**Mitotic recombination between homologous chromosomes drives genomic diversity in**

**diatoms**

Petra Bulánková[1,2,10]*, Mirna Sekulić[1,2,3,12], Denis Jallet[4,12], Charlotte Nef[5,12], Cock van

Oosterhout[6,] Tom O. Delmont[7], Ilse Vercauteren[1,2], Cristina Maria Osuna-Cruz[1,2,8], Emmelien

Vancaester[1,2,8,9], Thomas Mock[6], Koen Sabbe[3], Fayza Daboussi[4], Chris Bowler[5], Wim

Vyverman[3], Klaas Vandepoele[1,2,8] and Lieven De Veylder[1,2,10,11]*


**Affiliations:**

[1]    VIB Center for Plant Systems Biology, Technologiepark 71, 9052, Ghent, Belgium.

Petra Bulánková, Mirna Sekulić, Ilse Vercauteren, Cristina Maria Osuna-Cruz, Emmelien

Vancaester, Klaas Vandepoele and Lieven De Veylder


[2]    Department of Plant Biotechnology and Bioinformatics, Ghent University,

       Technologiepark 71, 9052, Ghent, Belgium.

Petra Bulánková, Mirna Sekulić, Ilse Vercauteren, Cristina Maria Osuna-Cruz, Emmelien

Vancaester, Klaas Vandepoele and Lieven De Veylder


[3]    Protistology and Aquatic Ecology, Department of Biology, Ghent University, 9000,

       Ghent, Belgium.

Mirna Sekulić, Koen Sabbe and Wim Vyverman


[4]    TBI , Université de Toulouse, CNRS, INRAE, INSA , 135 avenue de Rangueil F-31077,

       Toulouse , France.

Denis Jallet, Fayza Daboussi

25

26    [5]    Institut de Biologie de l'École Normale Supérieure (IBENS), École Normale Supérieure,

27          CNRS, INSERM, PSL Université Paris, 75005 Paris, France.

28    Charlotte Nef, Chris Bowler

29

30

31    [6]    School of Environmental Sciences, University of East Anglia, Norwich Research Park,

32          Norwich, NR4 7TJ, UK.

33    Cock Van Oosterhout, Thomas Mock

34

35    [7]    Génomique Métabolique, Genoscope, Institut François Jacob, CEA, CNRS, Univ Evry,

36          Université Paris-Saclay, 91000 Evry, France.

37    Tom O. Delmont

38

39    [8]    Bioinformatics Institute Ghent, Ghent University, Technologiepark 71, 9052, Ghent,

40          Belgium.

41    Cristina Maria Osuna-Cruz, Emmelien Vancaester and Klaas Vandepoele

42

43    [9]    Present address: Tree of Life, Wellcome Sanger Institute, Cambridge, CB10 1SA, UK

44

45    [10]   Senior author

46

47    [11]   Lead contact

48

49      [12]   These authors contributed equally to this work

50

51      *Correspondence to: lieven.deveylder@psb.vib-ugent.be or petra.bulankova@gmail.com

52

53 **SUMMARY**

54 Diatoms, an evolutionarily successful group of microalgae, display high levels of intraspecific

55 genetic variability in natural populations. However, the contribution of various mechanisms

56 generating such diversity is unknown. Here we estimated the genetic micro-diversity within a

57 natural diatom population and mapped the genomic changes arising within clonally

58 propagated diatom cell cultures. Through quantification of haplotype diversity by next-

59 generation sequencing and amplicon re-sequencing of selected loci, we documented a rapid

60 accumulation of multiple haplotypes accompanied by the appearance of novel protein

61 variants in cell cultures initiated from a single founder cell. Comparison of the genomic

62 changes between mother and daughter cells revealed copy number variation and copy-neutral

63 loss of heterozygosity leading to the fixation of alleles within individual daughter cells. The

64 loss of heterozygosity can be accomplished by recombination between homologous

65 chromosomes. To test this hypothesis, we established an endogenous read-out system and

66 estimated that the frequency of interhomolog mitotic recombination to be under standard

67 growth conditions 4.2 events per 100 cell divisions. This frequency is increased under

68 environmental stress conditions, including treatment with hydrogen peroxide and cadmium.

69 These data demonstrate that copy number variation and mitotic recombination between

70 homologous chromosomes underlie clonal variability in diatom populations. We discuss the

71 potential adaptive evolutionary benefits of the plastic response in the interhomolog mitotic

72    recombination rate, and we propose that this may have contributed to the ecological success

73    of diatoms.

74

**INTRODUCTION**

Diatoms, with as many as 100 000 estimated species, colonize a wide range of marine, freshwater and terrestrial environments[1]. Given their often vast census population size, diatoms possess high intraspecific genetic variation. Natural diatom population samples comprise between 87 to 100% clonal diversity as measured by microsatellite markers and a gene diversity ranging from 39 to 88%, suggesting that clonal lineages are significantly diverged [2]. However, many species reproduce asexually for long periods, and hence, some of this gene diversity is present as clonal diversity[2]. For example, the model diatoms *Phaeodactylum tricornutum* and *Thalassiosira pseudonana* have never been observed to produce F1 progeny[3,4] and in many other diatom species, sex is restricted by cell size[5]. Due to the mechanics of diatom cell wall formation from inorganic silica, the cell size of mitotically dividing cells in a population decreases over time with only cells under a species-specific size threshold being sexually competent. The period necessary to reach this threshold can last months to years[3,6-8]. So how do diatoms generate novel genetic variation that is required for adaptive evolution?

Genetic variation is generated by three evolutionary forces: mutation, recombination and gene flow. Ultimately, all novel genetic variation stems from mutations, whereas the two other forces merely shuffle this variation between haplotypes or genotypes (i.e. recombination), or between populations (gene flow)[9,10]. Mutation rates in diatoms appear to be at a similarly low level as in green algae[11] and coccoliths[12], but given their often large population size, considerable variation can be generated by the input of new mutations[13,14]. However, not all variation generated in a clonal lineage is readily available to selection due to recessivity and clonal interference, i.e. the competition between different beneficial mutations that occur in individuals within the same clonal lineage[15,16]. Potentially beneficial

99  variation can become "locked inside" in a poorly adapted genotype, and sexual reproduction

100  can relax this evolutionary constrained. The random assortment of alleles from both parents

101  during meiotic recombination[17,18] generates novel genotypes, which form the substrate for

102  natural selection.

103  In the absence of sex, mitotic interhomolog recombination can also generate novel

104  genotypes. Here, we define mitotic interhomolog recombination as the genetic exchange

105  between the haplotypes of the same individual that occurs in vegetative cells, when the

106  homologous chromosome is used as a template for homologous recombination, including

107  both crossing over and gene conversion events. As this can result in potentially harmful loss

108  of heterozygosity (LOH), chromosome rearrangements and consecutively mosaicism leading

109  to the onset of cancer in multicellular organisms, such recombination is typically strongly

110  suppressed in nearly all studied eukaryotes[19,20]. However, mitotic recombination is common

111  to many asexual and facultatively sexual species, including yeast[21,22], ascomycete fungi[23], and

112  oomycetes[24].

113  Here, we show that novel variation can be generated in diatom clonal populations

114  through mitotic recombination. We discuss how such variation may benefit adaptive evolution

115  under exposure to stress, and we hypothesize about the role this could play in the

116  evolutionary dynamics during clonal competition.

117

118  **RESULTS**

119  **Intraspecific SNV variability within a natural diatom population**

120  Previously, microsatellite-based approaches have demonstrated a high level of

121  intraspecific genetic variability in natural diatom populations [2]. However, as there is no good

122  evidence of genome-wide intraspecific diversity, we first sought to quantify the extent of

genome-wide intraspecific variability within natural diatom populations in situ using genome-wide metagenomic read recruitments from the Tara Oceans expeditions [1] to expand on previous knowledge. Diatom models commonly considered for fundamental genomic analyses and that are easily transformable, such as *Phaeodactylum tricornutum*, are however poorly retrieved in global environmental datasets such as the *Tara* Oceans metagenomes[1]. The most abundant diatom genera in terms of assigned 18S rRNA V9 rDNA reads in this dataset belong to *Chaetoceros* and *Fragilariopsis*, which account, respectively, for 23.1% and 15.5% of the total number of reads. While there are currently no whole-genome sequences from any *Chaetoceros* species, the genome of *Fragilariopsis cylindrus* is available and presents elevated genomic variability with around 25% of its diploid genome corresponding to highly divergent loci[25]. This genome was therefore chosen to explore diatom genomic variability in situ in the environment. We recruited metagenomic reads from *Tara* Oceans metagenomes using the *F. cylindrus* reference genome (mapping stringency >95% identity) and examined micro-diversity traits at the level of single nucleotide variants (SNVs)[26-28]. We were able to retrieve a large number of environmental sequences to this genome from Station 86 in the Southern Ocean (near the Antarctic peninsula, 64°30'88" S, 53°05'75" W), from both the surface (5m depth; mean coverage of 51.7X over genome) and deep chlorophyll maximum (DCM; 35m depth; mean coverage of 58.46X) layers. Overall, 89.64% (24,326) of *F. cylindrus* genes (total of 30.95 Mb) displayed coverage values similar to the entire genome in these two metagenomes and were considered for downstream analyses (see Methods). Within the scope of these genes, we identified 619,947 and 592,929 SNVs in the surface and DCM metagenomes, respectively (Data S1), which corresponds to an SNV density of ~2% (i.e., one SNV every 50 nucleotides). All the genes contained at least one SNV, with SNV density ranging from ~0.02% to ~10% (Data S1). Among these, 3,822 of these genes displayed SNV density <1% in both metagenomes. This

147   analysis suggests that the average nucleotide identity of the genomes considered is about

148   98%, supportive of the existence of a single population displaying thousands of micro-diversity

149   genomic traits. Parallel metabarcoding-based surveys based on 18S rRNA revealed sequences

150   most homologous to *F. cylindrus* in the samples from Station 86. Among the competing

151   nucleotides in the metagenomes, A-G and C-T transitions each contributed 30% of SNVs,

152   followed by the transversions A-C (13%), G-T (13%), A-T (10%) and C-G (4%). These statistics

153   were highly similar for the two metagenomes, with a comparable transition to transversion

154   ratio of around 1.5. Yet, of all the SNVs identified, only 429,530 (54.83%) were common to the

155   two metagenomes (Figure S1). We, therefore, conclude from this analysis that natural

156   populations of diatoms can harbour a large amount of micro-diversity, which is not restricted

157   to microsatellites but is present genome-wide.

158

159   **Genome-wide haplotype diversity**

160       In natural populations, the impact of sexual reproduction cannot be ruled out, and

161   given the vast census population size, their high nucleotide diversity is perhaps not surprising.

162   However, previous studies have hinted at genomic variability within laboratory clones, such

163   as extensive allelic diversity in the pennate polar diatom *F. cylindrus*[25] and differences

164   between *P. tricornutum* cultures belonging to the same strain derived originally from a single

165   cell[29]. Correspondingly, we noticed the presence of multiple haplotypes instead of the

166   expected two when sequencing various genomic loci in cultures clonally grown from a single

167   cell of *P. tricornutum* as well as of *Seminavis robusta*, grown under conditions that preclude

168   sexual reproduction. To map the distribution of loci with multiple haplotypes in these diatom

169   species, we took advantage of two available genome-wide datasets: short-read Illumina

170   sequencing was used to identify a set of reliable SNPs (single-nucleotide polymorphism,

171   present in at least 20% of reads) and PacBio and MinION long-read sequencing were used to

172   identify the number of haplotypes in *S. robusta* and *P. tricornutum,* respectively (Figure S2).

173   To decrease the error rate in long read sequencing, we used PacBio Circular Consensus

174   Sequences (CCS) reads and canu[30] for self-correction of both PacBio and MinION reads. Next,

175   we removed repeat regions and counted the number of combinations formed by confident

176   SNPs in individual reads in 1kb windows and selected loci with at least three haplotypes

177   supported each by a minimum of two reads. This analysis uncovered 1,405 of such loci in *S.*

178   *robusta* (125.6 Mb genome size) and 3,380 loci in *P. tricornutum* (27.4 Mb) (Figure 1A-D, Table

179   1, Figure S3, Data S2). To examine whether the number of uncovered haplotypes could be

180   caused by a high error rate in long-read sequencing datasets, we performed an equivalent

181   counting in available datasets from a haploid *Saccharomyces cerevisiae* (12 Mb) culture

182   derived from a single cell[31] and diploid *Arabidopsis thaliana* (135 Mb) datasets derived from

183   multiple inbred plants[32,33], where loci with multiple haplotypes are not expected. This yielded

184   only 3 and 83 loci with multiple haplotypes in *S. cerevisiae* and *A. thaliana*, respectively (Data

185   S2).

186

187   **Accumulation of novel haplotypes**

188       Genome-wide haplotype counting via long-read sequencing can suffer from increased

189   sequencing noise and as the datasets were derived from cultures with different cultivation

190   history, we could not conclude on the rate of the appearance of novel haplotypes. We,

191   therefore, validated the genome-wide data by observing selected loci from newly isolated,

192   single-cell cultures (Figure S4). For *S. robusta*, we profiled three loci identified in the genome-

193   wide haplotype analysis in three independent cultures, four months after single diploid cell

isolation. To overcome the potential problem of artefact generation during DNA amplification, we used emulsion PCR followed by Sanger sequencing of cloned PCR products. While the control mixture of two different alleles returned the two original haplotypes after PCR, we observed 2 to 6 haplotypes for the endogenous *S. robusta* loci (Figure 2A, Table 2) by manually examining the combinations of reliable SNPs in individual Sanger sequencing reads. Due to the low efficiency of emulsion PCR reactions, we were not able to sequence the founder cell. However, in every case, two prominent haplotypes were supported by a higher number of reads, possibly representing the haplotypes present in the founder cell, whereas the additional haplotypes presumably appeared during the four months in culture.

Independently, haplotype diversity and the rate at which new haplotypes appeared were analyzed for 62 *P. tricornutum* 2-kb loci using emulsion PCR followed by PacBio amplicon sequencing. Five loci (G32 to G36) were amplified at 1 month (T1) and 6 months (T6) after single-cell isolation, whereas the remaining 57 loci were amplified at T6 only. The heterozygosity of selected loci was profiled by SNP calling on the culture used for amplification at T1. Again, we used the short-read sequencing dataset to identify reliable SNPs and counted the number of haplotypes per locus formed by their combinations in corrected PacBio amplicon reads. The control reactions for random errors and artificial haplotype detection yielded the expected one and two haplotypes, respectively, demonstrating that the emulsion PCR, PacBio library preparation, and sequencing did not generate artefacts (Figure S4, Data S3). The number of recovered haplotypes varied between 1 and 15 (Figure 2C, Figure S4 and Data S3), with 6 loci displaying a single haplotype, 5 loci with two haplotypes, and 51 loci displaying at least three haplotypes. For four out of five loci amplified at both T1 and T6, an increase in the number of haplotypes was observed in the T6 sample (Figure 2B, Figure S4) despite deeper sequencing coverage of the T1 samples, suggesting that haplotypes

accumulate over time (Table 3). We analyzed the impact of haplotype variability on protein sequence in 20 genes fully covered by amplicon sequencing and found six for which the different haplotypes resulted in more than two putative protein variants, with up to six variants in the diatom-specific gene *Phatr3_J47122* (Figure 2D, Figure S4, Table S1).

Although only one locus was identified as homozygous by SNP calling in T1, five additional loci with single haplotypes were found in T6 sequencing (Figure 2C, Figure S4). These loci were identified as being heterozygous by SNP calling in T1, suggesting a loss of heterozygosity (LOH) [34]. Moreover, the novel haplotypes that accumulated in both *S. robusta* and *P. tricornutum* cultures were recombinants lacking *de novo* mutations. Such new combinations are typically generated during sexual reproduction through the meiotic recombination of homologous chromosomes[35].

**Genome-wide detection of loss of heterozygosity and copy number variation**

Although interhomolog recombination is rare in vegetative cells, we tested whether it could be the source of haplotype diversity in clonal diatom populations as sexual reproduction was excluded in our cultures. We sought to detect LOH and copy number variation (CNV) events in *P. tricornutum* under controlled conditions over a defined number of cell divisions. Three independent mother cultures (MC1 - MC3) were initiated from a single cell isolate and cultivated under conditions allowing approximately a single cell division per day (Figure 3A, Figure S5). After 30 days (T1), three single cells were again isolated from each mother culture to obtain nine daughter cultures (DC1.1 - DC3.3) that were harvested 30 days later (T2). At both T1 and T2, part of the mother cultures was also harvested. Following genome resequencing and SNP calling of all cultures, a pairwise comparison between the individual

241 daughter cultures and their respective mother cultures was performed to identify novel CNVs

242 and tracts of at least three consecutive SNPs that were lost in the daughter culture.

243 Changes in comparison with the mother cultures were found in four out of nine

244 daughter cells. One copy-neutral 8016 bp LOH, where one allele of the locus was replaced by

245 the other allele, was observed in DC1.2, three copy-neutral LOH events (296 bp, 614 bp and

246 1644 bp in length) and a 31.4 kb duplication covering 14 genes were observed in DC1.3 culture,

247 and one 30.9 kb and one 156.9 kb deletion were detected in cultures DC3.1 and DC2.1,

248 respectively, and were confirmed by Sanger re-sequencing or qPCR (Figure 3, Figure S5, Table

249 S2, Data S3). Besides the LOH events that were unique to a respective daughter culture, we

250 identified several regions with reduced SNP density common to all cultures. SNP density was

251 ten times lower than the genome average over almost the entire chromosome 19, and

252 seventeen and thirty-eight times lower at the extremities of chromosomes 27 and 28,

253 respectively, in comparison with the rest of the chromosome (Figure S5). These regions were

254 not found to be SNP poor when sequencing the same *P. tricornutum* strain from other

255 laboratories[36-38].

256 Profiling of the functional effect of 2914 SNPs in LOH regions revealed 59 SNPs

257 (0.362%) with possible high impact on gene function, 650 (3.984%) with low, 702 (4.303 %)

258 with moderate and 14,903 (91.351%) with modifier effect according to SnpEff

259 categorization[39]. Most SNPs with a high effect on protein function were found in the 156.9 kb

260 deletion on chromosome 26. This deletion was identified in the primary analysis as six LOH

261 regions and confirmed as a single deletion only after Sanger resequencing of LOH border

262 regions. Therefore, we were only able to analyze the effect of the 19 SNPs with high effect

263 found in the regions identified in the primary LOH analysis (Data S3) and found 3 SNPs that

264    caused a loss of function by introducing a premature stop codon in the respective gene (Data

265    S3).

266

**Copy-neutral loss of heterozygosity at the PtUMPS locus**

268    While the mechanism behind the observed deletions and duplication remains difficult

269    to interpret, the copy-neutral LOH events require an exchange of genetic information between

270    homologous chromosomes. To estimate the rate of interhomolog recombination in *P.*

271    *tricornutum*, we established a tractable endogenous readout system for copy-neutral LOH

272    detection, based on three strains containing two different mutant alleles of the *PtUMPS* gene,

273    generated through gene editing[40,41]. In strain *ptumps-1bp,* the 1 bp indel mutations in the two

274    alleles occur at a position only 1 bp apart, in strain *ptumps-320bp* they are separated by 320

275    bp and in strain *ptumps-1368bp* by 1368 bp (Figure 4A). As the *PtUMPS* protein is required for

276    uracil biosynthesis, cells with a wild-type (WT) allele can synthesize uracil, but also convert 5-

277    fluoroorotic (5-FOA) acid into the toxic 5-fluorouracil (5-FU), resulting in cell death. In contrast,

278    mutant cells are resistant to 5-FOA but are uracil auxotrophs. The *ptumps-/-* strains were

279    cultivated under non-selective conditions for 14 days (with uracil and without 5-FOA) to

280    permit potential recombination at the *PtUMPS* locus (Figure S6). Subsequently, $5 \times 10^7$ cells

281    from the culture were plated on a medium without uracil to select cells that underwent

282    recombination at the *PtUMPS* locus and restored the WT allele. We recovered no colonies in

283    strain *ptumps-1bp*, confirming that the WT allele was not restored by a random mutation, 12

284    colonies in strain *ptumps-320bp* and 83 colonies in strain *ptumps-1368bp* (Figure 4B, Data S4-

285    S5). Moreover, sequencing of *PtUMPS* alleles from ten *ptumps-1368bp* colonies and five

286    *ptumps-320bp* colonies corroborated the restoration of the WT allele through copy-neutral

287    LOH events (Figure 4C, Figure S6).

288    Next, the *PtUMPS* system was used to obtain an estimate of the interhomolog

289    recombination frequency. A total of $2 \times 10^7$ cells per replica from 5-FOA- and uracil-

290    supplemented medium (preventing recombination at the *PtUMPS* locus) were directly plated

291    onto medium without uracil to select only those cells that were in the process of interhomolog

292    recombination during a single round of cell division. The average frequency of interhomolog

293    recombination was 4.2 per 100 cell divisions per genome (Figure 4D, Data S4-S5),

294    approximately ten times higher than the rate reported for *S. cerevisiae* after recalculation per

295    cell division[42,43].

296    To test whether the rate of interhomolog recombination can be influenced by

297    environmental conditions, we employed the *PtUMPS* readout system to test the effect of the

298    DNA double-strand break inducing drug zeocin[44] and three physiologically relevant stresses:

299    hydrogen peroxide ($H_2O_2$), which is produced by various phytoplankton groups and can act as

300    a signalling molecule as well as cause oxidative damage[45], the trace metal cadmium[46], which

301    contaminates aquatic environments, and a polyunsaturated aldehyde (E,E)-2,4-Decadienal

302    that is involved in diatom intercellular signalling, stress surveillance, and defence against

303    grazers, but which can trigger lethality at high concentrations[47,48]. For each mock and stress

304    treatment, $25 \times 10^6$ cells per replica were transferred from 5-FOA- and uracil-supplemented

305    medium to medium containing uracil for 24 h, thus allowing a maximum of one cell division.

306    Next, cells were plated on a selective medium without uracil to recover cells that restored the

307    WT *PtUMPS* allele through interhomolog recombination. Only the zeocin treatment resulted

308    in the appearance of uracil prototrophic colonies in both *ptumps-1bp* and *ptumps-1368bp* in

309    a dose-dependent manner (Figure 4E). Sequencing of *ptumps-1bp* colonies revealed

310    restoration of the *PtUMPS* WT allele through *de novo* mutations (Figure S6). We thus suppose

311    that zeocin treatment induced robust DNA damage. No *ptumps-1bp* colonies were observed

312    in the other treatments hinting at a lack of *de novo* mutations. Whereas the (E, E)-2,4-

313    Decadienal treatment did not influence the rate of interhomolog recombination, we found a

314    positive, concentration-dependent effect of $H_2O_2$ and cadmium on the number of recovered

315    colonies (Figure 4F-G, Data S4-S5). These data illustrate that environmental stresses increase

316    the frequency of recombination between homologous chromosomes.

317

318    **DISCUSSION**

319    Analysis of environmental samples of subpopulations of the diatom *F. cylindrus*

320    showed that they harbour extensive genome-wide SNV diversity. However, as *F. cylindrus* is

321    not easily accessible to genome manipulation methods, we investigated the possible

322    underlying mechanism in commonly used diatom model species. By following the number of

323    haplotypes in cultures of *S. robusta* and *P. tricornutum* initiated from a single cell, we

324    documented that both diatom species rapidly accumulate recombined haplotypes throughout

325    the genome. The resulting novel SNP combinations in protein-coding genes can give rise to

326    novel protein variants that were not present in the founder cell, potentially contributing to

327    the physiological divergence of individual subclones in the clonal population. Additionally, a

328    comparison of genomic changes between mother cultures and their respective daughter

329    cultures revealed the appearance of copy-neutral LOH and CNV events over a brief period. We

330    hypothesize that CNVs arise from ectopic recombination or non-homologous end-joining[49-51].

331    In the copy-neutral LOH events, the information from the homologous chromosome either

332    replaces the original allele in case of a gene conversion event, or it leads to reciprocal

333    exchange in case of mitotic crossing over. Subsequent sister chromatid segregation during

334    mitosis may cause LOH tracts in the daughter cell(s), resulting in the fixation of polymorphisms

335  in a homozygous state[19,20], further contributing to phenotypic differences within the clonal

336  population.

337      Estimating the rate of the mitotic interhomolog recombination per cell division

338  revealed that the frequency in *P. tricornutum* exceeds by ten times the frequency in the yeast

339  *S. cerevisiae*, the key model in mitotic recombination research [42,43]. Although this comparison

340  does not take into account the possible differences between the diatom and yeast outcomes

341  of the interhomolog recombination, it suggests that such recombination is highly common in

342  *P. tricornutum* and that the constraints preventing the use of homologous chromosomes as a

343  template for homologous recombination might be relaxed in diatoms. The capability to rapidly

344  fix novel SNVs in a population through LOH could explain the differences observed in the

345  metagenomes of *F. cylindrus* subpopulations of the surface and DCM samples from the same

346  station.

347      We demonstrated that the rate of mitotic recombination increased under

348  environmental stress, which suggests that it has a degree of phenotypic plasticity. A similar

349  increase was documented for both meiotic and mitotic recombination in various organisms

350  including yeasts[52,53], plants[54-56] and metazoans[57-59] and has important implications for

351  evolution. Recombination related genomic changes were shown to shape the genomes of

352  pathogenic fungi such as *Candida albicans*, where the frequency increases under stress and

353  during host infections and contributes to the fitness advantage of resulting clones[60], as well

354  as in the oomycete *P. ramorum* where extensive runs of homozygosity (ROH) differentiate

355  individual invasive lineages[24]. However, in both pathogenic species, the recombination

356  involved preferentially the repetitive regions and transposons, and the exact frequency is not

357  known. By contrast, repetitive sequences were excluded in our analysis, and recombinant

358  haplotypes were found to be equally dispersed throughout the *P. tricornutum* genome.

359   Besides *de novo* LOH events, the analysis of mother and daughter cells revealed the

360   presence of regions with low SNP density common to all sequenced strains. These regions

361   were not detected as being low in SNP content in *P. tricornutum* Pt1 strains from other

362   laboratories [36,37] and a similar situation was reported for other genomic regions[29]. Low

363   heterozygosity regions can also arise due to inbreeding or purifying selection at a linked

364   genetic loci. However, *as P. tricornutum* has never been observed to reproduce sexually in

365   laboratory conditions, we propose that these might represent past LOH events in the ancestor

366   cell of the respective population. Low SNP density regions have also been observed in

367   presumably asexual isolates of the centric diatom *Thalassiosira pseudonana*[4]. It was

368   speculated that the loss of heterozygosity in these regions due to inbreeding resulted in the

369   fixation of mutations in genes required for sexual reproduction. In the light of high levels of

370   mitotic interhomolog recombination in *P. tricornutum*, an alternative cause of the decrease in

371   heterozygosity could be LOH accompanying such recombination.

372   Many diatom species accomplish rapid population expansion through clonal

373   reproduction[61]. During this phase of exponential growth, a small difference in fitness can have

374   a large effect on the eventual population number reached by each clonal lineage. However,

375   significant environmental changes are likely to deteriorate the fitness of any well-adapted

376   lineage. Yet, without sexual reproduction, each clonal lineage is limited in its adaptive

377   response by the variation contained within its genome. We hypothesize that mitotic

378   recombination can exploit the non-additive genetic variation (i.e. dominance and epistatic

379   variation) that is present within each genome but hidden from natural selection[62,63].

380   Plastic response in mitotic recombination could offer at least three important fitness

381   advantages during clonal competition. Firstly, the evolution of asexual microbes is generally

382   not limited by the number of beneficial single-point mutations, but rather, by overcoming

383 clonal interference and combining multiple mutations into a single genotype. For example,

384 beneficial mutations are readily available in yeast, but they compete with one another in the

385 population for fixation[64-66]. Mitotic recombination can relax the evolutionary constraints

386 imposed by clonal interference, by generating novel combinations of alleles. Alleles can thus

387 be 'tried and tested' against slightly different genomic backgrounds, which increases the

388 probability of finding a superior combination of multiple mutations. Mitotic recombination

389 thus can not only uncover hidden dominance variation by making loci homozygous, but it can

390 also reveal epistatic variation by creating novel allelic combinations that would otherwise not

391 have arisen. This could be particularly important during periods of environmental change and

392 stress, enabling the clonal lineage to discover other fitness peaks in a dynamic fitness

393 landscape[67].

394 Secondly, density and frequency-dependent processes are likely to regulate clonal

395 expansion. Lewontin put this succinctly: "a genotype is its own worst enemy, its fitness will

396 decrease as it becomes more common"[68]. Such negative frequency dependence is likely to

397 play an important role in asexual species, particularly during and after clonal expansion. The

398 ability to generate novel genotypes during mitotic recombination could mitigate this effect,

399 reducing the competition between clone-mates, and generating a more diffuse target for

400 antagonistically coevolving species, such as pathogens.

401 Thirdly, generating evolutionary novelty, either through mutation or recombination,

402 does impose a fitness cost to the individual or clonal lineage, i.e. a genetic load[69]. In a well-

403 adapted genotype, each mitotic recombination event is more likely to reduce fitness than to

404 increase it. However, occasionally, some mitotic recombination events could be selectively

405 advantageous, and this is more likely to occur if the clonal genotype is not optimally adapted

406 to its environment[22]. In diatoms such as *P. tricornutum*, natural selection thus trades off the

costs of the 'mitotic recombination load' against the potential benefits realized by such recombination. These benefits include reducing possible negative frequency-dependent effects, and uncovering hidden dominance and epistatic variance, enabling the genotype to climb or discover a fitness peak in the adaptive landscape. In other words, there may be an optimum level of mitotic recombination, depending on the stability of the environment, the match between the phenotype and the environment, and the amount of negative frequency-dependent selection. We hypothesize that phenotypic plasticity may enable lineages to track this optimal level of mitotic recombination, with natural selection favouring an increased rate under stressful environmental conditions. Alternatively, stressful environmental conditions could increase the mitotic recombination rate, for example by impairing DNA repair mechanisms. The question not answered by our experiments is whether the observed increase in mitotic recombination during stress is adaptive. We propose this is plausible, and that this hypothesis provides an interesting avenue for future research.

**AUTHOR CONTRIBUTIONS**

P.B. and L.D.V. conceptualized the study. P.B. initiated, designed and performed experiments and bioinformatics analysis on *P. tricornutum* and *S. robusta*. M.S. performed sequencing of uracil prototrophic colonies. D.J. generated *ptumps* mutant strains. C.N and T.D. performed the *F. cylindrus* metagenome assembly and analysis, I.V. helped with diatom culture maintenance. C.M.O-C. and E.M. provided bioinformatic datasets. P.B. and C.V.O. wrote the manuscript, generated all figures and data visualizations. L.D.V., K.V., F.D., C.B., W.V., K.S. supervised the research. P.B., L.D.V, K.V., W.V., F.D, K.S., C.B., C.V.O., T.M. reviewed and edited the manuscript.

**DECLARATION OF INTERESTS**

The authors declare no competing interests.

**MAIN FIGURE TITLES AND LEGENDS**

**Figure 1. Genome-wide distribution of haplotypes in *P. tricornutum* and *S. robusta*.** (A-B) Distribution of the detected number of haplotypes per 1 kb loci in *S. robusta* contigs above 20 kb (A) and *P. tricornutum* chromosomes 1-33 (B). From outside to inside: loci with more than two haplotypes (red), loci with two haplotypes (blue), loci with a single haplotype (grey), gene density*, SNP density* (black), GC content*; * per 10 kb. (C-D) Example of a representative genomic region from *S. robusta* (C) and *P. tricornutum* (D). The chromosome is represented by a grey rectangle. Above the chromosome: loci with single haplotype (grey bars), loci with two haplotypes (blue bars) and loci with more than two haplotypes (red bars). Below

454　chromosome from top to bottom: gene density**, SNP density** (grey line), GC content**;

455　** per 1 kb. See also Figures S1-S3 and Data S1-S2.

456

457　**Figure 2. Accumulation of novel haplotypes in cultures freshly initiated from a single cell.** (A)

458　Quantification of haplotypes on three loci in three *S. robusta* cultures (SR1-3) four months

459　after cultivation from a single cell. (B-C) Quantification of the number of haplotypes in *P.*

460　*tricornutum* at 1 month (T1) and 6 months (T6) after cultivation from a single cell. (B)

461　Quantification of the number of haplotypes at loci G32-G36 detected at T1 and T6. (C)

462　Quantification of the number of haplotypes detected at T6 (orange outline). Categories of the

463　shift in the number of haplotypes from founder cell to the number of haplotypes detected at

464　T6 are on the x-axis in the format: expected in founder cell > observed at T6, the number of

465　haplotypes on the y-axis. The size of the circle corresponds to the number of cases in a given

466　category. In (A-C) the expected two haplotypes in the founder cells are indicated in blue. (D)

467　Schematic representation of predicted proteins variants in Phatr3_J47122 gene. The top line

468　shows the position of amino acid variants on the protein indicated by red flags. Green regions

469　depict conserved domains according to the CDD/SPARKLE database [70]. The lines below

470　represent individual predicted variants. See also Figures S4, Table S1 and Data S3.

471

472　**Figure 3. Genome-wide detection of LOH in *P. tricornutum* mother and daughter cultures**

473　**after 30 days.** (A) Position of copy-neutral LOHs (orange), duplication (dark blue) and deletions

474　(red) in individual daughter cultures. Heterozygous regions are in light blue, blank space – no

475　SNPs. Nominators at the left refer to the daughter cell (DC) culture, whereas the lowercase

476　digit indicates the chromosome number. (B) Zoomed-in regions with detected LOH,

477 duplication or deletion events in the respective mother and daughter cultures (DC). Blue dots

478 – heterozygous SNPs, orange dots – homozygous SNPs in copy-neutral LOHs, red dots - SNPs

479 in deletions, grey area – sequence coverage. (C) Confirmation of a DNA duplication event on

480 chromosome 23 by qPCR in daughter culture DC1.3. The position of target loci (red boxes) is

481 shown in the upper part. The bar chart depicts the fold change in comparison to the MC1T1

482 sample on control loci D, E and F. Blue dots – heterozygous SNPs, dark blue - SNPs in

483 duplication. See also Figure S5, Table S2 and Data S3.

484

485 **Figure 4. Detection of LOH events at the *P. tricornutum PtUMPS* locus.** (A) Schematic of

486 alleles in the *PtUMPS* strains. Homologous chromosomes are depicted as grey bars with exons

487 in blue and green, loss-of-function indel mutations are in red, purple and orange bars

488 represent silent SNPs between the two original alleles. (B-C) Recombination in *PtUMPS*

489 mutant strains during 14 days of cultivation under non-selective conditions. (B) The number

490 of recovered uracil prototrophic colonies per strain. (C) Examples of sequenced recombinant

491 alleles in one *ptumps-320bp* and two *ptumps-1368bp* colonies. (D) Estimation of the

492 interhomolog recombination frequency per thousand cell divisions. Each dot represents one

493 replica. (E–H) Recombination events in response to stress-induced by (E) zeocin; (F) $H_2O_2$; (G)

494 cadmium and (H) 2,4-Decadienal. *ptumps-1bp* replicas are depicted in shades of grey, *ptumps-*

495 *1368bp* replicas are depicted in shades of blue. See also Figure S6 and Data S4-S5.

496

497

**TABLES**

**Table 1. Characteristics of loci with multiple haplotypes found in *S. robusta* and *P.**

***tricornutum***

| | Whole-genome average | | Loci with multiple haplotypes | |
|---|---|---|---|---|
| | **Total number** | **Percent** | **Total number** | **Percent** |
| ***Seminavis robusta* (1405 loci)** | | | | |
| GC content | | 48.5 % | | 48.8 % |
| SNPs total | 489799 | 100 % | 7714 | 100% |
| SNPs in intergenic regions | 149782 | 30.58 % | 2662 | 34.50 % |
| SNPs in protein coding genes | 339890 | 69.39 % | 5050 | 65.47 % |
| Intron | 17300 | 3.53 % | 233 | 3.03 % |
| Exon | 322590 | 65.86 % | 4817 | 62.44 % |
| Functional RNAs | 127 | 0.03 % | 2 | 0.03 % |
| ***Phaeodactylum tricornutum* (3380 loci)** | | | | |
| GC content | | 48.77 % | | 49.04 % |
| SNPs total | 290164 | 100% | 22531 | 100% |
| SNPs in intergenic regions | 110727 | 38.16 % | 6767 | 30.03 % |
| SNPs in protein coding genes | 178906 | 61.65 % | 15735 | 69.83 % |

| | | | | |
|---|---|---|---|---|
| Intron | 16271 | 5.61 % | 1242 | 5.51 % |
| Exon | 162635 | 56.04 % | 14493 | 64.32 % |
| Pseudogenes | 425 | 0.15 % | 27 | 0.12 % |
| Functional RNAs | 106 | 0.04 % | 2 | 0.01 % |

501

502 **Table 2. Verification of haplotype diversity in *S. robusta* at three selected loci in three**

503 **cultures (Sr1 – Sr3) through emulsion PCR amplification, cloning and Sanger sequencing of**

504 **individual clones**

| Locus | Number of SNPs | Culture | Number of haplotypes at 4 months after single cell isolation | Number of supporting reads for each haplotype |
|---|---|---|---|---|
| Sro_contig211: 7509-8241 | | | | |
| | 11 | Sr1 | 5 | 26; 15; 1; 1; 1 |
| | 11 | Sr2 | 4 | 23; 18; 8; 1 |
| | 11 | Sr3 | 3 | 26; 25; 2 |
| Sro_contig2103: 8397-9162 | | | | |
| | 5 | Sr1 | 3 | 19; 12; 1 |
| | 5 | Sr2 | 3 | 8; 6; 1 |
| | 5 | Sr3 | 2 | 24; 13 |
| Sro_contig872: 16034-16975 | | | | |
| | 4 | Sr1 | 3 | 24; 16; 1 |
| | 4 | Sr2 | 5 | 19; 14; 1; 1; 1 |
| | 4 | Sr3 | 6 | 29; 19; 3; 2; 1; 1 |
| Sro_contig556:54453-55487 - control mix of plasmids containing two alleles of the locus | | | | |
| | 3 | - | 2 | 44; 19 |

505

**Table 3. Change in number of recovered haplotypes over time in *P. tricornutum***

| Locus name | Coordinates | Number of SNPs | Samples harvested at 1 month after single cell isolation | | Samples harvested at 6 months after single cell isolation | |
|---|---|---|---|---|---|---|
| | | | Number of haplotypes | Coverage | Number of haplotypes | Coverage |
| G32 | 13:103145-105183 | 10 | 8 | 3903 | 12 | 878 |
| G33 | 27:205731-207770 | 18 | 6 | 5712 | 8 | 988 |
| G34 | 20:101923-103985 | 28 | 4 | 3140 | 8 | 654 |
| G35 | 12:519921-521994 | 12 | 9 | 2992 | 10 | 706 |
| G36 | 2:961633-963692 | 12 | 8 | 5225 | 8 | 688 |

507

508

509 **STAR METHODS**

510

511 **RESOURCE AVAILABILITY**

512 **Lead Contact**

513 Further information and requests for resources and reagents should be directed to and will

514 be fulfilled by the Lead Contact, Lieven De Veylder (lieven.deveylder@psb.vib-ugent.be)

515

516 **Materials Availability**

517 Material generated in this study is available upon request from the lead contact

518

519 **Data and Code Availability**

520 Raw sequencing data were deposited to the Sequence Read Archive (SRA) under BioProject

521 accessions PRJNA658511 and PRJNA658224. SRA accession numbers for individual samples

522 are listed in Data S3. Processed datasets were uploaded to zenodo: Aligned and processed

523 long-read sequencing datasets *S. robusta* PacBio, *P. tricornutum* MinION reads and SNP

524 selected for haplotype counting for both species are available at

525 https://doi.org/10.5281/zenodo.4005721. Aligned PacBio amplicon sequencing reads,

526 reference file and selected biallelic SNPs used in haplotype counting are available at

527 https://doi.org/10.5281/zenodo.4005643. Processed datasets from LOH detection in mother

528 and daughter cultures including ILLUMINA reads aligned to the reference P. tricornutum

529 genome used for SNP calling, SNP calls for individual samples and jointly called SNPs on all

530 samples are available at https://doi.org/10.5281/zenodo.4006016. **Code availability**

531    The haplotype coding script to count haplotypes in long-read sequencing datasets is available

532    on zenodo, at https://doi.org/10.5281/zenodo.4001752. A version of the haplotype counting

533    script that outputs the combination of bases at selected SNP sites is available on zenodo, at

534    https://doi.org/10.5281/zenodo.4173002. All other data are available from the authors upon

535    request.

536

537    **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

538

539    <u>Diatoms datasets and strains</u>

540    Datasets and strains used in this study are summarized in Data S3. The *S. robusta* D6 reference

541    strain (accession number DCG 0498) is available from the BCCM/DCG diatom culture collection

542    at Ghent University (http://bccm.belspo.be/about-us/bccm-dcg). Publicly available genomes

543    of *S. robusta* strain D6 [71] https://www.ebi.ac.uk/ena/browser/view/CAICTM010000000 and *P.*

544    *tricornutum*        Pt1        8.6        (CCMP2561)        strain        [91]

545    https://www.ebi.ac.uk/ena/browser/view/GCA_000150955.2        and        next-generation

546    sequencing datasets  were used for our analysis.

547

548    <u>Diatom cultivation conditions</u>

549    Both *S. robusta* strain D6 and *P. tricornutum* strain Pt1 subculture MC2 were cultivated in 1x

550    TMB medium consisting of 34.5 g/L of Tropic Marin Bio-Actif sea salt (Tropic Marin, Germany)

551    and 0.08g/L sodium bicarbonate (Sigma-Aldrich) supplemented with 1x Guillard's (F/2) Marine

552    Water Enrichment Solution (Sigma-Aldrich), 100 µg/mL ampicillin, 50 µg/mL gentamycin and

553    100 µg/mL streptomycin in 12 h/12 h light/dark cycle. *P. tricornutum* cultures were cultivated

554    at 20°C, under photosynthetic LED light with an intensity of 160 µmol photons $m^{-2}$ $s^{-1}$ and with

555  100 rpm shaking. *S. robusta* cultures were cultivated at 18°C with approximately 85 μmol

556  photons m$^{-2}$ s$^{-1}$ from cool-white fluorescent lights.

557

558  **METHOD DETAILS**

559  <u>Estimation of intra-specific variability in *Fragilariopsis cylindrus* metagenomes</u>

560  *Tara* Oceans metagenomic reads from 0.8-5 μm, 5-20 μm, 20-180 μm, 180-2000 μm and 0.8-

561  2000 μm size fractions were mapped against the FASTA file of the *Fragilariopsis cylindrus*

562  CCMP 1102 genome [25] (available at http://genome.jgi-psf.org/Fracy1/Fracy1.home.html)

563  using Bowtie2 v2.3.4.332 with a 95% identity filter. Two depths, surface (5m depth) and deep

564  chlorophyll maximum (DCM; 35m depth), both located in the epipelagic mixed layer [92] from

565  Station 86 situated in the Southern Ocean (near the Antarctic peninsula, 64°30'88" S,

566  53°05'75" W), displayed vertical sequence coverage superior or equal to 10X for three size

567  fractions (0.8-5, 5-20 and 0.8-2000 μm) and were thus selected for further analysis. Using

568  SAMtools v1.10 33, the resulting SAM files were converted into BAM files and for each sample

569  the BAM files of the three different size fractions were merged to increase the coverage (final

570  mean coverage 51.75X and 58.46X for surface and DCM respectively). Downstream analyses

571  were performed with the anvi'o platform[73] to generate profile databases based on the BAM

572  files that were combined into a merged profile database. Genes were imported into anvi'o at

573  the level of individual exons. Then, the program "anvi-summarize" was used with the "init-

574  gene-coverages" flag to characterize the mean coverage of each gene in the surface and DCM

575  samples. Genes that were considered for downstream analyses (n = 24,326) were invariably

576  detected within a population niche (here the metagenome) [93]. These genes had to occur in

577  the two samples and their mean coverage in each sample had to remain within a factor 3 of

578  the mean coverage of all 27,137 genes in the same metagenome. The filtering step based on

579    gene level coverage values is critical to remove outlier genes that may recruit reads from other

580    related genera or species that potentially co-occurred in the samples (e.g., the 18S rRNA gene

581    will recruit reads from other genera due to its high evolutionary stability, and genes from

582    closely related species will display higher coverage values compared to the species-specific

583    genes). Additionally, it allows to remove genes with hypervariable regions that will not recruit

584    reads, preventing the subsequent analysis of single nucleotide variants (hereafter referred to

585    as SNVs) [94]. Finally, the intra-population variability of *F. cylindrus* was analysed across the

586    selected genes and in the two samples using the programme "anvi-gen-variability-profile",

587    which provided tables reporting SNVs and their nucleotide frequencies in the recruited reads

588    (Data S1). We defined SNVs as positions displaying at least 10% variation from the consensus

589    nucleotide and with a mean vertical coverage ≥ 20X in the two samples. The variability tables

590    were imported into R v4.0.1 to compute the number of variable positions and SNV density (i.e.

591    the number of positions with SNVs for each exon in the selected genes divided by the

592    corresponding exon length) for each exon. Gene-level mean coverage, number of variable

593    positions and SNV density were computed using the information from the individual exons.

594

595    <u>Genome-wide haplotype counting in *S. robusta* and *P. tricornutum* next-generation</u>

596    <u>sequencing data</u>

597    For both *S. robusta* and *P. tricornutum* genome-wide haplotype counting, a reliable single

598    nucleotide polymorphism (hereafter referred to as SNP; in contrast to SNVs found in

599    metagenomes from natural populations, SNP had been supported by at least 20% of reads in

600    the sample from laboratory single strain) set was first identified in ILLUMINA short-read

601    sequencing datasets and then used for counting of the number of haplotypes in the PacBio RS

602    II and MinION long reads. The ILLUMINA and PacBio data of *S. robusta* from

603 https://www.ebi.ac.uk/ena/browser/view/PRJEB36614 and ILLUMINA and Minion data of *P.*

604 *tricornutum* from https://www.ebi.ac.uk/ena/browser/view/PRJNA487263. Because in long-

605 read sequencing the error rate for indels is higher than for SNPs, indels were ignored in our

606 analysis.

607

608 *SNP calling*: SNP calling on ILLUMINA short-read sequencing was done using GATK

609 HaplotypeCaller 3.7.0[80]. In short, adapters and reads with a quality score below 20 were

610 removed from ILLUMINA reads using BBduk2 [75] with minlen=35 qtrim=rl trimq=20 hdist=1 tbo

611 tpe options and custom adapter reference file. Next, the respective reads were aligned to *S.*

612 *robusta* v1 assembly CAICTM010000001-CAICTM010004752 (European Nucleotide Archive)

613 or *P. tricornutum* v2 assembly2 GCA_000150955.2 (European Nucleotide Archive) using

614 Burrows-Wheeler Alignment Tool (BWA)[76] algorithm BWA-MEM with -M option. Unmapped

615 and multi-mapped reads were removed using SAMtools [77] view with -h -F 4 -q 1 options.

616 Aligned reads were then sorted using picard-tools 1.8.0 [78] SortSam and duplicate reads were

617 marked with MarkDuplicates and indexed with BuildBamIndex. Read base quality scores were

618 adjusted by two round of recalibration. Here, SNPs and indels were called by GATK

619 HaplotypeCaller[79] and filtered with a set of hard filters using SelectVariants; QD < 2.0, FS >

620 60.0, MQ < 40.0, MQRankSum < -12.5, ReadPosRankSum < -8.0 for SNPs and QD < 2.0, FS >

621 200.0, ReadPosRankSum < -20.0 for indels. Recalibration table was generated with

622 BaseRecalibrator and recalibrated reads were printed with PrintReads. After the second round

623 of recalibration, germline SNPs were called using HaplotypeCaller with --genotyping_mode

624 DISCOVERY. Next, reliable biallelic SNPs were selected using SelectVariants with --

625 restrictAllelesTo BIALLELIC -selectType SNP and QD < 2.0, QUAL < 30.0, SOR > 3.0, FS > 60.2,

626 MQ < 40.0, MQRankSum < -12.5, ReadPosRankSum < -8.0, AF > 0.2 and DP < 10 options.

627　Repeat regions and low complexity DNA sequences in *S. robusta* and *P. tricornutum* were

628　identified using RepeatModeler 1.0.9[81] and masked using RepeatMasker 4.0.5 [82], and SNPs in

629　these regions were removed from the dataset using BEDtools [83] subtract algorithm. Finally,

630　selected fields (CHROM, POS, REF, ALT) from the SNP dataset were extracted from the vcf file

631　to a table and split into independent files by contig/chromosome using awk.

632

633　*S. robusta PacBio reads processing:* Circular Consensus Sequences (CCS) were obtained with

634　smrtanalysis 2.3.0 (PacBio) with minFullPasses 0 option. CCS reads were then self-corrected

635　using canu 1.4[30] with canu_correct genomeSize=136.0m errorRate=0.035 -pacbio-raw

636　options and trimmed with canu_trim genomeSize=136.0m errorRate=0.035 -pacbio-corrected

637　options. Corrected reads were mapped to the reference genome using BLASR[89] with -sam -

638　clipping soft options. The CIGAR string was corrected with samfixcigar, soft-clipped bases were

639　removed with biostar84452 from jvarkit[84] and uniquely mapped reads with mapping quality

640　>20 were selected using SAMtools[77]. Coverage was estimated using GATK 3.7.0

641　DepthOfCoverage. SAMtools view and awk were used to split the PacBio reads to separate

642　files per contig.

643

644　*P. tricornutum MinION reads processing*: MinION reads were self-corrected using canu 1.4[30]

645　with canu_correct genomeSize=30m errorRate=0.144 -nanopore-raw options. Reads were

646　aligned to the genome using GraphMap[85] with default settings and uniquely mapped reads

647　were selected using SAMtools view. The CIGAR string was corrected with samfixcigar, soft-

648　clipped bases were removed with biostar84452 from jvarkit. Coverage was estimated using

649　GATK 3.7.0 DepthOfCoverage. SAMtools view and awk were used to split the PacBio reads to

650　separate files per contig.

651

*Haplotype counting:* The haplotype counting was done with a custom script based on bash, awk and Sam2Tsv from jvarkit ([https://doi.org/10.5281/zenodo.4001752)](https://doi.org/10.5281/zenodo.4001752). In short, a record for every base in each processed PacBio/MinION read with the position, reference and the actual base was obtained using Sam2Tsv from jvarkit[84]. Next, only positions of SNPs selected in ILLUMINA reads were retained. The record was divided into fixed windows of max 1 kb from the first SNP and haplotypes for selected sites were written for each read separately. Reads containing an indel or another base than the reference or the alternative base at the selected SNP position or not covering the 1 kb region were removed and the number of haplotypes and number of supporting reads for each haplotype was counted using awk. Loci with multiple haplotypes were selected from the record of the number of haplotypes per 1 kb with the number of supporting reads with the following conditions: at least three haplotypes had to be supported each by at least 2 reads, the locus had to be at least 100 bp long and the coverage had to be below 100x to remove repeat regions that were not masked. Visualization of haplotype counting data was done using Circos [86] and karyoploteR [88].

666

*Haplotype counting in control genomes*: As a control for haplotype counting overcounting due to high error rate in long-read sequencing data, the genome assembly, ILLUMINA and PacBio sequencing data publicly available for haploid yeast *Saccharomyces cerevisiae* GLBRCY22-3 [31] grown from single colony [https://www.ncbi.nlm.nih.gov/bioproject/PRJNA279877](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA279877) and for several diploid *Arabidopsis thaliana* Ler plants [32,33] from [https://www.ncbi.nlm.nih.gov/bioproject/PRJNA311266](https://www.ncbi.nlm.nih.gov/bioproject/PRJNA311266) and [https://www.ncbi.nlm.nih.gov/bioproject/237120](https://www.ncbi.nlm.nih.gov/bioproject/237120) were used. The SNP calling, PacBio data

674  processing and haplotype counting were done as described above and are summarized in Data

675  S2.

676

677  <u>Calculation of the error rate in self-corrected long-read sequencing data used for haplotype</u>

678  <u>counting.</u>

679  The long-read sequencing technology is known to have a higher error rate than the short-read

680  sequencing technology[95]. To improve the read quality, we generated circular consensus

681  sequencing (CCS) reads for the PacBio datasets and used self-correction for both PacBio and

682  MinION reads (detailed description is available in the Methods). The self-correction was

683  preferred over the correction by ILLUMINA reads, as the ILLUMINA reads were used for calling

684  an SNP dataset and this could introduce bias in subsequent haplotype counting. The long-

685  sequencing reads were further processed after alignment to the genome by removal of soft

686  clipping, correction of CIGAR string and selection of uniquely mapped reads. The error rate in

687  corrected and aligned PacBio and MinION files was estimated using Alfred [90] (Table S8).

688  The haplotype counting script processes only positions of reliable SNPs selected in

689  ILLUMINA sequencing data and removes all reads containing insertions or deletions or other

690  bases than the expected reference and alternative allele at the selected sites. Therefore, only

691  the mismatch rate is relevant for the haplotype counting. Each position with a mismatch can

692  result in one of the other three bases different from the reference. If the probability of

693  mismatch to the reference is considered identical for each possible mismatched base, it is

694  equal to one third. As only one of these bases is accepted by the haplotype counting script,

695  the final probability of a mismatch at positions of reliable SNPs is equal to one-third of the

696  mismatch rate. Thus, the probability of error at selected SNP sites used for haplotype

697    determination ranged between 1.46% for *S. robusta* PacBio genome-wide sequencing and

698    0.05% for *P. tricornutum* T1 amplicon re-sequencing (Table S8).

699

700    <u>Resequencing of *S. robusta* and *P. tricornutum* loci with multiple haplotypes</u>

701    *DNA extraction:* DNA for deep sequencing was harvested and isolated by the CTAB DNA

702    extraction method. Cells from approximately 500 mL of exponentially growing *S. robusta* and

703    *P. tricornutum* cultures were harvested by centrifugation at 1216 x g for 5 min. The

704    supernatant was discarded and the cell pellet was resuspended in 400 µl of CTAB buffer (1%

705    (w/v) CTAB, 100 mM Tris-HCl pH 7.5, 10 mM EDTA pH 8, 700 mM NaCl and freshly added 4 µg

706    RNase A). *S. robusta* cells were disrupted by agitation with glass/zirconium beads (0.1-mm

707    diameter; Biospec) on a bead mill (Retsch) for three times 1 min at frequency 20 Hz. Samples

708    were incubated for 30 min at 60°C and afterwards let to cool down on the ice for 15 min. Next,

709    250 µl of chloroform:isoamylalcohol 24:1 was added and the samples were mixed manually

710    for 1 min. Phases were separated by centrifugation at 20 000 x g for 10 min. The upper

711    aqueous phase was transferred to a new tube and DNA was precipitated by the addition of an

712    equal volume of isopropanol followed by centrifugation for 15 min at 20 000 x g. The DNA

713    pellet was washed with 70% ethanol, air-dried and resuspended in 50 µl of 10 mM Tris-HCl

714    pH 8.5.

715    *Emulsion PCR:* Loci for re-sequencing of haplotypes were selected from the list of loci with

716    multiple haplotypes obtained through genome-wide haplotype detection. Three loci were

717    selected in *S. robusta* for Sanger sequencing verification (Table 2) and 62 loci for PacBio

718    amplicon sequencing verification in the case of *P. tricornutum* (Table S3). Primers for

719    amplification were designed manually (Data S3). To avoid PCR recombination artefacts [96,97],

720    selected loci were amplified by emulsion PCR using the MICELLULA DNA Emulsion &

Purification Kit (roboklon, Germany) according to the manufacturer's instructions. The DNA

concentration was measured on NanoDrop (ThermoFisher Scientific) and the number of DNA

template copies per µg of DNA was calculated according to the genome size of the respective

diatom. A maximum of $10^7$ DNA molecules was used per single emulsion PCR reaction. The

PCR reaction mix consisted of 1x OptiTaq PCR buffer B (roboklon), 200 µM dNTP mix, 2 µM of

each forward and reverse primer, DNA template with $10^6$-$10^7$ molecules, 1 mg/mL acetylated

BSA and 2.5U Opti Taq DNA polymerase (roboklon) in 50 µL of total volume. The emulsion mix

was prepared separately by mixing 220 µL of emulsion component 1, 20 µL of emulsion

component 2 and 60 µL of emulsion component 3 per PCR reaction. The 50 µL PCR reaction

was mixed with 300 µL of emulsion mix and emulsion was created by continuous vortexing at

1400 rpm at 4°C for 5 min. Each emulsion PCR reaction was split into three PCR tubes and run

with the following parameters: 94°C initial denaturation for 2 min, 26 cycles of 94°C

denaturation for 15 s, 56°C annealing for 30 s and 72°C extension with 1 kb/min relative to the

amplified fragment length, followed by a final extension at 72 °C for 10 min. The emulsion was

broken by the addition of 1 mL of isobutanol and vortexing. Next, 400 µL of Orange-DX solution

was added and reactions were gently mixed and centrifuged for 2 min at 20 000 x g. The

organic phase was removed and the aqueous phase was transferred to a Micellula spin column

activated by 40 µL of DX buffer. Columns were centrifuged at 11 000 x g for 1 min, washed

first with 500 µL of Wash-DX1 buffer, and then with 650 µL of Wash-DX2 buffer and the

leftovers of buffer were removed by an additional centrifugation for 2 min. PCR products were

eluted in 50 µL of Elution-DX buffer (all components: roboklon).


Sanger sequencing of *S. robusta* amplicons

744   *S. robusta* emulsion PCR products were cloned into the pGEM-T vector (Promega) according

745   to the manufacturer's instructions. In brief, the A overhangs were added by incubation of PCR

746   product with 10 µM dNTP mix and 1U of Taq DNA polymerase (Invitrogen) at 72°C for 10 min.

747   Then, 3.5 µL of PCR product was mixed with 5 µL of 2x ligase buffer, 0.75 µL of pGEM-T vector

748   and 2.25 U of T4 DNA ligase (all Promega) and incubated for 12 h at 4°C. Ligation mixtures

749   were transformed through electroporation into *E. coli* DH5alpha cells and transformants were

750   selected on LB supplemented with 100 µg/mL ampicillin (Duchefa). Clones containing cloned

751   PCR products were selected by Sanger sequencing with pGEM-5 and pGEM-6 primers (Data

752   S3). Sequencing results were aligned with the reference using Clustal Omega [98], and

753   haplotypes were manually assembled for each clone. Two alleles of Sro_contig556:54453-

754   55487 (Data S3) were cloned into the pGEM-T vector and an equimolar mix of these two

755   plasmids was used for emulsion PCR as a control for artefact generation. To simulate

756   conditions similar to emulsion PCR reactions on *S. robusta* genomic DNA, $10^6$-$10^7$ molecules

757   of *S. robusta* genomic DNA were added and the control samples were amplified with pGEM-3

758   and pGEM-4 primers (Data S3). The PCR products were again cloned into the pGEM-T vector

759   and sequenced with pGEM-5 and pGEM-6 primers. For counting the number of found

760   haplotypes, we considered only SNPs that were found in SNP call (so only the reference and

761   alternative allele at a selected position). Therefore, if the base at the SNP position was neither

762   reference nor an alternative base identified in the SNP call, the read was discarded. We

763   considered the probability of mismatch to the reference identical for each possible

764   mismatched base, and therefore one third. With 99.99% accuracy of the Sanger sequencing

765   (1 error per 10,000 sequenced base pairs) and the probability of mismatch turning the base

766   to the wrong reference/alternative allele being one third, the final probability of error at

767   selected sites is 1: 30,000.

768

769     *P. tricornutum* cultures and PacBio amplicon sequencing and haplotype counting

770     13 intergenic loci and 59 loci overlapping with coding regions (Table S3) were selected based

771     on an SNP call on T1 cell culture and the list of loci with more than two haplotypes for PacBio

772     Sequel amplicon sequencing. All loci were amplified by emulsion PCR as described above with

773     primers listed in Data S3. In the case of low amplification efficiency, the emulsion PCR was

774     repeated. Purified PCR products were concentrated using Genomic DNA Clean & Concentrator

775     (Zymo Research) according to the manufacturer's instructions. Amplifications of *CFP*, *GFP* and

776     *YFP* genes were used for control reactions for random mistakes and control reactions for

777     artificial haplotypes detection. The *CFP*, *GFP* and *YFP* sequences (Table S2 and Data S3) were

778     amplified by emulsion PCR with primers binding to vector backbone (Data S3) from GK-333-

779     CFP, GK-359 and GK-333-YFP plasmids respectively (GK-333, GenBank[72] accession MW934548

780     and GK-359, GenBank accession MW934549 were a gift from Dr. Nicole Poulsen, Center for

781     Molecular Bioengineering at TU Dresden). Control reactions for random errors consisted of

782     separately amplified GFP and YFP and control reactions for PCR-mediated recombination

783     consisted of mixed amplification of CFP+GFP and CFP+YFP (Table S2). Amplicons were pooled

784     together into two samples. Sample 1 contained 63 *P. tricornutum* endogenous 63 amplicons

785     from DNA harvested at a T6 time point, YFP amplified separately and CFP+GFP amplified in

786     one reaction. Sample 2 contained 5 *P. tricornutum* endogenous 5 amplicons from DNA

787     harvested at a T1 time point, GFP amplified separately and CFP+YFP amplified in one reaction.

788     Samples were barcoded, mixed in 9:1 ratio and sequenced on 1 PacBio Sequel SMRT cell at

789     Novogene (UK).

790     Circular Consensus Sequence (CCS) reads were obtained with SMRT Link 8.0.0 software

791     (PacBio) with --min-length 500 --max-length 2400 --min-passes 4 --by-strand and mapped to

792 the reference loci with BLASR[89] with default settings. The CIGAR string was corrected with

793 samfixcigar, soft-clipped bases were removed with biostar84452 from jvarkit[84] and reads with

794 mapping quality >20 were selected using SAMtools [77]. Next, the haplotype number per locus

795 was counted as described above. Only haplotypes supported either by at least 1% of valid

796 reads or at least two reads if read count was lower than 200 were selected (Table S3). The

797 genes fully covered by PacBio amplicons were manually annotated to detect putative variants.

798

799 <u>LOH and CNV detection</u>

800 *Cultures started from a single cell*: To isolate single cells from *P. tricornutum* cultures, an

801 aliquot from the respective culture was diluted to $10^6$, $10^9$ and $10^{12}$ in the growth medium.

802 200 µl of diluted culture were spread on a 120 x 120 mm Petri dish (Corning Gosselin) with

803 solid medium prepared with 17.25 g/L of Tropic Marin Bio-Actif sea salt solid and 10 g/L of

804 Plant Tissue Culture Agar (Neogen) supplemented with 1x Guillard's (F/2) Marine Water

805 Enrichment Solution (Sigma-Aldrich), 100 µg/L ampicillin, 50 µg/L gentamycin and 100 µg/L

806 streptomycin, and the plates were incubated at room temperature with a 12-h/12-h light/dark

807 cycle. Plates were checked for the presence of colonies after 14 days and single colonies were

808 transferred from the plate with the highest dilution factor that contained colonies into liquid

809 TMB medium using a pipette tip. This procedure was used to isolate and start three colonies

810 from a single cell from the Pt1 culture to obtain mother cultures MC1, MC2 and MC3. Thirty

811 days after mother culture isolation (T1 time point), daughter cells DC1.1-DC3.3 were isolated

812 from the respective mother cultures (Figure S5). Cultures for deep sequencing were harvested

813 twice; at time point T1 part of the mother cultures was harvested at time point T1, and 30

814 days later at time point T2 all cultures in the experiment were harvested.

815 _Estimation of cell division rate:_ To pre-test the number of cell divisions per 6 days in our culture

816 conditions, 1 x $10^6$ cells were inoculated on day 1 in 6 replicas and the number of cells was

817 counted after 6 full days). The average number of cells after 6 days was 8.03 x $10^6$. Using the

818 following formula: $N_t = N_0 2^{tf}$ , where N(t) is the number of cells at time t, $N_0$ is the initial

819 number of cells, t is the time in days and f is the frequency of cell cycles per unit time, the cell

820 division rate was estimated to be around 0.501 cell division per day. The actual rate of cell

821 division can be higher if cell mortality is counted.

822 _Illumina sequencing, SNP calling, LOH and CNV detection:_ DNA for deep sequencing was

823 harvested and isolated by the CTAB DNA extraction method as described above. Paired-end

824 libraries were prepared with the NEBNext Ultra DNA Library Prep Kit for Illumina (NEB) with a

825 500-bp insert size and sequencing was performed on a 2× 150bp Illumina NextSeq500 Medium

826 at the VIB Nucleomics Core (Leuven, Belgium). Adapters and reads with a quality score below

827 20 were removed using BBduk2[75] with minlen=35 qtrim=rl trimq=20 hdist=1 tbo tpe and

828 custom adapter reference file. Trimmed reads were aligned to _P. tricornutum_ v2 assembly

829 GCA_000150955.2[91] (European Nucleotide Archive) using Burrows-Wheeler Alignment Tool

830 (BWA)[76] algorithm BWA-MEM and processed and re-calibrated as described above. After the

831 second round of recalibration, SNPs were called in three different ways:

832 _Germline SNP calling with GATK HaplotypeCaller  a) joint genotyping with GATK 4.2.1 and b)_

833 _for individual samples with GATK 3.7.0_: First, germline SNPs were called with HaplotypeCaller

834 with either a) -ERC GVCF option and gVCF files were combined with CombineGVCFs and then

835 jointly genotyped with GenotypeGVCFs b) or without the -ERC GVCF option and joint

836 genotyping. SNPs were filtered with QD < 2.0,  QUAL < 30.0,  SOR > 3.0,  FS > 60.2,  MQ < 40.0,

837 MQRankSum < -12.5, ReadPosRankSum < -8.0 and DP < 10 and indels with QD < 2.0, QUAL <

838 30.0, FS > 200.0, ReadPosRankSum < -20.0, DP < 10 using VariantFiltration and filtered SNPs

839 and indels were removed with SelectVariants. These set of germline SNPs were used to build

840 a Panel of Normals for following pairwise comparison of mother and daughter cell cultures

841 using Mutect2 and for cross-verification of LOH regions identified in Mutect2.

842 *Pairwise comparison of mother and daughter cultures with GATK Mutect2:* All GATK algorithms

843 were version 4.2.1 if not stated otherwise. First, SNPs were called on each sample with

844 Mutect2 in tumour-only mode. Next, a vcf file for each sample was created by moving the

845 sample-level AF allele-fraction annotation were moved into the INFO field for each sample

846 using VariantsToTable. The Panel of Normals was prepared using the germline SNPs dataset

847 for individual samples generated by HaplotypeCaller as a germline-resource by first calling the

848 SNPs for each sample with Mutect2 in tumour-only mode, then merging all files with

849 CombineVariants (GATK3.7.0) and finally creating the Panel of Normals with

850 CreateSomaticPanelOfNormals. A Panel of Normals file for each pairwise comparison was

851 prepared by masking germline variants from the respective individual samples using germline

852 SNPs called by HaplotypeCaller using CatVariants (GATK3.7.0) to merge SNPs from both

853 samples and then masking them using SelectVariants with -XL option. Finally, LOH regions and

854 *de novo* mutations were detected by Mutect2 SNP call with vcf file with sample-level AF allele-

855 fraction in the INFO field used as --germline-resource, the respective masked Panel of Normals

856 file used as –pon, --genotype-germline-sites true and either the mother culture used as

857 tumour sample and daughter culture used as normal to detect LOH events in the daughter

858 culture or vice-versa to detect *de novo* mutations in daughter culture (Table S5 and Data S3).

859 Called SNPs were filtered with FilterMutectCalls and filtered SNPs were removed with

860 SelectVariants –excludeFiltered. As a control, the T1 with T2 time point of each mother culture

861 were compared. A minimum of three consecutive SNPs missing in the daughter or T2 mother

862 culture was considered as a LOH event and all LOH regions were reexamined in the datasets

863  of germline SNPs obtained either by individual SNP calling or joint genotyping. The nature of

864  the LOH was judged by comparison of coverage of the LOH region and its surrounding

865  heterozygous borders and through Sanger resequencing of the LOH border (Tables S5 and

866  Data S3).

867  *CNV detection*: The copy-number variation was detected using GATK 4.1.7 CNV detection

868  pipeline. First, intervals list with bin length set to 100 bp and -interval-merging-rule

869  OVERLAPPING_ONLY was prepared with PreprocessIntervals and collect raw counts were

870  collected using CollectReadCounts and CNV panel of normals was generated by

871  CreateReadCountPanelOfNormal with --minimum-interval-median-percentile 5.0 setting.

872  Standardized copy ratios and denoised copy ratios were obtained using DenoiseReadCounts

873  against a panel of normals. Reference and alternative allele counts at common germline sites

874  called on all samples by HaplotypeCaller in GVCF mode and jointly genotyped were obtained

875  using CollectAllelicCounts for each sample. Segments of contiguous copy ratios were acquired

876  by ModelSegments in a paired analysis with --denoised-copy-ratios and --allelic-counts from

877  the daughter culture and --normal-allelic-counts from the respective mother culture.

878  Amplified, deleted and copy-neutral segments were called with CallCopyRatioSegments with

879  default settings and plotted using PlotModeledSegments. Detected CNV events were cross-

880  verified in the datasets of germline SNPs obtained by SNP calling and joint genotyping and in

881  estimated coverage counts per 10 bp obtained using bedtools 2.2.28 coverage. Further,

882  identified LOH in tandem repeat on chromosome 5 in DC1.2 was verified by Sanger sequencing

883  and duplication on chromosome 23 in DC1.3 was confirmed by qPCR quantification (see

884  below).

885  *Re-sequencing of identified LOH events in mother versus daughter cell culture comparison and*

886  *CNV analysis:* The nature of LOH events identified by pairwise comparison between mother

887    and daughter cell culture was first judged by visual comparison of the sample coverage at the

888    LOH region and the neighbouring region. Next, the region was amplified by emulsion PCR as

889    described in the Methods, PCR products were cloned into the pGEM-T vector and individual

890    clones were sequenced. If the LOH region sequence was found together with both alleles of

891    neighbouring heterozygous SNP(s), the region was considered as copy-neutral LOH. If only a

892    single allele was found and the coverage data corresponded,  the region was considered as a

893    deletion. Predicted heterozygous SNPs in regions between identified LOH events on

894    chromosome 26 in DC2.1 culture were amplified by a standard PCR using Phusion High-Fidelity

895    DNA Polymerases (Thermo Scientific) according to the manufacturers' instructions and

896    sequenced by Sanger sequencing to verify their heterozygosity (Data S3). The same approach

897    was used to verify the hetero/homozygosity of three SNPs on chromosome 27 that were

898    predicted as LOH in culture DC2.2, but was called as homozygous by independent SNP calling

899    also in the mother culture MC2. The LOH on chromosome 5 in DC1.2 was identified as CNV,

900    but PCR amplification confirmed the presence of two alleles with different length (differing in

901    1960 bp). The longer allele contained a tandemly duplicated region while in the short allele

902    the duplication was missing. Subsequent Sanger sequencing of alleles showed that culture

903    DC1.2 contains two short alleles with LOH of SNPs in the surrounding region (Figure S5 and

904    Data S3). The effect of SNPs found in was profiled using SnpEff[39]. First, pre-build SnpEff *P.*

905    *tricornutum* database was downloaded and records for SNPs in LOH regions were selected

906    using bedtools intersect[83].  The selected SNP variants were annotated using snpEff with

907    default settings. SNPs annotated with "HIGH" effect were selected using grep and their exact

908    effect was manually annotated.

909    *qPCR quantification of duplication at chromosome 23 in DC1.4 culture:*

The relative copy number variation on chromosome 23 was examined in daughter cultures DC1.1, DC1.2 and DC1.3 in comparison with mother culture MC1 by quantitative real-time PCR (qPCR). DNA was extracted as described above and concentration was adjusted to 48 pg/μL for each sample. Two primer pairs were designed into the region containing putative duplication on chromosome 23 in DC1.3, two primer pairs were located on chromosome 23 outside the duplicated region and two primers pairs were targeting loci outside of chromosome 23, one on chromosome 6 and one on chromosome 22 (Data S3). qPCR was performed using the SYBR Green kit (Roche) with 100 nM primers and 0.125 μL DNA in a total volume of 5 μL per reaction. qPCR amplification reactions were run and analyzed on the LightCycler 480 (Roche) with following cycling conditions: 10 min polymerase pre-incubation at 95°C and 45 cycles of amplification at 95°C for 10 s, 60°C for 15 s, and 72°C for 15 s. Melting curves were recorder after the last cycle by heating from 65 to 95°C. All qPCR amplicons were sequenced by Sanger sequencing to confirm their integrity. For each reaction, three technical repeats were performed. Data were analyzed using qbase+ (Biogazelle)[99] with a copy number analysis option and with the mother culture MC1 as a reference sample and loci on the distal arm of chromosome 23 (locus D) chromosome 6 (locus E) and chromosome 22 (locus F) as reference targets.

PtUMPS read-out system

*PtUMPS cultures*: Strains *ptumps-1bp* and *ptumps-1368bp* were generated based on a previously described protocol[41]. Briefly, Cas9 ribonucleoprotein (RNP) complexes were assembled to target the *PtUMPS* locus, at either the gUMPS1 site or the gUMPS4 site (Data S3). An equimolar mixture of RNP gUMPS1 and RNP gUMPS4 (4 μg each) was bombarded into wild-type *P. tricornutum* Pt1 8.6 (CCMP2561) cells. Two rounds of selection were made on

934    silicate-free F/2 medium (Sigma) plates supplemented with 50 µg/mL uracil (Sigma) and 300

935    µg/mL 5-fluoroorotic acid (5-FOA; ThermoFisher). Cell lysates were then prepared to serve as

936    a template for genotyping. PCRs using the Q5 High Fidelity DNA polymerase (New England

937    Biolabs) and primers UMPS_5UTR_F and UMPS_3UTR_R (Data S3) were performed to amplify

938    the *PtUMPS* locus. The generated amplicons were subcloned employing the CloneJET PCR

939    cloning kit (Thermo Scientific) and analyzed through Sanger sequencing. The *ptumps-320bp*

940    strain was prepared and described previously under the name UA17[41]. *PtUMPS* mutant strains

941    were maintained in conditions described above in 1xTMB medium supplemented with

942    50 µg/mL uracil and 100 µg/mL 5-FOA to prevent the restoration of the wild-type allele.

943    *PtUMPS 14 day cultivation in non-selective conditions:* Cells densities of *ptumps-1bp*, *ptumps-*

944    *320bp* and *ptumps-1368bp* were estimated using Bürker counting chamber and $20 \times 10^6$ cells

945    from were harvested by centrifugation at 1216 x g for 5 minutes. The cell pellet was washed

946    four times with 50 mL of 1xTMB medium, then resuspended in 750 mL of 1x TMB medium

947    supplemented with 50 µg/mL uracil and resulting cultures were cultivated in 12h/12h

948    light/dark cycle. After 7 days, another 750 mL of fresh medium supplemented with uracil were

949    added. After 14 days, cultures were harvested by centrifugation, cell density was estimated

950    and $50 \times 10^6$ cells from the culture were washed four times with TMB medium and plated on

951    1% agar ½ TMB medium 245 x 245 mm Nunc™ Square BioAssay plates (ThermoFisher) and

952    incubated in 18h/6h light/dark cycle at 20°C for 6 weeks. The resulting colonies were manually

953    counted (Data S4 and S5). Colonies selected for sequencing of *PtUMPS* locus were transferred

954    to fresh 1xTMB medium and grown cell cultures were harvested as described above. The

955    *PtUMPS* locus was amplified with primers PtUMPS-1 and PtUMPS-2  (Data S3) using two

956    consecutive rounds of emulsion PCR. PCR products were cloned to pGEM-T vector and

957    sequenced by Sanger sequencing.

958

959 <u>Estimation of interhomolog recombination frequency:</u>

960 ~50x10$^6$ cells of three *ptumps-1bp* and five independent *ptumps-1368bp* cell subcultures

961 started from a single cell were harvested by centrifugation at 1216 x g for 5 minutes. The cell

962 pellet was washed four times with 50 mL of 1xTMB medium, then resuspended in 1xTMB

963 medium. The cell density was determined and 2 to 6 replicas of either 25x10$^6$ or 20x10$^6$ cells

964 were immediately plated on 1% agar ½ TMB medium and incubated as described above and

965 incubated in 18h/6h light/dark cycle at 20°C for 6 weeks. The resulting colonies were manually

966 counted (Data S4 and S5).

967 The frequency if interhomolog recombination was calculated based on the number of uracil

968 prototrophic colonies after the immediate transfer of *ptumps-1368bp* strain from 5-FOA- and

969 uracil-supplemented medium (only mutant cells survive) on plates without uracil (only cells

970 that restored the wild-type allele survive).  As the recombination that restored the wild-type

971 *PtUMPS* allele had to occur within 1368-bp region, first, the incidence of LOH events per 1000

972 bp was estimated for each replica by dividing the number of colonies by 1.38. The exact

973 nuclear genome size of *P. tricornutum* was determined from the genome fasta file using awk

974 as 27,450,724 bp. Taking into account the genome length and number of plated cells, the

975 recombination rate was recalculated per whole genome of 100 cells. As the copy-neutral LOH

976 is detectable in half of the cases of mitotic interhomolog recombination due to random

977 segregation of sister chromatids during mitotic metaphase, the number was multiplied by 2

978 to obtain the rate of mitotic recombination. The average value obtained from data for all

979 replicas was 4.2 x 10$^{-2}$. The rate of reciprocal cross-overs on a 120-kb region in *S. cerevisiae*

980 was estimated at 4 x 10$^{-5}$ per cell division and the rate of gene conversion events as 3.5 × 10$^{-6}$

981 per cell division [42] or 3.3 x 10$^{-3}$ per cell division for interstitial LOH, 1.4 x 10$^{-3}$ for terminal LOH

982　genome-wide [43]. The frequency of reciprocal cross-overs and gene conversion on a 120-kb

983　region were combined and recalculated first per 1 kb and subsequently per 11.89 Mb *S.*

984　*cerevisiae* genome. The resulting rate of interhomolog recombination in *S. cerevisiae* was

985　calculated as 4.3 x $10^{-3}$ in the case of 120 kb region or 4.7 x $10^{-3}$ in the case of the genome-

986　wide studies.

987

988　<u>Effect of treatment with cadmium, $H_2O_2$ and (E,E)-2,4-decadienal on interhomolog</u>

989　<u>recombination at PtUMPS locus</u>

990　The ranges of concentrations of chemicals used for treatments were first surveyed in

991　literature, then selected concentrations were tested by treatment of Pt1 *P. tricornutum* strain

992　for seven days. Afterwards, the cell survival of mock and treated cells was compared and the

993　maximal dose that did not cause a decrease in cell density was selected as the maximal dose

994　in the respective experiment. *ptumps-1bp* and *ptumps-1368bp* cell cultures were harvested

995　by centrifugation at 1216 x g  for 5 minutes and the cell pellet was washed four times with 50

996　mL of 1xTMB medium, then resuspended in 50 mL of 1xTMB medium. Cell density per mL was

997　estimated and $25 \times 10^6$ cells per replica were transferred to 200 mL of 1xTMB medium

998　supplemented with 50 µg/mL uracil and respective treatment or mock treatment. For zeocin

999　treatment, zeocin was added to the final concentration (InvivoGen) was added to final

1000　concentrations of 1 µg/mL and 10 µg/mL from 1000× stock solution. For cadmium treatment,

1001　CdCl₂ (Sigma-Aldrich) was added to final $Cd^{2+}$ concentrations of 5 µg/L and 50 µg/L from 4000×

1002　stock solution. For $H_2O_2$ treatment, a 30% solution of $H_2O_2$ (Merck) was added to a final

1003　concentration of 5 µM or 50 µM $H_2O_2$. For (E,E)-2,4-Decadienal treatment, 200 µL of DMSO

1004　was added to the mock-treated cell cultures and (E,E)-2,4-Decadienal (Sigma-Aldrich) was

1005　added to 0.1 µM and 1 µM final concentration from 1000x and 100x concentrated stock

1006      solution respectively. Cultures were incubated for 24h in 12h/12h light/dark cycle at 20°C,

1007      under photosynthetic LED light with an intensity of 160 µmol photons m$^{-2}$ s$^{-1}$ and with 100 rpm

1008      shaking. Afterwards, all samples were harvested by centrifugation at 1216 x g  for 5 minutes

1009      and cell pellets were washed four times with 50 mL of 1xTMB medium and plated on 1% agar

1010      ½ TMB medium as described above and incubated in 18h/6h light/dark cycle at 20°C for 6

1011      weeks. The resulting colonies were manually counted (Data S4 and S5).

1012

1013      **Quantification and statistical analysis**

1014      No sample-size calculations were performed. Sample sizes were determined to be adequate

1015      based on preliminary experiments and feasibility. The interhomolog mitotic recombination

1016      rate measured in three ptumps-1bp and five ptumps-1368bp independent biological

1017      replicates, each with two to six technical replicas. The influence of environmental stresses on

1018      interhomolog mitotic recombination rate was performed in three biological replicas. The

1019      number of biological replicates for each data panel is indicated in the figure panel, in

1020      Supplementary data and the source data files. No randomization was performed (not

1021      applicable). Data exclusions: A threshold was set to count haplotypes in long-read sequencing

1022      for genome-wide and re-sequencing experiments to remove false-positive data as described

1023      above. No other data were excluded from this study.

1024

1025      **DATA SX TITLES AND LEGENDS**:

1026

1027      **Data S1. Details of the SNV positions and densities in the selected genes in *F. cylindrus***

1028      **metagenomes. Related to STAR Methods and Figure S1.**

A) SNV positions of the exons present in the core genes, for the surface metagenome (10% divergence from consensus, at least 20X coverage in both samples). B) SNV positions of the exons present in the core genes, for the DCM metagenome (10% divergence from consensus, at least 20X coverage in both samples). C) SNV densities (%) of the exons present in the core genes. D) SNV densities (%) of the core genes (filtered by mean coverage).

**Data S2. List of loci with more than two haplotypes detected genome-wide in profiled organisms. Related to STAR Methods and Figures 1 and S2-S3.**

A) List of loci with more than two haplotypes detected in genome-wide analysis in *S. robusta*. B) List of loci with more than two haplotypes detected in genome-wide analysis in *P. tricornutum*. C) List of loci with more than two haplotypes detected in genome-wide analysis in *A. thaliana*. and D) List of loci with more than one haplotype detected in the genome-wide haplotype analysis in *S. cerevisiae*. E) Overview of genome-wide haplotype counting in diatoms *S. robusta* and *P. tricornutum*, yeast *S. cerevisiae* and plant *A. thaliana*. F) Error rate in processed long-read sequencing data used for haplotype counting estimated by Alfred.

**Data S3. Details of re-sequencing of loci with multiple haplotypes, analysis of genetic changes in daughter cells and materials used throughout the study. Related to STAR Methods and Figures 1-4 and S2-6.**

A) List of *S. robusta* loci used for Sanger sequencing verification of haplotype diversity. B) List of *P. tricornutum* 2-kb loci used for PacBio amplicon sequencing verification of haplotype diversity. C) The number of haplotypes at control amplicons obtained by PacBio amplicon sequencing. Only haplotypes supported by at least 1% of reads from the total read count or at least 2 reads were added to the haplotype count.  D) The number of haplotypes of *P.*

*tricornutum* loci used for the verification of haplotype diversity by PacBio amplicon sequencing, six months after single-cell isolation. Only haplotypes supported by at least 1% of reads from the total read count or at least 2 reads were added to the haplotype count. E) Position of alleles in genes used as for a control amplification in emulsion PCR experiments. F) De novo SNPs identified by pairwise comparison of mother and daughter cultures in the genome-wide LOH analysis. G) SNPs with high effect on protein function fixed by LOHs and deletions. H) Genes covered by duplication on chromosome 23 in DC1.3. I) List of diatom datasets and strains. J) List of oligonucleotides used throughout the study.

**Data S4. Photographs of plates with uracil prototrophic colonies resulting from mitotic interhomolog recombination. Related to STAR Methods and Figures 4 and S6.**

**Data S5. Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies. Related to STAR Methods and Figures 4 and S6.**

A) Counts of *ptumps-1bp*, *ptumps-320bp* and *ptumps-1368bp* uracil prototrophic colonies after a 14-day cultivation in non-selective conditions. B) Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies after immediate transfer from 5-FOA- and uracil-supplemented medium to medium without uracil. C) Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies after 24-h zeocin treatment in non-selective conditions. D) Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies after 24-h $H_2O_2$ treatment in non-selective conditions. E) Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies after 24-h cadmium treatment in non-selective conditions. F) Counts of *ptumps-1bp* and *ptumps-1368bp* uracil prototrophic colonies after 24-h (E,E)-2,4-Decadienal treatment in non-selective conditions.

1077

**REFERENCES**

1079 1. Malviya, S., Scalco, E., Audic, S., Vincenta, F., Veluchamy, A., Poulain, J., Wincker, P., Iudicone,

1080 D., de Vargas, C., Bittner, L., et al. (2016). Insights into global diatom distribution and

1081 diversity in the world's ocean. P Natl Acad Sci USA *113*, E1516-E1525.

1082 10.1073/pnas.1509523113.

1083 2. Godhe, A., and Rynearson, T. (2017). The role of intraspecific variation in the ecological and

1084 evolutionary success of diatoms in changing environments. Philos T R Soc B *372*. ARTN

1085 20160399 10.1098/rstb.2016.0399.

1086 3. Chepurnov, V.A., Mann, D.G., von Dassow, P., Vanormelingen, P., Gillard, J., Inze, D., Sabbe,

1087 K., and Vyverman, W. (2008). In search of new tractable diatoms for experimental biology.

1088 Bioessays *30*, 692-702. 10.1002/bies.20773.

1089 4. Koester, J.A., Berthiaume, C.T., Hiranuma, N., Parker, M.S., Iverson, V., Morales, R., Ruzzo,

1090 W.L., and Armbrust, E.V. (2018). Sexual ancestors generated an obligate asexual and globally

1091 dispersed clone within the model diatom species Thalassiosira pseudonana. Sci Rep-Uk *8*.

1092 ARTN 10492 10.1038/s41598-018-28630-4.

1093 5. Lewis, W.M. (1984). The Diatom Sex Clock and Its Evolutionary Significance. American

1094 Naturalist *123*, 73-80. Doi 10.1086/284187.

1095 6. Davidovich, N.A., Davidovich, O.I., Podunay, Y.A., Gastineau, R., Kaczmarska, I., Poulickova,

1096 A., and Witkowski, A. (2017). Ardissonea crystallina has a type of sexual reproduction that is

1097 unusual for centric diatoms. Scientific Reports *7*. ARTN 14670 10.1038/s41598-017-15301-z.

1098 7. Jewson, D.H. (1992). Size-Reduction, Reproductive Strategy and the Life-Cycle of a Centric

1099 Diatom. Philosophical Transactions of the Royal Society of London Series B-Biological

1100 Sciences *336*, 191-213. DOI 10.1098/rstb.1992.0056.

1101    8.    Fuchs, N., Scalco, E., Kooistra, W.H.C.F., Assmy, P., and Montresor, M. (2013). Genetic

1102          characterization and life cycle of the diatom Fragilariopsis kerguelensis. European Journal of

1103          Phycology *48*, 411-426. 10.1080/09670262.2013.849360.

1104    9.    Lynch, M., Ackerman, M.S., Gout, J.F., Long, H., Sung, W., Thomas, W.K., and Foster, P.L.

1105          (2016). Genetic drift, selection and the evolution of the mutation rate. Nat Rev Genet *17*,

1106          704-714. 10.1038/nrg.2016.104.

1107    10.   Hedrick, P.W. (2010). Genetics of populations.

1108    11.   Krasovec, M., Sanchez-Brosseau, S., and Piganeau, G. (2019). First Estimation of the

1109          Spontaneous Mutation Rate in Diatoms. Genome Biol Evol *11*, 1829-1837.

1110          10.1093/gbe/evz130.

1111    12.   Krasovec, M., Rickaby, R.E.M., and Filatov, D.A. (2020). Evolution of Mutation Rate in

1112          Astronomically Large Phytoplankton Populations. Genome Biol Evol *12*, 1051-1059.

1113          10.1093/gbe/evaa131.

1114    13.   Bürger, R. (2000). The mathematical theory of selection, recombination, and mutation

1115          (Wiley).

1116    14.   Ewens Warren, J. (2004). Mathematical population genetics [Texte imprimé] . I, Theoretical

1117          introduction / Warren J. Ewens, Second edition Edition (Springer).

1118    15.   Lang, G.I., Rice, D.P., Hickman, M.J., Sodergren, E., Weinstock, G.M., Botstein, D., and Desai,

1119          M.M. (2013). Pervasive genetic hitchhiking and clonal interference in forty evolving yeast

1120          populations. Nature *500*, 571-574. 10.1038/nature12344.

1121    16.   Maddamsetti, R., Lenski, R.E., and Barrick, J.E. (2015). Adaptation, Clonal Interference, and

1122          Frequency-Dependent Interactions in a Long-Term Evolution Experiment with Escherichia

1123          coli. Genetics *200*, 619-631. 10.1534/genetics.115.176677.

1124    17.   Fisher, R.A. (1930). The genetical theory of natural selection (Clarendon Press).

1125          10.5962/bhl.title.27468.

1126  18.  Muller, H.J. (1932). Some Genetic Aspects of Sex. The American Naturalist *66*, 118-138.

1127      10.1086/280418.

1128  19.  Johnson, R.D., and Jasin, M. (2001). Double-strand-break-induced homologous

1129      recombination in mammalian cells. Biochemical Society Transactions *29*, 196-201. Doi

1130      10.1042/Bst0290196.

1131  20.  Kadyk, L.C., and Hartwell, L.H. (1992). Sister Chromatids Are Preferred over Homologs as

1132      Substrates for Recombinational Repair in Saccharomyces-Cerevisiae. Genetics *132*, 387-402.

1133  21.  Aguilera, A., Chavez, S., and Malagon, F. (2000). Mitotic recombination in yeast: elements

1134      controlling its incidence. Yeast *16*, 731-754. Doi 10.1002/1097-

1135      0061(20000615)16:8<731::Aid-Yea586>3.0.Co;2-L.

1136  22.  James, T.Y., Michelotti, L.A., Glasco, A.D., Clemons, R.A., Powers, R.A., James, E.S., Simmons,

1137      D.R., Bei, F.Y., and Ge, S.H. (2019). Adaptation by Loss of Heterozygosity in Saccharomyces

1138      cerevisiae Clones Under Divergent Selection. Genetics *213*, 665-683.

1139      10.1534/genetics.119.302411.

1140  23.  Schoustra, S.E., Debets, A.J.M., Slakhorst, M., and Hoekstra, R.F. (2007). Mitotic

1141      recombination accelerates adaptation in the fungus Aspergillus nidulans. Plos Genet *3*. ARTN

1142      e68 10.1371/journal.pgen.0030068.

1143  24.  Dale, A.L., Feau, N., Everhart, S.E., Dhillon, B., Wong, B., Sheppard, J., Bilodeau, G.J., Brar, A.,

1144      Tabima, J.F., Shen, D., et al. (2019). Mitotic Recombination and Rapid Genome Evolution in

1145      the Invasive Forest Pathogen Phytophthora ramorum. Mbio *10*. ARTN e02452-18

1146      10.1128/mBio.02452-18.

1147  25.  Mock, T., Otillar, R.P., Strauss, J., McMullan, M., Paajanen, P., Schmutz, J., Salamov, A.,

1148      Sanges, R., Toseland, A., Ward, B.J., et al. (2017). Evolutionary genomics of the cold-adapted

1149      diatom Fragilariopsis cylindrus. Nature *541*, 536-540. 10.1038/nature20803.

1150  26.  Madoui, M.A., Poulain, J., Sugier, K., Wessner, M., Noel, B., Berline, L., Labadie, K., Cornils, A.,

1151      Blanco-Bercial, L., Stemmann, L., et al. (2017). New insights into global biogeography,

1152    population structure and natural selection from the genome of the epipelagic copepod

1153    Oithona. Molecular Ecology *26*, 4467-4482. 10.1111/mec.14214.

1154  27.  Olm, M.R., Crits-Christoph, A., Bouma-Gregson, K., Firek, B.A., Morowitz, M.J., and Banfield,

1155    J.F. (2021). inStrain profiles population microdiversity from metagenomic data and

1156    sensitively detects shared microbial strains. Nat Biotechnol. 10.1038/s41587-020-00797-0.

1157  28.  Quince, C., Delmont, T.O., Raguideau, S., Alneberg, J., Darling, A.E., Collins, G., and Eren, A.M.

1158    (2017). DESMAN: a new tool for de novo extraction of strains from metagenomes. Genome

1159    Biol *18*. ARTN 181 10.1186/s13059-017-1309-9.

1160  29.  Russo, M.T., Cigliano, R.A., Sanseverino, W., and Ferrante, M.I. (2018). Assessment of

1161    genomic changes in a CRISPR/Cas9 Phaeodactylum tricornutum mutant through whole

1162    genome resequencing. Peerj *6*. ARTN e5507 10.7717/peerj.5507.

1163  30.  Koren, S., Walenz, B.P., Berlin, K., Miller, J.R., Bergman, N.H., and Phillippy, A.M. (2017).

1164    Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat

1165    separation. Genome Res *27*, 722-736. 10.1101/gr.215087.116.

1166  31.  McIlwain, S.J., Peris, D., Sardi, M., Moskvin, O.V., Zhan, F.J., Myers, K.S., Riley, N.M., Buzzell,

1167    A., Parreiras, L.S., Ong, I.M., et al. (2016). Genome Sequence and Analysis of a Stress-

1168    Tolerant, Wild-Derived Strain of Saccharomyces cerevisiae Used in Biofuels Research. G3-

1169    Genes Genomes Genetics *6*, 1757-1766. 10.1534/g3.116.029389.

1170  32.  Kim, K.E., Peluso, P., Babayan, P., Yeadon, P.J., Yu, C., Fisher, W.W., Chin, C.S., Rapicavoli,

1171    N.A., Rank, D.R., Li, J., et al. (2014). Long-read, whole-genome shotgun sequence data for five

1172    model organisms. Scientific Data *1*. ARTN 140045 10.1038/sdata.2014.45.

1173  33.  Zapata, L., Ding, J., Willing, E.M., Hartwig, B., Bezdan, D., Jiao, W.B., Patel, V., James, G.V.,

1174    Koornneef, M., Ossowski, S., and Schneeberger, K. (2016). Chromosome-level assembly of

1175    Arabidopsis thaliana Ler reveals the extent of translocation and inversion polymorphisms. P

1176    Natl Acad Sci USA *113*, E4052-E4060. 10.1073/pnas.1607532113.

1177    34.    Hiraoka, M., Watanabe, K., Umezu, K., and Maki, H. (2000). Spontaneous loss of

1178           heterozygosity in diploid Saccharomyces cerevisiae cells. Genetics *156*, 1531-1548.

1179    35.    Hunter, N. (2015). Meiotic Recombination: The Essence of Heredity. Cold Spring Harb

1180           Perspect Biol *7*. 10.1101/cshperspect.a016618.

1181    36.    Rastogi, A., Vieira, F.R.J., Deton-Cabanillas, A.F., Veluchamy, A., Cantrel, C., Wang, G.H.,

1182           Vanormelingen, P., Bowler, C., Piganeau, G., Hu, H.H., and Tirichine, L. (2020). A genomics

1183           approach reveals the global genetic polymorphism, structure, and functional diversity of ten

1184           accessions of the marine model diatom Phaeodactylum tricornutum. Isme Journal *14*, 347-

1185           363. 10.1038/s41396-019-0528-3.

1186    37.    National Library of Medicine (US) (2018). Sequence Read Archive (SRA) SRR7762337.

1187           Bethesda (MD): National Center for Biotechnology Information.

1188    38.    (US), N.L.o.M. (2018). Sequence Read Archive (SRA) SRR7762336. Bethesda (MD): National

1189           Center for Biotechnology Information.

1190    39.    Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X.Y., and

1191           Ruden, D.M. (2012). A program for annotating and predicting the effects of single nucleotide

1192           polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain w(1118); iso-

1193           2; iso-3. Fly *6*, 80-92. 10.4161/fly.19695.

1194    40.    Sakaguchi, T., Nakajima, K., and Matsuda, Y. (2011). Identification of the UMP Synthase Gene

1195           by Establishment of Uracil Auxotrophic Mutants and the Phenotypic Complementation

1196           System in the Marine Diatom Phaeodactylum tricornutum. Plant Physiol *156*, 78-89.

1197           10.1104/pp.110.169631.

1198    41.    Serif, M., Dubois, G., Finoux, A.L., Teste, M.A., Jallet, D., and Daboussi, F. (2018). One-step

1199           generation of multiple gene knock-outs in the diatom Phaeodactylum tricornutum by DNA-

1200           free genome editing. Nat Commun *9*. ARTN 3924 10.1038/s41467-018-06378-9.

1201    42.    Barbera, M.A., and Petes, T.D. (2006). Selection and analysis of spontaneous reciprocal

1202        mitotic cross-overs in Saccharomyces cerevisiae. P Natl Acad Sci USA *103*, 12819-12824.

1203        10.1073/pnas.0605778103.

1204    43.    Sui, Y., Qi, L., Wu, J.K., Wen, X.P., Tang, X.X., Ma, Z.J., Wu, X.C., Zhang, K., Kokoska, R.J.,

1205        Zheng, D.Q., and Petes, T.D. (2020). Genome-wide mapping of spontaneous genetic

1206        alterations in diploid yeast cells. Proc Natl Acad Sci U S A. 10.1073/pnas.2018633117.

1207    44.    Povirk, L.F. (1996). DNA damage and mutagenesis by radiomimetic DNA-cleaving agents:

1208        bleomycin, neocarzinostatin and other enediynes. Mutat Res *355*, 71-89. 10.1016/0027-

1209        5107(96)00023-1.

1210    45.    D'Autreaux, B., and Toledano, M.B. (2007). ROS as signalling molecules: mechanisms that

1211        generate specificity in ROS homeostasis. Nat Rev Mol Cell Bio *8*, 813-824. 10.1038/nrm2256.

1212    46.    Brembu, T., Jorstad, M., Winge, P., Valle, K.C., and Bones, A.M. (2011). Genome-Wide

1213        Profiling of Responses to Cadmium in the Diatom Phaeodactylum tricornutum. Environ Sci

1214        Technol *45*, 7640-7647. 10.1021/es2002259.

1215    47.    Ianora, A., Miralto, A., Poulet, S.A., Carotenuto, Y., Buttino, I., Romano, G., Casotti, R.,

1216        Pohnert, G., Wichard, T., Colucci-D'Amato, L., et al. (2004). Aldehyde suppression of copepod

1217        recruitment in blooms of a ubiquitous planktonic diatom. Nature *429*, 403-407.

1218        10.1038/nature02526.

1219    48.    Vardi, A., Formiggini, F., Casotti, R., De Martino, A., Ribalet, F., Miralto, A., and Bowler, C.

1220        (2006). A stress surveillance system based on calcium and nitric oxide in marine diatoms. Plos

1221        Biol *4*, 411-419. ARTN e60 10.1371/journal.pbio.0040060.

1222    49.    Lieber, M.R. (2008). The mechanism of human nonhomologous DNA end joining. J Biol Chem

1223        *283*, 1-5. 10.1074/jbc.R700039200.

1224    50.    Alves, I., Houle, A.A., Hussin, J.G., and Awadalla, P. (2017). The impact of recombination on

1225        human mutation load and disease. Philos Trans R Soc Lond B Biol Sci *372*.

1226        10.1098/rstb.2016.0465.

1227   51.   Symington, L.S., Rothstein, R., and Lisby, M. (2014). Mechanisms and Regulation of Mitotic

1228       Recombination in Saccharomyces cerevisiae. Genetics *198*, 795-835.

1229       10.1534/genetics.114.166140.

1230   52.   Abdullah, M.F.F., and Borts, R.H. (2001). Meiotic recombination frequencies are affected by

1231       nutritional states in Saccharomyces cerevisiae. P Natl Acad Sci USA *98*, 14524-14529. DOI

1232       10.1073/pnas.201529598.

1233   53.   Forche, A., Abbey, D., Pisithkul, T., Weinzierl, M.A., Ringstrom, T., Bruck, D., Petersen, K., and

1234       Berman, J. (2011). Stress Alters Rates and Types of Loss of Heterozygosity in Candida

1235       albicans. Mbio *2*. ARTN e00129-11 10.1128/mBio.00129-11.

1236   54.   Modliszewski, J.L., Wang, H.K., Albright, A.R., Lewis, S.M., Bennett, A.R., Huang, J.Y., Ma, H.,

1237       Wang, Y.X., and Copenhaver, G.P. (2018). Elevated temperature increases meiotic crossover

1238       frequency via the interfering (Type I) pathway in Arabidopsis thaliana. Plos Genet *14*. ARTN

1239       e1007384 10.1371/journal.pgen.1007384.

1240   55.   Lloyd, A., Morgan, C., Franklin, F.C.H., and Bomblies, K. (2018). Plasticity of Meiotic

1241       Recombination Rates in Response to Temperature in Arabidopsis. Genetics *208*, 1409-1420.

1242       10.1534/genetics.117.300588.

1243   56.   Lucht, J.M., Mauch-Mani, B., Steiner, H.Y., Metraux, J.P., Ryals, J., and Hohn, B. (2002).

1244       Pathogen stress increases somatic recombination frequency in Arabidopsis. Nature Genetics

1245       *30*, 311-314. 10.1038/ng846.

1246   57.   Stevison, L.S., Sefick, S., Rushton, C., and Graze, R.M. (2017). Recombination rate plasticity:

1247       revealing mechanisms by design. Philos T R Soc B *372*. ARTN 20160459

1248       10.1098/rstb.2016.0459.

1249   58.   Jackson, S., Nielsen, D.M., and Singh, N.D. (2015). Increased exposure to acute thermal stress

1250       is associated with a non-linear increase in recombination frequency and an independent

1251       linear decrease in fitness in Drosophila. Bmc Evol Biol *15*. ARTN 175 10.1186/s12862-015-

1252       0452-8.

1253    59.    Lim, J.G.Y., Stine, R.R.W., and Yanowitz, J.L. (2008). Domain-Specific Regulation of

1254           Recombination in Caenorhabditis elegans in Response to Temperature, Age and Sex.

1255           Genetics *180*, 715-726. 10.1534/genetics.108.090142.

1256    60.    Gusa, A., and Jinks-Robertson, S. (2019). Mitotic Recombination and Adaptive Genomic

1257           Changes in Human Pathogenic Fungi. Genes (Basel) *10*. 10.3390/genes10110901.

1258    61.    Krueger-Hadfield, S.A., Balestreri, C., Schroeder, J., Highfield, A., Helaouet, P., Allum, J.,

1259           Moate, R., Lohbeck, K.T., Miller, P.I., Riebesell, U., et al. (2014). Genotyping an Emiliania

1260           huxleyi (prymnesiophyceae) bloom event in the North Sea reveals evidence of asexual

1261           reproduction. Biogeosciences *11*, 5215-5234. 10.5194/bg-11-5215-2014.

1262    62.    Wright, S. (1929). Fisher's Theory of Dominance. The American Naturalist *63*, 274-279.

1263    63.    Wright, S. (1934). Physiological and Evolutionary Theories of Dominance. The American

1264           Naturalist *68*, 24-53.

1265    64.    Desai, M.M., Fisher, D.S., and Murray, A.W. (2007). The speed of evolution and maintenance

1266           of variation in asexual populations. Curr Biol *17*, 385-394. 10.1016/j.cub.2007.01.072.

1267    65.    Kao, K.C., and Sherlock, G. (2008). Molecular characterization of clonal interference during

1268           adaptive evolution in asexual populations of Saccharomyces cerevisiae. Nat Genet *40*, 1499-

1269           1504. 10.1038/ng.280.

1270    66.    Lang, G.I., Botstein, D., and Desai, M.M. (2011). Genetic variation and the fate of beneficial

1271           mutations in asexual populations. Genetics *188*, 647-661. 10.1534/genetics.111.128942.

1272    67.    Bajic, D., Vila, J.C.C., Blount, Z.D., and Sanchez, A. (2018). On the deformability of an

1273           empirical fitness landscape by microbial evolution. Proc Natl Acad Sci U S A *115*, 11286-

1274           11291. 10.1073/pnas.1808485115.

1275    68.    Lewontin, R.C. (1974). The genetic basis of evolutionary change (Columbia University Press).

1276    69.    Crow, J.F. (1970). Genetic Loads and the Cost of Natural Selection. In Mathematical Topics in

1277           Population Genetics, K.-i. Kojima, ed. (Springer Berlin Heidelberg), pp. 128-177. 10.1007/978-

1278           3-642-46244-3_5.

1279 70. Lu, S.N., Wang, J.Y., Chitsaz, F., Derbyshire, M.K., Geer, R.C., Gonzales, N.R., Gwadz, M.,

1280 Hurwitz, D.I., Marchler, G.H., Song, J.S., et al. (2020). CDD/SPARCLE: the conserved domain

1281 database in 2020. Nucleic Acids Research *48*, D265-D268. 10.1093/nar/gkz991.

1282 71. Osuna-Cruz, C.M., Bilcke, G., Vancaester, E., De Decker, S., Bones, A.M., Winge, P., Poulsen,

1283 N., Bulankova, P., Verhelst, B., Audoor, S., et al. (2020). The Seminavis robusta genome

1284 provides insights into the evolutionary adaptations of benthic diatoms. Nat Commun *11*,

1285 3320. 10.1038/s41467-020-17191-8.

1286 72. Sayers, E.W., Cavanaugh, M., Clark, K., Pruitt, K.D., Schoch, C.L., Sherry, S.T., and Karsch-

1287 Mizrachi, I. (2021). GenBank. Nucleic Acids Res *49*, D92-D96. 10.1093/nar/gkaa1023.

1288 73. Eren, A.M., Esen, O.C., Quince, C., Vineis, J.H., Morrison, H.G., Sogin, M.L., and Delmont, T.O.

1289 (2015). Anvi'o: an advanced analysis and visualization platformfor 'omics data. Peerj *3*. ARTN

1290 e1319 10.7717/peerj.1319.

1291 74. Team, R.C. (2017). R: A language and environment for statistical computing. R Foundation for

1292 Statistical Computing.

1293 75. Bushnell, B. BBMap. sourceforge.net/projects/bbmap/.

1294 76. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler

1295 transform. Bioinformatics *25*, 1754-1760. 10.1093/bioinformatics/btp324.

1296 77. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G.,

1297 Durbin, R., and Genome Project Data Processing, S. (2009). The Sequence Alignment/Map

1298 format and SAMtools. Bioinformatics *25*, 2078-2079. 10.1093/bioinformatics/btp352.

1299 78. Institute, B. (Accessed: 2019, version 2.6.0). Picard Tools. Broad Institute, GitHub repository.

1300 79. Poplin, R., Ruano-Rubio, V., DePristo, M.A., Fennell, T.J., Carneiro, M.O., Van der Auwera,

1301 G.A., Kling, D.E., Gauthier, L.D., Levy-Moonshine, A., Roazen, D., et al. (2018). Scaling

1302 accurate genetic variant discovery to tens of thousands of samples. 201178. 10.1101/201178

1303 %J bioRxiv.

1304  80.  Van der Auwera, G.A., Carneiro, M.O., Hartl, C., Poplin, R., Del Angel, G., Levy-Moonshine, A.,

1305        Jordan, T., Shakir, K., Roazen, D., Thibault, J., et al. (2013). From FastQ data to high

1306        confidence variant calls: the Genome Analysis Toolkit best practices pipeline. Curr Protoc

1307        Bioinformatics *43*, 11 10 11-11 10 33. 10.1002/0471250953.bi1110s43.

1308  81.  Smit, A., Hubley, R. (2008-2015). RepeatModeler Open-1.0. <http://www.repeatmasker.org>.

1309  82.  Smit, A., Hubley, R & Green, P (2013-2015). RepeatMasker Open-4.0.

1310        <http://www.repeatmasker.org>.

1311  83.  Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing

1312        genomic features. Bioinformatics *26*, 841-842. 10.1093/bioinformatics/btq033.

1313  84.  Lindenbaum, P. (2015). JVarkit: java-based utilities for Bioinformatics. figshare.

1314  85.  Sovic, I., Sikic, M., Wilm, A., Fenlon, S.N., Chen, S., and Nagarajan, N. (2016). Fast and

1315        sensitive mapping of nanopore sequencing reads with GraphMap. Nat Commun *7*. ARTN

1316        11307 10.1038/ncomms11307.

1317  86.  Krzywinski, M., Schein, J., Birol, I., Connors, J., Gascoyne, R., Horsman, D., Jones, S.J., and

1318        Marra, M.A. (2009). Circos: An information aesthetic for comparative genomics. Genome

1319        Research *19*, 1639-1645. 10.1101/gr.092759.109.

1320  87.  Team, R. (2019). RStudio: Integrated Development for R. RStudio.

1321  88.  Gel, B., and Serra, E. (2017). karyoploteR: an R/Bioconductor package to plot customizable

1322        genomes displaying arbitrary data. Bioinformatics *33*, 3088-3090.

1323        10.1093/bioinformatics/btx346.

1324  89.  Chaisson, M.J., and Tesler, G. (2012). Mapping single molecule sequencing reads using basic

1325        local alignment with successive refinement (BLASR): application and theory. Bmc

1326        Bioinformatics *13*. Artn 238 10.1186/1471-2105-13-238.

1327  90.  Rausch, T., Fritz, M.H.Y., Korbel, J.O., and Benes, V. (2019). Alfred: interactive multi-sample

1328        BAM alignment statistics, feature counting and feature annotation for long- and short-read

1329        sequencing. Bioinformatics *35*, 2489-2491. 10.1093/bioinformatics/bty1007.

1330  91.  Bowler, C., Allen, A.E., Badger, J.H., Grimwood, J., Jabbari, K., Kuo, A., Maheswari, U.,

1331      Martens, C., Maumus, F., Otillar, R.P., et al. (2008). The Phaeodactylum genome reveals the

1332      evolutionary history of diatom genomes. Nature *456*, 239-244. 10.1038/nature07410.

1333  92.  Pesant, S., Not, F., Picheral, M., Kandels-Lewis, S., Le Bescot, N., Gorsky, G., Iudicone, D.,

1334      Karsenti, E., Speich, S., Trouble, R., et al. (2015). Open science resources for the discovery

1335      and analysis of Tara Oceans data. Sci Data *2*, 150023. 10.1038/sdata.2015.23.

1336  93.  Delmont, T.O., and Eren, A.M. (2018). Linking pangenomes and metagenomes: the

1337      Prochlorococcus metapangenome. Peerj *6*, e4320. 10.7717/peerj.4320.

1338  94.  Delmont, T.O., Kiefl, E., Kilinc, O., Esen, O.C., Uysal, I., Rappe, M.S., Giovannoni, S., and Eren,

1339      A.M. (2019). Single-amino acid variants reveal evolutionary processes that shape the

1340      biogeography of a global SAR11 subclade. Elife *8*. ARTN e46497 10.7554/eLife.46497.

1341  95.  Amarasinghe, S.L., Su, S., Dong, X.Y., Zappia, L., Ritchie, M.E., and Gouil, Q. (2020).

1342      Opportunities and challenges in long-read sequencing data analysis. Genome Biol *21*. ARTN

1343      30 10.1186/s13059-020-1935-5.

1344  96.  Williams, R., Peisajovich, S.G., Miller, O.J., Magdassi, S., Tawfik, D.S., and Griffiths, A.D.

1345      (2006). Amplification of complex gene libraries by emulsion PCR. Nat Methods *3*, 545-550.

1346      10.1038/nmeth896.

1347  97.  Kalle, E., Kubista, M., and Rensing, C. (2014). Multi-template polymerase chain reaction.

1348      Biomol Detect Quantif *2*, 11-29. 10.1016/j.bdq.2014.11.002.

1349  98.  Madeira, F., Park, Y.M., Lee, J., Buso, N., Gur, T., Madhusoodanan, N., Basutkar, P., Tivey,

1350      A.R.N., Potter, S.C., Finn, R.D., and Lopez, R. (2019). The EMBL-EBI search and sequence

1351      analysis tools APIs in 2019. Nucleic Acids Res *47*, W636-W641. 10.1093/nar/gkz268.

1352  99.  Hellemans, J., Mortier, G., De Paepe, A., Speleman, F., and Vandesompele, J. (2007). qBase

1353      relative quantification framework and software for management and automated analysis of

1354      real-time quantitative PCR data. Genome Biol *8*, R19. 10.1186/gb-2007-8-2-r19.