GOLD
OPEN ACCESS

THE GEOLOGICAL SOCIETY OF AMERICA®

# rmacrostrat: An R package for accessing and retrieving data from the Macrostrat geological database

**Lewis A. Jones[1], Christopher D. Dean[1], William Gearty[2], and Bethany J. Allen[3,4]**

[1]Department of Earth Sciences, University College London, Gower Street, London WC1E 6BT, UK
[2]Division of Paleontology, American Museum of Natural History, 200 Central Park W, New York, New York 10024, USA
[3]Department of Biosystems Science and Engineering, ETH Zurich, Klingelbergstrasse 48, 4056 Basel, Switzerland
[4]Computational Evolution Group, Swiss Institute of Bioinformatics, Quartier Sorge, Bâtiment Amphipôle, 1015 Lausanne, Switzerland

## ABSTRACT

The geological record is a vast archive of information that provides the only empirical data about the evolution of the Earth. In recent years, concentrated efforts have been made to compile macrostratigraphic data into the online centralized database Macrostrat. Macrostrat is a global stratigraphic database containing information regarding surface and subsurface rock units and their respective ages, lithologies, geographic extents, and various other associated metadata. However, these raw data are currently directly accessible only through the Macrostrat application programming interface, which is a barrier to potential users that are less familiar with such services. This data accessibility hurdle currently prevents full capitalization of the value offered by Macrostrat, particularly its potential to improve understanding of the geological and biological evolution of the Earth. Here, we introduce rmacrostrat, an R package that interfaces with the Macrostrat database to access and retrieve a variety of geological, paleontological, and economic data directly into the R programming environment. In this article, we provide details about how the package can be installed, its implementation, and potential use cases. For the latter, we showcase how rmacrostrat can be used to visualize regional stratigraphic columns, produce regional geologic outcrop maps, and investigate temporal trends in macrostratigraphic units. We hope that this package will make geological data more readily accessible and in turn will facilitate new research utilizing Earth system data.

## ■ INTRODUCTION

Earth's geologic record provides a unique spatiotemporal archive of the evolutionary history of the planet (Ernst and Youbi, 2017; Tetley et al., 2019; Cao et al., 2019; Scotese et al., 2021). Historically, to understand macroscale

Lewis A. Jones  https://orcid.org/0000-0003-3902-8986
Christopher D. Dean  https://orcid.org/0000-0001-6471-6903
William Gearty  https://orcid.org/0000-0003-0076-3262
Bethany Allen  https://orcid.org/0000-0003-0282-6407

Earth system trends through geological time, researchers were required to synthesize local or regional quantitative studies, predominantly from data gathered in the form of regional geological maps, sections, and individual sedimentary logs or boreholes (e.g., Ronov et al., 1980; Seslavinskiy, 1991; Bosscher and Schlager, 1993; Miall, 2022). However, the introduction of large online open-access databases, in which a variety of complementary data sets are already digitized and synthesized, has facilitated the development of macroscale analyses through both time and space. One such database is Macrostrat (https://macrostrat.org/), a relational geospatial database that aims to aggregate and synthesize field-derived geological data from geological maps and regional geological columns into a data set that describes the spatial distribution of geological units within the Earth's upper crust (Peters et al., 2018). Macrostrat contains information regarding individual rock "units," linked by unique identification numbers to associated lithological, environmental, paleontological, and economic attributes, alongside information regarding their respective chronostratigraphic context. These units are organized spatially into "columns," representing a cross section of the upper crust within particular geological basins, and temporally by Macrostrat's internal chronostratigraphic age model (Fig. 1). Sequentially deposited units bounded by unconformities form geological "sections," which also have their own unique identification numbers. Additionally, Macrostrat units are linked by unique identification numbers to geological mapping data amalgamated from a variety of sources as well as data from other large geoscience databases such as the Paleobiology Database (PBDB; https://paleobiodb.org/) (Peters and McClennen, 2016; Uhen et al., 2023).

Since its initial compilation in 2005 from the American Association of Petroleum Geologists Correlation of Stratigraphic Units of North America (COSUNA) charts (Peters, 2006), Macrostrat has grown into a comprehensive and well-established database containing over 35,000 units and 1500 geologic columns, all of which are publicly accessible. Macrostrat aims to provide such data on a global scale, and while the abundance and resolution of available data are currently geographically variable, improving spatial coverage is one of the major aims of the project moving forward (Quinn et al., 2024). Data hosted by Macrostrat have been used for a wide variety of applications in scientific research as well as science communication and
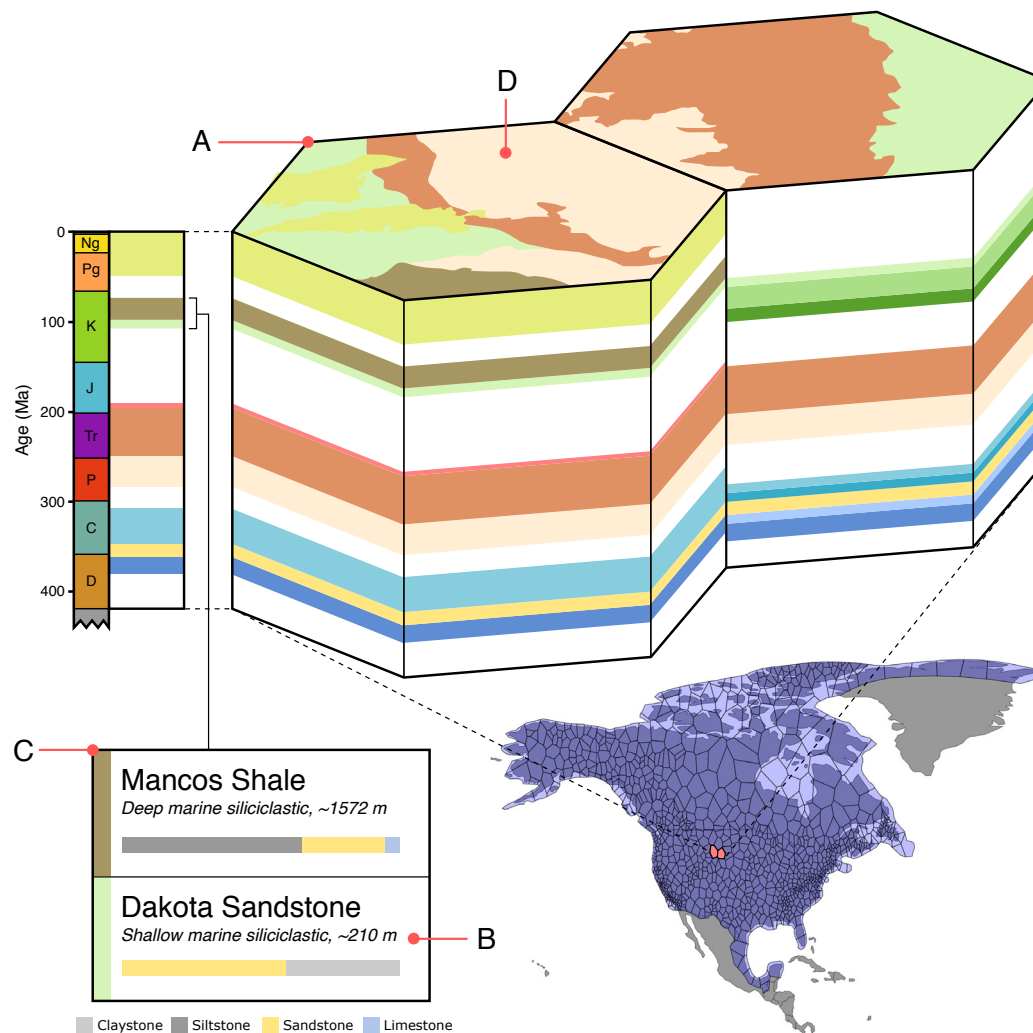
Figure 1. Schematic showing relationships between geological and stratigraphic data stored in Macrostrat. Data are arranged spatially as "columns" (A), which contain chronostratigraphic columns of distinct stratigraphic "units" (B), distinguished by a variety of attributes (e.g., lithology, environment, thickness, etc.) and organized temporally by Macrostrat's time-age model. Packages of continuously deposited units bounded by unconformities can be subset into "sections" (C). Macrostrat units can also be linked to Macrostrat's geological outcrop map (D), amalgamated from a variety of sources. Note that columns are idealized over the region and do not intersect accurately with surficial geological maps. Ng—Neogene; Pg—Paleogene; K—Cretaceous; J—Jurassic; Tr—Triassic; P—Permian; C—Carboniferous; D—Devonian.

education. The broad temporal and spatial scale of data hosted by Macrostrat has facilitated a diverse array of research related to Earth systems through time, including in the fields of sedimentology (Peters and Husson, 2017), stratigraphy (Tasistro-Hart and Macdonald, 2023), igneous petrology (Peters et al., 2021), geochemistry (Husson and Coogan, 2023), and paleobiology (Peters and Heim, 2010, 2011a, 2011b; Heim and Peters, 2011; Rook et al., 2013; Nelsen et al., 2016; Peters et al., 2017; Balseiro and Powell, 2020, 2023; Ye and Peters, 2023; Segessenman and Peters, 2024). Macrostrat has

also collaborated with the Extending Ocean Drilling Pursuits (eODP) project (https://eodp.github.io/) to integrate existing drill core data from sources such as the International Ocean Discovery Program (IODP) into the database (Sessa et al., 2023). Geologic map data held within Macrostrat are also displayed by a variety of web (e.g., Sift; https://macrostrat.org/sift/) and mobile applications (e.g., Rockd; https://rockd.org/) that aim to enable usage of geologic information by the wider scientific community, the general public, and university education platforms (Cohen et al., 2018). Macrostrat is also planning

to expand and integrate community-led validation of sections, ingestion of stratigraphic column data, and development of new software to facilitate data collaboration in the near future (Quinn et al., 2024). As such, Macrostrat is a vital resource for Earth scientists investigating a variety of issues related to both the geological history of our planet and the impacts of geological processes today.

Despite the apparent opportunities offered by Macrostrat, its hosted data can currently be directly accessed only via the database's application programming interface (API). Although a powerful resource, this single direct data access avenue means that familiarity with both the structure of the database and how to interact with APIs is necessary in order to use the database. Those able to overcome this data accessibility hurdle are still required to develop their own custom protocols to integrate Macrostrat data into coding-based scientific workflows; this can inherently lead to researchers "reinventing the wheel" and producing code that is case specific and difficult to repurpose, inhibiting the reproducibility of research conducted using Macrostrat data. Such processes are commonly carried out in the programming language R, which the Earth science community has broadly adopted to access, prepare, analyze, and plot data (e.g., Bell and Lloyd, 2015; Varela et al., 2015; Ortiz and Jaramillo, 2018; Barido-Sottani et al., 2019; Kocsis et al., 2019; Jones et al., 2023; Gearty, 2024). In particular, several R packages have been developed to interface with databases relevant to the geosciences through API services, supporting the generation of readable, reusable, and reproducible workflows (Varela et al., 2015; Gearty and Jones, 2023; Vidaña and Goring, 2023). However, until now, no such package has been available for interacting with the Macrostrat database.

Here, we present *rmacrostrat*, a dedicated R package for interfacing with the geological database Macrostrat. The package provides streamlined functionality for querying the database via its API service and retrieving various geological data (e.g., lithostratigraphic units), definitions, or metadata associated with the hosted data (e.g., lithological terms). First, we provide instructions for installing the package and details on its implementation. We then demonstrate the functionality available in *rmacrostrat* and provide typical usage examples. Finally, we provide details about the resources we have made available to support *rmacrostrat* users. By providing a programmatic solution to accessing the data hosted by Macrostrat, we endeavor to facilitate new research across the Earth sciences that is conducted in a streamlined, readable, reusable, and reproducible manner.

## ■ INSTALLATION

The *rmacrostrat* package can be installed from the Comprehensive R Archive Network (CRAN) using the install.packages() function in R (R Core Team, 2024):

```
install.packages("rmacrostrat")
```

If preferred, the development version of *rmacrostrat* can be installed from GitHub via the R package remotes (Csárdi et al., 2023):

```
remotes::install_github("palaeoverse/rmacrostrat")
```

Following installation, *rmacrostrat* can be loaded via the library() function in R:

```
library("rmacrostrat")
```

## Dependencies

The current version of *rmacrostrat* (ver. 1.0.0) was developed to fetch data from version 2 of Macrostrat's API. The package depends on R (≥4.0) (R Core Team, 2024) and imports functions from the R packages curl (Ooms, 2024, preprint), geojsonsf (Cooley, 2022), httr (Wickham, 2023), jsonlite (Ooms, 2014), and sf (Pebesma, 2018; Pebesma and Bivand, 2023). The package was developed with the support of the R packages devtools (Wickham et al., 2022), testthat (Wickham, 2011), and roxygen2 (Wickham et al., 2024).

## ■ IMPLEMENTATION

Functions are broadly grouped into two categories in *rmacrostrat*: (1) def_*, and (2) get_*. The def_* suite of functions provides access to the definitions (or metadata) associated with data stored in Macrostrat, such as lithologies [def_lithologies()], measurements [def_measurements()], or Macrostrat columns [def_columns()]. A summary of this suite of functions is provided in Table 1. The get_* suite of functions is for retrieving data from Macrostrat, such as Macrostrat columns [get_columns()], Macrostrat units [get_units()], or geological map outcrop objects [get_map_outcrop()]. Detailed descriptions of these functions are provided in Table 2.

## Definition Functions

Definitions (or metadata) of the various data stored in Macrostrat are retrieved from the Macrostrat API service via the def_* suite of functions (Table 1). The coverage of each of these functions should hopefully be immediately recognizable via their naming convention (e.g., def_lithologies() returns definitions of the lithologies used in Macrostrat). Data returned using the def_* suite of functions contain both categorical (and commonly hierarchical) information about data attributes of interest (e.g., def_lithologies() returns individual lithologies ["sandstone"], as well as the type ["siliciclastic"] and class ["sedimentary"] of the lithology) as well as unique identification numbers for individual attributes that can be used to query Macrostrat. Without

TABLE 1. SUMMARY OF THE DEFINITION SUITE OF FUNCTIONS (def_*) CURRENTLY AVAILABLE IN THE *rmacrostrat* R PACKAGE

| Function | Description |
|---|---|
| catalog() | Wrapper function to retrieve all definitions within a given definition set (e.g., lithologies) |
| def_columns() | Retrieve definitions for Macrostrat columns |
| def_drilling_sites() | Retrieve metadata for variables associated with the Extending Ocean Drilling Pursuits (eODP) project |
| def_econs() | Retrieve definitions for economic resources (e.g., coal) |
| def_environments() | Retrieve definitions for environments (e.g., dune) |
| def_grain_sizes() | Retrieve definitions for grain sizes (e.g., cobble) |
| def_intervals() | Retrieve definitions for time intervals (e.g., Cenozoic) |
| def_lithologies() | Retrieve definitions for lithologies (e.g., sandstone) |
| def_lithology_att() | Retrieve definitions for lithology attributes (e.g., tabular) |
| def_measurements() | Retrieve definitions for different measurements (e.g., porosity) |
| def_minerals() | Retrieve definitions for different minerals (e.g., agate) |
| def_plates() | Retrieve definitions for tectonic plates (e.g., Eurasia) |
| def_projects() | Retrieve definitions for Macrostrat projects (e.g., eODP) |
| def_references() | Retrieve definitions for published references |
| def_sources() | Retrieve definitions for geological maps (e.g., USGS) |
| def_strat_names() | Retrieve definitions for stratigraphic names (e.g., Hell Creek) |
| def_strat_name_concepts() | Retrieve definitions for stratigraphic name concepts (e.g., Dakota) |
| def_structures() | Retrieve definitions for geological structures (e.g., antiform) |
| def_timescales() | Retrieve definitions for timescales (e.g., international periods) |

*Notes:* USGS—U.S. Geological Survey.

TABLE 2. SUMMARY OF THE DATA RETRIEVAL SUITE OF FUNCTIONS (get_*) CURRENTLY AVAILABLE IN THE *rmacrostrat* R PACKAGE

| Function | Description |
|---|---|
| get_units() | Get data for Macrostrat units |
| get_sections() | Get data for Macrostrat sections |
| get_columns() | Get data for Macrostrat columns |
| get_age_model() | Get information about the age models for Macrostrat columns |
| get_map_outcrop() | Get spatial polygon data for geologic map outcrop |
| get_map_points() | Get spatial point data for geologic map measurements (e.g., strike, dip) |
| get_map_legends() | Get information from geologic map legends associated with outcrop and points |
| get_fossils() | Get Paleobiology Database collections data associated with Macrostrat entities |
| get_eodp() | Get Extending Ocean Drilling Pursuits data associated with various drilling programs |
| get_measurements() | Get measurements relevant to making geological inferences |
| get_paleogeography() | Get paleogeographic geometries based on the Wright et al. (2013) global plate model |
| get_stats() | Get statistics about Macrostrat projects |

user-specified arguments, all def_* functions return a data.frame object containing the entire data set of definitions associated with that function:

```
# Get all lithologies
def_lithologies()
# Get all minerals
def_minerals()
# Get all time intervals
def_intervals()
```

Alternatively, users can search for definitions of specific entities or groups of entities using the specific arguments for each def_* function. This can generally be achieved using specific unique identification numbers (integers) for those definitions or via a name (character strings):

```
# Get all marine environments by name
def_environments(environ_class = "marine")
# Get specific environment by ID
def_environments(environ_id = 2)
```

Jones et al. | *rmacrostrat*: Access and retrieve geological data

For convenience, we have also provided a wrapper around all def_* functions via the catalog() function. This function returns complete sets of definitions for each def_* function, which takes the suffix of an individual def_* function for its argument:

```
# Get all geological timescales
catalog(type = "timescales")
# Get all geological structures
catalog(type = "structures")
```

We strongly recommend using the def_* suite of functions prior to retrieving data from Macrostrat to better understand both the structure of the database and the utility offered by the functions available in *rmacrostrat*. Due to the wide variety of data available in Macrostrat, individual get_* functions include a large array of potential arguments that can differ substantially between functions (see Data Retrieval Functions section). By using the specific def_* functions related to potentially useful search criteria, users can efficiently identify arguments and parameters with which to query the database via the get_* suite of functions. Examples of the utility of the def_* functions are provided in the Application section below as well as in the available vignettes, which provide tutorials on how to use the package.

## Data Retrieval Functions

Data can be retrieved from the Macrostrat database API directly into the R environment using the get_* suite of functions (Table 2). These functions return either data related to specified Macrostrat entities (e.g., Macrostrat columns, units, sections, and age definitions), geologic map elements, or external data related to Macrostrat entities (e.g., PBDB collections, eODP data, paleogeographies); these data can be returned either as a standard data.frame or as a spatial simple features (i.e., sf) object (providing associated spatial geometries). In some instances where multiple values exist for a variable (e.g., proportions of lithologies within a unit), a hierarchical data.frame structure is employed (i.e., a data.frame within a data.frame). In accordance with the def_* suite of functions, the purpose of individual get_* functions is intended to be easily identifiable from their named suffix (e.g., get_columns() retrieves data for Macrostrat columns).

As opposed to the def_* functions, the get_* functions require at least one supplied argument for a valid database query. Although the array of possible arguments differs substantially between get_* functions, users can generally retrieve data based on several categories. Firstly, users can search by unique identification number for either the chosen data type to retrieve or based on another Macrostrat entity:

```
# Get specific column according to an ID
get_columns(column_id = 45)
# Get units and sections associated with a specific column ID
```

```
get_units(column_id = 45)
get_sections(column_id = 45)
# Get map outcrop related to specific unit ID
get_map_outcrop(unit_id = 1610)
```

It should be highlighted that for many get_* functions in *rmacrostrat* [e.g., get_columns(), get_units()], data can also be retrieved with associated spatial geometries that define the geographic extent or position of the retrieved data (see Fig. 1):

```
# Get specific column according to an ID
get_columns(column_id = 45, sf = TRUE)
# Get units associated with a specific column ID
get_units(column_id = 45, sf = TRUE)
```

Attribute information—such as lithostratigraphic name, lithology, environment, or economic source—can also be used independently, or in combination in some instances, to retrieve subsets of Macrostrat data. These attributes can be specified either using their unique identification number or by character string. Further information about each attribute to search by can be found in the respective def_* functions (e.g., lithology attribute information can be found in the def_lithologies() function):

```
# Get units inferred to be marine
get_units(environ_class = "marine")
# Get all sandstone units by name or ID
get_units(lithology = "sandstone")
get_units(lithology_id = 10)
```

Data can also be retrieved using temporal limits by specifying either a specific interval name as a character string (e.g., "Permian"), a unique identification number, or a numeric value (e.g., 275 Ma) or from providing constraints based on numerical limits (e.g., 251.9–298.9 Ma). All Macrostrat entities that overlap with the specified parameter(s) in terms of their chronostratigraphic range defined in the Macrostrat age model are returned:

```
# Get units by interval name
get_units(interval_name = "Aptian")
# Get units by interval ID
get_units(interval_id = 43)
# Get units by age
get_units(age = 200)
# Get units by age range
get_units(age_bottom = 250, age_top = 200)
```

Finally, some get_* functions allow the user to query the database using geographic or spatial information. This can be achieved either by specifying

coordinates in decimal latitude and longitude degrees or, if continental-scale resolution is desired, through the use of Macrostrat projects. Macrostrat data are split into regional projects, such as North America (project_id = 1) and New Zealand (project_id = 5); setting this argument returns all Macrostrat entities associated with that regional project. It should be noted that different Macrostrat projects currently have different levels of data completeness, ranging from virtually complete temporal and spatial coverage (e.g., North America, project_id = 1) to incomplete and limited coverage (e.g., Africa, project_id = 9):

```
# Get sections which appear at a specific longitude & latitude
get_sections(lng = -105.15, lat = 37.89)
# Get map outcrop which appears at a specific longitude & latitude
get_map_outcrop(lng = -105.15, lat = 37.89)
# Get all Macrostrat unit data for the North American continent
get_units(project_id = 1)
```

As aforementioned, it is recommended that these arguments be used in tandem with the def_* suite of functions to maximize search potential and data retrieval. For instance, a user interested in retrieving units deposited in a specific paleoenvironment may want to use the def_environments() function prior to their search to see the full variety of parameters by which to search. We reiterate the importance of always exploring the data fetched from Macrostrat and ensuring returned data are as expected. As an illustrative example of this, when seeking all carbonate-bearing marine units, the use of the lithology_type argument would return many more units than environ_type because higher-resolution paleoenvironment interpretations for all units are currently incomplete:

```
# Get Macrostrat "carbonate" units
get_units(lithology_type = "carbonate")
get_units(environ_type = "carbonate")
```

## ■ APPLICATION

Herein, we provide three example applications of the *rmacrostrat* package. These examples are greatly expanded in step-by-step vignettes provided alongside the package, available online via the associated package website (https://rmacrostrat.palaeoverse.org/articles/) and also bundled with the package, accessible via:
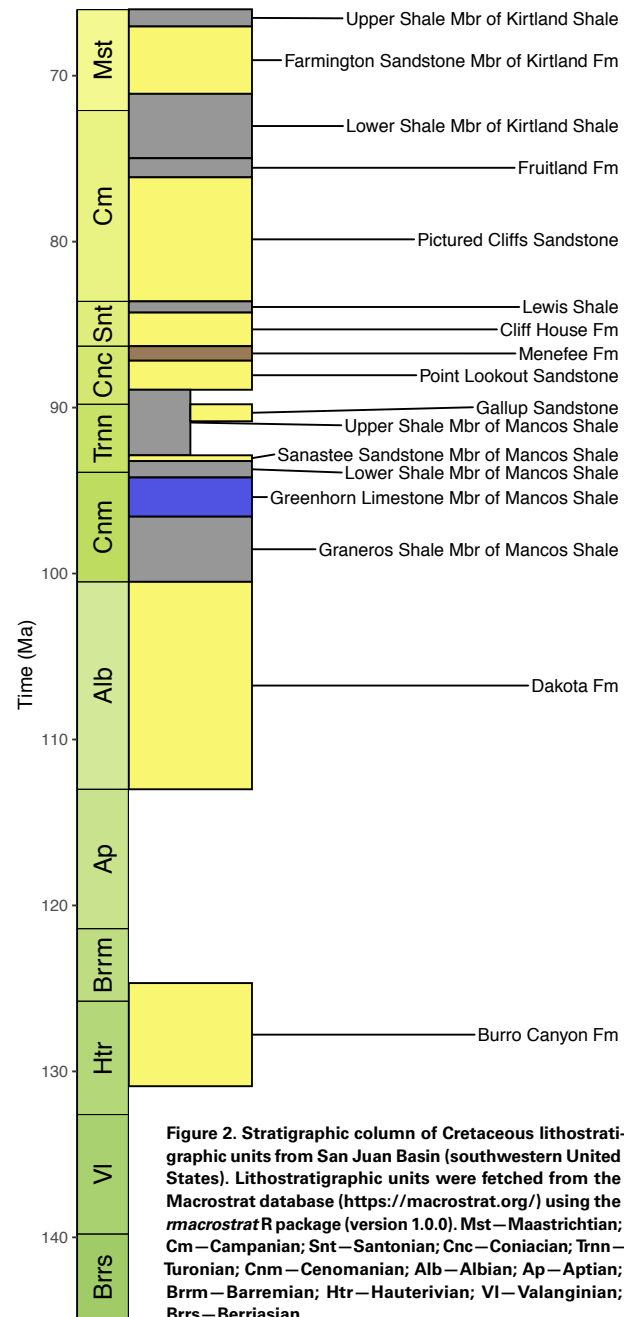
```
browseVignettes(package = "rmacrostrat")
```

### Constructing Stratigraphic Columns

An understanding of stratigraphy—that is, the relationships between adjacent geological units—is fundamental to accurately reading the geological record. This understanding enables researchers to put relative ages to lithological units and make temporal and spatial correlations with variables of interest. Using *rmacrostrat*, the geological data within the Macrostrat database can be easily retrieved and used to generate a stratigraphic column for a specific location and/or time interval. Below we provide an example showing how to retrieve and plot a stratigraphic column for the San Juan Basin, an asymmetric structural basin in northwestern New Mexico and southwestern Colorado (Four Corners region of the southwestern United States), containing sedimentary rocks ranging from Cambrian to Holocene in age (Fassett and Hinds, 1971). For this example, we restrict our column data to the Cretaceous, but this approach could equally be applied to any other basin or temporal interval. Columns are the most broad-scale geological entity available within Macrostrat, and by using the def_columns() function, the column associated with the San Juan Basin can be identified. The unique column identification number can then be used to get data for all appropriate units via get_units(). Given that the example focuses only on the Cretaceous, additional arguments available in get_units() can be used to further filter the queried data. With the returned data—Cretaceous lithostratigraphic units within the San Juan Basin—a stratigraphic column can be generated (Fig. 2):

```
# Load packages
library(rmacrostrat)
library(ggplot2)
library(ggrepel)
library(deeptime)
# Get the column definition of the San Juan Basin
column_def <- def_columns(column_name = "San Juan Basin")
# Using the column ID, retrieve all units of Cretaceous age
san_juan_units <- get_units(column_id = column_def$col_id,
                            interval_name = "Cretaceous")
# Specify x_min and x_max in dataframe
san_juan_units$x_min <- 0
san_juan_units$x_max <- 1
# Tweak values for overlapping units
san_juan_units$x_max[10] <- 0.5
san_juan_units$x_min[11] <- 0.5
# Add midpoint age for plotting
san_juan_units$m_age <- (san_juan_units$b_age +
                         san_juan_units$t_age) / 2
# Standardize and correct unit names according to USGS Geolex
san_juan_units$unit_name <- gsub(pattern = "Kirkland",
                                 replacement = "Kirtland",
                                 x = san_juan_units$unit_name)
san_juan_units$unit_name <- gsub(pattern = "Graneros Mbr",
                                 replacement =
                                 "Graneros Shale Mbr",
                                 x = san_juan_units$unit_name)
```

Figure 2. Stratigraphic column of Cretaceous lithostratigraphic units from San Juan Basin (southwestern United States). Lithostratigraphic units were fetched from the Macrostrat database (https://macrostrat.org/) using the *rmacrostrat* R package (version 1.0.0). Mst—Maastrichtian; Cm—Campanian; Snt—Santonian; Cnc—Coniacian; Trnn—Turonian; Cnm—Cenomanian; Alb—Albian; Ap—Aptian; Brrm—Barremian; Htr—Hauterivian; Vl—Valanginian; Brrs—Berriasian.

```
san_juan_units$unit_name <- gsub(pattern = "Sanostee Mbr",
                                 replacement =
                                 "Sanastee Sandstone Mbr",
                                 x = san_juan_units$unit_name)
# Plot stratigraphic column
ggplot(san_juan_units, aes(ymin = b_age, ymax = t_age,
                           xmin = x_min, xmax = x_max)) +
   # Plot units, colored by rock type
   geom_rect(fill = san_juan_units$color, color = "black") +
   # Add text labels
   geom_text_repel(aes(x = x_max, y = m_age, label =
              unit_name),
              box.padding = 0.1, nudge_x = 3,
              size = 3.5) +
   # Reverse direction of y-axis
   scale_y_reverse(limits = c(145, 66), n.breaks = 10,
              name = "Time (Ma)") +
   # Theming
   theme_classic() +
   theme(legend.position = "none",
       axis.line.x = element_blank(),
       axis.title.x = element_blank(),
       axis.text.x = element_blank(),
       axis.ticks.x = element_blank()) +
   # Add geological time scale
   coord_geo(pos = "left", dat = list("stages"), rot = 90)
```
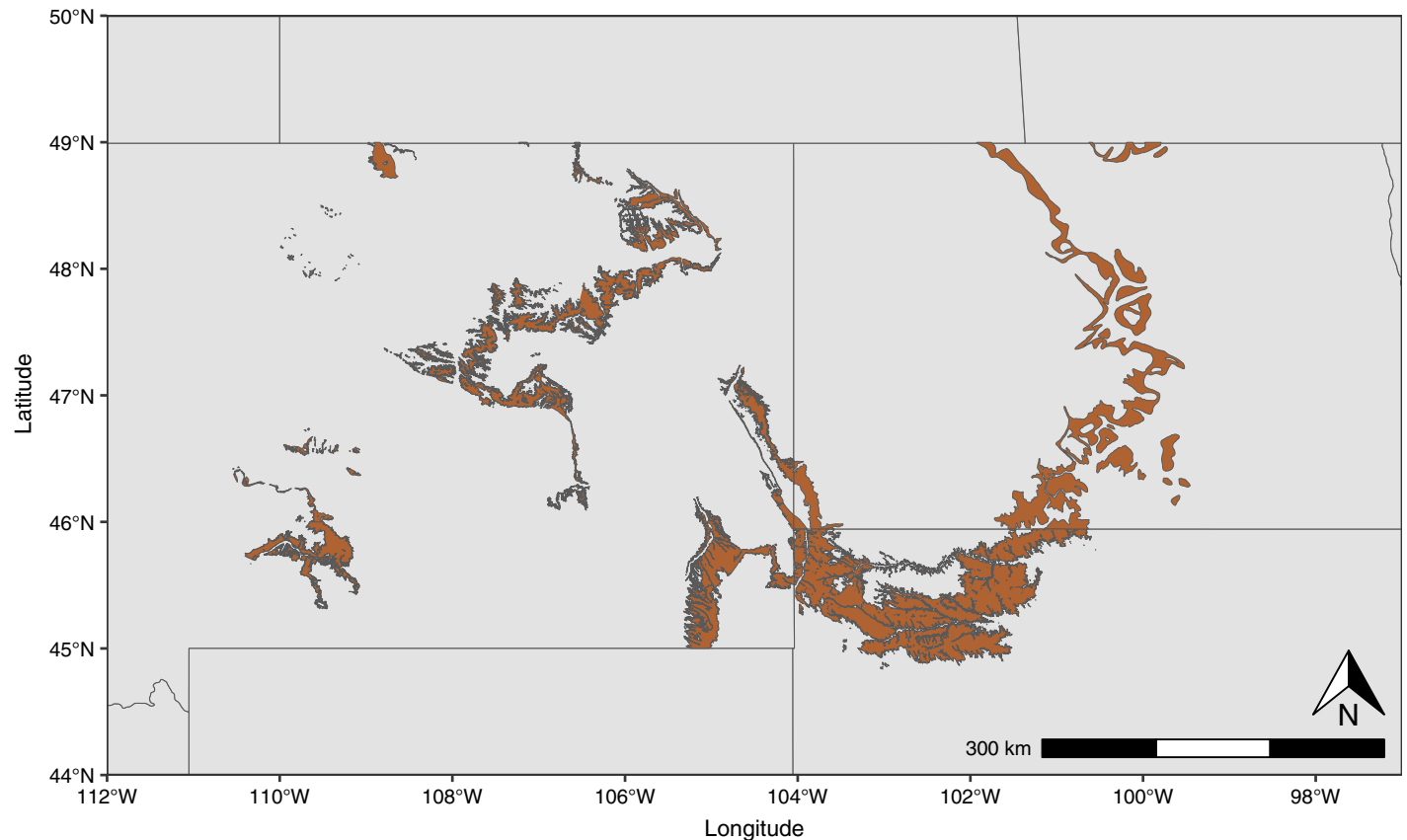
## Plotting Geologic Outcrop Maps

A commonly required figure across a range of disciplines within the geosciences is a geographic map of the outcrop for a specific geologic formation. Such a figure can be easily generated using the get_map_outcrop() function of *rmacrostrat*, which retrieves geospatial data associated with lithostratigraphic units. Below we provide an example for constructing a map of outcrop for the Hell Creek Formation, a geologic formation from the latest Cretaceous and early Paleogene of North America, which is found cropping out across Montana and North and South Dakota in the United States (Johnson et al., 2002; Fastovsky and Bercovici, 2016). Given that outcrop spatial data are compiled from various map sources, the definition function def_strat_names() is first used to find the appropriate identification numbers for any stratigraphic names of formations that include the Hell Creek. This information can then be used with the get_map_outcrop() function to retrieve geospatial data for the formation as a simple features (sf) object. These data can be plotted to produce a geological map (Fig. 3):

```
# Load libraries
library(rmacrostrat)
```

**Figure 3. Outcrop of Hell Creek Formation across Montana and North and South Dakota (United States). Outcrop data were fetched from the Macrostrat database (https://macrostrat.org/) using the *rmacrostrat* R package (version 1.0.0).**

```
library(rnaturalearth)
library(ggplot2)
library(ggspatial)
# Get data for chosen formation, specifying by stratigraphic rank
hc_def <- def_strat_names(strat_name = "hell creek,"
                                       rank = "Fm")
# Get spatial outcrop data for the chosen formation based on ID
hc <- get_map_outcrop(strat_name_id = hc_def$strat_name_id,
                      sf = TRUE)
# Load background maps
n_a <- ne_states(country = "united states of america",
                 returnclass = "sf")
ca <- ne_states(country = "canada",
                returnclass = "sf")
# Plot the map
ggplot() +
  geom_sf(data = n_a) +
  geom_sf(data = ca) +
  geom_sf(data = hc,
          fill = "#C7622B",
          lwd = 0) +
  coord_sf(xlim = c(-112, -97), ylim = c(44, 50), expand =
           FALSE) +
  labs(x = "Longitude", y = "Latitude") +
```

```
annotation_north_arrow(location = "br",
                        pad_y = unit(0.75, "cm"),
                        height = unit(1, "cm"),
                        width = unit(1, "cm")) +
annotation_scale(location = "br", width_hint = 0.3) +
theme_bw()
```

### Examining Macrostratigraphic Temporal Trends

Initial publications using data from the Macrostrat database quantified how the counts and proportion of Macrostrat entities, as well as different lithostratigraphic unit types associated with different paleoenvironments (e.g., marine, marginal, mixed, terrestrial), varied throughout the Phanerozoic (Peters and Heim, 2010). *rmacrostrat* facilitates access to these types of data and allows for similar analyses to be conducted. Below we provide an example of such an analysis, in this case estimating the number of igneous, metamorphic, and sedimentary units in North America throughout the Phanerozoic.

For this example, the relevant lithostratigraphic unit data from Macrostrat are first fetched using the get_units() function from *rmacrostrat*. For this query, several filters are applied to retrieve the appropriate data. First, the lithology_class argument is used to separate out the major rock classes. Second, the interval_name argument is used to filter to units only from the Phanerozoic. Finally, the project_id argument is used to filter results to units from the North American geological record:

```
# Load libraries
library(rmacrostrat)
# Get units by litholoy class, interval, and project ID
sedimentary <- get_units(lithology_class = "sedimentary",
                         interval_name = "Phanerozoic",
                         project_id = 1)
igneous <- get_units(lithology_class = "igneous",
                     interval_name = "Phanerozoic",
                     project_id = 1)
metamorphic<- get_units(lithology_class = "metamorphic",
                        interval_name = "Phanerozoic",
                        project_id = 1)
# Add column of sediment type
sedimentary$lithology_class <- "Sedimentary"
igneous$lithology_class <- "Igneous"
metamorphic$lithology_class <- "Metamorphic"
# Bind data
units <- rbind.data.frame(sedimentary, igneous, metamorphic)
```

With these data, the number of units for each rock type can be calculated for every international geological stage (i.e., time bin) through time. Functionality available in the palaeoverse R package can be used to retrieve relevant information about geological stages (Jones et al., 2023):

```
# Load libraries
library(palaeoverse)
# Generate stage-level time bins
bins <- time_bins(interval = c("Phanerozoic"), scale =
 "GTS2020")
# Rename age columns in units to be consistent with our bins
colnames(units)[which(colnames(units) == "b_age")] <- "max_ma"
colnames(units)[which(colnames(units) == "t_age")] <- "min_ma"
# Filter units with older maximum age than bins
units <- units[units$max_ma < 541,]
# Bin data
units <- bin_time(occdf = units, bins = bins,
                  min_ma = "min_ma", max_ma = "max_ma",
                  method = "all")
# Calculate the number of lithological classes per bin assignment
counts <- group_apply(occdf = units,
                      group = c("lithology_class",
                                "bin_assignment"),
                      fun = nrow)
# Rename columns to ease reading merging
colnames(counts) <- c("count", "lithology_class", "bin")
# Merge datasets by bin number
counts <- merge(x = bins, y = counts, by = "bin")
```

Using additional R packages for visualization, such as ggplot2 (Wickham, 2016) and deeptime (Gearty, 2024), Phanerozoic stage-level counts of North American lithostratigraphic units can be plotted by rock type (Fig. 4):

```
# Load libraries
library(ggplot2)
library(deeptime)
# Generate a plot of lithostratigraphic units through time
ggplot(counts, aes(fill = lithology_class, y = count, x =
 mid_ma)) +

  # Stacked bar chart with width specified by interval duration
  geom_bar(position = "stack", stat = "identity",
           width = counts$duration_myr,
           color = "black", linewidth = 0.1) +
  # Label y-axis
  scale_y_continuous("Counts of lithostratigraphic units") +
  # Label x-axis and reverse direction
  scale_x_reverse("Time (Ma)") +
```

```
# Data plotting colors
scale_fill_manual(values = c("#FF0000", "#00A08A",
 "#F2AD00")) +
# Theming
theme_bw() +
theme(legend.title = element_blank(),
      legend.position = c(0.9, 0.9)) +
# Add geological time scale
coord_geo()
```
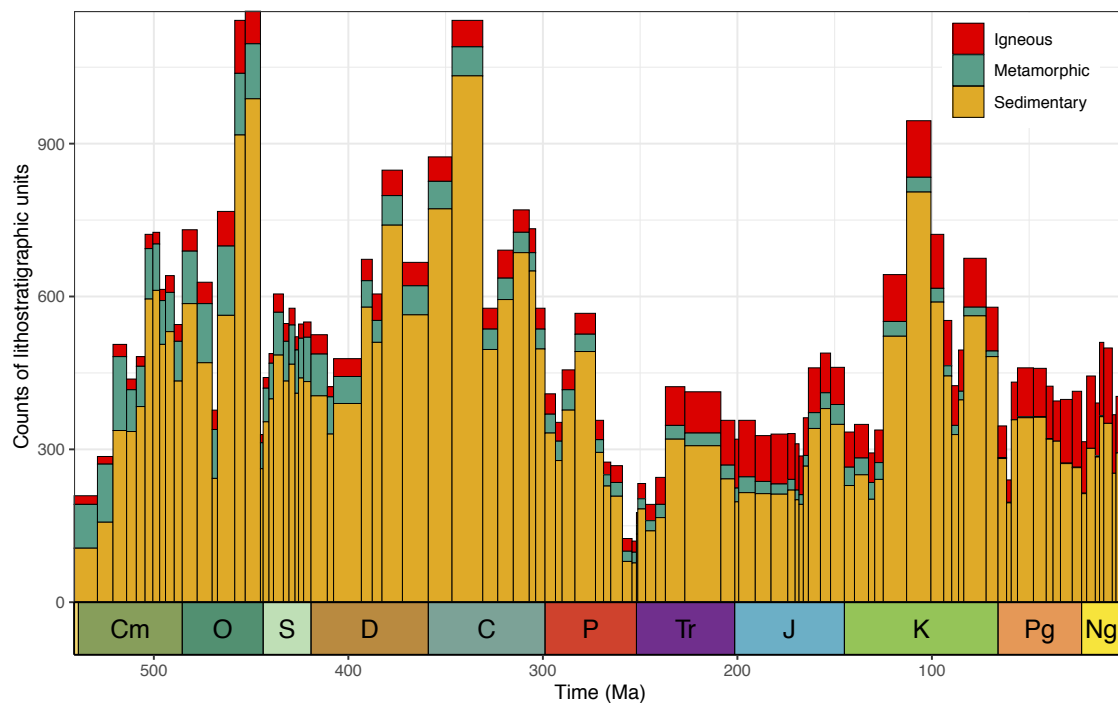
## ■ RESOURCES

We have made several resources available for our users. First, we have built a package website (https://rmacrostrat.palaeoverse.org) that provides information on how to use and contribute to *rmacrostrat*, how to report issues and bugs, and a contributor code of conduct. We have also made available three vignettes (i.e., tutorials) for the package, which provide user-friendly usage guides (https://rmacrostrat.palaeoverse.org/articles). Through *rmacrostrat*, we hope to further foster collaboration and the sharing of resources within the Earth science community. With this goal in mind, we warmly welcome the community to join and follow our community spaces, such as our GitHub organization page (https://github.com/palaeoverse) and Google Group (https://groups.google.com/g/palaeoverse), where users can share ideas and resources, advertise opportunities, and network with colleagues.

## ■ FUTURE PERSPECTIVES

The development of *rmacrostrat* expands upon the suite of software toolkits available within the Palaeoverse (https://palaeoverse.org/) "universe" (Jones, 2022; Gearty and Jones, 2023; Jones et al., 2023). The current version of *rmacrostrat* uses version 2 of Macrostrat's API to retrieve data, and it is our intention to track future versions of the API as updates become available, including the planned integration of eODP data into Macrostrat's data entities (e.g., columns). Through *rmacrostrat*, we hope to improve accessibility to the vast pool of geological data available within the Macrostrat database and facilitate new research across the Earth sciences. The *rmacrostrat* R package offers researchers the opportunity to streamline their research by providing a bridge between Macrostrat and the R environment as well as supporting the



**Figure 4. Number of Macrostrat lithostratigraphic units throughout the Phanerozoic (541–0 Ma) binned by stratigraphic stage-level bins and grouped by lithology class (igneous, metamorphic, and sedimentary). Units are counted for all time bins they overlap with. Lithostratigraphic units were fetched from the Macrostrat database (https://macrostrat.org/) using the *rmacrostrat* R package (version 1.0.0). Cm—Cambrian; O—Ordovician; S—Silurian; D—Devonian; C—Carboniferous; P—Permian; Tr—Triassic; J—Jurassic; K—Cretaceous; Pg—Paleogene; Ng—Neogene.**

capacity to generate fully reproducible pipelines. We hope that these benefits will encourage the community to further capitalize on the value offered by Macrostrat and may ultimately lead to higher data quality through peer review. As we have demonstrated with our example applications, *rmacrostrat* can be used to support the efficient plotting of stratigraphic columns, mapping of geological outcrop, and quantification of temporal dynamics in available macrostratigraphic units. However, we envision that *rmacrostrat* can also be used to support a wide range of additional analyses across the Earth sciences, such as economic resource exploration, comparisons between deep-time diversity dynamics and environmental change, and hazard mapping.

## REFERENCES CITED

Balseiro, D., and Powell, M.G., 2020, Carbonate collapse and the late Paleozoic ice age marine biodiversity crisis: Geology, v. 48, p. 118–122, https://doi.org/10.1130/G46858.1.

Balseiro, D., and Powell, M.G., 2023, Relative oversampling of carbonate rocks in the North American marine fossil record: Paleobiology, v. 49, p. 733–746, https://doi.org/10.1017/pab.2023.16.

Barido-Sottani, J., Pett, W., O'Reilly, J.E., and Warnock, R.C.M., 2019, FossilSim: An R package for simulating fossil occurrence data under mechanistic models of preservation and recovery: Methods in Ecology and Evolution, v. 10, p. 835–840, https://doi.org/10.1111/2041-210X.13170.

Bell, M.A., and Lloyd, G.T., 2015, strap: An R package for plotting phylogenies against stratigraphy and assessing their stratigraphic congruence: Paleontology, v. 58, p. 379–389, https://doi.org/10.1111/pala.12142.

Bosscher, H., and Schlager, W., 1993, Accumulation rates of carbonate platforms: The Journal of Geology, v. 101, p. 345–355, https://doi.org/10.1086/648228.

Cao, W., Williams, S., Flament, N., Zahirovic, S., Scotese, C., and Müller, R.D., 2019, Palaeolatitudinal distribution of lithologic indicators of climate in a palaeogeographic framework: Geological Magazine, v. 156, p. 331–354, https://doi.org/10.1017/S0016756818000110.

Cohen, P.A., Lockwood, R., and Peters, S., 2018, Integrating Macrostrat and Rockd into undergraduate earth science teaching: Cambridge University Press, https://doi.org/10.1017/9781108681445 (accessed May 2024).

Cooley, D., 2022, geojsonsf: GeoJSON to simple feature converter: R package, version 2.0.3, https://doi.org/10.32614/CRAN.package.geojsonsf (last accessed 19 September 2024).

Csárdi, G., Hester, J., Wickham, H., Chang, W., Morgan, M., and Tenenbaum, D., 2024, remotes: R Package Installation from Remote Repositories, Including 'GitHub', https://doi.org/10.32614/CRAN.package.remotes (accessed 19 September 2024).

Ernst, R.E., and Youbi, N., 2017, How Large Igneous Provinces affect global climate, sometimes cause mass extinctions, and represent natural markers in the geological record: Palaeogeography, Palaeoclimatology, Palaeoecology, v. 478, p. 30–52, https://doi.org/10.1016/j.palaeo.2017.03.014.

Fassett, J.E., and Hinds, J.S., 1971, Geology and fuel resources of the Fruitland Formation and Kirtland Shale of the San Juan Basin, New Mexico and Colorado: U.S. Geological Survey Professional Paper 676, 76 p., https://doi.org/10.3133/pp676.

Fastovsky, D.E., and Bercovici, A., 2016, The Hell Creek Formation and its contribution to the Cretaceous–Paleogene extinction: A short primer: Cretaceous Research, v. 57, p. 368–390, https://doi.org/10.1016/j.cretres.2015.07.007.

Gearty, W., 2024, deeptime: Plotting tools for anyone working in deep time: R package, version 1.1.1, https://doi.org/10.32614/CRAN.package.deeptime (last accessed 19 September 2024).

Gearty, W., and Jones, L.A., 2023, rphylopic: An R package for fetching, transforming, and visualising PhyloPic silhouettes: Methods in Ecology and Evolution, v. 14, p. 2700–2708, https://doi.org/10.1111/2041-210X.14221.

Heim, N.A., and Peters, S.E., 2011, Covariation in macrostratigraphic and macroevolutionary patterns in the marine record of North America: Geological Society of America Bulletin, v. 123, p. 620–630, https://doi.org/10.1130/B30215.1.

Husson, J.M., and Coogan, L.A., 2023, River chemistry reveals a large decrease in dolomite abundance across the Phanerozoic: Geochemical Perspectives Letters, v. 26, p. 1–6, https://doi.org/10.7185/geochemlet.2316.

Johnson, K.R., Nichols, D.J., and Hartman, J.H., 2002, Hell Creek Formation: A 2001 synthesis, *in* Hartman, J.H., Johnson, K.R., and Nichols, D.J., eds., The Hell Creek Formation and the Cretaceous-Tertiary Boundary in the Northern Great Plains: An Integrated Continental Record of the End of the Cretaceous: Geological Society of America Special Paper 361, p. 503–510, https://doi.org/10.1130/0-8137-2361-2.503.

Jones, L.A., 2022, sepkoski: Sepkoski's fossil marine animal genera compendium: R package, version 0.0.1, https://doi.org/10.32614/CRAN.package.sepkoski (19 September 2024).

Jones, L.A., et al., 2023, palaeoverse: A community-driven R package to support palaeobiological analysis: Methods in Ecology and Evolution, v. 14, p. 2205–2215, https://doi.org/10.1111/2041-210X.14099.

Kocsis, Á.T., Reddin, C.J., Alroy, J., and Kiessling, W., 2019, The R package divDyn for quantifying diversity dynamics using fossil sampling data: Methods in Ecology and Evolution, v. 10, p. 735–743, https://doi.org/10.1111/2041-210X.13161.

Miall, A.D., 2022, Stratigraphy: The modern synthesis, *in* Stratigraphy: A Modern Synthesis: Cham, Switzerland, Springer International Publishing, p. 341–417, https://doi.org/10.1007/978-3-030-87536-7_7.

Nelsen, M.P., DiMichele, W.A., Peters, S.E., and Boyce, C.K., 2016, Delayed fungal evolution did not cause the Paleozoic peak in coal production: Proceedings of the National Academy of Sciences of the United States of America, v. 113, p. 2442–2447, https://doi.org/10.1073/pnas.1517943113.

Ooms, J., 2014, The jsonlite package: A practical and consistent mapping between JSON data and R objects: arXiv, https://doi.org/10.48550/arXiv.1403.2805 (accessed June 2024).

Ooms, J., 2024, curl: A modern and flexible web client for r, https://doi.org/10.32614/CRAN.package.curl (last accessed 19 September 2024).

Ortiz, J.R., and Jaramillo, C.A., 2018, SDAR: A toolkit for stratigraphic data analysis in R, 10.32614/CRAN.package.SDAR (last accessed 19 September 2024).

Pebesma, E., 2018, Simple features for R: Standardized support for spatial vector data: The R Journal, v. 10, p. 439–446, https://doi.org/10.32614/RJ-2018-009.

Pebesma, E., and Bivand, R., 2023, Spatial Data Science: With Applications in R: New York, Chapman and Hall/CRC, 314 p., https://doi.org/10.1201/9780429459016.

Peters, S.E., 2006, Macrostratigraphy of North America: The Journal of Geology, v. 114, p. 391–412, https://doi.org/10.1086/504176.

Peters, S.E., and Heim, N.A., 2010, The geological completeness of paleontological sampling in North America: Paleobiology, v. 36, p. 61–79, https://doi.org/10.1666/0094-8373-36.1.61.

Peters, S.E., and Heim, N.A., 2011a, Macrostratigraphy and macroevolution in marine environments: Testing the common-cause hypothesis, *in* McGowan, A.J., and Smith, A.B., eds., Comparing the Geological and Fossil Records: Implications for Biodiversity Studies: Geological Society of London Special Publication 358, p. 95–104, https://doi.org/10.1144/SP358.7.

Peters, S.E., and Heim, N.A., 2011b, Stratigraphic distribution of marine fossils in North America: Geology, v. 39, p. 259–262, https://doi.org/10.1130/G31442.1.

Peters, S.E., and Husson, J.M., 2017, Sediment cycling on continental and oceanic crust: Geology, v. 45, p. 323–326, https://doi.org/10.1130/G38861.1.

Peters, S.E., and McClennen, M., 2016, The Paleobiology Database application programming interface: Paleobiology, v. 42, p. 1–7, https://doi.org/10.1017/pab.2015.39.

Peters, S.E., Husson, J.M., and Wilcots, J., 2017, The rise and fall of stromatolites in shallow marine environments: Geology, v. 45, p. 487–490, https://doi.org/10.1130/G38931.1.

Peters, S.E., Husson, J.M., and Czaplewski, J., 2018, Macrostrat: A platform for geological data integration and deep-time Earth crust research: Geochemistry, Geophysics, Geosystems, v. 19, p. 1393–1409, https://doi.org/10.1029/2018GC007467.

Peters, S.E., Walton, C.R., Husson, J.M., Quinn, D.P., Shorttle, O., Keller, C.B., and Gaines, R.R., 2021, Igneous rock area and age in continental crust: Geology, v. 49, p. 1235–1239, https://doi.org/10.1130/G49037.1.

Quinn, D.P., Idzikowski, C.R., and Peters, S.E., 2024, Building a multi-scale, collaborative, and time-integrated digital crust: The next stage of the Macrostrat data system: Geoscience Data Journal, v. 11, p. 11–26, https://doi.org/10.1002/gdj3.189.

R Core Team, 2024, R: A language and environment for statistical computing: Vienna, Austria, R Foundation for Statistical Computing, https://www.R-project.org/ (last accessed 19 September 2024).

Ronov, A.B., Khain, V.E., Balukhovsky, A.N., and Seslavinsky, K.B., 1980, Quantitative analysis of Phanerozoic sedimentation: Sedimentary Geology, v. 25, p. 311–325, https://doi.org/10.1016/0037-0738(80)90067-6.

Rook, D.L., Heim, N.A., and Marcot, J., 2013, Contrasting patterns and connections of rock and biotic diversity in the marine and non-marine fossil records of North America: Palaeogeography, Palaeoclimatology, Palaeoecology, v. 372, p. 123–129, https://doi.org/10.1016/j.palaeo.2012.10.006.

Scotese, C.R., Song, H., Mills, B.J.W., and van der Meer, D.G., 2021, Phanerozoic paleotemperatures: The earth's changing climate during the last 540 million years: Earth-Science Reviews, v. 215, https://doi.org/10.1016/j.earscirev.2021.103503.

Segessenman, D.C., and Peters, S.E., 2024, Transgression–regression cycles drive correlations in Ediacaran–Cambrian rock and fossil records: Paleobiology, v. 50, p. 150–163, https://doi.org/10.1017/pab.2023.31.

Seslavinskiy, K.B., 1991, Global transgressions and regressions during the Paleozoic: International Geology Review, v. 33, p. 107–114, https://doi.org/10.1080/00206819109465676.

Sessa, J.A., Fraass, A.J., LeVay, L.J., Jamson, K.M., and Peters, S.E., 2023, The Extending Ocean Drilling Pursuits (eODP) project: Synthesizing scientific ocean drilling data: Geochemistry, Geophysics, Geosystems, v. 24, https://doi.org/10.1029/2022GC010655.

Tasistro-Hart, A.R., and Macdonald, F.A., 2023, Phanerozoic flooding of North America and the Great Unconformity: Proceedings of the National Academy of Sciences of the United States of America, v. 120, https://doi.org/10.1073/pnas.2309084120.

Tetley, M.G., Williams, S.E., Gurnis, M., Flament, N., and Müller, R.D., 2019, Constraining absolute plate motions since the Triassic: Journal of Geophysical Research: Solid Earth, v. 124, p. 7231–7258, https://doi.org/10.1029/2019JB017442.

Uhen, M.D., et al., 2023, Paleobiology Database user guide version 1.0: PaleoBios, v. 40, no. 11, https://doi.org/10.5070/P9401160531.

Varela, S., González-Hernández, J., Sgarbi, L.F., Marshall, C., Uhen, M.D., Peters, S., and McClennen, M., 2015, paleobioDB: An R package for downloading, visualizing and processing data from the Paleobiology Database: Ecography, v. 38, p. 419–425, https://doi.org/10.1111/ecog.01154.

Vidaña, S.D., and Goring, S.J., 2023, neotoma2: An R package to access data from the Neotoma Paleoecology Database: Journal of Open Source Software, v. 8, 5561, https://doi.org/10.21105/joss.05561.

Wickham, H., 2011, testthat: Get started with testing: The R Journal, v. 3, p. 5–10, https://doi.org/10.32614/RJ-2011-002.

Wickham, H., 2016, ggplot2: Elegant Graphics for Data Analysis: Springer-Verlag, New York, 213 p., https://doi.org/10.1007/978-0-387-98141-3.

Wickham, H., 2023, httr: Tools for working with URLs and HTTP: R package, version 1.4.7, https://doi.org/10.32614/CRAN.package.httr (last accessed 19 September 2024).

Wickham, H., Hester, J., Chang, W., and Bryan, J., 2022, devtools: Tools to make developing R packages easier: R package, version 2.4.5, https://doi.org/10.32614/CRAN.package.devtools (last accessed 19 September 2024).

Wickham, H., Danenberg, P., Csárdi, G., and Eugster, M., 2024, roxygen2: In-line documentation for R: R package, version 7.3.2, https://doi.org/10.32614/CRAN.package.roxygen2 (last accessed 19 September 2024).

Wright, N., Zahirovic, S., Müller, R.D., and Seton, M., 2013, Towards community-driven paleogeographic reconstructions: Integrating open-access paleogeographic and paleobiology data with plate tectonics: Biogeosciences, v. 10, p. 1529–1541, https://doi.org/10.5194/bg-10-1529-2013.

Ye, S., and Peters, S.E., 2023, Bedrock geological map predictions for Phanerozoic fossil occurrences: Paleobiology, v. 49, p. 394–413, https://doi.org/10.1017/pab.2022.46.