

# popler: An R package for synthesis of population time series from long-term ecological research

Aldo Compagnoni<sup>\*a</sup>, Andrew J. Bibian<sup>a</sup>, Brad M. Ochocki<sup>a</sup>, Sam Levin<sup>b</sup>,  
Margaret O'Brien<sup>c</sup>, Kai Zhu<sup>d</sup> and Tom E.X. Miller<sup>a</sup>

<sup>a</sup>Department of BioSciences, Program in Ecology and Evolutionary Biology, Rice  
University, Houston, TX USA

<sup>b</sup>Sam's affiliation

<sup>c</sup>Margaret's affiliation

<sup>d</sup>Kai's affiliation

Running headline: The popler database and R package

---

<sup>\*</sup>email@aldo.edu

## Abstract

*[Very rough draft, will need to be improved.]*

1. Population dynamics has a central role in the historical and current development of both fundamental and applied ecological science. The nascent culture of open data promises to increase the value of population dynamics studies to the field of ecology. However, users interested in data synthesis using population data are stymied by the heterogeneity of the available datasets, and by how difficult it is to identify relevant datasets.
  2. To obviate these issues, we built popler, a relational database that can accommodate the vast majority of population data under a common structure, without the need for aggregating raw observations. Such common structure allows the user to more easily identify, compare, and manipulate datasets.
  3. Moreover, the availability of raw data confers maximum freedom when devising scientific analyses. We used popler to store population data originated by the United States Long Term Ecological Research (LTER) Network.
  4. Finally, we developed an R package to explore and query this resource.
- ...

## Keywords

## Introduction

1 Population dynamics – changes in species’ abundance and composition through time and space  
2 – are central to ecology for both applied and fundamental reasons. Population studies are a  
3 standard tool in conservation, particularly to perform risk assessment and viability analysis  
4 (Morris and Doak, 2002; Beissinger and McCullough, 2002). Population studies are also an  
5 integral part of fundamental ecological research. Populations are the building blocks of ecological  
6 dynamics at higher scales of organization, and examples abound showing how the study of  
7 population ecology improves understanding in evolution (Metcalf and Pavard, 2007), community  
8 ecology (Levine and HilleRisLambers, 2009), and ecosystem ecology (Medvigy et al., 2009; Fisher  
9 et al., 2018).

10 Given their central role in ecology, studies of population dynamics will be an essential com-  
11 ponent in the advances allowed by the flourishing culture of open access and data synthesis. The  
12 increase in freely available ecological data sets is poised to change ecological science (Laurance  
13 et al., 2016). The rising focus on open data is clear in changing publishing standards, in the de-  
14 sign of observational networks like the National Ecological Observatory Network (Schimel et al.,  
15 2007), and in the online availability of previously proprietary data (Kratz et al., 2003; Bechtold  
16 et al., 2005). This deluge of open data holds tremendous promise to facilitate comparative analy-  
17 ses and to test the generality of ecological hypotheses. For population dynamics in particular, it  
18 is the increasing availability of long-term data that will likely yield the most substantial scientific  
19 advances, as long time series are required to detect trends in abundance (Lindenmayer et al.,  
20 2012), quantify temporal variance (Compagnoni et al., 2016), and identify endogenous (Knape  
21 and de Valpine, 2012) or exogenous (Knape and de Valpine, 2011; Hampton et al., 2013) drivers  
22 of population fluctuations.

23 There is ample evidence demonstrating the importance of long-term data for population  
24 ecology to facilitate comparative and synthetic studies (Lindenmayer et al., 2012; Giron-Nava  
25 et al., 2017). To our knowledge, there is currently just one publicly accessible database focused on  
26 long-term population dynamics: the Global Population Dynamics Database (GPDD, Inchausti  
27 and Halley, 2001). The GPDD provides over 5000 time series of population size longer than 10  
28 years for over 1800 animal species. This database has been powerfully leveraged for comparative

analyses and syntheses (e.g., Knape and de Valpine, 2012) but it has some important limitations. GPDD time series are not spatially replicated – there is one observation of population size or density temporal replicate, with no estimate of uncertainty – making it difficult or impossible to isolate different sources of variability and requiring the assumption that population size is known without error. Retaining information on spatial replication would expand the variety of research questions that can be asked and improve estimation of uncertainty in the answers. Additionally, the GPDD focuses exclusively on single species dynamics, making it difficult or impossible to link the dynamics of multiple fluctuating populations within ecological communities (e.g., Ushio et al., 2018).

One of the best sources of publicly available long-term observations of population abundance of diverse plants and animals is the Long-Term Ecological Research (LTER) network supported by the U.S. National Science Foundation. The LTER was founded in 1980 and grew from the original six sites to the current 24 sites throughout North America plus one each in Puerto Rico and Antarctica [*current number of sites is 28*]. Synthetic and comparative studies from the LTER network have made valuable contributions to ecological understanding (Knapp et al., 2012). However, the majority of LTER synthesis research has focused on ecological dynamics at the community (e.g. Wilcox et al. (2017)) and ecosystem (e.g. Knapp and Smith (2001)) scales. Nevertheless, every LTER site collects population abundance data as one of its five core areas of continuous observations (Callahan, 1984). These population time series include both single- and multi-species studies. In our opinion, these data, which have been accumulating since 1980, are under-used. One issue that may limit the use of LTER population data in synthetic, comparative studies is their heterogeneity. Indeed, several authors have noticed that progress in ecology is hindered by difficulties in the way data is accessed, by heterogeneities in the way data is structured, and by failure to cite data originators [*Jones et al. 2004, Reichman et al. 2011*]. The structure of LTER data sets may be widely different, employing a variety of data types (counts of individuals, biomass estimates, percent cover, etc.), experimental designs driven by the priorities of particular PIs, and diverse replication schemes – idiosyncrasies that may be difficult to accommodate in a one-size-fits-all database. However, these challenges also present valuable opportunities. For example, the hierarchical replication structure of many LTER studies (e.g., subplots within plots within transects) can facilitate more sophisticated statistical investigation

59 than would be possible with simpler, aggregated, or unreplicated data.

60 To overcome the issues posed by heterogeneous data structures, we developed `popler` (POP-  
61 ulation dynamics in Long-term Ecological Research), an online database of LTER population  
62 studies. We also developed a companion R package to aid in discovery, querying, and synthesis.  
63 The `popler` database defines a common data structure to accommodate nearly every long-term  
64 population data set provided by the LTER network as of this writing. The common data struc-  
65 ture facilitates the identification, access, and manipulation of raw population data through a  
66 user-friendly R package. Our goals here are to provide introductions to the `popler` database  
67 and R package, which we believe will be useful resources in population, community, and macro  
68 ecology. Our focus here is on LTER time series but our database schema can in principle, ac-  
69 commodate any population level dataset and expanding `popler` beyond the LTER network is a  
70 priority for future development.

## 71 **The `popler` database**

72 `popler` aims to organize population data from the LTER network using a common structure.  
73 To achieve this, we identified a set of variables relevant to population studies (Table 1) and  
74 organized them into a relational database (Fig. 1). In `popler`, we store “raw” data, warts and  
75 all; we have not modified or aggregated the original observations. Our goal was to organize data  
76 under a common structure. If needed, this structure facilitates downstream aggregation and  
77 subsetting of the original data.

## 78 **Population data**

79 Time series of population size are the heart of the `popler` database. We defined ‘population data’  
80 as time-series of observations on the size or density of a population of a species or other taxonomic  
81 unit. Observations of population size are stored in a variable called `abundance_observation`  
82 and can be measured as a count, biomass, density, or cover. Counts and biomass are absolute  
83 measures, density is a measure of population size taken with respect to a one-, two-, or three-  
84 dimensional unit of observation, and cover is a relative (e.g., percentage) or absolute (e.g., area)  
85 measure. These four types of population data are stored in the homonymous tables of the

86 database (Fig. 1A).

87 The population datasets contained in `popler` are always replicated temporally. Temporal  
88 replicates are identified with up to three variables: `year`, `month`, and `day`. In theory, population  
89 data can be collected at a smaller temporal resolution (e.g. `hourly`). However, we did not  
90 encounter LTER datasets whose sampling was more frequent than daily.

91 Population data are also almost always spatially replicated. Hence, they provide measures  
92 of population abundance at multiple locations. Furthermore, spatial replicates are often nested,  
93 where for example a study might be replicated at separate sites (e.g. a watershed, a lake, an  
94 island, etc.), each of which contains intermediate spatial replicates (e.g. a transect, a block),  
95 which in turn contain the smallest spatial replicate at which observations are made (e.g. a plot,  
96 a quadrat). The hypothetical study described above would have three nested levels of spatial  
97 replication, identified by three numbered `spatial replication` variables. In `popler`, we  
98 accommodate data sets with up to five spatial replication levels (Table 1). For simplicity, we call  
99 the first, and therefore largest spatial replicate “study site” (Fig. 1C) – which is not the LTER  
100 site, which is identified elsewhere in the study’s metadata (Fig. 1D), but rather a location within  
101 it.

102 `popler` contains both observational and experimental studies. Experimental datasets con-  
103 tained information on one or more experimental treatments. `Popler` accommodates information  
104 on up to three experimental treatments, identified by three numbered `treatment type` vari-  
105 ables (Table 1).

106 Most datasets contain one or more variables in addition to the ones described above which we  
107 store in a list of variable called `covariates`. Covariates might report, for example, the hour of  
108 day or the temperature when population data was collected, the Cartesian coordinates of spatial  
109 replicates, and so on. Covariates are particularly useful for time series that contain information  
110 on population structure *[Would this be indicated in the metadata? How would someone search*  
111 *for such studies or know if they are working with one?]*. In these datasets, population size refers  
112 to subdivisions of the entire population. These subdivisions can be, for example, males and  
113 females, large and small individuals, developmental different stages, and so forth.

114 Finally, in addition to time series of abundance, `popler` contains datasets that provide  
115 information on the states or attributes of the individuals, or a subset thereof, that make up

116 a population. We store this information in a dedicated table (“Individual’””, Fig. 1A). As  
 117 individual attributes we consider variables that describe identity, size, sex, life stage (e.g. instar  
 118 versus adult in insects) or status (e.g. reproductive or non-reproductive, live or dead, etc.). We  
 119 refer to these individual attributes with the term “structure”: `popler` accommodates data sets  
 120 that measure up to four types of structure simultaneously. We store these data in up to four  
 121 numbered `structure_type` variables. While these data are not population time series, we  
 122 chose to include them in `popler` because they provide information on demographic transitions  
 123 that can be used to derive estimates of population growth. Moreover, in the cases of datasets  
 124 that sample all of the individuals in a population, individuals can be aggregated (i.e. summed)  
 125 as a measure of population size. *[Are ‘structure’ and ‘individual’ operationally the same thing?*  
 126 *Would there ever be individual data that is not somehow ‘structured’? Seems like no, especially*  
 127 *if identity is treated as structure. I am wondering if we can simplify this by just calling this*  
 128 *individual-level data, and saying that this will always be accompanied by one or more individual*  
 129 *attributes. We introduce the concept of structure above, as a possible covariates for time series.*  
 130 *Here ‘structure’ has a related but different interpretation. I would like to prevent confusion.]*

## 131 **Taxonomic information**

132 Each observation of abundance or individual state corresponds to a taxonomic unit (Fig. 1B).  
 133 Population data typically refer to a species or a genus, but we also include data that refer to  
 134 a higher taxonomic rank, such as family, or order. `Popler` provides 15 taxonomic ranks, and  
 135 two additional variables that refer to how taxonomic information is recorded in the original  
 136 datasets. The additional variables are `sppcode`, which are taxon-specific alphanumeric codes,  
 137 and `common_name`, the common name of each taxonomic unit (Table S1).

138 `Popler` stores the taxonomic information linked to each study in two tables: one containing  
 139 the original taxonomic information reported by the PI, the other containing the accepted tax-  
 140 onomic information derived from the former (Fig. 1B). Raw taxonomic data typically contains  
 141 ambiguities derived by the dynamic changes in species classifications (Chamberlain and Szöcs,  
 142 2013). The raw data also typically fail to include higher-level taxonomic information above the  
 143 genus level. To provide as much taxonomic information as possible, `popler` provides a second

144 table linking taxonomic units provided by the authors to accepted taxonomic units according to  
145 the algorithms provided by the R package `taxize` (Chamberlain and Szöcs, 2013). This package  
146 links taxon names to the unambiguous entry of an online taxonomic database.

## 147 Study site

148 We stored the locations of datasets by recording the latitude (`lat_study_site`) and longitude  
149 (`lng_study_site`) of study sites (Fig. 1C). Storing this information in a separate table allows for  
150 explicit connections between independent data sets collected at the same locations within LTER  
151 sites. We emphasize that “study sites” refer to the first and therefore largest level of nested  
152 spatial replication. These study sites are at a smaller scale than “LTER sites” (the 26 *get this*  
153 *number right* NSF-supported locations; Table A1).

## 154 Metadata

155 The metadata table describes 48 variables (Table S2) with general information about the dataset,  
156 its temporal and spatial replication, and its study design (Fig. 1D). *[It would be helpful to clarify*  
157 *how these 48 variables relate to the variables in Table 1.]* The general information includes a  
158 title designated by the original authors, link to online metadata, contact information for data  
159 originators, and the type of data provided by the dataset (i.e., which of the five tables in Fig. 1A  
160 the data is stored in). The metadata also contains variables on temporal replication *[elaborate*  
161 *on what this means. Is it total duration?]* (in years) and spatial replication (number of levels of  
162 nested spatial replicates, and the number of replicates for each spatially nested level). Finally, we  
163 provide basic information on the design of the study that originated the dataset. Studies could  
164 focus on a single species or on multiple species. *[Should we explicitly identify the community*  
165 *variable here? Also, this is not represented in Fig. 1D]* Studies can also be observational or  
166 experimental; if a study is experimental, we provide information on the type of treatments  
167 imposed by the study. Finally, while the objective of `popler` is to store raw population data,  
168 some publicly available datasets provide aggregated data. In this case, we consider these data as  
169 “derived”. *[Again, this is not represented in Fig. 1D. I think the correspondence between the text,*  
170 *Fig 1, and Table 1 are all much improved, but this can still be further tightened and cleaned up.]*



## 171 Creation and contents of the **popler** database

172 We implemented `popler` as an instance of PostgreSQL 9.2.18. We selected the candidate  
173 datasets online, pre-processed the associated files, and uploaded them in the database through a  
174 graphic user interface we developed in Python. We selected the candidate datasets by scraping  
175 information from the internet sites associated with each one of the 26 LTER sites (Table A1).  
176 We scraped this data using libraries `scrapy` in Python, and `rvest` and `RSelenium` in R.

177 We included or excluded studies based on several criteria. First, we only included studies that  
178 reported repeated observations of populations or individuals through time. Second, we included  
179 studies with at least five years of data. We performed this selection in early 2017, so datasets that  
180 are not currently stored in `popler` might now meet this criterion. We did not require population  
181 censuses to be consecutive (some studies have an irregular sampling frequency). Third, we se-  
182 lected only datasets for which the observations had corresponding taxonomic information, which  
183 implies that we ignored time series data on species richness or on the abundance of functional  
184 types.

185 Before uploading to the online database, we preprocessed datasets in python and R. We often  
186 had to combine the separate files associated with the same dataset, transform datasets from  
187 “wide” to “long” form, convert non-ASCII characters, and handling ambiguous study sites. We  
188 provide the details of this pre-processing in Appendix S1.

189 The `popler` database resulting from this process contains data from 215 studies (102 of  
190 which are experimental) representing 2574 cumulative years of observations with a mean study  
191 duration of 12.03 years. `popler` contains data from 691 plant species, 349 animal species, and  
192 1 fungal species. *[Note for Aldo: I thought we decided that `dictionary(variable)` would return  
193 the unique levels of factors and quantiles of continuous variables. This is not how it is working.  
194 `dictionary(duration_years)` returns unique levels of years. Also, the output shows studies less  
195 than 5 years, including zero years and one NA. This is not consistent with the criteria stated in  
196 the text. This will need some attention. Lastly, we should probably cut ‘Chromalveolata’ from the  
197 kingdom data.]*

## 198 The **popler** package

199 The **popler** R package consists of three core functions that allow users to browse and retrieve  
200 data from the database (Fig. 2). In order of intended use, these functions are: `dictionary()`,  
201 `browse()`, and `get_data()`. This order of use revolves around a summary table made up of  
202 the variables contained in the metadata (Fig. Fig. 1D) and taxonomic (Fig. Fig. 1B) tables of  
203 the online database. Each row of the summary table corresponds to a dataset and each column  
204 corresponds to one of the metadata variables that describe it. The `dictionary()` function  
205 explains what these variables (or columns) are, and it shows their content. This content provides  
206 the information needed to subset and visualize metadata through the `browse()` function, which  
207 returns studies (or rows) that satisfy the subset criteria. Finally, following exploration of the  
208 metadata, raw data can be downloaded with the `get_data()` function. Below we describe  
209 these functions and illustrate their use in greater detail.

## 210 The **dictionary()** function

211 The `dictionary` function is a good place for new users to begin working with **popler** (Fig. 2).  
212 With no arguments provided, this function returns a subset of the most useful metadata variables  
213 associated with each dataset (Fig. 1):

##	variable	description
----	----------	-------------

## 1	title	title of project
## 2	proj_metadata_key	unique project id
## 3	lterid	lter name
## 4	datatype	type of abundance data (e.g. count,biomass)
## 5	studytype	experimental or observational study?
## 6	duration_years	duration of project in years
## 7	community	does data set contain multiple taxa?
## 8	structure	types of individual structure
## 9	treatment	types of treatment
## 10	lat_lter	lter site latitude
## 11	lng_lter	lter site longitude
## 12	species	specific epithet of a taxonomic unit
## 13	kingdom	kingdom
## 14	phylum	phylum
## 15	class	class
## 16	order	order
## 17	family	family
## 18	genus	genus

214        Setting argument `full_tbl` to `TRUE` returns all 76 metadata variables. Each one of these  
215 variables name can be provided as an argument to `dictionary()`, which will then return the  
216 possible unique values of the variable. For example, `dictionary(lterid)` returns the three  
217 letter codes for all sites in the LTER network, `dictionary(genus)` returns all the genera  
218 represented in the database, etc. This output provides the starting point for developing criteria  
219 with which to identify relevant datasets. *[I cut the `durationyears` example because I don't think*  
220 *it makes sense to return unique values of continuous variables.]*

## 221 The `browse()` function

222 Once the user is familiar with the meaning and content of the variables that define popler  
223 datasets, they are ready to dig deeper using `browse()` (Fig. 2). Running `browse()` without  
224 arguments provides the metadata from all of the datasets currently stored in popler. This will be  
225 a 215\*19 data frame, with each row corresponding to a study and each column corresponding to  
226 a variable defined by `dictionary()`. *[Not sure this information is very meaningful here: Note  
227 that the taxonomic information associated with each dataset (potentially including many species)  
228 is contained in a list in the column `taxas`. This format allows `browse()` to return a relatively  
229 small, and therefore manageable data frame.]*

230 The full strength of `browse()` is achieved by subsetting studies according to desired criteria  
231 using standard logical expressions. For example, the user might want to consider only studies  
232 whose duration is 30 years or greater, which can be subsetted with:

```
LTER_30 <- browse( duration_years > 29)
```

233 This operation will create the object `LTER_30`, which provides metadata for the data sets  
234 that satisfy the specified criterion. Multiple criteria may be combined. For example, 30+ year  
235 studies of plants can be browsed with

```
LTER_30_plants <- browse( duration_years > 29 & kingdom == "Plantae")
```

236 It is at the `browse()` stage that users should vet the data sets that meet their criteria.  
237 To facilitate data exploration, `browse()` output can be printed in a more readable settings by  
238 providing `report = TRUE` as an argument, which opens up a formatted html document (Figure  
239 3). The metadata provided by `browse()` not only contains information on the characteristics  
240 of a study but also information on how to cite the study, its unique identifiers, including doi,  
241 Knowledge Network for Biocomplexity (<http://knb.ecoinformatics.org>) catalog system  
242 identifier, the contact information of study PIs, and a hyperlink to the url providing the original  
243 data.

## 244 The `get_data()` function

245 *[Random observation: the study titles in popler are sometimes in quotes, sometimes not. This*

246 *should be cleaned up. Also I would like to edit the default message of getdata() before we publish.]*

247 Once data sets of interest have been identified, `get_data()` downloads the data from a  
248 server that hosts the database. This function can take as its first argument a browse object, a  
249 logical expression, either, or both. For example,

```
LTER_30_dat <- get_data(LTER_30)
```

250 downloads the raw data from the studies described in the browse object `LTER_30` and

```
BNZ_30_dat <- get_data(LTER_30, lterid == "BNZ ")
```

251 *[Why is there a space after BNZ? Can we clean this up? It's a little embarrassing for a cherry-*  
252 *picked example.]* downloads the subset of data sets in `LTER_30` from the Bonanza Creek LTER  
253 site. The three letter code “BNZ” could have been located running `dictionary(lterid)` (we  
254 show the meaning of three letters codes in Table S3). The data downloaded from `popler` are  
255 in “long” form, meaning that each row of data reports a single measure of population size, and  
256 separate variables indicate the temporal and spatial replicate, taxonomic information, etc. This  
257 format makes it easy to further subset downloaded datasets with the aim of visualization and  
258 analysis. For example, to visualize the yearly population size of a single species across time and  
259 its spatial replicates, we can subset the study, the species, and flag its spatial replicates with  
260 a different color. We provide such an example using a 33-year long data set collected at the  
261 Bonanza Creek LTER. We select a dataset that measures the population size of plants as counts  
262 of observed individuals. We highlight one example species, *Viburnum edulis* *[Aldo wrote: ???]*  
263 and color code the highest level of spatial replication (the study sites of Fig. 1B). The following  
264 code draws the plot shown in Fig. ??.

```
plot(abundance_observation ~ year,  
     col = as.factor(spatial_replication_level_1),  
     data = subset(BNZ_30_dat, proj_metadata_key == 195 & sppcode == 'VIBUEDUL') )
```

265 This plot shows three idiosyncrasies of this particular dataset. First, because this study  
266 includes three nested spatial replication levels, each site contains several spatial replicates and  
267 therefore multiple points with the same color. Second, the temporal resolution of this dataset

is not consistent. Observations were collected at unequal intervals and, while the first and last observation are 33 years apart, the sampling was concentrated between 1983 and 2003, with only one or two sites sampled before and after. *[It would be nice to take this example a little further, drawing on the online study description - the type of digging one would have to do to really understand the data. For example, the abundance observation is a count, but per what sampling unit? Are these counts in a plot, counts on a transect, etc? We can't tell just from popler but we can probably get that information by reading the study design. It would be good to add something to this effect.]*

## Ancillary functions

popler also provides three additional functions to open the url of the original dataset, unpack covariates, and provide a citation for each dataset. First, we encourage users to consult the online study description associated with each dataset before starting their analyses. As described above, we have not modified the original datasets ; rather, we have rearranged their structure and added accepted taxonomies when possible. The function `metadata_url()` takes as its argument data objects produced by either the `browse()` or `get_data()` functions and launches a web browser. Second, as described above, popler stores all covariates associated with each observation but does not deliver these by default. To access covariates, the `cov_unpack()` function takes an object created by `get_data()` and extracts an R data frame with rows corresponding to abundance observations and columns corresponding to any covariates provided by the original PI. To interpret these covariates, users should read the documentation of the respective dataset using `metadata_url()`. Third we strongly encourage users of popler data to cite data originators and we provide the function `popler_citation()` to do so. The argument of this function is, again, a data object produced by either `browse()` or `get_data()`.

## Limitations and opportunities for development

*[I did not do too much to this section. I think it needs to be stronger but we can continue to polish this. Probably OK for sending to coauthors.]* Working with raw, spatially replicated, and non-aggregated data provides key advantages in quantitative analyses of population dynamics,

295 and these advantages were a driving force behind the development of `popler`. However, because  
296 we did not aggregate data, the user needs examine individual datasets and the associated online  
297 study descriptions to understand their peculiarities. This is important for several reasons. First,  
298 some studies do not define the difference between a missing value and an abundance observation  
299 of zero. Second, many datasets have gaps or changes in the design during the length of the  
300 study. Third, the covariates variable can hold key information which is best understood when  
301 consulting the online documentation of the original dataset.

302 In the future, there are opportunities to increase the size of `popler` and expand its scope.  
303 First, many of the studies included in `popler` are ongoing, so there will be opportunities to  
304 update these entries in the future with new observations. Second, because our schema (Fig.  
305 1) is very general, the database could be expanded to include population datasets outside of  
306 the LTER network. Third, it would be valuable to explicitly associate `popler`'s population-  
307 level data with environmental drivers, especially climate. It is our intention and hope that the  
308 resources provided by `popler` will advance ecological understanding of population dynamics  
309 within the LTER network, and more generally. *[I modified this paragraph to read 'X could be*  
310 *done' rather than 'We will do X'. Y'all have fun but I am getting out of the business of wrangling*  
311 *other people's messy data.]*

## 312 Acknowledgements

313 Support for database and package development was provided by the National Science Foundation  
314 (DEB-1543651). *[We should probably acknowledge NSF support of LTER, as well as coauthor-*  
315 *specific acknowledgements. ]*

## 316 References

317 W. A. Bechtold, P. L. Patterson, et al. *The enhanced forest inventory and analysis program:*  
318 *national sampling design and estimation procedures*, volume 80. US Department of Agriculture  
319 Forest Service, Southern Research Station Asheville, North Carolina, 2005.

320 S. R. Beissinger and D. R. McCullough. *Population viability analysis*. University of Chicago  
321 Press, 2002.

322 J. T. Callahan. Long-term ecological research. *BioScience*, 34(6):363–367, 1984.

323 S. A. Chamberlain and E. Szöcs. taxize: taxonomic search and retrieval in r. *F1000Research*, 2,  
324 2013.

325 A. Compagnoni, A. J. Bibian, B. M. Ochocki, H. S. Rogers, E. L. Schultz, M. E. Sneek, B. D.  
326 Elderd, A. M. Iler, D. W. Inouye, H. Jacquemyn, and T. E. X. Miller. The effect of demographic  
327 correlations on the stochastic population dynamics of perennial plants. 86:480–494, 2016. ISSN  
328 0012-9615. doi: 10.1002/ecm.1228.

329 R. A. Fisher, C. D. Koven, W. R. Anderegg, B. O. Christoffersen, M. C. Dietze, C. E. Farrior,  
330 J. A. Holm, G. C. Hurtt, R. G. Knox, P. J. Lawrence, et al. Vegetation demographics in earth  
331 system models: A review of progress and priorities. *Global change biology*, 24(1):35–54, 2018.

332 A. Giron-Nava, C. C. James, A. F. Johnson, D. Dannecker, B. Kolody, A. Lee, M. Nagarkar, G. M.  
333 Pao, H. Ye, D. G. Johns, et al. Quantitative argument for long-term ecological monitoring.  
334 *Marine Ecology Progress Series*, 572:269–274, 2017.

335 S. E. Hampton, E. E. Holmes, L. P. Scheef, M. D. Scheuerell, S. L. Katz, D. E. Pendleton, and  
336 E. J. Ward. Quantifying effects of abiotic and biotic drivers on community dynamics with  
337 multivariate autoregressive (mar) models. *Ecology*, 94(12):2663–2669, 2013.

338 P. Inchausti and J. Halley. Investigating long-term ecological variability using the global popu-  
339 lation dynamics database. *Science*, 293(5530):655–657, 2001.

340 J. Knappe and P. de Valpine. Effects of weather and climate on the dynamics of animal population  
341 time series. *Proceedings of the Royal Society of London B: Biological Sciences*, 278(1708):985–  
342 992, 2011.

343 J. Knappe and P. de Valpine. Are patterns of density dependence in the global population dy-  
344 namics database driven by uncertainty about population abundance? *Ecology letters*, 15(1):  
345 17–23, 2012.



346 A. K. Knapp and M. D. Smith. Variation among biomes in temporal dynamics of aboveground  
347 primary production. *Science*, 291(5503):481–484, 2001.

348 A. K. Knapp, M. D. Smith, S. E. Hobbie, S. L. Collins, T. J. Fahey, G. J. Hansen, D. A. Landis,  
349 K. J. La Pierre, J. M. Melillo, T. R. Seastedt, et al. Past, present, and future roles of long-term  
350 experiments in the Iter network. *BioScience*, 62(4):377–389, 2012.

351 T. K. Kratz, L. A. Deegan, M. E. Harmon, and W. K. Lauenroth. Ecological variability in space  
352 and time: Insights gained from the us Iter program. *AIBS Bulletin*, 53(1):57–67, 2003.

353 W. F. Laurance, F. Achard, S. Peedell, and S. Schmitt. Big data, big opportunities. *Frontiers*  
354 *in Ecology and the Environment*, 14(7):347–347, 2016.

355 J. M. Levine and J. HilleRisLambers. The importance of niches for the maintenance of species  
356 diversity. *Nature*, 461(7261):254, 2009.

357 D. B. Lindenmayer, G. E. Likens, A. Andersen, D. Bowman, C. M. Bull, E. Burns, C. R.  
358 Dickman, A. A. Hoffmann, D. A. Keith, M. J. Liddell, et al. Value of long-term ecological  
359 studies. *Austral Ecology*, 37(7):745–757, 2012.

360 D. Medvigy, S. Wofsy, J. Munger, D. Hollinger, and P. Moorcroft. Mechanistic scaling of ecosys-  
361 tem function and dynamics in space and time: Ecosystem demography model version 2. *Journal*  
362 *of Geophysical Research: Biogeosciences*, 114(G1), 2009.

363 C. J. E. Metcalf and S. Pavard. Why evolutionary biologists should be demographers. *Trends in*  
364 *Ecology & Evolution*, 22(4):205–212, 2007.

365 W. Morris and D. Doak. Quantitative conservation biology; theory and practice in conservation  
366 biology. *Sinauer, Sunderland, Massachusetts, USA*, 2002.

367 D. Schimel, W. Hargrove, F. Hoffman, and J. MacMahon. Neon: A hierarchically designed  
368 national ecological network. *Frontiers in Ecology and the Environment*, 5(2):59–59, 2007.

369 M. Ushio, C.-H. Hsieh, R. Masuda, E. R. Deyle, H. Ye, C.-W. Chang, G. Sugihara, and M. Kon-  
370 doh. Fluctuating interaction network and time-varying stability of a natural fish community.  
371 *Nature*, 554:360–363, Feb. 2018. ISSN 1476-4687. doi: 10.1038/nature25504.

372 K. R. Wilcox, A. T. Tredennick, S. E. Koerner, E. Grman, L. M. Hallett, M. L. Avolio, K. J.  
373 La Pierre, G. R. Houseman, F. Isbell, D. S. Johnson, J. M. Alatalo, A. H. Baldwin, E. W.  
374 Bork, E. H. Boughton, W. D. Bowman, A. J. Britton, J. F. Cahill, S. L. Collins, G. Du,  
375 A. Eskelinen, L. Gough, A. Jentsch, C. Kern, K. Klanderud, A. K. Knapp, J. Kreyling, Y. Luo,  
376 J. R. McLaren, P. Megonigal, V. Onipchenko, J. Prevéy, J. N. Price, C. H. Robinson, O. E.  
377 Sala, M. D. Smith, N. A. Soudzilovskaia, L. Souza, D. Tilman, S. R. White, Z. Xu, L. Yahdjian,  
378 Q. Yu, P. Zhang, and Y. Zhang. Asynchrony among local communities stabilises ecosystem  
379 function of metacommunities. *Ecology letters*, 20:1534–1545, Dec. 2017. ISSN 1461-0248. doi:  
380 10.1111/ele.12861.

Table 1: Variables used to store population or individual data in popler.

Variable	Description
abundance_observation	Measure of population abundance at a specific time and location. This variable measures abundance as a count, biomass, density, or cover. <i>[What is this value for individual data sets?]</i>
day	Day of observation
month	Month of observation
year	Year of observation
spatial_replicate_n	The $n^{th}$ level of spatial replication, where spatial_replicate_1 is the study site. popler accommodates up to five levels of spatial replication.
treatment_type_n	For datasets originating from an experimental study, the $n^{th}$ treatment. popler popler accommodates up to three treatments.
covariates	Ancillary observations that do not fall into the standard schema of popler.
structure_type_n	For individual data, these variables measure the $n^{th}$ attribute of individuals (identity, size, sex, status, stage). popler accommodates up to four structure types per dataset.

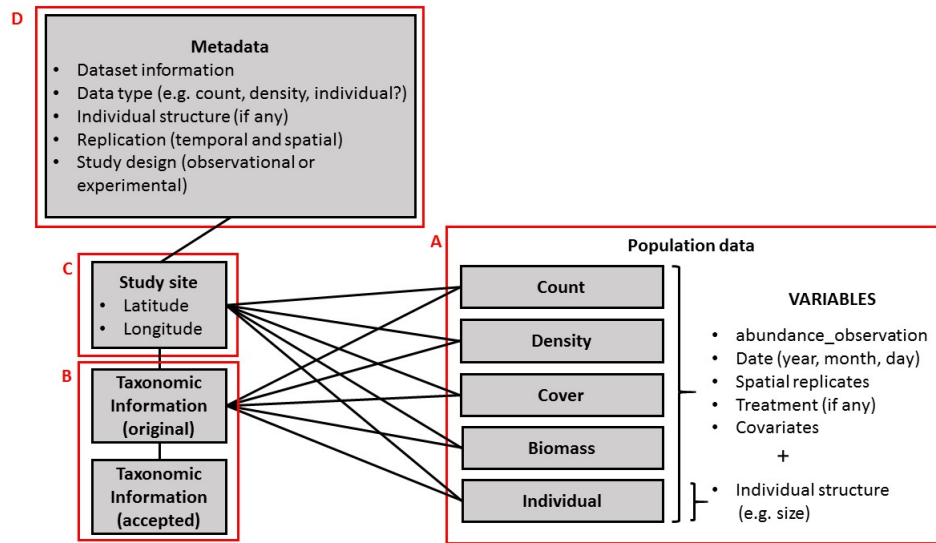


Figure 1: Schematic representation of the entity relationship diagram of the `popler` database. `popler` provides metadata on the studies that originated abundance data points (D). This metadata contains information on the unique identifiers of each study, on its design (observational or experimental), temporal and spatial replication. `Popler` stores the latitude and longitude of the study site (C). Each abundance data point corresponds to a specific taxonomic unit (B). Finally, the time series population data collected in a study can be of five different types: count, density, biomass, cover, and individual data (A).

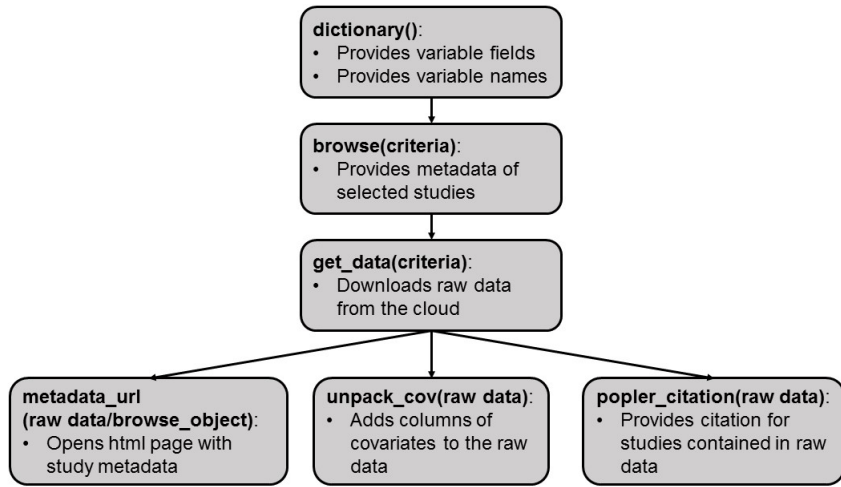


Figure 2: Suggested workflow when using the popler R package to interface with the homonymous online database. The function `dictionary()` refers to the variables of the metadata that describe the data sets contained in popler. `dictionary()` describes these variables and returns their possible values. This information advises which criteria to use when subsetting popler. The user can provide a criterion (that is, a logical statement) to subset the metadata, using `browse()`, and download dataset from the cloud using `get_data()`. Moreover, the output of `get_data()` (a data frame) can be the argument of three ancillary functions: `metadata_url()` opens the webpage containing the original dataset and their associated online metadata. `unpack_cov()` can be used to format the covariates contained in a raw data object into separate columns of a data frame. Finally, `popler_citation()` provides a citation for the downloaded data set(s).

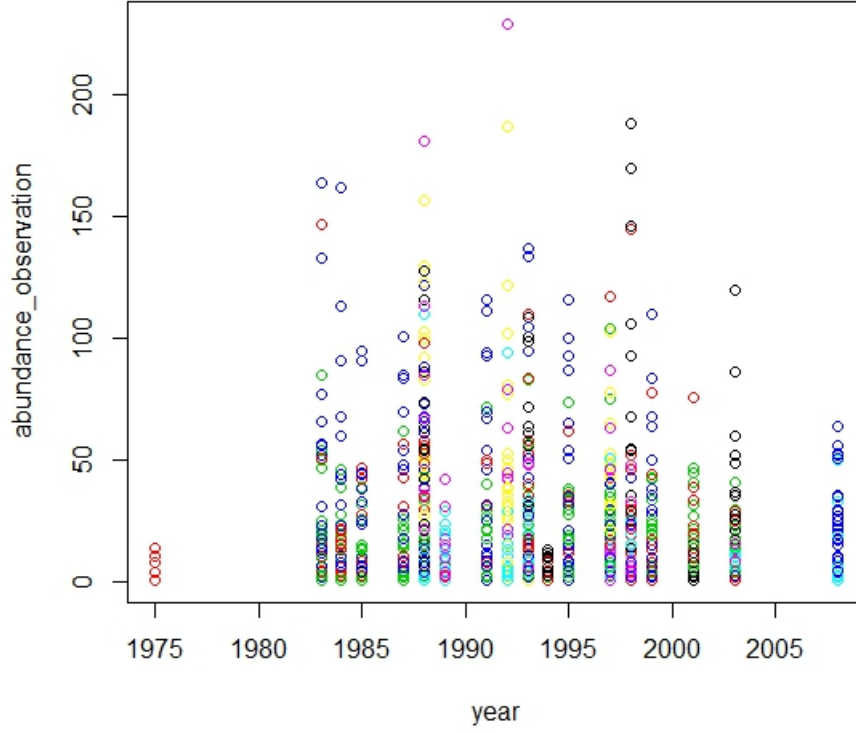


Figure 3: Time series of *Viburnum edulis* abundance counts at Bonanza Creek LTER. The abundance observation is a count and colors correspond to multiple study sites, the largest scale of spatial replication in popler. There are additional, smaller scales of spatial replication in this study, and hence multiple points per site per year.

## 381 Appendix S1: Pre-processing **popler** data

382 Before uploading datasets into the online **popler** database, we combined datasets, transformed  
383 datasets from wide to long form, converted non-ASCII characters, and modified ambiguous study  
384 site names.

385 The variables of many datasets were contained in two or more separate files, which we com-  
386 bined in a single file. When the original dataset provided data in wide form, we transformed  
387 it into long form. In wide form datasets, abundance data associated with different species was  
388 stored in separate columns. **popler** stores these datasets in long form, whereby each row of  
389 abundance data is related to a specific taxonomic unit in the table containing taxonomic infor-  
390 mation (Fig. 1B). We converted all data in ASCII format, because the encoding of the database  
391 is the UTF-8. We often re-defined study site names to unambiguously associate them with one  
392 of the 26 LTER sites. Many site names are alphanumeric codes (e.g. “U1”) which can overlap  
393 across several LTER sites. Hence, we changed site names following a standard formula (namely,  
394 from “U1” to “site\_sbc\_U1”, where “sbc” refers to the Santa Barbara coastal LTER site).

395 In a handful of cases, we removed single data rows from the original dataset. These data  
396 rows were associated with two types of typos in the original dataset. First, some abundance  
397 observations were not associated with a time of observation. We removed this data because  
398 **popler** can only accommodate population information associated with a time of observation.  
399 Second, a handful of abundance data points were clear typos (e.g. the letter “l” instead of a  
400 numeric value). We substituted these data points with a missing value (NULL in the database).  
401 We uploaded these pre-processed datasets in the **popler** database through a Graphic User  
402 Interface developed in Python using libraries **panda** and **pyqt5**.

Table S1: Taxonomic variables contained in the popler table on original taxonomic information.

Variable
sppcode
kingdom
subkingdom
infrakingdom
superdivision
division
subdivision
superphylum
phylum
subphylum
class
subclass
order
family
genus
species
common_name



Table S2: Metadata variables used to describe the datasets stored in popler.

Variable	Description
proj_metadat_key	Unique ID
lter_project_key	ID of LTER site
lter_project_key	ID of LTER site
title	Title of study
samplingunits	Unit of measure (if any) referred to population data.
datatype	Data type: count, biomass, cover, density, and individual. These correspond to the tables in Fig. 1A.
structured_type_n	If individual data, this shows what type of structure is stored. A study can contain up to $n = 4$ types of structure. <i>[is it 'structure' or 'structured']</i>
structured_type_n_units	Unit of measure (if any) referred to structure data.
studystartyr	Start year of the study
studyendyr	End year of the study
duration_years	Duration of the study in years. <i>[Is this derived as end-start?]</i>
samplefreq	Frequency of population census <i>[What is this value for studies with inconsistent frequency as in BNZ example?]</i>
studytype	Whether study is observational or experimental
community	Whether study includes single taxon (community = F) or multiple taxa (community = T)
spatial_replication_level	Extent of spatial replication level number $n$ . A dataset can have up to to 5 replication levels.
spatial_replication_level	Unit of spatial extent of the $n$ spatial replication level.
spatial_replication_level	Label of the spatial replication level (e.g. transect, plot, quadrat, ect.). The label of spatial replication level 1 is "site".
spatial_replication_level	The number of unique replicates for the $n$ th level of spatial replication.
treatment_type_n	The type of treatment. <i>[This is vague and does not correspond to level <math>n</math>. PROBLEM.]</i>
control_group	Is population size data raw, or is it derived (e.g. it is aggregated)? <i>[Problematic because we say the data are raw and non-aggregated.]</i>
authors	Author(s) of the original dataset
authors_contact	Contact of the author(s) associated with the original dataset. <i>[Is this always email address?]</i>
metalink	url of the original dataset
knbid	Knowledge Network for Biocomplexity identifier.

Table S3: LTER identification acronyms and their meaning as used in the popler database.

Variable	LTER name
SBC	Santa Barbara Coastal LTER
SEV	Sevilleta LTER
SGS	Shortgrass Steppe
VCR	Virginia Coastal Reserve LTER
AND	Andrew Forest LTER
NWT	Niwot Ridge LTER
BNZ	Bonanaza Creek LTER
CDR	Cedar Creek Ecosystem Science Reserve
GCE	Georgia Coastal Ecosystems LTER
ARC	Arctic LTER
CAP	Central Arizon - Phoneix LTER
FCE	Florida Coastal Everglades LTER
HFR	Harvard Forest LTER
KBS	Kellogg Biological Station LTER
CWT	Coweeta LTER
HBR	Hubbard Brook LTER
MCM	McMurdo Dry Valleys LTER
JRN	Jornada Basin LTER
CCE	California Current Ecosystem LTER
KNZ	Konza Prairie LTER