

1 **EukProt: a database of genome-scale predicted proteins across the diversity of** 2 **eukaryotic life**

3
4 Daniel J. Richter¹, Cédric Berney^{2,3}, Jürgen F. H. Strassert⁴, Fabien Burki⁵, Colomán de
5 Vargas^{2,3}

6
7 1. Institut de Biologia Evolutiva (CSIC-Universitat Pompeu Fabra), Passeig Marítim de la Barceloneta 37-49,
8 08003 Barcelona, Catalonia, Spain

9 2. Sorbonne Université, CNRS, Station Biologique de Roscoff, UMR 7144, ECOMAP, 29680 Roscoff, France

10 3. Research Federation for the study of Global Ocean Systems Ecology and Evolution, FR2022/Tara GOSEE,
11 Paris, France

12 4. Institute of Biology, Free University of Berlin, Königin-Luise-Straße 1-3, 14195 Berlin, Germany

13 5. Science for Life Laboratory and Department of Organismal Biology, Uppsala University, Norbyvägen 18D,
14 Uppsala 75236, Sweden

15

16

17 **Database Availability**

18

19 <https://doi.org/10.6084/m9.figshare.12417881.v2>

20

21

22 **Abstract**

23

24 EukProt is a database of published and publicly available predicted protein sets and unannotated
25 genomes selected to represent eukaryotic diversity, including 742 species from all major supergroups
26 as well as orphan taxa. The goal of the database is to provide a single, convenient resource for
27 studies in phylogenomics, gene family evolution, and other gene-based research across the spectrum
28 of eukaryotic life. Each species is placed within the UniEuk taxonomic framework in order to facilitate
29 downstream analyses, and each data set is associated with a unique, persistent identifier to facilitate
30 comparison and replication among analyses. The database is currently in version 2, and all versions
31 will be permanently stored and made available via FigShare. We invite the community to provide
32 suggestions for new data sets and new annotation features to be included in subsequent versions,
33 with the goal of building a collaborative resource that will promote research to understand eukaryotic
34 diversity and diversification.

35

36

37

38 **Introduction**

39

40 Over the past 15 years, the discovery of diverse novel microbial eukaryotes, coupled with methods to
41 reconstruct phylogenies based on hundreds of protein-coding genes (known as phylogenomics
42 (Eisen, 2003)) have led to a remarkable reshaping in our understanding of the eukaryotic tree of life,
43 the creation of new supergroups and the placement of enigmatic lineages in known supergroups
44 (Brown et al., 2018; Burki et al., 2012, 2020; Gawryluk et al., 2019; Janouškovec et al., 2017;
45 Kamikawa et al., 2014; Lax et al., 2018; Strassert et al., 2019; Yabuki et al., 2015). The phylogenomic
46 approach has also been used to investigate branching patterns within eukaryotic supergroups,
47 including those with implications for the early evolutionary events in the three major eukaryotic
48 lineages: land plants (Wickett et al., 2014), animals (King & Rokas, 2017) and fungi (Kiss et al., 2019).
49 Furthermore, analyses in diverse eukaryotes have the potential to reveal new genes, pathways and
50 mechanisms of function for biological processes currently characterized only in these three well-
51 studied lineages and their parasites (del Campo et al., 2014; Richter & Levin, 2019). In the ocean
52 alone, planetary-scale metagenomics studies across the full spectrum of life from viruses to animals
53 (Bork et al., 2015; Tara Oceans Coordinators et al., 2020) have already unveiled an extreme diversity

54 of eukaryotic genes (Carradec et al., 2018) whose phylogenetic origin and ecological function are
55 mostly unknown.

56
57 A critical prerequisite to all of these studies is the underlying database of predicted proteins from
58 which orthologs are extracted or other sequence analyses are performed. Because no single website
59 stores the protein predictions for the complete set of eukaryotic taxa that have been sequenced at a
60 genomic scale, each study has needed to assemble their own set, reducing reproducibility among
61 analyses (due to the inclusion of different taxa, or different data sets for the same taxon), and
62 producing a significant barrier to new researchers entering the field. In addition, because the large
63 majority of the protein data sets from diverse eukaryotes are not included in major databases (see
64 Table 1), researchers cannot easily access them via standard tools such as NCBI BLAST (Sayers et
65 al., 2020) in order to find diverse eukaryotic homologs of a protein of interest, nor are valuable
66 annotations such as protein domains (e.g., Pfam (El-Gebali et al., 2019), Interpro (Mitchell et al.,
67 2019)) or gene ontology (The Gene Ontology Consortium, 2019) easily accessible for searching.

68
69 To address this gap, we gathered the predicted proteins from a comprehensive set of species
70 representing known eukaryotic diversity. We placed each species within a unified taxonomic
71 framework, UniEuk (Berney et al., 2017), in order to ensure that the evolutionary relationships among
72 data sets are accurately and consistently described. Our database is designed to prioritize ease of
73 use, with unique, persistent identifiers assigned to each data set and a standard system of
74 nomenclature to facilitate repeatability of analyses.

75
76

77 **The EukProt Database**

78

79 The current version of EukProt (version 2) contains 742 eukaryotic species from 5 different types of
80 sequence data (Figure 1): genome (240 species), single-cell genome (25), transcriptome (453),
81 single-cell transcriptome (7) and expressed sequence tag (EST; 17). The data sets were downloaded
82 from 30 different sources (Table 1), with the two principal sources being NCBI (Sayers et al., 2020)
83 and the Moore Microbial Eukaryote Transcriptome Sequencing Project (MMETSP) (Keeling et al.,
84 2014). The representation among eukaryotic lineages is highly uneven, which is due to the difficulty of
85 discovering, culturing or sequencing species from many lineages, as well as a bias towards
86 sequencing either macroscopic, multicellular species or unicellular species that are parasites (e.g.,
87 Apicomplexa), photosynthetic (such as diatoms, which are part of the Ochrophyte lineage), or are
88 otherwise economically important (del Campo et al., 2014).

89

90 The EukProt database is organized around 5 guiding principles:

91

92 ***Breadth of species phylogenetic diversity:*** the objective of the database is to represent known
93 eukaryotic diversity as fully as possible. In practice, this means that all species with available data
94 sets are included for the majority of lineages. The exceptions are the land plants (Streptophyta),
95 animals (Metazoa) and fungi, and their parasites/pathogens (members of Apicomplexa,
96 Peronosporomycetes, Metakinetoplastina and Fornicata), for which we included a subset of species
97 selected to represent phylogenetic diversity within each group. Within animals, we also emphasized
98 the inclusion of marine planktonic species, as we anticipate these could serve as mapping targets for
99 data from large-scale metagenomic and metatranscriptomic ocean sequencing projects (e.g., *Tara*
100 *Oceans* (Carradec et al., 2018; Karsenti et al., 2011; Richter et al., 2019)).

101

102 ***Convenience of access:*** the full data set and its associated metadata can be downloaded from
103 FigShare with a single click. For 501 of the 724 species in the database, publicly available protein
104 sequences were ready to be included directly. For the remaining species, we processed the publicly
105 available data to produce protein sequences for inclusion in the database. The most common actions

106 were merging independent sets of protein predictions for the same species (105 species), translating
107 mRNA sequences from transcriptomes or EST projects (67), and *de novo* assembling and translating
108 transcriptomes from raw read data (47; we also provide these assemblies as part of the database, as
109 they are not publicly available). A full list of the different types of actions is described in the Methods,
110 and the actions taken for each species are available in the database metadata. We note that 16
111 species in the database are represented by genomic data lacking accompanying protein annotations;
112 15 of these are single-cell genomes. We considered annotation of these genomes to be outside the
113 scope of this release, so we include their nucleotide data in a separate file, which can be searched for
114 proteins of interest using translated homology search software (e.g., tblastn (Altschul, 1997) or BLAT
115 (Kent, 2002)).

116
117 **A unified taxonomic framework:** we placed all species into a unified taxonomic framework, UniEuk
118 (Berney et al., 2017), which is based on the most recent consensus eukaryotic taxonomy (Adl et al.,
119 2019). This will maximize interoperability of the database with other resources built upon to the
120 UniEuk taxonomic framework (e.g., EukBank and EukMap (Berney et al., 2017), EcoTaxa
121 (<https://ecotaxa.obs-vlfr.fr>), and EukRibo (Berney et al., in prep)), and ensure that results arising from
122 the database will be able to use a common set of identifiers to describe analyses at different
123 taxonomic levels. For example, this might include labeling groups in a phylogenetic tree, or
124 summarizing read placement data by taxonomic group. Knowledge of the taxonomy can also facilitate
125 the detection of mis-identified gene sequences within a data set, via comparison of the topology of
126 gene trees to the topology of the taxonomy. In addition to the full taxonomic lineage, we provide for
127 each species the “supergroup”, a very high-level taxonomic grouping, and “taxogroup”, a more fine-
128 grained grouping into evolutionarily or ecologically relevant lineages (see Methods). Finally, for each
129 species we include a list of names it was previously known by, to prevent issues resulting from
130 revisions to species names or lineages, which have occurred frequently with improvements in
131 sequencing and phylogenetic techniques and the discovery of new eukaryotic organisms.

132
133 **Appropriate references to publications and data sources:** to ensure that the researchers who
134 generated and provided each data set receive appropriate credit, we provide the DOI of the
135 publication describing each data set as well as the URL from which it was downloaded. The list of
136 URLs should also allow users of the database to download the original sequences for each data set
137 (although it is not possible to guarantee that the URLs for all data providers will be permanently
138 available).

139
140 **Reusability, persistence and replicability:** the database will be released in successive versions,
141 each of which will be permanently stored and accessible at FigShare. In that way, analyses using the
142 database will need only to specify which version was used, enabling follow-up analyses or replications
143 to begin with the identical database. In addition, each individual data set within the database is
144 assigned a unique, permanent identifier. When a new data set becomes available for a given species,
145 it is assigned a new unique identifier (and the identifier of the data set it replaces, if any, is indicated in
146 the database metadata). These and all other changes between versions of the database will be
147 logged as appropriate.

148
149

150 **Growing the EukProt Database with Community Involvement**

151

152 The core functionality of the database is the distribution of genome-scale protein sequences across
153 the diversity of eukaryotic life and within the UniEuk framework. However, we anticipate that
154 numerous other features might be useful to the community, and we hope to involve the community in
155 suggesting and adding new features in successive versions. These may include information or
156 analyses on full data sets, such as the sequencing technology that was used (e.g., Illumina, PacBio),
157 the 18S ribosomal DNA sequence corresponding to the data set, the completeness of the data set as

158 estimated by software such as BUSCO (Simão et al., 2015), or the estimation of potential
159 contamination levels with non-target species, as inferred using systematic sequence homology
160 searches. Additionally, we could consider adding features on the level of individual protein
161 sequences, such as protein domains from Pfam (El-Gebali et al., 2019)/Interpro (Mitchell et al., 2019),
162 gene ontology (The Gene Ontology Consortium, 2019)/eggNOG (Huerta-Cepas et al., 2019), or the
163 assignment of unique identifiers to individual protein sequences (currently, FASTA files are distributed
164 as is, without modification to headers/identifiers).

165
166 As new genome-scale eukaryotic protein data sets become available, we plan to add them to the
167 database. As yet, we do not have a formal mechanism to accomplish this, and will instead depend on
168 monitoring the literature and assistance from the community. We also plan to add any data sets we
169 may have inadvertently overlooked when building the current version of the database.

170
171 In the longer term, we hope the standardization of our database provides a path towards including all
172 data sets in a major sequence repository such as NCBI/EBI/DDBJ, so that they can be more broadly
173 accessible and integrated into the suites of tools available at these repositories.

174

175

176

177 **Methods**

178

179 *Species and strain identity*

180

181 We determined species and strain identities by reading the publications that described the data sets,
182 consulting the literature for naming revisions, and comparing 18S ribosomal DNA sequences to
183 reference sequence databases. For species that were previously known by other names, we recorded
184 them in the metadata for the data set, except in cases where a species was originally assigned to a
185 genus but not identified to the species level (e.g., *Goniomonas* sp., now identified as *Goniomonas*
186 *avonlea*, is not listed as a previous name).

187

188 *Supergroups and taxogroups from UniEuk*

189

190 Unlike the taxonomic system used in some other resources, the full taxonomic lineages we provide for
191 all species (which follow the framework developed in the UniEuk project) are not based on a fixed
192 number of ranks, but on a free, unlimited number of taxonomic levels, in order to match phylogenetic
193 evidence as closely as possible. This provides end-users more information and flexibility, but could
194 also make it more difficult to summarize results of downstream analyses. Therefore, we provide two
195 additional fields (“supergroup” and “taxogroup”) to help end-users whenever it is useful to distribute
196 eukaryotic diversity into a fixed number of taxonomic categories of equivalent phylogenetic depth or
197 ecological relevance. The 42 “supergroups” (of which 38 are included in EukProt) consist of strictly
198 monophyletic, deep-branching eukaryotic lineages of a phylogenetic depth equivalent to the “classic”
199 Alveolata, Rhizaria, and Stramenopiles. They are therefore highly variable in relative diversity, ranging
200 from clades consisting of a single, orphan genus (e.g. *Ancoracysta*, *Mantamonas*, *Palpitomonas*), to
201 Metazoa as a whole. The “taxogroup” level allows further subdivision of large supergroups into
202 lineages of relatively equivalent evolutionary or ecological relevance, based on current knowledge.
203 This level is more arbitrarily defined; it can include paraphyletic groupings in highly diversified
204 supergroups to accommodate minor lineages of similar ecology or phylogenetic depth. As illustrative
205 examples, diatoms are in the Diatomeae taxogroup within the Stramenopiles supergroup, and
206 coccolithophores are in the Prymnesiophyceae taxogroup within the Haptophyta supergroup. Small,
207 ecologically and morphologically homogeneous supergroups are not subdivided further; in such cases
208 the “taxogroup” level is the same as the “supergroup” level. The same approach will be used in

209 EukRibo, a database of reference ribosomal RNA gene sequences developed in parallel (Berney et
210 al., in prep) to help users link analyses of different types of genetic data.

211

212 *Merging strains from the same species*

213

214 In general, we only included data from a single strain/isolate per species. However, when only a
215 single transcriptome data set was available for a given strain of a species, and there were additional
216 published transcriptome data sets for other strains of the same species, we combined them using CD-
217 HIT (Li & Godzik, 2006) run with default parameter values, in order to guard against the possibility
218 that a single transcriptome might lack genes expressed only in one condition or experiment. When
219 multiple strains were merged to produce a species' data set (there were 25 such cases), this
220 information is indicated in the metadata for the data set.

221

222 *Processing steps applied to publicly available data*

223

224 EukProt metadata indicate the additional steps applied to each data set after downloading from the
225 data source (if any). All software parameter values were default, unless otherwise specified below or
226 in the metadata record for a given data set.

227

228 'assemble mRNA': *de novo* transcriptome assembly using Trinity v. 2.8.4, <http://trinityrnaseq.github.io/>
229 (Haas et al., 2013). We trimmed Illumina input reads for adapters and sequence quality using the
230 built-in '--trimmomatic' option. We trimmed 454 input reads prior to running Trinity with Trimmomatic v.
231 0.3.9, <http://www.usadellab.org/cms/?page=trimmomatic> (Bolger et al., 2014) with the directives
232 'ILLUMINACLIP:[454 adapters FASTA file]:2:30:10 SLIDINGWINDOW:4:5 LEADING:5 TRAILING:5
233 MINLEN:25'.

234

235 'translate mRNA': *de novo* translation of mRNA sequences with Transdecoder v. 5.3.0,
236 <http://transdecoder.github.io/>. When the number of predicted protein sequences for a given species
237 was less than half of the input mRNA sequences, we reduced the minimum predicted protein length to
238 50 (from the default of 100).

239

240 'CD-HIT': clustering of protein sequences to produce a non-redundant data set using CD-HIT v. 4.6,
241 <http://weizhongli-lab.org/cd-hit/> (Li & Godzik, 2006). We used this tool principally to combine protein
242 predictions for different strains of the same species, but also to reduce the size of very large predicted
243 protein sets (> 50,000 proteins) that showed evidence of redundancy.

244

245 'extractfeat', 'transeq', and 'trimseq': from the EMBOSS package v. 6.6.0.0,
246 <http://emboss.sourceforge.net/> (Rice et al., 2000). We used extractfeat to produce coding sequences
247 (CDS) from genomes with gene annotations but without publicly available protein sequences. We
248 used transeq to translate CDS directly into proteins. We used trimseq to trim EST sequences before
249 translation with Transdecoder.

250

251 *The EukProt database distribution on FigShare*

252

253 The database is distributed in a single archive containing five files. One file contains 726 protein data
254 sets, for species with either a genome (239) or single-cell genome (10) with predicted proteins, a
255 transcriptome (453), a single-cell transcriptome (7), or an EST assembly (17). A second file contains
256 16 genomes lacking predicted protein annotations, for 15 species with single-cell genomes, and 1
257 species with a genome sequence. These can be queried for proteins of interest with translated
258 sequence homology search software. A third file contains assembled transcriptome contigs, for 53
259 species with publicly available mRNA sequence reads but no publicly available assembly. The
260 proteins predicted from these assemblies are included in the proteins file. Finally, the database

261 metadata are distributed as two files: one file for the data sets included in the current version of the
262 database (742), and a second file for data sets not included (50), accompanied by the reason they
263 were not included (for example, if the sequences are published but not publicly available or if they
264 were replaced by a higher-quality data set for the same species).

265
266
267

268 **Acknowledgements**

269

270 We thank Michelle Leger for helpful advice, Núria Ros i Rocher for suggestions on the figure, and
271 members of the Multicellgenome Lab for discussions. This work was supported by the French
272 Government 'Investissement d'Avenir' program OCEANOMICS (ANR-11-BTBR-0008) and the ABiMS
273 computing cluster at the Station Biologique de Roscoff, France. DJR was supported by postdoctoral
274 fellowships from the Beatriu de Pinós programme of the Government of Catalonia's Secretariat for
275 Universities and Research of the Ministry of Economy and Knowledge, and from "la Caixa"
276 Foundation (ID 100010434), with the fellowship code LCF/BQ/PI19/11690008. CB was supported by
277 a grant from the Gordon and Betty Moore Foundation (GBMF5257 / UniEuk) and is grateful to the
278 International Society of Protistologists for additional support.

279

280

281

282 **References**

283

284 Adl, S. M., Bass, D., Lane, C. E., Lukeš, J., Schoch, C. L., Smirnov, A., Agatha, S., Berney, C.,
285 Brown, M. W., Burki, F., Cárdenas, P., Čepička, I., Chistyakova, L., del Campo, J., Dunthorn,
286 M., Edvardsen, B., Eglit, Y., Guillou, L., Hampl, V., ... Zhang, Q. (2019). Revisions to the
287 Classification, Nomenclature, and Diversity of Eukaryotes. *Journal of Eukaryotic Microbiology*,
288 66(1), 4–119. <https://doi.org/10.1111/jeu.12691>

289 Altschul, S. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search
290 programs. *Nucleic Acids Research*, 25(17), 3389–3402.
291 <https://doi.org/10.1093/nar/25.17.3389>

292 Berney, C., Ciuprina, A., Bender, S., Brodie, J., Edgcomb, V., Kim, E., Rajan, J., Parfrey, L. W., Adl,
293 S., Audic, S., Bass, D., Caron, D. A., Cochrane, G., Czech, L., Dunthorn, M., Geisen, S.,
294 Glöckner, F. O., Mahé, F., Quast, C., ... de Vargas, C. (2017). *UniEuk*: Time to Speak a
295 Common Language in Protistology! *Journal of Eukaryotic Microbiology*, 64(3), 407–411.
296 <https://doi.org/10.1111/jeu.12414>

297 Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: A flexible trimmer for Illumina sequence
298 data. *Bioinformatics*, 30(15), 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>

299 Bork, P., Bowler, C., de Vargas, C., Gorsky, G., Karsenti, E., & Wincker, P. (2015). Tara Oceans
300 studies plankton at planetary scale. *Science*, 348(6237), 873–873.
301 <https://doi.org/10.1126/science.aac5605>

302 Brown, M. W., Heiss, A. A., Kamikawa, R., Inagaki, Y., Yabuki, A., Tice, A. K., Shiratori, T., Ishida, K.-
303 I., Hashimoto, T., Simpson, A. G. B., & Roger, A. J. (2018). Phylogenomics Places Orphan
304 Protistan Lineages in a Novel Eukaryotic Super-Group. *Genome Biology and Evolution*, 10(2),
305 427–433. <https://doi.org/10.1093/gbe/evy014>

306 Burki, F., Okamoto, N., Pombert, J.-F., & Keeling, P. J. (2012). The evolutionary history of
307 haptophytes and cryptophytes: Phylogenomic evidence for separate origins. *Proceedings of*
308 *the Royal Society B: Biological Sciences*, 279(1736), 2246–2254.
309 <https://doi.org/10.1098/rspb.2011.2301>

310 Burki, F., Roger, A. J., Brown, M. W., & Simpson, A. G. B. (2020). The New Tree of Eukaryotes.
311 *Trends in Ecology & Evolution*, 35(1), 43–55. <https://doi.org/10.1016/j.tree.2019.08.008>

312 Carradec, Q., Pelletier, E., Da Silva, C., Alberti, A., Seeleuthner, Y., Blanc-Mathieu, R., Lima-Mendez,

- 313 G., Rocha, F., Tirichine, L., Labadie, K., Kirilovsky, A., Bertrand, A., Engelen, S., Madoui, M.-
314 A., Méheust, R., Poulain, J., Romac, S., Richter, D. J., Yoshikawa, G., ... Wincker, P. (2018).
315 A global ocean atlas of eukaryotic genes. *Nature Communications*, 9(1), 373.
316 <https://doi.org/10.1038/s41467-017-02342-1>
- 317 del Campo, J., Sieracki, M. E., Molestina, R., Keeling, P., Massana, R., & Ruiz-Trillo, I. (2014). The
318 others: Our biased perspective of eukaryotic genomes. *Trends in Ecology & Evolution*, 29(5),
319 252–259. <https://doi.org/10.1016/j.tree.2014.03.006>
- 320 Eisen, J. A. (2003). Phylogenomics: Intersection of Evolution and Genomics. *Science*, 300(5626),
321 1706–1707. <https://doi.org/10.1126/science.1086292>
- 322 El-Gebali, S., Mistry, J., Bateman, A., Eddy, S. R., Luciani, A., Potter, S. C., Qureshi, M., Richardson,
323 L. J., Salazar, G. A., Smart, A., Sonnhammer, E. L. L., Hirsh, L., Paladin, L., Piovesan, D.,
324 Tosatto, S. C. E., & Finn, R. D. (2019). The Pfam protein families database in 2019. *Nucleic
325 Acids Research*, 47(D1), D427–D432. <https://doi.org/10.1093/nar/gky995>
- 326 Gawryluk, R. M. R., Tikhonenkov, D. V., Hehenberger, E., Husnik, F., Mylnikov, A. P., & Keeling, P. J.
327 (2019). Non-photosynthetic predators are sister to red algae. *Nature*, 572(7768), 240–243.
328 <https://doi.org/10.1038/s41586-019-1398-6>
- 329 Haas, B. J., Papanicolaou, A., Yassour, M., Grabherr, M., Blood, P. D., Bowden, J., Couger, M. B.,
330 Eccles, D., Li, B., Lieber, M., MacManes, M. D., Ott, M., Orvis, J., Pochet, N., Strozzi, F.,
331 Weeks, N., Westerman, R., William, T., Dewey, C. N., ... Regev, A. (2013). De novo
332 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference
333 generation and analysis. *Nature Protocols*, 8(8), 1494–1512.
334 <https://doi.org/10.1038/nprot.2013.084>
- 335 Huerta-Cepas, J., Szklarczyk, D., Heller, D., Hernández-Plaza, A., Forslund, S. K., Cook, H., Mende,
336 D. R., Letunic, I., Rattei, T., Jensen, L. J., von Mering, C., & Bork, P. (2019). eggNOG 5.0: A
337 hierarchical, functionally and phylogenetically annotated orthology resource based on 5090
338 organisms and 2502 viruses. *Nucleic Acids Research*, 47(D1), D309–D314.
339 <https://doi.org/10.1093/nar/gky1085>
- 340 Janouškovec, J., Tikhonenkov, D. V., Burki, F., Howe, A. T., Rohwer, F. L., Mylnikov, A. P., & Keeling,
341 P. J. (2017). A New Lineage of Eukaryotes Illuminates Early Mitochondrial Genome
342 Reduction. *Current Biology*, 27(23), 3717–3724.e5. <https://doi.org/10.1016/j.cub.2017.10.051>
- 343 Kamikawa, R., Kolisko, M., Nishimura, Y., Yabuki, A., Brown, M. W., Ishikawa, S. A., Ishida, K.,
344 Roger, A. J., Hashimoto, T., & Inagaki, Y. (2014). Gene Content Evolution in Discobid
345 Mitochondria Deduced from the Phylogenetic Position and Complete Mitochondrial Genome
346 of *Tsukubamonas globosa*. *Genome Biology and Evolution*, 6(2), 306–315.
347 <https://doi.org/10.1093/gbe/evu015>
- 348 Karsenti, E., Acinas, S. G., Bork, P., Bowler, C., De Vargas, C., Raes, J., Sullivan, M., Arendt, D.,
349 Benzoni, F., Claverie, J.-M., Follows, M., Gorsky, G., Hingamp, P., Iudicone, D., Jaillon, O.,
350 Kandels-Lewis, S., Krzic, U., Not, F., Ogata, H., ... Wincker, P. (2011). A Holistic Approach to
351 Marine Eco-Systems Biology. *PLoS Biology*, 9(10), e1001177.
352 <https://doi.org/10.1371/journal.pbio.1001177>
- 353 Keeling, P. J., Burki, F., Wilcox, H. M., Allam, B., Allen, E. E., Amaral-Zettler, L. A., Armbrust, E. V.,
354 Archibald, J. M., Bharti, A. K., Bell, C. J., Beszteri, B., Bidle, K. D., Cameron, C. T., Campbell,
355 L., Caron, D. A., Cattolico, R. A., Collier, J. L., Coyne, K., Davy, S. K., ... Worden, A. Z.
356 (2014). The Marine Microbial Eukaryote Transcriptome Sequencing Project (MMETSP):
357 Illuminating the Functional Diversity of Eukaryotic Life in the Oceans through Transcriptome
358 Sequencing. *PLoS Biology*, 12(6), e1001889. <https://doi.org/10.1371/journal.pbio.1001889>
- 359 Kent, W. J. (2002). BLAT—The BLAST-Like Alignment Tool. *Genome Research*, 12(4), 656–664.
360 <https://doi.org/10.1101/gr.229202>
- 361 King, N., & Rokas, A. (2017). Embracing Uncertainty in Reconstructing Early Animal Evolution.
362 *Current Biology*, 27(19), R1081–R1088. <https://doi.org/10.1016/j.cub.2017.08.054>
- 363 Kiss, E., Hegedüs, B., Virágh, M., Varga, T., Merényi, Z., Kószó, T., Bálint, B., Prasanna, A. N.,
364 Krizsán, K., Kocsubé, S., Riquelme, M., Takeshita, N., & Nagy, L. G. (2019). Comparative

- 365 genomics reveals the origin of fungal hyphae and multicellularity. *Nature Communications*,
366 10(1). <https://doi.org/10.1038/s41467-019-12085-w>
- 367 Lax, G., Eglit, Y., Eme, L., Bertrand, E. M., Roger, A. J., & Simpson, A. G. B. (2018).
368 Hemimastigophora is a novel supra-kingdom-level lineage of eukaryotes. *Nature*, 564(7736),
369 410–414. <https://doi.org/10.1038/s41586-018-0708-8>
- 370 Letunic, I., & Bork, P. (2019). Interactive Tree Of Life (iTOL) v4: Recent updates and new
371 developments. *Nucleic Acids Research*, 47(W1), W256–W259.
372 <https://doi.org/10.1093/nar/gkz239>
- 373 Li, W., & Godzik, A. (2006). Cd-hit: A fast program for clustering and comparing large sets of protein
374 or nucleotide sequences. *Bioinformatics*, 22(13), 1658–1659.
375 <https://doi.org/10.1093/bioinformatics/btl158>
- 376 Marron, A. O., Ratcliffe, S., Wheeler, G. L., Goldstein, R. E., King, N., Not, F., de Vargas, C., &
377 Richter, D. J. (2016). The Evolution of Silicon Transport in Eukaryotes. *Molecular Biology and*
378 *Evolution*, 33(12), 3226–3248. <https://doi.org/10.1093/molbev/msw209>
- 379 Mitchell, A. L., Attwood, T. K., Babbitt, P. C., Blum, M., Bork, P., Bridge, A., Brown, S. D., Chang, H.-
380 Y., El-Gebali, S., Fraser, M. I., Gough, J., Haft, D. R., Huang, H., Letunic, I., Lopez, R.,
381 Luciani, A., Madeira, F., Marchler-Bauer, A., Mi, H., ... Finn, R. D. (2019). InterPro in 2019:
382 Improving coverage, classification and access to protein sequence annotations. *Nucleic Acids*
383 *Research*, 47(D1), D351–D360. <https://doi.org/10.1093/nar/gky1100>
- 384 Rice, P., Longden, I., & Bleasby, A. (2000). EMBOSS: The European Molecular Biology Open
385 Software Suite. *Trends in Genetics*, 16(6), 276–277. [https://doi.org/10.1016/S0168-](https://doi.org/10.1016/S0168-9525(00)02024-2)
386 [9525\(00\)02024-2](https://doi.org/10.1016/S0168-9525(00)02024-2)
- 387 Richter, D. J., & Levin, T. C. (2019). The origin and evolution of cell-intrinsic antibacterial defenses in
388 eukaryotes. *Current Opinion in Genetics & Development*, 58–59, 111–122.
389 <https://doi.org/10.1016/j.gde.2019.09.002>
- 390 Richter, D. J., Watteaux, R., Vannier, T., Leconte, J., Frémont, P., Reygondeau, G., Maillet, N.,
391 Henry, N., Benoit, G., Fernández-Guerra, A., Suweis, S., Narci, R., Berney, C., Eveillard, D.,
392 Gavory, F., Guidi, L., Labadie, K., Mahieu, E., Poulain, J., ... Tara Oceans Coordinators.
393 (2019). *Genomic evidence for global ocean plankton biogeography shaped by large-scale*
394 *current systems* [Preprint]. *Ecology*. <https://doi.org/10.1101/867739>
- 395 Sayers, E. W., Beck, J., Brister, J. R., Bolton, E. E., Canese, K., Comeau, D. C., Funk, K., Ketter, A.,
396 Kim, S., Kimchi, A., Kitts, P. A., Kuznetsov, A., Lathrop, S., Lu, Z., McGarvey, K., Madden, T.
397 L., Murphy, T. D., O'Leary, N., Phan, L., ... Ostell, J. (2020). Database resources of the
398 National Center for Biotechnology Information. *Nucleic Acids Research*, 48(D1), D9–D16.
399 <https://doi.org/10.1093/nar/gkz899>
- 400 Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V., & Zdobnov, E. M. (2015). BUSCO:
401 assessing genome assembly and annotation completeness with single-copy orthologs.
402 *Bioinformatics*, 31(19), 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>
- 403 Strasser, J. F. H., Jamy, M., Mylnikov, A. P., Tikhonenkov, D. V., & Burki, F. (2019). New
404 Phylogenomic Analysis of the Enigmatic Phylum Telonemia Further Resolves the Eukaryote
405 Tree of Life. *Molecular Biology and Evolution*, 36(4), 757–765.
406 <https://doi.org/10.1093/molbev/msz012>
- 407 Tara Oceans Coordinators, Sunagawa, S., Acinas, S. G., Bork, P., Bowler, C., Eveillard, D., Gorsky,
408 G., Guidi, L., Iudicone, D., Karsenti, E., Lombard, F., Ogata, H., Pesant, S., Sullivan, M. B.,
409 Wincker, P., & de Vargas, C. (2020). Tara Oceans: Towards global ocean ecosystems
410 biology. *Nature Reviews Microbiology*. <https://doi.org/10.1038/s41579-020-0364-5>
- 411 The Gene Ontology Consortium. (2019). The Gene Ontology Resource: 20 years and still GOing
412 strong. *Nucleic Acids Research*, 47(D1), D330–D338. <https://doi.org/10.1093/nar/gky1055>
- 413 Wickett, N. J., Mirarab, S., Nguyen, N., Warnow, T., Carpenter, E., Matasci, N., Ayyampalayam, S.,
414 Barker, M. S., Burleigh, J. G., Gitzendanner, M. A., Ruhfel, B. R., Wafula, E., Der, J. P.,
415 Graham, S. W., Mathews, S., Melkonian, M., Soltis, D. E., Soltis, P. S., Miles, N. W., ...
416 Leebens-Mack, J. (2014). Phylotranscriptomic analysis of the origin and early diversification

- 417 of land plants. *Proceedings of the National Academy of Sciences*, 111(45), E4859–E4868.
418 <https://doi.org/10.1073/pnas.1323926111>
- 419 Wideman, J. G., Lax, G., Leonard, G., Milner, D. S., Rodríguez-Martínez, R., Simpson, A. G. B., &
420 Richards, T. A. (2019). A single-cell genome reveals diplonemid-like ancestry of kinetoplastid
421 mitochondrial gene structure. *Philosophical Transactions of the Royal Society B: Biological*
422 *Sciences*, 374(1786), 20190100. <https://doi.org/10.1098/rstb.2019.0100>
- 423 Yabuki, A., Kamikawa, R., Ishikawa, S. A., Kolisko, M., Kim, E., Tanabe, A. S., Kume, K., Ishida, K., &
424 Inagki, Y. (2015). *Palpitomonas bilix* represents a basal cryptist lineage: Insight into the
425 character evolution in Cryptista. *Scientific Reports*, 4(1). <https://doi.org/10.1038/srep04641>

426 **Table 1.** Web hosts from which data sets were downloaded. The count for figshare.com includes 314
427 species from the MMETSP (Keeling et al., 2014) for which a procedure to remove cross-
428 contamination was applied (Marron et al., 2016).
429

Hostname	Number of Data Sets
figshare.com	351
ncbi.nlm.nih.gov	281
datadryad.org	32
ebi.ac.uk	16
rutgers.edu	9
genoscope.cns.fr	8
yadi.sk	6
oist.jp	5
jgi.doe.gov	3
genomics.cn	3
compagen.org	2
nrifs.fra.affrc.go.jp	2
obs-vlfr.fr	1
zenodo.org	1
licebase.org	1
ufl.edu	1
drive.google.com	1
liv.ac.uk	1
bioenergychina.org	1
algaegenome.org	1
giardiadb.org	1
nal.usda.gov	1
ovgu.de	1
umontreal.ca	1
ugent.be	1
bitbucket.org	1
dauidadlergold.com	1
sciencemag.org	1
malab.cn	1
treegenesdb.org	1

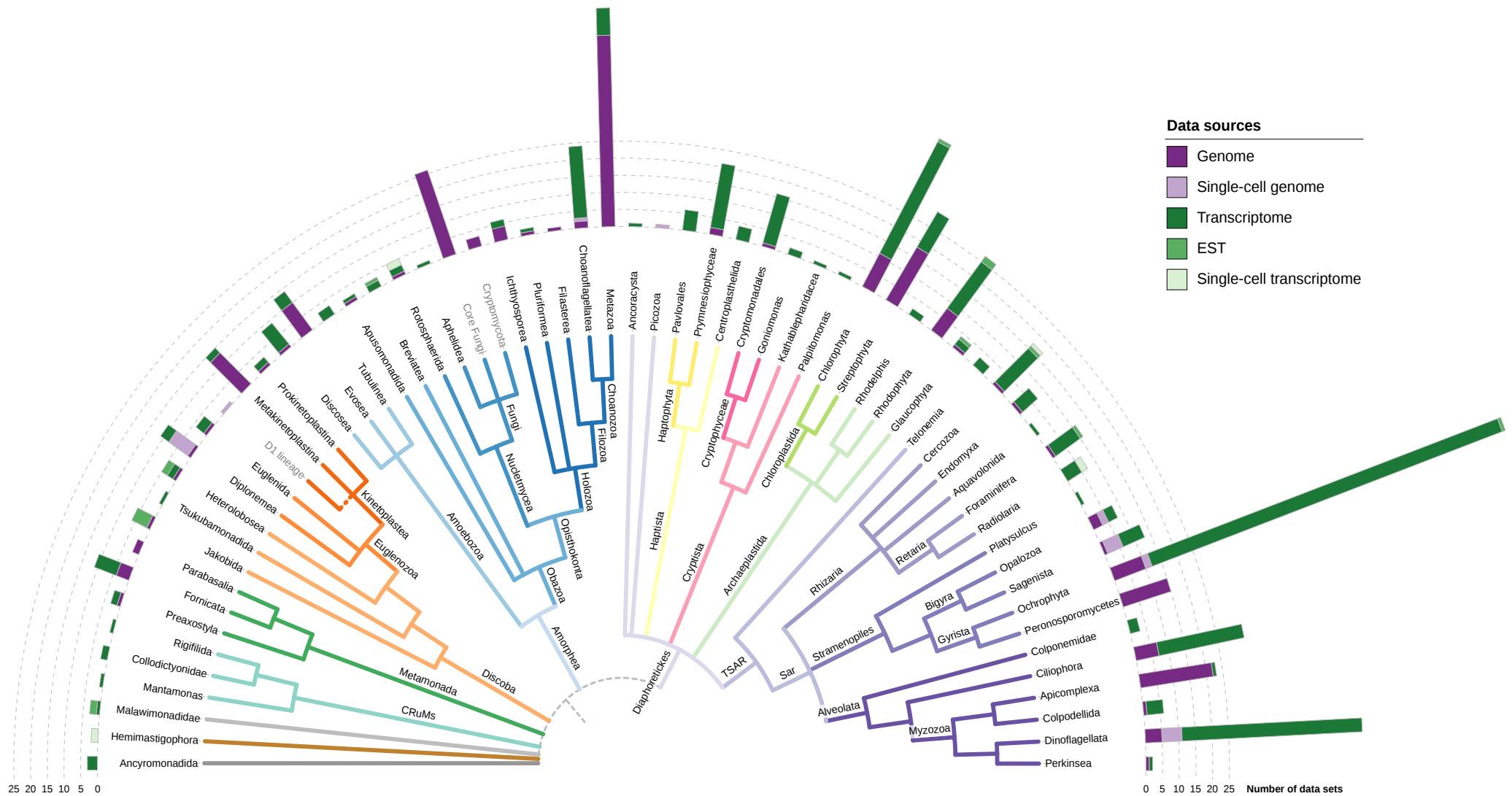


Figure 1. Distribution of 742 source data sets on the eukaryotic tree of life, separated by data set type, with taxonomy based on UniEuk (Adl et al., 2019; Berney et al., 2017). The position of the eukaryotic root, indicated with a dashed line, is currently unresolved. Group names shown in grey are not officially recognized: “D1 lineage” refers to a single-cell genome most closely related to Kinetoplastea (Wideman et al., 2019), whose affiliation with the group is indicated with a dashed line; “Cryptomycota” represents rozellids and microsporidia; “Core Fungi” encompasses all other fungal lineages from chytrids to Dikarya. The figure was created with iTOL (Letunic & Bork, 2019).