

1 *Technical Note*

2 **phylotaR: An automated pipeline for retrieving** 3 **orthologous DNA sequences from GenBank in R**

4 **Dominic J. Bennett** ^{1,2*}, **Hannes Hettling** ³, **Daniele Silvestro** ^{1,2}, **Alexander Zizka** ^{1,2}, **Christine D.**
5 **Bacon** ^{1,2}, **Søren Faurby** ^{1,2}, **Rutger A. Vos** ³ and **Alexandre Antonelli** ^{1,2,4,5}

6 ¹ Gothenburg Global Biodiversity Centre, Box 461, SE-405 30 Gothenburg, Sweden

7 ² Department of Biological and Environmental Sciences, University of Gothenburg, Box 461, SE-405 30
8 Gothenburg, Sweden

9 ³ Naturalis Biodiversity Center, P.O. Box 9517, 2300 RA Leiden, the Netherlands

10 ⁴ Gothenburg Botanical Garden, Carl Skottsbergs gata 22A, SE-413 19, Gothenburg, Sweden

11 ⁵ Department of Organismic and Evolutionary Biology, Harvard University, 26 Oxford St., Cambridge, MA
12 02138 USA

13 * Correspondence: dominic.john.bennett@gmail.com

14

15 **Abstract:** The exceptional increase in molecular DNA sequence data in open repositories is
16 mirrored by an ever-growing interest among evolutionary biologists to harvest and use those data
17 for phylogenetic inference. Many quality issues, however, are known and the sheer amount and
18 complexity of data available can pose considerable barriers to their usefulness. A key issue in this
19 domain is the high frequency of sequence mislabelling encountered when searching for suitable
20 sequences for phylogenetic analysis. These issues include the incorrect identification of sequenced
21 species, non-standardised and ambiguous sequence annotation, and the inadvertent addition of
22 paralogous sequences by users, among others. Taken together, these issues likely add considerable
23 noise, error or bias to phylogenetic inference, a risk that is likely to increase with the size of
24 phylogenies or the molecular datasets used to generate them. Here we present a software package,
25 phylotaR, that bypasses the above issues by using instead an alignment search tool to identify
26 orthologous sequences. Our package builds on the framework of its predecessor, PhyLoTa, by
27 providing a modular pipeline for identifying overlapping sequence clusters using up-to-date
28 GenBank data and providing new features, improvements and tools. We demonstrate our
29 pipeline's effectiveness by presenting trees generated from phylotaR clusters for two large
30 taxonomic clades: palms and primates. Given the versatility of this package, we hope that it will
31 become a standard tool for any research aiming to use GenBank data for phylogenetic analysis.

32 **Keywords:** BLAST; DNA, open source; phylogenetics; R; sequence orthology.

33

34 **1. Introduction**

35 The first step in any nucleotide-based phylogenetic analysis is the identification of sequence
36 homology. Without establishing homology, much like comparing apples and oranges, multiple
37 sequence alignment is meaningless. More precisely, in the context of species trees, sequences chosen
38 for phylogenetic analysis must represent the same ancestral region resulting from a speciation event,
39 i.e. they must be orthologous [1-2]. Sequence orthology is often determined through name-based
40 searches via large sequence databases, most commonly GenBank [3]. This approach, however, can
41 be problematic due to the possibility of sequences being mislabelled and differences in naming
42 conventions. For example, gene names can differ between working groups (e.g. COI, CO1, COX1
43 and COXI); different sections of a gene or region may be deposited under the same sequence name
44 [4]; and deposited sequences may represent multiple genes in what are termed chimeric sequences
45 [5]. In the best case scenario, these issues may lead to the failure to identify all relevant orthologous
46 sequences. Worst case, one or more of the downloaded sequences will represent different ancestral

47 regions, causing poor alignment and/or incorrect inference of phylogenetic trees. Without resolving
48 the problem of orthology in a programtic fashion, any large-scale attempt at self-updating,
49 automated pipelines and initiatives for constructing phylogenies, e.g. [6-7], are bound to fail.

50 In an early attempt to address these issues, Sanderson et al. [4] developed a pipeline, PhyLoTa, that
51 uses the Basic Local Alignment Search Tool (BLAST [8]) to identify orthologous sequences without
52 the need for gene name matching. For a given taxonomic group, PhyLoTa searches through
53 available sequences on GenBank and identifies orthologous sequence clusters. Users are then able to
54 survey the clusters via a web-interface [9] and select the ones that best suit their phylogenetic
55 analysis needs, e.g. by selecting the clusters that maximise coverage of their taxonomic groups of
56 interest. A downside of PhyLoTA is that the searching and clustering is performed via all-versus-all
57 BLASTing, the combinatorics of which become prohibitive above a certain taxonomic level – an ever
58 increasing barrier as public sequence databases grow. One solution is to perform the
59 BLASTing within taxonomic groups, leading to potentially shared clusters among taxonomic groups
60 remaining undiscovered by PhyLoTa.

61 More importantly, however, the current PhyLoTa release is outdated as it was built on a
62 GenBank release (representing 162,886,727 sequences, Release 194 [10]) in February 2013. Since then
63 over 44 million new sequences have been deposited in GenBank (Release 224, representing
64 207,040,555 sequences, [10]). Additionally, NCBI's taxonomic database has been updated as new
65 sequences are added, species are discovered and groupings are revised. Between March 2013 and
66 February 2018, 170,000 new nodes of the database's taxonomic tree were added [11], representing
67 30% of all current nodes. Clearly, more frequent and regular updates to the pipeline are needed for
68 phylogeneticists to make use of newly acquired and improved data.

69 To date, there have been just two alternatives for those who wish to discover orthologous
70 sequences from GenBank: rely on error-prone gene names, or make do with outdated information.
71 Here we present phylotaR, an R package comprising a pipeline for identifying and retrieving
72 orthologous sequence clusters directly from the latest GenBank release. Our pipeline recreates the
73 original PhyLoTa output for specific clades of interest to the user in a series of independent stages.
74 Additionally, a user has the option of a secondary cluster stage (*cluster*²) to identify and merge
75 clusters at higher taxonomic levels than is available with PhyLoTa. We demonstrate the capacity of
76 phylotaR by generating phylogenetic trees for two model clades, widely studied in evolutionary
77 biology: palms and primates.

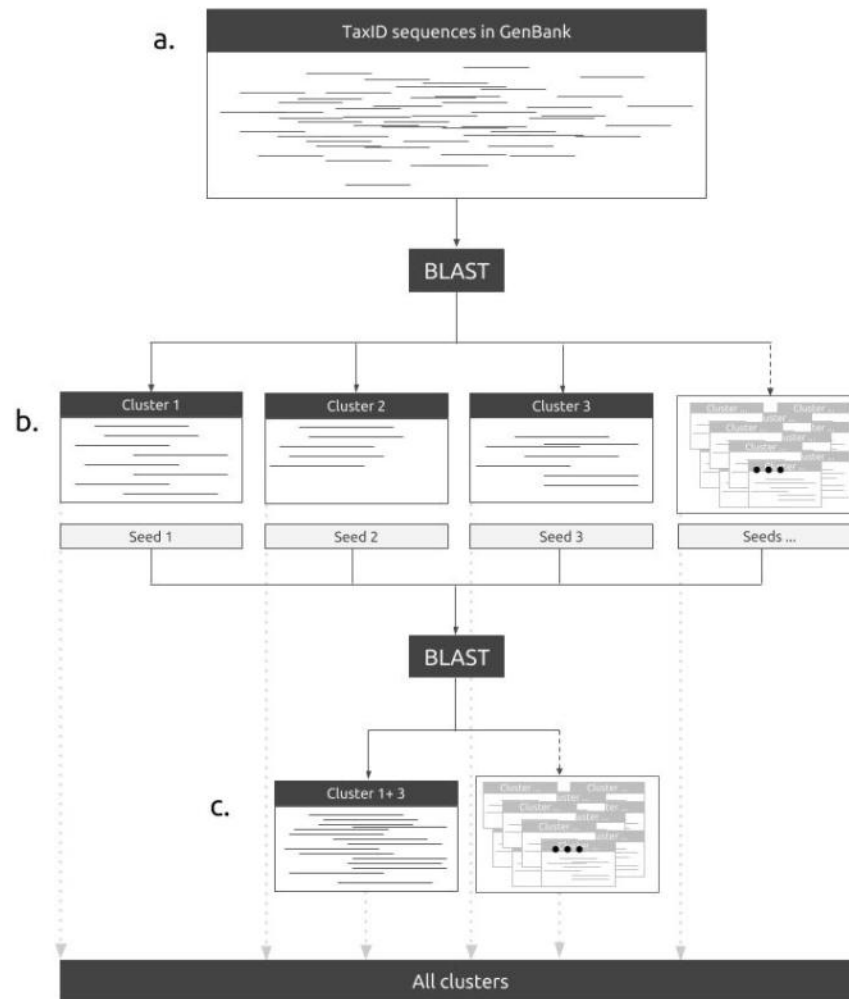
78 2. Implementation

79 2.1. The pipeline

80 The phylotaR pipeline consists of four automated, independent stages: **taxise** (identify all
81 descendant taxonomic nodes), **download** (hierarchically retrieve all sequences from across the
82 taxonomic tree), **cluster** (identify clusters from the downloaded sequences within nodes) and
83 **cluster**² (merge clusters identified within separate taxonomic nodes to identify clusters at higher
84 taxonomic levels) (see Figure S1 for a conceptual outline). At a minimum, all a user needs to do is
85 provide a taxonomic identity (name or NCBI ID at any taxonomic level), for which they would like
86 to generate sequence clusters, and then run the phylotaR pipeline. The pipeline mimics the original
87 PhyLoTa but with the following improvements: i) it makes use of sequence feature information to
88 break up large sequences which may have otherwise been discarded for being too long, ii) it can
89 generate paraphyletic clusters from nodes which are too small in themselves and iii) it has the
90 additional stage for matching sister clusters, *cluster*², which makes our method scalable to larger
91 groups of taxa with many sequences available. For more details on the pipeline see Figure 1 for an
92 outline of the process, refer to appendix A for a detailed description of each stage and Table S1 for a
93 description of all the parameters.

94 After the phylotaR pipeline stages have completed, the user can interrogate the identified
95 clusters using a series of functions supplied within the phylotaR package. For example, a user can

96 filter the sequences across the clusters to a given taxonomic rank, or select sequences with clusters
 97 using a range factors: sequence lengths, GC-ratios, sequence definitions, proportion of ambiguous
 98 nucleotides and/or maximum alignment density (MAD score, [4]). Additionally, plotting functions
 99 allow the user to see which taxa are covered by which clusters (for examples, see the Figures 3, S2
 100 and S3). After exploring, modifying and/or manipulating the clusters, the user can export them in
 101 tabular format as per the PhyLoTa database schema or as sequence files in FASTA format [12], which
 102 can be readily aligned by different software.
 103



104

105 **Figure 1.** The phylotaR pipeline identifies all sequences in GenBank associated with a user-specified
 106 taxonomic identity (a). The pipeline then performs all-vs-all BLAST across all the sequences to
 107 identify orthologous clusters (b). These searches are constrained to run within taxonomic groups up
 108 to a user-determined limit (default 50,000 sequences and 100,000 nodes). To generate higher
 109 taxonomic level clusters, an additional BLAST search is performed of the most connected sequences
 110 within clusters (i.e. the seed sequences) from the lower-level clusters. The clusters of overlapping
 111 seed sequences are then merged into larger clusters (c). All clusters, merged and unmerged, are then
 112 reported for inspection by the user. For more details on the pipeline, see appendix A.

113 2.2. Installation, features and future developments

114 The development version of the phylotaR package is currently available via GitHub
 115 (github.com/AntonelliLab/phylotaR) and can be installed through the R package devtools [13]. It
 116 will soon be available via CRAN [14]. The package depends on R (v 3+) and on a range of R packages
 117 [15-22], and requires stand-alone BLAST. Instructions for installing the BLAST+ suite can be found

118 via NCBI [23].

119 The entire pipeline can be run from an R console with just a few lines of code (Figure 2). The
120 package comes with tools that enable a user to kill the pipeline process at any point. All downloaded
121 and BLAST results are automatically cached, allowing a user to restart the pipeline after halting
122 without loss of data. A log file records any user interventions and all pipeline progress, thus,
123 increasing reproducibility of results and facilitating the methodological description in scientific
124 publications. Finally, the code is developed to be modular, allowing users to contribute additional
125 modules, functions and/or improvements. All internal functions and classes are documented to this
126 end. To maximise the transparency and stability of the phylotaR package we have submitted the
127 package to ROpenSci [24] – a community for ensuring good R coding practices in reproducible
128 research – for which it is currently under review.

```
library(phylotaR)
wd <- 'PATH TO WORKING DIRECTORY'
ncbi_dr <- 'PATH TO NCBI BLAST TOOLS'
# NCBI txid, e.g. primates
txid <- 9443
# Set-up working directory
setUp(wd=wd, txid=txid, ncbi_dr=ncbi_dr)
# Run pipeline
run(wd=wd)
```

129

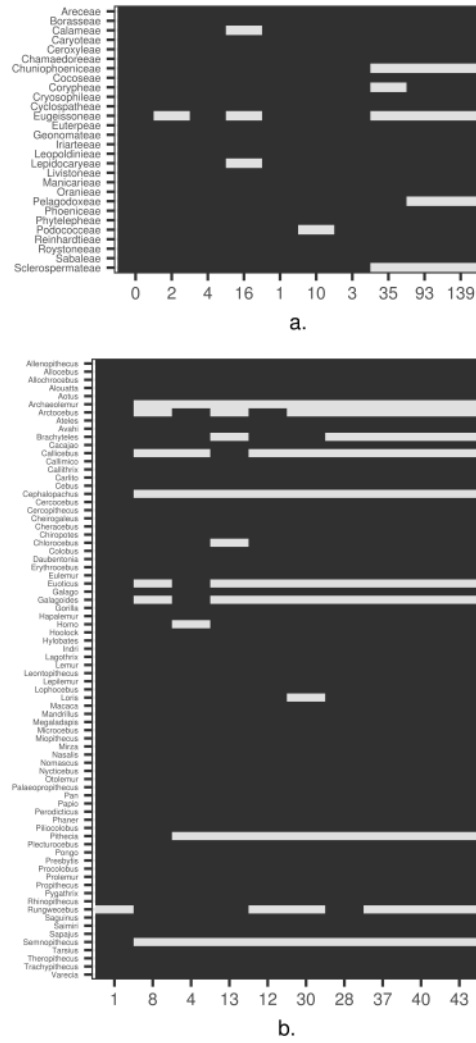
130 **Figure 2.** Initiating the phylotaR pipeline in R for primates (TaxID: 9443)

131 For future versions of phylotaR we envisage a range of additional features. For example,
132 the ability to identify clusters across disparate taxonomic groups using the cluster² stage; allow
133 users to incorporate their own unpublished sequences or specify their own taxonomy that is
134 independent of GenBank's; allow a user to download FTP copies of large sections of GenBank
135 to avoid slow-querying via Entrez [25]; and enable BLAST queries via NCBI's online BLAST
136 server [26] and other servers. Users can request new features and highlight bugs via the
137 phylotaR GitHub page (github.com/AntonelliLab/phylotaR).

138 3. Empirical Demonstration: Palms and Primates

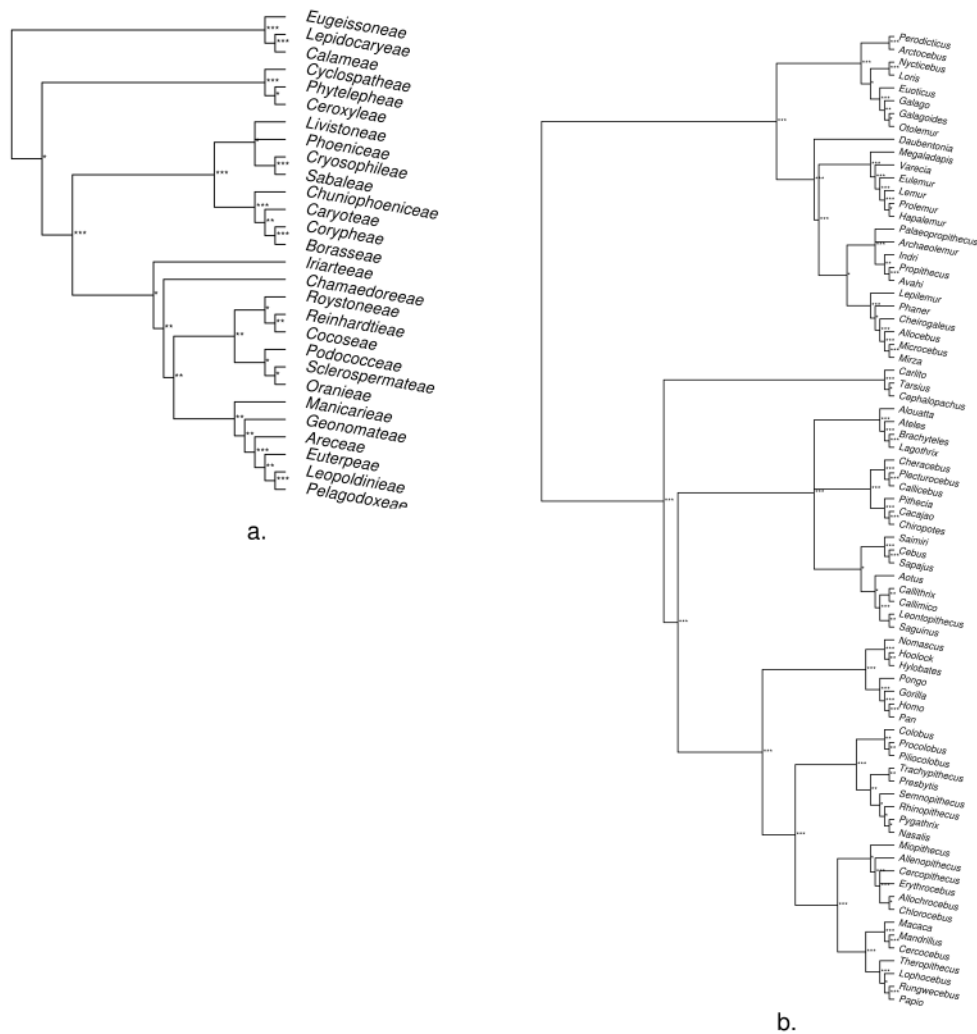
139 Here we demonstrate the phylotaR pipeline on palms (Arecaceae, TaxID: 4710) and primates
140 (Primates, TaxID: 9443) using default parameters, see Table S1. For palms and primates,
141 respectively, we identified 449 and 1,021 clusters, each containing over ten unique taxonomic nodes.
142 Taken together, these clusters included 1,238 and 653 unique taxonomic nodes and 13,344 and 56,112
143 sequences, respectively. See Figures S2 and S3 for visual representations of the relative distributions
144 of sequences/taxa/clusters among the different clusters and taxa.

145 To illustrate the ability of the pipeline to correctly identify orthologous sequence clusters and to
146 keep the demonstration fast, we generated small, representative trees for both palms and primates at
147 the tribe and genus levels, respectively. We reduced the number of sequences within each cluster to
148 just representatives within these taxonomic levels by filtering the sequences by proportion of
149 ambiguous nucleotides and length of sequence. We then selected the 'best' clusters for both palms
150 and primates from these reduced clusters, by dropping all clusters with sequence maximum
151 alignment densities (MAD, [4]) less than 0.75 and then selecting the top ten with the greatest number
152 of tribes/genera (see Tables S2a and S2b for detailed information on each of the selected clusters).
153 These top clusters were found to be representative of palms and primates as a whole (Figure 3). With
154 these clusters, we then generated alignments with MAFFT [27] and constructed phylogenetic trees
155 with RAxML [28]. We found the generated trees to be similar to those published for these groups
156 (Figures 4, S3 and S4) with few differences in topology. For further details on tree construction and
157 comparison methods, see appendix B. The methods described above and in the appendix can be
158 reproduced via scripts available from GitHub (github.com/AntonelliLab/phylotaR_demo).



159
160
161
162

Figure 3. Presence/absence of tribes and genera for palms (a) and primates (b), respectively, across the top ten best clusters. X-axis numbers are unique cluster Ids. For more details on each of these clusters, see Tables S2a and S2b.



163

164

165

166

167

168

169

Figure 4. Tribe- and genus-level trees for palms (a) and primates (b). Roots were determined manually by rooting with Strepsirrhini and Calamoideae for primates and palms respectively. Branch lengths have been removed. Support calculated from 100 rapid bootstraps: *** > .95, ** > .75 and * > .50. Complete tree construction methods are in appendix B. For tree comparisons with published trees for palms [29] and primates [30], see Figures S4 and S5.

170 5. Conclusions

171 The phylotaR package offers a user-friendly pipeline to obtain orthologous gene-sequences for
172 phylogenetic inferences from Genbank in R. Building on PhyLoTa, the pipeline makes use of
173 sequence feature information, is able to generate small paraphyletic clusters to prevent dropping
174 taxa and uses a subsequent cluster² to identify shared clusters across large taxonomic groups. The
175 phylotaR pipeline is modular, can be easily integrated into R workflow pipelines that, for example,
176 replicate the functionality of the SUPERSMART pipeline [6], and yields reliable results for
177 phylogenies of large taxonomic groups. The package is currently available via GitHub
178 (github.com/AntonelliLab/phylotaR) and comes with detailed vignettes containing documentation
179 and tutorials.

180 **Supplementary Materials:** The following are available online at www.mdpi.com/link, Figure S1: conceptual
181 outline of the phylotaR pipeline stages, Figure S2: relative number of sequences and clusters for each genus and
182 tribe, Figure S3: relative number of sequences and taxa for each cluster. Figures S4 and S5: comparison between
183 phylotaR-based trees to published trees, Table S1: default phylotaR pipeline parameters, Table S2(a)(b): details
184 on top clusters used for phylogeny constructions for palms and primates, primates.tre: resulting primate
185 Newick tree, palms.tre: resulting palms Newick tree.

186 **Acknowledgments:** We would like to thank Michael Sanderson for initiating the PhyLoTa project and
187 providing early feedback and advice. We would also like to thank ROpenSci for their feedback and support
188 regarding the development of the phylotaR package. In particular Scott Chamberlain for initial feedback and
189 Zebulun Arendsee and Naupaka Zimmerman for taking the time to review our code in thorough detail. D.S.
190 received funding from the Swedish Research Council (2015-04748). AA is supported by the Swedish Research
191 Council (B0569601), the Swedish Foundation for Strategic Research, a Wallenberg Academy Fellowship, the
192 Faculty of Sciences at the University of Gothenburg, the Wenner-Gren Foundations, and the David Rockefeller
193 Center for Latin American Studies at Harvard University.

194 **Author Contributions:** HH, RAV, DS, AZ and AA initiated, devised and developed the project. HH wrote the
195 initial pipeline code. DJB developed upon the pipeline, created the R package and wrote the manuscript. CDB
196 and SF reviewed the palm and primate trees. All authors contributed to the writing of the manuscript.

197 **Conflicts of Interest:** The authors declare no conflict of interest.

198 Appendix A

199 *Further pipeline details*

200 The first stage, **taxise**, looks up taxonomic information on all the descendant nodes from NCBI
201 taxonomy [31]. At the end of this stage a modifiable taxonomic dictionary is created containing
202 names, IDs, lineages and ranks of all descendant IDs. The taxonomic dictionary contains a
203 taxonomic tree (TreeMan class [20]), allowing fast querying of number of descendants and parent
204 taxa. In addition, the taxonomic nodes for which clusters can be generated are identified by counting
205 the number of possible sequences. Clades that contain too many sequences or too many
206 descendants, as defined by the user, are broken down into their subclades and these subclades are
207 analysed separately. A limit is required to prevent all-vs-all BLAST searches, at the clusters stage,
208 becoming too large.

209 In the second stage, **download**, sequences are hierarchically downloaded for each node
210 identified during the taxise stage. For ‘model organisms’ (sensu Sanderson et al. [4]: taxa for which
211 there are large numbers of sequence data available) a random subset of the available sequences are
212 downloaded. To prevent mega-clusters covering large numbers of different genes, all downloaded
213 sequences are broken down into their constituent annotated features. This stage queries GenBank [3]
214 through the rentrez package [16].

215 In the third stage, **cluster**, all-vs-all BLAST searches are performed within clades of
216 user-determined size using all the downloaded sequences. Clusters are identified for all nodes of the
217 taxonomy from sequences strictly associated with the node – direct – and all descendant sequences –
218 subtree. BLAST searches are performed externally to R through NCBI’s BLAST+ suite [23]. All
219 sequences that have BLAST e-values and coverages, respectively, less than and greater than the

220 user-determined maximum E-value and minimum coverage are considered part of a single cluster.
221 Entire clusters of sequences are then inferred from these BLAST results by identifying single-linkage
222 clusters, as per the original PhyLoTa pipeline [4], with the iGraph package [15]. For large clades that
223 have been broken down into many subclades for all-vs-all BLAST, an option in the phylotaR
224 parameters allows this stage to be run in parallel. Because the pipeline clusters hierarchically within
225 taxonomic clades, no clusters may be identified for small clades that are sister to very large clades
226 due to too few sequences. To prevent taxa from being excluded, paraphyletic clusters are generated
227 by non-hierarchically searching for clusters across all clades where no clusters were identified.

228 Because of the computational need to partition the cluster stage into subclades for very large
229 clades, a final, fourth phylotaR stage, **clusters**², is run to combine the subclade clusters into
230 higher-level clade clusters. For every cluster identified in the previous stage a 'seed sequence' is
231 determined as the sequence with the greatest number of BLAST hits with other sequences in its
232 cluster. This is slightly different from Sanderson et al.'s [4] conception of a 'seed sequence', which
233 was simply the starting point for single-linkage clustering. Because single-linkage clustering has a
234 tendency to wander – leading to stretched out clusters where the starting point may be far off from
235 the centre – we have opted to take the sequence with the highest connectivity instead. An all-vs-all
236 BLAST search is performed with these seed sequences. All subclade clusters where a valid BLAST
237 hit has occurred for their seed sequences are then merged into higher-level clusters.

238 At the end of the pipeline, a user will have identified clusters of the four different types
239 described above: direct, subtree, paraphyletic and merged.

240 **Appendix B**

241 *Further details on tree construction and assessment*

242 For each of the clusters, the sequences were written out in FASTA format and alignments were
243 constructed using mafft (v7.271, [27]) with its *--auto* argument. Supermatrices were then constructed
244 using the ten sets of alignments. RAxML (v8.2.4, [28]) was then used to construct the trees. We used
245 the GTRGAMMA model and partitioned the supermatrices based on the identified clusters. We ran
246 a rapid bootstrap analysis of 100 iterations.

247 To monitor for errors at any point in the demonstration pipeline two sequences were selected
248 instead of just one when reducing the clusters down to tribe/genus level. If in the resulting trees,
249 large numbers of representative sequences of the same taxon were not uniquely clustered, data and
250 methods were re-checked and the pipeline was run again. For palms and primates, sequences in the
251 same tribe/genus were sisters in, respectively, 100% and 94% of cases.

252 From the bootstrapped trees we calculated majority-rule consensus trees. These final trees
253 consisted of 28 and 80 tips for palms and primates, respectively representing 1.00 and 0.95 of all
254 known tribes and genera according to NCBI taxonomy (Federhen, 2012). Both trees were very
255 similar to published phylogenetic trees (Figures S3 and S4). We compared the final trees to
256 already-published ones (palms [29] and primates [30]) using the Robinson-Foulds distance [32] and
257 the triplet distance [33] through the R package treeman [20]. We found normalised Robinson-Foulds
258 and triplet distances between our trees and the published trees, respectively, of 0.083 and 0.002 for
259 palms and 0.189 and 0.016 for primates. For the primate tree, there were three key differences
260 between our tree and those published in the literature: the paraphyletic separation of the family
261 Lorisidae [(Perodicticus, Arctocebus), (Nyctiebus, Loris)], although this has been suggested before
262 [34]; a different branch ordering of the gibbons [35]; and the misplacement of the genus
263 Semenopithecus which has been suggested to be grouped with Trachypithecus [36]. For the palms
264 tree, the greatest problem was the misplacement of the subfamily Arecoideae (Iriarteae -
265 Pelagodoxeae) with the Coryphoideae (Livistoneae - Borasseae) rather than the Ceroxyloideae
266 (Cyclospatheae - Ceroxyleae), with which it is most often grouped [30, 37]. Additionally, there were
267 some potential errors and inconsistencies regarding the taxonomy: there is no representative of the
268 well-studied *Nypa* genus because there is no tribe name present in the NCBI taxonomy; also the
269 current accepted name for the Livistoneae is Trachycarpeae [38].

270 **References**

- 271 1. de Pinna, M.C.C. Concepts and tests of homology in the cladistics paradigm. *Cladistics*, **1991**, *7*(4),
272 367-394. DOI: 10.1111/j.1096-0031.1991.tb00045.x
- 273 2. Salemi, M.; Vandamme, A.-M.; Lemey, P. *The Phylogenetic Handbook: A Practical Approach to Phylogenetic*
274 *Analysis and Hypothesis Testing*. Cambridge University Press, Cambridge, UK, 2009.
- 275 3. Benson, D.A.; Karsch-Mizrachi, I.; Clark, K.; Lipman, D.J.; Ostell, J.; Sayers, E.W. GenBank. *Nucleic*
276 *Acids Research*, **2012**, *40*(Database issue), D48-53. DOI: 10.1093/nar/gkr1202
- 277 4. Sanderson, M.J.; Boss, D.; Chen, D.; Cranston, K.A.; Wehe, A. The PhyLoTA Browser: Processing
278 GenBank for molecular phylogenetics research. *Systematic Biology*, **2008**, *57*(3), 335-346.
279 DOI:10.1080/10635150802158688
- 280 5. Ashelford, K.E.; Chuzhanova, N.A.; Fry, J.C.; Jones, A.J.; Weightman, A.J. At least 1 in 20 16S rRNA
281 sequence records currently held in public repositories is estimated to contain substantial anomalies.
282 *Applied and Environmental Microbiology*, **2005**, *71*(12), 7724-7736. DOI:10.1128/AEM.71.12.7724-7736.2005
- 283 6. Antonelli, A.; Hettling, H.; Condamine, F.L.; Vos, K.; Nilsson, R.H.; Sanderson, M.J.; Sauquet, H.;
284 Scharn, R.; Silvestro, D.; Töpel, M.; Bacon, C.D.; Oxelman, B.; Vos, R.A. Toward a self-updating
285 platform for estimating rates of speciation and migration, ages, and relationships of Taxa. *Systematic*
286 *Biology*, **2017**, *66*(2), 153-166. DOI: 10.1093/sysbio/syw066
- 287 7. Pearse, W.D.; Purvis, A. phyloGenerator: An automated phylogeny generation tool for ecologists.
288 *Methods in Ecology and Evolution*, **2013**, *4*(7), 692-698. DOI: 10.1111/2041-210X.12055
- 289 8. Altschul, S.F.; Gish, W.; Miller, W.; Myers, E.W.; Lipman, D.J. Basic local alignment search tool.
290 *Journal of Molecular Biology*, **1990**, *215*(3), 403-410. DOI: 10.1016/S0022-2836(05)80360-2
- 291 9. PhyLoTa Browser. Available online: phylota.net (accessed on 28/03/18)
- 292 10. GenBank and WGS Statistics. Available online: www.ncbi.nlm.nih.gov/genbank/statistics (accessed
293 on 28/03/18)
- 294 11. Taxonomy Browser: Taxonomy Statistics. Available online:
295 www.ncbi.nlm.nih.gov/Taxonomy/taxonomyhome.html (accessed on 28/3/18)
- 296 12. Pearson, W.R.; Lipman, D.J. Improved tools for biological sequence comparison. *Proceedings of the*
297 *National Academy of Sciences*, **1988**, *85*(8), 2444-2448. DOI: 10.1073/pnas.85.8.2444
- 298 13. Wickham, H.; Hester, J.; Chang W.; Rstudio and R Core team. devtools: Tools to make developing R
299 packages easier. 2018. Available online: CRAN.R-project.org/package=devtools (accessed on 28/3/18)
- 300 14. The comprehensive R archive network. Available online: cran.r-project.org (accessed online 28/3/18)
- 301 15. Csardi G.; Nepusz T. The igraph software package for complex network research. *InterJournal,*
302 *Complex Systems*. **2018**, 1695.
- 303 16. Winter, D. rentrez: Entrez in R. R package version 1.1.0. 2017. Available online:
304 CRAN.R-project.org/package=rentrez (accessed on 28/3/18)
- 305 17. Lang, D.T.; the CRAN Team. XML: Tools for parsing and generating XML within R and S-Plus. 2018.
306 Available online: CRAN.R-project.org/package=XML (accessed on 28/3/18)
- 307 18. Wickham. H. *ggplot2: Elegant graphics for data analysis*. Springer-Verlag, New York, USA, 2009.

- 308 19. Ooms, J. sys: Portable system utilities. Available online: CRAN.R-project.org/package=sys (accessed
309 on 28/3/18)
- 310 20. Bennett, D.J.; Sutton, M.D.; Turvey, S.T. treeman: an R package for efficient and intuitive
311 manipulation of phylogenetic trees. *BMC Research Notes*, **2017**, *10*(1), 30.
- 312 21. Wilkins D. treemapify: Draw treemaps in 'ggplot2'. Available online:
313 CRAN.R-project.org/package=treemapify (accessed on 28/3/18)
- 314 22. Bengtsson, H. R.utils: Various Programming Utilities. Available online:
315 CRAN.R-project.org/package=R.utils (accessed on 28/3/18)
- 316 23. BLAST® Command Line Applications User Manual. Available online:
317 www.ncbi.nlm.nih.gov/books/NBK279690 (accessed on 28/3/18)
- 318 24. Transforming science through open data and software. Available online: ropensci.org (accessed on
319 28/3/18)
- 320 25. Entrez molecular sequence database system. Available online:
321 www.ncbi.nlm.nih.gov/Web/Search/entrezfs.html (accessed on 28/3/18)
- 322 26. Basic local alignment search tool. Available online: blast.ncbi.nlm.nih.gov/Blast.cgi (accessed on
323 28/3/18)
- 324 27. Katoh, K.; Standley, D.M. MAFFT multiple sequence alignment software version 7: Improvements in
325 performance and usability. *Molecular Biology and Evolution*, **2013**, *30*(4), 772-780.
326 DOI:10.1093/molbev/mst010
- 327 28. Stamatakis, A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large
328 phylogenies. *Bioinformatics*, **2014**, *30*(9), 1312-1313. DOI: 10.1093/bioinformatics/btu033
- 329 29. Baker, W.J.; Savolainen, V.; Asmussen-Lange, C.B.; Chase, M.W.; Dransfield, J.; Forest, F.; Harley,
330 M.M.; Uhl, N.W.; Wilkinson, M. Complete generic-level phylogenetic analyses of palms (Arecaceae)
331 with comparisons of supertree and supermatrix approaches. *Systematic Biology*, **2009**, *58*(2): 240-256.
332 DOI: 10.1093/sysbio/syp021
- 333 30. Perelman, P.; Johnson, W.E.; Roos, C.; Seuánez, H.N.; Horvath, J.E.; Moreira, M.A.M.; Kessing, B.;
334 Pontius, J.; Roelke, M.; Rumppler, Y.; Schneider, M.P.C.; Silva, A.; O'Brien, S.J.; Pecon-Slattery, J. A
335 molecular phylogeny of living primates. *PLoS Genetics*, **2011**, *7*(3), 1-17. DOI:
336 10.1371/journal.pgen.1001342
- 337 31. Federhen, S. The NCBI taxonomy database. *Nucleic Acids Research*, **2012**, *40*(Database issue), D136-43.
338 DOI: 10.1093/nar/gkr1178
- 339 32. Robinson, D.F.; Foulds, L.R. Comparison of phylogenetic trees. *Mathematical Biosciences*, **1981**, *53*(1-2),
340 131-147. DOI: 10.1016/0025-5564(81)90043-2
- 341 33. Critchlow D.E.; Pearl D.K.; Qian C. The triples distance for rooted bifurcating phylogenetic trees.
342 *Systematic Biology*, **1996**, *45*, 323-34.
- 343 34. Masters, J.C.; Anthony, N.M.; De Wit, M.J.; Mitchell, A. Reconstructing the evolutionary history of the
344 Lorisidae using morphological, molecular, and geological data. *American Journal of Physical*
345 *Anthropology*, **2005**, *127*(4), 465-480. DOI: 10.1002/ajpa.20149

- 346 35. Shi, C.M.; Yang, Z. Coalescent-based analyses of genomic sequence data provide a robust resolution
347 of phylogenetic relationships among major groups of gibbons. *Molecular Biology and Evolution*, **2018**,
348 35(1), 159-179. DOI: 10.1093/molbev/msx277
- 349 36. Osterholz, M.; Walter, L.; Roos, C. Phylogenetic position of the langur genera *Semnopithecus* and
350 *Trachypithecus* among Asian colobines, and genus affiliations of their species groups. *BMC*
351 *Evolutionary Biology*, **2008**, 8(1), 1-12. DOI: 10.1186/1471-2148-8-58
- 352 37. Couvreur, T.L.P.; Forest, F.; Baker, W.J. Origin and global diversification patterns of tropical rain
353 forests: Inferences from a complete genus-level phylogeny of palms. *BMC Biology*, **2011**, 9. DOI:
354 10.1186/1741-7007-9-44
- 355 38. Dransfield, J.; Uhl, N.W.; Asmussen-Lange, C.B.; Baker, W.J.; Harley, M.M.; Lewis, C.E. A new
356 phylogenetic classification of the palm family, *Arecaceae*. *Kew Bulletin*, **2005**, 60(September), 559-569.
357 DOI: 10.2307/25070242