1    **Metagenomics Accelerates Species Discovery and Unravel Great Biodiversity of Benthic**

2    **Invertebrates in Marine Sediment in Campos Basin, Brazil**

3

4    Milena MDP Schettini, Raony G C C L Cardenas, Marcella A A Detoni, Mauro F Rebelo.

5

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

7    Rio de Janeiro. Brasil.

8

9    **ABSTRACT**

10    Sediment fauna characterization and monitoring are mandatory requirements for obtaining oil

11    and gas (O&G) environmental licensing for exploration and production (E&P) activities.

12    Currently, for environmental characterizations and monitoring, biodiversity is assessed through

13    morphological taxonomy, a time-consuming process. Taxonomists are constantly failing to meet

14    the demands for biodiversity assessment required in monitoring programs. Thus, we combined

15    three different phylogenetic markers (rRNA 18S, rRNA 28S and COI), HTS and Bioinformatics to

16    identify benthic invertebrate organisms from sediment samples collected in five stations in the

17    Campos Basin in southeast Brazil, an important oil extraction area and one of the best-studied

18    marine biota in Brazil. Our results obtained with metagenomics were compared to morphology

19    data provided by the Habitats Project whereas the database *Global Biodiversity Information*

20    *Facility* was used for organism localization. We obtained around 4.83 µg of DNA from 15

21    samples. A total of 3.3 million sequences were clustered in Operational Taxonomic Units and

22    more than 1.6 million sequences (about 50% of all reads) were assigned to 957 prokaryotes and

23    577 eukaryotes. BLAST identified 23 phyla, 60 classes, 62 orders, 70 families, 67 genus and 46

24    species of eukaryotes. Our metagenomic analysis identified phyla that are traditionally found in

25    samples of marine benthos, such as Annelida, Arthropoda, Mollusca and Chordata, as well as

26    more rarely found phyla such as Bryozoa, Cnidaria, Echinodermata, Nematoda, Nemertea,

27    Platyhelminthes, Porifera and Priapulida; and even more rare phyla like Entoprocta and

Peer**J** Preprints

28    Gastrotricha. The low availability of genetic markers for Brazilian species in Genbank impaired

29    our ability to compare our findings with those obtained morphologically for which no sequences

30    were found in Genbank. Our study shows that metagenomics can be applied for environmental

31    characterization and monitoring programs and, with the possibility of automating the method,

32    may reduce from years to few months the time currently required for species identification and

33    biodiversity determination, which will certainly accelerate species discovery.

34

35    **INTRODUCTION**

36    Sediment fauna characterization and monitoring are mandatory requirements for obtaining oil

37    and gas (O&G) environmental licensing for exploration and production (E&P) activities. This

38    requirement is expected to remain a key element of environmental management in the future,

39    particularly in the frontiers of deep-sea offshore oil exploration areas: the Equatorial Margin and

40    Santos Pre-salt Basin in Brazil, or the Barents and Siberia seas in the Arctic Ocean.

41    Currently, for environmental characterizations and monitoring, biodiversity is assessed through

42    morphological taxonomy, a time-consuming process. As a general rule, taxonomic resolution at

43    species level is expected and for some fauna groups, the expertise required is so unique that only

44    a hand full of individuals in the world is fit for the task. Still, expert judgment is never 100%

45    accurate, with only 50% rate of identification success being shared among taxonomists

46    (Culverhouse et al., 2003). At last, fragile organisms that require special fixation procedures may

47    not be properly represented in the samples (Costa-Paiva; Paiva e Kautau, 2007). As a result,

48    invertebrate morphological identification efforts are often limited to few groups, including

49    Mollusca, Crustacea and Polichaeta (MMA, 2015) and some estimates suggest that more than

50    90% of all marine species have never been named (SCHEFFERS et al, 2012).

51    The typical number of sediment samples in a monitoring campaign is in the range of tenths, but

52    new areas to be explored can be as large as 300.000 km2, which can result in tenths of

53    thousands of samples for baseline environmental characterization.  Taxonomists are constantly

54    failing to meet the demands for biodiversity assessment required in monitoring programs. The

55 lack of experts is a major bottleneck in the process of identifying biodiversity (HEBERT et al,

56 2003; MORA; ROLLO; TITTENSOR, 2013), which delays operators execution of E&P projects to

57 reach 'first oil' and keep species from being identified.

58 In Brazil, which, according to the latest Report of the Convention on Biological Diversity (CBD,

59 2016), is the most biologically-diverse country, with more than 100,000 animal species been

60 accounted for, only 184 marine invertebrates had their conservation status accessed (MMA,

61 2015). It is possible that current risk estimates of environmental impact are based on

62 underestimated biodiversity inventories, representing a threat to species conservation.

63 Developing new technologies and approaches that accelerate species discovery and reveal

64 hidden biodiversity is crucial for setting conservation priorities and efforts.

65 Molecular methods use big data generated through high-throughput sequencing (HTS), which

66 greatly accelerates species discovery. This approach is particularly useful for marine sediment

67 analyses because the higher possibility of identifying minute organisms belonging to groups

68 such as Nematoda, Copepoda, Ostracoda, Rotifera, Kinohyncha, Loricifera, Tardigrada and of

69 species from debris and other sorts of environmental DNA (WANG et al, 2014), if compared to

70 morphology. To classify eukaryote organisms using DNA-based approaches, and which have not

71 yet been described morphologically, the concept of Operational Taxonomic Unit (OTU) can be

72 applied (Schmidt; Mafias Rodrigues; Von Mering, 2014; Stackebrandt; Goebel, 1994).

73

74 Since 2010, more than 600 papers have been published on the use of DNA-based identification

75 methods for species conservation (Bergman et al., 2016; Goldberg et al., 2014), biodiversity

76 inventory determination (Drummond et al. 2015); environmental monitoring (Bowman et al.,

77 2014; Brown et al., 2015; Chariton et al., 2015; Leray et al., 2015), DNA extraction/detection

78 (Eichmiller et al., 2014; Pedersen MW et al., 2015; Ficetola et al., 2016) and the technique has

79 been considered a major tool for Ocean's sustainability in the 21st century (Aricó, 2015).

80 In this study, we combined three different phylogenetic markers (rDNA 18S, rDNA 28S and COI),

81 HTS and Bioinformatics to identify benthic invertebrate organisms with metagenomes from

82   sediment samples collected in Campos Basin in southeast Brazil, an important oil extraction area

83   and one of the best-studied marine biota in Brazil (MILOSLAVICH et al, 2011).

84

85   **Material and Methods**

86   Sample collection and processing:

87   Samples were collected at Campos Basin in 2009 as part of 'Habitats Project – Campos Basin

88   Environmental  Heterogeneity' coordinated by CENPES/PETROBRAS. Table 1 presents

89   information (collection date, geographic coordinates and depth) on the five sampling stations:

90   B3, B4, C2, G2 and F5. Sediment samples were collected in triplicate, descending a Van Veen grab

91   in three different points around (150 m radius) each of the five stations, totaling 15 sediment

92   samples. At the time these samples were collected, no plans to have them genetically analyzed

93   had been set. Thus, they were kept at -20°C for 4 years until our analysis was done in 2013.

94   For each station, we manually homogenized 200 cm$^3$ of the muddy sediments and weighted 5g

95   for DNA extraction that was performed using the PowerMax Soil DNA Isolation (MoBio Inc),

96   according to manufacturer's instructions. DNA integrity was accessed by means of agarose gel

97   1.2%. Quantification was performed in Qibit 2.0 Fluorometer (Life Technologies).

98

99   Biogeography data:

100  Data on the organisms identified in this study were extracted from two main sources: the book

101  entitled "Biodiversidade bentônica da região central da Zona Econômica Exclusiva brasileira" by

102  Lavrado and Ignacio (2006) for the Cnidaria Crustacea, Echinodermata, Mollusca, Nematoda,

103  Polychaeta and Porifera groups, whereas the dada for organisms of the phyla Annelida,

104  Arthropoda, Brachiopoda, Bryozoa, Cnidaria, Echinodermata, Echiura, Foraminifera,

105  Haptophyte, Mollusca, Nematoda, Nemertea, Porifera, Priapula, Protozoa, Rodophyta were

106  identified by the Habitats Project and provided by Petrobras S.A. (unpublished data).

107  We also used the database *Global Biodiversity Information Facility* (www.gbif.org) for organism

108  localization.

109    In this study, we chose family as the taxonomic group to be used as reference in cladograms in

110    order to be able to compare our findings with those provided by morphological taxonomy.

111    Whenever species descriptions were available for both metagenomic and morphological

112    approach, they were also discussed.

113

114    PCR and high-throughput sequencing:

115    Information on PCR of COI, rDNA 18S and rDNA 28S genes is presented in Supplemented

116    material 1. We used the kit *Ion Xpress™ Plus Fragment Library* (Life Technologies) for preparing

117    the libraries for sequencing according to manufacturer's instructions of *Ion Xpress™ Plus gDNA*

118    *Fragment Library Preparation.* Template preparation and sequencing were done using the kit

119    Ion PGM™ Template OT2 400. Sequencing was done using the *Ion Personal Genome Machine*

120    *(PGM™) System* at the Life Technologies laboratories (São Paulo, SP), using *Chip* 318 v2.

121    Sequencing adapters were removed from reads using Torrent Suite *software* version 4.0.2 (Life

122    Technologies) and assigned to samples based on the combination primer tail-Ion Xpress

123    barcode. Prinseq version 0.20.4 (SCHMIEDER; EDWARDS, 2011) was used to remove either A/T

124    photopolymers bigger than 5 bases, reads with unidentified (N) bases, small length (<80bp) or

125    bad quality reads (Q<20). Remaining reads were clustered in OTUs using CD-HIT-EST version

126    4.6 (LI; GODZIK, 2006) (up to 97% identity under 100% coverage within a bigger read, word

127    size of 10 and 20 penalty points for gaps).

128    High quality and low redundancy sequences were compared to NCBI non-redundant nucleotide

129    repositories (NR) (http://www.ncbi.nlm.nih.gov/genbank/) using *Basic Local Alignment Search*

130    *Tool nucleotides* (BLASTn) version 2.3.0+ (Zhang et al, 2000). Max *e-value* was of $10^{-5}$ and the

131    number of events per query was limited to 100 (here called as *hits*).

132    Taxonomic names were attributed to each *read,* based on the reads group of BLAST hits, using

133    the 'Lowest Common Ancestor Assignment – LCA' algorithm in software MEGAN (MEta Genome

134    Analyzer; version 5.10.3; Huson et al., 2007) according to different parameters (Huson et al.,

Peer J Preprints

135    2011). Cladograms and rarefaction curves at family taxonomic level for each station were also

136    built using MEGAN.

137    The BLAST step was performed using the Elastic Compute Cloud (EC2) service of Amazon

138    (aws.amazon.com). The BLAST for each of the 15 sets of reads correspondent to the 15 samples,

139    run in a parallel scheme using eight threads on up to 96 AWS instances with 8 processors and 16

140    Gb of RAM each.

141

142    **RESULTS & DISCUSSION**

143    We obtained an average of 4.83 μg of DNA from each of the 15 samples. Sequencing generated

144    approximately 4.8 million sequences, which is within the expected values for the 318 v2 chip,

145    but with an average size of 155.1 bp, which is bellow the expected value for the OT2 400 kit.

146    Over 3.6 million sequences (75.35%) passed quality control and of these, around 3.3 million

147    were clustered in Operational Taxonomic Units by CD-HIT (Table 2). More than 1.6 million

148    sequences (about 50% of all reads) were assigned to 957 prokaryotes and 577 eukaryotes using

149    BLAST (Table 2). BLAST identified 23 phyla, 60 classes, 62 orders, 70 families, 67 genus and 46

150    species of eukaryotes (Supplementary Material 2 – Cladograms and Supplementary Material 3 –

151    list of species identified). Figure 1 shows the rate of OTU observed by metagenomics in each of

152    the stations distributed over the 13 invertebrate phyla (Figure 1A) and 38 invertebrate families

153    (Figure 2B). All other Prokaryota and non-invertebrate Eukaryota phyla observed in this study

154    are listed in the cladograms available in the supplementary material. A considerable number of

155    reads were assigned to the families Hominidea and Bovidae, increasing the number of reads

156    belonging to the Chordate phylum. However, these were read alignments generated against the

157    whole human and bovine genomes or chromosomes, as opposed to the three genetic markers

158    that we used in this study. Our results and discussion are focused on invertebrate families

159    belonging to marine benthos and no artifact findings on chordate will be further addressed.

160    One of the differentials of our study was that it was done using samples collected from the actual

161    areas were E&P activities are usually carried out. Several previous morphological taxonomic

162   studies were performed in these areas, either by the oil companies interested in obtaining their

163   licenses, or those involved in conservational programs (such as the Habitats Project) or by the

164   scientific community (the REVIZEE program).

165   The huge taxonomic effort of the Habitats Project generated a databank of 49,289 specimens. A

166   total of 17 phylum, 27 classes, 63 orders, 354 families, 768 genus and 749 species were

167   identified.

168   Out of the 1,773 macroinvertebrates *taxa* identified by morphological taxonomy,  1,211 or 68%

169   did not have any entry in Genbank found for any of the three markers (COI, rRNA 18S e 28S)

170   used in this study, indicating that Brazilian marine species are underrepresented in Genbank.

171   Thus, there is a need to increase efforts to have sequences from these three molecular markers

172   from more Brazilian species deposited in Genbank, as the limited number of sequences impairs

173   any parallel to be done between the findings obtained with molecular and those obtained with

174   morphological taxonomies.

175   Our metagenomic analysis identified phyla that are traditionally found in samples of marine

176   benthos, such as Annelida, Arthropoda, Mollusca and Chordata, as well as more rarely found

177   phyla such as Bryozoa, Cnidaria, Echinodermata, Nematoda, Nemertea, Platyhelminthes, Porifera

178   and Priapulida; and more rare phyla like, Entoprocta and Gastrotricha (Supplementary material

179   and Figure 1).

180   The great number of OTUs for Annelida, Arthropoda and Mollusca found by metagenomics

181   agrees with previous results for Campos Basin found by LAVRADO; IGNACIO, 2006 during the

182   REVIZEE project and also by those of the Habitats Project. Recent metagenomics study carried

183   out by Leray and Knowlton (2015) also identified Annelida and Arthropoda as the phyla with

184   more OTUs among the 22 phyla identified from approximately 0.09 m$^3$ sediments from coral reef

185   regions in Virginia and Florida, in the United States.

186   The Entoprocta (or Kamptozoa) phylum comprises about 170 aquatic and sessile species of sizes

187   between 0.5 and 5.0 mm and are mostly marine (Zhang, 2011). Until 2011, only 18 species of

188   Entoprocta were known on the Brazilian coast (Vieira; Migotto, 2011). In this study, all OTUs (6

189    in the C2 station and 24 in the G2 station) were attributed to the genus *Loxosomella* through the

190    marker rDNA 28S, with over 86% of sequence similarity. This result expands the distribution of

191    the genus that was previously limited to six species collected off the coast of São Paulo (VIEIRA;

192    MIGOTTO, 2011).

193    As for the cosmopolitan Gastrotricha phylum that comprises about 790 species of aquatic

194    organisms up to 1 mm in length (Zhang, 2011), all 22 OTUs assigned to the phylum (C2 station)

195    were in the *Tetranchyroderma* genus, with over 81% similarity with COI sequences found in the

196    Genbank. This occurrence also expands the distribution that had been previously limited to São

197    Paulo beaches (reported but not formally described – Garraffoni; ARAUJO, 2010), almost a

198    1000km away from the Campos Basin.

199    This is a pioneer study in which metagenomics results could be compared to those from a recent

200    comprehensive morphological taxonomy effort that worked with the same samples than those

201    used in our study. However, comparing results between studies should be taken with caution

202    because of the uncertainty on how much DNA is still available considering that samples have

203    been preserved at -20$^o$C for 40 years and the lack of available genetic markers for the Brazilian

204    marine species in the Genbank. It should also be noted that we analyzed 5g out of 200 gr of the

205    surface (0 to 2 cm) sediment for each of the 15 samples, while the morphological study worked

206    with 1000 cm$^3$ of sediment from each sample, comprising slices from 0 to 10 cm. Finally, for

207    many species, the sequences of the markers available in Genbank were partial and thus we

208    cannot ensure they properly aligned with the reads to attribute a taxonomic name. However,

209    these restrictions applies mostly to the families that we did not found and we believe that

210    observations made about the families that we actually found are valid.

211    Our analysis identified 38 families of invertebrates in the 15 samples from the 5 sampling

212    stations in Campos Basin. Figure 2 compares between the families from Annelida (9 families, fig.

213    2A), Arthropoda (10 families, fig. 2B) and Mollusca (7 families, fig. 2C) phyla identified by

214    metagenomics and morphology taxonomy in stations B3, B4, C2, F5 and G2.

215   Annelida families Amphinomidae, Enchytraeidae, Glyceridae, Orbiniidae, Serpulidae and

216   Spionidae were previously identified in Campos Basin by the Habitats Project while up to 28

217   annelida families previously reported by the Habitat project could not be identified by

218   metagenomics.  Family Hormogastridae found in our study is most likely a false positive since it

219   is not marine. The Arthropoda families Solenoceridae, Cylindroleberididae and Mysidae have

220   been previously identified in Campos Basin and in the Southeast of Brazil by other authors

221   (CARDOSO; SEREJO, 2007; GBIF, 2016; SEREJO et al, 2007; TÂMEGA; OLIVEIRA; FIGUEIREDO,

222   2013) while up to 29 arthropoda families previously reported by the Habitat Project were not be

223   identified by metagenomics. Families Miridae, Chalcididae and Formicidae found in our study

224   are most likely false positive since they are not marine insects. All Mollusca families identified by

225   metagenomics in Campos basin, except for Mytilidae, have not been identified by the Habitat

226   Projects, even though they have been previously found in the region (DORNELLAS; SIMONE,

227   2011; LAVRADO; IGNACIO, 2006; TÂMEGA; OLIVEIRA; FIGUEIREDO, 2013). Up to 15 mollusca

228   families previously reported by the Habitat Project could not be identified by metagenomics.

229   Moreover, metagenomics was also able to find several families not previously reported by

230   morphological taxonomy for a given station, suggesting that the family distribution could be

231   broader than anticipated. That is the case for Echiuridae, Hormogastridae and Pectinariidae

232   among the Annelidae; Desmosomatidae and  Hippolytidae for Arthropoda and Arcidae,

233   Mactridae and Pectinidae for Mollusca.

234   The Habitat Projects identified 749 organisms to the species level but only 64 had at least one

235   sequence of one of the three genetic markers (COI gene, 18S rRNA and 28S here studied)

236   deposited in Genbank and thus were 'eligible' for molecular identification. At first,

237   metagenomics identified 46 species. However, none of the 64 species previously identified by

238   morphological taxonomy by the Habitat Project were found by metagenomics. We believe that

239   these are false negative results that can be explained by the fact that samples were preserved at

240   -20$^{o}$C for 4 years, by the low sample volume and the fact that the genetic markers here studied

241   were missing in Genbank for a number of organisms that were identified by morphology.

242 However, we noticed that even after calibration of the parameters for the LCA algorithm (data

243 not shown), some incongruence in the attribution of the name of the species had happened. To

244 overcome that limitation, we manually searched for the 64 species names of found by the

245 Habitat Project among the names of the organisms generated by the BLAST hits for a given read.

246 We were able to identify more 45 species that had been previously described by morphological

247 taxonomy but were not picked by the LCA algorithm. The full list of species identified by

248 molecular and morphological taxonomies, together with the genetic markers available in

249 Genbank are listed in supplementary material 3 (or table). Other false negative results could

250 have been generated by the occurrence of synonymous names at the species level. For instance,

251 according to recent estimates, more than 80% of the algae of some genus and 38% of mollusca

252 have synonymous names. For marine species, this figure would reach 40% (COSTELLO; MAY;

253 STORK, 2013). An ongoing effort is dedicated to resolve synonymous names found in the GBIF

254 database.

255 Of the 46 invertebrate present in cladograms leaves (most specific possible position) that we

256 identified by the molecular taxonomy, 27 were invertebrates not previously described in the

257 region. These could represent new occurrence in Campos Basin. Because description in the main

258 biogeographic databases that we used (Habitats Project, Revizee and GBIF) usually goes no

259 further than the family taxonomic level, it is not possible to either claim or rule out that the

260 finding corresponds to the first description of the species in the region.

261 However, we wanted to calculate the likelihood that those events truly represented false

262 positive results, as oppose of being descriptions of new species. False positive results could

263 happen as an artifact due to similarities of genetic sequences shared among species belonging to

264 the same genus and the low representativeness of Brazilian species in Genbank. The high

265 similarity could have led BLAST to relate, with very low error probability, a read from one

266 species not present in the databank to another present in the Genbank and from the same genus

267 (phylogenetic similarity) but belonging to a completely different habitat. By using metadata on

268 the distribution of the species selected by BLAST, we managed to sort out at least one case

269 among our results. *Haliotis diversicolor*, which was identified in our study, is a small (25-85 mm)

270 gastropod form the Indo-Pacific Ocean, with georeferenced records on the coast of Japan,

271 Thailand, Australia, among other countries in the region (GBIF, 2016). Despite the geographical

272 distance, *Haliotis diversicolor* shares high sequence similarity with *H. aurantium*, which has been

273 identified in the Campos basin, and also with three other records corresponding to species found

274 in the Brazilian coast. The lack of genetic markers for these Brazilian species in Genbank may

275 have misled BLAST searches, which in this case erroneously classified the sequence of *H.*

276 *aurantium* as of *Haliotis diversicolor.*

277 To further remove false positive results, we wanted to find redundant identification done by

278 each of the three genetic markers for each of the 46 species found by molecular taxonomy,

279 hoping that a doubtful identification by one marker could be resolved by a positive confirmation

280 by the other two. Unfortunately, that was not the case. Out of the 64 species identified by the

281 Habitat Project, 16 species had sequences of all three markers available in Genbank and still

282 were not positively identified by metagenomics. Out of the 46 species identified by molecular

283 metagenomics, other 16 had sequences of all three genetic markers available in Genbank, but by

284 metagenomics they were identified only by one of the three markers and never by two or three.

285 We noticed that many times, even though the sequence for a genetic marker for a specific

286 organism was available in the Genbank, multiple names were attributed to the gene, only partial

287 sequences were available, or sequences were not validated experimentally. Genbank is the best

288 repository for genetic sequences yet available but still does not offer a high level of confidence

289 when it comes to the names attributed to genetic sequences. Our research team is currently

290 working on developing new algorithms to help overcome this limitation.

291 The problems related with having false positive and false negative results and with the

292 occurrence of synonymous names could be solved if we work only at the level of OTU to

293 compare taxa profile among samples seasonally. The frequency and abundance of OTUs could

294 then be related to environmental changes and could accelerate species discovery by showing

295 that genetic sequences vary according to environmental conditions. Further studies should be

296  done in which such strategy is adopted, as working with OTUs allows us to unravel the hidden

297  biodiversity of the thousands of 'no hit' OTUs and to relate their distribution to environmental

298  changes and activities.

299

300  **CONCLUSION**

301  Brazil has one of the strictest environmental laws and regulations for E&P activities of the O&G

302  sector in the world. Recent changes were made under resolution CONAMA 422/11 that

303  minimized bureaucracy required by the application process, increased transparency by sharing

304  information online and reduced liability for the O&G operators. In Brazil, the environmental

305  authority IBAMA (Brazilian Institute of the Environment and Renewable Natural Resources) is

306  responsible for issuing 'reference terms' that establish the guidelines and best practices for the

307  environmental licensing and monitoring.

308  Metagenomics can be applied for environmental characterization and monitoring programs and,

309  with the possibility of automating the method, may reduce from years to few months the time

310  currently required for species identification and biodiversity determination, which will certainly

311  accelerate species discovery.

312  Nevertheless, the fact that 68% of the organisms identified by morphology did not have

313  sequences of at least one of the three markers used in this study (COI, 18S and 28SrDNA)

314  deposited in the Genbank illustrates how low is the representation of molecular markers of

315  Brazilian marine species in the Genbank. Further studies should focus on sequencing organisms

316  and have their sequences deposited in the Genbank and other international databases.

317  We believe that metagenomic identification based on species' DNA overcomes several of the

318  limitations associated to morphological methodology. We have shown, as well as the studies

319  done by others, that metagenomics is a reliable approach for the identification of biodiversity,

320  that can be improved by adding more sequences of native species in public and proprietary

321  databanks. It is our opinion that metagenomics consists of the best available technique for

322    generating biodiversity inventories in marine sediments and should be acknowledged as such by

323    oil operators, environmental authorities and the scientific community at large.

324

325    **REFERÊNCIAS**

326    Aricò S. 2015. "Ocean Sustainability in the 21st Century" report. UNESCO Publishing /

327        Cambridge University Press . 324 pages.

328    Bergman PS, Schumer G, Blankenship S, Campbell E (2016) Detection of Adult Green Sturgeon

329        Using Environmental DNA Analysis. PLoS ONE 11(4): e0153500.

330        doi:10.1371/journal.pone.0153500

331    Bohmann et al., Environmental DNA for wildlife biology and biodiversity monitoring. Trends

332        Ecol. Evol. 29, 358–367 (2014).

333    Brown, Emily A; Chain, Frédéric J J; Crease, Teresa J; MacIsaac, Hugh J; Cristescu, Melania E. 2015.

334        Divergence thresholds and divergent biodiversity estimates: can metabarcoding reliably

335        describe zooplankton communities? Ecology and evolution vol. 5 (11) p. 2234-51

336    Cardoso, I. A.; Serejo, C. S. Deep Sea Caridea (Crustacea, Decapoda) from Campos Basin, Rj,

337        Brazil. **Brazilian Journal of Oceanography**, v. 55, n. 1, p. 39-50, 2007.

338    CBD, 2016 - https://www.cbd.int/countries/?country=br

339    Chariton et al., Metabarcoding of benthic eukaryote communities predicts the ecological condition of

340        estuaries. Environ. Pollut. 203, 165–174 (2015).

341    Costa-Paiva EE, Paiva PC and Klautau M. Anaesthetization and fixation effects on the morphology

342        of  sabellid polychaetes (Annelida: Polychaeta: Sabellidae). J. Mar. Biol. Ass. U.K. (2007),

343        87, 1127–1132.

344    Costello, M. J.; May, R. M.; Stork, N. E. Can We Name Earth's Species Before They Go Extinct?

345        **Science**, v. 339, n. 6118, p. 413-416, jan. 2013.

346    Culverhouse et al.: Do experts make mistakes? A comparison of human and machine

347        identification of dinoflagellates. Mar Ecol Prog Ser 247: 17–25, 2003.

348    Drummond et al. 2015. Evaluating a multigene environmental DNA approach for biodiversity

349        assessment. DOI: 10.1186/s13742-015-0086-1).

350    Eichmiller et al., 2014 The Relationship between the Distribution of Common Carp and Their

351        Environmental DNA in a Small Lake. DOI: 10.1371/journal.pone.0112611)

352    Ficetola G. F., P. Taberlet, E. Coissac, How to limit false positives in environmental DNA and

353        metabarcoding? Mol. Ecol. Resour. 16, 604–607 (2016).

354    Garraffoni, A. R. S.; Araújo, T. Q. Chave de identificação de Gastrotricha de águas continentais e

355        marinhas do Brasil. **Papéis Avulsos de Zoologia**, v. 50, n. 33, p. 535-552, 2010.

356    GBIF - Global Biodiversity Information Facility. Disponível em: <http://www.gbif.org/>. Acesso

357        em: 15 fev. 2016.

358    Goldberg C. S., K. M. Strickler, D. S. Pilliod, Moving environmental \{DNA\} methods from

359        concept to practice for monitoring aquatic macroorganisms. Biol. Conserv. 183, 1–3

360        (2015).

361    Hebert, P. D. N. et al. Biological identifications through DNA barcodes. **Proceedings of the Royal**

362        **Society of London Series B: Biological Sciences**, v. 270, p. 313-321, jan. 2003.

363    Huson, D. H. et al. Integrative analysis of environmental sequences using MEGAN4. **Genome**

364        **Research**, v. 21, n. 9, p. 1552-1560, set. 2011.

365    Huson, D. H. Et Al. MEGAN analysis of metagenomic data. **Genoma Research**, p. 1-10, jan. 2007.

366    Lavrado, H. P.; Ignacio, B. L. (Eds.). **Biodiversidade bentônica da região central da Zona**

367        **Econômica Exclusiva brasileira**. Rio de Janeiro: Museu Nacional - UFRJ, 2006.

368    Leray, M.; Knowlton, N. DNA barcoding and metabarcoding of standardized samples reveal

369        patterns of marine benthic diversity. **Proc. Natl. Acad. Sci. USA**, v. 112, n. 7, p. 2076-2081,

370        fev. 2015.

371    Li, W.; Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or

372        nucleotide sequences. **Bioinformatics**, v. 22, p. 1658-1659, 2006.

373    Miloslavich - Patricia Miloslavich1*, Eduardo Klein1, Juan M. Dı´az2, Cristia´n E. Herna´ndez3,

374        Gregorio Bigatti4, Lucia Campos5, Felipe Artigas6, Julio Castillo1, Pablo E. Penchaszadeh7,

375    Paula E. Neill8, Alvar Carranza9, Marı ´a V. Retana4, Juan M. Dı ´az de Astarloa10, Mirtha

376    Lewis4, Pablo Yorio4,11, Marı ´a L. Piriz4, Diego Rodrı ´guez10, Yocie Yoneshigue-

377    Valentin5, Luiz Gamboa12, Alberto Martı ´n1. Marine Biodiversity in the Atlantic and

378    Pacific Coasts of South America: Knowledge and Gaps. PLoS ONE, January 2011 | Volume 6

379    | Issue 1 | e14631.

380    MMA - Ministério do Meio Ambiente. Fifth National Report to the Convention on Biological

381    Diversity: Brazil. Brasília: Ministry of the Environment, 2015.

382    Mora, C.; Rollo, A.; Tittensor, D. P. Comment on "Can We Name Earth's Species Before They Go

383    Extinct?" **Science**, v. 341, n. 6143, p. 237, 2013.

384    Pedersen M. W.  et al., Ancient and modern environmental DNA. Philos. Trans. R. Soc. London B

385    Biol. Sci. 370 (2014), doi:10.1098/rstb.2013.0383.

386    Scheffers, B. R. et al. What we know and don't know about Earth's missing biodiversity. **Trends**

387    **in Ecology and Evolution**, v. 27, n. 9, p. 501-510, set. 2012a.

388    Scheffers, B. R. et al. Erratum to: ''What we know and don't know about Earth's missing

389    biodiversity''. **Trends in Ecology and Evolution**, v. 27, n. 12, p. 712-713, dez. 2012b.

390    Schmidt T. S., Matias Rodrigues J. F., Von Mering C. Ecological consistency of SSU rRNA-based

391    operational taxonomic units at a global scale. **PLOS Computational Biology**, v. 10, n. 4, p.

392    1-10, abr. 2014.

393    Schmieder, R.; Edwards, R. Quality control and preprocessing of metagenomic datasets.

394    **Bioinformatics**, v. 27, p. 863-864, 2011.

395    Serejo, C. S. et al. Abundância, diversidade e zonação dos crustáceos no talude da costa central do

396    Brasil (11° - 22° S) coletados pelo Programa REVIZEE/Score Central: prospecção

397    pesqueira. In: COSTA, R. A. S.; OLAVO. G.; MARTINS. A.S. (Eds.) **Biodiversidade da fauna**

398    **marinha profunda na costa central brasileira**. Rio de Janeiro: Museu Nacional, 2007, p.

399    133-162.

400  Stackebrandt, E.; Goebel, B. M. Taxonomic note: a place for DNA-DNA reassociation and 16S

401      rRNA sequence analysis in the present species definition in bacteriology. **International**

402      **Journal of Systematic Bacteriology**, v. 44, n. 4, p. 846-849, out. 1994.

403  Tâmega, F. T. S.; Oliveira, P. S.; Figueiredo, M. A. O. (Eds.) **Catalogue of the Benthic Marine Life**

404      **from Peregrino Oil Field, Campos Basin, Brazil**. Rio de Janeiro: Instituto Biodiversidade

405      Marinha, 2013.

406  Vieira, L. M.;  Migotto, A. E. Entoprocta Checklist of the State of São Paulo. **Biota Neotrop**., v. 11,

407      suplemento 1, p. 497-501, 2011.

408  Wang, Y., et al. Optimal Eukaryotic 18S and Universal 16S/18S Ribosomal RNA Primers and

409      Their Application in a Study of Symbiosis. **PLoS ONE**, v. 9, n. 3, 2014.

410  Zhang, Z.-Q. (ed.) Animal biodiversity: An outline of higher-level classification and survey of

411      taxonomic richness. **Zootaxa**, v. 3148, p. 1–237, 2011.

412  Zhang Z. et al. A greedy algorithm for aligning DNA sequences. **J Comput Biol.**, v. 7, n. 1-2, p. 203-

413  214,                                                                                        2000.

414

415 **Table 1** – Sampling date, location and depth Location of sampling stations B3, B4, C2, F5 and G2

416 in Campos Basin, southeast Brazil.

417

| | Sampling date | Latitude (SIRGAS2000) | Longitude (SIRGAS2000) | Depth (m) |
|---|---|---|---|---|
| Station B3 | 02/20/2009 | -22,997011 | -41,352583 | 77 |
| Station B4 | 02/21/2009 | -23,16851 | -41,052264 | 107 |
| Station C2 | 07/16/2009 | -22,625989 | -41,365082 | 54 |
| Station F5 | 02/24/2009 | -22,290999 | -40,110584 | 143 |
| Station G2 | 02/25/2009 | -21,98502 | -40,419918 | 56 |

418
419

420  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

421  of MgCl$_2$ (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

422  pure destilaed water (Invitrogen) to complete 50 µl final reaction volume.

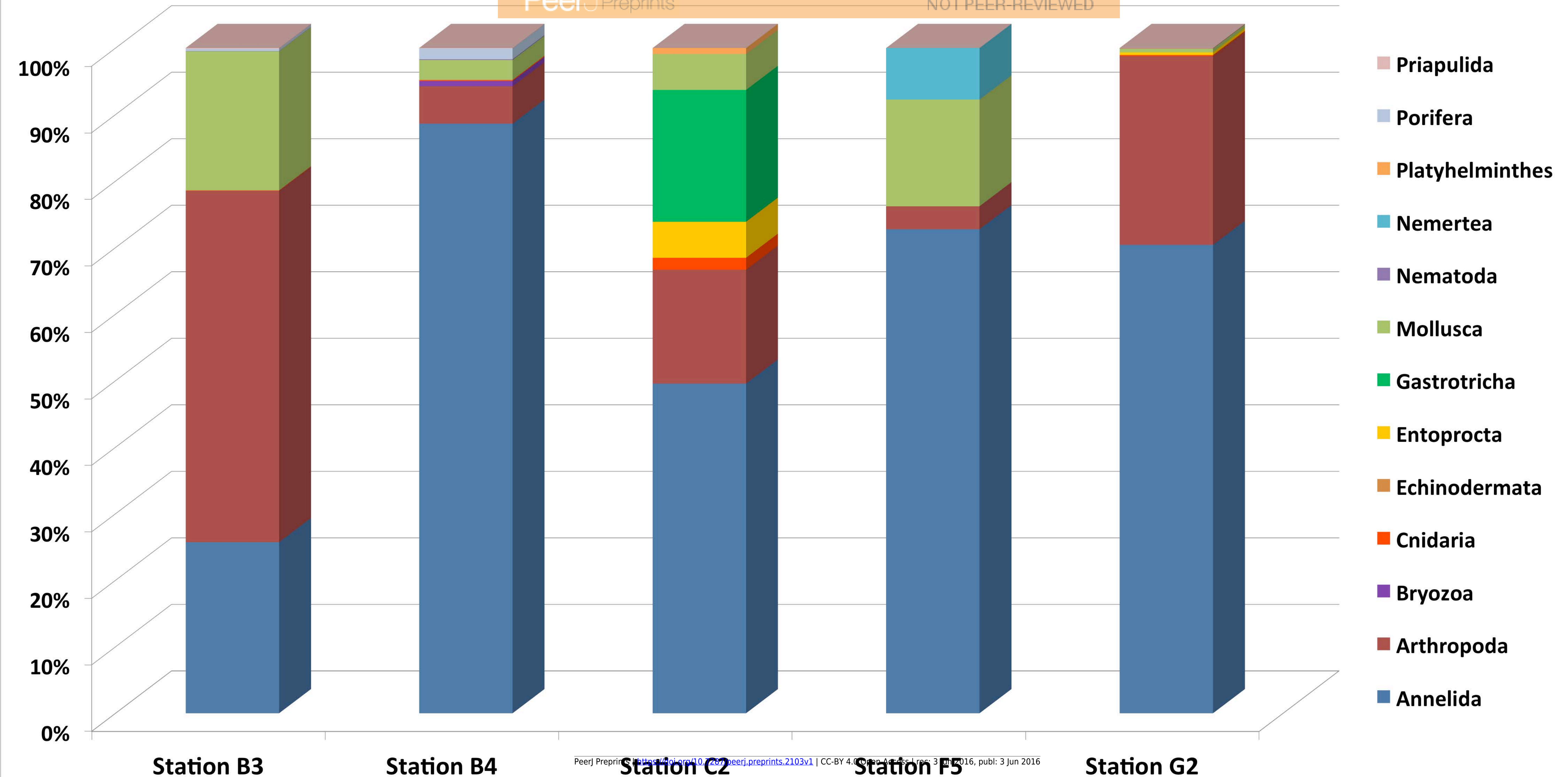| Target | Primer (F – Forward; R – reverse) | Denaturation | cycles | denaturation | anealing | Extension | Final extension | References |
|---|---|---|---|---|---|---|---|---|
| **COI** | TITCIAAYCAYAARGAYATTGG (F – jLCO1490); TAIACYTCIGGRTGICCRAARAAYCA (R – jHCO2198) | 1' @94oC | 10+30 | 30"@94oC | 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); GATCCTTCCGCAGGTTCACCTAC (R- 9R) | 2' @94oC | 40 | 30"@94oC | 30'@55oC | 1'@72oC | 5'@72oC | Whiting et al., 1997; Whiting, 2002 |
| **rRNA 28S** | ACCCGCTGAATTTAAGCAT (F – C1'); TGAACTCTCTCTTCAAAGTTCTTTTC (R- C2) | 2' @94oC | 40 | 30"@94oC | 30'@55oC | 1'@72oC | 5'@72oC | Van Le et al., 1993; Chen et al., 2003 |

423
424
425

426    **Table 2** – OTU per sample. OTU without a similar sequence on Genbank NR are under

427    'No Hits' fragments . OTU that did not comply with established LCA parameters (e.g.

428    score bellow 100) or do not add up to a node are under 'non attributed reads'. Also

429    under 'non-attributed' are Prokaryots attributed by rRNA16S, taxa attributed by genes

430    other than the 3 targets and taxa defined at Genbank as 'undefined'. They were also
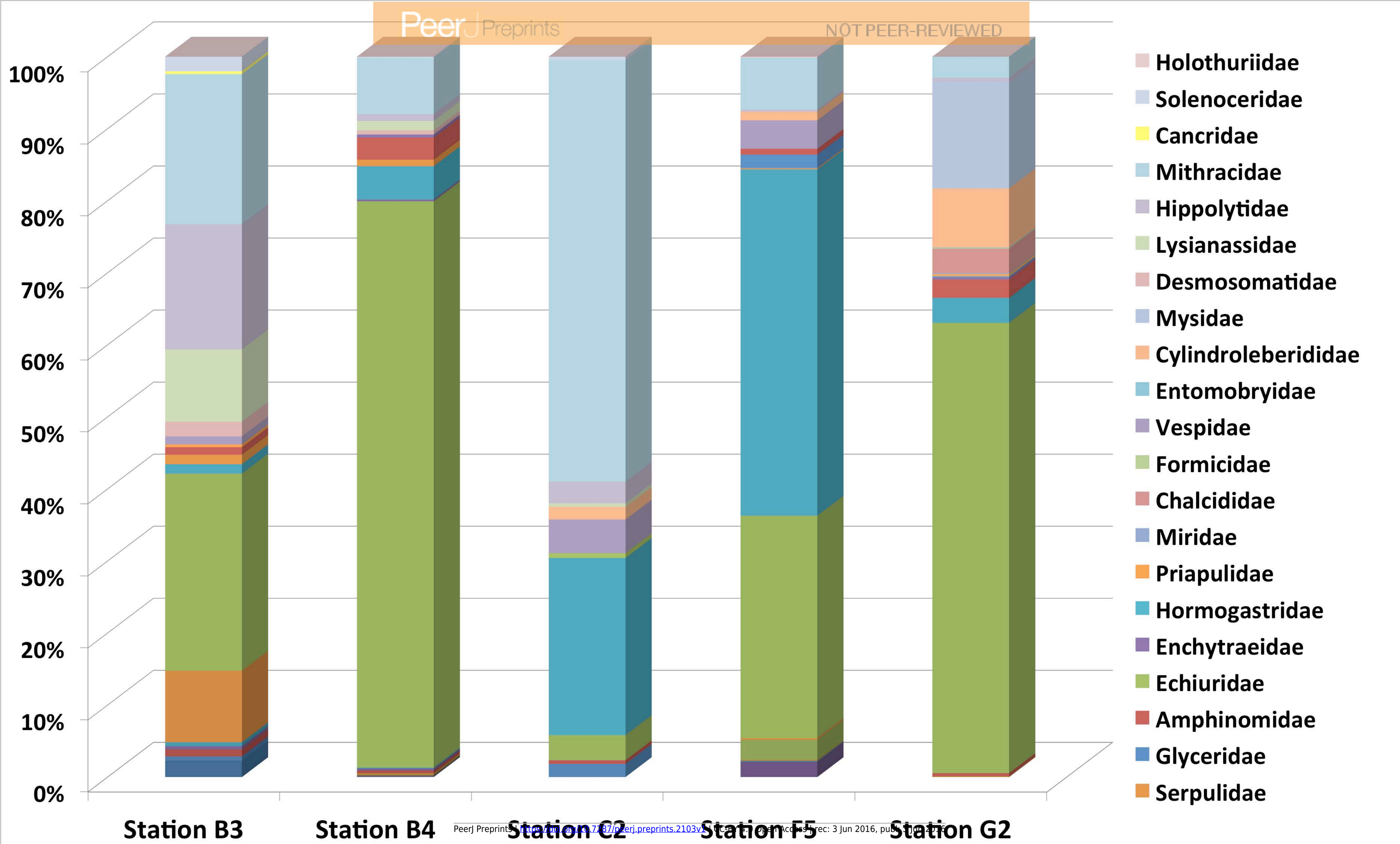
431    disabled at the cladograms.

432

| 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 |

433
434

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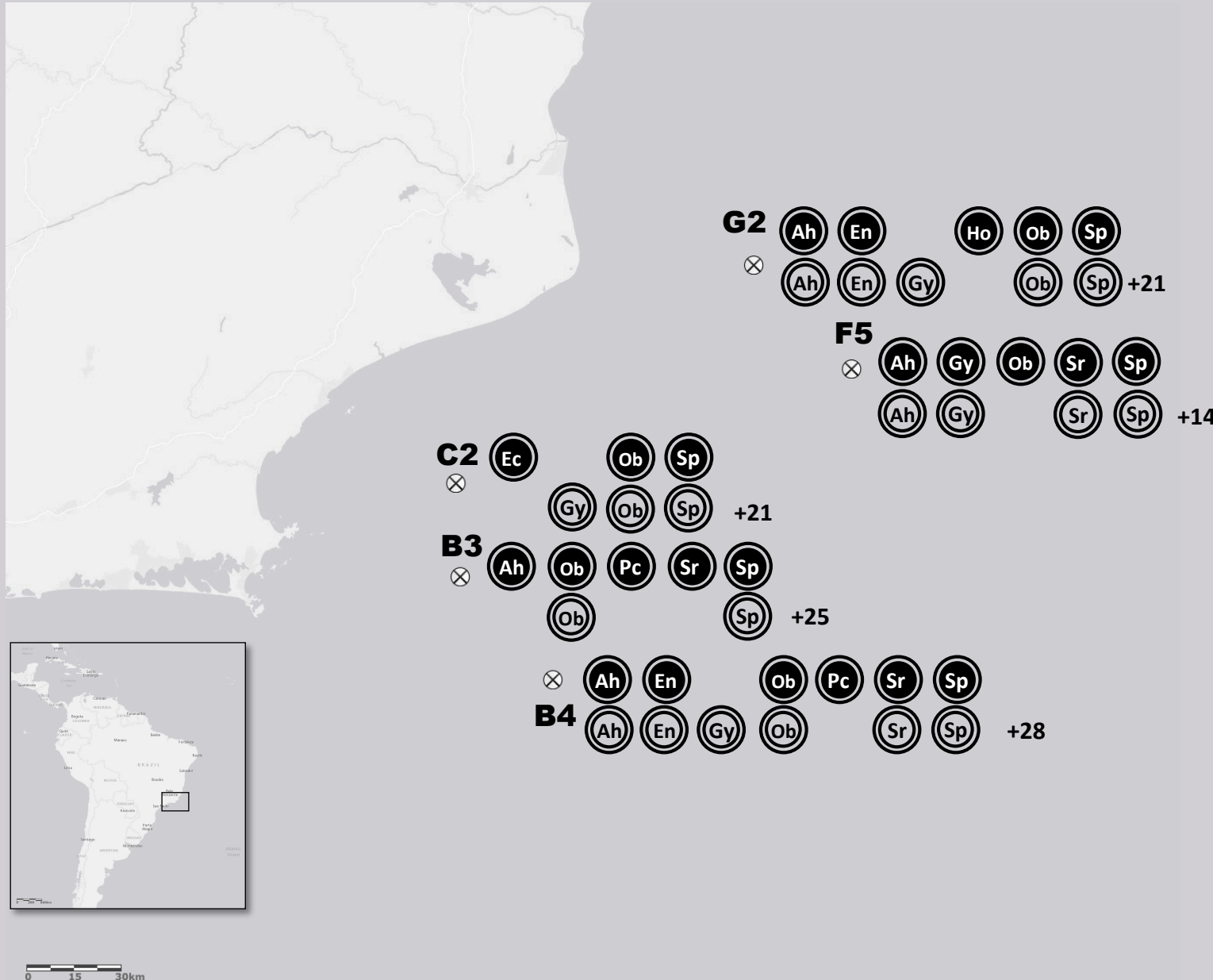
(B) in each station.

435

436

437 **Figure 2 – Distribution of the main invertebrate phylum identifyed by molecular**

438 **and morphological taxonomy in Campos Basin**. A) annelida distribution, b)

439 arthropoda distribution, C) mollusca distribution.

440
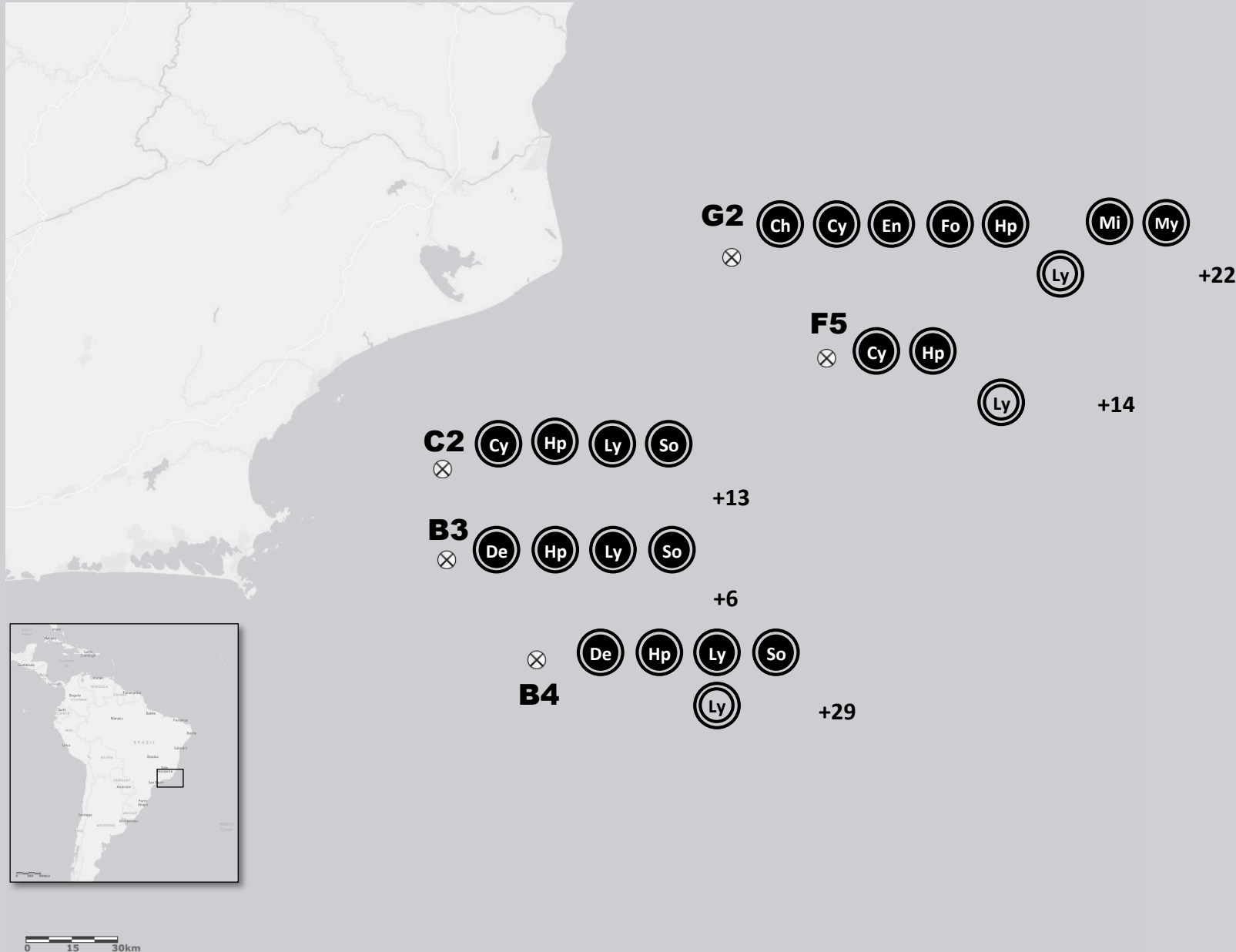
Annelida distribution

**Molecular** | Morphological

**Amphinomidae**

**Enchytraeidae**

**Echiuridae***

**Glyceridae**

**Hormogastridae****

**Orbiniidae**

**Pectinariidae***

**Serpulidae**

**Spionidae**

* Present in other stations of Habitats

** Non-marine family

# Mollusca distribution

**G2**

**+15**

**F5**

**+8**

**C2**

**+9**

**B3**

**+13**

**B4**

**+14**

Ar Ar **Arcidae\* \*\***

Ga Ga **Galeommatidae\*\***

Ha Ha **Haliotidae\*\***

Lo Lo **Lottiidae\*\***

Mt Mt **Mactridae\* \*\***

Ml Ml **Mytilidae**

Pe Pe **Pectinidae\* \*\***

\* Present in other stations of Habitats Project

\*\* Previous studies in Campos basin

0   15   30km

441    Suplementar material 2 – Family level Cladograms of the 5 sampling stations.

442    Cladograms were built using speciments identified with any of the 3 target genes. Bar

443    inside the squares represent the number of reads from each gene used to create the

444    node. A) Family cladogram for station B3; b) Family cladogram for station B4; C) Family

445    cladogram for station C2; D) Family cladogram for station G2; E) Family cladogram for

446    station F5.

447

*Family*

448    Suplementar material 3 – List of species identified by molecular and morphological

449    taxonomy

| Specie | 18S | 20S | COI | Study | |
|---|---|---|---|---|---|
| Cnemidocarpa verrucosa | + | + | + | Schettini | |
| Desmarestia dudresnayi | + | + | + | Schettini | |
| Erythrophyllum delesserioides | + | + | + | Schettini | |
| Eurythenes gryllus | + | + | + | Schettini | |
| Galeomma turtoni | + | + | + | Schettini | |
| Grifola frondosa | + | + | + | Schettini | |
| Haliotis diversicolor | + | + | + | Schettini | |
| Hormogaster redii | + | + | + | Schettini | |
| Lysmata seticaudata | + | + | + | Schettini | |
| Malassezia globosa | + | + | + | Schettini | |
| Marenzelleria arctia | + | + | + | Schettini | |
| Mimachlamys varia | + | + | + | Schettini | |
| Mysidium columbiae | + | + | + | Schettini | |
| Parotocinclus maculicauda | + | + | + | Schettini | |
| Pinctada imbricata | + | + | + | Habitats | and Hits |
| Platynereis dumerilii | + | + | + | Habitats | and Hits |
| Pontocaris lacazei | + | + | + | Habitats | |
| Praxillella affinis | + | + | + | Habitats | |
| Progoniada regularis | + | + | + | Habitats | and Hits |
| Protodorvillea kefersteini | + | + | + | Habitats | |
| Pteria colymbus | + | + | + | Habitats | |
| Scalibregma inflatum | + | + | + | Habitats | and Hits |
| Scapharca broughtonii | + | + | + | Schettini | |
| Serpula vermicularis | + | + | + | Schettini | |
| Syllis gracilis | + | + | + | Habitats | and Hits |
| Syllis variegata | + | + | + | Habitats | and Hits |
| Travisia brevis | + | + | + | Habitats | and Hits |
| Travisia forbesii | + | + | + | Habitats | and Hits |
| Travisia pupa | + | + | + | Habitats | and Hits |
| Aglaophamus circinata | | + | + | Habitats | and Hits |
| Alpheus formosus | | + | + | Habitats | |
| Amphipholis squamata | | + | + | Habitats | |
| Aricidea wassi | | + | + | Habitats | and Hits |
| Chelonia mydas | | + | + | Schettini | |
| Praxillella pacifica | | + | + | Habitats | and Hits |
| Priapulus caudatus | | + | + | Schettini | |
| Scolelepis bonnieri | | + | + | Schettini | |
| Scolelepis foliosa | | + | + | Schettini | |
| Amphimedon queenslandica | + | | + | Schettini | |
| Axiothella rubrocincta | + | | + | Habitats | and Hits |
| Bathyarca pectunculoides | + | | + | Habitats | |
| Bathyglycinde profunda | + | | + | Habitats | |
| Bathyglycinde sibogana | + | | + | Habitats | |
| Caprella equilibra | + | | + | Habitats | and Hits |
| Ceratocephale abyssorum | + | | + | Habitats | and Hits |

| Specie | 18S | 20S | COI | Study | |
|--------|-----|-----|-----|-------|---|
| Ciona intestinalis | + | | + | Schettini | |
| Clymenella torquata | + | | + | Habitats | and Hits |
| Pectinaria granulata | + | | + | Schettini | |
| Perna viridis | + | | + | Schettini | |
| Protaspis grandis | + | | + | Schettini | |
| Syllis hyalina | + | | + | Habitats | and Hits |
| Didemnum candidum | | | + | Schettini | |
| Leodamas rubra | | | + | Habitats | and Hits |
| Leodia sexiesperforata | | | + | Habitats | |
| Leptochelia dubia | | | + | Habitats | |
| Leucothoe urospinosa | | | + | Habitats | and Hits |
| Lumbrineris latreilli | | | + | Habitats | and Hits |
| Lysidice ninetta | | | + | Habitats | and Hits |
| Lysmata anchisteus | | | + | Schettini | |
| Macrochaeta clavicornis | | | + | Habitats | |
| Marphysa bellii | | | + | Habitats | and Hits |
| Mendicula ferruginosa | | | + | Habitats | and Hits |
| Mooreonuphis pallidula | | | + | Habitats | and Hits |
| Neanthes acuminata | | | + | Habitats | and Hits |
| Nereimyra punctata | | | + | Habitats | and Hits |
| Notomastus latericeus | | | + | Habitats | and Hits |
| Ophelina acuminata | | | + | Habitats | and Hits |
| Pyropia haitanensis | | | + | Schettini | |
| Scapharca kagoshimensis | | | + | Schettini | |
| Scoloplos armiger | | | + | Schettini | |
| Isolda pulchella | | | ++ | Habitats | and Hits |
| Apophlaea lyallii | + | + | | Schettini | |
| Chaetoceros curvisetus | + | + | | Schettini | |
| Coelomactra antiquata | + | + | | Schettini | |
| Crassinella lunulata | + | + | | Habitats | and Hits |
| Cryptococcus friedmannii | + | + | | Schettini | |
| Cyclaspis alba | + | + | | Habitats | |
| Cylichna alba | + | + | | Habitats | and Hits |
| Engraulis japonicus | + | + | | Schettini | |
| Euclymene oerstedi | + | + | | Habitats | and Hits |
| Eulalia viridis | + | + | | Habitats | and Hits |
| Eumida sanguinea | + | + | | Habitats | and Hits |
| Exogone dispar | + | + | | Habitats | and Hits |
| Galathowenia oculata | + | + | | Habitats | |
| Glycera americana | + | + | | Habitats | and Hits |
| Glycera southeastatlantica | + | + | | Habitats | and Hits |
| Goniada emerita | + | + | | Habitats | |
| Hesiospina aurantiaca | + | + | | Habitats | and Hits |
| Patelloida striata | + | + | | Schettini | |
| Scopelocheirus schellenbergi | + | + | | Schettini | |
| Subulatomonas tetraspora | + | + | | Schettini | |
| Ophelina cylindricaudata | | + | | Habitats | and Hits |
| Ophiactis lymani | | + | | Habitats | |
| Trypanosyllis zebra | | + | | Habitats | and Hits |

| Specie | 18S | 20S | COI | Study | |
|---|---|---|---|---|---|
| Ahnfeltiopsis leptophylla | + | | | Schettini | |
| Crucigera zygophora | + | | | Schettini | |
| Leitoscoloplos pugettensis | + | | | Schettini | |
| Malassezia nana | + | | | Schettini | |
| Ophiura ljungmani | + | | | Habitats | |
| Owenia fusiformis | + | | | Habitats | and Hits |
| Panthalis oerstedi | + | | | Habitats | and Hits |
| Paralacydonia paradoxa | + | | | Habitats | and Hits |
| Paramphinome jeffreysii | + | | | Habitats | and Hits |
| Pholoe minuta | + | | | Habitats | |
| Phtisica marina | + | | | Habitats | |
| Phyllodoce longipes | + | | | Habitats | and Hits |
| Solenocera crassicornis | + | | | Schettini | |
| Strombidium paracalkinsi | + | | | Schettini | |
| Phagomyxa odontellae | + | | | Schettini | |

450