

1 **MetaEuk – sensitive, high-throughput gene discovery and**  
2 **annotation for large-scale eukaryotic metagenomics**

3 Eli Levy Karin<sup>1\*</sup>, Milot Mirdita<sup>1</sup>, and Johannes Söding<sup>1\*</sup>

4

5

6 <sup>1</sup> Quantitative and Computational Biology, Max-Planck Institute for Biophysical Chemistry,  
7 37077 Göttingen, Germany.

8

9

10

11

12 \* To whom correspondence should be addressed:

13 Eli Levy Karin, Tel: +49 551 201-2881

14 E-mail: [eli.levy.karin@gmail.com](mailto:eli.levy.karin@gmail.com)

15 Johannes Söding, Tel: +49 551 201-2890

16 E-mail: [soeding@mpibpc.mpg.de](mailto:soeding@mpibpc.mpg.de)

17

18

19

20 Running title: gene discovery in eukaryotic metagenomics

21 Keywords: MetaEuk, eukaryotes, homology detection, prediction, annotation, contigs, marine;

## 22 **Abstract**

23 **Background:** Metagenomics is revolutionizing the study of microorganisms and their  
24 involvement in biological, biomedical, and geochemical processes, allowing us to investigate  
25 by direct sequencing a tremendous diversity of organisms without the need for prior cultivation.  
26 Unicellular eukaryotes play essential roles in most microbial communities as chief predators,  
27 decomposers, phototrophs, bacterial hosts, symbionts and parasites to plants and animals.  
28 Investigating their roles is therefore of great interest to ecology, biotechnology, human health,  
29 and evolution. However, the generally lower sequencing coverage, their more complex gene  
30 and genome architectures, and a lack of eukaryote-specific experimental and computational  
31 procedures have kept them on the sidelines of metagenomics.

32 **Results:** MetaEuk is a toolkit for high-throughput, reference-based discovery and annotation  
33 of protein-coding genes in eukaryotic metagenomic contigs. It performs fast searches with 6-  
34 frame-translated fragments covering all possible exons and optimally combines matches into  
35 multi-exon proteins. We used a benchmark of seven diverse, annotated genomes to show that  
36 MetaEuk is highly sensitive even under conditions of low sequence similarity to the reference  
37 database. To demonstrate MetaEuk's power to discover novel eukaryotic proteins in large-  
38 scale metagenomic data, we assembled contigs from 912 samples of the Tara Oceans project.  
39 MetaEuk predicted >12,000,000 protein-coding genes in eight days on ten 16-core servers.  
40 Most of the discovered proteins are highly diverged from known proteins and originate from  
41 very sparsely sampled eukaryotic supergroups.

42 **Conclusion:** The open-source (GPLv3) MetaEuk software  
43 (<https://github.com/soedinglab/metaeuk>) enables large-scale eukaryotic metagenomics  
44 through reference-based, sensitive taxonomic and functional annotation.

## 45 **Background**

46 Unicellular eukaryotes are present in almost all environments, including soil [1], oceans [2],  
47 and plant and animal-associated microbiomes [3,4]. They exhibit both autotrophic and  
48 heterotrophic lifestyles [5], exist in symbiosis with plants and animals [6], and interact with  
49 other microbial organisms [7]. They account for roughly half of the global primary productivity  
50 in the oceans, mostly by photosynthesis [8], are key contributors to the carbon and nitrogen  
51 cycles through carbon-dioxide fixation, organic matter degradation, and denitrification [9,10],  
52 and have been shown to be a source for chemically bioactive compounds [e.g., 11,12].

53 Since the advent of metabarcoding using 18S rRNA genes, the known evolutionary diversity  
54 of unicellular eukaryotes has increased by orders of magnitude [13], and novel phyla and  
55 supra-kingdoms are still being discovered [14,15]. Due to their vast diversity [16,17],  
56 unicellular eukaryotes are certain to hold invaluable secrets for biotechnology and  
57 biomedicine.

58 Protein-coding genes are major keys for understanding eukaryotic functions and activities [18].  
59 Metatranscriptomic and metagenomic studies provide unique means to reveal protein-coding  
60 genes. However, despite the great potential of studying uncultivable eukaryotes in their  
61 natural environment, they have received little attention in metatranscriptomic and  
62 metagenomic studies so far, with a few notable exceptions [e.g., 19,20]. The unique features  
63 of eukaryotic data, i.e., lower genomic coverage due to lower population densities in  
64 metagenomic samples, fewer reference genomes, increased genome sizes and higher  
65 complexity of gene structure negatively impact all stages of metagenomic analyses, from  
66 assembly, through binning, to protein prediction and annotation [as discussed by 21,22].

67 Specifically, identifying protein-coding genes in eukaryotes is inherently more challenging than  
68 in prokaryotes due to the exon-intron architecture of eukaryotic genes. To date, methods for  
69 eukaryotic gene calling [e.g., 23–25] consider two types of information when training models  
70 for gene prediction: intrinsic sequence signals (e.g., CpG islands) and extrinsic data, such as  
71 transcriptomics or an annotated genome from a closely-related organism. As splicing  
72 signatures are not well conserved throughout evolution, the predictive power of the trained  
73 models declines fast when applied to organisms that are phylogenetically distant from the  
74 organism on which the model was trained [26].

75 While these methods are very useful for genomics, their applicability to metagenomic data is  
76 severely limited. First, the transcriptomic or genomic data of annotated organisms that are  
77 sufficiently closely related are usually not available. Second, since the models need to be  
78 trained on a relatively narrow clade, the application of such methods to metagenomic data

79 requires to first bin the assembled contigs by their assumed genome of origin [as performed  
80 by 27], which is often quite inaccurate and slow, especially when the number of contigs is  
81 large, the coverage is low, the contigs are short, and the metagenomic data are species-rich  
82 [28–30]. Finally, model-training in itself is time consuming, taking hours to days per genomic  
83 bin [25,27], limiting this approach to the analysis of few genomic bins at a time.

84 Previously, methods that bypass or reduce the need to explicitly train models to detect protein-  
85 coding genes have been proposed in the context of genomics [e.g., 31,32]. These methods  
86 extract putative protein-coding fragments from the genome and join those that bear sequence  
87 similarity to available transcriptomic or protein sequence targets. Since the joined fragments  
88 can be separated by non-coding (intronic) regions, their match to the target is termed “spliced  
89 alignment”. Even at a genomic level, a brute force application of the spliced alignment  
90 approach poses a serious computational burden as it requires aligning each putative fragment  
91 to each target as well as recovering the set of putative fragments that best match a target.

92 Here, we developed MetaEuk, a novel and sensitive reference-based approach to identify  
93 single- and multi-exon protein-coding genes in eukaryotic metagenomic data. MetaEuk takes  
94 as input a set of assembled contigs and a reference database of target protein sequences or  
95 profiles. MetaEuk scans each contig in all six reading frames and extracts putative protein  
96 fragments between stop codons in each frame. Thus, MetaEuk makes no assumption about  
97 the splicing signal and does not rely on any preceding binning step. MetaEuk uses the  
98 MMseqs2 code library [33] for a very fast, yet sensitive identification of putative exons within  
99 the fragments. This step also discards the vast majority of fragments, which significantly  
100 reduces the computation time of all succeeding steps. The combinatorial task of considering  
101 all possible sets of putative exons to best match a given target is solved by means of dynamic  
102 programming. Since MetaEuk uses a homology-based strategy to identify protein-coding  
103 genes, it can directly confer annotations to its predictions from the matched target proteins.

104 We benchmarked MetaEuk by using annotated genomes and proteins of seven unicellular  
105 organisms from different parts of the eukaryotic tree of life under conditions of increasing  
106 evolutionary distance to sequences in the reference database. Despite its high speed and low  
107 false positive rates, MetaEuk is able to discover a large fraction of the known proteins in these  
108 benchmark genomes. We next applied MetaEuk to study marine eukaryotes. We assembled  
109 all Tara Oceans metagenomic samples [20] and focused on ~1,300,000 contigs of at least  
110 5kbp in length. We clustered more than 330,000,000 proteins to create a comprehensive  
111 catalog of over 87,000,000 protein profiles to serve as a reference database. We found the  
112 MetaEuk collection of >12,000,000 marine proteins is highly diverged, offering major  
113 eukaryotic lineage expansions.

## 114 **Results**

### 115 **The MetaEuk algorithm**

116 The main steps of the algorithm are presented schematically in Figure 1 and a detailed  
117 description is provided in the Methods section. For each input contig, all possible protein-  
118 coding fragments are translated in six reading frames and searched against a reference target  
119 database of protein sequences or profiles. Fragments from the same contig and strand that  
120 hit a reference target T are examined together. In each fragment, only the part that was aligned  
121 to the target protein T is considered as a putative exon. The putative exons are ordered  
122 according to their start position on the contig. Based on their contig locations and the locations  
123 of their aligned region on the target T, any two putative exons are either compatible or not. A  
124 dynamic programming procedure recovers the highest scoring path of compatible pairs of  
125 putative exons by computing the maximum scores of all paths ending with each putative exon.  
126 Since homologies among targets in the reference database can lead to multiple calls of the  
127 same protein-coding gene, redundancies are reduced by clustering the calls. To that end, all  
128 calls are ordered by their start position on the contig. The first call defines a new cluster and  
129 all calls that overlap it on the contig are assigned to its cluster if they share an exon with it.  
130 The next cluster is defined by the first unassigned call. After all calls are clustered, the best  
131 scoring call is selected as the representative of the cluster, termed a “prediction”. Finally, as  
132 overlaps of genes on the same strand are very rare [as reviewed by 34], gene predictions  
133 overlapping others on the same strand with a better E-value are removed.



156 protein alignment. We then computed the coverage of individual exons of the annotated  
157 proteins to which MetaEuk predictions were mapped. These mappings are fully described in  
158 the Methods section.

159

160 **Table 1 – Species used to benchmark MetaEuk.**

Species	Group	Genome size (Mbp)	# scaffolds	# annotated proteins	% multi-exon proteins	GC%	MetaEuk run time against UniRef90
<i>Schizosaccharomyces pombe</i>	Fungi	12.59	4	5,132	47%	36	35m
<i>Acanthamoeba castellanii</i> str. Neff	Amoebozoa	42.02	384	14,974	91%	57.8	59m
<i>Phytomonas</i> sp. isolate EM1	Excavata	17.78	138	6,381	0%	48	37m
<i>Babesia bigemina</i>	Alveolates	13.84	483	5,079	54%	50.6	35m
Nucleomorph of <i>Lotharella oceanica</i>	Rhizaria	0.68	4	668	39%	32.8	24m
<i>Phaeodactylum tricornutum</i>	Stramenopiles	27.45	88	10,408	46%	48.8	51m
<i>Aspergillus nidulans</i>	Eurotiomycetes	30.28	91	9,556	88%	50.3	52m

161

## 162 **Sensitivity at evolutionary distance**

163 Sequences from major eukaryotic clades, such as Rhizaria, Stramenopiles, and Dinoflagellata  
164 are poorly represented in public protein databases, despite their high abundance in the  
165 environment [17]. We therefore measured the ability of MetaEuk to identify homologous  
166 protein-coding genes in organisms, which have distant evolutionary relatives in the reference  
167 database, as would be the case in a typical metagenomic analysis. To that end, for each  
168 annotated organism, we considered five sets of MetaEuk predictions. The first is the base set,  
169 which consisted of all predictions. Since we worked with annotated species, their proteins are  
170 well represented in UniRef90. The base set therefore reflects ideal conditions, in which the  
171 queried organisms are close to the reference database. The other four sets reflect an  
172 increasing evolutionary distance and were generated by excluding MetaEuk gene calls whose  
173 Smith-Waterman alignment (computed using MMseqs2) to their UniRef90 target had more  
174 than 90%, 80%, 60% or 40% sequence identity. We measured sensitivity as the fraction of

175 annotated proteins from the query genome to which a MetaEuk prediction was mapped (see  
176 Methods). For all organisms, the sensitivity of the base set of predictions was at least 92%,  
177 and sensitivity decreased with the sequence identity threshold (Figure 2A). However, even at  
178 low thresholds (40% – 60%), a significant fraction of the annotated proteins were discovered.

179

### 180 **Annotated exon coverage**

181 We next assessed MetaEuk's performance at the level of individual exons. For each MetaEuk  
182 prediction from the base set and its mapped annotated protein, we computed the proportion  
183 of annotated exons that were covered by the prediction (see Methods). Overall, the majority  
184 of predictions covered the majority of exons and, as expected, the fraction of predictions that  
185 cover all annotated exons decreases with the number of exons in the annotated protein (Figure  
186 2B). For all organisms, most (77% – 91%) annotated exons were covered by MetaEuk  
187 predictions. In addition, we found that the fraction of multi-exon MetaEuk predictions was  
188 similar to that presented in Table 1 (average difference: 10%, Supp. Figure 1A) and that single-  
189 exon predictions tended to have longer exons than multi-exon predictions (Supp. Figure 1B).  
190 An additional measure of completeness of MetaEuk predictions is the coverage of the target  
191 UniRef90 protein based on which the prediction was made. We therefore aligned each  
192 predicted MetaEuk protein to its target and found that on average, > 83% of predictions  
193 covered > 90% of their target (Supp. Figure 2).

194

### 195 **Precision**

196 MetaEuk predictions that were mapped to annotated proteins were considered as true  
197 predictions. We first measured the precision of MetaEuk by using the NCBI annotations as  
198 gold standard and regarded all predictions in the base set that were not mapped to an  
199 annotated protein (8% – 35%, Supp. Figure 2) as false. We computed precision-recall curves  
200 by treating the predictions' E-values as a classifying score. We found good separation (AUC-  
201 PR > 0.7 in all cases) between predictions that mapped to annotated proteins and the rest  
202 (Figure 2C). However, a prediction that does not map to a known protein is not necessarily  
203 false as it might reflect an unannotated protein. We found that about 40% of the unmapped  
204 predictions overlap a protein-coding gene on the opposite strand or are on scaffolds that had  
205 no annotation at all (Supp. Figure 2), suggestive of post-hoc exclusion criteria in the NCBI  
206 annotation procedure. For this reason, we also measured the precision of MetaEuk  
207 independently of external annotations by using an inverted-sequence null model. For this  
208 annotation-free approach, we ran standard MetaEuk on the inverted sequences of the six



209 frame-translated putative fragments. Each prediction based on these inverted sequences can  
210 therefore be considered a false positive. We applied the same E-value cutoff for reporting  
211 predictions based on the original sequence data and based on the inverted set. For all  
212 organisms, the total number of false positive predictions produced by this approach was low  
213 (0 – 12), indicating very high precision (> 99.9%).

214

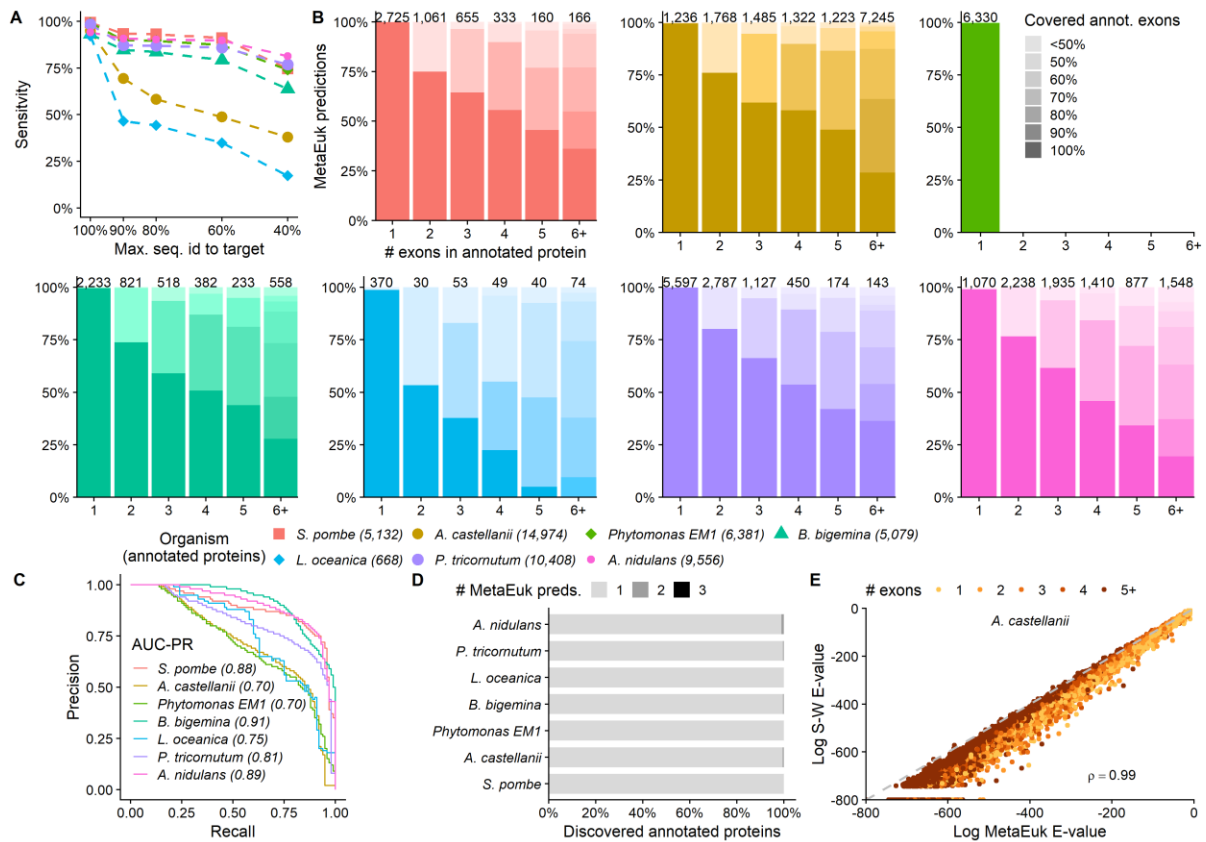
### 215 **Redundancy reduction**

216 MetaEuk's redundancy reduction procedure divides gene calls into disjoint clusters and retains  
217 a representative call as gene prediction for each cluster (see Methods). This reduces the  
218 number of potential protein-coding genes that need to be inspected. E.g., for *S. pombe*,  
219 MetaEuk produced over 1,100,000 calls that were reduced to a total of 5,564 predictions in  
220 the base set. A full reduction of redundancy is achieved when no two predictions correspond  
221 to same protein-coding gene. We thus identified cases in which two or more MetaEuk  
222 predictions were mapped to the same protein-coding gene. We found that for all benchmark  
223 organisms, redundancy is greatly reduced, as more than 99% of the annotated protein-coding  
224 genes in the benchmark scaffolds are only predicted once (Figure 2D).

225

### 226 **Statistical scores**

227 For each prediction, MetaEuk computes a bit-score between the set of translated and joined  
228 putative exons and the target protein. Based on this bit-score and the size of the reference  
229 database, an E-value is computed (see Methods). We evaluated MetaEuk's bit-scores and E-  
230 values by comparing them to those computed for each predicted protein and its target by the  
231 Smith-Waterman algorithm. Since MetaEuk penalizes missing and overlapping amino acids  
232 when joining putative exons, we expect the MetaEuk bit-score to be more conservative than  
233 the direct Smith-Waterman alignment bit-score. We found very high levels of agreement  
234 between the MetaEuk statistics and the Smith-Waterman statistics (Figure 2E, Supp. Figure  
235 3). This suggests a straightforward statistical interpretation of MetaEuk prediction scores.



236

237

238

239

240

241

242

243

244

245

246

247

### Effect of contig length

248

249

250

251

252

253

254

255

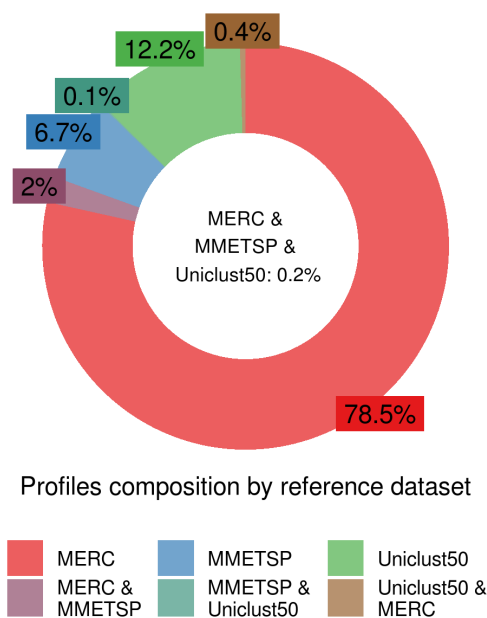
256

257

**Figure 2 – MetaEuk evaluation on benchmark.** MetaEuk predictions were mapped to annotated proteins. (A) Conditions of increasing evolutionary divergence were simulated by excluding gene calls based on their sequence identity to their target. Sensitivity is the fraction of annotated proteins from the query genome to which a MetaEuk prediction was mapped. (B) Fraction of exons covered by MetaEuk (color saturation). The number of MetaEuk predictions is indicated on top of each bar. (C) In an annotation-dependent precision estimation MetaEuk predictions that mapped to an annotated protein were considered as *true* and the rest as *false*. These sets of predictions are well separated by their E-values, as indicated by the high AUC-PR values. (D) Fraction of annotated protein-coding genes that were split by MetaEuk into two (dark grey) or three (black) different predictions. (E) Comparison of the E-values computed by MetaEuk and by the Smith-Waterman algorithm for *A. castellanii* proteins. The Spearman rho indicates high correlation for *A. castellanii* and the other organisms (Supp. Figure 3A).

## 258 Eukaryotic protein-coding genes in the ocean

259 To date, little is known about the biological activities of eukaryotes in the oceans [2,37]. We  
260 aimed to use MetaEuk to discover eukaryotic protein-coding genes in the Tara Oceans  
261 metagenomic dataset [20]. We first used MEGAHIT [38] to assemble all 912 samples of this  
262 project. We retained 1,351,204 contigs of at least 5kbp in length that were classified as  
263 potentially eukaryotic by EukRep [27]. We next constructed a comprehensive set of reference  
264 proteins by uniting over 21,000,000 representative sequences of the Uniclust50 database [39],  
265 the MERC dataset of over 292,000,000 protein sequence fragments assembled from  
266 eukaryotic Tara Oceans metatranscriptomic datasets [40], and over 18,500,000 protein  
267 sequences of MMETSP, the Marine Microbial Eukaryotic Transcriptome Sequencing Project  
268 [17,41]. We clustered the joint dataset of 331,913,793 proteins using the combined Linclust /  
269 MMseqs2 four-step cascaded clustering workflow [42] with a minimal sequence identity of  
270 20% and high sensitivity (-s 7). This resulted in 87,984,812 clusters, most of which (> 97%)  
271 contained proteins from a single reference dataset (Figure 3). For each cluster, a multiple  
272 sequence alignment was generated, based on which a sequence profile was computed.



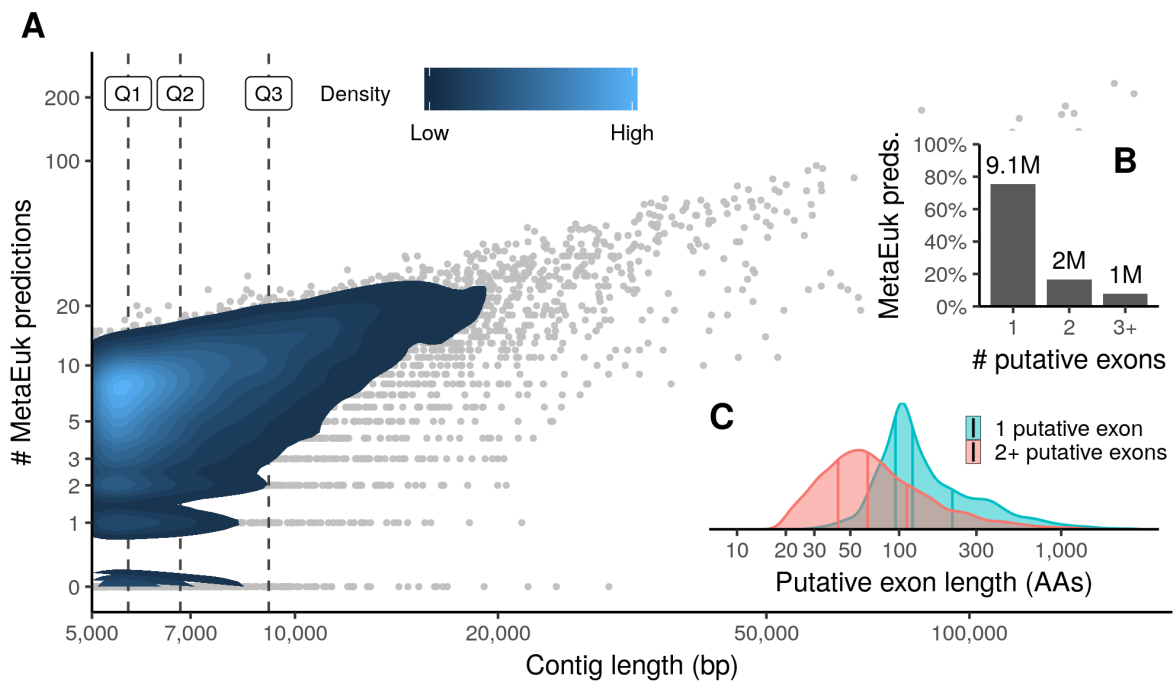
273

274 **Figure 3 – Reference profiles composition.** Proteins from three datasets: MERC (292 million), MMETSP (18.5 million) and  
275 Uniclust50 (21 million) were clustered into ~88 million clusters. Most clusters contained proteins from a single reference  
276 dataset. The profiles computed based on these clusters served as the reference database for the MetaEuk run on the Tara Oceans  
277 contigs.

278

279 MetaEuk's run using this reference database took eight days on ten 2x8-core servers and  
280 resulted in 12,111,301 predictions with no same-strand overlaps in 1,287,197 of the Tara  
281 Oceans contigs. Due to sequence similarities among the assembled contigs, some of these

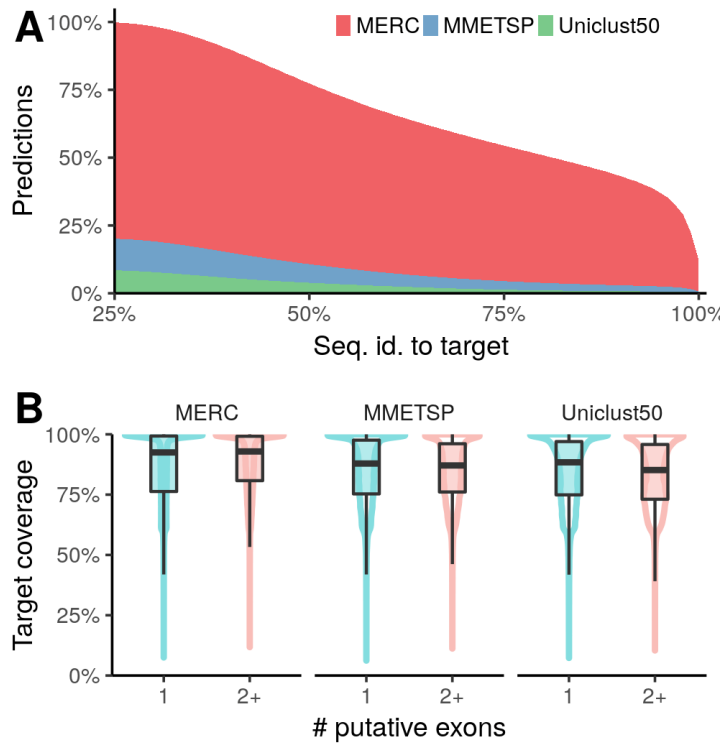
282 proteins are identical to each other, leaving a total of 6,158,526 unique proteins. We examined  
283 the distribution of predictions per contig, the number of putative exons in each prediction and  
284 the length of putative exons in single-exon and multi-exon predictions. We found that the  
285 number of predictions increases as a function of the contig length (Figure 4A), about 24% of  
286 predictions had more than one putative exon (Figure 4B) and multi-exon predictions tend to  
287 have shorter putative exons than single-exon predictions (4C). We analyzed the contribution  
288 of each reference dataset to the profiles based on which the MetaEuk predictions were made.  
289 MERC, MMETSP and Uniclust50 contributed 77.4%, 5.7% and 4.3% of the predictions,  
290 respectively. The rest of the predictions were based on mixed-dataset clusters (Supp. Figure  
291 5). We then used MMseqs2 to query the MetaEuk predicted proteins against their targets.  
292 Over 33% of the MetaEuk predictions have less than 60% sequence identity to their MERC,  
293 MMETSP or Uniclust50 target (Figure 5A). Finally, we found that 70% of the MetaEuk  
294 predicted proteins covered at least 80% of their reference target (Figure 5B).



295

296 **Figure 4 – MetaEuk predictions on Tara Oceans contigs.** MetaEuk was run on over 1.3 million contigs assembled from  
297 Tara Oceans metagenomic reads against a reference database of ~88 million protein profiles. (A) The number of MetaEuk  
298 predictions per contig increases with its length. Horizontal lines mark contig length quartiles. (B) Most (76%) MetaEuk  
299 predictions had a single putative exon. The absolute number of predictions is indicated above each bar. (C) Single-exon  
300 predictions tend to have longer putative exons than multi-exon predictions.

301

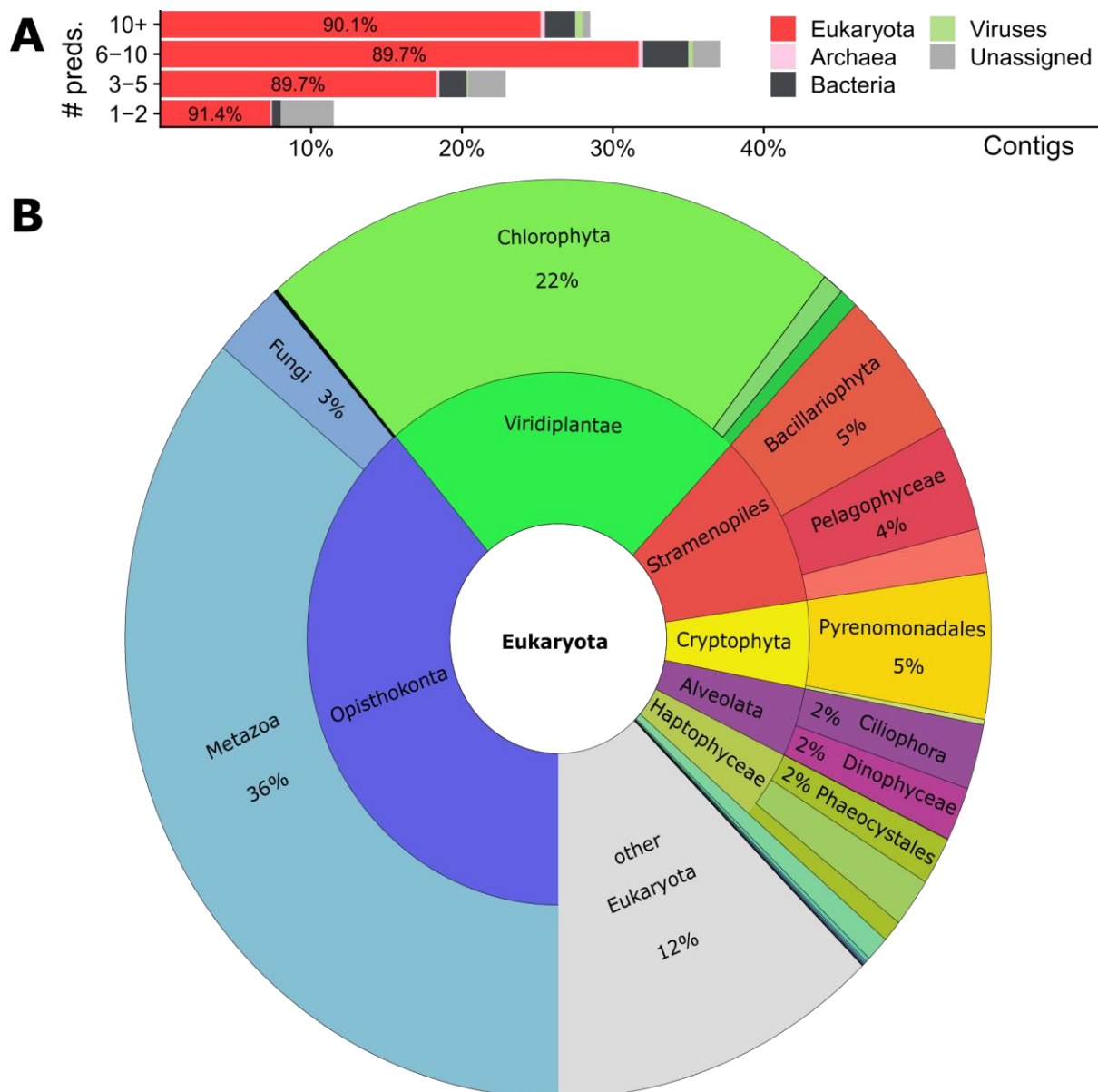


302

303 **Figure 5 – MetaEuk predictions compared to the reference datasets.** MetaEuk predicted proteins were queried against the  
304 representative sequence of their target reference cluster. (A) About one third of the predicted MetaEuk proteins had less than  
305 60% sequence identity to their target. (B) Targets are well covered by MetaEuk predicted proteins.

306

307 We next explored the taxonomic composition of the MetaEuk proteins. Since the majority  
308 (77%) of MetaEuk predictions were based on homologies to the MERC dataset, for which no  
309 taxonomic annotation is available, we queried the MetaEuk marine proteins collection against  
310 the Uniclust90 dataset [39] and the MMETSP dataset, both annotated using NCBI taxonomy  
311 (see Methods). We found that 63% of predictions based on homologies to the MERC dataset  
312 did not match any protein in either of the reference datasets, which means ~49% (63% of  
313 77%) of the MetaEuk marine proteins collection could not be assigned any taxonomy. This is  
314 in agreement with 52% of unassigned unigenes assembled from Tara Oceans  
315 metatranscriptomics [20]. We next assigned taxonomic labels to each assembled contig by  
316 conferring the taxonomic label with the best E-value of all MetaEuk predictions in the contig.  
317 This allowed us to annotate 92% of the contigs for which MetaEuk produced predictions (87%  
318 of all input contigs). We found that 82% of the contigs were assigned to the domain Eukaryota  
319 and 9% to non-eukaryotes, mostly bacteria (Figure 6A). We then examined the assigned  
320 eukaryotic supergroups below the domain level. About 12% of the eukaryotic contigs could  
321 not be assigned a supergroup. Among the most abundant eukaryotic supergroups are  
322 Metazoa and Chlorophyta (Figure 6B).



323

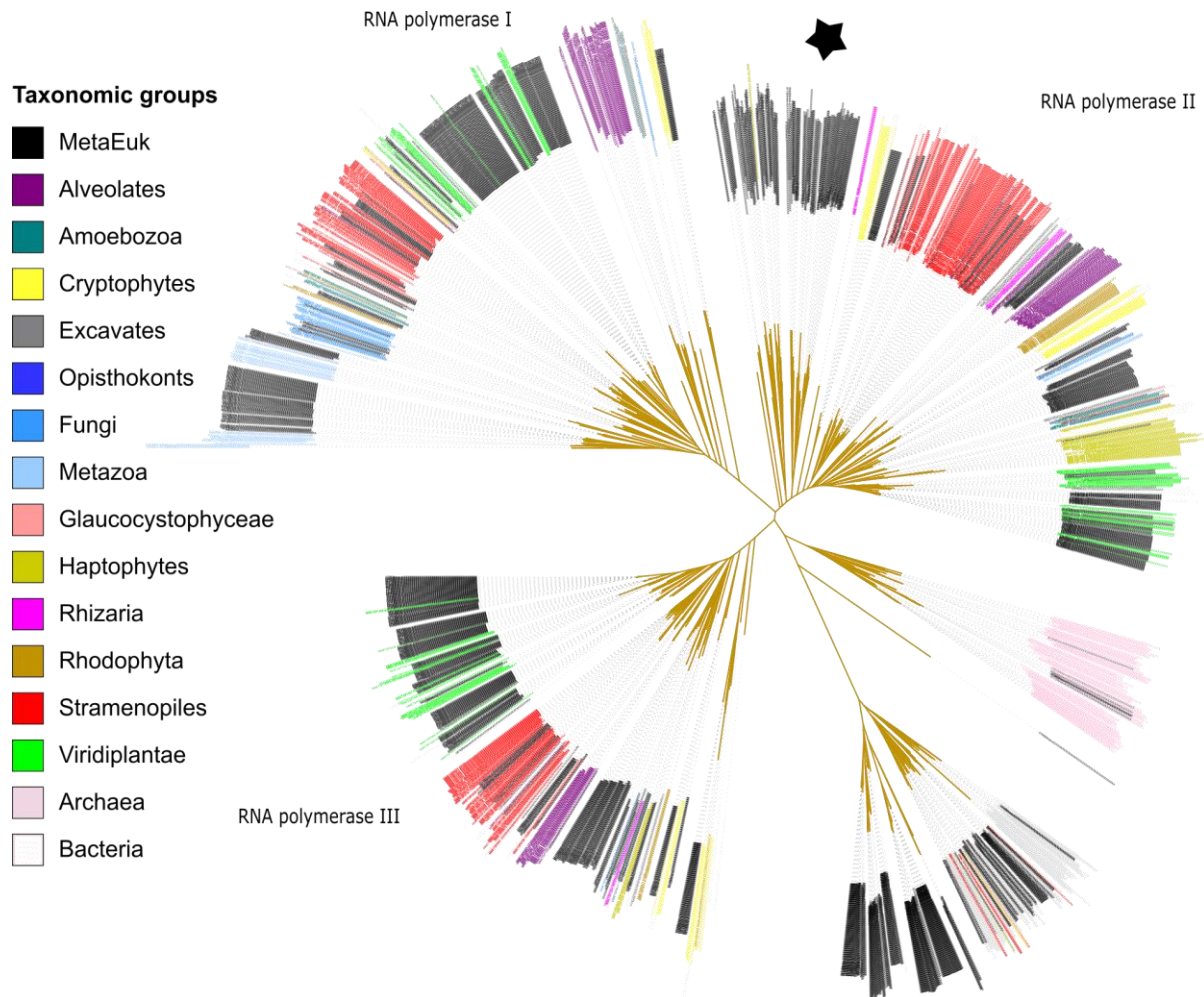
324 **Figure 6 – Taxonomy of Tara Oceans contigs with MetaEuk predictions.** The best-scoring taxonomic label of all  
 325 predictions on each contig was conferred to the contig. Contigs were divided into four categories according to their number of  
 326 MetaEuk predictions. Over 82% of the contigs were assigned to the domain Eukaryota. (A) The proportion of unassigned  
 327 contigs decreases with the number of MetaEuk predictions on the contig. The fraction of eukaryotic contigs out of all assigned  
 328 contigs is about 90% in all four categories. (B) Eukaryotic taxonomic labels below the domain level.

329

330 The high fraction of unassignable predictions (49%) prompted us to seek an additional way to  
 331 assess the diversity of the MetaEuk marine proteins. We thus collected orthologous  
 332 sequences of the large subunits of RNA polymerases, which are universal phylogenetic  
 333 markers [43] from 985 organisms for which we had taxonomic information, as well as 1,076  
 334 MetaEuk proteins, which consisted of all five Pfam domains of the large subunit in the right  
 335 order (see Methods). We aligned these sequences using MAFFT [44] and constructed the  
 336 maximum-likelihood phylogeny using RAxML [45]. The aim of this analysis was to delineate



337 the diversity of eukaryotic taxa of the MetaEuk marine proteins collection and not to resolve  
338 the exact phylogenetic relationships among them. As can be seen in Figure 7, MetaEuk  
339 proteins offer major lineage expansions in under-sampled eukaryotic supergroups.  
340 Importantly, the strict ortholog collection procedure performed for this analysis results in a  
341 conservative estimate of the diversity level of the MetaEuk marine proteins collection.  
342



343  
344 **Figure 7 – Diversity of MetaEuk marine eukaryotic proteins.** Homologous sequences of the large subunits of RNA  
345 polymerases of 985 species as well as 1,076 MetaEuk marine proteins were collected and a maximum-likelihood tree was  
346 computed based on their alignment. MetaEuk sequences (black) expand major eukaryotic lineages, including deeply rooted  
347 supergroups (denoted with star).

## 348 Discussion

349 We presented MetaEuk, an algorithm designed for large-scale analysis of eukaryotic  
350 metagenomic data. We demonstrated its utility for discovering proteins from highly diverged  
351 eukaryotic groups by analyzing assemblies of a huge set of 912 marine metagenomics  
352 samples. MetaEuk makes no assumption concerning splice site signatures and does not  
353 require a preceding binning procedure, which renders it suitable for the analysis of contigs  
354 from a mixture of highly diverged organisms. In order to achieve this, MetaEuk considers all  
355 possible putative protein-coding fragments from each input contig. Applying the spliced  
356 alignment dynamic programming procedure to recover the optimal set of putative exons  
357 directly on these fragments would result in a run time complexity per contig that is quadratic  
358 in the number of its fragments times the number of targets in the reference database. This is  
359 not feasible for metagenomics, as the number of fragments can be hundreds of millions (e.g.,  
360 from 1,351,204 Tara Oceans contigs, 152,519,258 fragments were extracted) and the  
361 reference database should be as comprehensive as possible (in this study, we used more  
362 than 87,000,000 protein profiles). To circumvent this limitation, MetaEuk takes advantage of  
363 the ultra-fast MMseqs2 search algorithm, which allows it to find putative exons matching a  
364 reference protein sequence with sufficient significance (in this study, a lenient E-value of 100).  
365 MetaEuk does not require significance at the exon level as it can combine sub-significant  
366 single exon matches to highly significant multi-exon matches. For example, two putative exons  
367 each with an E-value of 10 (corresponding to a bit-score of 25-40 in this study), are not  
368 individually significant but the sum of their bit-scores of at least 50 corresponds to a significant  
369 E-value of 1E-05.

370 MetaEuk is not designed to recover accurate splice sites, but rather to identify the protein-  
371 coding parts within exons. Indeed, we showed that MetaEuk predictions on the benchmark  
372 covered the majority (77% – 91%) of exons in annotated proteins. Since MetaEuk relies on  
373 local alignment at the amino acid level, it could potentially report pseudogenes, which still bear  
374 sequence similarity to reference proteins. However, we found that the majority of benchmark  
375 predictions (65% – 92%) mapped to NCBI annotated protein-coding genes, while the rest  
376 could be well separated from those that mapped by their E-values (AUC-PR > 0.7).  
377 Furthermore, unmapped predictions can reflect a missing annotation or post-hoc exclusion  
378 criteria (e.g., removal of annotations that overlap a better scoring one on the opposite strand).  
379 We therefore measured precision independently of annotations by running standard MetaEuk  
380 on the inverted sequences of the putative protein fragments extracted from the contigs. By  
381 using this annotation-free approach, we showed that MetaEuk's precision was greater than  
382 99.9% for all benchmark organisms. Put together, MetaEuk's strength is in describing the



383 protein-coding repertoire of versatile environments rather than in constructing statistical  
384 models of exon-intron transitions.

385 The Tara Oceans contigs analyzed in this study were assembled from Illumina HiSeq 2000  
386 short reads. High population diversity, repeat regions, and sequencing errors are among the  
387 major factors contributing to the computational challenge associated with metagenomic  
388 assembly [reviewed by 46]. These factors reduce the quality of the assembly as reflected, for  
389 example, in shorter contig lengths, chimeric contigs and contigs containing strand inversions.  
390 These in turn, directly and negatively impact MetaEuk. Shorter contigs limit its ability to  
391 discover multi-exon protein-coding genes as it searches for them within a contig. In addition,  
392 predictions on contig edges can be partial, which is more likely to happen in a highly  
393 fragmented assembly. By dividing each of the benchmark scaffolds to contigs whose lengths  
394 were drawn at random based on the length distribution of the Tara Ocean contigs, we showed  
395 that while MetaEuk retains its overall sensitivity to detect protein coding genes even under  
396 conditions of increasing evolutionary distance between the query organism and the target  
397 reference database, the completeness of its predictions is reduced. We thus expect MetaEuk  
398 to benefit from future improvements in assembly algorithms, higher sequencing coverage, and  
399 long-read sequencing technology [47–50].

400 In addition to developing MetaEuk, we generated two useful resources for the analysis of  
401 eukaryotes as part of this study. The first is the comprehensive protein profile database, which  
402 was computed using protein sequences from three sources: MERC, MMETSP and Uniclust50.  
403 With ~88 million records, it is the largest profile database focused on eukaryotes to date. Since  
404 MERC was assembled from the Tara Oceans metatranscriptomic data, we expected it to be  
405 a valuable resource for discovering protein-coding genes in the same environment. Indeed,  
406 we found that the majority of MetaEuk predictions (77%) were based on MERC protein  
407 profiles. Furthermore, the high fraction of MERC-based predictions that could not be assigned  
408 a taxonomic label (63%) demonstrates the uniqueness of this resource.

409 The second resource is the MetaEuk marine protein collection, which is available on our  
410 search webserver (<https://search.mmseqs.com/search>) for easy investigation [51]. Using a  
411 phylogenetic marker protein, we showed that this collection contains proteins spanning major  
412 eukaryotic lineages, including supergroups with very few available genomes. Over 33% of  
413 these proteins have less than 60% sequence identity to the representative reference proteins  
414 that were used to predict them, indicating their diversity with respect to the reference database.  
415 Unlike the MERC and MMETSP proteins, MetaEuk proteins are predicted in the context of  
416 genomic contigs. This allows us to learn of the number of putative exons that code for them  
417 as well as to examine them together with other proteins on the same contig. The latter is useful

418 for conferring taxonomic annotations to unlabeled predictions on the same contig as well as  
419 for detecting complex functional modules, by searching for co-occurrences of the module's  
420 proteins on the same contig.

421 As was demonstrated by the challenge of assigning taxonomy to highly diverged eukaryotic  
422 proteins, the paucity of eukaryotic sequences in reference databases is currently a major  
423 limitation in the study of eukaryotes. Thus, we expect the resources produced in this study  
424 and further analyses of eukaryotic metagenomic data using MetaEuk to produce a more  
425 comprehensive description of the tree of life [16,52–54].

426

## 427 **Conclusions**

428 MetaEuk is a sensitive reference-based algorithm for large-scale discovery of protein-coding  
429 genes in eukaryotic metagenomic data. Applying MetaEuk to large metagenomic datasets is  
430 expected to significantly enrich our databases with highly diverged eukaryotic protein-coding  
431 genes. By adding sequences from under-sampled eukaryotic lineages, we can improve  
432 sequence homology searches, protein profile computation and thereby homology-based  
433 function annotation, template-based and even de-novo protein structure prediction [55,56].  
434 These, in turn will allow for further exploration of eukaryotic activity in various environments  
435 [57].

## 436 **Methods**

### 437 **MetaEuk algorithm**

#### 438 **Code and resources availability**

439 The MetaEuk source code, compilation instructions and a brief user guide are available from  
440 <https://github.com/soedinglab/metaeuk> under the GNU General Public License v3.0. The  
441 resources produced during this study are available from  
442 <http://wwwuser.gwdg.de/~compbiol/metaeuk/>.

443

#### 444 **Putative exons compatibility**

445 In the first two stages of the MetaEuk algorithm all possibly coding protein fragments are  
446 translated from the input contigs. We scan each contig in six frames and extract the fragments  
447 between stop codons. These fragments are queried against the reference target database  
448 using MMseqs2. A set of fragments from the same contig and strand that have local matches  
449 to the same specific target  $T$  define a set of putative exons. We say two putative exons  $P_i$  and  
450  $P_j$  from the same set are compatible with each other if they can be joined together to a multi-  
451 exon protein.

452 Each  $P_i$  is associated with four coordinates: the amino-acid position on  $T$  from which the match  
453 to  $P_i$  starts ( $P_i^{ST}$ ) and ends ( $P_i^{ET}$ ); the nucleotide position on the contig from which the  
454 translation of  $P_i$  starts ( $P_i^{SC}$ ) and ends ( $P_i^{EC}$ ). We require a match of at least 10 amino acids (a  
455 minimal exon length). We consider putative exons  $P_i$  and  $P_j$  with  $P_i^{ST} < P_j^{ST}$  as compatible on  
456 the plus strand if:

- 457 (1) their order on the contig is the same as on the target:  $P_i^{SC} < P_j^{SC}$  ;  
458 (2) the distance between them on the contig is at least the length of a minimal intron but  
459 not more than the length of a maximal intron:  $15 \leq (P_j^{SC} - P_i^{EC}) \leq 10,000$ ;  
460 (3) their matches to  $T$  should not overlap. In practice we allow for a short overlap to  
461 account for alignment errors:  $(P_j^{ST} - P_i^{ET}) \geq -10$ .

462 In case  $P_i$  and  $P_j$  are on the negative strand, we modify conditions (1) and (2) accordingly:

- 463 (1)  $P_i^{SC} > P_j^{SC}$ ;  
464 (2)  $15 \leq (P_i^{EC} - P_j^{SC}) \leq 10,000$ .

465 Since the adjustment of conditions to the minus strand is straightforward, in the interest of  
466 brevity we focus solely on the plus strand in the following text.

467 We say a set of  $k > 1$  putative exons is compatible if, when ordered by their  $P_i^{ST}$  values, each  
468 pair of consecutive putative exons is compatible. (A set of a single exon is always compatible).

469

#### 470 **Bit-score and E-value computation**

471 A set of  $k$  compatible putative exons defines a pairwise protein alignment to the target  $T$ . This  
472 alignment is the concatenation of the ordered local alignments of all putative exons to  $T$ .  
473 Between each consecutive putative pair of exons  $P_i$  and  $P_{i+1}$  there might be unmatched amino  
474 acids in  $T$  or there might be a short overlap of their matches to  $T$ . We denote the number of  
475 unmatched amino acids between  $P_i$  and  $P_{i+1}$  as  $l_i$ , which can take a negative value in case of  
476 an overlap. MetaEuk computes the bit-score of the concatenated pairwise alignment  $S(P_{set}, T)$   
477 by summing the individual Karlin-Altschul [58] bit-scores  $S(P_i, T)$  of the putative exons to  $T$  and  
478 penalizing for unmatched or overlapping amino acids in  $T$  as follows:

$$479 \quad S(P_{set}, T) = \sum_{i=1}^k S(P_i, T) + \sum_{i=1}^{k-1} C(l_i) + \log_2(k!)$$

480 where the penalty function is  $C(l_i) = -|l_i|$  for  $l_i \neq 1$  and 0 if  $l_i = 1$ . The last term rewards the  
481 correct ordering of the  $k$  exons.

482 An E-value is the expected number of matches above a given bit-score threshold. Since for  
483 each contig, at most one gene call is reported per strand and target in the reference database,  
484 the E-value takes into account the number of amino acids in the reference database  $D$  and  
485 the search on two strands:

$$486 \quad E - Value(P_{set}, T) = 2 \times D \times 2^{-S(P_{set}, T)}$$

487

#### 488 **Dynamic programming**

489 Given a set of  $n$  putative exons and their target, MetaEuk finds the set of compatible exons  
490 with the highest combined bit-score. First, all putative exons are sorted by their start on the  
491 contig, such that  $P_1^{SC} \leq \dots \leq P_n^{SC}$ . The dynamic programming computation iteratively computes  
492 vectors  $S$ ,  $k$ , and  $b$  from their first entry 1 to their  $n^{th}$ . The entry  $S_i$  holds the score of the best  
493 exon alignment ending in exon  $i$  and  $k_i$  holds the number of exons in that set. Once the

494 maximum score is found, the exon alignment is back traced using  $b$ , in which entry  $b_i$  holds  
495 the index of the aligned exon preceding exon  $i$  (0 if  $i$  is the first aligned exon). Using the  
496 following values:

$$497 \quad S_0 = 0; k_0 = 0; b_0 = 0$$

498 all putative exons  $P_j$  are examined according to their order and the score vector is updated:

$$499 \quad S_j = \max_i (S_i + S(P_j, T) + C(l_j^i) + \log_2(k_i + 1) | 0 \leq i < j, i \text{ compatible with } j)$$

500  $k_j$  and  $b_j$  are updated accordingly. The terms  $\log_2(k_i + 1)$  add up to the score contribution  
501  $\sum_{i=1}^k \log_2(i) = \log_2(k!)$  and the transition 0 to  $j$  is defined as compatible with  $C(l_j^0) = 0$  for all  
502  $j$ . The optimal exon set is then recovered by tracing back from the exon with the maximal  
503 score. This dynamic programming procedure has time complexity of  $O(n^2)$ .

504

### 505 **Clustering gene calls to reduce redundancy**

506 MetaEuk assigns a unique identifier to each extracted putative protein fragment (stage 1 in  
507 Figure 1). A MetaEuk exon refers to the part of a fragment that matched some target  $T$  (stage  
508 2 in Figure 1, tinted background) and has the same identifier as the fragment. Two calls that  
509 have the same exon identifier in their exon set are said to share an exon. MetaEuk reduces  
510 redundancy by clustering calls that share an exon (stage 4 in Figure 1) and selecting a  
511 representative call as the gene prediction of each cluster. To that end, all  $N$  MetaEuk calls  
512 from the same contig and strand combination are ordered according to the contig start position  
513 of their first exon. Since this order can include equalities, they are sub-ordered by decreasing  
514 number of exons. The first cluster is defined by the first call, which serves as its tentative  
515 representative. Let  $m$  be the last contig position of the last exon of this representative. Each  
516 of the following calls is examined so long as its start position is smaller than  $m$  (i.e., it overlaps  
517 the representative on the contig). If that call shares an exon with the representative, it is  
518 assigned to its cluster. In the next iteration, the first unassigned call serves as representative  
519 for a new cluster and the following calls are examined in a similar manner, adding unassigned  
520 calls to the cluster in case they share an exon with the representative. The clustering ends  
521 with the assignment of all calls. At this stage, the final prediction is the call with the highest  
522 score in each cluster. This greedy approach has time complexity of  $O(N \times \log(N) + N \times A)$ ,  
523 where  $A$  is the average number of calls that overlap each representative on the contig. Since  
524 in practice,  $A \ll N$ , the expected time complexity is  $O(N \times \log(N))$ .

525

## 526 **Resolving same-strand overlapping predictions**

527 After the redundancy reduction step, MetaEuk sorts all predictions on the same contig and  
528 strand according to their E-value. It examines the sorted list and retains predictions only if they  
529 do not overlap any preceding predictions on the list.

530

## 531 **Benchmark datasets**

532 The UniRef90 database was obtained in March 2018. The annotated information of  
533 *Schizosaccharomyces pombe* (GCA\_000002945.2), *Acanthamoeba castellanii str. Neff*  
534 (GCA\_000313135.1), *Babesia bigemina* (GCA\_000981445.1), *Phytomonas sp. isolate EM1*  
535 (GCA\_000582765.1), Nucleomorph of *Lotharella oceanica* (GCA\_000698435.2),  
536 *Phaeodactylum tricornutum* (GCA\_000150955.2), and *Aspergillus nidulans*  
537 (GCA\_000149205.2) were downloaded from the NCBI genome assembly database (March –  
538 September 2018). This information included the genomic scaffolds, annotated protein  
539 sequences, and GFF3 files containing information about the locations of annotated proteins  
540 and other genomic elements. MetaEuk (Github commit  
541 47141068c171fcdd3d93411ac50958da0f2c4025, MMseqs2 submodule version  
542 ebb16f3631d320680a306c03aa7412c572f72ee3) was run with the following parameters: -e  
543 100 (a lenient maximal E-value of a putative exon against a target protein), --metaeuk-eval  
544 0.0001 (a stricter maximal cutoff for the MetaEuk E-value after joining exons into a gene call),  
545 --metaeuk-tcov 0.6 (a minimal cutoff for the ratio between the MetaEuk protein and the target)  
546 and --min-length 20, requiring putative exon fragments of at least 20 codons and default  
547 MMseqs2 search parameters.

548

## 549 **Mapping benchmark predictions to annotated proteins**

550 For each annotated protein, we listed the contig start and end coordinates of the coding part  
551 (CDS) of each of its exons. The lowest and highest of these coordinates were considered as  
552 the boundaries of the annotated protein, and the stretch between them as its “global” contig  
553 length. Similarly, we listed these coordinates and computed the boundaries and global contig  
554 length for each MetaEuk prediction. A MetaEuk prediction was globally mapped to an  
555 annotated protein if the overlap computed based on their boundaries was at least 80% of the  
556 global contig length of either of them and if, in addition, the alignment of their protein

557 sequences mainly consisted of identical amino acids or gaps (i.e., less than 10% mismatches).  
558 These criteria allow mapping MetaEuk predictions to an annotated protein, even if they miss  
559 some of its exons. Next, we computed the exon level mapping for all globally mapped pairs of  
560 MetaEuk predictions and annotated proteins. To that end, we compared their lists of exon  
561 contig coordinates. If an exon predicted by MetaEuk covered at least 80% of the contig length  
562 of an annotated protein's exon, we considered the annotated exon as "covered" by the  
563 MetaEuk prediction.

564

### 565 **Generating typical metagenomic contig lengths**

566 In order to evaluate MetaEuk's performance on contigs with a length distribution typical for  
567 assemblies from metagenomic samples, we recorded the lengths of the assembled contigs  
568 used for the analysis described in the "Tara Oceans dataset" section. The 1,351,204 contigs  
569 had a minimal length of 5,002 bps, 1<sup>st</sup> quartile of 5,661 bps, median of 6,763 bps, 3<sup>rd</sup> quartile  
570 of 9,020 bps and a maximal length of 1,524,677 bps. We divided each annotated scaffold into  
571 contigs of lengths that were randomly sampled from these recorded lengths. This resulted in  
572 1,392, 5,061, 1,816, 2,095, 80, 3,153 and 3,273 contigs for *S. pombe*, *A. castellanii*,  
573 *Phytomonas sp. isolate EM1*, nucleomorph of *L. oceanica*, *P. tricornutum*, and *A. nidulans*,  
574 respectively. MetaEuk was run on these contigs in the same way as on the original scaffolds.  
575 Since each of the new contigs corresponded to specific locations on the original scaffolds, we  
576 could carry out all benchmark assessments, which relied on mapping between MetaEuk  
577 predictions and annotated proteins.

578

### 579 **Tara Oceans dataset**

580 The 912 metagenomic SRA experiments associated with accession number PRJEB4352 were  
581 downloaded from the SRA (August – September 2018). The reads of each experiment were  
582 trimmed to remove adapters and low quality sequences using trimmomatic-0.38 [59] with  
583 parameters ILLUMINACLIP:TruSeq3-PE.fa:2:30:10 LEADING:3 TRAILING:3  
584 SLIDINGWINDOW:4:15 MINLEN:36 (SE for single-end samples). The resulting reads were  
585 then assembled with MEGAHIT [38] with default parameters. Contigs of at least 5kbp in length  
586 were classified as eukaryotic/non-eukaryotic using EukRep [27], which is trained to be highly  
587 sensitive to detecting eukaryotic contigs. MetaEuk was run on the contigs classified as  
588 eukaryotic with parameters: -e 100, --metaeuk-eval 0.0001, --min-ungapped-score 35, --min-



589 exon-aa 20, --metaeuk-tcov 0.6, --min-length 40, --slice-search (profile mode) and default  
590 MMseqs2 search parameters.

591

## 592 **Taxonomic assignment to predictions and contigs**

593 We used MMseqs2 to query the MetaEuk marine proteins collection against two taxonomically  
594 annotated datasets: Uniclust90 and the MMETSP protein dataset. Taxonomic labels  
595 associated with each of the MMETSP identifiers were downloaded from the NCBI website  
596 (BioProject PRJNA231566). We retained the hit with the highest bit-score value for each  
597 prediction if it had an E-value smaller than 1E-05. In addition, we examined the sequence  
598 identity between the MetaEuk prediction and the target in order to determine the rank of the  
599 taxonomic assignment. Similarly to [20], we used the following sequence identity cutoffs:  
600 >95% (species), >80% (genus), >65% (family), >50% (order), >40% (class), >30% (phylum),  
601 >20% (kingdom). Lower values were assigned at the domain level. The predictions on each  
602 contig were examined and the best-scoring one was used to confer taxonomic annotation to  
603 that contig. The assignment was visualized using Krona [60].

604

## 605 **Phylogenetic tree reconstruction**

606 We constructed the tree using the large subunit of RNA polymerases as a universal marker.  
607 This subunit contains five RNA\_pol\_Rpb domains (Pfam IDs: pf04997, pf00623, pf04983,  
608 pf05000, pf04998). As detailed below, protein sequences that contained all five domains in  
609 the right order were obtained in January-November 2019 from six sources to construct the  
610 multiple sequence alignment and tree. The sources were: (1) 75 sequences of the OrthoMCL  
611 [61] group OG5\_127924. The four-letter taxonomic codes of these sequences were converted  
612 to NCBI scientific names, based on information from the OrthoMCL website  
613 (<http://orthomcl.org/orthomcl/getDataSummary.do>). (2) 36 reviewed eukaryotic sequences  
614 were downloaded from UniProt [36]. These were used to distinguish between eukaryotic RNA  
615 Polymerase I (8 sequences), eukaryotic RNA Polymerase II (16 sequences) and eukaryotic  
616 RNA Polymerase III (12 sequences). We then ran an MMseqs2 profile search against the  
617 Pfam database (with parameters: -k 5, -s 7) with several query sets and retained results in  
618 which all five domains were matched in the right order with a maximal E-value of 0.0001. This  
619 allowed us to add the following sources: (3) 674 MMETSP proteins. (4) 100 archaeal proteins;  
620 (5) 100 bacterial proteins. For datasets (4) and (5), we first downloaded candidate proteins  
621 from the UniProt database by searching for the five domains and restricting taxonomy:archaea  
622 (bacteria). We then ran the previously described search procedure and randomly sampled



623 exactly 100 proteins from each group that matched the criterion. (6) 1,076 MetaEuk  
624 predictions. The joint set of 2,061 sequences was aligned using MAFFT v7.407 [44] and a  
625 phylogenetic tree was reconstructed by running RAxML v8 [45]. Tree visualization was  
626 performed in iTOL [62].

627 **Declarations**

628 **Ethics approval and consent to participate**

629 Not applicable

630

631 **Consent for publication**

632 Not applicable

633

634 **Availability of data and material**

635 The datasets generated and/or analyzed during the current study are available in  
636 <http://wwwuser.gwdg.de/~compbiol/metaeuk/>

637

638 **Competing interests**

639 The authors declare that they have no competing interests.

640

641 **Funding**

642 ELK is a recipient of a FEBS long-term fellowship and is an EMBO non-stipendiary long-term  
643 fellow. This work was supported by the EU's Horizon 2020 Framework Programme (Virus-X,  
644 grant 685778).

645

646 **Authors' contributions**

647 ELK and JS have designed the MetaEuk algorithm, benchmark and biological application. ELK  
648 and MM have developed the algorithm. ELK has analyzed the benchmark and Tara Oceans  
649 data. ELK and MM have generated the figures. ELK, JS and MM have drafted the manuscript.

650

651 **Acknowledgements**

652 We thank Dr. David Burstein from Tel Aviv University for his helpful insights concerning the  
653 phylogenetic analyses and Dr. Christian Woehle from Christian-Albrechts University Kiel for  
654 his comments on the manuscript.

## 655 **References**

- 656 1. Lentendu G, Hübschmann T, Müller S, Dunker S, Buscot F, Wilhelm C. Recovery of soil  
657 unicellular eukaryotes: an efficiency and activity analysis on the single cell level. *J Microbiol*  
658 *Methods*. 2013;95:463–9.
- 659 2. Keeling PJ, Campo J del. Marine protists are not just big bacteria. *Curr Biol*.  
660 2017;27:R541–9.
- 661 3. Parfrey LW, Walters WA, Knight R. Microbial eukaryotes in the human microbiome:  
662 ecology, evolution, and future directions. *Front Microbiol*. 2011;2:153.
- 663 4. Parfrey LW, Walters WA, Lauber CL, Clemente JC, Berg-Lyons D, Teiling C, et al.  
664 Communities of microbial eukaryotes in the mammalian gut within the context of  
665 environmental eukaryotic diversity. *Front Microbiol*. 2014;5.
- 666 5. Massana R. Eukaryotic picoplankton in surface oceans. *Annu Rev Microbiol*. 2011;65:91–  
667 110.
- 668 6. Flórez L V, Biedermann PHW, Engl T, Kaltenpoth M. Defensive symbioses of animals  
669 with prokaryotic and eukaryotic microorganisms. *Nat Prod Rep*. 2015;32:904–36.
- 670 7. Douglas AE. Symbiosis as a general principle in eukaryotic evolution. *Cold Spring Harb*  
671 *Perspect Biol*. 2014;6:a016113.
- 672 8. Field, Behrenfeld, Randerson, Falkowski. Primary production of the biosphere: integrating  
673 terrestrial and oceanic components. *Science*. 1998;281:237–40.
- 674 9. Jardillier L, Zubkov M V, Pearman J, Scanlan DJ. Significant CO<sub>2</sub> fixation by small  
675 prymnesiophytes in the subtropical and tropical northeast Atlantic Ocean. *ISME J*.  
676 2010;4:1180–92.
- 677 10. Woehle C, Roy A-S, Glock N, Wein T, Weissenbach J, Rosenstiel P, et al. A novel  
678 eukaryotic denitrification pathway in Foraminifera. *Curr Biol*. 2018;28:2536–2543.e5.
- 679 11. Michalak I, Chojnacka K. Algae as production systems of bioactive compounds. *Eng Life*  
680 *Sci*. 2015;15:160–76.
- 681 12. Falaise C, François C, Travers M-A, Morga B, Haure J, Tremblay R, et al. Antimicrobial  
682 compounds from eukaryotic microalgae against human pathogens and diseases in  
683 aquaculture. *Mar Drugs*. 2016;14:159.
- 684 13. Leray M, Knowlton N. DNA barcoding and metabarcoding of standardized samples

- 685 reveal patterns of marine benthic diversity. *Proc Natl Acad Sci*. 2015;112:2076–81.
- 686 14. Pawlowski J. The new micro-kingdoms of eukaryotes. *BMC Biol*. 2013;11:40.
- 687 15. Lax G, Eglit Y, Eme L, Bertrand EM, Roger AJ, Simpson AGB. Hemimastigophora is a  
688 novel supra-kingdom-level lineage of eukaryotes. *Nature*. 2018;564:410–4.
- 689 16. Burki F. The eukaryotic tree of life from a global phylogenomic perspective. *Cold Spring*  
690 *Harb Perspect Biol*. 2014;6:a016147.
- 691 17. Keeling PJ, Burki F, Wilcox HM, Allam B, Allen EE, Amaral-Zettler LA, et al. The Marine  
692 Microbial Eukaryote Transcriptome Sequencing Project (MMETSP): Illuminating the  
693 functional diversity of eukaryotic life in the oceans through transcriptome sequencing. *PLoS*  
694 *Biol*. 2014;12:e1001889.
- 695 18. Bik HM, Porazinska DL, Creer S, Caporaso JG, Knight R, Thomas WK. Sequencing our  
696 way towards understanding global eukaryotic biodiversity. *Trends Ecol Evol*. 2012;27:233–  
697 43.
- 698 19. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al. Structure  
699 and function of the global ocean microbiome. *Science*. 2015;348:1261359–1261359.
- 700 20. Carradec Q, Pelletier E, Da Silva C, Alberti A, Seeleuthner Y, Blanc-Mathieu R, et al. A  
701 global ocean atlas of eukaryotic genes. *Nat Commun*. 2018;9:373.
- 702 21. Escobar-Zepeda A, Vera-Ponce de León A, Sanchez-Flores A. The Road to  
703 Metagenomics: From Microbiology to DNA Sequencing Technologies and Bioinformatics.  
704 *Front Genet*. 2015;6:348.
- 705 22. Majaneva M, Hyytiäinen K, Varvio SL, Nagai S, Blomster J, Wachter R De. Bioinformatic  
706 amplicon read processing strategies strongly affect eukaryotic diversity and the taxonomic  
707 composition of communities. *PLoS One*. 2015;10:e0130035.
- 708 23. Stanke M, Morgenstern B. AUGUSTUS: a web server for gene prediction in eukaryotes  
709 that allows user-defined constraints. *Nucleic Acids Res*. 2005;33:W465–7.
- 710 24. Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. BRAKER1: unsupervised  
711 RNA-seq-based genome annotation with GeneMark-ET and AUGUSTUS. *Bioinformatics*.  
712 2016;32:767–9.
- 713 25. Hoff KJ, Stanke M. WebAUGUSTUS--a web service for training AUGUSTUS and  
714 predicting genes in eukaryotes. *Nucleic Acids Res*. 2013;41:W123–8.

- 715 26. Korf I. Gene finding in novel genomes. *BMC Bioinformatics*. 2004;5:59.
- 716 27. West PT, Probst AJ, Grigoriev I V, Thomas BC, Banfield JF. Genome-reconstruction for  
717 eukaryotes from complex natural microbial communities. *Genome Res*. 2018;28:569–80.
- 718 28. Sedlar K, Kupkova K, Provaznik I. Bioinformatics strategies for taxonomy independent  
719 binning and visualization of sequences in shotgun metagenomics. *Comput Struct Biotechnol*  
720 *J*. 2017;15:48–55.
- 721 29. Lu YY, Chen T, Fuhrman JA, Sun F, Sahinalp C. COCACOLA: Binning metagenomic  
722 contigs using sequence COmposition, read CoverAge, CO-alignment and paired-end read  
723 LinkAge. *Bioinformatics*. 2017;33:791–8.
- 724 30. Yu G, Jiang Y, Wang J, Zhang H, Luo H. BMC3C: binning metagenomic contigs using  
725 codon usage, sequence composition and read coverage. *Bioinformatics*. 2018;34:4172–9.
- 726 31. Gelfand MS, Mironov AA, Pevzner PA. Gene recognition via spliced sequence  
727 alignment. *Proc Natl Acad Sci U S A*. 1996;93:9061–6.
- 728 32. Gotoh O. Direct mapping and alignment of protein sequences onto genomic sequence.  
729 *Bioinformatics*. 2008;24:2438–44.
- 730 33. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the  
731 analysis of massive data sets. *Nat Biotechnol*. 2017;35:1026–8.
- 732 34. Kumar A. An overview of nested genes in eukaryotic genomes. *Eukaryot Cell*.  
733 2009;8:1321–9.
- 734 35. Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Ostell J, Pruitt KD, et al.  
735 GenBank. *Nucleic Acids Res*. 2018;46:D41–7.
- 736 36. Bateman A, Martin MJ, O'Donovan C, Magrane M, Alpi E, Antunes R, et al. UniProt: the  
737 universal protein knowledgebase. *Nucleic Acids Res*. 2017;45:D158–69.
- 738 37. Delmont TO, Quince C, Shaiber A, Esen ÖC, Lee ST, Rappé MS, et al. Nitrogen-fixing  
739 populations of Planctomycetes and Proteobacteria are abundant in surface ocean  
740 metagenomes. *Nat Microbiol*. 2018;3:804–13.
- 741 38. Li D, Liu C-M, Luo R, Sadakane K, Lam T-W. MEGAHIT: an ultra-fast single-node  
742 solution for large and complex metagenomics assembly via succinct de Bruijn graph.  
743 *Bioinformatics*. 2015;31:1674–6.
- 744 39. Mirdita M, von den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M. Uniclust

- 745 databases of clustered and deeply annotated protein sequences and alignments. *Nucleic*  
746 *Acids Res.* 2017;45:D170–6.
- 747 40. Steinegger M, Mirdita M, Söding J. Protein-level assembly increases protein sequence  
748 recovery from metagenomic samples manifold. *Nat Methods.* 2019;16:603–6.
- 749 41. Johnson LK, Alexander H, Brown CT. Re-assembly, quality evaluation, and annotation of  
750 678 microbial eukaryotic reference transcriptomes. *Gigascience.* 2019;8.
- 751 42. Steinegger M, Söding J. Clustering huge protein sequence sets in linear time. *Nat*  
752 *Commun.* 2018;9:2542.
- 753 43. Ren R, Sun Y, Zhao Y, Geiser D, Ma H, Zhou X. Phylogenetic resolution of deep  
754 eukaryotic and fungal relationships using highly conserved low-copy nuclear genes.  
755 *Genome Biol Evol.* 2016;8:2683–701.
- 756 44. Nakamura T, Yamada KD, Tomii K, Katoh K. Parallelization of MAFFT for large-scale  
757 multiple sequence alignments. *Bioinformatics.* 2018;34:2490–2.
- 758 45. Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of  
759 large phylogenies. *Bioinformatics.* 2014;30:1312–3.
- 760 46. Ghurye JS, Cepeda-Espinoza V, Pop M. Metagenomic assembly: overview, challenges  
761 and applications. *Yale J Biol Med.* 2016;89:353–62.
- 762 47. Frank JA, Pan Y, Tooming-Klunderud A, Eijsink VGH, McHardy AC, Nederbragt AJ, et  
763 al. Improved metagenome assemblies and taxonomic binning using long-read circular  
764 consensus sequence data. *Sci Rep.* 2016;6:25373.
- 765 48. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and  
766 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome*  
767 *Res.* 2017;27:722–36.
- 768 49. Driscoll CB, Otten TG, Brown NM, Dreher TW. Towards long-read metagenomics:  
769 complete assembly of three novel genomes from bacteria dependent on a diazotrophic  
770 cyanobacterium in a freshwater lake co-culture. *Stand Genomic Sci.* 2017;12:9.
- 771 50. Warwick-Dugdale J, Solonenko N, Moore K, Chittick L, Gregory AC, Allen MJ, et al.  
772 Long-read viral metagenomics captures abundant and microdiverse viral populations and  
773 their niche-defining genomic islands. *PeerJ.* 2019;7:e6800.
- 774 51. Mirdita M, Steinegger M, Söding J. MMseqs2 desktop and local web server app for fast,

- 775 interactive sequence searches. *Bioinformatics*. 2019;35:2856–8.
- 776 52. Mann DG, Droop SJM. Biodiversity, biogeography and conservation of diatoms.  
777 *Hydrobiologia*. 1996;336:19–32.
- 778 53. Norton TA, Melkonian M, Andersen RA. Algal biodiversity. *Phycologia*. 1996;35:308–26.
- 779 54. de Vargas C, Audic S, Henry N, Decelle J, Mahe F, Logares R, et al. Eukaryotic plankton  
780 diversity in the sunlit ocean. *Science*. 2015;348:1261605–1261605.
- 781 55. Ovchinnikov S, Park H, Varghese N, Huang P-S, Pavlopoulos GA, Kim DE, et al. Protein  
782 structure determination using metagenome sequence data. *Science*. 2017;355:294–8.
- 783 56. Söding J. Big-data approaches to protein structure prediction. *Science*. 2017;355:248–9.
- 784 57. Worden AZ, Allen AE. The voyage of the microbial eukaryote. *Curr Opin Microbiol*.  
785 2010;13:652–60.
- 786 58. Karlin S, Altschul SF. Methods for assessing the statistical significance of molecular  
787 sequence features by using general scoring schemes. *Proc Natl Acad Sci*. 1990;87:2264–8.
- 788 59. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence  
789 data. *Bioinformatics*. 2014;30:2114–20.
- 790 60. Ondov BD, Bergman NH, Phillippy AM. Interactive metagenomic visualization in a Web  
791 browser. *BMC Bioinformatics*. 2011;12:385.
- 792 61. Chen F, Mackey AJ, Stoeckert CJ, Roos DS. OrthoMCL-DB: querying a comprehensive  
793 multi-species collection of ortholog groups. *Nucleic Acids Res*. 2006;34:D363-8.
- 794 62. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new  
795 developments. *Nucleic Acids Res*. 2019;47:W256–9.