

A peer-reviewed version of this preprint was published in PeerJ on 9 October 2018.

[View the peer-reviewed version](https://doi.org/10.7717/peerj.5453) (peerj.com/articles/5453), which is the preferred citable publication unless you specifically need to cite this preprint.

Lazarus DB, Renaudie J, Lenz D, Diver P, Klump J. 2018. Raritas: a program for counting high diversity categorical data with highly unequal abundances. PeerJ 6:e5453 <https://doi.org/10.7717/peerj.5453>

Raritas and RaritasVox: Programs for counting high diversity categorical data with highly unequal abundances

David Lazarus ^{Corresp., 1}, Johan Renaudie ¹, Dorina Lenz ², Patrick Diver ³, Jens Klump ⁴

¹ Museum für Naturkunde, Berlin, Germany

² Leibniz-Institut für Zoo- und Wildtierforschung, Berlin, Germany

³ Divdat Consulting, Wesley, Arkansas, United States

⁴ CSIRO, Mineral Resources, Kensington, Australia

Corresponding Author: David Lazarus

Email address: david.lazarus@mfk.berlin

Acquiring data on the occurrences of many types of difficult to identify objects are often still made by human observation, e.g. in biodiversity and paleontologic research. Existing computer counting programs used to record such data have various limitations, including inflexibility and cost. We describe a pair of new open-source programs for this purpose - Raritas and RaritasVox, which share a similar graphical user interface for mouse based counting, and file output format. Raritas is written in Python and can be run as a standalone app for recent versions of either MacOS or Windows, or from the command line as easily customized source code. RaritasVox in addition supports voice based counting but is written in Java and is more complex to install or modify. Both programs explicitly support a rare category count mode which makes it easier to collect quantitative data on rare categories, e.g. rare species which are important in biodiversity surveys. Lastly, as to our knowledge no standards exist yet, we describe a new stratigraphic occurrence data (SOD) unitary file format which combines extensive metadata and a flexible structure for recording occurrence data of species or other categories in a series of samples.

1 Raritas and RaritasVox: programs for counting high diversity
2 categorical data with highly unequal abundances

3 David B. Lazarus¹, Johan Renaudie¹, Dorina Lenz², Patrick Diver³ and Jens Klump⁴

4 1 - Museum für Naturkunde - Leibniz-Institut für Evolutions- und Biodiversitätsforschung,
5 Berlin, Germany

6 2 - Leibniz-Institut für Zoo- und Wildtierforschung, Berlin, Germany

7 3 - Divdat Consulting, Wesley, Arkansas, USA

8 4 - CSIRO, Mineral Resources, Kensington, Australia

9 Corresponding author - David Lazarus, david.lazarus@mfn.berlin

10 Author contributions

11 DBL created the main program specifications, designed the GUI and wrote the paper. JR wrote
12 Raritas, DLenz and JK designed the voice functions and wrote RaritasVox. DBL and PD created
13 the SOD format.

14 **Abstract**

15 Acquiring data on the occurrences of many types of difficult to identify objects are often still
16 made by human observation, e.g. in biodiversity and paleontologic research. Existing computer
17 counting programs used to record such data have various limitations, including inflexibility and
18 cost. We describe a pair of new open-source programs for this purpose - Raritas and RaritasVox,
19 which share a similar graphical user interface for mouse based counting, and file output format.
20 Raritas is written in Python and can be run as a standalone app for recent versions of either
21 MacOS or Windows, or from the command line as easily customized source code. RaritasVox in
22 addition supports voice based counting but is written in Java and is more complex to install or
23 modify. Both programs explicitly support a rare category count mode which makes it easier to
24 collect quantitative data on rare categories, e.g. rare species which are important in biodiversity
25 surveys. Lastly, as to our knowledge no standards exist yet, we describe a new stratigraphic
26 occurrence data (SOD) unitary file format which combines extensive metadata and a flexible
27 structure for recording occurrence data of species or other categories in a series of samples.

28 Introduction

29 Human observations as a source of scientific data

30 Quantitative data about many aspects of the natural world are collected in modern science with
31 the use of instruments, but a substantial amount of observational data is still collected by human
32 observation. This is particularly common in ecology, organismal biology and behavioral sciences,
33 where the numeric data on the frequencies of occurrences of biologic phenomena are desired, but
34 the objects/phenomena to be counted are too complex to identify by instruments or fully
35 computerized image analysis systems. Up until the spread of desktop computers, such counts
36 were done mostly either with the aid of mechanical counter buttons (including arrays of several
37 buttons, to allow counting of multiple categories) or tallied by hand on printed list forms. Both
38 methods are slow and require re-entering the count values into a computer afterwards before
39 analysis, adding additional time and possibilities for error. Computer 'point-counting' programs
40 can in principle replace these methods and at the same time provide additional functions that
41 mechanical methods cannot, such as continuous statistical summaries of the data as it is being
42 collected, which provides useful feedback to the observer on how complete or accurate the
43 dataset being collected is.

44 Despite these obvious advantages counting programs have yet to fully replace manual methods.
45 There are many reasons for this including cost, inflexibility, compatibility and inadequate ease of
46 use. Numerous inexpensive or free simple tally counter programs are available that can replace
47 mechanical counter buttons (e.g. dozens of simple smartphone/tablet apps, or more sophisticated
48 desktop apps e.g. Versacount: (Kim & DeRisi, 2010). None of these however are well suited to
49 counting larger numbers of categories, which is common in ecology, and in related fields such as
50 paleontology. The need to count many objects in many categories is particularly acute in
51 biodiversity related disciplines, e. g. field surveys of species diversity; species counts of fossil
52 assemblages in micropaleontology. In such studies the diversity of objects and total numbers of
53 objects available for study are both very high. Several programs have been developed to assist in
54 biodiversity assessments (e.g. 'OrgaCount': www.aquaecology.de; 'Becam': www.avansee.com).
55 As many micropaleontologists work in commercial (oil industry) settings, there are also several
56 sophisticated counting programs available (many as commercial products) for counting large
57 numbers of microfossils: ; Polpal (Nalepka & Walanus, 2003); Foramsampler (McGann et al.,
58 2006); Counter (Zippi, 2007); Stratabug (Stratadata, 2014); Bugwin (Bugware, 2016). These
59 programs, whether for biologists or industrial micropaleontologists, however frequently are
60 limited in one or more ways. Many are embedded in larger, more specialized packages with
61 features for a single discipline, e.g. stratified ecologic sampling, biostratigraphic range charting,
62 petrologic thin section analyses. Programs are often complex to install, or are lacking in
63 flexibility, adaptability and/or ease of use. Many are also closed-source, expensive, and are
64 dependent on the commercial provider to maintain. There is thus a need for a program that is
65 relatively simple, free, open-source, less specialized and thus adaptable to counting a variety of
66 different types of objects, and that works with different operating systems. Most importantly, it
67 must be as easy to use as mechanical methods, since a program that is significantly slower will,
68 based on our experience, normally be rejected by users. Users often need to count thousands of
69 objects (see 'Rarity' below), and an even marginally slower data entry method will create an
70 unacceptable cumulative loss of the user's time. This is particularly true in counting objects such
71 as microfossils, or in field biodiversity surveys, where vast numbers of specimens are available
72 and can be quickly identified by the user, making data entry the time-limiting factor in data

73 collection.

74 **Rarity**

75 In addition to the general need for flexible, efficient counting programs, there is also a specific
76 need to count objects which have very different relative abundances. Many classes of objects in
77 the observable world show a characteristic pattern of unequal relative abundances that can be
78 approximated by power laws, including incomes, internet traffic, plankton sizes, and the sizes of
79 interstellar mineral grains (Mathis et al., 1977, Reed & Hughes, 2002, Buonassissi & Dierssen,
80 2010). Biologic entities, in particular species abundances in ecology and paleontology also
81 typically show such distributions, with a few species being relatively common, and the remainder
82 uncommon or quite rare (Preston, 1948, Brown et al., 2002). Counting objects at random from
83 such unevenly distributed populations results in many counts of the few common species, but
84 very few counts of rarer species. For example, in both the complete dataset, and in individual
85 samples, counts of fossil radiolarians in Neogene Southern Ocean sediments show a few very
86 common species, and many rare species (Figs. 1, 2). Even with >700,000 individuals, a
87 substantial fraction of the species are represented by 10 or fewer individuals. Thus, in order to
88 encounter at least one individual of all rare species very large numbers of specimens need to be
89 examined. For example, several thousand individuals needed to be examined in order to recover
90 95% of the estimated total species diversity (ca 200 species) in the single sample counted in Fig.
91 2 (Fig. 3).

92 Ecologists and paleontologists thus sometimes decide to base studies only on the small number of
93 species that are relatively common and thus whose abundances are easy to quantify. Many
94 applied micropaleontologic studies for example use the the environmental preferences of a
95 relatively small number of common species to reconstruct past environmental conditions (Imbrie
96 & Kipp, 1971, CLIMAP project members, 1976). Not all scientific questions can however be
97 addressed by examination of only a small number of common species. Unlike, e.g. mineral
98 grains, each biologic species is unique, with its own potential to contribute to ecosystem function
99 and, over the longer term, to evolutionary change. Biodiversity research in particular is concerned
100 about documenting total species richness and understanding threats to it, e.g. how current and
101 past environmental change affects it. The findings of such research feed into important decisions
102 on biodiversity conservation, land use and other global issues (i.e., the 'Rio' Convention on
103 Biological Diversity: www.cbd.int). Reasonably accurate estimates of total diversity - crucial in
104 biodiversity studies - can only be made when the majority of the diversity has been counted.
105 Extrapolations from less complete data tend to have unacceptably high error values (Colwell et
106 al., 2012). There is thus a major effort to understand the total species richness of modern and past
107 biologic systems (Mora et al., 2011), and consequently, the need to collect quantitative data on
108 many rare species (Roberts et al., 2016).

109 One approach to achieving this is based on the human ability to scan large populations to identify
110 a subset of target individuals much more rapidly than the same person could fully identify and
111 record the identity of each individual in the population. As a simple example, it is much faster to
112 scan a large crowd of people to identify a single category of persons of interest ('tall men with
113 beards'), than to identify each person in a crowd and record all of their names. Similarly, one can
114 quickly skip individuals belonging to a specific category to target other individuals. Biologists
115 and paleontologists collecting data on rare species make use of this ability by first counting all
116 individuals encountered to identify common species, then, mentally blocking out the common
117 species, continuing to count only species that are not in the 'common' group. In this 'rare category'

118 mode individuals of common species can be scanned over much more rapidly, and their counts
119 for the total area viewed estimated afterwards based on their abundances in 'all species' mode.
120 Larger total numbers of individuals are thereby examined, and a better estimate of total species
121 richness can be obtained (Gannon, 1971, Hinds, 1999, Stevenson et al., 2010). A good counting
122 program for such work should offer options that support this style of efficient counting of only
123 rare taxa. This ability is however, to our knowledge, normally not offered in currently available
124 counting programs, which are mostly designed to support counts of smaller numbers of species
125 and individuals in support of applied (paleo)environmental research.

126 **Materials and Methods**

127 Raritas and RaritasVox are two new programs for counting (tallying) multiple categories of
128 objects which meet these criteria. Both offer a flexible mouse-driven interface for counting
129 highly diverse lists of taxa, including both buttons for more common taxa, and hierarchical
130 menus to select rare taxa. An additional feature of the programs is the definition of a new file
131 format for storing such count data that uniquely combines the data and detailed metadata in a
132 user-friendly spreadsheet style layout. Compiled apps, source code, user guides, sample
133 configuration and output files are all publicly available at <https://github.com/plannapus/Raritas>.

134 The programs provide explicit support of dual-mode (all vs rare only) counting, and indeed this
135 feature is the basis for the program names. In standard mode, all individuals seen are counted. In
136 'rare only' mode, commonly occurring objects are no longer counted: only rare objects are. Not
137 having to pause to enter a count for the most frequently seen object types makes counting rare
138 object categories much faster. However, in order to be able to combine counts for common and
139 rare types together, it is also necessary to know the magnitude of observational effort made in
140 each counting mode, as the total frequencies of common objects are estimated for the 'rare objects
141 only' interval based on their frequency in 'all object' counting, and the observational effort spent
142 in 'rare' mode. A computer program that supports rare-only counting must therefore be able to
143 monitor observational effort in parallel to recording individual object counts. This is provided for
144 by a separate counter for observational effort, a 'track' counter which the user updates periodically
145 while counting.

146 The main program Raritas, is written in Python (van Rossum et al., 2010). The second -
147 RaritasVox - is written in Java, and was in fact the initial test development version. This older
148 version provides most, though not all of the features of the main Python version in mouse-based
149 counting. In addition it provides a unique option to register counts directly from voice input by
150 the user, who simply speaks the category names. Regardless of method or program variant, the
151 same type of output, setup and configuration files are used.

152 These programs' ease of use involve both ease of configuration as well as ease of use during
153 primary operation. Raritas and RaritasVox are configured almost entirely from the contents of a
154 simple tabular type file which can be created easily by users using a spreadsheet program. The
155 file contains list of which objects (e.g. species) are to be counted, how these are to be presented to
156 the user (button labels and other details). This also simplifies the program as there is no need to
157 write code for configuration, other than reading the configuration file.

158 Detailed metadata is captured for each dataset and saved with the data in the output files. This

159 often a weakness in other (e.g. commercial) programs where relatively little information is
160 captured. Reliance on program-external metadata capture such as embedding all metadata in
161 filenames is obviously limited in extent, not well structured and in our experience has not been
162 very reliable, particularly when metadata needs to be understandable over the long-term (i.e. by
163 other than the file creators).

164 Raritas been programmed in Python because it is a popular, well supported, and relatively easy to
165 learn multi-paradigm scripting computer language. It is more likely to be understandable to
166 workers in fields such as taxonomy/systematics than the more complex, object-oriented compiled
167 language Java. RaritasVox was programmed in Java in order to make use of specialized libraries
168 for voice recognition: the Sphinx open-source speech recognition engine (Walker et al., 2004)
169 (<http://www.speech.cs.cmu.edu/sphinx/doc/Sphinx.html>), and to insure speed, which is needed
170 for the complex task of voice recognition - Java code executes much faster than Python code.
171 Both programs run quickly on all hardware tested (desktop and laptop computers with Intel 'i'
172 series processors, running Windows 7-10; OS X 10.9-12). Raritas consists of ca 650 lines of
173 Python code; RaritasVox of nearly 4,000 lines of Java. The use of Python, plus the much smaller
174 size of the code, makes customization of the Raritas's features possible by technically savvy
175 users, without the need to employ a professional programmer. Python also provides excellent
176 packages for some functions such as plotting data that allow the program to produce better
177 outputs for the user without having to write additional code (e.g., matplotlib). Python is not
178 without problems - installing the various software modules (packages), including packages used
179 by other packages (dependencies) that an application needs can be very difficult for a non-
180 specialist, depending in part on the local python environment used. Raritas is therefore offered
181 both as a fully bundled program (double-clickable) with all needed packages included for Mac
182 OS X 10.11+ as well as for Windows 7 and 10; and also as source code: the former providing
183 ease-of-use for non specialists; the latter customizability. RaritasVox is also available either as a
184 bundled app (a .jar file) or as source code. The bundled versions are each ca 100 Mb in size.

185 **Installation**

186 No special installation procedure is needed for the Raritas program when used as the bundled
187 app. Using the source code version of Raritas (python) requires installing only two python
188 packages (and their dependencies): matplotlib and wxPython (Hunter, 2007, Dunn, 2014). These
189 must be installed using the appropriate python or OS package manager for the user's python
190 system, which will automatically install any dependencies. Some python distributions already
191 include both packages as part of their standard installation, thus requiring no special installations
192 by the user. RaritasVox requires a Java environment (available for free download, often installed
193 previously in many systems) in addition to the app itself. Installing the source code version of
194 RaritasVox is considerably more complicated: details are given in Appendix 1.

195 **Configuration file and starting the program**

196
197 Both programs read a single configuration file on starting - by default, the one previously used, or
198 a new one chosen by the user. The file (Fig. 4; Appendix 2) is in tab-text format and is just a list
199 of taxa names and how each should be presented to the user in the GUI interface. All names are
200 available by drop-down list by default. Names can also be shown as buttons (with abbreviations
201 to insure the button label fits). If a second set of names of higher level categories are provided
202 for the primary names, the name list is parsed into multiple list with multiple drop-down menus,

203 thus providing structure to longer name lists and more rapid access to taxa names.

204 Bundled versions of either program are started by the usual double-click of the app icon or other
205 standard GUI methods. The source code version of Raritas is started by a standard 'python
206 raritas.py' statement (optionally including a path name, if appropriate) at the command line. Once
207 the program starts all interaction takes place via the GUI interface that then appears. RaritasVox
208 cannot be run directly from the source code as Java is a compiled language - any customized
209 version of the RaritasVox Java code must first be compiled and linked either via the command
210 line or a programming tool such as an IDE.

211 **GUI interface for manual counting**

212 The main elements of the GUI interface for either version, once started, are: the metadata
213 window, the counting window, the rare count configuration window and the collector curve
214 window.

215 Metadata window (Fig. 5). When the program is first started a window appears which provides a
216 pop-up list of primary counting style options (file types), based on the SOD file specification
217 (described below). The next window collects the metadata appropriate for the file type, e.g. field
218 names that are used in the rest of the program for the material to be counted. At the moment the
219 program supports two types of primary data, both for microfossil occurrences: assemblages of
220 microfossils from deep-sea sediments obtained by the international deep-sea drilling programs, or
221 fossils from samples obtained from geologic sections on land, but other types can be defined. The
222 metadata window also provides a few run-time options for configuring the interface and behavior
223 during counting. Importantly, the user chooses which taxa name list configuration file they want
224 to use via a normal file open dialog at this time. When ready the 'start counting' button is clicked
225 and the counting window appears.

226 Counting window (Fig. 6). This is the main window that is used for most interaction with the
227 program. The upper part of the window is populated with the buttons for counting common
228 species, with labels as defined in the configuration file. Less common taxa are shown in the form
229 of popup lists, organized into higher level categories, again as defined in the configuration file.
230 Putting less common taxa into lists and common taxa on buttons allows most counts to be done
231 quickly with a button, while the comparatively slow process of selecting from a list is reduced to
232 a minimum. Lists are needed however as they can be of arbitrary length, while the number of
233 buttons is limited by screen size. Counting is active whenever the window is present. Clicking on
234 a button or selecting a taxa from the lists adds the species to the count data structures. A list of
235 recently counted objects is given in the sub-window (lower middle of main window). A button is
236 provided on the right to count observational effort ('Track', for number of 'tracks' scanned on a
237 microscope slide) and a counter shows the total tracks counted.

238 Clicking on 'Rare Count Mode' brings up a dialog (Fig. 7), where the counted objects are listed in
239 order of descending abundance, and the user can choose which to exclude from further counting.
240 When the dialog is dismissed counting resumes, with, for those taxa to be excluded, the taxa
241 buttons greyed out and pop-up list items inactivated.

242 Determining which species to exclude in rare count mode is not trivial. As this is a key feature of
243 Raritas we include the following suggestions, which are based on our experience of counting ca

244 700,000 total specimens (several thousand specimens per sample in over 100 samples) for the
245 study published in (Renaudie & Lazarus, 2013). The tally to use to trigger the switch to rare-only
246 counting, and the percentage threshold for species to be ignored during 'rare' count mode should,
247 as a rule of thumb, maximize the number of specimens to ignore while minimizing the error on
248 the abundant species percentages. In (Renaudie & Lazarus, 2013), we chose to stop the full count
249 mode when ca. 2,000 specimens were already counted and to ignore in 'rare' count mode species
250 with a percentage higher than ~5% of the community. Doing so allowed us to keep the error to ca.
251 10% of the investigated value. In other words, for a species that was present at 5% abundance in
252 full count mode, the theoretical standard error is slightly below 10% of this 5% value, i. e. a
253 theoretical percentage for the species between ca. 4.5 and ca. 5.5%; (Drooger, in (Zachariasse et
254 al., 1978) (Fig. 8a). These cut-off values eliminated 59.7% of the specimens during rare-only
255 mode (median of all samples counted, but varying from one sample to the other, black line on
256 Fig. 8b for median, dark grey area for interquartile range and light grey are for total range). An
257 additional, important criterion that was taken into consideration is that all samples encountered
258 had at least one species above the 'ignore in rare-only mode' percent threshold. Using an higher
259 threshold than 5% would have meant that some samples would have had to be counted entirely in
260 full count mode, as no species would have been abundant enough to exclude. In our study, there
261 were on average ca three (mean = 2.9) percent of the species above the cut-off threshold per
262 sample (blue and red lines of Fig. 7b).

263 The 'Show Collector's Curve' menu item (Raritas, or button, RaritasVox) brings up the fourth
264 main GUI element - a diversity accumulation plot (Fig. 9) showing the relationship to total
265 number of object types seen (species) vs total number of objects counted (specimens). For
266 typical biologic data these curves show a roughly logarithmic in shape - at first rising rapidly,
267 then, as increasingly species already seen previously are re-encountered, flattening out. The
268 curve's slope will eventually become zero when all object types in the sample have been detected
269 (compare to Fig. 2). The user can decide when the curve has become close enough to this state
270 for his/her purposes, and thus stop counting only when the data completeness quality is adequate.
271 If a series of samples are counted to the point where they have the same apparent slope at the end
272 of this dynamically generated diversity accumulation curve, they will share the property of being
273 'fairly' sampled, and relative differences in diversity will be shown without bias (Alroy, 2010,
274 Colwell et al., 2012). This type of feedback is important to insuring good quality observations
275 and is something that cannot be provided by simple mechanical count systems. It is however
276 rarely implemented in programs known to us.

277 **Voice interface**

278 RaritasVox has a similar GUI to Raritas, with only fairly minor differences in the layout of
279 elements or functional behavior (e.g., RaritasVox allows colors to be assigned to taxa names as an
280 aid to accurate name selection in the interface), and thus is not described separately here - details
281 are given in Appendix 1. The main difference in functionality is the ability to use a voice driven
282 counting mode, selected via a control button from the main counting window. The motivation
283 was the observation that, for some users, the constant change of focus between microscope and
284 counting program (or paper sheet) while counting microfossils under a microscope places a strain
285 on the user's vision. Some researchers affected by this problem had developed a voice-based
286 counting procedure: calling out species identifications and recoding the counts as audio
287 recordings, then later playing them back and transferring the species counts into their counting
288 sheets. RaritasVox was conceived as a way, by using speech recognition, to make this process

289 more efficient and ergonomic.

290 Since 2009 when RaritasVox was developed and today speech recognition has made tremendous
291 advances and has become a commonplace functionality in many everyday applications, e.g.
292 Apple's "Siri". Speech recognition systems can be classified into two categories. "Speaker
293 dependent" systems use "training" (also called "enrollment") where an individual speaker reads
294 text or isolated vocabulary into the system. The system analyzes the person's specific voice and
295 uses it to fine-tune the recognition of that person's speech, resulting in increased accuracy.
296 Systems that do not use training, including RaritasVox, are called "speaker independent" systems.
297 RaritasVox however makes use of the fact that the counting process uses an independent
298 vocabulary that is defined in a configuration file (Fig. 10; Appendix 2). The user may not only
299 use his or her own short terms for species rather than the full taxonomic name, e.g. "pachyleft"
300 instead of "*Globigerina pachyderma sinistral*", they can modify the configuration file so that the
301 program can better recognize an individual's normal pronunciation style. This is for example
302 useful for users with different native languages, as vowels in particular are often pronounced
303 differently, even for latin taxa names. For example "*Prunopyle*" is pronounced proo-no-peil by
304 English speakers, and proo-no-peel-ae by Germans.

305 At the time RaritasVox was first being planned (2009) only a few cross-platform packages were
306 available. The speech recognition software Sphinx and Java were chosen as the best combination
307 for an open-source, cross platform speech recognition package and language environment for our
308 purposes. For Sphinx the elemental components of speech sounds are interchangeably referred to
309 as "phones" or "phonemes" (see <http://www.speech.cs.cmu.edu/sphinx/doc/Sphinx.html> and
310 <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>). Only phonemes listed in the phoneme set of the
311 CMU Pronouncing Dictionary (around 40) can be used and it expects that the language used is
312 English. Only words consisting of one or more phonemes that are present in the customized
313 dictionary file (Fig. 10) can be recognized as "correct". The software will search for words
314 consisting of phonemes present in the dictionary which match best to the speech input. In
315 RaritasVox the spoken word is recognized, confirmation is shown on screen, and a count
316 command for that item is generated (Fig. 11).

317 RaritasVox was not used to collect research data and was only briefly tested for accuracy (Table
318 1).

319 Using a list of 18 words and 108 voice entries, four words were incorrectly identified (<4%),
320 resulting in 8 incorrect counts (7.5%). This is similar to accuracy in much more sophisticated,
321 general voice recognition systems [27], which is possible as RaritasVox uses a very limited
322 vocabulary. The count error rate may be too large for data collection where rare occurrences are
323 important (e.g. biostratigraphy) but adequate for others such as gross assemblage composition,
324 particularly when combined with statistical data reduction procedures such as factor analysis that
325 are insensitive to small amounts of random data scatter [13]. The accuracy is in any event
326 choosable by the user as they can, by monitoring the computer screen, correct errors before they
327 are counted using the spoken 'Remove' command to delete the last (incorrect) identified word.

328 **Output files**

329 **SOD File Format**

330 In addition to the diversity accumulation plots, which can be saved as graphics as often as desired
331 (the matplotlib library used in Raritas supports various file formats, e.g. png, pdf, jpg, tif), the
332 program saves the primary count data. This necessitates choosing, or creating a format for the
333 data files, as there is no universal community database which would allow a direct upload
334 solution. Despite a great deal of biostratigraphic or other data of the form of species by
335 samples/observations having been generated globally for many decades, no generally accepted or
336 even widely known file format exists for such data. Other fields have developed community data
337 formats for such data matrices, e.g. the BIOM format for biological observation matrices
338 (McDonald et al., 2012), as well as standard protocols to exchange information directly between
339 computer systems e.g. Darwin Core (Wieczorek et al., 2012). These formats are however of
340 limited use for paleontologic fossil occurrence matrices since they lack any way to store
341 metadata, general or individual sample, that is related to geologic age (sample position in section,
342 formation name, etc), and the metadata in general is optimized for biologic, not paleontologic
343 observations. One of the major biologic exchange protocols (ABCD: (Berendsohn, 2007),
344 <http://wiki.tdwg.org/ABCD/>) does have, via the EFG extension (<http://www.geocase.eu/efg>) the
345 ability to transmit both biologic and geologic data, but is a communication protocol, not storage
346 format, and the xml definition is not readable by normal users.

347 Within the field of paleontology, data on occurrences, outside of micropaleontology, are
348 dominated by simple taxa lists for a single locality (one sample). This is exemplified by the main
349 data input formats the most widely used paleontology community database PBDB (Alroy et al.,
350 2001), where data is entered, taxon by taxon, for one sample at a time. Within micropaleontology
351 taxa-by-sample data matrices are common (often referred to as 'range charts') but data is usually
352 given in the format of individual publications, without metadata in the files, in numerous
353 variations of a simple taxa-by-sample table. This is also the file format used by the deep-sea
354 drilling programs (DSDP, ODP, IODP), which have not generally captured micropaleontology
355 data except in a very limited form on-ship, using database entry forms, or simply archived data
356 copied from publications, with only minimal metadata stored separately from the data files.
357 Lastly there are several more comprehensive data file formats that are associated with
358 commercial micropaleontology, i.e. the oil industry. These formats include metadata, details of
359 stratigraphy etc, but are not compatible with each other and are mostly meant for internal use in
360 proprietary commercial programs, not for open file exchange. Most also tend to be quite user
361 unfriendly, giving sample and taxa names in separate definition blocks from the actual occurrence
362 data, and use a long, non-tabular, list type structure that makes comprehension difficult. There is
363 thus a need for a public (non-proprietary) file format that combines metadata and the taxa-by-
364 occurrences data in a single file, provides for geologic age or section information and which is
365 easy for scientists to read and use.

366 We have therefore adopted a new 'open file format': Stratigraphic Occurrence Data format, which
367 we abbreviate here simply as SOD format. This format originally was developed in response to
368 the need to merge metadata and occurrence data in user typed files, in order to manage a large
369 number of fossil occurrence matrix files that were being digitized from the literature for upload
370 into a database that provides a micropaleontologic equivalent to the PBDB: NSB (Lazarus, 1994,
371 Spencer-Cervato, 1999). This database reports occurrences of microfossils in deep-sea sediment
372 sections, and the data is mostly derived from studies that report the occurrences in the form of
373 simple samples by species tables, one table per section, per higher fossil group. The file format
374 itself is deliberately meant to be visually similar to the source publication data tables, being
375 essentially an enhanced version of the publication's tabular data matrix. This makes the file easily
376 read by users, and equally makes the transcription (keying-in) of data from publications into the

377 format relatively simple - in some cases, where a publication file is available in digital form,
378 simply by reformatting some of the fields, rather than re-entry of primary values. SOD format
379 however is significantly different from an 'ordinary' user data table in that it is based on a formal,
380 extendable definition of content. This definition adds more structure and detail for both taxa and
381 sample names, and uses the otherwise empty 'corner' of the matrix at the intersection of the row
382 and column labels to include, in a structured way, more general metadata about the occurrence
383 data in the file.

384 The file is laid out in 4 graphical blocks: general metadata: upper left corner block; taxa
385 metadata: left columns below metadata block; sample metadata: rows to right of corner metadata
386 block; and the occurrence data itself in the remaining lower right block (Fig. 12). Flexibility is
387 provided for in two ways. The individual fields in each block can be populated by different actual
388 data types, depending on the overall record type as determined by the 'File Type' field. Currently
389 there are only two defined file types, for deep-sea drilling data and more traditional land section
390 data (O and L, respectively). These differ both in general metadata (Site location vs geographic
391 name and geographic coordinates), and in the way in which sample names are structured: deep-
392 sea drilling samples ('O' files) use a consistent Site-Hole-Core-Section-Interval format, while land
393 sections are more variably defined, but usually include some combination of geologic formation,
394 vertical position in section and sample name (usually unique to each study); with additional
395 information often recorded on geologic age or biostratigraphic zone and lithology. SOD 'L'
396 formatted files include all these fields. Within the broad constraints on total fields available, the
397 number of file types using this layout is open to indefinite expansion. The SOD layout itself is
398 also extensible, as the version is written in the first metadata field in each file. The field
399 definitions and thus the data expected in each field are determined by these control fields, and
400 different layouts can be defined, for example with additional rows for sample name fields. This
401 flexibility however requires a separate source of information that defines, for the user and
402 programmer, what the field contents must be for each 'File Type' or SOD version number. These
403 definition requirements are the fundamental difference between regular data files as found in the
404 literature, and the SOD format. The definitions are given in two ways (which also allows cross
405 checking for data consistency). First, the tabular file definition requires full labeling - each cell,
406 row or column that holds data has an adjacent cell with fixed text content defining the data cell(s)
407 adjacent, so that the content resembles a simple key:value non-relational database structure. This
408 means the files are largely self documenting, and provides sufficient explanatory information to
409 users so that they can create new data files from a template file (containing labels but no data
410 values). Second, programs that read SOD files are expected to have a definition table of some
411 sort which gives the location and meaning of each cell for each file type and each SOD version.
412 Currently this is implemented in a table in the NSB database and used by programs (both a
413 python script and an R procedure at present) that read and upload SOD data into the NSB system.
414 This definition list could also be included (e.g. as a second 'page' in a spreadsheet file) with the
415 data files themselves. A full list of current SOD field definitions and additional details on the
416 format are given in Appendix 3.

417 Over 500 files have been created in SOD format, both typed or edited by users as described
418 above, or generated by the Raritas program during counting of microfossils. Raritas generates
419 only data for one sample at a time, but otherwise the output is identical to that used for complete
420 sample by taxa matrices in other SOD files. SOD formatted files are not intended to replace
421 more complex, formally controlled, computer-to-computer data exchange formats, defined in xml
422 or other systems. SOD is best viewed as complementary, providing a user accessible format that
423 encourages the capture of the metadata needed to adequately document stratigraphic occurrence

424 data, which until now has often not been done. It should also be noted that the SOD format is
425 much more flexible and can accommodate many more types of data than the current versions of
426 Raritas programs themselves, which are 'hard wired' to work e.g. with Taxa and Sample Names.
427 Future versions of these programs ideally should be modified to read the fields needed for the
428 metadata window, and output data file formats, directly from a SOD definition file.

429 **Diversity vs number of specimens**

430 The program outputs, in addition to the main count data, the cumulative diversity vs number of
431 counted objects history as a simple tab-text data file. This data can be useful for fitting
432 rarefaction curves in subsequent data analyses.

433 **Results**

434 The degree to which biodiversity assessments can be improved using our software depends on a
435 variety of factors - the distribution of taxon abundances (evenness) and absolute diversity of the
436 target population(s) being counted; and the ability of the user to mentally mask out taxa and focus
437 only on those not excluded. Most people can easily keep a 'skip' list of several taxa in mind when
438 counting, but not a much larger list, e.g. a dozen or more taxa. Thus the improvement in counting
439 with Raritas tends to be best when the abundances are significantly uneven and the total diversity
440 is less than a few hundred categories. In the example shown in Figures 1 and 7 of this paper,
441 from Antarctic Pleistocene radiolarian assemblages, by eliminating the 6 most common species
442 (cumulative abundance of >74% of the specimens in the sample) nearly 3/4 of the specimens can
443 be skipped, allowing an effective sampling of the rarer taxa that is 4X what would have been
444 possible by counting all specimens. In practice we have found that we more typically increase
445 our effective sample size by 2-3X by using rare count mode. These increased effective sample
446 sizes significantly improve the accuracy of diversity estimates, although the precise amount will
447 depend on total sample size, evenness and absolute diversity (Colwell et. al., 2012).

448 **Discussion and Conclusions**

449 The programs described here provide useful tools for counting populations with large numbers of
450 categories and unequal abundances of individuals in categories. They are, as programmed, best
451 suited to micropaleontology studies, but with only minor modification can be adapted to many
452 other uses in biodiversity research and other fields. The SOD definition provides a flexible,
453 internally documented yet easy to read file format for storing and exchanging occurrence data,
454 either for individual populations or matrices with multiple sets of observations. The Raritas
455 program described here has proved itself in actual use over several years in the junior author's
456 research group in Berlin. As noted above, it has been used to count >700,000 specimens
457 belonging to several hundred different species in >100 radiolarian microfossil assemblages, as
458 part of a study of biodiversity change in the Southern Ocean over the last 20 my (Renaudie &
459 Lazarus, 2013). It has been used by several individuals in other projects including students, on a
460 variety of computers.

461 **Acknowledgements**

462 The authors wish to thank the numerous individuals for the open-source software tools used to
463 create the Raritas programs.

464 **References**

- 465 Stevenson RJ, Pan Y, van Dam H. 2010. Assessing environmental conditions in rivers and
466 streams with diatoms, p. 57–85. In Smol JP, Stoermer EF (ed), *The Diatoms: Applications for*
467 *the Environmental and Earth Sciences*, Cambridge University Press, Cambridge.
- 468 Alroy J, Marshall CR, Bambach RK, Bezusko K, Foote M, Fürsich FT, Hansen TA, Holland SM,
469 Ivany LC, Jablonski D, Jacobs DK, Jones DC, Kosnik MA, Lidgard S, Low S, Miller AI,
470 Novack-Gottshall PM, Olszewski TD, Patzkowsky ME, Raup DM, Roy K, Sepkoski JJ, Jr.,
471 Sommers MG, Wagner PJ, Webber A. 2001. Effects of sampling standardization on estimates
472 of Phanerozoic marine diversification. *Proceedings of the National Academy of Sciences*
473 (USA) 98:6261–6266.
- 474 Alroy J. 2010. Fair sampling of taxonomic richness and unbiased estimation of origination and
475 extinction rates, p. 55–80. In Alroy, J, Hunt G (ed), *Quantitative Methods in Paleobiology*,
476 The Paleontological Society,
- 477 Berendsohn W. 2007. Access to Biological Collection Data. ABCD Schema 2.06 - ratified
478 TDWG standard. TDWG Task Group on Access to Biological Collection Data, BGBM,
479 Berlin <http://www.bgbm.org/TDWG/CODATA/Schema/default.htm>.
- 480 Brown JH, Gupta VK, Li BL, Milne BT, Restropo C, West GB. 2002. The fractal nature of
481 nature: power laws, ecological complexity and biodiversity. *Phil Trans R Soc* 357:619–626.
- 482 Bugware. 2016. Bugwin. <http://www.bugware.com>
- 483 Buonassissi CJ, Dierssen HM. 2010. A regional comparison of particle size distributions and the
484 power law approximation in oceanic and estuarine surface waters. *Journal of Geophysical*
485 *Research* 115:C10028 (1–12).
- 486 CLIMAP members. 1976. The surface of the ice-age earth. *Science* 191:1131–1137.
- 487 Colwell RK, Chao A, Gotelli NJ, Lin S-Y, Mao CX, Chazdon RL, Longino JT. 2012. Models and
488 estimators linking individual-based and sample-based rarefaction, extrapolation and
489 comparison of assemblages. *J Plant Ecol* 5:3–21.
- 490 Dunn R. 2014. wxPython, version. 3.0. wxpython.org.
- 491 Gannon JE. 1971. Two counting cells for the enumeration of zooplankton micro-crustacea. *Trans*
492 *Am Micros Soc* 90:486–490.
- 493 Hinds WC. 1999. *Aerosol Technology: Properties, Behavior, and Measurement of Airborne*
494 *Particles, 2nd Edition*, Wiley, Hoboken, NJ.
- 495 Hunter JD. 2007. matplotlib: A 2D graphics environment. *Computing in Science and Engineering*
496 9:90–95.
- 497 Imbrie J, Kipp NG. 1971. A new micropaleontological method for quantitative paleoclimatology:
498 application to a late Pleistocene Caribbean core, p. 71–181. In Turekian KK (ed), *Late*
499 *Cenozoic Glacial Ages*, Yale University Press, New Haven.
- 500 Kim CC, DeRisi JL. 2010. VersaCount: customizable manual tally software for cell counting.
501 *Source Code Biol Med* 5:web.
- 502 Lazarus DB. 1994. The Neptune Project - a marine micropaleontology database. *Math Geol*
503 26:817–832.
- 504 Mathis JS, Rimpl W, Nordsieck KH. 1977. The size distribution of interstellar grains. *The*
505 *Astrophysical Journal* 217:425–433.
- 506 McDonald D, Clemente JC, Kuczynski J, Rideout JR, Stombaugh J, Wendel D, Wilke A, Huse S,
507 Hufnagle J, Meyer F, Knight R, Caporaso JG. 2012. The Biological Observation Matrix
508 (BIOM) format or: how I learned to stop worrying and love the ome-ome. *GigaScience* 1:1–
509 6.
- 510 Mcgann M, McGann LB, Bonomassa O, Devries P, Luther J, Malmberg S, Nelson G, Pratt SIII.

- 511 2006. Foramsampler v. 3.0 - microfossil sample data management software. *Anuário do*
512 *Instituto de Geociências* 29:278–279.
- 513 Mora C, Tittensor DP, Adl SM, Simpson AGB, Worm B. 2011. How many species are there on
514 earth and in the ocean? *PLoS Biology* 9:1–8 (web).
- 515 Nalepka D, Walanus A. 2003. Data processing in pollen analysis. *Acta Paleobot* 43:125–134.
- 516 Preston FW. 1948. The commonness, and rarity, of species. *Ecology* 29:254–283.
- 517 Reed WJ, Hughes BD. 2002. From gene families and genera to incomes and internet file sizes:
518 why power-laws are so common in nature. *Physical Review E* 66:67103–67106.
- 519 Renaudie J, Lazarus D. 2013. On the accuracy of paleodiversity reconstructions: a case study in
520 Antarctic Neogene radiolarians. *Paleobiology* 39:491–509.
- 521 Roberts TE, Bridge TC, Caley MJ, Baird AH. 2016. The Point Count Transect Method for
522 Estimates of Biodiversity on Coral Reefs: Improving the Sampling of Rare Species. *PLoS*
523 *One* 11:e0152335.
- 524 Spencer-Cervato C. 1999. The Cenozoic deep sea microfossil record: explorations of the
525 DSDP/ODP sample set using the Neptune database. *Palaeontologica Electronica* 2:web.
- 526 Stratadata. 2014. Stratabugs biostratigraphic data management software.
527 <http://www.stratadata.co.uk>
- 528 van Rossum G, Drake J. 2010. Python Language Reference, version 2.7.
- 529 Walker W, Lamere P, Kwok P, Raj B, Singh R, Gouvea E. 2004. Sphinx-4: a Flexible Open
530 Source Framework for Speech Recognition, Sun Microsystems, Mountain View, CA.
- 531 Wieczorek J, Bloom D, Guralnick R, Blum S, Döring M, Giovanni R, Robertson T, Vieglaiss D.
532 2012. Darwin Core: An evolving community-developed biodiversity data standard. *PlosOne*
533 7:1–8.
- 534 Zachariasse WJ, Riedel WR, Sanfilippo A, Schmidt RR, Brolsma MJ, Schrader HJ, Gersonde R,
535 Drooger MM, Broekman JA. 1978. Micropaleontological counting methods and techniques-
536 an exercise on an eight metres section of the lower Pliocene of Capo Rossello, Sicily. *Utrecht*
537 *Micropaleontological Bulletins* 17:79-176.
- 538 Zippi P. 2007. Counter 4.5. PAZ Software. <http://www.pazsoftware.com>.

539 **Supporting Information Appendices**

- 540 S1 Appendix 1 - User Guides
541 S2 Appendix 2 - Sample Files
542 S3 Appendix 3 - SOD Definition

Figure 1

Assemblages with common and rare taxa

Microfossil assemblage as seen in the microscope (late Pleistocene, Southern Ocean, ODP Site 751). Specimens marked by black arrows all belong to *Antarctissa strelkovi* or *A. denticulata*. Other radiolarian species are marked by white arrows. Unmarked individuals are not targets for counting - broken radiolarians and diatom valves. Most individuals in this target assemblage belong to just a few species (particularly *A. strelkovi* and *A. denticulata*), making discovery of rarer taxa difficult.

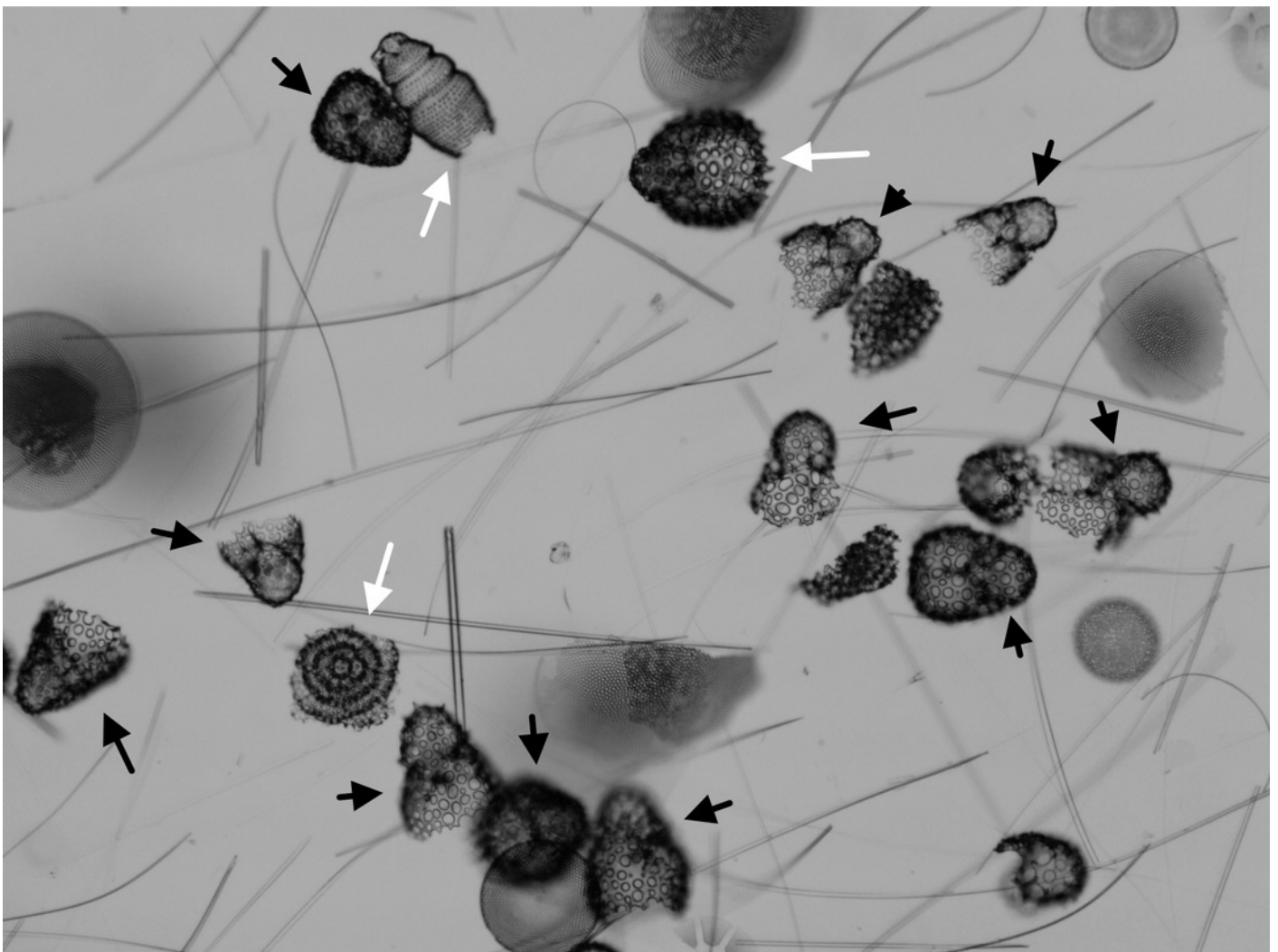


Figure 2(on next page)

Ranked relative abundances of fossil radiolarian species in single samples and combined multisample datasets.

Counts of species, sorted by abundance, of Neogene Southern Ocean radiolarian assemblages, showing total dataset (several dozen samples) and a single sample (Deep-sea drilling sample ODP 751A-6H-6, 98-100 cm). Despite a total count of 7071 specimens within the single sample, the majority of the species are represented by 6 or fewer individuals. From data in (Renaudie & Lazarus, 2013) SOM.

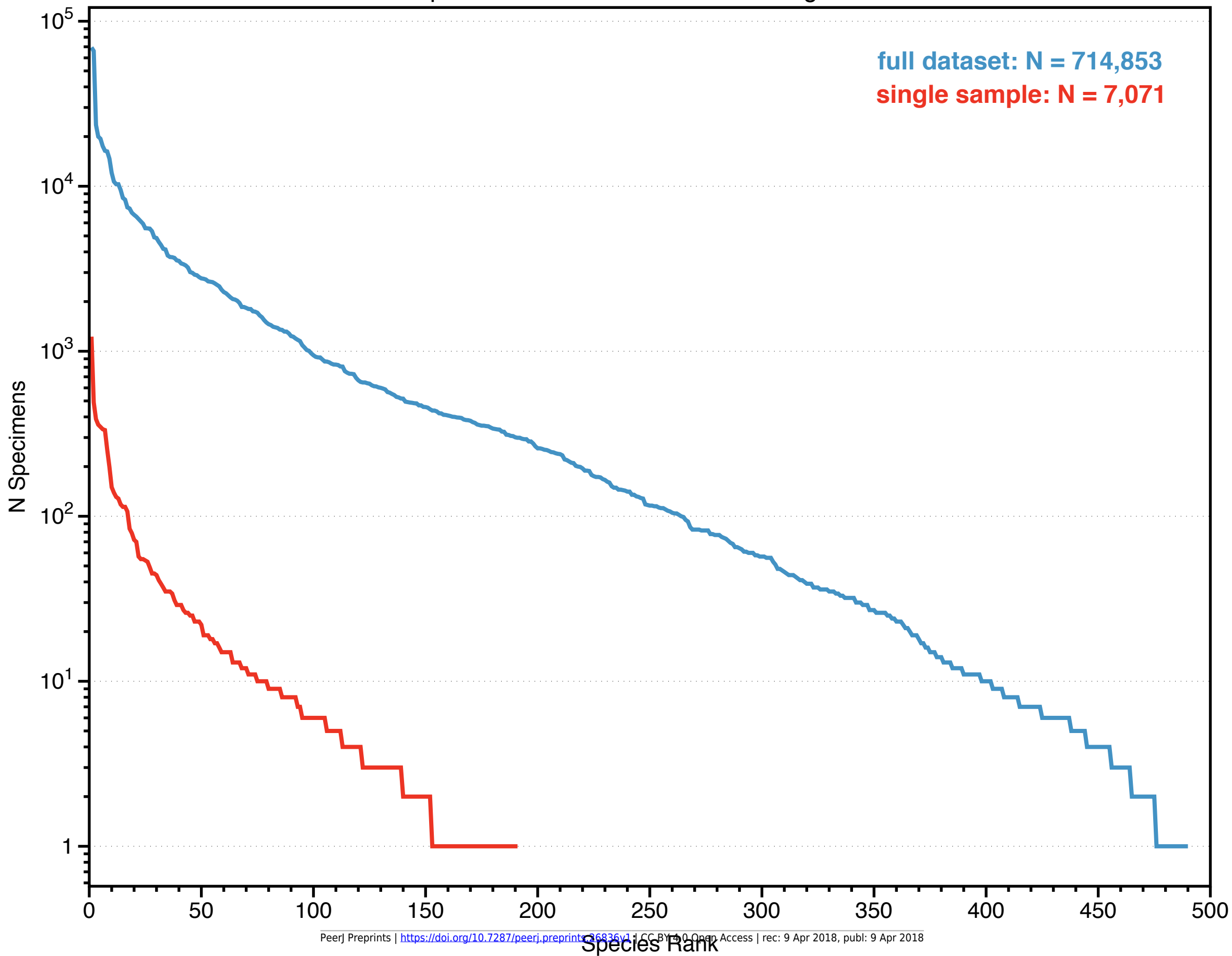


Figure 3(on next page)

Cumulative diversity vs sample size curve and estimated true diversity for a single sample.

Species-accumulation curve on a typical sample (sample ODP 751A-6H-6, 98-100 cm shown in Fig 1). Bold black curve is the species accumulation curve; light grey curve is a de Caprariis type curve-fit; dashed light grey line its asymptote (i.e. species diversity at infinite sample size). From (Renaudie & Lazarus, 2013).

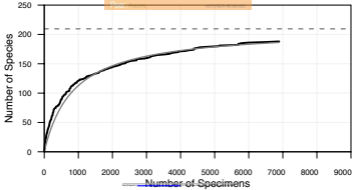


Figure 4

Configuration file to populate interface with category names.

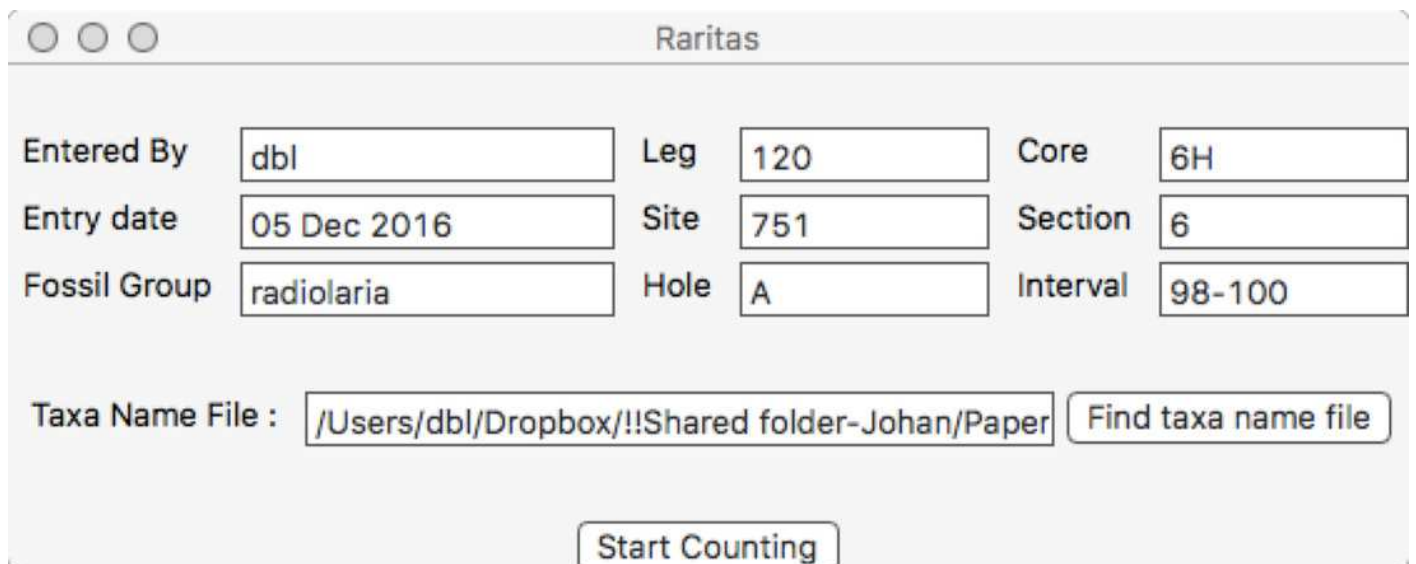
Configuration file format (a plain text file, here formatted for easier reading). Only a few fields - 'Genus' and 'Species' components of a taxonomic name, button (yes/no) are mandatory. A couple fields, e.g. 'Recognition Name' are used only by RaritasVox.

Genus	GQ	Species	SQ	Sub-species	Author	HigherTaxon	Comment	Color	Button	abbreviation	Recognition Name	listNr.
Acanthodesmia		micropora				Nassellaria			n	Aca mic		2
Acanthosphaera		actinota				Spumellaria			n	Aca act		3
Acanthosphaera		insignis				Spumellaria			n	Aca ins		3
Acrosphaera		australis				Collo/Entact/Phaeo			y	australis		1
Acrosphaera		cyrtodon				Collo/Entact/Phaeo			n	Acr cyr		1
Acrosphaera		labrata				Collo/Entact/Phaeo			y	labrata		1
Acrosphaera		lappacea				Collo/Entact/Phaeo			n	Acr lap		1
Acrosphaera		mercurius				Collo/Entact/Phaeo			n	Acr mer		1
Acrosphaera		murrayana				Collo/Entact/Phaeo			y	murrayana		1
Acrosphaera		cuniculauris				Collo/Entact/Phaeo			y	cuniculi		1
Acrosphaera		spinosa				Collo/Entact/Phaeo			n	Acr spi		1
Actinomma		arcadophorum				Spumellaria			n	Act arc		3
Actinomma		boreale				Spumellaria			n	Act bor		3
Actinomma		campilacantha				Spumellaria			n	Act cam		3
Actinomma		delicatulum				Spumellaria			y	Act del		3
Actinomma		golownini				Spumellaria			y	golownini		3
Actinomma		eldredgei				Spumellaria			n	Act spE		3
Actinomma		kerquelensis				Spumellaria			n	Act ker		3

Figure 5

Dialog to enter general sample metadata.

Metadata window used for Raritas. Information about the sample to be counted is entered here, including observer, date, class of objects being counted ('Fossil Group'), and sample identification information. RaritasVox has additional options (not shown), e.g. 'Save list of counted species with diversity' which, if checked, creates a second output file that gives the entire history of counting.



The image shows a screenshot of a software window titled "Raritas". The window contains several input fields for metadata, arranged in a grid. The fields are: "Entered By" (value: dbl), "Leg" (value: 120), "Core" (value: 6H), "Entry date" (value: 05 Dec 2016), "Site" (value: 751), "Section" (value: 6), "Fossil Group" (value: radiolaria), "Hole" (value: A), and "Interval" (value: 98-100). Below these fields is a "Taxa Name File" field with the path "/Users/dbl/Dropbox/!!Shared folder-Johan/Paper" and a "Find taxa name file" button. At the bottom center of the window is a "Start Counting" button.

Entered By	dbl	Leg	120	Core	6H
Entry date	05 Dec 2016	Site	751	Section	6
Fossil Group	radiolaria	Hole	A	Interval	98-100

Taxa Name File :

Figure 6

Main counting window with buttons, hierarchical category menus and count status information.

Main counting window. Objects to be counted are presented in two forms: an array of clickable buttons in the upper part of the window, and as a set of pop-up lists in the lower left and center part of the window. The number of lists and their contents is automatically built from the configuration file higher category labels for object entries. Button labels are also taken from this file on start-up. Other buttons or menu items control program behavior and call up other features e.g. voice recognition (RaritasVox only), show count plot, switch to Rare Count mode etc. A scrolling list of the most recently counted objects is shown in the lower middle. The 'Track' counter and clickable (large rectangular) button are on the lower right and are used to record observation effort in both regular and rare count modes. Note, in this image rare count mode has already been activated; thus some buttons are greyed out.

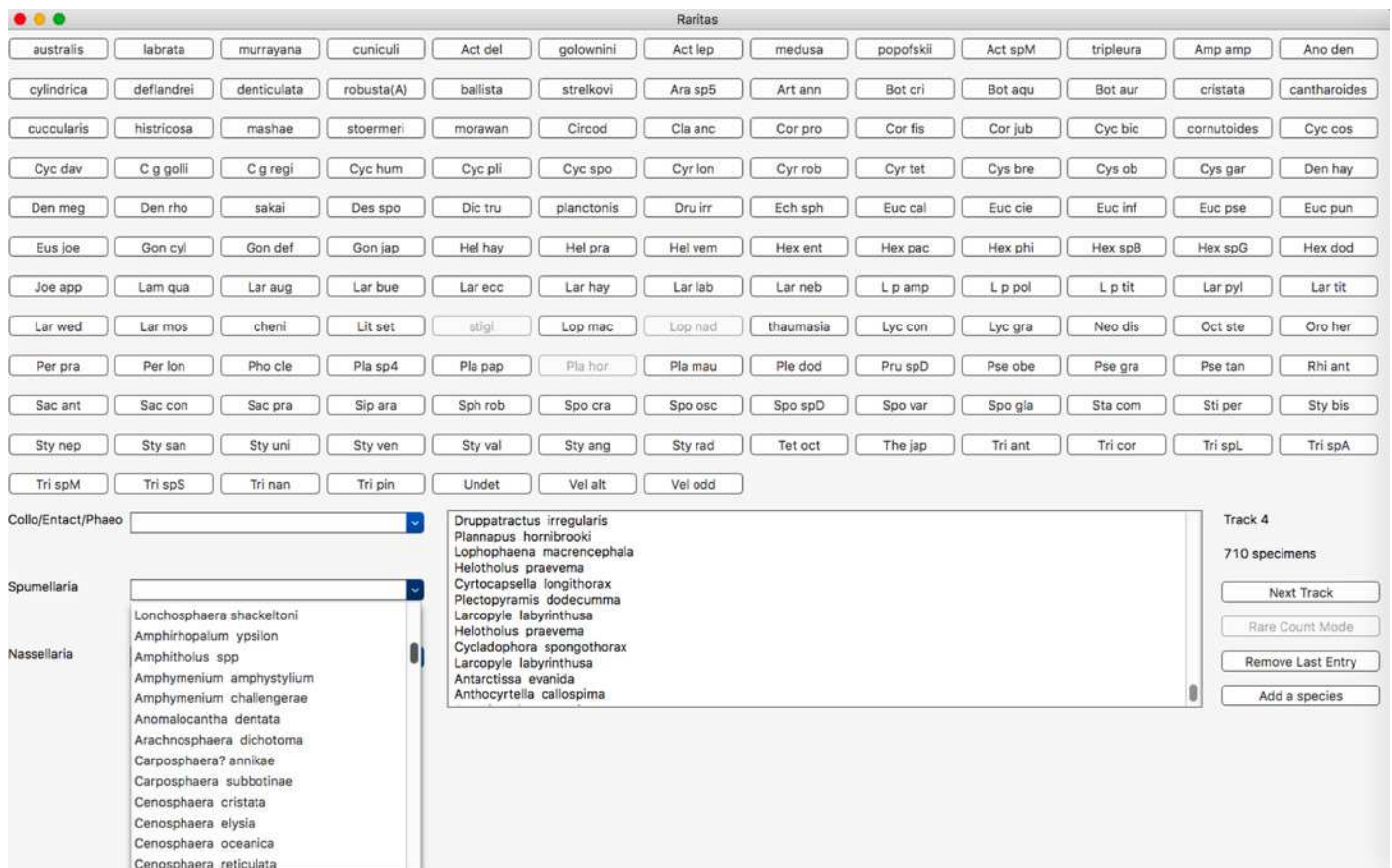


Figure 7

Dialog to configure rare count mode.

Configure rare count mode dialog. The object counts list, sorted by count frequencies, is presented and the user selects those objects (here, species names) that will in skipped and no longer counted in rare count mode.

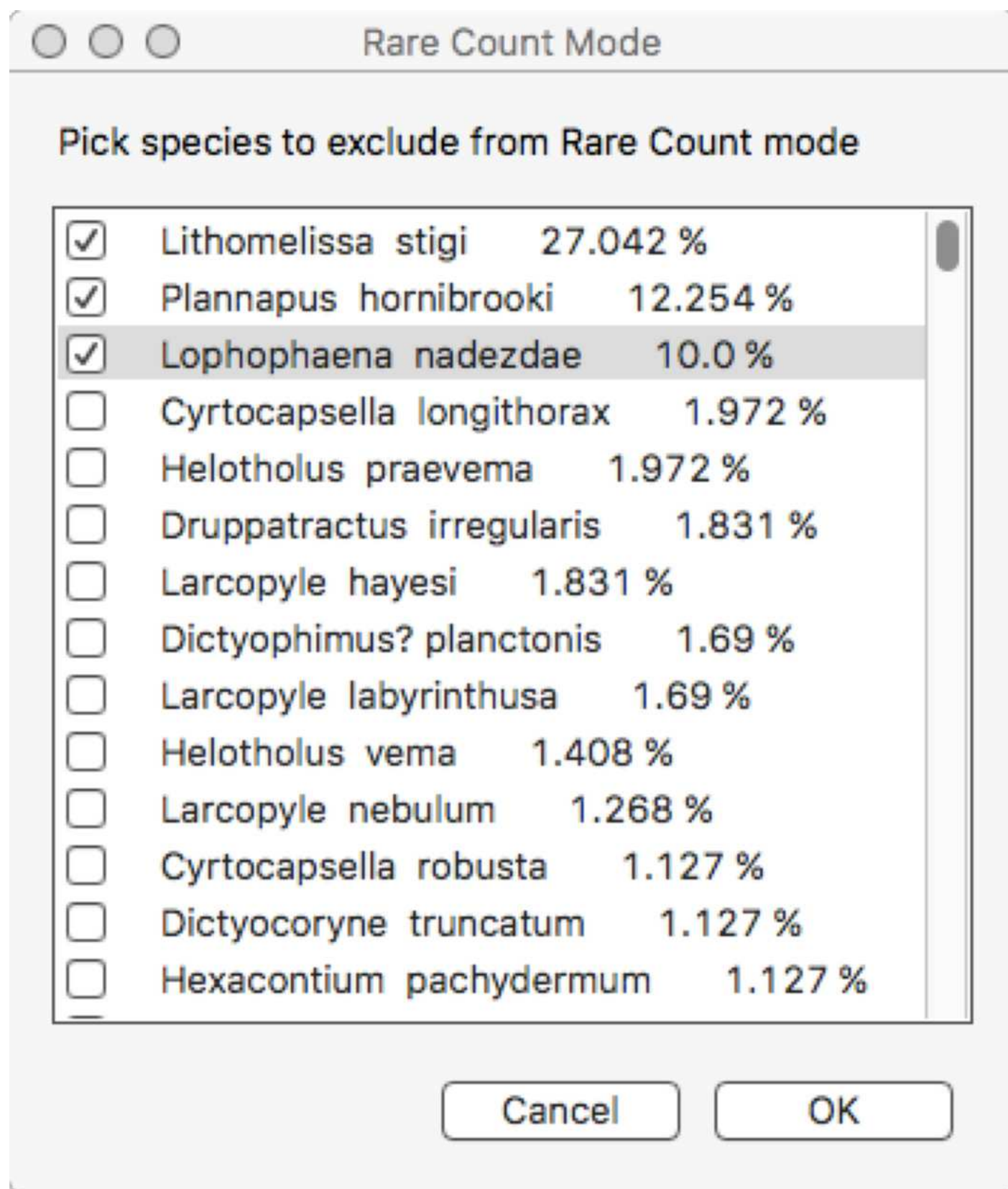


Figure 8(on next page)

Relationships between sample size and uncertainty of abundance estimates in generalized and actual biodiversity data.

Panel A (left) - Epsilon (size of confidence interval, relative to the abundance value, for a given species relative abundance in a population) plotted on a p (percent) vs N (number of specimen) landscape. Rule of thumb used in [12] marked by dashed lines (Renaudie & Lazarus, 2013) highlighted. Panel B (right) - Shows, for data reported in (Renaudie & Lazarus, 2013), red line: the percent of samples that have at least one species with percent higher than p ; blue line: the percent of species having a proportion higher than p in at least one sample, and black line with shading: the cumulative proportion of specimens of species with proportion higher than p (mean, inner-quartile range and total range over all 107 samples).

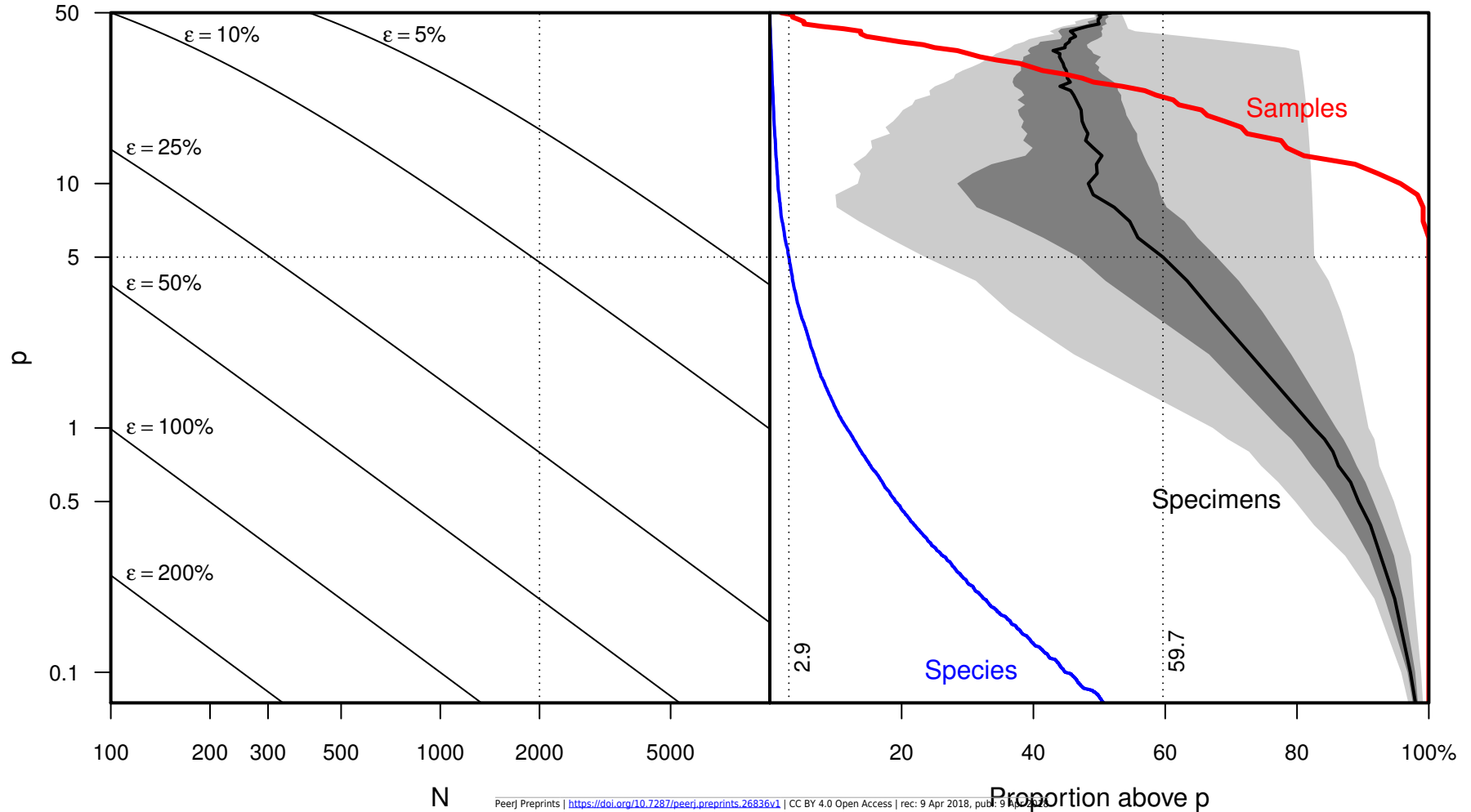


Figure 9(on next page)

Collecting curve, showing history of cumulative diversity vs sample size.

Count plot window, showing a simple graphic of how total diversity of objects ('species') is increasing with increased numbers of counted objects ('specimens'). The window appears whenever the user clicks the 'show count plot' button in the main counting window. This graphic is calculated and plotted anew with each invocation. The shape of the curve provides important feedback for the user, see text for details.

Collector's curve

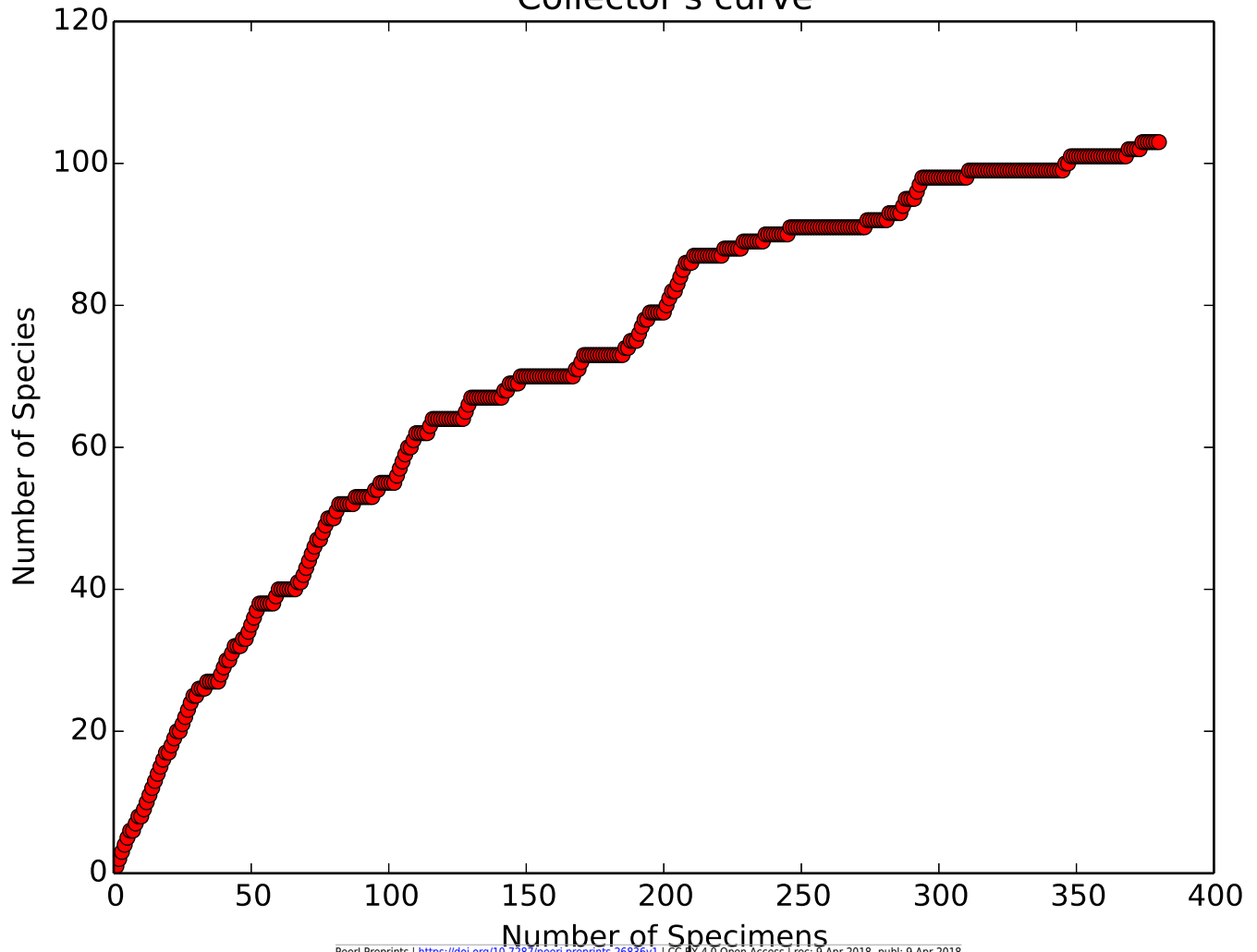


Figure 10(on next page)

RaritasVox defined vocabulary and pronunciation configuration file.

Configuration file for voice recognition using RaritasVox (extract only). Spoken words are on the left and the phoneme pronunciations on the right.

APROX	AH P R OW K S
STOCKI	S T OW K IY
AMPRADIOSA	AE M P R AE D IY OW S AH
ARTANNULATUS	AA R T AE N AH L EY T AH Z
AXIRREGULAR	AE K S IH R EH G Y AH L ER
BGRAN	B IY G R AE N
CRYPTBUSS	K R IH P T B AH S
GONDWANA	G AO N D W AE N AH
LOPHOHADRA	L OW F AH HH AE D R AH
MITA	M IY T AH
PODPAPILIS	P AA D PAE P IH L IH S
PSEUDODICT	S UW D OW D IH K T
ZYGO	S IY G ER
SPYRO	S P IY R ER
CORNUTELLA	K AO R N Y UH T EH L AH
CALOCYCLAS	K AE L OW S AY K L AH Z
BUNNYEARS	B AH N IY IH R Z
ZIGZAG	Z IH G Z AE G

Figure 11

Main counting window for RaritasVox

Screenshot of RaritasVox in voice-counting mode. A list of acceptable words is shown in the top window, the currently recognized word in large letters in the middle of the screen (to make it easy to see at a glance when e.g. working at a microscope), button controls below this and summary panes of count activity at the bottom.



Figure 12(on next page)

Example of SOD file format with data blocks framed.

Example of SOD file output (the main data output file produced by Raritas), with the 4 main areas (blocks) marked by bold lines. Metadata about the data file is stored in the upper left block, object labels and linked data such as author names, if known, are in the lower left block, sample information is in the upper right block, and the actual counting data in the lower right block. In output from the Raritas program only a single column of data is created but the SOD format definition permits the sample name and count values to repeat indefinitely (to the right of this figure). Note that only a few selected rows are shown here - the full file has ca 400 taxa names.

Table 1 (on next page)

Recognition accuracy in a simple test run of RaritasVox.

Accuracy of spoken entry using RaritasVox for a short list of species name abbreviations. Each name was spoken in random order 6 times. Note the independence of the spoken and data names e.g. zigzag for *L. robusta*. The spoken and formal names are linked in the Vox configuration file.

Genus	GQ	Species	SQ	spoken name	VOX count	Errors
Amphicraspedum		prolixum	gr.	aprox	5	1
Amphipyndax		stocki		stocki	6	0
Amphisphaera		radiosa		ampradiosa	6	0
Artostrobos		annulatus		artannulatus	6	0
Axoprimum		irregularis		axirregular	5	1
Buryella		granulata		bgran	6	0
Calocyclas		spp.		calocyclas	7	1
Cornutella		sp.		cornutella	6	0
Cryptocarpium		bussonii	gr.	cryptbuss	8	2
Gondwanaria	?	sp.		gondwana	6	0
Lithomelissa		robusta		zigzag	6	0
Lophocyrtis		hadra		lophohadra	5	1
Acrosphaera		cuniculiauris		bunyears	6	0
Mita	?	sp.		mita	6	0
Podocyrtis		papilis		podpapilis	5	1
Pseudodictyophimus		gracilipes		pseudodict	6	0
Spyrocyrtis		A	n.sp.	spyro	6	0
Zygocircus		buetschli		zygo	7	1
					108	8