# Environmental sequencing provides reasonable estimates of the relative abundance of specific picoeukaryotes

Caterina R. Giner, Irene Forn, Sarah Romac, Ramiro Logares, Colomban de Vargas, Ramon Massana

1    A Research Paper submitted to *Applied and Environmental Microbiology*

2

3    # Environmental sequencing provides reasonable estimates of the

4    # relative abundance of specific picoeukaryotes

5

6    Caterina R. Giner[a#], Irene Forn[a], Sarah Romac[b,c], Ramiro Logares[a], Colomban de Vargas[b,c],

7    and Ramon Massana[a#]

8

9    [a] Department of Marine Biology and Oceanography, Institut de Ciències del Mar (ICM-CSIC),

10    Barcelona, Spain.

11    [b] CNRS, UMR 7144, Station Biologique de Roscoff, Roscoff, France.

12    [c] Université Pierre et Marie Curie (UPMC) Paris 06, UMR 7144, Station Biologique de Roscoff,
13    France.

14

15    # Address correspondence to: Caterina R. Giner, caterina@icm.csic.es, and Ramon Massana,

16    ramonm@icm.csic.es

17    Institut de Ciències del Mar (CSIC), Passeig Marítim de la Barceloneta 37-49, 08003

18    Barcelona, Catalonia, Spain. Phone: 34-93-2309500; Fax: 34-93-2309555;

19

20    Running title: Correspondence between HTS and cell abundance

21  **Abstract**

23  High-throughput sequencing (HTS) is revolutionizing environmental surveys of microbial

24  diversity in the three domains of life by providing detailed information on which taxa are

25  present in microbial assemblages. However, it is still unclear how the relative abundance of

26  specific taxa gathered by HTS correlates with cell abundances. Here, we quantified the

27  relative cell abundance of 6 picoeukaryotic taxa in 13 planktonic samples from six European

28  coastal sites using epifluorescence microscopy on TSA-FISH preparations. These relative

29  abundance values were then compared with HTS data obtained in three separate molecular

30  surveys: 454 sequencing the V4 region of the 18S rDNA using DNA and RNA extracts

31  (DNA-V4 and cDNA-V4), and Illumina sequencing the V9 region (cDNA-V9). The

32  microscopic and molecular signals were generally correlated, indicating that a relative

33  increase in specific 18S rDNA was the result of a large proportion of cells in the given taxa.

34  Despite these positive correlations, the slopes often deviated from 1, precluding a direct

35  translation of sequences to cells. Our data highlighted clear differences depending on nucleic-

36  acid template or the 18S rDNA region targeted. Thus, the molecular signal obtained using

37  cDNA templates was always closer to relative cell abundances, while the V4 and V9 regions

38  gave better results depending on the taxa. Our data supports the quantitative use of HTS data

39  but warn about considering it as direct proxy of cell abundances.

42  Key words: 18S rDNA, 454 pyrosequencing, FISH, Illumina sequencing, picoeukaryotes,

43  specific abundances

46 **Importance**
47

48 Direct studies on marine picoeukaryotes by epifluorescence microscopy are problematic due

49 to the lack of morphological features, in addition to the limited number and poor resolution of

50 specific phylogenetic probes used in FISH routines. As a consequence, there is an increasing

51 use of molecular methods, including high-throughput sequencing (HTS), to study marine

52 microbial diversity. HTS can provide a detailed picture of the taxa present in a community,

53 and can reveal diversity not evident using other methods, but it is still unclear the meaning of

54 the sequence abundance in a given taxa. Our aim is to investigate the correspondence between

55 the relative HTS signal and relative cell abundances in selected picoeukaryotic taxa.

56 Environmental sequencing provides reasonable estimates of the relative abundance of specific

57 taxa. Better results are obtained when using RNA extracts as templates, while the region of

58 the 18S rDNA influenced differently depending on the taxa assayed.

59

## Introduction

Protists are key components of marine ecosystems, being major players in the global respiration and production budgets (1, 2) and playing central roles in marine food webs (3). Despite their importance and ubiquity, it was only during the past decade that environmental studies, based on molecular (i.e. culture-independent) techniques, revealed an unsuspected protist diversity in a large variety of marine ecosystems (4-13). These studies were based on the analysis of 18S ribosomal RNA (rRNA) genes retrieved directly from natural assemblages by PCR amplification, cloning and sequencing. Nowadays, the development and use of high-throughput sequencing tools (HTS), e.g. 454 or Illumina, which produce more than thousands of sequences from a single sample, has revolutionized the field, allowing deeper assessments of diversity (14), as well as better estimates of specific relative abundances. One of the main challenges of this approach, however, is to understand the correspondence between the relative abundances of sequences and cells. That is, how close is the specific diversity detected in molecular surveys to the true species composition of natural assemblages.

Few studies have analyzed the relationship between direct microscopic inspections and sequencing data in protists. One of the first studies compared cloning and sequencing results with an accurate list of protists species (5-100 μm size range) identified by microscopy (15). In that case, as the sequencing effort was very limited (less than 100 clones), few of the protists identified by morphology were detected in the sequencing set. In addition, the few sequences obtained did not represent the dominant observed species, a clear sign of the biases in this molecular approach. More recent comparative studies used HTS, therefore were not limited by the sequencing effort, but focused on specific taxa, in particular marine and freshwater ciliates (2, 16-18). Ciliate species have the advantage of having conspicuous morphological traits that allow proper identification by inverted microscopy. In most cases, the same species were found in microscopic and molecular datasets, but the relative

85    abundance of sequences and morphotypes were not in agreement, so each approach revealed a

86    different community structure. Other studies prepared mock communities and the results

87    obtained were similar: all individual taxa were detected, but the relative proportion of

88    sequence types was different from cell mixes (19, 20). Overall, the popularization of HTS

89    now allows a high-resolution exploration of protist richness present in natural samples, yet

90    when it comes to evenness, the picture obtained is still limited.

91        Among protists, picoeukaryotes (protists up to 3 µm in size) are known to be very

92    diverse, widely distributed, and ecologically important in the marine plankton realm (21).

93    Picoeukaryotes are counted as a group by epifluorescence microscopy using a general DNA

94    stain (22) or by flow cytometry (23), but due to their small size and lack of morphological

95    traits (24) they cannot be taxonomically identified by these tools. This can be achieved with

96    FISH (Fluorescence *in situ* Hybridization), which enables the visualization and quantification

97    of specific cells in natural assemblages by using oligonucleotide probes as phylogenetic stains

98    (25). FISH has served to identify the cells from novel environmental clades (11, 26, 27), and

99    has been applied in a few marine surveys (28-31). But this approach is relatively time

100   consuming and targets only one taxon at a time.

101       In this study, we assess the feasibility of using HTS data as a quantitative metric in

102   picoeukaryote diversity studies, by comparing relative HTS read abundances with relative

103   FISH cell counts in selected picoeukaryotic taxa. Differently to the previous studies in which

104   a single taxa (ciliates) or artificial communities were analyzed, here we focus in a set of

105   highly divergent lineages found in geographically separated and unrelated microbial

106   assemblages. Any pattern emerging from this heterogeneous and noisy dataset is expected to

107   be rather robust. We also investigate if there is a difference in community composition

108   assessed by using environmental DNA or RNA extracts as templates (DNA and cDNA reads,

109   respectively), sequencing different regions of the 18S rDNA (V4 versus V9), or using

110    different HTS platforms (454 versus Illumina). To address these questions we used published

111    sequencing datasets from several European coastal samples (Massana *et al.*, 2015 for

112    DNA/cDNA-V4 (32) and Logares *et al.*, 2014 for cDNA-V9 (33)) and chose 6 picoeukaryote

113    taxa (<3 µm) for which we had specific FISH probes for quantification.

114

115    **Materials and Methods**

116    *Sampling*

117         Samples were taken during the BioMarKs project (http://www.biomarks.org) in six

118    European coastal sites: Blanes (Spain, 41° 40' N, 2° 48' E), Gijon (Spain, 43° 40' N; 5° 35'

119    W), Naples (Italy, 40° 48' N, 14° 15' E), Oslo (Norway, 59° 16' N, 10° 43' E), Roscoff

120    (France, 48° 46' N, 3° 57' W) and Varna (Bulgaria, 43°10' N, 28° 50' E) (Table 1). Seawater

121    was collected with Niskin bottles attached to a CTD (conductivity-temperature-depth) rosette

122    at surface and deep chlorophyll maximum (DCM) depths. For molecular surveys, ~20 L of

123    seawater was pre-filtered through a 20 µm metallic mesh and then sequentially filtered

124    through 3 µm and 0.8 µm polycarbonate filters (142 mm diameter). The later filter contained

125    the picoplankton (0.8-3 µm size fraction) and was flash-frozen and stored at -80ºC. The

126    filtration time was less than 30 minutes to avoid RNA degradation.

127         Unfiltered seawater was taken for direct cell counts. For total microscopic counts,

128    seawater samples were fixed with glutaraldehyde (1% final concentration) and left for 1-24 h

129    at 4ºC. Then, aliquotes of 20 ml were filtered through 0.6 µm polycarbonate black filters and

130    stained with DAPI (4',6-diamidino-2-phenylindole, 5 µg ml$^{-1}$). Filters were mounted on a slide

131    and stored at -20ºC until processed. For TSA-FISH (Tyramide Signal Amplification-

132    Fluorescent *in situ* Hybridization) specific counts, aliquotes of 100 ml were fixed with filtered

133 formaldehyde (3.7% final concentration), incubated for 1-24 h in the dark at 4ºC and filtered

134 through 0.6 µm polycarbonate filters (25 mm diameter). Filters were kept at -80ºC until

135 processed. For flow cytometry counting of photosynthetic picoeukaryotes, aliquotes of 1.5 ml

136 were fixed with a mix of paraformaldehyde/glutaraldehyde (1%/0.25% final concentrations),

137 frozen in liquid nitrogen and stored at -80ºC until processed.

138 *Picoeukaryote cell abundance by DAPI staining and flow cytometry*

139 Total cell abundance of picoeukaryotes was estimated in DAPI-stained filters. Cells

140 were counted with an epifluorescence microscope (Olympus BX61) at 1000X under UV

141 excitation, changing to blue light excitation to verify the presence or absence of chlorophyll

142 autofluorescence (phototrophic and heterotrophic cells, respectively). A transect of about 13

143 mm was inspected and cells were classified in size classes: 2 µm, 3 µm, 4 µm, 5 µm and >5

144 µm. All data reported in the study refers to cells within the two smaller size classes (2-3 µm),

145 which account on average for 82% of the cells.

146 Cell abundance of photosynthetic picoeukaryotes was determined in a FACSort flow

147 cytometer by using the red fluorescence signal (chlorophyll) after exciting in a 488 nm laser

148 and the SSC (side-scattered light) of each particle. Fluorescent microspheres (0.95 µm beads)

149 were added as an internal standard (at $10^5$ beads ml$^{-1}$). Data was acquired for 2-4 minutes with

150 a flow rate of 50 to 100 µl min$^{-1}$ using the settings previously described (34).

151 *Cell abundance of specific picoeukaryote taxa by TSA-FISH*

152 The specific oligonucleotide probes used targeted several picoeukaryote taxa: NS4 and

153 NS7 targeted the uncultured clades MAST-4 and MAST-7; CRN02 and MICRO01 the

154 species *Minorisa minuta* and *Micromonas* spp.; PELA01 the class Pelagophyceae; and

155 ALV01 the environmental clade MALV-II (Table 2). These probes have been published in

156    other studies (see Table 2 for references) except NS7. Probe NS7 was designed here with

157    ARB (35) and targeted 91% of the 192 sequences from MAST-7 available in GenBank, had 1

158    mismatch with the remaining MAST-7 sequences, and at least 2 central mismatches with non-

159    target sequences. Probe NS7 gave better signal when combined with oligonucleotide helpers

160    contiguous to the probe region (NS7-HelperA: AACCAACAAAATAGCAC; NS7-HelperB:

161    CCCAACTATCCCTATTAA) that were added in the hybridization buffer at same

162    concentration than the probe. We tested a range of formamide concentration to find the best

163    hybridization condition, and checked that the probe gave negative signal with a variety of

164    non-target cultures. Finally, a probe targeting all eukaryotes (EUK502, 36) was also used. All

165    probes were labeled with horseradish peroxidase (HRP).

166         Hybridizations were performed as previously described (37). Filter pieces (about 1/10)

167    of the 0.6 μm polycarbonate filters were covered with 20 μl of hybridization buffer (40%

168    deionized formamide [except 30% for probe CNR01], 0.9 M NaCl, 20 mM Tris-HCl [pH 8],

169    0.01% sodium dodecyl sulfate [SDS]) and 2 μl of HRP-labeled probes (stock at 50 ng μl$^{-1}$),

170    and incubated overnight at 35ºC. After the hybridization, filter pieces were washed twice for

171    10 min at 37ºC with a washing buffer (37 mM NaCl [74mM NaCl when hybridizing with

172    20% formamide], 5 mM EDTA, 0.01% SDS, 20mM Tris-HCl [pH 8]). Tyramide signal

173    amplification (TSA) was carried out in a solution (1x PBS, 2 M NaCl, 1 mg ml$^{-1}$ blocking

174    reagent, 100 mg ml$^{-1}$ dextran sulfate, 0.0015% $H_2O_2$) containing Alexa 488-labeled tyramide

175    (4 μg ml$^{-1}$), by incubating in the dark at room temperature for 30-60 min. Filter pieces were

176    transferred twice to a phosphate buffer (PBS) bath in order to stop the enzymatic reaction and

177    air dried at room temperature. Cells were countersained with DAPI (5 μg ml$^{-1}$) and filter

178    pieces were mounted on a slide. Targeted FISH cells were counted by epifluorescence under

179    blue light excitation and checked with UV radiation (DAPI staining) for the presence of the

180    nucleus. Cells labeled with the probe EUK502 were counted using the same size classes as for

181    DAPI counts. Data reported refers to cells of 2-3 μm, which account on average for 84% of

182    the cells.

*High-throughput sequencing by 454 and Illumina*

184        HTS data derives from previously published papers taken during the BioMarKs

185    project (http://www.biomarks.org/). Total DNA and RNA from 13 picoplankton samples were

186    extracted simultaneously from the same filter. For RNA extracts, contaminating DNA was

187    removed and RNA was immediately reverse transcribed to cDNA. Data for the 454

188    sequencing derives from Massana *et al*. (32) and used the eukaryotic universal primers

189    TAReuk454FWD1 and TAReukREV3 (38), which amplified the V4 region of the 18S rDNA

190    (~ 380 bp). Amplicon sequencing from DNA and cDNA templates was carried out on a 454

191    GS FLX Titanium system (454 Life Sciences, USA) in Genoscope

192    (http://www.genoscope.cns.fr, France). The complete sequencing dataset is available at the

193    European Nucleotide Archive (ENA) under the accession number PRJEB9133

194    (http://www.ebi.ac.uk/ena/data/view/PRJEB9133). Data for the Illumina sequencing derives

195    from Logares *et al*. (33) and used the eukaryotic universal primers 1398f and 1510r (39),

196    which amplified the V9 region of the 18S rDNA (~130 bp). Paired-end 100 bp sequencing

197    was performed using a Genome Analyzer IIx (GAIIx) system located at Genoscope. Only

198    RNA (cDNA) samples were sequenced with Illumina. Sequences are publicly available at

199    MG-RAST (http://metagenomics.anl.gov) under accession numbers 4549958.3, 4549965.3,

200    4549959.3, 4549945.3, 4549943.3, 4549927.3, 4549941.3, 4549954.3, 4549922.3.

*Sequence analysis of HTS reads*

202        HTS reads by 454 and Illumina were quality checked following similar criteria as

203    detailed in the original papers (32, 33). After the quality control, chimera detection was run

204    with UCHIME (40) and ChimeraSlayer (41) using SILVA108 and $PR^2$ (42) as reference

205   databases. The final curated reads were clustered into OTUs (Operational Taxonomic Units)

206   by using UCLUST 1.2.22 (43) with similarity thresholds of 97% for V4-reads and 95% for

207   V9-reads. Representative reads of each OTU were taxonomically classified by using BLAST

208   against SILVA108, PR$^2$ and a marine microeukaryote database (44). After the taxonomic

209   assignment, metazoan OTUs were removed. From the complete OTU tables for 454 (32) and

210   Illumina datasets (33), the samples targeting the picoplankton were extracted: 13 samples for

211   DNA-V4, 13 samples for cDNA-V4 and 9 samples for cDNA-V9. Then, OTUs corresponding

212   to taxa typically larger than 3 μm (Dinophyceae, Ciliophora, Acantharia, Diatomea,

213   Polycystinea, Raphydophyceae, Ulvophyceae, Rodophyta and Xanthophyceae; in this order

214   of relative abundance) were removed. These groups accounted for 8.0% to 87.7% (average of

215   36.9%) of the 454 dataset and 11.5% to 73.5% (average of 33.9%) of the Illumina dataset.

216   The read number in the final OTU tables of picoeukaryotes was 110,258 for DNA-V4, 77,554

217   for cDNA-V4 and 1,753,600 for cDNA-V9.

218        The relative abundance of the picoeukaryotic groups of interest was retrieved from

219   these taxonomically classified OTU tables, by dividing the number of reads of the specific

220   OTUs corresponding to the groups of interest by the total number of reads in the sample.

221   Altogether, the six taxa of interest accounted for 36.4% of the DNA-V4 reads, 23.5% of the

222   cDNA-V4 reads, and 32.4% of the cDNA-V9 reads. Besides the taxonomic classification of

223   OTUs in the OTU table, we did an additional classification of the unclustered 454 and

224   Illumina reads, to obtain the raw reads for probe checking (see results) and to double-check

225   the taxonomic classification. For this second classification we downloaded GenBank

226   sequences representative of each picoeukaryotic group of interest and used this specific taxa-

227   database to retrieve HTS reads by local BLAST (sequence similarity >97%).

228

229    **Results**

230    *An overview of total picoeukaryote counts in marine coastal waters*

231         We estimated the total cell abundance of picoeukaryotes by epifluorescence

232    microscopy and flow cytometry in 13 planktonic samples taken in six geographically

233    separated European coastal sites and different depths (Table 1). Total picoeukaryote counts

234    (cells <3 μm) by epifluorescence microscopy of DAPI-stained samples revealed a wide range

235    of cell abundances, from 3,139 cells ml$^{-1}$ in Naples-2010 DCM to 24,346 cells ml$^{-1}$ in Oslo-

236    2010 DCM (average of 10,500 cells ml$^{-1}$ in all samples). Phototrophic and heterotrophic cells

237    were differentiated while counting the DAPI samples. The total abundance of phototrophic

238    cells was generally higher than heterotrophic cells (average of of 8,200 and 2,400 cells ml$^{-1}$,

239    respectively), with the exception of Naples-2010 Surface, where both assemblages have

240    similar abundances. In some cases (Blanes, Oslo-2010 DCM, Roscoff and Varna DCM)

241    phototrophic cells were >6 times more abundant than heterotrophic cells. Counts of

242    phototrophic picoeukaryotes obtained by flow cytometry correlated well with the microscopic

243    counts in the 10 samples analyzed (linear slope = 0.74, Pearson's r = 0.9, $P < 0.001$). When

244    forcing the regression line to intercept at 0, the slope was 0.90.

245         The general eukaryotic probe EUK502 was also used to estimate total picoeukaryotic

246    abundance. Cell counts by TSA-FISH were always lower than DAPI counts (60% on average)

247    (Fig. 1). In fact, the sample with the highest total cell abundance was different if estimated by

248    DAPI (Oslo-2010 DCM) or by TSA-FISH (Oslo-2009 Surface). The regression between both

249    datasets was significant, but with a slope very distant from 1 (linear slope = 0.26, Pearsons' r

250    = 0.74, $P < 0.05$). When forcing the line to intercept at 0, the slope was still very low, 0.43.

251    There was some tendency to this discrepancy, as TSA-FISH seemed to underestimate more

252    severely the total cell counts in samples dominated by very small cells. Clearly, DAPI counts

253    provided a better estimate of total picoeukaryotic abundance than TSA-FISH counts, and

254    therefore DAPI counts were used to calculate the relative cell abundances of each of the 6

255    specific picoeukaryotic groups: TSA-FISH counts of each group were at the numerator and

256    total DAPI counts at the denominator.

257    ***Abundance of specific picoeukaryotic taxa***

258    We used TSA-FISH to estimate the total abundance of six groups of picoeukaryotes,

259    chosen because they were well represented in the sequencing datasets of the picoplankton

260    from the studied samples (and poorly represented in the nanoplankton, Table S1). They

261    belonged to different eukaryotic supergroups: the Stramenopiles (MAST clades and

262    Pelagophyceae), Alveolates (the parasite clade MALV-II), Archaeplastida (*Micromonas* spp.)

263    and Rhizaria (*Minorisa minuta*). The taxonomic coverage of the used probes varied from

264    being very narrow targeting a species (*Minorisa minuta*) or a constrained phylogenetic clade

265    (*Micromonas* spp. and the MAST lineages), to being very wide targeting an algal class

266    (Pelagophyceae) or the diverse MALV-II group (formed by 44 phylogenetic clades). The sum

267    of heterotrophic cells (MASTs, *M. minuta* and MALV-II) represented on average 36% of

268    heterotrophic picoeukaryotes counted by DAPI, whereas the phototrophic cells targeted

269    (*Micromonas* and Pelagophyceae) represented on average only 22% of phototrophic

270    picoeukaryotes (Table 1).

271    The cell abundance of the six targeted groups varied strongly among the different

272    samples (Table S2). We found *Micromonas*, MAST-4, MAST-7 and MALV-II as the most

273    abundant taxa (averaged cell abundances of 1492, 279, 160, and 127 cells ml$^{-1}$, respectively),

274    detected in all samples. *Minorisa minuta* was very abundant in some sites, but absent in

275    others. By contrast, Pelagophyceae was the least abundant taxa (averaged cell abundances of

276    59 cells ml$^{-1}$). These cell counts pointed out that each sample contained a different

277    community. *Micromonas* was the most abundant taxa in 7 samples, MAST-4 in 4 samples and

278    *Minorisa* and MALV-II in the other two samples (Table S2).

279    ***In-silico validation of the FISH probes against raw V4-reads***

280        Before applying TSA-FISH, we evaluated the effectiveness of the probes against the

281    V4-reads obtained from the same samples. This analysis was done with raw reads (extracted

282    from the initial dataset by using GenBank sequences of each group as search templates) to

283    take into account all sequence variants. The number of raw reads per group obtained from this

284    way was very similar to the number derived from the OTU table (Table 2). About 1000 to

285    3000 reads were extracted per group (except for MALV-II, about 30,000 reads). Then, we

286    calculated the percentage of raw reads having a 100% match with the probes (Table 2). The

287    five specific probes validated this way retrieved a very high percentage of reads, more than

288    95% in all cases except in MALV-II (83%). Therefore, the vast majority of reads from these

289    five groups in our samples had the target region of the probes.

290        The probe targeting *Micromonas* was not designed at the V4 region of the 18S rDNA,

291    so it could not be directly evaluated with V4-reads from this study. Therefore, we took the

292    OTUs affiliating to *Micromonas* (7 OTUs and 11,166 reads), retrieved the closest GenBank

293    complete sequence from these OTUs (nearly identical at the V4 region), and verified the

294    efectiveness of the probe against these 7 GenBank sequences. Only 3 sequences (accounting

295    for 30% of the reads) exhibited a perfect match, whereas the remaining 4 sequences had a

296    mismatch in the first position of the probe. Thus, probe MICRO01 could be improved perhaps

297    removing the first base, but since this mismatch is located in the first position it likely does

298    not affect the FISH counts.

299

*Comparison of group specific read abundance and TSA-FISH counts*

301    The relative abundance of 454 V4-reads (from DNA and cDNA templates) and

302    illumina V9-reads (from cDNA templates) of each group of interest was compared with the

303    relative cell abundance assessed by epifluorescence microscopy (specific TSA-FISH counts

304    relative to total DAPI counts) in 13 samples for the V4-reads, and 9 samples for the V9-reads

305    (DCM samples from Naples and Oslo were excluded) (Fig. 2). The statistics of these plots are

306    shown in Table 3. For the DNA-V4, the correlation of the relative abundance of cells and

307    DNA reads was significant for all groups ($p < 0.05$) except for MAST-4 and Pelagophyceae,

308    and the goodness of these correlations varied among groups, being strongest for *Minorisa*

309    *minuta* ($R^2 = 0.97$) and weakest for MALV-II ($R^2 = 0.29$). Despite these good correlations,

310    linear slopes of the plots were always different from 1 except in MAST-7. In most cases

311    slopes were below 0.5, indicating an underestimation of cell abundance by 454 reads, while in

312    MALV-II the slope was very high (4.46), indicating a severe overestimation of the molecular

313    signal in this group.

314    By contrast, the correlations between relative cell and read abundances in the cDNA-

315    V4 survey were generally better for all groups, being also significant for Pelagophyceae and

316    MAST-4 (Table 3). Similar to the DNA-V4 survey, each group had a different slope, but in

317    this case there were three taxa (MAST-7, *M. minuta* and *Micromonas*) with slopes statistically

318    not different from 1, indicating that their relative abundances obtained by cell counts and 454

319    reads were comparable. In the six groups analysed, the slopes obtained in the cDNA survey

320    were closer to 1 than the slopes derived from the DNA survey, showing a better performance

321    of the cDNA approach.

322    For the Illumina cDNA-V9 survey, the correlations were slightly worse than for the

323    cDNA-V4 survey (Fig. 2, Table 3), as they were non-significant ($p > 0.05$) for MAST-4 and

324    MAST-7. Regarding the linear slopes, the three groups with a good performance at the

325    cDNA-V4 survey –*M. minuta*, Pelagophyceae and *Micromonas*– had slopes statistically

326    different from 1, indicating that in these groups the V4 region (and not the V9) could be used

327    as a proxy of cell counts. On the contrary, MALV-II had a better correlation with the V9-

328    cDNA reads than with the V4-reads, and its slope was not statistically different from one.

329    This highlights that there is not a "best region" that applies to all taxa.

330    ***Differences when targeting V4 and V9 regions of the 18S rDNA***

331        To discard that the differences observed between the V4 and the V9 regions were due

332    to the use of different sequencing platforms (454 for V4 and Illumina for V9), we sequenced

333    with Illumina (MiSeq platform) the V4 region of one sample of the dataset (Oslo-2009 DCM)

334    using both templates (DNA and cDNA). The relative abundance of ~60 taxonomic groups

335    inferred from the same targeted region (V4) in the two platforms displayed a very good

336    agreement, with $R^2$ of 0.97 and 0.91 (for DNA and cDNA, respectively), and linear slopes of

337    0.92 to 1.02. Both slopes were not significantly different from 1. Furthemore, this analysis

338    was also done in an additional set of 14 samples (from other planktonic size fractions and

339    sediments; data not shown) and both platforms performed similarly, with $R^2$ ranging from

340    0.57 to 1.00 (average of 0.91) and slopes ranging from 0.73 to 1.21 (average of 0.99).

341    Therefore, sequencing the same 18S rDNA region with 454 or Illumina (MiSeq) gave highly

342    consistent results.

343        Therefore, the differences outlined above between V4-454 and V9-Illumina

344    sequencing (Table 3) were due to targeting different 18S rDNA regions and not due to the

345    sequencing platform. In order to observe these differences in more detail, we compared the

346    relative abundance of cDNA-V4 reads and cDNA-V9 reads for the six picoeukaryotic taxa

347    studied here (Fig. 3). Clear and consistent differences were identified in each case. As before,

348  the correlations were good and significant, with $R^2$ ranging from 0.68 to 0.98 (being MALV-

349  II lower, 0.45), but the slopes deviated significantly from 1 ($p < 0.05$). The V9 analysis

350  increased significantly the relative abundance of the stramenopile groups (the two MAST

351  clades and Pelagophyceae), with slopes ranging from 2.3 to 3.4 while it was the opposite for

352  *Micromonas* and MALV-II (slope of 0.2 and 0.3, respectively) and the same for *Minorisa*

353  *minuta* (slope of 1.1).

354

**Discussion**

356      Identifying marine picoeukaryotes by direct microscopy is problematic because of

357  their small sizes, and as a consequence there is an increasing interest in using high-throughput

358  sequencing (HTS) technologies to explore their diversity. HTS surveys provide a detailed

359  picture of the taxa present in the community, including rare species in the assemblage (18,

360  33), and reveal diversity not evident using other methods. However, the interpretation of the

361  HTS signal in terms of total cell abundances is not straightforward. Interestingly, TSA-FISH

362  is able to bridge microscopic and sequencing approaches by using specific phylogenetic

363  probes to estimate true cell abundances (28, 45). FISH, besides being very laborious, is also

364  limited by the number of taxa-specific probes available as well as by their phylogenetic

365  resolution (46). Moreover, TSA-FISH could be inaccurate due to putative mismatches of the

366  probes with the target group, which would result in cell counts underestimates. We addressed

367  this issue by evaluating the six probes against sequences obtained from the same samples, and

368  found an acceptable performance (very good in four cases, 83% of reads for MALV-II and

369  only one terminal mismatch for *Micromonas*). This validated that the TSA-FISH cell counts

370  performed here were accurate and enabled the main objective of this study, which was to

16

371 evaluate how well the HTS signal estimates community structure in terms of specific

372 abundance.

373 *More sequences imply more cells*

374       Since the HTS signal is always relative (number of reads of a given taxa respect to the

375 total read number), we needed the total picoeukaryote abundance to calculate relative cell

376 abundances. In principle, using TSA-FISH with a universal eukaryotic probe would be

377 consistent with the study and would also provide an extra layer of certainty, since it allows an

378 easier differentiation of eukaryotic cells from fluorescent particles and large bacteria.

379 However, TSA-FISH counts systematically resulted in fewer cells than direct DAPI counts,

380 and we noticed protists that were not labeled with the EUK502 probe. Moreover, this

381 discrepancy was particularly critical in samples dominated by very small cells. The wide size

382 spectra of protist cells in natural samples implied a large variation in the fluorescent signal, so

383 small cells with dim fluorescence may remain unnoticed when close to large fluorescent cells,

384 and easily faded away while counting a field having many cells with diverse sizes and

385 morphologies. This problem did not happen when using specific probes, since then we

386 focused in counting a defined cell type (even with dim fluorescence). Therefore, we used the

387 direct DAPI counts to calculate relative cell abundances.

388       When comparing the relative abundance of HTS reads against the relative cell

389 abundance obtained by TSA-FISH for the different taxa, we generally found a good

390 correlation between both methods. The $R^2$ coefficients of each picoeukaryotic taxa were

391 similar in the three comparisons conducted (DNA-V4, cDNA-V4 and cDNA-V9 vs. TSA-

392 FISH), except a very poor correlation for Pelagophyceae in the DNA-V4 survey.

393 Nevertheless, the statistical significance was always better for the cDNA survey than for the

394 DNA. These correlations imply that relative read abundance was proportional to relative cell

395    abundance, i.e. an increase in the HTS signal from a particular taxon is the result of an

396    increase of the proportion of targeted cells in the sample. However, the correlation

397    coefficients were far from 1 in most cases, and this noisy signal was probably related to

398    molecular biases plus the large differences in the picoeukaryotic composition of each sample.

399          Molecular surveys based on a single gene are affected by the widely discussed PCR

400    biases (47). During PCR, some phylotypes can be amplified preferentially, some groups can

401    remain undetected due to primer mismatches (48) or there could be biases due to the number

402    of PCR cycles (49). So, it has been suggested that the relative read abundance can no longer

403    reflect the real composition of the original community, biasing diversity estimates and

404    producing over or underestimations of specific groups (2). Furthermore, sequencing errors

405    may create false or chimeric taxa (16, 50, 51). Our results indicate that PCR biases and

406    putative sequencing artifacts are not affecting proportionality between relative read and cell

407    abundance: more reads imply high proportion of cells. The significant correlations detected

408    here using this sample dataset, where each sample has large differences in the picoeukaryotic

409    composition because they were taken in distant sites and times of the year, justifies the use of

410    relative read abundance as a proxy of community composition for comparative purposes.

411    ***Relative abundances of sequences and cells often disagree***

412          Despite the significant correlations discussed above, HTS and TSA-FISH surveys did

413    not give the same quantitative information, as often the regression line was statistically

414    different from 1. Moreover, these slopes varied strongly among the three HTS surveys. In

415    order to compare these surveys, we analyzed the relative abundances of the six picoeukaryotic

416    groups (among themselves) in the different samples (Fig. 4). This showed a general

417    agreement between TSA-FISH and the two cDNA surveys, but depending on the composition

418    of the sample, the agreement was better using the V4 region or the V9 region. In samples

419    dominated by *Micromonas* (e.g Blanes, Oslo-2010, Roscoff, Varna DCM), the picture

420    obtained with the V4 region matched better the cell abundance, while in samples dominated

421    by stramenopiles (MAST-4, MAST-7, Pelagophyceae), the V9 region performed better. In

422    our samples, the cDNA-V4 survey gives a better representation of the true species

423    composition for 5 of the samples while cDNA-V9 performed better in 4 of the samples.

424         In all cases, the DNA survey gave a more biased perspective of the relative abundance

425    of the 6 picoeukaryotic taxa, being influenced by a very high abundance of MALV-II reads in

426    all samples. This is probably due to a particularly high number of rDNA-operon copies in

427    MALV groups (2, 30, 32). The SSU rDNA copy number can vary orders of magnitude among

428    protist taxa, from few copies per cell in some green algae (52) to about 30 copies in MAST-4

429    (53) or several thousand copies in some dinoflagellates (52), depending on the cell size and

430    genome size (54). Large differences in the copy number of the targeted gene will affect the

431    abundance estimates in DNA surveys (2). Moreover, reads retrieved in DNA surveys could

432    derive from dead organisms or dissolved extracellular DNA. It is known that dissolved DNA

433    is preserved in marine waters (55), escaping from degradation and persisting for different

434    periods of time, from hours to days (56). On the contrary, reads from cDNA surveys derive

435    from ribosomes and represent metabolically active taxa in the community, as ribosomes are

436    needed to perform the RNA translation in metabolically active cells (57, 58). This, in addition

437    to the SSU rDNA copy number, could explain the differences observed between DNA and

438    cDNA surveys. Moreover, our data also highlighted the impact of targeting different regions

439    of the 18S rDNA gene for estimating relative abundances. For example, the cDNA-V9 survey

440    showed a higher signal (more reads) for MAST taxa and a lower signal for *Micromonas* when

441    compared with cDNA-V4. It is known that the range of taxonomic groups detected by V4 and

442    V9 is different (38, 59, 60) and that some groups can be over- or underrepresented. In

443    particular, in our samples the V4 region gives good estimates of cell counts for MAST-7 and

444    *Micromonas* spp., the V9 for MALV-II, and both regions for *Minorisa minuta*. So, the region

445    targeted (and the primers used) is fundamental to interpret any existing molecular data.

446    **Concluding remarks**

447         To our knowledge, this is the first study investigating the correspondence between

448    HTS and cell counts for selected and relevant taxa of marine picoeukaryotes. Indeed, true cell

449    abundances of picoeukaryotic taxa require the TSA-FISH approach, but as this approach has

450    inherent limitations (is time consuming, few probes are available, fine resolution can not be

451    provided), we see the need to pursue with HTS studies. Our results indicate a good correlation

452    between both methods, implying that more cells results in more sequences, although they give

453    different quantitative information, i.e. the relative read abundance cannot be directly related to

454    relative cell abundance. The cDNA-V4 survey showed the best agreement with TSA-FISH

455    abundance, providing 1:1 relationships in half of the assayed taxa, but the cDNA-V9 was best

456    for other taxa. So, the region of the 18S rDNA gene targeted clearly affected the relative

457    abundance of specific taxa. Finally, based in the data mentioned here, we suggest that the

458    sequencing platform used (454 or Illumina) does not produce major biases in diversity. In

459    conclusion, the most quantitative option is to use cDNA templates rather than DNA, while the

460    choice of the targeted region will result in different relative abundances in each particular

461    taxa.

462

466

475

476

477 **References**

478 1.   **Boenigk J, Arndt H**. 2002. Bacterivory by heterotrophic flagellates. Antonie Van

479      Leeuwenhoek **81**:465–480.

480 2.   **Medinger R, Nolte V, Pandey RV, Jost S, Ottenwälder B, Schlötterer C, Boenigk**

481      **J**. 2010. Diversity in a hidden world: Potential and limitation of next-generation

482      sequencing for surveys of molecular diversity of eukaryotic microorganisms. Mol Ecol

483      **19**:32–40.

484 3.   **Sherr EB, Sherr BF**. 2002. Significance of predation by protists in aquatic microbial

485      food webs. Antonie Van Leeuwenhoek **81**:293–308.

486 4.   **López-García P, Rodríguez-Valera F, Pedrós-Alió C, Moreira D**. 2001.

487      Unexpected diversity of small eukaryotes in deep-sea Antarctic plankton. Nature

488      **409**:603–607.

489    5.    **Moon-van der Staay SY, De Wachter R, Vaulot D**. 2001. Oceanic 18S rDNA
490           sequences from picoplankton reveal unsuspected eukaryotic diversity. Nature **409**:607–
491           10.

492    6.    **Amaral-Zettler LA, Gómez F, Zettler E, Keenan BG, Amils R, Sogin ML**. 2002.
493           Microbiology: eukaryotic diversity in Spain's River of Fire. Nature **417**:137.

494    7.    **Dawson SC, Pace NR**. 2002. Novel kingdom-level eukaryotic diversity in anoxic
495           environments. Proc Natl Acad Sci USA **99**:8324–8329.

496    8.    **Stoeck T, Epstein S.** 2003. Novel eukaryotic lineages inferred from small-subunit
497           rRNA analyses of oxygen-depleted marine environments. Appl Environ Microbiol
498           **69**:2657–2663.

499    9.    **Berney C, Fahrni J, Pawlowski J**. 2004. How many novel eukaryotic "kingdoms"?
500           Pitfalls and limitations of environmental DNA surveys. BMC Biol. **2**:13.

501    10.    **Lovejoy C, Massana R, Pedrós-Alió C**. 2006. Diversity and distribution of marine
502           microbial eukaryotes in the Arctic Ocean and adjacent seas. Appl Environ Microbiol
503           **72**:3085–3095.

504    11.    **Not F, Gausling R, Azam F, Heidelberg JF, Worden AZ**. 2007. Vertical distribution
505           of picoeukaryotic diversity in the Sargasso Sea. Environ Microbiol **9**:1233–52.

506    12.    **Guillou L, Viprey M, Chambouvet A, Welsh RM, Kirkham AR, Massana R,**
507           **Scanlan DJ, Worden AZ**. 2008. Widespread occurrence and genetic diversity of
508           marine parasitoids belonging to Syndiniales (Alveolata). Environ Microbiol **10**:3349–
509           3365.

510    13.    **Massana R, Pedrós-Alió C**. 2008. Unveiling new microbial eukaryotes in the surface
511           ocean. Curr Opin Microbiol **11**:213–8.

512    14.    **Bik HM, Sung W, De Ley P, Baldwin JG, Sharma J, Rocha-Olivares A, Thomas**
513           **WK**. 2012. Metagenetic community analysis of microbial eukaryotes illuminates
514           biogeographic patterns in deep-sea and shallow water sediments. Mol Ecol **21**:1048–
515           1059.

516    15.    **Savin MC, Martin JL, LeGresley M, Giewat M, Rooney-Varga J**. 2004. Plankton
517           diversity in the bay of fundy as measured by morphological and molecular methods.
518           Microb Ecol **48**:51–65.

519    16.    **Bachy C, Dolan JR, López-García P, Deschamps P, Moreira D**. 2012. Accuracy of
520           protist diversity assessments: morphology compared with cloning and direct
521           pyrosequencing of 18S rRNA genes and ITS regions using the conspicuous tintinnid
522           ciliates as a case study. ISME J. **7**:244–255.

523    17.    **Santoferrara LF, Grattepanche JD, Katz LA, McManus GB**. 2014. Pyrosequencing
524           for assessing diversity of eukaryotic microbes: Analysis of data on marine planktonic
525           ciliates and comparison with traditional methods. Environ Microbiol **16**: 2752-2763

526    18.    **Stoeck T, Breiner HW, Filker S, Ostermaier V, Kammerlander B, Sonntag B**.
527           2014. A morphogenetic survey on ciliate plankton from a mountain lake pinpoints the
528           necessity of lineage-specific barcode markers in microbial ecology. Environ Microbiol
529           **16**:430–444.

530    19.    **Egge E, Bittner L, Andersen T, Audic S, de Vargas C, Edvardsen B**. 2013. 454
531           Pyrosequencing to describe microbial eukaryotic community composition, diversity
532           and relative abundance: A test for marine Haptophytes. PLoS One **8**: e74371.

533    20.    **Weber AAT, Pawlowski J**. 2013. Can abundance of Protists be inferred from
534           sequence data: A case study of Foraminifera. PLoS One **8**:1–8.

535    21.    **Massana R**. 2011. Eukaryotic picoplankton in surface oceans. Annu Rev Microbiol
536           **65**:91–110.

537    22.    **Porter KG, Feig YS**. 1980. The use of DAPI for identifying aquatic microfloral.
538           Limnol Oceanogr **25**:943–948.

539    23.    **Marie D, Partensky F, Simon N, Guillou L, Vaulot D.** 2000. Flow cytometry
540           analysis of marine picoplankton. In Living Colors: Protocols in Flow Cytometry and
541           Cell Sorting. Diamond, R.A. and DeMaggio, S. (eds). Springer Verlag, Berlin. pp.
542           421–454.

543    24.    **Potter D, Lajeunesse T**. 1997. Convergent evolution masks extensive biodiversity
544           among marine coccoid picoplankton. Biodiv Conserv **107**:99–108.

545    25.    **Delong EF, Wickham GS, Pace NR**. 1989. Phylogenetic stains: Ribosomal RNA-
546           based probes for the idenfication of single cells. Science **243**:1360–1363.

547    26.    **Massana R, Guillou L, Díez B, Pedrós-Alió C**. 2002. Unveiling the organisms behind
548           novel eukaryotic ribosomal DNA sequences from the ocean unveiling the organisms
549           behind novel eukaryotic ribosomal DNA sequences from the ocean. Appl Environ
550           Microbiol **68**:4554–4558.

551    27.    **Chambouvet A, Morin P, Marie D, Guillou L**. 2008. Control of toxic marine
552           dinoflagellate blooms by serial parasitic killers. Science **322**:1254–1257.

553    28.    **Not F, Latasa M, Marie D, Cariou T, Vaulot D, Simon N**. 2004. A single species,
554           *Micromonas pusilla* (Prasinophyceae), dominates the eukaryotic picoplankton in the
555           Western English Channel. Appl Environ Microbiol **70**:4064–4072.

556    29.    **Massana R, Terrado R, Forn I, Lovejoy C, Pedrós-Alió C**. 2006. Distribution and
557           abundance of uncultured heterotrophic flagellates in the world oceans. Environ
558           Microbiol **8**:1515–1522.

559    30.    **Siano R, Alves-de-Souza C, Foulon E, Bendif EM, Simon N, Guillou L, Not F.**.
560           2010. Distribution and host diversity of Amoebophryidae parasites across oligotrophic
561           waters of the Mediterranean Sea. Biogeosciences **8**:267–278.

562    31.    **Lin YC, Campbell T, Chung CC, Gong GC, Chiang KP, Worden AZ**. 2012.
563           Distribution patterns and phylogeny of marine stramenopiles in the North Pacific
564           Ocean. Appl Environ Microbiol **78**:3387–3399.

565    32.    **Massana R, Gobet A, Audic S, Bass D, Bittner L, Boutte C, Chambouvet A,**
566           **Christen R, Claverie JM,  Decelle J, Dolan JR, Dunthorn M, Edvardsen B, Forn I,**
567           **Forster D, Guillou L, Jaillon O, Kooistra W, Logares R, Mahé F, Not F, Ogata H,**
568           **Pawlowski J, Ernice MC, Probert I, Romac S, Richards T, Santini S, Shalchian-**
569           **Tabrizi K, Siano R, Simon N, Stoeck T, Valuot D, Zingone A, de Vargas C.** 2015.
570           Marine protist diversity in European coastal waters and sediments as revealed by high-
571           throughput sequencing. Environ Microbiol **17:**4035–4049.

572    33.    **Logares R, Audic S, Bass D, Bittner L, Boutte C, Christen R, Claverie JM, Decelle**
573           **J, Dolan JR, Dunthorn M, Edvardsen B, Gobet A, Kooistra WHCF, Mahé F, Not**
574           **F, Ogata H, Pawlowski J, Pernice MC, Romac S, Shalchian-Tabrizi K, Simon N,**
575           **Stoeck T, Santini S, Siano R, Wincker P, Zingone A, Richards TA, de Vargas C,**
576           **Massana R**. 2014. Patterns of rare and abundant marine microbial eukaryotes. Curr
577           Biol **24**:813–821.

578    34.    **Marie D, Partensky F, Vaulot D, Brussaard C.** 1999. Enumeration of phytoplankton,
579           bacteria, and viruses in marine samples. Current Protocols in Cytometry. John Wiley &
580           Sons, Inc.

581    35.    **Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner**
582           **FO**. 2013. The SILVA ribosomal RNA gene database project: Improved data
583           processing and web-based tools. Nucleic Acids Res **41**:590–596.

584    36.    **Lim EL, Dennett MR, Caron DA.** 1999. The ecology of *Paraphysomonas*
585           *imperforata* based on studies employing oligonucleotide probe identification in coastal
586           water samples and enrichment cultures. Limnol Oceanogr **44**:37–51.

587    37.    **Pernice MC, Forn I, Gomes A, Lara E, Alonso-Sáez L, Arrieta JM, del Carmen**
588           **Garcia F, Hernando-Morales V, MacKenzie R, Mestre M, Sintes E, Teira E,**
589           **Valencia J, Varela MM, Vaqué D, Duarte CM, Gasol JM, Massana R**. 2015.
590           Global abundance of planktonic heterotrophic protists in the deep ocean. ISME J.
591           **9**:782–792.

592    38.    **Stoeck T, Bass D, Nebel M, Christen R, Jones MDM, Breiner HW, Richards TA.**
593           2010. Multiple marker parallel tag environmental DNA sequencing reveals a highly
594           complex eukaryotic community in marine anoxic water. Mol Ecol **19**: 21–31.

595    39.    **Amaral-Zettler LA, McCliment EA, Ducklow HW, Huse SM**. 2009. A method for
596           studying protistan diversity using massively parallel sequencing of V9 hypervariable
597           regions of small-subunit ribosomal RNA genes. PLoS One **4**:1–9.

598    40.    **Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R**. 2011. UCHIME improves
599           sensitivity and speed of chimera detection. Bioinformatics **27**:2194–2200.

600    41.    **Haas BJ, Gevers D, Earl AM, Feldgarden M, Ward DV, Giannoukos G, Ciulla D,**
601           **Tabbaa D, Highlander SK, Sodergren E, Methé B, DeSantis TZ, Petrosino JF,**
602           **Knight R, Birren BW**. 2011. Chimeric 16S rRNA sequence formation and detection
603           in Sanger and 454-pyrosequenced PCR amplicons. Genome Res **21**:494–504.

604    42.    **Guillou L, Bachar D, Audic S, Bass D, Berney C, Bittner L, Boutte C, Burgaud G,**
605           **de Vargas C, Decelle J, del Campo J, Dolan JR, Dunthorn M, Edvardsen B,**
606           **Holzmann M, Kooistra WHCF, Lara E, Le Bescot N, Logares R, Mahé F,**
607           **Massana R, Montresor M, Morard R, Not F, Pawlowski J, Probert I, Sauvadet**
608           **AL, Siano R, Stoeck T, Vaulot D, Zimmermann P, Christen R**. 2013. The Protist
609           Ribosomal Reference database (PR2): A catalog of unicellular eukaryote small sub-
610           unit rRNA sequences with curated taxonomy. Nucleic Acids Res **41**:597–604.

611    43.    **Edgar RC**. 2010. Search and clustering orders of magnitude faster than BLAST.
612           Bioinformatics **26**:2460–2461.

613    44.    **Pernice MC, Logares R, Guillou L, Massana R**. 2013. General patterns of diversity
614           in major marine microeukaryote lineages. PLoS One **8**: e57170.

615    45.    **Not F, Simon N, Biegala IC, Vaulot D.**. 2002. Application of fluorescent in situ
616           hybridization coupled with tyramide signal amplification (FISH-TSA) to assess
617           eukaryotic picoplankton composition. Aquat Microb Ecol **28**:157–166.

618    46.    **Alonso-Sáez L, Balagué V, Sà EL, Sánchez O, González JM, Pinhassi J, Massana**
619           **R, Pernthaler J, Pedrós-Alió C, Gasol JM**. 2007. Seasonality in bacterial diversity in
620           north-west Mediterranean coastal waters: Assessment through clone libraries,
621           fingerprinting and FISH. FEMS Microbiol Ecol **60**:98–112.

622    47.    **Wintzingerode FV, Göbel UB, Stackebrandt E**. 1997. Determination of microbial
623           diversity in environmental samples: Pitfalls of PCR-based rRNA analysis. FEMS
624           Microbiol Rev **21**:213–229.

625    48.    **Hong S, Bunge J, Leslin C, Jeon S, Epstein S.** 2009. Polymerase chain reaction
626           primers miss half of rRNA microbial diversity. ISME J. **3**: 1365-1373.

627    49.    **Suzuki M, Rappé MS, Giovannoni SJ**. 1998. Kinetic bias in estimates of coastal
628           picoplankton community structure obtained by measurements of small-subunit rRNA
629           gene PCR amplicon length heterogeneity. Appl Environ Microbiol **64**:4522–4529.

630    50.    **Quince C, Lanzén A, Curtis TP, Davenport RJ, Hall N, Head IM, Read LF, Sloan**
631           **WT**. 2009. Accurate determination of microbial diversity from 454 pyrosequencing
632           data. Nat Methods **6**:639–641.

633    51.    **Kunin V, Engelbrektson A, Ochman H, Hugenholtz P**. 2010. Wrinkles in the rare
634           biosphere: Pyrosequencing errors can lead to artificial inflation of diversity estimates.
635           Environ Microbiol **12**:118–123.

636 52. **Zhu F, Massana R, Not F, Marie D, Vaulot D**. 2005. Mapping of picoeucaryotes in
637 marine ecosystems with quantitative PCR of the 18S rRNA gene. FEMS Microbiol
638 Ecol **52**:79–92.

639 53. **Rodríguez-Martínez R, Labrenz M, del Campo J, Forn I, Jürgens K, Massana R.**
640 2009. Distribution of the uncultured protist MAST-4 in the Indian Ocean, Drake
641 Passage and Mediterranean Sea assessed by real-time quantitative PCR. Environ
642 Microbiol **11**:397-408.

643 54. **Prokopowich CD, Gregory TR, Crease TJ**. 2003. The correlation between rDNA
644 copy number and genome size in eukaryotes. Genome **46**:48–50.

645 55. **Danovaro R, Corinaldesi C, Dell'Anno A, Fabiano M, Corselli C.** 2005. Viruses,
646 prokaryotes and DNA in the sediments of a deep-hypersaline anoxic basin (DHAB) of
647 the Mediterranean Sea. Environ Microbiol **7(4)**:586–592.

648 56. **Nielsen KM, Johnsen PJ, Bensasson D, Daffonchio D**. 2007. Release and persistence
649 of extracellular DNA in the environment. Environ Biosafety Res. **6**:37–53.

650 57. **Stoeck T, Zuendorf A, Breiner HW, Behnke A**. 2007. A molecular approach to
651 identify active microbes in environmental eukaryote clone libraries. Microb Ecol
652 **53**:328–339.

653 58. **Not F, del Campo J, Balagué V, de Vargas C, Massana R**. 2009. New insights into
654 the diversity of marine picoeukaryotes. PLoS One **4**:e7143.

655 59. **Dunthorn M, Klier J, Bunge J, Stoeck T**. 2012. Comparing the hyper-variable V4
656 and V9 regions of the small subunit rDNA for assessment of ciliate environmental
657 diversity. J. Eukaryot Microbiol **59**:185–7.

658    60.    **Decelle J, Romac S, Sasaki E, Not F, Mahé F**. 2014. Intracellular diversity of the V4
659          and V9 regions of the 18S rRNA in marine protists (radiolarians) assessed by high-
660          throughput sequencing. PLoS One **9**: e104297.

661

**Figure legends**

663 **Fig. 1.** Comparison of total picoeukaryotic abundance (cells <3 µm) by DAPI counts and

664 TSA-FISH counts using the eukaryotic probe EUK502 in all planktonic samples.

665 **Fig. 2.** Comparison of relative abundance of HTS reads against TSA-FISH cell counts in the

666 13 planktonic samples (9 samples for cDNA-V9 reads) for six picoeukaryotic taxa: MAST-4

667 (a), MAST-7 (b), *Minorisa minuta* (c), Pelagophyceae (d), *Micromonas* spp. (e) and MALV-

668 II (f). Dark blue symbols indicate DNA-V4 reads, light blue cDNA-V4 reads and green

669 cDNA-V9 reads. Regression lines are shown, and their statistics are presented in Table 3.

670 **Fig. 3.** Comparison of relative abundance of V9-Illumina reads and V4-454 reads (cDNA

671 surveys in both cases) in 9 planktonic samples for six picoeukaryote taxa: MAST-4 (a),

672 MAST-7 (b), *Minorisa minuta* (c), Pelagophyceae (d), *Micromonas* spp. (e) and MALV-II (f).

673 **Fig. 4.** Relative abundance of the different groups (among themselves) shown by the four

674 approaches (TSA-FISH, cDNA-V4, DNA-V4, cDNA-V9) in all planktonic samples. Gray

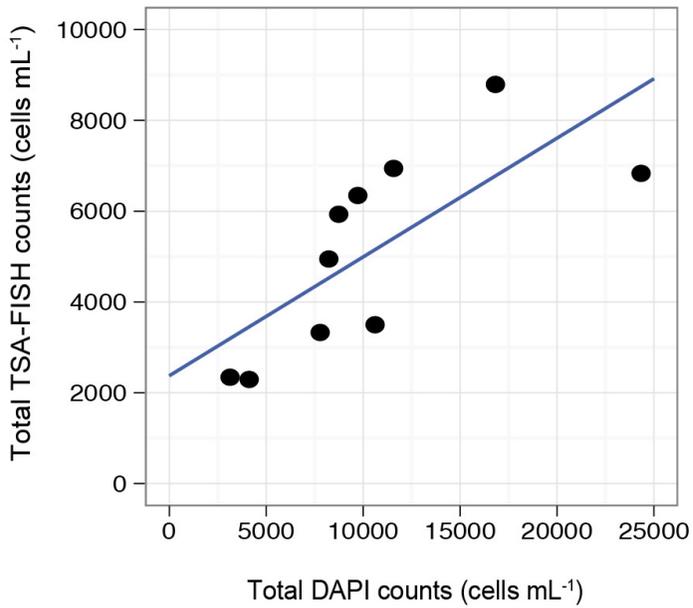675 bars indicate the absence of the sample.

676

Fig. 1- Comparison of total picoeukaryotic abundance (cells <3 μm) by DAPI counts and TSA-FISH counts using the eukaryotic probe EUK502 in all planktonic samples.
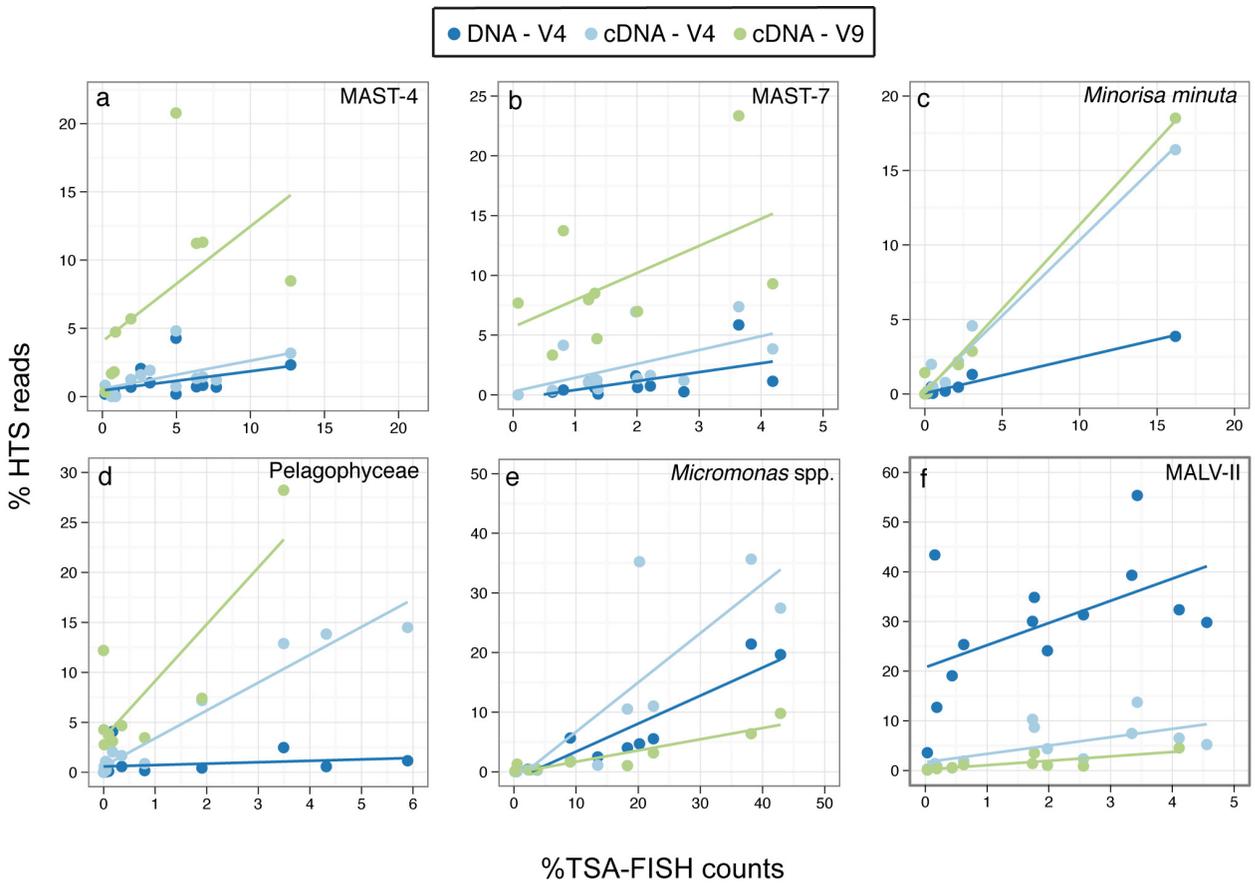
Fig. 2- Comparison of relative abundance of HTS reads against TSA-FISH cell counts
in the 13 planktonic samples (9 samples for cDNA-V9 reads), for six picoeukaryotic taxa:
MAST-4 (a), MAST-7 (b), *Minorisa minuta* (c), Pelagophyceae (d), *Micromonas* spp. (e) and MALV-II (f).
Dark blue symbols indicate DNA-V4 reads, light blue cDNA-V4 reads and green cDNA-V9 reads.
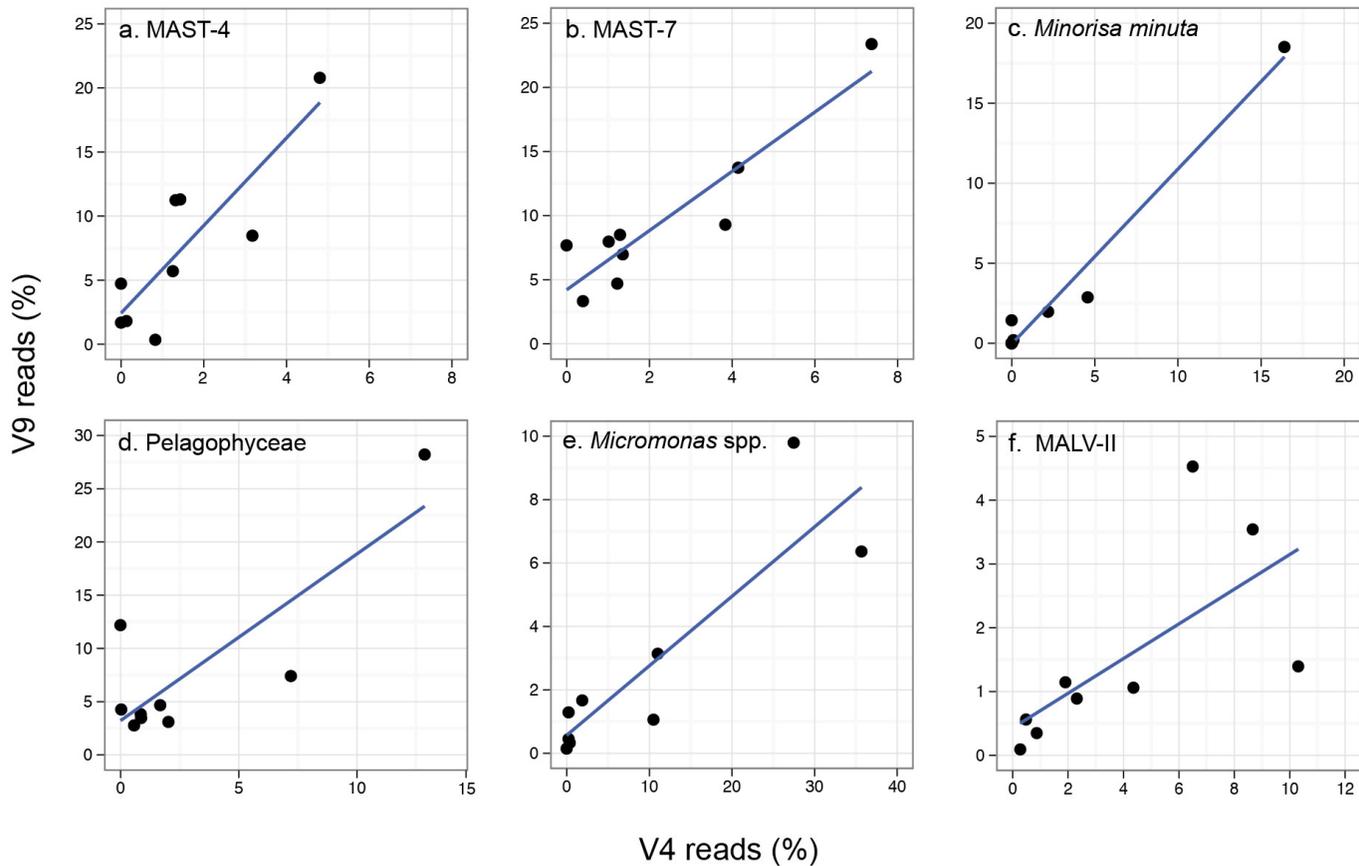Regression lines are shown, and their statistics are presented in Table 3.

Fig. 3 - Comparison of relative abundance of V9-Ilumina reads and V4-454 reads (cDNA surveys in both cases) in 9 planktonic samples for six picoeukaryote taxa: MAST-4 (a), MAST-7 (b), *Minorisa minuta* (c), Pelagophyceae (d), *Micromonas* spp. (e) and MALV-II (f).
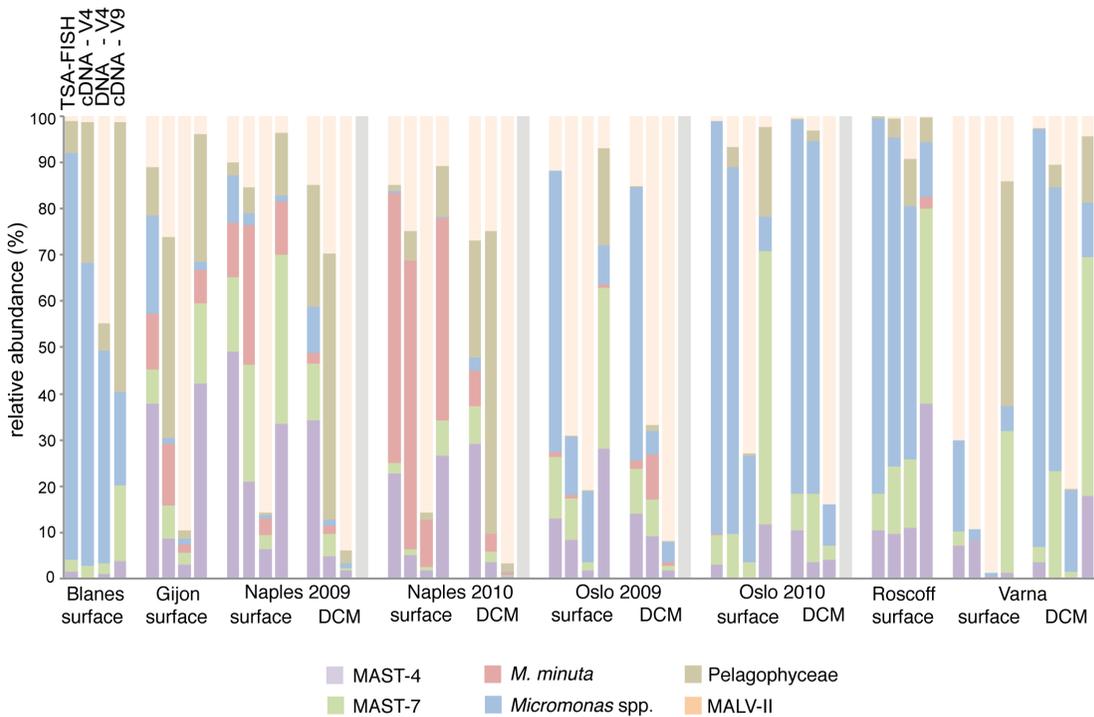
Fig. 4 - Relative abundance of the different groups (among themselves) shown by the four approaches (TSA-FISH, cDNA-V4, DNA-V4, cDNA-V9) in all planktonic samples. Gray bars indicate the absence of the sample.

**Table 1.** Planktonic samples analyzed (sampling site, date, depth and seawater temperature) and cell counts (cells ml$^{-1}$) in these samples: total picoeukaryote abundance (cells ≤3 µm) determined by DAPI (phototrophs and heterotrophs), and photosynthetic picoeukaryote abundance determined by flow cytometry (FC). The last two columns show the percentage of phototrophic and heterotrophic cells explained by the probes used.

| Sampling site | Date | Depth (m) | Temp. (ºC) | DAPI counts Phototr. | DAPI counts Heterotr. | FC counts Phototr. | % Phototr. | % Heterotr. |
|---|---|---|---|---|---|---|---|---|
| **Blanes** | Feb. 2010 | 1 (Surf.) | 12.5 | 9273 | 445 | 9215 | 48.6 | 53.7 |
| **Gijon** | Sep. 2010 | 1 (Surf.) | 20.2 | 1606 | 2503 | 2990 | 14.5 | 20.2 |
| **Naples** | Oct. 2009 | 1 (Surf.) | 22.8 | * | * | 2714 | - | - |
| | | 26 (DCM) | 22.4 | * | * | 2049 | - | - |
| | May 2010 | 1 (Surf.) | 19.2 | 4376 | 4372 | 4700 | 1.1 | 54.6 |
| | | 34 (DCM) | 15.5 | 1808 | 1331 | 1802 | 8.3 | 28.8 |
| **Oslo** | Sep. 2009 | 1 (Surf.) | 15.0 | 12342 | 4470 | 9540 | 12.4 | 21.9 |
| | | 20 (DCM) | 15.0 | 8773 | 2807 | 8930 | 17.9 | 38.4 |
| | Jun. 2010 | 1 (Surf.) | 15.0 | 7727 | 2893 | 13295 | 25.5 | 7.9 |
| | | 10 (DCM) | 12.5 | 21523 | 2823 | 17900 | 22.9 | 40.7 |
| **Roscoff** | Apr. 2010 | 1 (Surf.) | 9.9 | 7203 | 1034 | 8240 | 43.9 | 68.9 |
| **Varna** | May 2010 | 1 (Surf.) | 21.5 | * | * | 3861 | - | - |
| | | 40 (DCM) | 9.5 | 7043 | 731 | 9487 | 24.9 | 24.6 |

* DAPI counts were not performed, so picoeukaryotes could not be differentiated between phototrophs and heterotrophs. In these samples, total picoeukaryote counts were done on FISH filters and were: 4272 cells ml$^{-1}$ in Naples-2009 Surf, 1834 cells ml$^{-1}$ in Naples-2009 DCM, and 4656 cells ml$^{-1}$ in Varna Surf. These values were used in the correlations.

**Table 2.** List of oligonucleotide FISH probes used and effectiveness of the probes against reads from this study (% reads-probe). The table shows the number of 454 reads from each phylogenetic group extracted from the OTU table or from raw reads by local BLAST. The last column shows the percentage of raw reads in each group that have the probe target region with 0 mismatches.

| Probe Name | Target group | Probe sequence (5' – 3') | Probe reference | Num. of reads per Taxa | | % reads - probe |
| --- | --- | --- | --- | --- | --- | --- |
| | | | | In OTU table | From the raw reads | |
| NS4 | MAST-4 | TACTTCGGTCTGCAAACC | Massana *et al.,* 2002 | 2082 | 2082 | 98.0 |
| NS7 | MAST-7 | TCATTACCATAGTACGCA | This study | 2842 | 2833 | 95.7 |
| CRN02 | *Minorisa minuta* | TACTTAGCTCTCAGAACC | del Campo *et al.*, 2012 | 1853 | 1853 | 99.8 |
| PELA01 | Pelagophyceae | ACGTCCTTGTTCGACGCT | Not *et al.,* 2002 | 4440 | 3169 | 98.5 |
| MICRO01 | *Micromonas* spp. | AATGGAACACCGCCGGCG | Not *et al.,* 2004 | 11,166 | - | - |
| ALV01 | MALV-II | GCCTGCCGTGAACACTCT | Chambouvet *et al.*, 2008 | 35,359 | 29,894 | 83.0 |
| EUK502 | Eukaryotes | GCACCAGACTTGCCCTCC | Lim *et al.*, 1999 | - | - | - |

**Table 3.** Statistics ($R^2$, slope value, and p-value) of the correlations between relative abundance of reads and cells in the three molecular surveys: 454 DNA-V4 (Fig. 2, dark blue), 454 cDNA-V4 (Fig. 2, light blue) and Illumina cDNA-V9 (Fig. 2, green). The fourth statistics (p1) compares the slopes against the desired value of 1 (i.e. "ns" indicates that the slope is not significantly different from 1).

| | V4 - 454 survey | | | | | | | | V9 - Illumina survey | | | |
| | DNA | | | | cDNA | | | | cDNA | | | |
| | $R^2$ | slope | p-value | p1 | $R^2$ | slope | p-value | p1 | $R^2$ | slope | p-value | p1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| MAST-4 | 0.18 | 0.14 | ns | - | 0.31 | 0.21 | * | *** | 0.3 | 0.84 | ns | - |
| MAST-7 | 0.33 | 0.75 | * | ns | 0.31 | 1.16 | * | ns | 0.36 | 2.79 | ns | - |
| *Minorisa minuta* | 0.97 | 0.24 | *** | *** | 0.98 | 1.01 | *** | ns | 0.99 | 1.13 | *** | *** |
| Pelagophyceae | 0.06 | 0.14 | ns | - | 0.94 | 2.78 | *** | *** | 0.68 | 5.68 | ** | ** |
| *Micromonas* spp. | 0.87 | 0.47 | *** | *** | 0.73 | 0.83 | *** | ns | 0.87 | 0.2 | *** | *** |
| MALV-II | 0.29 | 4.46 | * | * | 0.39 | 1.68 | * | * | 0.60 | 0.89 | * | ns |

Significance codes: ***: <0.001; **: 0.001–0.01; *: 0.01–0.05; ns: no significant