

Manuscript version: Author's Accepted Manuscript

The version presented in WRAP is the author's accepted manuscript and may differ from the published version or Version of Record.

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/130153>

How to cite:

Please refer to published version for the most recent bibliographic citation information. If a published version is known of, the repository item page linked to above, will contain details on accessing it.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk.

1 **Long-read metabarcoding of the eukaryotic rDNA operon to phylogenetically and**
2 **taxonomically resolve environmental diversity**

3

4 Running title: Long-read metabarcoding of protists

5

6 Mahwash Jamy¹, Rachel Foster², Pierre Barbera³, Lucas Czech³, Alexey Kozlov³, Alexandros

7 Stamatakis^{3,4}, Gary Bending⁵, Sally Hilton⁵, David Bass^{2,6*}, Fabien Burki^{1,*}

8

9 ¹Science for Life Laboratory, Program in Systematic Biology, Uppsala University, Uppsala,

10 Sweden

11 ²Department of Life Sciences, Natural History Museum, London, UK

12 ³Computational Molecular Evolution Group, Heidelberg Institute for Theoretical Studies,

13 Heidelberg, Germany

14 ⁴Institute of Theoretical Informatics, Karlsruhe Institute of Technology, Karlsruhe, Germany

15 ⁵School of Life Sciences, The University of Warwick, Coventry, UK

16 ⁶Centre for Environment, Fisheries and Aquaculture Science (Cefas), Weymouth, Dorset, UK

17

18

19 **Corresponding authors:**

20 fabien.burki@ebc.uu.se

21 d.bass@nhm.ac.uk

22

23

24

25 **Abstract**

26 High-throughput DNA metabarcoding of amplicon sizes below 500 bp has revolutionized the
27 analysis of environmental microbial diversity. However, these short regions contain limited
28 phylogenetic signal, which makes it impractical to use environmental DNA in full
29 phylogenetic inferences. This lesser phylogenetic resolution of short amplicons may be
30 overcome by new long-read sequencing technologies. To test this idea, we amplified soil
31 DNA and used PacBio Circular Consensus Sequencing (CCS) to obtain a ~4500 bp region
32 spanning most of the eukaryotic SSU (18S) and LSU (28S) ribosomal DNA genes. We first
33 treated the CCS reads with a novel curation workflow, generating 650 high-quality OTUs
34 containing the physically linked 18S and 28S regions. In order to assign taxonomy to these
35 OTUs, we developed a phylogeny-aware approach based on the 18S region that showed
36 greater accuracy and sensitivity than similarity-based methods. The taxonomically-annotated
37 OTUs were then combined with available 18S and 28S reference sequences to infer a well-
38 resolved phylogeny spanning all major groups of eukaryotes, allowing to accurately derive
39 the evolutionary origin of environmental diversity. A total of 1019 sequences were included,
40 of which a majority (58%) corresponded to the new long environmental OTUs. The long-
41 reads also allowed to directly investigate the relationships among environmental sequences
42 themselves, which represents a key advantage over the placement of short reads on a
43 reference phylogeny. Altogether, our results show that long amplicons can be treated in a full
44 phylogenetic framework to provide greater taxonomic resolution and a robust evolutionary
45 perspective to environmental DNA.

46

47 Keywords: metabarcoding, taxonomy, phylogeny, protists, rDNA operon, PacBio

48

49 **Introduction**

50 Sequencing of environmental DNA (eDNA), here encompassing DNA contained in cells as
51 well as cell-free DNA, is a popular approach to study the diversity and ecology of microbial
52 eukaryotes, including small animals, fungi, and protists. eDNA has catalyzed the discoveries
53 of novel lineages at all taxonomic ranks from abundant to rare taxa, and revealed that most, if
54 not all, known groups of microbes are genetically much more diverse than anticipated (de
55 Vargas et al., 2015; Heger et al., 2018; Massana et al., 2015; Pawlowski et al., 2012). For
56 protists, recent global molecular surveys revealed that they can account for up to 80% of the
57 total diversity of eukaryotes in the environments (de Vargas et al., 2015; Logares et al., 2014;
58 Massana et al., 2015; Pawlowski et al., 2012). Initially, these molecular environmental studies
59 relied on cloning the small subunit ribosomal RNA gene (18S rDNA) followed by Sanger
60 sequencing, thereby generating reads of sufficient length to enable reasonably accurate
61 phylogenetic interpretation of the results (Amaral Zettler et al., 2002; Bass & Cavalier-Smith,
62 2004; Dawson & Pace, 2002; Diez et al., 2001; Edgcomb, Kysela, Teske, de Vera Gomez, &
63 Sogin, 2002; Lopez-Garcia, Philippe, Gail, & Moreira, 2003; López-García, Rodríguez-
64 Valera, Pedrós-Alió, & Moreira, 2001; Massana, Balagué, Guillou, & Pedrós-Alió, 2004;
65 Massana, Castresana, et al., 2004; Moon-Van Der Staay, De Wachter, & Vaulot, 2001; Stoeck
66 & Epstein, 2003; Stoeck, Taylor, & Epstein, 2003). Today, however, the overwhelming
67 majority of eDNA data corresponds to much shorter reads produced by Illumina, which
68 routinely generates several millions of reads (e.g. Bates *et al.*, 2013; de Vargas *et al.*, 2015;
69 Geisen, 2016). This enables sequencing a large fraction of the species present in an
70 environment, even including extremely rare organisms (de Vargas et al., 2015; Logares et al.,
71 2014). The drawback of this method is that only genetic regions limited to a few hundred
72 nucleotides (typically <500) can be sequenced at a time, for example the hypervariable V4 or

73 V9 regions of the 18S rDNA or the internal transcribed spacer (ITS) (Mahé et al., 2015;
74 Pawlowski et al., 2012; Stoeck et al., 2010).

75 Short amplicons contain relatively low phylogenetic signal (Dunthorn et al., 2014),
76 which complicates taxonomic identification especially when environmental reads are only
77 distantly related to reference sequences. To address the issue of low phylogenetic signal in
78 high-throughput data, a range of tools has been developed to provide reasonable taxonomic
79 identification of environmental OTUs (Operational Taxonomic Units). Given the mass
80 number of reads available, the most straightforward approach is to use pairwise sequence
81 similarity searches against reference databases (e.g. as done in de Vargas et al., 2015; Mahé et
82 al., 2017). While fast, this approach is highly sensitive to the taxon sampling and annotation
83 accuracy of the reference database. If a taxonomic group is absent or sequences are
84 misannotated in the reference database, the corresponding queries will be only approximately
85 annotated, remain unidentified, or worse, wrongly identified (Berger, Krompass, &
86 Stamatakis, 2011). Recognizing the limitations of similarity-based methods, new tools have
87 been developed that place short sequences into a phylogenetic context. The Evolutionary
88 Placement Algorithm (EPA; implemented in RAxML, or more recently in EPA-ng) (Barbera
89 et al., 2019; Berger et al., 2011) or pplacer (Matsen, Kodner, & Armbrust, 2010) are two such
90 tools. They are becoming popular methods that use a reference tree of carefully selected
91 (often long) sequences to successively score the optimal insertion position of every query
92 sequence or OTU. These methods perform well, and have contributed to the discovery of
93 novel eukaryotic lineages from environments where poor references exist (Bass et al., 2018;
94 Mahé et al., 2017). However, the phylogenetic placement of short reads still requires the
95 independent construction of a reference dataset, which by definition does not include the short
96 reads themselves. Thus, methods like EPA rely on the availability of reference sequences

97 generally produced by the less efficient and more expensive Sanger sequencing, or on genome
98 or transcriptome sequencing projects. Furthermore, references are often based on cell cultures,
99 which are available only for a small fraction of the diversity.

100 To better exploit the phylogenetic signal of the rDNA operon in environmental
101 metabarcoding studies, newer long-read sequencing technologies such as the Pacific
102 Biosciences platform (PacBio) hold great promise. PacBio has lower throughput and higher
103 error rates than Illumina but can produce reads that are over 20kb long at a fraction of the cost
104 of Sanger sequencing. In the last two years, PacBio sequencing has started to be applied to
105 metabarcoding studies, primarily on prokaryotic 16S rDNA (Mosher et al., 2014; Schloss,
106 Jenior, Koumpouras, Westcott, & Highlander, 2016; Wagner et al., 2016) and most recently
107 on larger amplicons also including the 23S rDNA (Martijn et al., 2017). For eukaryotes, the
108 18S rDNA was nearly fully sequenced for targeted microbial groups (Orr et al., 2018), whilst
109 longer regions also spanning the ITS and the 28S gene were used to analyze fungal diversity
110 (Heeger et al., 2018; Tedersoo & Anslan, 2019; Tedersoo, Tooming-Klunderud, & Anslan,
111 2018). These studies showed that in spite of the high error rates of PacBio, when applying a
112 corrective process based on multiple sequence passes (Circular Consensus Sequences - CCS)
113 together with rigorous quality filtering, long-amplicon sequencing is emerging as a robust
114 approach for studying environmental diversity.

115 Here, we used soil eDNA samples to generate broad eukaryote amplicons of about
116 4500 bp spanning the 18S rDNA, ITS1, 5.8S, ITS2, and the 28S rDNA regions. We used
117 PacBio-CCS to sequence these long-amplicons and applied several filtering steps to retain
118 only high-quality sequences. We then followed a full phylogenetic workflow to accurately
119 annotate long-sequences with taxonomy even in the absence of close references. These
120 annotated sequences were combined with available references to infer a well-resolved global

121 eukaryotic phylogeny from a concatenated 18S-28S alignment. Altogether, this study
122 represents an important step forward to use the full power of phylogenetics to derive the
123 accurate evolutionary origins of known and novel lineages present in the environment, as well
124 as expanding rDNA sequence databases for metabarcoding of eukaryotes.

125

126

127 **Materials and Methods:**

128

129 All new scripts listed below are available on Github (<https://github.com/Pbdas/long-reads>)

130

131 **Environmental samples and DNA extraction**

132 We used three environmental soil samples for this study: (1) soil from Tibet, China, collected
133 in summer 2011 from alpine meadows and coniferous forests; (2) rape seed rhizosphere
134 samples from Newbald, Nuneaton, York and Morden in the UK, collected in March 2015; and
135 (3) pooled set-aside agricultural soils from Wellesbourne, UK, collected in September 2010 as
136 described in in (Gosling, van der Gast, & Bending, 2017). Rhizosphere samples were
137 collected as follows: loosely adhering soil was removed from the roots leaving no more than 2
138 mm rhizosphere soil. Roots were washed sequentially in 4 x 25 ml sterile distilled water to
139 release the rhizosphere soil which was then centrifuged and the excess water drained to leave
140 a pellet of rhizosphere soil. All soil samples were extracted using PowerSoil DNA Isolation
141 Kit (MoBio Laboratories) following manufacturer's instructions with the following
142 modifications: (1) rhizosphere soil samples were homogenized in the TissueLyser II (Qiagen)
143 at 20 Hz for 2 x 10 minutes with a 180° rotation of the plates in-between; (2) set-aside

144 agricultural soils were processed using a Precellys 24 homogenizer (Bertin Technologies) for
145 the initial mechanical lysis step.

146

147 **PCR and PacBio sequencing**

148 We used two sets of eukaryotic universal primers to amplify a region covering the 18S, ITS1,
149 5.8S, ITS2, and 28S (Table 1). One 18S internal forward primer, 3NDf (which anneals to the
150 conserved region adjoining the 5' end of the V4 region, *E. coli* position 505) was used in
151 conjunction with two 28S internal reverse primers 21R (*E. coli* position 1926) and 22R (*E.*
152 *coli* position 1952) to amplify a ca. 4500 bp region. The forward and reverse primers are
153 described in (Cavalier-Smith *et al.*, 2009) and (Schwelm, Berney, Dixelius, Bass, &
154 Neuhauser, 2016), respectively and were chosen to maximize the eukaryotic diversity
155 obtained. Taxon coverage of the primers was checked *in silico* using SILVA TestProbe 3.0
156 (Quast *et al.*, 2013): primers 3NDf, 21R and 22R matched against 91.5%, 88.1% and 87.2%
157 of all eukaryotic sequences in SILVA release 132, respectively. For each sample, two PCRs
158 were carried out (one for each combination of forward and reverse primers), and the PCR
159 products were subsequently pooled.

160 PCRs were carried out using the Takara PrimeSTAR GXL high fidelity DNA
161 polymerase, selected for its capacity to amplify long fragments, in 25 µl reactions with 10-20
162 ng of template DNA. The following cycling conditions were used: denaturation at 98 °C for
163 10 s, primer annealing at 60 °C for 15 s, and extension at 68 °C for 90s. A final extension
164 time of 60 s was used after 30 cycles. This protocol corresponds to the rapid PCR protocol of
165 Takara GXL where extension time was shortened by adding twice as much polymerase. PCR
166 products were purified by polyethylene glycol and ethanol precipitation and were pooled and
167 concentrated using Amicon 0.5ml 50K columns (Merck, Germany). Amplicon sizes were

168 checked using TapeStation (Agilent Technologies) before SMRTbell library preparations.
169 Three SMRT cells (one per soil sample) on the PacBio Sequel instrument with v2 chemistry
170 were used for sequencing. Additionally, one RSII SMRT cell was used to sequence a
171 constructed sample with known diversity (see below). Sequencing and library preparation
172 were carried out at Uppsala Genome Center, Science for Life Laboratory, SE-75237 Uppsala.

173

174 **Sequencing of a known community**

175 To validate our curation pipeline and to assess error rates, we constructed a small community
176 of three fungal samples: two unidentified isolates of Agaricomycetes species (BOR77 and
177 BOR79) as well as the species *Phaeosphaeria luctuosa*. We amplified the 18S gene using two
178 sets of primers (Table 1): AU2 and AU4 for BOR77 and BOR79 (Vandenkoornhuysen,
179 Baldauf, Leyval, Straczek, & Young, 2002), and 3NDf and 1510R (Amaral-Zettler,
180 McCliment, Ducklow, & Huse, 2009) for *Phaeosphaeria*. All PCRs were conducted in 20 µl
181 final volumes with 1 µl of template DNA and a final concentration of 0.5 µM of each primer,
182 0.4 mM dNTPs, 2.5 mM of MgCl₂, 0.2 mg bovine serum albumin (BSA), 1x Promega Green
183 Buffer and 0.5 U of Promega GoTaq. Amplicons were sequenced with Sanger sequencing to
184 obtain reference sequences against which the PacBio sequences could be compared. To assess
185 error rate, curated PacBio sequences were searched against the 18S reference sequences with
186 VSEARCH v2.3.4 (Rognes *et al.*, 2016) using the --usearch_global option with the following
187 settings: --id 0.9 --strand both --maxaccepts 0 --top_hits_only --fulldp --userfields
188 query+target+id+alnlen+mism+gaps. The error rate was calculated as (mismatches +
189 indels)/length of alignment.

190

191 **Sequence curation and clustering pipeline**

192 To address PacBio's high error rate, we used a stringent sequence curation pipeline (Fig 1A;
193 Supp. Fig 1A). Circular Consensus Sequences (CCS) were generated from raw reads by
194 SMRT Link v4.0.0.190159 using a minimum number of two passes and Minimum Predicted
195 Accuracy of 0.99 with all other settings set to default. The latter was shown in (Schloss *et al.*,
196 2016) to be the most important factor in decreasing error rate. At this stage, we pooled
197 sequences from the three samples, resulting in one fastq file. A fasta file was generated using
198 the fastq.info command (pacbio=T) in mothur v1.39.5 (Schloss *et al.*, 2009). Sequences at this
199 step of the pipeline still include non-specific PCR amplicons, PCR artifacts such as chimeras
200 and some sequencing errors such as long homopolymer runs. These were filtered out using
201 the trim.seqs command in mothur using the following settings: minlength=2500,
202 maxlength=6000 (to discard non-specific and incomplete PCR amplicons), maxhomop=6 (to
203 stringently discard sequences with a homopolymer run of more 6 nucleotides), and
204 qwindowsize=50 and qwindowaverage=30 (to trim the few sequences with a stretch of low
205 quality sequence). The remaining non-specific PCR amplicons were filtered out by using
206 Barrnap v0.7 (--reject 0.4 --kingdom euk) (<https://github.com/tseemann/barrnap>), which
207 predicts the presence and location of 18S and 28S genes in the sequences. Reads with
208 unexpected structure (more than one 18S, 28S, 5.8S) or incomplete/non-specific reads
209 (missing 18S and/or 28S) were discarded. An in-house perl script was used to identify
210 sequences represented by reverse strand (using the Barrnap output) and subsequently reverse
211 complement them so that all sequences are in the same direction.

212 The sequences were then denoised by pre-clustering as described in Martijn *et al.*,
213 (2019) in order to curate the remaining sequencing errors that are randomly distributed.
214 Briefly, sequences were clustered at 99% similarity using VSEARCH v2.3.4 (Rognes, Flouri,
215 Nichols, Quince, & Mahé, 2016) (--cluster_fast --id 0.99). For each resulting pre-cluster with

216 three or more reads, we aligned the reads with mafft v7.271 (--auto) (Katoh & Standley,
217 2013) and generated a majority-rule consensus sequence using the consensus.seqs (cutoff=51)
218 option in mothur. Gaps were removed to yield final consensus sequences.

219 The denoised sequences as well as sequences from pre-clusters of size one and two
220 were subjected to *de novo* chimera detection using Uchime (Edgar, Haas, Clemente, Quince,
221 & Knight, 2011) (as implemented in mothur) (chunks=40, abskew=1; abundance of the
222 denoised sequences was taken as the number of sequences in their respective pre-clusters).
223 Our PCR primers amplified a few archaea ribosomal genes, and these were filtered out by
224 removing sequences with BLAST hits (Altschul, Gish, Miller, Myers, & Lipman, 1990) to
225 prokaryotic sequences in the SILVA SSU Ref NR 99 database v132 (Quast et al., 2013).
226 Finally, we used in-house perl scripts to extract the 18S and 28S sequences from the cleaned
227 reads, and aligned them with mafft-auto v7.271 (Katoh & Standley, 2013). Poorly aligned
228 sequences were removed after manual inspection.

229 We used the canonical 97% similarity threshold for 18S to cluster sequences into
230 Operational Taxonomic Units (OTUs) using an average-linkage hierarchical clustering
231 method. This was done by first generating a distance matrix using the dist.seqs (cutoff=0.2)
232 command in mothur and then clustering sequences using the cluster command. We used the
233 get.oturep command (label=0.03, method=distance) in mothur to obtain as representative
234 sequence of each OTU the sequence with the smallest distance to all other sequences in the
235 cluster, and extracted the same sequences from the 28S sequence set as OTU representatives.
236 From the total set of 1154 OTUs, we discarded all singletons to be conservative, and obtained
237 a final set of 650 OTUs (hereon referred to as queries).

238

239 **Taxonomic annotation**

240

241 Several datasets were constructed for phylogeny-aware taxonomic annotation and accuracy
242 assessment. These are summarized in Table 2 and described below.

243

244 *Phylogeny-aware annotation*: The 18S gene alone was used for taxonomic annotation as the
245 reference database for this gene is much more comprehensive than its 28S counterpart. The
246 basis for this pipeline is an 18S rDNA tree constructed with both labelled references and the
247 (yet unlabeled) queries. Known reference sequences (RS) were obtained from SILVA SSU
248 Ref NR 99 release 132 (Quast et al., 2013). The RS set comprised two subsets: (1) 504 RS
249 representative of global eukaryotic diversity—these were derived from the 512 taxa dataset
250 used in Mahe et al. 2017; and (2) two to five nearest neighbors of each query in the SILVA
251 database. To obtain these, each query sequence was aligned (mafft --auto) with the top 50
252 BLAST hits against high quality (pintail > 0) eukaryotic SILVA SSU sequences, and pairwise
253 ML distances were computed in RAxML (option -f x) (Stamatakis, 2014) under the
254 GTR+GAMMA model of substitution (Yang, 1994). RS with the lowest pairwise ML
255 distances with the query were selected as the nearest neighbors, resulting in 1157 RS after
256 removing duplicates. Combining the two subsets resulted in a total of 1661 RS, which
257 covered all major eukaryotic groups and, when available, included sequences closely related
258 to queries. The final dataset thus comprised the 650 queries plus the 1661 RS (2311
259 sequences in total; Table 2). These 2311 sequences were aligned with mafft (--retree 2 --
260 maxiterate 1000) and trimmed with trimal (-gt 0.3 -st 0.001), resulting in a multiple sequence
261 alignment (referred to as MSA) with 1589 alignment sites. The best unconstrained maximum
262 likelihood (ML) tree was selected from 20 tree searches run using RAxML-NG (v. 0.6.0)
263 (Kozlov, Darriba, Flouri, Morel, & Stamatakis, 2018). We assumed that the SILVA taxonomy

264 is correct and consistent with the exception of a few cases—preliminary tree searches
265 detected several potentially mislabeled RS, which were relabeled after careful inspection.

266 Based on this tree, a consensus taxonomy was derived using a combination of two
267 strategies (Fig 1B). Strategy 1: Use a custom program written with the Genesis library
268 (Czech, Barbera, & Stamatakis, 2019) to propagate the taxonomy of the closest related
269 reference to each query. Specifically, the program propagates the taxonomic annotation up the
270 tree (where one exists), solving conflicts at inner nodes by taking the intersection of the
271 taxonomic annotation (i.e. lowest common ancestor). Once complete, it propagates that
272 information down to the non-labeled taxa (queries). Strategy 2: Queries were first removed
273 from the tree before being placed back one at a time using EPA-ng (v0.2.1-beta; Barbera et
274 al., 2019). The location and likelihood weights of the placements are then used to compute the
275 taxonomic assignment and the confidence associated with each taxonomic rank as in SATIVA
276 (Kozlov *et al.*, 2016). This last step is implemented in the gappa tool "assign" (Czech et al.,
277 2019) (<https://github.com/lczech/gappa>). Finally, the consensus taxonomy for all queries was
278 produced by a perl script that calculates the intersection of taxonomic paths from strategies 1
279 and 2 (when the SATIVA-derived confidence score for a rank is 0.51 or above). Taxonomic
280 annotations assigned to each query were propagated to their 28S counterparts as they are
281 physically linked on the same molecule.

282

283 *Comparison with short reads*

284 We evaluated the effect of query sequence length by running the taxonomic annotation
285 pipeline with short Illumina reads that were generated *in silico*. We focused on the V4 region
286 (~ 500 bp) of the SSU gene, which is commonly used in barcoding studies, for example in
287 Mahé *et al.*, 2017. This dataset (MSA-V4) was derived from the original MSA by using the

288 V4 flanking primers (Table 1), TAREuk454FWD1 and TAREukREV3 (Stoeck *et al.*, 2010),
289 to trim only the query sequences (median length ~ 340 nucleotides), leaving the rest of the
290 MSA untouched (Table 2). After running the taxonomic annotation pipeline, we performed
291 the following analyses:

292 (1) The accuracy of placement is crucial for correct taxonomic annotation and we
293 compared that for the long and short queries using two metrics. (i) LWR (likelihood
294 weight ratio) of the most probable placement for each query—this is computed as the
295 ratio of the likelihood of the tree with the query at branch *x* to the sum over the
296 likelihoods of all other possible placements (Matsen *et al.*, 2010). (ii) EDPL (Expected
297 Distance between Placement Locations) shows how far the placements are spread
298 across the tree. It is computed as the sum of the distances between placements along
299 the branches of the tree, weighted by their probability (LWR) (Matsen *et al.*, 2010).

300 (2) We conducted pairwise comparisons of the taxonomic assignments and the confidence
301 for taxonomic ranks given to each query based on MSA and MSA-V4.

302

303 *Comparison with sequence similarity-based methods of taxonomic assignment*

304 To assess how our method compares with similarity-based methods, we initially constructed a
305 reference database consisting of high quality (pintail > 0) eukaryotic sequences in SILVA
306 SSU Ref NR99 release 132 and the 504 RS derived from (Mahé *et al.*, 2017). RS were
307 trimmed with the forward primer 3NDf and, both queries and RS were trimmed with the
308 reverse primer 1510R to ensure that they spanned the same region. Queries were searched
309 against this reference set using the global pairwise alignment strategy (--usearch_global
310 option) in VSEARCH v2.3.4 (Rognes *et al.*, 2016) with the following settings: --notrunclabels
311 --userfields query+id1+target --maxaccepts 0 --maxrejects 32 --top_hits_only --output_no_hits

312 --id 0.5 --iddef 1. Sequences were taxonomically assigned based on the top hit, and in case of
313 multiple top hits, the common ancestor of the hits was computed. For each query, the
314 percentage similarity to the closest reference sequence was recorded and the taxonomic
315 classification to deep-branching lineages was compared to that of the phylogeny-based
316 method.

317

318 **Phylogenetic analyses**

319

320 The information on the different alignments used for phylogenetic reconstruction can be
321 found in Table 2.

322

323 *18S+28S global phylogeny*

324 To phylogenetically resolve the biodiversity in our soil samples, we constructed a phylogeny
325 using a concatenated 18S + 28S dataset. For each of the two genes, queries were aligned with
326 their respective reference sequences using mafft v7.271 (--retree 2 --maxiterate 1000).

327 Alignments were filtered with trimal (-gt 0.3 -st 0.001) and a perl script

328 (https://github.com/iirisarri/phylogm/blob/master/concat_fasta.pl) was used to concatenate the

329 SSU and LSU alignments. The phylogeny was inferred with RAxML v8.2.10 as offered on

330 the Cipres web server (Miller, Pfeiffer, & Schwartz, 2010), with 20 tree searches under the

331 GTR+GAMMA model of substitution and 300 non-parametric bootstrap replicates. The

332 construction of the reference dataset is described below.

333 Reference sequences were included only when we could easily verify that the 18S and

334 28S genes originated from the same species or organism. These reference sequences were

335 derived from several public databases, as follows: (i) Searched NCBI nt using the following

336 search filters: ((ribosomal RNA) AND 4000:9000[Sequence Length]) AND
337 Eukaryota[Organism]. (ii) BLASTed whole queries (18S, ITS, 28S) against nt and retained
338 sequences with a minimum HSP of 2500 bp and 80% similarity. (iii) Obtained all 18S and
339 28S sequences from SILVA release 132 possessing the same accession number in the SSU
340 Ref NR 99 and LSU Ref databases. (iv) Included the 108 taxa dataset used in an article
341 studying the eukaryote tree with 18S+28S genes (Moreira et al., 2007). And lastly (v) used
342 barrnap to search all “protist” genomes available in Ensembl Release 92 (Zerbino et al.,
343 2018). This resulted in 3479 taxa after removing duplicates, from which we manually selected
344 sequences, in a best effort, to assemble the most representative dataset possible. Initial tree
345 building attempts placed certain cercozoan and apicomplexan lineages aberrantly among
346 Excavata and Amoebozoa due to long branch attraction. To mitigate this effect, we sorted
347 taxa by branch length (as in Heiss *et al.*, 2018) and removed the longest 118 (10.4 %)
348 branches from subsequent analyses. The final dataset contained 589 queries and 430 reference
349 sequences (1019 total) with 4304 alignment sites (Table 2).

350

351 *Apicomplexa phylogenies*

352 To investigate the effect of query length on resolving environmental diversity in more detail,
353 we constructed additional phylogenies of the fast-evolving group Apicomplexa (Table 2).
354 Reference sequences were obtained by downloading 40 GenBank 18S accessions of which 28
355 accessions had 28S sequences also available. We constructed concatenated and full-length
356 18S genes trees by aligning the references and queries separately using mafft v7.271 (--linsi)
357 and trimming alignments with trimal (-gt 0.3 -st 0.001). Trees were inferred with RAxML,
358 using the substitution model GTR+GAMMA from 20 searches and 100 bootstrap runs.
359 Finally, we constructed an 18S tree from reference sequences alone on which queries

360 shortened to the V4 region (trimmed with universal eukaryotic primers TAREuk454FWD1
361 and TAREukREV3; Stoeck et al., 2010), were placed with EPA-ng v0.3.5.

362

363

364 **Results**

365

366 **Sequence curation**

367 A total of 113,362 long rDNA Circular Consensus Sequences (CCS), all containing two or
368 more passes, were generated with PacBio Sequel. These CCS reads were filtered by a series
369 of stringent quality controls including the removal of non-specific amplicons and prokaryotic
370 sequences, as well as chimera detection (Fig 1A; Supp. Fig 1A). At the end of the curation
371 pipeline, the amplicons had on average 9.95 CCS passes (stdev=2.9) (Supp. Fig 1B). The
372 mean error rate was estimated to be 0.17% based on comparisons between CCS reads curated
373 by our pipeline and known Sanger sequences of the same species of fungi (see materials and
374 methods). OTUs were generated using a 97% similarity threshold based on the 18S region
375 only, leading to 650 high-quality clusters after removing singletons. These OTUs ranged in
376 length from 2501 to 5956 bp (Supp. Fig 1C). Most OTUs contained less than 10 reads, but
377 some were much larger and likely represented the most abundant organisms in the samples;
378 the largest OTU (6416 sequences) corresponded to *Brassica napus*, the main crop species
379 cultivated in one of the samples, while the second largest OTU (1322 sequences) belonged to
380 the gregarines (Apicomplexa), a group of parasites of various invertebrates that has been
381 shown to be particularly abundant in some soil environments (Mahé *et al.*, 2017).

382

383 **Phylogeny-aware taxonomic annotation**

384 In order to annotate the environmental queries with taxonomy, we developed a phylogeny-
385 aware approach that takes advantage of the increased sequence length (Fig 1B). We used only
386 the 18S part of the queries since the taxon sampling of 18S reference sequences is
387 considerably denser than that of 28S sequences. The taxonomic assignments were based on an
388 18S tree (Supp. Fig 2) inferred from the 650 queries together with 1661 full-length references
389 from the SILVA SSU database (i.e., a total of 2311 taxa; Table 2). A consensus taxonomy
390 was then derived from two strategies (Fig 1B; Materials and Methods).

391 Using this approach, we could confidently assign a majority of queries (627/650, or
392 96.5%) to deep-branching eukaryotic lineages (Supp. Fig 3), including queries with similarity
393 to references below 80%. The remaining 23 queries that were not assigned to any of the
394 recognized major lineages were all highly-divergent; of these, 18 could be classified with
395 confidence only to higher-rank assemblages that roughly correspond to the so-called
396 supergroups—the most inclusive established groups of eukaryotes (Burki, Roger, Brown &
397 Simpson, in press). For the remaining five queries, two were ambiguous even at the level of
398 supergroups, thus possibly representing novel deeply-branching lineages and/or sparsely
399 sampled taxonomic groups in the reference database, whilst three proved to be chimeras that
400 had escaped automated filtering. Interestingly, our method also performed well for low-rank
401 taxa, since 226 queries could be reliably annotated down to the genus and species levels
402 (Supp. Table 1).

403 We further investigated the performance of our method by comparing it to a
404 commonly used similarity-based taxonomic annotation tool (VSEARCH, Rognes *et al.*,
405 2016), which revealed several discrepancies. As expected, the most divergent sequences (i.e.
406 <80% similar to known references) showed the highest level of conflicts in taxonomic
407 assignment (Fig 2); 43.7% of these divergent sequences (21/48 queries) were assigned to

408 different deep-branching eukaryotic lineages, and even to different supergroups in four cases.
409 These conflicting assignments became less pronounced for more similar sequences (i.e.,
410 between 80 and 90% similarity), where we observed only 9 (1.6%) conflicts, while there was
411 no conflict for queries >90% similar to a reference (Fig 2).

412 To explore the conflicts between the different approaches in more detail, we focused
413 on the 10 most abundant lineages in our data (Fig 3A). For each lineage, taxonomic
414 assignments by VSEARCH was used as a reference and compared to the assignment derived
415 from our phylogeny-aware method, both using the full-length sequence or the V4 region.
416 Several differences were observed: false negatives, i.e. sequences assigned to the lineage by
417 phylogeny but not by similarity; false positives, i.e. sequences assigned to the lineage by
418 similarity but not by phylogeny; and higher-rank assignments (but not conflicting), i.e.
419 sequences assigned by phylogeny only to a more inclusive rank in the same taxonomic path.
420 The amount and type of differences were to a large extent group-specific. Three groups
421 contained no conflicting taxonomy (Ciliophora, Phytomyxea and Tubulinea), only a small
422 number of higher-rank assignments by V4 phylogenetic assignment. All other groups,
423 however, showed some levels of conflicts. Apicomplexa and Zoopagomycota displayed the
424 highest number of false negatives, with both phylogeny-based approaches classifying more
425 queries to these groups than VSEARCH (42.5% and 100% more queries respectively, blue
426 bars in Fig 3A). False positives were relatively more abundant in Colpodellida, where ~40%
427 of queries assigned to this group by VSEARCH was assigned to a different group by one or
428 the other phylogenetic method (pink bars in Fig 3A).

429 We found that a key difference between similarity and phylogeny is that the latter
430 approach is much more flexible in the level of taxonomic resolution without requiring
431 subjective decisions. For instance, ~23% of the queries with <90% similarity to known

432 sequences were conservatively classified to higher taxonomic ranks by our approach
433 compared to VSEARCH (Fig 2; Supp. Fig 4). A comparison of the lowest rank assignments
434 by all three methods illustrates this behavior (Fig 3B). The similarity method always
435 classified queries to the same predetermined rank, here corresponding to one of the 10 most
436 abundant lineages in our data. In sharp contrast, both phylogenetic methods displayed a
437 broader range of taxonomy, from higher to lower ranks (sometimes to species-level),
438 depending on the confidence in the assignment. Interestingly, the added information from the
439 longer sequences (long versus V4) translated into increased taxonomic resolution, i.e. more
440 assignments towards lower ranks (for example Phytomyxea in Fig 3B; Supp. Fig 5).
441 Furthermore, in the absence of closely-related references, our method can correctly propose
442 no specific annotation. A good test for this case was for the recently suggested supergroup-
443 level lineage Hemimastigophora was a good test case (Lax et al., 2018). One query with 85%
444 sequence similarity to Streptophyta in SILVA was labelled as an “unidentified eukaryote” by
445 our method, whilst it was logically annotated as a land plant by VSEARCH. When using
446 GenBank instead, this query revealed to be 98% similar to a newly added hemimastigote
447 sequence, thus not a land plant but indeed no specific grouping in the absence of that
448 sequence.

449

450 **Combined 18S-28S rDNA phylogeny of environmental DNA**

451 The availability of long queries allows, in principle, to better resolve the origin of
452 environmental sequences due to increased phylogenetic signal. We assembled a concatenated
453 18S-28S dataset including the annotated queries and reference sequences mined from various
454 public databases. The references were selected such that it could be verified that both the 18S
455 and 28S rDNA sequences originated from the same species (see material and method). We

456 included representatives of all major eukaryotic lineages where possible. In addition,
457 preliminary tree searches were used to identify long-branching taxa which were removed in
458 downstream analyses to reduce potential long branch attraction artifacts. This yielded a final
459 dataset of 1019 taxa, of which a majority (589 taxa = 58%) represented new environmental
460 queries. Importantly, because the 18S sequences are physically linked to their 28S counterpart
461 on the CCS reads, the taxonomic annotation inferred with our phylogeny-aware method could
462 be transferred to the combined 18S-28S reads. This provided a diverse set of taxonomically
463 annotated environmental queries in otherwise sparsely populated reference sequences (Fig 4).

464 Figure 4 shows a Maximum Likelihood (ML) tree of the 1019 taxa dataset. The
465 phylogenetic relationships were in general agreement with previous phylogenies based on the
466 18S and 28S (Moreira et al., 2007; Zhao et al., 2012), even recovering several well-
467 established supergroups that were first proposed based on substantially larger concatenated
468 protein datasets such as Sar (including the subclades Stramenopila, Alveolata, and Rhizaria)
469 or Opisthokonta (including the subclade Holomycota and Holozoa) (Baldauf, Roger, Wenk-
470 Siefert, & Doolittle, 2000; Burki et al., 2007). Overall, more than half of the newly sequenced
471 diversity (345 queries; 53% of all queries) corresponded to microbial taxa other than fungi or
472 animals (Fig 4). Members of Alveolata and Rhizaria accounted for nearly 70% of these protist
473 queries—the most dominant lineages in decreasing number of queries were Ciliophora,
474 Apicomplexa, Cercomonadida, Phytomyxea, Glissomonadida, and Vampyrellida. The
475 remaining sequenced diversity was dominated by fungal lineages, accounting for 203 queries
476 (31% of all queries) that equally represented dikarya (Ascomycota and Basidiomycota) and
477 the so-called early-diverging fungi (EDF). Of these EDF, Cryptomycota and Chytridiomycota
478 were particularly diverse. The remaining 16% of the queries corresponded to various animal
479 lineages as well as land plants.

480

481 **Comparison to 18S-only and V4-based phylogenetic classification of environmental**
482 **DNA**

483 The combined 18S-28S tree described above (Fig 4) provides a new solution for
484 obtaining a taxonomically annotated and well-resolved phylogenetic framework from high-
485 throughput environmental sequencing. To assess to which extent the added information of the
486 28S gene improved the phylogenetic resolution, we first compared the combined tree to the
487 18S-only tree constructed for the taxonomic assignment. Interestingly, both trees were largely
488 in agreement, suggesting that the ~1000bp-fragment sequenced for the 18S gene combined
489 with the substantially denser reference sampling available for this gene provided sufficient
490 phylogenetic signal to recover many groupings. However, the combined tree received
491 generally higher bootstrap support values: 54.3% of the bipartitions (552/1016) received \geq
492 75% bootstrap support in the combined tree compared to 43.1% bipartitions (994/2308) in the
493 18S tree. The combined tree also supported (bootstrap > 75%) more specific phylogenetic
494 position for a few queries that were taxonomically annotated only to high-rank taxa based on
495 the 18S tree. For example, two queries labelled as Opisthokonta could be assigned more
496 precisely to Aphelidea and as sister to nucleariid in the combined 18S-28S tree, respectively;
497 or one deep branching eukaryote in the 18S tree in fact corresponded to a long branch within
498 Ascomycota in the combined tree.

499 To investigate the benefits of long reads for phylogeny-based resolution of
500 environmental diversity in more detail, we constructed three additional datasets with varying
501 sequence lengths focusing on the Apicomplexa (Table 2). The sequence lengths corresponded
502 to i) the combined 18S-28S alignment ii) the full-length 18S-only alignment, and iii) an
503 alignment of full-length 18S reference sequences but with query sequences shortened to the

504 V4 region. The taxon-sampling across the three datasets was identical to facilitate
505 comparison, containing 67 queries and 40 reference sequences. The inspection of the
506 combined and the 18S trees (Supp. Figs 6-7) revealed no major discrepancies and placed all
507 56 queries among gregarines. As with the full eukaryotic tree, the bootstrap values were
508 globally higher in the combined tree but many relationships remained unsupported. However,
509 we observed several exceptions where the increased resolution of the combined tree allowed
510 for better interpretation. Most importantly, the monophyly of Apicomplexa was statistically
511 supported in the combined tree (83%) whereas it was unsupported in the 18S tree (23%). The
512 same was observed for the monophyly of other established groupings, such as the haematozoa
513 (79% vs. 55% in the combined and 18S trees, respectively) and haematozoa + coccidians
514 (92% vs. 47% in the combined and 18S tree respectively). Furthermore, the combined tree
515 (Supp. Fig 6) resolved the eugregarine superfamily, Actinocephaloidea, into two separate
516 clades with moderate to strong support (99% and 75%, respectively), while they did not form
517 separate clades in the 18S tree as previously noted based on this gene only. The comparison
518 with the phylogenetic placement of the V4 query sequences revealed that EPA-ng
519 successfully placed all apicomplexan queries among gregarines with the exception of one
520 query, which was placed close to *Plasmodium* instead. However, the reference-only 18S tree
521 had a different topology than the full 18S tree, presumably because the short queries were not
522 used for inferring the latter tree (Supp. Fig 7-8). Furthermore, a close inspection showed that
523 three queries were probably misplaced by EPA-ng on the branch leading to the
524 cephaloidophorids (parasites of marine invertebrates; Supp. Fig 8) when instead they formed
525 a robust monophyletic clade putatively representing novel lineages on both the 18S and
526 combined trees (Supp. Fig 6 and 7).

527

528

529 **Discussion**

530 In this study, we broadly sequenced the near-complete eukaryotic rDNA operon from
531 environmental soil samples, using PacBio sequencing. To our knowledge this is the first long
532 amplicon environmental sequencing study that uses a full phylogenetic approach to assess the
533 diversity of all eukaryotes. To reduce the inherently high error rate of PacBio, we combined
534 the Circular Consensus Sequencing (CCS) approach with a series of stringent filtering steps
535 and clustering. The final error rate of the CCS reads has a mean of 0.17%, which is
536 comparable with the error rate of Illumina (0.21%; Schirmer, D'Amore, Ijaz, Hall, & Quince,
537 2016), or in other PacBio-based studies (Schloss et al., 2016; Tedersoo et al., 2018; Wagner et
538 al., 2016). Even though the curating pipeline discarded the majority of CCS reads, many of
539 which might still be of high quality, the 650 OTUs that passed all filtering steps comprised a
540 large and broad diversity of eukaryotes. Almost all major microbial lineages were sampled,
541 from known abundant taxa in soils such as Ciliophora, Cercozoa, Apicomplexa, and fungi, to
542 rarer lineages in soil such as the mainly aquatic Bacillariophyceae (diatoms) and Chlorophyta
543 (green algae) (Bahram et al., 2018; Foissner & W., 1987; Stefan Geisen et al., 2018, 2015;
544 Stephen Geisen, Cornelia, Jörg, & Michael, 2014; Mahé et al., 2017). A few main protist
545 groups lacked new OTUs altogether, including Cryptista, Retaria, Rhodophyceae, and
546 Glaucophyta, but these are almost exclusively aquatic and thus less likely to be recovered
547 among soil sequences even if present in the environment at very low abundance (de Vargas et
548 al., 2015; Stefan Geisen et al., 2018; Lallias et al., 2015). However, not all major groups
549 typically widespread in soils were recovered with a correspondingly high sequence diversity.
550 This was for example the case for Amoebozoa, Excavata, or Centrohelida, whose low

551 diversity might be at least partially explained by primer bias (see materials and methods),
552 and/or by the ecological conditions represented by the samples.

553 The availability of longer environmental sequences opens up the possibility to
554 phylogenetically resolve environmental diversity with improved accuracy. Previous studies
555 employing both the 18S and 28S genes recovered many relationships within and between
556 major eukaryotic groups with greater resolution than that afforded by the 18S alone (Moreira
557 et al., 2007; Zhao et al., 2012). The use of both genes was proposed to more robustly derive
558 the origin of environmental sequences, particularly in the case of fast-evolving taxa, but this
559 was based on Sanger sequencing of clone libraries (Marande, López-García, & Moreira,
560 2009). Near full-length 18S amplicons and even longer fragments including parts of the 28S
561 have also recently been sequenced with PacBio for group-specific investigations,
562 demonstrating that long-read high-throughput sequencing is a promising complement to
563 Illumina for investigating the environmental diversity of eukaryotes (Heeger et al., 2018; Orr
564 et al., 2018). Here, we extended the approach to ~4500bp of the rDNA operon across the
565 whole phylogenetic diversity of eukaryotes. We built a combined 18S-28S tree of eukaryotes
566 that is globally well-resolved and can serve as a robust phylogenetic framework to describe
567 the environmental diversity in samples (Fig 4). Comparisons to the 18S region alone of the
568 queries (~1200 bp) provided a similar overall topology to the combined tree, but with lower
569 overall resolution (Supp. Fig 2). Furthermore, some key groups in the apicomplexan
570 phylogeny were either missing or not supported by the 18S-only tree, a pattern that was also
571 recovered by previous analyses (Simdyanov et al., 2017, 2018). Altogether, our phylogenetic
572 comparisons revealed that the 18S and 28S together provide increased resolution compared to
573 the 18S alone, but the differences between single and two-gene trees vary across groups.

574 In order to assign taxonomy to the long environmental reads, we applied a novel
575 phylogeny-aware approach that enables deriving robust annotation even in the absence of
576 closely related references. Most commonly, taxonomic annotation is conducted by similarity
577 comparison to reference databases (e.g. in de Vargas *et al.*, 2015; Mahé *et al.*, 2017).
578 Similarity works well when closely related references *are* available, however it requires the
579 use of arbitrary similarity cutoffs without biological grounding below which sequences are
580 considered of unknown origins (Bahram *et al.*, 2018; Stoeck *et al.*, 2010). To enable the use of
581 phylogenetics with short environmental reads, methods such as the Evolutionary Placement
582 Algorithm (EPA) have been recently developed and successfully applied to microbial
583 diversity (Bass *et al.*, 2018; Mahé *et al.*, 2017). Whilst the need for similarity cutoffs is
584 alleviated, the EPA still requires longer reference sequences to build a stable evolutionary
585 framework and thus does not fully overcome the limitations of short read sequencing when
586 references are lacking. Whilst our method relies partly on the EPA, it makes explicit use of
587 environmental sequences to build a reference tree and computes a confidence score for each
588 taxonomic rank. We show that it provides accurate taxonomic annotation with ranks
589 corresponding to the phylogenetic position of queries in the reference tree—higher ranks
590 correspond to deeper branches in the phylogeny—and that it performs better than similarity-
591 based methods for divergent sequences ($\leq 90\%$ similarity). Comparison with the classical use
592 of EPA with V4 reads revealed that whilst the overall annotations were similar, our approach
593 utilizing long queries led to higher confidence scores. It was also more informative than
594 placing short reads on a reference phylogeny, because the long queries directly contributed to
595 the phylogenetic inference by filling gaps between references. Thus, the relationships between
596 the queries themselves can be determined to reveal whether they cluster around known
597 sequences or form entirely new clades.

598 One of the main benefits of our approach is that it provides both the 18S and 28S
599 genes for the *same amplicon*. The 18S gene has long been the reference molecular marker for
600 environmental studies of protist diversity (de Vargas et al., 2015; Diez et al., 2001; López-
601 García et al., 2001; Massana, Balagué, et al., 2004; Massana et al., 2015; Moon-Van Der
602 Staay et al., 2001). With this approach, each 18S sequence should be paired with its 28S
603 counterpart (or ITS). As a result, we rapidly generated a massive increase of 28S sequence
604 diversity for which the attached 18S provides a direct link to the much larger availability of
605 18S sequences contained in databases such as SILVA, PR2, or GenBank. As a point of
606 comparison, the new sequences produced in this study alone represented the majority (58%)
607 of all broad eukaryote diversity for which we could gather reference sequences for both
608 genes. At lower taxonomic ranks, the increase in sequence diversity can be even more
609 significant. For example, we found a total of only nine species of gregarines (Apicomplexa)
610 that have both 18S and 28S genes in public databases. Here, we obtained 56 new gregarine
611 OTUs, corresponding to a 6-fold increase in diversity for this group. Thus, we suggest that the
612 newly generated long environmental sequences can be used in future studies as
613 taxonomically-annotated “anchor” sequences to fill phylogenetic gaps in addition to the more
614 traditional Sanger reference sequences.

615 In conclusion, we demonstrate several advantages of using high-throughput long
616 sequence metabarcoding for environmental studies of microbial eukaryote diversity. With
617 longer reads comes improved phylogenetic signal, and we show that it is possible to employ a
618 full phylogenetic approach to taxonomically classify sequences and obtain a robust
619 evolutionary framework of environmental diversity. This approach can be adapted for use
620 with other emerging long-read technologies, e.g. Nanopore sequencing, and may prove
621 particularly powerful in combination with even higher-throughput sequencing technologies

622 such as Illumina. Indeed, it will then be possible to map shorter but more abundant reads on a
623 much more comprehensive reference phylogeny obtained from the same environments. The
624 importance of eDNA studies continually grows in fields as varied as conservation biology,
625 evolutionary biology and ecology. Long metabarcoding of the eukaryotic rDNA operon will
626 undoubtedly play an increasingly important role in the close future.

627

628

629 **Acknowledgements.** We thank Stefan Geisen and Junling Zhang for providing soil samples
630 from Tibet. We further thank Thijs J.G. Ettema and Joran Martijn for their suggestions while
631 developing the read curation pipeline, and Vasily Zlatogursky for providing isolates for the
632 mock community. This work was supported by a grant from Science for Life Laboratory
633 available to FB, which covered salary of MJ, and experimental expenses. The work of DB and
634 RF was supported by the Standard Research Grant (NE/H009426/1), UK Department of
635 Environment, Food and Rural Affairs (Defra) under contract FC1214. The work of PB, LC
636 and AK were financially supported by the Klaus Tschira Foundation. The authors would like
637 to acknowledge support of the Uppsala Genome Center for providing assistance in massive
638 parallel sequencing. Work performed at Uppsala Genome Center has been funded by VR and
639 Science for Life Laboratory, Sweden.

640

641

642 **References**

643

- 644 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., & Lipman, D. J. (1990). Basic local
645 alignment search tool. *Journal of Molecular Biology*, *215*(3), 403–410. doi:
646 10.1016/S0022-2836(05)80360-2
647 Amaral-Zettler, L. A., McCliment, E. A., Ducklow, H. W., & Huse, S. M. (2009). A method for
648 studying protistan diversity using massively parallel sequencing of V9 hypervariable

649 regions of small-subunit ribosomal RNA Genes. *PLoS ONE*, 4(7), e6372. doi:
650 10.1371/journal.pone.0006372

651 Amaral Zettler, L. A., Gómez, F., Zettler, E., Keenan, B. G., Amils, R., & Sogin, M. L. (2002).
652 Eukaryotic diversity in Spain's River of Fire. *Nature*, 417(6885), 137–137. doi:
653 10.1038/417137a

654 Bahram, M., Hildebrand, F., Forslund, S. K., Anderson, J. L., Soudzilovskaia, N. A.,
655 Bodegom, P. M., ... Bork, P. (2018). Structure and function of the global topsoil
656 microbiome. *Nature*, 560(7717), 233–237. doi: 10.1038/s41586-018-0386-6

657 Baldauf, S. L., Roger, A. J., Wenk-Siefert, I., & Doolittle, W. F. (2000). A kingdom-level
658 phylogeny of eukaryotes based on combined protein data. *Science*, 290(5493), 972–977.
659 doi: 10.1126/science.284.5423.2124

660 Barbera, P., Kozlov, A. M., Czech, L., Morel, B., Darriba, D., Flouri, T., & Stamatakis, A.
661 (2019). EPA-ng: massively parallel evolutionary placement of genetic sequences.
662 *Systematic Biology*, 68(2), 365–369. doi: 10.1093/sysbio/syy054

663 Bass, D., & Cavalier-Smith, T. (2004). Phylum-specific environmental DNA analysis reveals
664 remarkably high global biodiversity of Cercozoa (Protozoa). *International Journal of*
665 *Systematic and Evolutionary Microbiology*, 54(6), 2393–2404. doi:
666 10.1099/ijs.0.63229-0

667 Bass, D., Czech, L., Williams, B. A. P., Berney, C., Dunthorn, M., Mahé, F., ... Williams, T. A.
668 (2018). Clarifying the Relationships between Microsporidia and Cryptomycota. *Journal*
669 *of Eukaryotic Microbiology*, 65(6), 773–782. doi: 10.1111/jeu.12519

670 Bates, S. T., Clemente, J. C., Flores, G. E., Walters, W. A., Parfrey, L. W., Knight, R., & Fierer,
671 N. (2013). Global biogeography of highly diverse protistan communities in soil. *The*
672 *ISME Journal*, 7(3), 652–659. doi: 10.1038/ismej.2012.147

673 Berger, S. A., Krompass, D., & Stamatakis, A. (2011). Performance, Accuracy, and Web Server
674 for Evolutionary Placement of Short Sequence Reads under Maximum Likelihood.
675 *Systematic Biology*, 60(3), 291–302. doi: 10.1093/sysbio/syr010

676 Burki, F., Shalchian-Tabrizi, K., Minge, M., Skjæveland, Å., Nikolaev, S. I., Jakobsen, K. S., &
677 Pawlowski, J. (2007). Phylogenomics reshuffles the eukaryotic supergroups. *PLoS ONE*,
678 2(8), e790. doi: 10.1371/journal.pone.0000790

679 Czech, L., Barbera, P., & Stamatakis, A. (2019). Genesis and Gappa: Processing, Analyzing
680 and Visualizing Phylogenetic (Placement) Data. *BioRxiv*, 647958. doi: 10.1101/647958

681 Dawson, S. C., & Pace, N. R. (2002). Novel kingdom-level eukaryotic diversity in anoxic
682 environments. *Proceedings of the National Academy of Sciences*, 99(12), 8324–8329.
683 doi: 10.1073/pnas.062169599

684 de Vargas, C., Audic, S., Henry, N., Decelle, J., Mahé, F., Logares, R., ... Romac, S. (2015).
685 Eukaryotic plankton diversity in the sunlit ocean. *Science*, 348(6237), 1261605. doi:
686 10.1007/s13398-014-0173-7.2

687 Diez, B., Pedros-Alio, C., Massana, R., Díez, B., Pedrós-Alió, C., & Massana, R. (2001). Study
688 of genetic diversity of eukaryotic picoplankton in different oceanic regions by small-
689 subunit rRNA gene cloning and sequencing. *Applied and Environmental Microbiology*,
690 67(7), 2932–2941. doi: 10.1128/AEM.67.7.2932-2941.2001

691 Dunthorn, M., Otto, J., Berger, S. A., Stamatakis, A., Mahé, F., Romac, S., ... Stoeck, T.
692 (2014). Placing Environmental Next-Generation Sequencing Amplicons from Microbial
693 Eukaryotes into a Phylogenetic Context. *Molecular Biology and Evolution*, 31(4), 993–
694 1009. doi: 10.1093/molbev/msu055

695 Edgar, R. C., Haas, B. J., Clemente, J. C., Quince, C., & Knight, R. (2011). UCHIME improves
696 sensitivity and speed of chimera detection. *Bioinformatics*, 27(16), 2194–2200. doi:
697 10.1093/bioinformatics/btr381

698 Edgcomb, V. P., Kysela, D. T., Teske, A., de Vera Gomez, A., & Sogin, M. L. (2002). Benthic
699 eukaryotic diversity in the Guaymas Basin hydrothermal vent environment. *Proceedings*
700 *of the National Academy of Sciences of the United States of America*, 99(11), 7658–
701 7662. doi: 10.1073/pnas.062186399

702 Foissner, & W. (1987). Soil Protozoa : fundamental problems, ecological significance.
703 adaptations in ciliates and tetaceans, bioindicators. and guide to the literature. *Prog.*
704 *Protistol.*, 2, 69–212. Retrieved from <https://ci.nii.ac.jp/naid/10019334967/>
705 Geisen, Stefan. (2016). Thorough high-throughput sequencing analyses unravels huge
706 diversities of soil parasitic protists. *Environmental Microbiology*, 18(6), 1669–1672.
707 doi: 10.1111/1462-2920.13309
708 Geisen, Stefan, Mitchell, E. A. D., Adl, S., Bonkowski, M., Dunthorn, M., Ekelund, F., ... Lara,
709 E. (2018). Soil protists: A fertile frontier in soil biology research. *FEMS Microbiology*
710 *Reviews*. doi: 10.1093/femsre/fuy006
711 Geisen, Stefan, Tveit, A. T., Clark, I. M., Richter, A., Svenning, M. M., Bonkowski, M., &
712 Urich, T. (2015). Metatranscriptomic census of active protists in soils. *The ISME*
713 *Journal*, 9(10), 2178–2190. doi: 10.1038/ismej.2015.30
714 Geisen, Stephen, Cornelia, B., Jörg, R., & Michael, B. (2014). Soil water availability strongly
715 alters the community composition of soil protists. *Pedobiologia*, 57(4–6), 205–213. doi:
716 10.1016/j.pedobi.2014.10.001
717 Gosling, P., van der Gast, C., & Bending, G. D. (2017). Converting highly productive arable
718 cropland in Europe to grassland: –a poor candidate for carbon sequestration. *Scientific*
719 *Reports*, 7(1), 10493. doi: 10.1038/s41598-017-11083-6
720 Heeger, F., Bourne, E. C., Baschien, C., Yurkov, A., Bunk, B., Spröer, C., ... Monaghan, M. T.
721 (2018). Long-read DNA metabarcoding of ribosomal RNA in the analysis of fungi from
722 aquatic environments. *Molecular Ecology Resources*, 18(6), 1500–1514. doi:
723 10.1111/1755-0998.12937
724 Heger, T. J., Giesbrecht, I. J. W., Gustavsen, J., del Campo, J., Kellogg, C. T. E., Hoffman, K.
725 M., ... Keeling, P. J. (2018). High-throughput environmental sequencing reveals high
726 diversity of litter and moss associated protist communities along a gradient of drainage
727 and tree productivity. *Environmental Microbiology*, 20(3), 1185–1203. doi:
728 10.1111/1462-2920.14061
729 Heiss, A. A., Kolisko, M., Ekelund, F., Brown, M. W., Roger, A. J., & Simpson, A. G. B. B.
730 (2018). Combined morphological and phylogenomic re-examination of malawimonads,
731 a critical taxon for inferring the evolutionary history of eukaryotes. *Royal Society Open*
732 *Science*, 5(4), 171707. doi: 10.1098/rsos.171707
733 Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7:
734 Improvements in performance and usability. *Molecular Biology and Evolution*, 30(4),
735 772–780. doi: 10.1093/molbev/mst010
736 Kozlov, A. M., Darriba, D., Flouri, T., Morel, B., & Stamatakis, A. (2018). RAxML-NG : A fast ,
737 scalable , and user-friendly tool for maximum likelihood phylogenetic inference.
738 *BioRxiv*. doi: 10.1101/447110
739 Lallias, D., Hiddink, J. G., Fonseca, V. G., Gaspar, J. M., Sung, W., Neill, S. P., ... Creer, S.
740 (2015). Environmental metabarcoding reveals heterogeneous drivers of microbial
741 eukaryote diversity in contrasting estuarine ecosystems. *The ISME Journal*, 9(5), 1208–
742 1221. doi: 10.1038/ismej.2014.213
743 Lax, G., Eglit, Y., Eme, L., Bertrand, E. M., Roger, A. J., & Simpson, A. G. B. (2018).
744 Hemimastigophora is a novel supra-kingdom-level lineage of eukaryotes. *Nature*,
745 564(7736), 410–414. doi: 10.1038/s41586-018-0708-8
746 Logares, R., Audic, S., Bass, D., Bittner, L., Boutte, C., Christen, R., ... Massana, R. (2014).
747 Patterns of rare and abundant marine microbial eukaryotes. *Current Biology*, 24(8),
748 813–821. doi: 10.1016/j.cub.2014.02.050
749 Lopez-Garcia, P., Philippe, H., Gail, F., & Moreira, D. (2003). Autochthonous eukaryotic
750 diversity in hydrothermal sediment and experimental microcolonizers at the Mid-
751 Atlantic Ridge. *Proceedings of the National Academy of Sciences*, 100(2), 697–702. doi:
752 10.1073/pnas.0235779100
753 López-García, P., Rodríguez-Valera, F., Pedrós-Alió, C., & Moreira, D. (2001). Unexpected
754 diversity of small eukaryotes in deep-sea Antarctic plankton. *Nature*, 409(6820), 603–

755 607. doi: 10.1038/35054537

756 Mahé, F., de Vargas, C., Bass, D., Czech, L., Stamatakis, A., Lara, E., ... Dunthorn, M. (2017).
757 Parasites dominate hyperdiverse soil protist communities in Neotropical rainforests.
758 *Nature Ecology & Evolution*, 1(4), 0091. doi: 10.1038/s41559-017-0091

759 Mahé, F., Mayor, J., Bunge, J., Chi, J., Siemensmeyer, T., Stoeck, T., ... Dunthorn, M. (2015).
760 Comparing high-throughput platforms for sequencing the V4 region of SSU-rDNA in
761 environmental microbial eukaryotic diversity surveys. *Journal of Eukaryotic*
762 *Microbiology*, 62(3), 338–345. doi: 10.1111/jeu.12187

763 Marande, W., López-García, P., & Moreira, D. (2009). Eukaryotic diversity and phylogeny
764 using small- and large-subunit ribosomal RNA genes from environmental samples.
765 *Environmental Microbiology*, 11(12), 3179–3188. doi: 10.1111/j.1462-
766 2920.2009.02023.x

767 Martijn, J., Lind, A. E., Schön, M. E., Spiertz, I., Juzokaite, L., Bunikis, I., ... Ettema, T. J. G.
768 (2019). Confident phylogenetic identification of uncultured prokaryotes through long
769 read amplicon sequencing of the 16S-ITS-23S rRNA operon. *Environmental*
770 *Microbiology*, 1462-2920.14636. doi: 10.1111/1462-2920.14636

771 Martijn, J., Lind, A. E., Spiers, I., Juzokaite, L., Bunikis, I., Pettersson, O. V., & Ettema, T. J. .
772 (2017). Amplicon sequencing of the 16S-ITS-23S rRNA operon with long-read
773 technology for improved phylogenetic classification of uncultured prokaryotes. *BioRxiv*,
774 234690. doi: 10.1101/234690

775 Massana, R., Balagué, V., Guillou, L., & Pedrós-Alió, C. (2004). Picoeukaryotic diversity in an
776 oligotrophic coastal site studied by molecular and culturing approaches. *FEMS*
777 *Microbiology Ecology*, 50(3), 231–243. doi: 10.1016/j.femsec.2004.07.001

778 Massana, R., Castresana, J., Balague, V., Guillou, L., Romari, K., Groisillier, A., ... Pedrós-
779 Alió, C. (2004). Phylogenetic and ecological analysis of novel marine stramenopiles.
780 *Applied and Environmental Microbiology*, 70(6), 3528–3534. doi:
781 10.1128/AEM.70.6.3528-3534.2004

782 Massana, R., Gobet, A., Audic, S., Bass, D., Bittner, L., Boutte, C., ... de Vargas, C. (2015).
783 Marine protist diversity in European coastal waters and sediments as revealed by high-
784 throughput sequencing. *Environmental Microbiology*, 17(10), 4035–4049. doi:
785 10.1111/1462-2920.12955

786 Matsen, F. A., Kodner, R. B., & Armbrust, E. V. (2010). pplacer: linear time maximum-
787 likelihood and Bayesian phylogenetic placement of sequences onto a fixed reference
788 tree. *BMC Bioinformatics*, 11(1), 538. doi: 10.1186/1471-2105-11-538

789 Miller, M. A., Pfeiffer, W., & Schwartz, T. (2010). Creating the CIPRES Science gateway.
790 *Proceedings of the Gateway Computing Environments Workshop (GCE)*, 1–7.
791 Retrieved from http://www.phylo.org/sub_sections/portal/sc2010_paper.pdf

792 Moon-Van Der Staay, S. Y., De Wachter, R., & Vaulot, D. (2001). Oceanic 18S rDNA
793 sequences from picoplankton reveal unsuspected eukaryotic diversity. *Nature*. doi:
794 10.1038/35054541

795 Moreira, D., von der Heyden, S., Bass, D., López-García, P., Chao, E., & Cavalier-Smith, T.
796 (2007). Global eukaryote phylogeny: Combined small- and large-subunit ribosomal
797 DNA trees support monophyly of Rhizaria, Retaria and Excavata. *Molecular*
798 *Phylogenetics and Evolution*, 44(1), 255–266. doi: 10.1016/j.ympev.2006.11.001

799 Mosher, J. J., Bowman, B., Bernberg, E. L., Shevchenko, O., Kan, J., Korlach, J., ... Kaplan, L.
800 A. (2014). Improved performance of the PacBio SMRT technology for 16S rDNA
801 sequencing. *Journal of Microbiological Methods*, 104, 59–60. doi:
802 10.1016/j.mimet.2014.06.012

803 Orr, R. J. S., Zhao, S., Klaveness, D., Yabuki, A., Ikeda, K., Watanabe, M. M., & Shalchian-
804 Tabrizi, K. (2018). Enigmatic Diphyllatea eukaryotes: culturing and targeted PacBio RS
805 amplicon sequencing reveals a higher order taxonomic diversity and global distribution.
806 *BMC Evolutionary Biology*, 18(1), 115. doi: 10.1186/s12862-018-1224-z

807 Pawlowski, J., Audic, S., Adl, S., Bass, D., Belbahri, L., Berney, C., ... de Vargas, C. (2012).

808 CBOL Protist Working Group: barcoding eukaryotic richness beyond the animal, plant,
809 and fungal kingdoms. *PLoS Biology*, 10(11), e1001419. doi:
810 10.1371/journal.pbio.1001419

811 Quast, C., Pruesse, E., Yilmaz, P., Gerken, J., Schweer, T., Yarza, P., ... Glöckner, F. O. (2013).
812 The SILVA ribosomal RNA gene database project: improved data processing and web-
813 based tools. *Nucleic Acids Research*, 41(D1), D590–D596. doi: 10.1093/nar/gks1219

814 Rognes, T., Flouri, T., Nichols, B., Quince, C., & Mahé, F. (2016). VSEARCH: a versatile open
815 source tool for metagenomics. *PeerJ*, 4, e2584. doi: 10.7717/peerj.2584

816 Schirmer, M., D'Amore, R., Ijaz, U. Z., Hall, N., & Quince, C. (2016). Illumina error profiles:
817 resolving fine-scale variation in metagenomic sequencing data. *BMC Bioinformatics*,
818 17(1), 125. doi: 10.1186/s12859-016-0976-y

819 Schloss, P. D., Jenior, M. L., Koumpouras, C. C., Westcott, S. L., & Highlander, S. K. (2016).
820 Sequencing 16S rRNA gene fragments using the PacBio SMRT DNA sequencing system.
821 *PeerJ*, 4, e1869. doi: 10.7717/peerj.1869

822 Schloss, P. D., Westcott, S. L., Ryabin, T., Hall, J. R., Hartmann, M., Hollister, E. B., ...
823 Weber, C. F. (2009). Introducing mothur: open-source, platform-independent,
824 community-supported software for describing and comparing microbial communities.
825 *Applied and Environmental Microbiology*, 75(23), 7537–7541. doi:
826 10.1128/AEM.01541-09

827 Schwelm, A., Berney, C., Dixelius, C., Bass, D., & Neuhauser, S. (2016). The large subunit
828 rDNA sequence of *Plasmodiophora brassicae* does not contain intra-species
829 polymorphism. *Protist*, 167(6), 544–554. doi: 10.1016/j.protis.2016.08.008

830 Simdyanov, T. G., Guillou, L., Diakin, A. Y., Mikhailov, K. V., Schrével, J., & Aleoshin, V. V.
831 (2017). A new view on the morphology and phylogeny of eugregarines suggested by the
832 evidence from the gregarine *Ancora sagittata* (Leuckart, 1860) Labbé, 1899
833 (Apicomplexa: Eugregarinida). *PeerJ*, 5, e3354. doi: 10.7717/peerj.3354

834 Simdyanov, T. G., Paskerova, G. G., Valigurová, A., Diakin, A., Kováčiková, M., Schrével, J., ...
835 Aleoshin, V. V. (2018). First Ultrastructural and Molecular Phylogenetic Evidence from
836 the Blastogregarines, an Early Branching Lineage of Plesiomorphic Apicomplexa.
837 *Protist*, 169(5), 697–726. doi: 10.1016/J.PROTIS.2018.04.006

838 Stamatakis, A. (2014). RAxML version 8: a tool for phylogenetic analysis and post-analysis of
839 large phylogenies. *Bioinformatics*, 30(9), 1312–1313. doi:
840 10.1093/bioinformatics/btu033

841 Stoeck, T., Bass, D., Nebel, M., Christen, R., Jones, M. D. M., Breiner, H. W., & Richards, T.
842 A. (2010). Multiple marker parallel tag environmental DNA sequencing reveals a highly
843 complex eukaryotic community in marine anoxic water. *Molecular Ecology*, 19(s1), 21–
844 31. doi: 10.1111/j.1365-294X.2009.04480.x

845 Stoeck, T., & Epstein, S. (2003). Novel Eukaryotic Lineages Inferred from Small-Subunit
846 rRNA Analyses of Oxygen-Depleted Marine Environments. *Applied and Environmental*
847 *Microbiology*, 69(5), 2657–2663. doi: 10.1128/AEM.69.5.2657-2663.2003

848 Stoeck, T., Taylor, G. T., & Epstein, S. S. (2003). Novel eukaryotes from the permanently
849 anoxic Cariaco Basin (Caribbean Sea). *Applied and Environmental Microbiology*, 69(9),
850 5656–5663. doi: 10.1128/AEM.69.9.5656-5663.2003

851 Tedersoo, L., & Anslan, S. (2019). Towards PacBio-based pan-eukaryote metabarcoding
852 using full-length ITS sequences. *Environmental Microbiology Reports*. doi:
853 10.1111/1758-2229.12776

854 Tedersoo, L., Tooming-Klunderud, A., & Anslan, S. (2018). PacBio metabarcoding of Fungi
855 and other eukaryotes: errors, biases and perspectives. *New Phytologist*, 217(3), 1370–
856 1385. doi: 10.1111/nph.14776

857 Vandenkoornhuyse, P., Baldauf, S. L., Leyval, C., Straczek, J., & Young, J. P. W. (2002).
858 Extensive fungal diversity in plant roots. *Science (New York, N.Y.)*, 295(5562), 2051.
859 doi: 10.1126/science.295.5562.2051

860 Wagner, J., Coupland, P., Browne, H. P., Lawley, T. D., Francis, S. C., & Parkhill, J. (2016).

861 Evaluation of PacBio sequencing for full-length bacterial 16S rRNA gene classification.
862 *BMC Microbiology*, 16(1), 274. doi: 10.1186/s12866-016-0891-4
863 Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences with
864 variable rates over sites: Approximate methods. *Journal of Molecular Evolution*, 39(3),
865 306–314. doi: 10.1007/BF00160154
866 Zerbino, D. R., Achuthan, P., Akanni, W., Amode, M. R., Barrell, D., Bhai, J., ... Flicek, P.
867 (2018). Ensembl 2018. *Nucleic Acids Research*, 46(D1), D754–D761. doi:
868 10.1093/nar/gkx1098
869 Zhao, S., Burki, F., Brate, J., Keeling, P. J., Klaveness, D., & Shalchian-Tabrizi, K. (2012).
870 Collodictyon--an ancient lineage in the tree of eukaryotes. *Molecular Biology and*
871 *Evolution*, 29(6), 1557–1568. doi: 10.1093/molbev/mss001
872
873

874 **Data accessibility.**

875 Raw PacBio Sequel reads have been submitted to the ENA database under accession number
876 PRJEB25197. Detailed software commands and custom scripts used in the read curation
877 pipeline are available on GitHub (<https://github.com/Pbdas/long-reads>).
878

879

880 **Author Contributions.**

881 FB and DB conceived the study. GB and SH provided soil samples from the UK and
882 performed DNA extraction. RF and DB performed the wet lab experiments. MJ performed
883 and/or coordinated most of the bioinformatic analyses, in close connection with PB, LC, AK,
884 and AS. FB, DB, and MJ wrote the first complete draft of the manuscript, and all authors
885 subsequently contributed to the final version.
886