

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 S. C. Nicholl⁶, Lucas Buffan¹⁰, Erin M. Dillon^{11,12}, Joseph T. Flannery-Sutherland¹³, and Alfio Alessandro Chiarenza¹

¹*Grupo de Ecología Animal, Departamento de Ecología e Biología 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



Abstract

1. The open-source programming language ‘R’ has become a standard tool in the palaeobiologist’s toolkit. Its popularity within the palaeobiology community continues to grow, with published articles increasingly citing the usage of R and R packages. However, there are currently a lack of agreed standards for data preparation and available frameworks to support implementation of such standards. Consequently, data preparation workflows are often unclear and not reproducible, even when code is provided. Moreover, due to a lack of code accessibility and documentation, palaeobiologists are often forced to ‘reinvent the wheel’ to find solutions to issues already solved by other members of the community.
2. Here, we introduce **palaeoverse**, a community-driven R package to aid data preparation and exploration for quantitative palaeobiological research. The package is freely available and has three core principles: (1) streamline data preparation and analyses; (2) enhance code readability; and (3) improve reproducibility of results. To develop these aims, we assessed the analytical needs of the broader palaeobiological community using an online survey, in addition to incorporating our own experiences.
3. In this work, we first report the findings of the survey which shaped the development of the package. Subsequently, we describe and demonstrate the functionality available in **palaeoverse** and provide usage examples. Finally, we discuss the resources we have made available for the community and the future plans for the broader **palaeoverse** project.
4. **palaeoverse** is a community-driven R package in palaeobiology, developed with the intention of bringing palaeobiologists together to establish agreed standards for high-quality quantitative research. The package provides a user-friendly platform for preparing data for analysis with well-documented open-source code to enhance transparency. The functionality available in **palaeoverse** improves code reproducibility and accessibility, which is beneficial for both the review process and future research.

Keywords

Analytical Palaeobiology, Computational Palaeobiology, R programming, Readable, Reusable, Reproducible

Resumen (Español)

1. El lenguaje de programación de código abierto ‘R’ se ha convertido en una herramienta común entre paleobiólogos. Su popularidad dentro de esta comunidad continúa creciendo, con cada vez más artículos que citan el uso de R y sus paquetes. Sin embargo, todavía faltan protocolos para la preparación de los datos, así como marcos de actuación para la aplicación de dichos protocolos. Por ello, los flujos de trabajo para la preparación de datos suelen resultar confusos y a menudo no reproducibles, incluso cuando se facilita el código. Además, debido a la falta de accesibilidad y documentación del código, los paleobiólogos se ven obligados frecuentemente a ‘reinventar la rueda’ para encontrar soluciones a problemas ya resueltos por otros miembros de la comunidad.
2. Por este motivo presentamos palaeoverse, un paquete colaborativo de R para facilitar la preparación y exploración de datos en paleobiología cuantitativa. El paquete es de acceso libre y posee tres principios básicos: (1) agilizar la preparación y análisis de datos; (2) facilitar la lectura del código y (3) mejorar la reproducibilidad de los resultados. Para alcanzar estos objetivos, evaluamos las necesidades analíticas de la comunidad paleobiológica a través de una encuesta en línea, además de incorporar nuestras propias experiencias.
3. En primer lugar, en el presente trabajo mostramos los resultados de la encuesta que determinaron el desarrollo de este paquete. A continuación, describimos y demostramos la funcionalidad disponible en palaeoverse, además de proporcionar ejemplos de utilización. Finalmente, discutimos los recursos disponibles para la comunidad y los planes futuros del proyecto palaeoverse.
4. palaeoverse es un paquete colaborativo de R en paleobiología, desarrollado con la intención de unir a paleobiólogos y establecer normas consensuadas para una investigación cuantitativa de alta calidad. Este paquete proporciona una plataforma de fácil utilización para la preparación de datos para su análisis, con un código abierto bien documentado para aumentar la transparencia. La funcionalidad de palaeoverse mejora la reproducibilidad y accesibilidad del código, lo que es beneficioso tanto para el proceso de revisión como para futuras investigaciones.

Palabras clave

Paleobiología analítica, paleobiología computacional, programación en R, legible, reutilizable, reproducible.

Introduction

Since the development of large palaeontological datasets from the 1970s onwards, palaeontologists have increasingly adopted computational approaches to address questions about the history of life on Earth (Benton & Harper, 1999; J. J. Sepkoski, 1978). Today, most sub-disciplines within palaeontology regularly use large datasets to perform experiments *in silico*. This has initiated a ‘Golden Age’ of palaeontology (D. Sepkoski & Ruse, 2009), where extensive datasets of various formats are used to test macroevolutionary and macroecological hypotheses (R. A. Close et al., 2020; Mannion et al., 2014; e.g. Quental & Marshall, 2013; Zaffos et al., 2017). The growth and increasing availability of such datasets has made coding an integral part of palaeobiological research. Today, palaeobiologists commonly use code to clean (Flannery-Sutherland, Raja, et al., 2022; e.g. Zizka et al., 2019), analyse (e.g. Guillerme, 2018; Kocsis et al., 2019), and visualise data (e.g. Bell & Lloyd, 2015), as well as build models (e.g. Silvestro et al., 2014; Starrfelt & Liow, 2016) and implement simulations (Barido-Sottani et al., 2019; e.g. Fraser, 2017; Furness et al., 2021; Jones et al., 2021). Whilst software has been developed in languages such as C++ (e.g. Garwood et al., 2019) and Python (e.g. Silvestro et al., 2014), the programming language R is currently the most popular in palaeobiology. This is due to the wide range of tools—in the form of R packages—available to help users work with their data. Many of these tools are often borrowed or repurposed from ecology (e.g. Chao et al., 2014; Oksanen et al., 2020), while others have been developed to specifically handle fossil data (Kocsis et al., 2019; e.g. Lloyd, 2016).

In spite of the growth of analytical tools, few packages explicitly focus on preparing data for analyses, forcing users to construct custom scripts. This can result in distinct differences in code style and practices amongst the community, including code legibility and documentation. Accordingly, custom scripts can be inaccessible to other users (Filazzola & Lortie, 2022). Although increasingly requested by journals, code is also not always provided as supplementary material nor made available in online repositories (e.g. GitHub, Zenodo, Dryad). A lack of available code can lead to research results being unreproducible, preventing future studies from extending the work. Even when code is available, it might be poorly documented or written in a way that is specific to the dataset being analysed, and as such it may require extensive reworking before it can be applied to other data. Consequently, researchers are often forced to ‘reinvent the wheel’, putting time and effort into writing code that already exists, but is unavailable, inaccessible, and/or difficult to repurpose (Filazzola & Lortie, 2022). Such issues are exacerbated by the absence of community standards for how data should be prepared for analyses; differing approaches utilised by different researchers result in a lack of consistency between studies, making comparison between results challenging. Thus, there is a well-established need for both protocols and tools for preparing palaeontological data for further analysis.

Here, we introduce the R package `palaeoverse`, a community-driven toolkit for streamlining palaeobiological analyses and improving code accessibility and reproducibility. Our approach differs from other palaeontological R packages in that it aims to bring the palaeobiological community together to establish consensus on the steps taken in data preparation for analysis, and how these steps should be implemented. The package contains functions that align with current researcher needs to cleanse, prepare, and explore occurrence datasets for further analysis. These needs were established via a survey conducted by members of a new working group. The functionality of `palaeoverse` is purposefully flexible and can be applied to a wide variety of occurrence datasets. In this paper, we report results from the survey, describe and detail the functionality of `palaeoverse`, and illustrate its features with usage examples.

Community survey

To assess the needs of the palaeobiological community, we conducted an online survey. The survey was distributed via social media (Twitter) and email, and included questions related to researchers' previous experience, pre-existing code (to identify potential contributions), and what functionality they consider to be useful in a new palaeobiological toolkit. We summarise the types of data participants typically work with, the tasks commonly carried out when working with this data, and the tools they would like to have access to in Figure 1. We found that survey participants ($n = 35$) work with a wide range of data (Figure 1) and the checking and transformation of data is the most commonly employed task. A wide variety of functions were requested by survey participants, with data plotting, time binning, and data access commonly suggested (Figure 1). Over 40% of participants also indicated that they were willing to contribute code to `palaeoverse`, highlighting the potential for a community-driven project. Specific details regarding the survey and responses can be found in the Supplementary Material.

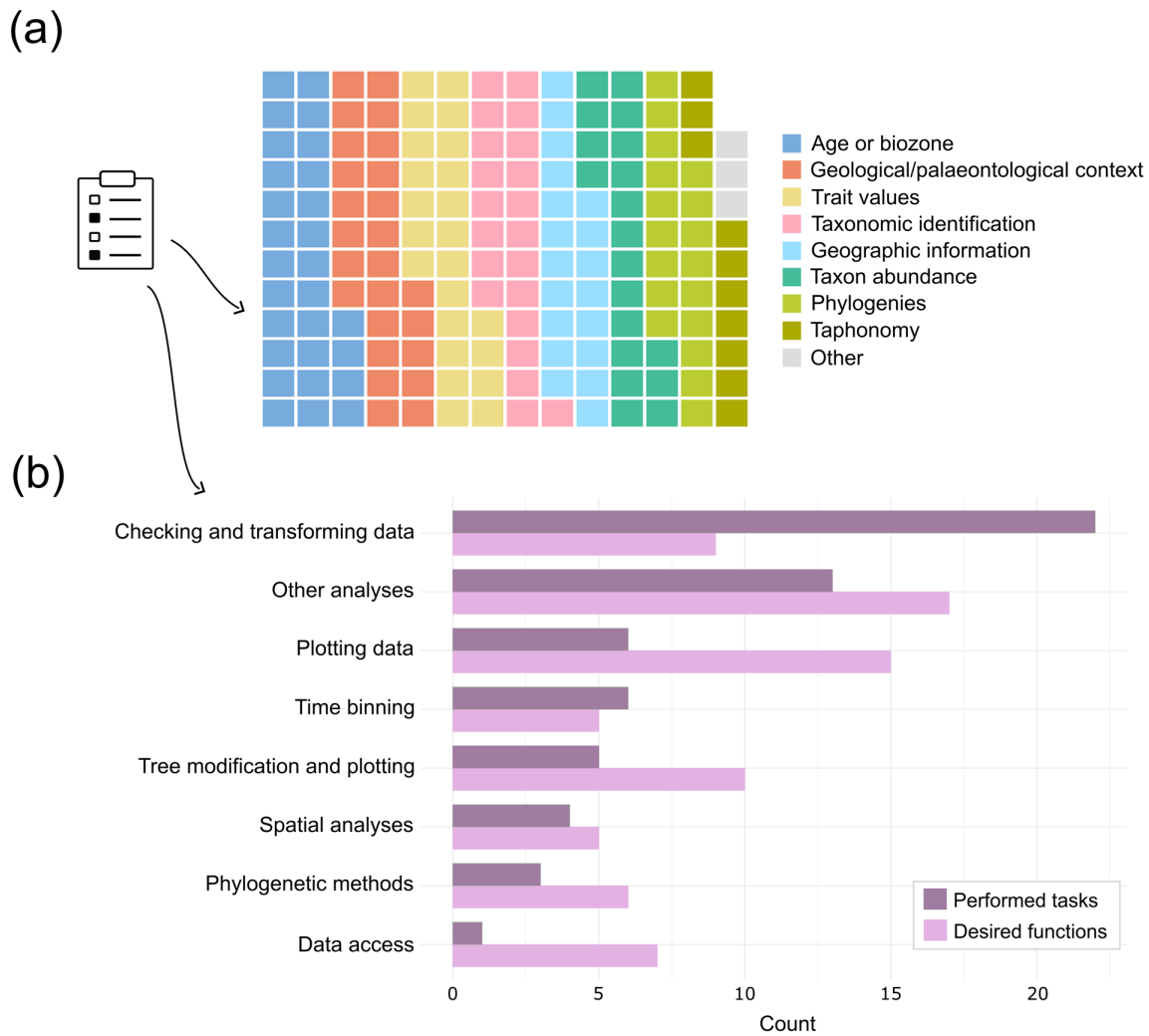


Figure 1: Summary of responses to the **palaeoverse** survey. (a) The types of palaeontological data that survey participants typically work with. Each box represents an individual check within a check-box list, in which participants could check multiple boxes. (b) Tasks that respondents routinely carry out in their own analyses (dark pink), and the functions they would find useful in the **palaeoverse** package (light pink).

Package description

After conducting the community survey, we combined participant input with our own experience to develop a toolkit for palaeobiologists, the **palaeoverse** R package. The package provides auxiliary functions to support data preparation and exploration for palaeobiological analysis. A summary of the functions currently available in **palaeoverse** is provided in Table 1, with further description provided in the

Features section. To demonstrate the functionality and versatility of the package, we also provide usage examples.

Installation

The `palaeoverse` package can be installed from CRAN using the `install.packages` function in R (R Core Team, 2022):

```
install.packages("palaeoverse")
```

If preferred, the development version of `palaeoverse` can be installed from GitHub via the `remotes` R package (Csárdi et al., 2021):

```
remotes::install_github("palaeoverse-community/palaeoverse")
```

Following installation, `palaeoverse` can be loaded via the `library` function in R:

```
library("palaeoverse")
```

Data

Functionality in `palaeoverse` was designed to be compatible with occurrence dataframes, such as those downloaded from the Paleobiology Database (<https://paleobiodb.org/#/>), the Geobiodiversity Database (<http://www.geobiodiversity.com>), or the Neptune Sandbox Berlin database (<https://nsb.mfn-berlin.de/>). Functionality is purposely flexible in `palaeoverse` and can be applied to various data sources with ease. In most instances, the returned object from a function is also a dataframe, which we consider the easiest data structure for most users to understand and work with. Although this might be undesirable for some advanced R users, transforming data structures should be straightforward for these users.

Functions

A summary of the functions available in `palaeoverse` along with their respective dependencies is provided in Table 1. All functions available in `palaeoverse` are novel in either their functionality or implementation, and collectively provide a flexible and versatile toolkit for palaeobiological research. Detailed descriptions of the functions are provided herein.

Table 1: A summary table of the functions currently available in the `palaeoverse` R package and respective dependencies. Base R dependencies are highlighted with an asterisk.

Function	Description	Dependency
<code>axis_geo</code>	Add a geological time scale axis to a plot	<code>deeptime</code> (Gearty, 2023)

Function	Description	Dependency
bin_lat	Bin fossil occurrences into latitudinal bins	-
bin_space	Bin fossil occurrences into spatial bins	h3jsr (O'Brien, 2023), sf (Pebesma, 2018)
bin_time	Bin fossil occurrences into time bins (choice of approaches)	stats* (R Core Team, 2022)
data	Datasets: 'tetrapods', 'reefs', 'interval_key', 'GTS2012', and 'GTS2020'	-
group_apply	Apply a function over user-defined groups	stats* (R Core Team, 2022)
lat_bins	Generate latitudinal bins	graphics* (R Core Team, 2022)
look_up	Link user-specified interval names to the International Geological Time Scale	-
palaeorotate	Reconstruct the palaeogeographic coordinates of fossil occurrences	curl (Ooms, 2023), geosphere (Hijmans, 2022), http (Wickham, 2022), h3jsr (O'Brien, 2023), pbapply (Solymos & Zawadzki, 2023), sf (Pebesma, 2018), stats* (R Core Team, 2022), utils* (R Core Team, 2022)
phylo_check	Check taxon names against tips in a phylogeny and/or remove tips from the tree	ape (Paradis & Schliep, 2019)
tax_check	Check for spelling mistakes in taxon names and flag potential issues	stats* (R Core Team, 2022), stringdist (van der Loo, 2014)
tax_range_space	Calculate the geographic range of taxa (choice of approaches)	geosphere (Hijmans, 2022), grDevices (R Core Team, 2022), h3jsr (O'Brien, 2023)

Function	Description	Dependency
<code>tax_range_time</code>	Calculate and plot the temporal range of taxa	graphics* (R Core Team, 2022)
<code>tax_expand_lat</code>	Convert taxon latitudinal ranges to bin-level pseudo-occurrences	-
<code>tax_expand_time</code>	Convert taxon temporal ranges to interval-level pseudo-occurrences	-
<code>tax_unique</code>	Calculate the number of unique taxa in a dataset of occurrences	stats* (R Core Team, 2022)
<code>time_bins</code>	Generate stratigraphic time bins or near-equal length time bins	-

Example datasets

Two occurrence datasets (`tetrapods` and `reefs`) are provided in `palaeoverse` to enable reproducible examples within function documentation. The `tetrapods` dataset is a compilation of Carboniferous–Early Triassic tetrapod occurrences ($n = 5,270$) from the Paleobiology Database. The dataset includes variables relevant to common palaeobiological analyses, covering the taxonomic identification of fossils and their geological, geographical and environmental context. The `reefs` dataset is a compilation of Phanerozoic reef occurrences ($n = 4,363$) from the PaleoReefs Database (Kiessling & Krause, 2022). This dataset includes information on the biological, geological, and geographical context of each reef. Except for the removal of superfluous columns and the renaming of some columns to improve clarity, both datasets are unaltered from their sources. Additional information on both datasets can be accessed via `?tetrapods` or `?reefs` once the package is loaded.

Time bins

We developed `time_bins` to enable access to two popular Geological Time Scales (GTS): GTS2012 and GTS2020 (Gradstein et al., 2012, 2020). Both GTS2012 and GTS2020 are included in the package as reference datasets. The `time_bins` function allows users to extract temporal bins at different temporal ranks (i.e. stage, epoch, period, era, or eon) using these datasets for a specified interval input:

```

189 # Get stage-level time bins
190 time_bins(interval = "Phanerozoic", rank = "stage", plot = TRUE)

```

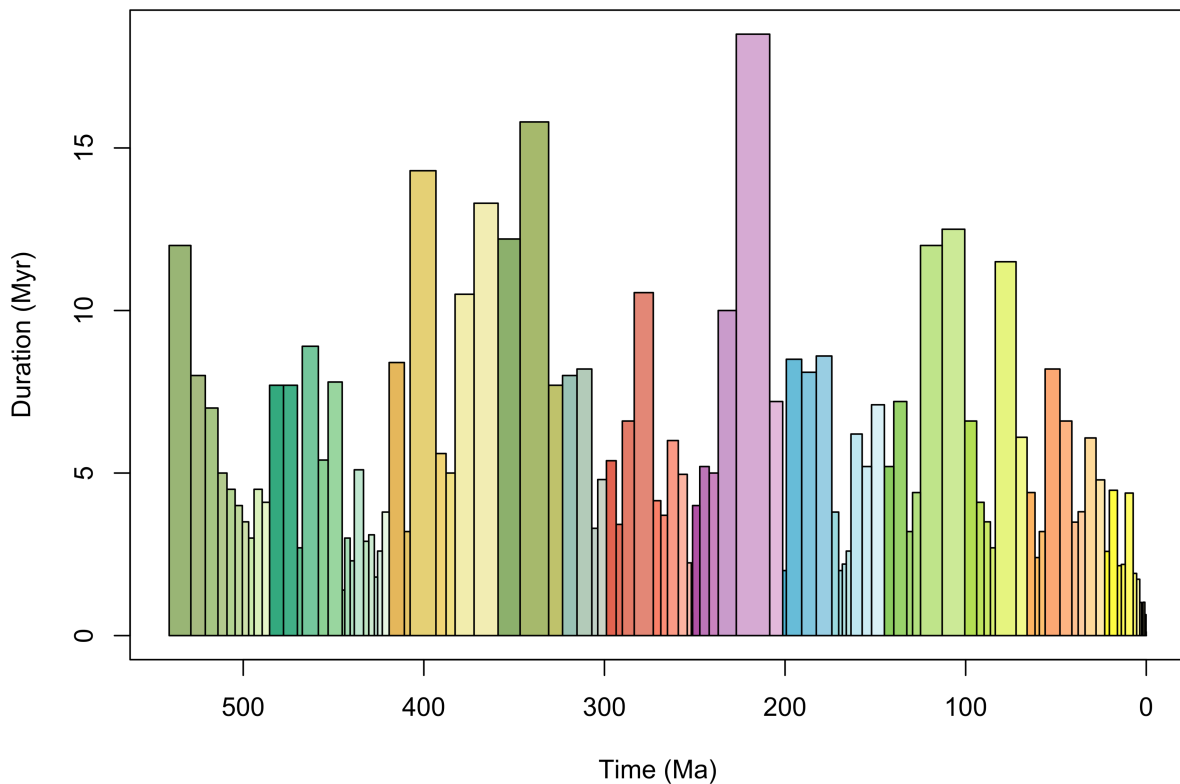


Figure 2: Phanerozoic stage-level time bins. Plot depicts the unevenness in duration of stratigraphic time bins. Bar colour filling follows the established colour scheme of the International Commission on Stratigraphy (<https://stratigraphy.org/>).

191 As is evident from Figure 2, GTS temporal bins are highly uneven in duration. Previous studies have
 192 attempted to circumvent this issue by generating near-equal-length time bins by grouping stages towards a
 193 target bin length (R. A. Close et al., 2020; e.g. Mannion et al., 2015). `time_bins` enables users to generate
 194 near-equal-length time bins following this approach (Figure 3) to a specified target size:

```

195 # Generate near-equal length time bins
196 time_bins(interval = "Phanerozoic", rank = "stage", size = 15, plot = TRUE)

```

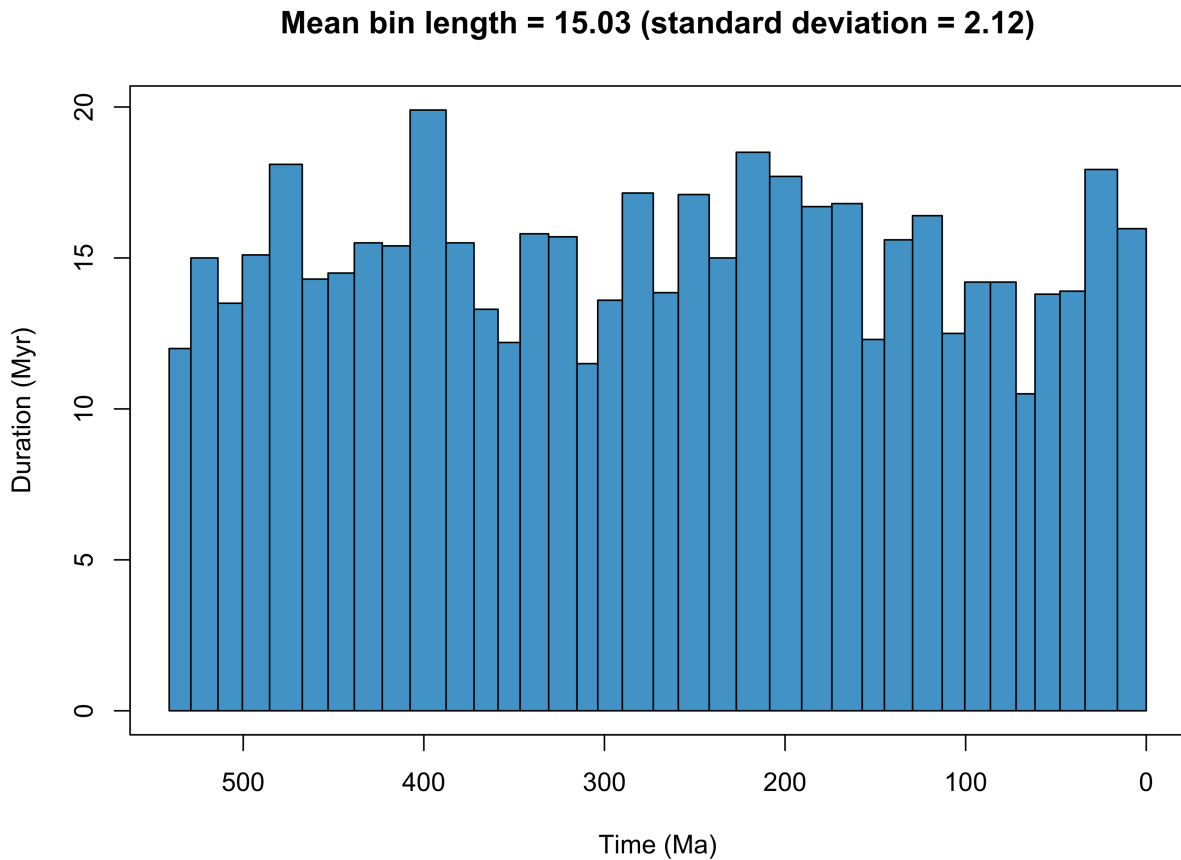


Figure 3: Phanerozoic near-equal-length time bins. Plot depicts composite stratigraphic bins (grouping stage-level bins) for the Phanerozoic of a target bin size of 15 million years. **Note:** time bins are still uneven but less so than stage-level bins.

Nevertheless, the appropriate set of time bins to use will depend upon the nature of subsequent analyses. Near-equal-length bins might be more desirable for calculating evolutionary rates through time, while GTS bins are defined on observed phenomena in the geological record, reflecting prior knowledge of cohesive biological units separated by some form of transition. Additional functionality in `time_bins` allows the user to assign occurrences to the generated bins if absolute ages are known (e.g. from radiometric dating). However, the bespoke `bin_time` function (discussed below) is likely to be the preferred option for most fossil occurrence data, which often have an age range.

Occurrence binning

Fossil occurrences are frequently ‘binned’ into distinct time intervals to enable quantification of changes (e.g. biodiversity or disparity) through geological time, as described by the respondents to our survey

207 (Figure 1). The function `bin_time` allows users to assign occurrences into time bins generated by the
208 function `time_bins`, or those defined by the user:

```
209 # Generate temporal bins  
210 bins <- time_bins()  
211 # Assign occurrences to bins  
212 bin_time(occdf = tetrapods, bins = bins, method = "mid")
```

213 Whilst binning occurrences with tightly defined temporal limits is straightforward and has been
214 implemented in other R packages (e.g. Lloyd, 2016), those with poorly constrained maximum and minimum
215 ages can span several intervals, and therefore cannot be easily assigned to a single bin. Palaeontologists
216 have identified numerous solutions to tackle this problem (Davies et al., 2017; Dean et al., 2020; Franeck
217 & Liow, 2020; e.g. Lloyd et al., 2012; Silvestro et al., 2016), but there is currently no consensus on the best
218 methodological approach or subsequent implementation. The `bin_time` function provides five approaches
219 defined by the ‘method’ argument: ‘mid’ (assigned based on the midpoint of the temporal range of the
220 occurrence), ‘majority’ (assigned to the bin which covers the majority of the temporal range of the
221 occurrence), ‘all’ (assigned to all bins within the temporal range of the occurrence), ‘random’ (assigned
222 randomly to bins with equal probability within the temporal range of the occurrence, repeated up to assigned
223 ‘reps’), and ‘point’ (assigned randomly using a user-defined probability distribution over the temporal range
224 of the occurrence, repeated up to assigned ‘reps’). We hope that formally including these options within the
225 `bin_time` function will encourage palaeontologists to routinely explore and compare the outcomes of
226 various binning approaches with ease.

227 In recent years, palaeobiologists have developed a heightened interest in the spatial structure of the fossil
228 record, with studies focused on understanding the spatial distribution of biodiversity and the processes that
229 drive them (Figure 1) (Antell et al., 2020; Chiarenza et al., 2022; R. Close et al., 2020; Flannery-Sutherland,
230 Silvestro, et al., 2022; Jones et al., 2022; Vilhena & Smith, 2013). In order to support such analyses,
231 `bin_space` has been developed for `palaeoverse`. The function allows the user to assign occurrence
232 data into equal-area grid cells using discrete hexagonal grids via the `h3jsr` package (O’Brien, 2023).
233 Additional functionality allows simultaneous assignation of occurrence data to cells of a finer-scale (i.e. a
234 ‘sub-grid’) within the primary grid. This might be desirable for users to evaluate differences in the amount
235 of area occupied by occurrences within their primary grid cells.

```
236 # Assign data to equal-area spatial bins  
237 bin_space(occdf = reefs, spacing = 250)  
238 bin_space(occdf = reefs, spacing = 250, sub_grid = 50)
```

239 Understanding the latitudinal distribution of biodiversity in deep time has also gained research interest in
240 recent years (Allen et al., 2020; Jones et al., 2021; Mannion et al., 2012, 2014; Powell, 2009; Song et al.,

2020). To ease implementation of such analyses, we have developed two functions, `lat_bins` and `bin_lat`, which can be used to generate latitudinal bins of a given size and assign occurrence data to those respective bins.

```
# Generate Latitudinal bins
bins <- lat_bins(size = 15)
# Assign occurrences to bins
bin_lat(occdf = tetrapods, bins = bins)
```

Palaeogeographic reconstruction

Using the present-day coordinates of fossil occurrences, plate rotation models can be used to reconstruct their location at the time of deposition. Existing fossil databases provide reconstructed coordinates for occurrences from only one or two of the many plate rotation models available (if any), and it is not always clear which model (or version of the model) has been used. This lack of transparency is reflected in some published articles that only cite the use of GPlates to reconstruct palaeocoordinates, yet lack specifics on which plate rotation model was used with the GPlates Web Service or desktop application (Müller et al., 2018). Furthermore, the uncertainty in palaeogeographic reconstructions is often underappreciated; reconstructed coordinates are treated as being well-established, rather than model-based estimates. Finally, online databases do not provide palaeocoordinates for all known samples. Both published and unpublished data (e.g. museum specimens) exists outside of online databases for which researchers might require palaeocoordinates.

We have developed the function `palaeorotate` to address these shortcomings. The function allows palaeocoordinates to be reconstructed within R using two different approaches: ‘point’ and ‘grid’. The first approach makes use of the GPlates Web Service and allows point data to be rotated to specific ages using the available models (see <https://gwsdoc.gplates.org>). The second approach uses reconstruction files of pre-generated palaeocoordinates to spatiotemporally link occurrences’ modern coordinates and age estimates with their respective palaeocoordinates. These reconstruction files were generated using an equal-area hexagonal grid (~100 km spacings) via the `h3jsr` package (O’Brien, 2023), and allow palaeocoordinates to be generated efficiently for large datasets. Furthermore, these reconstruction files allows the user to calculate the palaeolatitudinal range between reconstructed coordinates, as well as the great circle distance between the two most distant points (i.e. the palaeogeographic uncertainty). Finally, to encourage transparency in palaeobiological research, the function also reports additional information such as the plate rotation model used.

```
# Add midpoint age for rotation
tetrapods$age <- (tetrapods$max_ma + tetrapods$min_ma) / 2
```

```
274 # Palaeorotate occurrences and return uncertainty
275 palaeorotate(occdf = tetrapods, method = "grid", uncertainty = TRUE)
```

276 Taxon-related features

277 When working with large occurrence datasets, errors can easily creep into data. One frequently encountered
278 issue is spelling variations of the same taxon name. This can have undesirable consequences when
279 calculating metrics such as taxonomic richness or abundance. The `tax_check` function computes character
280 string distances between taxonomic names via the heuristic Jaro distance metric (Jaro, 1989). This metric
281 provides a measure of dissimilarity between character strings of 0 (exact match) to 1 (completely
282 dissimilar). During function call, the user defines a threshold for string dissimilarity to identify potential
283 synonyms. In `tax_check`, Jaro distances are calculated via the `stringdistmatrix` function from the
284 `stringdist` package (van der Loo, 2014). This function is provided to help researchers perform a spell
285 check on their dataset, with further similar functionality available in the `fossilbrush` package (Flannery-
286 Sutherland, Raja, et al., 2022). However, it should be made clear that this is no replacement for thorough
287 taxonomic vetting.

```
288 # Check for taxonomic errors
289 tax_check(taxdf = tetrapods, name = "genus")
```

290 The function `tax_unique` is provided to improve the accuracy of richness estimates from fossil occurrence
291 data. Palaeobiologists routinely discard occurrences not identified to their desired taxonomic resolution.
292 For example, if an analysis is conducted at species level, occurrences identified to the genus level (or above)
293 are discarded from the dataset. However, these occurrences can represent unique species, and their removal
294 can impact richness estimation. The `tax_unique` function reduces the number of unique taxa being
295 discarded by retaining fossils which are identified to a coarser taxonomic resolution than the desired level,
296 but must represent a clade not already in the filtered dataset. For instance, with three fossil occurrences
297 identified as *Tyrannosaurus rex*, *Spinosaurus aegyptiacus*, and *Diplodocidae* indet., the latter would be
298 discarded under species-level analysis (i.e. a species richness of two). However, this occurrence clearly
299 represents a different species to the two already present in the dataset. Using `tax_unique`, *Diplodocidae*
300 is treated as an additional species (i.e. a species richness of three) because this occurrence represents a
301 different species than the two already present in the dataset. Yet, the implementation is also conservative:
302 if multiple coarsely identified occurrences exist in the dataset, these are collapsed to the minimum number
303 of possible species (i.e. two occurrences of *Diplodocidae* indet. would be treated as only one species). This
304 method is similar to the ‘cryptic’ diversity measure introduced by Mannion et al. (2011).

```
305 # Evaluate unique taxa
306 tax_unique(occdf = tetrapods, genus = "genus", family = "family",
307            order = "order", class = "class", resolution = "genus")
```

Two functions exist in `palaeoverse` for computing taxon ranges. The first, `tax_range_time`, can be used to calculate and plot the temporal range of taxa. The function identifies all unique taxa provided in the occurrence dataframe and finds their first and last appearance dates. The second, `tax_range_space`, can be called to calculate the geographic range of taxa. This function allows the user to specify one of four different approaches (Darroch et al., 2020): (1) the area of a convex hull; (2) the (palaeo-)latitudinal range; (3) the maximum great-circle distance; and (4) the number and proportion of occupied equal-area grid cells. Similar to `tax_range_time`, the function will identify all unique taxa provided, and calculate these metrics based on the available occurrences of each taxon.

```
# Remove NA data
tetrapods <- subset(tetrapods, !is.na(order))
# Compute temporal range of orders
tax_range_time(occdf = tetrapods, name = "order", plot = TRUE)
```

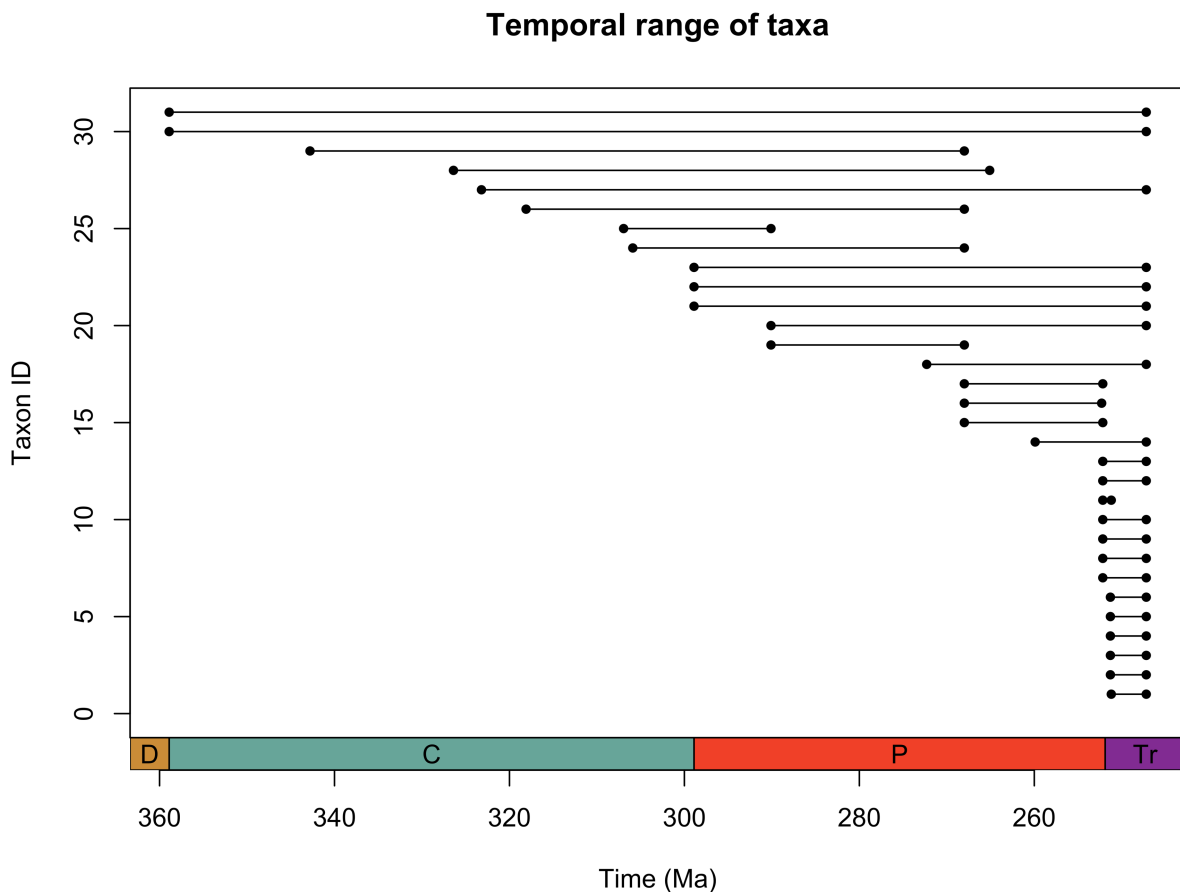


Figure 4: Temporal range of tetrapod orders in the `palaeoverse` example dataset.

```
# Compute latitudinal range of orders
tax_range_space(occdf = tetrapods, name = "order", method = "lat")
```

The provided `tax_expand_time` and `tax_expand_lat` functions are complementary to the taxonomic range functions. They convert temporal or latitudinal range data to bin-level pseudo-occurrences. These pseudo-occurrences serve to fill in ghost ranges, in which a taxon is presumed to be present, but no record exists. While these pseudo-occurrences should not be treated as equivalent to actual occurrence data, such data can be useful for performing statistical analyses where bin-level data is required.

Phylogeny wrangling

The function `phylo_check` compares a list of taxonomic names to the list of tip names in a user-provided phylogeny using the `ape` package (Paradis & Schliep, 2019). This comparison can be provided as a table describing the presence or absence of each taxon in the list and/or tips, or as counts of taxa present only in the list, only in the phylogeny, or in both. The function can also be used to trim the phylogeny to only include branches whose tip names are included within the list of taxonomic names.

Additional features

Datasets are frequently explored within groups in palaeobiology, such as time bins, collections or regions. The `group_apply` function has been included to allow users to run functions over a single, or multiple grouping variables, with ease.

```
# Compute the number of occurrences per collection  
group_apply(occdf = tetrapods, group = "collection_no", fun = nrow)
```

A common difficulty faced by palaeontologists is that the temporal information associated with fossil occurrence data is often asynchronous, and not directly comparable. Temporal data may be provided as either character-based interval names or numeric ages, and might conform to different time scales (e.g. international geological stages, or North American land mammal ages). Although interval names tend to be relatively stable over time, numerical age estimates are frequently updated with improved dating techniques, or the collection of new data. Consequently, where possible, interval names should be used to correlate occurrences from different stratigraphic time scales. The `look_up` function is provided to help assign a common time scale—typically international stages—to occurrence data. This is achieved with a user-defined table that links chosen interval names to corresponding stages on a common time scale (see example dataset `interval_key`). Numerical ages for the assigned stages can be provided by the user, or looked up in `GTS2012` or `GTS2020` (the default). This functionality therefore enables numerical ages to be assigned to datasets only containing character-based interval names (e.g. “Maastrichtian”).

```
reefs <- look_up(occdf = reefs,  
                early_interval = "interval",  
                late_interval = "interval",  
                int_key = interval_key)
```


355 Finally, a common feature request from our survey (Figure 1) was the ability to add the ‘Geological Time
 356 Scale’ to time-series plots in base R, with similar behaviour to the `deeptime` R package (Gearty, 2023) for
 357 `ggplot2` (Wickham, 2016). To address this request, the `axis_geo` function has been developed for the
 358 `palaeoverse` package (Figure 5).

```
359 # Palaeorotate reef dataset
360 reefs <- palaeorotate(occdf = reefs, age = "interval_mid_ma")

361 Warning in palaeorotate(occdf = reefs, age = "interval_mid_ma"): Palaeocoordi
362 nates could not be reconstructed for all points.
363 Either assigned plate does not exist at time of reconstruction or the plate r
364 otation model(s) does not cover the age of reconstruction.

365 # Plot palaeolatitudinal distribution through time
366 plot(x = reefs$interval_mid_ma, y = reefs$p_lat,
367      xlab = "Time (Ma)", ylab = "Palaeolatitude (°)",
368      xlim = c(541, 0), xaxt = "n", type = "p", pch = 20)
369 # Add Geological Time Scale
370 axis_geo(side = 1, intervals = "periods")
```

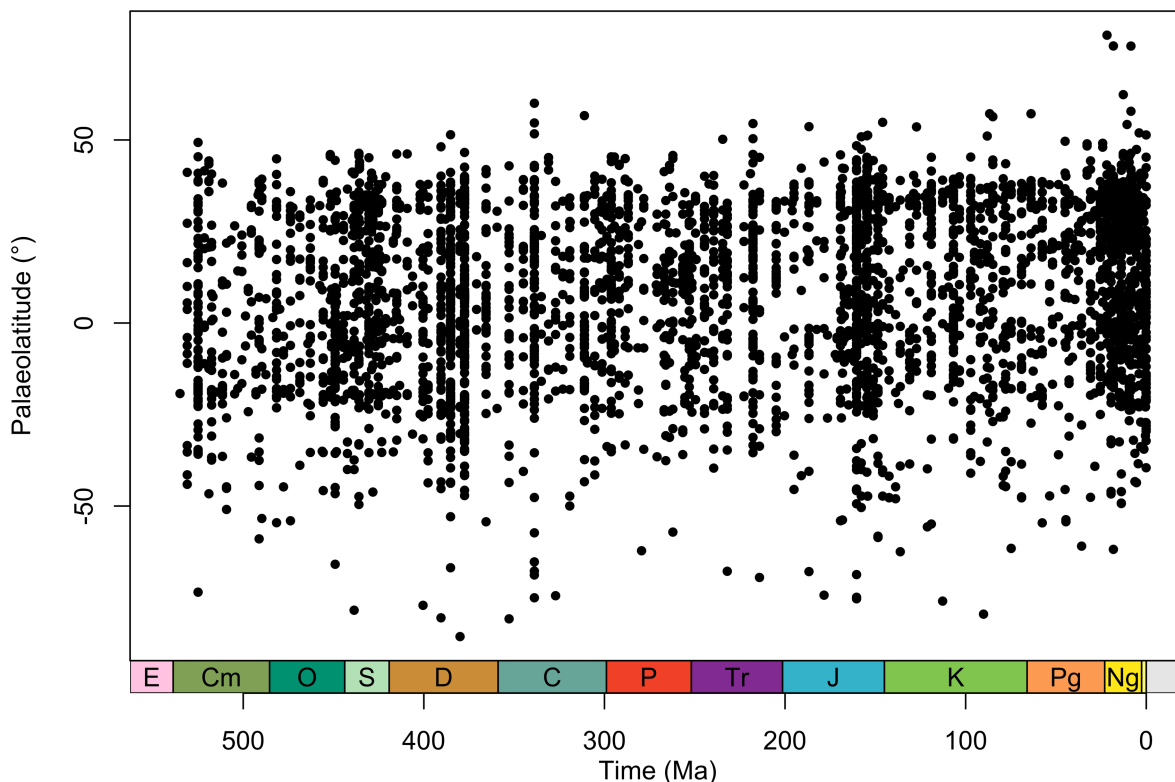


Figure 5: Example Phanerozoic plot of the palaeolatitudinal distribution of reefs through time. The plot demonstrates the usage of the `axis_geo` function for adding the Geological Time Scale to a base R plot.

Resources

To support the aims and use of `palaeoverse`, we have made several resources available to the palaeobiological community. Firstly, we have built a package website (<http://palaeoverse.palaeoverse.org>) which provides information on how to contribute to `palaeoverse`, how to report issues and bugs, and a general community code of conduct. Secondly, we have established a Google Group to foster collaboration and discussion on the issues faced by the community, such as establishing standards on data preparation (<https://groups.google.com/g/palaeoverse>).

Future perspectives

Palaeoverse is envisioned as a community project. While the initial development of the `palaeoverse` R package was led by the authors of this manuscript, it was also informed by the perspectives of 35 additional researchers (survey participants). Our hope is that `palaeoverse` will evolve into a community-driven package by welcoming contributions from the wider palaeontological community to broaden available functionality. To support this aim, we provide guidance on how the community can contribute to `palaeoverse` on the package website (<http://palaeoverse.palaeoverse.org>). Our working group also has the wider aim of establishing community standards and consensus in computational palaeobiological research and facilitating comparisons across studies. Through the `palaeoverse` R package, we hope to assist in making code more familiar and readable to fellow researchers, prevent researchers from ‘reinventing the wheel’ for common procedures, and improve the overall reproducibility of research through the use of computational tools which have been vetted and accepted by the broader community.

The development of the `palaeoverse` R package marks an initial effort to both streamline palaeobiological analysis pipelines and unite the computational palaeobiology community. Future efforts will see the expansion of the `palaeoverse` ‘universe’ with the development of Shiny applications to support non-R users and teaching exercises, tutorials to offer guidance for new researchers, and workshops to provide practical experience. In turn, we hope these efforts foster collaboration and the sharing of resources within the palaeobiology community. Finally, we warmly welcome the community to join these efforts and have established a community space accordingly to help facilitate the process (<https://groups.google.com/g/palaeoverse>).

Acknowledgements

The authors are extremely grateful to all survey respondents who helped to shape the development of `palaeoverse`. Special thanks are given to Emma M. Dunne whom participated in numerous discussions, and shared her experience with the development team. Thanks are also given to two anonymous reviewers that helped improve this manuscript. The contributions of LAJ, SG, and AAC were supported by the European Research Council under the European Union’s Horizon 2020 research and innovation program (grant agreement 947921; MAPAS project). LAJ was also supported by a Juan de la Cierva-formación 2021 fellowship (FJC2021-046695-I/MCIN/AEI/10.13039/501100011033) from the European Union “NextGenerationEU”/PRTR. AAC was also supported by a Juan de la Cierva-formación 2020 fellowship (FJC2020-044836-I/MCIN/AEI/10.13039/501100011033) from the European Union “NextGenerationEU”/PRTR. The contributions of WG were supported by the Population Biology Program of Excellence Postdoctoral Fellowship from the University of Nebraska-Lincoln School of Biological Sciences and the Lerner-Gray Postdoctoral Research Fellowship from the Richard Gilder Graduate School at the American Museum of Natural History. The contributions of BJA were supported by an ETH+ grant (BECCY). The contributions of CDD (RF_ERE_210013), MK (RGF_EA_180318) and CN (RGF_R1_180020) were supported by Royal Society grants. The contributions of PLG were supported by a FAPESP postdoctoral grant (2022/05697-9). This is Paleobiology Database publication no 450.

Conflict of Interest

We declare we have no conflict of interest.

Authors’ contributions

Lewis Jones conceived the project. All authors contributed to developing the project. Lewis A. Jones, Bethany J. Allen, William Gearty, Kilian Eichenseer, Christopher D. Dean, and Joseph Flannery-Sutherland contributed the code. All authors contributed to testing and reviewing the code. SG processed the survey results and produced the survey figures. All authors contributed to writing the manuscript.

Data accessibility

The `palaeoverse` R package is hosted on CRAN (<https://cran.r-project.org/web/packages/palaeoverse/>) and is available on GitHub (<https://github.com/palaeoverse-community/palaeoverse>). The code is also

archived in Zenodo through continuous integration (Jones et al., 2023). All example datasets are bundled with the R package. All code is released under a GPL (≥ 3) license.

References

- Allen, B. J., Wignall, P. B., Hill, D. J., Saupe, E. E., & Dunhill, A. M. (2020). The latitudinal diversity gradient of tetrapods across the permo-triassic mass extinction and recovery interval. *Proceedings of the Royal Society B*, 287(1929), 20201125.
- Antell, G. S., Kiessling, W., Aberhan, M., & Saupe, E. E. (2020). Marine biodiversity and geographic distributions are independent on large scales. *Current Biology*, 30(1), 115–121.e5. <https://doi.org/10.1016/j.cub.2019.10.065>
- Barido-Sottani, J., Pett, W., O'Reilly, J. E., & Warnock, R. C. (2019). FossilSim: An r package for simulating fossil occurrence data under mechanistic models of preservation and recovery. *Methods in Ecology and Evolution*, 10(6), 835–840.
- Bell, M. A., & Lloyd, G. T. (2015). Strap: An r package for plotting phylogenies against stratigraphy and assessing their stratigraphic congruence. In *Palaeontology* (Vol. 58, pp. 379–389). Wiley Online Library.
- Benton, M. J., & Harper, D. (1999). The history of life: Large databases in palaeontology. *Numerical Palaeobiology*, 249–283.
- Chao, A., Gotelli, N. J., Hsieh, T. C., Sande, E. L., Ma, K. H., Colwell, R. K., & Ellison, A. M. (2014). Rarefaction and extrapolation with hill numbers: A framework for sampling and estimation in species diversity studies. *Ecological Monographs*, 84, 45–67.
- Chiarenza, A. A., Mannion, P. D., Farnsworth, A., Carrano, M. T., & Varela, S. (2022). Climatic constraints on the biogeographic history of mesozoic dinosaurs. *Current Biology*, 32(3), 570–585.
- Close, R. A., Benson, R. B. J., Alroy, J., Carrano, M. T., Cleary, T. J., Dunne, E. M., Mannion, P. D., Uhen, M. D., & Butler, R. J. (2020). The apparent exponential radiation of phanerozoic land vertebrates is an artefact of spatial sampling biases. *Proceedings of the Royal Society B: Biological Sciences*, 287(1924), 20200372. <https://doi.org/10.1098/rspb.2020.0372>
- Close, R., Benson, R. B., Saupe, E., Clapham, M., & Butler, R. (2020). The spatial structure of phanerozoic marine animal diversity. *Science*, 368(6489), 420–424.
- Csárdi, G., Hester, J., Wickham, H., Chang, W., Morgan, M., & Tenenbaum, D. (2021). *Remotes: R package installation from remote repositories, including 'GitHub'*. <https://CRAN.R-project.org/package=remotes>
- Darroch, S. A., Casey, M. M., Antell, G. S., Sweeney, A., & Saupe, E. E. (2020). High preservation potential of paleogeographic range size distributions in deep time. *The American Naturalist*, 196(4), 454–471.
- Davies, T. W., Bell, M. A., Goswami, A., & Halliday, T. J. (2017). Completeness of the eutherian mammal fossil record and implications for reconstructing mammal evolution through the cretaceous/paleogene mass extinction. *Paleobiology*, 43(4), 521–536.

461 Dean, C. D., Chiarenza, A. A., & Maidment, S. C. (2020). Formation binning: A new method for
 462 increased temporal resolution in regional studies, applied to the late cretaceous dinosaur fossil record of
 463 north america. *Palaeontology*, 63(6), 881–901.

464 Filazzola, A., & Lortie, C. (2022). A call for clean code to effectively communicate science. *Methods in*
 465 *Ecology and Evolution*, 13(10), 2119–2128. <https://doi.org/10.1111/2041-210X.13961>

466 Flannery-Sutherland, J. T., Raja, N. B., Kocsis, Á. T., & Kiessling, W. (2022). Fossilbrush: An r package
 467 for automated detection and resolution of anomalies in palaeontological occurrence data. *Methods in*
 468 *Ecology and Evolution*, 13(11), 2404–2418. <https://doi.org/10.1111/2041-210X.13966>

469 Flannery-Sutherland, J. T., Silvestro, D., & Benton, M. J. (2022). Global diversity dynamics in the fossil
 470 record are regionally heterogeneous. *Nature Communications*, 13(1), 1–17.

471 Franeck, F., & Liow, L. H. (2020). Did hard substrate taxa diversify prior to the great ordovician
 472 biodiversification event? *Palaeontology*, 63(4), 675–687.

473 Fraser, D. (2017). Can latitudinal richness gradients be measured in the terrestrial fossil record?
 474 *Paleobiology*, 43(3), 479–494.

475 Furness, E. N., Garwood, R. J., Mannion, P. D., & Sutton, M. D. (2021). Evolutionary simulations clarify
 476 and reconcile biodiversity-disturbance models. *Proceedings of the Royal Society B*, 288(1949), 20210240.

477 Garwood, R. J., Spencer, A. R., & Sutton, M. D. (2019). REvoSim: Organism-level simulation of macro
 478 and microevolution. *Palaeontology*, 62(3), 339–355.

479 Gearty, W. (2023). *Deeptime: Plotting tools for anyone working in deep time*. [https://CRAN.R-](https://CRAN.R-project.org/package=deeptime)
 480 [project.org/package=deeptime](https://CRAN.R-project.org/package=deeptime)

481 Gradstein, F. M., Ogg, J. G., Schmitz, M. D., & Ogg, G. M. (2020). *Geologic time scale 2020*. Elsevier.

482 Gradstein, F. M., Ogg, J. G., Schmitz, M., & Ogg, G. (2012). *The geologic time scale 2012*. Elsevier.

483 Guillaume, T. (2018). dispRity: A modular r package for measuring disparity. *Methods in Ecology and*
 484 *Evolution*, 9(7), 1755–1763.

485 Hijmans, R. J. (2022). *Geosphere: Spherical trigonometry*. [https://CRAN.R-](https://CRAN.R-project.org/package=geosphere)
 486 [project.org/package=geosphere](https://CRAN.R-project.org/package=geosphere)

487 Jaro, M. A. (1989). Advances in record-linkage methodology as applied to matching the 1985 census of
 488 tampa, florida. *Journal of the American Statistical Association*, 84(406), 414–420.

489 Jones, L. A., Dean, C. D., Mannion, P. D., Farnsworth, A., & Allison, P. A. (2021). Spatial sampling
 490 heterogeneity limits the detectability of deep time latitudinal biodiversity gradients. *Proceedings of the*
 491 *Royal Society B*, 288(1945), 20202762.

492 Jones, L. A., Gearty, W., Eichenseer, K., Dean, C., Allen, B., & Flannery-Sutherland, J. (2023).
 493 *Palaeoverse-community/palaeoverse: v.1.1.1* (Version v.1.1.1) [Computer software]. Zenodo.
 494 <https://doi.org/10.5281/zenodo.7728639>

495 Jones, L. A., Mannion, P. D., Farnsworth, A., Bragg, F., & Lunt, D. J. (2022). Climatic and tectonic
 496 drivers shaped the tropical distribution of coral reefs. *Nature Communications*, 13(1), 1–10.

497 Kiessling, W., & Krause, C. (2022). *PaleoReefs database (PARED)* (Version 1.0) [Data set]. Zenodo.
 498 <https://doi.org/10.5281/zenodo.6037852>

499 Kocsis, Á. T., Reddin, C. J., Alroy, J., & Kiessling, W. (2019). The r package divDyn for quantifying
 500 diversity dynamics using fossil sampling data. *Methods in Ecology and Evolution*, 10(5), 735–743.

501 Lloyd, G. T. (2016). Estimating morphological diversity and tempo with discrete character-taxon
 502 matrices: Implementation, challenges, progress, and future directions. *Biological journal of the linnean*
 503 *society. Biological Journal of the Linnean Society*, 118, 131–151.

504 Lloyd, G. T., Pearson, P. N., Young, J. R., & Smith, A. B. (2012). Sampling bias and the fossil record of
 505 planktonic foraminifera on land and in the deep sea. *Paleobiology*, 38(4), 569–584.

506 Mannion, P. D., Benson, R. B., Carrano, M. T., Tennant, J. P., Judd, J., & Butler, R. J. (2015). Climate
 507 constrains the evolutionary history and biodiversity of crocodylians. *Nature Communications*, 6(1), 1–9.

508 Mannion, P. D., Benson, R. B., Upchurch, P., Butler, R. J., Carrano, M. T., & Barrett, P. M. (2012). A
 509 temperate palaeodiversity peak in mesozoic dinosaurs and evidence for late cretaceous geographical
 510 partitioning. *Global Ecology and Biogeography*, 21(9), 898–908.

511 Mannion, P. D., Upchurch, P., Benson, R. B., & Goswami, A. (2014). The latitudinal biodiversity
 512 gradient through deep time. *Trends in Ecology & Evolution*, 29(1), 42–50.

513 Mannion, P. D., Upchurch, P., Carrano, M. T., & Barrett, P. M. (2011). Testing the effect of the rock
 514 record on diversity: A multidisciplinary approach to elucidating the generic richness of sauropodomorph
 515 dinosaurs through time. *Biological Reviews*, 86(1), 157–181. [https://doi.org/10.1111/j.1469-](https://doi.org/10.1111/j.1469-185X.2010.00139.x)
 516 [185X.2010.00139.x](https://doi.org/10.1111/j.1469-185X.2010.00139.x)

517 Müller, R. D., Cannon, J., Qin, X., Watson, R. J., Gurnis, M., Williams, S., Pfaffelmoser, T., Seton, M.,
 518 Russell, S. H. J., & Zahirovic, S. (2018). GPlates: Building a virtual earth through deep time.
 519 *Geochemistry, Geophysics, Geosystems*, 19(7), 2243–2261. <https://doi.org/10.1029/2018GC007584>

520 O’Brien, L. (2023). *h3jsr: Access uber’s H3 library*. <https://CRAN.R-project.org/package=h3jsr>

521 Oksanen, J., Blanchet, F. G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., Minchin, P. R., O’Hara,
 522 R. B., Simpson, G. L., Solymos, P., Stevens, M. H. H., Szoecs, E., & Wagner, H. (2020). *Vegan:*
 523 *Community ecology package*. <https://CRAN.R-project.org/package=vegan>

524 Ooms, J. (2023). *Curl: A modern and flexible web client for r*. <https://CRAN.R-project.org/package=curl>

525 Paradis, E., & Schliep, K. (2019). Ape 5.0: An environment for modern phylogenetics and evolutionary
 526 analyses in R. *Bioinformatics*, 35, 526–528.

527 Pebesma, E. (2018). Simple Features for R: Standardized Support for Spatial Vector Data. *The R Journal*,
 528 10(1), 439–446. <https://doi.org/10.32614/RJ-2018-009>

529 Powell, M. G. (2009). The latitudinal diversity gradient of brachiopods over the past 530 million years.
 530 *The Journal of Geology*, 117(6), 585–594.

531 Quental, T. B., & Marshall, C. R. (2013). How the red queen drives terrestrial mammals to extinction.
 532 *Science*, 341(6143), 290–292.

533 R Core Team. (2022). *R: A language and environment for statistical computing*. R Foundation for
 534 Statistical Computing. <https://www.R-project.org/>

535 Sepkoski, D., & Ruse, M. (2009). *The paleobiological revolution: Essays on the growth of modern*
536 *paleontology*. University of Chicago Press.

537 Sepkoski, J. J. (1978). A kinetic model of phanerozoic taxonomic diversity i. Analysis of marine orders.
538 *Paleobiology*, 4(3), 223–251.

539 Silvestro, D., Salamin, N., & Schnitzler, J. (2014). PyRate: A new program to estimate speciation and
540 extinction rates from incomplete fossil data. *Methods in Ecology and Evolution*, 5(10), 1126–1131.

541 Silvestro, D., Zizka, A., Bacon, C. D., Cascales-Minana, B., Salamin, N., & Antonelli, A. (2016). Fossil
542 biogeography: A new model to infer dispersal, extinction and sampling from palaeontological data.
543 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371(1691), 20150225.

544 Solymos, P., & Zawadzki, Z. (2023). *Pbapply: Adding progress bar to '*apply' functions*.
545 <https://CRAN.R-project.org/package=pbapply>

546 Song, H., Huang, S., Jia, E., Dai, X., Wignall, P. B., & Dunhill, A. M. (2020). Flat latitudinal diversity
547 gradient caused by the permian–triassic mass extinction. *Proceedings of the National Academy of*
548 *Sciences*, 117(30), 17578–17583.

549 Starrfelt, J., & Liow, L. H. (2016). How many dinosaur species were there? Fossil bias and true richness
550 estimated using a poisson sampling model. *Philosophical Transactions of the Royal Society B: Biological*
551 *Sciences*, 371(1691), 20150219.

552 van der Loo, M. P. J. (2014). The stringdist package for approximate string matching. *The R Journal*, 6,
553 111–122. <https://CRAN.R-project.org/package=stringdist>

554 Vilhena, D. A., & Smith, A. B. (2013). Spatial bias in the marine fossil record. *PLoS One*, 8(10), e74470.

555 Wickham, H. (2016). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
556 <https://ggplot2.tidyverse.org>

557 Wickham, H. (2022). *Httr: Tools for working with URLs and HTTP*. [https://CRAN.R-](https://CRAN.R-project.org/package=httr)
558 [project.org/package=httr](https://CRAN.R-project.org/package=httr)

559 Zaffos, A., Finnegan, S., & Peters, S. E. (2017). Plate tectonic regulation of global marine animal
560 diversity. *Proceedings of the National Academy of Sciences*, 114(22), 5653–5658.

561 Zizka, A., Silvestro, D., Andermann, T., Azevedo, J., Duarte Ritter, C., Edler, D., Farooq, H., Herdean,
562 A., Ariza, M., Scharn, R., et al. (2019). CoordinateCleaner: Standardized cleaning of occurrence records
563 from biological collection databases. *Methods in Ecology and Evolution*, 10(5), 744–751.