# Bamboozle: A bioinformatic tool for identification and quantification of intraspecific barcodes

Matthew I M Pinder[†a]*, Björn Andersson[†a], Karin Rengefors[b], Hannah Blossom[b,c], Marie Svensson[b], and Mats Töpel[ad]

[†]These authors contributed equally to the manuscript


[a]Department of Marine Sciences, University of Gothenburg, Göteborg, Sweden
[b]Department of Biology, Lund University, Lund, Sweden
[c]Bigelow Laboratory for Ocean Sciences, Maine, USA
[d]IVL Swedish Environmental Research Institute, Göteborg, Sweden
*Corresponding author: mats.topel@ivl.se

**Keywords:** Metabarcoding, Population genomics, microbial evolution, microalgae, polymorphism,

**Competing Interest Statement:** The authors declare no competing financial interests.

25    **Abstract**

26

27    Evolutionary changes in populations of microbes, such as microalgae, cannot be traced using

28    conventional metabarcoding loci as they lack intraspecific resolution. Consequently, selection

29    and competition processes amongst strains of the same species cannot be resolved without

30    elaborate isolation, culturing, and genotyping efforts. Bamboozle, a new bioinformatic tool

31    introduced here, scans a species' entire genome and identifies allele-rich barcodes that enable

32    direct identification of different strains from a common population, and a single DNA sample,

33    using amplicon sequencing. We demonstrate its usefulness by identifying hypervariable

34    barcoding loci (<500 bp) from genomic data in two microalgal species, the diploid diatom

35    *Skeletonema marinoi*, and the haploid chlorophyte *Chlamydomonas reinhardtii*. Across the

36    genomes, only 26 loci capable of resolving all available strains' genotypes were identified, all of

37    which are within protein-coding genes of variable metabolic function. Single nucleotide

38    polymorphisms (SNPs) provided the most reliable genetic markers, and amongst 55 strains of *S.*

39    *marinoi,* three 500 bp loci contained, on average, 46 SNPs, 103 unique alleles, and displayed

40    100% heterozygosity. The prevalence of heterozygosity was identified as a novel opportunity to

41    improve strain quantification and detect false positive artefacts during denoising of amplicon

42    sequences. Finally, we illustrate how metabarcoding of a single genetic locus can be used to

43    track strain abundances of 58 strains of *S. marinoi* in an artificial selection experiment. As future

44    genomics datasets become available and DNA sequencing technologies develop, Bamboozle has

45    flexible user settings enabling optimal barcodes to be designed for other species and applications.

2

46 **Introduction**

47 Microalgae are a diverse paraphyletic group of aquatic microbes responsible for half of global

48 primary production (Falkowski et al. 1998; Field et al. 1998). Many species have enormous

49 population sizes compared with multicellular organisms and contain some of the highest

50 intraspecific genetic diversity observed within eukaryotes (Flowers et al. 2015). Despite their

51 ecological importance, knowledge about individual species' functional diversity and evolutionary

52 potential is limited (Godhe and Rynearson 2017). Microalgal blooms can contain thousands to

53 millions of different clones (Sassenhagen et al. 2021), but designs of evolution experiments are

54 currently constrained to only one or a few strains as they generally utilise clonal cultures, e.g.,

55 (Lohbeck et al. 2012; Schaum et al. 2017; Schaum et al. 2018; Sefbom et al. 2015; Wolf et al.

56 2019). If these strains are co-cultured in competition or *in situ*, as needed for many evolutionary

57 and ecological questions, the individual strains cannot readily be identified or quantified using

58 existing microscopic or genotyping technologies. Currently, genotyping of microalgae relies

59 heavily on microsatellites (Rengefors et al. 2017), although newer sequencing methods such as

60 RAD-seq provide a method for fine-scale genotyping and genomic information (Rengefors et al.

61 2021). Although suitable for population genetic and genomic studies, these methods still rely on

62 single DNA samples from each strain for genotype identification. This severely restricts these

63 methods' usefulness in quantitative evolutionary studies, where strain isolation is time-

64 consuming or even impossible if the species is difficult to culture (Rengefors et al. 2021), or if

65 the population's clonal diversity is large (Schaum et al. 2017; Scheinin et al. 2015).

66 An alternative approach to track cellular abundances involves using massively parallel amplicon

67 sequencing of a diverse locus. The main advantage of such an approach is that millions of

68 individual DNA molecules, and by extension the cells they originate from, can be quantified

3

69    from a single DNA sample. This simple principle - metabarcoding - has revolutionised the field

70    of microbial ecology in the past decade (Hebert et al. 2003; Taberlet et al. 2012). Depending on

71    the taxonomic group and the resolution required, various highly conserved genes are commonly

72    used in metabarcoding, including ribosomal RNA genes (16S, 18S, 23S, or 28S), mitochondrial

73    cytochrome c-oxidase subunit 1 (COI), and the chloroplastic ribulose-1,5-bisphosphate

74    carboxylase/oxygenase large subunit (*rbcL)* (Guo et al. 2015; Pujari et al. 2019; Yamada et al.

75    2017). Metabarcoding enables a relatively simple and standardised quantification of ecosystem-

76    wide parameters, such as diversity, presence of low abundance or cryptic taxa, and community

77    composition. However, the conserved nature of the available metabarcoding loci restricts their

78    usefulness as intraspecific markers (Canesi and Rynearson 2016; Godhe et al. 2006; Guo et al.

79    2015).

80    With a sufficiently diverse intraspecific barcode marker it should be possible to track the

81    abundances and quantify fitness of individual strains from a mixed microbial population.  Such

82    an approach would be analogous to 'barcoding analysis by sequencing', or Bar-seq experiments,

83    where unique barcodes are used to replace genes in knock-out mutants (Robinson et al. 2014).

84    Bar-seq experiments then uses amplicon sequencing to quantify fitness changes in thousands of

85    non-lethal gene deletion mutants simultaneously, using a co-culture experimental design (Kim et

86    al. 2010; Li et al. 2019; Robinson et al. 2014; Smith et al. 2009), and it has accelerated the

87    functional annotations of several microbial species' genomes (Dent et al. 2005; Giaever and

88    Nislow 2014; Han et al. 2010). Access to 'intraspecific' DNA barcode loci in microalgae and

89    other microbes would provide a powerful tool to address fundamental questions about their

90    evolution and population structures, analogous to how metabarcoding and Bar-seq have radically

4

91    changed our understanding of microbial ecology and phenotypic functions of genes in the last

92    decade.

93    Genome sequencing projects are beginning to resolve the genetic diversity between species of

94    microalgae (Armbrust et al. 2004; Bowler et al. 2008; Merchant et al. 2007; Mock et al. 2017;

95    Read et al. 2013; Worden et al. 2009), and population genomic datasets are becoming

96    increasingly available (Blanc-Mathieu et al. 2017; Flowers et al. 2015; Osuna-Cruz et al. 2020;

97    Rastogi et al. 2020; von Dassow et al. 2015), providing information about intraspecific genetic

98    diversity. From this data, it should be possible to identify novel barcoding loci with intraspecific

99    resolution. Here we introduce Bamboozle, a bioinformatic approach that uses whole genome

100   sequencing data to identify intraspecific barcodes by scanning a species' entire genome for the

101   most suitable sites. As a proof-of-concept, we identify and design primers for a set of barcodes in

102   the diploid marine diatom *Skeletonema marinoi* and the haploid model chlorophyte

103   *Chlamydomonas reinhardtii*. Furthermore, we demonstrate how one single barcode locus can be

104   used to track strain selection in an artificial evolution experiment incorporating 58 environmental

105   strains of *S. marinoi* from two Baltic Sea inlets on Sweden's east coast. Based on the data from

106   this experiment we identify opportunities, as well as challenges, in employing amplicon

107   sequencing to address questions regarding evolutionary processes in microbes. Although we see

108   obvious applications in microalgae and other microbes, Bamboozle can identify intraspecific

109   barcodes for both haploid and diploid species when provided with a reference genome and whole

110   genome sequencing data from multiple individuals.

5

**New approaches**

Our aim with the Bamboozle tool was to enable identification of novel intraspecific barcoding loci with resolution down to the level of genotypes originating from one common population. To this end, we needed to identify loci fulfilling three criteria that render them suitable as barcoding loci, as outlined by Kress and Erickson (2008): 1) they should be short enough to allow amplification and sequencing with current technologies, 2) they should have significant interspecific (in our case intraspecific) variability, and 3) they should have conserved flanking regions that allow the binding of primers. Our approach builds on existing approaches for identifying novel barcoding loci by scanning the entire genome (Angers-Loustau et al. 2016; Paracchini et al. 2017). In contrast to the aforementioned studies, we focus on identification of barcodes with intraspecific resolution, implement several filters to optimise allele richness, exclude unsuitable genomic regions, and minimise the amount of manual screening needed.

The steps of the Bamboozle approach for barcode identification are outlined in Fig. 1 and described in detail below. The required input data includes a reference genome for the organism of interest, and whole genome sequencing (WGS) data from the strains to be analysed. Throughout the paper we use the strain concept (Lakeman et al. 2009), which for most practical purposes is synonymous with genotype in multicellular organisms. The WGS data is required in two formats, one pair per strain: 1) a sorted BAM file generated by aligning the WGS reads to the reference, and 2) single nucleotide polymorphism (SNP) data in VCF format generated from the aforementioned BAM file using GATK (Van der Auwera and O'Connor 2020).

In the case of diploid species, phased BAM and VCF files (i.e., where data from each allele is in a separate file) are also required to analyse the sequences of individual alleles. The main

6

133    Bamboozle tool, written and tested in Python version 3.7.10, consists of six main steps, detailed

134    below.

135    *1. Identification of regions of unusual coverage.* When mapping whole genome shotgun

136        sequencing reads to a reference genome, one expects a similar read coverage across the

137        entirety of the reference. While short regions deviating from this expectation could result

138        from potentially informative deletions or insertions in the mapped strain versus the reference,

139        longer such regions could be due to misassemblies in the reference genome, repeat regions

140        with significant length differences between strains, or gene duplications. As these

141        phenomena could be problematic in identifying appropriate barcoding loci and downstream

142        amplicon sequencing, SAMtools (Danecek et al. 2021) and Bedtools (Quinlan and Hall 2010)

143        are called by Bamboozle and used to determine the median read depth of each contig in each

144        strain (`samtools depth`), and to identify regions which deviate heavily from this median

145        value (`bedtools genomecov`). By default, regions with coverage less than 50% or more than

146        200% of the contig median in that strain are flagged, in order to filter out large heterozygous

147        deletions and gene duplications, respectively.

148    *2. Identification of variant positions.* Using a custom algorithm written in Python version

149        3.7.10, the locations of all variants in each strain are extracted from their respective VCF

150        files and consolidated.

151    *3. Genome screening based on conserved sites.* As one of the core requirements for a suitable

152        barcode is conserved flanking regions for primer binding, windows containing such

153        conserved regions are identified in this step. Bamboozle traverses each contig/chromosome

154        (step size 1; window size defined by the user), saving for further analysis those windows with

7

155    conserved regions (i.e., regions without any SNP's) at each end, based on the variant list

156    from step 2. The size of the window and the conserved end regions can be defined by the

157    user with the `--window_size` and `--primer_size` options, respectively, to account for

158    different amplification and sequencing strategies.

159    *4. Merging of overlapping windows.* Where both conserved end regions of two windows

160    overlap, these two are merged. For example, a 500 bp window with conserved regions from

161    contig positions 1-21 and 480-500, and a second such window with conserved regions from

162    contig positions 2-22 and 481-501, would be merged to form a single 501 bp window with

163    conserved regions from contig positions 1-22 and 480-501. This is done to reduce the

164    number of overlapping windows to be screened manually after Bamboozle outputs the end

165    results. For a visual example, see step 4 of Fig. 1.

166    *5. Exclusion of windows with unusual read coverage.* The merged windows are now compared

167    to the coverage data generated in step 1, and windows that overlap with a region of read

168    coverage fluctuation by ten or more bases are excluded. As length variations can be

169    potentially valuable in terms of differentiating between strains, allowing short (<10 bp)

170    stretches of irregular coverage (indicative of short indel regions) is intended as a compromise

171    between the potential for informative windows and the complications such regions may

172    introduce, as noted in step 1.

173    *6. Identifying informative loci.* As a unique sequence is required to identify and quantify each

174    strain in, for example, a strain selection or phenotyping experiment, SAMtools (`samtools`

175    `faidx`) and BCFtools (`bcftools consensus`) (Danecek et al. 2021) are applied to each

176    remaining window to generate consensus sequences for each strain (and each allele, in the

177    case of a diploid organism) using the reference genome and the respective VCF files. The

178    individual allele sequences are then analysed to determine if each strain in the dataset

179    contains at least one unique allele, in which case the window is reported to the user as a

180    potentially suitable barcoding locus.

181   As output, Bamboozle reports the identified barcoding loci in 1) a tab-separated file, giving

182   coordinates and metadata for each locus, 2) a BED file of locus coordinates, intended for easy

183   visualisation, such as in a genome browser, and 3) a multi-FASTA file for each locus, containing

184   the sequence of each allele in each strain. This enables quick and easy manual quality control of

185   the output, and design of primers for PCR tests and amplicon sequencing. As a proof-of-concept,

186   we performed amplicon sequencing on one locus identified in the diploid *S. marinoi* and used a

187   custom-made script that merges amplicons from paired-end Illumina data using BBMerge

188   (Bushnell et al. 2017), denoises the data using DADA2 (Callahan et al. 2016), and takes

189   advantage of heterozygosity to quality control the output and translate allele abundances into

190   counts of individual strains. Bamboozle offers a complete pipeline from WGS data through to

191   amplicon sequencing analysis of the barcoding loci, and its performance has been fine-tuned and

192   benchmarked using the diploid diatom *S. marinoi*.

193 **Results**

194 *Performance of the bioinformatic pipeline: millions of potential barcoding loci reduced to*

195 *twenty-six*

196 As expected, common metabarcoding loci lacked intraspecific resolution in *S. marinoi*, with 80-

197 100% of strains having the same single allele (Table S1). We applied Bamboozle to WGS data

198 from 54 strains of *S. marinoi*, to identify novel barcoding loci for use with 2x301 bp paired-end

199 Illumina MiSeq sequencing. Two of the initial 56 datasets were excluded from the analysis, one

200 due to low WGS read coverage, and one for being a clone of another strain in the study. Initial

201 attempts to identify suitable barcoding loci of 300-400 bp failed in this species (data not shown),

202 so we scanned for loci of 500 bp. Within the *S. marinoi* genome, most (99.5%) of the 54,753,029

203 possible 500 bp windows were filtered out due to either coverage fluctuations (main problem,

204 Table 1) or unconserved flanking regions. Of the remaining 263,000 windows, the average

205 frequency of SNP positions was 12±11, and the average number of alleles was 14 (Table 1). Four

206 potential barcoding loci were identified as containing at least one unique allele per strain and

207 passing all filters (Table 1). Three of these, plus an additional hypervariable locus (*Sm_C2W24*)

208 from a previous iteration of Bamboozle that did not include a coverage filter (step 1, Fig. 1),

209 were selected for further evaluation. Compared to other 500 bp windows passing the read depth

210 and conserved flanking region filters, three barcodes proved especially rich in terms of SNPs and

211 unique alleles, falling in the 97.6th to 99.9th percentiles, respectively (Fig. 2A and B).

212 *Sm_C2W24* was predicted to contain an even higher SNP frequency (Fig. 2B). The *Sm_C2W24*

213 locus was retained through the analysis to assess the performance of Bamboozle in identifying

214 loci without coverage filter implementation.

10

215     The Bamboozle tool was also applied to seventeen haploid *C. reinhardtii* strains to identify loci

216     of ~300 bp. While the initial analysis on all strains returned no suitable barcoding loci,

217     separation by mating type (eight mt+ strains and nine mt- strains) was successful. From a

218     reference genome of 107,613,365 bp [nuclear genome, mitochondrion, plastid, and mt- locus

219     (Merchant et al. 2007)], with 107,161,335 potential 300 bp windows, only 137,000 loci (0.13%)

220     contained less than two-fold coverage deviation compared with the contig median (again the

221     main issue, Table 1) and had conserved flanking regions. About 40% of remaining windows

222     contained only one allele in the population, and less than 1% contained more than seven (Fig. 2C

223     and E). Only 22 of these loci were identified by Bamboozle as having one unique allele for all

224     strain genomes - fifteen for mt+ and seven for mt- (Table S2). In addition to being significantly

225     more diverse in allelic richness, these Bamboozle-identified loci also contained a much higher

226     number of SNPs than other loci across the genome (three out of four ranked in the 99.8th

227     percentile; Fig. 2D and F).

228     *In silico annotation of intraspecific barcodes and common features*

229     All 26 barcodes identified in *S. marinoi* and *C. reinhardtii* were inside predicted protein-coding

230     genes (Table S2). Several barcodes were also in close physical proximity within the same, or

231     adjacent genes, e.g., *Sm_C12W1*, *-2*, and *-3*, located within a 2,197 bp region containing two

232     gene models (Sm_t00009768-RA and Sm_t00009769-RA), as well as *Cr_Chr9W-12* and *-13*

233     inside *CHLRE_09g389134v5*, and *Cr_Chr11W-16, -17,* and *-18,* inside *CHLRE_11g467528v5*.

234     The predicted function of the proteins did not have a consistent pattern across the two species,

235     and ranged from ribosomal proteins (Sm_t00008465-RA [ribosome biogenesis protein WDR12]

236     and     *CHLRE_10g447800v5*     [60kDa     SS-A/Ro     ribonucleoprotein]),     to     enzymes

237     (*CHLRE_03g207250v5*     [putative     glutamine     synthetase]     and     *CHLRE_13g592050v5*

11

238     [allantoinase]) and ion transporters (*CHLRE_11g467528v5* [calcium channel]). Eleven barcodes

239     were located inside genes without conserved domains or with conserved Domains of Unknown

240     Function (DUF).

241     Four barcodes of each species were selected for primer design and further barcode development

242     (Fig. 3). Among these, no part of the ~500 bp barcoding loci was intronic in *S. marinoi*. Instead,

243     three out of four barcodes spanned domains with repetitive amino acid regions (*Sm_C2W24*,

244     *Sm_C12W1*, and *Sm_C16W4*) with the conserved primer sites anchored inside the same or

245     flanking conserved domains (Fig. 3A). In contrast, for *C. reinhardtii,* the major part of the 300

246     bp barcoding loci was intronic, while the conserved primer sites were located either in introns or

247     exons (Fig. 3C). A more detailed annotation of the barcode-containing genes is provided in the

248     Supplemental Information.

249     As future barcode applications may involve co-cultivation with other species, or *in situ*

250     experiments in natural microalgal communities, we evaluated to what extent the primers were

251     species-specific. In the target species, primers could be designed to amplify the selected loci

252     around the 21 bp conserved regions (Fig. 3B and D), although strain amplification appeared

253     uneven between strains for *Sm_C12W2, Sm_C16W4, Cr_Chr3W1*, and *Cr_Chr1W3*, suggesting

254     possible PCR amplification bias. Homologous genes for the *S. marinoi* barcode loci were

255     identified in three other species of centric diatoms within the order *Thalassiosirales*

256     (*Thalassiosira pseudonana* and *Thalassiosira oceanica,* and *Skeletonema subsalsum*) (Table S3).

257     In the two *Thalassiosira* species, each locus's primer sites contained mismatches in at least three

258     positions. However, in *S. subsalsum* the *Sm_C12W2* locus contained only one mismatch, and this

259     was also the only locus that yielded a PCR product using *S. subsalsum* DNA as template. None

260     of the primer pairs amplified in eleven other Baltic Sea phytoplankton species (Fig. S1A and B).

261    Although the 54 strains of *S. marinoi* all came from a remote population of *S. marinoi* inside the

262    Baltic Sea, the primer sites were conserved also in three strains from the American Atlantic

263    Coast and the Mediterranean Sea (Table S3), suggesting that the barcodes will function outside

264    the local Baltic Sea population. PCRs of the *C. reinhardtii* barcode loci *Cr_Chr9W12* and

265    *Cr_Chr3W10* did not amplify in two other chlorophytes, including another *Chlamydomonas* sp.

266    strain, suggesting they could also be species-specific (Fig. S1C).

267    *Accuracy of bioinformatic predictions and amplicon sequencing of barcodes in S. marinoi*

268    The allele sequences predicted by Bamboozle were validated using amplicon sequencing of

269    individual strains' PCR products, or Sanger sequencing for the haploid *C. reinhardtii*. This was

270    done for two reasons: firstly, to generate the accurate allele sequences in the strains and contrast

271    this with Bamboozle's predictions based on genomic shotgun sequences, and secondly, to assess

272    what PCR and sequencing artefacts were prevalent.

273    Although *Sm_C2W24* was identified by an earlier version of Bamboozle that did not employ a

274    coverage filter, we retained this locus throughout the analysis to illustrate issues that micro- or

275    minisatellite-like loci can produce in the pipeline. Bamboozle predicted that *Sm_C2W24*

276    contained 95 SNPs and five indels within our studied population, but manual inspection

277    indicated that some strains had low sequencing coverage across the most variable region.

278    Amplicon sequencing revealed length differences of 198 bp amongst alleles of *Sm_C2W24* in the

279    *S. marinoi* population, caused by variable copy number (2 to 24 copies) of a 9 bp repeat

280    encoding Alanine- Asparagine- Glutamic acid (AN(E)) (length difference also confirmed by gel-

281    electrophoresis, Fig. 3B). The variant calling algorithm of GATK, and by extension Bamboozle,

282    had misinterpreted these repeats as a high abundance of SNPs (Table 2). Consequently, almost

283    all allele-sequences in the different strains were incorrectly predicted, and per-base accuracy

284    across the locus was only 82.6%. Hence, as outlined in Fig. 1, adjustments were made to

285    Bamboozle to filter out similar regions, and *Sm_C2W24,* together with 90-99% of the genomes

286    (Table 1), was disregarded in subsequent analysis.

287    All strain-specific sequences of the three *S. marinoi* barcoding loci, as well as *Cr_Chr9W12* and

288    *Cr_Chr3W10* in *C. reinhardtii*, identified with the coverage filter were predicted with >99.9%

289    accuracy (Table 2). Looking into the reason for the remaining discrepancies, we discovered that

290    some of the strains were triploid (GP2-4_54 and VG1-2_78; Fig. 4), confounding the software's

291    expectations of diploidy. In other cases, the inaccuracies were caused by insufficient filtering of

292    the variant calling data, resulting in spurious base-calling errors.

293    Once the inaccuracies in the reference alleles were corrected using the confirmed sequences,

294    merged amplicons that still did not match the expected allele sequences were investigated to

295    determine the sources of error. This artifact search was limited to *S. marinoi* genotype samples

296    with coverage >100 amplicons per locus. A detailed description of the analysis and its findings

297    can be found in the Supplemental Information (Fig. S2, S3 and text). Briefly, 61-89% of

298    amplicons contained sequencing-errors, with 39-76% failing to merge. 11 to 40% of merged

299    amplicons mapped perfectly to known alleles, while remaining sequences were explained by

300    erroneous base-calling (10-40%), off-target amplification (0.1-12%), and PCR chimeras (1-3%).

301    Most amplicon sequence-errors could be either corrected or filtered out using the denoising

302    algorithm of DADA2 (Callahan et al. 2016). We finetuned the settings for the 523 bp

303    *Sm_C12W1* locus to minimize false negative (alleles not predicted as real Amplicon Sequence

304    Variants [ASV]) and false positive (any ASVs not corresponding to a biological allele)

14

305 observations (Table S6). From mixed DNA samples, the most stringent settings evaluated failed

306 to predict three out of 110 alleles, while 16 artefact sequences were retained as ASV, although

307 they comprised only 0.03% of all merged sequences. By analysing the patterns of the remaining

308 false positive ASVs across these samples, searching for linked alleles, we were able to iteratively

309 identify the heterozygous alleles of three strains for which we lacked genotype sequencing data

310 (GP2-4_32, VG1-2_65, VG1-2_99). The remaining false positives were annotated as either

311 chimeras that DADA2's *de novo* chimera filter failed to remove, or one bp sequencing errors

312 from an abundant allele. Relaxing filtering settings, resulted in higher false positive ASVs (346,

313 totalling 1.2% of all amplicons), but detection of all known alleles as ASVs (Table S6). Using

314 the later settings on a mixed DNA sample, the *Sm_C12W1* alleles could be used to enumerate the

315 abundance of 58 strains accurately (Fig. 4).

316 *Implementation of an intraspecific metabarcoding method in a microbial evolution experiment*

317 The quantitative performance of *Sm_C12W1* as a metabarcoding locus was further evaluated in

318 an artificial evolution experiment using *S. marinoi*. In the artificial evolution experiment, 30

319 strains from Gåsfjärden (VG) and 28 strains from Gropviken (GP) were density normalised

320 (relative cell abundance of each strain was 3.4±0.89% based on cell counts using microscopy),

321 pooled separately for each population, and subjected to 42 days (50-100 generations) selection

322 with and without toxic copper stress. This data was used to evaluate the quantitative performance

323 of *Sm_C12W1*'s ability to track strain abundances.

324 Despite our expectations of diploidy in the *S. marinoi* strains, there were signs of polyploidy in

325 the experimental data. Some strains exhibited a skewed allele ratio closer to 1:3 (GP2-4_43,

326 GP2-4_55, GP2-4_63, VG1-2_45, and VG1-2_61; Fig. 5) suggesting that they could be

15

327 polyploids or that the PCR reaction amplified alleles with slightly different efficiencies. The

328 combined allele abundance of strains was sometimes over-represented two- to three-fold (e.g.,

329 GP2-4_57, VG1-2_45, VG1-2_90; Fig. 4) while others were underrepresented at a similar

330 magnitude (VG1-2_103, VG1-2_105, VG1-2_56; Fig. 4). This could not be explained by

331 differences in cell density ($R^2$=0.0003, slope=-0.0059, N=58), again suggesting ploidy

332 differences, or different DNA extraction efficiencies between cells of strains. Despite these

333 subtle deviations from expectations, the barcode *Sm_C12W1* could be used to identify all strains

334 at close to the expected relative abundances, most often with strains two alleles at close to a 1:1

335 ratio (Fig. 4). As the populations evolved through selection on strain diversity, one or two strains

336 ultimately became dominant in each treatment, and during this process, their allelic ratio was

337 conserved at 1:1 (Fig. 5). These observations show that relative changes in allele frequencies

338 should reflect changes in the relative abundance of strains.

339 Throughout the experiment, most strains experienced negative selection. At around 10 ASV

340 observations per sample, DADA2's algorithm started filtering out observations (compare Fig. 5,

341 with non-denoised sequence matches in Fig. S4), indicating that false-negative artefacts were

342 created. However, DADA2 effectively removed all but four false-positive observations in the

343 experiment (allele of GP population seen in the VG population, and vice versa: Tables S6). The

344 remaining four false positives appeared to be chimeric, as did allele GP2-4_44#1, which

345 displayed a skewed allelic ratio representing up to 5% (>1,000 copies) of amplicons in several

346 replicate samples, while GP2-4_44#2 was not detected at all (Fig. 5). As the first 316 bp of this

347 allele are identical with GP2-4_27#1, and the last 354 bp with GP2-4_27#2 (the two alleles of

348 the most dominant strain in the same saples), there is a central 147 bp region where a PCR

349 chimera from these two alleles would produce GP2-4_44#1. In addition to enabling chimera

16

350    detection in situations like the one above, heterozygosity provided an additional analytical

351    advantage, enabling differential equations to accurately partition alleles that were shared between

352    two or three strains (Fig. 4 and 5). These results show that heterozygosity can be used to both

353    filter out chimeric observations, and to enable improved quantification of strains, even if they

354    lack a unique allelic marker.

355 **Discussion**

356 By scanning the entire genome for suitable loci, we illustrate how the Bamboozle approach can

357 identify hyper-variable barcodes with intraspecific resolution. In the case of both *S. marinoi* and

358 *C. reinhardtii*, most of the genome (>90%) was not suitable for barcoding and was effectively

359 filtered out by the algorithm. Of the remaining loci fulfilling the criteria for barcodes (Kress and

360 Erickson 2008), only a fraction contained sufficient allelic richness to provide allele-rich markers

361 for strain quantification, highlighting the power of a whole-genome scan approach for

362 development of novel intraspecific barcodes. Using a combination of original (54) and published

363 [17: (Flowers et al. 2015; Ness et al. 2016)] WGS datasets, we illustrate how such intraspecific

364 barcodes can be designed from either publicly available sequencing data, or through the

365 generation of new data for a species or population of interest.

366 Within our amplicon sequencing data, many amplicons (~40%) did not correspond to any known

367 allele in our database. Therefore, we attempted to identify the sources of these artefacts, and how

368 to address them. This was done in order to make recommendations for future studies using

369 workflows similar to ours.

370 First, we note that the 523 bp *Sm_C12W1* locus, paired with Illumina MiSeq's 2×300 bp

371 sequencing technology, creates some challenges that may not be encountered in shorter loci. For

372 example, we could not use DADA2 as a stand-alone tool to process reads since it truncated

373 100% of all amplicons during merging (Table S6). BBMerge performed better but still failed to

374 merge 78% of all reads, due to the short overlapping region and low base-call quality at the 3'

375 end of reads. By pairing BBMerge with DADA2, we found a flexible approach that could correct

376 99-99.97% of all sequencing errors in reads that had merged, with only spurious false-positive

18

377   allelic observations (Fig. 5). Although a shorter locus is preferable and will likely perform better

378   than this, our analysis shows that the tools available to make robust biological inferences from

379   noisy amplicon data extend to 500 bp barcode loci, albeit with a loss of four-fifths of the

380   sequence data. If future studies cannot identify window sizes smaller than 500 bp, accurate long-

381   read sequencing of barcodes (e.g., PacBio's Circular Consensus Sequences technique (Rhoads

382   and Au 2015)), or paired end 300 bp sequences from the newer platform Illumina NovaSeq 6000

383   could be options worth considering. Long-read sequencing, in combination with DADA2's error

384   correction, has enabled differentiation between pathogenic and benign *E. coli* strains using single

385   SNP markers and metabarcoding of the whole ~1450 bp 16S rRNA gene (Callahan et al. 2019).

386   Such a sequencing approach opens up for loci of several kb to be used as intra-specific barcodes.

387   In contrast to sequencing errors, PCR chimeras provided a more problematic type of artefact. We

388   found that PCR chimeras could not be removed effectively using software built for processing

389   metabarcoding data, such as Mothur (Rognes et al. 2016; Schloss et al. 2009), UCHIME2 (Edgar

390   2016), or DADA2s chimera removal algorithm (Callahan et al. 2016). Exacerbating this issue is

391   the potential for 'perfect fake' sequences (chimeras identical to another, non-chimeric sequence in

392   the dataset (Edgar 2016)). Out of the 110 alleles in *Sm_C12W1,* we identified several cases of

393   such 'perfect fakes' (e.g., allele 1 in GP2-4_44; Fig. 5), and these instances should become

394   progressively more numerous if more strains and alleles are included in the experiments and

395   analysis. However, we could identify these PCR chimeras in our experimental design primarily

396   by looking at deviation from the expected heterozygous allelic ratios of individual strain (Fig. 5).

397   Depending on the research question and experimental design, chimeras can be more or less

398   effectively filtered out using similar approaches in future studies. If chimeras risk causing

399   erroneous biological interpretations from the data, approaches to reduce the number of chimeras

400    produced includes optimising a low PCR cycle threshold (Smyth et al. 2010) or using chimera-

401    free PCR kits, such as emulsion PCR (Williams et al. 2006).

402    The allele phasing step (upstream of Bamboozle) also created some problems, as it could not

403    handle cases of triploidy (Fig. 4). Haplotype phasing of polyploids is non-trivial and benefits

404    greatly from long-read data (Abou Saada et al. 2021). However, as long as polyploidy is

405    accurately identified (e.g., by amplicon sequencing of each strain's barcode individually), it

406    should not provide any issues with downstream analysis of strain abundances and may even aid

407    in providing increased allelic diversity for quantification purposes. Consequently our evaluation

408    of Bamboozle as both a tool to identify intra-specific barcodes suggest it can accurately locate

409    such hypervariable regions. Our quantitative analyse show that with appropriate choice of

410    barcode lengths, PCR kits, sequencing platforms and denoising strategies, amplicon sequencing

411    of intra-specific barcodes can be used to track the abundances of large numbers of strains from

412    mixed DNA samples.

413
414    **Conclusions**

415    We built and employed the Bamboozle tool to develop metabarcoding loci that provide

416    intraspecific resolution in the diploid diatom *S. marinoi*. This allowed us to identify and quantify

417    the abundance of 58 environmental strains from mixed DNA sample. We illustrated the

418    usefulness of such an intraspecific barcode locus in an artificial evolution experiment where we

419    tracked the abundance of all strains during co-cultivation. Although further tests are needed to

420    assay performance in new populations and environmental samples, as well as the amplicon

421    sequencing performance of the *C. reinhardtii* loci, our preliminary tests suggest that intra-

422    specific metabarcoding can simplify or enable novel studies of microbial ecology and evolution

20

423    in natural and artificial settings. This is not limited to artificial selection experiments like ours

424    but possibly also to determine genotype diversity in natural populations, studies of ancient DNA

425    from sediment cores (Adams et al. 2019; Härnström et al. 2011), mesocosm incubations of

426    natural phytoplankton communities (Scheinin et al. 2015; Tatters et al. 2013), effects of

427    predation (Sjöqvist et al. 2014), nutrient competition within (Collins 2011) and between species

428    (Descamps-Julien and Gonzalez 2005), or other drivers that are challenging or impossible to

429    observe in mono-culture experiments (Baert et al. 2016; Collins and Schaum 2021). Using the

430    haploid chlorophyte *C. reinhardtii*, we illustrate how the Bamboozle tool should, if possible, be

431    able to identify loci in both haploid and diploid organisms with a reference genome and whole

432    genome sequencing data from multiple genotypes. Our comparative analysis of barcoding loci

433    identified in *S. marinoi* and *C. reinhardtii* suggests that much of the genome is unsuitable for

434    barcode design, and that different species have different optimal barcode loci. Consequently,

435    Bamboozle's strategy of scanning the entire genome of a species is a suitable approach to

436    identify novel barcodes for evolutionary studies.

437    **Acknowledgments**

442    **Data availability**

443 The code for Bamboozle is available at https://github.com/topel-research-group/Bamboozle and

444 the WGS and amplicon sequencing data has been deposited at NCBI under BioProject

445 PRJNA939970.

446

22

**Materials and Methods**

*Strain retrieval and whole genome sequencing*

Individual strains of *S. marinoi* were germinated and isolated from sediment resting stages using standard micropipetting techniques (Härnström et al. 2011). Surface sediment was collected from two semi-enclosed inlets of the brackish Baltic Sea, where one (Gåsfjärden [VG]; 57°34.35'N 16°34.98'E) has been exposed to historical copper mining (Söderhielm and Sundblad 1996) while the other (Gropviken [GP]; 58°19.92'N 16°42.35'E) has not. A total of 69 and 55 strains were isolated from VG and GP, with 88% and 94% survival, respectively. Strains were cultured in locally-sourced seawater (salinity 7), which had been sterile-filtered (Sarstedt's [Helsingborg, Sweden] 0.2 mm polyethersulfone membrane filter), and amended with f/2 nutrients (Guillard 1975) and 106 µM $SiO_2$. Cultures were continuously screened for contamination by other microalgae, and auxospore formation or bimodal cell sizes, which indicate sexual inbreeding in *S. marinoi* (Ferrante et al. 2019), and such cultures were discarded.

DNA extraction was performed within one month of revival from resting stages in a total of 28 strains from GP and 30 from VG. Cultures were made temporarily axenic via a combination of mechanical cleaning (triple washing in 20 µg mL$^{-1}$ Triton X-100 media and cell collection on a 3.0-µm polycarbonate filter), followed by a five-day antibiotic cocktail treatment (90 µg mL$^{-1}$ Paromomycin, ten µg mL$^{-1}$ Ciprofloxacin, and 40 µg mL$^{-1}$ Cefotaxime). On day six, cells were collected via centrifugation in 50 mL Falcon tubes (10 000 x g, for ten min.), flash-frozen in liquid nitrogen, and stored at -80ºC. DNA was extracted using a CTAB-phenol–chloroform protocol as described in (Godhe et al. 2001), but with additional RNA digestion during cell lysis (65ºC for 60 min. using 1 mg RNaseA mL$^{-1}$ CTAB buffer). DNA yield was quantified using

23

469 Qubit (Thermo Fisher Scientific). In two subsequential attempts, strains VG1-2_65 and VG1-

470 2_99 did not survive the antibiotic treatment, and extractions yielded insufficient amounts of

471 DNA for PCRs.

472 Whole genome sequencing was performed on the remaining 56 *S. marinoi* strains using ½ lane

473 of an Illumina NovaSeq S4 flowcell. Library preparation was done using Nextera's PCR-based

474 protocol with enzymatic fragmentation. This resulted in 722.82 Mreads (2x150 bp paired-end).

475 *Data preparation for the pipeline – S. marinoi*

476 Trimming of reads was performed on the whole genome sequencing data to remove the first 15

477 bases of each read. The trimmed reads were mapped to the *S. marinoi* reference genome version

478 1.1.2 (https://albiorix.bioenv.gu.se/Skeletonema_marinoi.html) using Bowtie2 version 2.3.4.3

479 (Langmead and Salzberg 2012), and the resultant SAM files were converted to sorted BAM files

480 and indexed using SAMtools version 1.9 (Danecek et al. 2021). Variant calling was performed

481 using GATK version 4.1.8.0 (Van der Auwera and O'Connor 2020), following the Best Practices

482 recommendations on the tool's website, and the resultant VCF files were indexed using

483 BCFtools version 1.10.2 (Danecek et al. 2021). In addition, SAMtools version 1.10 was used to

484 phase the BAM files, and GATK was re-run on each of these phased BAM files individually.

485 These phased files were indexed as described above.

486 *Data preparation for the pipeline – C. reinhardtii*

487 Sequencing data from the strains of interest was downloaded from either the NCBI Sequence

488 Read Archive (SRA; downloaded using the `prefetch` and `fastq-dump` tools from the SRA

489 toolkit version 2.9.6 (Leinonen et al. 2010)) or the European Nucleotide Archive (ENA;

24

490    downloaded    using    the    `enaDataGet`    tool    from    enaBrowserTools    version    1.6

491    [https://github.com/enasequence/enaBrowserTools]) (Table S4; retrieved 22 February 2021).

492    PCR    duplicates    were    removed    using    the    filterPCRdupl    Perl    script    version    2.3

493    (https://github.com/linneas/condetri), then trimmed with Cutadapt version 3.2 (Martin 2011),

494    using a quality threshold of 30 and a post-trimming length threshold of either 90 (ENA reads) or

495    45 (SRA reads).

496    The trimmed reads were mapped to a reference consisting of the following sequences from

497    *C. reinhardtii*: the contig-level assembly of the reference nuclear genome version 5.5 (accession

498    no. ABCN00000000.2); the mitochondrial genome (accession no. NC_001638.1); the plastid

499    genome (accession no. NC_005353.1); and the mating type minus locus (accession no.

500    GU814015.1). Mapping was performed using Bowtie2 version 2.3.4.3 (Langmead and Salzberg

501    2012), with the resultant SAM files being converted to sorted BAM files and indexed using

502    SAMtools version 1.10 (Danecek et al. 2021). Variant calling and VCF indexing were performed

503    as described above for *S. marinoi.*

504    *Bamboozle workflow*

505    Barcodes for *S. marinoi* were obtained using commit 75467da of the Bamboozle main branch

506    (except for *Sm_C2W24*, which was obtained using commit ae28b1b). Barcodes for *C. reinhardtii*

507    were obtained using commit db64632 of the `matt_improvement` branch. The latter commit has

508    the fourth and fifth steps of the pipeline reversed (i.e., windows with irregular coverage are

509    excluded prior to window merging), along with other minor improvements.

510    The input files (reference genome FASTA files, plus BAM and VCF files [and respective

511    indexes] for each strain) were prepared as described above.

25

512    Bamboozle was run in both species with a `-window_size` parameter of both 500 and 300, and a

513    `-primer_size` parameter of 21.

514    *Development of the bioinformatic pipeline*

515    The predicted barcodes were first visualised against the reference genomes using IGV version

516    2.6.3 (Robinson et al. 2011) to look for anomalies. In an early iteration of the pipeline (commit

517    ae28b1b) without coverage depth filters (steps 1 and 5, Fig. 1), many loci contained obvious

518    errors in the conserved primer sites, or spanned regions of low read coverage (e.g., *Sm_C2W24*)

519    suggesting the pipeline could not effectively process regions with large indels (e.g. micro- or

520    mini-satellites). This issue has largely been resolved by adding the coverage depth filter, as

521    evidenced by the reduction of identified barcodes in *S. marinoi* from 245 to 4.

522    Additional adjustments of the pipeline used for analysing both *S. marinoi* and *C. reinhardtii*,

523    included extending the code to allow analysis of haploid genomes, and steps 4 and 5 of the

524    pipeline were switched (i.e., windows with irregular coverage are excluded prior to window

525    merging), in order to avoid excluding potentially informative windows without coverage issues.

526    *In silico annotation of intraspecific barcodes*

527    The gene models of the predicted barcodes were retrieved from version 5.5 of the *C. reinhardtii*

528    reference genome (https://ensembl.gramene.org/Chlamydomonas_reinhardti) and version 1.1.2

529    of the *S. marinoi* reference genome (https://albiorix.bioenv.gu.se/Skeletonema_marinoi.html).

530    Functional domains were annotated using NCBI's CD-Search webtool (Marchler-Bauer and

531    Bryant 2004). When gene models lacked a functional annotation, putative functions were

532    assigned by a homology search using BLASTp against the NCBI database (Altschul et al. 1990),

533    and comparison of functional domain structure.

26

534  To assess primer specificity outside our Baltic Sea strains of *S. marinoi,* we mapped

535  transcriptional data from three strains from global culture collections (Keeling et al. 2014), and

536  manually looked for variants within the predicted conserved regions. A targeted BLASTn search

537  was used to identify orthologous genes in three species within the order *Thalassiosirales* with

538  reference genomes, namely: *Thalassiosira oceanica* (Lommer et al. 2012), *Thalassiosira*

539  *pseudonana* (Armbrust et al. 2004), and the closely related *Skeletonema subsalsum* (Sarno et al.

540  2005) using an in-house draft genome for the latter (Pinder et al. unpublished data).

*Confirmation of the bioinformatic predictions*

542  Illumina amplicon sequencing was used to assess the sequence accuracy of each of Bamboozle's

543  locus predictions in the diploid *S. marinoi,* and Sanger sequencing for the haploid *C. reinhardtii*.

544  Primers (standard de-salted oligos) were designed for the 21 bp conserved regions flanking the

545  variable region (Table S5) and extended with Illumina adapter sequences for *S. marinoi. rbcL*

546  was used as a positive control of amplification in all diatoms (Guo et al. 2015), and 16s or 18S

547  rRNA across other phytoplankton taxa (Stoeck et al. 2010; Sundberg et al. 2013). The barcodes

548  were amplified from 100 ng DNA and a final primer concentration of 0.1 µM using the

549  Phusion® High-Fidelity PCR Kit (Thermo Scientific) with individually optimised annealing

550  temperatures between 62 and 65ºC. A one-step PCR reaction, with adapter extended primers,

551  was run for 30 cycles, with 5 s denaturation (98ºC), 5 s annealing, 30 s extension (72ºC), and a

552  final 5 min extension. The evenness of amplification across strains was qualitatively assessed

553  using gel electrophoresis. Strains with poor barcode amplification were re-run with either

554  increased DNA concentration or up to 35 PCR cycles. Species specificity of the primers was

555  assessed using DNA from two strains of *Skeletonema subsalsum*, as well as a mixture of other

556  common phytoplankton species from the Baltic Sea (see Supplemental information). The *C.*

557    *reinhardtii* primers were only tested on two other green algae strains; *Chlamydomonas* sp and

558    *Microglena* sp.

559    For the 56 *S. marinoi* strains, the barcode alleles were sequenced individually using 2x301 bp

560    reads on 1/10 of a lane on the Illumina MiSeq platform with v3 chemistry (expected 1.8 million

561    reads in total). The four PCR reactions were run independently for each barcode, and the

562    products pooled based on volume. Further library preparation, including amplicon size selection

563    >300 bp, Nextera dual-indexing, and amplicon sequencing were performed at SciLifeLab (NGI

564    Stockholm), according to the manufacturer's instructions.

565    *Identifying sources of error in the amplicon sequencing results*

566    Paired end reads were first quality filtered and trimmed using Cutadapt version 3.2 to remove

567    adapter and primer sequences, with a quality threshold of 28 and a minimum read length of 180

568    bp. The resultant trimmed reads were merged using BBMerge version 38.86 (Bushnell et al.

569    2017). Amplicon sequences from the individually strains that did not match the two expected

570    allele sequences were investigated to determine the sources of error. Chimera detection was

571    initially performed using the `chimera.vsearch` function of Mothur version 1.47.0 (Schloss et al.

572    2009). However, as this highlighted fewer chimeras than we knew to exist through manual

573    inspection, we instead switched to the `uchime2_ref` function of Usearch version 11.0.667

574    (Edgar 2016). This was run with the options `'-strand plus -mode sensitive'` against a

575    database containing all alleles expected to appear in that strain.

576    Identification of truncated sequences, sequences containing Ns, off-targets, and localisation of

577    sequencing errors and SNPs along the amplicons, was performed using a custom Python script,

28

578    incorporating Bowtie2 version 2.3.4.3 (Langmead and Salzberg 2012), SAMtools version 1.12,

579    and BCFtools version 1.12 (Danecek et al. 2021).

580    *Performance of intraspecific metabarcoding in a selection experiment*

581    The quantitative performance of *Sm_C12W1* barcodes locus was validated in a selection

582    experiment using the 58 *S. marinoi* strains. A detailed description of this experimental design

583    and the results will be published elsewhere (Andersson et al, unpublished). Briefly, the single

584    strain cultures (non-axenic) were pooled back to their original populations (Gropviken: GP and

585    Gåsfjärden: VG) at even densities, and maintained in exponential growth phase via a semi-

586    continuous dilution scheme as outlined in Andersson et al. (2020) through dilutions every third

587    day for a total of 42 days. Five 100 mL replicates per population were grown with toxic copper

588    stress of 8.65 μM $CuSO_4$, and five replicates were grown without toxic stress. The dilution

589    bottleneck was ~10,000 chains of cells per replicate. Cultures were harvested at the start of the

590    experiment, after nine, and 42 days of selection. DNA was extracted from experimental samples

591    as for the WGS preparation outlined above, but without the axenic treatment. Barcodes were

592    PCR amplified from each replicate timepoint. For the single time 0 sample, two technical DNA

593    extractions replicates × two PCR replicates were processed and analysed in parallel (i.e., N=4

594    technical replicates). The selection experiment samples were indexed and pooled together with

595    the individual genotypes, on the same MiSeq run using 9/10 of the read capacity (i.e., an

596    expected 16.2 million reads).

597    The individual strains' genotype samples were used to create a database of all strains' alleles,

598    and amplicons from the selection experiment samples were compared against this database

599    (script available at: https://github.com/topel-research-group/Live2Tell). As the experimental

29

600     samples contained mixtures of strains at variable densities, and allele calling was expected to

601     require single SNP resolution across long amplicons (up to 523 bp), we explored avenues to

602     denoise the data. Initial attempts to use the nf-core/ampliseq pipeline version 2.4.0 (Straub et al.

603     2020) failed to detect all expected alleles. Running DADA2 version 1.16.0 (Callahan et al. 2016)

604     – a component of the nf-core/ampliseq pipeline – as a standalone program for identifying

605     amplicon sequence variants (ASVs) proved more successful. Different settings and combinations

606     of steps were iteratively explored (described in detail in Supplemental information). We used

607     four criteria to evaluated the quality of the output including: 1) how many known allele

608     sequences were <u>not</u> identified as ASVs (false negatives); 2) how many unexpected ASVs were

609     outputted (false positives); 3) what fraction of total raw reads assembled into amplicons

610     matching either known alleles or false positives; and 4) if false positive or false negative

611     observations risked affecting the biological interpretation of the artificial evolution experiment.

612     The relative abundances of strains were quantified based on the number of amplicons matching

613     its two alleles. In this step, a set of differential equations were used to parse out the observations

614     of alleles shared between strains. In general, the proportion of shared allele $i_1$ belonging to a

615     specific strain x, y, and z was computed individually for each sample based on strains second

616     unique allele ($x_2$, $y_2$, and $z_2$):

617     Eq. 1: Reads of $i_1$ belonging to strain $x = i_1 \times \dfrac{\text{strain } x \text{ allele } x_2}{(\text{strain } x \text{ allele } x_2 + \text{strain } y \text{ allele } y_2 + \text{strain } z \text{ allele } z_2)}$

618     The two heterozygous allele counts were then summed up to represent the abundance of each

619     strain, with the third allele of triploid strains omitted. Amplicon counts per strain were

620     simultaneously normalized to relative abundance (RA) as:

30

621     Eq. 2: $$\text{RA strain } x = \frac{\text{counts allele } x_1 + counts \text{ allele } x_2}{\sum counts\; allels\; _{a-z}}$$

622     To validate the approach, we assessed how closely the barcode enumeration of strains matched

623     with microscopic cell counts in the time zero samples and how closely the abundance of each

624     strain's two alleles correlated throughout the selection experiment.

## References

Abou Saada, O., A. Tsouris, C. Eberlein, A. Friedrich, and J. Schacherer. 2021. nPhase: an accurate and contiguous phasing method for polyploids. Genome Biol. **22:** 1-27. doi:10.1186/s13059-021-02342-x

Adams, C. I., M. Knapp, N. J. Gemmell, G.-J. Jeunen, M. Bunce, M. D. Lamare, and H. R. Taylor. 2019. Beyond biodiversity: can environmental DNA (eDNA) cut it as a population genetics tool? Genes **10:** 192. doi:10.3390/genes10030192

Altschul, S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. 1990. Basic local alignment search tool. J. Mol. Biol. **215:** 403-410. doi:10.1016/S0022-2836(05)80360-2

Andersson, B., A. Godhe, H. L. Filipsson, K. Rengefors, and O. Berglund. 2020. Differences in metal tolerance among strains, populations, and species of marine diatoms-importance of exponential growth for quantification. Aquat. Toxicol. **226:** 105551. doi:10.1016/j.aquatox.2020.105551.

Angers-Loustau, A., M. Petrillo, V. Paracchini, D. M. Kagkli, P. E. Rischitor, A. Puertas Gallardo, A. Patak and others 2016. Towards plant species identification in complex samples: a bioinformatics pipeline for the identification of novel nuclear barcode candidates. PLoS one **11:** e0147692. doi:10.1371/journal.pone.0147692

Armbrust, E. V., J. A. Berges, C. Bowler, B. R. Green, D. Martinez, N. H. Putnam, S. Zhou and others 2004. The genome of the diatom *Thalassiosira pseudonana*: ecology, evolution, and metabolism. Science **306:** 79-86. doi:10.1126/science.1101156

Baert, J. M., F. De Laender, K. Sabbe, and C. R. Janssen. 2016. Biodiversity increases functional and compositional resistance, but decreases resilience in phytoplankton communities. Ecology **97:** 3433-3440. doi:10.1002/ecy.1601

Blanc-Mathieu, R., M. Krasovec, M. Hebrard, S. Yau, E. Desgranges, J. Martin, W. Schackwitz and others 2017. Population genomics of picophytoplankton unveils novel chromosome hypervariability. Sci. Adv. **3:** e1700239. doi:10.1126/sciadv.1700239

Bowler, C., A. E. Allen, J. H. Badger, J. Grimwood, K. Jabbari, A. Kuo, U. Maheswari and others 2008. The *Phaeodactylum* genome reveals the evolutionary history of diatom genomes. Nature **456:** 239-244. doi:10.1038/nature07410

Bushnell, B., J. Rood, and E. Singer. 2017. BBMerge–accurate paired shotgun read merging via overlap. PloS one **12:** e0185056. doi:10.1371/journal.pone.0185056

Callahan, B. J., P. J. McMurdie, M. J. Rosen, A. W. Han, A. J. A. Johnson, and S. P. Holmes. 2016. DADA2: High-resolution sample inference from Illumina amplicon data. Nat. Methods **13:** 581-583. doi:10.1038/nmeth.3869

Callahan, B. J., J. Wong, C. Heiner, S. Oh, C. M. Theriot, A. S. Gulati, S. K. McGill and others 2019. High-throughput amplicon sequencing of the full-length 16S rRNA gene with single-nucleotide resolution. Nucleic Acids Res. **47:** e103-e103. doi:10.1093/nar/gkz569

Canesi, K. L., and T. A. Rynearson. 2016. Temporal variation of *Skeletonema* community composition from a long-term time series in Narragansett Bay identified using high-throughput DNA sequencing. Mar. Ecol. Prog. Ser. **556:** 1-16. doi:10.3354/meps11843

Collins, S. 2011. Competition limits adaptation and productivity in a photosynthetic alga at elevated CO2. Proc. R. Soc. B. **278:** 247-255. doi:10.1098/rspb.2010.1173

Collins, S., and C. E. Schaum. 2021. Growth strategies of a model picoplankter depend on social milieu and p CO2. Proc. R. Soc. B **288:** 20211154.

32

670    Danecek, P., J. K. Bonfield, J. Liddle, J. Marshall, V. Ohan, M. O. Pollard, A. Whitwham and
671          others 2021. Twelve years of SAMtools and BCFtools. Gigascience **10:** giab008.
672          doi:10.1093/gigascience/giab008

673    Dent, R. M., C. M. Haglund, B. L. Chin, M. C. Kobayashi, and K. K. Niyogi. 2005. Functional
674          genomics of eukaryotic photosynthesis using insertional mutagenesis of *Chlamydomonas*
675          *reinhardtii*. Plant Physiol. **137:** 545-556. doi:10.1104/pp.104.055244

676    Descamps-Julien, B., and A. Gonzalez. 2005. Stable coexistence in a fluctuating environment: an
677          experimental demonstration. Ecology **86:** 2815-2824. doi:10.1890/04-1700

678    Edgar, R. C. 2016. UCHIME2: improved chimera prediction for amplicon sequencing. BioRxiv**:**
679          074252. doi:10.1101/074252

680    Falkowski, P. G., R. T. Barber, and V. V. Smetacek. 1998. Biogeochemical Controls and
681          Feedbacks on Ocean Primary Production. Science **281:** 200-207.
682          doi:10.1126/science.281.5374.200

683    Ferrante, M. I., L. Entrambasaguas, M. Johansson, M. Töpel, A. Kremp, M. Montresor, and A.
684          Godhe. 2019. Exploring Molecular Signs of Sex in the Marine Diatom *Skeletonema*
685          *marinoi*. Genes **10:** 494. doi:10.3390/genes10070494

686    Field, C. B., M. J. Behrenfeld, J. T. Randerson, and P. Falkowski. 1998. Primary production of
687          the biosphere: integrating terrestrial and oceanic components. Science **281:** 237-240.
688          doi:10.1126/science.281.5374.237

689    Flowers, J. M., K. M. Hazzouri, G. M. Pham, U. Rosas, T. Bahmani, B. Khraiwesh, D. R. Nelson
690          and others 2015. Whole-genome resequencing reveals extensive natural variation in the
691          model green alga *Chlamydomonas reinhardtii*. Plant Cell **27:** 2353-2369.
692          doi:10.1105/tpc.15.00492

693    Giaever, G., and C. Nislow. 2014. The yeast deletion collection: a decade of functional
694          genomics. Genetics **197:** 451-465. doi:10.1534/genetics.114.161620

695    Godhe, A., M. R. McQuoid, I. Karunasagar, I. Karunasagar, and A. S. Rehnstam☐Holm. 2006.
696          Comparison of three common molecular tools for distinguishing among geographically
697          separated clones of the diatom *Skeletonema marinoi* Sarno et Zingone
698          (Bacillariophyceae) 1. J. Phycol. **42:** 280-291. doi:10.1111/j.1529-8817.2006.00197.x

699    Godhe, A., S. K. Otta, A. S. Rehnstam-Holm, I. Karunasagar, and I. Karunasagar. 2001.
700          Polymerase chain reaction in detection of *Gymnodinium mikimotoi* and *Alexandrium*
701          *minutum* in field samples from southwest India. Mar. Biotechnol. **3:** 152-162.
702          doi:10.1007/s101260000052

703    Godhe, A., and T. Rynearson. 2017. The role of intraspecific variation in the ecological and
704          evolutionary success of diatoms in changing environments. Phil. Trans. R. Soc. B. **372:**
705          20160399. doi:10.1098/rstb.2016.0399

706    Guillard, R. R. L. 1975. Culture of phytoplankton for feeding marine invertebrates., p. 29–60. *In*
707          W. L. Smith and M. H. Chanley [eds.], *Culture of marine invertebrate animals.* Springer.

708    Guo, L., Z. Sui, S. Zhang, Y. Ren, and Y. Liu. 2015. Comparison of potential diatom
709          'barcode'genes (the 18S rRNA gene and ITS, COI, rbcL) and their effectiveness in
710          discriminating and determining species taxonomy in the Bacillariophyta. Int. J. Syst.
711          Evol. Microbiol. **65:** 1369-1380. doi:10.1099/ijs.0.000076

712    Han, T. X., X.-Y. Xu, M.-J. Zhang, X. Peng, and L.-L. Du. 2010. Global fitness profiling of
713          fission yeast deletion strains by barcode sequencing. Genome Biol. **11:** 1-13.
714          doi:10.1186/gb-2010-11-6-r60

715   Hebert, P. D. N., A. Cywinska, S. L. Ball, and J. R. deWaard. 2003. Biological identifications
716         through DNA barcodes. Proc. R. Soc. B. **270:** 313-321. doi:10.1098/rspb.2002.2218
717   Härnström, K., M. Ellegaard, T. J. Andersen, and A. Godhe. 2011. Hundred years of genetic
718         structure in a sediment revived diatom population. Proc. Natl. Acad. Sci. USA **108:**
719         4252–4257. doi:10.1073/pnas.1013528108
720   Keeling, P. J., F. Burki, H. M. Wilcox, B. Allam, E. E. Allen, L. A. Amaral-Zettler, E. V.
721         Armbrust and others 2014. The Marine Microbial Eukaryote Transcriptome Sequencing
722         Project (MMETSP): illuminating the functional diversity of eukaryotic life in the oceans
723         through transcriptome sequencing. PLoS Biol. **12:** e1001889.
724         doi:10.1371/journal.pbio.1001889
725   Kim, D.-U., J. Hayles, D. Kim, V. Wood, H.-O. Park, M. Won, H.-S. Yoo and others 2010.
726         Analysis of a genome-wide set of gene deletions in the fission yeast
727         Schizosaccharomyces pombe. Nat. Biotechnol. **28:** 617-623. doi:10.1038/nbt.1628
728   Kress, W. J., and D. L. Erickson. 2008. DNA barcodes: genes, genomics, and bioinformatics.
729         Proc. Natl. Acad. Sci. USA **105:** 2761-2762. doi:10.1073/pnas.0800476105
730   Lakeman, M. B., P. von Dassow, and R. A. Cattolico. 2009. The strain concept in phytoplankton
731         ecology. Harmful Algae **8:** 746-758. doi:10.1016/j.hal.2008.11.011
732   Langmead, B., and S. L. Salzberg. 2012. Fast gapped-read alignment with Bowtie 2. Nat.
733         Methods **9:** 357-359. doi:10.1038/nmeth.1923
734   Leinonen, R., H. Sugawara, M. Shumway, and I. N. S. D. Collaboration. 2010. The sequence
735         read archive. Nucleic Acids Res. **39:** D19-D21. doi:10.1093/nar/gkq1019
736   Li, X., W. Patena, F. Fauser, R. E. Jinkerson, S. Saroussi, M. T. Meyer, N. Ivanova and others
737         2019. A genome-wide algal mutant library and functional screen identifies genes required
738         for eukaryotic photosynthesis. Nat. Genet. **51:** 627-635. doi:10.1038/s41588-019-0370-6
739   Lohbeck, K. T., U. Riebesell, and T. B. Reusch. 2012. Adaptive evolution of a key
740         phytoplankton species to ocean acidification. Nat. Geosci. **5:** 346. doi:10.1038/ngeo1441
741   Lommer, M., M. Specht, A. S. Roy, L. Kraemer, R. Andreson, M. A. Gutowska, J. Wolf and
742         others 2012. Genome and low-iron response of an oceanic diatom adapted to chronic iron
743         limitation. Genome Biol. **13:** R66. doi:10.1186/gb-2012-13-7-r66
744   Marchler-Bauer, A., and S. H. Bryant. 2004. CD-Search: protein domain annotations on the fly.
745         Nucleic Acids Res. **32:** W327-W331. doi:10.1093/nar/gkh454
746   Martin, M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads.
747         EMBnet. journal **17:** 10-12. doi:10.14806/ej.17.1.200
748   Merchant, S. S., S. E. Prochnik, O. Vallon, E. H. Harris, S. J. Karpowicz, G. B. Witman, A.
749         Terry and others 2007. The *Chlamydomonas* genome reveals the evolution of key animal
750         and plant functions. Science **318:** 245-250. doi:10.1126/science.1143609
751   Mock, T., R. P. Otillar, J. Strauss, M. McMullan, P. Paajanen, J. Schmutz, A. Salamov and
752         others 2017. Evolutionary genomics of the cold-adapted diatom *Fragilariopsis cylindrus*.
753         Nature **541:** 536-540. doi:10.1038/nature20803
754   Ness, R. W., S. A. Kraemer, N. Colegrave, and P. D. Keightley. 2016. Direct estimate of the
755         spontaneous mutation rate uncovers the effects of drift and recombination in the
756         *Chlamydomonas reinhardtii* plastid genome. Mol. Biol. Evol. **33:** 800-808.
757         doi:10.1093/molbev/msv272
758   Osuna-Cruz, C. M., G. Bilcke, E. Vancaester, S. De Decker, A. M. Bones, P. Winge, N. Poulsen
759         and others 2020. The *Seminavis robusta* genome provides insights into the evolutionary

760         adaptations of benthic diatoms. Nat. Commun. **11:** 3320. doi:10.1038/s41467-020-17191-
761         8
762 Paracchini, V., M. Petrillo, A. Lievens, A. P. Gallardo, J. T. Martinsohn, J. Hofherr, A. Maquet
763         and others 2017. Novel nuclear barcode regions for the identification of flatfish species.
764         Food Control **79:** 297-308. doi:10.1016/j.foodcont.2017.04.009
765 Pujari, L., C. Wu, J. Kan, N. Li, X. Wang, G. Zhang, X. Shang and others 2019. Diversity and
766         Spatial Distribution of Chromophytic Phytoplankton in the Bay of Bengal Revealed by
767         RuBisCO Genes (rbcL). Front. Microbiol. **10:** 1501. doi:10.3389/fmicb.2019.01501
768 Quinlan, A. R., and I. M. Hall. 2010. BEDTools: a flexible suite of utilities for comparing
769         genomic features. Bioinformatics **26:** 841-842. doi:10.1093/bioinformatics/btq033
770 Rastogi, A., F. R. J. Vieira, A.-F. Deton-Cabanillas, A. Veluchamy, C. Cantrel, G. Wang, P.
771         Vanormelingen and others 2020. A genomics approach reveals the global genetic
772         polymorphism, structure, and functional diversity of ten accessions of the marine model
773         diatom *Phaeodactylum tricornutum*. ISME J **14:** 347-363. doi:10.1038/s41396-019-0528-
774         3
775 Read, B. A., J. Kegel, M. J. Klute, A. Kuo, S. C. Lefebvre, F. Maumus, C. Mayer and others
776         2013. Pan genome of the phytoplankton *Emiliania* underpins its global distribution.
777         Nature **499:** 209-213. doi:10.1038/nature12221
778 Rengefors, K., R. Gollnisch, I. Sassenhagen, K. Harnstrom Aloisi, M. Svensson, K. Lebret, D.
779         Certnerova and others 2021. Genome-wide single nucleotide polymorphism markers
780         reveal population structure and dispersal direction of an expanding nuisance algal bloom
781         species. Mol. Ecol. **30:** 912-925. doi:10.1111/mec.15787
782 Rengefors, K., A. Kremp, T. B. Reusch, and A. M. Wood. 2017. Genetic diversity and evolution
783         in eukaryotic phytoplankton: revelations from population genetic studies. J. Plankton
784         Res. **39:** 165-179. doi:10.1093/plankt/fbw098
785 Rhoads, A., and K. F. Au. 2015. PacBio Sequencing and Its Applications. Genom. Proteom.
786         Bioinform. **13:** 278-289. doi:10.1016/j.gpb.2015.08.002
787 Robinson, D. G., W. Chen, J. D. Storey, and D. Gresham. 2014. Design and analysis of Bar-seq
788         experiments. G3-Genes Genom. Genet. **4:** 11-18. doi:10.1534/g3.113.008565
789 Robinson, J. T., H. Thorvaldsdóttir, W. Winckler, M. Guttman, E. S. Lander, G. Getz, and J. P.
790         Mesirov. 2011. Integrative genomics viewer. Nat. Biotechnol. **29:** 24-26.
791         doi:10.1038/nbt.1754
792 Rognes, T., T. Flouri, B. Nichols, C. Quince, and F. Mahé. 2016. VSEARCH: a versatile open
793         source tool for metagenomics. PeerJ **4:** e2584. doi:10.7717/peerj.2584
794 Sarno, D., W. H. Kooistra, L. K. Medlin, I. Percopo, and A. Zingone. 2005. Diversity in the
795         genus *Skeletonema* (Bacillariophyceae). II. An assessment of the taxonomy of *S.*
796         *costatum*□like species with the description of four new species 1. J. Phycol. **41:** 151-176.
797         doi:10.1111/j.1529-8817.2005.04067.x
798 Sassenhagen, I., D. Erdner, B. Lougheed, M. Richlen, and C. Sjöqvist. 2021. Estimating
799         proportion of clones and genotype richness in aquatic microalgae. Authorea Preprints.
800         doi:10.22541/au.161876383.32271114/v1
801 Schaum, C.-E., S. Barton, E. Bestion, A. Buckling, B. Garcia-Carreras, P. Lopez, C. Lowe and
802         others 2017. Adaptation of phytoplankton to a decade of experimental warming linked to
803         increased photosynthesis. Nat. Ecol. Evol. **1:** 1-7. doi:10.1038/s41559-017-0094

804  Schaum, C.-E., A. Buckling, N. Smirnoff, D. Studholme, and G. Yvon-Durocher. 2018.
805        Environmental fluctuations accelerate molecular evolution of thermal tolerance in a
806        marine diatom. Nat. Commun. **9:** 1719. doi:10.1038/s41467-018-03906-5
807  Scheinin, M., U. Riebesell, T. A. Rynearson, K. T. Lohbeck, and S. Collins. 2015. Experimental
808        evolution gone wild. J. R. Soc. Interface **12:** 20150056. doi:10.1098/rsif.2015.0056
809  Schloss, P. D., S. L. Westcott, T. Ryabin, J. R. Hall, M. Hartmann, E. B. Hollister, R. A.
810        Lesniewski and others 2009. Introducing mothur: open-source, platform-independent,
811        community-supported software for describing and comparing microbial communities.
812        Appl. Environ. Microbiol. **75:** 7537-7541. doi:10.1128/AEM.01541-09
813  Sefbom, J., I. Sassenhagen, K. Rengefors, and A. Godhe. 2015. Priority effects in a planktonic
814        bloom-forming marine diatom. Biol. Lett. **11:** 20150184. doi:10.1098/rsbl.2015.0184
815  Sjöqvist, C., A. Kremp, E. Lindehoff, U. Båmstedt, J. Egardt, S. Gross, M. Jönsson and others
816        2014. Effects of grazer presence on genetic structure of a phenotypically diverse diatom
817        population. Microb. Ecol. **67:** 83-95. doi:10.1007/s00248-013-0327-8
818  Smith, A. M., L. E. Heisler, J. Mellor, F. Kaper, M. J. Thompson, M. Chee, F. P. Roth and others
819        2009. Quantitative phenotyping via deep barcode sequencing. Genome Res. **19:** 1836-
820        1842. doi:10.1101/gr.093955.109
821  Smyth, R. P., T. E. Schlub, A. Grimm, V. Venturi, A. Chopra, S. Mallal, M. P. Davenport and
822        others 2010. Reducing chimera formation during PCR amplification to ensure accurate
823        genotyping. Gene **469:** 45-51. doi:10.1016/j.gene.2010.08.009
824  Stoeck, T., D. Bass, M. Nebel, R. Christen, M. D. Jones, H. W. Breiner, and T. A. Richards.
825        2010. Multiple marker parallel tag environmental DNA sequencing reveals a highly
826        complex eukaryotic community in marine anoxic water. Mol. Ecol. **19:** 21-31.
827        doi:10.1111/j.1365-294X.2009.04480.x
828  Straub, D., N. Blackwell, A. Langarica-Fuentes, A. Peltzer, S. Nahnsen, and S. Kleindienst.
829        2020. Interpretations of Environmental Microbial Community Studies Are Biased by the
830        Selected 16S rRNA (Gene) Amplicon Sequencing Pipeline. Front. Microbiol. **11**.
831        doi:10.3389/fmicb.2020.550420
832  Sundberg, C., W. A. Al-Soud, M. Larsson, E. Alm, S. S. Yekta, B. H. Svensson, S. J. Sørensen
833        and others 2013. 454 pyrosequencing analyses of bacterial and archaeal richness in 21
834        full-scale biogas digesters. FEMS Microbiol. Ecol. **85:** 612-626. doi:10.1111/1574-
835        6941.12148
836  Söderhielm, J., and K. Sundblad. 1996. The Solstad Cu☐Co☐Au mineralization and its relation
837        to post☐Svecofennian regional shear zones in southeastern Sweden. GFF **118:** 47-47.
838        doi:10.1080/11035899609546323
839  Taberlet, P., E. Coissac, F. Pompanon, C. Brochmann, and E. Willerslev. 2012. Towards
840        next☐generation biodiversity assessment using DNA metabarcoding. Mol. Ecol. **21:**
841        2045-2050. doi:10.1111/j.1365-294X.2012.05470.x
842  Tatters, A. O., M. Y. Roleda, A. Schnetzer, F. Fu, C. L. Hurd, P. W. Boyd, D. A. Caron and
843        others 2013. Short-and long-term conditioning of a temperate marine diatom community
844        to acidification and warming. Phil. Trans. R. Soc. B. **368:** 20120437.
845        doi:10.1098/rstb.2012.0437
846  Van der Auwera, G. A., and B. D. O'Connor. 2020. Genomics in the cloud: using Docker,
847        GATK, and WDL in Terra, 1st ed. Sebastopol, CA: O'Reilly Media Inc.

848 von Dassow, P., U. John, H. Ogata, I. Probert, M. Bendif el, J. U. Kegel, S. Audic and others
849       2015. Life-cycle modification in open oceans accounts for genome variability in a
850       cosmopolitan phytoplankton. ISME J **9:** 1365-1377. doi:10.1038/ismej.2014.221
851 Williams, R., S. G. Peisajovich, O. J. Miller, S. Magdassi, D. S. Tawfik, and A. D. Griffiths.
852       2006. Amplification of complex gene libraries by emulsion PCR. Nat. Methods **3:** 545-
853       550. doi:10.1038/nmeth896
854 Wolf, K. K., E. Romanelli, B. Rost, U. John, S. Collins, H. Weigand, and C. J. Hoppe. 2019.
855       Company matters: The presence of other genotypes alters traits and intraspecific selection
856       in an Arctic diatom under climate change. Global Change Biol. **25:** 2869-2884.
857       doi:10.1111/gcb.14675
858 Worden, A. Z., J.-H. Lee, T. Mock, P. Rouzé, M. P. Simmons, A. L. Aerts, A. E. Allen and
859       others 2009. Green evolution and dynamic adaptations revealed by genomes of the
860       marine picoeukaryotes *Micromonas*. Science **324:** 268-272. doi:10.1126/science.1167222
861 Yamada, M., M. Otsubo, Y. Tsutsumi, C. Mizota, Y. Nakamura, K. Takahashi, and M. Iwataki.
862       2017. Utility of mitochondrial□encoded cytochrome c oxidase I gene for phylogenetic
863       analysis and species identification of the planktonic diatom genus *Skeletonema*. Phycol.
864       Res. **65:** 217-225. doi:10.1111/pre.12179
865
866

867  Table 1. Statistical parameters from Bamboozles genomic scan for barcode windows. Shown are

868  average (standard deviation) for the variable barcoding loci identified by Bamboozle (selected

869  windows), with regard to the number of predicted SNPs, indels, and unique alleles per locus, in

870  contrast to average statistics for the genome as a whole, obtained using a sliding-window

871  approach (other windows). Where two sets of values are given, the first gives statistics are for

872  those windows where only the coverage filter is applied, and the second (in square brackets)

873  gives statistics for those windows were both the coverage filter and the conserved end region

874  filter are applied.

875

| | *S. marinoi* | *C. reinhardtii* Mating type + | *C. reinhardtii* Mating type - |
|---|---|---|---|
| Reference genome size (nuclear, mitochondrion, and plastid) | 54,794,446 bp | 107,613,365 bp | |
| Number of strains included | 54 | 8 | 9 |
| Length of genome conforming to coverage filtering | 4,740,972 | 630,835 | 464,867 |
| SNP density of selected windows | 45.0 (9.13)[†] | 17.4 (3.25) | 22.86 (4.74) |
| SNP density of other windows | 26.26 (14.53) [12.36 (10.83)] | 5.30 (6.02) [2.17 (3.42)] | 4.68 (5.87) [1.92 (2.95)] |
| Indel density of selected windows | 2.75 (4.27) | 1.47 (0.83) | 1.71 (1.98) |
| Indel density of other windows | 3.27 (4.55) [1.38 (2.76)] | 0.53 (1.03) [0.27 (0.73)] | 0.46 (0.98) [0.23 (0.67)] |
| Number of unique alleles in selected windows | 100 (4.27) | 8 (0) [††] | 9 (0) [††] |
| Number of unique alleles in other windows | 30.56 (20.41) [14.30 (13.89)] | 2.87 (1.58) [2.08 (1.23)] | 2.99 (1.85) [2.11 (1.31)] |

876  [†]Note that in the case of *S. marinoi*, these statistics exclude *Sm_C2W24* (identified using an older version
877  of the pipeline) but include *Sm_C12W3* (identified in the more recent version).
878  [††]When comparing all *C. reinhardtii* selected windows among all 17 strains (i.e., not separated by mating
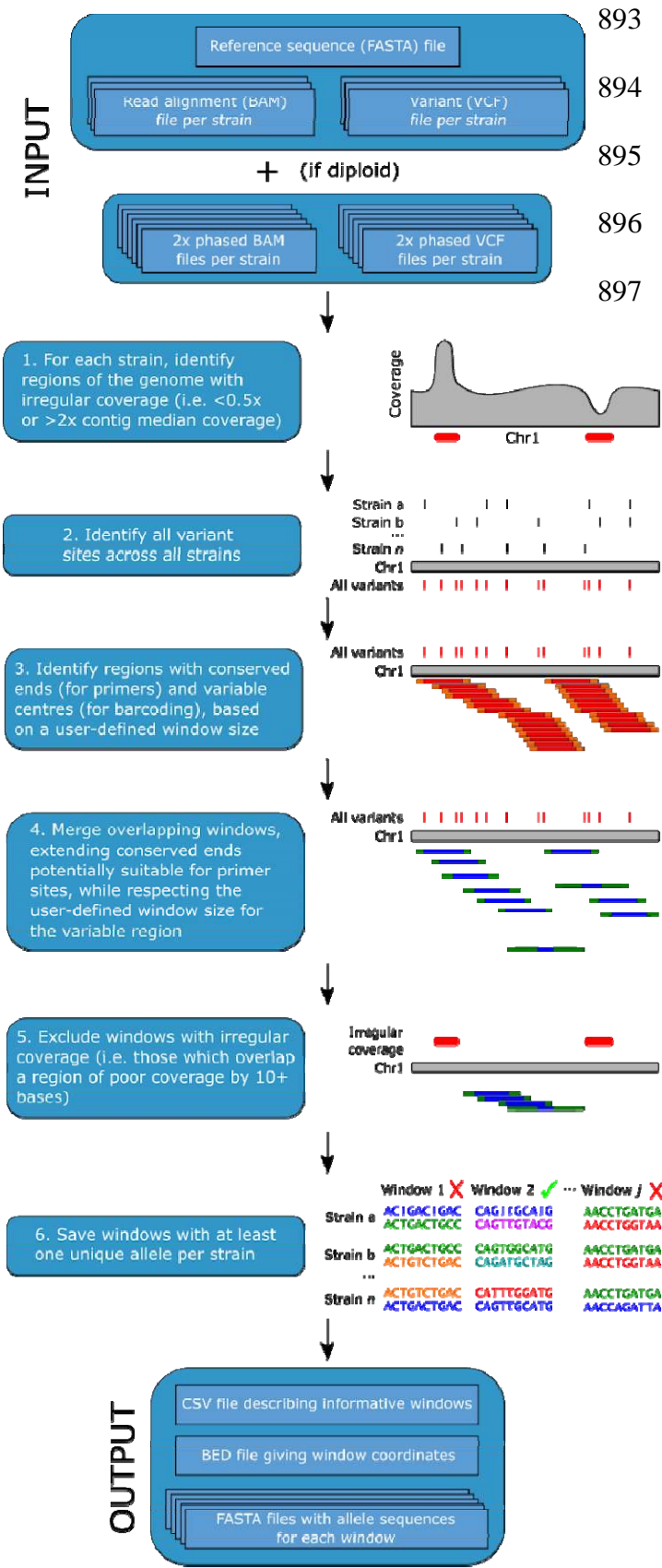879  type), there are 12.55 (1.79) alleles per locus.
880

38

881 Table 2: Accuracy of Bamboozle's allele sequence predictions for selected barcodes in *S. marinoi* and *C. reinhardtii*. 'Predictions' corresponds to Bamboozle's output sequences based on WGS data, and 'observations' are confirmations based on targeted PCR-based sequencing. *C. reinhardtii* loci with poor amplification and lack of Sanger sequencing data are shown as n/a. *Sm_C2W24* results marked with * indicate results from an earlier iteration of the pipeline, which didn't take coverage or allele phasing into account, and was variant-called using BCFtools (cf. GATK in the most recent iteration).

888

| *Skeletonema marinoi*[†] | *Sm_C2W24* | *Sm_C12W1* | *Sm_C12W2* | *Sm_C16W4* |
|---|---|---|---|---|
| Gene model | Sm_t00004275-RA | Sm_t00009768-RA | Sm_t00009769-RA | Sm_t00008465-RA |
| Barcode length in reference (2 alleles in ref. genome) | 471 (471/471) | 523 (523/523) | 494 (494/494) | 504 (501/504) |
| Predicted SNP positions | 95* | 37 | 44 | 58 |
| Predicted indel positions | 5* | 2 | 0 | 9 |
| Predicted unique alleles | 56* | 98 | 99 | 106 |
| Observed SNP positions | 19 | 38 | 44 | 54 |
| Observed indel positions | 198 | 0 | 0 | 18 |
| Mean WGS sequencing depth (range, N=55 ) | 49 (13-101) | 60 (17-140) | 59 (17-131) | 63 (19-131) |
| Mean amplicon sequencing depth (range, N=55) | 5,391 (5-35,434) | 4,156 (15-10,934) | 1,693 (5-6,960) | 3,591 (6-19,143) |
| Accuracy of prediction across all strains and bases | 82.65% | 99.91% | 99.93% | 99.88% |

889

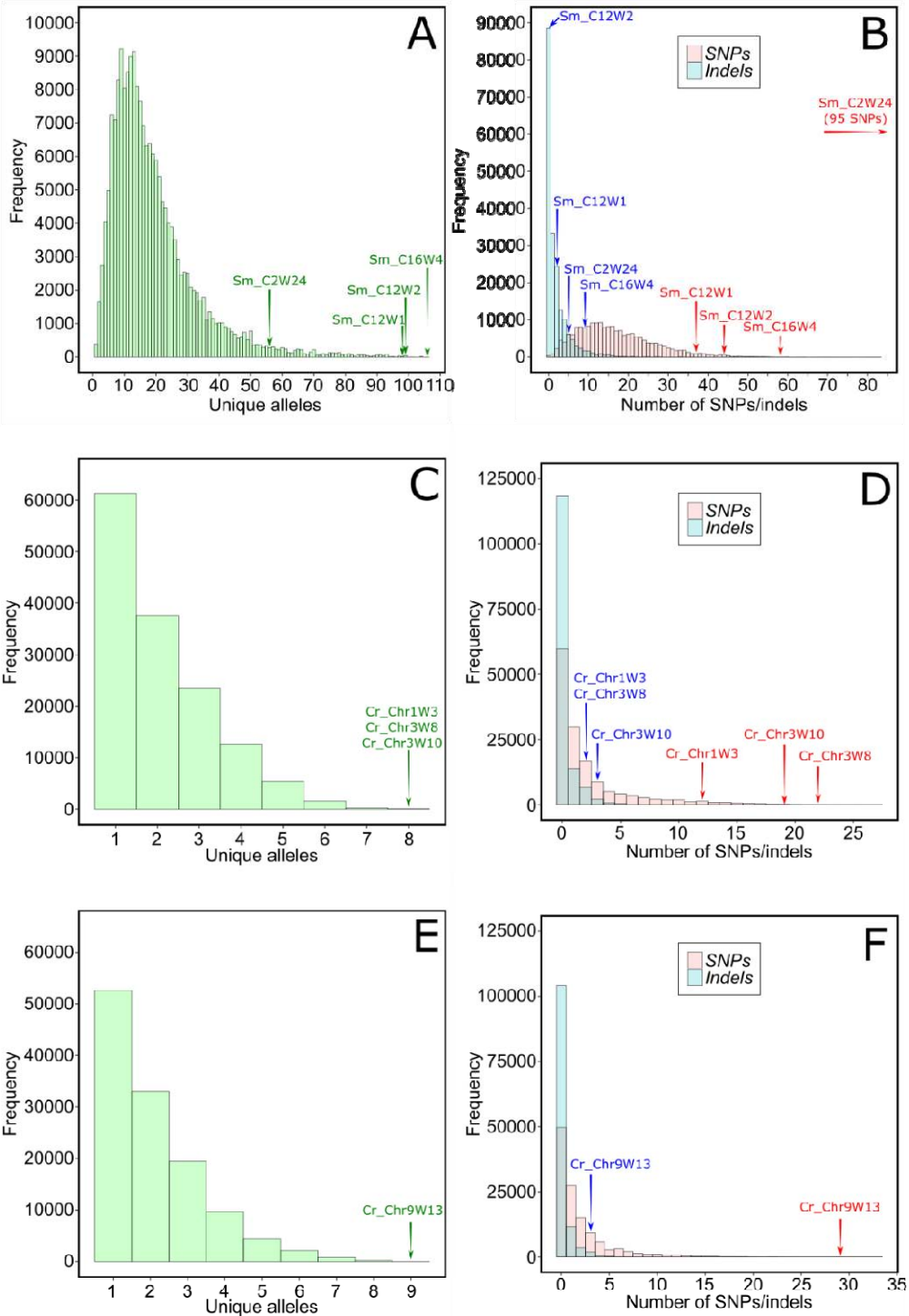| *Chlamydomonas reinhardtii*[††] | mt+ | | mt- | |
|---|---|---|---|---|
| | Cr_Chr1W3 | Cr_Chr3W8 | Cr_Chr3W10 | Cr_Chr9W12 |
| Gene model | CHLRE_01g040900v5 | CHLRE_03g181100v5 | CHLRE_03g207250v5 | CHLRE_09g389134v5 |
| Barcode length | 298 | 297 | 299 | 290 |
| Predicted SNP positions | 21 | 23 | 28 | 32 |
| Predicted indel positions | 4 | 2 | 8 | 3 |
| Observed SNP positions | n/a | n/a | 27 | 32 |
| Observed indel positions | n/a | n/a | 8 | 3 |
| Mean WGS coverage (range N=19) | 29 (11-48) | 37 (17-58) | 32 (16-52) | 33 (14-53) |
| Accuracy of prediction across all strains and bases | n/a | n/a | 99.98 | 99.88 |

890 [†]One strain (GP2-4_32) was excluded from this analysis due to a failed sequencing.
891 [††]One *strain* (CC-2938) was excluded from this analysis due to failed DNA extraction.

892

Figure 1: Schematic summarising the steps of the Bamboozle pipeline, as well as input and output files. Each step is described in more detail in the New Approaches section.
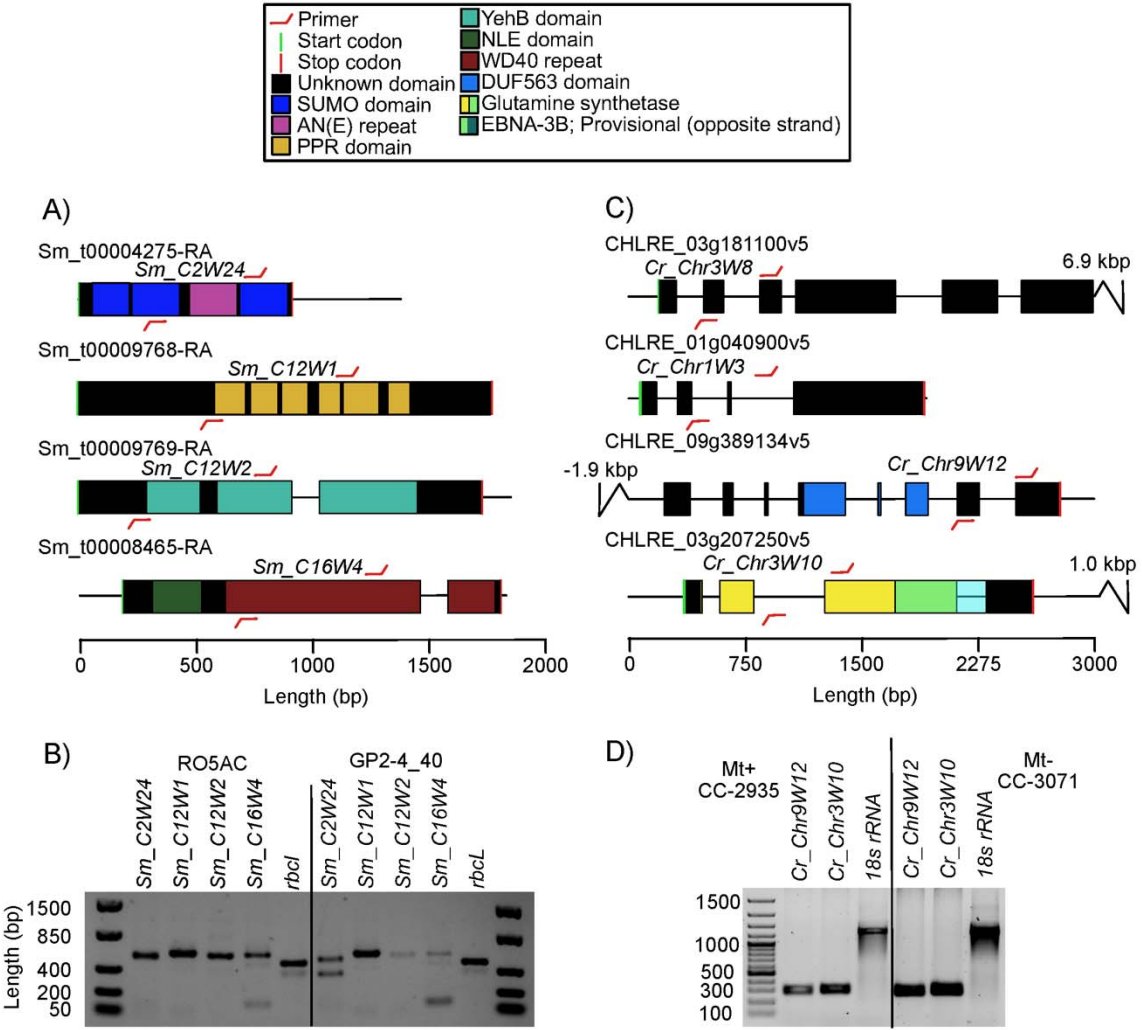
40

898    Figure 2: Histograms showing the predicted number of unique alleles (panels A, C, and E), and

899    predicted frequency of SNPs and indels (panels B, D, and F), at each window of 500 bp (*S.*

900    *marinoi*) or 300 bp (*C. reinhardtii*). Shown are data from all genomic windows that passed

901    Bamboozle's filters for coverage depth and flanking region conservation, as described in the New

902    Approaches section. In the case of *S. marinoi*, only nuclear contigs are included. Panels A and B

903    depict figures from *S. marinoi*; panels C and D depict figures from *C. reinhardtii* mating type +;

904    panels E and F depict figures from *C. reinhardtii* mating type -. Arrows indicate the relative

905    locations of the selected barcoding loci in terms of predicted SNP (red) and indel (blue)

906    frequency, and unique alleles (green). Note that *Sm_C2W24* was predicted with an earlier

907    version of Bamboozle, resulting in deviant values compared to other loci.

908
909

42

910  Figure 3: Genomic location of the four barcoding loci identified in *S. marinoi* and selected four

911  loci in *C. reinhardtii,* corresponding to those that primers were designed and evaluated for (for

912  complete list of all loci, and brief gene model annotations, see Table S2). A and C) Gene models

913  with functional domains, primer sites, and variable regions in *S. marinoi* and *C. reinhardtii,*

914  respectively. B and D) Gel electrophoresis of PCR products in two strains of *S. marinoi* and one

915  strain each of mating type + and - in *C. reinhardtii,* of loci with good amplification results.

916  Abbreviations – SUMO: Small Ubiquitin-like Modifier; AN(E): the three amino acids Alanine,

917  Asparagine, and Glutamic acid [occasionally replaced by Glycine]; PPR: pentatricopeptide

918  [RNA binding]; YehB: Uncharacterised membrane protein [partial, UPF0754 family]; NLE:

919  NUC135 domain located in N terminal of WD40 repeats; WD40: domain involved in various

920  cellular function including signal transduction, pre-mRNA processing and cytoskeleton

921  assembly; DUF563: Domain of Unknown Function; EBNA-3B: PHA03378 superfamily domain

922  (located on the reverse strand). Glutamine synthetase is represented by its PLN02284

923  superfamily domain and overlap with EBNA-3B is shown in green. Note that the primer symbols

924  are not drawn to scale but the tip of each symbol indicates the exact location where they end and

925  the variable region begins.

926  Department of Biology, Lund University, Lund, Sweden

927
928

Figure 4: Metabarcoded relative abundance of strains alleles using the barcode loci *Sm_C12W1*.

Shown are two mixed populations of *S. marinoi*. A) a mixture of 28 strains from Gropviken

(GP), and B) 30 strains from Gåsfjärden (VG). Legend numbers indicate whether this is the first,

second, or third allele in each strain, with colours/asterisks indicating if alleles are homologous

in two (*) or three (**) other strains. For asterisk-marked alleles, the abundance has been

partitioned between strains using differential equations. The horizontal line indicates the

expected fraction of amplicons (allele 1+2) per strain, assuming even cell densities, Error bars

show standard deviation, with N=3 technical PCR replicates for panel A, and N=4 for B. Data

has been denoised using the settings for 'Relaxed DADA2 on amplicons' as outlined in the
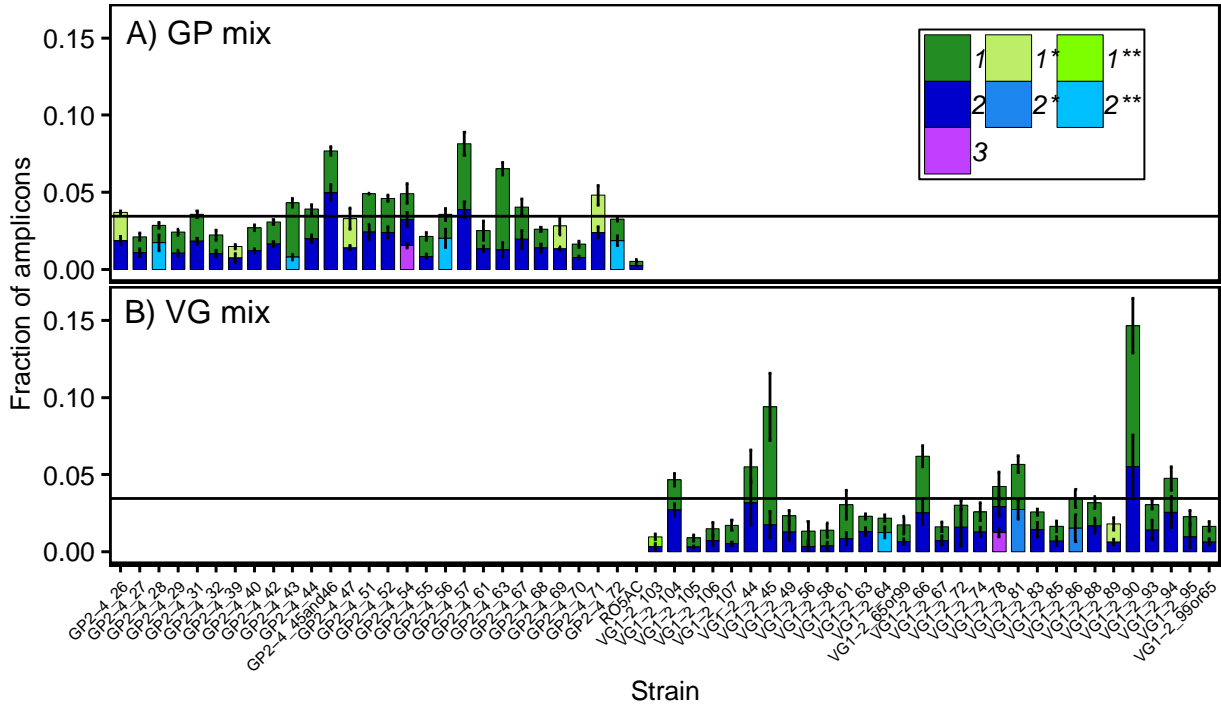
Supplementary Information.

941    Figure 5: Correlations between heterozygous allele of 58 strains of *S. marinoi* across the 42-day

942    selection experiment. For visibility, data has been transformed (x +1) with removal of 0+0 cases.

943    Shared homologous alleles are indicated by alphabetical indexes in the corner of each allele (i.e.,

944    lower right equals allele one, upper left allele two), and their abundances have been allocated

945    between strains using differential equations. All such equations assume a 1:1 ratio between

946    alleles, except allele two in GP2-4_43, where a ratio of 0.15 is used based on the genotype

947    sample. The black line across panels corresponds to the expected 1:1 ratio between heterozygous

948    alleles. The blue line is a second order polynomial function fitted to the data, with shaded area

949    corresponding to 95% confidence interval. Data has been denoised using the settings for

950    'Relaxed DADA2 on amplicons' as outlined in the Supplemental Information. Asterisks

951    highlight false-positive observations that the denoising pipeline failed to remove. Fig. S4 shows

952    the same figure without denoising (i.e. 'exact matches' settings in Table S6).

953