

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

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

7
8 ¹*Grupo de Ecoloxía Animal, Departamento de Ecoloxía e Bioloxía Animal, Universidade de Vigo, 36310*
9 *Vigo, Spain.*

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

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

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

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

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

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

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

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

21 ¹⁰*Département de Biologie, École Normale Supérieure de Lyon, Université Claude Bernard Lyon 1,*
22 *69342 Lyon Cedex 07, France.*

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

24 ¹²*Department of Ecology, Evolution, and Marine Biology, University of California, Santa Barbara, CA*
25 *93106, USA.*

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

27
28 **Corresponding author:** LewisAlan.Jones@uvigo.es

29

30

Abstract

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

55

Keywords

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

58

59 **Introduction**

60 Since the development of large palaeontological datasets from the 1970s onwards, palaeontologists have
61 increasingly adopted computational approaches to address questions about the history of life on Earth
62 (Sepkoski, 1978; Benton and Harper, 1999). Today, most sub-disciplines within palaeontology regularly
63 use large datasets to perform experiments *in silico*. This has initiated a ‘Golden Age’ of palaeontology
64 (Sepkoski and Ruse, 2009), where extensive datasets of various formats are used to test macroevolutionary
65 and macroecological hypotheses (Quental and Marshall, 2013; Mannion et al., 2014; Zaffos, Finnegan and
66 Peters, 2017; Close et al., 2020a). The growth and increasing availability of such datasets has made coding
67 an integral part of palaeobiological research. Today, palaeobiologists commonly use code to clean (Zizka
68 et al., 2019; Flannery-Sutherland et al., 2022), analyse (Guillerme, 2018; Kocsis et al., 2019), and visualise
69 data (Bell and Lloyd, 2015), as well as build models (Silvestro, Salamin and Schnitzler, 2014; Starrfelt and
70 Liow, 2016) and implement simulations (Fraser, 2017; Barido-Sottani et al., 2019; Furness et al., 2021;
71 Jones et al., 2021). Whilst software has been developed in languages such as C++ (Garwood, Spencer and
72 Sutton, 2019) and Python (Silvestro et al., 2014), the programming language R is currently the most popular
73 in palaeobiology. This is due to the wide range of tools—in the form of R packages—available to help users
74 work with their data. Many of these tools are often borrowed or repurposed from ecology (e.g. Chao et al.,
75 2014; Oksanen et al., 2020), while others have been developed to specifically handle fossil data (e.g. Lloyd,
76 2016; Kocsis et al., 2019).

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

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

101 **Community survey**

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

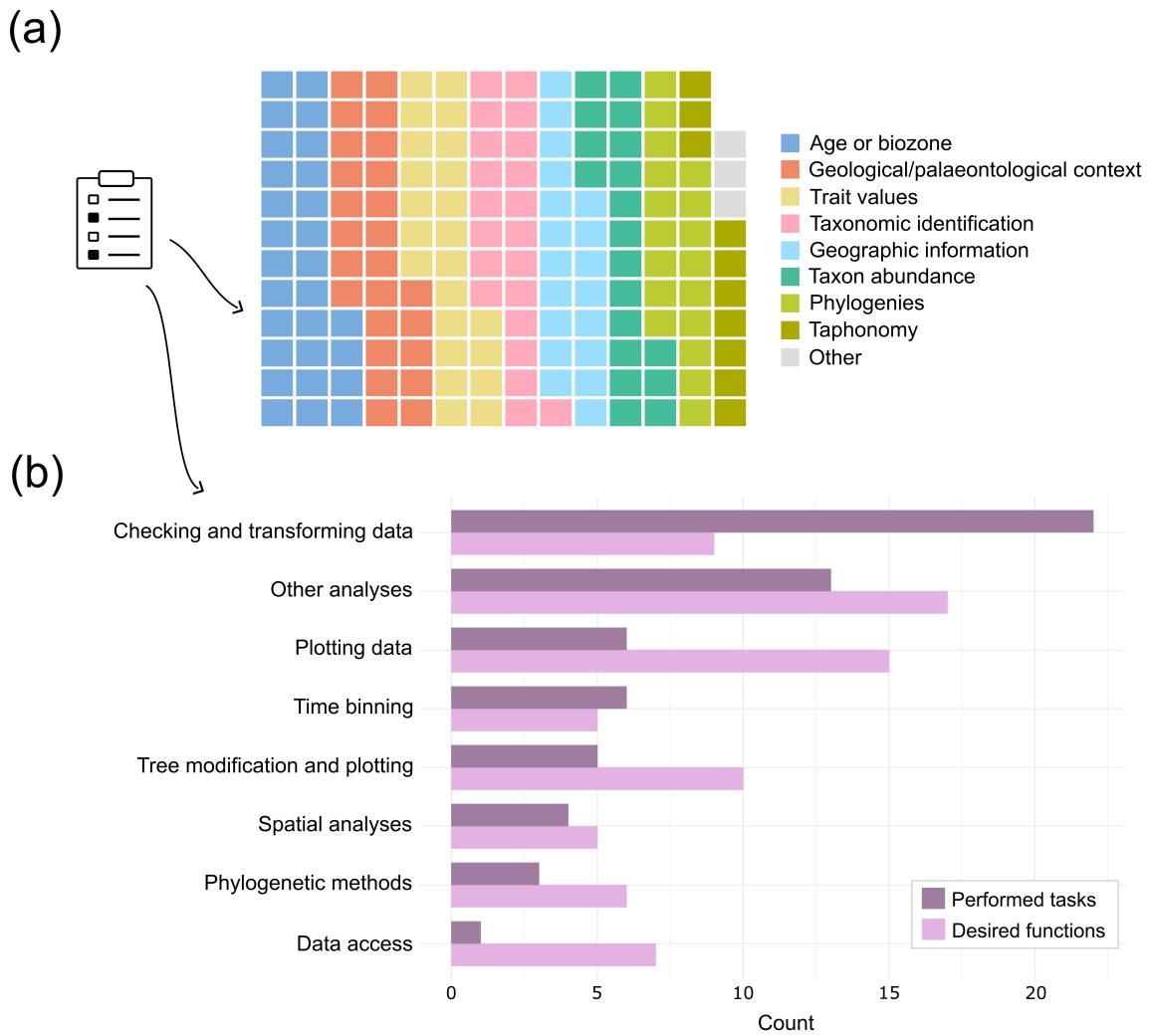


Figure 1: Summary of responses to the `palaeoverse` survey. (a) The types of palaeontological data that survey participants typically work with. Each box represents an individual check within a checkbox list, in which participants could check multiple boxes. (b) Tasks that respondents routinely carry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

113 **Package description**

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

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

120 **Installation**

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

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

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

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

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

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

129 **Data**

130 Functionality in `palaeoverse` was designed to be compatible with occurrence dataframes, such as those
131 downloaded from the Paleobiology Database (<https://paleobiodb.org/#/>), the Geobiodiversity Database
132 (<http://www.geobiodiversity.com>), or the Neptune Sandbox Berlin database (<https://nsb.mfn-berlin.de/>).

133 Functionality is purposely flexible in `palaeoverse` and can be applied to various data sources with ease.

134 In most instances, the returned object from a function is also a dataframe, which we consider the easiest
135 data structure for most users to understand and work with. Although this might be undesirable for some
136 advanced R users, transforming data structures should be straightforward for these users.

137 **Functions**

138 A summary of the functions available in `palaeoverse` is provided in Table 1. Detailed descriptions of the
139 functions are provided herein.

140 Table 1: A summary table of the functions currently available in the `palaeoverse` R package

Function	Description
<code>axis_geo</code>	Add a geological time scale axis to a plot
<code>bin_lat</code>	Bin fossil occurrences into latitudinal bins
<code>bin_space</code>	Bin fossil occurrences into spatial bins
<code>bin_time</code>	Bin fossil occurrences into time bins (choice of approaches)

Function	Description
data	Datasets: ‘tetrapods’, ‘reefs’, ‘interval_key’, ‘GTS2012’, and ‘GTS2020’
group_apply	Apply a function over user-defined groups
lat_bins	Generate latitudinal bins
look_up	Link user-specified interval names to the International Geological Time Scale
palaeorotate	Reconstruct the palaeogeographic coordinates of fossil occurrences
phylo_check	Check taxon names against tips in a phylogeny and/or remove tips from the tree
tax_check	Check for spelling mistakes in taxon names and flag potential issues
tax_range_space	Calculate the geographic range of taxa (choice of approaches)
tax_range_time	Calculate and plot the temporal range of taxa
tax_expand_lat	Convert taxon latitudinal ranges to bin-level pseudo-occurrences
tax_expand_time	Convert taxon temporal ranges to interval-level pseudo-occurrences
tax_unique	Calculate the number of unique taxa in a dataset of occurrences
time_bins	Generate stratigraphic time bins or near-equal length time bins

141 Example datasets

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

152 Time bins

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

```

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

```

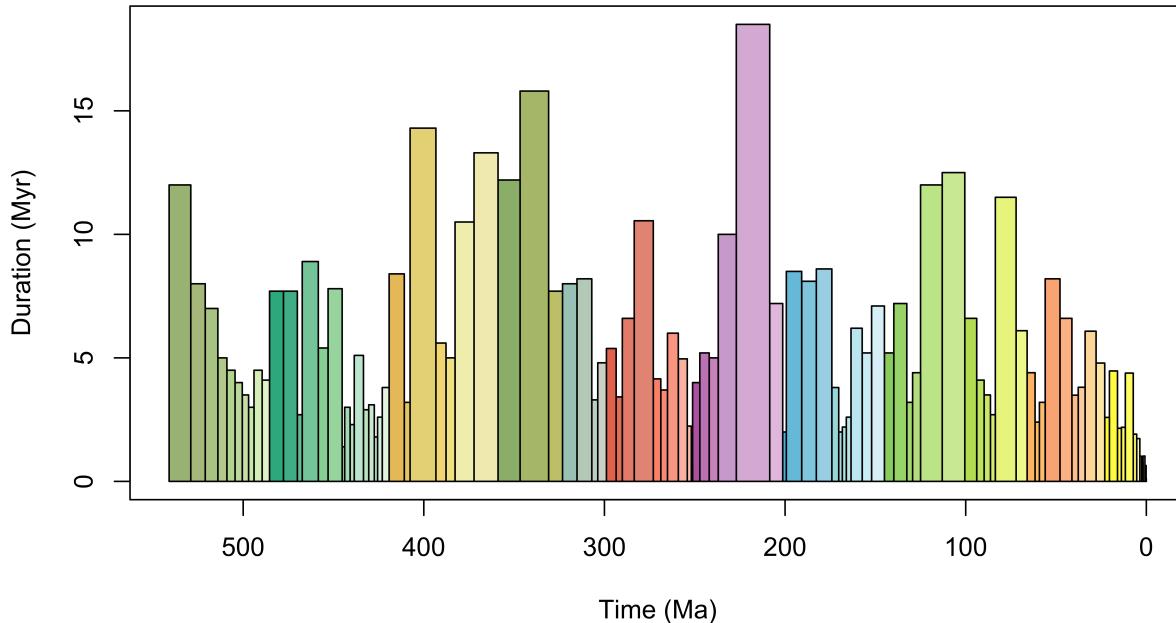


Figure 2: Phanerozoic stage-level time bins. Plot depicts the unevenness in duration of stratigraphic time bins. Bar colour filling follows the established colour scheme of the International Commission on Stratigraphy (<https://stratigraphy.org/>).

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

```

163 # Generate near-equal Length time bins
164 time_bins(interval = "Phanerozoic", rank = "stage", size = 15, plot = TRUE)

```

Mean bin length = 15.03 (standard deviation = 2.12)

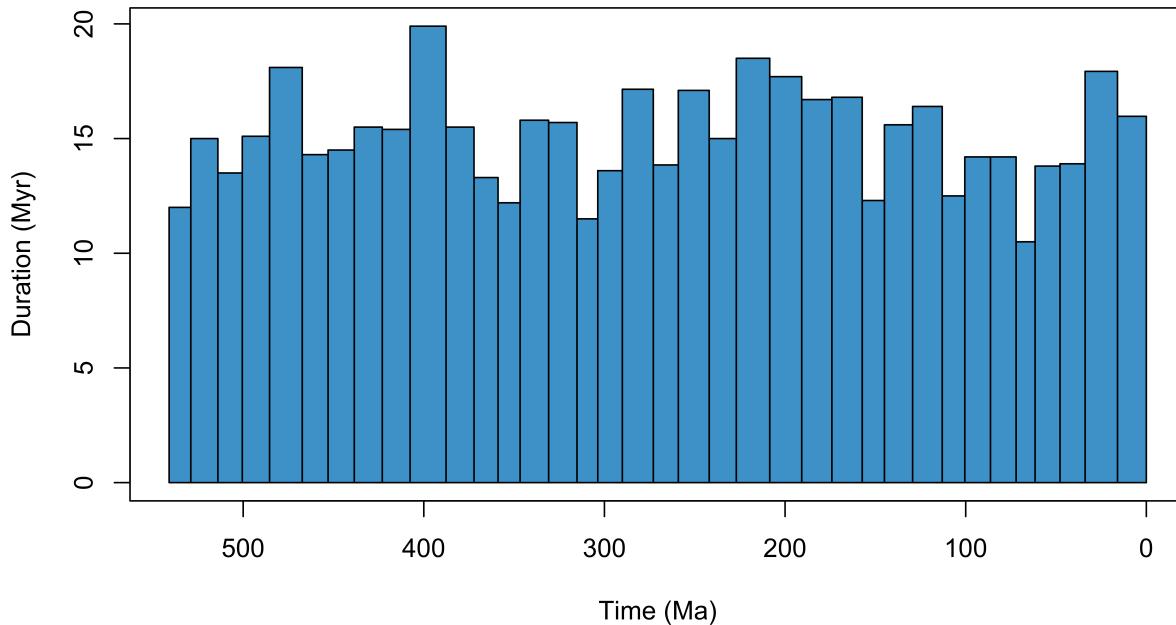


Figure 3: Phanerozoic near-equal-length time bins. Plot depicts composite stratigraphic bins (grouping stage-level bins) for the Phanerozoic of a target bin size of 15 million years. **Note:** time bins are still uneven but less so than stage-level bins.

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

172 Occurrence binning

173 Fossil occurrences are frequently ‘binned’ into distinct time intervals to enable quantification of changes
174 (e.g. biodiversity or disparity) through geological time. The function `bin_time` allows users to assign
175 occurrences into time bins generated by the function `time_bins`, or those defined by the user:

```
176 # Generate temporal bins
177 bins <- time_bins()
178 # Assign occurrences to bins
179 bin_time(occdf = tetrapods, bins = bins, method = "mid")
```

Whilst binning occurrences with tightly defined temporal limits is straightforward, those with poorly constrained maximum and minimum ages can span several intervals, and therefore cannot be easily assigned to a single bin. Palaeontologists have identified numerous solutions to tackle this problem (Lloyd et al., 2012; Silvestro et al., 2016; Davies et al., 2017; Dean, Chiarenza and Maidment, 2020; Franeck and Liow, 2020), but there is currently no consensus on the best methodological approach or subsequent implementation. The `bin_time` function provides five approaches defined by the ‘method’ argument: ‘mid’ (assigned based on the midpoint of the temporal range of the occurrence), ‘majority’ (assigned to the bin which covers the majority of the temporal range of the occurrence), ‘all’ (assigned to all bins within the temporal range of the occurrence), ‘random’ (assigned randomly to bins with equal probability within the temporal range of the occurrence, repeated up to assigned ‘reps’), and ‘point’ (assigned randomly from a uniform distribution over the temporal range of the occurrence, repeated up to assigned ‘reps’). We hope that formally including these options within the `bin_time` function will encourage palaeontologists to routinely explore and compare the outcomes of various binning approaches with ease.

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

```
202 # Assign data to equal-area spatial bins  
203 bin_space(occdf = reefs, spacing = 250)  
204 bin_space(occdf = reefs, spacing = 250, sub_grid = 50)
```

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

```
210 # Generate latitudinal bins  
211 bins <- lat_bins(size = 15)  
212 # Assign occurrences to bins  
213 bin_lat(occdf = tetrapods, bins = bins)
```

214 **Palaeogeographic reconstruction**

215 Using the present-day coordinates of fossil occurrences, plate rotation models can be used to reconstruct
216 their location at the time of deposition. Existing fossil databases provide reconstructed coordinates for
217 occurrences from only one or two of the many plate rotation models available (if any), and it is not always
218 clear which model (or version of the model) has been used. This lack of transparency is reflected in some
219 published articles that only cite the use of GPlates to reconstruct palaeocoordinates, yet lack specifics on
220 which plate rotation model was used with the GPlates Web Service or desktop application (Müller et al.,
221 2018). Furthermore, the uncertainty in palaeogeographic reconstructions is often underappreciated;
222 reconstructed coordinates are treated as being well-established, rather than model-based estimates. Finally,
223 online databases do not provide palaeocoordinates for all known samples. Both published and unpublished
224 data (e.g. museum specimens) exists outside of online databases for which researchers might require
225 palaeocoordinates.

226 We have developed the function `palaeorotate` to address these shortcomings. The function allows
227 palaeocoordinates to be reconstructed within R using two different approaches: ‘point’ and ‘grid’. The first
228 approach makes use of the GPlates Web Service and allows point data to be rotated to specific ages using
229 the available models (see <https://gwsdoc.gplates.org>). The second approach uses reconstruction files of pre-
230 generated palaeocoordinates to spatiotemporally link occurrences’ modern coordinates and age estimates
231 with their respective palaeocoordinates. These reconstruction files were generated using a $1^\circ \times 1^\circ$ spatial
232 grid and allows palaeocoordinates to be generated efficiently for large datasets. Furthermore, these
233 reconstruction files allows the user to calculate the palaeolatitudinal range between reconstructed
234 coordinates, as well as the great circle distance between the two most distant points (i.e. the
235 palaeogeographic uncertainty). Finally, to encourage transparency in palaeobiological research, the
236 function also reports additional information such as the plate rotation model used.

```
237 # Add midpoint age for rotation
238 tetrapods$age <- (tetrapods$max_ma + tetrapods$min_ma) / 2
239 # Palaeorotate occurrences and return uncertainty
240 palaeorotate(occdf = tetrapods, method = "grid", uncertainty = TRUE)
```

241 **Taxon-related features**

242 When working with large occurrence datasets, errors can easily creep into data. One frequently encountered
243 issue is spelling variations of the same taxon name. This can have undesirable consequences when
244 calculating metrics such as taxonomic richness or abundance. The `tax_check` function computes character
245 string distances between taxonomic names via the heuristic Jaro distance metric (Jaro, 1989). This metric
246 provides a measure of dissimilarity between character strings of 0 (exact match) to 1 (completely

247 dissimilar). During function call, the user defines a threshold for string dissimilarity to identify potential
248 synonyms. In `tax_check`, Jaro distances are calculated via the `stringdistmatrix` function from the
249 `stringdist` package (van der Loo, 2014). This function is provided to help researchers perform a spell
250 check on their dataset. However, it should be made clear that this is no replacement for taxonomic vetting.

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

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

```
268 # Evaluate unique taxa
269 tax_unique(occdf = tetrapods, genus = "genus", family = "family",
270             order = "order", class = "class", resolution = "genus")
```

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

```

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

```

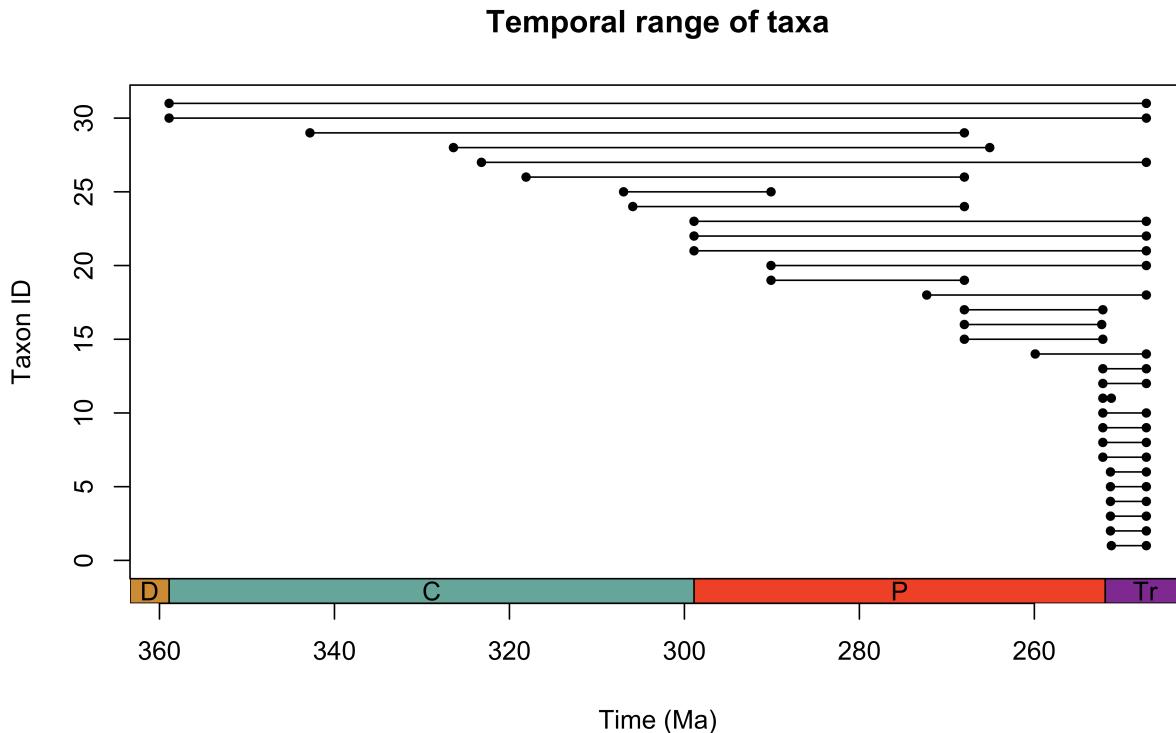


Figure 4: Temporal range of tetrapod orders in the `palaeoverse` example dataset.

```

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

```

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

290 Phylogeny wrangling

291 The function `phylo_check` compares a list of taxonomic names to the list of tip names in a user-provided
 292 phylogeny. This comparison can be provided as a table describing the presence or absence of each taxon in
 293 the list and/or tips, or as counts of taxa present only in the list, only in the phylogeny, or in both. The
 294 function can also be used to trim the phylogeny to only include branches whose tip names are included
 295 within the list of taxonomic names.

296 Additional features

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

```
300 # Compute the number of occurrences per collection  
301 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”).

```
314 reefs <- look_up(occdf = reefs,  
315                 early_interval = "interval",  
316                 late_interval = "interval",  
317                 int_key = interval key)
```

Finally, a common feature request from our survey was the ability to add the ‘Geological Time Scale’ to time-series plots in base R, with similar behaviour to the `deeptime` R package (Gearty, 2022) for `ggplot2` (Wickham, 2016). To address this request, the `axis_geo` function has been developed for the `palaeoverse` package (Figure 5).

```

322 # Palaeorotate reef dataset
323 reefs <- palaeorotate(occdf = reefs, age = "interval_mid_ma")
324 # Plot palaeolatitudinal distribution through time
325 plot(x = reefs$interval_mid_ma, y = reefs$p_lat,
326       xlab = "Time (Ma)", ylab = "Palaeolatitude (\u00B0)",
327       xlim = c(541, 0), xaxt = "n", type = "p", pch = 20)
328 # Add Geological Time Scale
329 axis geo(side = 1, intervals = "periods")

```

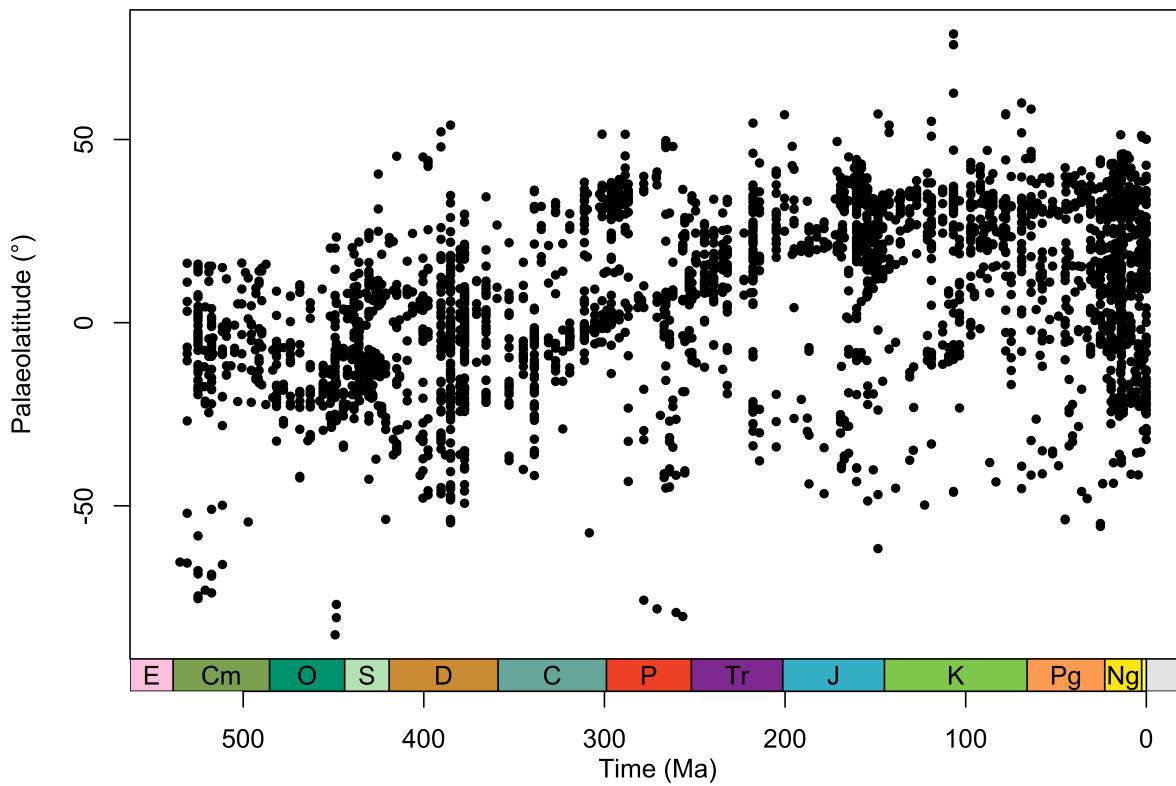


Figure 5: Example Phanerozoic plot of the palaeolatitudinal distribution of reefs through time. The plot demonstrates the usage of the `axis_geo` function for adding the Geological Time Scale to a base R plot.

330 Resources

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

337 **Future perspectives**

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

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

357 **Acknowledgements**

358 The authors are extremely grateful to all survey respondents who helped to shape the development of
359 `palaeoverse`. Special thanks are given to Emma M. Dunne whom participated in numerous discussions,
360 and shared her experience with the development team. The contributions of LAJ, SG, and AAC were
361 supported by the European Research Council under the European Union’s Horizon 2020 research and
362 innovation program (grant agreement 947921; MAPAS project). AAC was also supported by a Juan de la
363 Cierva-formación 2020 fellowship funded by FJC2020-044836-I / MCIN /AEI / 10.13039 /501100011033
364 from the European Union “NextGenerationEU”/PRTR. The contributions of WG were supported by the
365 Population Biology Program of Excellence Postdoctoral Fellowship from the University of Nebraska-
366 Lincoln School of Biological Sciences and the Lerner-Gray Postdoctoral Research Fellowship from the
367 Richard Gilder Graduate School at the American Museum of Natural History. The contributions of BJA

368 were supported by an ETH+ grant (BECCY). The contributions of CDD (RF_ERE_210013), MK
369 (RGF_EA_180318) and CN (RGF_R1_180020) were supported by Royal Society grants. The contributions
370 of PLG were supported by a FAPESP postdoctoral grant (2022/05697-9). This is Paleobiology Database
371 publication no XXX.

372 **Authors' contributions**

373 LAJ conceived the project. All authors contributed to developing the project. LAJ, BJA, WG, KE, CD, and
374 JFS contributed the code. All authors contributed to testing and reviewing the code. SG processed the
375 survey results and produced the survey figures. All authors contributed to writing the manuscript.

376 **Data accessibility**

377 The `palaeoverse` R package is hosted on CRAN (<https://cran.r-project.org/web/packages/palaeoverse/>)
378 and is available on GitHub (<https://github.com/palaeoverse-community/palaeoverse>). All example datasets
379 are bundled with the R package. All code is released under a GPL (>= 3) license.

380 **References**

- 381 Allen, B.J., Wignall, P.B., Hill, D.J., Saupe, E.E. and Dunhill, A.M. (2020) The latitudinal diversity
382 gradient of tetrapods across the permo-triassic mass extinction and recovery interval. *Proceedings of the
383 Royal Society B*, **287**, 20201125.
- 384 Antell, G.S., Kiessling, W., Aberhan, M. and Saupe, E.E. (2020) Marine biodiversity and geographic
385 distributions are independent on large scales. *Current Biology*, **30**, 115–121.e5.
- 386 Barido-Sottani, J., Pett, W., O'Reilly, J.E. and Warnock, R.C. (2019) FossilSim: An r package for
387 simulating fossil occurrence data under mechanistic models of preservation and recovery. *Methods in
388 Ecology and Evolution*, **10**, 835–840.
- 389 Bell, M.A. and Lloyd, G.T. (2015) *Strap: An r Package for Plotting Phylogenies Against Stratigraphy
390 and Assessing Their Stratigraphic Congruence*. Wiley Online Library.
- 391 Benton, M.J. and Harper, D. (1999) The history of life: Large databases in palaeontology. *Numerical
392 palaeobiology*, 249–283.
- 393 Chao, A., Gotelli, N.J., Hsieh, T.C., Sande, E.L., Ma, K.H., Colwell, R.K., et al. (2014) Rarefaction and
394 extrapolation with hill numbers: A framework for sampling and estimation in species diversity studies.
395 *Ecological Monographs*, **84**, 45–67.
- 396 Chiarenza, A.A., Mannion, P.D., Farnsworth, A., Carrano, M.T. and Varela, S. (2022) Climatic
397 constraints on the biogeographic history of mesozoic dinosaurs. *Current Biology*, **32**, 570–585.

- 398 Close, R.A., Benson, R.B.J., Alroy, J., Carrano, M.T., Cleary, T.J., Dunne, E.M., et al. (2020a) The
399 apparent exponential radiation of phanerozoic land vertebrates is an artefact of spatial sampling biases.
400 *Proceedings of the Royal Society B: Biological Sciences*, **287**, 20200372.
- 401 Close, R., Benson, R.B., Saupe, E., Clapham, M. and Butler, R. (2020b) The spatial structure of
402 phanerozoic marine animal diversity. *Science*, **368**, 420–424.
- 403 Csárdi, G., Hester, J., Wickham, H., Chang, W., Morgan, M. and Tenenbaum, D. (2021) *Remotes: R*
404 *Package Installation from Remote Repositories, Including 'GitHub'*.
- 405 Darroch, S.A., Casey, M.M., Antell, G.S., Sweeney, A. and Saupe, E.E. (2020) High preservation
406 potential of paleogeographic range size distributions in deep time. *The American Naturalist*, **196**, 454–
407 471.
- 408 Davies, T.W., Bell, M.A., Goswami, A. and Halliday, T.J. (2017) Completeness of the eutherian mammal
409 fossil record and implications for reconstructing mammal evolution through the cretaceous/paleogene
410 mass extinction. *Paleobiology*, **43**, 521–536.
- 411 Dean, C.D., Chiarenza, A.A. and Maidment, S.C. (2020) Formation binning: A new method for increased
412 temporal resolution in regional studies, applied to the late cretaceous dinosaur fossil record of north
413 america. *Palaeontology*, **63**, 881–901.
- 414 Filazzola, A. and Lortie, C. (2022) A call for clean code to effectively communicate science. *Methods in*
415 *Ecology and Evolution*, **13**, 2119–2128.
- 416 Flannery-Sutherland, J.T., Raja, N.B., Kocsis, Á.T. and Kiessling, W. (2022) Fossilbrush: An r package
417 for automated detection and resolution of anomalies in palaeontological occurrence data. *Methods in*
418 *Ecology and Evolution*.
- 419 Flannery-Sutherland, J.T., Silvestro, D. and Benton, M.J. (2022) Global diversity dynamics in the fossil
420 record are regionally heterogeneous. *Nature Communications*, **13**, 1–17.
- 421 Franeck, F. and Liow, L.H. (2020) Did hard substrate taxa diversify prior to the great ordovician
422 biodiversification event? *Palaeontology*, **63**, 675–687.
- 423 Fraser, D. (2017) Can latitudinal richness gradients be measured in the terrestrial fossil record?
424 *Paleobiology*, **43**, 479–494.
- 425 Furness, E.N., Garwood, R.J., Mannion, P.D. and Sutton, M.D. (2021) Evolutionary simulations clarify
426 and reconcile biodiversity-disturbance models. *Proceedings of the Royal Society B*, **288**, 20210240.
- 427 Garwood, R.J., Spencer, A.R. and Sutton, M.D. (2019) REvoSim: Organism-level simulation of macro
428 and microevolution. *Palaeontology*, **62**, 339–355.
- 429 Gearty, W. (2022) *Deeptime: Plotting Tools for Anyone Working in Deep Time*.
- 430 Gradstein, F.M., Ogg, J.G., Schmitz, M. and Ogg, G. (2012) *The Geologic Time Scale 2012*. Elsevier.
- 431 Gradstein, F.M., Ogg, J.G., Schmitz, M.D. and Ogg, G.M. (2020) *Geologic Time Scale 2020*. Elsevier.
- 432 Guillerme, T. (2018) dispRity: A modular r package for measuring disparity. *Methods in Ecology and*
433 *Evolution*, **9**, 1755–1763.

- 434 Jaro, M.A. (1989) Advances in record-linkage methodology as applied to matching the 1985 census of
435 tampa, florida. *Journal of the American Statistical Association*, **84**, 414–420.
- 436 Jones, L.A., Dean, C.D., Mannion, P.D., Farnsworth, A. and Allison, P.A. (2021) Spatial sampling
437 heterogeneity limits the detectability of deep time latitudinal biodiversity gradients. *Proceedings of the*
438 *Royal Society B*, **288**, 20202762.
- 439 Jones, L.A., Mannion, P.D., Farnsworth, A., Bragg, F. and Lunt, D.J. (2022) Climatic and tectonic drivers
440 shaped the tropical distribution of coral reefs. *Nature communications*, **13**, 1–10.
- 441 Kiessling, W. and Krause, C. (2022) PaleoReefs database (PARED).
- 442 Kocsis, Á.T., Reddin, C.J., Alroy, J. and Kiessling, W. (2019) The r package divDyn for quantifying
443 diversity dynamics using fossil sampling data. *Methods in Ecology and Evolution*, **10**, 735–743.
- 444 Lloyd, G.T. (2016) Estimating morphological diversity and tempo with discrete character-taxon matrices:
445 Implementation, challenges, progress, and future directions. *Biological journal of the Linnean Society*.
446 *Biological Journal of the Linnean Society*, **118**, 131–151.
- 447 Lloyd, G.T., Pearson, P.N., Young, J.R. and Smith, A.B. (2012) Sampling bias and the fossil record of
448 planktonic foraminifera on land and in the deep sea. *Paleobiology*, **38**, 569–584.
- 449 Mannion, P.D., Benson, R.B., Carrano, M.T., Tenant, J.P., Judd, J. and Butler, R.J. (2015) Climate
450 constrains the evolutionary history and biodiversity of crocodylians. *Nature communications*, **6**, 1–9.
- 451 Mannion, P.D., Benson, R.B., Upchurch, P., Butler, R.J., Carrano, M.T. and Barrett, P.M. (2012) A
452 temperate palaeodiversity peak in mesozoic dinosaurs and evidence for late cretaceous geographical
453 partitioning. *Global Ecology and Biogeography*, **21**, 898–908.
- 454 Mannion, P.D., Upchurch, P., Benson, R.B. and Goswami, A. (2014) The latitudinal biodiversity gradient
455 through deep time. *Trends in Ecology & Evolution*, **29**, 42–50.
- 456 Mannion, P.D., Upchurch, P., Carrano, M.T. and Barrett, P.M. (2011) Testing the effect of the rock
457 record on diversity: A multidisciplinary approach to elucidating the generic richness of sauropodomorph
458 dinosaurs through time. *Biological Reviews*, **86**, 157–181.
- 459 Müller, R.D., Cannon, J., Qin, X., Watson, R.J., Gurnis, M., Williams, S., et al. (2018) GPlates: Building
460 a virtual earth through deep time. *Geochemistry, Geophysics, Geosystems*, **19**, 2243–2261.
- 461 O'Brien, L. (2022) *H3jsr: Access Uber's H3 Library*.
- 462 Oksanen, J., Blanchet, F.G., Friendly, M., Kindt, R., Legendre, P., McGlinn, D., et al. (2020) *Vegan:*
463 *Community Ecology Package*.
- 464 Powell, M.G. (2009) The latitudinal diversity gradient of brachiopods over the past 530 million years. *The*
465 *Journal of Geology*, **117**, 585–594.
- 466 Quental, T.B. and Marshall, C.R. (2013) How the red queen drives terrestrial mammals to extinction.
467 *Science*, **341**, 290–292.
- 468 R Core Team. (2022) *R: A Language and Environment for Statistical Computing*. R Foundation for
469 Statistical Computing, Vienna, Austria.

- 470 Sepkoski, J.J. (1978) A kinetic model of phanerozoic taxonomic diversity i. Analysis of marine orders.
471 *Paleobiology*, **4**, 223–251.
- 472 Sepkoski, D. and Ruse, M. (2009) *The Paleobiological Revolution: Essays on the Growth of Modern*
473 *Paleontology*. University of Chicago Press.
- 474 Silvestro, D., Salamin, N. and Schnitzler, J. (2014) PyRate: A new program to estimate speciation and
475 extinction rates from incomplete fossil data. *Methods in Ecology and Evolution*, **5**, 1126–1131.
- 476 Silvestro, D., Zizka, A., Bacon, C.D., Cascales-Minana, B., Salamin, N. and Antonelli, A. (2016) Fossil
477 biogeography: A new model to infer dispersal, extinction and sampling from palaeontological data.
478 *Philosophical Transactions of the Royal Society B: Biological Sciences*, **371**, 20150225.
- 479 Song, H., Huang, S., Jia, E., Dai, X., Wignall, P.B. and Dunhill, A.M. (2020) Flat latitudinal diversity
480 gradient caused by the permian–triassic mass extinction. *Proceedings of the National Academy of*
481 *Sciences*, **117**, 17578–17583.
- 482 Starrfelt, J. and Liow, L.H. (2016) How many dinosaur species were there? Fossil bias and true richness
483 estimated using a poisson sampling model. *Philosophical Transactions of the Royal Society B: Biological*
484 *Sciences*, **371**, 20150219.
- 485 van der Loo, M.P.J. (2014) The stringdist package for approximate string matching. *The R Journal*, **6**,
486 111–122.
- 487 Vilhena, D.A. and Smith, A.B. (2013) Spatial bias in the marine fossil record. *PLoS One*, **8**, e74470.
- 488 Wickham, H. (2016) *Ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- 489 Zaffos, A., Finnegan, S. and Peters, S.E. (2017) Plate tectonic regulation of global marine animal
490 diversity. *Proceedings of the National Academy of Sciences*, **114**, 5653–5658.
- 491 Zizka, A., Silvestro, D., Andermann, T., Azevedo, J., Duarte Ritter, C., Edler, D., et al. (2019)
492 CoordinateCleaner: Standardized cleaning of occurrence records from biological collection databases.
493 *Methods in Ecology and Evolution*, **10**, 744–751.