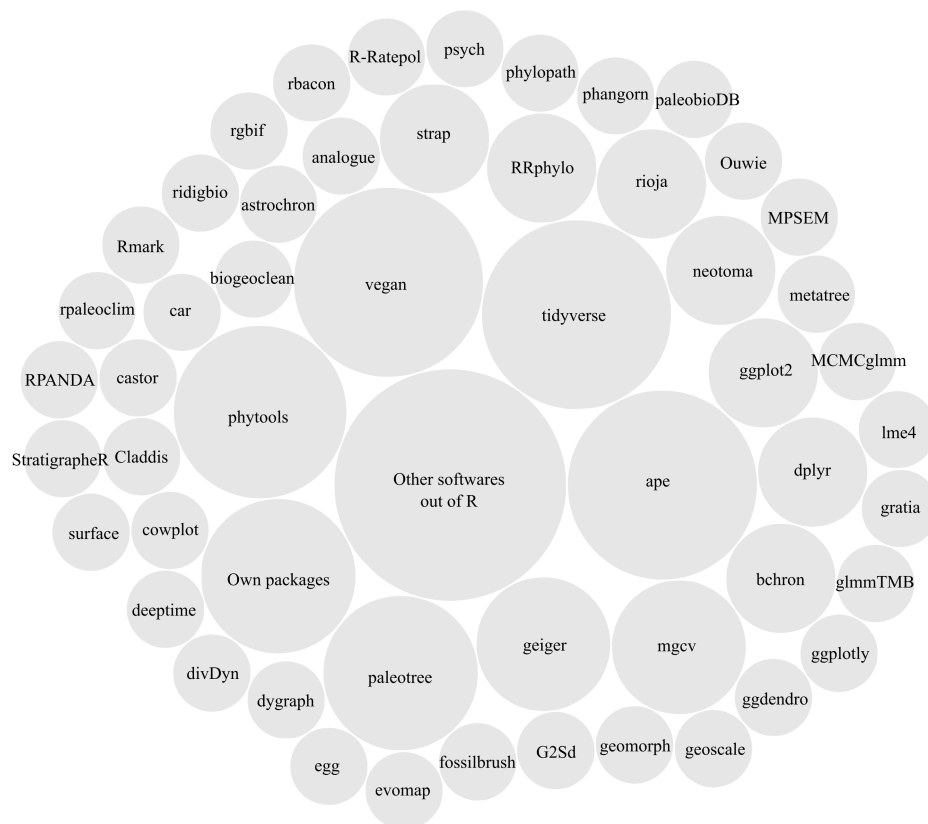


Figure 1: Summary of responses to the palaeoverse survey. Preferred tools for processing and analysing paleontological data. Both resources inside and outside the R environment and R packages are included as categories.