LIMNOLOGY AND
OCEANOGRAPHY

# Endemicity and climatic niche differentiation in three marine ciliated protists

| | |
|---|---|
| Abstract: | The biogeographic pattern of single-celled eukaryotes (protists), including ciliates, is poorly understood. Most marine species are believed to have a relatively high dispersal potential, such that both globally-distributed and geographically-isolated taxa exist. Primary occurrence data for three large, easily identified ciliate species, Parafavella gigantea, Schmidingerella serrata and Zoothamnium pelagicum, and environmental data drawn from the National Oceanic and Atmospheric Administration's (NOAA) World Ocean Atlas were used to estimate each species' spatial and environmental distributions using Maxent v3.3.3k. The predictive power of the models was tested with a series of spatial stratification studies, which were evaluated using partial receiver operating characteristic statistics. Differences between niches occupied by each taxon were evaluated using background similarity tests. All predictions showed significant ability to anticipate test points. The null hypotheses of niche similarity were rejected in all background similarity tests comparing the niches among the three species. This paper provides a first quantitative assessment of environmental conditions associated with three species of ciliates and a first estimate of their spatial distributions in the North Atlantic, which can serve as a benchmark against which to document distributional shifts. These species follow consistent, predictable patterns related to climate and environmental biochemistry; the importance of climatic conditions as regards protist distributions is noteworthy considering the effects of global climate change. |

SCHOLARONE™
Manuscripts

1  **TITLE**
2  Endemicity and climatic niche differentiation in three marine ciliated protists
3
4  **AUTHORS**
5  Richard AJ Williams[1,2,3*], Hannah L Owens[2,4], John Clamp[5]†, A Townsend Peterson[2], Alan
6
7  Warren[6], Mercedes Martín Cereceda[7]
8
9
10
11
12  [1]   Centre for Ecology and Evolution in Microbial Model Systems, Linnaeus University,
13       SE-391 82, Kalmar, Sweden
14  [2]   Biodiversity Institute, University of Kansas, Lawrence, KS 66045, USA
15  [3]   Department of Biodiversity, Ecology and Evolution, Universidad Complutense de
16       Madrid, 28040, Madrid, Spain
17  [4]   Department of Biology, University of Florida, Gainesville, FL 32611, USA
18  [5]   Department of Biological and Biomedical Sciences, North Carolina Central
19       University, Durham, NC 27707, USA
20  [6]   Department of Life Sciences, Natural History Museum, London SW7 5BD, UK
21  [7]   Department of Genetics, Physiology and Microbiology, Universidad Complutense de
22       Madrid, 28040, Madrid, Spain
23  *   Corresponding author: richard.williams@lnu.se
24  †   John Clamp passed away on Feb. 6th, 2018, following a short illness
25
26
27
28
29  **KEYWORDS**
30  Ciliates, Ecological Niche Models, Niche Differentiation, *Parafavella gigantea*,
31  *Schmidingerella serrata, Zoothamnium pelagicum*
32
33  **RUNNING HEAD**
34  Marine ciliate niche and distribution
35
36
37  **Abstract:** 213 words
38  **Paper:** 4445 words
39  **References:** 53
40
41

**ABSTRACT**

The biogeographic pattern of single-celled eukaryotes (protists), including ciliates, is poorly

understood. Most marine species are believed to have a relatively high dispersal potential,

such that both globally-distributed and geographically-isolated taxa exist. Primary occurrence

data for three large, easily identified ciliate species, *Parafavella gigantea*, *Schmidingerella*

*serrata* and *Zoothamnium pelagicum*, and environmental data drawn from the National

Oceanic and Atmospheric Administration's (NOAA) World Ocean Atlas were used to

estimate each species' spatial and environmental distributions using Maxent v3.3.3k. The

predictive power of the models was tested with a series of spatial stratification studies, which

were evaluated using partial receiver operating characteristic statistics. Differences between

niches occupied by each taxon were evaluated using background similarity tests. All

predictions showed significant ability to anticipate test points. The null hypotheses of niche

similarity were rejected in all background similarity tests comparing the niches among the

three species. This paper provides a first quantitative assessment of environmental conditions

associated with three species of ciliates and a first estimate of their spatial distributions in the

North Atlantic, which can serve as a benchmark against which to document distributional

shifts. These species follow consistent, predictable patterns related to climate and

environmental biochemistry; the importance of climatic conditions as regards protist

distributions is noteworthy considering the effects of global climate change.

61 **INTRODUCTION**

62 Marine microorganisms are essential to the planet's health, marine biodiversity, and the

63 fishing industry. They carry out nearly 50% of the world's photosynthesis, 70% of which is

64 locked up in long-term carbon storage (Field et al. 1998). Microbes are the base of the marine

65 ecological pyramid, and are vital to the marine food web that transfers energy from primary

66 producers to higher trophic levels (Azam et al. 1983; Fenchel 2008). Understanding the

67 current spatial distributions of marine microorganisms is thus important to improving

68 understanding of their functions, and is fundamental for impact assessment of effects of

69 climate change.

70

71 Microorganisms—used here to refer to all microscopic organisms (including protists,

72 archaea, bacteria and fungi)—are generally considered to be cosmopolitan in distribution

73 (Beijerinck 1913): "everything is everywhere, but the environment selects". They are thought

74 to be ubiquitous thanks to huge population sizes and a consequent low probability of

75 extirpation; also, they are limitlessly dispersible, have negligible rates of allopatric speciation

76 compared to larger organisms, and can potentially occur at any place that meets their

77 biological requirements (Finlay et al. 1996). These assumptions suggest that processes

78 driving biogeographic patterns of microbe diversity are fundamentally different from

79 macroscopic organisms, with profound implications for understanding mechanisms driving

80 microbial distributions and evolution. A variant on this view states that whereas most free-

81 living protist species are globally distributed, about one-third of them may be

82 biogeographically restricted (Foissner 1999; Foissner 2006).

83

84 Nonetheless, microbial geographic patterns are still poorly understood. Description of the

85 biogeography of protists is hampered by a scarcity of taxonomists, high frequency of

86   misidentification and dramatic under-sampling (Foissner 2006). Traditionally, protists have

87   been identified by morphological features, though this approach risks lumping organisms

88   with discrete biogeography, ecology and genetics into single morphospecies (Gentekaki and

89   Lynn 2010; Katz et al. 2011). The high rates of synonymy and existence of polymorphic life

90   stages (phenotypic variability) in ciliates can exacerbate the problem (Dolan 2016). Protist

91   biogeography has benefited from molecular studies, as environmental sequencing can provide

92   abundant distributional data. Several contributions have described geographic

93   circumscription in 18 S rDNA sequences from diverse aquatic protist groups—diatoms

94   (Evans et al. 2009), Cercozoa (Bass et al. 2007), Amoebozoa (Aguilar et al. 2014),

95   heterotrophic flagellates (Boenigk et al. 2006), and ciliates (Bass et al. 2009). Several recent

96   morphological and molecular studies have found that many protists, particularly the rarer

97   ones, follow discrete biogeographic patterns (Filker et al. 2016; Logares et al. 2015; Segovia

98   et al. 2017). Further studies relating morphology, functionality and molecular information are

99   needed to describe the dispersal ability of these and other protist species.

100

101   Ciliophora is a diverse phylum of heterotrophic or mixotrophic protists—the ciliates—with

102   ~4500 taxonomically valid, free-living species (Foissner et al. 2008), characterized by the

103   presence of hair-like organelles called cilia, used for locomotion, attachment, feeding,

104   sensation, etc., and nuclear dimorphism (Gao et al. 2016). They have been sampled from

105   marine environments; indeed, the Ocean Biogeographic Information System (OBIS) has

106   locality data for nearly 200,000 Ciliophora observations, albeit with a clear bias towards

107   sampling along major trade routes. Some microbial species are more conspicuous and easier

108   to identify with confidence than others, leading some authors to propose these "flagship"

109   species as the ultimate proof of endemism (Foissner 2006). Such species, the "elephants" of

110   the microbial world, cannot be missed if they are present, because of their distinctive

4

111　morphological features and/or significantly large size, with the caveat that traditional

112　morphological-based identification of microbial species risks lumping diverse ecological and

113　genetic species into a single morphospecies complex and underestimating cryptic diversity

114　(Gentekaki and Lynn 2010).

115

116　Ecological niche models (ENM) infer suitable abiotic habitat conditions for non-model

117　organisms by generating a correlational model that unites occurrence information and

118　environmental data for the taxon of interest to determine the geographic distribution of

119　habitat conditions correlated with species occurrences (Soberón and Peterson 2005). They

120　have been used to explore diverse topics in distributional ecology, including species'

121　geographic distributions, niche conservatism, spread of invasive macrospecies and diseases,

122　and effects of climate change on species distributions; see Peterson et al. (2011) for a

123　summary. ENMs have long been used to infer distributions of species in the marine

124　environment (Wiley et al. 2003). To date, marine ENMs have been used mainly to create

125　models for macrofaunal distributions (Bentlage et al. 2013; Saupe et al. 2014), although some

126　authors have applied ENM approaches to characterize phytoplankton (Brun et al. 2015) and

127　Foraminifera (Langer et al. 2013). Relatively few attempts have been made to describe the

128　ecological dimensions of microbial distributions for free-living terrestrial microbes (Aguilar

129　and Lado 2012). One noteworthy recent contribution explored the phylogeography of 18S

130　rDNA variants of the myxomycete *Badhamia melanospora* (phylum Amoebozoa) (Aguilar et

131　al. 2014).

132

133　The aim of this contribution is to show, for the first time in ciliates, the utility of ENMs in

134　putting ciliate morphospecies on the map, that is, in analysing and understanding their

135　biogeography. We focus on three flagship species from the phylum Ciliophora: *Parafavella*

136　　*gigantea* (Brandt, 1896), *Schmidingerella* (= *Favella*) *serrata* (Möbius, 1887), and

137　　*Zoothamnium pelagicum* Du Plessis, 1891. These species are particularly obvious thanks to

138　　their large body size. *Parafavella gigantea* (up to 750 μm in length) and *Schmidingerella*

139　　*serrata* (up to 350 µm long) are relatively obvious within plankton samples as they are much

140　　larger than many other protists. *Zoothamnium pelagicum*, though somewhat smaller (up to

141　　120 μm), is colonial, and often forms discoidal clusters 2-3 mm thick that are clearly visible

142　　to the human eye (Laval 1968). These were the only species for which sample size and

143　　taxonomic clarity were sufficient to permit developing ecological niche models (ENMs).

144

145　　**METHODS**

146　　**Data on ciliate distributions and ecological dimensions**

147　　Primary occurrence data for Ciliophora species were downloaded from the Global

148　　Biodiversity Information Facility (GBIF; http://www.gbif.org/) and Ocean Biogeographic

149　　Information System (OBIS; http://www.iobis.org/) on 10 December 2015. We obtained more

150　　than 200,000 observations, but most were not usable for development of ENMs: many were

151　　not supported by voucher specimens or sequences (which we consider to be necessary

152　　documentation), and taxonomic resolution for many taxa was poor. Observations were often

153　　duplicated from the same locality; and sample sizes were low for most taxa.

154

155　　We found only three ciliate species with sufficient sample size, taxonomic clarity, and

156　　geographic spread to develop ENMs:  *P. gigantea* ($N = 368$), *S. serrata* ($N = 135$), and *Z.*

157　　*pelagicum* ($N = 279$). Records for these taxa were collected between 1903 and 2008 (see Fig.

158　　1); duplicate records falling within the same 1° grid square were eliminated. The area of

159　　coverage was limited to the North Atlantic (20°E-80°W; 0°-90°N) centred on the best-

160　　sampled ocean regions worldwide. Occurrence data were visualized in ESRI ARCGIS 9.3.1

161   (ESRI, Redlands, CA, USA), and showed broadly consistent distributional patterns. The few

162   data (<2% for each taxon) that fell outside this region were nonetheless included in

163   development of models to avoid prejudging limits of unknown distributional patterns (see

164   Fig. 1).

165

166   Environmental data layers were drawn from National Oceanic and Atmospheric

167   Administration's (NOAA) World Ocean Atlas (NOAA 1999), at a native spatial resolution of

168   1°: ~110 km at the Equator. These data layers represent long-term annual means from

169   oceanographic datasets covering 1900-1997. We developed two models: Model 1, based on

170   six annual surface temperature and salinity variables (mean, maximum and minimum, for

171   each), and Model 2, based on the same six variables for surface temperature and salinity, as

172   well as six annual average surface biochemical variables (Supplementary Table 15). The data

173   for each environmental coverage layer were converted into standard normal variates in

174   ArcGIS 9.3.1 prior to analysis.

175

176   **Ecological Niche Modelling**

177   ENMs were calibrated for each species using the maximum entropy algorithm Maxent

178   v3.3.3k (Phillips et al. 2006). Maxent is a correlational algorithm that uses presence-only

179   species occurrence data to estimate ecological niche parameters (Phillips et al. 2006). Maxent

180   fits a suitability surface for the species of interest to the set of pixels across the study region,

181   maximising the entropy of the probability distribution but constrained to return higher per-

182   pixel suitability scores for pixels with environmental variable combinations most like those at

183   which the species has been detected. That is, the resulting suitability surface is a raster map in

184   which each pixel is scored regarding its similarity to climatic and biochemical conditions at

185   sites where the species has been observed.

186

187    From Maxent suitability surfaces, distributions of suitable areas for each species were defined

188    using a minimum training occurrence threshold (Pearson et al. 2007). We defined two

189    thresholds for suitable and unsuitable areas for each species: the least presence threshold

190    (LPT)—equal to the lowest probability at any occurrence location—and a more conservative

191    10th percentile training presence threshold, which is the highest suitability value that includes

192    90% of the calibration data. LPT is appropriate when there is no error in occurrence data that

193    may be introduced during geo-referencing or identification, whereas the 10th percentile

194    approach anticipates up to 10% error among occurrences. More complex methods of

195    thresholding have demonstrated better skill in classifying suitable and unsuitable

196    environments (Jímenez-Valverde and Lobo 2007; Liu et al. 2005), but our dataset lacks true

197    absence data, so a simple omission-based technique was preferable (Bean et al. 2012). In

198    addition, under this approach, the map of suitable conditions for a species is defined to be

199    inclusive of all habitable conditions, and as such should include all known populations

200    (Peterson 2014).

201

202    **Testing predictive power of ENMs**

203    This study examines whether the occurrences of three ciliate taxa follow a consistent and

204    predictable environmental pattern in the Atlantic Ocean. As such, we tested whether models

205    could predict independent subsets of occurrence data reliably. These tests covered the entire

206    study area, using replicate random subsets of the occurrence data, and spatially stratified

207    subsets of the occurrence data, 50% of localities for model calibration and 50% for testing for

208    both models under default parameters, and choosing logistic output format with suitability

209    values from 0 (unsuitable) to 1 (suitable). To avoid extrapolation in model features, no

210   clamping or extrapolation was permitted (Owens et al. 2013); to avoid overly complex

211   models, no hinge or threshold features were permitted.

212

213   For prediction across the study area, 10 bootstrap replicate runs with a maximum of 10,000

214   iterations were conducted using a random seed with 70% of occurrence points. For spatial

215   stratification challenges, occurrence data were split spatially into quadrants above and below

216   the median longitude and latitude of the occurrence data. From this spatial stratification, we

217   developed the following three pairs of quadrants: west versus east of the median longitude,

218   north versus south of the median latitude, and on-diagonal (upper left-hand and lower right-

219   hand quadrants) versus off-diagonal (lower left-hand and upper right-hand quadrants). In

220   each case, we developed both reciprocal predictions, testing the ability of ENMs to anticipate

221   the spatial distribution of occurrence data in areas for which no sampling is available. Models

222   were evaluated by applying partial receiver operating characteristic (ROC) statistics to the

223   50% subset of occurrences withheld from model development for testing. Area under the

224   curve (AUC) ratios were calculated using the Partial ROC function available in CONABIO's

225   NicheToolBox (http://shiny.conabio.gob.mx:3838/nichetoolb2/). Final models were

226   developed using 50% of available data.

227

228   One concern in comparisons of models based on different sets of environmental data is that

229   increased model complexity may lead to overfitting (Peterson et al. 2007; Radosavljevic and

230   Anderson 2014). A recent tendency, as a consequence, has been to use the Akaike

231   information criterion (AIC) as a means of comparing model likelihood values while

232   penalizing complex (and ostensibly overfit) models (Warren and Seifert 2011). Specifically,

233   we calculated the sample-size-corrected AICc statistic using ENMTools version 1.3.3

234   (Warren et al. 2010), and chose as the "best" model the one that had the lowest AICc values.

235    Results based on AIC (without the sample-size correction) and the alternative Bayesian

236    information criterion (BIC) metric were similar, and so are not presented.

237

238    **Niche difference and range restriction**

239    To assess whether niches differed among taxa (which would imply range limits), we

240    evaluated differences in niches occupied by each taxon using background similarity tests

241    (Warren et al. 2008) available in the ENMTools R package version 0.1

242    (https://github.com/danlwarren/ENMTools). This test considers similarity between predicted

243    geographic distributions among species, using two statistics, Schoener's *D* and a modified

244    Hellinger's *I* metric, to quantify similarity. This test evaluates whether ENMs generated from

245    two species are more different than expected when occurrences are drawn from the same

246    underlying distribution across the region accessible to each species. The test allows

247    specification of an area of analysis ("the background"), which we equate with the area

248    accessible to a species over relevant time periods (Barve et al. 2011; Soberón and Peterson

249    2005). Numbers of points sampled from the background were set at observed sample sizes. In

250    each test, 100 replicate analyses were performed to estimate probabilities associated with null

251    hypotheses of niche similarity. The hypothesis that species were no more like each other than

252    if points had randomly been drawn from the study area was rejected if observed similarity

253    between models fell below the $5^{th}$ of the null distribution.

254

255    **RESULTS**

256    The area under the curve (AUC) ratios for the independent *Parafavella gigantea* testing data

257    were 1.77 and 1.79 for Model 1 and Model 2, respectively ($P < 0.05$; Table 1). The 10%

258    threshold for both models predicted a broad potential distribution across the cold temperate

259    zone of the North Atlantic north to about 79 / 81° N (Model 1 / Model 2) in the eastern

260    Atlantic, and about 67° N (both models) off the Greenland coast (Fig. 2) and the Labrador

261    Sea. The southerly limit was at about 48° / 47° N in the eastern Atlantic, and 36° / 38° N off

262    the US coast. Both models showed a more northerly prediction in the eastern Atlantic than in

263    the western Atlantic, owing to the ameliorating effect of the Gulf Stream / North Atlantic

264    Drift / Norwegian Current. Locality points with a borderline prediction for both models were

265    principally from the poorly represented southeastern quadrant of the study area, the

266    Greenland Sea in the northeast, and Labrador Sea and Hudson Bay in the northwest. The 10th

267    percentile training presence fractional area predicted was about 22% of the study area for

268    each model, whereas that predicted under least presence threshold (LPT) was about 43 / 44%

269    of the study area.

270

271    The area under the curve ratios for the independent *Schmidingerella serrata* testing data were

272    1.83 and 1.81 for Model 1 and Model 2, respectively (Table 1). Both models predicted a

273    broad distribution in the North Atlantic from about 74° / 71° N in the eastern Atlantic, around

274    65° N off the Greenland coast, and 52° N off the Canadian coast (Fig. 2**)**. The southern

275    predicted limit extended from about 45° / 46° N in the European coast, curving slightly to the

276    north in the mid-Atlantic, and about 36° N off the US coast. Most "borderline" points were

277    either in the Labrador Sea and the northwest Atlantic or the southeast part of the distribution,

278    including the Mediterranean. The least presence threshold prediction was similar between the

279    two models, covering about 30% and 28% of the study area, respectively. The proportion of

280    the study area predicted by the 10th percentile area was 18.3% by Model 1, compared to

281    18.9% in Model 2.

282

283    The area under the curve ratios for the independent *Z. pelagicum* testing data were 1.76 and

284    1.75 for Model 1 and Model 2, respectively (Table 1). Both predicted (at the 10th percentile

11

285  training presence fractional area) a broad distribution in the North Cold Temperate Atlantic

286  from about 72° / 71° N in the eastern Atlantic, respectively, to around 52° N off the Canadian

287  coast (Fig. 2). The southern predicted limit included much of the North Temperate Atlantic,

288  to about 36° N for both models off the US coast, almost all the Mediterranean, and even some

289  Tropical areas (to 18° / 15° N in the southeast Atlantic). Model 1's least presence threshold

290  prediction was broader than that of Model 2 (66.7% / 61.4% of the study area, respectively).

291  By contrast, the $10^{th}$ percentile prediction was rather small for each, and, again, more

292  restricted for Model 2: 30.0% / 24.5% of the study area. In each case, it was restricted to a

293  central band (most of the North Temperate Atlantic). All 10 replicates of spatial stratification

294  tests for each taxon and for both models showed significant ability to predict test points (*P* <

295  0.05; see Table 1 and Supplementary Material).

296

297  All tests comparing niches of the three species rejected the null hypothesis of niche similarity

298  between pairs of species when compared to a null distribution generated from the background

299  region (*P* < 0.05; Figure 3, Supplementary Material Table 13). *Parafavella gigantea* was

300  predicted to range in the cold temperate zone of the North Atlantic, extending to the Arctic

301  Ocean, whereas the predicted distribution of *Z. pelagicum* was more southerly, extending to

302  the Caribbean and the Equator; the predicted distribution of *S. serrata* was intermediate. The

303  sample-size corrected Akaike information criterion (AICc) statistics, in all three cases,

304  indicated that the simpler models based only on six annual surface temperature and salinity

305  variables (i.e., Model 1) were preferable to the more complex models that included six

306  biochemical variables as well (Model 2), as differentials in AICc were >90 in all three cases

307  (Supplementary Material Table 14).

308

309    Because all major areas of the predicted-suitable area for each of the three species appeared

310    to be inhabited--at least as far as limited sampling permitted us to conclude (e.g., in parts of

311    the North Sea)--we conclude that these species likely have quite-excellent dispersal abilities.

312

313    **DISCUSSION**

314    Documentation of spatial patterns of microbial species remains undefined and contentious,

315    owing to the complicated nature of detecting them over wide areas.  Evidence for

316    cosmopolitan distributions has been demonstrated for several protistan lineages (Cermeño

317    and Falkowski 2009; Richards et al. 2005). However, many microbial species have broad but

318    restricted distributional patterns (Bass et al. 2009; Bass et al. 2007; Foissner 2006); see Bass

319    and Boenigk (2011) for a comprehensive review. This pattern of moderate endemicity is also

320    seen in prokaryotic microbes (Noguez et al. 2005; Tamames et al. 2010). Some species are

321    even endemic to discrete geographic areas and ecosystems (Foissner 2006). It is unclear to

322    what extent microbial biogeography is obscured by poor understanding of species limits.

323    Only one study (Aguilar et al. 2014) explored the spatial distribution of a single

324    morphospecies of amoeba, *Badhamia melanospora*, using molecular genotyping and ENMs.

325    The authors detected two geographically-structured groups of ribotypes for *B. melanospora*,

326    each of which showed limited distributions, and concluded that this species is not

327    cosmopolitan. Thus, it may be the case that morphospecies mask microbial diversity and

328    biogeographic patterns (Gentekaki and Lynn 2010; Katz et al. 2011).

329

330    Our ecological niche models for three ciliophoran species clearly detected an environmental

331    signal unique to each, such that each species occupies a distinct fundamental ecological

332    niche. The area identified as suitable for each species varied greatly on broad spatial scales,

333    ranging from about 18% of the study area for *S. serrata* to 25-30% for *Z. pelagicum* (10%

13

334  threshold). None of the species is likely to be ubiquitous across the entire study area and they

335  showed statistically significant differences in their environmental characteristics. Our models

336  demonstrated that the distribution of each species is constrained by environment, particularly

337  maximum and minimum temperatures.

338

339  Our locality data were drawn from specimens from the North Atlantic Ocean, the region for

340  which sampling is most dense and complete, although we are conscious that these species'

341  ranges may extend more broadly. Our models were consistently able to predict non-North

342  Atlantic locality data, with the following exception: *S. serrata* beyond the study area are

343  found off the coast of Ecuador and Peru, the North Pacific, and the eastern Mediterranean.

344  Although the latter two sets of points were well predicted by our models, the Ecuador-Peru

345  points were in an area with higher temperature and salinity than those used to train our *S.*

346  *serrata* model, and thus fell outside the ecological niche estimated in our models. Still, the *S.*

347  *serrata* models showed a good fit for the North Atlantic locality data, and provide a

348  parsimonious prediction. The non-Atlantic *S. serrata* points may form part of the natural

349  distribution of the species, or represent recent, perhaps human-mediated, invasion. For

350  instance, the points from the Ecuador-Peru coast might be explained by natural dispersal and

351  colonization by the species following the opening of the Panama Canal (similar to the

352  Lesseps immigrant Indo-Pacific species that entered the Mediterranean via the Suez Canal),

353  or direct transport in a ship's ballast water (Foissner 2011). However, we know of no

354  evidence that can clarify which is the case.

355

356  Additional documented populations for the other two ciliates comprised occurrences that

357  were highly consistent with our predictions. *Parafavella gigantea* has been detected in the

358  North Pacific and along the Russian northern coastal areas, consistent with our North Atlantic

359    prediction. The *P. gigantea* models showed a good fit for the North Atlantic locality data, in

360    that they recovered the areas where the species has been detected, without predicting an area

361    too broad to be credible; however, the North Pacific was not included in the model

362    predictions. Additional locality data for *P. gigantea* compiled in a recent study (Dolan et al.

363    2017) comprised 38 additional locality points for the species, and all were anticipated in our

364    least presence threshold model (significant *P* < 0.05). Few locality points for *Z. pelagicum*

365    were available from outside the North Atlantic (*N* = 6), all from the North Pacific. An

366    exhaustive literature search for additional records revealed no additional locality data for *Z.*

367    *pelagicum*. Our model does not provide a prediction for the North Pacific distributions of

368    these three species. However, climatic and biochemical characteristics of the North pacific

369    are broadly similar to those of the North Atlantic, and we anticipate that these locality points

370    would be predicted by broader analyses.

371

372    The *S. serrata* populations occurring between the Galapagos and the Peruvian coast occupy a

373    tropical niche, quite distinct from the cold temperate distribution of populations detected in

374    the North Atlantic. The true environmental tolerances of this species may indeed be broader

375    than the model we have developed here. If so, our ENM requires locality data from areas with

376    higher temperature and salinity to reflect this ecological tolerance, and our model would thus

377    fail to represent the full fundamental ecological niche. Alternately, the tropical Pacific

378    population may have ecological tolerances distinct from those of the North Atlantic

379    population, and may even represent a distinct species within the morphospecies, suggesting a

380    need for investigation of the ecophysical constraints and genetic distinctiveness of the two

381    populations. Multiple phylogeographic studies on diverse protists, but not ciliate groups

382    (Aguilar et al. 2014; Bass et al. 2007; Evans et al. 2009), have shown distinct geographic

383    ribotypes in morphological species. To the best of our knowledge, no studies have assessed

384    phylogeographic patterns in any of the taxa investigated here. This study demonstrates that

385    three common, morphologically conspicuous and widespread morphospecies occupy distinct

386    geographic distributions and ecological niches in the North Atlantic.

387

388    Our occurrence data were collected over a relatively long time period, just over a century

389    (1903-2008). Global climate change and short-term regional climatic phenomena, including

390    the North Atlantic Oscillation, the Atlantic Multidecadal Oscillation, and the El Niño-

391    Southern Oscillation, may have significant effects on distributions of the three ciliate species

392    analyzed. The environmental data layers used as input to our models were averaged over

393    much of the twentieth century, which overlaps well with the temporal provenance of the

394    occurrence data. These considerations may introduce minor biases, nonetheless, owing to a

395    rise in global mean surface temperature of about 0.1°C toward the end of the century (Hansen

396    et al. 2010). More generally, the environmental signature for a given pixel was assigned an

397    average value in our analyses, rather than the values for the year in which the sample was

398    collected. Such models with finer temporal resolution can and should be developed to resolve

399    this issue, considering the appropriate environmental regime over the period of specimen

400    collection. We aim to take this "next step" when we can develop both the relevant data layers

401    and occurrence data that are sufficiently rich and dense to permit such analyses. For now,

402    however, our models provide a first estimate of environmental envelopes for these three

403    ciliate species, with the caveat that predictions may involve a slight northern bias that is

404    unlikely to affect the argument regarding the ubiquity hypothesis versus moderate

405    endemicity.

406

407    Predictions across the study area based on different environmental data sets were similar for

408    each species. In all models, minimum and maximum temperature contributed the most to the

409  prediction, although the percent contribution of these two temperature coverages was reduced

410  by the introduction of additional data layers for the Model 2 series. The additional coverages

411  of phosphate and silicate content were most important for *P. gigantea* and *S. serrata*, whereas

412  oxygen saturation and silicate content were most important for *Z. pelagicum*. In contrast, the

413  salinity layers and the apparent oxygen usage layer were relatively unimportant in all models.

414  Model 2 provided finer resolution for each taxon, predicting less of the study area and

415  omitting fewer of the independent evaluation points at the least presence threshold,

416  suggesting a better model. However, the model complexity evaluation indicated clearly that

417  the increased detail of Model 2 for each species did not outweigh the negative effects of

418  increasing model complexity and dimensionality, such that Model 1 was preferable for each

419  species. A clear link exists between the biogeography of each of the species and climatic

420  conditions, so distributions of the species will likely change with climate change. This

421  possibility merits further investigation, considering the current scientific focus on

422  anthropogenic global climate change, and other short-term regional climatic phenomena. This

423  study is a first reference point for the distribution of each species, and can serve as a

424  benchmark against which to compare future distributional patterns.

425

426  **CONCLUSIONS**

427  This study represents a first effort to describe the spatial and environmental distributions,

428  using ecological niche models, of three species of ciliated protists, a group of microorganisms

429  that are essential to marine carbon cycling and trophic chains, marine biodiversity, and even

430  to the fishing industry (Caron and Countway 2009; Lom and Dyková 1992) but for which

431  geographic distributions are poorly documented. This study thus serves as a first reference for

432  the distribution of each species, and as a benchmark against which to compare potential

433  future distributional patterns. Although future work remains to be done to refine these

434    models, particularly to consider climate variability, our findings point clearly to a situation in

435    which each of the studied species has a unique environmental signature and geographic

436    distribution.

437

18

**REFERENCES**

Aguilar, M., A. M. Fiore-Donno, C. Lado, and T. Cavalier-Smith. 2014. Using environmental niche models to test the 'everything is everywhere' hypothesis for *Badhamia*. ISME J. **8:** 737-745.

Aguilar, M., and C. Lado. 2012. Ecological niche models reveal the importance of climate variability for the biogeography of protosteloid amoebae. ISME J **6:** 1506-1514.

Azam, F., T. Fenchel, J. G. Field, J. S. Gray, L. A. Meyerreil, and F. Thingstad. 1983. The ecological role of water-column microbes in the sea. Mar. Ecol.-Prog. Ser. **10:** 257-263.

Barve, N. and others 2011. The crucial role of the accessible area in ecological niche modeling and species distribution modeling. Ecol. Model. **222:** 1810-1819.

Bass, D., and J. Boenigk. 2011. Everything is everywhere: a twenty-first century de/reconstruction with respect to protists., p. 88-110. *In* D. Fontaneto [ed.], Biogeography of microscopic organisms: is everything small everywhere? Cambridge University Press.

Bass, D. and others 2009. A molecular perspective on ecological differentiation and biogeography of Cyclotrichiid Ciliates. J. Eukaryot. Microbiol. **56:** 559-567.

Bass, D., T. A. Richards, L. Matthai, V. Marsh, and T. Cavalier-Smith. 2007. DNA evidence for global dispersal and probable endemicity of protozoa. BMC Evol. Biol. **7:** e162.

Bean, W. T., R. Stafford, and J. S. Brashares. 2012. The effects of small sample size and sample bias on threshold selection and accuracy assessment of species distribution models. Ecography **35:** 250-258.

Beijerinck, M. W. 1913. De infusies en de ontdekking der backteriën. Müller (Reprinted in Verzamelde geschriften van M.W. Beijerinck, vijfde deel, pp. 119–140. Delft, 1921).

Bentlage, B., A. T. Peterson, N. Barve, and P. Cartwright. 2013. Plumbing the depths: extending ecological niche modelling and species distribution modelling in three dimensions. Glob. Ecol. Biogeogr. **22:** 952-961.

Boenigk, J., K. Pfandl, T. Garstecki, H. Harms, G. Novarino, and A. Chatzinotas. 2006. Evidence for geographic isolation and signs of endemism within a protistan morphospecies. **72:** 5159-5164.

Brun, P. and others 2015. Ecological niches of open ocean phytoplankton taxa. Limnol. Oceanogr. **60:** 1020-1038.

Caron, D. A., and P. D. Countway. 2009. Hypotheses on the role of the protistan rare biosphere in a changing world. Aquat. Microb. Ecol. **57:** 227-238.

Cermeño, P., and P. G. Falkowski. 2009. Controls on diatom biogeography in the ocean. Science **325:** 1539-1541.

Dolan, J. R. 2016. Planktonic protists: little bugs pose big problems for biodiversity assessments. **38:** 1044-1051.

Dolan, J. R., R. W. Pierce, and E. J. Yang. 2017. Tintinnid ciliates of the marine microzooplankton in Arctic Seas: a compilation and analysis of species records. Polar Biol. **40:** 1247-1260.

Evans, K. M., V. A. Chepurnov, H. J. Sluiman, S. J. Thomas, B. M. Spears, and D. G. Mann. 2009. Highly differentiated populations of the freshwater diatom *Sellaphora capitata* suggest limited dispersal and opportunities for allopatric speciation. Protist **160:** 386-396.

Fenchel, T. 2008. The microbial loop-25 years later. J. Exp. Mar. Biol. Ecol. **366:** 99-103.

Field, C. B., M. J. Behrenfeld, J. T. Randerson, and P. Falkowski. 1998. Primary production of the biosphere: integrating terrestrial and oceanic components. Science **281:** 237-240.

487  Filker, S., R. Sommaruga, I. Vila, and T. Stoeck. 2016. Microbial eukaryote plankton
488      communities of high-mountain lakes from three continents exhibit strong
489      biogeographic patterns. Mol. Ecol. **25:** 2286-2301.
490  Finlay, B. J., J. O. Corliss, G. Esteban, and T. Fenchel. 1996. Biodiversity at the microbial
491      level: the number of free-living ciliates in the biosphere. Q. Rev. Biol. **71:** 221-237.
492  Foissner, W. 1999. Protist diversity: estimates of the near-imponderable. Protist **150:** 363-
493      368.
494  ---. 2006. Biogeography and dispersal of micro-organisms: a review emphasizing protists.
495      Acta Protozool. **45:** 111-136.
496  ---. 2011. Dispersal of protists: the role of cysts and human introductions, p. 61 – 87. *In* D.
497      Fontaneto [ed.], Biogeography of microscopic organisms: is everything small
498      everywhere? Cambridge University Press.
499  Foissner, W., A. Chao, and L. A. Katz. 2008. Diversity and geographic distribution of ciliates
500      (Protista: Ciliophora). Biodivers. Conserv. **17:** 345-363.
501  Gao, F. and others 2016. The all-data-based evolutionary hypothesis of ciliated protists with a
502      revised classification of the phylum Ciliophora (Eukaryota, Alveolata). Sci. Rep. **6:**
503      e24874.
504  Gentekaki, E., and D. H. Lynn. 2010. Evidence for cryptic speciation in *Carchesium
505      polypinum* Linnaeus, 1758 (Ciliophora: Peritrichia) inferred from mitochondrial,
506      nuclear, and morphological markers. J. Eukaryot. Microbiol. **57:** 508-519.
507  Hansen, J., R. Ruedy, M. Sato, and K. Lo. 2010. Global surface temperature change. Rev.
508      Geophys. **48:** RG4004.
509  Jímenez-Valverde, A., and J. M. Lobo. 2007. Threshold criteria for conversion of probability
510      of species presence to either-or presence-absence. Acta Oecol.-Int. J. Ecol. **31:** 361-
511      369.
512  Katz, L. A., J. DeBerardinis, M. S. Hall, A. M. Kovner, M. Dunthorn, and S. V. Muse. 2011.
513      Heterogeneous rates of molecular evolution among cryptic species of the ciliate
514      morphospecies *Chilodonella uncinata*. J. Mol. Evol. **73:** 266-272.
515  Langer, M. R., A. E. Weinmann, S. Lotters, J. M. Bernhard, and D. Rodder. 2013. Climate-
516      driven range extension of *Amphistegina* (Protista, Foraminiferida): models of current
517      and predicted future ranges. PLoS ONE **8:** e54443.
518  Laval, M. 1968. Zoothamnium pelagicum du Plessis. Cilié péritricheplanctonique:
519      morphologie, croissance et comportement. Protistologica **4:** 333-363.
520  Liu, C. R., P. M. Berry, T. P. Dawson, and R. G. Pearson. 2005. Selecting thresholds of
521      occurrence in the prediction of species distributions. Ecography **28:** 385-393.
522  Logares, R., J. F. Mangot, and R. Massana. 2015. Rarity in aquatic microbes: placing protists
523      on the map. Res Microbiol **166:** 831-841.
524  Lom, J., and I. Dyková. 1992. Protozoan parasites of fishes, p. 315. Developments in
525      aquaculture and fisheries science. Elsevier Science.
526  NOAA. 1999. World Ocean Atlas 1998.
527  Noguez, A. M., H. T. Arita, A. E. Escalante, L. J. Forney, F. García-Oliva, and V. Souza.
528      2005. Microbial macroecology: highly structured prokaryotic soil assemblages in a
529      tropical deciduous forest. Glob. Ecol. Biogeogr. **14:** 241-248.
530  Owens, H. L. and others 2013. Constraints on interpretation of ecological niche models by
531      limited environmental ranges on calibration areas. Ecol. Model. **263:** 10-18.
532  Pearson, R. G., C. J. Raxworthy, M. Nakamura, and A. T. Peterson. 2007. Predicting species
533      distributions from small numbers of occurrence records: a test case using cryptic
534      geckos in Madagascar. J. Biogeogr. **34:** 102-117.
535  Peterson, A. and others 2011. Ecological niches and geographic distributions. Princeton
536      University Press.

537 Peterson, A. T. 2014. Mapping disease transmission risk: enriching models using
538     biogeography and ecology. Johns Hopkins University Press.
539 Peterson, A. T., M. Papeş, and M. Eaton. 2007. Transferability and model evaluation in
540     ecological niche modeling: a comparison of GARP and Maxent. Ecography **30:** 550-
541     560.
542 Phillips, S. J., R. P. Anderson, and R. E. Schapire. 2006. Maximum entropy modeling of
543     species geographic distributions. Ecol. Model. **190:** 231-259.
544 Radosavljevic, A., and R. P. Anderson. 2014. Making better Maxent models of species
545     distributions: complexity, overfitting and evaluation. J. Biogeogr. **41:** 629-643.
546 Richards, T. A., A. A. Vepritskiy, D. E. Gouliamova, and S. A. Nierzwicki-Bauer. 2005. The
547     molecular diversity of freshwater picoeukaryotes from an oligotrophic lake reveals
548     diverse, distinctive and globally dispersed lineages. Environ. Microbiol. **7:** 1413-
549     1425.
550 Saupe, E. E., J. R. Hendricks, A. T. Peterson, and B. S. Lieberman. 2014. Climate change and
551     marine molluscs of the western North Atlantic: future prospects and perils. J.
552     Biogeogr. **41:** 1352-1366.
553 Segovia, B. T. and others 2017. Common and rare taxa of planktonic ciliates: influence of
554     flood events and biogeographic patterns in neotropical floodplains. Microb. Ecol. **74:**
555     522-533.
556 Soberón, J., and A. T. Peterson. 2005. Interpretation of models of fundamental ecological
557     niches and species' distributional areas. Biodiversity Informatics **2:** 1-10.
558 Tamames, J., J. J. Abellán, M. Pignatelli, A. Camacho, and A. Moya. 2010. Environmental
559     distribution of prokaryotic taxa. BMC Microbiol. **10:** article 85.
560 Warren, D. L., R. E. Glor, and M. Turelli. 2008. Environmental niche equivalency versus
561     conservatism: quantitative approaches to niche evolution. Evolution **62:** 2868-2883.
562 Warren, D. L., R. E. Glor, and M. Turelli. 2010. ENMTools: a toolbox for comparative
563     studies of environmental niche models. Ecography **33:** 607-611.
564 Warren, D. L., and S. N. Seifert. 2011. Ecological niche modeling in Maxent: the importance
565     of model complexity and the performance of model selection criteria. Ecol. Appl. **21:**
566     335-342.
567 Wiley, E. O., K. M. McNyset, A. T. Peterson, C. R. Robins, and A. M. Stewart. 2003. Niche
568     modeling and geographic range predictions in the marine environment using a
569     machine-learning algorithm. Oceanography **16:** 120–127.

570

571    **ACKNOWLEDGEMENTS**

583

584    **SUPPLEMENTARY MATERIAL**

585    This document provides a summary for all data management and spatial stratification

586    procedures followed in developing predictive models for three ciliate species; also provided

587    are heuristic estimate of relative contributions of the environmental variables, maps of all 36

588    spatial stratification predictions, results of the tests to determine range restriction using

589    ENMTools background similarity, and a summary environmental layers used for model

590    development.

591

592    **DATA ACCESSIBILITY**

593    Primary occurrence data are fully and openly accessible via the Global Biodiversity

594    Information Facility (GBIF; http://www.gbif.org/) and Ocean Biogeographic Information

595    System (OBIS; http://www.iobis.org/). Environmental coverages were drawn from National

596     Oceanic and Atmospheric Administration's (NOAA) World Ocean Atlas (NOAA 1999). Full

597     results are available in Supplementary Material.

598

599  **FIGURE LEGENDS**

600  **Figure 1**. Global map showing distribution records of *Parafavella gigantea* (A; $N$ = 4495),

601  *Schmidingerella serrata* (B; $N$ = 1536), and *Zoothamnium pelagicum* (C; $N$ = 1738)] based

602  on data drawn from the Ocean Biogeographic Information System (OBIS;

603  http://www.iobis.org), and Global Biodiversity Information Facility (GBIF;

604  http://www.gbif.org/), both accessed 10.12.2015. The sample sizes that we used for model

605  development were greatly reduced as we included only the locality data from vouchered

606  reference material, and removed all duplicate records of a species at any given site.

607

608  **Figure 2**.  Predictions of suitable areas across the study area for each of three species, based

609  on models calibrated using two environmental coverage sets.

610

611  **Figure 3**.  Model 1 background similarity tests comparing the niche occupied by each taxon

612  using the background similarity tests available in the ENMTools R package version 0.1.

613  Schoener's $D$ and a modified Hellinger's $I$ metric, are used to quantify the similarity of two

614  probability distributions. In each case, the niche occupied by each taxon is more different

615  than expected from the study region. Model 2 background similarity tests are shown in

616  Supplementary Material Figure 14.

617

24

618 **TABLES**
619 **Table 1. Summary table of ecological niche modelling results for both models of three**
620 **species of ciliates**. Results of spatial stratification models are shown in Supplementary
621 Material.

| Species | $N$ training / testing points [1] | Testing AUC ratio [2] | Fractional predicted area LPT [3] | Fractional predicted area 10% PT [4] |
|---|---|---|---|---|
| All data predictive models | | | | |
| *P. gigantea* 1 | 184 / 184 | 1.77 | 0.430 | 0.221 |
| *P. gigantea* 2 | 184 / 184 | 1.79 | 0.442 | 0.220 |
| *S. serrata* 1 | 68 / 67 | 1.83 | 0.304 | 0.183 |
| *S. serrata* 2 | 68 / 67 | 1.81 | 0.283 | 0.189 |
| *Z. pelagicum* 1 | 140 / 139 | 1.76 | 0.667 | 0.299 |
| *Z. pelagicum* 2 | 140 / 139 | 1.75 | 0.614 | 0.245 |

622

---

[1] $N$, number of points used for model training / testing.
[2] Training/ testing AUC data, area under the curve of the receiver operating characteristic calculated using NicheToolBox (http://shiny.conabio.gob.mx:3838/nichetoolb2/).
[3] LPT, Least Prediction Threshold.
[4] 10% PT, predictive threshold that excludes the 10% most outlying points.

A. *P.gigantea*

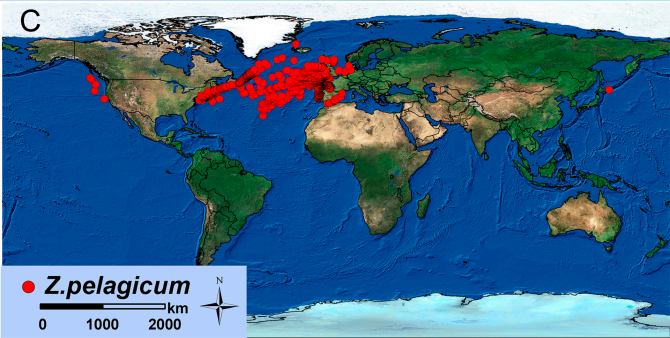B. *S.serrata*

C. *Z.pelagicum*

**Figure 2**.

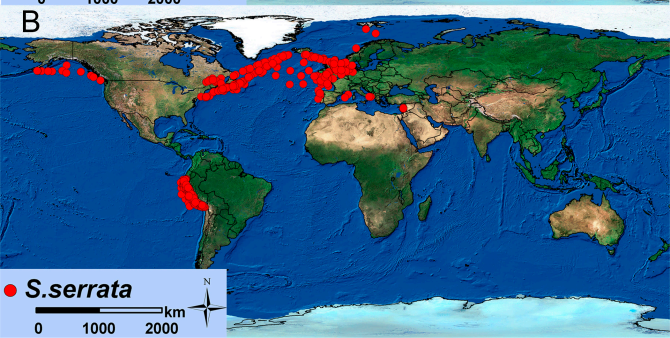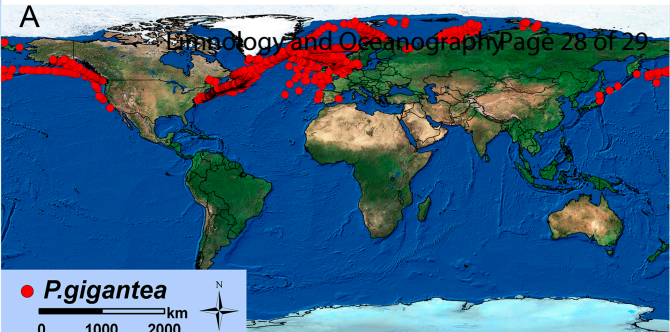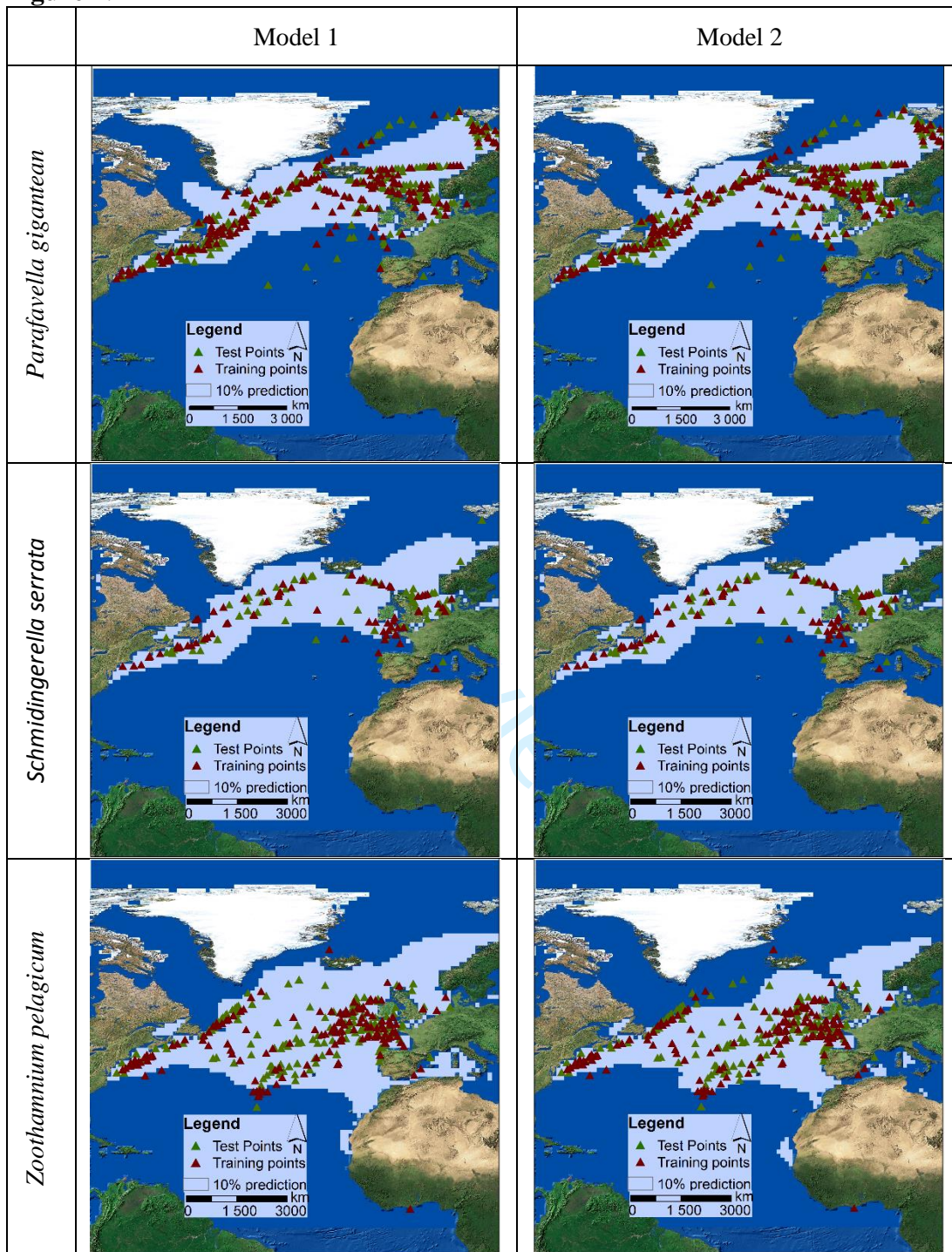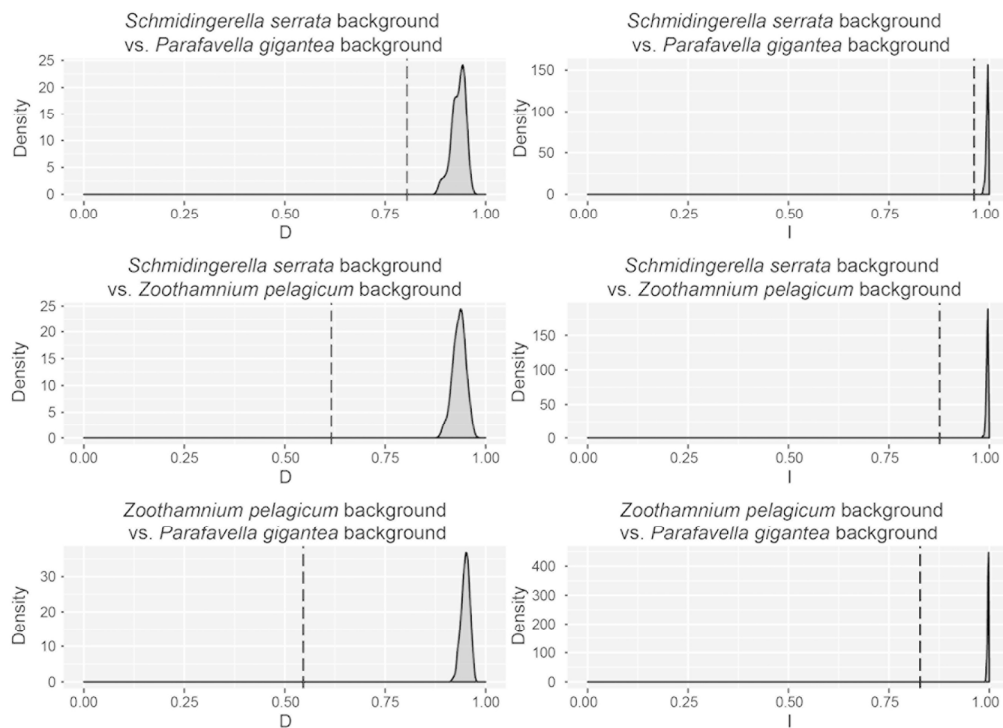| | Model 1 | Model 2 |
|---|---|---|
| *Parafavella gigantean* | | |
| *Schmidingerella serrata* | | |
| *Zoothamnium pelagicum* | | |

Figure 3.  Model 1 background similarity tests comparing the niche occupied by each taxon using the background similarity tests available in the ENMTools R package version 0.1. Schoener's D and a modified Hellinger's I metric, are used to quantify the similarity of two probability distributions. In each case, the niche occupied by each taxon is more different than expected from the study region. Model 2 background similarity tests are shown in Supplementary Material Figure 14.