

1

2 **Metabarcoding analysis on European coastal samples**
3 **reveals new molecular metazoan diversity**

4

5 David López-Escardó¹, Jordi Paps², Colomán de Vargas^{3,4}, Ramon Massana⁵, Iñaki
6 Ruiz-Trillo^{1,6,7*}, Javier del Campo^{1,5*}

7

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

10 ²*School of Biological Sciences, University of Essex, Wivenhoe Park, Colchester, CO4*
11 *3SQ, UK*

12 ³*CNRS, UMR 7144, Adaptation et Diversité en Milieu Marin, Station Biologique de*
13 *Roscoff, Roscoff, France*

14 ⁴*UPMC Univ. Paris 06, UMR 7144, Station Biologique de Roscoff, Roscoff, France*

15 ⁵*Department of Marine Biology and Oceanography, Institut de Ciències del Mar*
16 *(CSIC), Barcelona, Catalonia, Spain*

17 ⁶*ICREA, Pg. Lluís Companys 23, 08010 Barcelona, Catalonia, Spain*

18 ⁷*Departament de Genètica, Microbiologia i Estadística, Universitat de Barcelona,*
19 *Barcelona, Catalonia, Spain*

20 **Correspondence and requests for materials should be addressed to JdC (email:*
21 *jdelcampo@icm.csic.es) or IR-T (email: inaki.ruiz@ibe.upf-csic.es).*

22 **Abstract**

23 Although animals are among the best studied organisms, we still lack a full
24 description of their diversity, especially for microscopic taxa. This is partly due to the
25 time-consuming and costly nature of surveying animal diversity through
26 morphological and molecular studies of individual taxa. A powerful alternative is the
27 use of high-throughput environmental sequencing, providing molecular data from all
28 organisms sampled. We here address the unknown diversity of animal phyla in marine
29 environments using an extensive dataset designed to assess eukaryotic ribosomal
30 diversity among European coastal locations. A multi-phylum assessment of marine
31 animal diversity that includes water column and sediments, oxic and anoxic
32 environments, and both DNA and RNA templates, revealed a high percentage of
33 novel 18S rRNA sequences in most phyla, suggesting that marine environments have
34 not yet been fully sampled at a molecular level. This novelty is especially high among
35 Platyhelminthes, Acoelomorpha, and Nematoda, which are well studied from a
36 morphological perspective and abundant in benthic environments. We also identified
37 based on molecular data a potentially novel group of widespread tunicates. Moreover,
38 we recovered a high number of reads for Ctenophora and Cnidaria in the smaller
39 fractions suggesting their gametes might play a greater ecological role than previously
40 suspected.

41

42 **Introduction**

43 The animal kingdom is one of the best-studied branches of the tree of life ¹, with more
44 than 1.5 million species described in around 35 different phyla ². Some authors have
45 suggested there may be more than 10 million species of animals, indicating that there
46 is an extensive unknown animal diversity. This hidden diversity may vary according
47 to the animal phyla considered. Not surprisingly, those animal phyla with microscopic
48 representatives (i.e., those animals with a size below 2mm ³, also known as
49 micrometazoans ⁴) are suggested to contain most of this potential unknown diversity
50 ³.

51
52 Marine environments cover most of the earth's surface. More importantly, all
53 metazoan phyla, except onychophorans, have marine representatives, with up to 60%
54 including microscopic members ⁵. Copepods, for instance, are the most abundant
55 multicellular group of organisms on earth ⁶, highlighting the key role of microbial
56 animals in marine ecosystems. Given that the marine benthic meiofauna is also one of
57 the hot spots of alpha-diversity in the biosphere, marine environments thus appear to
58 be ideal sites in which to analyze animal diversity across phyla.

59
60 Classical methods to survey animal diversity, such as isolation and morphological
61 identification, might be ineffective to comprehensively analyze
62 micro/mesozooplanktonic ⁷ and meiofaunal diversity ⁸. The microscopic size of the
63 organisms and the wide variety of morphologies makes the identification process
64 tedious and slow, requiring taxonomists with experience in different groups to
65 properly assess the composition of the community and describe new species or
66 groups. Molecular techniques, and especially high-throughput environmental

67 sequencing (HTES), have recently provided a more efficient method to assess and
68 understand ecological patterns in the microbial world ⁹, including metazoans ^{8,10–12}.
69 Although, these studies have mainly focused on richness patterns in marine benthic
70 communities or in zooplanktonic communities, with special attention on copepods ^{7,13}.
71 Studies of microbial eukaryotes ^{14–16} and even some animal clades ¹⁷ suggest that
72 HTES could also be used to detect novel lineages. However, such an approach has yet
73 to be applied across the whole animal kingdom.

74 To obtain a better understanding of the genetic diversity of the different metazoan
75 phyla, and the potential of HTES to quantify diversity and novelty levels, we analyzed
76 a large dataset of ribosomal small subunit (18S rRNA) V4 region tags from European
77 coastal sampling sites in the context of the BioMarKs project, which was designed to
78 analyze the diversity of unicellular eukaryotes. The BioMarKs dataset is based on 137
79 RNA and DNA samples from six locations ^{14,18} (Fig. S1; Table S1). The use of RNA
80 in this dataset allows analysis that goes beyond the detection of cells or DNA material
81 in the environment, as it provides a window on biological activity. For each sampling
82 site, there is data from both pelagic and benthic environments, with the pelagic
83 samples being divided into different depths and size fractions (Table S2). The large
84 quantity of data, together with the use of a phylogenetically curated taxonomic
85 assignment has provided a global view of genetic diversity across all metazoan phyla.
86 Our data show that 18S rRNA HTES approaches can be used to infer diversity and
87 novelty. Furthermore, we provide evidence that many unsampled lineages remain
88 among animals, and that there are even some potential novel groups. Consequently,
89 greater efforts should be made to sample specific animal groups, especially in benthic
90 environments.

91

92 **Results**

93 *Metazoan 18S rRNA reference database*

94 An important point to consider when analyzing diversity by metabarcoding is how the
95 taxonomic assignment is done. It is known that the use of GenBank or SILVA as
96 reference databases to perform the taxonomic assignment ^{7,8,12,13,19,20} can be
97 problematic ²¹. The reason is that those databases contain numerous missannotations
98 that affect the final taxonomic assignment. To avoid this problem and to have the best
99 possible taxonomic assignment, we manually constructed a novel phylogenetically
100 curated metazoan 18S rRNA reference dataset.

101 Our database included 19,364 18S rRNA sequences retrieved from GenBank. The
102 database was curated in a phylogenetic-wise manner, so that each animal phylum had
103 the widest possible representation of internal groups and that each sequence had a
104 clear taxonomic assignment. The resulting database was subsequently used to assign a
105 taxonomic identity to the approximately 1.5 million reads analyzed, providing a
106 holistic and phylogenetically accurate view of the metazoan diversity.

107

108 *General abundance and richness patterns of microbial animals*

109 We first analyzed the relative abundance of metazoan reads within the whole
110 eukaryotic dataset. We found that metazoans reads were quite abundant compared to
111 other eukaryotic groups in both the DNA and RNA samples (Fig. 1; Fig. S2). This
112 high percentage of metazoan reads was especially notable in anoxic pelagic

113 environments and in oxic sediments (Fig.1B). Interestingly, metazoan reads were not
114 only abundant in the micro/mesoplankton fraction (68% DNA, 49% RNA of the total
115 eukaryotic reads), but also in the smaller fractions (i.e., the pico/nano fractions which
116 are less than 20um). The presence of a high percentage of metazoan reads in the
117 smaller fractions is especially relevant in the anoxic environment, with 75% of the
118 DNA reads (and 33% of the RNA) being assigned to metazoans.

119 The clustering of reads into OTUs yielded 1067 OTUs from 23 different metazoan
120 phyla (Fig.2, Table S4). 469 OTUs were found to be exclusive to benthic
121 environments, 505 to pelagic environments and 102 OTUs were present in both
122 (Fig.2A). Crustacea appeared as the richest clade (246 OTUs) within the pelagic-
123 exclusive dataset, followed by Polychaeta (45). Within the benthic (sediment)-specific
124 samples, the largest number of OTUs were from Nematoda (227), followed by
125 Crustacea (101). Polychaeta (31) and Crustacea (23) dominated the OTUs present in
126 both environments (Fig.2A).

127 The largest proportion of animal reads in oxic water column environments were from
128 Crustacea, which represented up to the 89% of DNA and 53% of RNA in the overall
129 metazoans reads from the micro/meso fractions (Fig. 1A). More than 80% of the
130 crustacean RNA reads, however, corresponded to 8 specific OTUs that were assigned
131 to copepods (Table S5). Besides crustaceans, there was also a high abundance of
132 reads from tunicates (5% DNA only, but 28% RNA) within the oxic
133 micro/mesoplanktonic samples, most of them corresponding to appendicularians
134 (Table S5). On the other hand, benthic samples were dominated by polychaetes (30%
135 DNA, 23% RNA) and crustaceans (19% DNA, 23% RNA) (Fig. 1B). Within benthic
136 Crustacea, ostracods and copepods were the most abundant groups (Table S6).

137

138 *Community structure across environments and size fractions*

139 To determine the biogeographical patterns of the microbial animals in our dataset, we
140 analyzed the presence/absence of OTUs in all five sites (discarding the anoxic
141 samples). A large fraction of the OTUs (668 out of 1076) were present in just one
142 single location. However, the number of reads of these "endemic" OTUs (around
143 $4 \cdot 10^4$) was three times lower than the 8 OTUs present in all sampling sites (around
144 $1.2 \cdot 10^5$ reads) (Fig. 2B). The taxonomic composition of the cosmopolitan OTUs (Fig.
145 2B) differed greatly from the complete dataset except for the crustacean dominance
146 (Fig. 2B). In particular, there were no nematodes or polychaetes among the
147 cosmopolitan OTUs, whereas a cnidarian and a craniate OTU appeared to be present
148 over the 5 sampling sites. Our analysis also showed that all the cosmopolitan OTUs
149 belonged to the water column, whereas more than half (56%) of the "endemic" ones
150 belonged to the sediments. These endemic OTUs represented 80% of the total benthic
151 OTUs.

152 RNA reads indicate metabolically active cells²². Interestingly, we found a relatively
153 high percentage of RNA reads assigned to metazoans in the smaller fractions (from
154 0.8 to 20 μm): 2.4 % in oxic and 32.4 % in anoxic samples (Fig. 1A). Therefore, we
155 decided to analyze the potential source of those RNA reads. Most of the reads were
156 crustaceans (36% RNA reads), followed by tunicates, ctenophores, cnidarians and
157 polychaetes (Fig. 1A). Ctenophores (85% RNA pico/nano fractions) and cnidarians
158 (16% RNA pico/nano fractions) dominated the reads assigned to metazoans in the
159 anoxic waters of Varna, Black Sea (Fig. 1B).

160 To understand whether the reads from the smaller fractions were directly derived from
161 the larger ones, we filtered the data based on their co-occurrence between the
162 pico/nano fraction and the micro/meso fractions. We observed that OTUs present in
163 both smaller and larger fractions had a clearly different proportion of reads (Fig. 3).
164 Most of the reads in the smaller fractions belonged to the ctenophores (58%), whereas
165 crustaceans dominated (52%) the micro/mesoplanktonic fractions. In this regard,
166 OTUs corresponding to *Pleurobrachia pileus* (a ctenophore) and *Aurelia aurita* (a
167 cnidarian) were especially enriched in the smaller fraction (Fig. 3), representing 57%
168 of all metazoan RNA reads, and up to 33% of all eukaryotic RNA reads in the anoxic
169 samples (Table S7) (Fig.1A).

170

171 *Sequence novelty*

172 We performed BLAST searches against the NCBI nt nr database to interrogate the
173 level of novelty in our molecular dataset across all animal phyla. The results revealed
174 a high degree of sequence novelty (Fig. 4A). In particular, 35.5% of our OTUs
175 (representing 10.5% of the reads) had a BLAST identity lower than 97% compared to
176 NCBI sequences (Fig. 4B). Moreover, up to 10% of the OTUs, which accounts for 5%
177 of the metazoan reads, had BLAST identities lower than 90%. The putative novelty
178 was especially high among platyhelminthes, acoelomorphs, and nematodes, in which
179 most of their OTUs (75%) had a BLAST identity lower than 97%. Gastrotrichs and
180 crustaceans also had significant novelty (40-50% of their OTUs had a BLAST identity
181 below 97%).

182 Interestingly, the OTUs that appear to be most abundant within the water column
183 (Table S5) and sediments (Table S6) correspond either to already known sequences or

184 with high similarity to known sequences. The level of novelty is also different
185 between benthic and pelagic environments. Thus, 70% of the OTUs found in benthic
186 environments had a BLAST identity of less than 97% (Fig. 2A), while this percentage
187 decreased to 21% of OTUs in the water column or to 11% of OTUs present in both
188 water column and benthos. This suggests that benthic marine environments are a
189 potential hot-spot to find new metazoan taxa or lineages.

190 Among the potential novelty, we detected a group of three OTUs that had a relatively
191 large number of RNA reads in the water column (1.8%). (Fig. 1, labelled as "MAME
192 1"; MArine METazoan group 1), and with BLAST identities around 95% against two
193 unclassified environmental sequences from GenBank (KC582969 and HQ869055).
194 Analysis of this group of OTUs in other HTES studies based on the 18S rDNA gene
195 revealed 66 more OTUs retrieved from SRA (14 OTUs) and Tara Oceans⁹ (52 OTUs)
196 that are potentially from the same MAME 1 clade. Those 69 OTUs from BioMarks,
197 SRA and Tara Oceans represent 389,703 reads in total, an indication that OTUs
198 assigned to this group are relatively common in marine environments. Indeed, we
199 found that MAME 1 was present in coastal and open waters with a widespread
200 distribution across the world's oceans (except for the Arctic) in both the surface and
201 the deep chlorophyll maximum (Fig. S5B).

202 To have a better understanding of its phylogenetic position, we performed
203 phylogenetic trees. Our trees placed the MAME 1 GenBank sequence within tunicates
204 by both maximum likelihood and Bayesian inference (Table S3), and with good nodal
205 support (79% bootstrap support and 0.99 Bayesian posterior probability), although
206 with relatively longer branches than the rest of the metazoans. To determine its
207 specific phylogenetic position within the tunicates, we inferred an additional tree with

208 most of the available 18S rRNA sequences of tunicates, representing most of the
209 known diversity of this phylum. In this tunicate-focused tree, the MAME 1 sequence
210 clustered with thaliaceans as sister-group to the genus *Doliolum*, although with low
211 nodal support (Fig. 5). Finally, we ran a RAxML-EPA analysis to place the 69 OTUs
212 plus the other NCBI sequence within the reference tree of metazoans and the tree of
213 tunicates. In both cases, the 69 OTUs clustered together, with the reference MAME 1
214 sequences forming a monophyletic clade. Thus, our phylogenetic analysis suggests
215 that MAME 1 represents a novel, previously undescribed group of tunicates. Given
216 their extremely long-branches, however, additional molecular data will be needed to
217 further confirm this relationship.

218

219 **Discussion**

220 *High-throughput sequencing, a powerful methodology to assess diversity*

221 HTES is a useful method, but it also has some caveats. For example, it is well known
222 that it may be misleading to directly translate reads and OTU numbers into biomass
223 and number of species, respectively. In particular, the use of amplicon data as a proxy
224 for metazoan biomass abundance has been disputed, also with RNA data²³. Different
225 number of rRNA copies in the genomes of different taxa, PCR primer mismatches and
226 amplicon lengths can all affect the correlation between morphological and molecular
227 data^{7,24}. However, some studies have indeed shown positive correlations between
228 read abundances and biomass patterns in bivalve and decapod larvae¹⁹ and within
229 copepod groups⁷. Thus, we believe our approach to biomass abundance, although not
230 perfect, is useful enough to report the most abundant groups. A good indication of our
231 approach is that we recovered the general patterns previously described in

232 micro/mesoplanktonic communities based on morphological observations ^{25,26}, in
233 which copepods were found to be predominant within micro/mesoplanktonic
234 communities ⁶ followed by appendicularians ²⁶. Moreover, we found a more
235 heterogenic distribution in benthic habitats, which is to be expected considering that
236 sediments are known to harbor most of the metazoan diversity ⁵.

237 Overall, our data confirms that, although with some caveats, HTES is a powerful tool
238 to assess diversity. In this regard, the construction of a phylogenetically curated
239 database to assign the OTU taxonomy has proven to be crucial for our analysis aimed
240 at describing novelty in different metazoan phyla. Our clustering of OTUs at 97% is
241 likely a conservative approach for metazoans ²⁷, and some of our OTUs may indeed
242 represent more than one species. This largely depends on each metazoan lineage and
243 its specific 18S rRNA evolution rate. Moreover, primer bias can affect the detection
244 of some groups, meaning that some taxa can be present in the environment but
245 missing in our dataset ²⁸. However, by clustering at 97% we can directly compare the
246 results with the rest of the eukaryotes and get a more stringent output avoiding
247 polymorphisms effects and an overrepresentation of the retrieved diversity.

248

249 *Benthic-Pelagic relationship*

250 Analysis of benthic and pelagic metazoan communities in our dataset revealed that
251 most OTUs are exclusively pelagic or benthic, showing few overlaps between the two
252 communities, in agreement with our beta-diversity analyses (Fig. S3, Fig. S4A) and
253 the literature available ^{29,30}. Only 10% of OTUs from our dataset were present in both
254 benthic and pelagic communities, and these mainly corresponded to polychaetes,
255 crustaceans, molluscs and cnidarians (Fig. 2A). Among the shared OTUs Polychaeta

256 and Mollusca water column reads probably represent juvenile pelagic stages^{31,32}

257 while the benthic reads from crustaceans and cnidarians, that are predominantly

258 pelagic, come likely from death organisms or debris.

259 In addition, our data clearly shows that the pelagic OTUs tend to be present in more

260 sites, while most of the benthic OTUs are restricted to one location. The restricted

261 presence of meiofaunal OTUs has been described previously²⁰. Thus, the distribution

262 in the water column fits more with the consideration that “everything is everywhere”

263³³, probably because pelagic animals have fewer dispersal barriers than do benthic

264 ones³⁴.

265

266 *An ecological role for gametes?*

267 Somewhat surprisingly, we observed a high percentage of metazoan reads in the

268 smaller size fractions of most water column samples (Figure 1). This includes, as

269 well, the samples derived from RNA templates, probably indicating a significant

270 biological activity of metazoans in those smaller fractions. We believe it is unlikely

271 that those metazoan RNA reads could come from an extracellular origin because RNA

272 is fragile and quickly degraded by ribonucleases, and its structure is easily affected by

273 both oxygen and water³⁵. Furthermore, the RNA reads from pico/nanoplanktonic

274 fractions contain a different taxonomic distribution compared to the extracellular

275 DNA samples and the micro/mesoplanktonic RNA samples (Fig. 1A and Fig. 3A).

276 Thus, and taking into account the small size reported for certain animal gametes, we

277 hypothesize that a large part of those metazoan reads from the smaller fractions most

278 likely come from metazoan gametes.

279 This is the case, for example, of the reads from smaller fractions assigned to tunicates,
280 ctenophores, cnidarians and polychaetes, since they all use external fertilization.
281 Ctenophora and Cnidaria, which are not only abundant in DNA reads but also have a
282 relatively high number of RNA reads in the smaller fractions (Fig. 3B), might be a
283 particularly notable example of the importance of gametes in the environment. The
284 co-occurrence of reads in both smaller and larger fractions, the overrepresentation in
285 the smaller ones and the fact that their sperm size is smaller than $5\ \mu\text{m}$ ^{36,37} are good
286 indicators that at least the RNA signal of cnidarians and ctenophores might
287 corresponds to gametes. That will not be the case for the reads assigned to copepods
288 in the smaller fractions. They cannot come from gametes, since copepods use internal
289 fertilization and release eggs larger than $50\ \mu\text{m}$ ³⁸. Therefore, the crustacean RNA
290 reads observed in smaller fractions (from 0.8 to $20\ \mu\text{m}$) are probably the result of cell
291 breakage from larger fractions (Fig. 3A). Finally, we note that some of the OTUs that
292 are exclusively retrieved from smaller fractions could also correspond to sperm from
293 organisms that are larger than 2mm or from benthic fauna with external fertilization
294 and gamete sizes less than $10\ \mu\text{m}$, such as certain ctenophores and polychaetes (Table
295 S7).

296 It is worth mentioning that metazoan RNA reads corresponding to germline cells
297 could account, in our data, for as much as 3.2% of the total eukaryotic RNA reads in
298 the smaller fractions (Table S7), and up to 33% of eukaryotic reads in anoxic samples.
299 Thus, their numbers are comparable to those from the unicellular heterotrophic
300 flagellates, which usually reach abundances of up to the 40% of eukaryotic RNA
301 reads in pico and nano plankton³⁹. Thus, and considering those abundances, sperm
302 may play an important ecological role in those environments, particularly in the Black
303 Sea anoxic waters. Further research is needed to assess the effect of sperm in

304 microbial nutrient fluxes, especially during spawning events, when it may represent a
305 passive member of the community eaten by other metazoans or protists from micro-
306 scale fractions.

307

308 *Novelty in different metazoan phyla*

309 We performed an analysis on novelty by plotting the pairwise identities of the first
310 BLAST hit against NCBI non-redundant database. This provided a distribution of the
311 "novel" OTUs (those with sequence identities lower than 97% to any NCBI sequence)
312 along different environments (Fig. 2) and for different metazoan phyla (Fig. 4).

313 Interestingly, we found that 45% of our metazoan OTUs had less than 97% identity
314 against the NCBI nt nr database. Why a threshold of 97% for novelty? We believe it
315 is the safest one to detect novelty, although we probably miss a lot of intra-genera or
316 intra-class variation, depending in the animal group. It is worth mentioning, however,
317 that by having a threshold of pair-wise identities below 97%, we avoid any potential
318 intra-individual polymorphic variants⁴⁰. Therefore, we follow the rationale that OTUs
319 that do not have 100% identities but close (98% or higher) against the first BLAST hit
320 from NCBI non-redundant database, are probably the same taxa (maybe representing
321 intraindividual variations) or very closely related species. In contrast, the OTUs that
322 have a BLAST identity under 97% represent much deeper changes, and so, they
323 clearly represent, at least, different taxa than the ones represented in Genbank. Some
324 OTUs, especially those 10% of our OTUs with pairwise identities against GenBank
325 under 90%, may even represent new clades.

326 Although one could argue that this degree of novelty might reflect sequencing
327 artifacts, we are confident it is not the case because 1) we have followed a stringent

328 chimera and singletons removal process, 2) the reads are distributed across different
329 samples, and 3) they are not homogeneously distributed among taxonomic groups. In
330 addition, around 80% of our OTUs have RNA reads and their taxonomic distribution
331 is almost identical to the DNA OTUs. So, these novel variants present in the RNA
332 subset are transcribed by active organisms and are less prone to be artifacts or rare
333 variants ⁴¹.

334 We are aware that detection of novelty in metazoans just with molecular data is
335 challenging, given that the number of described animal species is larger than the
336 number of 18S rRNA sequences available in public databases (Fig. S7B). Therefore, a
337 novel sequence might belong to a species that has already been described but not yet
338 sequenced. A complete database linking morphological and molecular data is needed
339 to fully solve this issue. However, the 18S rRNA data so far available certainly is a
340 good representation of known animal diversity (Fig. S7B), and we believe our study
341 does indicate which metazoan lineages contain the higher levels of hidden molecular
342 diversity, and so, which are the animal groups needed for a more extensive sampling.

343 Those animal groups with the higher levels of novelty are not others than crustaceans,
344 nematodes, platyhelminthes, gastrotrichs and acoelomorphs. With the exception of
345 crustaceans, these groups occupy early branching phylogenetic positions within the
346 Ecdysozoa or the Lophotrochoa/Spiralia, or even within the Bilateria ⁴². Moreover,
347 the high genetic diversity in often neglected groups such as Acoelomorpha ¹⁷ and
348 Gastrotricha ¹⁰ reveals that these groups need a deeper exploration. We cannot rule
349 out the possibility that the relatively fast evolutionary rates of the 18S sequences from
350 nematodes, acoelomorphs and chaetognaths may have an effect on these low
351 similarity values. In addition, intragenomic variability of the 18S rRNA gene, already

352 described in some metazoan groups such as Platyhelminthes⁴³ or Chaetognaths⁴⁴,
353 can also contribute to these novelty values. Nevertheless, those are specific, isolated
354 cases. There is certainly extensive genetic novelty in our dataset, suggesting that most
355 acoelomorph, platyhelminth, chaetognath, and nematode species have not yet been
356 sequenced. Some of these hidden animal OTUs occupy key phylogenetic positions,
357 which can help to better reconstruct the metazoan tree of life and unravel the
358 evolution of extant species from the Urmetazoan¹⁷.

359

360 *A potential novel group of tunicates revealed by HTES*

361 We also recovered and genetically described a potential novel group of tunicates, here
362 named as “MAME 1”. It could be argued that this group represents an already
363 described Thaliacean related to the genus *Doliolum* that happens to have never been
364 sequenced or rare variants of the 18S gene belonging to known species. However, we
365 consider these two options unlikely for several reasons. First, the group seems to be
366 well populated (69 OTUs between our data and public repositories) and present in
367 many environments worldwide, not only in coastal waters (Fig. S5). Moreover, the
368 pairwise identity of the two MAME 1 sequences retrieved from NCBI is about 89%,
369 suggesting is not a single species, but rather an entire group of sequences with high
370 genetic variability, forming an independent clade related to Thaliaceans (Fig. 5). In
371 fact, the nucleotide identity among MAME 1 OTUs is similar as the observed among
372 distant *Aplousobranchia* species (for example, there is an 88% of identity between the
373 18S rRNA of *Distaplia dubia* and *Diplosoma virens*). Finally, different classes of the
374 18S rRNA gene have not been reported yet in Tunicates (there are 628 tunicate 18S
375 ribosomal sequences available at Genbank) and the percentage of identity of MAME

376 1 sequences against described Tunicate species seems too low (78% of identity with
377 the best BLAST hit *Thalia democratica*) for a different 18S rRNA type. In animal
378 groups in which different classes of 18S rRNA gene have been described, such as in
379 chaetognaths, the intra-individual variation among 18S classes lies around 90-93% of
380 identity⁴⁴. Therefore, we suggest that MAME 1 might corresponds to a new group of
381 tunicates that contains a large number of RNA reads within micro/mesoplankton
382 environments and is present in different habitats. However, without morphological
383 data, we cannot truly discard the possibility that those sequences belong to a
384 molecular divergent group of Thaliacean species, already morphologically described,
385 but without genetic data available. Although this emphasizes the powerful of HTES to
386 assess biodiversity and detect novelty, it also highlights its limitations. Thus, it is
387 crucial to continue and improve the classical screenings of marine diversity, with the
388 aim to link altogether morphological and genetic information in order to better
389 understand the metazoan biodiversity of our oceans.

390 **Conclusions**

391 We have reported an analysis of micrometazoan diversity in the European coast based
392 on HTES that includes, for the first time, both water column and sediments, oxic and
393 anoxic environments, and both DNA and RNA templates. To assess taxonomy, we
394 constructed a novel reference dataset comprising all animal phyla, which was
395 manually and phylogenetically curated. Our data show general read abundance and
396 richness patterns that partially corroborate previous morphological^{5,6,25,26} and
397 molecular studies^{8,10,13,19,20,45}. Our data showed a high relative abundance of
398 metazoan RNA reads within pico-nano size fractions (0.8-20 μm), suggesting that the
399 sperm of Ctenophores and Cnidarians plays a relevant ecological role as part of the

400 microbial food network. These results show the potential of HTES techniques as a fast
401 and exhaustive method to approach the study of micrometazoan biomass and diversity
402 patterns.

403 This kind of data has allowed us to describe novelty values found in different animal
404 phyla. We observed that some animal phyla have much genetic novelty that is yet to
405 be unraveled, including novelty in several well sampled groups such as Crustacea,
406 Platyhelminthes or Nematoda. Our finding of a potential new group of widespread
407 tunicates (MAME 1) highlights the value of phylogenetic approaches to identify novel
408 groups within phyla. The finding of MAME 1 in several HTES datasets could be
409 considered the first step in a reverse taxonomic process ⁴⁶ potentially leading to
410 isolation and detailed description. Overall, our data show that, if we truly want to
411 understand the biodiversity of marine environments, it is important to further sample
412 animal taxa within those environments. To achieve that, we need to have better tools
413 for the genetic screening, and especially for the isolation and morphological
414 characterization of these organisms.

415 **Materials and Methods**

416 *Sampling, 454 sequencing, curation of the sequences and diversity analysis*

417 During the BioMarKs project (biomarks.eu), samples were collected in six European
418 coastal sites (Fig. S1; Table S1). For sampling collection details, DNA/RNA
419 extraction methods, PCR amplifications, 454 sequencing details and read filtering
420 process see the electronic supplementary material. Processed reads allowed to build a
421 OTU (Operational Taxonomic Unit) table (reads per sample) with usearch v8.1.861 ⁴⁷,
422 using the UPARSE OTU clustering algorithm ⁴⁸, at a threshold of 97% similarity.
423 Afterwards, we used our own metazoan reference dataset (available at figshare

424 <https://dx.doi.org/10.6084/m9.figshare.3475007.v1>) to assign a taxonomical
425 affiliation to our OTUs. Finally, we removed the putative chimeric metazoan
426 sequences using Mothur's Chimera Slayer ⁴⁹ and discarded all the singletons. We
427 determined the degree of novelty of our dataset, by blasting the OTU sequences
428 against NCBI nt nr (September 23 2014). The metazoan OTU table obtained was
429 processed for alpha and beta-diversity analyses using QIIME ⁵⁰. See the electronic
430 supplementary material for details on this section.

431 *Analysis of the RNA reads from the small fractions*

432 Using QIIME scripts, we binned the OTUs that contain RNA reads within the water
433 column of each sampling site into three different groups: 1) OTUs containing the
434 small fractions (pico/nano), 2) OTUs containing the larger fraction (micro/meso), and
435 3) OTUs present in both small and large size classes. OTUs representing less than 10
436 RNA reads per site were discarded.

437 *Phylogenetic analysis of MAME1 sequence tags*

438 In order to phylogenetically place the short reads assigned to the novel metazoan
439 group (MAME 1) within an animal and tunicate backbone, we performed a RAxML-
440 EPA analysis ⁵¹ using a metazoan and a tunicate reference tree using the longest
441 putative MAME 1 sequence found by BLAST at NCBI nt nr database (*KC582969*), as
442 a unique MAME 1 representative. Using the MAME1 tree and alignment as a
443 reference we recruited environmental 18S rDNA short reads from SRA and Tara
444 Oceans and used them to perform abundance and distribution analyses (see the
445 electronic supplementary material).

446

447 **Data accessibility**

448

449 Electronic supplementary material that accompanies the online version of this article
450 includes materials and methods and supplementary figures and tables. The complete
451 BioMarks sequencing dataset is available at European Nucleotide Archive (EMBL-
452 EBI) <http://www.ebi.ac.uk/ena>, under project accession number PRJEB9133. OTU
453 tables, 18S metazoan database, MAME 1 group OTU table and phylogenetic trees
454 data (alignments, sequences and trees) are available at Figshare:
455 <https://dx.doi.org/10.6084/m9.figshare.3475007.v1>.

456

457 **Acknowledgment**

458 This work was supported by an Institució Catalana de Recerca i Estudis Avançats
459 contract, two grants (BFU-2011-23434 and BFU2014-57779-P) from the Ministerio de
460 Economía y Competitividad (MINECO), one of which (BFU2014-57779-P) was co-
461 funded by the European Regional Development Fund (FEDER), and a European
462 Research Council Consolidator Grant (ERC-2012-Co -616960) to IR-T. We also
463 acknowledge financial support from the Secretaria d'Universitats i Recerca del
464 Departament d'Economia i Coneixement de la Generalitat de Catalunya (Project 2014
465 SGR 619). JdC is supported by a Marie Curie International Outgoing Fellowship grant
466 (FP7-PEOPLE-2012-IOF - 331450 CAARL). JP acknowledges support from the
467 European Research Council under the European Union's Seventh Framework Program
468 (FP7/2007- 2013) / ERC grant [268513]. The work is part of the EU ERA-Net program
469 BiodivERsA, under the project BioMarKs (Biodiversity of Marine euKaryotes).

470

471 **Author Contributions**

472 JdC and IR-T designed and coordinated the study. RM and CdV provided the data. DL-
473 E, JdC and JP prepared the 18S metazoan database. DL-E and JdC analyzed the data
474 and prepared the figures. DL-E, JdC, JP and IR-T interpreted the data. DL-E, JdC and
475 IR-T wrote the manuscript, while all authors commented the manuscript.

476

477 **Additional Information**

478 Competing financial interests: The authors declare no competing financial interests

479

480 **References**

- 481 1. del Campo, J. *et al.* The others: our biased perspective of eukaryotic genomes.
482 *Trends Ecol. Evol.* **29**, 252–259 (2014).
- 483 2. Zhang, Z.-Q. Q. Animal biodiversity: An update of classification and diversity
484 in 2013. *Zootaxa* **3703**, 5 (2013).
- 485 3. Blaxter, M. L. *et al.* Defining operational taxonomic units using DNA barcode
486 data. *Philos. Trans. R. Soc. B Biol. Sci.* **360**, 1935–1943 (2005).
- 487 4. Guil, N. Molecular approach to micrometazoans. Are they here, there and
488 everywhere? in *Biogeography of Microscopic Organisms* (ed. Fontaneto, D.)
489 284–306 (Cambridge University Press, 2011).
490 doi:10.1017/CBO9781107415324.004
- 491 5. Snelgrove, P. V. R. Getting to the bottom of marine biodiversity: sedimentary
492 habitats. *Bioscience* **49**, 129 (1999).
- 493 6. Humes, A. How many copepods? *Hydrobiologia* **292/293**, 1–7 (1994).
- 494 7. Hirai, J., Kuriyama, M., Ichikawa, T., Hidaka, K. & Tsuda, A. A metagenetic
495 approach for revealing community structure of marine planktonic copepods.
496 *Mol. Ecol. Resour.* **15**, 68–80 (2015).
- 497 8. Fonseca, V. G. *et al.* Second-generation environmental sequencing unmask
498 marine metazoan biodiversity. *Nat. Commun.* **1**, 98 (2010).
- 499 9. de Vargas, C. *et al.* Eukaryotic plankton diversity in the sunlit ocean. *Science*
500 **348**, 1261605–1261605 (2015).
- 501 10. Chariton, A. A., Court, L. N., Hartley, D. M., Colloff, M. J. & Hardy, C. M.
502 Ecological assessment of estuarine sediments by pyrosequencing eukaryotic
503 ribosomal DNA. *Front. Ecol. Environ.* **8**, 233–238 (2010).
- 504 11. Lallias, D. *et al.* Environmental metabarcoding reveals heterogeneous drivers
505 of microbial eukaryote diversity in contrasting estuarine ecosystems. *ISME J.*
506 **9**, 1208–1221 (2015).
- 507 12. Bik, H. M. *et al.* Metagenetic community analysis of microbial eukaryotes
508 illuminates biogeographic patterns in deep-sea and shallow water sediments.
509 *Mol. Ecol.* **21**, 1048–1059 (2012).
- 510 13. Pearman, J. K., El-Sherbiny, M. M., Lanzén, A., Al-Aidaros, A. M. &
511 Irigoien, X. Zooplankton diversity across three Red Sea reefs using
512 pyrosequencing. *Front. Mar. Sci.* **1**, 1–11 (2014).
- 513 14. del Campo, J. *et al.* Diversity and distribution of unicellular opisthokonts along
514 the European coast analysed using high-throughput sequencing. *Environ.*
515 *Microbiol.* **17**, 3195–3207 (2015).

- 516 15. Richards, T. A. *et al.* Molecular diversity and distribution of marine fungi
517 across 130 European environmental samples. *Proc. R. Soc. B Biol. Sci.* **282**,
518 20152243–20152243 (2015).
- 519 16. Pan, J., del Campo, J. & Keeling, P. J. Reference Tree and Environmental
520 Sequence Diversity of Labyrinthulomycetes. *J. Eukaryot. Microbiol.* **64**, 88–96
521 (2017).
- 522 17. Arroyo, A. S., López-Escardó, D., de Vargas, C. & Ruiz-Trillo, I. Hidden
523 diversity of Acoelomorpha revealed through metabarcoding. *Biol. Lett.* **12**,
524 20160674 (2016).
- 525 18. Massana, R. *et al.* Marine protist diversity in European coastal waters and
526 sediments as revealed by high-throughput sequencing. *Environ. Microbiol.* **17**,
527 4035–4049 (2015).
- 528 19. Lindeque, P. K., Parry, H. E., Harmer, R. a, Somerfield, P. J. & Atkinson, A.
529 Next generation sequencing reveals the hidden diversity of zooplankton
530 assemblages. *PLoS One* **8**, e81327 (2013).
- 531 20. Fonseca, V. G. *et al.* Metagenetic analysis of patterns of distribution and
532 diversity of marine meiobenthic eukaryotes. *Glob. Ecol. Biogeogr.* **23**, 1293–
533 1302 (2014).
- 534 21. Bik, H. M. *et al.* Sequencing our way towards understanding global eukaryotic
535 biodiversity. *Trends Ecol. Evol.* **27**, 233–43 (2012).
- 536 22. Felske, A. *et al.* Phylogeny of the Main Bacterial 16S rRNA Sequences in
537 Drentse A Grassland Soils. *Appl. Environ. Microbiol.* **64**, 871–879 (1998).
- 538 23. Lejzerowicz, F. *et al.* High-throughput sequencing and morphology perform
539 equally well for benthic monitoring of marine ecosystems. *Sci. Rep.* **5**, 13932
540 (2015).
- 541 24. Porazinska, D. L., Sung, W., Giblin-Davis, R. M. & Thomas, W. K.
542 Reproducibility of read numbers in high-throughput sequencing analysis of
543 nematode community composition and structure. *Mol. Ecol. Resour.* **10**, 666–
544 676 (2010).
- 545 25. Beaugrand, G., Brander, K. M., Alistair Lindley, J., Souissi, S. & Reid, P. C.
546 Plankton effect on cod recruitment in the North Sea. *Nature* **426**, 661–664
547 (2003).
- 548 26. Bouquet, J.-M. *et al.* Culture optimization for the emergent zooplanktonic
549 model organism *Oikopleura dioica*. *J. Plankton Res.* **31**, 359–370 (2009).
- 550 27. Tang, C. Q. *et al.* The widely used small subunit 18S rDNA molecule greatly
551 underestimates true diversity in biodiversity surveys of the meiofauna.
552 *Proceedings of the National Academy of Sciences* **109**, 16208–16212 (2012).

- 553 28. Creer, S. *et al.* Ultrasequencing of the meiofaunal biosphere: Practice, pitfalls
554 and promises. *Mol. Ecol.* **19**, 4–20 (2010).
- 555 29. Johnson, C. L. *et al.* Biodiversity and ecosystem function in the Gulf of Maine:
556 pattern and role of zooplankton and pelagic nekton. *PLoS One* **6**, e16491
557 (2011).
- 558 30. Giere, O. *Meiobenthology: the microscopic motile fauna of aquatic sediments.*
559 *Meiobenthology* (Springer-Verlag Berlin Heidelberg, 2008).
- 560 31. Ellien, C., Thiébaud, E., Dumas, F., Salomon, J.-C. & Nival, P. A modelling
561 study of the respective role of hydrodynamic processes and larval mortality on
562 larval dispersal and recruitment of benthic invertebrates: example of *Pectinaria*
563 *koreni* (Annelida: Polychaeta) in the Bay of Seine (English Channel). *J.*
564 *Plankton Res.* **26**, 117–132 (2004).
- 565 32. Andresen, H., Strasser, M. & van der Meer, J. Estimation of Density-
566 Dependent Mortality of Juvenile Bivalves in the Wadden Sea. *PLoS One* **9**,
567 e102491 (2014).
- 568 33. Beaugrand, G., Reid, P. C., Ibañez, F., Lindley, J. A. & Edwards, M.
569 Reorganization of North Atlantic marine copepod biodiversity and climate.
570 *Science* **296**, 1692–1694 (2002).
- 571 34. Thornhill, D. J., Mahon, A. R., Norenburg, J. L. & Halanych, K. M. Open-
572 ocean barriers to dispersal: A test case with the Antarctic Polar Front and the
573 ribbon worm *Parborlasia corrugatus* (Nemertea: Lineidae). *Mol. Ecol.* **17**,
574 5104–5117 (2008).
- 575 35. Lindahl, T. Instability and decay of the primary structure of DNA. *Nature* **362**,
576 709–715 (1993).
- 577 36. Gaino, E. & Scoccia, F. Sperm ultrastructure of a member of the black coral
578 family Aphanipathidae: *Rhipidipathes reticulata* (Anthozoa, Antipatharia).
579 *Tissue Cell* **42**, 391–394 (2010).
- 580 37. Fischer, A. H., Pang, K., Henry, J. Q. & Martindale, M. Q. A cleavage clock
581 regulates features of lineage-specific differentiation in the development of a
582 basal branching metazoan, the ctenophore *Mnemiopsis leidyi*. *Evodevo* **5**, 4
583 (2014).
- 584 38. Castellani, C. & Lucas, I. A. N. Seasonal variation in egg morphology and
585 hatching success in the calanoid copepods *Temora longicornis*, *Acartia clausi*
586 and *Centropages hamatus*. *J. Plankton Res.* **25**, 527–537 (2003).
- 587 39. Logares, R. *et al.* Diversity patterns and activity of uncultured marine
588 heterotrophic flagellates unveiled with pyrosequencing. *ISME J.* **6**, 1823–1833
589 (2012).
- 590 40. Stoeck, T. *et al.* Multiple marker parallel tag environmental DNA sequencing

- 591 reveals a highly complex eukaryotic community in marine anoxic water. *Mol.*
592 *Ecol.* **19**, 21–31 (2010).
- 593 41. Pawlowski, J., Esling, P., Lejzerowicz, F., Cedhagen, T. & Wilding, T. A.
594 Environmental monitoring through protist next-generation sequencing
595 metabarcoding: Assessing the impact of fish farming on benthic foraminifera
596 communities. *Mol. Ecol. Resour.* **14**, 1129–1140 (2014).
- 597 42. Telford, M. J., Budd, G. E. & Philippe, H. Phylogenomic insights into animal
598 evolution. *Current Biology : CB* **25**, (2015).
- 599 43. Carranza, S., Giribet, G., Ribera, C., Bagnuà, J. & Riutort, M. Evidence that
600 two types of 18S rDNA coexist in the genome of *Dugesia* (Schmidtea)
601 mediterranea (Platyhelminthes, Turbellaria, Tricladida). *Mol. Biol. Evol.* **13**,
602 824–832 (1996).
- 603 44. Gasmi, S. *et al.* Evolutionary history of Chaetognatha inferred from molecular
604 and morphological data: a case study for body plan simplification. *Front. Zool.*
605 **11**, 84 (2014).
- 606 45. Pearman, J. K. & Irigoien, X. Assessment of Zooplankton Community
607 Composition along a Depth Profile in the Central Red Sea. *PLoS One* **10**,
608 e0133487 (2015).
- 609 46. Markmann, M. & Tautz, D. Reverse taxonomy: an approach towards
610 determining the diversity of meiobenthic organisms based on ribosomal RNA
611 signature sequences. *Philos. Trans. R. Soc. B Biol. Sci.* **360**, 1917–1924 (2005).
- 612 47. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.
613 *Bioinformatics* **26**, 2460–2461 (2010).
- 614 48. Edgar, R. C. UPARSE: highly accurate OTU sequences from microbial
615 amplicon reads. *Nat. Methods* **10**, (2013).
- 616 49. Schloss, P. D. *et al.* Introducing mothur: Open-source, platform-independent,
617 community-supported software for describing and comparing microbial
618 communities. *Appl. Environ. Microbiol.* **75**, 7537–7541 (2009).
- 619 50. Caporaso, J. G. *et al.* QIIME allows analysis of high-throughput community
620 sequencing data. *Nat. Methods* **7**, 335–336 (2010).
- 621 51. Berger, S. A., Krompass, D. & Stamatakis, A. Performance, accuracy, and Web
622 server for evolutionary placement of short sequence reads under maximum
623 likelihood. *Syst. Biol.* **60**, 291–302 (2011).
- 624

625 **Figure Legends**

626 **Fig. 1: Relative abundances of different metazoan groups and metazoan relative**
627 **abundance compared to the eukaryotes.** Relative abundances of different metazoan
628 groups (colored columns) and metazoan relative abundance compared to total
629 eukaryotes (black columns) in **(a)** oxic fractions and anoxic fractions, and **(b)**
630 different depths, separated by DNA and RNA templates. The number above each
631 column represents the number of metazoan reads in the fraction/environment for the
632 given template (RNA or DNA).

633 **Fig. 2: Metazoan richness. (a)** The OTU distribution for each metazoan group
634 divided into pelagic specific, sediment specific and those present in both
635 environments. BLAST identities are also plotted against NCBI nr nt in dark/light blue.
636 On the right, there is a representation of the number of OTUs (blue line) and number
637 of reads (red line) based on their environment. **(b)** Environmental distribution of
638 OTUs is shown based on prevalence: In blue, pelagic-specific OTUs (i.e., OTU with
639 more than 90% of the reads within the water column); in green, OTUs present both in
640 the water column and the sediments; in brown, OTUs present only in sediments (i.e.,
641 OTUs with more than 90% of the reads within the sediments). In addition, BLAST
642 identities are shown against NCBI nr nt in dark/light blue. The number of OTUs (blue
643 line) and number of reads (red line) based on their occurrence in 1 or more (up to 5)
644 geographical site is shown to the right.

645 **Fig. 3: Analysis of the small (pico and nano) and large (micro/meso) fractions,**
646 **and extracellular DNA. (a)** Taxonomic distribution of the OTU reads in the smaller
647 and larger fractions and within the extracellular DNA. **(b)** Ratio of the numbers of
648 reads from the smaller fractions and large fraction for these OTUs.

649 **Fig. 4: Sequence novelty plus summary of OTUs/read numbers of the main**
650 **Metazoan phyla in our dataset. (a)** Distribution of OTU BLAST identities against
651 NCBI nt nr for the main phyla of our dataset. **(b)** Summary of the number of OTUs
652 (blue) and the number of reads (red) of the given phyla.

653 **Fig. 5: Tunicate 18S rRNA phylogenetic tree placing the novel metazoan group**
654 **MAME 1.** The tree was inferred using RaxML-EPA from the 18S rRNA gene
655 nucleotide sequence and including representatives from all sequenced tunicate groups.
656 The nodal support values marked with a dot correspond to maximum likelihood 100-
657 replicate bootstrap support and Bayesian posterior probabilities.

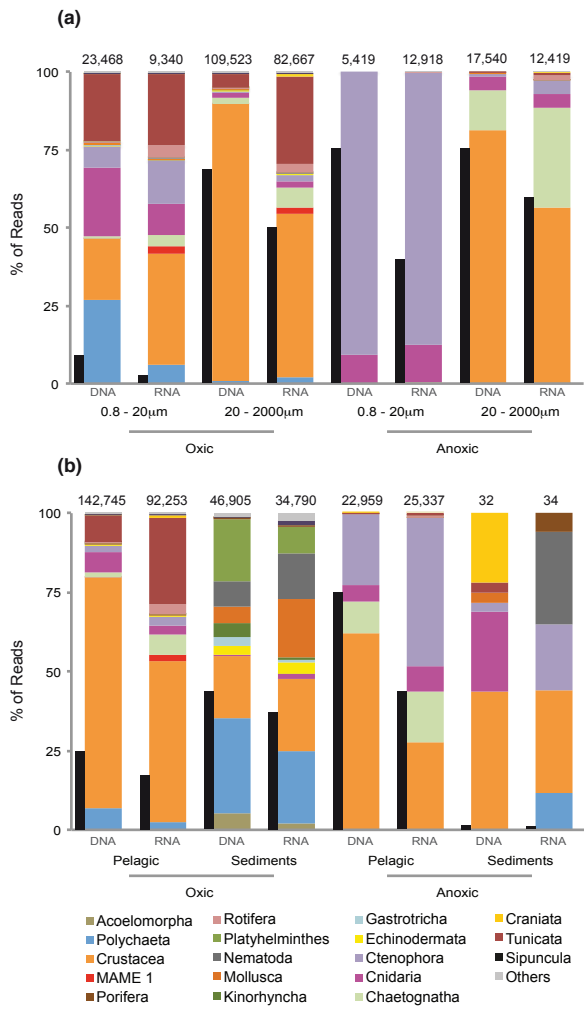


Fig 1

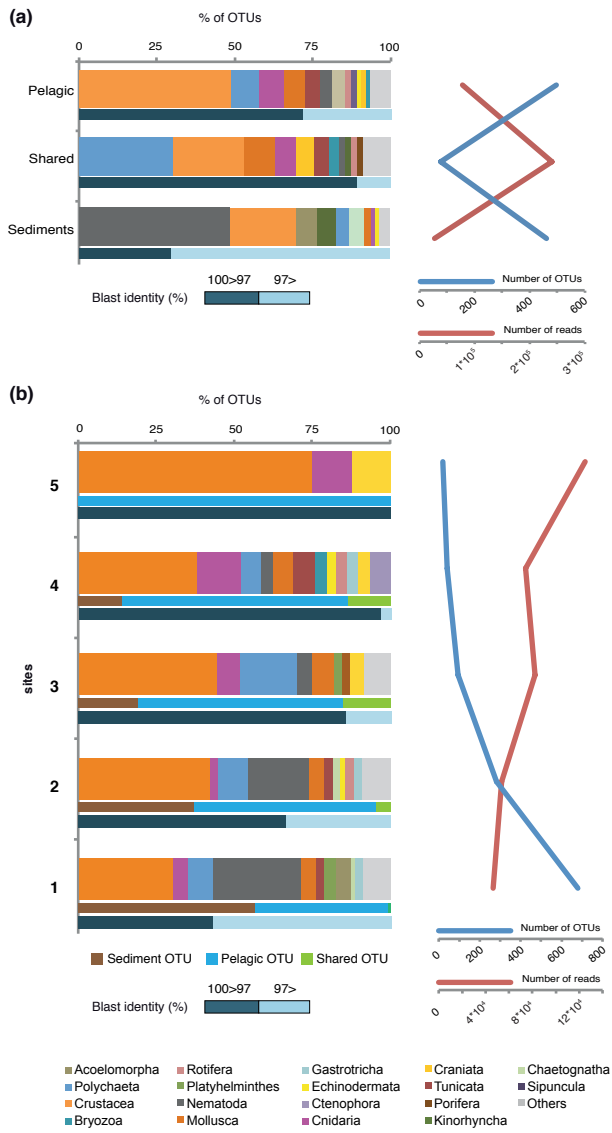


Fig 2

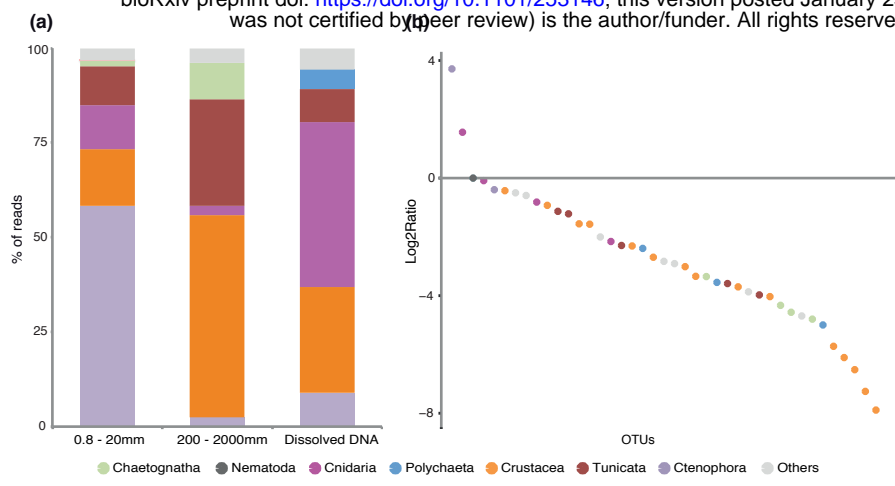


Fig 3

(a)

(b)

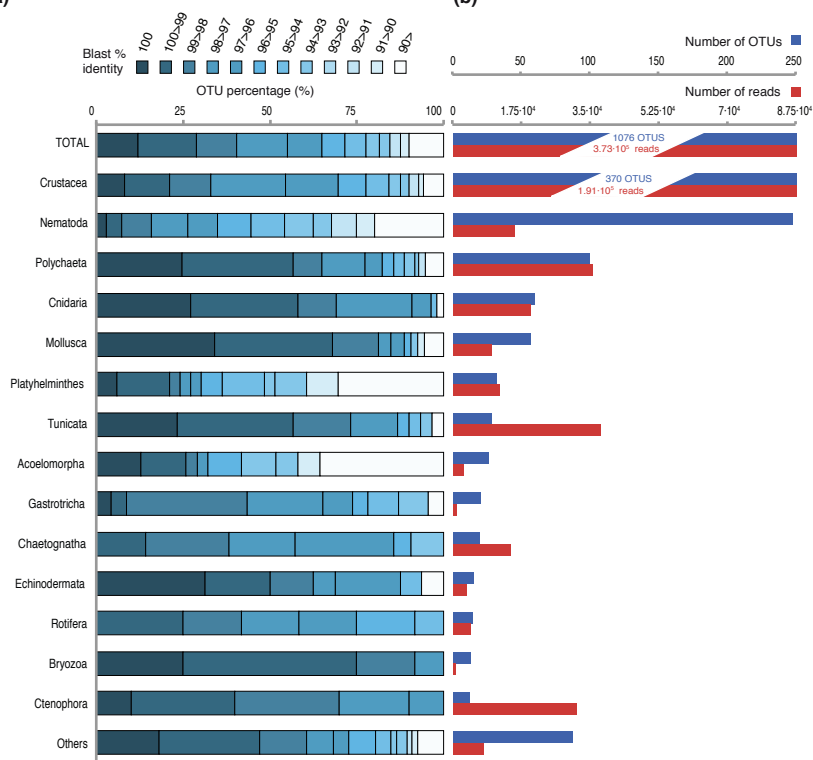


Fig 4

