

1 **Meta-Barcoding Accelerates Species Discovery and Unravels**
2 **the High Biodiversity of Benthic Invertebrates in Marine**
3 **Sediments of the Campos Basin, Brazil**

4

5 Milena MDP **Schettini**, Raony GCCL **Cardenas**, Marcella AA **Detoni**, Mauro F **Rebello**.

6 Instituto de Biofísica Carlos Chagas Filho. Universidade Federal do Rio de Janeiro. Rio de Janeiro, Rio de Janeiro. Brasil.

7

8 **KEYWORDS:** molecular taxonomy, 18S rRNA, 28 S rRNA, COI, metagenomics

9

10 **ABSTRACT**

11 Biodiversity is currently assessed for characterization and monitoring of the environment through
12 morphological taxonomy, a laborious and time-consuming process. We used 18S rRNA, 28S rRNA, and
13 cytochrome *c* oxidase I together with next-generation sequencing and bioinformatics to identify benthic
14 invertebrate organisms from sediment samples collected at five stations in the Campos Basin in southeast
15 Brazil, an important oil extraction area with one of the best-studied marine biota in Brazil. A total of 3.3
16 million sequences were clustered in operational taxonomic units, and more than 1.6 million sequences
17 (about 50% of all reads) were assigned to 957 prokaryotes and 577 eukaryotes. BLAST identified 23
18 phyla, 60 classes, 62 orders, 70 families, 67 genera, and 46 species of eukaryotes. Using meta-barcoding,
19 we identified phyla that are traditionally found in samples of marine benthos, including Annelida,
20 Arthropoda, Mollusca, and Chordata, as well as rare phyla such as Entoprocta and Gastrotricha. Taxa
21 identified through meta-barcoding were compared to data obtained through morphology from previous
22 studies in the area (REVIZEE, Habitats Project) and geo-validated with the Global Biodiversity Information
23 Facility database. This is the first report of a number of taxa in the Campos area, and the large number of
24 operational taxonomic units that were identified reveal a high level of benthic biodiversity in the Campos
25 Basin that has not been previously reported. Our study supports the application of meta-barcoding for
26 environmental characterization and monitoring programs, which could greatly reduce the time currently
27 required for species identification and biodiversity determination.

28 **INTRODUCTION**

29 Sediment fauna characterization and monitoring are mandatory requirements for obtaining oil
30 and gas (O&G) environmental permits for exploration and production (E&P) activities. This
31 requirement is expected to remain a key element for environmental management in the future,
32 particularly in the frontiers of deep-sea offshore oil exploration areas such as the Equatorial
33 Margin and the Santos Pre-salt Basin in Brazil.

34 Biodiversity identification, which is required for environmental characterization and
35 monitoring, is commonly carried out by morphological taxonomy, a laborious and time-
36 consuming process. As a general rule, taxonomic resolution at the species level is expected, but
37 for some fauna groups, the expertise required for this task is so specialized that only a handful of
38 individuals in the world are qualified to perform it. In addition, expert judgment is never 100%
39 accurate, and a 50% rate of identification consistency among taxonomists has been reported
40 (Culverhouse et al., 2003). The occurrence of pseudo-absence is frequent, especially for fragile
41 organisms that require special fixation procedures (Costa-Paiva et al., 2007). As a result,
42 invertebrate morphological identification efforts are often limited to few groups, including
43 Mollusca, Crustacea, and Polychaeta (Scaramuzza, 2015), while some estimates suggest that
44 more than 90% of all marine species have never been named (Scheffers et al., 2012).

45 The typical number of sediment samples in a monitoring campaign is in the range of tens, but in
46 sedimentary basins as large as 300,000 km², this number can extend to tens of thousands of
47 samples for baseline environmental characterization. The lack of expertise for morphological
48 taxonomy of some groups is a major bottleneck in the process of identifying biodiversity (Hebert
49 et al., 2003; Mora et al., 2013), and as a result, taxonomists are frequently unable to meet the
50 demands for biodiversity assessment required for monitoring programs, delaying the
51 development of economic activities and the discovery of new species.

52 According to the latest Report of the Convention on Biological Diversity (Diversity, 2016), Brazil
53 is the most biologically-diverse country in the world, with more than 100,000 animal species
54 described. However, only 184 marine invertebrates have had their conservation status assessed

55 (Scaramuzza, 2015). It is possible that current risk estimates of environmental impact are based
56 on underestimated biodiversity inventories, representing a threat to species conservation (Wu,
57 1982). Developing new technologies and approaches that accelerate species discovery and
58 reveal hidden biodiversity is crucial for setting conservation priorities and efforts.
59 Meta-barcoding uses genetic marker data generated through high-throughput next-generation
60 sequencing (NGS) of environmental samples (Leray and Knowlton, 2015) to greatly accelerates
61 species discovery and assess biodiversity. Since 2010, more than 600 papers have been
62 published on the use of DNA-based identification methods for species conservation (Goldberg et
63 al., 2015; Bergman et al., 2016), biodiversity inventory determination (Drummond et al., 2015),
64 environmental monitoring (Bohmann et al., 2014; Chariton et al., 2015; Leray and Knowlton,
65 2015; Brown et al., 2015), and DNA extraction/detection (Pedersen et al., 2014; Eichmiller et al.,
66 2014; Ficetola et al., 2016), and the technique has been considered a major tool for ocean
67 sustainability in the 21st century (Aricò, 2015). This approach is particularly useful because of
68 its sensitivity in identifying minute organisms and species in sediment (Wang et al., 2014). For
69 eukaryotic organisms that have not yet had their genetic markers sequenced or have not yet
70 been described morphologically, the concept of operational taxonomic unit (OTU) can be applied
71 (Stackebrandt and Goebel, 1994; Pedersen et al., 2014).

72 In this study, we combined three different phylogenetic markers (18S rRNA, 28S rRNA, and
73 cytochrome *c* oxidase subunit I [COI]) together with NGS and bioinformatics to identify benthic
74 invertebrate organisms using metagenomes from sediment samples collected in the Campos
75 Basin in southeast Brazil, an important oil extraction area with one of the best-studied marine
76 biota in Brazil (Miloslavich et al., 2011).

77

78 MATERIAL AND METHODS

79 Sample collection and processing:

80 Samples were collected in a survey in 2009 as part of the Habitats Project – Campos Basin
81 Environmental Heterogeneity coordinated by CENPES/PETROBRAS. Table 1 presents

82 information (collection date, geographic coordinates, and depth) on the five sampling stations
83 (B3, B4, C2, G2, and F5) in the Campos Basin. Sediment samples were collected in triplicate by
84 lowering a Van Veen grab at three different points around (150 m radius) each of the five
85 stations, resulting in a total of 15 sediment samples. At the time of collection, there were no
86 plans to have them genetically analyzed, and thus they were kept at -20°C for 4 years until our
87 analysis was done in 2013.

88 A 200-g subsample of the 0–2 cm slice of sediment of each sample was manually homogenized,
89 and DNA was extracted from 5 g of this subsample using the PowerMax Soil DNA Isolation kit
90 (MoBio Inc) according to the manufacturer's instructions. DNA integrity was checked on a 1.2%
91 agarose gel. Quantification was performed using a Qubit 2.0 Fluorometer (Life Technologies).

92 **Biogeography data:**

93 Biogeography data on the organisms identified in this study were extracted from previous
94 studies. Data on Cnidaria, Crustacea, Echinodermata, Mollusca, Nematoda, Polychaeta, and
95 Porifera groups were taken from the Brazilian REVIZEE (Living Resources in the Exclusive
96 Economic Zone) program (Lavrado and Ignacio, 2006), whereas the data for organisms of the
97 phyla Annelida, Arthropoda, Brachiopoda, Bryozoa, Cnidaria, Echinodermata, Echiura,
98 Foraminifera, Haptophyte, Mollusca, Nematoda, Nemertea, Porifera, Priapulida, Protozoa, and
99 Rodophyta were taken from the Habitats Project and provided by CENPES/PETROBRAS
100 (unpublished data). We also used the Global Biodiversity Information Facility (GBIF) database
101 (www.gbif.org) for organism geo-localization.

102

103 **PCR and high-throughput sequencing:**

104 Information on PCR of the COI, 18S rRNA, and 28S rRNA genes is presented in Figure 1 in
105 Supplementary Material. We used the Ion Xpress™ Plus Fragment Library kit (Life Technologies)
106 for preparing the libraries for sequencing according to the manufacturer's instructions for gDNA
107 fragment library preparation. Template preparation and sequencing were done using the Ion

108 PGM™ Template OT2 400 kit. Sequencing was done using the Ion Personal Genome Machine
109 (PGM™) System at the Life Technologies laboratories (São Paulo, SP), using Chip 318 v2.

110 **Bioinformatics and Taxonomic Name Attribution:**

111 Reads were prefiltered using the Torrent Suite software version 4.0.2 (Life Technologies) and
112 assigned to samples based on a combination primer tail-Ion Xpress barcode. PRINSEQ version
113 0.20.4 (Schmieder and Edwards, 2011) was used to remove poly A/T tails longer than 5 bases,
114 reads with unidentified (N) bases, reads shorter than 80 bp, and bad quality reads (Q<20). The
115 remaining reads were clustered in operational taxonomic units (OTUs) using CD-HIT-EST
116 version 4.6 (Li and Godzik, 2006) (up to 97% identity under 100% coverage within a bigger
117 read, word size of 10, and 20 penalty points for gaps).

118 High quality and low redundancy sequences were compared to NCBI
119 (<http://www.ncbi.nlm.nih.gov>) non-redundant nucleotide repositories (NR) using the
120 Nucleotide Basic Local Alignment Search Tool (BLASTn) version 2.3.0+ (Zhang et al., 2000). Max
121 e-value was of 10^{-5} and the number of events (referred to henceforth as 'hits') per query was
122 limited to 100.

123 Taxonomic names were attributed to each read based on the read's group of BLAST hits using
124 the Lowest Common Ancestor Assignment (LCA) algorithm in the MEGAN software (MEta
125 Genome Analyzer v. 5.10.3; Huson et al., 2007) with parameter adjustment (Huson et al., 2011).
126 Cladograms and rarefaction curves at the family taxonomic level were also built for each station
127 using MEGAN.

128 The BLAST step was performed using the Elastic Compute Cloud (EC2) service of Amazon
129 (aws.amazon.com). The BLAST for each of the 15 sets of reads corresponding to the 15 samples
130 was run in a parallel scheme using eight threads on up to 96 AWS instances with 8 processors
131 and 16 Gb of RAM each.

132

133 **RESULTS**

134 We obtained an average of 4.83 μg of DNA from each of the 15 samples. Sequencing generated
135 approximately 4.8 million sequences with an average size of 155.1 bp. Over 3.6 million
136 sequences (75.35%) passed quality control, and of these around 3.3 million were clustered in
137 OTUs by CD-HIT. Table 2 shows (1) the total number of OTUs, (2) the number of OTUs with no
138 hits in BLAST, (3) the number of OTUs with reads not attributed to any taxa by LCA, and (4) the
139 number of OTUs with attributed reads. For the five sampling stations, more than 1.6 million
140 sequences (about 50% of all reads) were assigned to 957 prokaryotes and 577 eukaryotes by
141 the LCA algorithm in MEGAN using hits produced by the BLAST similarity algorithm with any of
142 the three molecular markers (18S rRNA, 28S rRNA, or COI). LCA further identified 23 phyla, 60
143 classes, 62 orders, 70 families, 67 genera, and 46 species of eukaryotes. Figure 1A shows the
144 distribution of the 13 invertebrate phyla OTU identified by meta-barcoding for each of the five
145 stations, and Figure 1B shows the same for the 38 invertebrate families OTU identified. All other
146 prokaryotes and eukaryotes identified in this study by any of the three molecular markers and
147 classified to the taxonomic level of family are listed in the cladograms available in Supplementar
148 Material 2 in for each of the five sampling stations.

149 Our analysis identified 38 families of invertebrates in the 15 samples from the five sampling
150 stations in the Campos Basin. Figure 2 compares the spatial distribution of families identified by
151 meta-barcoding from phyla with most abundant frequencies (Annelida, 9 families, Fig. 2A;
152 Arthropoda, 10 families, Fig. 2B; and Mollusca, 7 families, Fig. 2C) in relation to previously
153 published morphologic taxonomy results for stations B3, B4, C2, F5 and G2.

154 Initially, the LCA algorithm identified 46 species, of which 27 were invertebrates not previously
155 described in the region. A text search of the list of BLAST hits allowed for identification of an
156 additional 45 species of invertebrates that had previously been identified in the Campos Basin.

157 The full list of species identified in this study is in Supplementary Material 3.

158

159 **DISCUSSION**

160 In this study we report the first meta-barcoding description of eukaryote biodiversity in the
161 deep-sea Brazilian continental shelf. Of the 3.3 million sequences that were classified as OTUs,
162 more than 1.6 million were assigned to 957 prokaryotes and 577 eukaryotes. Even though the
163 association between a given OTU and a given species should be established carefully, the
164 remaining 1.6 million OTUs that were not identified in the current study based on the sequences
165 of genetic markers available in Genbank suggests that the benthic biodiversity of the Campos
166 Basin could be orders of magnitude higher than that established by previous studies based on
167 morphological taxonomy.

168 One of the differentiating characteristics of our study is the fact that we used samples collected
169 from areas where E&P activities have been carried out and where several previous morphological
170 taxonomic studies were performed, either by oil companies as part of the process to obtain
171 environmental permits or those involved in conservation programs (such as the Habitats
172 Project) or by the scientific community (specially the REVIZEE program).

173 The approximately 4.8 million sequences we found are within the expected range for the 318 v2
174 chip. Even though the average read size of 155.1 bp was below the expected number for the OT2
175 400 kit, it did not compromise our analysis.

176 When further analyzing the OTUs distributed in the 23 phyla, we found that a considerable
177 number of reads were assigned to the families Hominidea and Bovidae, increasing the number of
178 reads belonging to the Chordate phylum. However, these were the result of read alignments
179 generated against the whole human and bovine genomes or chromosomes, as opposed to the
180 three specific genetic markers. As the focus of this study was on benthic invertebrates because of
181 their significance for legally mandated environmental characterization and monitoring in
182 offshore areas, these artifact findings in the Chordata were ignored.

183 Our meta-barcoding analysis identified phyla that are traditionally found in samples of marine
184 benthos, including Annelida, Arthropoda, Mollusca and Chordata, as well as more rarely found
185 phyla such as Bryozoa, Cnidaria, Echinodermata, Nematoda, Nemertea, Platyhelminthes,

186 Porifera, and Priapulida, and even rarer phyla such as Entoprocta and Gastrotricha (Figure 1 and
187 Figure 2 in Supplementary Material).

188 The large number of OTUs for Annelida, Arthropoda, and Mollusca found by meta-barcoding
189 agrees with previous results for the Campos Basin (Lavrado and Ignacio, 2006) obtained by the
190 REVIZEE project and also with those from the Habitats Project. A recent meta-barcoding study
191 (Leray and Knowlton, 2015) also identified Annelida and Arthropoda as the phyla with the most
192 OTUs among the 22 phyla identified in approximately 0.09 m³ of sediment from coral reef
193 regions in Virginia and Florida in the United States.

194 The Entoprocta (or Kamptozoa) phylum comprises about 170 aquatic and sessile species that
195 are between 0.5 and 5.0 mm in size and are mostly marine (Zhang, 2011). Until 2011, only 18
196 species of Entoprocta were known on the Brazilian coast (Vieira and Migotto, 2011). In this
197 study, all Entoprocta OTUs (6 at the C2 station and 24 at the G2 station) were attributed to the
198 genus *Loxosomella* via the 28S rRNA marker and had over 86% similarity to sequences found in
199 Genbank. This result expands the distribution of the genus, which was previously limited to six
200 species collected off the coast of São Paulo (Vieira and Migotto, 2011). As for the cosmopolitan
201 Gastrotricha phylum that comprises about 790 species of aquatic organisms up to 1 mm in
202 length (Zhang, 2011), all 22 OTUs assigned to the phylum (C2 station) were in the
203 *Tetranchyroderma* genus, with over 81% similarity to COI sequences found in Genbank. This
204 finding also expands the distribution of this genus, which was previously limited to São Paulo
205 beaches located approximately 1000 km away from the Campos Basin (Garraffoni and Araújo,
206 2010).

207 Our meta-barcoding results were compared to those from a recent comprehensive
208 morphological taxonomy effort (the Habitats Project coordinated by CENPES/PETROBRAS) that
209 analyzed the same samples used in our study. The Habitats Project generated a databank with
210 data for almost 50,000 specimens and identified 17 phyla, 27 classes, 63 orders, 354 families,
211 768 genera, and 749 species. The comparison between the findings obtained with molecular and
212 morphological taxonomies was restricted, however, since 1211 (68%) of the 1773

213 macroinvertebrate taxa identified by morphological taxonomy did not have an entry in Genbank
214 for any of the three markers (18S rRNA, 28S rRNA, and COI) used in this study (indicating a
215 major underrepresentation of Brazilian marine species in Genbank and a need to increase efforts
216 to have sequences from more Brazilian species deposited in that database). Other factors that
217 limited the comparison between the two taxonomic approaches included uncertainty about how
218 much DNA was still available in the sediment samples that had been preserved at -20°C for 4
219 years as well as the limited amount of sample analyzed at each station (5 g of the 200 g of the 0–
220 2 cm slice of sediment that was homogenized, compared to 4 L of the 0–10 cm slice for the
221 morphological study). Finally, for many species, only partial sequences of the markers were
222 available in Genbank, and thus it was not possible to ensure that they were sufficiently aligned
223 with a given read to permit assignment to a taxon. These restrictions suggest that some families
224 that were apparently absent may not actually have been so. The continuing effort to add more
225 sequences to Genbank should clarify this issue.

226 Out of the 70 families identified by meta-barcoding, 21 were invertebrate. The families
227 Amphinomidae, Enchytraeidae, Glyceridae, Orbiniidae, Serpulidae, and Spionidae, belonging to
228 the phylum Annelida, were previously identified in the Campos Basin by the Habitats Project,
229 which also identified 28 other Annelida families not found by meta-barcoding. Hormogastridae,
230 which was found in our study, is most likely a false positive since it is not a marine family.

231 The families Solenoceridae, Cylindroleberididae, and Mysidae, belonging to the phylum
232 Arthropoda, have previously been identified in the Campos Basin and in southeastern Brazil by
233 other authors (Cardoso, 2007; Serejo et al., 2007; Tâmega et al., 2013), while 29 arthropod
234 families previously reported by the Habitats Project were not identified by meta-barcoding. The
235 families Miridae, Chalcididae, and Formicidae, all of which were found in our study, are most
236 likely false positives since they are non-marine insects.

237 All Mollusca families identified by meta-barcoding in the Campos basin except for Mytilidae have
238 previously been found in the region (Lavrado and Ignacio, 2006; Dornellas and Simone, 2011;

239 Tâmeaga et al., 2013), although not by the Habitats Project. The latter identified 15 Mollusca
240 families not identified by meta-barcoding.
241 Meta-barcoding was also able to find families at every sampling station that had not been
242 previously reported by the Habitats Project, such as Echiuridae and Pectinariidae in the
243 Annelida; Desmosomatidae and Hippolytidae in the Arthropoda; and Arcidae, Mactridae, and
244 Pectinidae in the Mollusca. This indicates that the distribution of these families may be broader
245 than that suggested by morphological taxonomy.

246 Of the 46 species identified by meta-barcoding, none of the 24 benthic invertebrates had been
247 previously described by the Habitats Project and may represent new occurrences in the region,
248 with the exception of the arthropod *Eurythenes gryllus*, previously identified in REVIZEE. We
249 must remember that even though species level resolution is desirable, taxonomic penetration to
250 the family level is accepted by environmental agencies, and most specimens in previous studies
251 have been identified only to this level. We found records of the families of all newly observed
252 species in the Habitats Project, REVIZEE and GBIF databases, which supports the argument that
253 these are new occurrences, even though we cannot rule out the possibility of false positive.

254 The comparison between our data and that of the Habitats Projects was limited by the
255 availability of sequences of the three genetic markers (18S rRNA, 28S rRNA and COI) deposited
256 in Genbank. Only 64 out of the 749 organisms identified to the species level by the Habitats
257 Project had at least one genetic marker sequence found in Genbank and thus were eligible for
258 molecular identification. However, none of those 64 species were identified by meta-barcoding.
259 We believe that these 64 missing species are pseudo-absence results, and that they were not
260 found because the samples had been preserved at -20°C for 4 years. However, another
261 explanation is possible. When analyzing our data, we noticed that even after calibration of the
262 parameters for the LCA algorithm (data not shown), some incongruence in the attribution of the
263 taxonomic name to a species could happen due to the selection of an unlikely BLAST hit to name
264 the query OTU. To overcome this problem, we text-searched the names of the organisms of all
265 BLAST hits that were associated to a given OTU and compared the names found for those of the

266 64 species found by the Habitat Project. We were thus able to identify 45 additional species that
267 had previously been described by morphological taxonomy but were not picked by the LCA
268 algorithm. The full list of species identified by molecular and morphological taxonomies,
269 together with their genetic markers that are available in Genbank, are listed in Supplementary
270 Material 3. Other pseudo-absence results could have been generated by the occurrence of
271 synonymous names at the species level. For instance, according to recent estimates, more than
272 80% of the algae in some genera and 38% of Mollusca have synonymous names. For marine
273 species, this percentage could reach 40% (Costello et al., 2013). An ongoing effort is dedicated to
274 resolving synonymous names found in the GBIF database.

275 The use of biogeographic databases (Habitats Project, REVIZEE and GBIF) to verify and adjust
276 the meta-barcode observations has proven to be a good strategy. False positive results can be an
277 artifact of the low representativeness of Brazilian species in Genbank. Due to similarities of
278 genetic sequences shared among species belonging to the same genus, BLAST can associate, with
279 very low error probability, a read from a species not present in the Genbank to a sequence from
280 a phylogenetically similar species from a different habitat. By using metadata on the distribution
281 of the species selected by BLAST, we managed to identify at least one instance of this occurring
282 in our analysis. The small (25–85 mm) gastropod *Haliotis diversicolor*, which was identified in
283 our study, is native to the Indo-Pacific Ocean, with geo-referenced records on the coast of Japan,
284 Thailand, and Australia (GBIF, 2016). Thus this result might have represented a new occurrence
285 of this species in a completely new environment or, alternatively, a false positive. *Haliotis*
286 *aurantium*, a small gastropod belonging to the same genus as *H. diversicolor*, has previously been
287 identified in the Campos Basin, and we believe that because Genbank contained no sequence of
288 *H. aurantium* corresponding to any of the three genetic markers, the LCA algorithm may have
289 erroneously assigned an OTU from *H. aurantium* to *H. diversicolor*. A system able to resolve these
290 types of incongruences would greatly improve meta-barcoding analysis.

291 In an attempt to prevent false positive results, we tested whether redundant identification by a
292 second or third genetic marker could confirm potential positive results of species that were

293 identified by a single marker. Unfortunately, that was not the case. Of the 46 species identified
294 by meta-barcoding, 16 had sequences of all three genetic markers available in Genbank but were
295 identified by just one of the three markers and not by the others. It was frequently the case that,
296 even though the sequence for a genetic marker for a specific organism was available in Genbank,
297 multiple names were attributed to the gene, only partial sequences were available, or sequences
298 had not been validated experimentally. Genbank is the best repository for genetic sequences yet
299 available, but it still does not offer a high level of confidence when it comes to the names
300 attributed to genetic sequences. Our research team is currently working on developing new
301 algorithms to help overcome this limitation.

302 The problems related to the presence of false positives and pseudo absences could be solved if
303 biodiversity characterization was done in a taxonomic-free context, looking only at OTUs to
304 compare biodiversity profiles among samples. The frequency and abundance of OTUs could then
305 be related to environmental changes, either spatial or seasonal, and species discovery would be
306 accelerated by identification of OTUs that vary according to environmental conditions. Were
307 such a strategy to be adopted, not only would it allow us to work with the hidden biodiversity of
308 the thousands of 'no hit' OTUs, but OTU profiles and their distribution could inform us of
309 environmental changes. Species names are fundamental to ecology, and in spite of all the
310 uncertainty that they entail, a lot of the accumulated knowledge in biology is associated with
311 these units. Although we may never give up on the idea of naming species, the ease of gathering
312 OTU data is unprecedented, making the prospect of a taxonomic-free ecology complementary to
313 the traditional one more and more likely in the years to come.

314

315 **CONCLUSION**

316 This study confirms that meta-barcoding can be a reliable, fast, and low cost tool for
317 environmental characterization and monitoring. It may be the only suitable method for
318 producing information critical for decision making about extensive, little-explored areas within
319 a short time period, thus safeguarding the environment without delaying economic activities. It

320 may also accelerate species discovery and contribute to ecological knowledge in ways that have
321 not yet been fully explored.

322 The methodology could be improved by adding more sequences of native species to public and
323 proprietary databases, but it is our opinion that meta-barcoding can already be considered to be
324 the best available technique for generating biodiversity inventories in marine sediments, and it
325 should be acknowledged as such by oil operators, environmental authorities, and the scientific
326 community at large.

327 Brazil has one of the strictest set of environmental laws and regulations for the O&G sector in
328 the world, and one that is constantly being improved. Recent changes made under resolution
329 CONAMA 422/11 minimized bureaucracy in the application process, increased transparency by
330 sharing information online, and reduced liability for the O&G operators. The Brazilian
331 environmental authority IBAMA (Brazilian Institute of the Environment and Renewable Natural
332 Resources) establishes guidelines and best practices for environmental licensing and monitoring
333 by means of 'reference terms.' By becoming an early adopter of meta-barcoding, IBAMA could
334 play a leading role in the implementation of this innovative methodology that can greatly
335 contribute to the conservation of deep-sea environments worldwide.

336

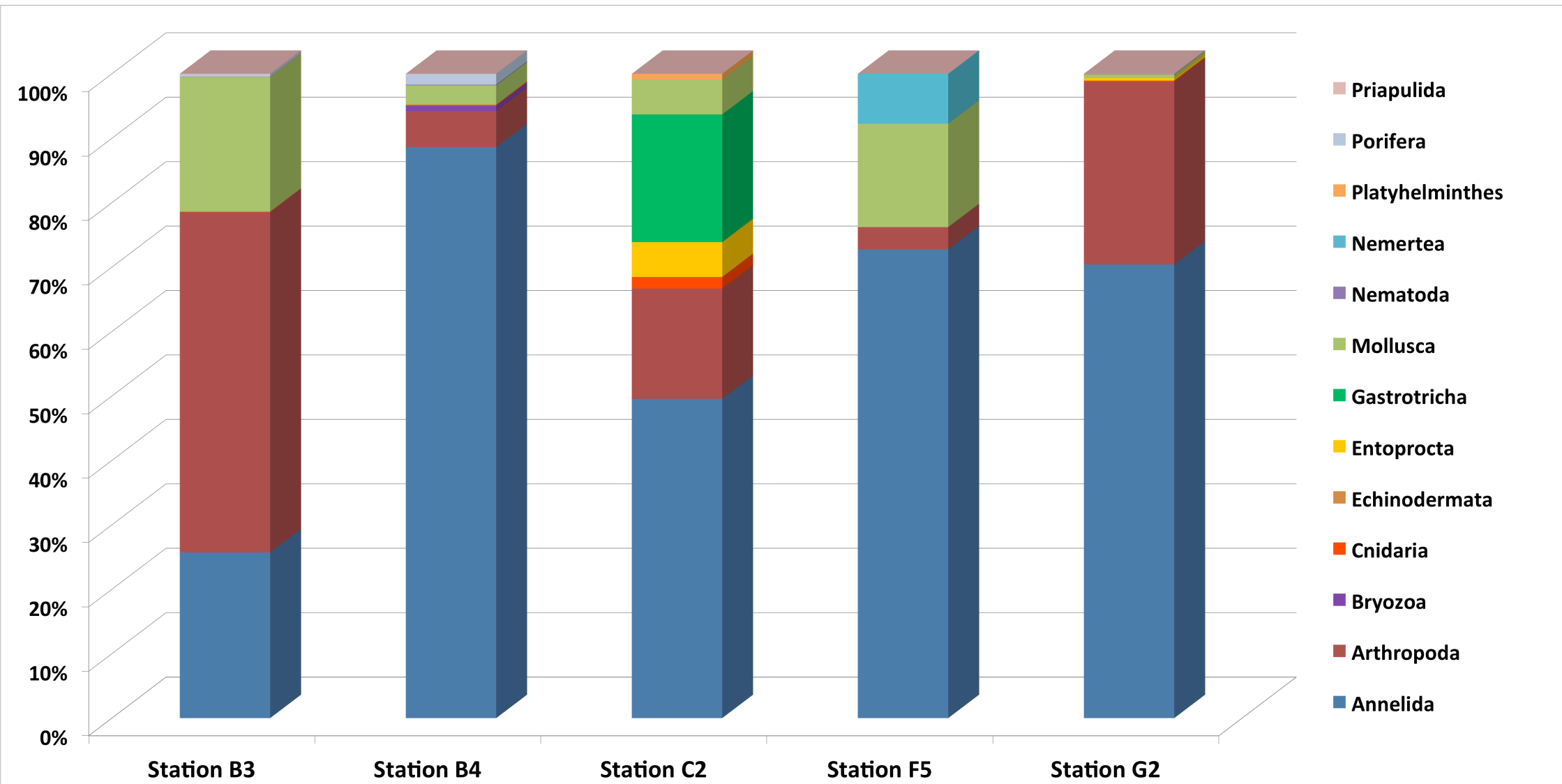
Table 1 – Sampling date, location and depth Location of sampling stations B3, B4, C2, F5 and G2 in Campos Basin, southeast Brazil.

	Sampling date	Latitude (SIRGAS2000)	Longitude (SIRGAS2000)	Depth (m)
Station B3	02/20/2009	-22,99701	-41,352583	77
Station B4	02/21/2009	-23,16851	-41,052264	107
Station C2	07/16/2009	-22,62599	-41,365082	54
Station F5	02/24/2009	-22,29010	-40,110584	143
Station G2	02/25/2009	-21,98502	-40,419918	56

Table 2 – OTU per sample. OTU without a similar sequence on Genbank NR are under ‘No fragments’. OTU that did not comply with established LCA parameters (e.g. score below threshold) or OTU that did not add up to a node are under ‘non attributed reads’. Also under ‘non-attributed’ are OTU not attributed by rRNA16S, taxa attributed by genes other than the 3 targets and taxa defined in Genbank as ‘undefined’. They were also disabled at the cladograms.

Sample	Total OTU	No Hits	Non attributed	Attributed
St. B3 rep. #1	101,966	20,505	73,653	7,808
St. B3 rep. #2	379,812	65,557	97,849	222,406
St. B3 rep. #3	84,180	12,167	57,290	14,723
St. B4 rep. #1	103,053	25,721	57,290	14,723
St. B4 rep. #2	332,953	35,384	64,066	236,503
St. B4 rep. #3	302,290	50,143	65,134	187,013
St. C2 rep. #1	245,233	34,452	40,687	170,094
St. C2 rep. #2	307,780	59,289	60,866	187,625
St. C2 rep. #3	249,969	56,247	81,114	112,608
St. F5 rep. #1	139,992	50,900	35,349	53,743
St. F5 rep. #2	105,435	32,435	47,684	25,316
St. F5 rep. #3	83,962	43,377	34,877	5,708
St. G2 rep. #1	173,740	71,230	60,632	41,780
St. G2 rep. #2	312,446	88,627	79,156	144,663
St. G2 rep. #3	347,494	32,832	120,519	194,143
TOTAL	3,270,206	678,866	959,986	1,631,453

Figure 1 – OTU occurrence in each station. Percentage of OTU for phyla (A) and Family (B) in each station.



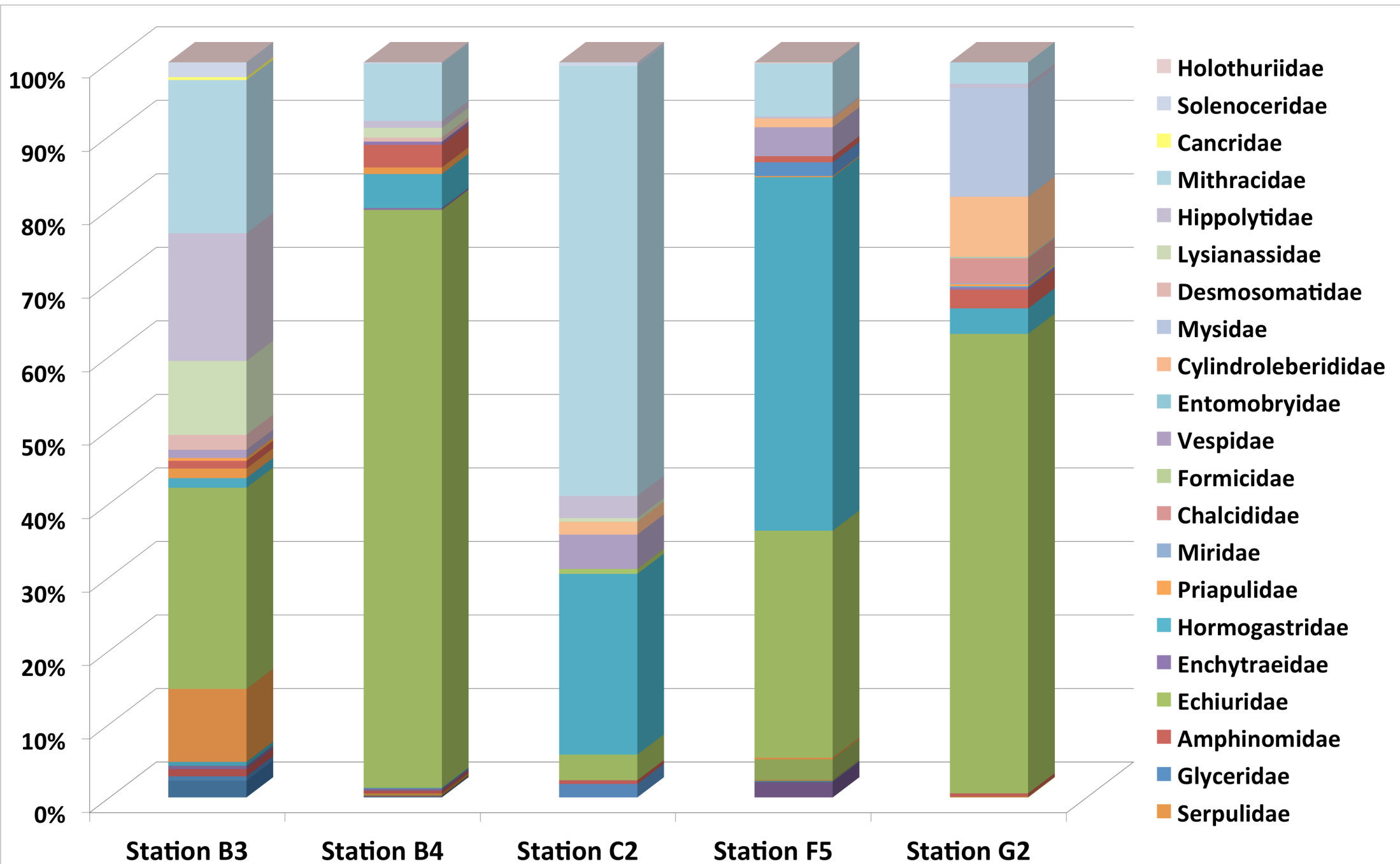
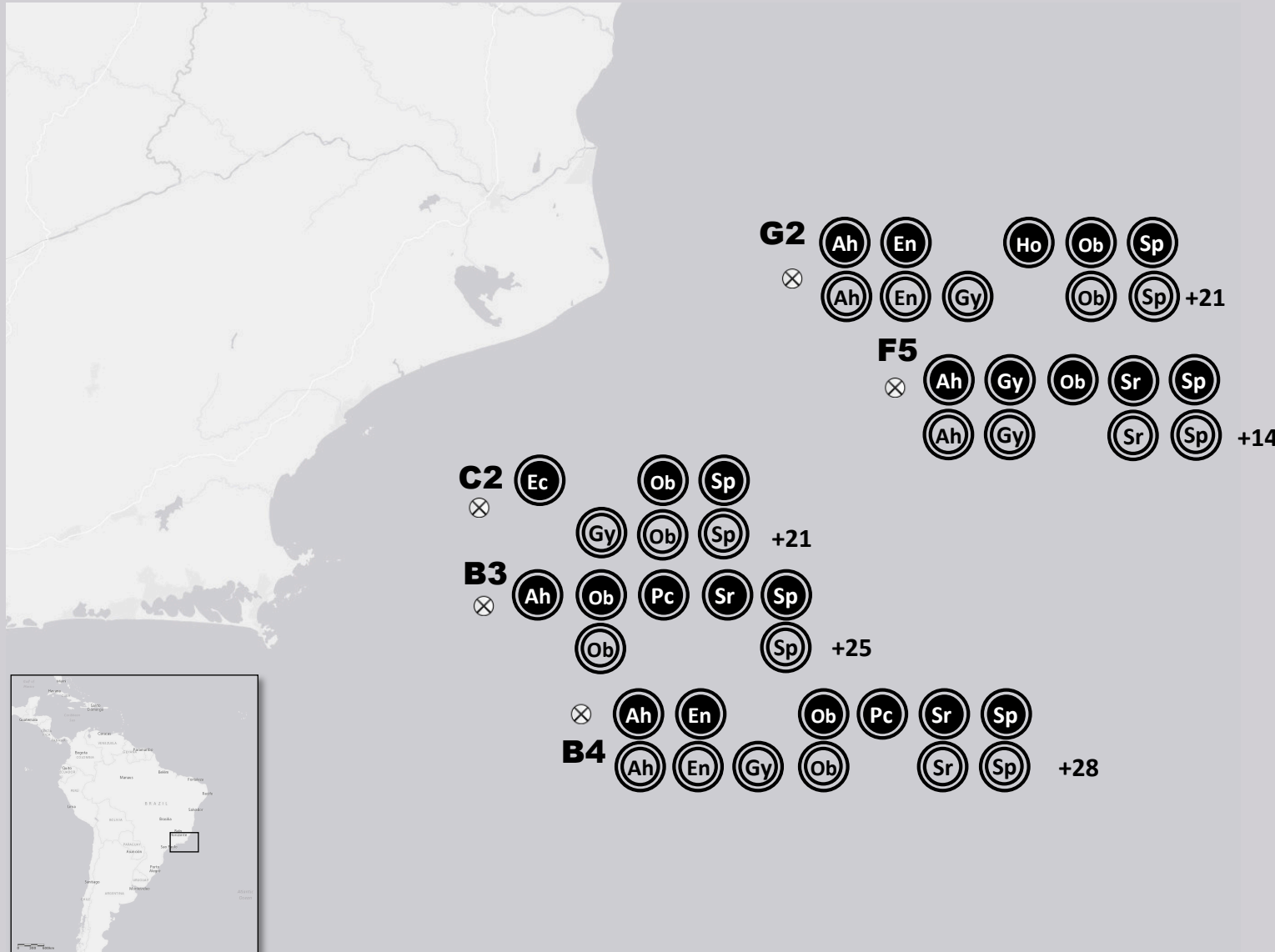


Figure 2 – Distribution of the main invertebrate phylum identified by molecular and morphological taxonomy in Campos Basin. A) annelida distribution, b) arthropoda distribution, C) mollusca distribution.

Annelida distribution

Molecular

Morphological



- Ah** **Ah** Amphinomidae
- En** **En** Enchytraeidae
- Ec** **Ec** Echiuridae*
- Gy** **Gy** Glyceridae
- Ho** **Ho** Hormogastridae**
- Ob** **Ob** Orbiniidae
- Pc** **Pc** Pectinariidae*
- Sr** **Sr** Serpulidae
- Sp** **Sp** Spionidae

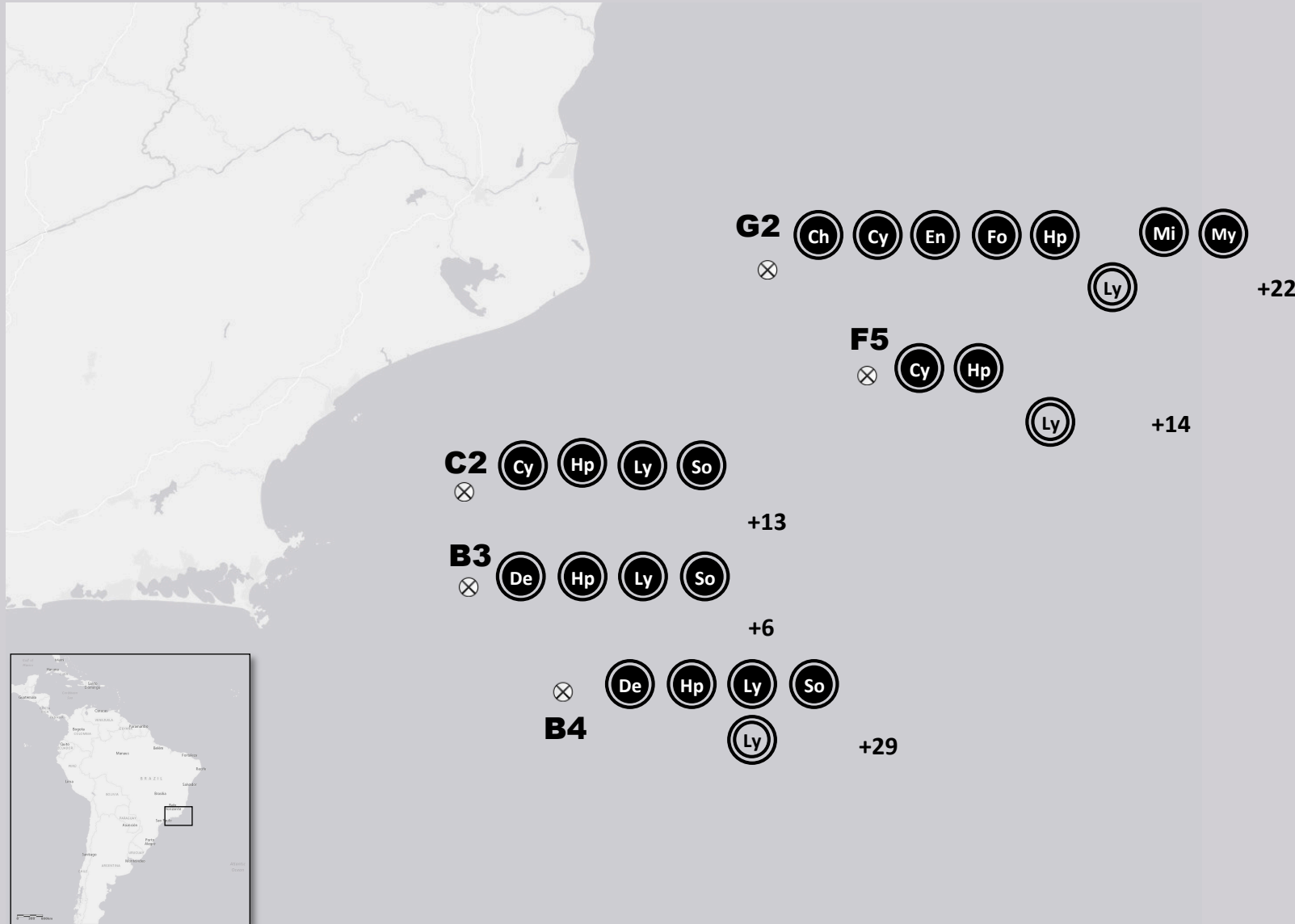
* Present in other stations of Habitats

** Non-marine family

Arthropoda distribution

Molecular

Morphological



- Ch Ch Chalcididae***
- Cy Cy Cyndroleberididae**
- De De Desmosomatidae*
- En En Entomobryidae
- Fo Fo Formicidae***
- Hp Hp Hippolytidae*
- Ly Ly Lysianassidae
- Mi Mi Miridae***
- My My Mysidae**
- So So Solenoceridae**

* Present in other stations of Habitats Project

** Previous studies in Campos basin

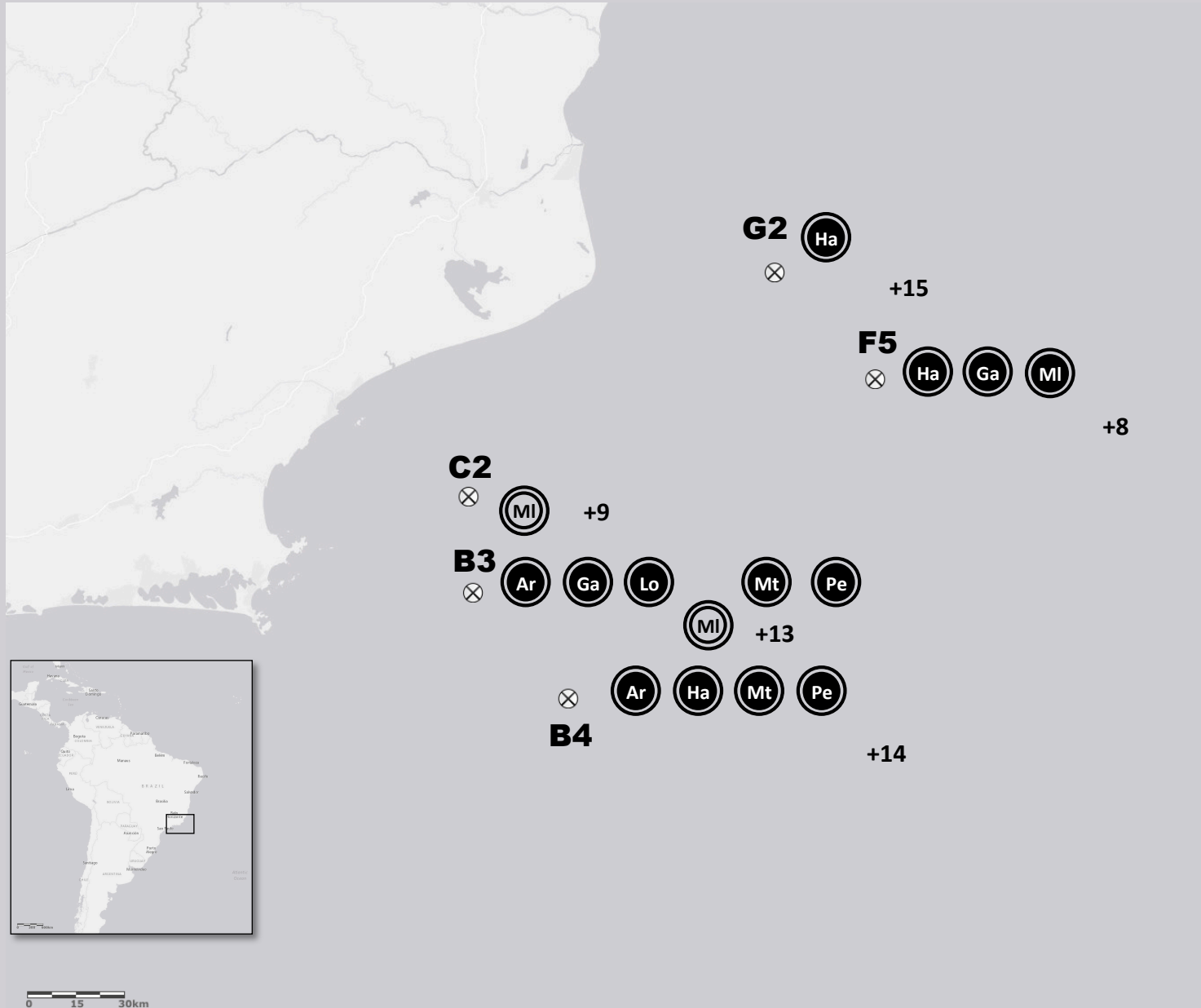
*** Non-marine family



Mollusca distribution

Molecular

Morphological



- Ar Ar Arcidae* **
- Ga Ga Galeommatidae**
- Ha Ha Haliotidae**
- Lo Lo Lottiidae**
- Mt Mt Mactridae* **
- MI MI Mytilidae
- Pe Pe Pectinidae* **

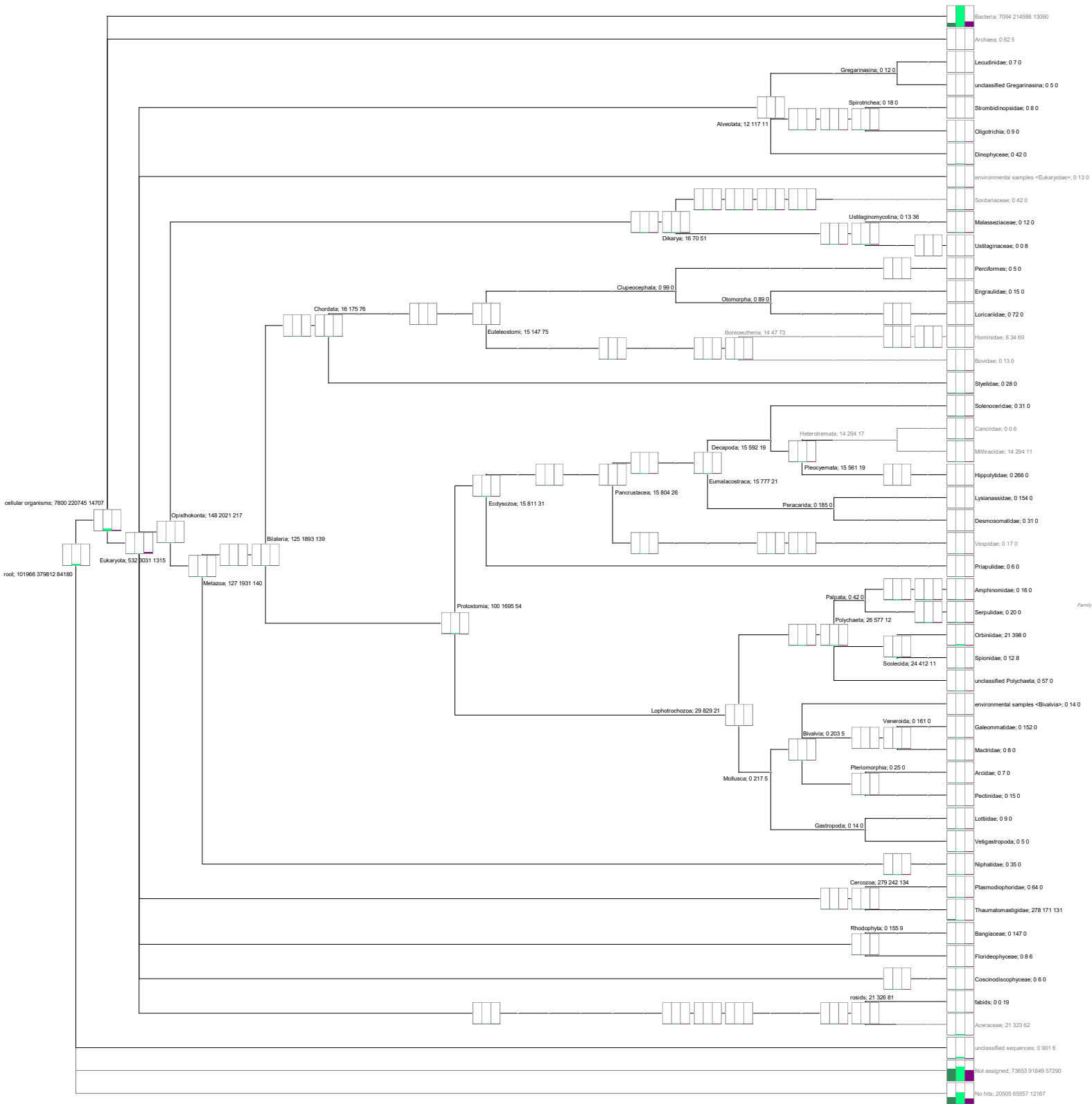
* Present in other stations of Habitats Project

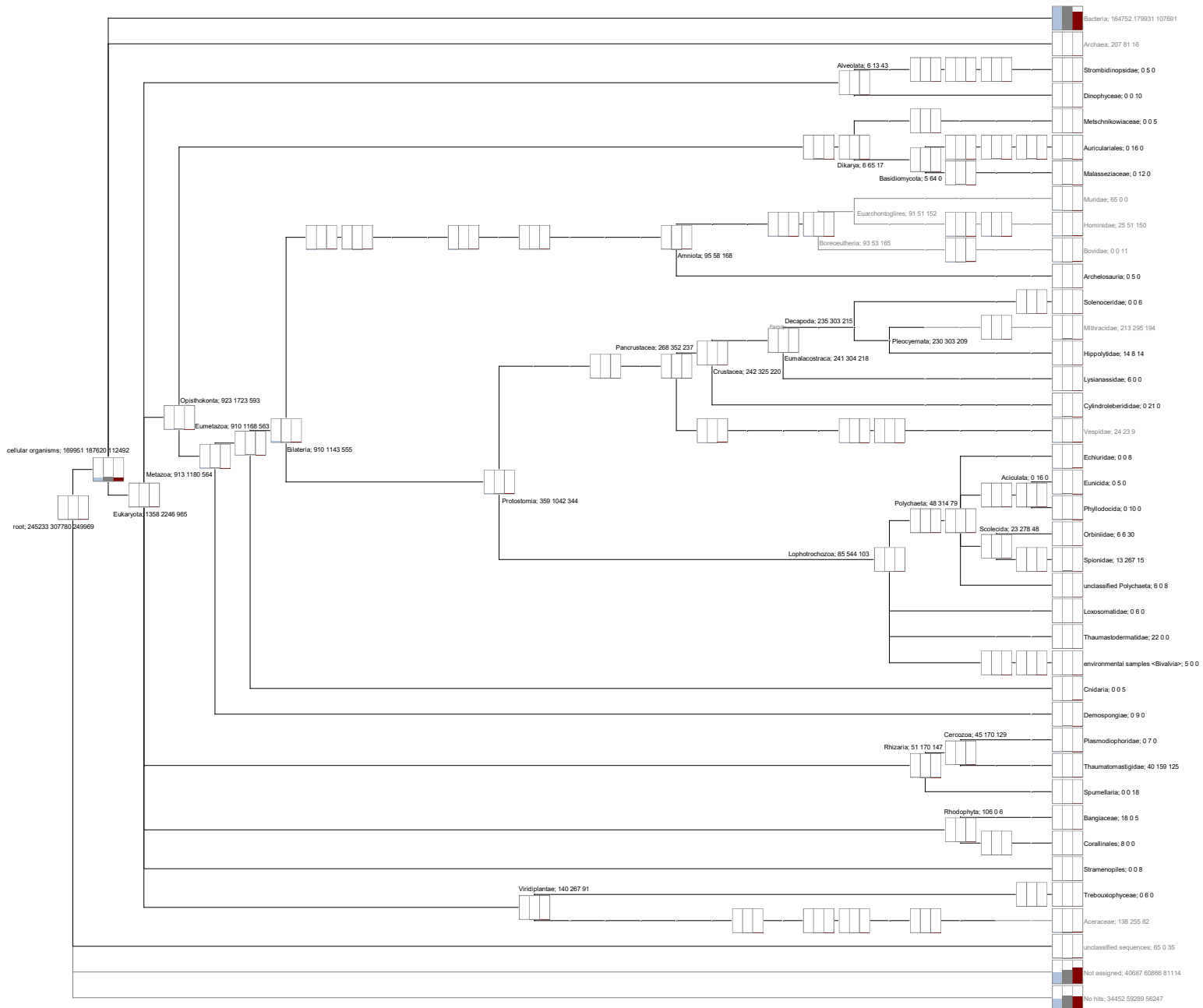
** Previous studies in Campos basin

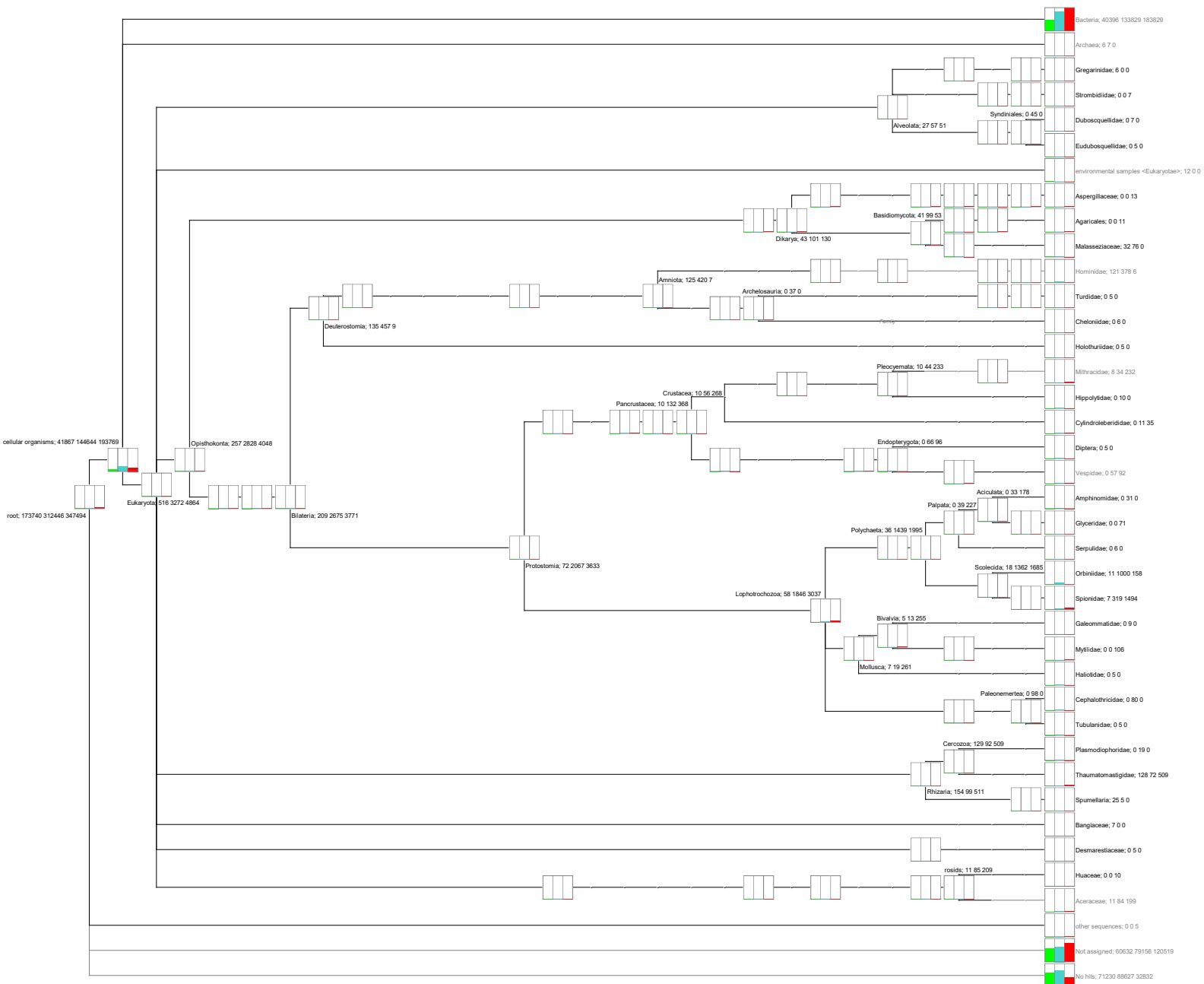
Suppelementar Material 1 – PCR primers and conditions. 1-5 µL of DNA template, 1 µL (5µM) of primers Forward and reverse), 5 µl of 10X buffer, 2 µl of MgCl₂ (25 mM), 1 µl of dNTP 10 µM (Fermentas), 0.2 µl de Platinum® Taq DNA Polymerase High Fidelity 5 U.µL-1 (Thermo Scientific) and ultra pure destilaed water (Invitrogen) to complete 50 µl final reaction volume.

Target	Primer (F – Forward; R – reverse)	Denaturation	cycles	denaturation	anealing	Extension	Final extension	References
COI	TITCIAAYCAYAARGAYATTGG (F – jLCO1490); TAIACYTCIGGRTGICRAARAAYCA (R – jHCO2198)	1' @94oC	10+30	30''@94oC	1'30''@61-52oC (-1oC per cycle) + 1'30''@61-52oC	1'@72oC	5'@72oC	Geller et al., 2013
rRNA 18S	ATGGTTGCAAAGCTGAAC (F – a2.0); GATCCTCCGCAGGTTACCTAC (R- 9R)	2' @94oC	40	30''@94oC	30'@55oC	1'@72oC	5'@72oC	Whiting et al., 1997; Whiting, 2002
rRNA 28S	ACCCGCTGAATTTAAGCAT (F – C1'); TGA ACTCTCTCTCAAAGTTCTTTTC (R- C2)	2' @94oC	40	30''@94oC	30'@55oC	1'@72oC	5'@72oC	Van Le et al., 1993; Chen et al., 2003

Supplementar material 2 – Family level Cladograms of the 5 sampling stations. Cladograms were built using specimens identified with any of the 3 target genes. Bar inside the squares represent the number of reads from each gene used to create the node. A) Family cladogram for station B3; b) Family cladogram for station B4; C) Family cladogram for station C2; D) Family cladogram for station G2; E) Family cladogram for station F5.







Supplementar material 3 – List of species identified by molecular and morphological taxonomy.

Specie	18S	20S	COI	Study
<i>Cnemidocarpa verrucosa</i>	+	+	+	Schettini
<i>Desmarestia dudresnayi</i>	+	+	+	Schettini
<i>Erythrophyllum delesserioides</i>	+	+	+	Schettini
<i>Eurythenes gryllus</i>	+	+	+	Schettini
<i>Galeomma turtoni</i>	+	+	+	Schettini
<i>Grifola frondosa</i>	+	+	+	Schettini
<i>Haliotis diversicolor</i>	+	+	+	Schettini
<i>Hormogaster redii</i>	+	+	+	Schettini
<i>Lysmata seticaudata</i>	+	+	+	Schettini
<i>Malassezia globosa</i>	+	+	+	Schettini
<i>Marenzelleria arctica</i>	+	+	+	Schettini
<i>Mimachlamys varia</i>	+	+	+	Schettini
<i>Mysidium columbiae</i>	+	+	+	Schettini
<i>Parotocinclus maculicauda</i>	+	+	+	Schettini
<i>Pinctada imbricata</i>	+	+	+	Habitats and Hits
<i>Platynereis dumerilii</i>	+	+	+	Habitats and Hits
<i>Pontocaris lacazei</i>	+	+	+	Habitats
<i>Praxillella affinis</i>	+	+	+	Habitats
<i>Progoniada regularis</i>	+	+	+	Habitats and Hits
<i>Protodorvillea kefersteini</i>	+	+	+	Habitats
<i>Pteria colymbus</i>	+	+	+	Habitats
<i>Scalibregma inflatum</i>	+	+	+	Habitats and Hits
<i>Scapharca broughtonii</i>	+	+	+	Schettini
<i>Serpula vermicularis</i>	+	+	+	Schettini
<i>Syllis gracilis</i>	+	+	+	Habitats and Hits
<i>Syllis variegata</i>	+	+	+	Habitats and Hits
<i>Travisia brevis</i>	+	+	+	Habitats and Hits
<i>Travisia forbesii</i>	+	+	+	Habitats and Hits
<i>Travisia pupa</i>	+	+	+	Habitats and Hits
<i>Aglaophamus circinata</i>		+	+	Habitats and Hits
<i>Alpheus formosus</i>		+	+	Habitats
<i>Amphipholis squamata</i>		+	+	Habitats
<i>Aricidea wassi</i>		+	+	Habitats and Hits
<i>Chelonia mydas</i>		+	+	Schettini
<i>Praxillella pacifica</i>		+	+	Habitats and Hits
<i>Priapulus caudatus</i>		+	+	Schettini
<i>Scolelepis bonnierii</i>		+	+	Schettini
<i>Scolelepis foliosa</i>		+	+	Schettini
<i>Amphimedon queenslandica</i>	+		+	Schettini
<i>Axiothella rubrocincta</i>	+		+	Habitats and Hits
<i>Bathycarca pectunculooides</i>	+		+	Habitats
<i>Bathyglycinde profunda</i>	+		+	Habitats
<i>Bathyglycinde sibogana</i>	+		+	Habitats
<i>Caprella equilibra</i>	+		+	Habitats and Hits
<i>Ceratocephale abyssorum</i>	+		+	Habitats and Hits
<i>Ciona intestinalis</i>	+		+	Schettini
<i>Clymenella torquata</i>	+		+	Habitats and Hits
<i>Pectinaria granulata</i>	+		+	Schettini
<i>Perna viridis</i>	+		+	Schettini
<i>Protaspis grandis</i>	+		+	Schettini
<i>Syllis hyalina</i>	+		+	Habitats and Hits
<i>Didemnum candidum</i>			+	Schettini
<i>Leodamas rubra</i>			+	Habitats and Hits
<i>Leodia sexiesperforata</i>			+	Habitats
<i>Leptocheilia dubia</i>			+	Habitats
<i>Leucothoe urospinosa</i>			+	Habitats and Hits
<i>Lumbrineris latreilli</i>			+	Habitats and Hits
<i>Lysidice ninetta</i>			+	Habitats and Hits
<i>Lysmata anchisteus</i>			+	Schettini
<i>Macrochaeta clavicornis</i>			+	Habitats
<i>Marphysa bellii</i>			+	Habitats and Hits
<i>Mendicula ferruginosa</i>			+	Habitats and Hits
<i>Mooreonuphis pallidula</i>			+	Habitats and Hits
<i>Neanthes acuminata</i>			+	Habitats and Hits
<i>Nereimyra punctata</i>			+	Habitats and Hits

Specie	18S	20S	COI	Study
<i>Notomastus latericeus</i>			+	Habitats and Hits
<i>Ophelina acuminata</i>			+	Habitats and Hits
<i>Pyropia haitanensis</i>			+	Schettini
<i>Scapharca kagoshimensis</i>			+	Schettini
<i>Scoloplos armiger</i>			+	Schettini
<i>Isolda pulchella</i>			+	Habitats and Hits
<i>Apophlaea lyallii</i>	+	+		Schettini
<i>Chaetoceros curvisetus</i>	+	+		Schettini
<i>Coelomactra antiquata</i>	+	+		Schettini
<i>Crassinella lunulata</i>	+	+		Habitats and Hits
<i>Cryptococcus friedmannii</i>	+	+		Schettini
<i>Cyclaspis alba</i>	+	+		Habitats
<i>Cylichna alba</i>	+	+		Habitats and Hits
<i>Engraulis japonicus</i>	+	+		Schettini
<i>Euclymene oerstedii</i>	+	+		Habitats and Hits
<i>Eulalia viridis</i>	+	+		Habitats and Hits
<i>Eumida sanguinea</i>	+	+		Habitats and Hits
<i>Exogone dispar</i>	+	+		Habitats and Hits
<i>Galathowenia oculata</i>	+	+		Habitats
<i>Glycera americana</i>	+	+		Habitats and Hits
<i>Glycera southeastatlantica</i>	+	+		Habitats and Hits
<i>Goniada emerita</i>	+	+		Habitats
<i>Hesiospina aurantiaca</i>	+	+		Habitats and Hits
<i>Patelloida striata</i>	+	+		Schettini
<i>Scopelocheirus schellenbergi</i>	+	+		Schettini
<i>Subulatomonas tetraspora</i>	+	+		Schettini
<i>Ophelina cylindricaudata</i>		+		Habitats and Hits
<i>Ophiactis lymani</i>		+		Habitats
<i>Trypanosyllis zebra</i>		+		Habitats and Hits
<i>Ahnfeltiopsis leptophylla</i>	+			Schettini
<i>Crucigera zygophora</i>	+			Schettini
<i>Leitoscoloplos pugettensis</i>	+			Schettini
<i>Malassezia nana</i>	+			Schettini
<i>Ophiura ljunghmani</i>	+			Habitats
<i>Owenia fusiformis</i>	+			Habitats and Hits
<i>Panthalis oerstedii</i>	+			Habitats and Hits
<i>Paralacydonia paradoxa</i>	+			Habitats and Hits
<i>Paramphinome jeffreysii</i>	+			Habitats and Hits
<i>Pholoe minuta</i>	+			Habitats
<i>Phtisica marina</i>	+			Habitats
<i>Phyllodoce longipes</i>	+			Habitats and Hits
<i>Solenocera crassicornis</i>	+			Schettini
<i>Strombidium paracalkinsi</i>	+			Schettini
<i>Phagomyxa odontellae</i>	+			Schettini